

Energy-Aware Data-Driven Model Selection in LLM-Orchestrated AI Systems

Daria Smirnova, Hamid Nasiri, Marta Adamska, Zhengxin Yu, and Peter Garraghan

School of Computing and Communications, Lancaster University

Email: {d.smirnova1, h.nasiri, m.adamska, z.yu8, p.garraghan}@lancaster.ac.uk

Abstract—As modern artificial intelligence (AI) systems become more advanced and capable, they can leverage a wide range of tools and models to perform complex tasks. Today, the task of orchestrating these models is often performed by Large Language Models (LLMs) that rely on qualitative descriptions of models for decision-making. However, the descriptions provided to these LLM-based orchestrators do not reflect true model capabilities and performance characteristics, leading to suboptimal model selection, reduced accuracy, and increased energy costs. In this paper, we conduct an empirical analysis of LLM-based orchestration limitations and propose GUIDE, a new energy-aware model selection framework that accounts for performance–energy trade-offs by incorporating quantitative model performance characteristics in decision-making. Experimental results demonstrate that GUIDE increases accuracy by 0.90%–11.92% across various evaluated tasks, and achieves up to 54% energy efficiency improvement, while reducing orchestrator model selection latency from 4.51 s to 7.2 ms.

I. INTRODUCTION

Modern AI systems have seen widespread adoption. Conventional AI systems are designed to complete a singular or narrow set of tasks such as object detection [1] or speech recognition [2]. AI systems – particularly those using Large Language Models (LLMs) – can facilitate complex tasks by leveraging an assortment of tools and models in tandem to expand system versatility and capability. Orchestration – responsible for determining model selection, planning, and execution to complete tasks [3] – is key in AI systems. As system scale and the number of models grow, there is an urgent need to design AI system orchestrators that operate at high speed, accuracy, and low cost.

In response, researchers have recently identified that LLMs due to their strong reasoning capabilities can be instructed to operate as AI system orchestrators. These *LLM-based orchestrators* leverage foundational models to reason, make decisions, and act to invoke specific tools and models [4, 5, 6]. The effectiveness of these orchestrators has resulted in their increased adoption within agentic tool-augmented systems [7].

However, existing LLM-based orchestrators encounter many challenges pertaining to their performance and efficiency:

(i) *Orchestrator decision making for model selection is limited to qualitative data* of a model’s capability, specifically model cards and textual descriptions (e.g., recommended tasks or number of user likes on HuggingFace) [8]. While such information helps provide task context to the orchestrator, it

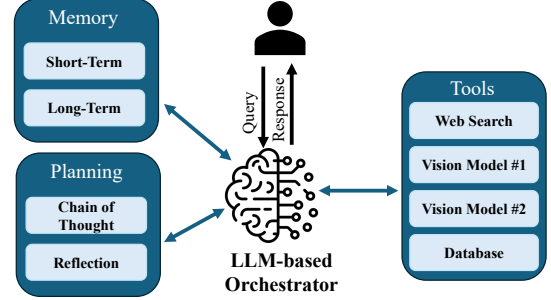


Fig. 1. Overview of an LLM-orchestrated AI system.

does not capture task-level performance (i.e. accuracy, energy cost). Reliance solely on qualitative model information results in suboptimal selection of models that are neither the most accurate nor cost-efficient.

(ii) *LLM-based orchestration requires extensive communication with an LLM to operate* — i.e., textual descriptions of models must be sent to the LLM to provide context, such as API descriptions and usage examples for each tool [9]. This incurs significant token usage and latency (especially at greater AI system scale), both of which increase energy consumption and reduce system throughput. Moreover, the LLM’s context window limits the number of tools the orchestrator can consider.

Hence, we assert that the current design of LLM-based orchestrators for model selection significantly impairs AI system scalability, performance, and efficiency.

To address these limitations, we postulate that incorporating quantitative data (accuracy, energy cost, etc.) into orchestrator decisions would result in more effective model selection, enabling the AI system to balance task accuracy and cost-efficiency (specifically energy use, given the growing concerns of AI sustainability and datacenter growth [10]). There have been extensive discussions on performance–energy trade-offs in adjacent domains such as LLM-serving [11, 12, 13, 14] and IoT [15]. However, to the best of our knowledge, there is no prior work that has analyzed the above limitations or the effectiveness of incorporating quantitative data within LLM-based orchestration for AI systems.

In this paper, we present a comprehensive experimental analysis of LLM-based orchestrator operation and highlight recurrent shortcomings. We identify that LLM orchestrators

systematically select underperforming models that increase costs and energy use with minimal accuracy gain, while introducing overhead, and highlight the need for a quantitative, cost-aware selection policy. Derived from these findings, we present GUIDE – an energy-aware approach to model selection within LLM-orchestrated AI systems. GUIDE is designed to incorporate system-level model performance characteristics into orchestrator decision-making during selection, and clearly defined policies enable the maximization of task quality whilst minimizing cost. Our framework incorporates performance–energy trade-offs in decision-making and couples it with real-time system energy metrics to ensure decisions account for current energy usage.

The key contributions of this paper are:

- An experimental analysis of state-of-the-art LLM-based model selection methods, demonstrating current performance challenges. Our findings show that LLMs invoke more models than needed (in over 99% of cases for some tasks) and make suboptimal model choices.
- A quantitative energy-aware LLM-orchestration framework GUIDE, comprising (i) an energy budget tracker monitoring system energy levels, and (ii) a model selector that prioritizes performance-energy trade-off using Pareto-optimization, informed by the budget tracker. GUIDE attains higher accuracy across tasks, improved performance-energy trade-off, and orders-of-magnitude lower latency.

The rest of the paper is organized as follows: Section II reviews background concepts in LLM-orchestrated AI systems. Section III presents an empirical evaluation of existing LLM-orchestrated frameworks, highlighting their limitations in task identification and model selection. Section IV introduces our proposed energy-aware model selection framework and its main components. Section V presents the experimental setup and results demonstrating the effectiveness of our proposed approach. Section VI discusses related work, and Section VII concludes the paper.

II. BACKGROUND

LLM-orchestrated systems have emerged as a compelling new paradigm for integrating diverse computational capabilities into a single pipeline (Figure 1). In these systems, an LLM functions as the central coordinator across various tasks. Typically, an LLM-orchestrated system includes (i) a task identification module, (ii) a task planner that manages task dependencies, (iii) a tool or model selector that is connected to a registry of available ML models and tools [16, 17], and (iv) a response generator that consolidates outputs from all models and delivers a coherent response to the user.

