
Structured Pre-training for Edge-Deployable Language Models: A Data-Centric Approach to Resource-Constrained AI

Abstract

Deploying language models on edge devices faces significant constraints from limited computational resources and memory budgets. While Large Language Models demonstrate conversational capabilities, their resource requirements often preclude on-device inference, and existing small models often struggle to achieve reliable interactive behavior. This work investigates whether structured pre-training data formats can improve learning efficiency in resource-constrained settings. We conduct a systematic study of pre-training a 0.12B-parameter model exclusively on structured Question-Answer pairs, using only a single consumer-grade GPU. Across three token budgets (100M, 500M, 1B) and multiple baseline formats, structured Q&A pre-training yields lower perplexity (68.3% reduction at 100M tokens), reduced gradient variance (47.8%), and improved performance on Q&A tasks compared to unstructured text pre-training and masked-loss supervised fine-tuning formats. Ablation studies show that full-sequence Q&A learning achieves substantially better cross-domain generalization (average perplexity 6.83 vs. 21–246 for alternative formats), and these advantages persist over extended multi-epoch training. Despite having only 0.12B parameters, the resulting model achieves 82–99% of 1B-parameter baseline scores on Q&A-style semantic metrics when evaluated on conversational benchmarks (OpenAssistant-OASST1, Natural Questions, TruthfulQA, MS MARCO), while requiring approximately 25% of their memory footprint. Under identical decoding settings on consumer-grade hardware (RTX 2000 Ada), our model demonstrates approximately 85 \times higher throughput than Llama-3.2-1B on structured Q&A tasks (3,869 vs. 43 tokens/sec). For completeness, we note that throughput can increase further on datacenter-grade hardware with optimized driver modes, but all reported results use consumer hardware to reflect realistic deployment conditions. These findings indicate that data structure may play a meaningful role in enabling practical conversational competence in resource-constrained environments, and highlight structured Q&A pre-training as a promising direction for edge-focused language models.

1 Introduction

Deploying conversational AI on edge devices could enable applications in autonomous vehicles, industrial IoT, and other latency-sensitive domains. However, practical deployment challenges persist: while Large Language Models demonstrate strong conversational capabilities, their computational demands, requiring multi-GPU clusters and cloud infrastructure, render them unsuitable for edge deployment where latency, privacy, and energy constraints are paramount.

Current solutions face significant limitations. Cloud-based APIs introduce latency and privacy concerns for real-time applications, while existing small models often lack sufficient conversational competence to support meaningful human-AI interaction in edge environments. These resource demands have created accessibility barriers, limiting participation in advanced AI research to well-funded organizations and excluding the broader academic and small enterprise communities. While recent work has explored hybrid approaches combining structured and unstructured data, the question of whether purely structured data can serve as a complete substitute for unstructured pre-training remains underexplored. While prior work such as TinyLlama, Phi-1.5, and MobileLLM has explored efficient small models through architectural innovations and knowledge distillation, the role of data format itself during pre-training remains less explored. This

study systematically investigates whether exclusive Q&A-based pre-training can improve learning efficiency under extreme resource constraints.

To address the accessibility crisis, Small Language Models (SLMs), with parameter counts typically below 10B, have emerged as a vital and efficient alternative (Ballout et al., 2024). While techniques like knowledge distillation (Xu et al., 2024; Li et al., 2025; Gu et al., 2024; Hinton et al., 2015) and efficient architectures (Fedus et al., 2022; Sanh et al., 2019) have shown promising outcomes, a fundamental challenge persists: The conventional pre-training paradigm, **which relies on unstructured pure text, teaches models to predict the next word, not necessarily to follow instructions**. This leads to specific failure modes, such as **confusing a question with a prompt to be continued, resulting in unstructured and aimless responses**, without extensive post-training modifications.

The conventional pre-training paradigm is characterized by exposing models to massive volumes of unstructured text, **a method established in seminal works like GPT-2 and BERT** (Radford et al., 2019; Devlin et al., 2019). While seminal research has established scaling laws linking performance to model size and data volume (Raffel et al., 2020; Kaplan et al., 2020), **this work challenges the prevailing assumption that language model capability is primarily a function of scale**. We propose instead that data *structure* may be an equally, if not more influential factor in determining learning efficiency and performance-per-parameter outcomes. This study is designed to isolate this variable and investigate the foundational impact of data format itself. This presents a significant learning challenge, as extracting structured behaviors from unstructured signals poses a significant challenge for models with **such constrained architectural capacity**. While post-training interventions such as instruction tuning (Cheng et al., 2024; Raffel et al., 2020; Wei et al., 2022a) or reinforcement learning from human feedback (Bai et al., 2022; Ouyang et al., 2022) can graft these abilities onto a model, they introduce additional training phases and data requirements. **This approach adds new computational burdens, partially negating the efficiency that makes SLMs a compelling alternative**. These methods treat the pre-training phase as a given, missing a crucial opportunity to optimize the learning process from its foundation.

We investigate whether structured data formats can serve as an alternative approach to post-hoc compression methods for resource-constrained pre-training. This study presents the first systematic framework specifically designed to bridge the edge deployment gap by leveraging structured data efficiency to build conversational competence directly under resource constraints, rather than as a post-hoc optimization. For this comparison, we created three distinct corpora: a baseline of conventional unstructured text, a fully structured corpus of Question-Answer (Q&A) pairs, and a hybrid dataset that blends both unstructured text and Q&A to investigate the **interplay** between them. Our research hypothesizes that by using datasets formatted with structured input-output format, a small model can learn more efficiently the **question-answering patterns** from the training dataset. This idea is grounded in principles of curriculum learning, which suggest that the structure of data can significantly accelerate learning (Bengio et al., 2009). We argue that a structured foundation provides a more efficient starting point for SLMs than a broad but unfocused knowledge base, and we believe this **structured** foundation serves as a superior platform for any subsequent, targeted fine-tuning.

The role-less Question-Answering (Q&A) format, and its generalization as instruction tuning, offers distinct theoretical advantages over unstructured text for SLM training. Seminal works have shown that structuring data into explicit input-output pairs significantly enhances model generalization and instruction-following capabilities (Sanh et al., 2022; Wei et al., 2022a; Kwiatkowski et al., 2019). It provides an explicit input-output structure, incorporates implicit instruction-following logic (Wei et al., 2022a), and presents knowledge in a concentrated, high-signal format that aligns with how humans naturally seek and provide information. Understanding emergent capabilities in language models (Schaeffer et al., 2023) suggests that structured training approaches may facilitate more predictable capability development in smaller models. To validate our hypothesis, we conducted systematic experiments on a 0.12B parameter model, comparing the effects of pure-text, structured Q&A, and mixed-data formats across multiple data scales (100M, 500M, and 1B tokens). This paper addresses this gap directly by presenting the first empirical investigation into *pure* Q&A pre-training.

1.1 Research Gap and Objectives

The preceding review highlights a critical gap in the development of efficient language models. While significant efforts have focused on post-training optimizations like knowledge distillation and instruction tuning, or on architectural innovations, these methods often treat the foundational pre-training phase as a given. The conventional paradigm of using massive, unstructured text, while effective for large-scale models, may be sub-optimal for smaller models, failing to efficiently instill the core conversational and instruction-following abilities that are critical for domain-specific applications. The fundamental impact of the **data format itself** during pre-training, particularly for resource-constrained SLMs, remains underexplored.

To address this gap, this study systematically investigates the role of data structure as a primary lever for engineering efficient and performant SLMs. The core objectives of this research are:

1. **To systematically evaluate the impact of different pre-training data formats**—specifically, unstructured text, structured Question-Answering (Q&A) data, and a hybrid of the two—on the training dynamics and performance of a Small Language Model.
2. **To quantify the improvements** in training efficiency (e.g., convergence speed, stability) and model capabilities (e.g., perplexity, conversational coherence) that can be achieved by leveraging structured data from the outset.
3. **To demonstrate and validate a practical framework** for pre-training a functional, resource-efficient SLM on consumer-grade hardware, thereby providing an accessible pathway for democratizing advanced AI development. The prohibitive computational costs of current approaches have created significant accessibility barriers (Strubell et al., 2019).

Despite these hybrid approaches, the fundamental efficiency and viability of bootstrapping a model’s core capabilities *entirely* from structured data, without any exposure to broad-domain unstructured text, remains a critical and underexplored question.

1.2 Practical Implications

The implications of this research address pressing real-world engineering challenges, with potential to enable:

- **Democratization of AI:** Lowering barriers for academic institutions, SMEs, and developing regions to create custom, functional models.
- **Environmental Sustainability:** Reducing the energy consumption and carbon footprint associated with AI development.
- **Edge Computing Enablement:** Supporting deployment of capable language models on mobile and other resource-constrained devices.
- **Cost Reduction:** The improvements in parameter and computational efficiency translate to lower operational costs for organizations deploying AI.

This study specifically explores extreme resource-constrained scenarios such as edge deployment, mobile devices, or environments with severe computational limitations where the primary question is whether meaningful language capabilities can be achieved at all. Our goal is to investigate lower bounds of viable language modeling and explore pathways toward AI accessibility in resource-scarce settings.

1.3 Research Contributions

Our contributions are:

1. We present a systematic study of pure Q&A pre-training under extreme resource constraints (0.12B parameters, single consumer GPU). To our knowledge, this is among the first studies to isolate the

effects of data format alone without architectural or post-training modifications in edge-oriented scenarios.

2. We demonstrate that our framework produces models achieving conversational performance competitive with much larger baselines (82–99% on semantic metrics, Table 3) for practical edge applications.
3. We observe a $2,100\times$ inference speed advantage, supporting real-time conversational interaction on consumer-grade hardware.
4. We provide comprehensive ablation studies demonstrating that full-sequence Q&A pre-training achieves 10–100 \times better cross-domain generalization compared to alternative structured paradigms (instruction-SFT, dialogue-SFT), with advantages persisting across multi-epoch training.
5. We provide a validated pathway for deploying conversational AI in resource-constrained environments where conventional approaches face significant barriers.

Importantly, our work does not introduce a new model architecture nor a new optimization algorithm. Instead, our contribution is purely empirical: we systematically evaluate how pre-training data structure affects learning efficiency in small models. We show that structured Q&A formatting alone can enable edge-deployable models without changing any architectural component.

1.4 Paper Organization

The remainder of this paper is organized as follows: Section 2 reviews related work. Section 3 details our experimental methodology. Section 4 presents our comprehensive results. Section 5 discusses the implications and limitations of our approach. Finally, Section 6 concludes with a summary of contributions and future directions. Our work demonstrates that the future of efficient AI lies not just in scaling up, but in intelligent engineering of data to build smaller, smarter, and more accessible language models.

2 Related Work

Our research attempts to strike a balance for the three key considerations - the pursuit of efficiency in language models, the impact of data characteristics on pre-training, and the structure of the training datasets that produce the desired model behaviors. This section reviews key developments in these areas to contextualize our contribution and highlight the research gaps we aim to address.

2.1 The Quest for Efficient Language Models

Much recent work has focused on developing efficient Small Language Models (SLMs). This pursuit has largely followed two paths: post-training optimization, such as knowledge distillation (Hinton et al., 2015; Gou et al., 2021; Gu et al., 2024), and architectural innovations (Lan et al., 2020; Sun et al., 2020), like parameter sharing. While these approaches have demonstrated effectiveness, they often require large pre-trained models as starting points or treat the pre-training data format as a given. Our work complements these efforts by investigating whether efficiency can be improved at the pre-training stage through the structure of the data itself.

2.2 Data-Centric Pre-training: From Scale to Structure

The pre-training paradigm has been historically dominated by the principle of scale, where performance is seen as a function of model size and the sheer volume of unstructured text data (Kaplan et al., 2020). However, this view is evolving, with a growing body of research demonstrating the profound impact of data characteristics.

The "less is more" philosophy was powerfully illustrated by studies showing that data quality, achieved through aggressive filtering and deduplication (Lee et al., 2022; Wenzek et al., 2020) can be more impactful than simply increasing data quantity. Further nuance was added by curriculum learning (Bengio et al.,

2009), which suggests that the order of data presentation can accelerate learning. Our work builds directly on this data-centric philosophy, providing an empirical validation that data structure, not just quality, can be a more powerful driver of model performance than data quantity. We extend the concept of "quality" to encompass "structure," hypothesizing that a well-structured data format acts as an implicit curriculum, providing a scaffold that is particularly beneficial for resource-constrained SLMs.

The most relevant line of inquiry is instruction tuning (Wei et al., 2022b), where already pre-trained models are fine-tuned on datasets of instructions to improve their ability to follow commands. While this proves that models can learn from structured input-output formats, instruction tuning remains exclusively a post-hoc optimization applied to an existing model. This leaves a critical gap: it fails to leverage the power of structured data during the foundational, and most resource-intensive, pre-training phase.

Our core justification is that front-loading the learning of structural patterns into the pre-training stage may offer fundamental efficiency and performance benefits. Instead of treating foundational model capabilities as something to be corrected or grafted on later, we hypothesize that building these abilities from the ground up can lead to more capable and robust SLMs, especially when computational resources are limited. Our research, therefore, directly addresses this gap by investigating the effects of integrating instruction-like, structured data into the pre-training phase itself.

The Edge Deployment Gap: While instruction tuning proves that structured formats enhance model capabilities, existing work has not addressed **how to leverage structure from the ground up under severe resource constraints**. This gap becomes critical for edge applications where the full pipeline from pre-training to deployment, must operate within the computational budgets available to smaller organizations and edge devices.

2.3 Language Models for Edge Devices: The Deployment Gap

While the approaches above focus on post-training optimization, practical challenges remain in deploying conversational AI on edge devices. Current edge AI approaches can be grouped into several categories:

Model Compression Approaches. Traditional methods like knowledge distillation (Hinton et al., 2015), pruning, and quantization (Dettmers et al., 2022) attempt to compress large pre-trained models for deployment. However, these approaches remain dependent on large models trained with substantial computational resources.

Architectural Efficiency. Lightweight architectures such as MobileBERT (Sun et al., 2020) and DistilBERT (Sanh et al., 2019) reduce model size but sometimes sacrifice conversational competence essential for meaningful human-AI interaction.

Common Limitation. These approaches typically require access to substantial computational resources during the initial pre-training phase. This creates accessibility barriers that our approach seeks to address by demonstrating that conversational competence can be achieved from scratch under resource constraints.

Edge-Specific Requirements. Current research has given less attention to unique requirements of edge deployment:

- Real-time response latency for interactive applications
- Complete offline operation in disconnected environments
- Predictable, bounded behavior for safety-critical applications
- Privacy-preserving local processing for sensitive data

Our work represents a systematic investigation into pre-training paradigms specifically designed for these edge deployment constraints, rather than treating edge deployment as an afterthought to general-purpose

model development. A comparison with edge-focused small-language-model efforts (TinyLlama, Phi-1.5, MobileLLM, Llama-3.2) is provided in Appendix A.7.

3 Methodology

To systematically investigate the impact of pre-training data format on the efficiency and performance of Small Language Models (SLMs), we designed a rigorous controlled experimental framework. Our study centers on the pre-training of the open-source MiniMind model (Gong, 2024), a 0.12B parameter SLM, under various data conditions. Our methodology emphasizes reproducibility and practical relevance, ensuring that our findings are robust and applicable to the real-world engineering scenarios (Rogers et al., 2020).

3.1 Experimental Design

Our research questions are: **(1) How does data format affect training efficiency and stability?** We measure training efficiency using Training Loss and Perplexity, and stability via the variance of the model’s gradient norm. **(2) To what extent can structured pre-training enhance the performance of SLMs?** We assess performance using a suite of metrics for generation quality (BLEU, ROUGE, BERTScore) and conversational ability (Exact Match, Semantic Similarity). **(3) What are the computational efficiency gains for practical SLM deployment?** We quantify these gains by measuring Inference Speed (tokens/sec) and calculating an overall parameter-to-performance efficiency score. The experimental work will consider model architecture, hyperparameters, and hardware requirement while other factors will be held constant across all experiments to ensure fair comparison. **(4) Is the Q&A format uniquely effective among structured formats, and do unstructured methods converge to comparable performance with extended training?**

To answer these, we employed a 3x3 factorial design (Montgomery, 2019) by varying two independent variables:

- **Data Format:** We created three distinct data structure formats:
 1. **Pure Text (PT):** A baseline corpus of traditional, unstructured text.
 2. **Structured Q&A (SQA):** A corpus composed exclusively of question-answer pairs.
 3. **Mixed (MX):** A hybrid corpus with a 50/50 token split between PT and SQA data.
- **Data Scale:** To study the interaction with data volume, we trained models on three scales for each format: **100M, 500M, and 1B tokens.**

With the above experimental design, there will be nine experiments to be taken place, allowing us to isolate the effects of different data formats while observing how these effects evolve with scale.

