

# Focus Area: Computational Efficiency

ProjectX 2023

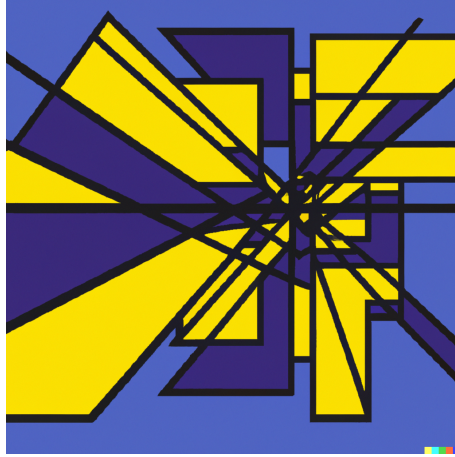


Figure 1: "Computational Efficiency in the form of modern art" - DALL-E 2

## Introduction

With the rapid rise in Artificial Intelligence (AI) innovations within the last decade, a major hurdle keeping AI from solving the world's problems is the computational cost of training models. With recent advancements in computational capabilities, AI has seen an unprecedented amount of focus for research and development in the industry.

However, the current methods of model training often result in prolonged training times, excessive energy consumption, and inefficient hardware utilization [MIT, ]. As such, the central challenge still remains - to improve the computational efficiency of AI model training, making the process more resource-efficient, cost-effective, and environmentally sustainable. The key issues identified are as follows:

- **Training Time:** Current AI model training processes can take days or even weeks to complete, leading to delays in research, development, and deployment of new models and applications.
- **Energy Consumption:** The energy consumed during AI model training contributes to a significant carbon footprint and operational costs.
- **Hardware Utilization:** Inefficient use of hardware resources such as GPUs and TPUs can result in underutilization, wasting potential processing power and increasing training time.
- **Scalability:** As AI models grow in complexity and size, data centers struggle to scale their infrastructure to accommodate the training needs, resulting in performance bottlenecks.

These improvements can be through either algorithms or hardware. An algorithmic approach might include finetuning a model's parameters while a hardware approach may include the use of AI specific hardware to perform specialized tasks. Together, these improvements unlock the potential for AI adaptations in a variety of fields, regardless of the resource constraints [Hernandez and Brown, 2020].

## Background

As technology continues to evolve, the vast influx of data and the increasing utilization of AI across various domains have highlighted the need for improved computational efficiency. By optimizing the allocation and utilization of resources, AI systems can operate more effectively, accelerating model training, reducing inference latency, and enhancing scalability. This optimization extends beyond raw computational power, encompassing efficient algorithms, model compression, and parallel processing strategies [Deering, 1984]. Diverse strategies are being employed to achieve this goal.

As you conduct your research, you'll find that one prominent avenue for improving computational efficiency is the development and optimization of specialized hardware. New developments in AI specialized hardware played a major role in creating the Machine Learning landscape of today. For instance, the creation of Application-Specific Integrated Circuits (ASICs) and Tensor Processing Units (TPUs) tailored to the demands of AI tasks drastically accelerated computations while minimizing power consumption [Jouppi et al., 2017]. As shown in Figure 2, ASIC's provide a significant increase in computation speed as compared to CPU's.

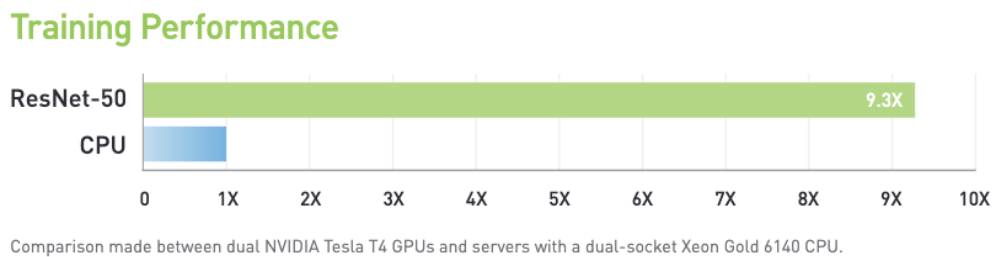


Figure 2: NVIDIA Tesla T4 vs Xeon Gold CPU [noa, ]

Furthermore, architectural improvements within AI systems are yielding substantial gains in efficiency. Techniques such as knowledge distillation involve training compact models to replicate the behavior of larger, more resource-intensive models. This knowledge transfer not only reduces the computational requirements for inference but also aids in compressing models for deployment on resource-constrained devices like smartphones and Internet of Things (IoT) devices. Moreover, advancements in algorithms and model designs can significantly cut down on memory and processing requirements without substantially compromising performance.

The importance of computational efficiency expands beyond just speed. Such systems also consume less power, leading to reduced carbon emissions to improve sustainability. Efficient AI computation enables real-time applications such as autonomous vehicles, medical diagnostics, and industrial automation, where low latency and rapid decision-making are critical [Schwartz et al., 2020].