LLMs are typically at the center of all of these aforementioned subtasks. Given a few in-context examples, they can successfully coordinate them. For example, for model selection, an LLM is provided with a predefined list of candidate models, which often includes extensive textual descriptions and APIs, and in-context examples of how to perform model selection and invoke these models [9].

Despite LLMs being versatile tools that can be instructed to perform many system subtasks, each LLM-based decision incurs substantial costs, including energy and token usage, as well as latency [18]. Hence, using LLMs for simple subtasks, such as model selection in LLM-orchestrated systems, is not always practical. Moreover, LLMs can hallucinate, fail to follow in-context examples, or choose incompatible models. Therefore, using LLMs for orchestration can be impractical and unreliable.

III. LIMITATIONS OF CURRENT LLM-ORCHESTRATED SYSTEMS

A. Evaluation setup

Our objective is to empirically analyze how current LLM-based orchestrators perform model selection, to what extent they consider performance and cost, and what overhead can arise from this approach.

To ground this analysis, we evaluate a representative LLM-orchestrated framework that follows a typical design pattern for an LLM-orchestrated AI system: a registry of models and an orchestrator (LLM) that can build workflows and select and invoke models. We use JARVIS (Hugging GPT) as this case study as it is a popular framework that connects an LLM with a multitude of models hosted on HuggingFace via APIs [4]. The goal of this evaluation is to assess the efficiency of LLM-based decision-making used in model selection procedures across various tasks, using JARVIS as a representative case study.

We evaluated tasks commonly used in AI systems:

- **ICapt** - Image Captioning. The dataset (jpawan33/fkr30k-image-captioning-dataset) contains approximately 30,000 512x512 images with accompanying textual descriptions.
- **VQA** - Visual Question Answering. The dataset (Immslab/OK-VQA) consists of images of varying sizes, textual questions about them, and a list of potential answers provided for each image.
- **OD** - Object Detection. The dataset (COCO2017) contains over 200,000 annotated images depicting everyday objects in different real-world scenarios.
- **IGen** - Image Generation. The dataset (Falah/image_generation_prompts_SDXL) consists of 1M textual descriptions of scenes and objects, to be used as prompts for image generation.

For this analysis, we used subsets of these datasets. In particular, for VQA and OD, we used 5000 prompts, and for ICapt and IGen, 1000. This is due to low observed variability in model selection across the ICapt and IGen tasks (Table I). Each prompt was sent via an HTTP request to the JARVIS server, set up locally on our Ubuntu machine, equipped with an NVIDIA RTX 6000 Ada GPU. On this machine, we had 26 models of these four task types locally hosted.

The data collected in this experiment include details on the identified task types for each prompt, the selected models, and the corresponding time, CPU, and GPU energy consumption. For CPU and GPU energy data collection, we used Turbostat and Zeus Python libraries, respectively.

B. Analysis findings

Upon the analysis of the results, two main issues with LLM orchestrated decision making were discovered:

Task misclassification. We observed that, for the majority of the prompts in the VQA and OD datasets, an LLM either selected the wrong task type or decided that more than one task was required to fulfill the user’s request (Table I). Using more than one task led to lower overall accuracy and significantly increased energy cost, while incorrectly identifying a task type simply does not answer a user query. Such task misclassification is problematic because different task combinations exhibit distinct energy profiles and accuracies.

For example, in 99.2% of VQA prompts, the LLM built a workflow with at least 2 tasks. The most popular task combination for VQA prompts is *ICapt, VQA* (70% of all VQA prompts). JARVIS exhibited similar behaviour on the OD task, where 99.2% of prompts were classified as *OD, VQA*, instead of *OD*.

As demonstrated in Table I, different task combinations exhibit distinct energy profiles and accuracies. Whenever an LLM identifies a VQA prompt correctly as a *VQA* task, the average GPU energy consumption is only 65 J (corresponding to 0.3% of all VQA prompts), whereas the most popular combination *ICapt, VQA* (70% of prompts) has a mean energy of 224 J. Similarly, correctly identifying an OD prompt can use approximately two times less energy than the predominant *OD, VQA* combination (18.2 J compared to 37.5 J). Accuracy follows the same pattern: it is the highest (92.3%) for VQA prompts that are identified correctly as a *VQA* task, and in the OD experiment, the predominant *OD, VQA* combination achieves 20.1%, compared to 68.2% for the correct *OD* combination.

This means that using multiple models to answer a user’s prompt can paradoxically lead to lower-quality responses, likely because the LLM loses some information in the process of consolidating model outputs into one coherent response for the user. This exposes a resource-efficiency issue in LLM-orchestrated systems that must be addressed to improve their efficiency.

Suboptimal Model Selection. Even when an LLM correctly identifies a task type, it can still make suboptimal decisions when selecting a model from the candidate list. We observed two main issues with LLM-based model selection: (1) a strong bias toward more popular models when model metadata contains popularity metrics, such as HuggingFace model card likes, and (2) overreliance on the LLM’s internalized knowledge about models when such metadata is missing from model descriptions. We discuss both issues in the following subsections.

1) *Popularity-Based Selection Bias:* Our analysis shows that, in most cases, the LLM favors the model with the greatest recognition on HuggingFace (measured by the number of likes). This behavior is problematic because popularity does not necessarily correlate with model efficiency or accuracy. As a result, the LLM often selects models that are neither the most accurate nor the most energy-efficient among the available options.

TABLE I
PROMPT-LEVEL TASK COMBINATIONS WITH ASSOCIATED ENERGY USAGE AND ACCURACY.

Task Combination	Proportion of Prompts	Avg. Energy (J)	Task Accuracy
VQA Dataset			
ICapt, VQA	70%	224	62.5%
ICapt, VQA, OD	17.5%	581	65.0%
ICapt, DocVQA	8.3%	576	22.0%
VQA, IClass	2.3%	129	76.8%
ICapt	0.5%	167	55.6%
VQA	0.3%	65	92.3%
ICapt, OD	0.2%	435	50.0%
VQA, OD	0.1%	334	85.7%
OD Dataset			
OD, VQA	99.2%	37.5	20.1%
OD	0.7%	18.2	68.2%
VQA	0.1%	19.4	1.7%
ICapt Dataset			
ICapt	100%	64.1	—
IGen Dataset			
IGen	100%	512.8	—

TABLE II
PERCENTAGE OF MODEL SELECTIONS PER TASK CATEGORY USING JARVIS DEFAULT MODEL SELECTION METHOD.

