

Computing Within Limits: An Empirical Study of Energy Consumption in ML Training and Inference

Ioannis Mavromatis*, Kostas Katsaros, and Aftab Khan†,

*Digital Catapult, London, UK

†Bristol Research & Innovation Laboratory, Bristol, Toshiba Europe Ltd., UK

Emails: {ioannis.mavromatis, kostas.katsaros}@digicatapult.org.uk, aftab.khan@toshiba-bril.com

Abstract—This study presents an empirical investigation into the energy consumption of Machine Learning (ML) in immersive media applications. Acknowledging the growing environmental impact of ML, we examine various model architectures and hyperparameters in both training and inference phases to identify energy-efficient practices. Our study leverages software-based power measurements for ease of replication across different configurations, models and datasets. In this paper, we examine multiple models and hardware configurations to identify correlations across the different measurements and metrics and key contributors to energy reduction. Our analysis offers practical guidelines for constructing sustainable ML operations, emphasising energy consumption and carbon footprint reductions while maintaining performance. As identified, short-living profiling can quantify the long-term expected energy consumption. Moreover, model parameters can be used to accurately estimate the expected total energy without the need for extensive experimentation.

Index Terms—Machine Learning, Power Profiling, Energy Consumption, Sustainable AI, Green AI

I. INTRODUCTION

In recent years, media capture and processing technologies enabled new forms of true 3-D media content that increase the degree of user immersion. We currently see applications around Virtual Reality (VR) and Augmented Reality (AR) gaming, interactive art installations, education, etc. [1]. Machine Learning (ML) plays a pivotal role in all the above applications, acting as a Quality of Experience (QoE) predictor [2], recognising and classifying images [3], optimising the Content Delivery Networks (CDNs) [4], and more.

As immersive media applications progressively become the norm [5], the content creation (e.g., videos for 360° experiences, games for VR/AR platforms, etc.) and content delivery pipelines (e.g., ICT infrastructure - data centres, cloud computing, CDNs, etc.) become a key contributor to global energy consumption and the sector’s emissions [6]. ML is again at the forefront, providing ways for reducing the energy consumed across the entire pipeline [7].

ML is undoubtedly one of the key factors for innovation in the area of immersive media. However, its extensive use across the entire pipeline makes ML a significant sustainability concern [8]. The intensive computation required for training and deploying Deep Learning (DL) models contributes to substantial energy consumption and, thus, carbon emissions [8]. Projections estimate that ML pipelines will produce 2% of the global carbon emissions by 2030 [9].

In the immersive media literature, energy reductions have become a priority [7], [10], but the power consumed by the use of ML models is often overlooked. Driven by the above, in this paper, we present an empirical study of the energy consumption of a simple image classification example using ML. With this work, we aim to offer practical guidelines and best practices that researchers and practitioners can adopt in their ML pipelines.

Current trends reveal an increased interest in Green and Sustainable ML [11], [12]. Sustainable ML practices [11] encompass efficient use of computational resources and holistic optimisation of ML pipelines that collectively lead to reduced energy consumption, minimised carbon footprints, and economic benefits. As ML become increasingly integral nowadays, its sustainability will be crucial for its overall impact and acceptability [12].

Our paper aims to provide insights into how ML lifecycles can be optimised for lower energy consumption without compromising performance. We analyse various model architectures and hyperparameters, both for training and inference, to identify areas where energy consumption can be reduced. Based on our findings, we will critically comment on the key contributors to energy reduction and provide ways for estimating the expected energy consumption based on various model parameters. Our work can be leveraged by ML practitioners aiming for energy-aware optimisations in their ML pipelines.

The remainder of this paper is structured as follows: Sec. II presents recent activities around sustainable ML and discusses their limitations. Green ML is described in III outlining the energy consumed within an MLOps pipeline. The methodology used for our extensive investigation is illustrated in Sec. IV. Secs V and VI present our results and lessons learned, respectively. Finally, the paper is concluded in Sec. VII

II. RELATED WORK

Many works have discussed Green and Sustainable ML. Some notable examples are [11]–[13], where the authors provide statistics on how ML’s energy consumption will increase over time. Authors in [11] also compare transformer models running in Google’s data centres. All papers discuss the potential benefits of energy reduction from good practices (e.g., early existing, knowledge transfer, etc.) but do not systematically assess those. Our work contributes towards that by conducting an empirical study on real-world hardware.

