

# HOLISTICALLY EVALUATING THE ENVIRONMENTAL IMPACT OF CREATING LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

As the performance of artificial intelligence systems has dramatically increased, so too has the environmental impact of creating these systems. While many model developers release estimates of the power consumption and carbon emissions from the final training runs for their latest models, there is comparatively little transparency into the impact of model development, hardware manufacturing, and total water usage throughout. In this work, we estimate the real-world environmental impact of developing a series of language models, ranging from 20 million to 7 billion active parameters, trained on up to 5 trillion tokens each. When accounting for hardware manufacturing, model development, and our final training runs, we find that our series of models released **270 metric tons** of carbon emissions, equivalent to powering about 53 homes in the United States for one year, and consumed **1.137 million liters of water**, equivalent to about 10 years of water usage by a person in the United States, even though our data center is extremely water-efficient. We measure and report the environmental impact of our model development; to the best of our knowledge we are the first to do so for LLMs, and we find that model development, the impact of which is generally not disclosed by most model developers, amounted to  $\sim 80\%$  of that of training. By looking at detailed time series data for power consumption, we also find that power usage throughout training is not consistent, fluctuating between  $\sim 15\%$  and  $\sim 85\%$  of our hardware’s maximum power draw, with negative implications for grid-scale planning as demand continues to grow. We close with a discussion on the continued difficulty of estimating the environmental impact of AI systems, and key takeaways for model developers and the public at large.

## 1 INTRODUCTION

In recent years, the field of artificial intelligence has progressed at an unprecedented pace, driven in large part by the development and deployment of large language and multimodal models. However, the development of these models comes with significant environmental costs (Schwartz et al., 2020; Strubell et al., 2020; Wu et al., 2022). Training these models requires massive computational resources, which, in turn, require large amounts of energy. Powering training both emits carbon (by burning fossil fuels) and consumes water (by evaporating or polluting it in power plants, data centers, and hardware manufacturing processes; Li et al. (2023)). There is a growing demand for energy to power AI workloads – for instance, Microsoft recently signed a deal to purchase the next 20 years of energy generated by re-opening a nuclear power plant<sup>1</sup>, and meanwhile energy providers are extending the life of aging fossil fuel energy plants to keep up with demand<sup>2</sup>. As such, especially as increasing numbers of stakeholders become involved in the development and use of AI systems, it is imperative to carefully characterize the true cost of building and deploying state-of-the-art models, to inform more effective strategies for mitigating potential harms, and planning for future demand.

In this paper, we estimate the energy use and environmental impacts caused by training a series of dense transformer language models<sup>3</sup> ranging in size from 20 million to 7 billion active parameters, trained on 1.7 to 5 trillion tokens. To do this, we calculate Scope 2 CO<sub>2</sub> emissions in accordance with

<sup>1</sup><https://www.technologyreview.com/2024/09/26/1104516/three-mile-island-microsoft/>

<sup>2</sup><https://www.wsj.com/business/energy-oil/electricity-demand-coal-gas-retirement-charts-dd07029a>

<sup>3</sup>Details are currently omitted to preserve anonymity but will be added upon publication.

054 the Greenhouse Gas Protocol’s definitions,<sup>4</sup> and Scope 1 and Scope 2 water consumption following  
055 Li et al. (2023); in addition, we calculate “upstream” embodied carbon and water consumption, and  
056 provide “downstream” estimates from use of our models (which are part, but not all, of Scope 3).

057 Importantly, we calculate (i) electricity consumption, (ii) carbon emissions, and (iii) water consump-  
058 tion at three points in the machine learning pipeline: early model development (e.g., hyperparameter  
059 tuning and experiments before the final training run), training of the main model, and inference. To  
060 the best of our knowledge, we are the first to report this information for model development of large  
061 language models, and we find the environmental impact of developing even our relatively small  
062 models (only up to 7B parameters) is equivalent to burning 1.5 gasoline tanker trucks of fuel, or the  
063 amount of water consumed by one average person in the United States in 4.5 years. We encour-  
064 age the reader to consider larger models released by other organizations to have equivalently larger  
065 environmental impacts.

066 Our methodology draws upon best practices from recent publications, aiming to provide the most  
067 thorough reporting yet of the environmental impact of LLMs. For example, unlike previous works  
068 that assume GPUs operate at 100% of their theoretical maximum power draw (Dubey et al., 2024)  
069 and report only the cost to train a small set of released models, we measure power consumption  
070 at sub-second intervals throughout training. We focus our efforts on a wide range of model sizes,  
071 optimized for widespread deployment (Dubey et al., 2024; Mehta et al., 2024; Team et al., 2024),  
072 and estimate what the environmental impact would be if our models were deployed in a variety of  
073 different scenarios. We find that in some scenarios, our models would only need to run inference  
074 on 200,000 instances to match the electricity consumed, carbon emitted, and water consumed of the  
075 *entire* training process.

076 We conclude that more transparency is needed across the industry in reporting the environmental  
077 impact of AI systems. AI systems orders of magnitude larger than those in this paper are being  
078 built, and put into production at a global scale, leading to emissions 10s or 100s of times larger than  
079 what we report. This work is a step in the right direction, but responsibility of calculating, reporting,  
080 and reducing the environmental impact should fall on those training the largest models, as they are  
081 having the largest impact.

## 082 2 RELATED WORK

083 While most publicly available models do not report any climate impact, including CO<sub>2</sub> emissions,  
084 water usage, or embodied carbon, a few reports recently have included some estimates. For example,  
085 Luccioni et al. (2023) reported estimates for emissions from the manufacturing process (embodied  
086 emissions), from electricity consumption during training, and from electricity consumption of the  
087 cluster while it was idle (see their Table 2). Dodge et al. (2022) measured electricity consump-  
088 tion and carbon emissions for training language models and computer vision models with granular  
089 timesteps with region-specific carbon intensity, but didn’t measure development costs, water con-  
090 sumption, or inference. Similarly, developers of the Llama models (Touvron et al., 2023a;b; Dubey  
091 et al., 2024) reported electricity consumption and carbon emissions estimates of training their fi-  
092 nal models; they did not estimate development cost or water consumption, and their approach to  
093 carbon intensity varied<sup>5</sup>. Gemma developers (Team et al., 2024) only report a single number: the  
094 total emissions from pretraining their models, not broken down by model or by different stages of  
095 training, or by electricity consumption and carbon intensity. The OLMo report (Groeneveld et al.,  
096 2024) documents electricity consumption per model, and uses region-specific carbon intensity to  
097 estimate emissions for two regions, but does not estimate other environmental impacts. Energy use  
098 and environmental impacts are not typically documented for proprietary models.

