# Why LLM Benchmarking is Broken and How to Fix It

## Guanhua Zhang

Social Foundations of Computation

# Ranking Is All You Need

At the core of applied machine learning are *model rankings*

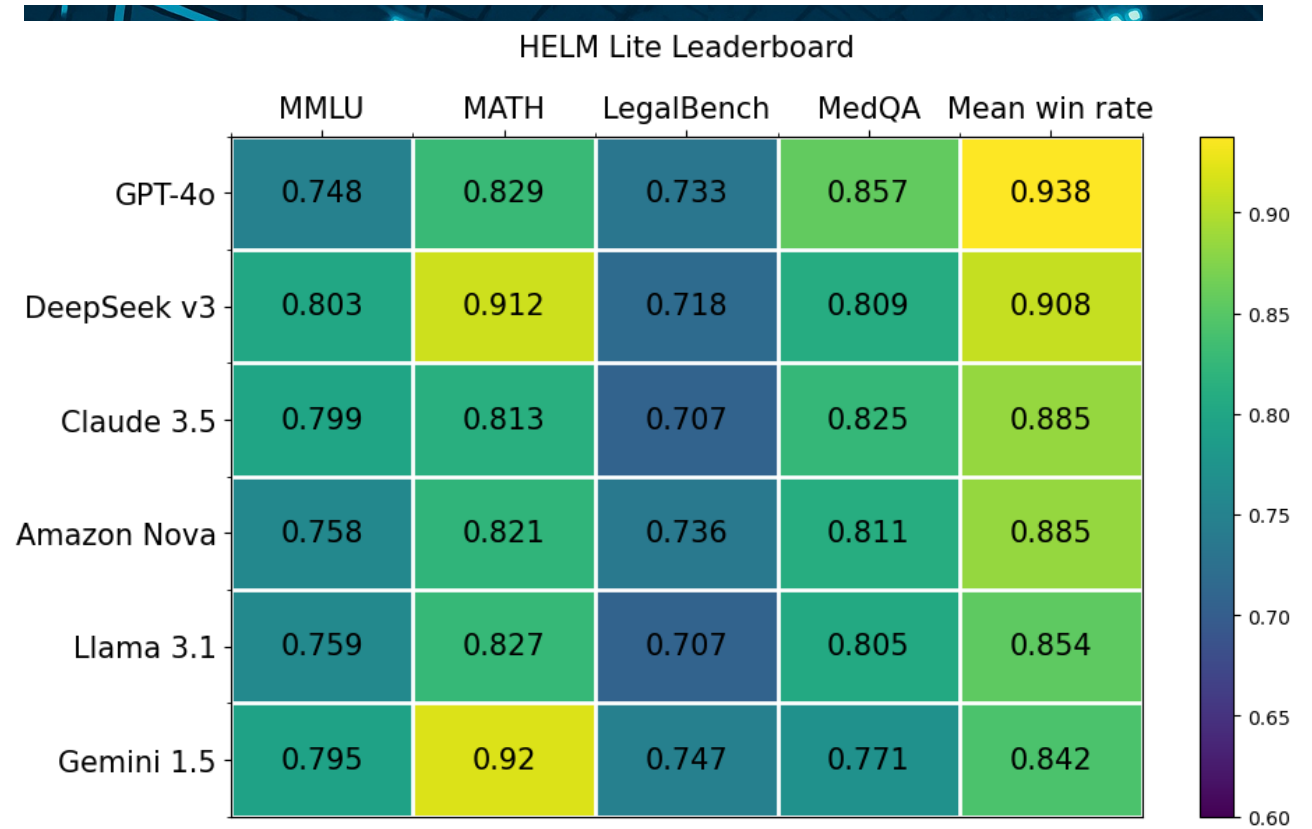*Good model rankings are the goal of benchmarking*