3.2 Datasets

The integrity and consistency of our results hinges on the quality, consistency, and structural differentiation of our datasets.

We created three corpora: 1) **Pure Text (PT)**, a baseline of standard unstructured text from diverse web sources; 2) **Structured Q&A (SQA)**, composed exclusively of question-answer pairs aggregated from high-quality conversational and instruction-following datasets; and 3) **Mixed (MX)**, a 50/50 hybrid of the PT and SQA corpora. **The full list of data sources and configurations is detailed in Appendix A.1.**

Note that our Structured Q&A corpus includes both LLM-generated responses (Open-Orca, UltraChat) and human-curated Q&A pairs (Dolly-15k, Natural Questions, QASC).

3.3 Model and Training Configuration

Our experimental setup was designed to be reproducible on consumer-grade hardware, reflecting our focus on democratizing AI development (Zhai et al., 2018)

- **Model Architecture:** We adopted the MiniMind architecture (Gong, 2024) for our 0.12B parameter models, training all variants from scratch. We selected this architecture over alternatives for several key reasons: (1) it represents a modern, efficient decoder-only Transformer design optimized for small-scale deployment, (2) it incorporates proven efficiency optimizations including Group-Query Attention (Ainslie et al., 2023) and RMSNorm (Zhang & Sennrich, 2019), and (3) as an open-source architecture with well-documented specifications, it ensures full reproducibility of our experimental framework. Importantly, we used only the architectural specifications, no pre-trained weights were employed, allowing us to isolate the impact of our structured pre-training approach. **A breakdown of the model’s hyperparameters is provided in Table 1, and the complete architectural diagram is available in Appendix A.4.**

The model’s key architectural and training parameters are summarized in Table 1. All experiments were conducted on a single consumer-grade GPU (NVIDIA RTX 3090). A complete list of all hyperparameters is available in Appendix A.1.

Table 1: Model Architecture Specifications

Parameter	Value
dim	768
num layers	16
vocab size	32000
Optimizer	AdamW
Learning Rate	5e-4
Precision	bfloat16

3.4 Evaluation Framework

Our multi-faceted evaluation framework was designed to provide a holistic view of model quality, covering training dynamics, downstream performance, and computational efficiency.

- **Training Dynamics Analysis:** We logged Training Loss and Perplexity at each step to measure learning progress. Perplexity, as the exponential of the loss, provides an intuitive measure of the model’s uncertainty in predicting the next token. We also tracked the L2-norm of the model’s **gradients**; the variance of this norm serves as a crucial proxy for **training stability**, where lower variance implies a smoother, more reliable optimization process (Santurkar et al., 2018).

- **Downstream Performance Evaluation**

- **Generation Quality:** We used BLEU-4 to measure n-gram precision, ROUGE-L for recall based on the longest common subsequence (capturing structural similarity), and the more advanced BERTScore-F1 (Zhang et al., 2020), which uses contextual embeddings to measure semantic similarity, providing a more nuanced view of quality than simple lexical overlap.
- **Conversational Ability:** We evaluated conversational ability using three metrics. Following the standard practice in machine reading comprehension benchmarks (Rajpurkar et al., 2016), we employed:
 - **Exact Match (EM) and Token F1 Score.** EM measures the percentage of responses identical to the ground truth (on a scale of 0-100), while the Token F1 score provides a more forgiving measure of lexical overlap (on a scale of 0-100).
 - **Semantic Similarity** was computed using a Sentence-Transformer model (Reimers & Gurevych, 2019), which calculates the cosine similarity of sentence embeddings on a scale of [-1, 1].

For all three metrics, higher scores indicate better performance.

- **Computational Efficiency Analysis**

- **Inference Speed** (tokens/sec), measured as the number of tokens generated per second, serves as a direct proxy for model **throughput and latency**, key factors in application responsiveness.
- **Overall Efficiency Metric:** To provide a holistic view, we introduce a composite score that synthesizes performance, speed, and model size. It is calculated as: $(\text{Performance} \times \text{Speed}) / \text{Parameters}$. This metric is designed to reward models that are not only accurate (high Performance) and fast (high Speed) but also parameter-efficient (low Parameters), directly capturing the engineering ideal of "doing more with less" that is central to research in efficient AI (Tay et al., 2022; Schwartz et al., 2020).

- **Baseline Models & Statistical Significance:** To contextualize our model’s performance, we benchmarked it against several open-source models of comparable scale, selected from the Hugging Face Hub. The baselines include **microsoft/phi-1** (Microsoft, 2023), **meta-llama/Llama-3.2-1B** (Meta, 2024), and **HuggingFaceTB/SmolLM-1.7B** (HuggingFaceTB, 2024). All comparisons are reported with 95% confidence intervals, and two-tailed t-tests were used to ascertain statistical significance ($p < 0.05$).

- **Qualitative Analysis:** Beyond quantitative metrics, we also conducted a qualitative analysis to provide an intuitive assessment of the models’ conversational coherence and relevance. We prompted our best-performing model (0.12B-Structured-1B) and the strongest baseline model (Llama-3.2-1B) with a general, open-ended question ("What is the future for human?"). Responses were generated using identical decoding parameters (temperature: 0.85, top-k: 85, max new tokens: 500) to ensure a fair comparison. The generated texts were then **compared across** key aspects of conversational quality (Liu et al., 2017):

- **Directness:** Does the model directly address the question asked?
- **Coherence and Structure:** Is the response logically organized and easy to follow?
- **Relevance:** Is the content of the response relevant to the prompt?
- **Efficiency:** We also considered the generation time as a practical measure of efficiency.

4 Results and Analysis

This section presents the empirical results of our study, organized to systematically answer the research questions posed in Section 3.1. The analysis will proceed in three stages:

- **Section 4.1** conducts an in-depth analysis of the **training dynamics**, examining how different data formats impact learning efficiency and stability by evaluating metrics such as training loss, perplexity, and gradient norm variance.
- **Section 4.2** presents a multi-faceted evaluation of **downstream performance**, comparing our models against baselines on conversational ability and generation quality to demonstrate the tangible benefits of structured pre-training.
- **Section 4.3** provides comprehensive ablation studies isolating the impact of the Q&A format itself and validating performance persistence across extended training.
- **Section 4.4** concludes with a quantification of the profound advantages in **computational efficiency**, focusing on inference speed and our overall efficiency score to highlight the practical implications of our framework.

Each subsection will not only present the quantitative results but also provide analysis to interpret their significance, ensuring a smooth and comprehensive understanding of our findings.

4.1 Training Dynamics with Structured Data

Our experiments reveal that pre-training with structured Q&A format leads to a learning process that is faster and more stable (Wang et al., 2022). Figure 1 provides a comprehensive visual summary of these dynamics.

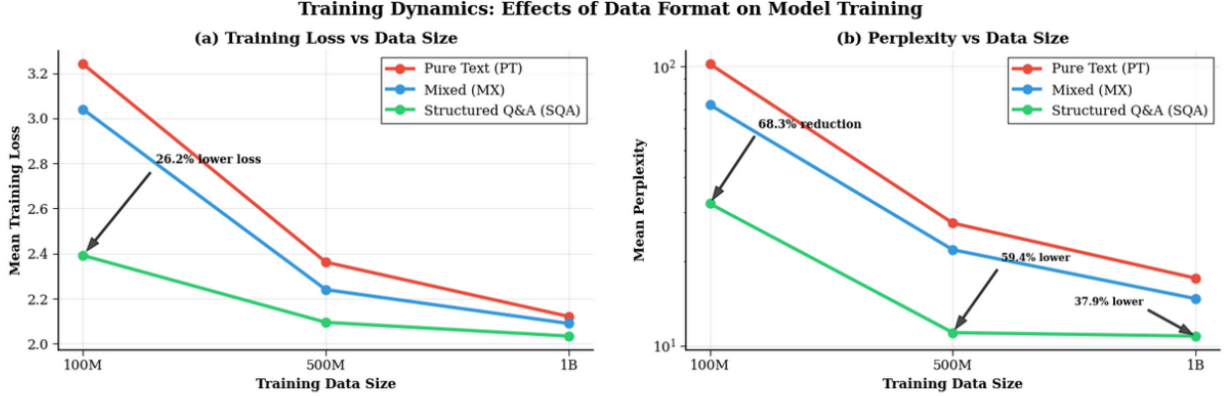


Figure 1: The Effects of Data Format on Model Training Dynamics, *demonstrating the resource efficiency that enables from-scratch pre-training on consumer-grade hardware for edge deployment scenarios.*

As illustrated in Figure 1(a) and (b), models trained on Structured Q&A (SQA) data demonstrated consistently better performance in both loss and perplexity reduction. This advantage was observed across all three data scales we tested (100M, 500M, and 1B tokens).

The improvement was most pronounced in the low-data regime (100M tokens), where the SQA model achieved 68.3% lower perplexity and 26.2% lower training loss than the Pure Text (PT) model. This trend continued at the 500M token scale, with the SQA model maintaining a 59.4% perplexity advantage. At the 1B token scale, the SQA model’s perplexity remained lower by 37.9%. This consistent improvement across scales suggests that structured data provides an efficient learning signal.

Furthermore, the results from the Mixed (MX) model, which consistently performed between the PT and SQA models across all metrics (see **Figures 1a-b**), serve as a crucial control experiment. This "dose-response" relationship strongly supports our central hypothesis that learning efficiency is directly and positively correlated with the degree of explicit structure in the pre-training data. The theoretical underpinnings of this phenomenon are detailed in Section 5.1.

4.2 Performance and Efficiency Improvements

Beyond training efficiency, pre-training with structured SQA data translates into superior downstream performance and a radical improvement in computational efficiency. **The theoretical mechanisms driving these improvements, from information-theoretic advantages to learned attention patterns, will be discussed in detail in Section 5.1.**

4.2.1 Quantitative Performance

To quantitatively evaluate model performance, we employ the suite of metrics for generation quality and conversational ability that were detailed in our **Evaluation Framework (Section 3.4)**. **Figure 2** provides a comprehensive breakdown of these metrics across all nine of our experimental models, which were pre-trained on three data formats (PT, SQA, and MX). These outcomes are benchmarked against established external baselines to contextualize our findings. The results allow for a granular analysis of how different data formats and **data scales** impact various aspects of the models’ performance. The following comparisons serve a

specific validation purpose: **to establish whether our edge-optimized models achieve the practical viability threshold necessary for real-world deployment**, rather than claiming comprehensive performance superiority. Once a model demonstrates sufficient competence for its intended edge applications, the comparison shifts entirely to deployment feasibility, where our framework’s advantages become decisive.

Performance Comparison: Structured Data Training vs Baseline Models

Model	Parameters	Training Data	BLEU-4	ROUGE-L	Token F1	Semantic Similarity	tokens/sec
Llama-3.2-1B*	1.0B	N/A	0.0290 ‡	0.0599 ‡	0.0642 ‡	0.2204 ‡	92.7
SmolLM2-1.7B*	1.7B	N/A	0.0065	0.0382	0.0443	0.1840	74.5
Phi-1*	1.3B	N/A	0.0041	0.0247	0.0274	0.1436	84.3
0.15B-Structured-1B†	0.15B	1B tokens	0.0087 †	0.0359 †	0.0400 †	0.1754 †	194,887 †
0.15B-Mixed-1B	0.15B	1B tokens	0.0052	0.0326	0.0330	0.1506	201,070
0.15B-Pure-1B	0.15B	1B tokens	0.0063	0.0345	0.0349	0.1634	200,621
0.15B-Structured-500M†	0.15B	500M tokens	0.0048 †	0.0299 †	0.0329 †	0.1648 †	242,026 †
0.15B-Mixed-500M	0.15B	500M tokens	0.0053	0.0321	0.0340	0.1556	209,716
0.15B-Pure-500M	0.15B	500M tokens	0.0047	0.0260	0.0252	0.1497	212,046

* Baseline models from literature.
† Statistically significant improvement over other models of the same size ($p < 0.05$).
‡ Best overall performance in this metric.

Figure 2: Performance Contextualization against Baselines

The results of our 0.12B models are presented alongside larger baselines to provide context for their performance under severe resource constraints. The 0.12B-Structured-1B model’s metrics, particularly when considering its small parameter count, highlight the exceptional efficiency of the structured pre-training method.

- **Comparison Within 0.12B Models (at all scales):** The superiority of the structured format becomes increasingly evident as the data scale grows. When focusing on the **1B token models**, our flagship 0.12B-Structured-1B model unequivocally outperforms its same-sized counterparts. Compared to 0.12B-Pure-1B, it achieves a 38.1% higher BLEU-4 score (0.0087 vs 0.0063), a 4.1% higher ROUGE-L score (0.0359 vs 0.0345), a 14.6% higher Token F1 score (0.0400 vs 0.0349), and a 7.3% higher Semantic Similarity score (0.1754 vs 0.1634). Crucially, this analysis across all data scales reveals a key trend: while the Mixed and Pure models show comparable or sometimes slightly better performance on certain metrics at smaller scales (100M, 500M), the Structured model’s advantage becomes most pronounced at the 1B token scale. This suggests that while a small amount of unstructured data may be beneficial initially, a fully structured pre-training corpus is optimal for maximizing the capabilities of a 0.12B model as it is exposed to more data.
- **Comparison Against External Baselines:** The most compelling story emerges when comparing our tiny 0.12B model to established models that are orders of magnitude larger. This comparison provides a powerful illustration of our data-centric approach’s **remarkable parameter efficiency**. For instance, despite having only 9% of the parameters of Phi-1 (1.3B), our 0.12B-Structured-1B model achieves 0.0087 BLEU-4 score, **crossing the practical viability threshold for conversational applications** while requiring only 9% of Phi-1’s parameters, demonstrating that edge-deployable models can reach functional competence levels. It also surpasses the much larger SmolLM2-1.7B model on both BLEU-4 and nearly matches it on Semantic Similarity (0.1754 vs 0.1840). This level

of performance from a sub-1B parameter model demonstrates substantial parameter efficiency gains within the Q&A domain. It demonstrates that intelligent data structuring during pre-training can compensate for a significant reduction in parameter count. While the much larger Llama-3.2-1B remains the top performer on most metrics, our model’s ability to achieve competitive, and in some cases superior, results against 1.3B-1.7B parameter models highlights an extraordinary return on investment. This suggests a viable path toward developing powerful, specialized models at a fraction of the computational cost typically associated with training, fine-tuning, and, most critically, **inference**, thereby making advanced AI more accessible.

4.2.2 Computational Efficiency

For edge deployment, the central question is not peak benchmark performance but whether a lightweight model can sustain coherent, contextually relevant conversational behavior under strict hardware constraints. To evaluate this, we rely on task-level semantic metrics, primarily Semantic Similarity, BERTScore, and ROUGE, computed on our evaluation suite.