In the dynamic landscape of evolving technology and burgeoning AI integration, the pursuit of computational efficiency stands as an essential guiding progress. The synthesis of hardware innovation and algorithmic optimization collectively forges a path toward better AI systems. This journey towards efficiency is not only a journey of speed and power conservation but also a journey of sustainability and real-world applicability. As we stride ahead, the fusion of these strategies promises to reshape industries, advance transformative applications, and pave the way for a future where intelligent systems operate seamlessly, swiftly, and sustainably in harmony with the world around us.

# Focus Area: Data Efficiency

ProjectX 2023

## Introduction

In the realm of artificial intelligence (AI), the effective utilization of data plays a pivotal role in achieving accurate, reliable, and scalable AI models. However, the current AI landscape is plagued by challenges related to data inefficiency, hindering the development and deployment of AI systems.

The overarching goal of addressing data efficiency in AI is to advance the field by:

- **Reducing Data Dependency:** Developing AI models that require fewer data samples for training while improving performance.
- **Optimizing Resource Usage:** Designing techniques to minimize memory, computational, and energy requirements during both training and deployment.
- **Improving Generalization:** Enhancing model generalization by efficiently leveraging diverse datasets and learning from limited samples.
- **Enabling Edge AI:** Facilitating the deployment of AI on edge devices with limited resources, enabling real-time, on-device processing.
- **Enhancing Adaptability:** Creating AI systems that can adapt efficiently to new data distributions and tasks, enabling lifelong learning.

## Background

At its heart, the pursuit of data optimization seeks to find the delicate balance between the computational resources required and the performance of the model. This challenge stems not only from practical concerns such as reducing energy consumption and training time but also from economic and environmental considerations. Effective resource management has the potential to democratize AI by making it more accessible, encouraging innovation, and facilitating the integration of AI systems into a wide array of fields.

Training models on large amounts of data is costly and may not always lead to the most accurate or precise model. In fact, too much data can result in overtraining in which the model is no longer generalizable to solve other problems.

One of the primary methods to optimize training data is distillation. Dataset distillation involves transforming a large dataset into a smaller synthetic dataset to reduce the required memory by aggregating the most important information [noa, 2021]. Such distillation can be done through a variety of algorithms and is particularly useful when dealing with large amounts of data that may have unnecessary detail. Concepts like active learning describe the idea of the model “asking” for the information they need to improve performance [noa, 2022]. In this sense, the model will only require the essential data it needs to perform its task. The focus is then shifted onto the quality rather than the quantity of data. Thus, having metrics to grade data quality is important.

Evaluating data efficiency involves considering two fundamental elements: quantity and quality.

## Quantity

Quantity pertains to either the memory footprint occupied by a dataset or the number of observations it encompasses. A larger quantity of data can be beneficial for training robust models and enhancing their generalization capabilities.

## Quality

On the other hand, the second key metric, quality, involves aspects such as completeness and accuracy within the dataset. A high-quality dataset is characterized by consistent, complete values with minimal noise. Inaccurate or missing values can introduce bias and significantly impact a model's performance. Data noise, referring to random errors or variations, can mislead models by detecting patterns that don't truly exist, leading to erroneous predictions. To mitigate these challenges, data preprocessing is essential. Techniques such as imputation and exclusion can address missing values, while identifying and managing outliers, along with appropriate transformations, can reduce data noise.

Ensuring data accuracy and consistency requires validation through cross-referencing reliable sources and incorporating human validation. It's not uncommon for datasets to be outdated, repetitive, inconsistent, or poorly organized, necessitating preprocessing before they are usable.

Notably, data quality and quantity are interconnected, and their effective management is pivotal for successful model development and deployment.

As an illustrative example, recent research conducted at Cornell University highlights the significance of data quality and quantity. This research introduced a new large language model with notable size reductions compared to competitors, achieved through the utilization of "textbook quality" data obtained from the web and synthetically generated textbooks and exercises, as demonstrated in their work on GPT-3.5 [Gunasekar et al., 2023]. Table 1 compares their model (phi-1) against other prominent language models today, highlighting their competitiveness with a fraction of a dataset size.

Model Name	Model Size	Dataset Size	HumanEval	MBPP
<b>phi-1 [Gunasekar et al., 2023]</b>	<b>1.3B</b>	<b>7B</b>	<b>50.6%</b>	<b>55.5%</b>
WizardCoder [Luo et al., 2023]	16B	1T	57.3%	51.8%
GPT-3.5 [OpenAI, 2023]	175B	N.A.	47%	-
StarCoder [Li et al., 2023]	15.5B	1T	33.6%	52.7%

Table 1: phi-1 evaluated against other prominent language models.