Task Type	Model	Chosen	Number of likes	Mentions likes
VQA Dataset				
DocVQA	LayoutLM-DQA	100%	174	98.07%
IClass	ViT-B16	100%	169	83.19%
	ViT-GPT2	78.98%	219	53.70%
ICapt	BLIP-Capt-L	13.89%	52	38.24%
	BLIP2-6.7B	4.49%	24	0%
	BLIP2-2.7B	2.09%	25	0%
	TrOCR-B-P	0.55%	56	0%
OD	DETR-R-50	100%	129	96.05%
VQA	ViLT-B32-VQA	100%	86	99.00%
OD Dataset				
OD	DETR-R-50	100%	129	99.98%
VQA	ViLT-B32-VQA	100%	86	99.86%
ICapt Dataset				
ICapt	ViT-GPT2	100%	219	87.21%
IGen Dataset				
IGen	SD-1.5	100%	6367	99.83%

For example, Table II shows that 100% of prompts from the ICapt dataset sent to JARVIS resulted in the LLM choosing the *ViT-GPT2* model out of five candidate ICapt models available on our server. Moreover, in 87.21% of responses, the LLM explicitly mentions the number of likes in its reasoning for this model choice. After a quick profiling of all candidate models on

items from the same dataset, we found that *ViT-GPT2* (accuracy: 0.284, energy: 12.7) is neither the most energy-efficient model nor the most accurate among the five candidates (Table III). In fact, the most accurate model for the ICapt task is *BLIP2-6.7B* (accuracy: 0.320, energy: 110.8), while the best trade-off between accuracy and energy efficiency is achieved by *BLIP-Capt-B* (accuracy: 0.315, energy: 12.5). Notably, *BLIP-Capt-B* was never selected due to its low number of likes (i.e., 44).

TABLE III
PROFILING OF ICAPT MODELS.

Model	ROUGE-L	Avg. Energy (J)	Time (s)	Number of likes
ViT-GPT2	0.284	12.7	0.07	219
BLIP-Capt-L	0.286	17.2	0.10	52
BLIP-Capt-B	0.315	12.5	0.09	44
BLIP2-2.7B	0.301	43.2	0.19	25
BLIP2-6.7B (8-bit)	0.320	110.8	0.41	24

A similar issue is observed in the VQA dataset experiment. As shown in Table II, the LLM selected *ViLT-B32-VQA* (accuracy: 0.414, energy: 6.5) for 100% of all VQA tasks, and in at least 99% of responses, the number of likes was mentioned in the reasoning. Although *ViLT-B32-VQA* is the least costly model, it is not the most accurate (Table IV). If the LLM instead selected *BLIP-VQA-B* (accuracy: 0.531, energy: 6.9), accuracy could have increased by 28.3%, adding only 6.2% in energy cost.

These results clearly show that the number of likes and downloads is not a reliable proxy for model quality or efficiency. High popularity may reflect good community recognition, but does not necessarily indicate outstanding performance. Therefore, qualitative metrics such as likes should not guide model selection; they should only provide context for understanding model use.

TABLE IV
PROFILING OF VQA MODELS.

Model	ROUGE-L	Avg. Energy (J)	Time (s)	Number of likes
ViLT-B32-VQA	0.414	6.5	0.03	86
BLIP-VQA-L	0.531	7.3	0.06	52
BLIP-VQA-B	0.531	6.9	0.06	44
GIT-L-VQAv2	0.048	6.6	0.05	3
GIT-B-TextVQA	0.048	6.6	0.04	1
GIT-B-VQAv2	0.069	6.8	0.05	1

2) *Opaque Model Selection in the Absence of Popularity Metrics*: Removing the popularity data from the model metadata did not solve these issues—the LLM still made suboptimal choices. Without the popularity data, the LLM exhibited a strong preference for certain models, driven by its

internal knowledge, which can be too general to reflect models’ actual performance.

To investigate this further, we conducted an additional set of experiments using the same setup, but excluded the number of likes and downloads from the model descriptions. Without these metrics, the LLM was provided only with the model IDs and short textual descriptions (up to 100 characters) extracted from HuggingFace model cards. In most cases, these model descriptions only included the model’s full name. For clarity, from this point onward, we call the default LLM-based model selection “JARVIS”, and the variant without popularity information as “Name-Only”. The results of Name-Only experiments are presented in Table V.

On the OD task, Name-Only showed a strong preference for *DETR-R-101*, selecting it in 99.82% of cases, stating in the reasoning that this model has higher performance capabilities than *DETR-R-50*. However, this information was not included in the provided model description; therefore, this leads us to the assumption that the LLM relies mostly on LLM’s prior internal knowledge about the performance of these models. Indeed, *DETR-R-101* has shown improved performance over *DETR-R-50* (+1.6 AP on COCO), with some increase in computational cost (GFLOPs 187 → 253, +35%), as reported by [19]. Another good model choice would have been *YOLOv8-S*. This lightweight object detection model achieves an AP of 44.9 with only 11.2M parameters and 28.6 GFLOPs (approximately 89% fewer GFLOPs than *DETR-R-101*), enabling faster, more efficient computation. Despite this, *YOLOv8-S* is not picked by LLM-based methods. Hence, relying on LLMs’ judgment calls is highly unreliable, as they do not reflect the model’s true capabilities and can exhibit biases due to the data they were trained on.

In our VQA experiment, the LLM consistently selected the vision transformer model *ViLT-B32-VQA* (accuracy: 0.414, energy: 6.5), and in its reasoning it states that this model is “specifically finetuned for visual question answering” and has “the vision and language transformer vilt architecture”. However, choosing *BLIP-VQA-B* (accuracy: 0.531, energy: 6.9) would result in a better accuracy-energy trade-off. For ICapt tasks, *BLIP-Capt-L* (accuracy: 0.531, energy: 7.3) is consistently chosen with reasoning stating that it is “known for its strong performance in unified vision-language understanding”. However, as previously mentioned, our quick profiling shows that *BLIP-Capt-B* (accuracy: 0.531, energy: 6.9) is just as accurate but consumes less energy, yet it was never selected.

These suggest that the internal knowledge of LLMs about ML models is too generalized and biased, and does not accurately reflect the models’ true capabilities when compared with other candidates in the list. Hence, relying on black-box LLM decision-making for model selection consistently leads to suboptimal outcomes. This means that neither popularity metrics nor LLM’s internal knowledge about the models should be used as primary decision-making factors.

