
BUDGET-EFFICIENT LLM SELECTION VIA AUTOMATED SKILL PROFILING

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054
Anonymous Authors¹

ABSTRACT

Introduction. The surge in large language models (LLMs) has introduced a wide range of deployment options—from lightweight open-source models to high-performing proprietary APIs. While larger models typically yield better accuracy and reasoning, they come with significantly higher monetary and energy costs. This raises a core systems challenge: how can we select the most suitable LLM for a given task, under real-world constraints such as latency, energy, or budget?

Benchmark metrics alone are insufficient to guide this decision. Choosing the most accurate model often ignores which capabilities are required, why failures occur, or whether a cheaper model might suffice. We introduce **BELLA** (Budget-Efficient LLM Selection via Automated Skill Profiling), a system that decomposes LLM outputs into interpretable skills and weaknesses, builds structured model-task-skill matrices, and recommends models that maximize utility within resource limits. While we demonstrate BELLA on financial reasoning tasks, its architecture generalizes to any domain exhibiting cost-performance tradeoffs.

Background and Related Work. Conventional LLM evaluation pipelines emphasize aggregate task metrics (e.g., accuracy, F1) or compute costs (e.g., FLOPs, latency), but provide limited insight for deployment decisions. Recent profiling tools offer finer granularity: EvalTree (Zeng et al., 2025) builds hierarchical skill trees but assumes one skill per instance; Skill-Slices (Moayeri et al., 2024) cluster rationales to derive skill-level performance but require detailed reasoning traces. QualEval (Murahari et al., 2024) and CheckList (Ribeiro et al., 2020) provide test-based diagnostics but lack per-instance attribution and do not incorporate deployment-time constraints like cost or latency.

Separately, model routing frameworks like GraphRouter (Feng et al., 2025), AS-LLM (Wu et al., 2024), and ModelSwitch (Chen et al., 2025) improve efficiency via learned policies or model agreement, and in some cases explicitly factor in cost. However, they do not reason over skill-level capabilities or offer interpretable insight into why a model is selected for a given task. BELLA complements these systems by aligning required and demonstrated skills while optimizing for real-world resource limits—enabling cost-aware selection grounded in interpretable, per-instance capability profiling.

Methodology. BELLA is a modular pipeline for skill-based model selection under deployment constraints. It consists of five stages:

1. **Benchmarking.** We evaluate a suite of LLMs on diverse multi-skill reasoning tasks, collecting standardized performance metrics (e.g., accuracy, F1) and operational costs (e.g., API usage, latency, energy). These serve as the empirical foundation for downstream profiling and selection.
2. **Skill Profiling via Critic LLM.** For each model-task-instance, a critic LLM is prompted with the input, reference output, and model output. It returns a structured list of demonstrated and missing skills, grounded in the model’s reasoning trace. This yields a fine-grained, per-instance capability profile.
3. **Skill Clustering and Canonicalization.** We embed and cluster all annotated skill/weakness phrases to derive a compact, interpretable set of $\sim 15\text{--}20$ canonical skills. Each instance is relabeled with binary indicators for skill presence and impact on outcome.
4. **Skill-Aware Model Selection.** For a new task, BELLA infers the required skills (via prompting or user input), filters models that cover those skills, and applies multi-objective optimization to select one that maximizes performance within budget constraints.
5. **Evaluation.** We evaluate BELLA using **leave-one-task-out cross-validation**. For each held-out task, we infer its required skills, apply BELLA’s selection process, and compare performance and cost against oracle and naive baselines. We also test alternative skill profilers (e.g., EvalTree) to assess robustness to representation quality.

BELLA enables plug-and-play skill extraction, interpretable model routing, and cost-aware selection—bridging the gap between benchmark metrics and deployment-ready decisions. Its modular design supports transparent, efficient use of LLMs across diverse tasks, models, and constraints, making it a deployable, domain-agnostic solution for cost-efficient inference.

055 **REFERENCES**056 Chen, A. et al. Do we truly need so many samples? multi-
057 llm repeated sampling efficiently scales test-time com-
058 pute. *arXiv preprint arXiv:2504.00762*, 2025.
059060 Feng, T., Shen, Y., and You, J. Graphrouter: A
061 graph-based router for llm selections. *OpenReview*,
062 2025. URL <https://openreview.net/forum?id=eU39PDsZtT>.
063064 Moayeri, M., Gehrmann, S., and Beutel, A. Unearthing skill-
065 level insights for understanding trade-offs of foundation
066 models. *arXiv preprint arXiv:2410.13826*, 2024.
067068 Murahari, V., Zhang, T., and Wang, Y. Qualeval: Auto-
069 matic capability profiling and skill discovery for llms. In
070 *NAACL*, 2024.
071072 Ribeiro, M. T., Wu, T., Guestrin, C., and Singh, S. Beyond
073 accuracy: Behavioral testing of nlp models with checklist.
074 In *ACL*, 2020.
075076 Wu, X., Zhong, Y., Wu, J., and Tan, K. C. As-llm: When
077 algorithm selection meets large language model. *Open-
078 Review*, 2024. URL <https://openreview.net/forum?id=17aD9VMQUq>.
079080 Zeng, Z., Wang, Y., Hajishirzi, H., and Koh, P. W. Evaltree:
081 Profiling language model weaknesses via hierarchical
082 capability trees. *arXiv preprint arXiv:2503.08893*, 2025.
083084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109