# Comparing energy consumption and accuracy in text classification inference

Johannes Zschache and Tilman Hartwig

Application Lab for AI and Big Data, German Environment Agency, Alte Messe 6, Leipzig, 04103, Saxony, Germany.

*Corresponding author(s). E-mail(s): tilman.hartwig@uba.de;
Contributing authors: johannes.zschache@uba.de;

## Abstract

The increasing deployment of large language models (LLMs) in natural language processing (NLP) tasks raises concerns about energy efficiency and sustainability. While prior research has largely focused on energy consumption during model training, the inference phase has received comparatively less attention. This study systematically evaluates the trade-offs between model accuracy and energy consumption in text classification inference across various model architectures and hardware configurations. Our empirical analysis shows that the best-performing model in terms of accuracy can also be energy-efficient, while larger LLMs tend to consume significantly more energy with lower classification accuracy. We observe substantial variability in inference energy consumption ($<$mWh to $>$kWh), influenced by model type, model size, and hardware specifications. Additionally, we find a strong correlation between inference energy consumption and model runtime, indicating that execution time can serve as a practical proxy for energy usage in settings where direct measurement is not feasible. These findings have implications for sustainable AI development, providing actionable insights for researchers, industry practitioners, and policymakers seeking to balance performance and resource efficiency in NLP applications.

# 1 Introduction

Artificial intelligence (AI) systems, particularly large language models (LLMs), have driven remarkable progress in Natural Language Processing (NLP) applications. This

development has been enabled by the Transformer architecture (Vaswani et al., 2017) and exemplified by the emergence of large-scale models such as GPT-3 (Brown et al., 2020), which have significantly advanced task performance. However, this progress has come at a cost: the escalating energy demands of AI systems pose significant environmental and computational challenges. Data centers that support AI computations are major electricity consumers, often dependent on fossil fuels, thereby contributing to greenhouse gas emissions (Lacoste et al., 2019; Axenbeck et al., 2025). This increasing energy demand challenges global climate objectives such as the Paris Agreement (United Nations, 2015a) and the United Nations' Sustainable Development Goals (SDGs), specifically Goal 13 on climate action (United Nations, 2015b). Consequently, designing energy-efficient AI systems is imperative for aligning technological advancements with sustainability goals. Moreover, beyond sustainability, energy-efficient models offer additional advantages, including reduced hardware requirements, lower financial costs, and faster inference times.

When evaluating machine learning models, most studies concentrate on the quality of the model responses by tracking e.g. the accuracy, the RMSE, or other measures. And even if the energy consumption is taken into account, prior research has mainly focused on the training phase (Strubell et al., 2019; Patterson et al., 2021; Luccioni and Hernandez-Garcia, 2023). The inference phase, which is repeatedly executed in real world deployments, has received comparatively less attention. However, energy efficiency during the operational phase is an increasingly relevant topic as LLM applications become ubiquitous and LLM models are trained to use additional test-time compute to improve performance (OpenAI, 2024; DeepSeek-AI, 2025). Addressing this gap, we present a systematic study on the energy consumption of language models during inference, providing actionable insights for balancing accuracy with efficiency.

A particularly popular machine learning task is text categorization, a task that lightweight models have been shown to handle effectively. For instance, Joulin et al. (2017) show that a simple classifier built on word embeddings is often as accurate as deep learning classifiers. Despite this, some authors argue for the use of pre-trained LLMs for text classification because it reduces the need for model training and simplifies data preprocessing (Wang et al., 2024). Additionally, popular software tutorials promote LLMs for classification tasks (LangChain Team, 2023; Lamini Team, 2023), further encouraging their use even when more efficient alternatives exist. In order to justify the usage of LLM in relatively simple tasks such as text categorization, we advocate a consequent comparison of a model's response quality to its energy efficiency.

Given a practical use case that is occurring in public administration, our study empirically analyzes trade-offs between model accuracy and energy consumption across various language models and hardware configurations. We find that the best performing model is energy efficient while LLMs show higher energy usage with lower accuracy. Generally, we see significant variability in inference energy consumption, influenced by model type, model size, and hardware specifications. Additionally, the energy consumption during inference is shown to highly correlate with the model's runtime. This makes the duration of computations a valuable proxy measure for energy consumption in settings where the latter cannot be traced. Our findings have implications for researchers, industry practitioners, and policymakers advocating for sustainable

AI development (Kaack et al.; Luccioni et al., 2025). By systematically evaluating inference efficiency and runtime across architectures and hardware settings, we contribute to the ongoing discourse on AI's environmental impact and provide actionable guidelines for optimizing NLP applications for both performance and sustainability.

## 2 Previous research

Research on the environmental impact of machine learning (ML) has primarily focused on the energy consumption and carbon emissions produced during the training phase of large-scale models. Most famously, Strubell et al. (2019) quantify the carbon footprint of NLP models, revealing that the training of a single large-scale transformer model can emit as much carbon as five cars over their entire lifetimes (their measurements include thousands of hyperparameter tuning jobs, which makes it difficult to disentangle model-inherent efficiency from experimental setup). This seminal work spurred further investigations into the environmental costs of training neural networks, including large language models (Patterson et al., 2021; Luccioni and Hernandez-Garcia, 2023; Patterson et al., 2022).

While training remains a significant contributor to energy consumption, recent studies have begun to focus on the inference phase. Samsi et al. (2023) highlighted the substantial energy demands of LLM inference but did not explore the relationship between energy consumption and task-specific performance. Liu et al. (2022) underscore the importance of evaluating NLP models not just on efficiency metrics but also on accuracy by introducing the Efficient Language Understanding Evaluation (ELUE) benchmark. ELUE aims to establish a Pareto frontier that balances performance and efficiency. It includes various language understanding tasks, facilitating fair and comprehensive comparisons among models. However, the framework adopts number of parameters and FLOPs as the metrics for model efficiency, disregarding hardware specific factors. Similarly, Chien et al. (2023) estimate the energy consumption associated with the inference phase of generative AI applications based on the output word count and several assumptions about the application such as the number of FLOPS per inference and the sampling rate.

In contrast, we promote energy-efficient NLP models by the direct measurement of the power consumed during inference. Hence, our work follows the approach of the SustaiNLP 2020 shared task (Wang and Wolf, 2020). SustaiNLP demonstrated that substantial energy savings are achievable with minimal performance loss. While this study was limited to the performance of a couple of small language models on a single benchmark, we extend these efforts to a greater number of partially very large models deployed to a practical inference scenario.

This makes our study very similar to the one by Alizadeh et al. (2025), who investigated the trade-offs between accuracy and energy consumption when deploying large language models (LLMs) for software development tasks. Besides the finding that larger LLMs with higher energy consumption do not always yield significantly better accuracy, the authors demonstrated that architectural factors, such as feedforward layer size and transformer block count, directly correlate with energy usage.

3

Finally, Luccioni et al. (2024) provide one of the most comprehensive analyses of energy consumption during ML model inference. Their study systematically compared the energy costs of 88 models across 10 tasks and 30 datasets, including both smaller task-specific and larger multi-purpose models. They found that the larger models are orders of magnitude more energy-intensive than smaller task-specific ones, especially for tasks involving text and image generation. Furthermore, their research underscores the variability in energy consumption across tasks and model architectures. The authors advocate for increased transparency and sustainable deployment practices, emphasizing that the environmental costs of deploying large, multi-purpose AI systems must be carefully weighed against their utility.

# 3 Data and methods

Our experiments are inspired by an occasionally occurring use case in public administration: the management of objections that are submitted by the population. Due to a potentially very large amount of submissions, an automatic preprocessing of the objections is of high value. One of the possible steps of an automated workflow is to categorize each submission for optimal forwarding to the responsible department.

The data of our study originates from the process of selecting a repository site for high-level radioactive waste in Germany. During the first phase, sub-areas were identified and discussed in a process called FKTG (Fachkonferenz Teilgebiete). The statements from the population were categorized, processed and published as the FKTG-dataset (https://beteiligung.bge.de/index.php). The text of the submission is given by the column 'Beitrag' (input). The column 'Themenkomplex' (topic) contains the category of the text.