Machine Learning Model Development and Deployment Phases

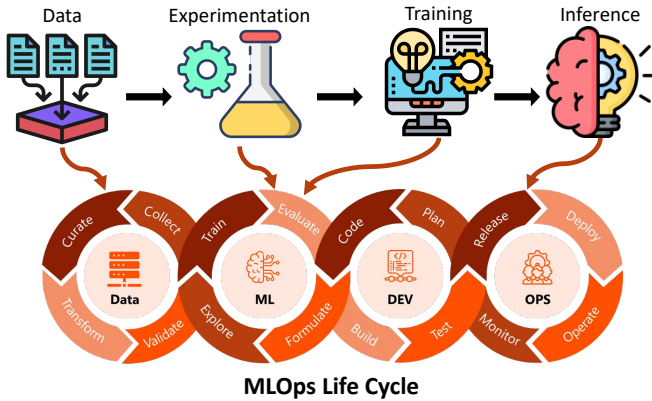


Fig. 1: ML model development and deployment phase and the associated MLOps life cycle.

Researchers have explored various energy reduction algorithms, e.g., pruning [14] or quantisation [15], etc. These works are smaller-scale investigations and focus on methods that affect the accuracy of a given model. On the contrary, in our large-scale study, we explore ways for energy reduction without changes in accuracy.

A notable work presenting various measurement campaigns is outlined in [16]. The authors focus on how various modifications in an ML pipeline can reduce the environmental impact, targeting system-level holistic optimisations. However, the individual measurements or the models used are not detailed. Our work studies a set of well-known models and datasets to enable readers to understand the differences between distinct hyperparameters and models. Moreover, open-sourcing our code will also enable other researchers to replicate our study with different models, datasets or hardware.

The recent literature includes two relevant studies to our evaluation based on real-world measurements [17], [18]. The authors of the first study [17] focused primarily on shallow single-layer models. The authors of the second study [18] investigated larger transformer-based models. However, neither includes a deep exploration of how different model characteristics or hyperparameters affect energy consumption. This is a key contribution of our paper.

III. GREEN MLOps: A STRATEGIC IMPERATIVE

DevOps merges software development with IT operations to speed up development time using automation and integration tools. MLOps, an extension of DevOps for ML pipelines, focuses on managing ML model lifecycles efficiently, tackling issues like data management and reproducibility. All production systems supporting ML-enabled immersive media applications will usually integrate some sort of MLOps framework [19]. Green MLOps extends the idea of MLOps, providing a framework that streamlines ML operations in an energy-aware and cost-effective fashion [12].

A. Energy Consumption in MLOps

MLOps (Fig. 1) usually comprises a **Data Processing** phase where data are collected, curated and labelled, and weights are applied to individual features based on their importance. This is followed by an **Experimentation** phase, where practitioners implement and evaluate potential algorithms, model architectures and training techniques. Various hyperparameter combinations are tested during this phase, most frequently in a grid-search fashion, to achieve reasonably robust functionality.

Once one / many solutions are determined as promising, during the **Training** phase, the chosen models are trained on extensive – larger quantity and feature-rich – datasets, “productising” them. Further hyperparameter tuning may be deemed necessary during this phase as well. Finally, when the model is ready, it is deployed in production, and the **Inference** phase starts. The model is fed real-time data and takes on-the-fly decisions. During this phase, the model’s performance is continuously monitored, and re-training cycles can be triggered when required.

Evidently [18], significant energy is consumed while training, experimenting or inferring on a model. Facebook’s AI research team [16] discusses that the total compute cycles for inference predictions exceed the corresponding training cycles, having a split of 10 : 20 : 70 (in %) between **Experimentation**, **Training** and **Inference**, respectively. Moreover, for the end-to-end pipeline, the energy footprint is roughly 31 : 29 : 40 (in %) for the **Data**, **Experimentation/Training**, and **Inference** phases. What is more, as described in [18], a poor hyperparameter tuning strategy can increase the total energy consumed by a typical Natural Language Model (NLP) by $\times 2000$ and a poorly managed neural network (NN) architecture by up to $\times 3000$ for a transformer-based NLP.

B. A Comparison between Computation and Data Exchange

In an MLOps pipeline, energy consumption involves computational resources and data exchange. Unlike typical ICT systems like video streaming, where computation and data transfer energy use are roughly equal [20], [21], ML pipelines are expected to demand more energy for computation. This is because while activities like **Experimentation** and **Training** may occur once (e.g., no retraining of the model is required), the **Inference** and **Data** phases will always need continuous computation. Additionally, trends like Federated Learning (FL), which distributes training or inference to edge nodes, show even higher energy consumption compared to centralised learning approaches (particularly considering complex ML models) despite the data exchange being decreased [22]. This is due to the difficulty of parallelising computation as training and inference are executed across a large number of clients on smaller datasets. Currently, no study compares the energy consumed by the computation against the data exchange in large-scale ML deployments (e.g., as the one presented from Facebook in [16]), highlighting a research gap.

The above motivated our investigation, i.e., to identify different model characteristics and hyperparameters that impact

the energy consumed within an ML pipeline. We focus our investigation on the **Experimentation, Training and Inference** phases of an MLOps pipeline. More specifically, we will compare parameters such as the model size, the batch size, the time required for training and inference, the Multiply–Accumulate (MAC) operations, the hardware utilisation and the model parameters as a function of the energy consumed.