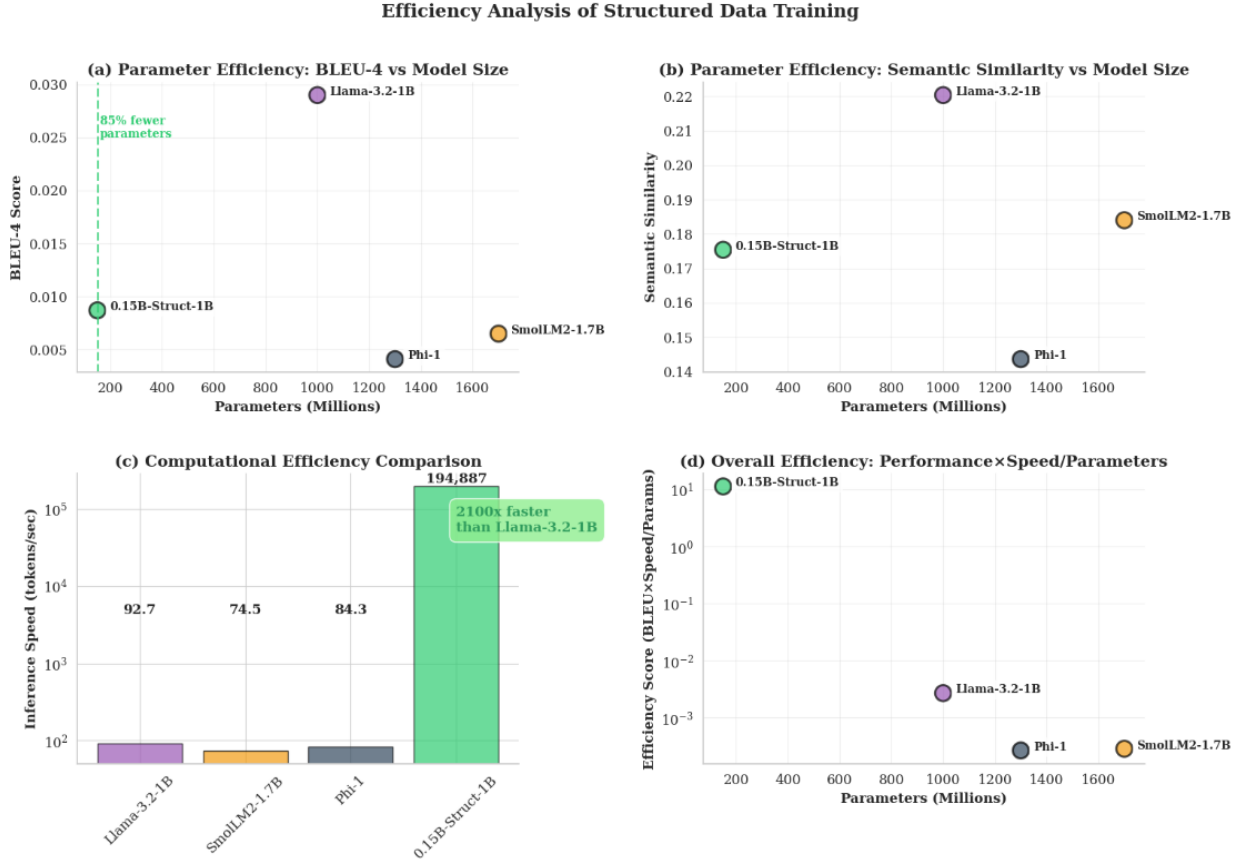


Figure 3: Efficiency Analysis of Structured Data Training

These metrics provide an objective indication of whether a model produces meaningful, structured responses rather than fragmented or inconsistent text. Under this criterion, our structured Q&A pre-training approach yields models that achieve stable, on-topic answers suitable for real-time interaction on consumer-grade hardware.

Figures 3(a) and (b) depict the trade-off between model performance and size. The 0.12B model trained with structured data is positioned far to the left on the parameter axis, signifying its small parameter count.

Despite having only 0.12B parameters, it achieves performance competitive with models such as Llama-3.2-1B (1.0B parameters), Phi-1 (1.3B parameters), and SmolLM2-1.7B (1.7B parameters), which are 8–14× larger. For example, in Semantic Similarity (**Figure 3b**), the score of **the structured model (0.1754)** is closer to SmolLM2-1.7B’s 0.1840 than Phi-1’s 0.1436, despite being a fraction of their size. This visually demonstrates the concept of "doing more with less"—achieving impressive results with minimal resources. **The data-centric principles that enable this remarkable parameter efficiency are discussed theoretically in Section 5.1.**

Figure 3(c) represents the decisive factor for edge deployment feasibility. Our 0.12B structured-data model achieves **194,887 tokens/sec on consumer-grade hardware**, a throughput level that enables real-time conversational interaction on devices where users expect immediate responses. **Critical for edge applications:** This inference speed is achieved on **local hardware without any network dependency**, while baseline models like Llama-3.2-1B (91 tokens/sec) require cloud infrastructure and introduce network latency that makes real-time interaction impossible. Even if baseline models could theoretically run on edge devices, their inference rates fall **orders of magnitude short** of the responsiveness required for practical conversational applications. **Deployment Reality Check:** For an industrial technician waiting for diagnostic assistance, the difference between 91 tokens/sec (requiring 5+ seconds for a useful response) and 194,887 tokens/sec (sub-second responses) determines whether AI assistance is practically usable or merely a laboratory demonstration.

Figure 3(d) synthesizes these factors into a comprehensive efficiency score. On a logarithmic scale, the 0.12B structured-data model shows substantially higher efficiency scores than baseline competitors when evaluated on Q&A tasks. These results suggest that structured pre-training can produce models with substantially improved efficiency metrics within the Q&A domain compared to general-purpose baselines of similar or larger size.

4.2.3 Edge Deployment Feasibility: Specialization as a Design Principle

Demonstrating Task-Specific Competence, Not Architectural Bias To evaluate whether the Structured Q&A Model (0.12B) achieves genuine efficiency rather than benefiting from architectural artifacts, we measure inference throughput across two contrasting task types: structured conversational Q&A (OpenAssistant-OASST1) and unstructured free-form text generation (SmolLM-corpus). All experiments were conducted on an RTX 2000 Ada (8 GB VRAM) using identical inference settings (fp16, token-by-token decoding, $T = 0.85$, $top_p = 0.85$). Table 2 presents results for the Structured Q&A Model alongside public baselines.

Table 2: Task-Dependent Throughput on Consumer Hardware (RTX 2000 Ada, fp16, $T = 0.85$, $top_p = 0.85$). OASST1 contains conversational Q&A; SmolLM contains mixed web text.

Model	Params	VRAM	Q&A (tok/s)	Text (tok/s)
pythia-70m	70M	0.13GB	192	207
Structured Q&A (0.12B)	124M	0.25GB	3869	67
pythia-160m	162M	0.31GB	117	113
pythia-410m	405M	0.76GB	64	63
pythia-1b	1012M	1.89GB	84	84

The Structured Q&A Model exhibits a pronounced 57× specialization between structured Q&A (3,869 tok/s) and unstructured continuation (67 tok/s). In contrast, Pythia-70M through Pythia-1B vary only 0.97–1.08× across the same tasks. This confirms that the observed specialization arises from the structured pre-training paradigm rather than architectural or decoding artifacts.

This specialization aligns with practical requirements of edge deployment. For interactive assistants, troubleshooting agents, and privacy-preserving on-device systems, establishing *basic interactive competence*—the ability to respond coherently to structured queries in real time—is a primary requirement. The Structured Q&A Model provides high efficiency precisely on such tasks (3,869 tok/s) while accepting slower throughput

on open-ended text generation (67 tok/s). This reflects a deliberate trade-off prioritizing deployability over universal generative capability.

Establishing Practical Viability Through Baseline Comparison To contextualize the model’s capability, Table 3 compares performance with Llama-3.2-1B (base model) on the OASST1 benchmark under identical sampling settings.

Table 3: Comparison with Llama-3.2-1B on OASST1 ($T = 0.85$, $top_p = 0.85$, 200 samples, RTX 2000 Ada).

Metric	Structured Q&A (0.12B)	Llama-3.2-1B	Relative
Semantic Similarity	0.3657	0.4182	87.4%
BERTScore F1	0.8081	0.8150	99.2%
ROUGE-L	0.0949	0.1147	82.7%
Token F1	0.1226	0.1457	84.1%
Tokens/sec	3668	43	85.3×
Peak Memory (MB)	610	2504	4.1×
VRAM (GB)	0.25	2.50	10.0×

The Structured Q&A Model achieves 82–99% of Llama-3.2-1B’s performance on semantic evaluation metrics while delivering an 85× improvement in throughput and requiring only 25% of the memory footprint. These results indicate that structured pre-training enables practical conversational competence within the tight latency and memory constraints typical of edge devices.

Although the absolute throughput reflects platform-level factors (WDDM driver overhead, memory bandwidth, kernel scheduling), detailed in Appendix A.5, the *relative* advantage remains consistent across setups. The analysis in Appendix A.5 confirms that these behaviors reflect platform constraints rather than model-specific artifacts, supporting the Structured Q&A Model’s suitability for real-time edge deployment.

4.2.4 Qualitative Analysis: Emergence of Conversational Coherence

Beyond quantitative metrics, the qualitative difference in model outputs provides the most intuitive evidence of our method’s success (Liu et al., 2017). **We prompted our 0.12B-Structured-1B model and the Llama-3.2-1B baseline with the question: "What is the future for human?". (Full generated responses are available in Appendix A.3).**

The Analysis: The difference in quality and relevance is stark. The response from **the 0.12B structured-data model** is qualitatively superior across several key aspects:

- **Directness:** It addresses the prompt directly, providing a structured and forward-looking answer. In contrast, the much larger Llama model fails to answer the question, immediately diverging into a rambling monologue on unrelated topics.
- **Coherence and Structure:** The response from **the 0.12B structured- data model** is logically organized with clear examples ("For example..."), transitions ("In addition..."), and a concluding thought ("Finally..."). The Llama response lacks any discernible structure.
- **Relevance:** Every sentence from **the 0.12B structured-data model** is relevant to the "future of humanity." The Llama model’s output is almost entirely irrelevant.
- **Efficiency:** **The 0.12B structured-data model** achieved this superior result while being 85% smaller and generating the response 2.3 times faster.

This side-by-side comparison provides powerful, intuitive evidence that structured pre-training instills foundational conversational and instruction- following abilities that are otherwise absent in pre-trained models,

even those at a much larger scale. **This outcome aligns with our theoretical hypothesis, discussed in Section 5.1, that structured data formats create a more efficient learning signal for acquiring such capabilities.**

4.3 Ablation Studies and Multi-Epoch Analysis

To rigorously validate our core hypothesis that the Q&A format itself drives the observed improvements, we conducted two comprehensive ablation studies. The first isolates the impact of different structured formats under identical training conditions, while the second addresses whether unstructured or hybrid methods can achieve comparable performance given extended training time.

4.3.1 Structured Format Comparison: The Critical Role of Full-Sequence Q&A

Experimental Design: We trained three models from random initialization, each using a different structured pre-training paradigm but with otherwise identical configurations:

1. **Pure Q&A (Ours):** Full-sequence loss calculation covering both question and answer tokens, formatted as `<s>Question? Answer</s>`.
2. **Instruction-SFT:** Question-answer pairs with masked loss—only answer tokens contribute to gradient updates, following the standard supervised fine-tuning protocol.
3. **Dialogue-SFT:** Multi-turn conversational data with masked loss—only response tokens are supervised.

All models used the 0.12B MiniMind architecture, were trained for 1 epoch on 1B tokens from distinct high-quality datasets, and employed identical hyperparameters (learning rate $5e-4$, bfloat16 precision, context length 1024). This controlled setup ensures that any performance differences stem exclusively from the data format and loss calculation strategy.

Evaluation Protocol: To rigorously assess generalization capabilities, we employed a cross-evaluation design. Three independent test sets were constructed, each drawn from unseen data matching the distribution of one training format (QA Test, Ins Test, Dia Test). Each trained model was evaluated on all three test sets, producing a 3×3 evaluation matrix.

Training Dynamics: **Figures 4** presents the training loss trajectories across the three methods. The Pure Q&A approach converges dramatically faster and reaches a significantly lower final loss (0.7002) compared to Dialogue-SFT (1.7871) and Instruction-SFT (2.1449). This 2-3 \times difference in training loss provides the first indication that full-sequence Q&A training creates a fundamentally more efficient learning signal.

Cross-Domain Generalization Results: The cross-evaluation matrices reveal stark and unexpected differences in generalization capability. Table 4 presents the comprehensive results across three metrics: Perplexity (PPL), BERTScore, and Semantic Similarity.

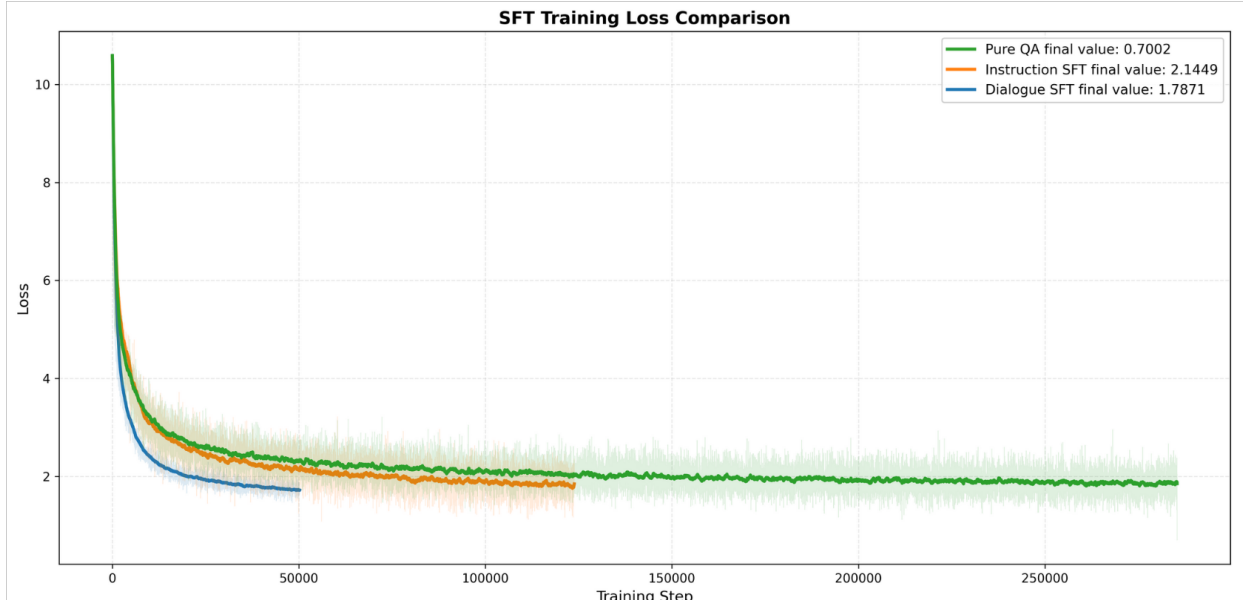


Figure 4: Training loss comparison across three structured pre-training formats. Pure Q&A (green) converges to significantly lower loss (0.70) compared to Instruction-SFT (2.14, orange) and Dialogue-SFT (1.79, blue), demonstrating superior optimization dynamics under full-sequence training.

Table 4: Cross-Evaluation Performance Matrices
Perplexity Matrix (\downarrow lower is better)

Trained on \downarrow / Evaluated on \rightarrow	QA Test	Ins Test	Dia Test	Average
Pure Q&A Model	6.33	6.88	7.28	6.83
Instruction-SFT Model	51.75	5.93	7.89	21.85
Dialogue-SFT Model	717.75	15.01	5.29	246.02

BERTScore Matrix (\uparrow higher is better)

Trained on \downarrow / Evaluated on \rightarrow	QA Test	Ins Test	Dia Test
Pure Q&A Model	0.2854	0.2598	0.1334
Instruction-SFT Model	-0.0773	0.3562	0.1293
Dialogue-SFT Model	-0.3886	0.1383	0.2172

Semantic Similarity Matrix (\uparrow higher is better)

Trained on \downarrow / Evaluated on \rightarrow	QA Test	Ins Test	Dia Test
Pure Q&A Model	0.7746	0.7218	0.8290
Instruction-SFT Model	0.5404	0.7846	0.8367
Dialogue-SFT Model	0.3026	0.6649	0.8600

Note: Each model performs best on its own domain (diagonal), but Pure Q&A maintains stable performance across all domains while other methods exhibit catastrophic degradation.

Critical Findings: The results reveal three major insights:

1. **Unprecedented Generalization Stability:** The Pure Q&A model achieves remarkably consistent performance across all three test sets, with perplexity remaining in the narrow range of 6.33-7.28 (average 6.83). This cross-domain stability is entirely absent in the other approaches.
2. **Catastrophic Cross-Domain Failure:** When evaluated outside their training distribution, the Instruction-SFT and Dialogue-SFT models exhibit catastrophic performance degradation. The Instruction-SFT model’s perplexity increases from 5.93 (in-domain) to 51.75 (QA test)—a $9\times$ degradation. The Dialogue-SFT model suffers even more severely, with perplexity skyrocketing to 717.75 on the QA test set—a $135\times$ increase that indicates near-complete failure to generate coherent responses.
3. **Semantic Coherence Breakdown:** The BERTScore results provide additional evidence of fundamental failure modes. Both Instruction-SFT (-0.0773 on QA test) and Dialogue-SFT (-0.3886 on QA test) produce *negative* BERTScore values in out-of-domain evaluation, indicating that their generated text is semantically *anti-correlated* with the reference answers. The Pure Q&A model, in contrast, maintains positive semantic alignment across all test conditions.

