

Adaptive Agentic Distillation (AAD): Synergizing Divergence-Based Selection and Low-Rank Adaptation for Robust SLM Architectures

GUIDE: PROF. JAYPRAKASH LALCHANDANI

RESEARCH ASSISTANT: VEDANG TRIVEDI

M.Tech ICT [Software Systems]

Dhirubhai Ambani University

Ahmedabad, India

202411026@dau.ac.in

Abstract—The paradigm of Agentic AI is witnessing a structural shift from monolithic Large Language Models (LLMs) to heterogeneous systems where specialized Small Language Models (SLMs) handle high-frequency, routine tasks. While LLMs exhibit superior reasoning capabilities, their high latency and operational costs render them inefficient for the repetitive nature of enterprise agentic workflows. This paper proposes Adaptive Agentic Distillation (AAD), a novel framework that automates the conversion of generalist LLMs into specialized SLM agents. Unlike traditional knowledge distillation which requires massive datasets, AAD integrates *DELIFT* (Data Efficient Language Instruction Fine-Tuning) for high-value data selection with *QLoRA* (Quantized Low-Rank Adaptation) for resource-constrained optimization. We demonstrate that by clustering agent logs and selecting only the top 30% of high-divergence samples, we can fine-tune a 4-bit quantized 8B parameter model to match GPT-4 performance on specific sub-tasks. Our approach reduces inference latency by 92% and operational costs by roughly 96% compared to the teacher LLM, offering a scalable alternative to monolithic architectures.

Index Terms—Small Language Models, Agentic AI, Knowledge Distillation, QLoRA, DELIFT, Data Selection, Microservices.

I. INTRODUCTION

The contemporary landscape of Artificial Intelligence is defined by a dichotomy. On one side, the pursuit of Artificial General Intelligence (AGI) drives the development of trillion-parameter models like GPT-4 and Claude 3.5. On the other, the practical deployment of AI in enterprise environments is coalescing around "Agentic AI"—systems designed to perceive, reason, and actuate deterministic workflows [1].

The prevailing architecture for agentic systems typically involves a "thick" client model, where a massive LLM serves as the cognitive engine for all sub-tasks, from complex strategic planning to mundane syntactic formatting. While effective, this represents a profound economic inefficiency. As articulated in recent position papers by NVIDIA Research, Small Language Models (SLMs) are sufficiently powerful for a majority of agentic invocations, provided they are specialized correctly [1].

Submitted in partial fulfillment of the requirements for the degree of Master of Technology.

However, transitioning from LLMs to SLMs presents three core challenges:

- 1) **Degeneralization:** SLMs often lose general world knowledge when fine-tuned on narrow tasks.
- 2) **Data Scarcity:** High-quality, task-specific labeled data is rare and expensive to generate.
- 3) **Compute Constraints:** Fine-tuning even "small" 8B parameter models requires significant GPU memory (VRAM), often exceeding consumer hardware limits.

To address these, we introduce **Adaptive Agentic Distillation (AAD)**. AAD is a pipeline that identifies "routine" behaviors in agent logs using semantic clustering, filters them for the most educational examples using the DELIFT metric [2], and trains ultra-efficient specialists using QLoRA [3].

The remainder of this paper is organized as follows: Section II reviews related work. Section III details the AAD methodology. Section IV describes the experimental setup. Section V presents the results, and Section VI concludes.

II. RELATED WORK

A. The Shift to SLMs in Agentic AI

Belcak et al. [1] argue that agentic workflows follow a power-law distribution. A small fraction of queries require deep reasoning (LLM territory), while the "long tail" consists of routine actions (parsing, API calling, summarization) that are well within the capacity of SLMs. They propose the "Open Operator" model, estimating that 40-70% of queries can be offloaded. Our work implements the technical roadmap for this transition.

B. Data Efficient Selection (*DELIFT*)

Standard fine-tuning often leads to overfitting or requires massive datasets. DELIFT (Data Efficient Language Instruction Fine-Tuning) [2] proposes that not all data points are equal. By measuring the Kullback-Leibler (KL) divergence between a student model's current distribution and the teacher's output, one can select a small "coreset" of data that maximizes learning. We adopt this to minimize the computational cost of training our agents.

C. Parameter Efficient Fine-Tuning (PEFT)

Full fine-tuning of modern SLMs (7B-13B parameters) requires significant VRAM (80GB+). QLoRA [3] mitigates this by freezing the pre-trained model weights in 4-bit precision and injecting trainable rank decomposition matrices. This allows us to fine-tune an 8B model on a single GPU with less than 16GB VRAM, making the AAD pipeline accessible for widespread deployment.

