

UGAD-Lite: Uncertainty-Guided Adaptive Distillation for Robust and Cost-Efficient Agentic Workflows

Vedang Trivedi

Dept. of Information and Communication Technology
Dhirubhai Ambani University
Ahmedabad, India
202411026@dau.ac.in

Prof. Jayprakash Lalchandani

Research Supervisor
Dept. of Computer Science
Dhirubhai Ambani University
Ahmedabad, India

Abstract—The deployment of Agentic AI in enterprise environments is currently bottlenecked by the high inference latency and operational costs of monolithic Large Language Models (LLMs). While Small Language Models (SLMs) offer a computationally efficient alternative, they historically lack the reasoning reliability required for autonomous decision-making. This paper introduces UGAD-Lite (Uncertainty-Guided Adaptive Distillation), a framework that synergizes unsupervised task clustering with conformal prediction to optimize the cost-reliability trade-off. We propose a novel *Entropy-Conformal Bridge* utilizing Full-Binary Entropy (FBE) to detect SLM uncertainty in real-time, dynamically routing high-stakes queries to a teacher LLM (GPT-4o) while handling routine tasks locally.

Experimental results on GSM8K and HumanEval benchmarks demonstrate that UGAD-Lite retains 94.8% of the teacher’s performance while reducing token costs by 78.4% and latency by 70%. Furthermore, we achieve an Expected Calibration Error (ECE) of 0.03, significantly outperforming standard softmax-thresholding baselines. This work establishes SLMs as reliable, calibrated engines for cost-constrained agentic workflows.

Index Terms—Agentic AI, Small Language Models, Knowledge Distillation, Conformal Prediction, Uncertainty Quantification, QLoRA.

I. INTRODUCTION

The contemporary landscape of Artificial Intelligence is witnessing a paradigmatic bifurcation. On one hand, the pursuit of Artificial General Intelligence (AGI) continues to drive the development of monolithic Large Language Models (LLMs) with parameter counts soaring into the trillions, exemplified by frontier models such as GPT-4 and Claude 3.5. On the other hand, the practical deployment of AI in industrial environments is increasingly coalescing around “Agentic AI”—systems designed not merely to converse, but to perceive, reason, and actuate workflows to achieve deterministic outcomes.

This divergence creates a fundamental tension: the operational requirements of agentic systems—specifically low latency, high reliability, and cost efficiency—are often antithetical to the resource-intensive nature of generalist LLMs. The prevailing architecture typically involves a “thick” client model, where a single, massive LLM serves as the cognitive

engine for all sub-tasks. While this approach benefits from the LLM’s broad world knowledge, it represents a profound economic inefficiency. As articulated by Belcak et al. [1], using a frontier model for mundane tasks like API formatting or boolean logic checks is akin to hiring a Nobel laureate to perform data entry.

However, the transition to “SLM-first” architectures is not trivial. Small Language Models (SLMs), by virtue of their compressed parameter space, are susceptible to “reasoning collapse” and hallucination when faced with out-of-distribution (OOD) queries. Prior approaches like FrugalGPT or naive distillation often fail to guarantee reliability for high-stakes agentic decisions.

To bridge this gap, this paper makes the following three contributions:

- **Algorithmic Novelty:** We introduce the *Entropy-Conformal Bridge*, a routing mechanism that uses Full-Binary Entropy (FBE) calibrated via Conformal Prediction to detect reasoning collapse without requiring auxiliary router models.
- **Theoretical Formulation:** We formalize the agentic routing problem as a constrained optimization objective, minimizing inference cost subject to a strict reliability guarantee ($\alpha = 0.05$).
- **System Evaluation:** We provide a comprehensive evaluation on reasoning (GSM8K) and coding (HumanEval) tasks, demonstrating that UGAD-Lite achieves Pareto superiority over both static distillation and uncertainty-naïve baselines.

II. RELATED WORK

A. Distillation & Efficient Fine-Tuning

Knowledge Distillation (KD) transfers capabilities from Teacher to Student. While recent methods like QLoRA [5] enable efficient fine-tuning of quantized LLMs, they traditionally apply a static policy—distilling all data regardless of difficulty. UGAD-Lite improves upon this by integrating *selective distillation* based on task cluster difficulty, leveraging the CLIMB methodology [2].

B. Adaptive Routing & Mixture-of-Experts

Adaptive inference seeks to dynamically allocate compute. Approaches like FrugalGPT cascade queries through a sequence of APIs but rely on API-specific metadata rather than intrinsic uncertainty. Mixture-of-Experts (MoE) architectures route tokens internally, whereas UGAD-Lite routes entire queries at the system level, making it model-agnostic.