099 Comparably little transparency has been provided on the water consumption of AI systems. Li et al.  
100 (2023) estimate the water consumption of some closed models like GPT-3, but these estimates are  
101 based on speculation about location of training, energy consumption, etc., as there is very little  
102 public information about GPT-3’s training. Similarly, there are few estimates of embodied carbon  
103 for AI systems, as the manufacturing process is notoriously opaque. In addition, almost all reporting  
104

105  
106 <sup>4</sup><https://ghgprotocol.org/sites/default/files/standards/ghg-protocol-revised.pdf>

107 <sup>5</sup>Llama 1 did not use the data center location’s carbon intensity, instead using US national average carbon  
intensity; Llama 2 did not specify the carbon intensity; Llama 3 used a region-specific carbon intensity

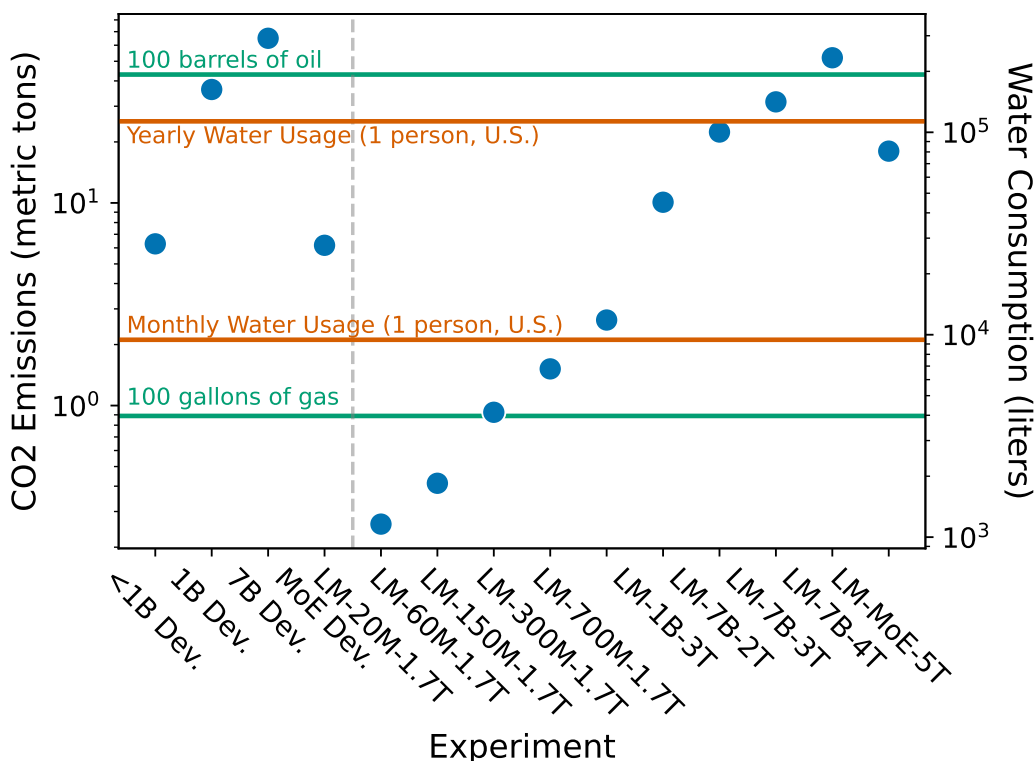


Figure 1: The environmental impact for model development and final training runs described in Section 4.1, where we plot each experiment by both its total CO<sub>2</sub> emissions and water consumption. We see that development costs are substantial, and comparable to that of the most expensive full training runs. We also see that environmental impact is log-linear for both the size of the model (keeping the size of the dataset consistent), and the size of the training dataset (keeping the model size consistent), highlighting the multi-dimensional factors that dictate total environmental impact.

of environmental impact is based on *training* of the *final* model that is released. Instead of only focusing on training, Luccioni et al. (2024) estimate the impact of inference of deployed AI systems. To the best of our knowledge our work provides the first public estimates of environmental impact of development of an LLM, i.e. hyperparameter tuning and ablations before the main training run.

### 3 METHODOLOGY

Our goal in this work is to characterize the holistic environmental impacts of large language models in as much detail as possible, enabling assessment of key challenges and future directions towards reducing those impacts. Typically studies documenting language model training and development methodology will address this concern by reporting the cost to train the final, deployed model measured in GPU hours, kWh energy, and/or CO<sub>2</sub>e emissions. However, this calculation provides an incomplete characterization of the factors leading to environmental degradation due to LLMs that under-estimates impacts and provides insufficient information to inform strategies for developing and deploying LLMs in a more environmentally conscious way.

Following the more comprehensive analysis provided for the BLOOM model (Luccioni et al., 2023), we expand our measurement to include both *operational* GHG emissions arising from the energy required for the development, training, and inference phases of the ML model lifecycle, as well as *embodied* emissions attributed to manufacturing of the hardware supporting those operations. We also go beyond previous work to report non-GHG externalities such as water use, and finer-grained data such as variance in energy use throughout training. We describe our methodology for measuring and estimating these impacts in more detail below.

### 3.1 OPERATIONAL IMPACTS

Operational environmental impacts of LLMs are those that arise directly from the development and use of models, and include the GHG emissions arising from energy sources used to power model training and deployment, including servers and data center cooling. We base our analysis of operational emissions around the following equation introduced by Schwartz et al. (2020) to describe the amount of computation required to produce a machine learning artifact, such as an LLM:

$$Cost(R) \propto E \cdot D \cdot H \quad (1)$$

where the cost of a scientific result  $R$  (e.g. a claim that a particular training setup reaches X accuracy on benchmark Y) is proportional to the product of the cost of processing a single example  $E$ , the size of the training dataset  $D$ , and the number of hyperparameter experiments  $H$ . In previous work,  $E \cdot D$ , the cost of training on the training dataset, is what is most commonly reported, and  $H$ , the total number of experiments, is most often excluded.

In our analysis, we calculate the total power consumption during model training, development, and inference, and use this to estimate the total carbon emissions and water consumption during each stage. We follow previous work (Luccioni et al., 2023; Dubey et al., 2024; Team et al., 2024) to calculate CO<sub>2</sub> emissions from power consumption:

$$CO_2Emissions = P \cdot PUE \cdot CI \quad (2)$$

where the total carbon emissions is equal to the power usage  $P$ , multiplied by the power usage effectiveness  $PUE$ <sup>6</sup> of the data center, multiplied by the carbon intensity  $CI$  of the local power grid. We ran every experiment in the same data center, and our data center provider informed us that their  $PUE$  is between 1.1 and 1.2 depending on the current total utilization, so we conservatively assume a consistent value of 1.2 for our calculations. The power provider, which will be deanonymized upon publication, last reported a carbon intensity of 0.332 kg CO<sub>2</sub>e per kWh in 2021.

We follow Li et al. (2023) to calculate water consumed onsite and through power generation:

$$Consumption = P \cdot PUE \cdot (WUE_{\text{onsite}} + WUE_{\text{offsite}}) \quad (3)$$

where  $WUE_{\text{onsite}}$  is the water usage effectiveness of the data center, dictated by the cooling hardware used, and  $WUE_{\text{offsite}}$  is the water usage effectiveness of the local power provider, dictated by the precise mixture of sources of power generation, as thermo- and hydro-electric power plants lead to evaporated water that is lost and will not re-enter circulation in the local environment.

As our data center uses an efficient closed-loop cooling system with no evaporative cooling, we assume a  $WUE_{\text{onsite}}$  of 0.2 liters per kWh following Li et al. (2023). Our data center is in Texas, so we use the reported average for Texas' power generation for our  $WUE_{\text{offsite}}$ , or 1.29 L per kWh (Reig et al., 2020). Together, these lead to a total  $WUE$  of 1.49 L per kWh.