IV. METHODOLOGY

To investigate the above, measuring the absolute power at frequent intervals and the time required for each experiment is essential. Hardware statistics like the utilisation of resources and the model characteristics should also be captured as part of our experimentation. We implemented a framework able to capture all the above and produce the results found in the paper. Our codebase can be found at github.com/ioannismavromatis/sustainable-ai.

A. Gathering Software-Based Energy Consumption Data

Monitoring energy consumption in computing environments can be achieved through hardware or software tools. Hardware methods offer precision [23] but face challenges in synchronisation and control [24], particularly for brief measurements like testing a shallow NN. These methods often require external clocks and costly equipment, limiting accessibility for all ML practitioners. Our investigation employs a software-based approach to measure energy consumption to overcome these issues. This choice not only reduces cost and complexity but also enhances consistency and scalability. Moreover, it allows for parallel evaluations of multiple devices and facilitates testing in complex scenarios, such as FL deployments.

In software-based measurements, power consumption is typically assessed in two ways. The first method estimates power based on a hardware component’s Thermal Design Power (TDP) and its utilisation (in a linear relationship). TDP, measured in Watts (W), indicates the maximum power consumption under theoretical full load. However, this method oversimplifies the relationship between power consumption and utilisation [25], as modern hardware can dynamically adjust the frequency and deactivate entire cores to save energy. A more nuanced approach is based on the hardware’s capacitance C , voltage V , and frequency f , as $P = 1/2 CV^2 f$, but obtaining these values for all components is rather challenging.

As a workaround, manufacturers offer a solution by providing access to energy data through Model Specific Registers (MSRs), like Nvidia’s Management Library (NVML) for GPUs and Intel’s Running Average Power Limit (RAPL) for CPU and DRAM usage. These methods are reliable with a reported variance of about ± 5 W in absolute values while following consistent trends in relative measurements [26], [27]. For consumer CPUs that MSRs do not provide DRAM measurements, DRAM’s energy consumption is approximated using $P_{\text{DRAM}} = \sum N_{\text{DIMM}} \times P_{\text{DIMM}}$, where N_{DIMM} is the number of DIMMs and $P_{\text{DIMM}} = 1/2 CV^2 f$. The operational V and f are accessible from the OS, and C is fixed for all our experiments. This equation is a good approximation as voltage

variations during DRAM operations are almost negligible, and operational frequency does not change over time [28].

B. Calculating Energy Usage in Machine Learning Processes

As discussed, our investigation will focus on the **Experimentation, Training and Inference** phases. **Training and Experimentation** phases are very similar (a model is trained using a set of preconfigured hyperparameters) and thus can be approached in a similar way in our investigation. To measure the energy consumption we define two metrics, i.e., E_{tr} , which is the total energy consumed during one training session (i.e., for a given model and dataset, with a pre-defined set of hyperparameters and a fixed number of epochs), and E_{in} , which is the total energy during inference (i.e., for a given model and dataset, inferring across all samples with a given batch size). They are as follows:

$$E_{\text{tr}} = \int_{t=0}^{T_{\text{tr}}} P_{\text{tr}}(t) dt - \int_{t=0}^{T_m} P_{\text{idle}}(t) dt \quad (1)$$

$$E_{\text{in}} = \int_{t=0}^{T_{\text{in}}} P_{\text{in}}(t) dt - \int_{t=0}^{T_m} P_{\text{idle}}(t) dt \quad (2)$$

where T_{tr} and T_{in} are the training and inference times, T_m is a hardcoded time interval used for the idle experiment, and P_{tr} , P_{in} and P_{idle} are the power measurements during training, testing and when the system is idle. Our framework captures the power consumption at frequent intervals Δt . Denoting t_i as the i -th time interval, the power $P(t_i)$ (this could be either for training or inference) is:

$$P(t_i) = P_{\text{CPU}}(t_i) + P_{\text{GPU}}(t_i) + P_{\text{DRAM}}(t_i) \quad (3)$$

where P_{CPU} , P_{GPU} and P_{DRAM} are the power consumption, taken in real-time for the CPU, GPU and DRAM, respectively. The energy within i -th interval can be calculated as the $E(t_i) = P(t_i) \Delta t$. Based on that, the Eqs. (1) and (2) can be approximated with the cumulative sum of all intervals, i.e.:

$$E_{\text{tr}} = \sum_{i=0}^{N_{\text{tr}}} P_{\text{tr}}(t_i) \Delta t - \sum_{t=0}^{N_m} P_{\text{idle}}(t_i) \Delta t \quad (4)$$

$$E_{\text{in}} = \sum_{t=0}^{N_{\text{in}}} P_{\text{in}}(t_i) \Delta t - \sum_{t=0}^{N_m} P_{\text{idle}}(t_i) \Delta t \quad (5)$$

where N_{tr} and N_{in} are the total number of intervals during training or inference, respectively. As discussed, data exchange and processing, even though they play a significant role in the energy consumed, will not be considered at this stage.