We scraped the dataset from the website and restricted it to entries for which the topic occurs at least 10 times. The remaining 378 entries were split into half: 189 entries for training and 189 entries for testing. This unusual 50:50 split was done so that the test set should be sufficiently representative by containing enough examples of each of the 14 categories. Each of the following experiments was repeated 10 times with different train-test-splits. To increase comparability, every experiment was run with the same 10 train-test-splits.

An experiment run consists of a training phase and a testing phase. Since large language models have been argued to be applicable to text categorization without training (zero-shot), we omit the training phase for these models and apply LLMs without fine-tuning. We report the energy consumption and accuracy only for the test phase as averages over all runs.

## 3.1 Traditional models

Besides LLMs, we initially run the experiments with lightweight NLP models that we call traditional because they have been used for categorization tasks long before LLMs existed. Specifically, we use a linear model (logistic regression) and a gradient boosting algorithm (xgboost). Logistic regression is a simple, interpretable model that estimates the probability of a class based on a linear combination of input features. XGBoost

(Extreme Gradient Boosting) is an efficient, scalable machine-learning algorithm that combines predictions from multiple decision trees to improve accuracy.

We consider three different types of features: bag-of-words (BoW), term frequency-inverse document frequency (TF-IDF), and a pretrained multilingual sentence embedding. BoW represents text by counting word occurrences without considering order, while TF-IDF adjusts word counts by their importance across documents, capturing rare but informative terms. The TF-IDF features are calculated on all 2-gram and 3-gram character sequences, which capture local patterns in the text. The multilingual sentence embedding (https://huggingface.co/sentence-transformers/paraphrase-multilingual-mpnet-base-v2) provides dense vector representations of text, preserving semantic meaning across languages. This embedding is not fine-tuned on the training data. Both models are trained using the default parameters provided by sklearn.linear_model.LogisticRegression and xgboost.XGBClassifier.

## 3.2 Large language models

Large language models (LLMs) were applied without training (zero-shot) using the test set only. Table 1 gives the names and sources of the models used. The LLMs were selected by the following criteria:

- availability on Huggingface
- support of german language
- capability of processing the `dspy`-prompt (see appendix A)

Additionally, Jamba Mini 1.5 was chosen as model with an alternative architecture that includes next to Transformer also Mamba layers (a state-space model). The Deepseek distillations (DS) were added to include models with reasoning capabilities (test-time compute).

| Model | Link |
|---|---|
| Llama 3.1 8B | https://huggingface.co/meta-llama/Meta-Llama-3.1-8B-Instruct |
| Llama 3.1 70B | https://huggingface.co/meta-llama/Meta-Llama-3.1-70B-Instruct |
| Qwen 2.5 7B | https://huggingface.co/Qwen/Qwen2-7B-Instruct |
| Qwen 2.5 72B | https://huggingface.co/Qwen/Qwen2-72B-Instruct |
| Phi 3.5 Mini | https://huggingface.co/microsoft/Phi-3.5-mini-instruct |
| Phi 3.5 MoE | https://huggingface.co/microsoft/Phi-3.5-MoE-instruct |
| Jamba Mini 1.5 | https://huggingface.co/ai21labs/AI21-Jamba-1.5-Mini |
| DS Qwen 14B | https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-14B |
| DS Qwen 32B | https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-32B |
| DS Llama 8B | https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-8B |
| DS Llama 70B | https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-70B |

**Table 1** Selection of large language models

## 3.3 Computing Resources

We used different computing systems for a comparative analysis of energy efficiency across diverse hardware architectures. This enables the assessment of how architectural

differences - especially GPU tensor core capabilities - affect the inference speed and power usage. A diversity in computational infrastructure is crucial for generalizing findings across different environments and ensuring the validity and replicability of experimental results in machine learning research. Furthermore, insights gained from using multiple platforms contribute to optimizing resource allocation strategies and improving cost-effectiveness in large-scale machine learning projects.

To run our experiments, we were granted access to the high-performance computing (HPC) systems of TUD Dresden University of Technology (https://doc.zih.tu-dresden.de/) and Leipzig University (https://www.sc.uni-leipzig.de/). For GPU-accelerated computing, three different systems are available named `Capella`, `Paula`, and `Clara` (see Table 2). The main difference for our study is the GPU: while a node on the `Capella` cluster is equipped with 4 x H100, there are 8 x A30 on each node on `Paula` and 4 x V100 on `Clara`. This means that a large model such as Llama 3.1 70B or Qwen 2.5 72B fits on a single node of `Capella` (requiring 2 GPUs) or `Paula` (requiring all 8 GPUs) but takes up two nodes of the `Clara` cluster (assuming a 16-bit floating point representation of the parameters).

| Cluster | Capella | Paula | Clara |
|---|---|---|---|
| HPC center | TUD Dresden University of Technology | Leipzig University | Leipzig University |
| number of nodes | 144 | 12 | 6 |
| CPU per node | 2 x AMD (32 cores) 2.7GHz | 2 x AMD (64 cores) 2.0GHz | 1 x AMD (32 cores) 2.0GHz |
| RAM per node | 768 GB | 1 TB | 512 GB |
| GPU per node | 4 x NVIDIA H100 (94GB) | 8 x NVIDIA A30 (24 GB) | 4 x NVIDIA V100 (32GB) |
| single GPU max power consumption | 700W | 165W | 250W |

**Table 2** HPC Resources

LLMs were deployed using the `vllm` library (https://github.com/vllm-project/vllm), which runs on a ray cluster (https://www.ray.io/) for multi-node computations. If a model is too large to be deployed on a single GPU, the model weights are distributed over multiple GPUs, which allow for a parallel computation of the activations (c.f. tensor model parallelism (TMP) in Bai et al., 2024, pp.16). In cases where two computing nodes are needed, the model is split into two parts and executed sequentially (c.f. pipeline model parallelism (PMP) in Bai et al., 2024, p.17): first the model part on the first node and then the model part on the second node.

The energy consumption and the runtime of the inference phase were measured by the CodeCarbon package (https://github.com/mlco2/codecarbon). This package uses the NVIDIA Management Library (NVML) and the Intel RAPL files to track the power usage of GPU and CPU (https://mlco2.github.io/codecarbon/methodology.html#power-usage). The power consumption of the memory is flatly added with 0.375W/GB of memory used. In settings where the model is deployed on more than one node, the inference duration is taken as the maximum and the energy as the sum over all nodes.

Various software tools have been created to monitor energy consumption during the application of machine learning models (https://github.com/tiingweii-shii/Awesome-Resource-Efficient-LLM-Papers?tab=readme-ov-file#%EF%B8%8F-energy-metrics). Similar to CodeCarbon, Carbontracker (Anthony et al., 2020) and experiment-impact-tracker (Henderson et al., 2020) estimate energy consumption by monitoring hardware usage. In some settings, CodeCarbon is considered more accurate, yielding values closer to those obtained via physical wattmeters (Bouza et al., 2023). Comparing different tools of energy monitoring is beyond the scope of our paper.

# 4 Results

For each model, we report accuracy, energy consumption, and inference duration. The energy consumption and duration were measured only for the inference step, i.e., after the model and data were already loaded. One inference run involves classifying 189 text samples from a test set. All tables and figures present the average results over 10 runs on different test sets, with the same 10 test sets used for each model. Measurement variance was generally low: $< 0.002$ for accuracy, and $< 0.2$ dex for both energy consumption and duration (logarithmically scaled to base 10).

Figure 1 illustrates the trade-off between energy consumption and accuracy across all models. For these experiments, a single node of the `Capella` system was used. The minimum number of H100 GPUs required varies by model (see Table B1).

The highest accuracy was achieved by a traditional linear model using pre-trained sentence embeddings. Notably, even the most energy-efficient model - a linear model with TF-IDF features - outperformed several large language models (LLMs). Among LLMs with relatively high accuracy, the best small model (Qwen 2.5 7B) consumes seven times less energy than the most accurate model (Qwen 2.5 72B), with only a minor accuracy reduction of 0.07 points. Deepseek models, despite their extensive reasoning processes during inference, exhibit lower accuracy than non-reasoning LLMs while consuming significantly more energy and taking longer to complete inference.

