

Energy-Aware Machine Learning Algorithm Design

Dheeraj Vaddepally

dheeraj.vaddepally@gmail.com

Abstract

The exponential increase in machine learning (ML) use on mobile and edge devices indicated a necessity to adopt efficient algorithm design to conserve energy for future consumption and sustainability. Power reduction for energy-constrained platforms like smartphones, Internet of Things devices, and autonomous cars, at training and inference, is critical of importance. This book discusses design techniques for energy-conscious machine learning algorithms, specifically CPU and GPU energy profiling and reducing the power usage with techniques. Profiling techniques and tools are discussed to find out the energy requirements of various algorithms, and model pruning, quantization, knowledge distillation, and low-precision inference are discussed for minimizing inference power usage. For training, efficient backpropagation, energy-conscious optimizers, and distributed training are taken into account. The work also discusses energy efficiency-performance trade-offs and the promise of energy-aware NAS and dynamic resource management. The influence of energy-aware algorithm design is shown through examples of mobile and IoT device, edge computing, and data center applications. Last but not least, hardware constraints and scalability issues are presented, and future directions for designing more energy-efficient ML systems are provided.

Keywords: Energy-Aware Algorithms, Machine Learning, CPU/GPU Profiling, Power Consumption, Inference Optimization, Training Efficiency, Model Pruning, Quantization, Low-Power AI, Resource-Constrained Devices

I. INTRODUCTION

The increasing prevalence of mobile, IoT, and embedded devices has generated the need for energy-efficient machine learning algorithms. These energy-consuming devices are increasingly becoming crucial building blocks across numerous applications, from real-time analytics to autonomous systems and edge computing. Machine learning models bring unprecedented capabilities, but they come with equally gigantic energy demands, particularly during training and inference. As these devices are battery-limited, it has become a necessity to make sure that machine learning algorithms are designed keeping energy efficiency in mind.

Machine learning models, especially large-scale deep learning models, are computationally expensive to train and execute, requiring large resources for both training and inference. On low-resource devices, high power usage can quickly drain battery life and limit the capability of the device. With the model's complexity increasing, the energy required also increases proportionally. Thus, it is important to create machine learning algorithms that are not only accurate and reliable but also energy sensitive.[1]

The focus of this paper is to talk about techniques for CPU/GPU usage and power optimization profiling of algorithms during training and inference stages of machine learning models. These techniques include advanced methods such as hardware profiling, model pruning, and quantization that can be utilized to minimize energy usage without degrading accuracy or performance. It is essential to understand how various design choices impact energy consumption to deploy machine learning models in real-time and low-power applications.[2]

II. BACKGROUND

Machine learning models typically have a few energy-intensive phases such as data pre-processing, model training, and inference. Training is very well-known to be the energy-intensive phase because it involves learning from large amounts of data in an iterative manner. The operation would likely involve various passes over data, gradient computation, and optimization steps for models that need heavy computational loads, particularly when dealing with deep neural networks. The process of inference, although relatively energy-low-intensive, is energy-intensive in the nature that it actually makes real-time prediction or classification of data streams, often with energy-constrained systems.[2]

An important amount of research has been accomplished in reducing the energy expense of machine learning before this. Hardware optimization approaches such as offloading to GPU and low-energy CPU architecture have been studied for the purpose of saving power. Model compression methods such as pruning, quantization, and distillation have also been created to compress model size and computational requirements of machine learning models. These techniques attempt to achieve a model performance vs. efficiency trade-off in a bid to make machine learning models deployable in a more scalable manner in energy-constrained environments. Apart from that, a number of ways of optimization to the software include adjusting learning rates and batch size in efforts to lower energy consumption cost as per model training.