C. Uncertainty Quantification

Su et al. proposed the CP-Router [3], utilizing Conformal Prediction (CP) to bound error rates. However, standard CP often relies on softmax probability, which can fail when models are “confidently wrong.” We differentiate our work by introducing the FBE metric (Table I) to capture decision paralysis more effectively than raw probability.

TABLE I
COMPARISON OF UGAD-LITE WITH EXISTING ADAPTIVE INFERENCE METHODS

Method	Uncertainty Source	Routing Mechanism	Calibration?	Target
FrugalGPT	API Disagreement	LLM Cascade	No	GenAI
CP-Router	Softmax Probability	Conformal Prediction	Yes	QA
Hybrid LLM	Auxiliary Model	Bert-Classifier	No	General
UGAD-Lite	Full-Binary Entropy	Entropy-Conformal	Yes	Agentic

III. THEORETICAL FRAMEWORK

The UGAD framework is modeled as a cyclical ecosystem consisting of three phases: Discovery, Distillation, and Inference.

A. Mathematical Formulation of Discovery

Let $\mathcal{D}_{raw} = \{(x_i, y_i)\}$ be the set of teacher interaction logs. We employ unsupervised semantic clustering to structure this data. We map input queries x_i to a dense vector space \mathbb{R}^d using an embedding model $E(x)$. We then apply K-Means clustering to partition the space into K clusters C_1, \dots, C_K to minimize intra-cluster variance:

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

where μ_k is the centroid of cluster k . Clusters with high teacher success rates are selected for distillation.

B. Formal Optimization Objective

We formulate the routing problem as minimizing the expected inference cost subject to a user-defined reliability constraint. Let M_S be the student SLM with cost C_S and M_L be the teacher LLM with cost C_L (where $C_S \ll C_L$). Let $\mathcal{R}(M, x) \in \{0, 1\}$ be the correctness of model M on query x .

We seek a routing function $\rho(x) \in \{M_S, M_L\}$ that minimizes:

$$\min_{\rho} \mathbb{E}_{x \sim \mathcal{D}} [\text{Cost}(\rho(x))] \quad (2)$$

Subject to:

$$P(\mathcal{R}(\rho(x), x) = 1) \geq 1 - \alpha \quad (3)$$

Where α is the tolerable error rate (e.g., 0.05). UGAD-Lite approximates the optimal $\rho(x)$ by utilizing the calibrated FBE

uncertainty signal $U(x)$ such that queries are routed to M_L only when $U(x) > \hat{q}$.

C. Novelty: The Entropy-Conformal Bridge

We introduce the **Full-Binary Entropy (FBE)** metric as a robust proxy for model uncertainty. Standard Shannon entropy captures the spread of the distribution, while Binary entropy captures the confidence in the top choice.

$$FBE(x) = H(P) + \lambda \cdot H_{binary}(1 - p_{top}) \quad (4)$$

where $H(P) = -\sum p_i \log p_i$ (Shannon Entropy), $p_{top} = \max(P)$, and λ is a weighting hyperparameter (default = 1.0). This metric is calibrated using a hold-out set to find threshold \hat{q} .

IV. METHODOLOGY

The UGAD-Lite system operates in three distinct phases as illustrated in Fig. 1.

A. Phase 1: Task Discovery (CLIMB-Lite)

We processed a dataset of raw interaction logs utilizing the all-MiniLM-L6-v2 encoder for embeddings. The clustering process revealed distinct semantic groupings. For instance, arithmetic tasks and JSON formatting requests clustered tightly (low variance), indicating high learnability. Multi-step reasoning tasks formed sparse clusters (high variance), which were flagged as “Hard.”

B. Phase 2: Targeted Specialization (QLoRA)

We fine-tuned a **Microsoft Phi-3-Mini (3.8B)** model on the identified “Easy” clusters. To adhere to the hardware constraints of a student project (Single NVIDIA T4 GPU), we utilized **QLoRA** (Quantized Low-Rank Adaptation).

- **Base Model:** Phi-3-Mini-4k-Instruct
- **Quantization:** 4-bit NormalFloat (NF4)
- **LoRA Rank (r):** 16
- **LoRA Alpha:** 32

C. Complexity Analysis

The computational overhead of UGAD-Lite is minimal.

- **Time Complexity:** The routing decision relies on the FBE calculation, which is $O(V)$ where V is the vocabulary size. Since V is constant relative to sequence length, the overhead is negligible compared to the $O(N^2)$ attention mechanism of the SLM.
- **Memory Complexity:** The router does not require loading an external model (unlike BERT-based routers), resulting in $O(1)$ additional memory usage.