## 4.1 Analysis of hardware settings

This section analyzes the impact of different hardware configurations (see Tab. 2) on energy consumption. We focus on GPU usage due to its dominant role in machine learning inference.

As shown in Figure 2, GPU consumption accounts for the largest share of total energy usage in all experiments. The only exceptions are traditional models without embeddings, which do not use the GPU during inference.

### 4.1.1 Varying the Number of GPUs

We examined the effect of varying the number of GPUs on energy consumption and inference duration. Most LLMs were tested on 1, 2, or 4 GPUs on a single `Capella` system node. Larger models (Qwen 72B, Phi MoE, Llama 70B, Jamba Mini, and DS Llama 70B) required either 2 or 4 GPUs. Increasing the number of GPUs consistently
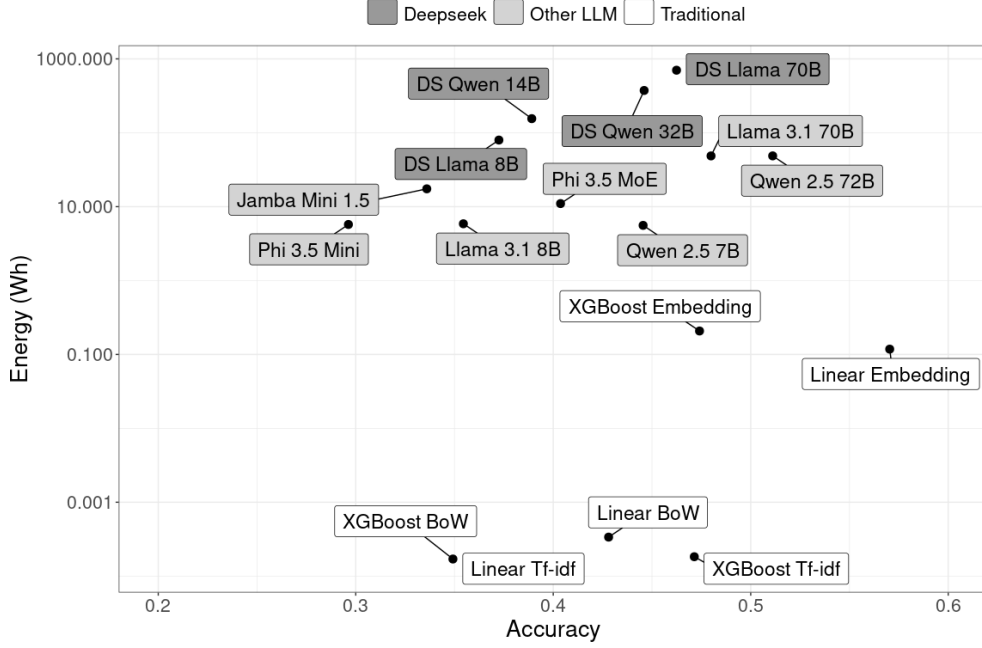
**Fig. 1** Accuracy-energy-trade-off of all models for the inference task on the `Capella` system (single node). The energy consumption for the same task spans over six orders of magnitude with traditional models being the most energy-efficient models and reasoning models are most energy-consuming. The best model for this specific task is a traditional model (Linear Embedding) with moderate energy consumption.

reduced inference duration but did not reduce energy consumption. In some cases, energy consumption increased due to the additional GPUs in operation (see Figure 3).

### 4.1.2 Varying the Number of Nodes

While large models can often be executed on a single computing node, certain hardware limitations or shared high-performance computing (HPC) environments may necessitate using multiple nodes. In shared systems, it is often easier to access two nodes with half the number of available GPUs than a single node with all its GPUs, due to scheduling constraints and resource allocation policies. However, deploying models across multiple nodes increases network communication overhead and significantly raises energy consumption.

We evaluated this effect for the largest models on the `Capella` system by comparing a 'single-node' configuration (2 GPUs on one node) with a 'double-node' configuration (1 GPU on each of two nodes). For the double-node configuration, energy consumption was summed across both nodes and averaged over 10 runs, while the reported duration reflects the average of the maximum value between the two nodes.

As shown in Figure 4, using two nodes increased energy consumption by a factor that depends on the model (see also Table B2). This increase stems from the overhead
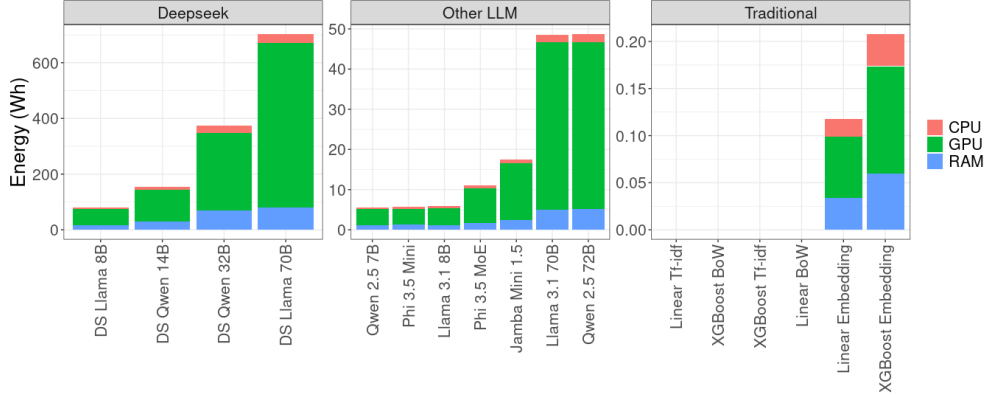
**Fig. 2** Energy consumption of all models for the inference task on the `Capella` system (single node)
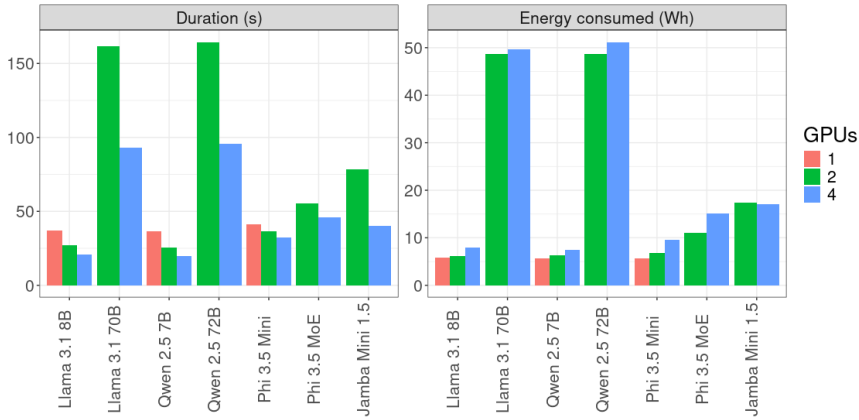


**Fig. 3** Effects of the number of GPUs on the runtime and consumed energy (`Capella`, single node). Deepseek models are not shown.

of coordinating across nodes. Inference duration also increased by the same factor due to the sequential execution of model components and the required inter-node communication.