Both calculations rely on total power usage. To calculate power usage during development and training, we analyze detailed time series data for a single node throughout each run, logging power data at sub-second intervals, and extrapolate to the total number of nodes. As we only measure GPU power consumption, our estimates should be viewed as a lower bound on the true amount of power consumed during development and training.

### 3.2 EMBODIED IMPACTS

Embodied impacts are those arising from the production of physical elements required to support LLM development and use, such as hardware manufacturing and data center construction. To calculate embodied emissions, we follow Luccioni et al. (2023) by amortizing the carbon emissions from manufacturing over the lifetime of the hardware to get an estimate of the per hour cost, and multiplying by the number of GPU hours used throughout model development and training. We

<sup>6</sup><https://www.techtarget.com/searchdatacenter/definition/power-usage-effectiveness-PUE>

216 extend this to include water consumption as well, by amortizing estimates of water consumption  
217 during manufacturing over the lifetime of the hardware.  
218

### 219 3.3 MODELS, DATA, AND HARDWARE 220

221 Most of the models we evaluate are standard dense transformers, with an architecture similar to  
222 Llama (Touvron et al., 2023a;b; Dubey et al., 2024), OLMo (Groeneveld et al., 2024), and other  
223 recent popular models, ranging in size from 20 million to 7 billion active parameters. Each of the  
224 sub-billion parameter models was trained on 1.7 trillion tokens, the 1 billion parameter model was  
225 trained to 3 trillion tokens, and the 7 billion parameter models were trained to 2, 3 and 4 trillion  
226 tokens. We additionally evaluate a mixture-of-experts (MoE) model with 1 billion active and 7  
227 billion total parameters, trained to 5 trillion tokens.

228 Each model was trained on the same compute cluster, using standard HGX servers with 8 NVIDIA  
229 H100 GPUs per server, with high speed InfiniBand interconnect between each node, and we used  
230 between 2 and 64 nodes concurrently per training run.  
231

### 232 3.4 SIMULATING INFERENCE 233

234 Because we do not deploy our models, we do not collect or report data about real usage of our mod-  
235 els. We instead report estimated costs associated with deployment of a subset of our models, along  
236 with comparison models, with varying inference configurations. Though in reality causal language  
237 models can have a variety of use cases and be deployed on a variety of hardware infrastructure,  
238 we collect measurements assuming models are served via SGLang (Zheng et al., 2024) on a single  
239 H100 GPU that users interact with the models via chat. All three inference configurations used can  
240 be mapped to a previously proposed realistic online inference scenario (Reddi et al., 2020; Peng  
241 et al., 2023). Specifically, other than the “batching” scenario where all requests are sent instant-  
242 aneously, the requests follow a Poisson distribution, albeit at different rates that influence different  
243 batch sizes. The requests themselves come from the ShareGPT dataset<sup>7</sup>, and each inference scenario  
244 involves the same sample of 2400 prompts (same random seed). Input and output lengths, therefore,  
245 are the same in theory for a given model, but due to differences in tokenization and model context  
246 length, there are slight variations in mean input/output lengths across models, 225-250 and 190-230  
247 tokens respectively.

248 In our inference experiments, we measure cumulative energy consumption using CodeCarbon  
249 Courty et al. (2024) tracking, which was checked against the same time series monitoring used  
250 throughout training. Notably, we measure total power and energy consumption associated with only  
251 the relevant processes, excluding the overhead associated with, for example, holding the model in  
252 memory or listening for requests.

253 The hardware used for our inference simulations is from the same cluster as that used in training,  
254 but we use only a single H100 GPU at a time. See Appendix A for details about our inference  
255 methodology and assumptions.

## 256 4 RESULTS 257

### 258 4.1 BUILDING OUR MODELS 259

260 In this section, we aim to report a full accounting of the environmental impact of training our series  
261 of models, from hardware manufacturing, to development, and the final training runs. We follow the  
262 methodology outlined in Section 3.1 and Section 3.2.

263 When calculating environmental impact, we use information from our data center provider and their  
264 power provider to assume a carbon intensity of 0.332 kg CO<sub>2</sub> emitted per kWh, a power usage  
265 effectiveness (*PUE*) of 1.2, and a total water usage effectiveness (*WUE*) of 1.49 liters per kWh.  
266

267 **Hardware manufacturing** NVIDIA does not release the embodied carbon emissions or water  
268 consumption about the hardware it produces, so we assume the same embodied carbon emissions  
269

<sup>7</sup>[https://huggingface.co/datasets/anon8231489123/ShareGPT\\_Vicuna\\_unfiltered/resolve/main/ShareGPT\\_V3\\_unfiltered\\_cleaned\\_split.json](https://huggingface.co/datasets/anon8231489123/ShareGPT_Vicuna_unfiltered/resolve/main/ShareGPT_V3_unfiltered_cleaned_split.json), anon8231489123/ShareGPT\_Vicuna\_unfiltered

Table 1: We developed our models in four distinct groups, based on parameter count and architecture: less than 1 billion, 1 billion, and 7 billion parameters, and our mixture-of-experts model with 1 billion active and 7 billion total parameters. We found that  $\sim 55\%$  of our total environmental impact came from developing the 7B models, and the total impact was emissions equivalent to 1.5 tanker trucks’ worth of gasoline, and water consumption equal to 4 and a half years of water used by the average person in the United States.

|               | GPU Hours | Total MWh | # Runs | Carbon Emissions (tCO <sub>2</sub> eq) | Equivalent to ...             | Water Consump. (kL) | Equiv. to... (water usage, 1 person) |
|---------------|-----------|-----------|--------|--|-------------------------------|---------------------|--------------------------------------|
| <b>&lt;1B</b> | 29k       | 19        | 20     | 6                                      | 675 gallons of gasoline       | 28                  | 3 mo                                 |
| <b>1B</b>     | 164k      | 109       | 227    | 36                                     | 40x NY ↔ SF flights, 1 person | 163                 | 1 yr, 5 mo                           |
| <b>7B</b>     | 269k      | 196       | 375    | 65                                     | 150 oil barrels               | 291                 | 2 yr, 7 mo                           |
| <b>MoE</b>    | 27k       | 19        | 35     | 6                                      | 3 tons of coal                | 28                  | 3 mo                                 |
| <b>Total</b>  | 490k      | 342       | 657    | 114                                    | 1.5 gasoline tanker trucks    | 510                 | 4 yr, 6 mo                           |

as Luccioni et al. (2023), or 3700 kg of CO<sub>2</sub>eq per 8x server node, equal 463 kg per GPU. There is little public information on how much water is required to produce a single GPU, though chip manufacturing facilities require millions of liters per day<sup>8</sup>. Some estimates<sup>9</sup> place TSMC water usage at 12.33 liters per square centimeter of hardware, which equals 100.4 liters per H100, which we use for our analysis.

We additionally estimate the environmental impact from mining rare earth metals used during manufacturing, assuming an H100 is 0.1% rare earth metal by mass. Mining 1 kg rare earth materials consume about 11 kL of water, and releases 65.4 kg CO<sub>2</sub>eq (Browning et al., 2016), and one silicon wafer weighs 125 grams and produces about 63 H100s<sup>10</sup>. Together, these add an additional 2.2 liters consumed and 0.013 kg CO<sub>2</sub>eq per GPU.

