The Price of Prompting: Profiling Energy Use in Large Language Models Inference

Erik Johannes Husom, SINTEF.

erik.johannes.husom@sintef.no

Lwin Khin Shar,

Singapore Management University, lkshar@smu.edu.sg

Arda Goknil,

SINTEF,

arda.goknil@sintef.no

Sagar Sen, SINTEF.

sagar.sen@sintef.no

Abstract

In the rapidly evolving realm of artificial intelligence, deploying large language models (LLMs) poses increasingly pressing computational and environmental challenges. This paper introduces MELODI – Monitoring Energy Levels and Optimization for Data-driven Inference – a multifaceted framework crafted to monitor and analyze the energy consumed during LLM inference processes. MELODI enables detailed observations of power consumption dynamics and facilitates the creation of a comprehensive dataset reflective of energy efficiency across varied deployment scenarios. The dataset, generated using MELODI, encompasses a broad spectrum of LLM deployment frameworks, multiple language models, and extensive prompt datasets, enabling a comparative analysis of energy use. Using the dataset, we investigate how prompt attributes, including length and complexity, correlate with energy expenditure. Our findings indicate substantial disparities in energy efficiency, suggesting ample scope for optimization and adoption of sustainable measures in LLM deployment. Our contribution lies not only in the MELODI framework but also in the novel dataset, a resource that can be expanded by other researchers. Thus, MELODI is a foundational tool and dataset for advancing research into energy-conscious LLM deployment, steering the field toward a more sustainable future. The released code and dataset are available at https://github.com/ejhusom/MELODI.

1 Introduction

The swift progression of artificial intelligence (AI) has precipitated the emergence of advanced natural language processing (NLP) systems. Large language models (LLMs) have ascended to prominence among these systems, proving indispensable across a vast spectrum of applications. Their utility ranges from simple text predictions to managing complex dialogues. Concurrently, as these models grow in computational complexity, their energy demands escalate correspondingly [25, 32, 31]. In an epoch that prioritizes sustainable computing, the energy expenditure of LLMs is a consideration of mounting importance and one that commands focused attention [5, 17].

In the current literature, several endeavors aim to quantify and mitigate the environmental impact of AI. Prior studies have predominantly concentrated on the training phase of traditional machine learning (ML) models, with the inference phase often relegated to a secondary focus despite its cumulative energy footprint over the lifespan of a model [28]. Tools such as Green Algorithms [16] provide estimations of carbon emissions for computational tasks and yet lack real-time monitoring capabilities. Furthermore, while initiatives like CodeCarbon [6] and its integration to ML pipelines [13] offer

a more automated tracking of emissions, they do not cater specifically to the nuanced demands of LLMs nor provide granular data concerning inference-related energy consumption.

This paper introduces the MELODI framework designed to monitor and analyze energy usage during LLM inference. The urgent need for such a framework stems from the growing environmental and economic costs associated with the operation of LLMs. MELODI leverages two specialized power usage monitoring tools (i.e., Scaphandre [27] and Nvidia-smi [22]). Scaphandre tracks the CPU's power consumption for the LLM process, while nvidia-smi measures the GPU's total power usage. For accurate measurements, the GPU must be exclusively used for the LLM inference, with no concurrent processes. MELODI records the energy consumption for each inference process, in addition to the prompt and the LLM's response. Furthermore, it facilitates the collection of metadata associated with LLM inference tasks, enabling a deeper analysis of the relationship between task complexity and energy utilization.

An important part of our research involved creating a robust dataset essential for evaluating MELODI's capabilities. This included gathering a diverse set of prompts to monitor the energy consumption of various LLMs across different types of instructions. We deployed these models on different hardware to ensure a thorough analysis. The data collection includes not only energy consumption measurements but also related metadata, such as the size and complexity of prompts, the responses, and the time metrics for each inference process. This rich dataset helps us explore the relationships between prompt complexity and the energy expenditure they incur.

Contributions. In summary, we make three main contributions. Firstly, we have developed a framework integrating open-source tools to monitor power and energy consumption during the inference phase of LLMs. Secondly, we have undertaken an extensive data collection initiative involving a variety of language model deployment frameworks, multiple LLMs, and diverse prompt datasets to facilitate a comprehensive comparison of energy consumption across these dimensions. Thirdly, we offer a generalized analysis that elucidates the relationships between prompt characteristics, such as length and features, and the energy consumption patterns observed within the collected data.