### 4.1.3 Comparing GPU Architectures

Finally, we compared the energy efficiency of different GPU architectures (see Figure 5 and Table B3). Interestingly, the expected efficiency gains from using the more powerful H100 instead of V100 or A30 GPUs were only observed for the Deepseek models. This discrepancy is likely to arise because Deepseek models engage in extended reasoning by generating a larger output of words before making a classification decision. Consequently, the efficiency of H100 GPUs becomes evident only when substantial
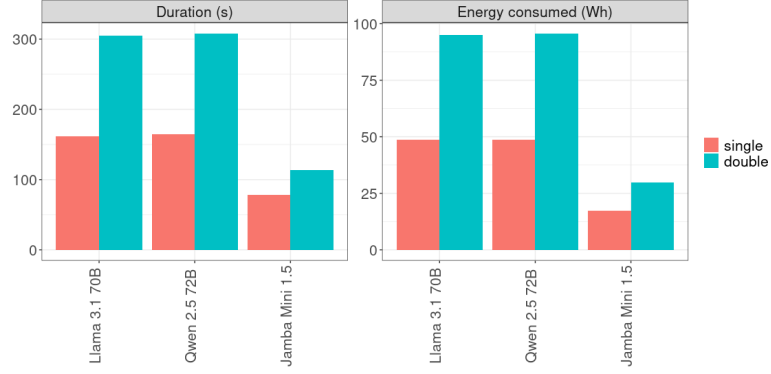
**Fig. 4** Comparison single node vs. double node deployment (`Capella`).

text is generated. For models generating a single token per inference, a V100 or even a A30 GPU is more efficient in inference.
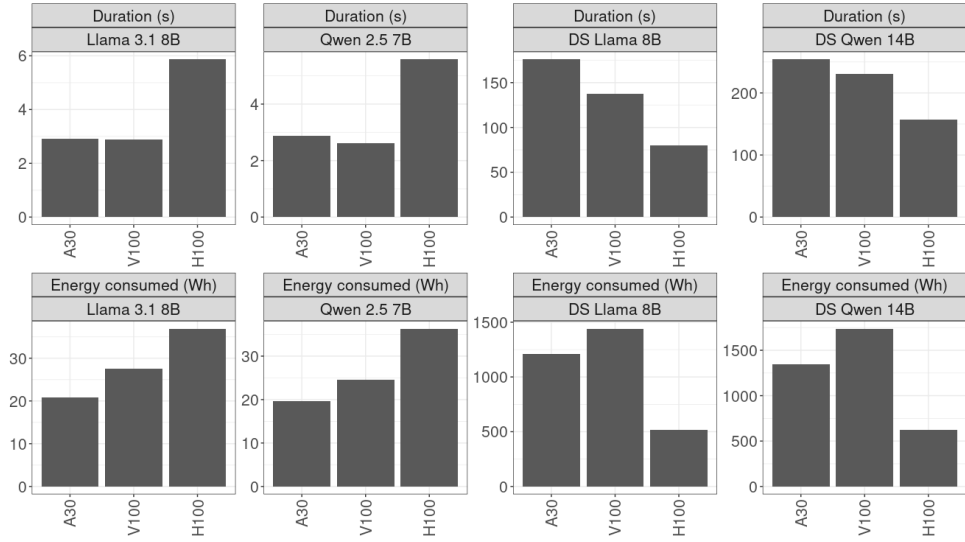


**Fig. 5** Comparison of different GPU cards: four exemplary LLMs. Single node deployment.

## 4.2 Linear relationship between duration and energy

In most of the tables in appendix B, we report both the duration of each inference run and its corresponding energy consumption. Since energy is the integral of power over time, these two measures exhibit a strong correlation. If the power is constant over time, this correlation should be linear. Figure 6 illustrates this relationship for

all experiments conducted on a single node of the `Capella` cluster. When controlling for the number of GPUs used for model deployment, the relation between duration and energy is approximately linear. Therefore, the duration appears to serve as a good proxy for the energy consumed.
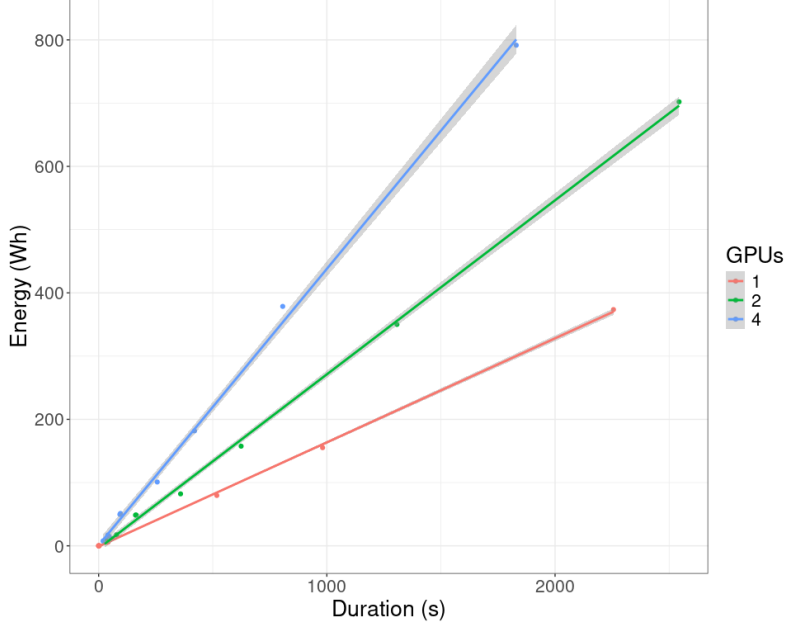


**Fig. 6** Plotting the relationship between duration and energy consumption (single node on `Capella`). The lines are added by running a linear regression model.

To further quantify the relationship between duration and energy consumption, we performed a linear regression analysis for each hardware configuration (see Table 3). This analysis includes all experiments, regardless of the number of nodes used for model deployment. The consistently high $R^2$ values across all configurations indicate that, for a given hardware setup, duration and energy consumption are nearly interchangeable as measures of computational effort.

Moreover, when the regression coefficients are known for a specific computing system, energy consumption can be reliably estimated from the duration and the number of GPUs. Only the coefficients of duration ($a$) and of the interaction term duration:GPUs ($c$) are statistically significant. The other coefficients ($b$ and $d$) are omitted from the approximation:

$$\text{Energy} \approx (a + c \cdot \text{GPUs}) \cdot \text{Duration}. \qquad (1)$$

For instance, on the `Capella` system, the following approximation holds for any computation:

11

$$\frac{\text{Energy}}{1\,\text{Wh}} \approx (0.1 + 0.09 \cdot \text{GPUs}) \cdot \frac{\text{Duration}}{1\,\text{s}}. \tag{2}$$

This relationship suggests that, under fixed hardware conditions, monitoring the duration of computations provides an efficient means of estimating energy usage with minimal additional measurement overhead.

|  | Dependent variable: Energy | | |
|---|---|---|---|
|  | Capella | Clara | Paula |
|  | (1) | (2) | (3) |
| Duration ($a$) | 0.097*** | 0.061*** | 0.079*** |
|  | (0.008) | (0.002) | (0.026) |
| GPUs ($b$) | −0.500 | 0.048 | −2.195 |
|  | (2.297) | (0.339) | (3.472) |
| Duration:GPUs ($c$) | 0.090*** | 0.036*** | 0.054*** |
|  | (0.004) | (0.0002) | (0.004) |
| Constant ($d$) | −6.205 | −0.826 | 3.328 |
|  | (5.725) | (1.368) | (17.220) |
| Observations | 44 | 19 | 23 |
| $R^2$ | 0.998 | 1.000 | 0.989 |
| Adjusted $R^2$ | 0.998 | 1.000 | 0.987 |

*Note:*      $^*$p<0.1; $^{**}$p<0.05; $^{***}$p<0.01

**Table 3** Linear regression of energy consumption on duration (table format by Hlavac, 2022). The numbers (coefficients) give the estimated effects of each predictor on the dependent variable. A positive coefficient means the variable increases the outcome, while a negative coefficient means it decreases the outcome. The standard error (in parenthesis) estimates the variability of the coefficient estimate. The p-value (given by the asterisks) indicates whether the predictor is statistically significant (different from zero).

# 5 Discussion

We would like to mention the limitations of our study, which also point to the areas of future research. First, while traditional models were trained on approximately 200 examples, the large language models (LLMs) were applied in a zero-shot setting, meaning they had no access to labeled examples. Previous research has shown that few-shot prompting - where representative examples are included in the prompt - can improve performance (Brown et al., 2020). For the present study, we kept the prompt

as simple as possible (see section A). But in an actual application, we would add background information about the data and the categories. In general, prompt engineering, the addition of representative examples to the prompt, or even fine-tuning an LLM could yield higher accuracy rates. On the other hand, energy efficiency in LLMs can be improved through model quantization, which reduces computational demands by compressing model parameters (Jacob et al., 2017).

