

**ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA**

**DEPARTMENT OF COMPUTER SCIENCE
AND ENGINEERING**

ARTIFICIAL INTELLIGENCE

MASTER THESIS

in

Intelligent Systems

**HARDWARE DIMENSIONING FOR
ENVIRONMENTAL SUSTAINABILITY:
BENCHMARK OF AI ALGORITHMS AND
ENVIRONMENTAL IMPACT**

CANDIDATE

Enrico Morselli

SUPERVISOR

Prof. Andrea Borghesi

CO-SUPERVISOR

Allegra De Filippo, PhD.

Academic year 2024-2025

Session 5th

dedicated(X) :- friend(X).

Contents

1	Introduction	1
1.1	Background and Rationale	1
2	Related Works	3
2.1	Sustainability in AI	3
2.2	Tools for tracking Carbon Emissions	5
3	Metodology	7
3.1	Empirical Model Learning in HADA	7
3.2	Integration of CodeCarbon in HADA	10
4	Experimental Analysis	12
4.1	Benchmarking on Different Hardware Platforms	12
5	HADA-as-a-Service	13
5.1	HADA Web Application	13
6	Conclusions	15
	Bibliography	16
	Acknowledgements	19

List of Figures

List of Tables

Chapter 1

Introduction

1.1 Background and Rationale

In recent years we have witnessed a dramatic increase in the performance of Artificial Intelligence technologies. Even if AI still fails to exceed human ability in some complex cognitive tasks, as of 2023 it has surpassed human capabilities in a range of tasks, such as image classification, basic reading comprehension, visual reasoning and natural language inference [10]. Not to mention the astonishing results achieved by Generative AI in tasks as Image and Video Generation [10]. This great advances in performance were made possible by a massive upscale of model sizes and computational resources ("compute" in short) dedicated to training state-of-the-art AI models. Research shows that for frontier AI models (i.e. those that were in the top 10 of training compute when they were released), the training compute has grown by a factor of 4-5x/year since 2010 [13]. This surge in required compute has driven a corresponding spike in energy consumption for AI, and consequently, an higher environmental impact due to CO₂ emissions. For instance, for training their LLaMA models, Meta AI researchers have estimated a period of approximately 5 months of on 2048 A100 80GB GPUs, resulting in a total of 2,638MWh of energy and a total emission of 1,015 tCO₂eq [16]. Given the widespread application, the steep increase in model size and complexity and

the crescent energy requirements for AI applications, the Carbon Footprint of AI has become a growing concern in the context of the current climate emergency.

In this work, we will explore an approach for addressing the issue of AI sustainability, by means of HADA (HARdware Dimensioning for AI Algorithms), which is a framework that uses ML to learn the relationship between an algorithm configuration and performance metrics, like total runtime, solution cost and memory usagem and then uses Optimization to find the best Hardware architecture and its configuration to run an algorithm under required performances and budget limits which is the problem known as Hardware Dimensioning [3]. What we will do is to extend this framework in order to consider also the performance of the algorithms in terms of Energy consumption and Carbon Emissions, so that, ideally, we could find the best Algorithm and Hardware configuration that reduces the Carbon Footprint of computation. We will then proceed to test this approach on some small-scale algorithms that we could easily execute in a timely manner on local machines and HPC clusters.

The rest of the work is structured as follows:

- **Chapter 2** Introduces Related works that addressed the issue of AI's carbon footprint, and how Carbon Footprint is determined
- **Chapter 3** Introduces some theoretical background about HADA, and explains the integration of the new metrics.
- **Chapter 4** Presents the experimental setup and the results of the experiments
- **Chapter 5** Presents the HADA framework, providing an overview of the tool
- **Chapter 6** Presents the conclusions