2 Background

2.1 Large Language Models (LLMs)

LLMs represent a significant advancement in natural language processing (NLP) and deep learning (DL). These models, such as OpenAI's GPT-3.5 architecture, are based on transformer neural networks, which have revolutionized NLP tasks by capturing long-range dependencies in text [29]. LLMs, trained on vast text datasets, learn intricate language patterns and relationships through unsupervised learning, predicting subsequent text from context. The integration of LLMs into digital ecosystems is underpinned by advanced APIs and deployment frameworks. These APIs are essential tools that allow developers to use LLM features for prompt submission and response generation, making it possible to add LLM capabilities to various applications. Among the plethora of tools, several APIs and frameworks, such as Ollama [23], OpenAI's GPT API [24], Huggingface's Transformers [33], Llama.cpp [18], GPT4All [2], and vLLM [15], have emerged as pivotal to the LLM deployment landscape.

2.2 Monitoring Tools and Technologies

Advanced monitoring tools and technologies are essential for tracking and optimizing energy consumption, thereby mitigating the environmental impact of AI operations. In this context, MELODI employs the following two tools for their complementary features:

Scaphandre [27] is one of the few tools that provide detailed insights into the energy usage of computing systems at the process level. Its ability to deliver real-time power metrics makes it invaluable for assessing the energy efficiency of LLM deployments.

NVIDIA System Management Interface (nvidia-smi) [22] is a command-line utility providing the real-time monitoring of GPU, including power consumption, for systems equipped with NVIDIA GPUs. It can be used to evaluate the energy efficiency of LLMs leveraging GPU acceleration.

3 Dataset Collection: MELODI Framework for Monitoring Energy Usage

To understand the energy consumption patterns of LLMs, we created the MELODI framework, a tool designed to monitor energy usage during the inference processes of LLMs. Figure 1 provides an architecture of this framework. The core of MELODI's approach involves systematically sampling prompts from a curated collection, submitting these prompts to an LLM, and recording the prompts, their responses, metadata, and the energy metrics during inference. This data is cataloged to allow for a detailed analysis of energy consumption trends during the operational phase of LLMs.

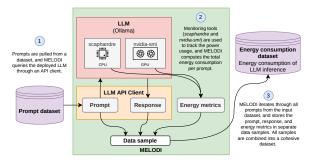


Figure 1: Overview of the MELODI framework.

- Input Data (Prompt Dataset): MELODI extracts prompts from a pre-existing dataset (prompt dataset in Figure 1) and uses them to query a deployed LLM through an API client.
- Power Usage Monitoring and Data Collection: MELODI integrates two power monitoring tools, Scaphandre [27] and nvidia-smi [22], to track the power consumption of the deployed LLM, such as the Ollama model depicted in Figure 1. These tools provide complementary functionalities: Scaphandre measures the CPU's power usage specific to the LLM process, while nvidia-smi assesses the total power usage of the GPU. It is crucial that no other processes operate on the GPU during LLM inference to ensure the accuracy of nvidia-smi's measurements. MELODI then calculates the energy consumption for each prompt, compiling data samples that encompass the prompt, the LLM's response, and the associated energy metrics, thus providing a granular view of energy utilization.
- **Dataset Generation:** MELODI iterates through all prompts from the dataset, capturing the prompt, its response, and energy metrics. These data samples are aggregated into a comprehensive dataset, encapsulating the energy consumption of the LLM during inference.

The MELODI framework is contingent upon LLM services that can be queried for responses. While MELODI is designed to be compatible with any service offering an API endpoint for response generation, our current implementation specifically caters to services that conform to either the Ollama- or OpenAI API specifications, encapsulating the majority of prevalent LLM deployment services. These API endpoints furnish not only the LLM's natural language responses but also critical metadata, including timestamps and token lengths for both prompts and responses. The landscape of LLM deployment tools is diverse, featuring platforms such as Llama.cpp [18], GPT4All [2], vLLM [15], Llamafile [19], and Ollama [23]. For the scope of our experiments, we chose to utilize Ollama [23] as our model deployment mechanism, attributing this choice to its straightforward installation process, dual compatibility with CPU and GPU architectures, and the expansive repository of open-source LLMs it supports for facile download and deployment.

Scaphandre [27] operates with process-level granularity, monitoring and recording the power usage of active processes on a computing device. We apply a regular expression to the process names to isolate the specific process associated with the LLM service. The frequency at which Scaphandre samples power metrics is configurable with granularity fine enough to reach the nanosecond scale. Although a higher sampling frequency increases the resolution and accuracy of the data, it also requires more storage for data retention and increases the energy consumption of Scaphandre itself. Therefore, we also monitor the energy consumption of Scaphandre, enabling an assessment of the energy overhead incurred by the measurement process. In contrast, nvidia-smi [22] tracks the power usage of the entire GPU and lacks the capability to distinguish among multiple processes concurrently utilizing the GPU. To guarantee precise measurements, we limit the GPU to solely operating the LLM inference process. Our monitoring tools measure the power consumption P of a process in microwatts (μ W). We then compute the energy consumption E by integrating the power consumption over the duration E0, applying the trapezoidal rule. With E1 expressed in seconds, we convert E2 into kilowatt-hours (kWh) by dividing by E3.6 × E10.12 facilitating a standardized energy usage assessment.