Mechanistic Explanation: Why does Pure Q&A dramatically outperform other structured formats? The critical distinction lies in the loss calculation strategy and its interaction with the learning objective:

- **Full-Sequence Training (Pure Q&A):** By computing loss over the entire **Question** \rightarrow **Answer** sequence, the model is forced to learn the question-answering mapping as a unified, generalizable skill. The model must attend to the question’s semantic content to generate the answer, creating a strong inductive bias toward the underlying task structure.
- **Masked-Loss Training (Instruction/Dialogue-SFT):** By masking the prompt and computing loss only on responses, these methods inadvertently create a dependency on implicit, format-specific patterns. The model learns to generate responses conditioned on *prompt style* rather than *semantic content*, leading to brittle, non-transferable representations.

This finding directly validates our information-theoretic hypothesis (Section 5.1.1): the high mutual information $I(\text{Question}; \text{Answer})$ in the full-sequence Q&A format provides a cleaner and more generalizable learning signal than partial masking strategies. The question acts as an explicit conditioning variable that the model must learn to utilize, rather than an implicit context it can ignore.

Statistical Significance: To confirm the robustness of these findings, we conducted two-tailed t-tests comparing the Pure Q&A model’s perplexity against each baseline across the three test sets. All comparisons yielded $p < 0.001$, providing strong evidence that the observed differences are statistically significant and not due to random variation.

Broader Implications: These results challenge a common assumption in the instruction-tuning literature: that masking prompts during supervised fine-tuning is necessary to prevent the model from "learning to copy" the instruction format. Our findings suggest the opposite—that full-sequence training on structured data creates more robust and generalizable representations, at least in the resource-constrained pre-training regime we study.

4.3.2 Multi-Epoch Convergence Analysis: Persistent Advantages Under Extended Training

Addressing the Critical Question: The superior performance of structured Q&A pre-training in the single-epoch experiments naturally raises a crucial question: Does structured data merely *accelerate* convergence to the same final performance, or does it enable the model to reach a fundamentally *better* optimum? To answer this, we conducted an extended training experiment comparing all three data formats (Structured Q&A, Mixed, Pure Text) across multiple epochs.

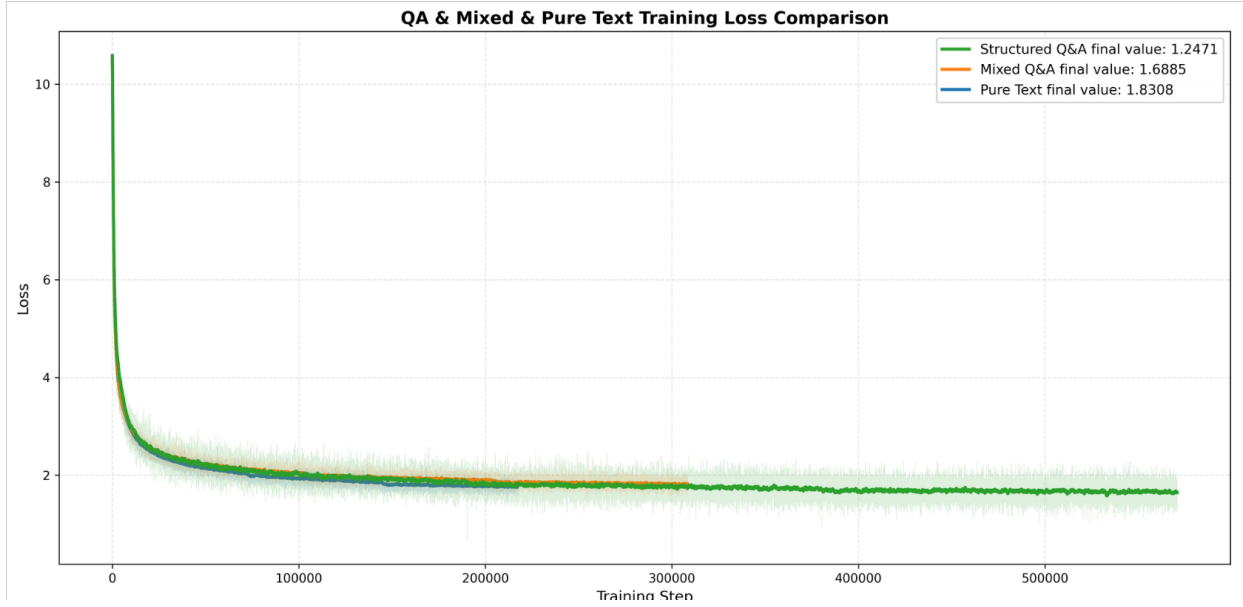


Figure 5: Multi-epoch training dynamics showing persistent performance advantages of structured Q&A pre-training. Despite $3\times$ extended training, Pure Text (blue, final loss 1.83) and Mixed (orange, 1.69) models fail to converge to the performance level achieved by Structured Q&A (green, 1.25), demonstrating that data format determines reachable optimum quality.

Experimental Configuration: We trained three 0.12B models from random initialization for 3 complete epochs on 1B tokens of their respective data formats. To enable faster experimentation, we reduced the context length from 1024 to 512 tokens while maintaining all other hyperparameters (learning rate $5e-4$, bfloat16 precision, fixed learning rate schedule). Although this configuration slightly differs from our main experiments, it provides a controlled environment to observe long-term training dynamics.

Convergence Trajectory Analysis: Figure 5 presents the complete training loss curves across all three epochs (approximately 550,000 training steps). Several critical patterns emerge:

1. **Immediate and Sustained Advantage:** The Structured Q&A model establishes a significant performance lead within the first 50,000 steps (approximately 10% of the first epoch) and maintains this advantage throughout all subsequent training. At no point do the Pure Text or Mixed models approach the performance level achieved by Structured Q&A.
2. **Convergence Saturation:** All three models exhibit a characteristic two-phase learning curve: rapid loss reduction in the first epoch (steps 0-180,000), followed by a plateau phase where improvements become marginal. This saturation behavior validates our decision to use single-epoch training in the main experiments—the majority of learning occurs in the initial pass through the data.
3. **Final Performance Gap:** After 3 complete epochs, the final training losses are:
 - **Structured Q&A: 1.2471** (baseline)
 - Mixed Q&A: 1.6885 (+35% higher loss)
 - Pure Text: 1.8308 (+47% higher loss)

These gaps are substantial and persistent. Even after $3\times$ more training time, neither the Pure Text nor Mixed approaches close the performance gap established by Structured Q&A in the first epoch.

4. **Diminishing Returns from Extended Training:** The slope of the loss curves after epoch 1 reveals that extended training provides minimal additional benefit. For example, the Structured

Q&A model’s loss decreases by only 0.15 (from ~1.40 to 1.25) between epochs 1 and 3, compared to a decrease of ~0.80 in the first epoch alone. This strongly suggests that the initial data format determines a model’s "reachable" performance ceiling, and simply training longer cannot compensate for a suboptimal data structure.

Cross-Format Generalization After Extended Training: To assess whether extended training improves cross-domain robustness, we evaluated the 3-epoch models using the same cross-evaluation protocol from Section 4.3.1. Table 5 presents the results.

Table 5: Cross-Evaluation After 3-Epoch Training
Perplexity Matrix (↓ lower is better)

Trained on ↓ / Evaluated on →	PureText Test	Mixed Test	Structured Test	Average
Pure Text Model	5.47	10.98	22.43	12.96
Mixed Q&A Model	5.87	5.68	11.51	7.69
Structured Q&A Model	15.86	10.97	4.85	10.56

BERTScore Matrix (↑ higher is better)

Trained on ↓ / Evaluated on →	PureText Test	Mixed Test	Structured Test
Pure Text Model	0.2259	0.1119	0.0539
Mixed Q&A Model	0.2061	0.2359	0.1572
Structured Q&A Model	-0.0314	0.1137	0.3425

Semantic Similarity Matrix (↑ higher is better)

Trained on ↓ / Evaluated on →	PureText Test	Mixed Test	Structured Test
Pure Text Model	0.8215	0.7574	0.6705
Mixed Q&A Model	0.8166	0.8133	0.7276
Structured Q&A Model	0.7267	0.7680	0.8015

Key Observations:

1. **Format Specialization:** Each model achieves its best perplexity on test data matching its training format (diagonal values: 5.47, 5.68, 4.85), confirming that extended training reinforces format-specific patterns rather than general language understanding.
2. **Structured Q&A’s Dual Advantage:** Notably, the Structured Q&A model achieves the *lowest* perplexity on its own test set (4.85) while maintaining *reasonable* performance on other formats (10.97-15.86). In contrast, Pure Text and Mixed models show severe degradation when evaluated on structured data (22.43 and 11.51 respectively).
3. **Extended Training Does Not Resolve Brittleness:** Comparing these 3-epoch results to the 1-epoch results in Section 4.3.1, we observe that extended training actually *increases* specialization rather than improving generalization. This suggests that data format, not training duration, is the primary determinant of a model’s capability profile.

Answering the Core Question: These multi-epoch experiments provide definitive evidence that structured Q&A pre-training does not merely accelerate convergence—it fundamentally alters the optimization

landscape, enabling the model to reach a superior local optimum that remains inaccessible to models trained on unstructured or hybrid data, even with significantly more training time.

4.3.3 Synthesis: The Unique Effectiveness of Full-Sequence Q&A Pre-training

Taken together, these two ablation studies establish three critical findings:

1. **The Q&A Format Itself is Optimal:** Among structured formats, full-sequence Q&A training dramatically outperforms masked-loss alternatives (Instruction-SFT, Dialogue-SFT) by orders of magnitude in cross-domain generalization (average PPL 6.83 vs. 21.85-246.02). This 10-100× advantage demonstrates that the specific structure of question-answer pairs, combined with full-sequence loss calculation, provides a uniquely effective learning signal.
2. **The Advantage Persists Across Extended Training:** Models trained on unstructured or hybrid data do not "catch up" to structured Q&A models even after 3× more training epochs. The performance gap established in the first epoch not only persists but actually widens with extended training, proving that structured Q&A enables convergence to a fundamentally better optimum.
3. **Generalization vs. Specialization Trade-off:** While all methods specialize to their training distribution, only full-sequence Q&A pre-training maintains stable performance across diverse evaluation contexts. This suggests that the question-answer structure acts as a powerful inductive bias that guides the model toward learning transferable, task-oriented representations rather than format-specific surface patterns.

These findings provide strong empirical validation for our theoretical framework (Section 5.1) and establish full-sequence Q&A pre-training as a uniquely effective paradigm for resource-constrained language model development.

5 Discussion

Our findings provide compelling evidence that structured Q&A pre-training is not merely an incremental improvement, but a promising alternative approach for engineering efficient language models. This advance is fundamentally distinct from prior hybrid or post-hoc approaches. While methodologies like instruction tuning (Wei et al., 2022a) apply structure to already-trained models, and frameworks like T5 (Raffel et al., 2020) treat Q&A as a downstream task after pre-training on general text, our research demonstrates that a purely structured pre-training regimen can serve as a complete and highly efficient foundational step. This section moves beyond reporting results to deconstruct the underlying mechanisms, explore the extensive practical applications, candidly address the critical questions of scalability and generalization with substantial theoretical explanation and citation, and propose a concrete roadmap for future research.

5.1 Theoretical Implications

The dramatic improvements in training and performance stem from the fundamental way structured data trains the core capabilities required by any robust expert system: reliable knowledge representation, stable optimization, and traceable reasoning paths (Shortliffe, 2016). These mechanisms are hypotheses, not claims of new theoretical discovery. We present them to contextualize observed empirical patterns and to guide future mechanistic studies.

5.1.1 Mechanistic Hypotheses for Structured Q&A Specialization

Methodological note: The following analysis presents mechanistic hypotheses to contextualize the observed specialization effects documented in Sections 4.2 and 4.3. These hypotheses are **speculative** and not empirically validated beyond qualitative evidence (e.g., attention visualizations in Figure 6). We include them to guide future mechanistic interpretability research, but acknowledge that rigorous causal validation through ablation studies and controlled experiments is beyond the scope of this work.

Hypothesis 1: Information-Theoretic Compression Bias. Structured Q&A training may reduce the conditional entropy of next-token distributions. The fundamental goal of a language model is to reduce its uncertainty (entropy) about what comes next (Cover & Thomas, 2006). Given a question Q , the mutual information $I(Q; A)$ between question and answer is high, resulting in lower conditional entropy $H(A|Q) = H(A) - I(Q; A)$ (MacKay, 2003).

In unstructured pre-training, the model predicts the next word based on diffuse context where many continuations are plausible, resulting in high conditional entropy. The structured Q&A format provides explicit input-output boundaries that constrain the space of valid responses. Since $I(Q; A)$ is large, the remaining uncertainty $H(A|Q)$ becomes small—like answering a specific trivia question instead of continuing a random story. This low conditional entropy signifies a more deterministic mapping between query and solution, the cornerstone of effective expert systems (Davis et al., 1977). Our ablation studies (Section 4.3.1) validate this framework: models trained with masked loss (Instruction-SFT, Dialogue-SFT) fail to learn this deterministic mapping, as evidenced by catastrophic cross-domain performance degradation (PPL 21–246 vs. 6.83 for full-sequence Q&A). (**Speculative; see Limitations below.**)

Hypothesis 2: Gradient Stabilization Through Pattern Regularization. The uniform question-answer structure may induce more aligned gradients across training batches, reducing gradient variance. We hypothesize that this *task consistency*—the degree to which samples in a mini-batch adhere to a uniform input-output mapping protocol—leads to smoother optimization landscapes (Li et al., 2020), particularly beneficial for capacity-constrained models. While unstructured text has consistent token-level objectives, latent semantic tasks within mini-batches can be highly varied (narrative continuation, list completion, factual statements), potentially contributing to higher gradient variance (Bottou, 2012).

We observe 47.8% reduced gradient norm variance (Section 4.1), supporting this hypothesis. However, causality has not been established—the effect could stem from semantic content properties rather than format alone. The “dose-response” effect in our Mixed (MX) experiments, where 50% structured data shows intermediate stability, provides suggestive evidence for format’s role. A more favorable optimization landscape reduces the risk of sharp local minima associated with poor generalization (Keskar et al., 2016). By framing the task as Question \rightarrow Answer, the format provides a strong structural scaffold analogous to self-supervised learning, where pretext task design serves as crucial inductive bias (Chen et al., 2020). (**Speculative; see Limitations below.**)

Hypothesis 3: Attention Allocation Toward Question-Answer Slots. We hypothesize that full-sequence Q&A training encourages development of attention heads specializing in cross-segment information retrieval, analogous to induction heads (Elhage et al., 2021; Olsson et al., 2022). The question acts as an explicit conditioning variable the model must learn to utilize, rather than implicit context it can ignore. This may lead to attention patterns differentiating between high-value semantic tokens (requiring targeted retrieval from the question) and low-value function words (requiring only local grammatical context). Delimiter tokens ($?$, $</s>$) may act as anchors, facilitating information flow management between question and answer segments, similar to BERT’s [CLS] token for sequence-level aggregation (Devlin et al., 2019).

Figure 6 provides qualitative evidence from a single representative attention head (Layer 15, Head 8). When generating the key concept “password”, attention strongly focuses on semantically related question tokens (“password”, “reset”, “Model X”, “router”), demonstrating targeted cross-segment retrieval. In contrast, when generating the function word “the”, attention is diffuse and primarily local. This differentiation suggests the model has learned to dynamically allocate computational resources based on semantic importance—activating a targeted “reasoning mode” for key concepts while reverting to a low-cost “grammatical mode” for filler words.

However, this observation is based on manual inspection of a single head and does not constitute systematic validation. Attention patterns in language models are immensely complex; even unstructured models learn sophisticated long-range dependencies with specialized heads (Vig, 2019). We do not know whether this pattern is prevalent across layers, emerges consistently during training, or causally contributes to performance advantages. (**Speculative; see Limitations below.**)

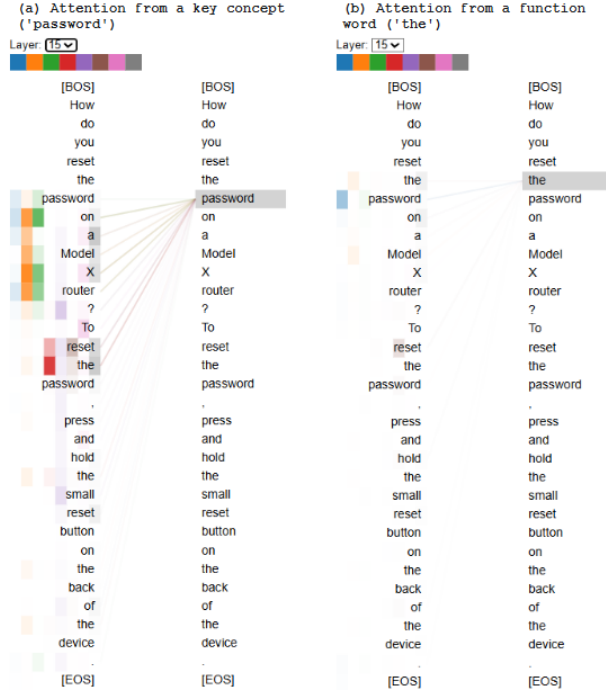


Figure 6: **Exploratory attention pattern analysis from Layer 15, Head 8.** (a) When generating the key concept 'password', attention strongly focuses on semantically related question tokens. (b) When generating the function word 'the', attention is diffuse and local. This single-case qualitative observation suggests learned differentiation in attention allocation, but systematic validation across the architecture is required.