Chapter 2

Related Works

2.1 Sustainability in AI

The environmental impact of training large AI models is illustrated by recent empirical assessments. Strubell et al. (2019) quantified the CO₂ emissions of several NLP models and found that training a big transformer with extensive hyperparameter tuning (including neural architecture search) emitted roughly 626,000 pounds of CO₂ - about the same as the lifetime emissions of five cars [14] [4]. These findings brought attention to the fact that accuracy gains in AI often come at a steep energy and carbon price. In response, researchers have begun to systematically reporting energy use and consequent CO₂ emissions of model training to raise awareness [4] [11]. Beyond individual models, broader studies have examined AI's total energy and environmental footprint across the industry. Henderson et al. (2020) [7] introduced a framework for tracking real-time energy consumption and carbon emissions during ML experiments, encouraging researchers to include these metrics in publications. Their work underscored that transparent reporting is essential for understanding and ultimately reducing AI's climate impacts. Industry-scale analyses also reveal sobering trends. Gupta et al. (2021) [6] analyzed the end-to-end footprint of computing and found that while operational emissions (from running hardware) have been partly curbed by efficiency improvements, the overall

carbon footprint of computing continues to grow due to increased scale. Notably, they showed that for modern data centers and mobile devices, manufacturing and infrastructure (embodied carbon) now account for the majority of emissions. In other words, as data centers adopt cleaner power, the emissions “hidden” in hardware supply chains (chip fabrication, server manufacturing, etc.) become a dominant concern. Similarly, Wu et al. (2022) [17] present a holistic study of AI at a large tech company (Meta/Facebook), examining the entire AI model lifecycle - from data processing and training to inference and hardware lifecycle. They report super-linear growth in AI workloads and infrastructure: for instance, daily inference operations doubled over a recent 3-year span, forcing a 2.5x expansion in server capacity. Crucially, Wu et al. also highlight that embodied carbon is an increasing fraction of AI’s total footprint, echoing that improvements in hardware efficiency alone cannot eliminate AI’s impact. Their analysis argues for looking beyond training alone - considering data center construction, supply chains, and the frequency of model retraining - to truly grasp AI’s environmental impact.

A recurring theme in these studies is the diminishing return on energy investment for AI model improvements. As models get larger and more complex, the incremental accuracy gains often require disproportionately more compute (and thus energy). Schwartz et al. (2020) [12] dub the status quo “Red AI,” where researchers prioritize accuracy at almost any computational cost, and they note this trend is unsustainable both environmentally and even economically. Thompson et al. (2021) [15] similarly observed that progress in benchmarks was coming with exponentially increasing computing cost, warning of diminishing returns and calling the situation unsustainable. Another challenge is equitable access: massive energy requirements make cutting-edge AI research expensive, potentially concentrating it in wealthy institutions and regions with robust infrastructure. This raises concerns that AI’s growing energy hunger not only harms the planet but also exacerbates inequalities in who can afford to do top-tier research. These concerns have prompted calls

for a paradigm shift toward “Green AI”, where efficiency and sustainability are treated as primary goals in model development. Several studies have addressed the issue of AI’s carbon footprint.

2.2 Tools for tracking Carbon Emissions

The growing awareness about AI’s Carbon Footprint also motivated the development of dedicated tools and methodologies to monitor the carbon footprint of AI workloads. A number of open-source tools and frameworks have been created to help practitioners measure the energy consumption and CO₂ emissions of their code. ML CO₂ Impact (Machine Learning Emissions Calculator), introduced by Lacoste et al. (2019), [8] was one of the early tools for estimating emissions from model training. It is a web-based calculator where users input information about their training run - such as hardware type (CPU/GPU), runtime, cloud provider, and location - and it computes the energy consumed and corresponding carbon emissions. The calculator leverages known power draw profiles of hardware and regional carbon intensity factors. According to [1] the latest incarnation of this tool has evolved under the umbrella of the CodeCarbon initiative, integrating its functionality into CodeCarbon’s open-source package.

CodeCarbon [2], which is the tool we used to expand HADA in this work, is an open-source Python package for tracking the carbon footprint of computing projects. It integrates into ML code to log the resources used (CPU, GPU, etc.) and estimates the CO₂ emissions produced by the workload. Uniquely, CodeCarbon accounts for the location of the computation - using region-specific electricity carbon intensity data - to provide location-dependent emission estimates. This allows developers to see, for instance, that running the same training job on a low-carbon grid (e.g. hydro-powered Montreal) results in a much smaller footprint than running it on a coal-heavy

grid. The goal is to inform and incentivize researchers and engineers to optimize or relocate their workloads to reduce emissions.