4 Dataset Composition

This section outlines the energy consumption dataset collected with MELODI, using diverse computing systems (from high-end servers to laptops) and multiple prompt datasets for several LLMs. Table 1 lists the prompt datasets in dataset collection, i.e., Alpaca [1]

Table 1: Prompt datasets.

	Prompt dataset	Avg. energy cons. (kWh)	Avg. response token length	
	Alpaca	4.09e-04	229.34	
	Code-feedback	6.93e-04	372.51	
-				

and Code-Feedback [12], and their average energy consumption and response token length. The Alpaca dataset, produced using OpenAI's text-davinci-003 engine, contains 52,000 prompts designed to enhance instruction-following capabilities in language models. dataset primarily includes text generation tasks and is tailored to train pre-trained language models to follow complex instructions more effectively. The Code-Feedback dataset supports the OpenCodeInterpreter model [35] designed to refine code by integrating code execution and human feedback. This dataset features 68,000 multi-turn interactions, combining user instructions with compiler responses to enhance model training in coding scenarios.

Table 2 summarizes the LLMs with their average energy Table 2: LLMs used in data collection. consumption and response token lengths. The models range from codelama-7b to llama3-70b, with energy consumption varying from as low as 8.30×10^{-5} kWh for gemma-2b to 2.26×10^{-3} kWh for llama3-70b. The average response token lengths also vary, with codelama-7b generating the longest responses at 451.52 tokens, while gemma-7b records shorter lengths at 249.41 tokens.

Model name/size	Avg. energy	Avg. response
WIOGCI Hallic/SIZC	cons. (kWh)	token length
codellama-7b	2.05e-04	451.52
codellama-70b	4.40e-03	330.04
gemma-2b	8.30e-05	279.32
gemma-7b	1.46e-04	249.41
llama3-8b	1.34e-04	255.20
llama3-70b	2.26e-03	251.46

Table 3 details the hardware used in our data collection. The diversity of the hardware used during data collection covers a wide range of processing power, memory capacities, and graphics capabilities, crucial for a comprehensive analysis of energy consumption across different computing environments.

Table 3: Hardware used in the data collection.

Machine	CPU	Memory	GPU	GPU memory
Server	AMD EPYC 7643 48-Core Processor	528GB	NVIDIA RTX A5000	24GB
Workstation	Intel Xeon W-2223 8-core @ 3.6GHz	128GB	NVIDIA RTX A2000	12GB
Laptop 1	Intel i5 11th Gen 12-core @ 2.4GHz	16GB	None	None
Laptop 2	Intel i7 10th Gen 12-core @ 2.7GHz	32GB	NVIDIA Quadro RTX 4000	8GB

Table 4 provides a detailed summary of the energy consumption dataset obtained using MELODI, capturing data from different hardware setups and LLMs under varying operational conditions. The dataset encompasses an array of data points per sample. These include the prompt and response; token lengths for both prompt and response; the timestamp of the API query; time series data on the power usage of the LLM service during inference and of Scaphandre during the monitoring period; and the aggregated energy consumption of the LLM service during inference.

Table 4: Overview of the energy consumption dataset.

	Tuest it ever the energy consumption duties.						
ID	Prompt dataset	Model	Model size	Hardware	No. of prompts	Avg energy cons.	Avg response
110						per response (kWh)	token length
1	codefeedback	codellama	7b	workstation	3084	1.83e-04	431.13
2	codefeedback	codellama	7b	laptop1	5295	1.85e-04	403.63
3	codefeedback	codellama	7b	laptop2	3555	2.47e-04	520.32
4	codefeedback	codellama	70b	workstation	161	4.40e-03	330.04
5	alpaca	gemma	2b	laptop1	5295	1.85e-04	403.63
6	alpaca	gemma	2b	workstation	11828	3.65e-05	187.91
7	codefeedback	gemma	2b	workstation	9897	7.30e-05	318.22
8	codefeedback	gemma	2b	laptop2	4972	7.36e-05	305.29
9	alpaca	gemma	2b	laptop2	5101	4.70e-05	181.52
10	alpaca	gemma	7b	laptop2	5099	9.81e-05	160.60
11	codefeedback	gemma	7b	workstation	5885	1.81e-04	338.23
12	alpaca	gemma	7b	workstation	8735	1.05e-04	165.09
13	codefeedback	gemma	7b	laptop2	3387	2.01e-04	333.73
14	alpaca	llama3	8b	laptop2	5101	1.34e-04	255.20
15	alpaca	llama3	70b	server	1026	2.26e-03	251.46