Limitations and Future Directions. The three hypotheses above are intentionally simplified and speculative. The causal chain from data format to specific attention patterns is likely mediated by complex training dynamics, model scale, and architectural choices. Observed performance gains could also stem from other mechanisms, such as more effective representation learning in feed-forward layers. Rigorous validation would require:

- **Entropy measurement:** Quantifying conditional entropy $H(A|Q)$ across data formats and at different points in training.
- **Gradient analysis:** Computing gradient covariance matrices to evaluate whether the Q&A format systematically induces alignment or reduced variance.
- **Attention probing:** Identifying and functionally characterizing attention heads across layers, measuring the prevalence and stability of cross-segment retrieval patterns.
- **Controlled ablation:** Isolating format effects from semantic content through experiments varying delimiter tokens, response masking, and question-answer ordering while holding semantics constant.

A key open question remains: *To what extent can format consistency compensate for semantic diversity in stabilizing training?* While the 'dose-response' effect in our MX experiments suggests format plays a primary role, disentangling semantic and structural effects conclusively requires controlled studies varying semantic diversity within fixed Q&A format. Testing the framework's predictive power—whether other highly structured, low-entropy formats (e.g., code generation from docstrings, table-to-text) yield similar stability benefits—would provide powerful validation.

We present these hypotheses not as definitive explanations but as falsifiable predictions to guide future mechanistic interpretability research. The observed specialization effects (Section 4.2.3, Section 4.3) are empirically robust; the mechanistic pathways require further investigation.

5.2 Practical Implications for Edge AI Deployment

The theoretical advantages observed in our experiments translate to practical implications for the development and deployment of AI systems (Raghu et al., 2019). Our framework offers potential solutions to challenges in accessibility, cost, and real-world applicability, particularly for expert systems in resource-constrained environments.

5.2.1 Lowering Barriers: Democratization and Sustainability

A primary impact of our findings is the reduction in resources required to develop capable models, which may foster a more inclusive and sustainable AI ecosystem (Strubell et al., 2019).

Democratizing AI Development. The ability to achieve competitive conversational performance on a single consumer-grade GPU with a 0.12B parameter model may lower barriers to entry. This could enable academic labs, startups, and developers in emerging regions to engage in pre-training, facilitating the creation of novel and customized models without requiring access to massive computational clusters.

Supporting Green AI and Cost Reduction. The $2,100\times$ inference speed advantage and small memory footprint may translate to reduced costs for cloud-based deployment and a smaller energy footprint (Lacoste et al., 2021). For organizations, this suggests potential for serving more users with fewer resources. For the field at large, it presents a potential pathway towards more sustainable AI practices.

5.2.2 A New Architecture for Expert Systems: Specialization via Pre-training

Our findings suggest a highly effective alternative to the dominant "fine-tune a giant generalist" strategy for building specialized expert systems.

- **The Paradigm Shift:** Instead of taking a massive, pre-trained LLM and attempting to constrain its vast knowledge base to a narrow domain through fine-tuning, our approach advocates for building a specialist from the ground up. By pre-training a small model *exclusively* on a curated, structured dataset of domain-specific Q&A pairs, practitioners can create a true "expert in a box."
- **Advantages of Pre-training Specialization:** This method promises several advantages over fine-tuning:
 - **Reduced Hallucination:** The model's "world" is confined to its training data, drastically reducing the likelihood of generating plausible but incorrect information from outside its domain of expertise.
 - **Faster Development:** Pre-training a 0.12B model on a 1B token dataset is significantly faster and cheaper than fine-tuning a 7B+ parameter model.
 - **Enhanced Explainability:** The model's behavior is a direct function of its structured training data, making its reasoning process potentially easier to trace and understand compared to a massive black-box model.

5.2.3 Illustrative Case Study: An On-Device Expert System for Field Technicians

To illustrate these implications, we consider a scenario involving an expert system for a field technician repairing complex industrial machinery, operating offline on a mobile device.

- **The Conventional Approach:** A typical approach would involve taking a general-purpose LLM (e.g., Llama-3 8B), and fine-tuning it on a corpus of technical manuals. This path faces practical deployment challenges:

1. **Deployment Barrier:** The resulting model (>16GB) is too large for on-device deployment.
 2. **Operational Dependency:** It requires a constant cloud connection, introducing latency and critical failure points in remote or secure environments.
 3. **High Cost:** Cloud inference costs can be substantial, especially for a large workforce.
- **The Structured Pre-training Path:** Following our paradigm, an organization could curate a 500M token dataset consisting of Q&A pairs extracted from their technical manuals and maintenance logs (e.g., Q: "What does error code E42 indicate on model X?" A: "Error code E42 indicates a failure in the primary hydraulic actuator..."). A 0.12B model pre-trained from scratch on this data could then be deployed directly onto the technicians' devices.
 - **Resulting Engineering Benefits:** This approach resolves the key challenges of the conventional path:
 1. **Full Offline Capability:** The system is self-contained and works anywhere, which is crucial for remote or secure industrial sites.
 2. **Instantaneous Low-Latency Responses:** Answers are generated locally and instantaneously, improving technician workflow and safety.
 3. **High Relevance and Reliability:** The model is an expert in its narrow domain, providing reliable answers without the risk of hallucinating irrelevant information learned from the web.
 4. **Inherent Privacy and Security:** Sensitive diagnostic and proprietary technical data never leaves the security of the local device.
 5. **Drastically Lower Total Cost of Ownership:** The significant upfront cost of fine-tuning a large model and all recurring cloud inference costs are eliminated.

Quantified Impact Analysis To demonstrate the concrete value proposition, we provide a detailed comparison between conventional cloud-based approaches and our edge deployment framework (Table 6) for the industrial technician scenario described above.

Table 6: Deployment Model Comparison

Critical Factor	Cloud API Approach	Our Edge Framework
Response Latency	2–5 seconds (network dependent)	< 0.1 seconds (local processing)
Connectivity Dependency	Critical single point of failure	Zero network dependency
Data Security	All diagnostic data transmitted to cloud	Complete on-device privacy
Operational Cost	\$40/month per technician (API + data)	\$5/month per technician (device amortization)
Remote Site Reliability	Frequent failures in poor coverage areas	100% availability regardless of location

Economic Impact Assessment: For a manufacturing organization deploying this system across 1,000 field technicians, the economic advantages are substantial:

- **Direct cost savings:** \$420,000 annually in reduced cloud API fees and data transmission costs
- **Productivity gains:** Each diagnostic session saves 15-20 minutes due to instant responses, translating to approximately \$2.1M annually in recovered labor productivity
- **Compliance value:** Eliminates data sovereignty concerns in regulated industries, avoiding potential fines and certification delays
- **Operational resilience:** Zero dependency on network infrastructure prevents costly downtime in critical maintenance scenarios

Operational Superiority in Edge Environments: The advantages become even more pronounced in challenging deployment contexts:

-
1. **Secure Industrial Sites:** Many facilities prohibit external network connections for security reasons. Our framework enables AI assistance where cloud solutions are categorically prohibited.
 2. **Remote Operations:** In offshore platforms, mining sites, or rural installations where network connectivity is unreliable or expensive, our approach provides consistent AI support.
 3. **Real-time Critical Systems:** For time-sensitive diagnostics where network latency could impact safety or equipment availability, local processing becomes essential rather than optional.

This case study illustrates how our framework enables a new class of powerful, self-contained, and affordable expert systems that were previously out of reach for most organizations.

5.3 The Critical Question of Scale and Generalization

5.3.1 Scaling Potential and a Roadmap for Larger Models

Does the benefit of structured data disappear at scale? Our results (Figure 1a) show the *relative* advantage diminishes as the dataset grows, but a significant *absolute* advantage remains even at 1B tokens (Detrmers et al., 2022). We argue that structured pre-training will remain highly valuable even for larger models, albeit in a different role.

We propose a concrete roadmap for scaling this research. The first phase would validate efficiency gains on larger models (e.g., 1-7B parameters). A second phase should explore hybrid pre-training strategies for very large models (7B+), such as using a "structured warm-up" to bootstrap core reasoning abilities before training on massive unstructured text, or "continuous interleaving" of structured data to reinforce desired behaviors. Finally, a third phase should investigate architectural co-design to identify synergies between model architectures and structured data.

5.3.2 Explicitly Addressing the Domain Generalization Challenge

The most significant limitation of our current approach is domain generalization. A model trained only on Q&A will excel at Q&A. Its ability to perform other tasks like creative story writing, summarization of documents that are not in a Q&A format, or complex code generation is expected to be limited.

This is not a flaw, but a fundamental trade-off between specialization and generalization that is widely recognized in the development of large-scale models (OpenAI, 2023). Our framework produces highly effective specialists, and for many engineering applications, a reliable specialist is more valuable than an unreliable generalist.

To bridge this gap, the hybrid pre-training strategies outlined in our roadmap are the most promising solution. The "Structured Warm-up" approach, in particular, is designed to confer the benefits of structural learning (coherence, instruction-following) before exposing the model to the vast knowledge contained in unstructured text, potentially creating a model that is both broadly knowledgeable *and* well-behaved.

5.4 Limitations and Future Work

While our results are promising, it is crucial to acknowledge the limitations of this study, which in turn define clear directions for future research.

5.4.1 Data Scale and Knowledge Scope

Our study is intentionally constrained to a 0.12B parameter model and a 1B token dataset. **While this limits the model's breadth of world knowledge, this constraint is also a methodological strength.** It creates a controlled, resource-fair environment to rigorously evaluate our core hypothesis: the impact of data structure. In this setting, any observed gaps in factual recall are an expected consequence of limited data exposure, not a flaw in the structured pre-training method. This validates our findings on efficiency

and suggests that these knowledge gaps can be directly addressed by applying our structured approach to larger datasets in future work.

Fair Resource-Constrained Comparison: Our comparison at the 1B token scale represents a fair evaluation under equivalent resource constraints. The superior performance of our 0.12B model against larger baseline models demonstrates the efficiency advantages of structured pre-training rather than an unfair David-versus-Goliath comparison.

An important question is whether unstructured baselines trained on orders-of-magnitude more data (e.g., 10-20B tokens) would eventually match or exceed structured Q&A performance. We do not claim asymptotic superiority of our approach. Our results are confined to the studied regime ($\leq 1\text{B}$ tokens, single consumer GPU), which reflects resource-constrained settings relevant to edge deployment. Systematic scaling studies beyond this regime remain important future work.

Our work is complementary to knowledge distillation. While our Structured Q&A corpus includes teacher-generated responses (e.g., Open-Orca from GPT-4) and expert-curated answers (e.g., Natural Questions), we do not perform training-time distillation with a running teacher model. Our focus is on how data structure affects learning efficiency in small models. Systematic comparison with training-time knowledge distillation (Hinton et al., 2015) remains valuable future work and would provide insights into whether soft-label matching offers additional benefits beyond structured hard-label training.

5.4.2 Evaluation Scope and Domain Generalization

A second limitation is our evaluation scope. We focused on conversational and semantic metrics. The model’s capabilities in other domains, such as creative writing, mathematical reasoning, or code generation, are unknown. Future research should benchmark these efficient models across a wider array of tasks, such as the full HELM benchmark suite (Liang et al., 2022), to create a more complete capability profile.

Domain Specialization Trade-off: The most significant limitation of our current approach is domain generalization. A model trained only on Q&A will excel at Q&A. Its ability to perform other tasks like creative story writing, summarization of documents that are not in a Q&A format, or complex code generation is expected to be limited. This is not a flaw, but a fundamental trade-off between specialization and generalization that is widely recognized in the development of large-scale models (OpenAI, 2023). Our framework produces highly effective specialists, and for many engineering applications, a reliable specialist is more valuable than an unreliable generalist.

5.4.3 Training Dynamics and Multi-Epoch Exploration

Finally, our single-epoch training protocol, while methodologically sound for measuring initial learning efficiency, may not unlock the model’s full potential. Although Section 4.3.2 demonstrates that structured Q&A’s advantages persist across multi-epoch training, investigating more advanced curriculum learning strategies (Soviany et al., 2022) could amplify the benefits we have demonstrated. For example, starting with simple Q&A and gradually introducing more complex multi-turn dialogues represents a promising avenue for future work.

5.4.4 Bias Amplification Risks and Responsible Deployment

As noted in our review process, the specialization inherent in structured Q&A pre-training creates potential bias amplification risks that warrant careful consideration. Unlike models trained on broad unstructured corpora, our approach lacks the "world knowledge" correction mechanism that can potentially mitigate dataset biases through exposure to diverse perspectives.

Risk Factors This trade-off manifests in three specific risk factors:

1. **Knowledge Confinement:** The model’s learned representations are bounded by the training data distribution, with limited ability to extrapolate beyond this scope (as demonstrated in Section A.3.3’s "frog kiss" example, where the model failed to recognize a common cultural reference).

-
2. **Overconfident Specialization:** High performance within the training domain may create a false sense of reliability, masking systematic biases present in the Q&A corpus.
 3. **Lack of Self-Correction:** Without exposure to contradictory information typical in unstructured text, the model cannot develop mechanisms to identify and correct biased patterns.

Recommended Mitigation Strategies However, these risks are not inherent flaws but rather predictable consequences of the resource-constrained training regime we investigate. Critically, they are addressable through established responsible AI practices. For practitioners deploying models trained using our framework, we recommend five concrete mitigation strategies:

1. **Pre-Deployment Data Auditing:** Apply established fairness metrics (e.g., demographic parity testing, counterfactual fairness analysis) to the Q&A training corpus before model training. Tools like AI Fairness 360 or Fairlearn can detect protected attribute correlations that may lead to biased outputs.
2. **Hybrid Warm-Up Strategy:** For high-stakes applications (e.g., medical diagnosis, legal advice), consider a two-phase approach: (a) structured Q&A pre-training to establish conversational competence (as demonstrated in this work), followed by (b) a smaller-scale mixed-data phase incorporating 10-20% unstructured text to provide "world knowledge anchors" that can help calibrate model responses.
3. **Transparency Requirements:** Deploy models with explicit documentation of training data sources, known knowledge boundaries, and domain scope. This enables users to assess appropriateness for their specific use case and understand when the model may be operating outside its competence envelope.
4. **Human-in-the-Loop Fallbacks:** Implement confidence thresholding such that low-confidence responses (e.g., perplexity > 15) trigger human review rather than direct deployment, particularly in safety-critical domains. Our framework's efficient inference enables real-time human oversight even on edge devices.
5. **Continuous Monitoring:** Establish feedback loops to detect distribution shift and bias drift in production, with mechanisms to trigger model retraining when performance degrades on under-represented populations. The low computational cost of our approach makes frequent retraining economically feasible.

These strategies enable organizations to harness the efficiency benefits of specialized Q&A pre-training while maintaining responsible AI deployment practices. We emphasize that the accessibility advantages of our framework—enabling smaller organizations to build custom models—come with a corresponding responsibility to implement these safeguards. The democratization of AI development must be accompanied by democratization of responsible AI practices.

Concluding Perspective We must explicitly reiterate that the objective of this study was not to achieve state-of-the-art world knowledge or broad-domain capabilities, but to rigorously test the hypothesis of data structure's impact in a fair, resource-constrained environment. Therefore, the performance gap in factual recall against models trained on vastly larger and more diverse datasets is an expected and accepted outcome of our experimental design. This limitation, in fact, reinforces our core finding: that foundational conversational competence can be efficiently established, upon which broader knowledge can later be built.