Machine learning operations call for energy profiling with the goal to identify and get accustomed to learning operation usage and consumption behaviors. Profiling instruments can be used to monitor the amount of energy a given CPU or GPU operation consumes, and this allows scientists to identify stages where optimization is critical. This can translate into more effective algorithmic fine-tuning that can result in brain-bending power savings. Profiling information can also inform energy-aware algorithms by identifying the energy cost of given hyperparameters, model architectures, and data processing methods. Overall, energy profiling is a central method to construct machine learning algorithms that are performance-optimized and energy-optimized, especially to be run on embedded and mobile platforms.[3]

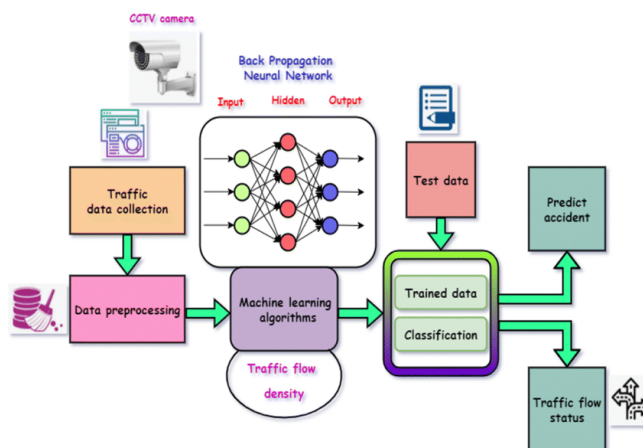


Fig. 1. Machine Learning Algorithm Design

III. PROFILING ALGORITHMS FOR CPU/GPU USAGE

Profiling is the systematic quantification of resource utilization—time, memory, and energy—during the execution of machine learning algorithms. Profiling in the context of energy-conscious machine learning entails measuring the energy expenditure of models at different stages, including data preprocessing, training, and inference. Profiling's main objective is to locate bottlenecks and inefficiencies that lead to excessive power usage and resource consumption. Profiling allows developers to see the energy consumption patterns of their models on various hardware platforms (e.g., CPU versus GPU) and design to balance performance and energy efficiency.

For mobile and embedded devices that have limited power supplies (such as batteries), profiling becomes essential. These devices tend to execute real-time applications, so energy minimization at the cost of model performance is a key design necessity. Profiling enables quantifying the energy expense of different model elements—such as forward passes, matrix multiplications, or memory access patterns—to guide specific optimizations aimed at saving power without sacrificing the overall effectiveness of the model.[3]

A. Energy Profiling Tools and Techniques

Several tools and techniques exist for energy profiling the power consumption of machine learning models. These tools allow for quantifying power usage on both GPUs and CPUs, enabling insights into energy consumption at the hardware level. Some of the most popularly used profiling tools are:

- **PyJoules:** A Python tool that is used to quantify the energy usage of different system components (e.g., CPU, GPU, RAM) in a fine-grained fashion. PyJoules can be incorporated into machine learning pipelines to track power consumption at different phases of the algorithm.

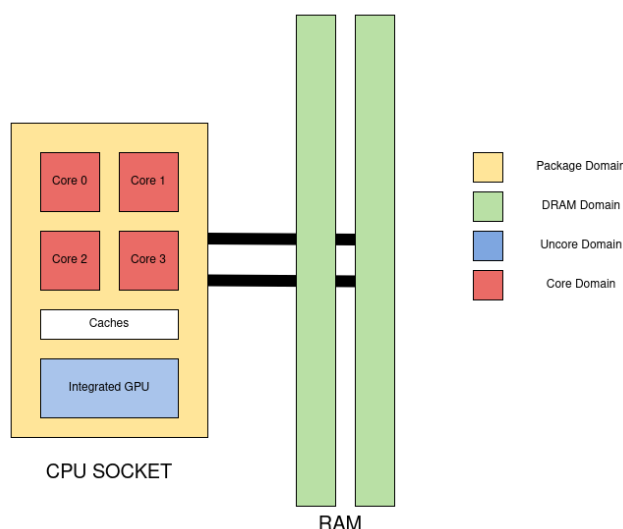


Fig. 2. Pyjoules

- **PowerAPI:** An open-source framework that facilitates energy monitoring and power modeling. It offers readings on the power usage of individual applications and hardware devices, providing detailed insights into energy consumption patterns.
- **NVIDIA's NVML (NVIDIA Management Library):** A monitoring API that offers GPU metrics like power usage, temperature, and utilization in real time. It's popularly used for profiling GPU-intensive tasks, including training deep learning models.
- **Intel VTune:** An Intel performance profiling tool used by developers to detect performance bottlenecks and energy inefficiencies in CPU-intensive programs. VTune offers granular CPU energy consumption information to fine-tune CPU-based machine learning workloads.[4]

B. CPU vs. GPU Power Consumption

Machine learning models have a very different energy consumption profile based on whether they are run on CPUs or GPUs. CPUs are built for generic computing operations with flexibility and accuracy but at a loss of efficiency in terms of energy when executed for large-scale parallel computations. GPUs, being built for parallel processing, fit best for resource-intensive operations such as matrix arithmetic in deep models. GPUs take more energy, however, at the overall level, especially training, because they are very highly powered and intensive in their processing and usage.