Dataset Analysis

We analyze the energy consumption dataset to address three Research Questions (RQ)s:

• RQ1. How does the energy consumption of LLM inference vary across different hardware, models, and prompt datasets?

- RQ2. What is the relationship between prompt complexity, response characteristics, and the energy consumption of LLMs during the inference process?
- RQ3. Can we develop a predictive model that accurately forecasts the energy consumption of LLMs based on prompt features and response characteristics?

5.1 Analysis Setup

Data Preparation. Our raw data consists of power usage over time, which we aggregate to calculate the energy consumption for each inference process as outlined in Section 3. These aggregated values are then compiled into distinct energy consumption datasets, categorized by the specific hardware, model, and prompt dataset utilized during the inference.

Statistical Analysis. Our analysis of the energy consumption datasets focuses on two metrics: energy per token and energy per response. Energy consumption per token is crucial as it evaluates resource use across different inference setups regardless of text volume. Meanwhile, energy per response is beneficial for estimating resource usage relative to the input prompt. We conduct statistical analyses on these metrics to compare across model types, sizes, hardware setups, and the specific prompt datasets used during inference, providing a comprehensive view of energy dynamics.

Visualization. We employ box plots to visualize the distribution of energy consumption for each model and hardware configuration (see Figure 2). The box plots depict the median energy consumption as a red line, the interquartile range with the main box, and extend to 1.5 times the interquartile range with whiskers, providing a clear summary of variability and central tendency.

Interpretation. We interpret results through a comparative analysis of energy consumption across different scenarios. The relationship between prompt complexity and energy consumption is assessed using the Pearson correlation coefficient, which ranges from 0 to 1 to indicate variable correlation. Additionally, predictive models for energy consumption are evaluated using the \mathbb{R}^2 score.

5.2 Results

5.2.1 RQ1 (energy consumption of LLM inference across different setups)

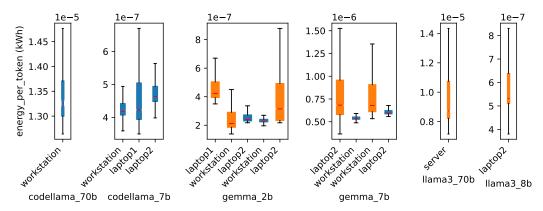
To address RQ1, we performed a comparative analysis of energy consumption across various LLMs operating on different hardware. The box plots in Figure 2 graphically represent this analysis, showing the energy usage of LLMs during inference on diverse hardware setups.

Model Size. Figure 2a categorizes energy consumption per token by the model utilized for inference. Notably, the largest models (codellama-70b with Code-Feedback and llama3-70b with Alpaca) exhibit energy usage approximately 100 times greater than their smallest counterparts (codellama-7b with Code-Feedback and llama3-8b with Alpaca). Additionally, the larger models (gemma-7b with both Alpaca and Code-Feedback) show an energy consumption that is roughly ten times higher than their smaller versions (gemma-2b with Alpaca and Code-Feedback).

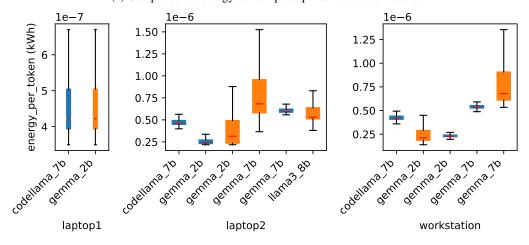
Hardware. Figure 2b organizes results by the hardware used for model execution. Energy consumption is noticeably higher when LLM tasks are run on the laptops than on the workstation. Specifically, Laptop1 (CPU-only) shows similar energy usage levels when operating different models (codellama-7b with Code-Feedback and gemma-2b with Alpaca), suggesting inefficiencies in CPU-based processing for LLM tasks since we would expect a significant difference when running a 2b model vs a 7b model. This observation could also reflect the limitations of the CPU-monitoring tool Scaphandre, as the power metrics are only estimates [14]. Comparisons between Laptop2 and the workstation, despite the latter's more robust GPU (12GB vs. 8GB), reveal roughly equivalent energy usage, underscoring potential disparities in hardware efficiency.