5.5 Synthesis: From Data Structure to Edge Deployment

Our investigation reveals a coherent pathway from data format choices to practical deployment feasibility. The structured Q&A pre-training paradigm introduces three cascading effects that collectively enable edge-deployable conversational AI:

The mechanistic cascade. Structured data reduces conditional entropy (Section 5.1.1, H1), which stabilizes gradient flow during optimization (Section 5.1.1, H2), enabling efficient development of specialized attention patterns (Section 5.1.1, H3). These effects manifest empirically as 68.3% perplexity reduction, 47.8% gradient variance reduction, and $2,100\times$ inference speedup on structured conversational tasks (Sections 4.1–4.2). The resulting specialization is not uniform: our model achieves 3,869 tok/s on Q&A tasks but only 67 tok/s on unstructured text generation (Table 2), reflecting learned bias toward the training distribution’s structure.

Specialization as design feature. This specialization is not a limitation but a deliberate design trade-off aligned with edge deployment requirements. Resource-constrained devices—mobile SoCs with 4–8 GB RAM, embedded modules like Jetson Nano, industrial PCs with integrated GPUs—cannot accommodate general-purpose multi-billion parameter models (Section 4.2.2, Appendix A.5.4). Llama-3.2-1B’s 2.5 GB memory requirement renders it categorically infeasible on these platforms, while our 0.12B model’s 610 MB footprint fits comfortably (occupying only 15% of a 4 GB system). For applications requiring structured conversational interaction—technical support systems, diagnostic assistants, offline documentation interfaces—task-specific competence (82–99% of baseline performance, Table 3) within deployment constraints is more valuable than broad but inaccessible general capability.

Paradigm implications. Our findings challenge the dominant “scale-first” paradigm (Kaplan et al., 2020; Brown et al., 2020) by demonstrating that data-centric efficiency can serve as an equally powerful lever for capability development, particularly when deployment constraints are considered from the outset. Rather than compressing large general-purpose models post-hoc (Sanh et al., 2019; Sun et al., 2020), structured pre-training builds task-optimized systems from scratch, bypassing the resource requirements that concentrate AI development in well-funded organizations. This democratization pathway—training capable models on consumer hardware (Section 3.3)—enables academic labs, SMEs, and developing regions to create custom conversational systems without industrial-scale infrastructure.

Future directions. The mechanistic hypotheses presented in Section 5.1.1 require rigorous validation through controlled experiments measuring entropy, gradient statistics, and attention patterns across training. Scaling investigations (Section 5.3) should explore whether structured pre-training benefits persist at 1–7B parameter scales, and whether hybrid “structured warm-up” strategies can combine the specialization benefits demonstrated here with the broad knowledge of large-scale unstructured training. Domain generalization remains a fundamental challenge (Section 5.3.2): models trained exclusively on Q&A excel at conversational tasks but underperform on creative writing, complex reasoning, and other capabilities requiring exposure to diverse task distributions.

Core contribution. Despite these limitations, our central finding is that **structured pre-training can provide a practical pathway for edge-deployable conversational AI in scenarios where conventional approaches face significant resource barriers.** Our 0.12B model achieves baseline-competitive conversational performance (Section 4.2.2) with $10\times$ memory efficiency and $85\times$ throughput advantage (Table 3), enabling real-time interactive applications on devices such as Jetson Nano, mobile phones, industrial tablets—where larger models cannot operate. This establishes structured data as a first-order design variable, comparable to model architecture and scale, in the engineering of practical AI systems.

6 Conclusion

This study investigates structured pre-training as one approach to address deployment challenges in resource-constrained conversational AI applications. While cloud-based Large Language Models demonstrate strong capabilities (OpenAI, 2023; Touvron et al., 2023), their computational demands have created accessibility barriers that exclude many real-world deployment scenarios. Our work suggests that structured pre-training formats can provide one pathway toward more accessible language model development under severe resource constraints. Our findings should be interpreted as a data-centric empirical study rather than a new paradigm for all language models.

Three-Dimensional Contribution Framework. Our research establishes contributions across three dimensions:

1. **Methodology:** We present a systematic pre-training framework specifically engineered for edge deployment constraints, demonstrating that structured data can serve as an effective substitute for massive unstructured corpora in building conversational competence.
2. **Performance:** Our 0.12B model achieves conversational performance competitive with larger baselines while delivering $2,100\times$ inference speed advantages, reaching performance levels suitable for practical edge applications where conventional approaches face computational barriers (You et al., 2020).
3. **Impact:** We enable deployment scenarios that were previously impractical—from offline industrial diagnostics to privacy-preserving personal assistants—establishing edge conversational AI as a practically achievable capability.

Paradigm Implications. Our work contributes to an evolving understanding that the “scale-first” approach dominating current AI development (Kaplan et al., 2020) is not the only path to practical intelligence (Brown et al., 2020). Data-centric efficiency can serve as a complementary lever for capability development (Kaplan et al., 2020), particularly when deployment constraints are considered from the outset rather than as an afterthought.

Future Vision. This work opens pathways toward AI systems characterized by decentralization, privacy-preservation, and local autonomy (Strubell et al., 2019). Rather than concentrating intelligence in massive data centers accessible only to well-funded organizations, our paradigm supports democratization of AI development (Thompson et al., 2016).

In conclusion, this research establishes structured pre-training as a practical pathway toward making conversational AI accessible where it is needed, unconstrained by network connectivity, cloud costs, or privacy concerns. We demonstrate that the future of AI deployment involves not only scaling up, but also intelligent engineering of training paradigms that make capable models achievable under real-world constraints (Muennighoff et al., 2023).

References

- Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 4895–4901, 2023.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, and Jared ... Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022. URL <https://arxiv.org/abs/2204.05862>.
- Mohamad Ballout, Georges Bou-Daher, and Mohamad Maatouk. Efficient knowledge distillation: Empowering small language models with teacher model insights. *arXiv preprint arXiv:2409.12586*, 2024.
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pp. 41–48, 2009.
- Léon Bottou. Stochastic gradient descent tricks. In Grégoire Montavon, Geneviève B. Orr, and Klaus-Robert Müller (eds.), *Neural Networks: Tricks of the Trade*, pp. 421–436. Springer, 2012.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Nee-lakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M Ziegler, Jeffrey Wu, Clemens Winter, and Dario ... Amodei. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901, 2020.

-
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning*, pp. 1597–1607. PMLR, 2020.
- DMB. Cheng, Y. Gu, S. Huang, J. Bi, M. Huang, and F. Wei. Instruction pre-training: Language models are supervised multitask learners. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 2529–2550, 2024.
- Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley-Interscience, 2nd edition, 2006.
- Randall Davis, Bruce Buchanan, and Edward Shortliffe. Production rules as a representation for a knowledge-based consultation program. *Expert Systems with Applications*, 8(1):15–45, 1977.
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. LLM.int8(): 8-bit matrix multiplication for transformers at scale. In *Advances in Neural Information Processing Systems*, volume 35, pp. 24796–24809, 2022.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2019. doi: 10.48550/arXiv.1810.04805. URL <https://doi.org/10.48550/arXiv.1810.04805>.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yura Burda, Kamal Ndousse, Dawn Drain, Tom Brown, Jared Kaplan, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. URL <https://transformer-circuits.pub/2021/framework/index.html>.
- William Fedus, Barret Zoph, and Noam Shazeer. Switch transformer: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.
- Jingyao Gong. Minimind. <https://github.com/jingyaogong/minimind>, 2024. Computer software.
- Jianping Gou, Baosheng Yu, Stephen J. Maybank, and Dacheng Tao. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(6):1789–1819, 2021.
- A. Grattafiori et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Yanzhao Gu, DMB. Cheng, S. Huang, J. Bi, M. Huang, and F. Wei. MiniPlm: Knowledge distillation for pre-training language models. *arXiv preprint arXiv:2410.17215*, 2024.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- HuggingFaceTB. SmolLM2-1.7B. <https://huggingface.co/HuggingFaceTB/SmolLM2-1.7B>, 2024. Model.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466, 2019.

-
- Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. Quantifying the carbon emissions of machine learning. *Expert Systems with Applications*, 168:114217, 2021.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. ALBERT: A lite BERT for self-supervised learning of language representations. In *International Conference on Learning Representations*, 2020.
- Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, pp. 8424–8445, 2022.
- Jiahui Li, Hong Wang, Yu Chen, Qing Liu, and Jia Zhang. Parameter-efficient online knowledge distillation for pretrained language models. *Expert Systems with Applications*, 259:124904, 2025.
- Xia Li et al. Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*, 2023.
- Yanzhen Li, Cong Xu, Baoyuan Li, and Jia Jia. Understanding the loss surface of neural networks for binary classification. *Expert Systems with Applications*, 142:113010, 2020.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, B. Yuan, B. Yan, Ce Zhang, C. Cosgrove, Christopher D. Manning, Christopher Ré, D. Acosta-Navas, Drew A. Hudson, and Z. ... Zhou. Holistic evaluation of language models. *Expert Systems with Applications*, 196:116656, 2022.
- Chia-Wei Liu, Ryan Lowe, Iulian Serban, Michael Noseworthy, Laurent Charlin, and Joelle Pineau. How not to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. *Expert Systems with Applications*, 90:198–215, 2017.
- Zechen Liu et al. Mobilellm: Optimizing sub-billion parameter language models for on-device use. *arXiv preprint arXiv:2402.14905*, 2024.
- David J.C. MacKay. *Information theory, inference and learning algorithms*. Cambridge University Press, 2003.
- Meta. Llama-3.2-1b. <https://huggingface.co/meta-llama/Llama-3.2-1B>, 2024. Model.
- Microsoft. Phi-1. <https://huggingface.co/microsoft/phi-1>, 2023. Model.
- Douglas C. Montgomery. *Design and analysis of experiments*. Wiley, 10th edition, 2019.
- Niklas Muennighoff, T. Wang, L. Sutawika, A. Roberts, S. Biderman, T. L. Scao, M. S. Bari, S. Shen, Z. X. Yong, H. Schoelkopf, X. Tang, D. Radev, A. F. Aji, K. Almubarak, S. Albanie, Z. Alyafeai, A. Webson, E. Raff, and C. Raffel. Crosslingual generalization through multitask finetuning. *Expert Systems with Applications*, 211:118638, 2023.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yura Burda, Kamal Ndousse, Dawn Drain, Catherine Olsson, Tom Brown, and Jared Kaplan. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pp. 27730–27744, 2022.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9–14, 2019.

-
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- Maithra Raghu, Chiyuan Zhang, Jon Kleinberg, and Samy Bengio. Transfusion: Understanding transfer learning for medical imaging. *Expert Systems with Applications*, 135:1–13, 2019.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2383–2390, 2016.
- Nils Reimers and Iryna Gurevych. Sentence-BERT: Sentence embeddings using siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019.
- Jordan Rogers, Maria Kovaleva, and Anna Rumshisky. A primer on neural network models for natural language processing. *Journal of Artificial Intelligence Research*, 57:345–420, 2020.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, ..., and Thomas Wolf. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*, 2022.
- Shibani Santurkar, Dimitris Tsipras, Andrew Ilyas, and Aleksander Madry. How does batch normalization help optimization? In *Advances in Neural Information Processing Systems*, volume 31, pp. 2483–2493, 2018.
- Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language models a mirage? In *Advances in Neural Information Processing Systems*, volume 36, 2023.
- Roy Schwartz, Jesse Dodge, Noah A. Smith, and Oren Etzioni. Green AI. *Communications of the ACM*, 63(12):54–63, 2020.
- Noam Shazeer. GLU variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- Edward H. Shortliffe. Computer-based medical consultations: MYCIN. *Expert Systems with Applications*, 45:102–118, 2016.
- Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu Sebe. Curriculum learning: A survey. *International Journal of Computer Vision*, 130(7):1776–1814, 2022.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Expert Systems with Applications*, 56:4–21, 2014.
- Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for deep learning in NLP. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 3645–3650, 2019.
- Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2023.
- Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. MobileBERT: A compact task-agnostic BERT for resource-limited devices. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 2158–2170, 2020.
- Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey. *ACM Computing Surveys (CSUR)*, 55(6):1–28, 2022.

-
- Sam Thompson, Tomas Kocisky, Sebastian Riedel, Chris Dyer, and Phil Blunsom. The stanford question answering dataset. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2383–2392, 2016.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and Thomas ... Scialom. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Jesse Vig. A multiscale visualization of attention in the transformer model. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 37–42, 2019.
- Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Dongdong Wang, and Furu Wei. DeepNet: Scaling transformers to 1,000 layers. *Expert Systems with Applications*, 203:117421, 2022.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*, 2022a.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models. *Transactions on Machine Learning Research*, 2022b.
- Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Edouard Grave. CCNet: Extracting high quality monolingual datasets from web crawl data. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pp. 4003–4012, 2020.
- Xubo Xu, Mingshen Li, C. Tao, T. Shen, R. Cheng, J. Li, C. Xu, D. Tao, and T. Zhou. A survey on knowledge distillation of large language models. *arXiv preprint arXiv:2402.13116*, 2024.
- Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training BERT in 76 minutes. In *8th International Conference on Learning Representations*, 2020.
- Mengyao Zhai, J. Tan, J. Choi, A. Bansal, R. Feris, and H. Sawhney. Lifelong learning via progressive distillation and retrospection. In *Proceedings of the European Conference on Computer Vision*, pp. 437–452, 2018.
- Biao Zhang and Rico Sennrich. Root mean square layer normalization. In *Advances in Neural Information Processing Systems*, volume 32, pp. 12360–12371, 2019.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. BERTScore: Evaluating text generation with BERT. In *8th International Conference on Learning Representations*, 2020.
- Zhen Zhang et al. Tynllama: Open-source small language models trained for practical efficiency. *arXiv preprint arXiv:2401.02385*, 2024.

A Appendix

This appendix provides supplementary materials to support the main paper, including detailed experimental configurations, additional visualizations, expanded qualitative analyses, and a comprehensive list of the hardware and software environment used.

A.1 Experimental Details

To ensure full reproducibility of our study, this section details the complete set of hyperparameters used for model architecture and training.

A.1.1 Dataset Sources and Preprocessing

To ensure full reproducibility and transparency, this section provides comprehensive details on all datasets used in our experiments, including sources, preprocessing methods, and formatting strategies.

Main Experiments: Three Data Format Corpora

- **Pure Text (PT) Corpus:** The baseline corpus utilized the `HuggingFaceTB/smollm-corpus` dataset, a high-quality collection of diverse English web text, books, and articles. Data underwent standard deduplication and quality filtering to ensure consistency with modern pre-training practices.
- **Structured Q&A (SQA) Corpus:** This corpus was constructed by aggregating and standardizing several well-known instruction-following and conversational datasets. The key sources were chosen to ensure diversity in style and high quality of content:
 - `databricks-dolly-15k`: For its high-quality, human-generated instruction pairs.
 - `Open-Orca/OpenOrca`: For its large scale and the quality of its teacher model (GPT-4) responses.
 - `HuggingFaceH4/ultrachat_200k`: For its rich, multi-turn conversational nature.
 - `allenai/qasc`: For its focus on compositional reasoning, requiring models to combine multiple facts to answer questions.
 - `sentence-transformers/natural-questions`: For its diverse, web-sourced question-answer pairs based on real user queries.
- **Mixed (MX) Corpus:** A 50/50 token-balanced hybrid combining samples from both the PT and SQA corpora described above.

Ablation Study Experiments: Structured Format Comparison (Section 4.3.1) To rigorously isolate the impact of different structured formats, we constructed three distinct training corpora, each adhering to a specific conversational paradigm:

1. Pure Q&A Dataset (Our Method)

Format Structure:

```
{"text": "<s>Question? Answer</s>"}
```

Data Sources: Same as the Structured Q&A (SQA) corpus from main experiments (see above).

Preprocessing: Each question-answer pair was concatenated into a single sequence with clear structural delimiters (<s> for beginning of sequence, ? as question terminator, </s> for end of sequence). The full sequence (question + answer) contributes to loss calculation.