Second, we do not account for the energy costs associated with training the traditional models because it is infeasible to compare them to the training costs of LLMs. The LLMs used in this study were pre-trained by external organizations and made publicly available. As a result, the energy costs of training are distributed among all users, making it difficult to estimate per-user energy consumption. Even if training energy costs for an LLM were known, the number of users remains uncertain. Additionally, hosting LLMs (e.g., on Hugging Face) and managing network traffic also contribute to energy consumption. Deploying an LLM on a dedicated server (e.g., using vLLM) requires setup time and additional energy. Beyond inference, significant time and computational resources are also required for development tasks, including data processing, testing different models and prompts, parameter tuning, and debugging - workloads that apply to both traditional models and LLMs. The measurement of additional related energy consumptions (such as network traffic or disk storage) is beyond the scope of this paper.

Third, energy consumption was measured using CodeCarbon, a tool recognized for providing reliable estimates of a machine's total energy use (Bouza et al., 2023). However, it does not allow for precise measurement of energy consumption at the level of individual processes. Moreover, power intake was recorded at 15-second intervals, meaning the accuracy of energy estimates improves with longer-running processes. Another limitation of CodeCarbon is that RAM energy consumption is approximated at 0.375W per GB of memory used. While the Running Average Power Limit (RAPL) framework can directly measure RAM power consumption, it is not supported on all CPUs (https://github.com/mlco2/codecarbon/issues/717#issuecomment-2589805160). Additionally, in shared computing environments such as high-performance computing (HPC) clusters, measurements may be affected by other users' activities. Especially when an LLM was deployed across multiple nodes, variations in network traffic at different times may have influenced energy measurements. A more precise assessment of energy efficiency would benefit from using dedicated computing resources with physical wattmeters and high-resolution energy measurement tools(e.g. Ilsche et al., 2019).

In the following, we assess further limitations of the present study in more detail. More specifically, we address our focus on a single dataset in section 5.1 and the limitation to the text categorisation task in section 5.2. Subsequently, we contextualise our work in the broader context of planet-centered LLMs (section 5.3).

## 5.1 Analysis on other datasets

Our analysis was conducted on a highly specialized dataset. To assess the generalizability of our findings, we replicated the experiments using four additional, widely used datasets (see table 4). These datasets were selected from the HuggingFace platform

based on popularity and had to meet two criteria: suitability for text classification and inclusion of two columns - `text` and `label`. To maintain comparability with our initial analysis, we randomly sampled 200 training examples and 200 test examples from each dataset. Using a slightly larger training set might have provided an advantage to traditional models, as the LLMs were applied in a zero-shot setting without fine-tuning. Each model experiment was repeated 10 times with different samples, ensuring that each model was tested on the same 10 sets.

| Name | Classification Task | ID on https://huggingface.co/datasets |
|------|---------------------|----------------------------------------|
| news | news topics: World, Sports, Business, Sci/Tech | fancyzhx/ag_news |
| yelp | sentiment: 1-5 stars | Yelp/yelp_review_full |
| tomatoes | sentiment: pos, neg | cornell-movie-review-data/rotten_tomatoes |
| emotion | emotion: anger, fear, joy, love, sadness, surprise | dair-ai/emotion |

**Table 4** Selection of datasets for text classification tasks.

Figure 7 visualizes the relationship between accuracy and energy consumption for these additional text classification tasks. For clarity, we restricted the visualization to the models with the three highest accuracy scores and included the linear model with sentence embeddings for comparison (see Tables B4 and B5 for details).

Similar to our findings with the FKTG dataset, the DeepSeek models do not outperform the best non-reasoning models in most cases. The only exception is the emotion dataset, where DeepSeek Llama 70B achieves an accuracy of 0.61, slightly surpassing the best non-reasoning model, Phi 3.5 MoE (0.60). However, unlike in the previous analysis, for every dataset, at least one LLM outperforms the best traditional model. For the news dataset, Llama 3.1 70B achieves an accuracy 0.05 points higher than the best linear model (0.88 vs. 0.83). However, this comes at the cost of significantly higher energy consumption (34.15 Wh vs. 0.0021 Wh), highlighting the need for careful trade-off considerations.

In the case of sentiment analysis on the Yelp dataset, traditional models perform considerably worse than LLMs, justifying the energy costs of LLM deployment. In some cases, a smaller model, such as Qwen 2.5 7B, may be sufficient. While its accuracy is slightly lower than the version with 72B parameters (0.60 vs. 0.68), it consumes only one-eighth of the energy. A similar pattern is observed for sentiment analysis on the Rotten Tomatoes dataset, where traditional models fail to match LLM performance. Among the larger models, Jamba Mini 1.5 stands out as one of the most efficient choices, offering strong accuracy while consuming significantly less energy. Notably, despite having nearly as many parameters as Llama 3.1 70B and Qwen 2.5 72B (51.6B vs. 70B/72B), Jamba Mini 1.5 requires only a quarter of the energy for the same task.

Finally, for emotion classification, the linear model with sentence embeddings is among the top-performing models. In this case, a traditional model provides the most efficient solution. Hence, accuracy-energy trade-offs must be assessed on a case-by-case basis. In some scenarios, traditional models are sufficient, while in others, LLMs offer justifiable benefits despite higher energy consumption. However, a reason for the
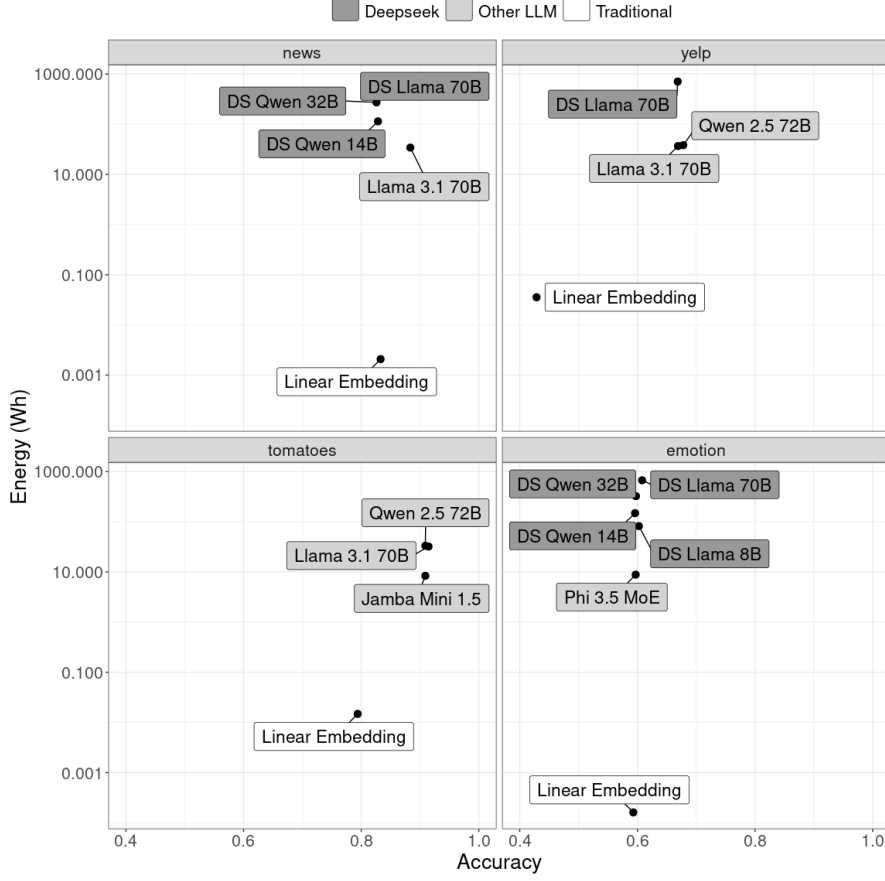
**Fig. 7** Accuracy-energy-trade-off of the best models for the inference task on different datasets (the Linear Embedding model was added for comparison), `Capella` system, single node. See Tables B4 and B5 for results of all models.

superior performance of LLMs on some datasets might be that the data were included in the model's training data. Our study uses data that are probably not part of any LLM training set. Nevertheless, test-time compute, as featured by the Deepseek models, has no benefits in text classification tasks, and the linear relationship between computation runtime and energy consumption holds across all datasets (see Table B6).