Internally, we assume a 4 year lifespan for our GPUs, leading to an embodied emissions of 0.013 kg of CO<sub>2</sub>eq and 0.003 liters of water consumed per GPU hour. We used 1.17 million GPU hours in total, leading to a total of **16 tCO<sub>2</sub>eq** emitted and **3.4 kL** of water consumed during manufacturing.

**Development** Before launching our final training runs for each model, we ran a series of controlled experiments to improve and stabilize our training setup, and to determine our final hyperparameters and data mixtures. We ran these in four distinct groups: small models (less than 1 billion parameters), 1 billion parameter models, 7 billion parameter models, and our mixture-of-experts model. We report detailed development costs for each group in Table 1.

Unsurprisingly, we find that the majority of development costs ( $\sim 55\%$ ) were incurred at the 7 billion parameter scale, due to both the relative size of the model and our own prioritization, and we see this both in the total environmental impact and the number of individual runs per category. Using our data center’s efficiency factors, we find that our development runs led to **114 tCO<sub>2</sub>eq** emitted and **510 kL** of water consumed.

**Final training runs** Finally, we fully trained our series of models, ranging from 20 million to 7 billion active parameters, with detailed information provided in Table 2. As we saw during development, the majority of the cost incurred came from training our 7B models, which we trained to 2, 3, and 4 trillion tokens. We also see that the 1B dense model required about as much energy per trillion tokens as the MoE model with 1B active parameters, though the MoE model was slightly less efficient, most likely due to the extra compute required for routing tokens. In summary, we find that our training runs led to **140 tCO<sub>2</sub>eq** emitted and **627 kL** of water consumed.

<sup>8</sup><https://www.azcentral.com/story/opinion/op-ed/joannaallhands/2024/06/12/tsmc-arizona-water-use-recycling/74059522007/>

<sup>9</sup><https://www.semiconductor-digest.com/water-supply-challenges-for-the-semiconductor-industry/>

<sup>10</sup><https://anysilicon.com/die-per-wafer-formula-free-calculators/>

Table 2: We list the estimated power usage, carbon emissions, and water consumption from training our dense transformers, ranging from 20 million to 7 billion parameters, trained on 1.7 to 4 trillion tokens, and a mixture-of-experts model with 1 billion active and 7 billion total parameters, trained to 5 trillion tokens. We find that the environmental impact is quite high, even for our relatively small models. Training our series of models emitted equivalent carbon to over 27 years of electricity use by the average household in the U.S., and consumed equivalent water to the average person in the U.S. for 5 and a half years.

|                     | Power Usage (MWh) | Carbon Emissions (tCO <sub>2</sub> eq) | Equiv. to... (energy usage, 1 home, U.S.) | Water Consumption (kL) | Equiv. to... (water usage, 1 person, U.S.) |
|---------------------|-------------------|--|---|------------------------|--|
| <b>BLOOM-176B</b>   | 520               | 30                                     | 4 years                                   | -                      | -  |
| <b>Llama 2 7B</b>   | 81                | 31                                     | 6 yrs, 1 mo                               | -                      | -  |
| <b>Llama 3 8B</b>   | -                 | 420                                    | 83 years                                  | -                      | -  |
| <b>Llama 3.2 1B</b> | -                 | 107                                    | 14 years                                  | -                      | -  |
| <b>OLMo 7B</b>      | 149               | 0*                                     | -   | -                      | -  |
| <b>OLMo 7B</b>      | 114               | 70                                     | 13 yrs, 10 mo                             | -                      | -  |
| <b>LM-20M-1.7T</b>  | 0.8               | 0.3                                    | 3 weeks                                   | 1                      | 4 days                                     |
| <b>LM-60M-1.7T</b>  | 1.2               | 0.4                                    | 1 month                                   | 2                      | 6 days                                     |
| <b>LM-150M-1.7T</b> | 2.4               | 1                                      | 2 mo, 1 wk                                | 4                      | 13 days                                    |
| <b>LM-300M-1.7T</b> | 5                 | 2                                      | 5 months                                  | 7                      | 22 days                                    |
| <b>LM-700M-1.7T</b> | 8                 | 3                                      | 7 months                                  | 12                     | 38 days                                    |
| <b>LM-1B-3T</b>     | 30                | 10                                     | 2 years                                   | 45                     | 5 months                                   |
| <b>LM-7B-2T</b>     | 67                | 22                                     | 4 yrs, 4 mo                               | 100                    | 11 months                                  |
| <b>LM-7B-3T</b>     | 95                | 32                                     | 6 yrs, 4 mo                               | 141                    | 1 yr, 3 mo                                 |
| <b>LM-7B-4T</b>     | 157               | 52                                     | 10 yrs, 4 mo                              | 234                    | 2 yr, 1 mo                                 |
| <b>LM-MoE-5T</b>    | 54                | 18                                     | 3 yrs, 7 mo                               | 81                     | 9 months                                   |
| <b>Total (Ours)</b> | <b>421</b>        | <b>140</b>                             | <b>27 yrs, 7 mo</b>                       | <b>627</b>             | <b>5 yr, 6 mo</b>                          |

**Putting it in perspective** In total, our series of models led to at least **270 tCO<sub>2</sub>eq** emitted. Using the U.S. Environmental Protection Agency’s Greenhouse Gas Equivalencies Calculator<sup>11</sup>, this is equivalent to 3.6 tanker trucks’ worth of gasoline burned, emissions from the average yearly energy use for 35.2 homes in the U.S., or the amount of carbon sequestered by 315 acres of U.S. forests in one year. We additionally estimate we consumed at least **1,137 kL** of water, which is equivalent to about 10 years of water consumption by the average person in the U.S.<sup>12</sup>.

**Other Costs** In this work we strive to provide a thorough accounting of the total cost of developing our models. However, there remain a number of sources of emissions and water consumption that are difficult, if not impossible to comprehensively measure without access to proprietary information across a range of industries, such as transportation and end of life hardware disposal. While the costs we report above represent a large portion of the total development process, more transparency is needed to understand the full impact of model training.

## 4.2 SIMULATING DEPLOYMENT & INFERENCE

We report *simulated* inference costs; that is, we explore the question of what our models’ impact might be if they were put into production. In contrast to §4.1, where we reported the actual impact from our actions, this section reports partial estimates of Scope 3 carbon emissions and water consumption: the impact from the downstream actions of others using our models. We include comparisons with recent instruction-tuned models as well.

In Table 3, we display 1) power and energy costs, 2) carbon and water consumption, and 3) the time to complete 100 requests. We additionally report “breakeven” points, that is the number of inferences in each scenario required for inference costs to be equal or greater to training costs.