Model. In energy assessments of different models of the same size (codellama-7b with Code-Feedback vs. gemma-7b with Code-Feedback) run on both a workstation and Laptop2, gemma-7b consistently displayed higher energy consumption than codellama-7b. This observation suggests that codellama-7b is a slightly more energy-efficient model than gemma-7b.

Prompt Datasets. To evaluate the energy consumption per token across two prompt datasets, we analyzed the box plots for gemma-2b and gemma-7b in Figure 2a, the only models tested with both prompt datasets. The analysis reveals that, except in one instance, the Alpaca dataset consistently led



(a) Comparison of energy consumption per token across models.



(b) Comparison of energy consumption per token across hardware setups (Codellama-70b is excluded for the workstation to not skew the plot).

Figure 2: Comparative energy consumption per token across inference setups.

to higher energy consumption per token. Additionally, the interquartile range, representing variability in energy consumption, is notably wider in all cases, suggesting greater fluctuation in energy use when using the Alpaca dataset.

Table 4 facilitates a comparison of energy consumption per prompt, revealing that the median consumption for the Code-Feedback prompt dataset exceeds that of the Alpaca dataset during inference. This discrepancy likely arises from the longer average response length associated with the Code-Feedback prompts, which at 373 words surpasses the 229 words of the Alpaca prompts.

RQ1 Conclusion. Our analysis indicates marked disparities in energy consumption among different LLMs and hardware configurations. Larger models like codellama-70b and llama3-70b use roughly 100 times more energy per token than smaller ones. Energy demands are higher for LLMs running on laptops than workstations, particularly due to inefficiencies in CPU processing. Models of the same size exhibit varying energy efficiencies. Notably, the Code-Feedback dataset leads to greater energy use per token than the Alpaca dataset, likely due to longer response lengths.

5.2.2 RQ2 (the relationship between prompt complexity, response characteristics and the energy consumption)

To address RQ2, we used Python frameworks—Spacy [11], nltk [4], textblob [20], and text-stat [30]—to derive 57 text-based features. We integrated these features with our energy consumption dataset and calculated the Pearson correlation coefficient between energy consumption per prompt and each feature. Table 5 presents the top 20 features from this analysis.

Our analysis revealed that significant correlations with energy consumption per prompt are primarily linked to response characteristics rather than the complexity of the prompt. Key factors such as response token length, response duration, and total inference duration show strong positive correlations with energy consumption, with coefficients of 0.846, 0.625, and 0.618, respectively. This observation suggests that energy consumption escalates with an increase in response

Table 5: Correlations between features and energy consumption per response, calculated across all samples of our datasets. All features are derived from the input prompt, except those marked with a star (*), which are obtained from the response.

Feature	Correlation	Feature	Correlation
*response_token_length	0.846	reading_time	0.080
*response_duration	0.625	lexicon_count	0.078
*total_duration	0.618	prompt_token_length	0.075
adjectives	0.113	adverbs	0.073
adj_count	0.113	adverb_count	0.073
polysyllabcount	0.095	stop_word_count	0.072
long_word_count	0.090	noun_count	0.070
syllable_count	0.083	monosyllabcount	0.068
letter_count	0.082	word_count	0.068
char_count	0.080	verb_count	0.067

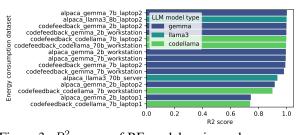
tokens due to more extensive processing within the model, highlighting that longer and more timeintensive responses significantly drive higher energy usage.

Upon analyzing correlations with prompt complexity features, we observe only low positive correlations. Attributes like adjective count, syllable count, and long word count demonstrate modest correlations with energy consumption, with coefficients of 0.113, 0.083, and 0.090, respectively. These results indicate that prompt complexity (given our set of features) has a minimal impact on energy consumption. Instead, the length and duration of responses are more significant factors. This insight suggests that optimizing the response generation process could be more effective in reducing energy consumption than merely simplifying the input prompts.

RQ2 Conclusion. The results show that response characteristics such as token length and duration are strongly correlated with energy usage, indicating higher consumption with longer responses. Conversely, prompt complexity features like adjective and syllable counts exhibit low correlations, suggesting minimal impact on energy use. These findings emphasize that managing response generation may offer more significant energy savings than simplifying prompts.