## 5.2 Transferability to other tasks

Another limitation of the present study is its exclusive focus on the categorization task, which confines the analysis to a narrow subset of machine learning challenges. While this focus allows for a straightforward measurement of a model's performance (using the accuracy metric), it neglects the applicability of the results to other tasks. Recent studies suggest that similar comparisons in terms of efficiency and accuracy can be insightful in a variety of domains beyond categorization. For instance, Clavié et al.

(2025) demonstrate that simple encoder-based models can effectively tackle generative tasks, expanding the potential applications of smaller, less energy-hungry models.

Moreover, a growing body of research highlights the advantages of fine-tuned small models for specialized tasks, where they often outperform larger models (Savvov, 2024). This trend is evident in studies such as Wei et al. (2024), where an diabetes-specific LLM - despite having significantly fewer parameters - outperforms both GPT-4 and Claude-3.5 in processing various diabetes tasks. Similarly, Lu et al. (2023) report that their fine-tuned models achieve performance levels comparable to GPT-4 on domain-specific annotation tasks, yet with hundreds of times fewer parameters and significantly reduced computational costs. Zhan et al. (2025) further emphasize the superior performance of fine-tuned small models over zero-shot LLMs, particularly in in-domain content moderation tasks.

The study by Luccioni et al. (2024) provides additional insights into the balance between model size and efficiency while looking at ten different machine learning tasks including image classification and captioning, question answering, summarization, as well as image and text generation. The authors demonstrate that smaller models can achieve high performance with considerably less resource consumption. Their initiative resulted into the AI Energy Score (https://huggingface.co/AIEnergyScore), a tool designed to assess the environmental impact of AI models on a range of tasks, and reinforces the growing importance of considering energy efficiency in model evaluation.

## 5.3 Further Requirements of Planet-Centered LLMs

While energy consumption and the associated carbon footprint remain crucial considerations for sustainable AI, truly planet-centered LLMs must meet a broader set of requirements that go beyond mere efficiency. These include other limited resources (water, rare-earth metals, landuse,...), transparency, accessibility, ethical considerations, and technical adaptability to ensure responsible and sustainable AI deployment.

Transparency in AI models is essential for trust and reproducibility (Raji et al., 2020). The predictions of traditional LLM models are generally more transparent than those of LLMs. Open-source LLMs, where both model architectures and training data are publicly available, contribute to scientific progress, allow for direct model comparisons such as this present study, and reduce dependency on proprietary technologies (Wei et al., 2023). Furthermore, the ability to inspect training data is crucial to assess potential biases and copyright compliance (Bender et al., 2021). Many proprietary models, such as GPT-4, lack such transparency, making it difficult to evaluate their fairness and ethical considerations. The EU AI Act will require providers of general-purpose AI models to publish a sufficiently detailed summary of their training data starting in August 2025, which further highlights the call for transparency.

LLMs vary significantly in size, ranging from lightweight models such as fast-Text (Joulin et al., 2017) to massive architectures like BLOOM-176B, which require substantial GPU memory and network bandwidth (Luccioni et al., 2023). These computational demands translate into high operational costs and environmental impacts. Moreover, some models require proprietary hardware, limiting their accessibility and

long-term sustainability. Future AI systems should prioritize modularity and adaptability, enabling efficient integration into diverse infrastructures without excessive resource demands.

The relevance and fairness of AI-generated outputs depend on the quality and recency of training data. Stale or biased datasets can lead to misleading results and reinforce harmful stereotypes (Bender et al., 2021; Gehman et al., 2020). In particular, the presence of toxic content or hate speech in training data can result in models generating harmful or discriminatory outputs, which poses serious challenges for their deployment in sensitive contexts such as education, healthcare, or public administration. Moreover, safety concerns—such as the risk of models producing factually incorrect, manipulative, or otherwise harmful content—are especially critical in public-sector applications, where accountability and trust are paramount (Weidinger et al., 2021). Addressing these challenges requires robust bias-mitigation strategies and transparent documentation of model behavior.

To align with global sustainability and ethical AI principles, future research should emphasize the development of adaptable, transparent, and energy-efficient LLMs. By integrating principles of openness, fairness, and regulatory compliance, we can foster AI systems that not only minimize environmental impact but also promote responsible and equitable usage across sectors.

**Author contribution statements.** T.H. conceived the study, initiated the project, led the research effort, and contributed to the literature review and manuscript writing. J.Z. designed and implemented the experiments, developed the codebase, conducted data analysis, and contributed to drafting the manuscript.

**Competing interests.** There are no competing interests.

**Availability of data and code.** All underlying data will be shared upon reasonable request to the corresponding author. The source code will be made public.

# Appendix A   LLM prompt

For the zero-shot classification, we prompted the LLM with the following instruction (originally in German):

```
Classify the text as one of the following categories:
— <category 1>
— <category 2>
—  ...
```

The categories were a fixed set of 14 options that occurred in the training as well as the test dataset: 'geoWK', 'Tongestein', 'Aktive Störungszonen', 'Öffentlichkeitsbeteiligung', 'Kristallingestein', 'FEP/Szenarien/Entwicklungen des Endlagersystems', 'Anwendung geoWK', 'Mindestanforderungen', 'Steinsalz in steiler Lagerung', 'Datenverfügbarkeit', 'Modellierung', 'Referenzdatensätze', 'Bereitstellung der Daten', 'Ausschlusskriterien'.

Since we deployed the `dspy` framework (https://dspy.ai/) to query the LLMs, the final prompt was automatically extended to the following:

```
— role: system
  content: |—
    Your input fields are:
    1. 'text' (str)

    Your output fields are:
    1. 'category' (str)

    All interactions will be structured in the following way,
    with the appropriate values filled in.

    [[ ## text ## ]]
    {text}

    [[ ## category ## ]]
    {category}

    [[ ## completed ## ]]

    In adhering to this structure, your objective is:
            Classify the text as one of the following categories:
            — <category 1>
            — <category 2>
            —  ...
— role: user
  content: |—
    [[ ## text ## ]]
    <text>

    Respond with the corresponding output fields, starting with
    the field '[[ ## category ## ]]', and then ending with the
    marker for '[[ ## completed ## ]]'.
```

# Appendix B   Tables

| Model | GPUs | Energy (Wh) | Accuracy | Duration (s) | Average Power (W) |
|---|---|---|---|---|---|
| Linear BoW | 1 | <0.01 | 0.43 | 0.01 | 139.96 |
| Linear Tf-idf | 1 | <0.01 | 0.41 | 0.01 | 43.72 |
| Linear Embedding | 1 | 0.12 | 0.57 | 1.64 | 259.41 |
| XGBoost BoW | 1 | <0.01 | 0.35 | 0.01 | 63.32 |
| XGBoost Tf-idf | 1 | <0.01 | 0.47 | 0.01 | 67.77 |
| XGBoost Embedding | 1 | 0.21 | 0.47 | 2.87 | 259.94 |
| Llama 3.1 8B | 1 | 5.86 | 0.35 | 36.88 | 572.49 |
| Llama 3.1 70B | 2 | 48.60 | 0.48 | 161.59 | 1082.82 |
| Qwen 2.5 7B | 1 | 5.58 | 0.45 | 36.28 | 553.84 |
| Qwen 2.5 72B | 2 | 48.66 | 0.51 | 164.44 | 1065.31 |
| Phi 3.5 Mini | 1 | 5.74 | 0.30 | 41.45 | 498.46 |
| Phi 3.5 MoE | 2 | 11.00 | 0.40 | 55.51 | 713.34 |
| Jamba Mini 1.5 | 2 | 17.42 | 0.34 | 78.61 | 797.94 |
| DS Llama 8B | 1 | 79.64 | 0.37 | 517.83 | 553.67 |
| DS Llama 70B | 2 | 702.06 | 0.46 | 2543.47 | 993.68 |
| DS Qwen 14B | 1 | 155.20 | 0.39 | 981.35 | 569.33 |
| DS Qwen 32B | 1 | 373.56 | 0.45 | 2255.99 | 596.11 |

