

11. Green AI

Sustainable Software Engineering
CS4295



Luís Cruz
L.Cruz@tudelft.nl

1. Overview of Green AI
2. Green data-centric AI
3. Green AI at Facebook
4. Tim Yarally on Green AI

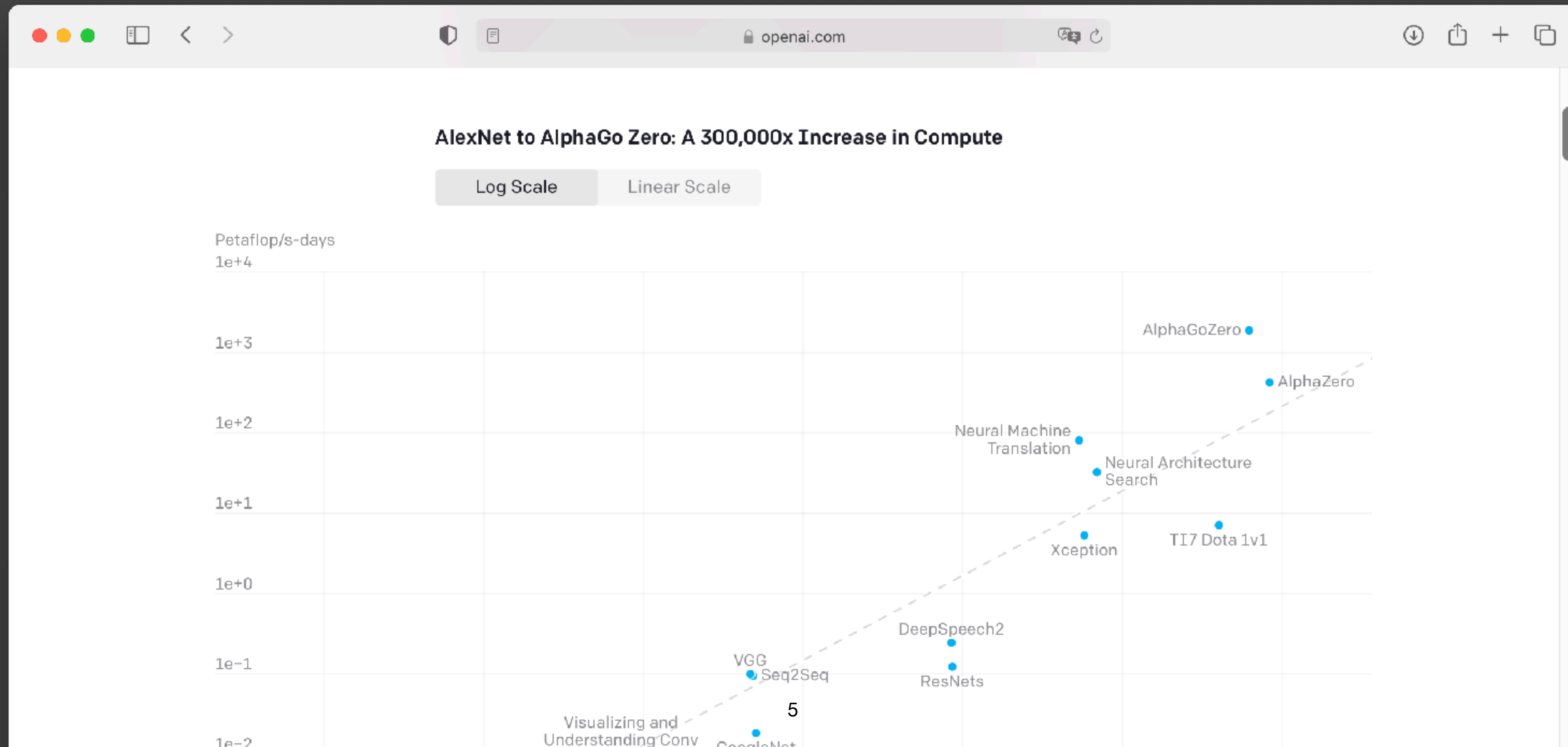
Tomorrow's class

AI

- Artificial Intelligence (AI) is the branch of computer science that deals with **automating** tasks that typically require **human intelligence**.
- In the past years AI has been widely applied across different domains. Health care.
- To deploy AI systems, we test them against **benchmarks** (or validation sets).
 - The goal is to outperform the previous existing models.
 - E.g., in Machine Learning we usually resort to accuracy metrics. The highest the accuracy, the better the model.

Since 2012, the amount of computing used for AI training **has been doubling every 3.4 months**

- <https://openai.com/blog/ai-and-compute/>



- To create better AI systems we are currently adding
 - **More data**
 - **More experiments**
 - **Larger models**

The Equation of Red AI

$$Cost(R) \propto E \cdot D \cdot H$$

Cost of a single (**E**)xample

Number of (**H**)yperparameters

Size of (**D**)ataset

Issues of Red AI

- High costs (hardware, electricity, data access, etc.)
- Limited reproducibility.
- Energy consumption.
- Carbon emissions.
- SMEs can hardly be competitive.
- Groundbreaking AI research is mostly done by tech giants.

A few examples of Red AI

- Google's BERT-large
 - 350 million features
 - Trained for 2.5 days using 512 TPU chips, costing \$60K+
- Open-GPT3
 - 175 billion features
- AlphaGo
 - 1920 CPUs, 280 GPUs, costing \$35M

Red AI



Accuracy: 0.9999999999

Green AI



- Energy
- Time
- Reproducibility
- Reusage

Research on Green AI

- Most literature revolves around **position papers**. (?)

On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?

Emily M. Bender*
ebender@uw.edu
University of Washington
Seattle, WA, USA

Timnit Gebru*
timnit@blackinai.org
Black in AI
Palo Alto, CA, USA

Angelina McMillan-Major
aymm@uw.edu
University of Washington
Seattle, WA, USA

Shmargaret Shmitchell
shmargaret.shmitchell@gmail.com
The Author

ABSTRACT

The past 3 years of work in NLP have been characterized by the development and deployment of ever larger language models, especially for English. BERT, its variants, GPT-2/3, and others, most recently Switch-C, have pushed the boundaries of the possible both through architectural innovations and through sheer size. Using these pretrained models and the methodology of fine-tuning them for specific tasks, researchers have extended the state of the art on a wide array of tasks as measured by leaderboards on specific benchmarks for English. In this paper, we take a step back and ask: How big is too big? What are the possible risks associated with this technology and what paths are available for mitigating those risks? We provide recommendations including weighing the environmental and financial costs first, investing resources into curating and carefully documenting datasets rather than ingesting everything on the web, carrying out pre-development exercises evaluating how the planned approach fits into research and development goals and supports stakeholder values, and encouraging research directions beyond ever larger language models.

