

Environmental Footprint of GenAI Research: Insights from the Moshi Foundation Model

Anonymous authors
Paper under double-blind review

Abstract

New multi-modal large language models (MLLMs) are continuously being trained and deployed, following rapid development cycles. This generative AI frenzy is driving steady increases in energy consumption, greenhouse gas emissions, and a plethora of other environmental impacts linked to datacenter construction and hardware manufacturing. Mitigating the environmental consequences of GenAI remains challenging due to an overall lack of transparency by the main actors in the field. Even when the environmental impacts of specific models are mentioned, they are typically restricted to the carbon footprint of the final training run, omitting the research and development stages.

In this work, we explore the impact of GenAI research through a fine-grained analysis of the compute spent to create Moshi, a 7B-parameter speech-text foundation model for real-time dialogue developed by Kyutai, a leading privately funded open science AI lab. For the first time, our study dives into the anatomy of compute-intensive MLLM research, quantifying the GPU-time invested in specific model components and training phases, as well as early experimental stages, failed training runs, debugging, and ablation studies. Additionally, we assess the environmental impacts of creating Moshi from beginning to end using a life cycle assessment methodology: we quantify energy and water consumption, greenhouse gas emissions, and mineral resource depletion associated with the production and use of datacenter hardware.

Our detailed analysis allows us to provide actionable guidelines to reduce compute usage and environmental impacts of MLLM research, paving the way for more sustainable AI research.

1 Introduction

The environmental footprint of modern artificial intelligence systems has become a growing concern (Zhuk, 2023), as the rapid scaling of its energetic requirements unfolds on a planet already under significant ecological strain (Rockström et al., 2024). Large text and multimodal models now require millions of GPU-hours for training alone (Zhao et al., 2025; Le Scao et al., 2022), raising questions about the sustainability of current development practices (Varoquaux et al., 2025) and MLLM research itself.

In this work, we present a detailed analysis of the full development life cycle of Moshi, a state-of-the-art speech-to-text foundation model developed by Kyutai, a leading research organization in large language models and speech technologies. Rather than focusing exclusively on the final training run or hyperparameter search, we study the complete sequence of research activities that led to the released system, from early exploratory experimentation through final training.

This broader perspective is essential. Indeed, most research works in machine learning only report the cost of the training of the final model, implicitly assuming it to be the dominant contributor to overall cost. A few environmental impact studies have also considered the cost of hyperparameter search, and others the impact of inference, but research practices largely remain terra incognita. In practice, MLLM research involves extensive experimentation, debugging, hyperparameter tuning, architectural exploration, benchmarking, ablation studies, and discarded runs. These activities can collectively account for a substantial

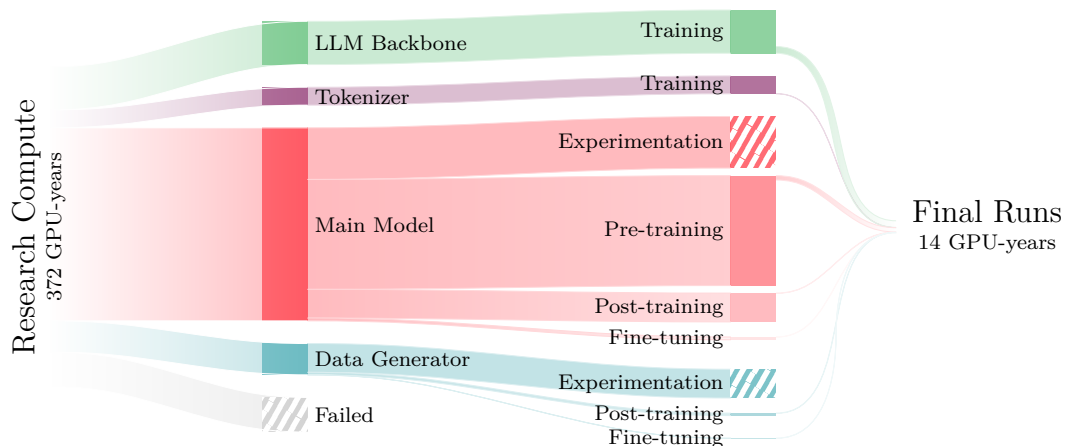


Figure 1: **From research to final compute.** *Research Compute* is split among individual model components and their respective training phases. *Failed* reflects the cost of failed experiments, and *Experimentation* gathers early versions that differ significantly from the definitive architecture and training scheme choices for specific components. *Final Runs* isolates the compute of training only one definitive version of each model component.

share of both compute usage and environmental impact, yet they are rarely measured or reported. Our study is enabled by an unprecedented level of access to internal training logs provided by Kyutai. This access allows us, for the first time, to quantify the computational and environmental impact of *all* stages of a specific industry-scale MLLM research project, beyond standard hyperparameter search.

In addition, we conduct a comprehensive life cycle assessment (LCA) of the Moshi research project. Beyond operational energy consumption, we estimate operational and embodied impacts in four impact categories.

Our analysis yields several key findings. First, we show that training the final deployed model accounts for only a small fraction of the total environmental impact, around 4%, while the environmental impact of experimentation, debugging, failed runs, model evaluation and ablation studies is significant. Second, we find that the ratio between the full cost and the final training highly depends on the novelty of the approach, being around $6.5\times$ for the relatively standard LLM backbone but around $40\times$ for the main Moshi model. Finally, we observe a strong decoupling between the number of runs and their intensity: 13% of runs account for nearly 89% of the total compute. Together, these results shed new light on current development workflows and highlight concrete opportunities to reduce the environmental footprint of future AI systems.

2 Related Work

We first review the core literature on LCA methodology and discuss the specific challenges of applying it to AI systems (sec. 2.1). We then survey prior work that applies these methodologies to AI systems (sec. 2.2). Finally, we provide an overview of Moshi, the system analyzed in this paper (sec. 2.3).

2.1 Methodologies for Environmental Assessment

Life Cycle Assessment (LCA) Methodology. Life cycle assessment (LCA) (Hauschild et al., 2018; Heijungs & Suh, 2002) is an environmental impact assessment methodology formalized in the 14040/44 ISO standards. It was proposed as a holistic approach to evaluating the environmental impacts incurred throughout the *life cycle* of a product system: raw material extraction, manufacturing, transport, use, and end of life. LCA considers a variety of *impact categories* with effects on human health, natural resources, and the natural environment, measured via *indicators* such as global warming potential, water consumption,

or human toxicity. These impacts are assigned with respect to a *functional unit*: a quantitative description of a function provided by the system at a desired level of performance.

LCA is essential to diagnose impacts shifting between life cycle phases and impact categories. For example, using more efficient hardware reduces energy consumption, but increases production impacts due to semiconductor manufacturing.

LCA Methodology for AI Systems An AI system encompasses an AI model and the tangible infrastructure involved in its creation and deployment: sensors used for data collection, hyperscale datacenters for training, servers for hosting the model and running inference, etc. Initiatives for adapting the LCA methodology to AI systems have recently been proposed (Ligozat et al., 2022; OECD.AI Expert Group on AI Compute and Climate, 2022; Kaack et al., 2022), based on frameworks specific to information and communications technology (ICT) (Hilty & Hercheui, 2010; ITU, 2024). These developments have led to tools such as MLCA (Morand et al., 2024) or Boavizta (Simon et al., 2025), which we leverage in this work.

Several distinctions in the type of impact are important for our work. First, the impact can be attributed to three types of effects (Horner et al., 2016; Hilty & Hercheui, 2010; Kaack et al., 2022): (i) *first-order effects* are directly related to the development and operation of AI systems; (ii) *second-order effects* result from changes in industry when using an AI system; (iii) *third-order effects* are large-scale changes in lifestyle and economic structures following the widespread use of AI. Second, first-order impacts can be divided into *operational impacts* incurred directly while using the hardware, and *embodied impacts* corresponding to the remaining life cycle phases of the hardware, such as manufacturing, transport and end of life (Horner et al., 2016; Gupta et al., 2022; Kaack et al., 2022). Third, the *AI system development life cycle* (Wu et al., 2022) can be divided into four stages: (i) data collection, processing, and storage; (ii) research and development; (iii) model training; (iv) model deployment (inference). We highlight the distinction between *research*, which involves free-range experimentation on modeling choices, model architecture design, training techniques, and other aspects; and *development*, which entails more structured hyperparameter searches and scaling law experiments in preparation for final model training.

In this work, we consider the operational and embodied first-order impacts of the research and development and training stages of a speech-text AI foundation model. Detailing the impact of the research phase is the key novelty of our work compared to existing assessments, that we detail in the next section.

2.2 Environmental Assessments of AI Systems.

In this section, we give an overview of existing environmental reports for AI systems. We first list some works that assess the impact of AI services without considering the research and development stage. We then outline works that focus on research and development impacts and are closer to our study. Finally, we mention company reports on their LLM development costs, which are relevant but often remain very high-level.

AI System Deployment. As the use of commercial AI solutions has become widespread, a body of work has focused on assessing the growing impacts of model deployment, considering both embodied and operational impacts. Such works assess impacts ranging from just global warming potential (Gupta et al., 2022; Chien et al., 2023; Li et al., 2025b) up to several environmental impact indicators, including abiotic depletion potential (Berthelot et al., 2024) or water consumption (Elsworth et al., 2025). Jegham et al. (2025) extend their assessment of model deployment to three impact indicators, but do not consider embodied impacts. Opposite to these works, we focus on research and development costs.

AI System Research and Development. A few published studies assess the impacts of training suites of LLMs while also considering *development* overheads. Among these, some report the impact of development activities as a single number (Lakim et al., 2022), whereas others provide breakdowns by model size (Morrison et al., 2025).

Strubell et al. (2019) were arguably the first to quantify the *research and development* cost of a novel NLP model. In their work, they compare the impacts of training a single instance of the model, of tuning the

model, and of the full research and development process. They also estimate the impacts of running a neural architecture search to improve model architecture, which was revisited by Patterson et al. (2021). Other more recent works also consider research costs, with varying levels of granularity: Wu et al. (2022) merge research, development, and training costs into a single “offline training” cost, whereas Luccioni et al. (2023) provide a breakdown by model size and high-level activity, i.e. model evaluation and miscellaneous processes. In contrast, we provide a detailed breakdown of research costs for a foundation model approaching a task in a radically novel way.

Company Reports on LLM Development. Well-known LLM developers have reported the environmental footprint of training specific models, but, to the best of our knowledge, none of these reports give detailed insights of the research and development process: Google and Meta AI estimate the final training carbon footprints of T5, GPT-3 (Patterson et al., 2021), Gemma (Gemma Team et al., 2024), the Llama family (Touvron et al., 2023a;b; Meta, 2024), OPT (Zhang et al., 2022), and other models, without considering research and development costs. The OPT model report (Zhang et al., 2022) provides a logbook registering informal comments made during development, but does not quantify the impact of the registered incidences. Similarly, the model report of Gopher (Rae et al., 2022), which focuses on compute, omits “compute arising from development, pre-emption, or other sources of inefficiency”, although it does offer a detailed account of the compute spent on evaluation across several benchmarks.

These same entities do, however, quantify research and development costs at the company level: Wu et al. (2022) mention the share of infrastructure power capacity devoted to experimentation at Facebook AI, and indicate the compute intensity of typical experimental runs. On the other hand, Patterson et al. (2022) report the energy consumption spent on machine learning at Google -including research, development, testing, and production- but do not isolate the consumption of each of these activities.

Other companies are more transparent regarding their environmental impacts: Allen AI provides holistic assessments of the OLMo model family (Groeneveld et al., 2024; Team OLMo et al., 2025), quantifying the impacts of development training runs for different model sizes (Morrison et al., 2025); and Mistral reports the results of a comprehensive life cycle assessment of two of its LLMs (Mistral AI, 2025), but without clear mention of the research and development stage.

Contrary to these reports, we analyze in detail all the compute that was used in the research and development of the Moshi speech-text foundation model.