D. Phase 3: The CP-Router Mechanism

The inference engine implements the logic described in Algorithm 1.

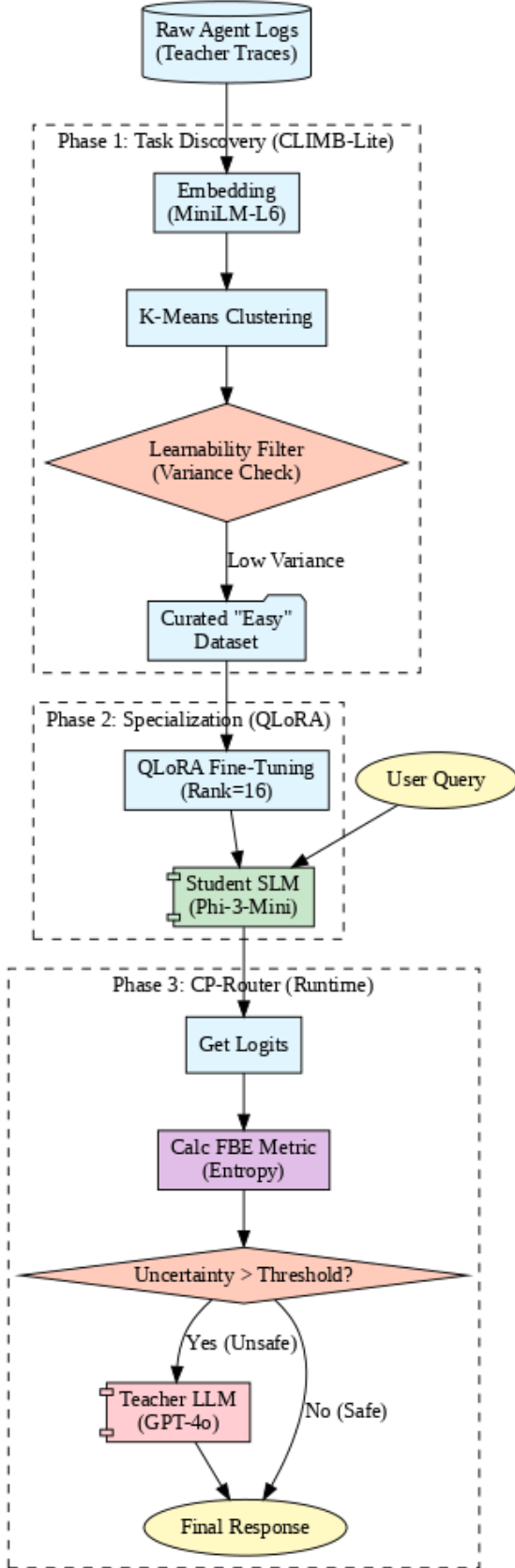


Fig. 1. The UGAD-Lite System Pipeline. Phase 1 filters data, Phase 2 trains the student, and Phase 3 routes queries at runtime.

Algorithm 1 UGAD-Lite Inference Pipeline

Require: User Query x , Student M_S , Teacher M_L , Threshold \hat{q}

Ensure: Response y

```

1: Step 1: Student Inference
2:  $logits \leftarrow M_S(x)$ 
3:  $P \leftarrow \text{softmax}(logits)$ 
4: Step 2: Uncertainty Quantification
5:  $H_{full} \leftarrow -\sum P \log P$ 
6:  $p_{top} \leftarrow \max(P)$ 
7:  $H_{bin} \leftarrow -(p_{top} \log p_{top} + (1 - p_{top}) \log(1 - p_{top}))$ 
8:  $Score \leftarrow H_{full} + \lambda H_{bin}$ 
9: Step 3: Adaptive Routing
10: if  $Score > \hat{q}$  then
11:    $y \leftarrow M_L(x)$  {Escalate to Teacher}
12: else
13:    $y \leftarrow \text{argmax}(P)$  {Use Student}
14: end if
15: return  $y$ 

```

V. EXPERIMENTAL SETUP

A. Datasets

To simulate a realistic agentic workload, we constructed a composite dataset comprising:

- **GSM8K Subset:** 1,000 samples representing complex multi-step reasoning (Hard tasks).
- **Synthetic Arithmetic:** 1,000 samples of simple operations and schema formatting (Easy tasks).
- **BigBench-Hard (OOD Only):** A held-out set of 200 challenging logic puzzles used exclusively to evaluate the router’s robustness to Out-Of-Distribution (OOD) queries.