For conventional machine learning algorithms (e.g., decision trees, k-nearest neighbors), CPUs tend to be more energy-efficient than GPUs. These algorithms tend to have fewer matrix operations and are more amenable to serial execution, so CPU use is more power-efficient. Deep learning models—particularly those with large neural networks—take advantage of GPU parallelism, but at increased energy expense. Profiling these models can determine if a change from GPU to CPU (or vice versa) might provide improved energy efficiency based on the task complexity and target hardware.

C. Case Study: Profiling Energy Consumption for CPU and GPU

Suppose a convolutional neural network (CNN) model is trained on an image classification task. Upon modeling the energy draw of this model on CPU as well as on GPU through PyJoules and NVML, outputs could be such that GPU training initially uses more power due to it having greater draw but, based on its capacity to execute this task at lightning speed compared to CPUs, ultimately might have an even lower (power * time) total for the energy. On the other hand, while inferring, the CPU can provide more energy efficiency, particularly for models with smaller size and lower complexity. [4]

IV. TECHNIQUES TO REDUCE INFERENCE POWER CONSUMPTION

A. Model Pruning and Quantization

Model pruning and quantization are popular techniques that are used in an attempt to compress the size and computations of the neural networks and, in turn, decrease the energy used at inference. Model pruning involves the elimination of the redundant or smaller weights of the neural network. Pruning reduces the computational burden for inference by reducing the active connections and weights, thereby saving energy. For example, structured pruning removes entire filters or neurons entirely, reducing the workload of operations for CPUs and GPUs heavily but not heavily impacting accuracy.

Quantization is another effective technique that projects high-precision floating-point computation (e.g., FP32) onto lower-precision, less accurate representations (e.g., FP16 or INT8). This reduces memory and computation, especially during inference. With low-bit representation utilization, quantization reduces the energy required for multiply-accumulate operations—one of the largest drivers of inference energy cost in deep learning. Quantization is a post-training technique and thus a good technique for energy conservation in pre-trained models running on low-end devices.

B. Knowledge Distillation

Knowledge distillation is a method by which a smaller, but better-performing model (the "student model") is trained to mimic a larger, but more complex model (the "teacher model"). The smaller model is easy to deploy and uses a lot less computation during inference time, conserving energy. Knowledge distillation proves useful when used on mobile phones with minimal computing power and where it would not be practical to use a large model. Knowledge distillation enables the smaller model to learn from the larger model's output probabilities or feature representations and have the same accuracy level while consuming less energy.[5]

C. Low-Precision Inference

Low-precision inference pushes energy efficiency to the next level by constraining computation to fewer bits. Low-precision arithmetic methods such as FP16 (16-bit floating-point) and INT8 (8-bit integer) have the potential to reduce power consumption by an enormous factor without sacrificing much accuracy. The majority of contemporary machine learning hardware accelerators such as NVIDIA's Tensor Cores and Google's TPUs already support low-precision arithmetic for speeding up and lowering power consumption for inference. Low-precision arithmetic reduces the memory access burden and the arithmetic unit burden, therefore minimizing the power to perform inference on embedded or mobile platforms.