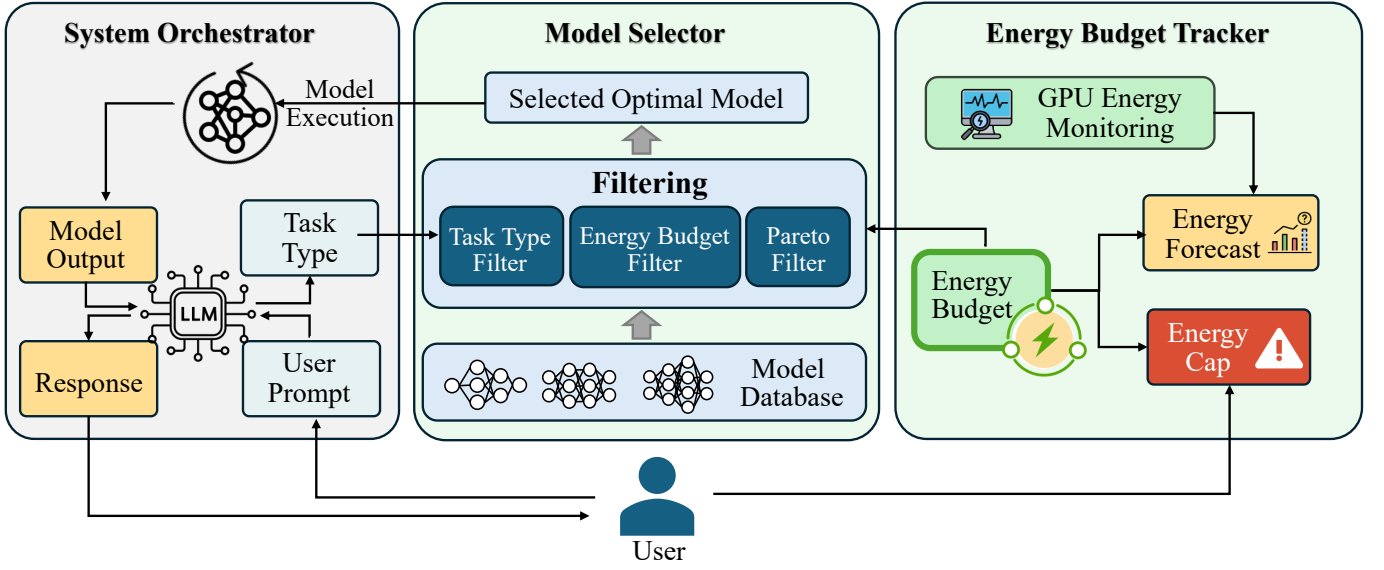


Fig. 2. Overview of the energy-aware LLM-orchestrator model selection framework. The Energy Budget Tracker (right) estimates the current per-slot energy budget based on GPU energy deltas using EMA and the user-defined energy cap. The Model Selector (middle) filters models by task and budget, applies a Pareto filter on accuracy–energy, and returns the most accurate model from the Pareto efficient subset. The System Orchestrator (left) executes the selected model and returns the response to the user.

IV. ENERGY-AWARE DATA-DRIVEN MODEL SELECTION

We propose that model selection should rely on explicit quantitative performance and efficiency metrics, such as accuracy and energy consumption. To test this hypothesis, we developed a data-driven model selection methodology and incorporated it into JARVIS to evaluate how the proposed framework affects the efficiency and performance of such LLM-orchestrated AI systems.

We design a framework for energy-adaptive model orchestration, in which model selection is explicitly constrained by a user-defined energy budget per time slot (Figure 2). The main idea is that an LLM-orchestrated system should not only determine which model is most suitable for a given task but also verify that the chosen model does not exceed the user-defined energy target. To achieve this, we implement an energy tracker that continuously monitors GPU energy consumption and estimates the amount of energy we can afford to use in this time slot without overshooting the target, using an exponential moving average (EMA). Based on this EMA estimate, the orchestrator selects the best model among the candidates that will not cause the energy to overshoot the user-defined target. We use Pareto optimization to balance accuracy and average energy usage when selecting models. In this way, the system balances quality and efficiency based on the real-time system state, rather than assuming that computational resources are unlimited.

The proposed framework consists of two main components: the energy budget tracker and the model selector. Based on the user-defined soft per-window energy target, our framework aims to select the most accurate Pareto-efficient model that fits within this energy constraint.

Algorithm 1 Energy Budget Tracker

Input: Slot duration S , energy level C , polling interval Δt , EMA weight α
Output: E_{usable}
 $E_{\text{used}} = 0$, $P_{\text{EMA}} = 0$
for $k = 1$ **to** $\lfloor S/\Delta t \rfloor$ **do**
 $E_k \leftarrow \text{energy}(\Delta t)$
 $P_k \leftarrow E_k/\Delta t$ ▷ instantaneous power (W)
 $P_{\text{EMA}} \leftarrow \alpha P_k + (1 - \alpha)P_{\text{EMA}}$
 $E_{\text{used}} \leftarrow E_{\text{used}} + E_k$
 $t_{\text{rem}} \leftarrow S - k\Delta t$
 $E_{\text{rem}}^{\text{pred}} \leftarrow P_{\text{EMA}} \times t_{\text{rem}}$
 $E_{\text{tot}}^{\text{pred}} \leftarrow E_{\text{used}} + E_{\text{rem}}^{\text{pred}}$
 $E_{\text{usable}} \leftarrow \max(0, C - E_{\text{tot}}^{\text{pred}})$
end for

A. Energy Budget Tracker

The tracker runs in the background, performing GPU metering online and forecasting its energy consumption (Algorithm 1). The GPU time is split into fixed slots, and every 100 ms, it pulls the GPU’s energy reading and converts it into instantaneous power. Slot length is fixed at 2 s to match typical per-request latencies and to ensure fast adaptation. The GPU meter readings are then used to compute an EMA of power, smoothing out noise and making sure the tracker can still react to sudden spikes in energy consumption.

Given time slot duration S , user-defined energy target C , polling interval Δt , and EMA weight $\alpha \in (0, 1)$, the tracker estimates how much energy will be used in the current time slot, and exposes a usable per-slot energy budget that we can

TABLE V
PERCENTAGE OF MODEL SELECTIONS PER TASK CATEGORY USING
NAME-ONLY MODEL SELECTION METHOD.