CCS CONCEPTS

- Computing methodologies → Natural language processing

ACM Reference Format:
Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? . In *Conference on Fairness, Accountability, and Transparency*.

Energy and Policy Considerations for Deep Learning in NLP

Emma Strubell Ananya Ganesh Andrew McCallum
College of Information and Computer Sciences
University of Massachusetts Amherst
{strubell, aganesh, mcallum}@cs.umass.edu

Abstract

Recent progress in hardware and methodology for training neural networks has ushered in a new generation of large networks trained on abundant data. These models have obtained notable gains in accuracy across many NLP tasks. However, these accuracy improvements depend on the availability of exceptionally large computational resources that necessitate similarly substantial energy consumption. As a result these models are costly to train and develop, both financially, due to the cost of hardware and electricity or cloud compute time, and environmentally, due to the carbon footprint required to fuel modern tensor processing hardware. In this paper we bring this issue to the attention of NLP researchers by quantifying the approximate financial and environmental costs of training a variety of recently successful neural network models for NLP. Based on these findings, we propose actionable recommendations to reduce costs and improve equity in NLP research and practice.

Introduction

Advances in techniques and hardware for training deep neural networks have recently enabled impressive accuracy improvements across any fundamental NLP tasks (Bahdanau et al., 2015; Luong et al., 2015; Dzout and Manning, 2017; Vaswani et al., 2017), with the most computationally-hungry models obtaining the highest scores (Peters et al., 2018; Devlin et al., 2019; Radford et al., 2019; So et al., 2019). As a result, training a state-of-the-art model now requires substantial computational resources which demand considerable energy, along with the associated financial and environmental costs. Re-

Consumption	CO ₂ e (lbs)
Air travel, 1 passenger, NY→SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000

Training one model (GPU)	CO ₂ e (lbs)
NLP pipeline (parsing, SRL) w/ tuning & experimentation	39
Transformer (big) w/ neural architecture search	78,468
	192
	626,155

Table 1: Estimated CO₂ emissions from training common NLP models, compared to familiar consumption.¹

NLP models could be trained and developed on a commodity laptop or server, many now require multiple instances of specialized hardware such as GPUs or TPUs, therefore limiting access to these highly accurate models on the basis of finances.

Even when these expensive computational resources are available, model training also incurs a substantial cost to the environment due to the energy required to power this hardware for weeks or months at a time. Though some of this energy may come from renewable or carbon credit-offset resources, the high energy demands of these models are still a concern since (1) energy is not currently derived from carbon-neutral sources in many locations, and (2) when renewable energy is available, it is still limited to the equipment we have to produce and store it, and energy spent training a neural network might better be allocated to heating a family's home. It is estimated that we must cut carbon emissions by half over the next decade to deter escalating rates of natural disaster, and based

Green AI

contributed articles

DOI:10.1145/3381831

Creating efficiency in AI research will decrease its carbon footprint and increase its inclusivity as deep learning study should not require the deepest pockets.

BY ROY SCHWARTZ, JESSE DODGE, EN ETZIONI

This trend is driven by the strong focus of the AI community on obtaining “state-of-the-art” results,² as exemplified by the popularity of leaderboards,^{3,4} which typically report accuracy (or other similar measures) but omit any mention of cost or efficiency (see, for example, leaderboards.allenai.org).⁵ Despite the clear benefits of improving model accuracy, the focus on this single metric ignores the economic, environmental, and social cost of reaching the reported results.

We advocate increasing research activity in Green AI—AI research that is more environmentally friendly and inclusive. We emphasize that Red AI research has been yielding valuable scientific contributions to the field, but it has been overly dominant. We want to shift the balance toward the Green AI option—to ensure any inspired undergraduate with a laptop has the opportunity to write high-quality papers that could be accepted at premier research conferences. Specifically, we propose making efficiency a more common evaluation criterion for AI papers alongside accuracy and related measures.

key insights

- The computational costs of state-of-the-art AI research has increased 300,000x in recent years. This trend, denoted **Red AI**, stems from the AI community's focus on accuracy while paying attention to efficiency.
- Red AI leads to a surprisingly large carbon footprint, and makes it difficult for academics, students, and researchers to engage in deep learning research.
- An alternative is Green AI, which treats efficiency as a primary evaluation criterion alongside accuracy. To measure efficiency, we suggest reporting the number of floating-point operations required to generate a result.
- Green AI research will decrease AI's environmental footprint and increase its inclusivity.

of artificial intelligence (AI) has progress on a broad range of object recognition, game playing, and machine translation.¹³ Much of this progress has been achieved by increasingly large and intensive deep learning models.⁴ Figure 1 plots training energy for state-of-the-art deep learning models. AlexNet in 2012²⁴ to AlphaZero in 2017²⁵ represents an overall increase of 300,000x, or roughly doubling every few months. An even more recent example, observed in NLP word-embedding models at ELMo²⁴ followed by BERT,⁹ Megatron-LM,⁴² T5,³⁶ and GPT-3,⁴ has estimated the carbon footprint of training such models as Red AI.