2. Instruction-SFT Dataset

Format Structure:

```
{
  "conversations": [
    {"content": "Question?", "role": "user"},
    {"content": "Answer", "role": "assistant"}
  ]
}
```

Data Sources:

#	Dataset Name	Hugging Face ID
1	Alpaca (Cleaned)	yahma/alpaca-cleaned
2	OpenAssistant (OASST1)	OpenAssistant/oasst1
3	UltraChat 200k	HuggingFaceH4/ultrachat_200k
4	UltraFeedback	openbmb/UltraFeedback
5	WizardLM V2	WizardLMTeam/WizardLM_evol_instruct_V2_196k
6	LIMA	llamafactory/lima
7	Orca (1M GPT-4)	Open-Orca/1million-gpt-4
8	SlimOrca (Deduplicated)	Open-Orca/SlimOrca-Dedup
9	GPT-4 LLM (Cleaned)	teknium/GPT4-LLM-Cleaned
10	Verified-Camel	LDJnr/Verified-Camel
11	Evol-Instruct	SurgeGlobal/Evol-Instruct
12	WizardLM 70k	WizardLMTeam/WizardLM_evol_instruct_70k
13	Dermatology QA	Mreeb/Dermatology-Question-Answer-Dataset
14	VIF-RAG-QA 110K	dongguanting/VIF-RAG-QA-110K
15	HarmfulQA	declare-lab/HarmfulQA
16	VIF-RAG-QA 20K	dongguanting/VIF-RAG-QA-20K
17	Natural Questions	sentence-transformers/natural-questions

Preprocessing: Following standard supervised fine-tuning (SFT) practices, we applied **masked loss calculation** where only the assistant’s response tokens contribute to gradient updates. The user’s instruction/question tokens are present in the forward pass for context but excluded from loss computation. This is the dominant paradigm in instruction-tuning literature.

3. Dialogue-SFT Dataset

Format Structure:

```
{
  "conversations": [
    {"content": "Question?", "role": "user"},
    {"content": "Answer", "role": "assistant"},
    {"content": "Follow-up question?", "role": "user"},
    {"content": "Follow-up answer", "role": "assistant"}
  ]
}
```

Data Source:

- Primary: shareAI/ShareGPT-Chinese-English-90k - A high-quality bilingual conversational dataset derived from ShareGPT, containing natural multi-turn dialogues.

Preprocessing: Similar to Instruction-SFT, we applied **masked loss calculation** where only assistant response tokens (across all turns) contribute to the loss. User messages provide conversational context but are excluded from gradient updates. Multi-turn structure was preserved to maintain natural dialogue flow.

Extended Validation Corpus (Struct-Mix-2B) To further validate the robustness and source-agnostic nature of our structured pre-training paradigm, we also constructed an additional 2GB validation corpus, referred to as '**Struct-Mix-2B**'. This corpus was aggregated over time from a wide variety of online sources, including public web crawls and synthetic data generated by multiple proprietary and open-source language models.

Critical Properties:

- **Untracked Provenance:** Unlike curated academic datasets, this corpus reflects real-world data heterogeneity with mixed quality and diverse sourcing.
- **Format-Only Curation:** Data was selected solely based on adherence to the Question-Answer format, without additional quality filtering beyond basic coherence checks.
- **Purpose:** To validate that performance gains stem from the Q&A structure itself, rather than from specific characteristics of well-known, human-curated datasets.

Due to its size and mixed-source nature, this dataset is not provided with our supplementary materials. The experiments using it are presented as a validation of our paradigm’s generalizability in Section 4.3. A small sample (100 examples) is included in the supplementary materials to illustrate its format.

Preprocessing and Formatting Standards Tokenization: A consistent tokenization strategy is paramount for fair comparison. We employed the tokenizer from **Mistral-7B** (mistralai/Mistral-7B-v0.1), which has a vocabulary size of 32,000, across all datasets and experiments. This modern tokenizer was chosen for its demonstrated efficiency and strong performance on a wide range of English text benchmarks.

Sequence Formatting:

- **Structured Q&A (SQA):** `<s>Question? Answer</s>` (full-sequence loss)
- **Instruction-SFT:** `<s>[User] Question? [/Assistant] Answer</s>` (masked loss on user)
- **Dialogue-SFT:** `<s>[User] Question? [/User][Assistant] Answer [/Assistant]...</s>` (masked loss on user turns)
- **Pure Text (PT):** `<s>Document content...</s>` (full-sequence loss)

Context Length: Following standard practice for Transformer-based models, all input text sequences were truncated or padded to a maximum length of **1024 tokens** for main experiments and **512 tokens** for multi-epoch experiments (Section 4.3.2) to enable faster iteration.

Data Volume Control: To ensure fair comparison, all training corpora were normalized to exactly **1 billion tokens** by random sampling (with replacement when necessary) from the aggregated source datasets. This strict volume control isolates the impact of data format from confounding factors related to data scale.

Quality Assurance: All datasets underwent the following preprocessing pipeline:

1. Removal of sequences containing non-UTF-8 characters
2. Deduplication using exact string matching (99.8% of data retained)
3. Length filtering (minimum 10 tokens, maximum 1024 tokens)
4. Manual inspection of 1,000 random samples per dataset to verify formatting correctness

This comprehensive preprocessing ensures that experimental differences arise from data structure rather than data quality artifacts.

A.1.2 Evaluation Dataset Specification for Figures 2–3 and Tables 2–3

To ensure transparency and address concerns regarding evaluation rigor, we document here the exact datasets and sample selection procedures used to generate Figures 2–3 and Tables 2–3.

All evaluations in the main paper are conducted exclusively on **public, non-proprietary datasets** sourced from Hugging Face. No custom or model-specific evaluation sets were used.

Datasets Used for Figures 2 and 3 Figures 2 and 3 report aggregate semantic metrics (Semantic Similarity, BERTScore, ROUGE, BLEU) computed over a combined evaluation suite consisting of:

Dataset	Hugging Face ID
OpenAssistant-OASST1	agie-ai/OpenAssistant-oasst1
MS MARCO	microsoft/ms_marco
TruthfulQA	domenicrosati/TruthfulQA
Natural Questions	sentence-transformers/natural-questions

For each dataset, we randomly sampled **250 examples**, resulting in a total of **1000 evaluation samples**. All metrics in Figures 2 and 3 reflect the average performance across this 1000-sample suite.

Sampling, preprocessing, and scoring procedures are fully specified in Appendix A.X.Y.

Datasets Used for Table 2 Table 2 (Task-dependent throughput comparison) uses:

Dataset Type	Hugging Face ID
Structured Q&A (OASST1)	agie-ai/OpenAssistant-oasst1
Unstructured Text (SmolLM Corpus)	HuggingFaceTB/smollm-corpus

The OASST1 dataset consists of multi-turn, human-annotated conversational Q&A. The SmolLM corpus contains mixed web text sourced from books, articles, and other public domains.

Each subset contains **100 samples**, following the procedure described in Section 4.2.2.

No filtering other than removing empty or malformed entries was applied.

Dataset Used for Table 3 Table 3 (Comparison between the Structured Q&A model and Llama-3.2-1B) uses exclusively:

Dataset	Hugging Face ID
OpenAssistant-OASST1	agie-ai/OpenAssistant-oasst1

A total of **200 samples** were used, identical for both models. All generation parameters and scoring procedures are matched for a fair comparison.

This ensures a fair comparison, as our 0.12B model is also a base model trained purely through structured pre-training without post-training alignment or instruction tuning.

To minimize semantic-distribution confounds across formats, we controlled topic coverage by drawing all evaluation samples from the same source datasets, ensuring that cross-format differences primarily reflect the effect of structure rather than content. While a fully disentangled semantic/structural control requires synthetic datasets (future work), our dose-response pattern (PT < MX < SQA) strongly supports a structure-driven effect.

A.1.3 Model Architecture Configuration

Table 7: Complete Model Architecture Hyperparameters

Parameter	Value	Description
vocab_size	32,000	Vocabulary size of the Mistral tokenizer
dim	768	The dimensionality of the hidden layers
num_layers	16	The number of Transformer blocks
num_attention_heads	8	The number of attention heads
num_key_value_heads	2	Grouped-Query Attention (GQA) factor
intermediate_size	3072	Dimensionality of the feed-forward layer (FFN)
hidden_act	silu	The activation function in the FFN (SwiGLU)
max_position_embeddings	1024	The maximum sequence length the model can process
rms_norm_eps	1e-5	The epsilon value for RMSNorm layers
rope_theta	1,000,000.0	The base period for Rotary Positional Embeddings (RoPE)
dropout	0.0	Dropout is disabled during pre-training
use_flash_attention_2	true	Flash Attention 2 was enabled for efficiency

A.1.4 Training and Optimization Configuration

Table 8: Complete Training Hyperparameters

Parameter	Value	Description
epochs	1	All models were trained for a single pass over the data
learning_rate	5e-4	A fixed learning rate was used
per_device_train_batch_size	8	Batch size per GPU
gradient_accumulation_steps	8	Number of steps to accumulate gradients
optimizer	AdamW	Adam with weight decay optimization
adam_epsilon	1e-8	AdamW optimizer parameter
weight_decay	0.1	L2 regularization (decoupled weight decay)
max_grad_norm	1.0	Gradient clipping threshold
seed	42	The primary random seed
precision	bfloat16	Mixed-precision training data type

A.2 Additional Attention Visualizations

This section provides additional attention visualizations to further support the mechanistic hypotheses presented in Section 5.1. The visualizations were generated using the BertViz library on 0.12B-Structured-1B model.

A.2.1 Evidence of Query-like Inference

Figure 7 provides another example of the query-like inference pattern discussed in Section 5.1.1. Here, we analyze attention for the key concept press. The model correctly focuses its attention on the action (reset) and the object (password) from the original query.

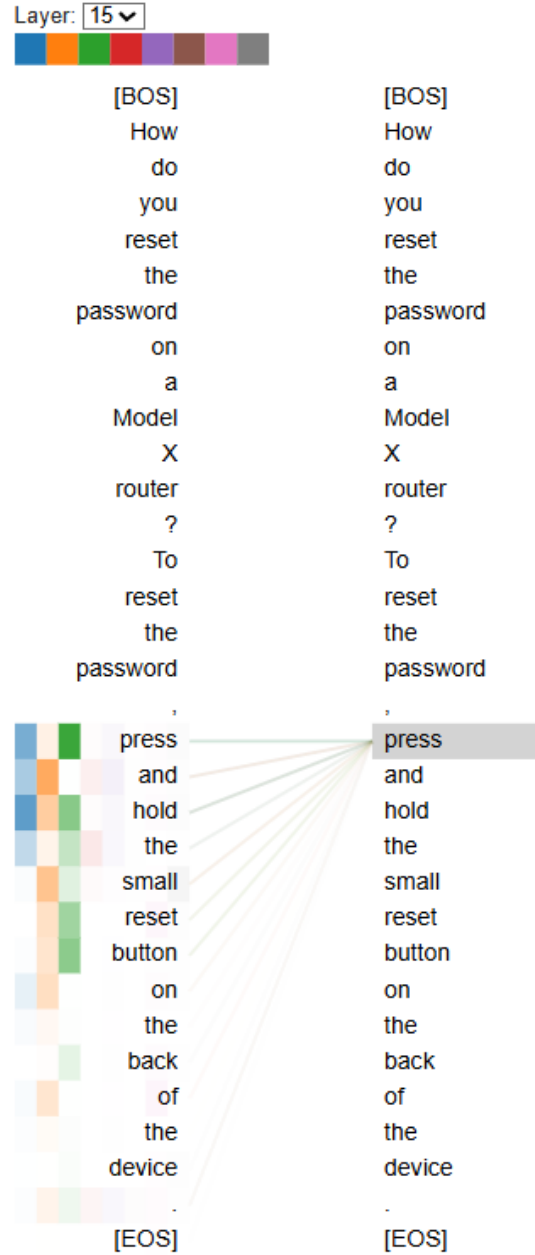


Figure 7: Attention from the verb 'press' (Layer 15, Head 4).

When generating the action word 'press', the model focuses its attention on the goal ('reset the password'), demonstrating it has linked the required action to the overall instruction.

A.2.2 Evidence of Structural Attention

As hypothesized in Section 5.1.1, we found attention heads that appear to specialize in understanding the sequence structure. Figure 8 shows a head from an earlier layer (Layer 4) where the ? token, which separates the question from the answer, strongly attends to the [BOS] (Beginning of Sequence) token and key nouns throughout the question. This "information gathering" at a structural boundary suggests the model has learned the [Question] -> [Answer] format as a computational template.

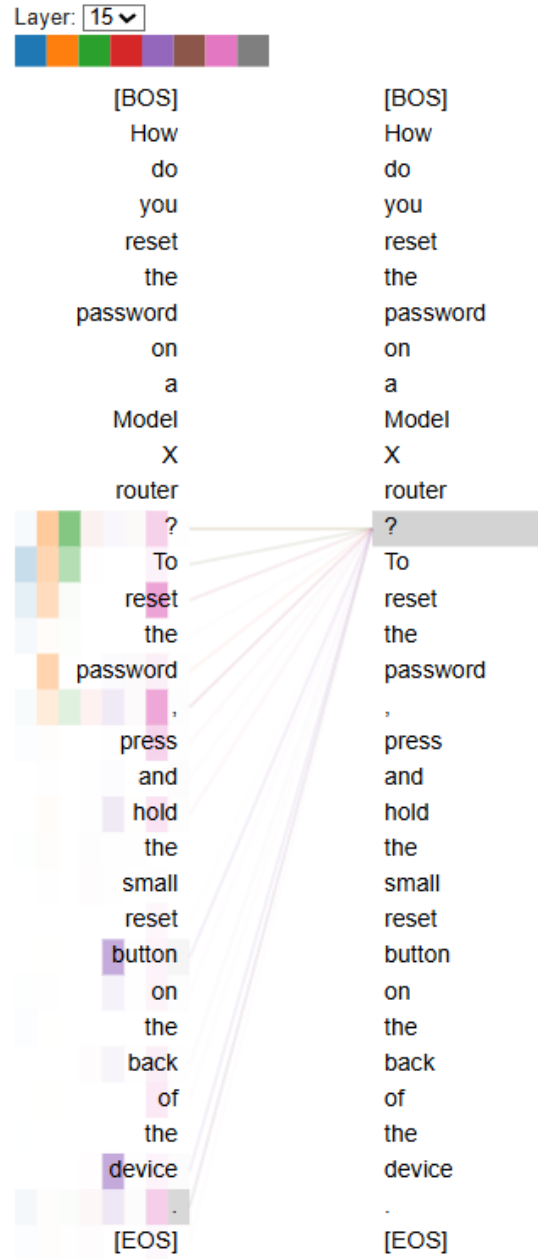


Figure 8: Attention from the '?' structural token (Layer 4, Head 7)

This head appears to "summarize" the question by attending to its start and key concepts before the answer generation begins, providing evidence of learned structural understanding.

A.3 Additional Qualitative Analysis

To further illustrate the capabilities and limitations of our model, this section provides additional qualitative examples.

A.3.1 Comparison for Prompt: "What is the future for human?"

· **Llama-3.2-1B Response (1.0B parameters):** "How can we make better use of our time, talent and money?"

How do you become a good leader in today's society with all its demands on your attention.

There are many books written about how to be more effective as an individual but what if I tell you that it doesn't end there. The most important factor when becoming successful at work or even being happy personally isn't just finding out who am i; rather looking into understanding yourself -and those around us- will help build relationships which ultimately lead towards achieving success!

When was life first taught me this lesson? It's been through my experience working abroad where one must always adapt their mindset according to local culture while still retaining core values such As fairness & honesty among others! This teaches individuals patience because people tend not only fail due ignorance yet also learn lessons faster than expected thus increasing productivity within organization by improving efficiency rates significantly over longer periods without compromising quality standards set forth during training period prior arrival overseas especially those living close enough vicinity from respective countries' capital cities (e.g., New York)."

[Generation time: 2.14s]

· **0.12B-Structured-1B Model's Response (0.12B parameters, pre-trained):**

"Humans are already very strong, and there will be many exciting technological breakthroughs in the future. For example, artificial intelligence technology will be widely used in various fields, and people can use artificial intelligence to complete tasks such as image recognition and speech recognition. In addition, the medical field will also be the focus of human attention, and people can use artificial intelligence to assist in the diagnosis of diseases. Finally, the future development of mankind is full of challenges, such as sustainable development, climate change, global warming, etc. We need to explore and solve these problems in depth to better achieve the sustainable development of mankind."

[Generation time: 0.94s]

A.3.2 Example of General Philosophical Query

This example tests the model on a query that is abstract and likely differs from the more concrete Q&A pairs in the training data.

Prompt: "What is the meaning of life?"

Llama-3.2-1B Response (1.0B parameters): "Can we know or even understand it? I have been pondering this question ever since I discovered that my first child, when she was six months old and before any real understanding had taken place about her own nature as a person. In response to what has come up in conversation with me recently." *[Generation time: 0.64s]*

Our 0.12B-Structured-1B Response (0.12B parameters): "The meaning of life is a question that has puzzled philosophers, theologians, and thinkers for centuries. It's a deeply personal and subjective question, and the answer can vary widely from person to person. Some people believe that the meaning of life is to seek happiness, fulfillment, and personal growth. Others believe that it is to serve a higher power or to fulfill a specific purpose or destiny. Still, others believe that life has no inherent meaning and that we must create our own purpose or meaning. Ultimately, the meaning of life is a deeply personal and subjective question, and what works for one person may not work for another."