5.2.3 RQ3 (predictive model for the energy consumption of LLMs)

To address RQ3, we developed ML models using the energy consumption dataset, focusing on either prompt or response characteristics as input features. Our exploratory analysis identified Random Forest (RF) as the optimal model for effectively handling both types of input features. This choice was based on their performance and suitability for the data structure encountered in our dataset. Hence, we developed RF models (using



data structure encountered in our dataset. Figure 3: R^2 scores of RF models using only response length for forecasting energy consumption per response. either response token length or prompt characteristics) across hardware setups, models, and prompt datasets (one RF model for one LLM utilizing one prompt dataset on one hardware setup).

Table 5 shows that energy consumption highly correlates with response characteristics such as token length and duration. Figure 3 presents the R^2 scores of RF models using only response token length, demonstrating high predictive performance in most cases, with the majority achieving R^2 greater than 0.97 and an average R^2 score of 0.94. However, the RF models for laptop1 with only a CPU show a lower R^2 of 0.74, suggesting potential inaccuracies in CPU-based energy consumption estimates. The findings reveal that constraining the length of responses—by explicitly specifying the maximum text length in the prompts—could significantly reduce energy consumption during model operations.

Predictive models relying solely on prompt characteristics exhibit limited performance, as depicted in Figure 4. Analysis of the \mathbb{R}^2 scores across different models reveals variability in predictive accuracy. For example, the RF models for LLMs like gemma-7b, utilized with the Code-Feedback dataset on both a workstation and Laptop2, as well as gemma-2b using the Code-Feedback dataset on the

workstation, show relatively higher R^2 scores, ranging between 0.549 and 0.580. This demonstrates some levels of predictive effectiveness, particularly under specific hardware and dataset conditions.

Figure 4 illustrates that RF models for some LLMs, such as codellama-7b using Code-Feedback on laptop1 and laptop2, demonstrate notably poor predictive performance, with negative R^2 scores signaling ineffective model fitting and potential discrepancies in model predictions. Specifically, the RF model for codellama-7b utilizing Code-Feedback on laptop2 shows an R^2 score of -205.897, indicating a sig-

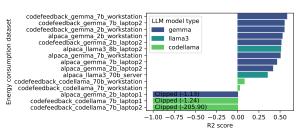


Figure 4: \mathbb{R}^2 scores of RF models using prompt characteristics for forecasting energy consumption per response.

nificant divergence between predicted and actual energy consumption. Generally, the RF models for gemma models show more predictable energy consumption patterns than llama models. Nonetheless, the small dataset size limits the robustness of these conclusions.

RQ3 Conclusion. Random Forest models showed strong predictive abilities, especially when using response characteristics, with most achieving an \mathbb{R}^2 score over 0.97. However, models for CPU-only laptops performed poorly, indicating potential inaccuracies in CPU-based energy estimates. Overall, models trained on response features were more effective than those using prompt characteristics, highlighting areas for energy optimization in LLM operations.

5.3 Threats to validity

Threats to internal validity. The primary concern regarding the internal validity of our data collection and analysis is the accuracy of the power usage monitoring tools. The precision of the nvidia-smi tool, as noted in the documentation of NVIDIA [7], has a stated accuracy within ± 5 Watts. However, a study by Yang et al. [34] suggests the error margin might reach up to $\pm 5\%$. Accurate GPU power usage readings also necessitate that no other processes utilize the GPU concurrently. While we ensured this in most experiments, some base processes remained active on the GPU of the workstation, though their impact on baseline energy consumption was minimal. Scaphandre only provides an estimate of the power usage based on the RAPL interface of Intel processors [27] without quantifying accuracy, potentially affecting the reliability of our energy consumption data. Additionally, power efficiency variations in identical hardware models could introduce more variability to our results. Future research should utilize more precise power measurement tools and standardized tests across different hardware to better understand and quantify these discrepancies.

Threats to external validity. The generalizability of our results is constrained by the limited range of inference scenarios analyzed, particularly in terms of hardware setups and LLMs utilized. The small sample size makes it challenging to draw robust conclusions, and the extent to which our findings can be generalized depends on the representativeness of the selected hardware and model configurations relative to the broader array of LLM deployments.

6 Discussion and Future Work

Our study provides insights into the energy consumption of LLM inference while identifying areas needing further research. One limitation identified is the accuracy of CPU-based power monitoring, which may compromise our results. Despite using the Scaphandre tool for measuring CPU power, we observed similar energy consumption across models of significantly different sizes, suggesting possible measurement imprecision. This was particularly evident in the reduced predictive performance of energy consumption on CPU-only setups, as opposed to configurations equipped with GPUs.