## References

- [MIT, ] The computing power needed to train AI is now rising seven times faster than ever before. Available at <https://tinyurl.com/bdc3wnocr>.
- [noa, ] Nvidia t4 tensor core gpus for accelerating inference. Available at <https://www.nvidia.com/en-us/data-center/tesla-t4/>.
- [noa, 2021] (2021). Training machine learning models more efficiently with dataset distillation. Available at <https://ai.googleblog.com/2021/12/training-machine-learning-models-more.html>.
- [noa, 2022] (2022). How to train ml models more efficiently with active learning. Available at <https://venturebeat.com/ai/how-to-train-ml-models-more-efficiently-with-active-learning/>.
- [Deering, 1984] Deering, M. F. (1984). Hardware and software architectures for efficient AI. In *Proceedings of the Fourth AAAI Conference on Artificial Intelligence, AAAI'84*, pages 73–78, Austin, Texas. AAAI Press.
- [Gunasekar et al., 2023] Gunasekar, S., Zhang, Y., Aneja, J., Mendes, C. C. T., Del Giorno, A., Gopi, S., Javaheripi, M., Kauffmann, P., de Rosa, G., Saarikivi, O., Salim, A., Shah, S., Behl, H. S., Wang, X., Bubeck, S., Eldan, R., Kalai, A. T., Lee, Y. T., and Li, Y. (2023). Textbooks are all you need. arXiv:2306.11644 [cs].
- [Hernandez and Brown, 2020] Hernandez, D. and Brown, T. B. (2020). Measuring the algorithmic efficiency of neural networks. arXiv:2005.04305 [cs, stat].
- [Jouppi et al., 2017] Jouppi, N. P., Young, C., Patil, N., Patterson, D., Agrawal, G., Bajwa, R., Bates, S., Bhatia, S., Boden, N., Borchers, A., Boyle, R., Cantin, P.-l., Chao, C., Clark, C., Coriell, J., Daley, M., Dau, M., Dean, J., Gelb, B., Ghaemmaghami, T. V., Gottipati, R., Gulland, W., Hagmann, R., Ho, C. R., Hogberg, D., Hu, J., Hundt, R., Hurt, D., Ibarz, J., Jaffey, A., Jaworski, A., Kaplan, A., Khaitan, H., Killebrew, D., Koch, A., Kumar, N., Lacy, S., Laudon, J., Law, J., Le, D., Leary, C., Liu, Z., Lucke, K., Lundin, A., MacKean, G., Maggiore, A., Mahony, M., Miller, K., Nagarajan, R., Narayanaswami, R., Ni, R., Nix, K., Norrie, T., Omernick, M., Penukonda, N., Phelps, A., Ross, J., Ross, M., Salek, A., Samadiani, E., Severn, C., Sizikov, G., Snelham, M., Souter, J., Steinberg, D., Swing, A., Tan, M., Thorson, G., Tian, B., Toma, H., Tuttle, E., Vasudevan, V., Walter, R., Wang, W., Wilcox, E., and Yoon, D. H. (2017). In-datacenter performance analysis of a tensor processing unit. In *Proceedings of the 44th Annual International Symposium on Computer Architecture*, pages 1–12, Toronto ON Canada. ACM.
- [Li et al., 2023] Li, R., Allal, L. B., Zi, Y., Muennighoff, N., Kocetkov, D., Mou, C., Marone, M., Akiki, C., Li, J., Chim, J., Liu, Q., Zheltonozhskii, E., Zhuo, T. Y., Wang, T., Dehaene, O., Davaadorj, M., Lamy-Poirier, J., Monteiro, J., Shliazhko, O., Gontier, N., Meade, N., Zebaze, A., Yee, M.-H., Umapathi, L. K., Zhu, J., Lipkin, B., Oblokulov, M., Wang, Z., Murthy, R., Stillerman, J., Patel, S. S., Abulkhanov, D., Zocca, M., Dey, M., Zhang, Z., Fahmy, N., Bhattacharyya, U., Yu, W., Singh, S., Luccioni, S., Villegas, P., Kunakov, M., Zhdanov, F., Romero, M., Lee, T., Timor, N., Ding, J., Schlesinger, C., Schoelkopf, H., Ebert, J., Dao, T., Mishra, M., Gu, A., Robinson, J., Anderson, C. J., Dolan-Gavitt, B., Contractor, D., Reddy, S., Fried, D., Bahdanau, D., Jernite, Y., Ferrandis, C. M., Hughes, S., Wolf, T., Guha, A., von Werra, L., and de Vries, H. (2023). Starcoder: may the source be with you! arXiv:2305.06161 [cs].
- [Luo et al., 2023] Luo, Z., Xu, C., Zhao, P., Sun, Q., Geng, X., Hu, W., Tao, C., Ma, J., Lin, Q., and Jiang, D. (2023). Wizardcoder: empowering code large language models with evol-instruct. arXiv:2306.08568 [cs].
- [OpenAI, 2023] OpenAI (2023). Gpt-4 technical report. arXiv:2303.08774 [cs].
- [Schwartz et al., 2020] Schwartz, R., Dodge, J., Smith, N. A., and Etzioni, O. (2020). Green ai. *Communications of the ACM*, 63(12):54–63.