2.3 Background on Moshi

Moshi (Défossez et al., 2024) is an open-source speech-text multimodal foundation model designed for natural, expressive, and real-time interaction. It was developed by Kyutai¹ over a 9-month period and is arguably the first end-to-end speech-to-speech model, making it a good case study for assessing the impacts of innovative research at an industrial scale, which go beyond the hyperparameter tuning or dataset validation common in well-established LLM development pipelines.

Kyutai kept and shared detailed logs of the 3,534 individual training runs necessary for the development of the most innovative parts of their model, enabling detailed analysis of research costs. They also provided us with the global development and final training cost for their more standard LLM backbone.

As illustrated in fig. 2, the development of Moshi relies on the following four modules:

- **LLM backbone:** Helium, a pure-text LLM trained from scratch and used to initialize Moshi’s main transformer.
- **Data generator:** a text-to-speech module used to create custom fine-tuning datasets.
- **Tokenizer:** Mimi, a neural audio codec that converts waveforms to speech tokens and back.
- **Main model:** a 7B-parameter transformer that consumes and produces tokenized speech.

¹<https://kyutai.org/>

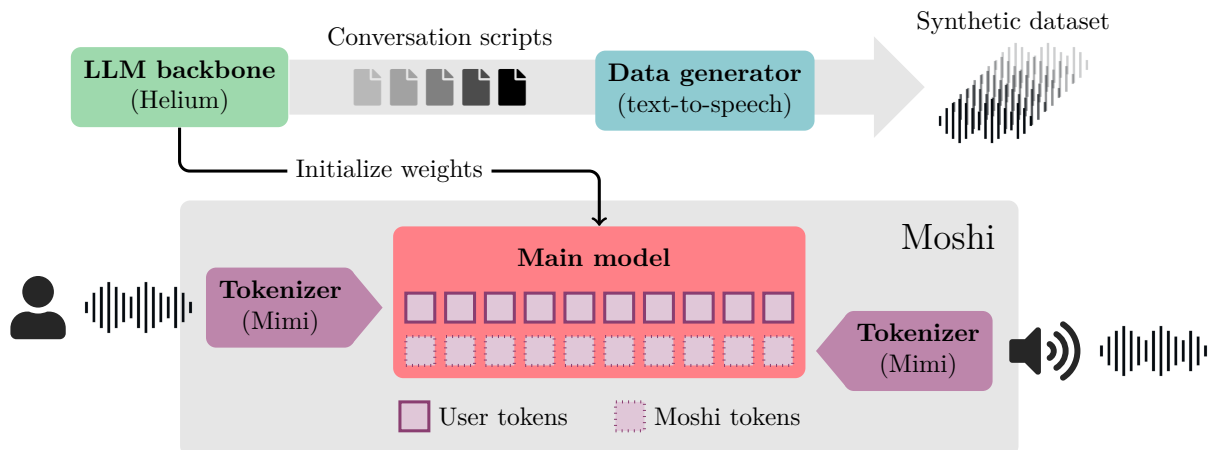


Figure 2: **Moshi modules.** Mimi ■ tokenizes input waveforms and feeds them to the main transformer model ■, whose predictions are converted back to waveform by Mimi. The transformer is initialized with the weights of the custom LLM Helium ■, and a data generator ■ converts synthetic conversation scripts into a fine-tuning speech dataset.

As all the runs took place on identical compute nodes on the Scaleway cloud-computing platform², we define the *compute* associated to a run as its duration multiplied by the number of GPUs, and use it as the main unit of comparison in our analysis. For example, a run executed on 8 GPUs for 10 hours has a compute intensity of 80 GPU-hours.

3 Compute Distribution Analysis

In this section, we analyze how compute is allocated across high-level objectives, including training, evaluation, hyperparameter search, ablations, and final model training (sec. 3.1). We then examine Moshi in detail, quantifying the compute cost of its modules and training phases (sec. 3.2). Finally, we characterize the distribution of compute by analyzing the intensity of the training runs across these categories (sec. 3.3).

3.1 Compute Distribution by Goal

In this section, we first quantify the compute spent on each run phase (training, validation, and evaluation), and then analyze how compute is distributed across research phases (from debugging to ablation studies).

Compute per Run Phase. Within a run—the training of a module with a fixed hyperparameter configuration—we distinguish three phases:

- **optimization** per se, where gradients are computed through backpropagation, and the parameters are updated,
- **validation**, where, periodically during training, inference is performed on a held-out subset of the data to monitor optimization progress and detect overfitting,
- **evaluation**, where, periodically during training, inference is performed on the test set. Model outputs are saved for human assessment, and more metrics are computed, potentially requiring the generation of large volumes of data and scoring with external models.

Fig. 3 shows the breakdown of compute along these three run phases. Core training accounts for 90% of the total compute and is by far the dominant computational driver, but validation and evaluation still account for 2.8% and 7.2% of the compute budget respectively, i.e., over 35 GPU-years, which is far from negligible.

²<https://www.scaleway.com/>

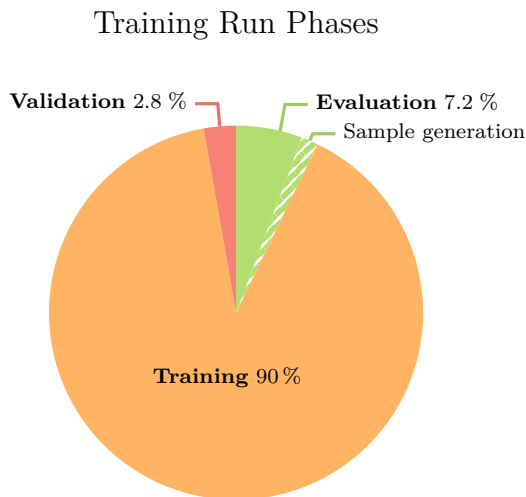


Figure 3: **Compute per run phase.** Runs are split into training, validation, and evaluation. We aggregate the compute for each phase across all runs, excluding LLM development.

Research Phases

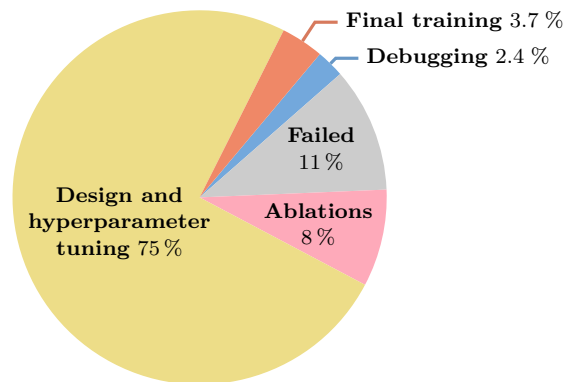


Figure 4: **Compute per research phase.** AI research includes design, tuning, and final training, but also debugging, failed runs and ablations.

Note that approximately one quarter of the evaluation compute is dedicated to sample generation for human assessment.

Compute per Research Phase. We distribute the compute cost across the main research phases that we identify from the logs:

- **debugging,**
- **failed runs,** corresponding to unusable training runs that yielded very low performance and were therefore not used for hyperparameter selection,
- **core development, model design, and hyperparameter tuning,** including experimental exploration of directions that were later discarded,
- **final model training,**
- **ablation studies and safety analyses,** performed for rigorously validating design choices, writing the scientific paper, and making the model ready to be released.

As shown in fig. 4, 75% of the compute is devoted to architecture design and hyperparameter tuning, whereas training the final models accounts for less than 4% of the total compute. The ablation studies and safety analyses reported in Moshi’s white paper alone represent 8% of the compute budget. Notably, debugging and failed runs together account for over 13% of the total compute, underscoring their substantial contribution to overall resource usage.

Key Takeaways

- **Periodic evaluation during training adds a significant overhead:** over 7% of the compute is spent on evaluating models in case the need for deeper analysis or human assessment arises, which calls into question the common practice of performing these evaluations regularly, and invites to consider using less expensive and less frequent performance tracking.
- **Final model training is a small part of the total research and development budget,** accounting for less than 4% of the total compute, emphasizing the importance of research costs.
- **Debugging and failed runs are costly:** together, they account for more than 13% of compute, suggesting that improved debugging practices and early-phase experiment diagnostics could significantly reduce overall resource consumption.
- **Ablation studies are costly:** while they are key to rigorous design choice validation and research publications, ablation studies represent 8% of the total computation budget, again questioning common research practices.

3.2 Compute Distribution by Module and Training Phase.

Training Phase Definition. In this section, we identify the training runs corresponding to each module of Moshi described in sec. 2.3. To better understand the research process, we also classify the runs of the most innovative parts of Moshi, namely the data generator and the main transformer model, into sequential training phases that could be applied to many AI development workflows:

- **Experimentation (Exp):** An exploratory phase used to test alternative architectures and functional modes; nothing is finalized at this point, and proof-of-concept models are evaluated through numerous runs, typically with moderate compute.
- **Pre-training (Pre):** Once the general pipeline and architecture type are fixed, models are trained on a large corpus. This phase typically involves fewer but substantially more compute-intensive runs. *Example:* Moshi learns speech representations from a dataset of 7M hours of audio, mostly containing English speech.
- **Post-training (Post):** The model’s weights are refined to handle specific input/output formats using datasets tailored for this purpose. *Example:* Moshi learns conversational turn-taking, i.e., when to speak and when to listen, from a dialogue dataset with separate audio tracks.
- **Fine-tuning (FT):** The model is adapted to a specific application by training on smaller, specialized datasets. *Example:* Moshi learns to behave like a useful conversational assistant from 20k hours of synthetic instruction data.

For the more standard modules, namely the tokenizer and the LLM backbone, which do not require exploration and which are trained in a single phase, we simply refer to their training as **Train**. Note that this does not mean that a single training occurs, since more standard development activities such as hyperparameter search are still needed.

We also keep the **Failed** run category described in the previous section, common to all modules and training phases, and corresponding to runs that produced under-performing models due to factors such as bugs, misconfigured hyperparameters, or inadequate architectural variants.

A visualization of the development costs for the different phases can be seen in fig. 1, and we analyze them below in more detail.

Project Timeline. In fig. 5, we visualize the number of runs and GPU in use (top) as well as the cumulative compute attributed to the different modules and training phases (bottom). It outlines that early phases of the project require a large number of short runs on few GPUs, and that the brunt of the computational load is taken by few expensive runs, especially long pre-training runs.

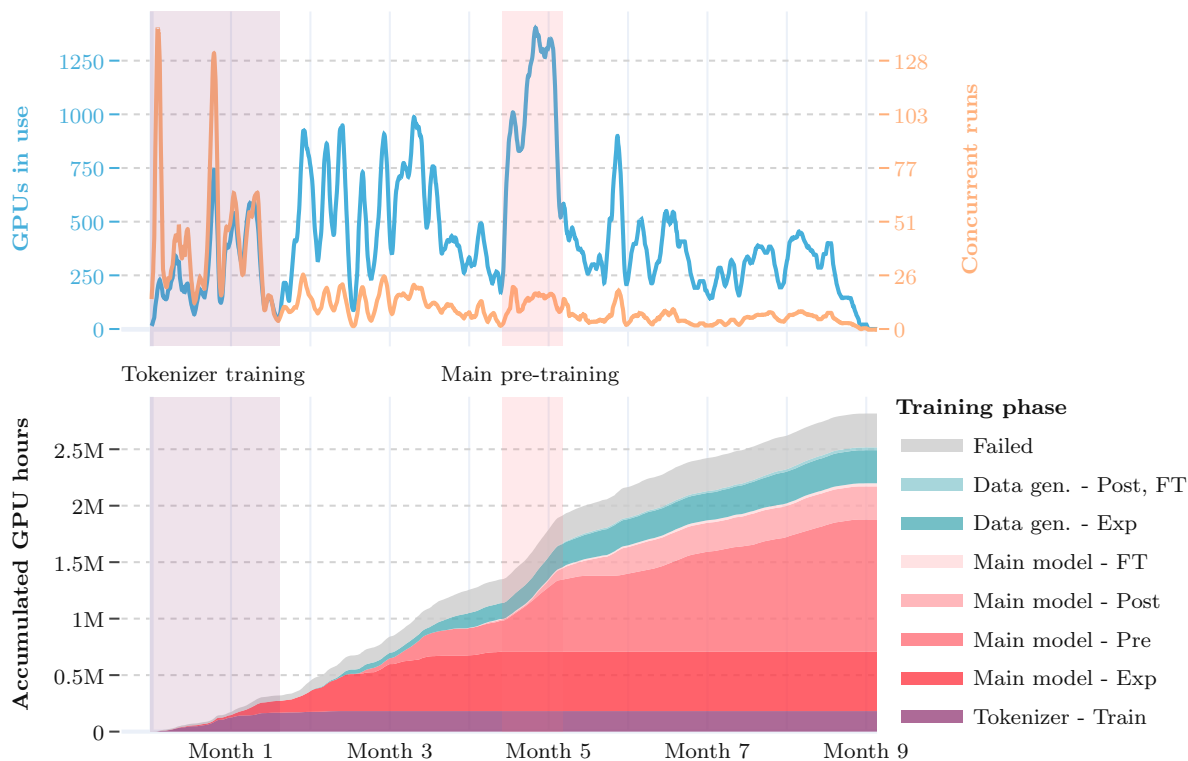


Figure 5: **Project compute intensity timeline.** *Top:* Number of GPUs in use (blue) and number of concurrent runs (orange) over the duration of the project. *Bottom:* Accumulated GPU hours per module and training phase: experimentation (Exp), pre-training (Pre), post-training (Post), and fine-tuning (FT). All quantities are sampled every 30 minutes, and smoothed with a sliding-window average over 100 steps. The plots do not include LLM backbone runs.