# Multi-Task Benchmarking for LLMs

LLMs can solve many tasks
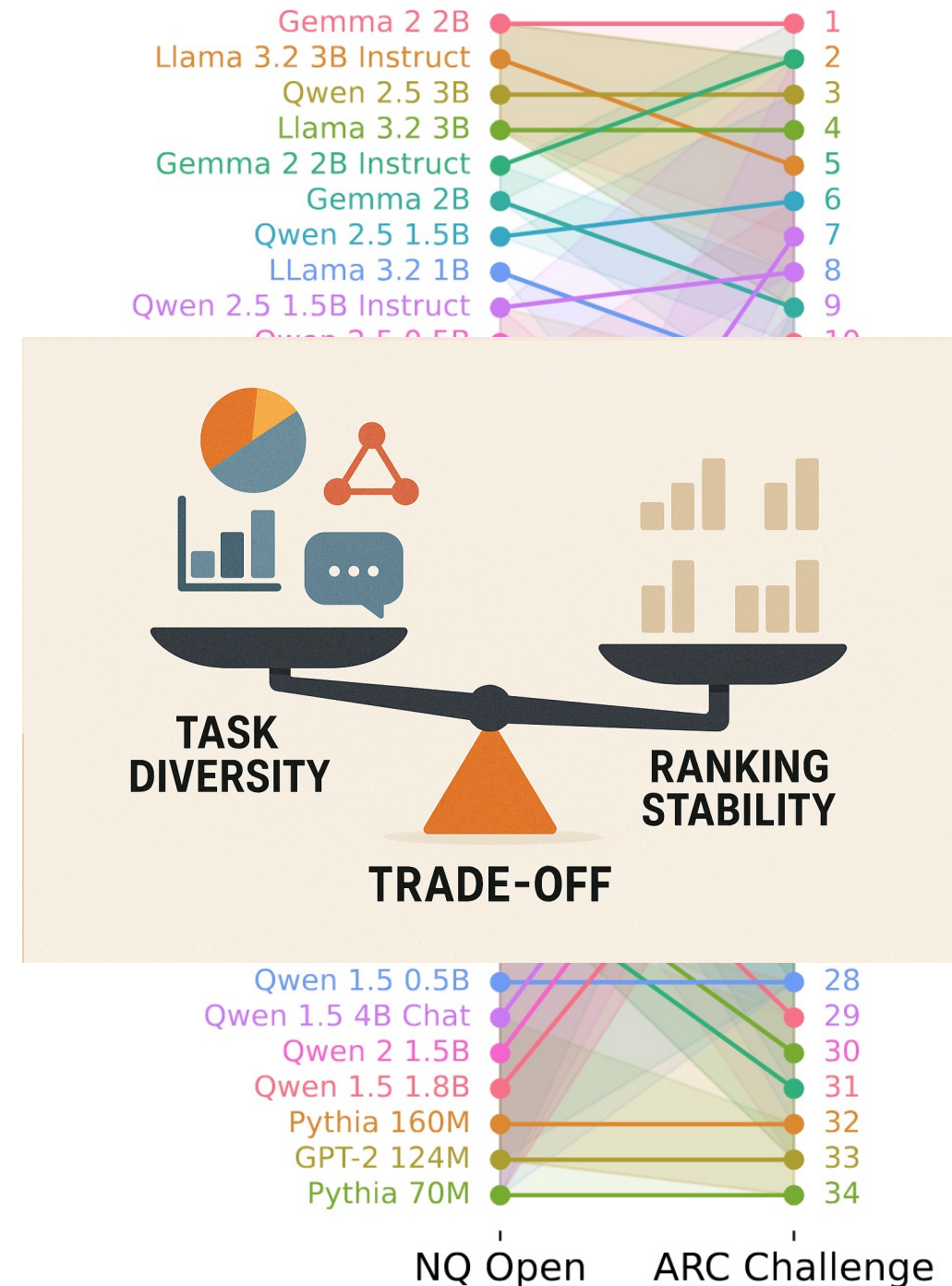
Which ranking should we look at?

Multi-task benchmarks: Just evaluate them on everything!