C. Hardware Stats and Model Characteristics

Our framework also collects utilisation statistics for all resources and the model characteristics. The NVML library provides the GPU (and its VRAM) utilisation. For the CPU, the utilisation metrics were directly collected from the OS as a function of each CPU core. The CPU utilisation is calculated as the average utilisation at a given time between all cores. Similarly, DRAM’s utilisation was also provided by the OS.

Concerning the model features, several key characteristics emerge as critical indicators when considering the operational efficiency of the model. These include the *model size*, the number of *total and trainable parameters*, *buffer size*, and *MACs*. The model size, measured when the model is decompressed and loaded in the VRAM, includes both the parameters and buffers and represents the overall footprint of the model in memory. The model size is measured in bytes (B).

The total number and the trainable parameters are key indicators of a model’s complexity. These parameters are different when layers in the model are frozen (i.e., are not updated). A larger number of parameters typically implies a more complex model, which can potentially achieve higher accuracy but at the cost of increased computational resources and memory usage. This complexity can lead to longer training times and may require more powerful hardware.

The buffer size indicates the additional data structures often used for storing intermediate outputs and constants that do not change during training, such as batch normalisation parameters. While they do not directly contribute to the learning capacity of the model, they impact the overall memory footprint. A large buffer size can lead to inefficiencies in systems with limited memory.

Finally, the MAC is the fundamental operation in NNs, especially in convolutional layers. The number of MACs provides an estimate of the computational complexity of the model. Higher MACs generally indicate increased computation for both training and inference, leading to longer processing times and increased energy consumption. All the above-mentioned model characteristics are calculated when the model is loaded in the GPU before the execution of each experiment. For our investigation, either independently or as a combination, these parameters will be explored towards total energy consumption.

V. RESULTS

For our investigation, we consider a simple image classification task. This task was chosen due to the ample models and datasets available in the literature. We conducted a thorough investigation to observe the behaviours of different ML models. In brackets, we present the model variant chosen for our experimentation. We picked: SimpleDLA, DPN (26), DenseNet (121), EfficientNet (B0), GoogLeNet, LeNet, MobileNet, MobileNetV2, PNASNet, PreActResNet (18), RegNet (X_200MF), ResNet (18), ResNeXt (29_2x64d), SENet (18), ShuffleNetV2, and VGG (16), to capture a diverse range of architectures and sizes. All experiments were conducted with the same hyperparameters (batch size of 128, learning rate 0.001, SGD optimiser, categorical cross-entropy loss and weight decay 5×10^{-4}). We also fixed the seed to ensure consistency across different runs.

We used three different Hardware Configurations (HCs) summarised in Tab. I. These three HCs provide diverse playgrounds to explore and identify their differences or similarities. As the space in the paper is limited, we will present a subset of the results and discuss the rest in the text. All results can be found in our GitHub repository for further analysis.

TABLE I: Hardware Configurations (HCs). In brackets is the TDP for each hardware component.

	HC-1	HC-2	HC-3
CPU*	i7-8700K (95 W)	i9-11900KF (125 W)	i5-12500 (65 W)
DRAM	4x16 GB DDR4 3600 MHz	4x32 GB DDR4 3200 MHz	2x16 GB DDR5 3200 MHz
GPU ¹	RTX 3080 (320 W) 10 GB	RTX 3090 (350 W) 24 GB	RTX A2000 (70 W) 12 GB