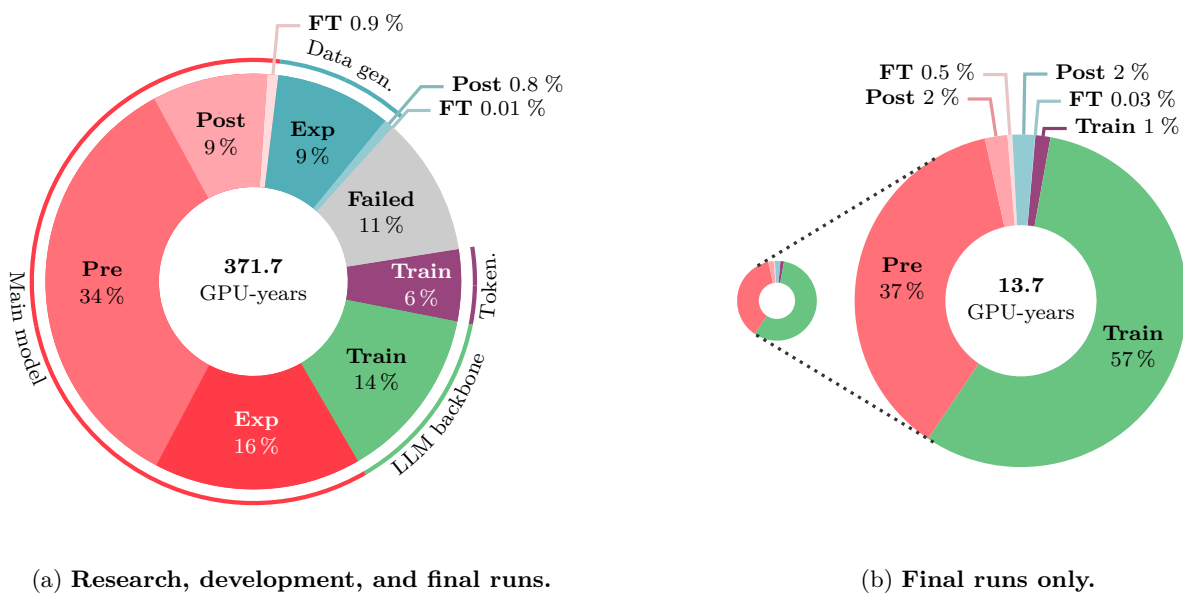


Figure 6: **Compute per training phase.** We distribute the compute among experimentation (Exp), pre-training (Pre), post-training (Post), and fine-tuning (FT) for each module. Fig. 6a considers all research and development runs plus the final training runs, and fig. 6b isolates the final runs. The area of each chart is proportional to the compute it represents.

Table 1: **Final vs. total compute.** We present the percentage of total research and development compute spent on the final training run, for each module and training phase independently: pre-training (Pre), post-training (Post), and fine-tuning (FT).

Module	Tokenizer	Main model			Data generator		LLM backbone
	Train	Pre	Post	FT	Post	FT	Train
Final-to-total compute ratio (%)	1.0	4.0	0.9	2.0	10.6	8.2	15.5

Distribution across Training Phases. In fig. 6a, we report the distribution of the research and development compute between the training phases. Researching, developing, and training the main model dominates the budget (60%), with pre-training alone contributing 34%, while fine-tuning accounts for less than 1%. In contrast, developing and training the LLM represents 14% of the overall compute. Early experimentation is responsible for a significant share of compute, accounting for 25% of the total.

In fig. 6b, we report the same breakdown restricted to the final runs only, i.e. a single run for each training phase whose resulting model is deployed in production. Under this setting, training the LLM and pre-training the main model account for 57% and 37% of the final compute respectively, largely outweighing all other modules and training phases. The notable increase in the share of LLM compute is explained by the (proportionally) lower development costs of the more standard LLM backbone compared to other modules. On the contrary, the main model requires almost 60% of the total compute, but only 37% of the final training compute. We believe that this difference is actually related to the novelty of the different modules, the development cost of more standard modules being limited to hyperparameter search, while the most novel require more research and exploration.

Relative Cost between Final Runs and Total Compute. To better analyze this effect, we report in tab. 1 the ratios of the final training cost to the total research, development, and training cost for each module and training phase. The main model and tokenizer have the lowest ratios corresponding to high research and development costs. The case of the data generator is interesting: this module is in fact a variant of the main model, with slightly modified post-training and fine-tuning schemes. After designing and training the main model, training the data generator was thus easier, which results in higher ratios, still however below the more standard LLM.

Key Takeaways

- **Modules and phases whose final costs are negligible might have significant research costs**, emphasizing the importance of not only analyzing the final training and considering all training phases.
- **Experimentation on novel elements has a significant cost, 25%.**
- **The final-to-total cost ratios differ significantly by module**, the most standard module (LLM) having a much higher ratio, emphasizing the importance of assessing research costs.

3.3 Compute Distribution by Run Compute Intensity

We analyze the different runs according to their compute intensity, i.e., their GPU-time. We group runs into eight compute intensity categories with thresholds at one GPU-hour, one GPU-day, one GPU-week, one GPU-month, one GPU-year, three GPU-years, and five GPU-years.

Compute Intensity of All Runs. In fig. 7, we report both the number of runs per intensity category and their contribution to total compute. We observe a pronounced 90/10 effect, with a frontier around one GPU-month: 13% of runs account for 89% of the total compute. Low-intensity runs (below one GPU-day) represent 42% of all runs but contribute negligibly to total compute, as they are primarily associated with

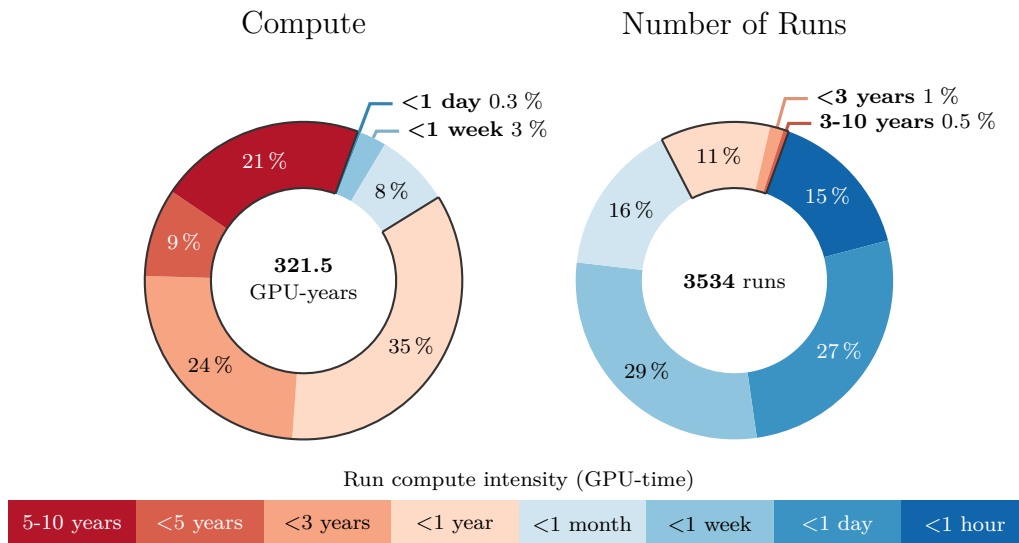


Figure 7: **Run compute intensity categories.** Distribution of runs across compute-intensity categories, excluding LLM runs. The highlighted sectors correspond to 90% of the compute concentrated in 10% of the runs.

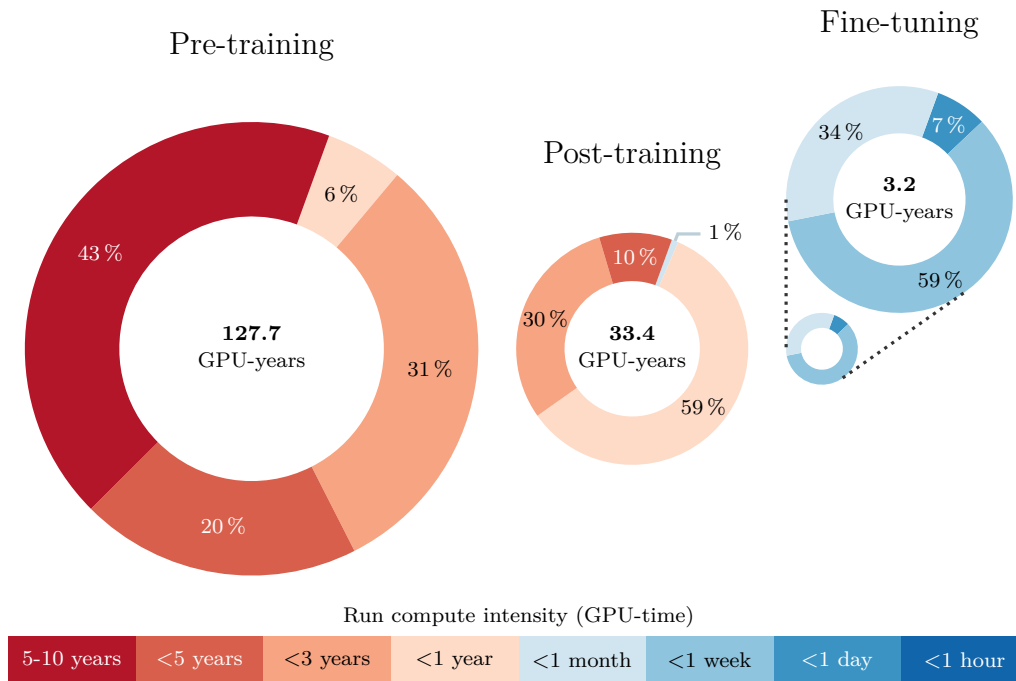


Figure 8: **Run compute intensity by training phase.** Compute-intensity distribution for pre-training, post-training, and fine-tuning runs of the main model. The area of each chart is proportional to the compute of the corresponding phase.

debugging, fine-tuning, and tokenizer training. At the opposite end of the spectrum, only 19 runs (0.5%) with intensities exceeding three GPU-years account for 30% of the total compute.

Compute Intensity of Failed and Ablation Runs. As observed in fig. 4, failed runs and ablation studies each contribute around one tenth of the development compute. However, failed runs are much

less compute-intensive than ablation studies: out of 1,479 failed runs, only five have an intensity over one GPU-year; whereas ablation studies only comprise 139 runs, twelve of them over the one GPU-year mark.

