# LLM efficiency challenge: 1 LLM + 1 GPU + 1 Day

Weiwei Yang -  Microsoft Research
Mark Saroufim - Meta

NEURAL INFORMATION
PROCESSING SYSTEMS

# Our goal

**Current LLMs:**

1. Lack transparency and community
2. Lack meaningful evaluation metrics
3. Lack hardware

# Rules

1. Only approved base models
2. Most public datasets are allowed
3. Open source

4090 and A100 track

# Everything is OSS

winner submissions https://llm-efficiency-challenge.github.io/leaderboard

evaluation code https://github.com/llm-efficiency-challenge
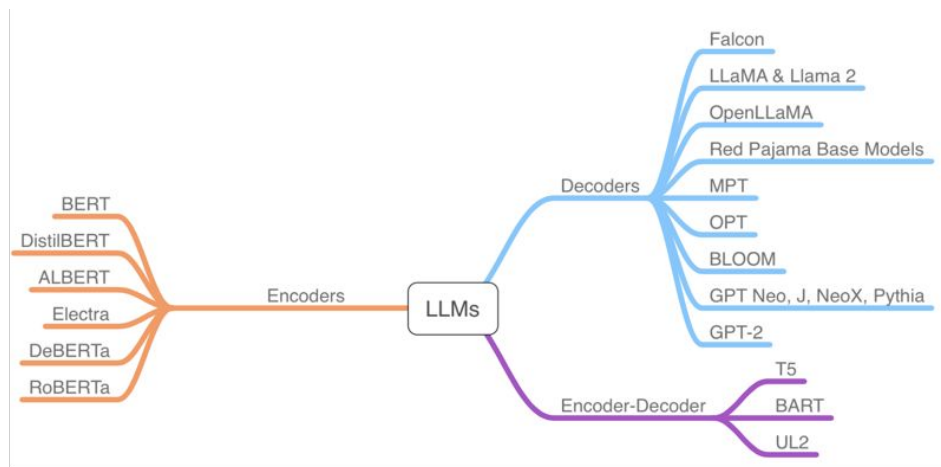
# Participation

1,400+ people on Discord

700+ successful submissions on Discord Leaderboard

182 teams

# What is a good base model?

1. Scores well on standard eval datasets
2. Powers a popular consumer app
3. Generalizes better after seeing more data



Visualization source https://lightning.ai/pages/community/tutorial/neurips2023-llm-efficiency-guide/

# Open vs closed eval

## ⅓ Open

- BigBench
- MMLU
- TruthfulQA
- CNN/DailyMail
- GSM8k
- BBQ

## ⅔ Closed

- ETHICS
- MATH
- Samsum
- corr2cause

# We want to do it again!

Public open leaderboards are not enough

Communities are necessary for evaluation

Vibrant community of educators and tools

Docker and pip are not that reproducible of a combo

Eval and data tools are very early

Open to variants

# Seeding the community

## Optimizing LLMs From a Dataset Perspective

Sep 15, 2023
by Sebastian Raschka

This article focuses on improving the modeling performance of LLMs by finetuning them using carefully curated datasets. Specifically, this article highlights strategies that involve modifying, utilizing, or manipulating the datasets for instruction-based finetuning rather than altering the model architecture or training algorithms (the latter will be topics of a future article). This article will also explain how you can prepare your own datasets to finetune open-source LLMs.

Note that the NeurIPS LLM Efficiency Challenge is currently underway, aiming to train a Large Language Model on a single GPU within a 24-hour period, which is super interesting for practitioners and researchers interested in LLM efficiency. The techniques discussed in this article have direct relevance to this competition, and we will delve into how these dataset-centric strategies could potentially be applied within the challenge setting. Additionally, the article will offer suggestions for new experiments you might consider trying.