Our results reveal that response token length is a reliable indicator of energy consumption, suggesting that managing response lengths could effectively control energy use. However, the connection between prompt complexity and energy consumption is less definitive. While models based on prompt characteristics show potential, their performance varies. Future research could explore the potential of using LLMs or other NLP models to analyze prompts for energy consumption prediction, balancing the benefits against the potential energy costs of such analyses.

Our current analysis does not address the quality of responses generated. Our dataset, however, provides a rich foundation for future research to explore the relationship between response quality and energy consumption. Understanding whether higher-quality responses inherently require more energy could inform the development of more efficient models and prompt designs.

To enhance our study's generalizability, we plan to test a wider array of hardware configurations and model types, allowing us to better understand energy consumption variations. Comparing the energy efficiency of task-specific models against general-purpose models in different applications would also be valuable. This analysis could reveal if specialized models yield significant energy savings for specific tasks, ensuring that comparisons maintain a consistent base model for accuracy.

Future research could explore patterns in energy consumption across models, focusing beyond mere size to include architectural differences. This could reveal how specific design choices impact energy demands and help optimize models for better energy efficiency without sacrificing performance.

7 Related Work

The merging of NLP and sustainable computing, emphasized in studies like the investigation of Strubell et al. [28] into the energy use and carbon footprint of NLP model training, is a pivotal aspect of AI. These studies stress the need for energy-efficient practices across AI model lifecycles and have spurred extensive research into sustainable AI. Meanwhile, initiatives by Henderson et al. [10] advocate for transparent reporting of energy consumption, enhancing accountability, especially during the training phase. Tools like Green Algorithms [16], CodeCarbon [6], and CarbonTracker [3] have been developed to estimate computational carbon emissions, yet they often provide general rather than detailed analyses necessary for large-scale LLM deployment.

Luccioni et al. [21] and Desislavov et al. [8] found that general-purpose ML models are more energy-intensive than task-specific ones and that efficiency gains in NLP and computer vision do not necessarily correlate with increased energy consumption. In contrast, MELODI, focusing on the inference stage, delivers real-time energy assessments of LLMs by utilizing Scaphandre and nvidia-smi to directly measure power usage. This method provides detailed energy metrics at the API level, offering a more nuanced evaluation of energy use than broader estimates and enhancing MELODI's role in sustainable AI.

Recent studies, including those by Liang et al. [17] and Samsi et al. [26], focus on sustainable AI by analyzing the LLMs' energy usage and carbon footprint. These studies highlight energy demands across various LLM configurations. However, MELODI enhances these efforts by providing a comprehensive framework that measures LLM energy consumption by leveraging Scaphandre and nvidia-smi for real-time energy tracking of both CPUs and GPUs. This capability enables a more dynamic energy efficiency optimization in LLM deployment, distinguishing MELODI with its predictive and monitoring strengths.

Everman et al. [9] analyzed the carbon footprint of open-source LLMs during inference, using the Software Carbon Intensity (SCI) to gauge environmental impact. Their findings debunked the notion that higher carbon footprints correlate with better model quality and highlighted GPUs' role in reducing emissions. In contrast, MELODI supports real-time power monitoring during inference, generating data that serves as a benchmark and facilitates the development of ML models for predicting energy consumption. Thus, while Everman et al. focus on comparative carbon impact analysis, MELODI enables continuous monitoring and predictive energy management.

8 Conclusion

In conclusion, MELODI represents a noteworthy leap forward in advancing energy-aware practices within the domain of LLMs. By formulating a precise methodology and unveiling an open-source instrument tailored for real-time monitoring of energy consumption throughout the LLM inference process, MELODI effectively bridges a vital chasm in sustainable computing. Our data collection, encompassing various LLM architectures, hardware setups, and prompt datasets, provides a robust platform for rigorous comparison and nuanced analysis. This endeavor facilitates a deeper understanding of the energy dynamics associated with LLM deployment scenarios, offering critical insights that contribute to the optimization and sustainability of LLM applications.

Acknowledgments and Disclosure of Funding

The work has been conducted as part of the ENFIELD project (101120657) funded by the European Commission within the HEU Programme.