Run Compute Intensity by Training Phase. We focus on runs from the pre-training, post-training, and fine-tuning phases of the main model and analyze their compute intensity (fig. 8). Pre-training concentrates the most compute-intensive runs and is the only phase containing runs exceeding five GPU-years, which make up 43% of the pre-training compute. In contrast, no fine-tuning run exceeds one GPU-month; instead, 66% of the fine-tuning compute corresponds to runs lasting under one GPU-week. Post-training occupies an intermediate regime, with 60% of its compute coming from runs under one GPU-year.

Key Takeaways

- **A small fraction of runs dominates compute usage:** 13% of runs account for nearly 90% of total compute.
- **Failed runs and ablation studies have distinct compute-intensity profiles,** with failed runs being less intensive, but ten times more frequent.
- **Pre-training concentrates the most compute-intensive runs,** including all individual runs exceeding five GPU-years.

4 Environmental Assessment

In this section, we quantify the environmental impact of developing the Moshi model from the first experiments up to the last training run.

Cluster Configuration. All training runs took place on the Scaleway Nabuchodonosor supercomputer (Scaleway, 2024), an NVIDIA DGX SuperPOD (NVIDIA, 2025b) made up of 127 NVIDIA DGX H100 (NVIDIA, 2025a) nodes and located in Paris ³.

Impact Indicators. We estimate operational and embodied impacts across the following environmental impact indicators:

- **Primary energy (PE)** measures the consumption of renewable and non-renewable energy resources extracted from nature, expressed in megajoules (MJ) (Boavizta, 2023).
- **Global warming potential (GWP)** quantifies the contribution of greenhouse gas emissions to climate change, expressed in kilograms of carbon dioxide equivalent (kgCO₂eq) (Boavizta, 2023).
- **Water consumption (WC)** measures the volume of water used and not returned to its original source (through evaporation, incorporation into products, or migration), expressed in liters (L) (Li et al., 2025a).
- **Abiotic resource depletion (ADP)** quantifies the depletion of non-renewable mineral and metal (ADPe) and fossil (ADPf) resources, expressed in kilograms of antimony equivalent (kgSbeq) (Boavizta, 2023).

Scope. The object of our assessment (or *functional unit*) covers the complete research and development process of the Moshi model, from early experiments to final trainings and ablations. We exclude data acquisition, processing, and storage due to a lack of detailed information, and we do not account for the environmental costs associated with deployment and inference after public release.

As summarized in tab. 2, we omit end-of-life impacts because of the general lack of reliable data (Ficher et al., 2025; Baldé et al., 2024). We further restrict water consumption estimates to the use phase only. Although studies on the water consumption of hardware manufacturing exist (Falk et al., 2025; Boyd, 2012),

³Kyutai rented additional nodes in an unspecified location during a period of three months, but we omit this fact due to a lack of detailed data.

Table 2: **Life cycle assessment scope.** Impact indicators considered for each life cycle phase of the hardware. We group raw material extraction and manufacturing into a single *production* phase. (✓) means that the impacts are only partially accounted for (Simon et al., 2025).

Impact indicator	Life cycle phase			
	Production	Transport	Use	End of life
Primary energy (PE)	✓	(✓)	✓	
Global warming potential (GWP)	✓	(✓)	✓	
Water consumption (WC)			✓	
Abiotic depletion potential (ADP)	✓	✓	✓	

extrapolating these results to individual hardware components would require additional assumptions and is therefore outside the scope of this work.

4.1 Methodology

We provide here a high-level description of our methodology for estimating the environmental impacts of developing and training Moshi, starting from measured compute usage. We refer the reader to sec. B for the full set of equations and parameter values. We distinguish between *operational impacts*, which arise from the use phase of the hardware during model training (i.e., electricity consumption while compute nodes are operating), and *embodied impacts*, which correspond to impacts incurred during hardware production, transport, and end of life.

Operational Impacts. GPU energy consumption is estimated as the product of the number of concurrently used GPUs, the maximum rated GPU power, and a 95% utilization factor, which accounts for brief periods of non-GPU-intensive work. Based on Kyutai’s observations, we assume a CPU utilization of 5%. Following prior analyses of similar compute nodes (Spetko et al., 2020), we assume the power consumption of RAM and other node components (including fans, SSDs, network cards, and the motherboard) to be constant.

We account for energy overheads associated with storage, management, and communication at the SuperPOD level, as well as datacenter infrastructure overheads. We do not include the consumption of idle compute nodes, assuming that nodes not used for developing Moshi were allocated to other workloads by Scaleway.

To compute primary energy impacts, we follow the methodology of Boavizta (2026) and multiply energy consumption by the consumption of fossil fuel resources per kilowatt-hour, and adding an overhead to also account for renewable energy sources. We similarly multiply by an abiotic depletion factor per kilowatt-hour to obtain use phase resource depletion impacts⁴.

Following prior work (Luccioni et al., 2023; Lannelongue et al., 2021), we estimate greenhouse gas emissions as the product of total energy consumption and the yearly average carbon intensity of electricity at the location of computation. Finally, we estimate water consumption using the methodology proposed by Li et al. (2025a).

Embodied Impacts. We estimate embodied impacts for GPUs, CPUs, RAM, SSDs, power supplies, motherboards, and cases, as well as for the assembly of the compute nodes. Production and transport impacts of individual hardware components are estimated using per-component impact factors provided by Boavizta (Simon et al., 2025). For GPU production and transport impacts, we refer to a recent report by ADEME (Lees-Perasso et al., 2026)⁵. These embodied impacts are then allocated proportionally to

⁴Like Boavizta, we only contemplate mineral and metal resource depletion (ADPe) in the use phase, excluding fossil resource depletion (ADPf).

⁵For GPU abiotic depletion potential (ADP) impacts, we only report mineral and metal depletion (ADPe). For the remaining hardware, ADP also includes fossil resource depletion (ADPe + ADPf).

the duration of hardware use during Moshi’s development relative to its typical service lifetime, following established practice in prior work (Luccioni et al., 2023; Morand et al., 2024; Falk et al., 2025).

4.2 Analysis

We first report the total environmental impacts across the four indicators under study, comparing them, when possible, to yearly per-capita impacts. We then disaggregate each indicator by hardware component and by scope (computation, datacenter overheads, and embodied impacts), before focusing specifically on hardware production. Finally, we conclude with a simulation of how the environmental impacts of model development would vary across different geographic locations. Detailed numerical results are reported in sec. A.

Environmental Impacts of Research. The research, development, and training compute of the model ascended to 3M GPU-hours, or **372 GPU-years**. This translates into:

- An energy consumption of **5 gigawatt-hours**, equivalent to the yearly consumption of 727 people in France (RTE, 2025) and resulting in **68 terajoules** of primary energy extracted from the environment.
- A global warming potential of **319 tonnes of carbon dioxide equivalent**, that of 39 people in a year in France (Baude & Larrieu, 2025), or of 132 round trip flights between Paris and San Francisco (Sustainable Travel International, 2024).
- A water consumption of **19 megaliters**, that of 342 people in a year in France (OECD, 2025).
- An abiotic depletion potential of **8 kilograms of antimony equivalent**, that of 6,566 smartphones (Sánchez et al., 2024) or 483 laptops (Baur et al., 2023).

Fig. 9 breaks down the total environmental impacts by hardware component and by impact scope: *computation* ■ due to compute node operation; *datacenter* ■ due to datacenter management, cooling, ventilation, and other overheads; and *embodied* ■ due to hardware production. For water consumption (fig. 9c), operational impacts are split between *datacenter cooling* ■ and power-plant cooling for *electricity production* ■.

When it comes to primary energy (fig. 9a) and global warming potential (fig. 9b), operational impacts due to computation ■ and datacenter overheads ■ are directly proportional. In contrast, embodied ■ impacts make up a larger share of the total global warming potential, whereas their contribution to total primary energy is small. This difference is very pronounced in the case of RAM. These larger embodied global warming potential impacts are mainly explained by the emission of fluorinated gases and wet chemicals during the manufacturing process (Hess & Nowicka, 2026).

Focusing on water consumption (fig. 9c), datacenter cooling ■ consumes a negligible amount of water in comparison to cooling the power plants that generate the electricity to power the datacenter ■. We do not estimate embodied water consumption, but it should be considerable due to the ultrapure water and the electricity required for semiconductor manufacturing (Boyd, 2012; Li et al., 2025a).

GPUs are responsible for the majority of the impact across all indicators except abiotic depletion potential (fig. 9d), where node components other than CPUs and GPUs require the most material resources by a large margin. Notably, these other node components, including fans and network cards, contribute significantly to primary energy (fig. 9a) and global warming potential (fig. 9b), yet they are seldom taken into account in related work. Conversely, CPUs have the lowest impacts overall, which is explained by their low utilization during training.

Embodied Impact Details. We now break down embodied impacts for each impact category across the hardware components of a compute node. As shown in fig. 10, GPUs contribute the most to embodied primary energy (PE) and global warming potential (GWP), closely followed by RAM. The impacts of producing one RAM module are around five times lower than those of a GPU, but each compute node contains thirty-two RAM modules as opposed to only eight GPUs. The distribution changes for abiotic depletion

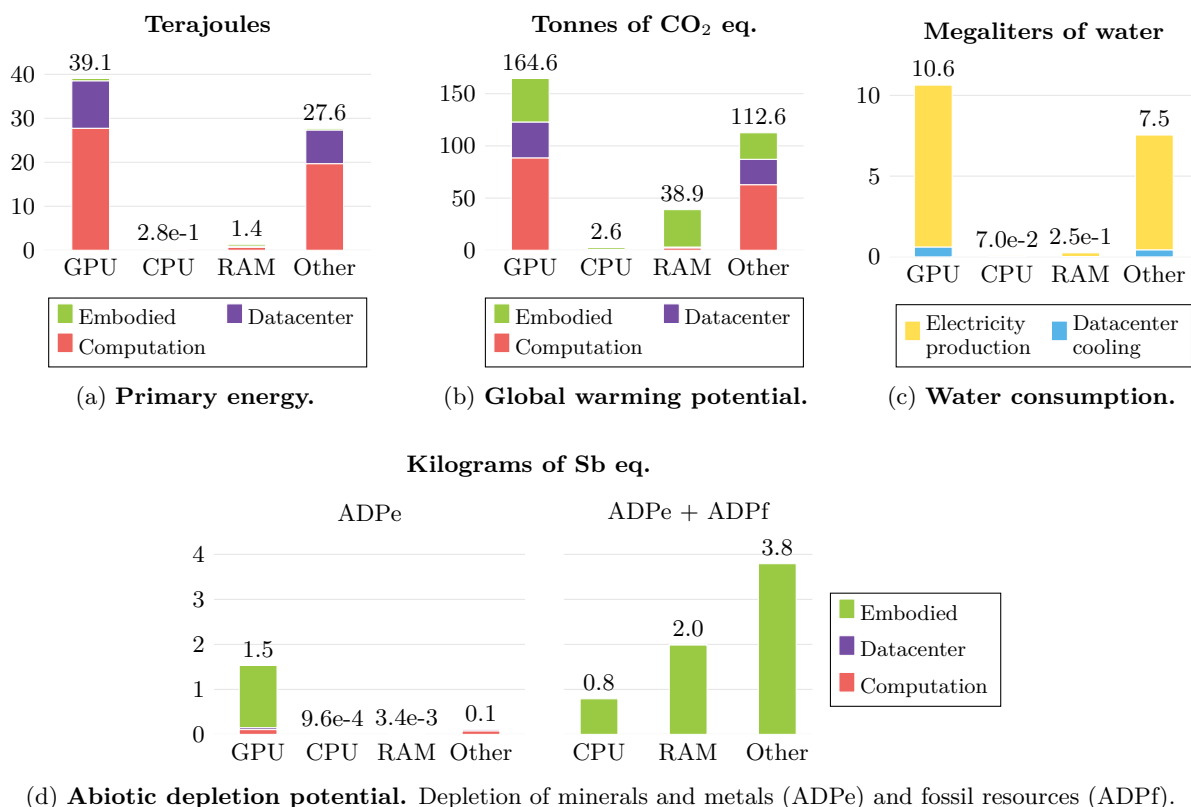


Figure 9: **Environmental impacts of research.** Each impact indicator (primary energy, global warming potential, water consumption, abiotic depletion potential) is disaggregated by hardware component (GPU, CPU, RAM, Other), and by scope.