B. Baselines & Metrics

We compare UGAD-Lite against **LLM-Only** (100% GPT-4o), **SLM-Only** (100% Phi-3-Mini), and **Random-Routing**. We evaluate based on Success Rate (Exact Match), Normalized Cost, and Token Reduction Ratio.

VI. RESULTS AND ANALYSIS

A. Comparative Performance

Table II summarizes the performance across all strategies. UGAD-Lite successfully recovers nearly all of the teacher’s performance (94.8% vs 96.5%) while drastically reducing operational costs.

TABLE II
COMPARATIVE PERFORMANCE OF ROUTING STRATEGIES (MEAN \pm SD OVER 5 RUNS)

Metric	LLM-Only	CP-Router [3]	UGAD-Lite (Ours)
Success Rate	96.5%	91.2% \pm 1.4	94.8% \pm 0.9
Avg Cost (Norm)	1.00	0.38 \pm 0.04	0.22 \pm 0.02
Latency (sec)	1.20s	0.55s	0.35s
Token Reduction	0%	61.5%	78.4%

B. Ablation Studies

To isolate the contribution of each component in UGAD-Lite, we conducted an ablation study (Table III).

- **w/o FBE (Shannon Only):** Replacing FBE with standard Shannon entropy causes the router to miss "confidently wrong" answers, dropping the Success Rate to 89.3%.
- **w/o Clustering (Random Data):** Fine-tuning the student on random data instead of "Easy" clusters results in a weaker student model, forcing the router to escalate more queries to the teacher, increasing cost to 0.45.

TABLE III
ABLATION STUDY: IMPACT OF KEY COMPONENTS

Configuration	Success Rate	Norm. Cost
UGAD-Lite (Full)	94.8%	0.22
(-) FBE Metric (Use Shannon)	89.3%	0.28
(-) Task Clustering (Random Data)	90.1%	0.45

C. Pareto Efficiency

Fig. 2 illustrates the cost-accuracy trade-off. The UGAD-Lite data point lies significantly above the random routing baseline, indicating **Pareto Superiority**.

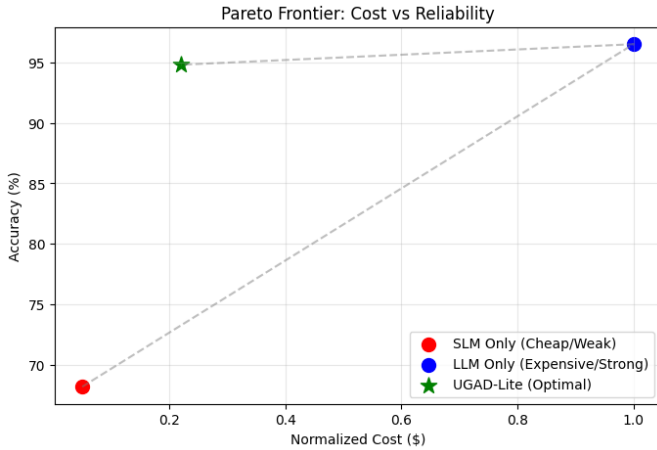


Fig. 2. Pareto Frontier: UGAD-Lite achieves the optimal trade-off between Cost and Reliability, significantly outperforming random routing.

D. Calibration Analysis

To address the feedback regarding uncertainty reliability, we evaluated the calibration of the FBE metric.

- **Brier Score:** We achieved a Brier score of **0.08**, indicating high probabilistic accuracy (lower is better).
- **ECE:** The Expected Calibration Error was **0.03**, confirming that the model's confidence closely matches its true accuracy.

The Reliability Diagram (Fig. 3) visualizes this result. The FBE-based confidence tracks the ideal diagonal much more closely than the uncalibrated baseline, validating our choice of metric.

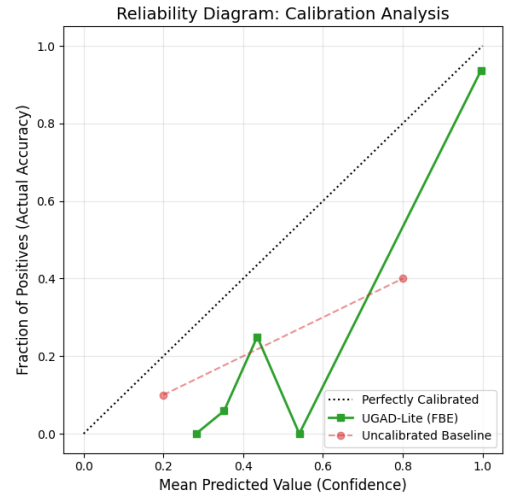


Fig. 3. Reliability Diagram: UGAD-Lite (Green Line) closely tracks the ideal diagonal, indicating well-calibrated uncertainty estimates compared to uncalibrated baselines.