Task Type	Model	Chosen
VQA Dataset		
DocVQA	LayoutLM-DQA	100%
IClass	BEiT-B16	78.85%
	DETR-R-50	21.15%
ICapt	BLIP-Capt-L	96.66%
	BLIP2-6.7B	2.75%
	TrOCR-B-P	0.59%
OD	DETR-R-101	95.37%
	YOLOv8-S	3.02%
	DETR-R-50	1.61%
VQA	ViLT-B32-VQA	100%
OD Dataset		
OD	DETR-R-101	99.82%
	DETR-R-50	0.18%
VQA	BLIP-Capt-L	100%
ICapt Dataset		
ICapt	BLIP-Capt-L	96.79%
	ViT-GPT2	3.21%
IGen Dataset		
IGen	SD-1.5	99.62%
	OpenJourney	0.36%
	SD-2.1	0.02%

use to do additional work (e.g. run an inference). All quantities are computed at discrete times $t_k = k\Delta t$.

Within each time slot, we collect the cumulative energy E_{used} used up until this moment in time, and calculate the remaining time in the slot: t_{rem} . Given t_{rem} , we then estimate the remaining energy usage for this slot using EMA: $E_{\text{rem}}^{\text{pred}} = P_{\text{EMA}} \times t_{\text{rem}}$. We then calculate usable energy budget available for new work (e.g., a new ML inference), based on the user-defined energy level: $E_{\text{usable}} = \max(0, C - (E_{\text{used}} + E_{\text{rem}}^{\text{pred}}))$.

B. Model Selector

For each new task, this component first filters candidate models based on task type, energy usage, and Pareto efficiency, aiming to avoid scheduling an inference that would cause the slot to exceed the defined energy level (Algorithm 2). The selector then chooses the model with the best accuracy.

The Model Selector filters given set of models \mathcal{M} by task type τ and then, using the budget E_{usable} , retrieved from the Tracker, it finds only the models that fit this budget ($E_{\text{avg}}(m) \leq E_{\text{usable}}$). After that, the Selector performs Pareto-Efficient filtering to obtain a subset of models on the Pareto Frontier of $(\text{Acc}, E_{\text{avg}})$. If the resulting subset is empty, the selector waits for Δt_{retry} s, obtains an updated budget E_{usable} from the tracker, and retries to obtain a Pareto-efficient subset of models. From a Pareto-efficient subset, the model selector picks a model that has the highest accuracy on the given task. We set $\Delta t_{\text{retry}} = 0.5$ s.

Algorithm 2 Model Selector

Input: Model set $\mathcal{M} = \{(m, \tau(m), E_{\text{avg}}(m), \text{Acc}(m))\}$, usable energy budget E_{usable} , current task type τ , Δt_{retry}
Output: Selected model m^*

repeat
 $\mathcal{M}_\tau \leftarrow \{m \in \mathcal{M} \mid \tau(m) = \tau\}$
 $E_{\text{usable}} \leftarrow \text{Tracker.pull()}$ \triangleright read latest usable energy from Energy Budget Tracker
 $\mathcal{M}_{\text{budget}} \leftarrow \{m \in \mathcal{M}_\tau \mid E_{\text{avg}}(m) \leq E_{\text{usable}}\}$
if $|\mathcal{M}_{\text{budget}}| = 0$ **then**
Wait Δt_{retry} seconds
end if
until $|\mathcal{M}_{\text{budget}}| > 0$
 $\mathcal{P} \leftarrow \{m \in \mathcal{M}_{\text{budget}} \mid \nexists m' \in \mathcal{M}_{\text{budget}} : E_{\text{avg}}(m') < E_{\text{avg}}(m) \wedge \text{Acc}(m') > \text{Acc}(m)\}$ \triangleright Pareto-efficient subset
 $m^* \leftarrow \arg \max_{m \in \mathcal{P}} \text{Acc}(m)$

This process allows our framework to ensure a balance of accuracy and energy, and adherence to user-defined energy targets.

V. PERFORMANCE EVALUATION

A. Experiment Setup

We evaluate our energy-adaptive model selection framework in comparison to a prominent popularity-driven method, JARVIS, and a knowledge-based method, Name-Only. Using JARVIS as a representative case study is sufficient to demonstrate the inefficiency of popularity- and knowledge-driven model selection and to evaluate the benefits of introducing quantitative, energy-aware decision-making.

Hardware. All experiments ran on a single NVIDIA RTX 6000 Ada GPU (50 GB VRAM), on an Ubuntu 20.04.6 LTS. The system uses NVIDIA driver 535.247.01 (CUDA runtime 12.2). The CUDA toolchain is nvcc 11.8 (11.8.89) with GCC 9.4.0.

Software. We evaluate within the JARVIS framework, used as a representative LLM-orchestrated AI system that leverages multiple ML models. The model selector and the energy budget tracker are implemented in Python. The energy budget tracker is a separate, parallel-running process, while the model selector that communicates with the tracker makes per-request model choices. Inference is executed by the framework’s built-in execution engine. To reduce variance, we disable background services unrelated to the experiment and maintain a fixed software environment across all runs.

Model Selection Methods. We compare three model selection policies within the identical JARVIS pipeline:

- JARVIS: The framework’s default model selection method that prioritizes models with the highest number of likes.
- Name-Only: The knowledge-based model selection method that makes decisions based on only the model’s textual descriptions, such as its full name.
- GUIDE: Our proposed model selection framework with four user-defined energy targets of 100 J, 150 J, 400

J, and 600 J per slot (GUIDE-100/150/400/600). These levels were chosen based on our model profiling. Lower energy levels are suited to small-scale tasks (ICapt, VQA, and OD), whereas higher energy levels target the more demanding IGen workload.

LLM-based policies (JARVIS and Name-Only) use OpenAI’s GPT-4o-mini as the model orchestrator, chosen for its high performance at an affordable cost.

Models. We use the same 26 models as in the analysis in Section III. For each model, we measure energy consumption for (i) the forward pass only and (ii) the full model lifecycle, which includes loading into and unloading from device memory. The default built-in execution engine of JARVIS loads and unloads a model for every request; therefore, we changed this logic to keep models in memory for the whole experiment to better reflect real-world system dynamics, where models are typically kept loaded in device memory to serve multiple requests.

Datasets. For a fair comparison, we evaluate all three policies on fixed-size subsets from the datasets mentioned in Section III. Specifically, we use 100 prompts per task from these sources. Each prompt can result in the execution of models of different task types, not limited to the dataset’s task type.

Evaluation metrics. To assess the quality of models’ performance on VQA and ICapt tasks, we report the ROUGE-L score, as it is an effective measure of similarity between two texts. To evaluate the performance of OD models, we use the standard mAP@0.5 metric, which measures how accurately a model identifies objects in an image at an Intersection over Union (IoU) threshold of 0.5. To report the response quality of IGen models, we use CLIP (openai/clip-vit-base-patch32), which assesses semantic alignment between prompts (textual descriptions) and generated images. To evaluate the energy usage of models and processes, we use the standard unit of Joules (J). We also report Accuracy per Joule (Acc/J) as an energy-efficiency metric that captures the trade-off between energy consumption and accuracy per request. This metric, commonly used in prior work [20, 21], measures the amount of accuracy achieved per unit of energy consumed. We adopt this efficiency metric as it allows a fair comparison of model selections across different policies.