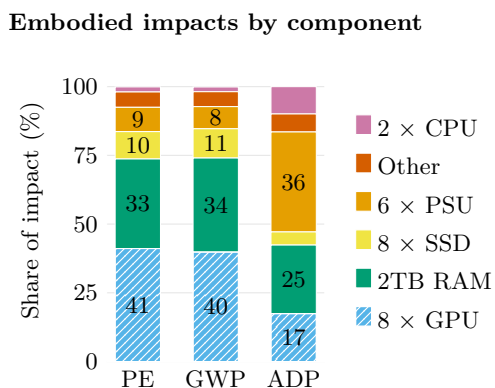


Figure 10: **Embodied impacts by component.** Share of embodied primary energy (PE), global warming potential (GWP), and abiotic depletion potential (ADP) for each hardware component in a node. Solid fill impacts are estimated using Boavizta (Simon et al., 2025) and include ADPe + ADPf; dashed impacts come from ADEME (Lees-Perasso et al., 2026) and include ADPe only.

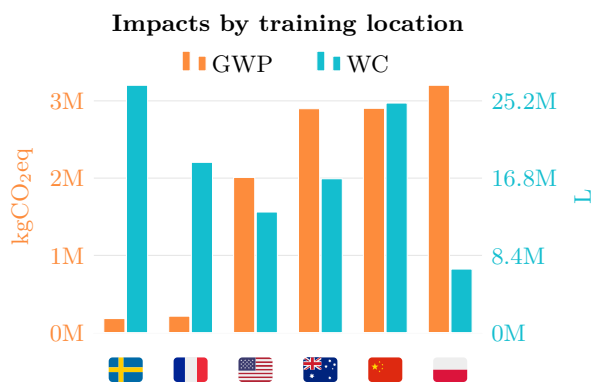


Figure 11: **Operational impacts by training location.** Hypothetical global warming potential (GWP) and water consumption (WC) of developing the model in different locations, excluding embodied impacts.

potential (ADP), with power supplies being responsible for most of the resource depletion, followed by RAM modules and GPUs. The impacts of producing a single power supply are the highest among all components, and there are six power supplies in each compute node. We provide details on the per-unit impacts of each hardware component in tab. 5 and fig. 12 (sec. A).

Importance of Location. Factors such as carbon intensity and water consumption for power plant cooling vary by location, depending on the fuel mix powering the electricity grid. We explore the hypothetical operational impacts of developing Moshi in different locations in fig. 11. Overall, the figure shows no correlation between global warming potential and water consumption impacts; however, there is a marked trade-off between both impacts in Sweden, France, and Poland.

Carbon intensity in Sweden and France is low thanks to their reliance on hydroelectric and nuclear energy, as opposed to fossil fuels. However, hydropower and nuclear are among the most water-intensive energy sources (Reig et al., 2020). Conversely, Poland has one of the most carbon-intensive electricity grids worldwide (Electricity Maps, 2026), yet its water consumption is low: the Polish grid relies extensively on coal and gas, with almost no hydropower or nuclear power plants (EMBER, 2025b).

5 Discussion

We conclude this study by discussing our main results, how they compare to those available in the literature, and possible mitigation strategies for the growing impact of Gen-AI research.

Research vs. Development Computational Costs. The final training of Moshi represents less than 4% of the total compute and energy consumption. This is much lower than the numbers reported in the literature. For example, according to Luccioni et al. (2023), training the final BLOOM model represented 31% of the compute (considering only a part of the project, the one performed the cluster used for the final training), and Morrison et al. (2025) report that training the open source versions of the OLMo suite of LLMs accounted for 50% of the environmental impacts of the project. We believe that this is due to three main effects. First, while developed by experienced researchers, Moshi was built completely from scratch, following the creation of Kyutai, which we believe enables us to better account for the exploratory research phase than previous studies. Indeed, we argue that the boundary of a specific development project in institutes that have already performed LLM research is more blurry, as previous or related projects are likely leveraged. Second, we believe that this large difference is also due to the originality of the speech-to-speech Moshi model. While large language models benefit from years of empirical optimization—such as scaling laws (Hoffmann et al., 2022; Bhagia et al., 2025), which enable controlled experimentation on smaller proxy models, speech-to-speech modeling remains comparatively underexplored. As a result, a larger portion of the development process must be carried out at scale, increasing the relative cost of experimentation. This is reflected in our results: when restricting the analysis to the LLM backbone alone, the share of compute attributed to final training rises to 15%, bringing it closer to previously reported values. Third, a substantial fraction of the total compute is spent on experimentation, debugging, and other development stages that are often not accounted for in prior work.

Reducing Unnecessary Compute. Failed experiments make up 11% of Moshi’s total compute. Most of these training runs were quickly canceled, but their contribution is related to their large number: 42% of the runs were discarded due to poor hyperparameter combinations, mistaken configurations, or bugs. While failed experiments are unavoidable in research, these figures invite to take special care when launching compute-intensive experiments, e.g. above one GPU-month, and monitor them closely.

Debug runs, which are also inexpensive individually, still make up 2.4% of the research and development compute, almost as much as training the final model. Debugging is naturally necessary, but these figures invite to perform it as much as possible on infrastructures with a low power consumption instead of a production environment, potentially using downscaled versions of the models and datasets.

Periodic evaluation and validation during model training account for 10% of the total compute. We believe that this important cost should be taken into account to modify standard practices, performing evaluation and validation at a lower frequency, and on smaller datasets.

Questioning Research Practices and Expectations. Ablation studies are often at the core of a machine learning research article, validating the findings and claims. However, these ablation studies have an important cost: 8% of the compute of Moshi being spent on ablation studies and safety analyses, mostly carried out while the final model was already trained. This relatively large share stems mainly from ablations on pre-training. We believe that this should lead to more questioning of the necessity and practices of ablation studies. For example, comparisons could be made at early stages of training, or fine-tuning a given model for a short time, or systematically on smaller versions of the models and datasets.

In a similar vein, we argue that given the strong environmental impact of AI research, the computational budget should be more systematically included in the evaluation, for example, evaluating expected performance as a function of the computational budget, including hyperparameter search (Dodge et al., 2019). Beyond reducing environmental impact, we believe that such considerations would also make comparisons more meaningful by factoring out the important impact of compute scaling on results (Mertens et al., 2026).

Reducing Environmental Impacts for a Given Compute Budget. The location of the computational resources has a significant effect on operational impacts. A natural direction to reduce AI research impact is thus to select data centers based on their power and water usage, and taking into account the carbon and water intensities of the local electrical grid. However, low-carbon grids might consume large amounts of water for power plant cooling, which is problematic in regions under water stress, and adapting the computational load to local resources demand is still a rare practice.

The largest share of embodied impacts stems from GPU and RAM manufacturing, as well as power supply production in the case of resource depletion. The impacts of GPU manufacturing are considerable due to both their high per-unit impact and the amount of GPUs required for training. Although producing a single RAM module has a lower impact, compute nodes designed for AI training may easily contain up to thirty or sixty of these modules. These observations should serve as an incentive to boost research into smaller models and training schemes with low memory footprints, as well as to extend the lifetimes of GPUs by reducing compute usage.

Measurements and Publicity. As a first step toward these evolutions and to better question the impact of GenAI, we argue that measuring and publicly reporting the computational and environmental costs of not only the final training, but also development, and even complete research projects, with breakdown per project stage, should become common practice. For operational impacts, tools such as CodeCarbon⁶ are both easily accessible and accurate. Similar tools are also available for embodied impacts, for example, Boavizta (Simon et al., 2025) or MLCA (Morand et al., 2024), whose methodology we outline in sec. B.2.

⁶<https://codecarbon.io>

References

- ADEME. Base IMPACTS®. <https://base-empreinte.ademe.fr/donnees/download-data>, 2023. Accessed: 2026-03-16.
- Cornelis P Baldé, Ruediger Kuehr, Tales Yamamoto, Rosie McDonald, Elena D’Angelo, Shahana Althaf, Garam Bel, Otmar Deubzer, Elena Fernandez-Cubillo, Vanessa Forti, et al. Global e-waste monitor 2024. <https://ewastemonitor.info/the-global-e-waste-monitor-2024/>, 2024. Accessed: 2025-11-25.
- Manuel Baude and Sylvain Larrieu. L’empreinte carbone de la France de 1990 à 2024. <https://www.statistiques.developpement-durable.gouv.fr/lempreinte-carbone-de-la-france-de-1990-2024>, 2025. Accessed: 2025-11-18.
- Sarah-Jane Baur, Marina Proske, and Erik Poppe. Life cycle assessment of the Framework Laptop 2022. LCA report (ISO 14044 and ISO 14067). Technical report, Fraunhofer IZM, Berlin, 2023.
- Adrien Berthelot, Eddy Caron, Mathilde Jay, and Laurent Lefèvre. Estimating the environmental impact of Generative-AI services using an LCA-based methodology. *Procedia CIRP*, 122:707–712, 2024. ISSN 2212-8271. doi: 10.1016/j.procir.2024.01.098. 31st CIRP Conference on Life Cycle Engineering.
- Akshita Bhagia, Jiacheng Liu, Alexander Wettig, David Heineman, Oyvind Tafjord, Ananya Harsh Jha, Luca Soldaini, Noah A. Smith, Dirk Groeneveld, Pang Wei Koh, Jesse Dodge, and Hannaneh Hajishirzi. Establishing task scaling laws via compute-efficient model ladders, 2025.
- Boavizta. Impacts criteria - Boavizta API documentation. <https://doc.api.boavizta.org/Explanations/impacts/>, 2023. Accessed: 2025-12-19.
- Boavizta. Electrical impact factors - Boavizta API documentation. <https://doc.api.boavizta.org/Explanations/impacts/>, 2026. Accessed: 2026-03-16.
- Sarah B. Boyd. *Life-Cycle Assessment of Semiconductors*. Springer, New York, NY, 2012. ISBN 978-1-4419-9987-0 978-1-4419-9988-7. doi: 10.1007/978-1-4419-9988-7.
- Andrew A Chien, Liuzixuan Lin, Hai Nguyen, Varsha Rao, Tristan Sharma, and Rajini Wijayawardana. Reducing the carbon impact of generative AI inference (today and in 2035). In *Proceedings of the 2nd Workshop on Sustainable Computer Systems*, pp. 1–7, Boston MA USA, 2023. ACM. ISBN 979-8-4007-0242-6. doi: 10.1145/3604930.3605705.
- Jeongdong Choe. Micron B47R 3D CTF CuA NAND Die, world’s first 176L (195T). <https://semitechnology.com/micron-b47r-3d-ctf-cua-nand-die-worlds-first-176l-195t>, 2021. Accessed: 2025-11-13.
- Jeongdong Choe. Industry-leading DDR5 technology: Micron vs. Samsung vs. SK Hynix. <https://www.techinsights.com/blog/industry-leading-ddr5-technology>, 2022. Accessed: 2025-11-13.
- CodeCarbon. Power usage methodology — CodeCarbon 3.2.1 documentation. <https://mlco2.github.io/codecarbon/methodology.html#power-usage>, 2026. Accessed: 2026-01-20.
- Alexandre Défossez, Laurent Mazaré, Manu Orsini, Amélie Royer, Patrick Pérez, Hervé Jégou, Edouard Grave, and Neil Zeghidour. Moshi: A speech-text foundation model for real-time dialogue, 2024.
- Clément Desroches, Martin Chauvin, Louis Ladan, Caroline Vateau, Simon Gosset, and Philippe Cordier. Exploring the sustainable scaling of AI dilemma: A projective study of corporations’ AI environmental impacts, 2025.
- Jesse Dodge, Suchin Gururangan, Dallas Card, Roy Schwartz, and Noah A. Smith. Show your work: Improved reporting of experimental results. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2185–2194, Hong Kong, China, 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1224.