<sup>11</sup><https://www.epa.gov/energy/greenhouse-gas-equivalencies-calculator>

<sup>12</sup><https://www.epa.gov/watersense/statistics-and-facts>

Table 3: Measurements and estimates of resource costs from SGLang benchmarking on 2400 prompts from ShareGPT at varying request rates. Since the models were served on machines from the same cluster that our models were trained on, we have the same PUE and WUE coefficients of 1.2 and 1.49 L / kWh respectively, and carbon intensity of 0.332 kg CO<sub>2</sub>e / kWh – note the difference in units for energy consumption and carbon emissions, namely MWh → kWh, tons → grams CO<sub>2</sub>eq, and kL → L. The measurements reported in this table account for the processes associated with active inference, but not server startup time or overhead. Thus, these numbers can be considered as strictly lower bounds on usage in similar settings. Also of note is the relatively small variability in carbon emissions and water consumption across different model sizes in cases where batches are not saturated, despite faster inference in smaller models in fully saturated batching scenarios – greater peak efficiency does not guarantee efficient deployment, as resource consumption can be heavily influenced by total uptime of a service.

|                      | Request freq. | GPU Power Usage (kWh) | Total Process Energy (kWh) | Carbon Emissions (g CO <sub>2</sub> eq) | Water consump. (L) | Seconds per 100 req. | # Inf. for CO <sub>2</sub> equiv. w/ training |
|----------------------|---------------|-----------------------|----------------------------|---|--------------------|----------------------|---|
| <b>Llama 3.2 1B</b>  | ∞             | 0.003                 | 0.007                      | 2.3                                     | 0.010              | 1.02                 | 110.5 bil                                     |
|                      | 8 / sec       | 0.032                 | 0.084                      | 27.9                                    | 0.126              | 12.65                | 9.2 bil                                       |
|                      | 1 / sec       | 0.154                 | 0.662                      | 219.8                                   | 0.986              | 100.59               | 1.2 bil                                       |
| <b>Llama 2 7B</b>    | ∞             | 0.020                 | 0.036                      | 12.0                                    | 0.053              | 4.20                 | 6.2 bil                                       |
|                      | 8 / sec       | 0.052                 | 0.106                      | 35.2                                    | 0.158              | 12.87                | 2.1 bil                                       |
|                      | 1 / sec       | 0.331                 | 0.855                      | 283.9                                   | 1.274              | 100.64               | 262.1 mil                                     |
| <b>Llama 3 8B</b>    | ∞             | 0.011                 | 0.021                      | 7.0                                     | 0.032              | 2.44                 | 144.6 bil                                     |
|                      | 8 / sec       | 0.050                 | 0.107                      | 35.5                                    | 0.160              | 12.81                | 28.4 bil                                      |
|                      | 1 / sec       | 0.330                 | 0.856                      | 284.2                                   | 1.276              | 100.64               | 3.6 bil                                       |
| <b>LM-1B-3T</b>      | ∞             | 0.004                 | 0.009                      | 3.0                                     | 0.013              | 1.26                 | 8.0 bil                                       |
|                      | 8 / sec       | 0.034                 | 0.084                      | 27.9                                    | 0.125              | 12.64                | 860.6 mil                                     |
|                      | 1 / sec       | 0.165                 | 0.676                      | 224.4                                   | 1.008              | 100.58               | 106.9 mil                                     |
| <b>LM-7B-4T</b>      | ∞             | 0.019                 | 0.033                      | 11.0                                    | 0.049              | 4.10                 | 11.4 bil                                      |
|                      | 8 / sec       | 0.049                 | 0.096                      | 31.9                                    | 0.144              | 12.80                | 3.9 bil                                       |
|                      | 1 / sec       | 0.321                 | 0.818                      | 271.6                                   | 1.219              | 100.60               | 459.5 mil                                     |
| <b>LM-1BA-7BT-5T</b> | ∞             | 0.007                 | 0.017                      | 5.6                                     | 0.025              | 2.11                 | 7.7 bil                                       |
|                      | 8 / sec       | 0.037                 | 0.097                      | 32.2                                    | 0.144              | 12.82                | 1.3 bil                                       |
|                      | 1 / sec       | 0.146                 | 0.650                      | 215.8                                   | 0.969              | 100.60               | 200.2 mil                                     |

Surprisingly, we find that for most models tested, the number of inferences required to outweigh training costs is in the hundreds of millions to tens of billions, except for the most over-trained models. As many of these models were created to be efficient in deployment-focused scenarios – such as on edge devices, or in popular online products – it is important to consider inference costs in addition to training costs. The largest model providers are producing up to hundreds of billions of tokens per day<sup>13</sup>, highlighting that deployed models can quickly reach this tipping point.

### 4.3 POWER FLUCTUATIONS DURING TRAINING

One problem caused by training AI models at large scales is that the power demand starts and stops suddenly (Dubey et al., 2024), which power grids can struggle to handle. When demand sharply rises, generation sources that can be quickly started and stopped – generally powered by fossil fuels, such as coal and natural gas – must be brought online quickly, increasing the marginal carbon intensity of the grid and potentially negatively impacting other consumers in cases where demand rises more quickly than generation can handle. When demand sharply drops, excess power is discarded – by grounding the power or venting steam – until generation sources can spin down. Power grids can generally manage some large variations (for example, when communities experience a sudden power outage), but as we add more variability to the system, it becomes more difficult to maintain this delicate balance, and infrastructure is not set up to handle frequent, large fluctuations.

In Figure 2, we show a snapshot of our model’s GPU power consumption during pre-training. We find that power consumption is not consistent – instead, power is consistent *while the model is train-*

<sup>13</sup><https://x.com/sama/status/1756089361609981993>



ing, but drops quickly while saving checkpoints. Though our models are relatively small, and we have since improved checkpointing performance, other model developers have experienced similar issues caused by checkpointing and synchronization between nodes (Dubey et al., 2024).

## 5 DISCUSSION

### 5.1 MORE TRANSPARENCY IS (STILL) NEEDED

While many model developers—including some of the largest for profit entities operating in this space—make best efforts to report at least part of the cost of building their AI systems (Dubey et al., 2024; Team et al., 2024), more transparency is still needed throughout the development pipeline. Proposed legislation, such as the Artificial Intelligence Environmental Impacts Act<sup>14</sup> in the United States, would start the process for defining voluntary environmental impact reporting standards for model developers, but until such standards are created and accepted in the community, improved transparency can only come through voluntary efforts by companies and research organizations. Policy action is needed to ensure there is public visibility into environmental impacts across the entire supply chain, from hardware manufacturing, data center construction, and energy production, all the way through to model deployment and inference.

#### Embodied emissions are still an enigma

Though a vital piece of all model development pipelines, the environmental impact of manufacturing the GPUs used to train models is essentially unknown. In previous work, Luccioni et al. (2023) highlighted the fact that researchers focused on AI’s environmental impact are forced to use unreliable estimates of the cost of manufacturing state of the art computational hardware, and the situation is no better now, nearly two years later. Many companies that manufacture other pieces of data center hardware disclose estimates of the lifetime environmental impact,<sup>15</sup> and until GPU manufacturers release similar information—on a voluntary or compulsory basis—this will not improve.

**Development costs are substantial, and unreported** As reported in Section 4.1, we present detailed information on the cost of developing our training pipeline, in contrast with previous work. We found that development costs—associated with failed runs, hyperparameter searches, testing architecture changes, and more—are responsible for a substantial portion of the total environmental impact of creating our systems, highlighting a need for more transparency from model developers. This is especially important in light of AutoML tools, where many models may be automatically trained while searching for a solution, and scaling law experiments, where many smaller models are trained to predict the performance of larger models, and then discarded (Li et al., 2024).