References

- [1] Alpaca. https://huggingface.co/datasets/tatsu-lab/alpaca, Visited in 2024.
- [2] Yuvanesh Anand, Zach Nussbaum, Adam Treat, Aaron Miller, Richard Guo, Ben Schmidt, GPT4All Community, Brandon Duderstadt, and Andriy Mulyar. Gpt4all: An ecosystem of open source compressed language models. *arXiv preprint arXiv:2311.04931*, 2023.
- [3] Lasse F Wolff Anthony, Benjamin Kanding, and Raghavendra Selvan. Carbontracker: Tracking and predicting the carbon footprint of training deep learning models. *arXiv* preprint *arXiv*:2007.03051, 2020.
- [4] Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python: analyzing text with the natural language toolkit.* "O'Reilly Media, Inc.", 2009.
- [5] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [6] CodeCarbon. https://codecarbon.io/, Visited in 2024.
- [7] NVIDIA Corporation. Nvidia system management interface, 2023. 2024-05-03.
- [8] Radosvet Desislavov, Fernando Martínez-Plumed, and José Hernández-Orallo. Trends in ai inference energy consumption: Beyond the performance-vs-parameter laws of deep learning. *Sustainable Computing: Informatics and Systems*, 38:100857, 2023.
- [9] Brad Everman, Trevor Villwock, Dayuan Chen, Noe Soto, Oliver Zhang, and Ziliang Zong. Evaluating the carbon impact of large language models at the inference stage. In *IEEE international performance, computing, and communications conference (IPCCC)*, pages 150–157. IEEE, 2023.
- [10] Peter Henderson, Jieru Hu, Joshua Romoff, Emma Brunskill, Dan Jurafsky, and Joelle Pineau. Towards the systematic reporting of the energy and carbon footprints of machine learning. *Journal of Machine Learning Research*, 21(248):1–43, 2020.
- [11] Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear, 2017.
- [12] https://huggingface.co/datasets/m-a-p/Code-Feedback. https://huggingface.co/datasets/tatsu-lab/alpaca, Visited in 2024.
- [13] Erik Johannes Husom, Sagar Sen, and Arda Goknil. Engineering carbon emission-aware machine learning pipelines. In *IEEE/ACM 3th International Conference on AI Engineering—Software Engineering for AI (CAIN)*, 2024.
- [14] Mathilde Jay, Vladimir Ostapenco, Laurent Lefèvre, Denis Trystram, Anne-Cécile Orgerie, and Benjamin Fichel. An experimental comparison of software-based power meters: focus on cpu and gpu. In 2023 IEEE/ACM 23rd International Symposium on Cluster, Cloud and Internet Computing (CCGrid), pages 106–118. IEEE, 2023.
- [15] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pages 611–626, 2023.
- [16] Loïc Lannelongue, Jason Grealey, and Michael Inouye. Green algorithms: quantifying the carbon footprint of computation. *Advanced science*, 8(12):2100707, 2021.
- [17] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*, 2022.
- [18] Llama.cpp. https://github.com/ggerganov/llama.cpp, Visited in 2024.

- [19] Llamafile. https://github.com/Mozilla-Ocho/llamafile, Visited in 2024.
- [20] Steven Loria. textblob documentation. Release 0.15, 2, 2018.
- [21] Alexandra Sasha Luccioni, Yacine Jernite, and Emma Strubell. Power hungry processing: Watts driving the cost of ai deployment? *arXiv preprint arXiv:2311.16863*, 2023.
- [22] NVIDIA. https://developer.nvidia.com/nvidia-system-management-interface, Visited in 2024.
- [23] Ollama. https://github.com/ollama/ollama, Visited in 2024.
- [24] OpenAI. https://platform.openai.com/, Visited in 2024.
- [25] Matthias C Rillig, Marlene Ågerstrand, Mohan Bi, Kenneth A Gould, and Uli Sauerland. Risks and benefits of large language models for the environment. *Environmental Science & Technology*, 57(9):3464–3466, 2023.
- [26] Siddharth Samsi, Dan Zhao, Joseph McDonald, Baolin Li, Adam Michaelas, Michael Jones, William Bergeron, Jeremy Kepner, Devesh Tiwari, and Vijay Gadepally. From words to watts: Benchmarking the energy costs of large language model inference. In *IEEE High Performance Extreme Computing Conference (HPEC)*, pages 1–9. IEEE, 2023.
- [27] Scaphandre. https://github.com/hubblo-org/scaphandre, Visited in 2024.
- [28] Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for deep learning in nlp. *arXiv preprint arXiv:1906.02243*, 2019.
- [29] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. Advances in neural information processing systems, 27, 2014.
- [30] textstat. https://github.com/textstat/textstat, Visited in 2024.
- [31] Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359*, 2021.
- [32] Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, Po-Sen Huang, John Mellor, Amelia Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh, et al. Taxonomy of risks posed by language models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 214–229, 2022.
- [33] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- [34] Zeyu Yang, Karel Adamek, and Wesley Armour. Part-time power measurements: nvidia-smi's lack of attention, 2024.
- [35] Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhu Chen, and Xiang Yue. Opencodeinterpreter: Integrating code generation with execution and refinement. https://arxiv.org/abs/2402.14658, 2024.