*Intel Core, ¹Nvidia driver v530.30.02, CUDA v12.1

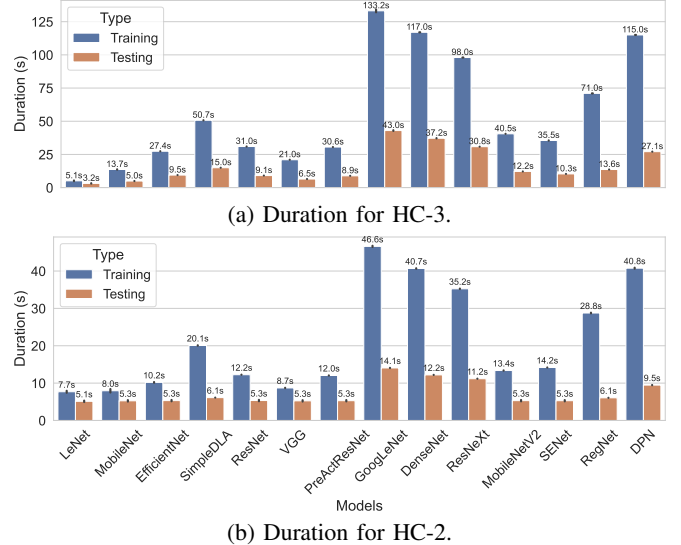


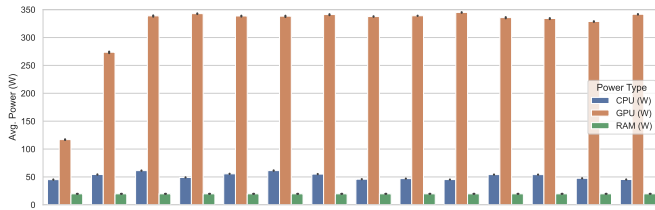
Fig. 2: Training and inference duration (for 50k samples).

Our investigation was based on the CIFAR-10 dataset [29]. The dataset consists of 60000 32×32 colour images in 10 classes, with 6000 images per class. The split between the training and testing set is 50000 : 10000. For our evaluation, we replicated the testing set 5 times (i.e., to 50k samples), so there are consistent samples between training and testing.

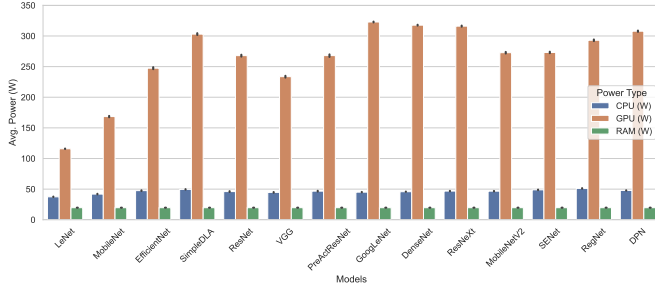
A. Initial Statistics

Starting with the maximum accuracy achieved, most models achieved around 87% – 91% after 100 epochs. The shallower LeNet underperformed as expected, reaching only around 68%, whereas MobileNet and EfficientNet reached 81% and 83%, respectively. Comparing the time required for training and inference (one epoch of training and 50k samples of inference), we see the results in Fig. 2. For most models, training takes three times longer than inference due to back-propagation and parameter updating. However, as seen, models like DPN, RegNet, etc. do not adhere to this rule of thumb. As shown, there is no direct correlation between the training and inference across different models, so each model should be independently investigated.

Comparing now HC-3 (Fig. 2a) with HC-2 (Fig. 2b), there is no observable difference during the training phase in terms of time required (relatively – between models). However, for the inference, we observe that a more powerful GPU (HC-2) parses the dataset in roughly equal time across most models. With the inference dictating the energy consumption (as discussed in Sec. III-A), as a rule of thumb, models



(a) Power usage by model (training).



(b) Power usage by model (inference).

Fig. 3: Average power usage with HC-2.

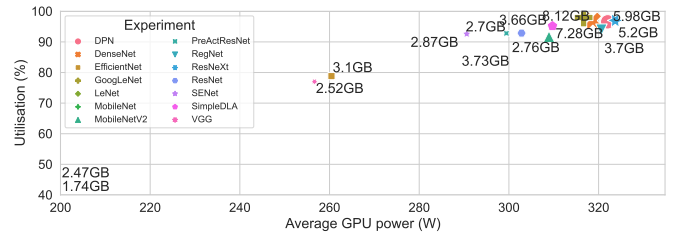
achieving similar accuracy but inferring quicker can introduce significant energy benefits in the long term, even if they require more time to train.

B. Power Consumption Measurements

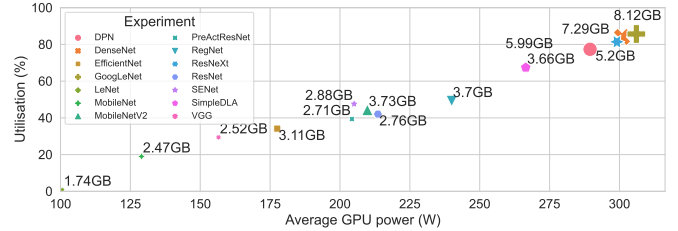
Fig. 3 shows the average power consumed for HC-2. A comparison between the training and inference phases is also shown. All the larger models force the GPU to operate close to its TDP (Fig. 3a). As expected, CPU and DRAM, being underutilised, show roughly equal and not significantly high average power consumption across all models. However, this is not the case for the inference (Fig. 3b). As shown, many models operate $\geq 30\%$ lower than the GPU’s TDP (e.g., VGG), with the CPU and DRAM showing similar results with the training. This is the case for the other two HCs, with the difference being more prominent for HC-1 and less prominent for HC-3.