Cite other works which provides tools to track carbon emissions, e.g. LLM-Carbon [5].

Chapter 3

Metodology

3.1 Empirical Model Learning in HADA

As we mentioned in the introduction, the main focus of this work is Hardware Dimensioning, i.e. the problem of finding the right Hardware configuration to run the desired algorithm under required performance and budget limits. More specifically, we are given a target algorithm, and we would like to know what is the optimal settings of Algorithm Hyperparameters and Hardware Architecture to run the algorithm under a set of constraints in terms of performance requirements. This is not a trivial problem to model, owing to the complexity of knowing beforehand the performance of an AI algorithm on different HW architectures and evaluating the effect of different HW and algorithm configurations. The idea behind HADA is to integrate the domain knowledge held by experts (time constraints, required solution quality, budget limits) with data-driven models. The basis of HADA is the **Empirical Model Learning (EML)**, [9] framework, which uses Machine Learning (ML) techniques to learn the model of an Optimisation problem. The idea is to learn, rather than directly express, complex relationships between algorithm performance and HW resources via ML models, using *empirical* knowledge. These models are then embedded in the optimisation problem. In general, this approach is useful

whenever the problem to solve using optimisation is very complex and therefore not trivial to derive a model to describe it. Broadly speaking, EML deals with solving declarative optimisation models with a complex component, h , which represents the relation between variables which can be acted upon (the decision variables, also called *decidables*) x and the *observables* y related to the system considered; the function $h(x) = y$ describes this relationships. Since the h component is complex and we cannot optimise directly over it, we exploit empirical knowledge to build a surrogate model h_θ learned from data, where θ is the parameter vector.

We present here the general mathematical formulation of EML:

$$\min \quad f(x, y) \quad (3.1)$$

$$\text{s.t.} \quad h_\theta(x) = y \quad (3.2)$$

$$g_j(x, y) \quad \forall j \in J \quad (3.3)$$

$$x_i \in D_i \quad \forall x_i \in x \quad (3.4)$$

where x is the decision variables vector (each x_i has its own domain D_i) and y is the vector of observed variables. The goal is to minimize the objective function $f(x, y)$, which might depend on both decision and observed variables. Both decision and observed variables can be subject to a set of constraints $g_j(x, y)$, such as classical inequalities from Mathematical Programming and combinatorial predicates from Constraint Programming (e.g., global constraints). The function $h_\theta(x)$ represents the approximate behaviour of the complex relation between decision and observed variables, and it is, in practice, the encoding of a ML model.

The surrogate model $h_\theta(x)$ is trained as in the typical Supervised Learning setting; EML requires a training set $\mathcal{S} = (x_i, h(x_i))_{i=1}^m$ which is then used to find the parameter vector θ minimizing a loss function:

$$\arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m L(h_{\theta}(x_i), y_i^*) \quad (3.5)$$

where y_i^* are the ground-truth labels (or targets, in case of regression) of each data point in the training set \mathcal{S} and L is the loss function (e.g., L_1 or L_2 loss for scalar y 's).

The way in which HADA works can then be summarized in three main phases:

1. *Data Set Collection* - an initial phase to collect the data set \mathcal{S} by running multiple times the target algorithms, under different configuration. This implies running the algorithms with different configurations on different Hardware platforms;
2. *Surrogate Model Creation* - once a training set is available, a set of ML models is then trained on such data. These are the h_{θ} surrogate models. Then, these models are encoded as a set of variables and constraints following the EML paradigm;
3. *Optimization* - post the user-defined constraints and objective function on top of the combinatorial structure formed by the encoded ML models and the domain-knowledge constraints, and finally solve the optimization model (either until an optimal solution or a time limit is reached).

In the HADA paper [3] is described how they gathered a dataset containing records describing:

- Algorithm configuration;
- Time required to find a solution;
- Cost of the solution found;
- Average memory used by the algorithm.

Our goal is to extend this data with metrics that will allow us to assess the impact of the Algorithm Configuration and the Hardware Architecture on Energy Consumption, and consequently, on the Carbon Footprint of the algorithm.