B. Experiment Results and Analysis

Relying on qualitative descriptions of model popularity degrades model selection performance in terms of both accuracy and energy efficiency. Table VI shows that model selection based on quantitative data (GUIDE) improves accuracy across all four tasks (with absolute gains ranging from 0.90% to 11.92%). The accuracy of GUIDE is higher than that of both JARVIS and Name-Only on ICapt and VQA tasks. Moreover, for the ICapt, VQA, and OD tasks, it achieves better overall accuracy–energy trade-offs, as shown in Figure 3. GUIDE-100 demonstrates the largest efficiency gains, improving Accuracy-per-Joule by 54%, 14.4%, and 18.4% over JARVIS on ICapt, OD, and VQA, respectively. It also outperforms Name-Only on this metric by 67% on ICapt, 6.9% on OD, and 18.4% on

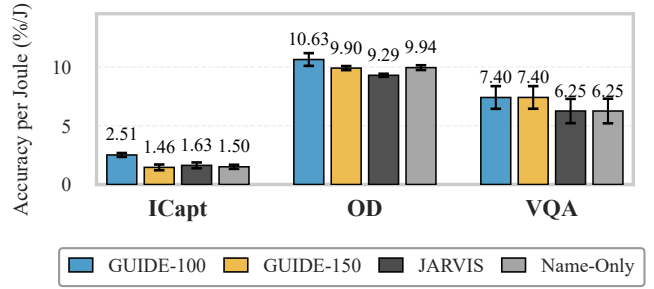


Fig. 3. Accuracy per Joule, calculated on weighted energy and accuracy results, on ICapt, VQA and OD task types for each method: GUIDE model selection (with two GUIDE level settings) and two comparable model selection methods - JARVIS and Name-Only. Error bars represent 95% CI.

TABLE VI
SYSTEMS COMPARISON BY TASK, AGGREGATED OVER ALL REQUESTS FROM FOUR DATASETS (SORTED BY ACC/J WITHIN EACH GROUP). ABSOLUTE VALUES SHOWN. ACC = ACCURACY; ACC/J = ACCURACY PER JOULE.

System	Acc (%)	Energy (J)	Acc/J (%/J)
ICapt			
GUIDE-100	31.50	12.54	2.51
Jarvis	28.60	17.58	1.63
Name-Only	28.67	19.07	1.50
GUIDE-150	31.55	21.68	1.46
OD			
GUIDE-100	70.03	6.59	10.63
Name-Only	71.11	7.15	9.94
GUIDE-150	71.26	7.20	9.90
Jarvis	65.00	7.00	9.29
VQA			
GUIDE-150	52.54	7.10	7.40
GUIDE-100	52.50	7.09	7.40
Jarvis	40.62	6.50	6.25
Name-Only	40.62	6.50	6.25
IGen			
GUIDE-400	34.20	193.33	0.18
Jarvis	34.20	193.33	0.18
Name-Only	34.19	193.40	0.18
GUIDE-600	35.10	431.64	0.08

VQA. GUIDE-100 both increases accuracy and reduces energy usage on ICapt and OD tasks. For the IGen tasks, we observe an accuracy improvement as well, but this is likely due to the use of a more relaxed energy level (600 J), which allows the use of larger models.

Data-driven model selection methodology allows selections to concentrate on a Pareto-efficient subset of models, leading to a better accuracy–energy balance. Data-driven model selection methodology, GUIDE, places 100% of model selections on the Pareto frontier of (accuracy, energy). In contrast, the two comparable methods, JARVIS and Name-

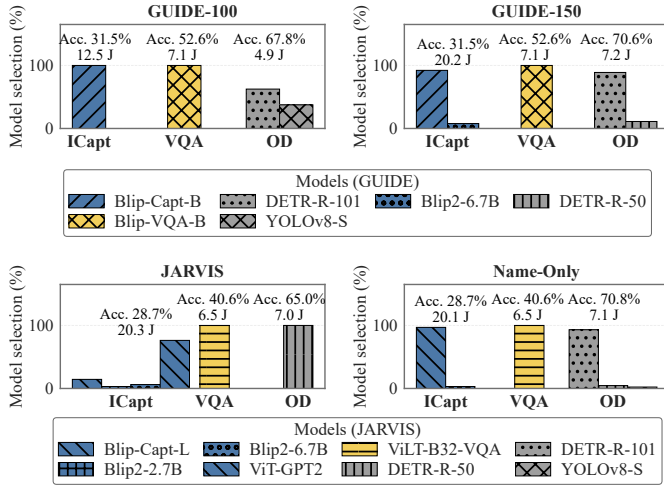


Fig. 4. Model selection performance on the VQA dataset.

Only, selected Pareto-efficient models in 72.3% and 72.7% of cases, respectively, across all datasets and task types. Hence, in at least 28% of all model selections, these policies chose a model that was not on the Pareto frontier, meaning that there existed at least one other model that was better on both metrics (i.e., higher accuracy and lower energy consumption). Our model selection approach avoids making such suboptimal choices, ensuring a more efficient use of available resources.

Our proposed selector operates at a small fraction of the overhead imposed by LLM-based model selectors. The energy-aware model selector GUIDE introduces a mean latency of 7.2 ms per request and consumes on average 21.9 mJ on the CPU per decision (data collected from 10,000 requests). The NVML-based energy tracker that runs in the background contributes on average 1 mJ of CPU energy per selection. Additionally, due to constant polling of GPU readings (NVML queries) at 100 ms intervals, the tracker adds 0.316 W overhead to the GPU power consumption (about 18.96 J per minute). In contrast, for the LLM-based selectors (based on 10,000 selections using JARVIS), the selection step uses on average 2,081 tokens per request (input: 1,977, output: 104) and incurs a mean selection time of 4.51 s. The input token number is large because for each decision, JARVIS sends a few in-context examples and a list of candidate models with their textual descriptions to the LLM. Due to this, each LLM-based decision halts the system for over 4 s, which can significantly impair quality of service. Overall, GUIDE achieves 6.4×10^2 lower latency (4.51 s vs. 7.2 ms) and negligible energy overhead, in comparison to LLM-based decision-making.