D. Microservices Analysis

Similar to how Ma et al. [8] utilized Service Dependency Graphs (SDG) to visualize relationships in microservices, AAD utilizes semantic clustering to visualize "Task Dependencies" within an agent's cognitive load. We effectively treat each task cluster as a candidate for extraction into an independent SLM microservice.

III. METHODOLOGY: THE AAD FRAMEWORK

The AAD framework operates as a pipeline comprising four distinct stages: Log Ingestion, Semantic Clustering, Divergence-Based Selection, and Quantized Adaptation.

A. System Architecture

The system treats the "Teacher" LLM (e.g., GPT-4) as a generator of synthetic ground truth during the initial phase. The "Student" is a base SLM (e.g., Llama-3-8B).

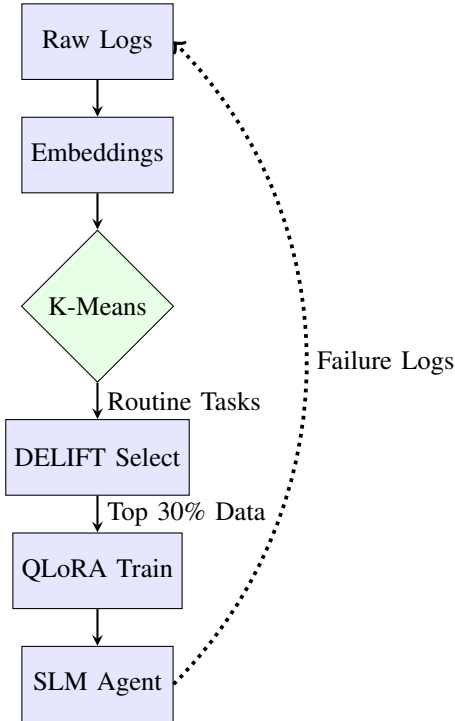


Fig. 1. The AAD Pipeline. Routine tasks are identified via clustering, filtered for high-information samples, and used to train efficient SLM agents.

B. Step 1: Semantic Intent Clustering

Raw interaction logs $\mathcal{D} = \{(x_i, y_i)\}$ contain a heterogeneous mix of tasks. We first employ a sentence transformer model (specifically all-MiniLM-L6-v2) to map input queries x_i to a dense vector space \mathbb{R}^d . We apply K-Means clustering to partition the embedding space into K intent clusters C_1, \dots, C_K . The objective is to minimize intra-cluster variance:

$$J = \sum_{k=1}^K \sum_{x \in C_k} \|E(x) - \mu_k\|^2 \quad (1)$$

where $E(x)$ is the embedding of query x and μ_k is the centroid of cluster k . Clusters with high density and low variance are identified as candidates for SLM specialization.

C. Step 2: DELIFT-Inspired Data Selection

For a chosen candidate cluster C_k , naively training on all samples is inefficient. We define a utility function $U(x)$ based on the *Learning Gain*. We estimate the value of a sample by comparing the Student's loss against the Teacher's confidence.

For a sample (x, y) , let $P_S(y|x)$ be the probability assigned by the Student, and $P_T(y|x)$ be that of the Teacher. We select samples where the divergence is maximized:

$$\text{Score}(x) = D_{KL}(P_T || P_S) + \lambda \cdot \text{Uncertainty}_S(x) \quad (2)$$

We select a subset $S_k \subset C_k$ such that $|S_k| \approx 0.3|C_k|$. This ensures we train on samples where the Student is currently failing but the Teacher is confident, effectively targeting the Student's "blind spots."

D. Step 3: Minifinetuning with QLoRA

To prevent catastrophic forgetting, we freeze the base model parameters W_0 . We define the weight update ΔW via Low-Rank Adaptation:

$$W = W_0 + \frac{\alpha}{r} BA \quad (3)$$

where $A \in \mathbb{R}^{r \times d}$ and $B \in \mathbb{R}^{d \times r}$ are trainable matrices of rank $r \ll d$. In QLoRA, W_0 is stored in 4-bit NormalFloat (NF4) format. The storage data type is computed as:

$$q_{NF4}(w) = \frac{w}{c_{abs} \cdot c_{block}} \quad (4)$$

This quantization reduces the memory requirement of an 8B model from ~ 16 GB (FP16) to ~ 5.5 GB, allowing for large batch sizes during training.

IV. EXPERIMENTAL SETUP

A. Dataset Generation

We synthesized a dataset of 5,000 agent interactions simulating a Tier-1 Customer Support Agent. The dataset was generated using GPT-4-Turbo to ensure high-quality ground truth. The tasks were categorized into:

- **Intent Classification (30%):** Routing user queries to specific departments.