3.2 Integration of CodeCarbon in HADA

The first step to extend HADA with Carbon Intensity is therefore to gather a dataset which includes data about carbon footprint of the algorithms execution. To do so, we can extend the code for running our algorithms with Codecarbon [2]. The Carbon Footprint of an algorithm depends on two factors: the energy needed to run it and the pollutants emitted when producing such energy. The former depends on the computing resources used (e.g. number of cores, running time, data centre efficiency) while the latter, called carbon intensity, depends on the location and production methods used (e.g. nuclear, gas or coal). The Carbon Intensity is expressed in terms of carbon dioxide equivalent (CO_2e), and summarises the global warming effect of the Greenhouse Gases (GHG) emitted in the determined timeframe. [Green Algorithms]. According to [insert codecarbon docs ref], the carbon footprint can then be computed with the following formula:

$$C = E \times CI \quad (3.6)$$

Where:

- E is the energy consumed by the computational infrastructure: quantified as kilowatt-hours (kWh).
- CI is the Carbon Intensity of the electricity consumed for computation: quantified as g of CO_2e emitted per kilowatt-hour of electricity.

Codecarbon directly measures the energy consumption of the CPU, GPU and RAM on which the code is executed at intervals of 15 sec (by default,

but it can also be modified). It also monitor the duration of code execution to compute the total electricity consumption. Then it retrieves information about the carbon intensity of the electricity in your geographic area based on the location. Carbon Intensity of the consumed electricity is calculated as a weighted average of the emissions from the different energy sources that are used to generate electricity, including fossil fuels and renewables. Based on the mix of energy sources in the local grid, this package calculates the Carbon Intensity of the electricity consumed.

To extend HADA for sustainable AI, we integrate CodeCarbon to track emissions in:

- ANTICIPATE and CONTINGENCY algorithms.
- MaxFlow Algorithms:
 - Boykov-Kolmogorov (BK)
 - Excess Incremental Breadth First Search (EIBFS)
 - Hochbaum's Pseudo Flow (HPF)

Chapter 4

Experimental Analysis

4.1 Benchmarking on Different Hardware Platforms

Experiments were conducted on:

- MacBook Pro (2019)
- Leonardo Supercomputer (CINECA HPC)

Each algorithm was executed on 30 instances with hyperparameter values ranging from 1 to 100, generating datasets with 6,000 records per algorithm per hardware platform.

Chapter 5

HADA-as-a-Service

5.1 HADA Web Application

Benchmark data was integrated into the HADA web application, requiring:

- Creation of JSON configuration files for each algorithm-hardware combination.
- Specification of hyperparameters and performance targets.

Example JSON structure:

```
{
  "name": "anticipate",
  "HW_ID": "macbook",
  "hyperparams": [
    {"ID": "num_scenarios", "type": "int", "LB": 1, "UB": 100}
  ],
  "targets": [
    {"ID": "time", "LB": null, "UB": null},
    {"ID": "memory", "LB": null, "UB": null},
    {"ID": "emissions", "LB": null, "UB": null}
  ]
}
```

}

Chapter 6

Conclusions

This work extends HADA by integrating carbon emission constraints, enhancing its applicability for sustainable AI hardware selection. Through experimental benchmarks on laptops and HPC systems, we validated the framework's ability to balance performance and environmental impact. The web-based prototype enables users to make informed decisions when configuring AI workloads under sustainability constraints.

Bibliography

- [1] N. Bannour et al. “Evaluating the Carbon Footprint of NLP Methods: A Survey and Analysis of Existing Tools”. In: *Proceedings of the Second Workshop on Simple and Efficient Natural Language Processing*. Virtual: Association for Computational Linguistics, 2021, pp. 11–21. DOI: 10.18653/v1/2021.sustainlp-1.2. URL: <https://aclanthology.org/2021.sustainlp-1.2/>.
- [2] B. Courty et al. *mlco2/codecarbon: v2.4.1*. Version v2.4.1. 2024. DOI: 10.5281/zenodo.11171501. URL: <https://doi.org/10.5281/zenodo.11171501>.
- [3] A. De Filippo et al. “HADA: An automated tool for hardware dimensioning of AI applications”. In: *Knowledge-Based Systems* 251 (2022), p. 109199. ISSN: 0950-7051. DOI: <https://doi.org/10.1016/j.knosys.2022.109199>. URL: <https://www.sciencedirect.com/science/article/pii/S0950705122005974>.
- [4] J. Dodge et al. “Measuring the Carbon Intensity of AI in Cloud Instances”. In: *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. FAccT ’22. Seoul, Republic of Korea: Association for Computing Machinery, 2022, pp. 1877–1894. ISBN: 9781450393522. URL: <https://doi.org/10.1145/3531146.3533234>.