In this article, but our focus is on AI research that relies on deep

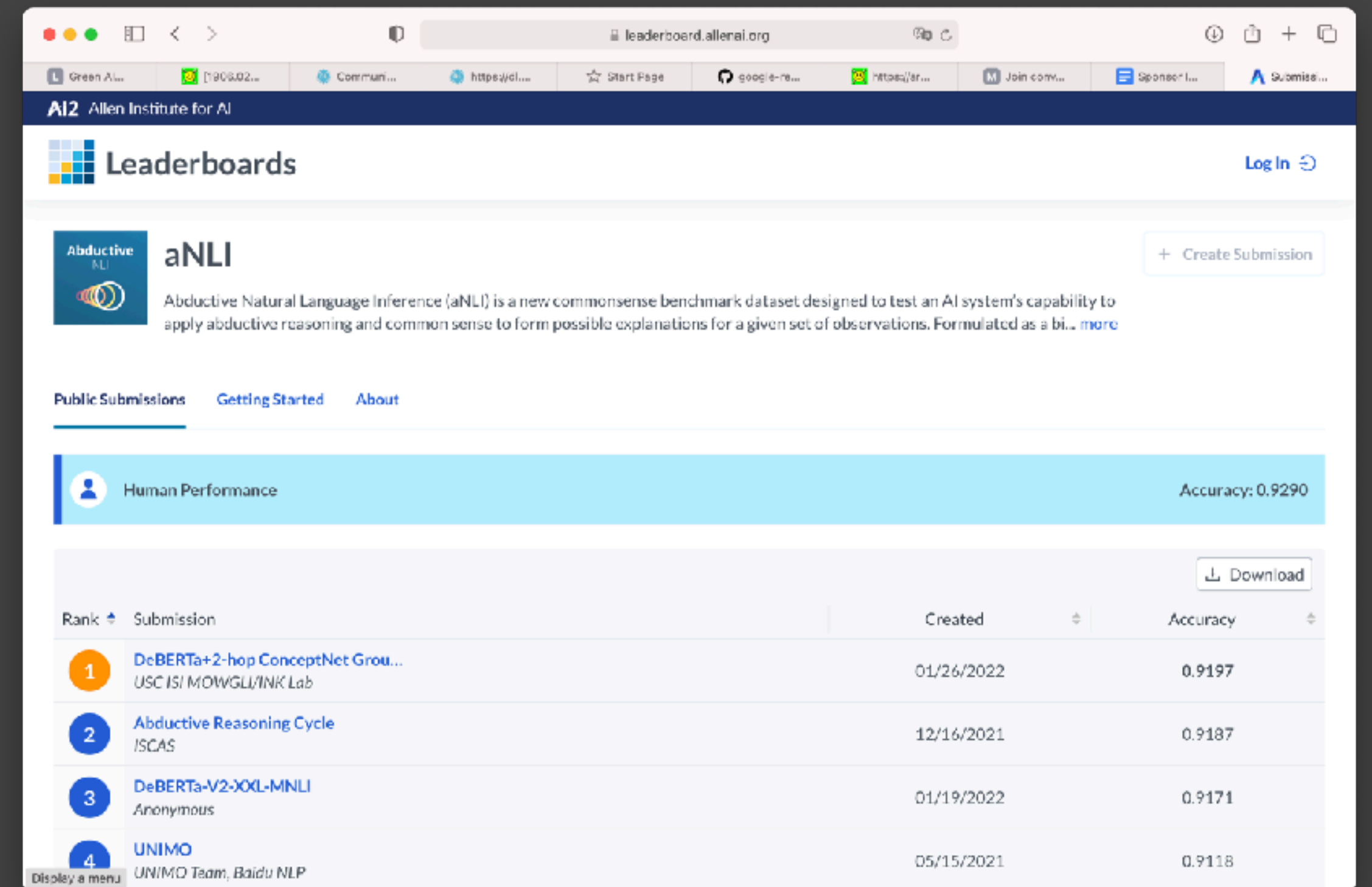
DECEMBER 2020 | VOL. 03 | NO. 12

How can we adopt **Green AI**

- **Check whether AI is needed.**
- Select green datacenters.
- Run on **low carbon intensity** hours.
- Opt for **GPU-optimised** solutions
- Opt for **low-power hardware** (e.g., Nvidia Jetson boards)
 - Or GPUs that provide energy metrics (e.g., NVIDIA GPUs via the **nvidia-smi** tool)
- **Report** energy/carbon metrics (e.g., embed in MLFlow?)
- Use pre-trained models (Transfer Learning)
- Preprocess dataset to reduce size.
- Improve parameter-tuning strategy.

Reporting energy/carbon footprint

- We need **benchmarks**.
- AllenAI leaderboard <https://leaderboard.allenai.org>
 - **No carbon metrics**, yet
- Report comparable proxies for energy consumption.
 - **!** Learning algorithms behave in a non-deterministic
 - **!** Different data-points lead to different energy consumption

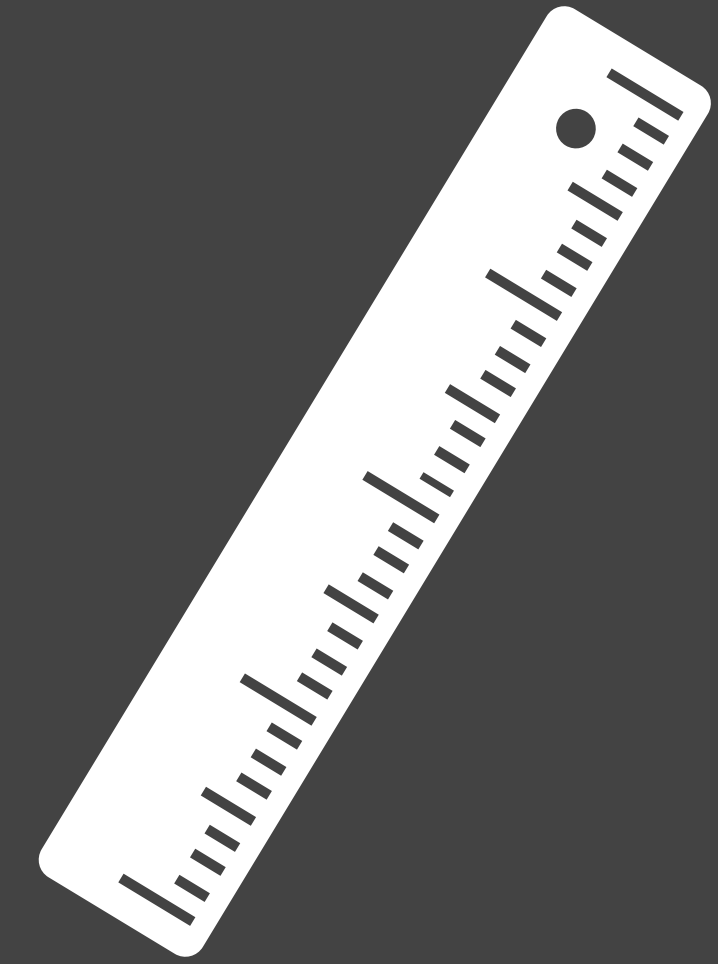


The screenshot shows the AllenAI Leaderboards website for the aNLI benchmark. The page displays the top human performance and a list of public submissions. The human performance is 0.9290. The top submissions are:

Rank	Submission	Created	Accuracy
1	DeBERTa+2-hop ConceptNet Group USC ISI MOWGLI/INK Lab	01/26/2022	0.9197
2	Abductive Reasoning Cycle ISCAS	12/16/2021	0.9187
3	DeBERTa-V2-XXL-MNLI Anonymous	01/19/2022	0.9171
4	UNIMO UNIMO Team, Baidu NLP	05/15/2021	0.9118

Reporting energy/carbon footprint

- Reporting **measured energy consumption**
 - + Accurate
 - + Easy to map to carbon emissions
 - - Hard to measure
- Reporting **time** / estimation based on **time & hardware**
 - + Easy to measure
 - + Correlates with energy consumption in most cases.
 - - Difficult to compare with measurements from other setups
- E.g., **floating point operations** (FPOs) (?)
 - + comparable across different setups
 - + cheap
 - - does not factor in energy consumption in memory
 - - does not reflect carbon emissions

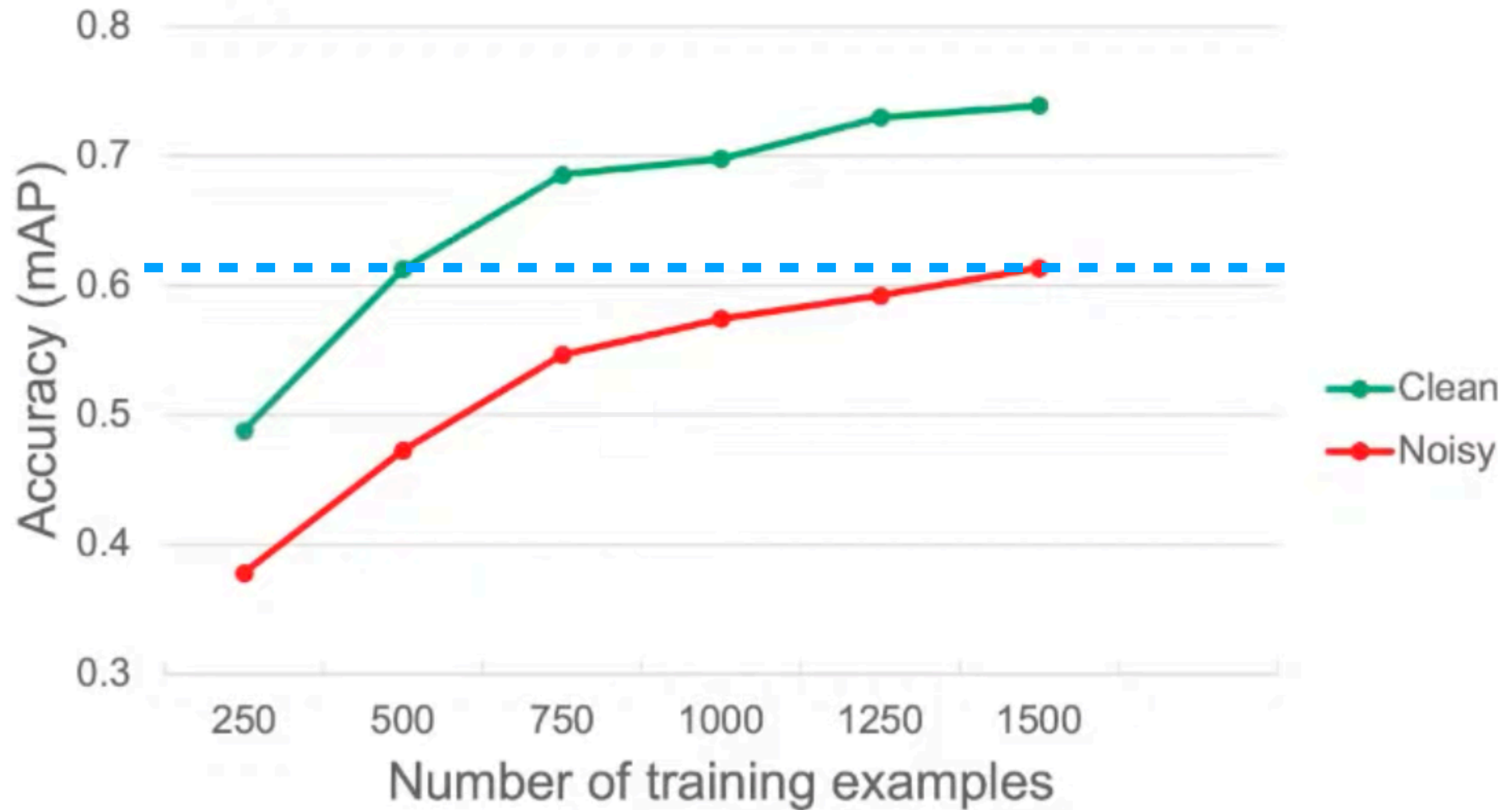


Data-centric AI

Data-centric AI

- Emerging discipline that deals with systematically engineering data to build AI systems.
 - Shift from **improving the training strategy** to **improving the data**.
 - It is better to have **small but reliable** datasets than **large but noisy** datasets.
 - => Improve **data collection**, **data labelling**, and **data preprocessing**.
- More about data-centric AI by Andrew Ng:
<https://www.youtube.com/watch?v=06-AZXmwHjo>