- Electricity Maps. Real-time carbon intensity map. https://app.electricitymaps.com/map/live/fifteen_minutes, 2026. Accessed: 2026-01-19.
- Cooper Elsworth, Keguo Huang, David Patterson, Ian Schneider, Robert Sedivy, Savannah Goodman, Ben Townsend, Parthasarathy Ranganathan, Jeff Dean, Amin Vahdat, Ben Gomes, and James Manyika. Measuring the environmental impact of delivering AI at Google Scale, 2025.
- EMBER. Electricity data explorer - Carbon intensity in 2024. https://ember-energy.org/data/electricity-data-explorer/?data=co2_intensity&fuel=total&chart=single_year, 2025a. Accessed: 2025-11-10.
- EMBER. Electricity data explorer - Electricity generation in 2024. https://ember-energy.org/data/electricity-data-explorer/?metric=pct_share&chart=single_year, 2025b. Accessed: 2025-11-19.
- EMBER. Electricity data explorer - renewable electricity generation in france in 2024. https://ember-energy.org/data/electricity-data-explorer/?data=generation&fuel=res&chart=single_year&entity=France&tab=main&metric=pct_share&date=2024-01-01, 2026. Accessed: 2026-03-16.
- Sophia Falk, David Ekchajzer, Thibault Pirson, Etienne Lees-Perasso, Augustin Wattiez, Lisa Biber-Freudenberger, Sasha Luccioni, and Aimee van Wynsberghe. More than carbon: Cradle-to-grave environmental impacts of GenAI training on the Nvidia A100 GPU. <http://arxiv.org/abs/2509.00093>, 2025.
- Marion Ficher, Tom Bauer, and Anne-Laure Ligozat. A comprehensive review of the end-of-life modeling in LCAs of digital equipment. *The International Journal of Life Cycle Assessment*, 30(1):20–42, 2025. ISSN 1614-7502. doi: 10.1007/s11367-024-02367-x.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikula, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L. Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu-hui Chen, Zafarali Ahmed, Zhitaogong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. Gemma: Open models based on Gemini research and technology, 2024.
- Dirk Groeneveld, Iz Beltagy, Evan Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, William Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah Smith, and Hannaneh Hajishirzi. OLMo: Accelerating the science of language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15789–15809, Bangkok, Thailand, 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.841.

- Udit Gupta, Young Geun Kim, Sylvia Lee, Jordan Tse, Hsien-Hsin S. Lee, Gu-Yeon Wei, David Brooks, and Carole-Jean Wu. Chasing carbon: The elusive environmental footprint of computing. *IEEE Micro*, 42(4): 37–47, 2022. ISSN 0272-1732. doi: 10.1109/MM.2022.3163226.
- Michael Z. Hauschild, Ralph K. Rosenbaum, and Stig Irving Olsen (eds.). *Life Cycle Assessment: Theory and Practice*. Springer International Publishing, Cham, 2018. ISBN 978-3-319-56474-6 978-3-319-56475-3. doi: 10.1007/978-3-319-56475-3.
- Reinout Heijungs and Sangwon Suh. *The Computational Structure of Life Cycle Assessment*, volume 11 of *Eco-Efficiency in Industry and Science*. Springer Netherlands, Dordrecht, 2002. ISBN 978-90-481-6041-9 978-94-015-9900-9. doi: 10.1007/978-94-015-9900-9.
- Julia Christina Hess and Maria Nowicka. Direct emissions in semiconductor manufacturing are increasing again – what is behind the shift? <https://www.interface-eu.org/publications/semiconductor-emissions-data-2026>, January 2026. Accessed: 2026-03-23.
- Lorenz M. Hilty and Magda David Hercheui. ICT and sustainable development. In Jacques Berleur, Magda David Hercheui, and Lorenz M. Hilty (eds.), *What Kind of Information Society? Governance, Virtuality, Surveillance, Sustainability, Resilience*, volume 328, pp. 227–235. Springer Berlin Heidelberg, Berlin, Heidelberg, 2010. ISBN 978-3-642-15478-2 978-3-642-15479-9.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Oriol Vinyals, Jack W. Rae, and Laurent Sifre. Training compute-optimal large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22*, pp. 30016–30030, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 978-1-7138-7108-8.
- Nathaniel C Horner, Arman Shehabi, and Inês L Azevedo. Known unknowns: Indirect energy effects of information and communication technology. *Environmental Research Letters*, 11(10):103001, 2016. ISSN 1748-9326. doi: 10.1088/1748-9326/11/10/103001.
- ITU. Recommendation ITU-T L.1410 (11/2024) - Methodology for environmental life cycle assessments of information and communication technology goods, networks and services. <https://www.itu.int/epublications/publication/itu-t-l-1410-2024-11-methodology-for-environmental-life-cycle-assessments-of-information-and-communication-technology-goods-networks-and-services>, 2024.
- Nidhal Jegham, Marwan Abdelatti, Lassad Elmoubarki, and Abdeltawab Hendawi. How hungry is AI? Benchmarking energy, water, and carbon footprint of LLM inference, 2025.
- Yi Jin, Paul Behrens, Arnold Tukker, and Laura Scherer. Water use of electricity technologies: A global meta-analysis. *Renewable and Sustainable Energy Reviews*, 115:109391, 2019. ISSN 1364-0321. doi: 10.1016/j.rser.2019.109391.
- Lynn H. Kaack, Priya L. Donti, Emma Strubell, George Kamiya, Felix Creutzig, and David Rolnick. Aligning artificial intelligence with climate change mitigation. *Nature Climate Change*, 12(6):518–527, 2022. ISSN 1758-6798. doi: 10.1038/s41558-022-01377-7.
- Imad Lakim, Ebtesam Almazrouei, Ibrahim Abualhaol, Merouane Debbah, and Julien Launay. A holistic assessment of the carbon footprint of Noor, a very large Arabic language model. In Angela Fan, Suzana Ilic, Thomas Wolf, and Matthias Gallé (eds.), *Proceedings of BigScience Episode #5 – Workshop on Challenges & Perspectives in Creating Large Language Models*, pp. 84–94, virtual+Dublin, 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.bigscience-1.8.
- Loïc Lannelongue, Jason Grealey, and Michael Inouye. Green Algorithms: Quantifying the carbon footprint of computation. *Advanced Science*, 8(12):2100707, 2021. ISSN 2198-3844. doi: 10.1002/advs.202100707.

- Teven Le Scao, Thomas Wang, Daniel Hesslow, Stas Bekman, M Saiful Bari, Stella Biderman, Hady Elsahar, Niklas Muennighoff, Jason Phang, Ofir Press, Colin Raffel, Victor Sanh, Sheng Shen, Lintang Sutawika, Jaesung Tae, Zheng Xin Yong, Julien Launay, and Iz Beltagy. What language model to train if you have one million GPU hours? In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 765–782, Abu Dhabi, United Arab Emirates, 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-emnlp.54.
- Etienne Lees-Perasso, David Ekchajzer, Gauthier Roussilhe, and Thomas De Latour. Analyses de cycle de vie de gpu pour l’intelligence artificielle. Technical report, ADEME, 2026.
- Pengfei Li, Jianyi Yang, Mohammad A. Islam, and Shaolei Ren. Making AI less ‘thirsty’. *Communications of the ACM*, 68(7):54–61, 2025a. ISSN 0001-0782. doi: 10.1145/3724499.
- Yueying Li, Zhanqiu Hu, Esha Choukse, Rodrigo Fonseca, G. Edward Suh, and Udit Gupta. EcoServe: Designing carbon-aware AI inference systems. <https://arxiv.org/abs/2502.05043>, 2025b.
- Anne-Laure Ligozat, Julien Lefevre, Aurélie Bugeau, and Jacques Combaz. Unraveling the hidden environmental impacts of AI solutions for environment life cycle assessment of AI solutions. *Sustainability*, 14(9): 5172, 2022. ISSN 2071-1050. doi: 10.3390/su14095172.
- Alexandra Sasha Luccioni, Sylvain Viguier, and Anne-Laure Ligozat. Estimating the carbon footprint of BLOOM, a 176B parameter language model. *Journal of Machine Learning Research*, 24(253):1–15, 2023. ISSN 1533-7928.
- María Castrillo Melguizo, Jaime Iglesias Blanco, and Álvaro López García. Wattnet: Matching electricity consumption with low-carbon, low-water footprint energy supply, 2026.
- Matthias Mertens, Natalia Fischl-Lanzoni, and Neil Thompson. Is there “secret sauce” in large language model development?, 2026.
- Meta. Hugging face - meta-llama/Meta-Llama-3-70B. <https://huggingface.co/meta-llama/Meta-Llama-3-70B>, 2024. Accessed: 2025-11-21.
- Mistral AI. Our contribution to a global environmental standard for AI. <https://mistral.ai/news/our-contribution-to-a-global-environmental-standard-for-ai>, 2025. Accessed: 2025-12-04.
- Ki-Il Moon, Ho-Young Son, and Kangwook Lee. Advanced packaging technologies in memory applications for future generative AI era. In *2023 International Electron Devices Meeting (IEDM)*, pp. 1–4, San Francisco, CA, USA, 2023. IEEE. ISBN 979-8-3503-2767-0. doi: 10.1109/IEDM45741.2023.10413890.
- Clément Morand, Anne-Laure Ligozat, and Aurélie Névéol. MLCA: A tool for machine learning life cycle assessment. In *2024 10th International Conference on ICT for Sustainability (ICT4S)*, pp. 227–238. IEEE, 2024. doi: 10.1109/ICT4S64576.2024.00031.
- Jacob Morrison, Clara Na, Jared Fernandez, Tim Dettmers, Emma Strubell, and Jesse Dodge. Holistically evaluating the environmental impact of creating language models, 2025.
- NVIDIA. Introduction to NVIDIA DGX H100/H200 systems - NVIDIA DGX H100/H200 user guide. <https://docs.nvidia.com/dgx/dgxh100-user-guide/introduction-to-dgxh100.html>, 2025a. Accessed: 2025-11-10.
- NVIDIA. NVIDIA DGX SuperPOD: Next generation scalable infrastructure for AI leadership reference architecture featuring NVIDIA DGX H100. <https://docs.nvidia.com/dgx-superpod/reference-architecture-scalable-infrastructure-h100/latest/abstract.html>, 2025b. Accessed: 2025-11-17.
- NVIDIA. Planning a data center deployment - NVIDIA DGX SuperPOD: Data center design featuring NVIDIA DGX H100 systems. <https://docs.nvidia.com/dgx-superpod/design-guides/dgx-superpod-data-center-design-h100/latest/planning.html>, 2025c. Accessed: 2025-11-10.