Given that CPU and DRAM do not change drastically across different models, moving on, in Fig. 4, we present an example of the power consumption as a function of the utilisation and the usage of the GPU’s VRAM. For training, a larger GPU VRAM reflects, most of the time, higher utilisation and increased power consumption. This is even more prominent during inference. Moreover, the results indicate a high correlation between utilisation and power consumption, but up to a certain point. Beyond a power draw of ~ 300 W, any further increase did not translate to a dramatic increase in the GPU utilisation. At this point, utilisation was almost 100%, so performance was pretty much at its maximum. This is more clear in Fig. 4a, where, as said earlier, most models push the GPU to operate close to its TDP.

We observe a linear relationship by investigating the time and energy consumption (e.g., per epoch for training or per X number of samples for inference). Our results for that can be found in our repository. However, we observe the



(a) Power usage by model (training).



(b) Power usage by model (testing).

Fig. 4: Utilisation and power consumption (considering the GPU RAM usage) - HC-1.

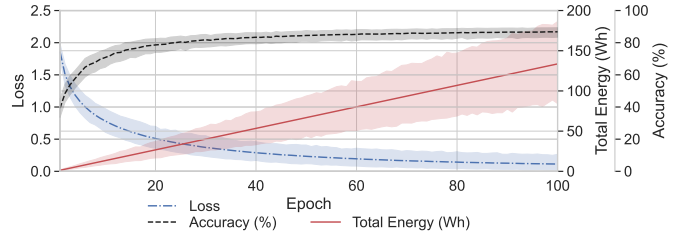


Fig. 5: Loss, energy and accuracy per epoch - HC-3.

following by comparing the model loss, accuracy, and total energy accumulated as the number of epochs increases while we train a model (Fig. 5). This figure shows the average across all models and the range of values for a given epoch. As observed, even though there is no correlation between accuracy and total energy consumed, as the number of epochs increases, the energy benefits that can be observed when replacing a model can significantly outnumber the change in accuracy.

Comparing the MACs of each model as a function of the total energy, we see a high correlation between them (results in our repository). An even higher correlation is shown if considering a new metric: the number of MACs per model parameter (Fig. 6). For both training (Fig. 6a) and inference (Fig. 6b), we see a strong correlation across them.

Finally, we compare the batch size for training and inference (Fig. 7). As expected, smaller batch sizes increase the power consumption. There is a direct correlation with the GPU utilisation for each model. For all setups, there is a batch size that minimises the power consumption, with no further improvements shown if the batch size is increased.

VI. DISCUSSION

The previous section presented a subset of our results from our extensive investigation across several models and hardware setups. This section will summarise our findings and critically comment on them. Starting with our initial observations (Sec. V-A), it is clear that each model’s unique

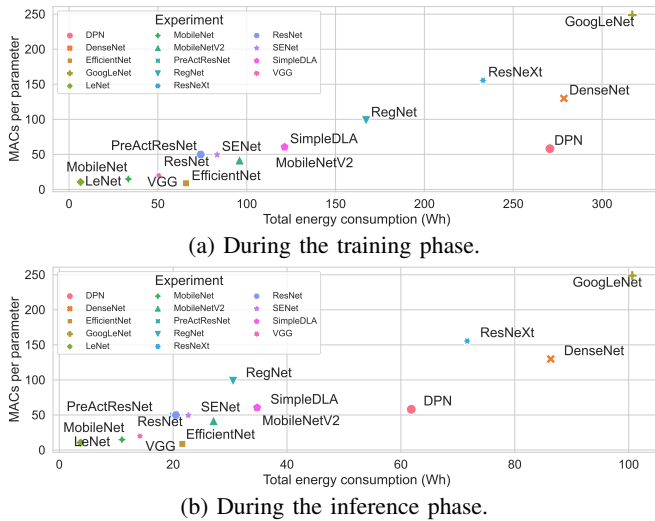


Fig. 6: Total energy consumption as a function of the MACs per parameter - HC-3.

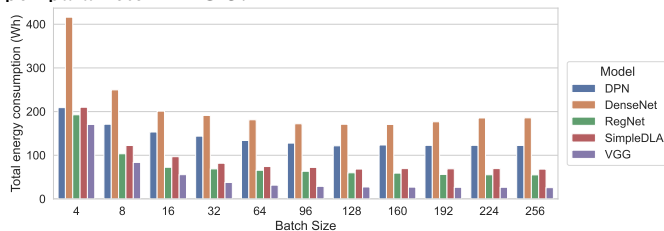


Fig. 7: Total energy consumption as a function of the batch size - HC-2.

architecture does not allow room for cross-model observations (i.e., if a model’s energy consumption is low, there is no obvious way to say that another model will have an equally low energy consumption). Further investigation of the specific architectures and model layers and how they affect energy consumption could identify more similarities. However, it was outside this work’s scope and can be considered in the future.