HELM Lite Leaderboard

|  | MMLU | MATH | LegalBench | MedQA | Mean win rate |
|---|---|---|---|---|---|
| GPT-4o | 0.748 | 0.829 | 0.733 | 0.857 | 0.938 |
| DeepSeek v3 | 0.803 | 0.912 | 0.718 | 0.809 | 0.908 |
| Claude 3.5 | 0.799 | 0.813 | 0.707 | 0.825 | 0.885 |
| Amazon Nova | 0.758 | 0.821 | 0.736 | 0.811 | 0.885 |
| Llama 3.1 | 0.759 | 0.827 | 0.707 | 0.805 | 0.854 |
| Gemini 1.5 | 0.795 | 0.92 | 0.747 | 0.771 | 0.842 |

# Tasks disagree with each other



- The model rankings in different tasks often differ, even if the two tasks are similar

- *Analogy with voting system*:

    Each task is a voter; each model is a candidate. Each voter ranks candidates

    Social choice theory: It's hard to aggregate many rankings into one good ranking.

- Our result: Inherent trade-off between task diversity and ranking stability in multi task benchmarks
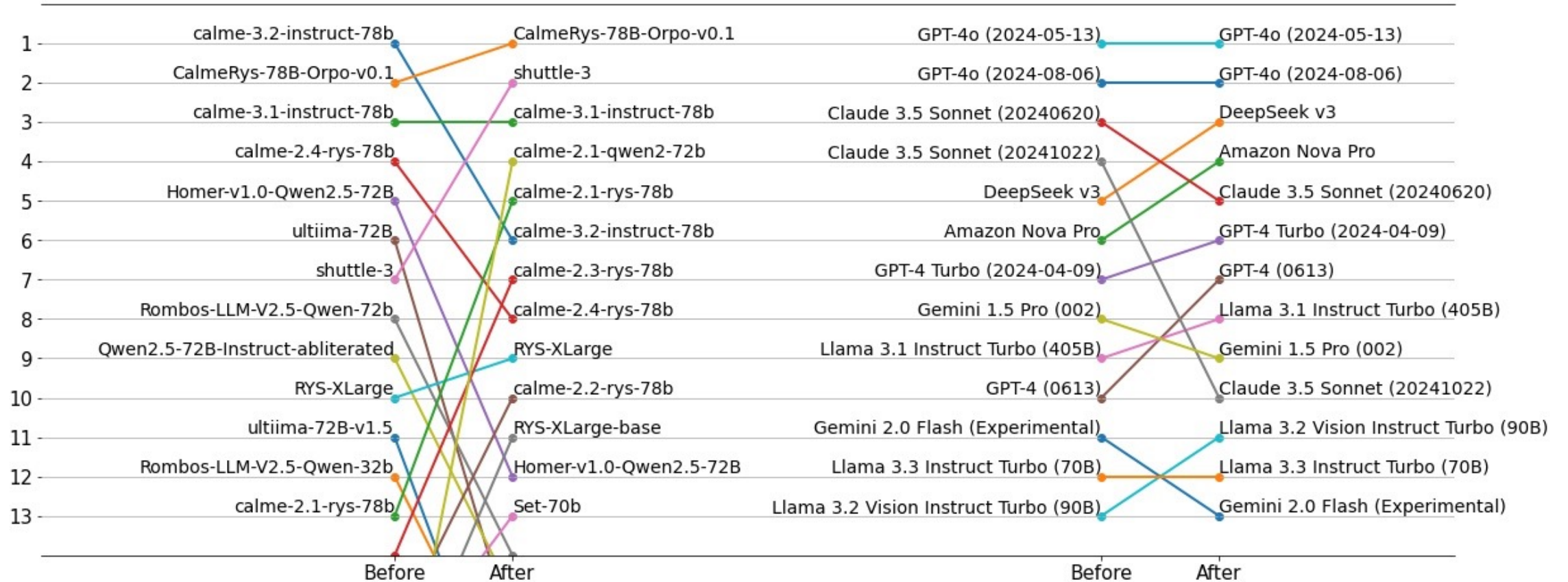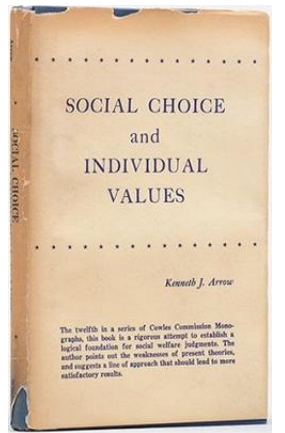