Example: Clean vs. noisy data



Green Data-centric AI

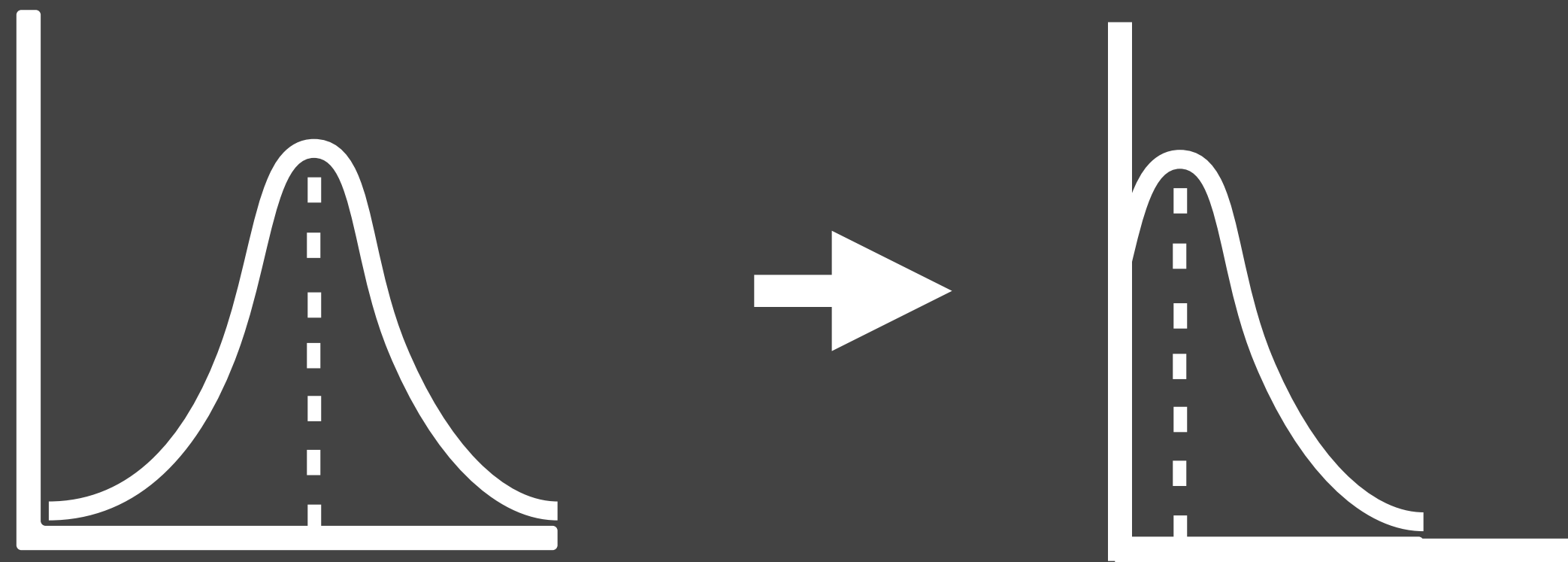
- How do different ML algorithms compare in terms of energy consumption?
- How does **number of rows** relate to the energy consumption of ML models?
- How does **number of features** relate to the energy consumption of ML models?
- What is the impact of reducing data in the **performance** of the model?
- Method -> results -> discussion



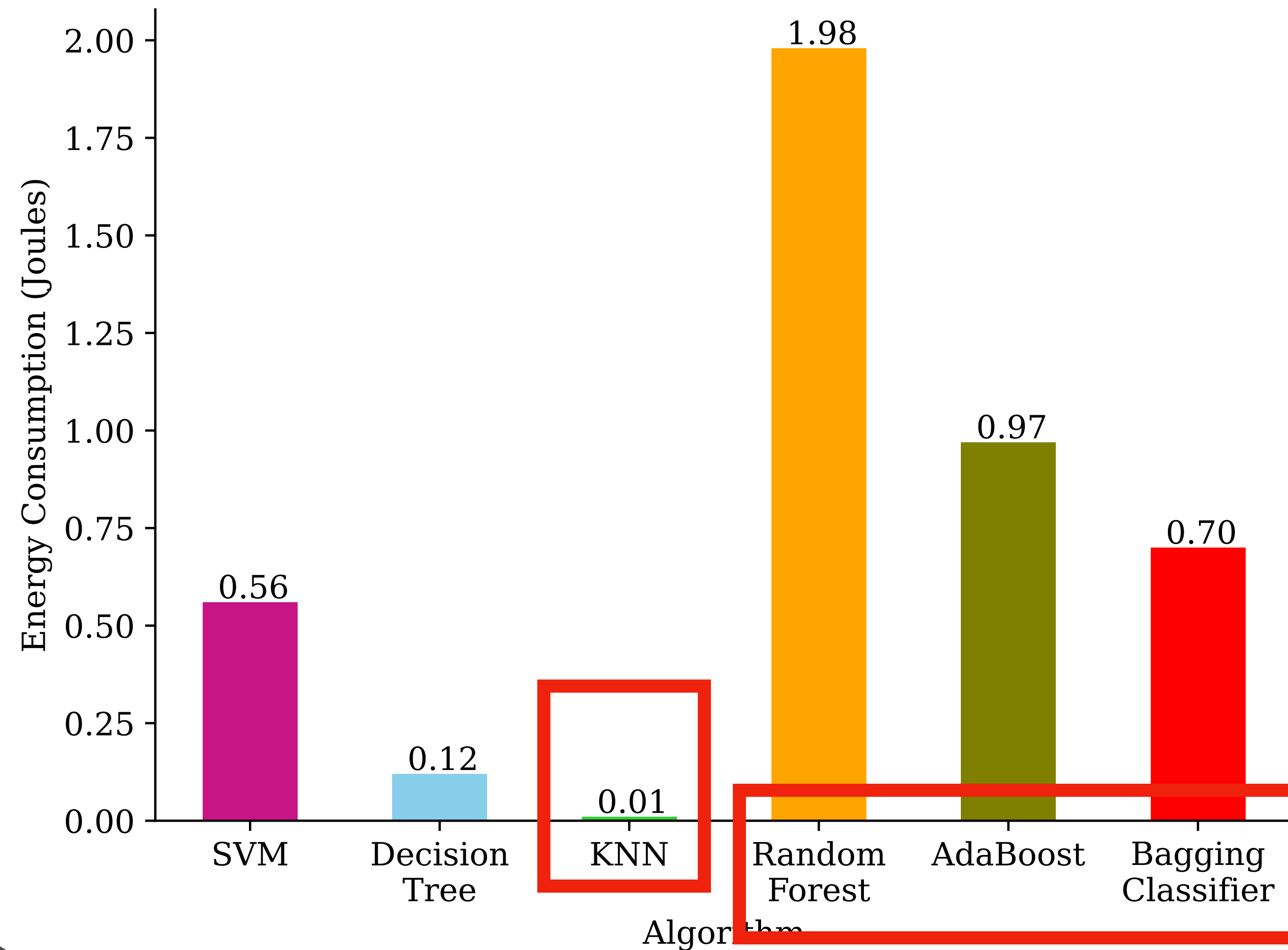
Method

- Single object of study: natural language model to **detect spam messages**.
- 6 machine learning algorithms: **SVM, Decision Tree, KNN, Random Forrest, AdaBoost, Bagging Classifier**.
- Reduce the number of rows. 10%, 20%, .., 100%
 - **Stratified random sampling** (?)
- Reduce the number of features. 10%, 20%, .., 100%
 - **Feature importance** metric based on the Chi-Square Test (Chi2)
- Estimate energy consumption using a RAPL-based tool. (?)