The GUIDE model selection methodology adheres closely to user-defined energy targets. As summarized in Table VII, compliance is strongest at higher energy levels (e.g., 150 J or above), where for most time slots the realized energy stays below our target. When the target is set to lower levels, such as 100 J, we observe that the rate of overshoots increases. This issue comes up because the user-set energy target is

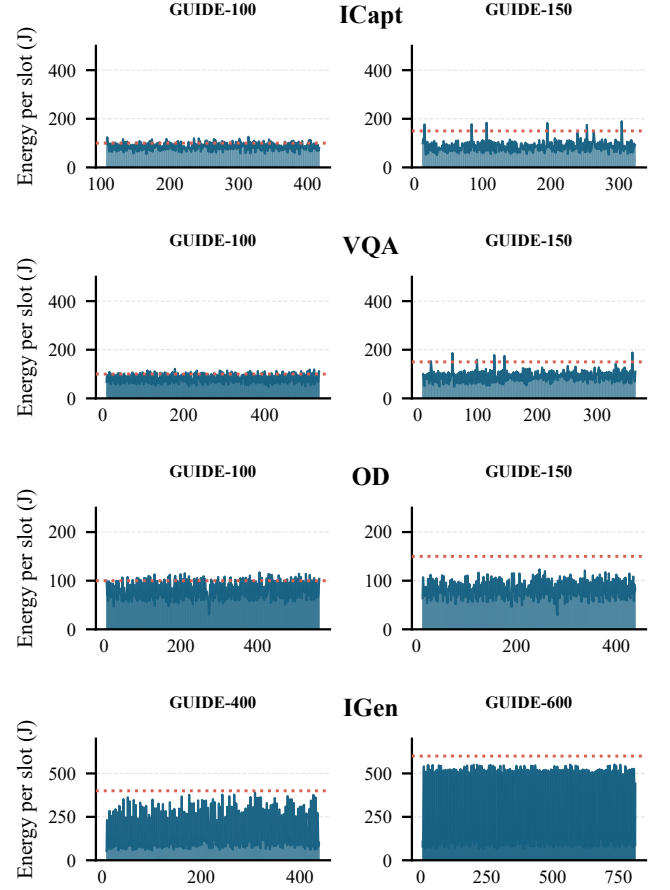


Fig. 5. Per-slot energy usage. The dotted red line indicates the user-defined energy target; actual realized energy usage in each time slot is shown as blue bars. Includes GPU base draw of 40-50 J/s. The x-axis denotes time slots. Each prompt in a dataset of a certain type can result in the execution of many models, not limited to the dataset’s primary type, due to JARVIS multi-step reasoning and task misclassifications, as discussed in Section III.

nearly equal to the system’s baseline fluctuation. As a result, it becomes harder for the energy budget tracker to make stable predictions about current energy availability. Overall, the model selector adapts its model choices to match user-defined energy targets while still maintaining performance. In most cases, the selector successfully follows these targets, demonstrating sufficient compliance.

In the absence of popularity signals, an LLM-based model selector outperforms a popularity-driven baseline on OD tasks. JARVIS, with popularity cues encoded in the model metadata, performs worse on the OD task than Name-Only. Figure 3 shows that, while on the VQA task the two baselines perform comparably, the data for the OD task indicate a clear advantage for the Name-Only method. Compared with JARVIS, it achieved 6.11% higher accuracy on the OD task, while increasing energy cost by only 0.15 J (Table VI). On the OD task, Name-Only has achieved better energy efficiency in Accuracy-per-Joule (+0.65%/J). Across the other three tasks, the

TABLE VII
PERCENTAGE OF TIME SLOTS THAT STAYED UNDER THE ENERGY TARGET.

Dataset	Level (J)	% Violations rate
OD	100	15.3%
OD	150	0%
VQA	100	16.7%
VQA	150	1.7%
ICapt	100	24.2%
ICapt	150	1.9%
IGen	400	0%
IGen	600	0%

differences in their performance are insignificant. This confirms that popularity metadata in model descriptions degrades system performance.

GUIDE replaces LLM-driven selection with a local, data-driven policy that boosts accuracy, cuts costs, and avoids both token usage and response delays. Our model selection policy integrates energy awareness into model orchestration, delivering higher accuracy and more predictable energy consumption than popularity- or knowledge-driven LLM-based selectors. GUIDE improves Accuracy-per-Joule by +54% (ICapt), +18.4% (VQA), and +14.4% (OD), while concentrating its choices on only Pareto-efficient models. GUIDE also saves both tokens and time (7.2 ms vs 4.51 s for JARVIS), avoiding extensive LLM communication, resulting in an average use of 2,081 tokens. GUIDE has higher accuracy at a similar cost in most cases, avoiding inefficient overreliance on LLMs. It allows for better balance of performance and efficiency and introduces negligible overhead.

VI. RELATED WORK

Inefficient use of LLMs and the costs they incur have been extensively studied in other domains. For example, LLM-serving systems can switch between different LLM variants based on query complexity, resulting in lower costs with negligible accuracy loss. These systems invoke larger LLMs only when necessary, and prioritize smaller ones otherwise [11, 12, 14]. An offline energy-optimal framework also showed that profiling LLMs to derive energy and runtime characteristics enables selecting cheaper models to lower system energy usage [22]. However, these energy- and cost-aware approaches focus on LLM serving and do not consider LLM-orchestrated systems that must choose among many ML models.

Current LLM-orchestrated AI systems still make decisions largely from text-based priors or popularity-driven signals and ignore quantitative performance-energy trade-offs and real system energy usage. Similar energy-aware adaptation has been explored in ML for resource-constrained settings, such as IoT: [15] use available energy to choose between variants of tinyML models, and [23] choose a smaller model or apply early exiting to ensure the system stays within a predefined energy budget. Our framework is inspired by such adaptation but targets LLM-orchestrated AI systems with many heterogeneous models.

In large-scale computing, data centers commonly schedule workloads under energy constraints and use mechanisms such as

dynamic GPU frequency scaling and carbon-aware scheduling [24]. Energy-aware scheduling is also used in model-serving systems such as InFaas that choose among multiple model variants to meet accuracy and cost requirements [25, 26]. A cost-aware tool selection framework, CATP-LLM, reduces invocation cost, maintaining accuracy [27]. These efforts show the value of explicit performance–cost trade-offs in decision-making, yet none use real-time energy data with quantitative model performance when selecting models in LLM-orchestrated AI systems. GUIDE closes this gap by making energy-aware, data-driven selections directly in the orchestration layer.