**Water costs are real, and under-explored** While under-explored in previous work, AI’s growing water consumption is beginning to receive more and more attention<sup>16</sup> (Li et al., 2023), though not as much as it may deserve. As shown in Section 4.1, even training a series of comparatively small models uses an enormous amount of water, the amount of which is also drastically impacted by both the cooling systems used in data centers as well as the power generation methods used. Without

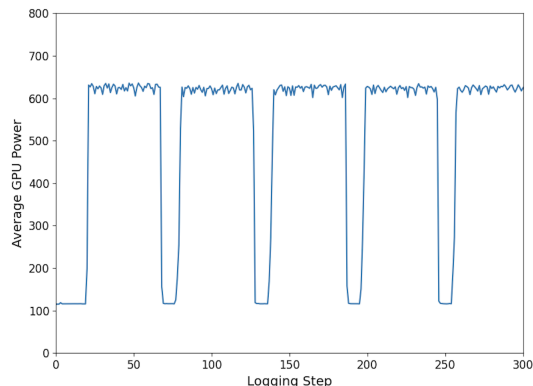


Figure 2: Average GPU power for a single node for the first 300 logging steps during LM-7B-4T training. The first spike is the beginning of training, and each drop happens when a model checkpoint is saved. When actively training, the average GPU power is over 600W, over 85% of an H100’s maximum power draw of 700W, and during checkpointing, power usage drops to just over 100W, or about 15% maximum.

<sup>14</sup>[https://www.markey.senate.gov/imo/media/doc/artificial\\_intelligence\\_environmental\\_impacts\\_act\\_of\\_2024\\_-\\_020124pdf.pdf](https://www.markey.senate.gov/imo/media/doc/artificial_intelligence_environmental_impacts_act_of_2024_-_020124pdf.pdf)

<sup>15</sup><https://www.hpe.com/psnow/doc/a50005151enw>

<sup>16</sup><https://www.washingtonpost.com/technology/2024/09/18/energy-ai-use-electricity-water-data-centers/>

486 more transparency from developers on when, where, and how they are training their models, it will  
 487 continue to be difficult to quantify the scale of the issue, stymieing efforts to address it.  
 488

## 490 5.2 SMALL CHOICES DURING TRAINING CAN HAVE LARGE IMPACTS

491 While many issues relating to transparency require action from corporations and large research  
 492 groups, choices made during training have a large effect downstream.  
 493

494  
 495 **Smaller models are cheaper to train and use, but at what cost?** Until recently, to achieve high  
 496 model performance, a large model was needed. Compute-optimal scaling laws for neural network  
 497 training (Hoffmann et al., 2022; Kaplan et al., 2020) imply that it is more efficient to put more data  
 498 into a larger model, because of diminishing returns from “over-training” a small model. This meant  
 499 that models were expensive to both train and deploy, limiting how widespread they could become,  
 500 and how financially feasible they were to be used in a variety of scenarios.

501 Recently, however, continuing to train models on more and more tokens beyond the “compute-  
 502 optimal” limit<sup>17</sup> has been extremely successful in making “deployment-optimized” models that can  
 503 be substantially cheaper to perform inference with. This has led to an explosion in both training  
 504 cost for small models, and total inference compute cost, as API-based models become cheaper to  
 505 use<sup>18,19</sup> and small models are deployed on-device (Gunter et al., 2024; Abdin et al., 2024). This may  
 506 be an instance of Jevons’ Paradox (Jevons, 1865): when a resource’s efficiency increases, overall  
 507 consumption of that resource tends to increase, rather than decrease. In other words, as the financial  
 508 and environmental cost of training models decreases, the downstream impact may continue to grow.

509 This is especially clear in context of our results in Section 4.2, showing that though the raw num-  
 510 ber of inferences required to outweigh training is objectively quite large, smaller models are being  
 511 deployed in many new scenarios that will drastically increase their total usage. Many inference use  
 512 cases are also not able to be batched (e.g. generating text on a phone for immediate use), meaning  
 513 that deployers cannot schedule many of these requests to take advantage of cheaper and/or cleaner  
 514 energy, and instead must make use of immediately available power. Given that this trend will most  
 515 likely only accelerate, it is vital that we quickly improving transparency into the total cost of de-  
 516 ployment in all deployment scenarios.  
 517

### 518 **Power fluctuations reveal inefficiencies at best, challenges to power grid control at worst**

519 While it is known that the dramatic spike in power consumption at the beginning of training and the  
 520 subsequent drop at the end are problematic for power grid operators at large scales, little has been  
 521 discussed publicly about how power consumption changes throughout training. We found that our  
 522 models, using an optimized code base and publicly available tooling, sees rapid power fluctuations  
 523 throughout training caused by the commonplace practice of frequently saving model checkpoints.  
 524 This means that without careful engineering, one training run can cause thousands of rapid power  
 525 fluctuations, which poses an immediate challenge for large-scale LLM training in data centers, which  
 526 typically source energy directly from power providers. Generated power needs to go somewhere,  
 527 and rapid, large drops in consumption during training breaks common assumptions about data center  
 528 supply and demand, leading to significant control challenges in power systems. While some frame-  
 529 works have begun to implement workarounds to manage this issue,<sup>20</sup> more awareness is needed on  
 530 the part of researchers and engineers as training runs scale to tens of thousands of GPUs<sup>21</sup> or more,  
 531 as even some of the largest model developers encounter difficulties from regularly shifting power  
 532 demand throughout training (Dubey et al., 2024). We emphasize that addressing this will require  
 533 more comprehensive solutions such as parallelized checkpointing, improved demand response in  
 534 data centers running large AI workloads, and new, heterogeneous methods for distributed training  
 535 spanning software, hardware, and scheduling.

536 <sup>17</sup>e.g. scaling from 1 to 2 to 15T tokens for Llama 1, 2, and 3 (Touvron et al., 2023a;b; Dubey et al., 2024)

537 <sup>18</sup><https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>

538 <sup>19</sup><https://developers.googleblog.com/en/gemini-15-flash-updates-google-ai-studio-gemini-api/>

539 <sup>20</sup>E.g. the new `PYTORCH_NO_POWERPLANT_BLOWUP` environment variable in PyTorch.