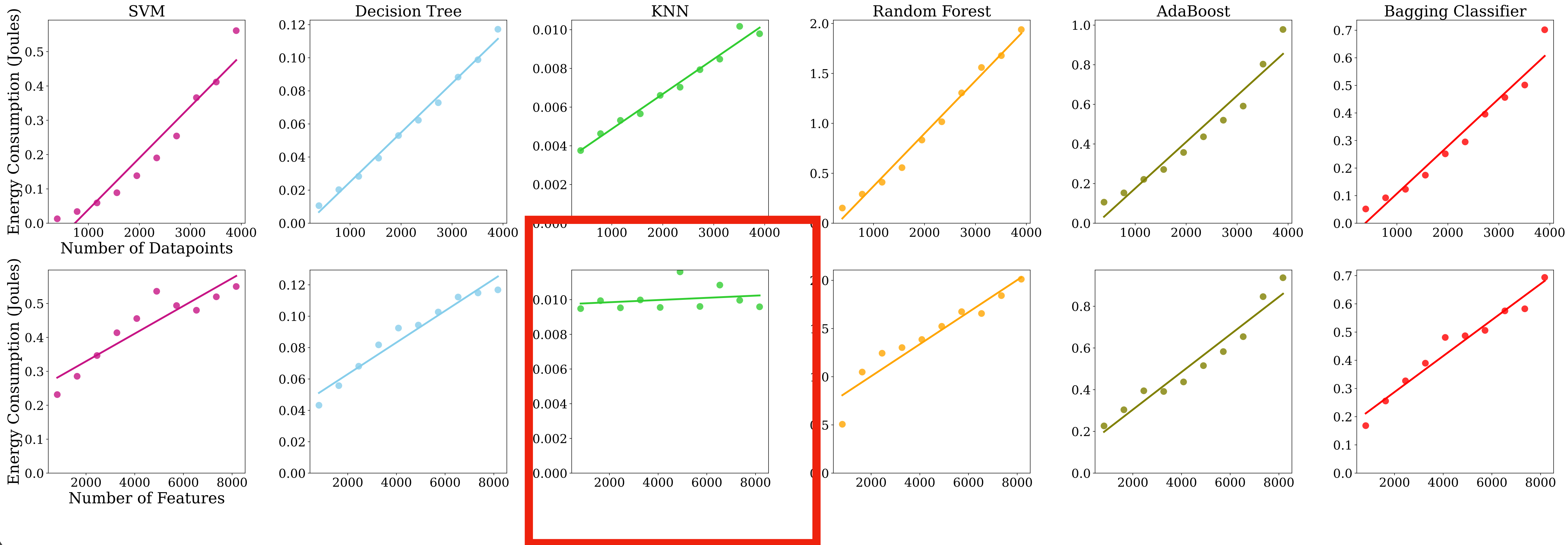
- Repeat 30 times
- Fix random seeds
- ...
- Data was **not Normal** => tailed Normal distribution.



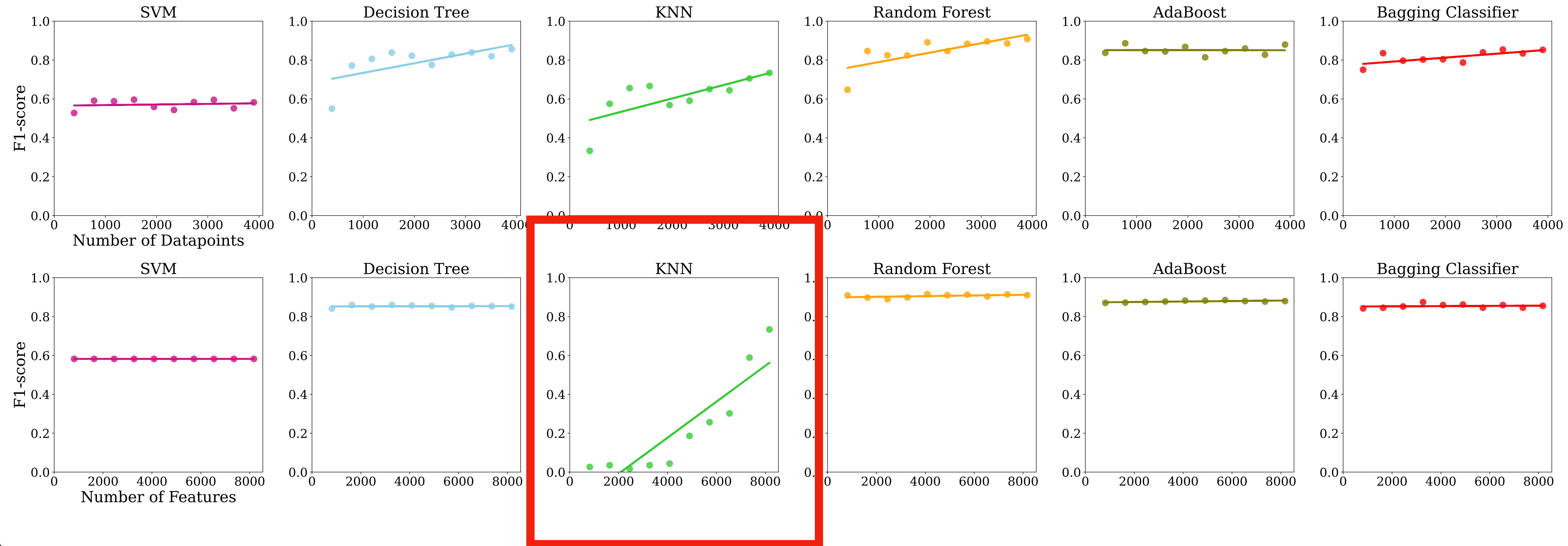
Results: energy consumption of algorithms