Broadly speaking (Fig. 5), benefits in energy reduction can outperform the gains in accuracy on many occasions. Moreover, training and inference durations are not correlated; thus, cross-phase or cross-hardware estimations are not promising. Even though a rule of thumb could say that training will require three times more time for the same number of samples, this is not always the case.

Considering that time and total energy are linear, short-living profilers (e.g., training for one epoch or inferring for a small number of samples) can be used to extrapolate the total energy for larger scenarios. In addition, for accuracy and duration, it was evident that models achieving comparable accuracy but “running faster” can introduce huge energy benefits in the long term. From what was shown in Fig. 3 and taking into account Facebook’s energy split presented in Sec. III-A, a less power-hungry model during inference should be prioritised for a pipeline over a less energy-intensive model during training. This observation can be combined with the total energy and time to get even more accurate estimations for both training and inference across different models.

As shown in Fig. 4, hardware devices’ power profiles are

not exactly linear. Usually, manufacturers push their devices to the limit to squeeze a narrow increase in performance. Smart strategies like the one introduced in [30] (power capping optimisations) can exploit that and significantly reduce the total energy consumed. Finally, if an estimation of the model’s expected energy is required, contrary to the literature that proposes using the model’s MACs (with Spearman correlation of 0.8 – HC-3), we identified the MACs per model parameter as a more suitable candidate (with Spearman correlation of 0.9 – HC-3). Similar correlations are observed across all setups, proving that as a uniform solution for estimating the expected energy consumption with high accuracy.

VII. CONCLUSIONS

This work presented an extensive analysis across multiple ML models and hardware setups to uncover techniques for improving sustainability without sacrificing effectiveness. The investigation methodology combined software-based power measurements with tracking of hardware utilisation and model characteristics. The experiments demonstrated that for many models, reductions in energy consumption can outpace marginal accuracy improvements, highlighting the need to balance performance and efficiency. Additionally, assumptions about energy use cannot be reliably made across training and inference or across hardware due to a lack of consistent correlations. However, normalising model MACs by the number of parameters provides an excellent indicator of energy consumption in most cases. The insights from this study can guide decisions when constructing ML pipelines, whether choosing architectures and hyperparameters or provisioning hardware resources. There remains ample opportunity for future work to further improve sustainability through novel architectures optimised for efficiency and adoption of best practices around selective retraining, power capping, and inference-focused model selection. Overall, the evidence clearly shows that with careful planning, ML can continue advancing while aligning with environmental responsibility.

ACKNOWLEDGMENT

This work is a contribution by Project REASON, a UK Government funded project under the Future Open Networks Research Challenge (FONRC) sponsored by the Department of Science Innovation and Technology (DSIT). This work was also funded in part by Toshiba Europe Ltd. and Bristol Research and Innovation Laboratory (BRIL).