<sup>21</sup><https://time.com/7021709/elon-musk-xai-grok-memphis/>

## REFERENCES

- 540  
541  
542 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen  
543 Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko,  
544 Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dong-  
545 dong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang  
546 Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit  
547 Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao,  
548 Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin  
549 Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim,  
550 Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden,  
551 Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong  
552 Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro  
553 Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norrick, Barun Patra, Daniel Perez-  
554 Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo  
555 de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim,  
556 Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla,  
557 Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua  
558 Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp  
559 Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Ji-  
560 long Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan,  
561 Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan  
562 Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your  
563 phone, 2024. URL <https://arxiv.org/abs/2404.14219>.
- 564 Callum Browning, Stephen Northey, Nawshad Haque, Warren Bruckard, and Mark Cooksey. *Life*  
565 *Cycle Assessment of Rare Earth Production from Monazite*, pp. 83–88. Springer International  
566 Publishing, Cham, 2016. ISBN 978-3-319-48768-7. doi: 10.1007/978-3-319-48768-7\_12. URL  
[https://doi.org/10.1007/978-3-319-48768-7\\_12](https://doi.org/10.1007/978-3-319-48768-7_12).
- 567 Benoit Courty, Victor Schmidt, Goyal-Kamal, MarionCoutarel, Luis Blanche, Boris Feld, inimaz,  
568 Jérémy Lecourt, LiamConnell, SabAmine, supatomic, Mathilde Léval, Alexis Cruveiller, oumi-  
569 nasara, Franklin Zhao, Aditya Joshi, Christian Bauer, Amine Saboni, Patrick LLORET, Alexis  
570 Bogroff, Niko Laskaris, Hugues de Lavoreille, Alexandre Phiev, Edoardo Abati, rosekelly6400,  
571 Douglas Blank, Ziyao Wang, Lucas Otávio, and Armin Catovic. mlco2/codecarbon: v2.7.1,  
572 September 2024. URL <https://doi.org/10.5281/zenodo.13744486>.
- 573 Jesse Dodge, Taylor Prewitt, Remi Tachet Des Combes, Erika Odmark, Roy Schwartz, Emma  
574 Strubell, Alexandra Sasha Luccioni, Noah A. Smith, Nicole DeCario, and Will Buchanan. Mea-  
575 suring the carbon intensity of ai in cloud instances, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2206.05229)  
576 [2206.05229](https://arxiv.org/abs/2206.05229).
- 577 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
578 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
579 *arXiv preprint arXiv:2407.21783*, 2024.
- 580 Dirk Groeneveld, Iz Beltagy, Evan Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya  
581 Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Au-  
582 thur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel,  
583 Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crys-  
584 tal Nam, Matthew Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh  
585 Shah, William Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi,  
586 Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini,  
587 Noah Smith, and Hannaneh Hajishirzi. OLMO: Accelerating the science of language models. In  
588 Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meet-*  
589 *ing of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15789–15809,  
590 Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/  
591 2024.acl-long.841. URL <https://aclanthology.org/2024.acl-long.841>.
- 592  
593 Tom Gunter, Zirui Wang, Chong Wang, Ruoming Pang, Andy Narayanan, Aonan Zhang, Bowen  
Zhang, Chen Chen, Chung-Cheng Chiu, David Qiu, Deepak Gopinath, Dian Ang Yap, Dong

- 594 Yin, Feng Nan, Floris Weers, Guoli Yin, Haoshuo Huang, Jianyu Wang, Jiarui Lu, John Pee-  
595 bles, Ke Ye, Mark Lee, Nan Du, Qibin Chen, Quentin Keunebroek, Sam Wiseman, Syd Evans,  
596 Tao Lei, Vivek Rathod, Xiang Kong, Xianzhi Du, Yanghao Li, Yongqiang Wang, Yuan Gao,  
597 Zaid Ahmed, Zhaoyang Xu, Zhiyun Lu, Al Rashid, Albin Madappally Jose, Alec Doane, Alfredo  
598 Bencomo, Allison Vanderby, Andrew Hansen, Ankur Jain, Anupama Mann Anupama, Areeba  
599 Kamal, Bugu Wu, Carolina Brum, Charlie Maalouf, Chinguun Erdenebileg, Chris Dulhanty, Do-  
600 minik Moritz, Doug Kang, Eduardo Jimenez, Evan Ladd, Fangping Shi, Felix Bai, Frank Chu,  
601 Fred Hohman, Hadas Kotek, Hannah Gillis Coleman, Jane Li, Jeffrey Bigham, Jeffery Cao, Jeff  
602 Lai, Jessica Cheung, Jiulong Shan, Joe Zhou, John Li, Jun Qin, Karanjeet Singh, Karla Vega,  
603 Kelvin Zou, Laura Heckman, Lauren Gardiner, Margit Bowler, Maria Cordell, Meng Cao, Nicole  
604 Hay, Nilesh Shahdadpuri, Otto Godwin, Pranay Dighe, Pushyami Rachapudi, Ramsey Tantawi,  
605 Roman Frigg, Sam Davarnia, Sanskruti Shah, Saptarshi Guha, Sasha Sirovica, Shen Ma, Shuang  
606 Ma, Simon Wang, Sulgi Kim, Suma Jayaram, Vaishaal Shankar, Varsha Paidi, Vivek Kumar,  
607 Xin Wang, Xin Zheng, Walker Cheng, Yael Shrager, Yang Ye, Yasu Tanaka, Yihao Guo, Yun-  
608 song Meng, Zhao Tang Luo, Zhi Ouyang, Alp Ayyar, Alvin Wan, Andrew Walkingshaw, Andy  
609 Narayanan, Antonie Lin, Arsalan Farooq, Brent Ramerth, Colorado Reed, Chris Bartels, Chris  
610 Chaney, David Riazati, Eric Liang Yang, Erin Feldman, Gabriel Hochstrasser, Guillaume Seguin,  
611 Irina Belousova, Joris Pelemans, Karen Yang, Keivan Alizadeh Vahid, Liangliang Cao, Mah-  
612 yar Najibi, Marco Zuliani, Max Horton, Minsik Cho, Nikhil Bhendawade, Patrick Dong, Piotr  
613 Maj, Pulkit Agrawal, Qi Shan, Qichen Fu, Regan Poston, Sam Xu, Shuangning Liu, Sushma  
614 Rao, Tashweena Heeramun, Thomas Merth, Uday Rayala, Victor Cui, Vivek Rangarajan Sridhar,  
615 Wencong Zhang, Wenqi Zhang, Wentao Wu, Xingyu Zhou, Xinwen Liu, Yang Zhao, Yin Xia,  
616 Zhile Ren, and Zhongzheng Ren. Apple intelligence foundation language models, 2024. URL  
<https://arxiv.org/abs/2407.21075>.
- 617 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
618 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom  
619 Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aure-  
620 lia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and  
621 L. Sifre. Training compute-optimal large language models. *ArXiv*, abs/2203.15556, 2022. URL  
622 <https://api.semanticscholar.org/CorpusID:247778764>.
- 623 William Stanley Jevons. *The Coal Question; An Inquiry Concerning the Progress of the Nation, and*  
624 *the Probable Exhaustion of Our Coal Mines*. London: Macmillan and Co, 1865.
- 625  
626 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,  
627 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language  
628 models. *arXiv preprint arXiv:2001.08361*, 2020.
- 629  
630 Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Gadre, Hritik Bansal,  
631 Etash Guha, Sedrick Keh, Kushal Arora, Saurabh Garg, Rui Xin, Niklas Muennighoff, Rein-  
632 hard Heckel, Jean Mercat, Mayee Chen, Suchin Gururangan, Mitchell Wortsman, Alon Al-  
633 balak, Yonatan Bitton, Marianna Nezhurina, Amro Abbas, Cheng-Yu Hsieh, Dhruva Ghosh,  
634 Josh Gardner, Maciej Kilian, Hanlin Zhang, Rulin Shao, Sarah Pratt, Sunny Sanyal, Gabriel Il-  
635 harco, Giannis Daras, Kalyani Marathe, Aaron Gokaslan, Jieyu Zhang, Khyathi Chandu, Thao  
636 Nguyen, Igor Vasiljevic, Sham Kakade, Shuran Song, Sujay Sanghavi, Fartash Faghri, Se-  
637 woong Oh, Luke Zettlemoyer, Kyle Lo, Alaaeldin El-Nouby, Hadi Pouransari, Alexander Toshev,  
638 Stephanie Wang, Dirk Groeneveld, Luca Soldaini, Pang Wei Koh, Jenia Jitsev, Thomas Kol-  
639 lar, Alexandros G. Dimakis, Yair Carmon, Achal Dave, Ludwig Schmidt, and Vaishaal Shankar.  
640 Datacomp-1m: In search of the next generation of training sets for language models, 2024. URL  
641 <https://arxiv.org/abs/2406.11794>.
- 642 Pengfei Li, Jianyi Yang, Mohammad A. Islam, and Shaolei Ren. Making ai less "thirsty": Uncover-  
643 ing and addressing the secret water footprint of ai models, 2023. URL <https://arxiv.org/abs/2304.03271>.
- 644  
645  
646 Alexandra Sasha Luccioni, Sylvain Viguier, and Anne-Laure Ligozat. Estimating the carbon foot-  
647 print of bloom, a 176b parameter language model. *Journal of Machine Learning Research*, 24  
(253):1–15, 2023. URL <http://jmlr.org/papers/v24/23-0069.html>.