Results: energy vs data shape



Results: performance vs data shape



Discussion

- Other data properties should be investigated.
 - E.g., data types
- **Reporting energy data** is essential. It can lead to different model selection without hindering model performance.
- There is a big opportunity in:
 - **Model distillation.** Distill the knowledge of a complex model into a simpler one.
 - **Core set extraction.** (extracting the smallest subset that keeps the 623 key dataset properties)
 - **Dataset distillation.** (Derive/distill a smaller dataset into a synthetic dataset that yields similar models)

Green AI at Facebook

Sustainable AI: Environmental Implications, Challenges and Opportunities (2022)

arXiv:2111.00364v2 [cs.LG] 9 Jan 2022

Sustainable AI: Environmental Implications, Challenges and Opportunities

Carole-Jean Wu, Ramya Raghavendra, Udit Gupta, Bilge Acun, Neerasha Ardalani, Kiwan Macng, Gloria Chang, Fiona Aga Behram, James Huang, Charles Bai, Michael Gschwind, Anurag Gupta, Mylic Ott, Anastasia Melnikov, Salvatore Candido, David Brooks, Geeta Chauhan, Benjamin Lee, Hsien-Hsin S. Lee, Rugra Akyildiz, Maximilian Balandat, Joe Spisak, Ravi Jain, Mike Rabhat, Kim Hazelwood

Facebook AI

Abstract—This paper explores the environmental impact of the super-linear growth trends for AI from a holistic perspective, spanning *Data*, *Algorithms*, and *System Hardware*. We characterize the carbon footprint of AI computing by examining the model development cycle across industry-scale machine learning use cases and, at the same time, considering the life cycle of system hardware. Taking a step further, we capture the operational and manufacturing carbon footprint of AI computing and present an end-to-end analysis for *what* and *how* hardware-software design and at-scale optimization can help reduce the overall carbon footprint of AI. Based on the industry experience and lessons learned, we share the key challenges and chart out important development directions across the many dimensions of AI. We hope the key messages and insights presented in this paper can inspire the community to advance the field of AI in an environmentally-responsible manner.

I. INTRODUCTION

Artificial Intelligence (AI) is one of the fastest growing domains spanning research and product development and significant investment in AI is taking place across nearly every industry, policy, and academic research. This investment in AI has also stimulated novel applications in domains such as science, medicine, finance, and education. Figure 1 analyzes the number of papers published within the scientific disciplines, illustrating the growth trend in recent years¹.

AI plays an instrumental role to push the boundaries of knowledge and sparks novel, more efficient approaches to conventional tasks. AI is applied to predict protein structures radically better than previous methods. It has the potential to revolutionize biological sciences by providing in-silico methods for tasks only possible in a physical laboratory setting [1]. AI is demonstrated to achieve human-level conversation tasks, such as the Blender Bot [2], and play games at superhuman levels, such as AlphaZero [3]. AI is used to discover new electrocatalysts for efficient and scalable ways to store and utilize renewable energy [4], predicting renewable energy availability in advance to improve energy utilization [5], operating hyperscale data centers efficiently [6], growing plants using less natural resources [7], and, at the same time, being used to tackle climate changes [8], [9]. It is projected that, in the next five years, the market for AI will increase by 10× into hundreds of billions of dollars [10]. All of these investments

¹Based on monthly counts, Figure 1 estimates the cumulative number of papers published per category on the arXiv database.

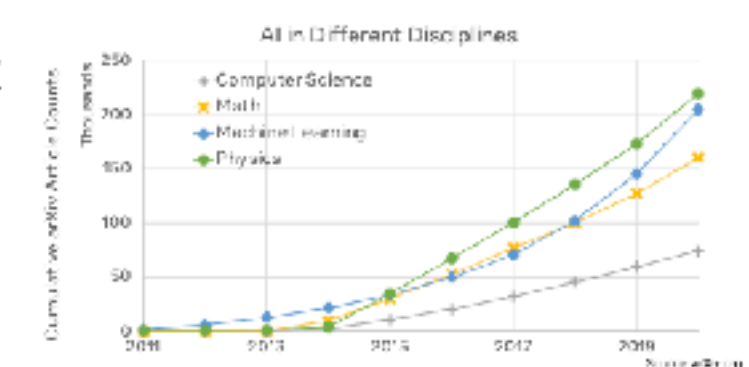


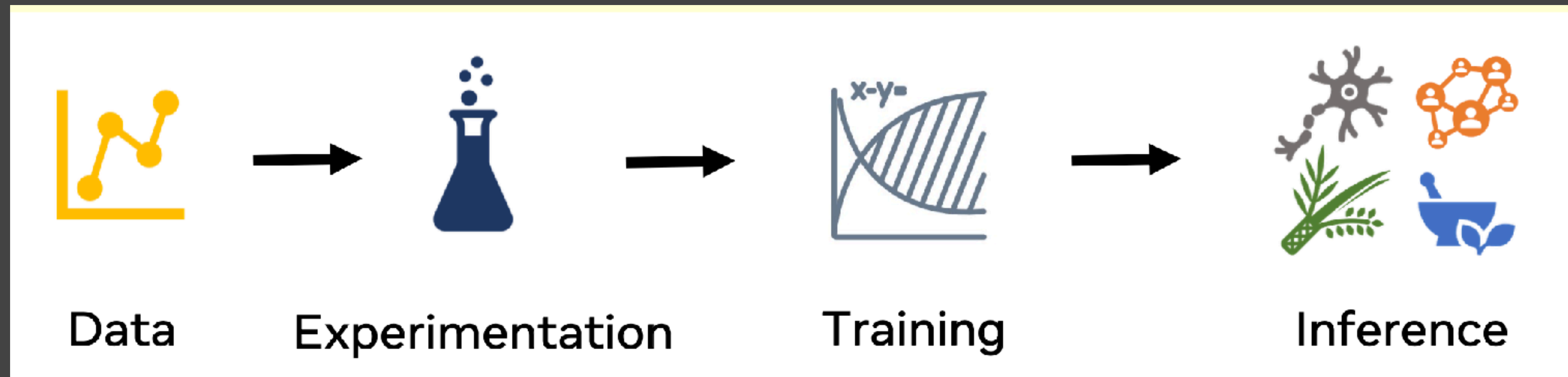
Fig. 1. The growth of AI is exceeding that of many other scientific disciplines. Significant research growth in machine learning is observed in recent years as illustrated by the increasing cumulative number of papers published in machine learning with respect to other scientific disciplines based on the monthly count (y-axis measures the cumulative number of articles on arXiv).

in research, development, and deployment have led to a super-linear growth in AI data, models, and infrastructure capacity. With the dramatic growth of AI, it is imperative to understand the environmental implications, challenges, and opportunities of this nascent technology. This is because technologies tend to create a self-accelerating growth cycle, putting new demands on the environment.