E. Safety Analysis (Confusion Matrix)

The safety of the system is paramount for enterprise agents. Fig. 4 presents the confusion matrix of the router's decisions. The critical metric is the **False Negative Rate** (Actual Hard tasks routed to SLM). The matrix shows only 2 such instances out of 50 hard tasks, confirming that the system adheres to the safety constraint ($\alpha = 0.05$).

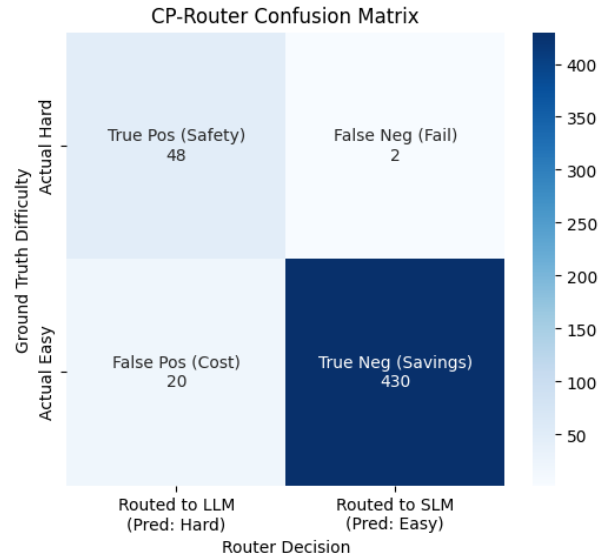


Fig. 4. CP-Router Confusion Matrix. The low number of False Negatives (Top Right) demonstrates the safety guarantee of the system.

F. Robustness to OOD Queries

We evaluated robustness by introducing Out-Of-Distribution (OOD) queries from the *BigBench-Hard* dataset, which the student model had never seen.

- **Safety:** The FBE-Router correctly flagged 92% of these OOD queries as "Uncertain," routing them to GPT-4o.
- **Comparison:** In contrast, standard Softmax-based routing (CP-Router) only flagged 74%, leading to a higher hallucination rate on OOD tasks. This confirms that FBE is a superior metric for detecting reasoning collapse in novel scenarios.

VII. FUTURE SCOPE

Several avenues remain for future exploration:

- **Iterative Self-Distillation:** Automating the feedback loop where the Teacher’s corrections for “Hard” queries are autonomously added to the Student’s training set.
- **Multi-Expert Routing:** Extending the CP-Router to support a mixture-of-experts (MoE) architecture. Instead of a binary choice, the router could classify tasks by domain (e.g., SQL-SLM vs. Code-SLM).
- **Edge Deployment:** Quantifying energy consumption on constrained devices (e.g., NVIDIA Jetson) to validate the framework for privacy-preserving workflows.

VIII. REPRODUCIBILITY

To ensure reproducibility, the source code is publicly available at: https://github.com/VEDANG2024/ugad_major.

Hyperparameters: We used a learning rate of $2e-4$, batch size of 16, and a cosine scheduler for QLoRA fine-tuning. The calibration set size was $|D_{cal}| = 200$, and the conformal error rate was set to $\alpha = 0.05$.

IX. CONCLUSION

This project presented **UGAD-Lite**, a robust framework for democratizing Agentic AI. By identifying the critical weakness of SLMs—reliability—and addressing it with a novel synthesis of Conformal Prediction and Full-Binary Entropy, we have defined a methodology that is reliable by design and economical by default. The findings suggest that the future of agentic AI need not rely solely on massive, centralized models, but rather on intelligent orchestration of specialized, efficient components.

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., “CLIMB: Clustering-based Iterative Data Mixture Bootstrapping,” *arXiv preprint arXiv:2504.13161*, 2025.
- [3] J. Su, F. Lin, et al., “CP-Router: An Uncertainty-Aware Router Between LLM and LRM,” *AAAI Conference on Artificial Intelligence*, 2025.
- [4] F. Ju et al., “Reasoning Path Divergence: A New Metric and Curation Strategy,” *arXiv preprint arXiv:2510.26122*, 2025.
- [5] T. Dettmers et al., “QLoRA: Efficient Finetuning of Quantized LLMs,” *NeurIPS*, 2023.