- OECD. Adapting the Paris metropolitan area to a water-scarce future. Technical report, OECD Publishing, Paris, 2025.
- OECD.AI Expert Group on AI Compute and Climate. Measuring the environmental impacts of artificial intelligence compute and applications: The AI footprint. OECD Digital Economy Papers 341, OECD Publishing, Paris, 2022.
- David Patterson, Joseph Gonzalez, Quoc Le, Chen Liang, Lluis-Miquel Munguia, Daniel Rothchild, David So, Maud Texier, and Jeff Dean. Carbon emissions and large neural network training, 2021.
- David Patterson, Joseph Gonzalez, Urs Hölzle, Quoc Le, Chen Liang, Lluis-Miquel Munguia, Daniel Rothchild, David R. So, Maud Texier, and Jeff Dean. The carbon footprint of machine learning training will plateau, then shrink. *Computer*, 55(7):18–28, 2022. ISSN 1558-0814. doi: 10.1109/MC.2022.3148714.
- Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d’Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Ed Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorraine Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training Gopher, 2022.
- Paul Reig, Tianyi Luo, Eric Christensen, and Julie Sinistore. Guidance for calculating water use embedded in purchased electricity. *World Resources Institute*, 2020.
- Johan Rockström, Jonathan F. Donges, Ingo Fetzer, Maria A. Martin, Lan Wang-Erlandsson, and Katherine Richardson. Planetary boundaries guide humanity’s future on Earth. *Nature Reviews Earth & Environment*, 5(11):773–788, 2024. ISSN 2662-138X. doi: 10.1038/s43017-024-00597-z.
- RTE. Analyses et données de l’électricité - Bilan électrique 2024 - émissions. <https://analysesetdonnees.rte-france.com/bilan-electrique-2024/emissions#Introduction>, 2025. Accessed: 2025-11-18.
- David Sánchez, Sarah-Jane Baur, and Lara Eguren. Life cycle assessment of the Fairphone 5. Technical report, Fraunhofer IZM, Berlin, 2024.
- Scaleway. Infrastructures for LLMs in the cloud. <https://www.scaleway.com/en/blog/infrastructures-for-llms-in-the-cloud/>, 2024. Accessed: 2025-11-10.
- Scaleway. Scaleway impact report 2025. https://www-uploads.scaleway.com/Impact_Report2025_22ee3a8232.pdf, 2025. Accessed: 2025-10-30.
- Ian Schneider, Hui Xu, Stephan Benecke, David Patterson, Keguo Huang, Parthasarathy Ranganathan, and Cooper Elsworth. Life-cycle emissions of AI hardware: A cradle-to-grave approach and generational trends, 2025.
- Thibault Simon, David Ekchajzer, Adrien Berthelot, Eric Fourboul, Samuel Rince, and Romain Rouvoy. BoaviztAPI: A bottom-up model to assess the environmental impacts of cloud services. *ACM SIGEnergy Energy Informatics Review*, 4(5):84–90, 2025. doi: 10.1145/3727200.3727213.
- Matej Spetko, Lubomir Riha, and Branislav Jansik. Performance, power consumption and thermal behavioral evaluation of the DGX-2 platform. In *Parallel Computing: Technology Trends*, pp. 614–623. IOS Press, 2020. doi: 10.3233/APC200091.

- Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for deep learning in NLP. In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 3645–3650, Florence, Italy, 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1355.
- Sustainable Travel International. Carbon footprint calculator for travel. <https://sustainabletravel.org/our-work/carbon-offsets/calculate-footprint/>, 2024. Accessed: 2025-11-18.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, Allyson Ettinger, Michal Guerquin, David Heineman, Hamish Ivison, Pang Wei Koh, Jiacheng Liu, Saumya Malik, William Merrill, Lester James V. Miranda, Jacob Morrison, Tyler Murray, Crystal Nam, Jake Poznanski, Valentina Pyatkin, Aman Rangapur, Michael Schmitz, Sam Skjonsberg, David Wadden, Christopher Wilhelm, Michael Wilson, Luke Zettlemoyer, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2 OLMo 2 furious, 2025.
- TechPowerUp. NVIDIA H100 SXM5 80 GB. <https://www.techpowerup.com/gpu-specs/h100-sxm5-80-gb.c3900>, 2022. Accessed: 2025-11-10.
- TechPowerUp. Intel Xeon Platinum 8480+. <https://www.techpowerup.com/cpu-specs/xeon-platinum-8480.c2958>, 2023. Accessed: 2025-11-10.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. LLaMA: Open and efficient foundation language models, 2023a.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutvi Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023b.
- Gaël Varoquaux, Alexandra Sasha Luccioni, and Meredith Whittaker. Hype, sustainability, and the price of the bigger-is-better paradigm in AI, March 2025.
- Carole-Jean Wu, Ramya Raghavendra, Udit Gupta, Bilge Acun, Newsha Ardalani, Kiwan Maeng, Gloria Chang, Fiona Aga, Jinshi Huang, Charles Bai, Michael Gschwind, Anurag Gupta, Myle Ott, Anastasia Melnikov, Salvatore Candido, David Brooks, Geeta Chauhan, Benjamin Lee, Hsien-Hsin Lee, Bugra Akyildiz, Maximilian Balandat, Joe Spisak, Ravi Jain, Mike Rabbat, and Kim Hazelwood. Sustainable AI: Environmental implications, challenges and opportunities. In D. Marculescu, Y. Chi, and C. Wu (eds.), *Proceedings of Machine Learning and Systems*, volume 4, pp. 795–813, 2022. doi: 10.48550/arXiv.2111.00364.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. OPT: Open pre-trained transformer language models, 2022.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A survey of large language models, 2025.

Alesia Zhuk. Artificial intelligence impact on the environment: Hidden ecological costs and ethical-legal issues. *Journal of Digital Technologies and Law*, 1(4):932–954, 2023. ISSN 2949-2483. doi: 10.21202/jdtl.2023.40.

A Additional Results

In this appendix, we gather the numerical values represented in fig. 9 and fig. 11 of the main text: tab. 3 compiles all the results of our environmental impact assessment, disaggregated by hardware component and impact scope; and tab. 4 summarizes the hypothetical global warming potential and water consumption impacts of developing Moshi in different locations.

Table 3: **Environmental impacts of research.** For each environmental impact indicator, we estimate embodied impacts due to hardware production and operational impacts due to computation and datacenter energy consumption overheads. In the case of water consumption, we discern between datacenter cooling and power plant cooling for electricity production. We abbreviate *Motherboard* as *MoBo*. Represented in fig. 9.

Primary energy (MJ)										
Scope	Component									Total
	GPU	CPU	RAM	SSD1	SSD2	PSU	MoBo	Case	Assembly	
Embodied	5.71e5	2.56e4	4.52e5	1.87e4	1.38e5	1.23e5	1.62e4	4.26e4	1.33e3	1.39e6
Datacenter	1.08e7	7.09e4	2.55e5				7.64e6			1.87e7
Computation	2.77e7	1.82e5	6.56e5				1.97e7			4.82e7
Total	3.91e7	2.79e5	1.36e6				2.76e7			6.83e7

Global warming potential (kgCO ₂ eq)										
Scope	Component									Total
	GPU	CPU	RAM	SSD1	SSD2	PSU	MoBo	Case	Assembly	
Embodied	4.19e4	1.81e3	3.60e4	1.52e3	1.12e4	8.47e3	1.28e3	2.90e3	1.29e2	1.05e5
Datacenter	3.44e4	2.26e2	8.14e2				2.44e4			5.98e4
Computation	8.83e4	5.81e2	2.09e3				6.27e4			1.54e5
Total	1.65e5	2.62e3	3.89e4				1.13e5			3.19e5

Abiotic depletion potential (kgSbeq)										
Scope	Component									Total
	GPU	CPU	RAM	SSD1	SSD2	PSU	MoBo	Case	Assembly	
Embodied	1.38e0	7.90e-1	1.99e0	5.84e-2	3.80e-1	2.89e0	7.14e-2	3.91e-1	2.73e-05	7.95e0
Computation	1.05e-1	6.92e-4	2.49e-3				7.46e-2			1.83e-1
Datacenter	4.09e-2	2.69e-4	9.69e-4				2.90e-2			7.12e-2
Total	1.53e0	7.91e-01	1.99e0				3.90e0			8.21e0

Water consumption (L)					
Scope	Component				Total
	GPU	CPU	RAM	Other	
Datacenter cooling	1.00e7	6.60e4	2.38e5	7.12e6	1.75e7
Electricity production	6.02e5	3.96e3	1.42e4	4.27e5	1.05e6
Total	1.06e7	7.00e4	2.52e5	7.54e6	1.85e7

Table 4: **Operational impacts by training location.** Hypothetical global warming potential and water consumption of developing the model in different locations, excluding embodied impacts. Represented in fig. 11.

Operational impact	Location					
	Sweden	France	USA	Australia	China	Poland
Global warming potential (kgCO ₂ eq)	1.83e5	2.13e5	2.01e6	2.90e6	2.90e6	3.20e6
Water consumption (L)	2.69e7	1.85e7	1.31e7	1.67e7	2.49e7	6.91e6

In fig. 10 of the main text, we show how the embodied impacts of a compute node are distributed across component types (GPUs, CPUs, RAM modules, etc.). We compliment this information in tab. 5 and fig. 12, which illustrate the impacts of producing *a single unit* of each component type, plus the assembly process of the full compute node.

Table 5: **Embodied impacts of one component.** Production impacts for *a single unit* of each hardware component, plus the assembly process of the compute node: primary energy (PE), global warming potential (GWP), and abiotic depletion potential (ADP). We abbreviate *Motherboard* as *MoBo*. Represented in fig. 12.

Impact indicator	Component								
	GPU	CPU	RAM	SSD1	SSD2	PSU	MoBo	Case	Assembly
PE (MJ)	3.69e3	6.62e2	7.30e2	4.83e2	8.93e2	1.06e3	8.36e2	2.20e3	6.86e1
GWP (kgCO ₂ eq)	2.70e2	4.67e1	5.81e1	3.93e1	7.23e1	7.29e1	6.61e1	1.50e2	6.68e0
ADP (kgSbeq)	8.94e-3	2.04e-2	3.20e-3	1.51e-3	2.45e-3	2.49e-2	3.69e-3	2.02e-2	1.41e-6

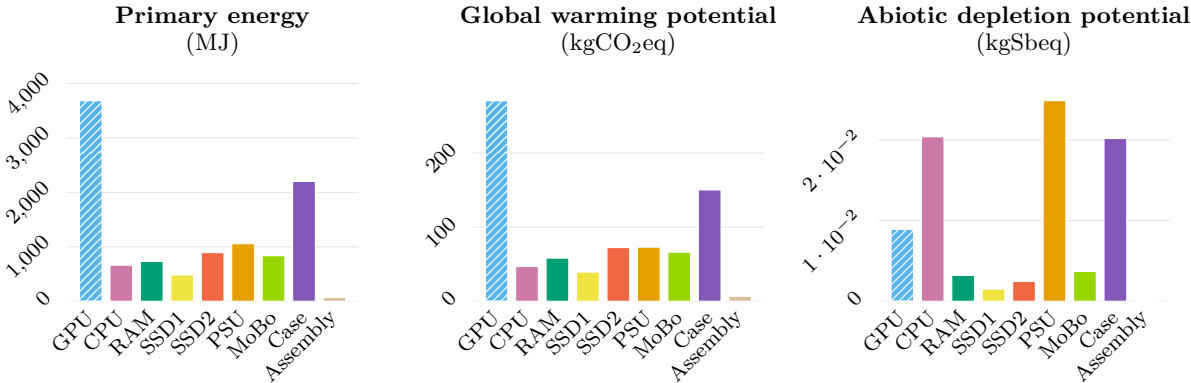


Figure 12: **Embodied impacts of one component.** Production impacts for *a single unit* of each hardware component. Solid fill impacts are estimated using Boavizta (Simon et al., 2025) and combine mineral and metal (ADPe) and fossil resource (ADPf) depletion in the case of abiotic depletion potential (ADP); dashed impacts come from ADEME (Lees-Perasso et al., 2026) and include ADPe only. We abbreviate *Motherboard* as *MoBo*. Values in tab. 5.

B Environmental Assessment Methodology

This appendix details all the formulas and values that we use to carry out our environmental assessment, thus providing complete transparency of our results and outlining sources of uncertainty. Sec. B.1 and sec. B.2 present the formulas used to derive, respectively, operational and embodied impacts.