**This article is a cross-post that originally appeared here on the Lightning AI Blog**.

https://sebastianraschka.com/blog/2023/optimizing-LLMs-dataset-perspective.html

# A place to clarify rules

# A place to run evaluation jobs

# Data is all you need (from the 4090 track winners)

# Thank you!

Sponsors

aws    Lightning AI    Microsoft

moz://a ai    λ Lambda

Special Thanks

- Stanford for creating and supporting HELM
- Weights and Biases for promoting our competition in their LLM finetuning course
- Greg Bowyer for donating 4090 GPUs for our eval infra

https://llm-efficiency-challenge.github.io/sponsors

# Thank you!

## Organizers

- Mark Saroufim
- Weiwei Yang
- Joe Isaacson
- Luca Antiga
- Driss Guessous
- Greg Bowyer
- Christian Puhrsch
- Geeta Chauhan
- Supriya Rao
- Artidoro Pagnoni
- Vicki Boykis
- Aaron Gonzales
- Davide Eynard

## Advisors

- Sebastian Raschka
- Yifan Mai
- Yotam Perlitz
- Leshen Choshen
- Danylo Baibak
- Jean Schmidt
- Eli Uriegas
- Fahd Husain
- Kaleab Kinfu
- Matthias Reso

https://llm-efficiency-challenge.github.io/organizers

https://llm-efficiency-challenge.github.io/advisors

## Schedule

| | | | |
|---|---|---|---|
| Fri 1:30 p.m. - 1:45 p.m. | 🔖 | **Kick-Off to Efficiency: Welcoming statement for the organizers** ( Speak ) | *Mark Saroufim · Weiwei Yang* |
| Fri 1:45 p.m. - 2:00 p.m. | 🔖 | **Invited Speaker: Jeremy Howard-Lessons from 25 years of machine learning competitions** ( In-person presentation ) | |
| Fri 2:00 p.m. - 2:15 p.m. | 🔖 | **Invited Speaker: Sebastian Raschka (lightning.ai) - LoRA in Action: Insights from Finetuning LLMs with Low-Rank Adaptation** ( In-person presentation ) | *Sebastian Raschka* |
| Fri 2:15 p.m. - 2:30 p.m. | 🔖 | **Unveiling Success: A100 track Team percent_bdf's Winning Strategies** ( Zoom presentation ) | *Ao Liu* |
| Fri 2:30 p.m. - 2:45 p.m. | 🔖 | **Invited Speaker: Tim Dettmers QLoRA** ( In-person presentation ) | *Tim Dettmers* |
| Fri 2:45 p.m. - 3:00 p.m. | 🔖 | **Invited Speaker: Sourab Mangrulka -- Generative AI for AI: 🤗 PEFT: Finetuning made simple, efficient and extendable** ( Zoom presentation ) | *SOURAB MANGRULKAR* |
| Fri 3:00 p.m. - 3:15 p.m. | 🔖 | **Coffee break** | |
| Fri 3:15 p.m. - 3:30 p.m. | 🔖 | **Unveiling Success: 4090 Track winning team's strategies** ( In-person presentation ) | |
| Fri 3:30 p.m. - 3:45 p.m. | 🔖 | **Invited Speaker: Keming Lu (Alibaba Research) - Qwen: Towards a Generalist Model** ( In-person presentation ) | *Keming Lu* |
| Fri 3:45 p.m. - 4:00 p.m. | 🔖 | **Invited Speaker: Mojan Javaheripi (Microsoft Research) - Unleashing the power of Small Language Models** ( Zoom presentation ) | *Mojan Javaheripi* |
| Fri 4:00 p.m. - 4:15 p.m. | 🔖 | **Invited Speaker: Leshem Choshen (IBM Research) - Efficient Evaluation for Efficient Training** ( In-person presentation ) | *Leshem Choshen* |
| Fri 4:15 p.m. - 4:30 p.m. | 🔖 | **Award celemony and open floor discussion** ( panel discussion ) | *Weiwei Yang · Mark Saroufim · Christian Puhrsch · Joseph Isaacson · Vicki Boykis* |