D. Algorithm Optimization

Optimization of the algorithm itself is another technique to reduce power consumption during inference. Early-exit architecture provides an opportunity to reduce computation by allowing the model to early-exit inference after becoming certain about the prediction. This reduces the need for full forward pass throughout the entire network, saving computational power and resources. Another optimization technique is to design thin neural network structures such as MobileNet, EfficientNet, or SqueezeNet that are particularly designed for mobile phones. These models are optimized to be light, efficient, and low-power but accuracy comparable.[6]

V. BALANCING ENERGY EFFICIENCY AND PERFORMANCE

A. Trade-off between Model Accuracy and Energy

In model building to obtain proper machine learning and energy use, the most common issue is sacrificing model quality for energy use. Lowering energy use has the secondary effect of lowering model complexity, performing lower accuracy calculations, or eliminating redundant parameters. This may be at the cost of model quality or precision, though. The problem is how to do so while keeping energy usage to a minimum without sacrificing too much in terms of prediction quality. Small models, for instance, are energy efficient, but the model size would be something to cap in order to account for complexity within data patterns, i.e., sacrificed accuracy.

B. Energy-Efficient Model Architecture Search (NAS)

Neural Architecture Search (NAS) is an algorithmic procedure to find the most effective neural network structures from a pre-specified set of constraints, i.e., energy usage.

Energy-efficient NAS allows one to discover model architectures that balance energy consumption, computational expense, and accuracy. The method applies search algorithms to explore an exponentially large space of potential architectures, testing each in terms of energy consumption and accuracy. Incorporating energy efficiency as a factor, NAS can assist in designing models that are resource-efficient for resource-constrained settings like in embedded and mobile systems without sacrificing performance. [10]

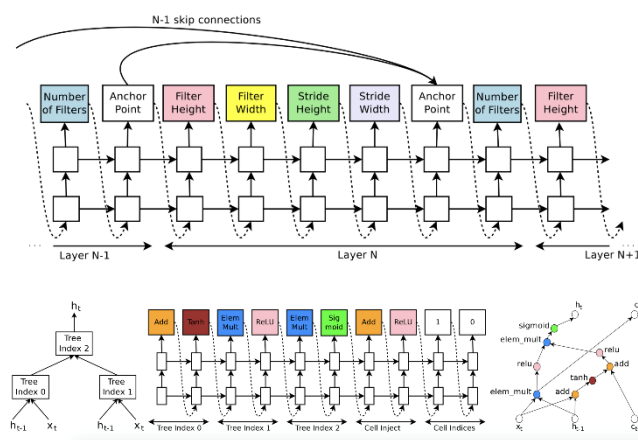


Fig. 3. NAS

C. Dynamic Resource Allocation

Dynamic resource allocation are methods that allocate computing resources as a function of workload and available energy.

For instance, the active number of cores or the clock rate of a processor can be reduced by machine learning algorithms when the computational demand is low, hence conserving energy. Dynamic workload distribution between GPUs and CPUs as a function of available energy or real-time power cap is possible. This efficient approach is used for optimizing energy consumption under maintaining the performance constraints, particularly in energy-constrained platforms like mobile apps, IoT, or edge computing. [7]

VI. USE CASES AND APPLICATIONS

A. Mobile and IoT Devices

Energy-conscious machine learning is very important in wearables, IoT devices, and mobile phones as they are power-limited. Machine learning models deployed on them need to walk a tight rope between performing inference in real time and saving energy. For example, models for applications such as speech recognition, image classification, or monitoring a person's health need to be lightweight and low-power. Methods such as model quantization, pruning, and edge computing make it possible for machine learning platforms in these systems to provide low power consumption-based efficient insights. [9]

B. Data Centers and Cloud-Based ML

Low-power machine learning methods are the pillars in big cloud systems and data centers in providing lower operating costs and environmental effects.

Since artificial intelligence workloads are increasingly utilized within data centers for deployment in recommendation systems, natural language applications, and training of deep neural networks at scale, power consumption increases exponentially. Hardware accelerators and energy-efficient algorithms, workload consolidation and hardware optimization and end-of-line high-end cooling hardware, can decouple the end-to-end machine learning process of the cloud by lowering overall machine learning energy consumption. [7]

C. Edge Computing and Autonomous Systems

Power-hungry machine learning processes draw power to a significant extent in edge computing environments where computation power is distributed close to sources of data (e.g., sensors, cameras) rather than centralizing it in the cloud servers.