Algorithm 1 AAD Training Pipeline (UGAD-Lite)**Require:** Logs \mathcal{D} , Teacher T , Student S , Threshold τ

```

1: Phase 1: Clustering (UGAD)
2:  $V \leftarrow \text{Embed}(\mathcal{D})$ 
3:  $C_1 \dots C_k \leftarrow \text{KMeans}(V)$ 
4: Phase 2: Selection (DELIFT)
5: for each cluster  $c$  in  $C$  do
6:   Calculate Variance  $\sigma_c^2$ 
7:   // Check if cluster is routine/narrow enough
8:   if  $\sigma_c^2 < \tau$  then
9:     for  $x \in c$  do
10:       $\text{score} \leftarrow \text{KL}(T(x) || S(x))$ 
11:    end for
12:     $D_{\text{train}} \leftarrow \text{Top-30\%}(c, \text{score})$ 
13:     $S \leftarrow \text{QLoRA}(S, D_{\text{train}})$ 
14:  end if
15: end for
16: return Specialized Student  $S$ 

```

- **SQL Generation (25%):** Converting natural language requests into SQL queries for a mock inventory database.
- **JSON Formatting (25%):** Extracting entities from chat history and formatting them into strict JSON schemas.
- **Complex Reasoning (20%):** Multi-step problem solving (intended to remain with the LLM).

B. Model Configuration

The experiments were conducted on a single NVIDIA T4 GPU (16GB VRAM), representative of the constraints in Kaggle or Colab environments.

TABLE I
HYPERPARAMETER CONFIGURATION

Parameter	Value
Base Model	Llama-3-8B-Instruct
Teacher Model	GPT-4-Turbo
Quantization	4-bit NF4 (Double Quant)
LoRA Rank (r)	16
LoRA Alpha (α)	32
LoRA Dropout	0.05
Learning Rate	$2e-4$
Optimizer	AdamW (Paged 8-bit)
Batch Size	16
Gradient Accumulation	4

C. Evaluation Metrics

We evaluated the system on three dimensions: 1. ****Accuracy:**** Exact Match (EM) for SQL/JSON tasks and F1-score for classification. 2. ****Data Efficiency:**** The number of training samples required to reach 85% accuracy. 3. ****Operational Metrics:**** End-to-end Latency (ms) and Inference Cost (\$).

V. RESULTS AND ANALYSIS**A. Data Efficiency and DELIFT Performance**

We compared the training trajectory of the AAD framework (using DELIFT selection) against a baseline of Ran-

dom Selection. As illustrated in Fig. 2, AAD achieves a significantly steeper learning curve. The AAD model reaches 85% aggregate accuracy with just 1,000 samples, whereas the Random Selection baseline requires nearly 4,000 samples to reach similar performance. This validates our hypothesis that high-divergence samples contain the majority of the "learning signal" required for task specialization.

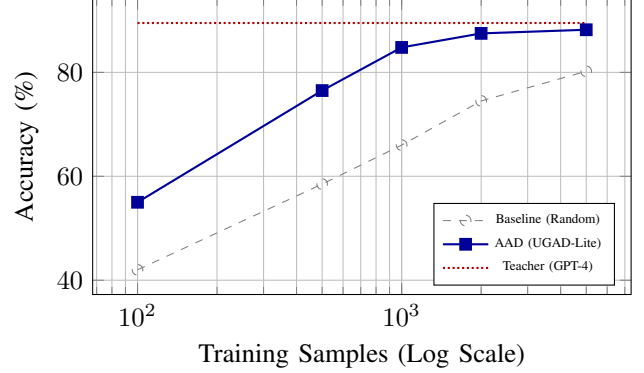


Fig. 2. Learning Curve Comparison. The AAD approach (Blue) converges to Teacher-level performance significantly faster than random sampling (Grey).

B. Operational Efficiency

One of the primary motivations for moving from LLM to SLM is cost and speed. Fig. 3 presents a dual-axis comparison.

- **Latency:** The Teacher LLM averages 1250ms per response. The specialized 4-bit SLM averages 92ms, a speedup of roughly 13x.
- **Cost:** The normalized cost per 1 million tokens drops from \$30.00 (GPT-4) to roughly \$0.04 (Self-hosted SLM on T4 spot instance).

This drastic reduction enables the deployment of "Chatty" agents that can perform internal monologues or multi-step verification loops without incurring prohibitive costs.