-
- [5] A. Faiz et al. *LLMCarbon: Modeling the end-to-end Carbon Footprint of Large Language Models*. 2024. arXiv: 2309.14393 [cs.CL]. URL: <https://arxiv.org/abs/2309.14393>.
 - [6] U. Gupta et al. “Chasing Carbon: The Elusive Environmental Footprint of Computing”. In: *CoRR* abs/2011.02839 (2020). arXiv: 2011.02839. URL: <https://arxiv.org/abs/2011.02839>.
 - [7] P. Henderson et al. “Towards the Systematic Reporting of the Energy and Carbon Footprints of Machine Learning”. In: *Journal of Machine Learning Research* 21.248 (2020), pp. 1–43. URL: <http://jmlr.org/papers/v21/20-312.html>.
 - [8] A. Lacoste et al. “Quantifying the Carbon Emissions of Machine Learning”. In: *arXiv preprint* (2019). arXiv: 1910.09700. URL: <https://arxiv.org/abs/1910.09700>.
 - [9] M. Lombardi, M. Milano, and A. Bartolini. “Empirical decision model learning”. In: *Artificial Intelligence* 244 (2017). Combining Constraint Solving with Mining and Learning, pp. 343–367. ISSN: 0004-3702. DOI: <https://doi.org/10.1016/j.artint.2016.01.005>. URL: <https://www.sciencedirect.com/science/article/pii/S0004370216000126>.
 - [10] N. Maslej et al. *The AI Index 2024 Annual Report*. Tech. rep. AI Index Steering Committee, Stanford University, 2024.
 - [11] D. A. Patterson et al. “Carbon Emissions and Large Neural Network Training”. In: *CoRR* abs/2104.10350 (2021). arXiv: 2104.10350. URL: <https://arxiv.org/abs/2104.10350>.
 - [12] R. Schwartz et al. “Green AI”. In: *CoRR* abs/1907.10597 (2019). arXiv: 1907.10597. URL: <http://arxiv.org/abs/1907.10597>.

-
- [13] J. Sevilla and E. Roldán. *Training Compute of Frontier AI Models Grows by 4-5x per Year*. Accessed: 2025-03-03. 2024. URL: <https://epoch.ai/blog/training-compute-of-frontier-ai-models-grows-by-4-5x-per-year>.
 - [14] E. Strubell, A. Ganesh, and A. McCallum. “Energy and Policy Considerations for Deep Learning in NLP”. In: *CoRR* abs/1906.02243 (2019). arXiv: 1906.02243. URL: <http://arxiv.org/abs/1906.02243>.
 - [15] N. Thompson et al. “Deep Learning’s Diminishing Returns: The Cost of Improvement is Becoming Unsustainable”. In: *IEEE Spectrum* 58 (2021), pp. 50–55. ISSN: 0018-9235.
 - [16] H. Touvron et al. *LLaMA: Open and Efficient Foundation Language Models*. 2023. arXiv: 2302.13971 [cs.CL]. URL: <https://arxiv.org/abs/2302.13971>.
 - [17] C. Wu et al. “Sustainable AI: Environmental Implications, Challenges and Opportunities”. In: *CoRR* abs/2111.00364 (2021). arXiv: 2111.00364. URL: <https://arxiv.org/abs/2111.00364>.

Acknowledgements

I'm very grateful to the inventor of the Prolog language, without whom this thesis couldn't exist. I'd also like to acknowledge my advisor Prof. Mario Rossi by tail-recursively acknowledging my advisor.