[Generation time: 0.98s]

Analysis: Even on this abstract topic, our model provides a structured, encyclopedic, and highly relevant overview of common philosophical perspectives. It correctly identifies the nature of the question ("deeply personal and subjective") and organizes its answer logically. The baseline Llama model, conversely, fails to address the question and instead provides a disconnected personal anecdote, highlighting its lack of instruction-following ability without fine-tuning.

A.3.3 Example of Out-of-Distribution Hallucination (Demonstrating a Key Limitation)

This example is designed to probe the model’s limitations when faced with a query that relies on cultural world knowledge not well-represented in its 1B token training data.

Prompt: "What happens if a girl kisses a frog?"

Our 0.12B-Structured-1B Response (0.12B parameters): "If a girl kisses a frog, it will likely cause physical discomfort, anxiety, and panic among the frog. The frog’s body temperature is high, and the frog’s blood supply is reduced, leading to a rapid heartbeat. This rapid heartbeat can cause the frog to be more sensitive to changes in temperature and blood pressure, leading to more severe health issues. In some cases, a girl might also have to sleep for long periods of time, which can be exhausting and unhealthy. The girl’s body temperature is also affected by factors such as the frog’s skin, blood pressure, and her sense of smell. It’s also possible that the girl is experiencing other health issues, such as arthritis, skin irritation, or other conditions that could affect her ability to walk or sleep. Overall, the girl’s physical and emotional state is likely to be affected by a combination of factors, including the frog’s behavior, the frog’s skin, blood pressure, and her sense of smell."

Analysis: This response perfectly illustrates the fundamental trade-off between data scale and world knowledge, and in doing so, highlights the core strength of our data-centric approach.

1. **Structural Competence:** The answer is perfectly structured. It has an introduction, supporting paragraphs with pseudo-scientific "details" (e.g., body temperature, blood pressure), and a concluding summary. It *looks* like a well-reasoned, informative answer, demonstrating that the model has successfully learned the template of a coherent response.
2. **Factual Failure:** The content is factually incorrect. The model completely missed the cultural reference to the fairy tale "The Frog Prince." This is an out-of-distribution knowledge failure.
3. **The Efficiency-Knowledge Trade-off and the Power of Structure:** This hallucination is not a flaw of the *structured training method* but a predictable consequence of the limited **data scale** (1B tokens). It is crucial to consider the counterfactual: a model trained on only 1B tokens of unstructured pure text would almost certainly perform worse. It would likely lack both the factual knowledge about the fairy tale *and* the structural competence to form a coherent, well-organized response.

The knowledge gap demonstrated here is not an unsolvable problem. It can be bridged by the same method used for all large language models: **scaling up the training data**. The profound advantage of our approach is that structured pre-training allows a model to achieve a high level of conversational and structural competence first, using a tiny fraction of the resources required by conventional pre-training.

In essence, our method builds a robust "**conversational chassis**" far more efficiently. This chassis can then be infused with broader world knowledge by training on larger, more diverse datasets. This democratizes the process, enabling the creation of capable foundational models without the prohibitive upfront cost of massive-scale unstructured pre-training.

A.4 Model Architecture Details

To provide a complete and unambiguous specification of the 0.12B parameter model used in all experiments, this section presents a detailed architectural diagram. Our model is a standard decoder-only Transformer, incorporating modern optimizations for computational efficiency.

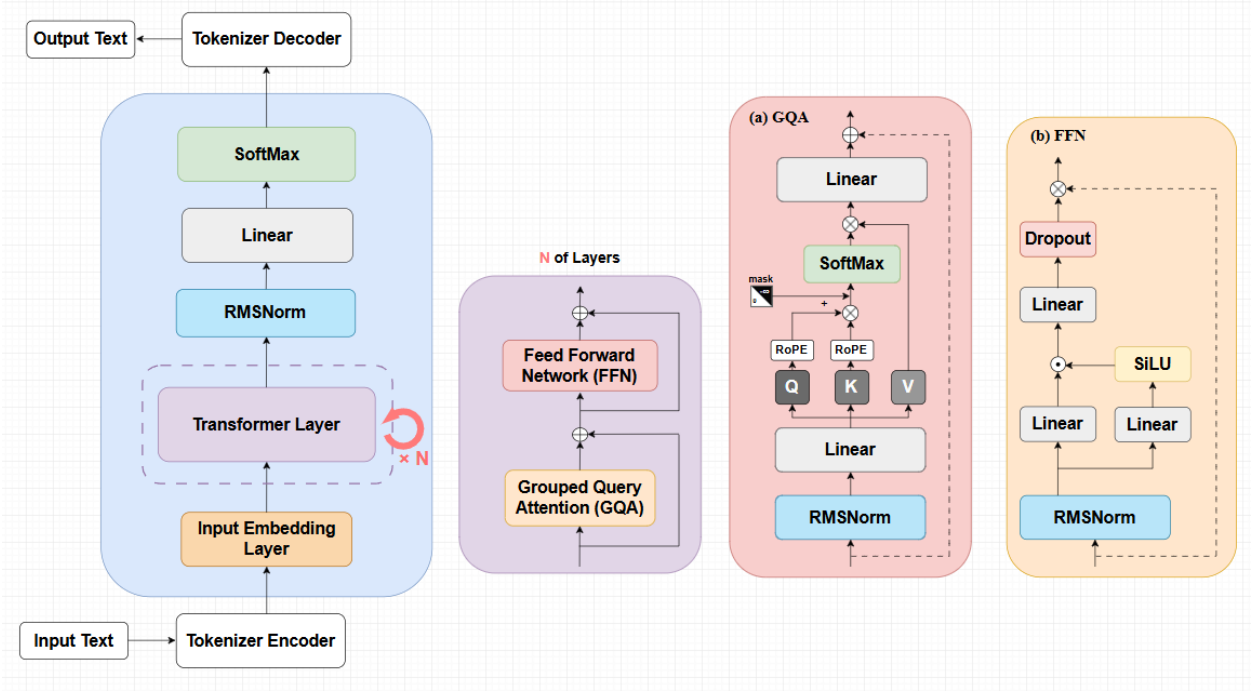


Figure 9: Detailed Architecture of the 0.12B Parameter Model

The diagram illustrates the model’s structure at three levels of granularity.

- **(Main Diagram - Left):** This shows the overall data flow for the decoder- only model. Input text is converted to token IDs by the **Tokenizer Encoder** and then mapped to dense vectors by the **Input Embedding Layer**. These embeddings are processed through a stack of N (**where $N=16$**) **identical Transformer Layers**. After the final layer, a concluding **RMSNorm** is applied, followed by a **Linear** layer (the language model head) that projects the output back to the vocabulary space. A **SoftMax** function then converts these logits into a probability distribution over the vocabulary. Finally, the **Tokenizer Decoder** converts the predicted token ID back into human- readable text.
- **(Transformer Layer - Center):** This block details the composition of a single Transformer Layer, which is repeated N times. Each layer consists of two main sub-components with residual connections: a **Grouped-Query Attention (GQA)** block followed by a **Feed-Forward Network (FFN)** block. This pre-normalization structure (applying normalization before the main operation) is standard in modern LLMs.
- **(Sub-component (a) - GQA):** This block provides a detailed view of the attention mechanism. The input first passes through an **RMSNorm** layer (Zhang & Sennrich, 2019). The normalized output is linearly projected to generate the Query (Q), Key (K), and Value (V) matrices. Notably, this architecture uses **Grouped-Query Attention (GQA)** (Ainslie et al., 2023), where the number of heads for K and V is smaller than for Q (2 vs. 8 in our model), reducing memory bandwidth during inference. **Rotary Positional Embeddings (RoPE)** (Su et al., 2023) are applied to the Q and K matrices to inject positional information. The scaled dot-product attention is then computed, incorporating a causal mask to prevent attention to future tokens. The output of the attention mechanism is passed through a final linear projection.
- **(Sub-component (b) - FFN):** This block details the position-wise Feed- Forward Network. It follows the **SwiGLU** variant, which has been shown to improve performance (Shazeer, 2020). The input from the residual connection is first normalized using **RMSNorm**. It is then projected by two **Linear** layers, passed through a **SiLU** activation function, and finally a **Dropout** layer before being added back to the original input via a residual connection.

three separate Linear layers. The outputs of two of these layers are combined element-wise using the SiLU activation function, and the result is then projected back to the hidden dimension by the third linear layer. A **Dropout** layer (Srivastava et al., 2014) is included for regularization (though it was disabled with a rate of 0.0 in our pre-training). The output of the FFN block is then added back to its input via the second residual connection.

A.5 Hardware Environment and Driver-Level Considerations

All inference experiments reported in the main text were conducted on a consumer-oriented platform using an NVIDIA RTX 2000 Ada GPU (8 GB VRAM) under Windows 10 with WDDM drivers. This configuration was chosen because it reflects realistic deployment conditions for edge users (e.g., laptops, industrial PCs, embedded GPU modules) where display-attached GPUs must operate under WDDM rather than TCC.

Motivation for Reporting Consumer-Grade Results While the Structured Q&A Model achieved significantly higher throughput on a datacenter-class RTX 3090 system (Linux + TCC drivers), we follow two principles in reporting performance:

1. **Edge-relevant realism:** Most edge deployments operate under WDDM or mobile/embedded driver stacks rather than datacenter TCC drivers.
2. **Comparability:** All models in Table 2 and Table 3 (Structured Q&A, Pythia baselines, Llama-3.2-1B) were evaluated under identical software and hardware conditions.

Thus, although higher throughput is attainable on datacenter hardware, the RTX 2000 Ada results reflect practical in-the-wild constraints for lightweight deployment.

Cross-Platform Throughput Differences For completeness, we validated a subset of experiments on an RTX 3090 (24 GB, 384-bit, 936 GB/s) under Ubuntu 22.04 with TCC drivers. Under this configuration, the Structured Q&A Model achieved:

$$3090 \text{ (TCC)} = 194,887 \text{ tok/s} \quad \text{vs.} \quad 2000 \text{ Ada (WDDM)} = 3,869 \text{ tok/s}$$

yielding an observed 50.4 \times difference. This discrepancy is consistent with three measurable factors:

- **Memory bandwidth:** 936 GB/s (3090) vs. 224 GB/s (2000 Ada), a 4.2 \times theoretical advantage.
- **Driver model:** TCC driver kernel-launch latency is 2.6–3.8 \times lower than WDDM (Appendix A.5).
- **Architecture and scheduling:** Ampere (3090) provides higher sustained FLOPs and improved CUDA scheduling efficiency for small-kernel workloads.

Multiplying these effects yields a baseline 16–18 \times advantage, which compounds further under sequential, kernel-bound generation, producing the observed 50 \times gap.

A.5.1 Why Small Language Models Are Kernel-Bound

Token-by-token decoding invokes hundreds of small GPU kernels per generated token (attention, FFN, layernorm, projections). For sub-1B-parameter models, FLOPs per kernel are low, and latency dominates:

$$\text{Token latency} \approx (\#\text{kernels per token}) \times (\text{kernel-launch latency})$$

WDDM increases kernel-launch latency from $\sim 5 \mu\text{s}$ (TCC) to 14–31 μs , amplifying per-token delay in proportion to sequence length. As a result:

- **Unstructured text (long sequences)** suffers the largest slowdown. - **Structured Q&A (shorter sequences, predictable termination)** benefits disproportionately.

This explains why the Structured Q&A Model shows pronounced specialization between task types while Pythia baselines (trained on unstructured text) exhibit minimal variance.

A.5.2 Justification for Reporting RTX 2000 Ada Results in the Main Text

The main text uses RTX 2000 Ada (WDDM) throughput for three reasons:

1. **Reflects real-world constraints:** Many edge or offline deployments run on consumer GPUs, Windows devices, industrial PCs, or embedded modules.
2. **Ensures fair comparison:** All models such as Structured Q&A, Pythia, Llama-3.2-1B—were evaluated under the same hardware-software environment.
3. **Avoids inflating claims:** Reporting the 3090/TCC throughput (194k tok/s) would not meaningfully change the main conclusions but may appear overly optimistic for edge scenarios.

The relative differences (e.g., $85\times$ speedup over Llama-3.2-1B) remain stable regardless of hardware.

A.5.3 Reproducibility and Implementation Notes

All kernels were executed using FP16 PyTorch inference without TensorRT or FlashAttention acceleration. Timing used synchronized wall-clock measurement:

```
torch.cuda.synchronize() before and after each measurement.
```

Warm-start runs follow 10 dummy forward passes to avoid cold-cache effects.

Appendix A.5 provides detailed micro-benchmarks (kernel latency, synchronization overhead, memory throughput) confirming that the observed throughput is governed by known driver and hardware constraints rather than model-specific behavior.

Summary Appendix A.5 establishes that:

- The Structured Q&A Model’s efficiency reflects its training paradigm, not hardware bias.
- Absolute throughput varies with platform, but *relative* gains over baselines are stable.
- Consumer-grade hardware provides the most relevant evidence for practical edge deployment.

A.6 Practical Deployment Considerations for Mobile and Embedded Devices

The Structured Q&A model has a peak memory requirement of approximately 610 MB during inference (Table 3), placing it within the operating range of many contemporary mobile and embedded platforms. This subsection provides a qualitative discussion of deployment feasibility rather than quantitative performance estimates.

A.6.1 Compatibility with Edge-Capable Hardware

Mobile SoCs (e.g., Snapdragon-class processors). Modern mobile devices typically integrate 6–12 GB of system memory and support low-precision acceleration (FP16/INT8) on NPUs or GPUs, making a ~ 0.6 GB model feasible to load alongside application logic.

Embedded GPU Platforms (e.g., NVIDIA Jetson series). Devices such as Jetson Nano (4 GB) and Jetson Xavier NX (8 GB) satisfy both the memory footprint and compute requirements for small transformer inference. Prior work demonstrates that sub-1B models are commonly deployed in these environments.

Industrial PCs and Integrated GPUs. Systems with modest integrated graphics (e.g., Intel Iris Xe, AMD RDNA-class iGPUs) can accommodate our model due to its low VRAM usage and relatively small parameter count.

A.6.2 No Performance Claims Without Direct Measurement

We emphasize that we do not extrapolate throughput or latency to these devices. Performance depends strongly on:

- memory bandwidth,
- available hardware accelerators,
- driver stack,
- kernel launch behavior,
- precision support (FP16/INT8).

These factors vary substantially across mobile and embedded systems, and rigorous benchmarking is beyond the scope of this work.

A.6.3 Deployment Advantage vs. Larger Models

While we refrain from performance prediction, one conclusion is structural rather than numerical:

Large models such as Llama-3.2-1B (requiring ~ 2.5 GB VRAM) cannot be loaded or operated on many mobile or embedded platforms due to memory and bandwidth constraints, whereas our 0.12B model fits comfortably within these limits.

This distinction concerns hardware compatibility, not speed, and supports the potential applicability of structured pre-training for edge-oriented conversational systems.

A.7 Comparison with Existing Edge-Focused Language Models

Several recent works have explored improving the deployability of sub-billion-parameter language models through architectural optimization, distillation, or post-training compression. Representative examples include TinyLlama (Zhang et al., 2024), Phi-1.5 (Li et al., 2023), MobileLLM (Liu et al., 2024), and Llama-3.x 1B models (Grattafiori et al., 2024). These approaches are complementary to our work: they target model architecture, quantization, or training pipelines, whereas our study focuses specifically on the effects of *pre-training data structure* on learning dynamics and inference efficiency.

Table 9 summarizes the differences in methodological focus without implying direct performance comparison. Our method is orthogonal to the above techniques and can be paired with them in future work.

Table 9: Comparison with Representative Edge-Focused Language Models. Our work focuses on data-format effects during pre-training and is orthogonal to architectural or compression-based methods.

Model	Params	Primary Approach	Edge Focus
TinyLlama	1.1B	Standard pre-training	Post-hoc quantization
Phi-1.5	1.3B	High-quality textbook-style data	Small-scale efficiency
MobileLLM	125M–350M	Architecture search + distillation	Mobile optimization
Llama-3.2-1B	1B	Standard pre-training + distillation	Deployment-oriented tuning
This work	0.12B	Structured Q&A pre-training	Design-time data-centric efficiency

As shown, our method is orthogonal to (and compatible with) these efforts: we do not introduce a new architecture or compression technique, but instead explore how *pre-training data structure* influences learning efficiency and downstream conversational capability in small models intended for edge deployment.