VII. CONCLUSION

In this work, we have analyzed LLM-based model selection methods and their limitations. Our results show that these methods often invoke more models than necessary, misclassify tasks, and make suboptimal model choices. These issues significantly degrade system accuracy and energy efficiency. To address these shortcomings, we propose an energy-aware data-driven model selection framework, GUIDE, that considers model energy and accuracy metrics and uses Pareto optimization to enable explicit performance-cost trade-offs. GUIDE achieves significant accuracy improvements on four vision tasks (with absolute gains ranging from 0.90% to 11.92%) and up to 54% in Accuracy-per-Joule on three tasks, while operating at 7.2 ms latency per request (vs 4.51 s for LLM-based methods).

VIII. ACKNOWLEDGMENTS

This work was supported by the Engineering and Physical Sciences Research Council (Fellowship number EP/V007092/1).

REFERENCES

- [1] Ajantha Vijayakumar and Subramaniaswamy Vairavasundaram. YOLO-based object detection models: A review and its applications. *Multimedia Tools and Applications*, 83(35):83535–83574, 10 2024.
- [2] Sadeen Alharbi, Muna Alrazgan, Alanoud Alrashed, Turkiayh Alnomasi, Raghad Almojel, Rimah Alharbi, Saja Alharbi, Sahar Alturki, Fatimah Alshehri, and Maha Almojel. Automatic speech recognition: Systematic literature review. *IEEE Access*, 9:131858–131876, 2021.
- [3] Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng Li, Tianhe Ren, Xueyan Zou, Jianwei Yang, Hang Su, Jun Zhu, Lei Zhang, Jianfeng Gao, and Chunyuan Li. LLaVA-Plus: Learning to use tools for creating multimodal agents, 2023.
- [4] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. HuggingGPT: Solving AI tasks with ChatGPT and its friends in Hugging Face, 2023.
- [5] Dídac Surís, Sachit Menon, and Carl Vondrick. ViperGPT: Visual inference via python execution for reasoning, 2023.
- [6] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual ChatGPT:

- Talking, drawing and editing with visual foundation models, 2023.
- [7] Joykirat Singh, Raghav Magazine, Yash Pandya, and Akshay Nambi. Agentic reasoning and tool integration for LLMs via reinforcement learning, 2025.
 - [8] Xiangyan Liu, Rongxue Li, Wei Ji, and Tao Lin. Towards robust multi-modal reasoning via model selection, 2024.
 - [9] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools, 2023.
 - [10] Alex De Vries. The growing energy footprint of artificial intelligence. *Joule*, 7(10):2191–2194, 2023.
 - [11] Surya Narayanan Hari and Matt Thomson. Tryage: Real-time, intelligent routing of user prompts to large language models, 2023.
 - [12] Dujian Ding, Ankur Mallick, Chi Wang, Robert Sim, Subhabrata Mukherjee, Victor Ruhle, Laks V. S. Lakshmanan, and Ahmed Hassan Awadallah. Hybrid LLM: Cost-efficient and quality-aware query routing, 2024.
 - [13] Dimitris Stripelis, Zijian Hu, Jipeng Zhang, Zhaozhuo Xu, Alay Dilipbhai Shah, Han Jin, Yuhang Yao, Salman Avestimehr, and Chaoyang He. TensorOpera Router: A multi-model router for efficient LLM inference, 2024.
 - [14] Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E. Gonzalez, M Waleed Kadous, and Ion Stoica. RouteLLM: Learning to route LLMs with preference data, 2025.
 - [15] Adnan Sabovic, Jaron Fontaine, Eli De Poorter, and Jeroen Famaey. Energy-aware tinyML model selection on zero energy devices. *Internet of Things*, 30:101488, 2025.
 - [16] Orr Zohar, Shih-Cheng Huang, Kuan-Chieh Wang, and Serena Yeung. LOVM: Language-only vision model selection, 2023.
 - [17] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. ToolLLM: Facilitating large language models to master 16000+ real-world APIs, 2023.
 - [18] Siddharth Samsi, Dan Zhao, Joseph McDonald, Baolin Li, Adam Michaleas, Michael Jones, William Bergeron, Jeremy Kepner, Devesh Tiwari, and Vijay Gadepally. From words to watts: Benchmarking the energy costs of large language model inference, 2023.
 - [19] Yian Zhao, Wenyu Lv, Shangliang Xu, Jinman Wei, Guanzhong Wang, Qingqing Dang, Yi Liu, and Jie Chen. DETRs beat YOLOs on real-time object detection, 2024.
 - [20] Vanessa Mehlin, Sigurd Schacht, and Carsten Lanquillon. Towards energy-efficient deep learning: An overview of energy-efficient approaches along the deep learning lifecycle, 2023.
 - [21] Zeyu Yang and Wesley Armour. The hidden joules: Evaluating the energy consumption of vision backbones for progress towards more efficient model inference. In *Forty-second International Conference on Machine Learning*, 2025.
 - [22] Qunyou Liu, Darong Huang, Marina Zapater, and David Atienza. GreenLLM: SLO-Aware dynamic frequency scaling for energy-efficient LLM serving, 2025.
 - [23] Marcello Bullo, Seifallah Jarak, Pietro Carnelli, and Deniz Gündüz. Energy-aware dynamic neural inference, 2024.
 - [24] Adam Lechowicz, Rohan Shenoy, Noman Bashir, Mohammad Hajiesmaili, Adam Wierman, and Christina Delimitrou. Carbon- and precedence-aware scheduling for data processing clusters. In *Proceedings of the ACM SIGCOMM 2025 Conference*, page 1241–1244, New York, USA, 2025.
 - [25] Francisco Romero, Qian Li, Neeraja J. Yadwadkar, and Christos Kozyrakis. INFaaS: Automated model-less inference serving. In *Proceedings of the 2021 USENIX Annual Technical Conference (USENIX ATC 2021)*, pages 397–411, 2021.
 - [26] Haoran Qiu, Weichao Mao, Archit Patke, Shengkun Cui, Saurabh Jha, Chen Wang, Hubertus Franke, Zbigniew T. Kalbarczyk, Tamer Başar, and Ravishankar K. Iyer. Power-aware deep learning model serving with μ -serve. In *Proceedings of the 2024 USENIX Annual Technical Conference (USENIX ATC 2024)*, pages 75–93, Santa Clara, USA, 2024.
 - [27] Duo Wu, Jinghe Wang, Yuan Meng, Yanning Zhang, Le Sun, and Zhi Wang. CATP-LLM: Empowering large language models for cost-aware tool planning, 2025.