The algorithms are especially critical in autonomous systems, such as drones or autonomous vehicles, which must make real-time decisions using a limited power supply. Power efficiency in this case directly influences the system's lifespan and performance. Power-efficient algorithms provide more run time and faster response, making them the key to the success of autonomous and edge-based applications. [8]

VII. CONCLUSION

This paper covered some of the most important power profiling techniques and power reduction techniques for machine learning models. We described how energy profiling tools such as PyJoules and NVML are needed to achieve the energy usage of models run on GPUs as well as CPUs. Techniques such as model pruning, quantization, knowledge distillation, and low-precision inference were highlighted as good means of decreasing power consumption at inference time. We also discussed

performance vs. energy efficiency trade-offs in the models, the NAS role for optimizing the models, and power consumption and workload control dynamic resource management techniques.

As machine learning models are being pushed to a wide variety of devices that range from cloud infrastructure to mobile and IoT at unprecedented scales, energy efficiency is increasingly becoming a focus area of interest.

The need for models that can perform real-time computation without draining resources is paramount in battery-constrained and resource-constrained systems. Through the application of energy-efficient design techniques in algorithms, engineers are able to embed machine learning products with efficiency and scalability in a way that allows them to be used across a variety of applications from data centers to edge computing.

In the long term, continued emphasis on energy-efficient model design will have consequences well down into the sustainability of AI systems.

As more machine learning is incorporated into common technology, the ability to construct models that are low-power consuming as well as highly performing will become crucial to minimizing the environmental cost of AI deployments. Low power consumption AI also has the ability to unlock new innovation in autonomous systems, smart cities, and low-power wearables. By making energy-awareness center stage, next-generation AI can be more environmentally friendly, scalable, and accessible on a broad base of platforms and industries.

REFERENCES

- [1] Berral, J. L., Goiri, I., Nou, R., Julià, F., Fitó, J. O., Guitart, J., ... & Torres, J. (2012). *Toward Energy-Aware Scheduling Using Machine Learning*. *Energy-Efficient Distributed Computing Systems*, 215-244.
- [2] Hribar, J., Marinescu, A., Chiumento, A., & DaSilva, L. A. (2021). Energy-aware deep reinforcement learning scheduling for sensors correlated in time and space. *IEEE Internet of Things Journal*, 9(9), 6732-6744.
- [3] Li, R., Gong, W., Wang, L., Lu, C., & Dong, C. (2023). *Co-evolution with deep reinforcement learning for energy-aware distributed heterogeneous flexible job shop scheduling*. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 54(1), 201-211.
- [4] Lazzaro, D., Cinà, A. E., Pintor, M., Demontis, A., Biggio, B., Roli, F., & Pelillo, M. (2023, September). *Minimizing energy consumption of deep learning models by energy-aware training*. In *International Conference on Image Analysis and Processing* (pp. 515-526). Cham: Springer Nature Switzerland.
- [5] Hoffmann, J. L. C., & Fröhlich, A. A. (2021). *Online machine learning for energy-aware multicore real-time embedded systems*. *IEEE Transactions on Computers*, 71(2), 493-505.
- [6] Turkkan, B. O., Dai, T., Raman, A., Kosar, T., Chen, C., Bulut, M. F., ... & Sow, D. (2022, June). *GreenABR: Energy-aware adaptive bitrate streaming with deep reinforcement learning*. In *Proceedings of the 13th ACM Multimedia Systems Conference* (pp. 150-163).
- [7] Yang, T. J., Chen, Y. H., & Sze, V. (2017). *Designing energy-efficient convolutional neural networks using energy-aware pruning*. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 5687-5695).

- [8] *Janjani, H., Agarwal, T., Gopinath, M. P., Sharma, V., & Raja, S. P. (2024). Designing Energy-Aware Scheduling and Task Allocation Algorithms for Online Reinforcement Learning Applications in Cloud Environments. IEEE Transactions on Computational Social Systems.*
- [9] *Xu, C., Wang, K., Li, P., Xia, R., Guo, S., & Guo, M. (2018). Renewable energy-aware big data analytics in geo-distributed data centers with reinforcement learning. IEEE Transactions on Network Science and Engineering, 7(1), 205-215.*
- [10] *Jafarzadeh, S. Z., & Moghaddam, M. H. Y. (2014, October). Design of energy-aware QoS routing algorithm in wireless sensor networks using reinforcement learning. In 2014 4th International Conference on Computer and Knowledge Engineering (ICCCKE) (pp. 722-727). IEEE.*