- 648 Sasha Luccioni, Yacine Jernite, and Emma Strubell. Power hungry processing: Watts driving the  
649 cost of ai deployment? In *The 2024 ACM Conference on Fairness, Accountability, and Trans-*  
650 *parency*, pp. 85–99, 2024.
- 651
- 652 Sachin Mehta, Mohammad Hossein Sekhavat, Qingqing Cao, Maxwell Horton, Yanzi Jin, Chen-  
653 fan Sun, Seyed Iman Mirzadeh, Mahyar Najibi, Dmitry Belenko, Peter Zatloukal, and Moham-  
654 mad Rastegari. OpenELM: An efficient language model family with open training and inference  
655 framework. In *Workshop on Efficient Systems for Foundation Models II @ ICML2024*, 2024.  
656 URL <https://openreview.net/forum?id=XNMbTkxroF>.
- 657 Hao Peng, Qingqing Cao, Jesse Dodge, Matthew E. Peters, Jared Fernandez, Tom Sherborne, Kyle  
658 Lo, Sam Skjonsberg, Emma Strubell, Darrell Plessas, Iz Beltagy, Evan Pete Walsh, Noah A.  
659 Smith, and Hannaneh Hajishirzi. Efficiency pentathlon: A standardized arena for efficiency eval-  
660 uation, 2023. URL <https://arxiv.org/abs/2307.09701>.
- 661
- 662 Vijay Janapa Reddi, Christine Cheng, David Kanter, Peter Mattson, Guenther Schmuelling, Carole-  
663 Jean Wu, Brian Anderson, Maximilien Breughe, Mark Charlebois, William Chou, Ramesh  
664 Chukka, Cody Coleman, Sam Davis, Pan Deng, Greg Diamos, Jared Duke, Dave Fick, J. Scott  
665 Gardner, Itay Hubara, Sachin Idgunji, Thomas B. Jablin, Jeff Jiao, Tom St. John, Pankaj Kanwar,  
666 David Lee, Jeffery Liao, Anton Lokhmotov, Francisco Massa, Peng Meng, Paulius Micikevicius,  
667 Colin Osborne, Gennady Pekhimenko, Arun Tejusve Raghunath Rajan, Dilip Sequeira, Ashish  
668 Sirasao, Fei Sun, Hanlin Tang, Michael Thomson, Frank Wei, Ephrem Wu, Lingjie Xu, Koichi  
669 Yamada, Bing Yu, George Yuan, Aaron Zhong, Peizhao Zhang, and Yuchen Zhou. Mlperf in-  
670 ference benchmark. In *2020 ACM/IEEE 47th Annual International Symposium on Computer*  
*Architecture (ISCA)*, pp. 446–459, 2020. doi: 10.1109/ISCA45697.2020.00045.
- 671
- 672 Paul Reig, Tianyi Luo, Eric Christensen, and Julie Sinistore. Guidance for calculating water  
673 use embedded in purchased electricity, 2020. URL [https://www.wri.org/research/](https://www.wri.org/research/guidance-calculating-water-use-embedded-purchased-electricity)  
674 [guidance-calculating-water-use-embedded-purchased-electricity](https://www.wri.org/research/guidance-calculating-water-use-embedded-purchased-electricity).
- 675
- 676 Roy Schwartz, Jesse Dodge, Noah A. Smith, and Oren Etzioni. Green ai. *Commun. ACM*, 63(12):  
677 54–63, November 2020. ISSN 0001-0782. doi: 10.1145/3381831. URL [https://doi.org/](https://doi.org/10.1145/3381831)  
[10.1145/3381831](https://doi.org/10.1145/3381831).
- 678
- 679 Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for  
680 modern deep learning research. In *Proceedings of the AAAI conference on artificial intelligence*,  
681 volume 34, pp. 13693–13696, 2020.
- 682
- 683 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya  
684 Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard  
685 Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex  
686 Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, An-  
687 tonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo,  
688 Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric  
689 Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Hen-  
690 ryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski,  
691 Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu,  
692 Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee,  
693 Lucas Dixon, Machel Reid, Maciej Mikula, Mateo Wirth, Michael Sharman, Nikolai Chinaev,  
694 Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko  
695 Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruiho  
696 Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree  
697 Pandya, Siamak Shakeri, Soham De, Ted Klimentko, Tom Hennigan, Vlad Feinberg, Wojciech  
698 Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh  
699 Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin  
700 Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah  
701 Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. Gemma: Open models based on  
gemma research and technology, 2024. URL <https://arxiv.org/abs/2403.08295>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Ar-

- 702 mand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation  
703 language models, 2023a. URL <https://arxiv.org/abs/2302.13971>.  
704
- 705 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
706 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
707 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 708 Carole-Jean Wu, Ramya Raghavendra, Udit Gupta, Bilge Acun, Newsha Ardalani, Kiwan Maeng,  
709 Gloria Chang, Fiona Aga, Jinshi Huang, Charles Bai, et al. Sustainable ai: Environmental impli-  
710 cations, challenges and opportunities. *Proceedings of Machine Learning and Systems*, 4:795–813,  
711 2022.
- 712 Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao,  
713 Christos Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark Barrett, and Ying Sheng. Sglang:  
714 Efficient execution of structured language model programs, 2024. URL <https://arxiv.org/abs/2312.07104>.  
715  
716  
717

## 718 A INFERENCE SIMULATION DETAILS

719

720 Additional details, currently omitted in order to preserve anonymity, will be shared upon publication.  
721

### 722 A.1 LIMITATIONS

723

724 We present only a limited set of inference simulations following a number of simplistic assumptions.

725 Specifically, we simulate only settings where a deployed model is ingesting input tokens and gen-  
726 erating output tokens following default parameters defined in SGLang (Zheng et al., 2024) – as  
727 opposed to, for instance, evaluating only the likelihood of a given text.  
728

729 Additionally, we note that practitioners frequently quantize LLMs before deploying them, and/or  
730 deploy to and run inference on edge device, sometimes even without GPUs. We do not account for  
731 these scenarios in our experiments.  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755