B.1 Operational Impacts

We employ the following methodology to estimate operational primary energy, global warming potential, and water consumption, using the values provided in tab. 6 and starting from the energy consumed by the datacenter:

Energy Consumption. We estimate operational energy consumption as the sum of energy consumed for computation ($\text{oper}_E^{\text{computation}}$) and datacenter overheads ($\text{oper}_E^{\text{datacenter}}$). We consider cluster management, networking, and storage overheads (o_{cluster}) (NVIDIA, 2025c); as well as cooling, ventilation, and other datacenter overheads (power usage effectiveness or PUE):

$$\text{oper}_E = \text{oper}_E^{\text{computation}} + \underbrace{((\text{PUE} - 1) \times o_{\text{cluster}} + (o_{\text{cluster}} - 1)) \times \text{oper}_E^{\text{computation}}}_{\text{oper}_E^{\text{datacenter}}}.$$

Energy consumption for computation is the aggregate of the consumption of all hardware components in a compute node:

$$\text{oper}_E^{\text{computation}} = \text{oper}_E^{\text{GPU}} + \text{oper}_E^{\text{CPU}} + \text{oper}_E^{\text{RAM}} + \text{oper}_E^{\text{other}},$$

where the consumption for each hardware component is estimated as:

$$\text{oper}_E^{\text{hw}} = \mathcal{C} \times \frac{q_{\text{hw}}}{q_{\text{GPU}}} \times \begin{cases} u_{\text{hw}} \times \text{TDP}_{\text{hw}} & \text{if hw} \in \{\text{GPU, CPU}\} \\ P_{\text{hw}} & \text{if hw} \in \{\text{RAM, other}\} \end{cases},$$

with \mathcal{C} the total development compute in GPU-hours, q_{hw} the quantity of component hw per node, u_{hw} the hardware utilization, TDP_{hw} the thermal design power of the hardware, and P_{hw} the constant power consumption of the hardware.

We establish average utilization factors for CPU and GPU (u_{CPU} and u_{GPU}) based on Kyutai’s observations during training.

To obtain the power consumption of RAM, we apply CodeCarbon’s methodology with efficiency scaling (CodeCarbon, 2026), which results in $5 \times (4 + 4 \times 0.9 + 8 \times 0.8 + 16 \times 0.7) = 126$ watts for the 32 RAM modules of a compute node, and therefore 3.94 watts per RAM module (P_{RAM}). We define the power consumption of the remaining node hardware, P_{other} , as the difference between the total power consumption of a compute node (10,200 watts (NVIDIA, 2025a)) and the consumption of GPUs, CPUs, and RAM.

Primary Energy (PE). We compute operational primary energy as done in Boavizta (2026), starting from energy consumption:

$$\text{oper}_{\text{PE}} = \text{PE}_{\text{kWh}} \times \text{oper}_E,$$

with

$$\text{PE}_{\text{kWh}} = \frac{\text{ADP}_{\text{f}_{\text{kWh}}}}{1 - \% \text{renewable}},$$

where $\text{ADP}_{\text{f}_{\text{kWh}}}$ is the fossil resource depletion per kilowatt-hour of the electrical grid, and $\% \text{renewable}$ is the percentage of electricity generated from renewable energy sources.

Global Warming Potential (GWP). We estimate operational global warming potential by multiplying the operational energy consumption by the carbon intensity (CI) of the electrical grid:

$$\text{oper}_{\text{GWP}} = \text{oper}_E \times \text{CI},$$

making the appropriate unit conversions. Although not shown in the equation, we split operational global warming potential into impacts due to computation and datacenter overheads, as we do with primary energy.

Abiotic Depletion Potential (ADP). Similarly to primary energy, we compute operational mineral and metal resource depletion from energy consumption as follows:

$$\text{oper}_{\text{ADP}} = \text{ADPe}_{\text{kWh}} \times \text{oper}_{\text{E}},$$

with ADPe_{kWh} the mineral and metal resource depletion per generated kilowatt-hour.

Water Consumption (WC). Following the methodology of Li et al. (2025a), we estimate operational water consumption as:

$$\text{oper}_{\text{WC}} = \underbrace{\text{WUE} \times o_{\text{cluster}} \times \text{oper}_{\text{E}}^{\text{computation}}}_{\text{oper}_{\text{WC}}^{\text{datacenter cooling}}} + \underbrace{\text{EWIF} \times \text{oper}_{\text{E}}}_{\text{oper}_{\text{WC}}^{\text{electricity production}}},$$

where WUE is the water usage effectiveness of the datacenter, and EWIF is the electricity water intensity factor of the local electrical grid. Once again, we make the appropriate unit conversions.

We obtain the electricity water intensity factor (EWIF) as a weighted average of the water consumption per energy source reported by Reig et al. (2020), weighted by the energy mix in the target location in 2024 as reported in the EMBER database (EMBER, 2025b). Since Reig et al. (2020) and EMBER (2025b) do not use the same terminology to name energy sources, we establish correspondences as follows: *solar - photovoltaic*, *bioenergy - biomass*, *other renewables - geothermal*⁷, *gas - natural gas*, *coal - hard coal*, and *other fossil - heavy fuel oil*. It should be noted that Reig et al. (2020) consider the water consumption of solar and wind energy to be zero, whereas other sources do not (Jin et al., 2019). For future work, we recommend referring to platforms such as Wattnet (Melguizo et al., 2026).

B.2 Embodied Impacts

We estimate the production impacts for a single unit of hardware as follows:

$$\begin{aligned} \text{emb}_{\text{imp}}^{\text{CPU unit}} &= \text{base}_{\text{imp}}^{\text{CPU}} + \text{die_size}^{\text{CPU}} \times \text{die}_{\text{imp}}^{\text{CPU}} \\ \text{emb}_{\text{imp}}^{\text{mem unit}} &= \text{base}_{\text{imp}}^{\text{mem}} + \frac{\text{capacity}_{\text{mem}}^{\text{mem}}}{\text{density}_{\text{mem}}^{\text{mem}}} \times \text{die}_{\text{imp}}^{\text{mem}} \quad \forall \text{mem} \in \{\text{RAM}, \text{SSD}\} \\ \text{emb}_{\text{imp}}^{\text{PSU unit}} &= \text{weight}^{\text{PSU}} \times \text{base}_{\text{imp}}^{\text{PSU}}, \end{aligned}$$

where $\text{imp} \in \{\text{PE}, \text{GWP}, \text{ADP}\}$ for primary energy, global warming potential, and abiotic depletion potential respectively. The impacts of the remaining hardware, and of the assembly process of the compute node, are constant values from Boavizta (Simon et al., 2025) and the ADEME report on GPU production impacts (Lees-Perasso et al., 2026). The base and die impact factors that we use are gathered in tab. 7, and tab. 8 lists the specifications of each hardware component. The allocated embodied impact for the duration of use is:

$$\text{emb}_{\text{imp}}^{\text{hw}} = \frac{\mathcal{C}}{\mathcal{D}} \times \frac{q_{\text{hw}}}{q_{\text{GPU}}} \times \text{emb}_{\text{imp}}^{\text{hw unit}},$$

where \mathcal{C} is the total development compute in GPU-hours; q_{hw} is the quantity of component hw per node, with hw one of: GPU, CPU, RAM, SSD1, SSD2, PSU, motherboard, case, or assembly; and where $\mathcal{D} = \text{lifespan} \times \text{utilization_rate}$ is the total duration of use of the hardware equipment throughout its lifespan, in hours. We assume an equipment lifespan of four years, in line with values employed in related work (Morand et al., 2024; Schneider et al., 2025; Falk et al., 2025; Desroches et al., 2025), and a reasonable average utilization rate of 0.6 (Luccioni et al., 2023; Wu et al., 2022). The values of \mathcal{C} and q_{hw} can be found in tab. 6.

⁷Only applies to the United States, where the share of *other renewables* is 0.4%, and France, where it is 0.1%.

Table 6: **Environmental assessment variables.** Definitions, values, and sources of the main variables used in the environmental assessment of Moshi. For sources marked with *, values are computed instead of taken directly.

Variable	Notation	Value	Unit	Source
Total compute	C	3.26e6	GPU-hours	Kyutai logs
Hardware quantity per node	q_{GPU}	8	-	(NVIDIA, 2025a)
	q_{CPU}	2	-	
	q_{RAM}	32	modules	
	q_{SSD1}	2	disks	
	q_{SSD2}	8	disks	
	q_{PSU}	6	-	
	$q_{\text{motherboard}}$	1	-	
	q_{case}	1	-	
	q_{assembly}	1	-	
	q_{other}	1	-	
Average GPU utilization	u_{GPU}	9.50e-1	-	Kyutai estimates
Average CPU utilization	u_{CPU}	5.00e-2	-	
GPU thermal design power	TDP_{GPU}	7.00e2	W	(TechPowerUp, 2022)
CPU thermal design power	TDP_{CPU}	3.50e2		(TechPowerUp, 2023)
RAM module power	P_{RAM}	3.94e0		(CodeCarbon, 2026; NVIDIA, 2025a)*
Other node hardware power	P_{other}	3.77e3		(NVIDIA, 2025a)*
Power usage effectiveness	PUE	1.25e0	-	(Scaleway, 2025)
Water usage effectiveness	WUE	2.50e-1	L/kWh	(Scaleway, 2025)
Cluster overheads	α_{cluster}	1.11e0	-	(NVIDIA, 2025c)*
Carbon intensity (2024)	CI_{SE}	3.50e1	gCO ₂ eq/kWh	(EMBER, 2025a)
	CI_{FR}	4.10e1		
	CI_{US}	3.84e2		
	CI_{AU}	5.54e2		
	CI_{CN}	5.55e2		
	CI_{PL}	6.12e2		
Electricity water intensity factor (2024)	EWIF_{SE}	4.94e0	L/kWh	(EMBER, 2025b; Reig et al., 2020)*
	EWIF_{FR}	3.34e0		
	EWIF_{US}	2.30e0		
	EWIF_{AU}	2.99e0		
	EWIF_{CN}	4.57e0		
	EWIF_{PL}	1.12e0		
Fossil depletion per kWh (FR)	$\text{ADP}_{\text{kWh}}^{\text{f}}$	9.31e0	MJ/kWh	(ADEME, 2023)
Mineral and metal depletion per kWh (FR)	$\text{ADP}_{\text{kWh}}^{\text{m}}$	4.86e-8	kgSbeq/kWh	(ADEME, 2023)
Renewable-generated electricity (FR, 2024)	%renewable	2.72e-1	-	(EMBER, 2026)

Table 7: **Embodied impact factors.** Impact factors provided by Boavizta (Simon et al., 2025) and the ADEME agency (Lees-Perasso et al., 2026) to estimate global warming potential (GWP), primary energy (PE), and abiotic depletion potential (ADP) impacts of hardware production and transport. We abbreviate *Motherboard* as *MoBo*.

	Boavizta									ADEME
	RAM, SSD die	RAM base	SSD base	CPU die base		PSU	MoBo	Case	Assembly	GPU
GWP kgCO ₂ eq	2.20e0 per cm ²	5.22e0	6.34e0	1.97e0 per cm ²	9.14e0	2.43e1 per kg	6.61e1	1.50e2	6.68e0	2.70e2
PE MJ	2.73e1 per cm ²	7.40e1	7.40e1	2.65e1 per cm ²	1.56e2	3.52e2 per kg	8.36e2	2.20e3	6.86e1	3.69e3
ADP kgSbeq	6.30e-5 per cm ²	1.69e-3	5.63e-4	5.87e-7 per cm ²	2.04e-2	8.30e-3 per kg	3.69e-3	2.02e-2	1.41e-6	8.94e-3

Table 8: **Compute node component specifications.** Hardware specifications of the NVIDIA DGX H100 compute node (NVIDIA, 2025a). We omit fans, network cards, and NVSwitches. For SSD memory density, we select a value from the Boavizta repository reflecting high-end SSDs before the release of the DGX H100.

Component	Model/Type	Specification	Value	Unit	Source
CPU	Intel Xeon Platinum 8480C	die size	19.08	cm ²	(TechPowerUp, 2023)
		die size	8.14	cm ²	(TechPowerUp, 2022)
GPU	NVIDIA H100 SXM HBM3	VRAM capacity	80	GB	(NVIDIA, 2025a)
		VRAM density	1.65	GB/cm ²	(Moon et al., 2023)
RAM	DDR5	capacity	64	GB	(NVIDIA, 2025a)
		density	2.66	GB/cm ²	(Choe, 2022)
SSD1	NVMe M.2	capacity	1920	GB	(NVIDIA, 2025a)
		density	128	GB/cm ²	(Simon et al., 2025; Choe, 2021)
SSD2	NVMe U.2	capacity	3840	GB	(NVIDIA, 2025a)
		density	128	GB/cm ²	(Simon et al., 2025; Choe, 2021)
PSU	-	weight	3	kg	(Simon et al., 2025)