REFERENCES

- [1] M.-A. Moïnnereau, A. A. de Oliveira, and T. H. Falk, "Immersive Media Experience: A Survey of Existing Methods and Tools for Human Influential Factors Assessment," *Quality and User Experience*, vol. 7, no. 1, p. 5, Jun 2022.
- [2] A. Doumanoglou, D. Griffin, J. Serrano *et al.*, "Quality of Experience for 3-D Immersive Media Streaming," *IEEE Trans. Multimedia.*, vol. 64, no. 2, pp. 379–391, 2018.
- [3] J. Teo and J. T. Chia, "Deep Neural Classifiers For Eeg-Based Emotion Recognition In Immersive Environments," in *Proc. of Int. Conf. on ICSCCE*, jul 2018, pp. 1–6.
- [4] V. Vasilev, J. Leguay, S. Paris *et al.*, "Predicting QoE Factors with Machine Learning," in *Proc. of IEEE ICC*, May 2018, pp. 1–6.
- [5] A. Feldmann, O. Gasser, F. Lichtblau *et al.*, "Implications of the COVID-19 Pandemic on the Internet Traffic," in *Proc. of ACM Int. Meas. Conf.*, Oct. 2021.
- [6] Greenpeace, "Clicking Clean: A Guide to Building the Green Internet," 2015. [Online]. Available: <https://www.greenpeace.org/static/planet4-international-stateless/2015/05/153e0823-2015clickingclean.pdf>
- [7] R. Kumar, S. K. Gupta, H.-C. Wang *et al.*, "From Efficiency to Sustainability: Exploring the Potential of 6G for a Greener Future," *Sustainability*, vol. 15, no. 23, 2023.
- [8] D. Buragohain, S. Chaudhary, G. Pungeng *et al.*, "Analyzing the Impact and Prospects of Metaverse in Learning Environments Through Systematic and Case Study Research," *IEEE Access*, vol. 11, pp. 141 261–141 276, Dec. 2023.
- [9] A. Luccioni, A. Lacoste, and V. Schmidt, "Estimating Carbon Emissions of Artificial Intelligence [Opinion]," *IEEE Technol. Soc. Mag.*, vol. 39, no. 2, pp. 48–51, Jun. 2020.
- [10] A. Katsenou, J. Mao, and I. Mavromatis, "Energy-Rate-Quality Trade-offs of State-of-the-Art Video Codecs," in *Proc. of PCS*, Dec. 2022.
- [11] D. Patterson, J. Gonzalez, U. Hölzle *et al.*, "The Carbon Footprint of Machine Learning Training Will Plateau, Then Shrink," *Computer*, vol. 55, no. 7, pp. 18–28, Jun. 2022.
- [12] R. Schwartz, J. Dodge, N. A. Smith *et al.*, "Green AI," *Commun. ACM*, vol. 63, no. 12, p. 54–63, nov 2020.
- [13] R. Verdecchia, J. Sallou, and L. Cruz, "A Systematic Review of Green AI," *WIREs Data Mining and Knowledge Discovery*, vol. 13, no. 4, p. e1507, 2023.
- [14] T.-J. Yang, Y.-H. Chen, and V. Sze, "Designing Energy-Efficient Convolutional Neural Networks Using Energy-Aware Pruning," in *Proc. of IEEE CVPR*, Jul. 2017, pp. 6071–6079.
- [15] N. S. Eliezer, R. Banner, H. Ben-Yaakov *et al.*, "Power Awareness In Low Precision Neural Networks," in *Proc. of ECCV Workshop*, Feb. 2023, p. 67–83.
- [16] C.-J. Wu, R. Raghavendra, U. Gupta *et al.*, "Sustainable AI: Environmental Implications, Challenges and Opportunities," in *Proc. of MLS*, vol. 4, Oct. 2022, pp. 795–813.
- [17] M. S. Islam, S. N. Zisad, A.-L. Kor *et al.*, "Sustainability of Machine Learning Models: An Energy Consumption Centric Evaluation," in *Proc. of ECCE*, Feb. 2023, pp. 1–6.
- [18] E. Strubell, A. Ganesh, and A. McCallum, "Energy and Policy Considerations for Modern Deep Learning Research," in *Proc. of AAAI*, vol. 34, no. 09, Apr. 2020, pp. 13 693–13 696.
- [19] M. Testi, M. Ballabio, E. Frontoni *et al.*, "MLOps: A Taxonomy and a Methodology," *IEEE Access*, vol. 10, pp. 63 606–63 618, jun 2022.
- [20] E. Gelenbe, "Electricity Consumption by ICT: Facts, Trends, and Measurements," *Ubiquity*, vol. 2023, no. August, Aug. 2023.
- [21] Carbon Trust, "Carbon Impact of Video Streaming," Tech. Rep., 2021.
- [22] X. Qiu, T. Parcollet, J. Fernandez-Marques *et al.*, "A First Look into the Carbon Footprint of Federated Learning," *J. of MLR*, vol. 24, no. 129, pp. 1–23, 2023.
- [23] G. Conti, D. Jimenez, A. del Rio *et al.*, "A Multi-Port Hardware Energy Meter System for Data Centers and Server Farms Monitoring," *Sensors*, vol. 23, no. 1, Dec. 2023.
- [24] S. Rinaldi, F. Bonafini, P. Ferrari *et al.*, "Software-based Time Synchronization for Integrating Power Hardware in the Loop Emulation in IEEE1588 Power Profile Testbed," in *Proc. of IEEE ISPCS*, Sep. 2019, pp. 1–6.
- [25] W. Lin, T. Yu, C. Gao *et al.*, "A Hardware-aware CPU Power Measurement Based on the Power-exponent Function model for Cloud Servers," *Information Sciences*, vol. 547, pp. 1045–1065, Oct. 2021.
- [26] NVIDIA Corporation, "nvidia-smi.txt," 7 2016. [Online]. Available: <https://developer.download.nvidia.com/compute/DCGM/docs/nvidia-smi-367.38.pdf>
- [27] H. David, E. Gorbatov, U. R. Hanebutte *et al.*, "RAPL: Memory power estimation and capping," in *Proc. of ACM/IEEE ISLPED*, Aug. 2010, pp. 189–194.
- [28] T. Vogelsang, "Understanding the Energy Consumption of Dynamic Random Access Memories," in *Proc. of IEEE/ACM MICRO*, Dec. 2010, pp. 363–374.
- [29] A. Krizhevsky, "Learning Multiple Layers of Features from Tiny Images," Apr. 2009. [Online]. Available: <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>
- [30] I. Mavromatis, S. De Feo, P. Carnelli *et al.*, "FROST: Towards Energy-efficient AI-on-5G Platforms - A GPU Power Capping Evaluation," in *Proc. of IEEE CSCN*, Nov. 2023.