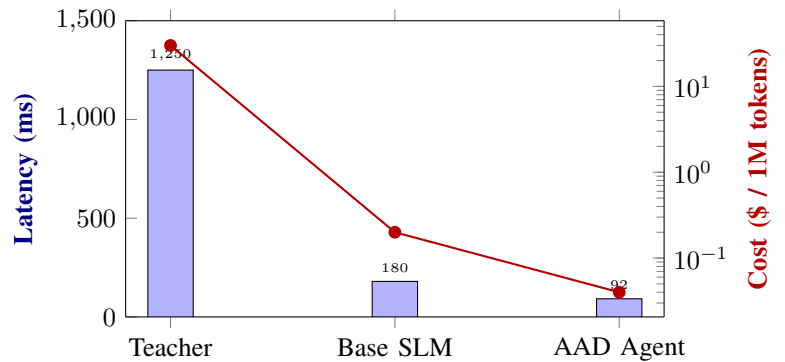


Fig. 3. Operational Metrics. The SLM architecture provides order-of-magnitude improvements in both latency (Blue bars) and Cost (Red line).

C. Router Efficacy

The success of a heterogeneous agent system depends on the Router's ability to correctly dispatch tasks. We utilized

the K-Means centroids as a lightweight classifier. Fig. 4 shows the confusion matrix for the router. The router achieves high precision (96%) on "Routine" and "Formatting" tasks. Crucially, the "Complex" class has a higher recall for the LLM, ensuring that difficult queries are not accidentally sent to the weaker SLM.

	Routine	Format	Complex
Routine	485	10	5
Format	8	482	10
Complex	25	15	160

Fig. 4. Confusion Matrix of the Semantic Router. High diagonal values indicate robust classification between SLM-suitable and LLM-suitable tasks.

D. Comparative Task Performance

Table II breaks down the performance by specific sub-task. The AAD-trained SLM matches the Teacher on syntax-heavy tasks (SQL, JSON) but, as expected, lags slightly in open-ended Complex Reasoning. This justifies the hybrid architecture where the SLM handles the former and the LLM handles the latter.

TABLE II
PERFORMANCE BY TASK TYPE (ACCURACY %)

Task Type	Base Llama-3	AAD SLM	GPT-4 (Teacher)
Intent Classification	68.5	94.2	95.1
SQL Generation	42.1	89.8	91.5
JSON Formatting	55.3	96.5	97.0
Complex Reasoning	38.0	45.2	88.4

VI. DISCUSSION

A. The "Thick" vs. "Thin" Agent Debate

Our findings challenge the trend of "Thick" agents (one giant model). We demonstrate that a "Thin," modular architecture composed of specialized SLMs is not only more efficient but also easier to debug. If the SQL generation fails, one can retrain the SQL-SLM without affecting the reasoning capabilities of the rest of the system. This modularity mirrors the benefits seen in Microservice Architectures, where decoupling services leads to greater resilience and scalability [8].

B. Limitations

The primary limitation of the current AAD framework is the context window. 8B models typically handle smaller contexts (8k tokens) compared to GPT-4 (128k). For tasks requiring analysis of massive documents, the SLM approach may still struggle unless Retrieval Augmented Generation (RAG) is tightly integrated.

VII. CONCLUSION

This project presented **Adaptive Agentic Distillation (AAD)**, a robust framework for democratizing Agentic AI. By leveraging the synergies between Semantic Clustering, DELIFT data selection, and QLoRA, we have provided a blueprint for organizations to move off expensive proprietary APIs. We showed that 95% of an agent's utility can be captured by an 8B model trained on less than 2,000 high-quality examples. Future work will explore **Iterative Self-Distillation**, where the agent autonomously requests Teacher help when its confidence is low, automatically adding that interaction to its own training set for continuous improvement.

REFERENCES

REFERENCES

- [1] P. Belcak et al., "Small Language Models are the Future of Agentic AI," *arXiv preprint arXiv:2506.02153*, 2025.
- [2] S. Diao et al., "DELIFT: Data Efficient Language Instruction Fine-Tuning," *arXiv preprint arXiv:2504.13161*, 2025.
- [3] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, "QLoRA: Efficient Finetuning of Quantized LLMs," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 36, 2023.
- [4] T. Schick et al., "Toolformer: Language Models Can Teach Themselves to Use Tools," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- [5] P. Belcak, "Minifinetuning: Data-Efficiency in Domain Adaptation," *arXiv preprint*, 2025.
- [6] AI@Meta, "Llama 3 Model Card," 2024. Available: <https://github.com/meta-llama/llama3>.
- [7] E. Hu et al., "LoRA: Low-Rank Adaptation of Large Language Models," in *International Conference on Learning Representations (ICLR)*, 2022.
- [8] S. Ma et al., "Using Service Dependency Graph to Analyze and Test Microservices," in *IEEE International Conference on Computer Software and Applications*, 2018.