**Table B1**  Measurements of all models for the inference task on the FKTG dataset, `Capella` system, single node, shown are averages over 10 runs

| Model | Duration (s) | | | Energy consumed (Wh) | | |
|---|---|---|---|---|---|---|
| | single | double | ratio | single | double | ratio |
| Llama 3.1 70B | 161.59 | 304.77 | 1.89 | 48.60 | 94.88 | 1.95 |
| Qwen 2.5 72B | 164.44 | 308.16 | 1.87 | 48.66 | 95.70 | 1.97 |
| Jamba Mini 1.5 | 78.61 | 113.88 | 1.45 | 17.42 | 29.81 | 1.71 |
| DS Llama 70B | 2543.47 | 6792.54 | 2.67 | 702.06 | 1899.86 | 2.71 |

**Table B2**  Comparison single vs. double node deployment, `Capella` system

| Model | Duration (s) | | | Energy consumed (Wh) | | |
|---|---|---|---|---|---|---|
| | A30 | V100 | H100 | A30 | V100 | H100 |
| Llama 3.1 8B | 20.78 | 27.52 | 36.88 | 2.91 | 2.88 | 5.86 |
| Qwen 2.5 7B | 19.58 | 24.64 | 36.28 | 2.87 | 2.63 | 5.58 |
| Phi 3.5 Mini | 19.18 | 25.02 | 41.45 | 2.65 | 2.50 | 5.74 |
| Phi 3.5 MoE | 77.60 | 32.53 | 45.93 | 17.77 | 6.04 | 15.04 |
| DS Llama 8B | 1210.90 | 1439.58 | 517.83 | 175.83 | 137.90 | 79.64 |
| DS Qwen 14B | 1348.09 | 1736.21 | 624.38 | 254.01 | 230.72 | 157.58 |
| DS Qwen 32B | 1688.23 | 2192.53 | 806.68 | 444.67 | 457.60 | 378.58 |

**Table B3**  Comparison of different GPU cards, single node deployment.

19

| Dataset | news | | yelp | |
|---|---|---|---|---|
| Model | Energy (Wh) | Accuracy | Energy (Wh) | Accuracy |
| Linear BoW | <0.01 | 0.65 | <0.01 | 0.36 |
| Linear Tf-idf | <0.01 | 0.65 | <0.01 | 0.34 |
| Linear Embedding | <0.01 | 0.83 | 0.04 | 0.43 |
| XGBoost BoW | <0.01 | 0.48 | <0.01 | 0.31 |
| XGBoost Tf-idf | <0.01 | 0.52 | <0.01 | 0.29 |
| XGBoost Embedding | 0.03 | 0.74 | 0.01 | 0.40 |
| Llama 3.1 8B | 4.31 | 0.71 | 4.73 | 0.58 |
| Llama 3.1 70B | 34.15 | 0.88 | 36.71 | 0.67 |
| Qwen 2.5 7B | 4.21 | 0.01 | 4.52 | 0.60 |
| Qwen 2.5 72B | 33.75 | 0.79 | 38.20 | 0.68 |
| Phi 3.5 Mini | 3.30 | 0.53 | 15.55 | 0.58 |
| Phi 3.5 MoE | 8.53 | 0.78 | 8.32 | 0.58 |
| Jamba Mini 1.5 | 9.34 | 0.78 | 11.45 | 0.56 |
| DS Llama 8B | 60.58 | 0.82 | 97.18 | 0.62 |
| DS Llama 70B | 483.73 | 0.83 | 707.03 | 0.67 |
| DS Qwen 14B | 113.81 | 0.83 | 177.41 | 0.63 |
| DS Qwen 32B | 271.92 | 0.83 | 358.62 | 0.63 |

**Table B4** Measurements of all models for the inference task on the news and yelp datasets, `Capella` system, single node, shown are averages over 10 runs

| Dataset | tomatoes | | emotion | |
|---|---|---|---|---|
| Model | Energy (Wh) | Accuracy | Energy (Wh) | Accuracy |
| Linear BoW | <0.01 | 0.59 | <0.01 | 0.36 |
| Linear Tf-idf | <0.01 | 0.59 | <0.01 | 0.40 |
| Linear Embedding | 0.01 | 0.79 | <0.01 | 0.59 |
| XGBoost BoW | <0.01 | 0.54 | <0.01 | 0.30 |
| XGBoost Tf-idf | <0.01 | 0.55 | <0.01 | 0.33 |
| XGBoost Embedding | <0.01 | 0.76 | <0.01 | 0.53 |
| Llama 3.1 8B | 4.12 | 0.87 | 4.46 | 0.56 |
| Llama 3.1 70B | 32.07 | 0.91 | 34.12 | 0.58 |
| Qwen 2.5 7B | 4.04 | 0.73 | 4.17 | 0.37 |
| Qwen 2.5 72B | 33.25 | 0.91 | 34.81 | 0.58 |
| Phi 3.5 Mini | 7.20 | 0.87 | 5.13 | 0.53 |
| Phi 3.5 MoE | 7.72 | 0.89 | 8.82 | 0.60 |
| Jamba Mini 1.5 | 8.37 | 0.91 | 10.22 | 0.56 |
| DS Llama 8B | 72.15 | 0.83 | 81.82 | 0.60 |
| DS Llama 70B | 510.86 | 0.90 | 670.40 | 0.61 |
| DS Qwen 14B | 134.02 | 0.89 | 148.20 | 0.60 |
| DS Qwen 32B | 246.48 | 0.89 | 323.48 | 0.60 |

**Table B5** Measurements of all models for the inference task on the tomatoes and emotion datasets, `Capella` system, single node, shown are averages over 10 runs

| | Dependent variable: Energy | | | |
|---|---|---|---|---|
| | tomatoes | emotion | news | yelp |
| | (1) | (2) | (3) | (4) |
| Duration | 0.040*** | 0.043*** | 0.052*** | 0.045*** |
| | (0.002) | (0.002) | (0.003) | (0.003) |
| GPUs | −0.079 | −0.052 | 0.536 | 0.810 |
| | (0.950) | (1.011) | (1.470) | (1.545) |
| Duration:GPUs | 0.122*** | 0.120*** | 0.115*** | 0.120*** |
| | (0.001) | (0.001) | (0.002) | (0.002) |
| Constant | −0.397 | −0.464 | −1.300 | −1.773 |
| | (1.290) | (1.372) | (1.985) | (2.103) |
| Observations | 17 | 17 | 17 | 17 |
| $R^2$ | 1.000 | 1.000 | 1.000 | 1.000 |
| Adjusted $R^2$ | 1.000 | 1.000 | 1.000 | 1.000 |

*Note:* $^{*}$p<0.1; $^{**}$p<0.05; $^{***}$p<0.01

**Table B6** Linear regression of energy consumption on duration for the datasets of section 5.1 (table format by Hlavac, 2022).