This work explores the environmental impact of AI from a holistic perspective. More specifically, we present the challenges and opportunities to designing sustainable AI computing across the key phases of the machine learning (ML) development process — *Data*, *Experimentation*, *Training*, and *Inference* — for a variety of AI use cases at Facebook, such as vision, language, speech, recommendation and ranking. The solution space spans across our fleet of datacenters and on-device computing. Given particular use cases, we consider the impact of AI *data*, *algorithms*, and *system hardware*. Finally, we consider emissions across the life cycle of hardware systems, from manufacturing to operational use.

AI Data Growth. In the past decade, we have seen an exponential increase in AI training data and model capacity. Figure 2(b) illustrates that the amount of training data at Facebook for two recommendation use cases — one of the fastest growing areas of ML usage at Facebook — has increased by 2.4× and 1.9× in the last two years, reaching exabyte scale. The increase in data size has led to a 3.2× increase in data ingestion bandwidth demand. Given this increase, data storage and the ingestion pipeline accounts for a significant portion of

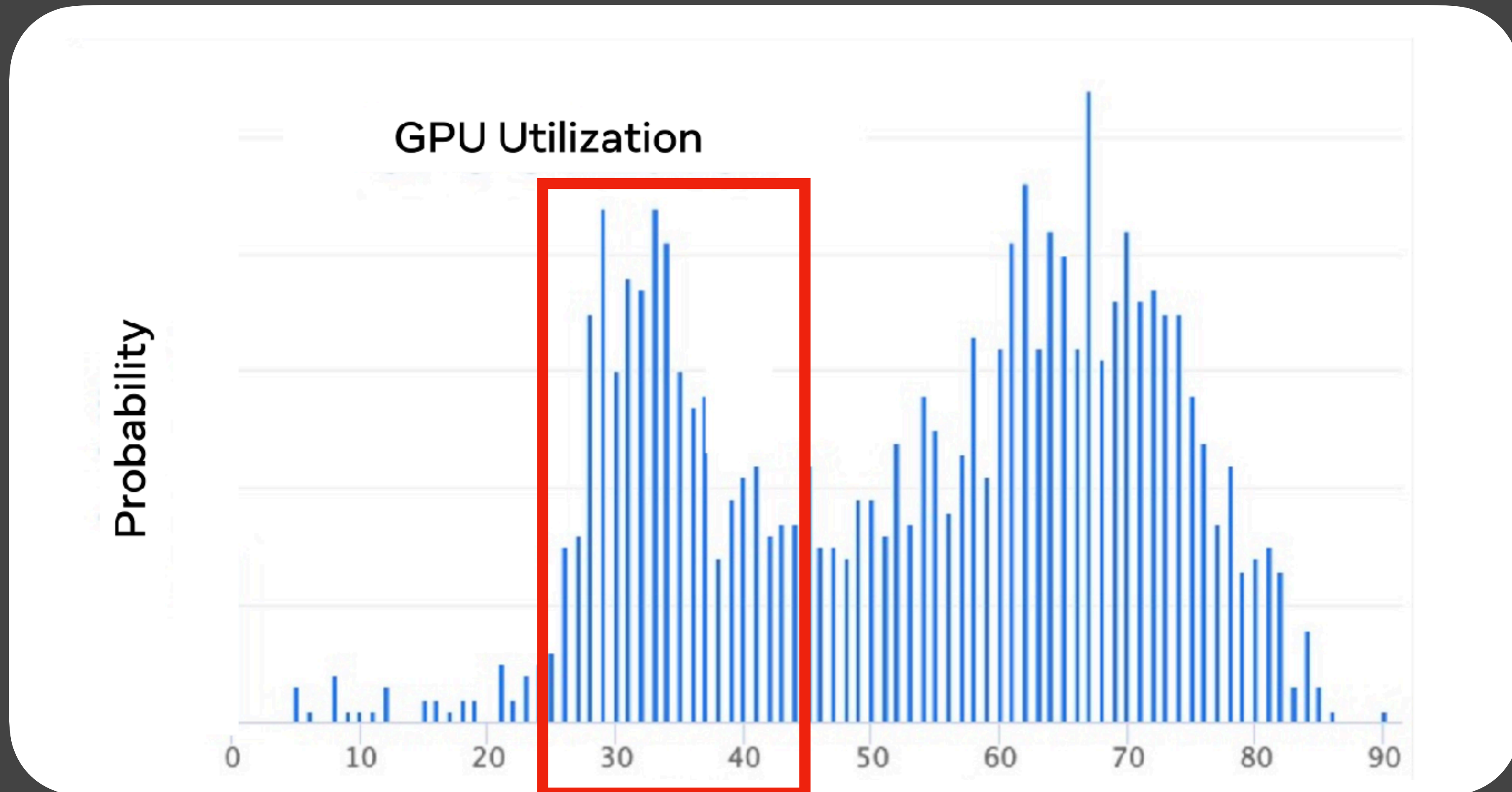
Carbon footprint mapped to the AI lifecycle



- There are 4 main overarching stages where carbon emissions need to be isolated: **data collection, experimentation, training, inference**.
- At Facebook, recommendation systems split energy consumption **evenly between training and inference**; text translation models have a **35%/65%** split. (Operational cost)
- Operational/embodied cost split: **30%/70%**

Open issues according to Meta

- A vast portion of projects only use **GPUs at 30%**. Should be higher to attenuate embodied carbon.



Federated learning

- Federated learning consists of training a ML model across **multiple decentralized edge devices** holding local data samples.
- Federated-learning is a nice solution for **devices with limited energy resources**. E.g., IoT.

Federated learning might not be the solution

- Most of the carbon footprint stems from **communications**