# References

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Advances in Neural Information Processing Systems, vol. 30 (2017)

Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., *et al.*: Language models are few-shot learners. In: Advances in Neural Information Processing Systems, vol. 33, pp. 1877–1901 (2020)

Lacoste, A., Luccioni, A., Schmidt, V., Dandres, T.: Quantifying the Carbon Emissions of Machine Learning (2019). https://arxiv.org/abs/1910.09700

Axenbeck, J., Kunkel, S., Blain, J., et al.: Global Embodied Emissions of Digital Technologies: The Hidden 42%. Research Square. Preprint (Version 1) available at Research Square (2025). https://doi.org/10.21203/rs.3.rs-6479454/v1 . https://www.researchsquare.com/article/rs-6479454/v1

United Nations: Paris Agreement (2015). https://unfccc.int/sites/default/files/english_paris_agreement.pdf

United Nations: Transforming our world: the 2030 Agenda for Sustainable Development (2015). https://sustainabledevelopment.un.org/post2015/

transformingourworld

Strubell, E., Ganesh, A., McCallum, A.: Energy and Policy Considerations for Deep Learning in NLP (2019). https://arxiv.org/abs/1906.02243

Patterson, D., Gonzalez, J., Le, Q., Liang, C., Munguia, L.-M., Rothchild, D., So, D., Texier, M., Dean, J.: Carbon Emissions and Large Neural Network Training (2021). https://arxiv.org/abs/2104.10350

Luccioni, A.S., Hernandez-Garcia, A.: Counting Carbon: A Survey of Factors Influencing the Emissions of Machine Learning (2023). https://arxiv.org/abs/2302.08476

OpenAI: Learning to reason with LLMs (2024). https://openai.com/index/learning-to-reason-with-llms/

DeepSeek-AI: DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning (2025). https://arxiv.org/abs/2501.12948

Joulin, A., Grave, E., Bojanowski, P., Mikolov, T.: Bag of tricks for efficient text classification. In: Lapata, M., Blunsom, P., Koller, A. (eds.) Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers, pp. 427–431. Association for Computational Linguistics, Valencia, Spain (2017). https://aclanthology.org/E17-2068/

Wang, Z., Pang, Y., Lin, Y., Zhu, X.: Adaptable and Reliable Text Classification using Large Language Models (2024). https://arxiv.org/abs/2405.10523

LangChain Team: Classification Tutorial – LangChain Documentation. https://python.langchain.com/docs/tutorials/classification/. Accessed: 2025-02-23 (2023)

Lamini Team: CAT Documentation – Lamini. https://docs.lamini.ai/cat/. Accessed: 2025-02-23 (2023)

Kaack, L.H., Donti, P.L., Strubell, E., Kamiya, G., Creutzig, F., Rolnick, D.: Aligning artificial intelligence with climate change mitigation **12**(6), 518–527 https://doi.org/10.1038/s41558-022-01377-7

Luccioni, A.S., Strubell, E., Crawford, K.: From Efficiency Gains to Rebound Effects: The Problem of Jevons' Paradox in AI's Polarized Environmental Debate (2025). https://arxiv.org/abs/2501.16548

Patterson, D., Gonzalez, J., Hölzle, U., Le, Q., Liang, C., Munguia, L.-M., Rothchild, D., So, D., Texier, M., Dean, J.: The Carbon Footprint of Machine Learning Training Will Plateau, Then Shrink (2022). https://arxiv.org/abs/2204.05149

Samsi, S., Zhao, D., McDonald, J., Li, B., Michaleas, A., Jones, M., Bergeron, W., Kepner, J., Tiwari, D., Gadepally, V.: From Words to Watts: Benchmarking the

Energy Costs of Large Language Model Inference (2023). https://arxiv.org/abs/2310.03003

Liu, X., Sun, T., He, J., Wu, J., Wu, L., Zhang, X., Jiang, H., Cao, Z., Huang, X., Qiu, X.: Towards Efficient NLP: A Standard Evaluation and A Strong Baseline (2022). https://arxiv.org/abs/2110.07038

Chien, A.A., Lin, L., Nguyen, H., Rao, V., Sharma, T., Wijayawardana, R.: Reducing the carbon impact of generative ai inference (today and in 2035). In: Proceedings of the 2nd Workshop on Sustainable Computer Systems. HotCarbon '23. Association for Computing Machinery, New York, NY, USA (2023). https://doi.org/10.1145/3604930.3605705 . https://doi.org/10.1145/3604930.3605705

Wang, A., Wolf, T.: Overview of the sustainlp 2020 shared task. In: SUSTAINLP (2020). https://api.semanticscholar.org/CorpusID:226283937

Alizadeh, N., Belchev, B., Saurabh, N., Kelbert, P., Castor, F.: Language Models in Software Development Tasks: An Experimental Analysis of Energy and Accuracy (2025). https://arxiv.org/abs/2412.00329

Luccioni, S., Jernite, Y., Strubell, E.: Power hungry processing: Watts driving the cost of ai deployment? In: Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency. FAccT '24, pp. 85–99. Association for Computing Machinery, New York, NY, USA (2024). https://doi.org/10.1145/3630106.3658542

Bai, G., Chai, Z., Ling, C., Wang, S., Lu, J., Zhang, N., Shi, T., Yu, Z., Zhu, M., Zhang, Y., Song, X., Yang, C., Cheng, Y., Zhao, L.: Beyond Efficiency: A Systematic Survey of Resource-Efficient Large Language Models (2024). https://arxiv.org/abs/2401.00625

Anthony, L.F.W., Kanding, B., Selvan, R.: Carbontracker: Tracking and Predicting the Carbon Footprint of Training Deep Learning Models (2020). https://arxiv.org/abs/2007.03051

Henderson, P., Hu, J., Romoff, J., Brunskill, E., Jurafsky, D., Pineau, J.: Towards the systematic reporting of the energy and carbon footprints of machine learning. J. Mach. Learn. Res. **21**(1) (2020)

Bouza, L., Bugeau, A., Lannelongue, L.: How to estimate carbon footprint when training deep learning models? a guide and review. Environmental Research Communications **5**(11), 115014 (2023) https://doi.org/10.1088/2515-7620/acf81b

Hlavac, M.: stargazer: Well-Formatted Regression and Summary Statistics Tables. R package version 5.2.3 (2022). https://CRAN.R-project.org/package=stargazer

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger,

23

G., Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., Amodei, D.: Language models are few-shot learners. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M.F., Lin, H. (eds.) Advances in Neural Information Processing Systems, vol. 33, pp. 1877–1901. Curran Associates, Inc., ??? (2020)

Jacob, B., Kligys, S., Chen, B., Zhu, M., Tang, M., Howard, A., Adam, H., Kalenichenko, D.: Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference (2017). https://arxiv.org/abs/1712.05877

Ilsche, T., Hackenberg, D., Schöne, R., Bielert, M., Höpfner, F., E. Nagel, W.: Metricq: A scalable infrastructure for processing high-resolution time series data. In: 2019 IEEE/ACM Industry/University Joint International Workshop on Data-center Automation, Analytics, and Control (DAAC), pp. 7–12 (2019). https://doi.org/10.1109/DAAC49578.2019.00007

Clavié, B., Cooper, N., Warner, B.: It's All in The [MASK]: Simple Instruction-Tuning Enables BERT-like Masked Language Models As Generative Classifiers (2025). https://arxiv.org/abs/2502.03793

Savvov, S.: Your Company Needs Small Language Models (2024). https://towardsdatascience.com/your-company-needs-small-language-models-d0a223e0b6d9/

Wei, L., Ying, Z., He, M., Chen, Y., Yang, Q., Hong, Y., Lu, J., Li, X., Huang, W., Chen, Y.: An adapted large language model facilitates multiple medical tasks in diabetes care (2024). https://arxiv.org/abs/2409.13191

Lu, Y., Yao, B., Zhang, S., Wang, Y., Zhang, P., Lu, T., Li, T.J.-J., Wang, D.: Human Still Wins over LLM: An Empirical Study of Active Learning on Domain-Specific Annotation Tasks (2023). https://arxiv.org/abs/2311.09825

Zhan, X., Goyal, A., Chen, Y., Chandrasekharan, E., Saha, K.: SLM-Mod: Small Language Models Surpass LLMs at Content Moderation (2025). https://arxiv.org/abs/2410.13155

Raji, I.D., Gebru, T., Mitchell, M., Buolamwini, J., Lee, J., Denton, E.: Saving face: Investigating the ethical concerns of facial recognition auditing. Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (2020)

Wei, J., Bosma, M., Zhao, V.Y., et al.: Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023)

Bender, E.M., Gebru, T., McMillan-Major, A., Shmitchell, S.: On the dangers of stochastic parrots: Can language models be too big? Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (2021)

Luccioni, A.S., Viguier, S., Ligozat, A.-L.: Estimating the carbon footprint of bloom, a 176b parameter language model. J. Mach. Learn. Res. **24**(1) (2023)

Gehman, S., Gururangan, S., Sap, M., Choi, Y., Smith, N.A.: Realtoxicityprompts: Evaluating neural toxic degeneration in language models. Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP) (2020)

Weidinger, L., Mellor, J., et al.: Ethical and social risks of harm from language models. arXiv preprint arXiv:2112.04359 (2021)