- **Sensitivity:**
  1. *Add different label noises to different tasks*
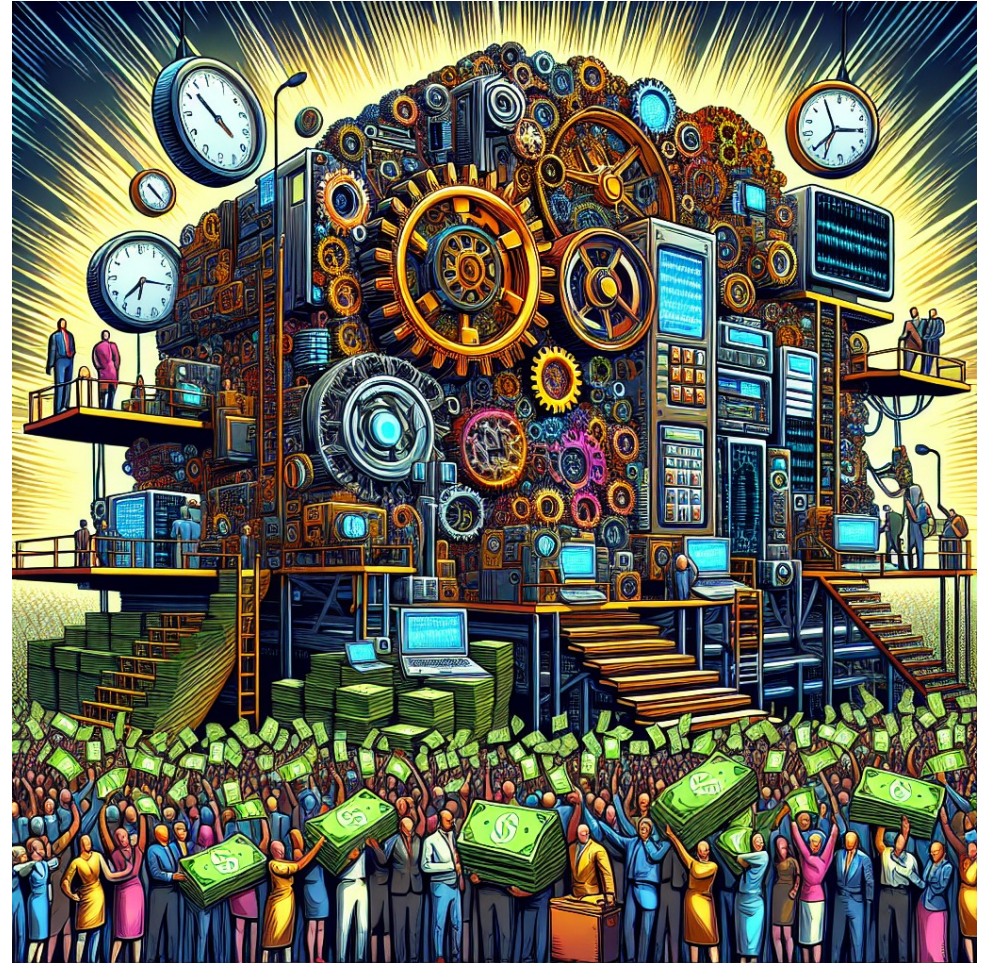  2. *Add some irrelevant weak models*
- **Diversity:**
  - *Ranking disagreement measured by Kendall's W*

# It gets worse: LLM Benchmarking is Costly

- Evaluating a single 176B parameter model, Bloom, on the HELM multi-task benchmark required 4,200 GPU hours

- People have proposed methods for benchmark performance prediction to speed up evaluation

- Our result: These methods fail at the frontier, where models are better than old models

So, it seems we're in a pinch:

1. Rankings are inconsistent

2. Computing many rankings is costly



But there's good news:

Ranking inconsistency is an artifact of how LLMs were trained

We can remove this artifact and recover highly consistent rankings

# Ranking inconsistency is due to training on the test task

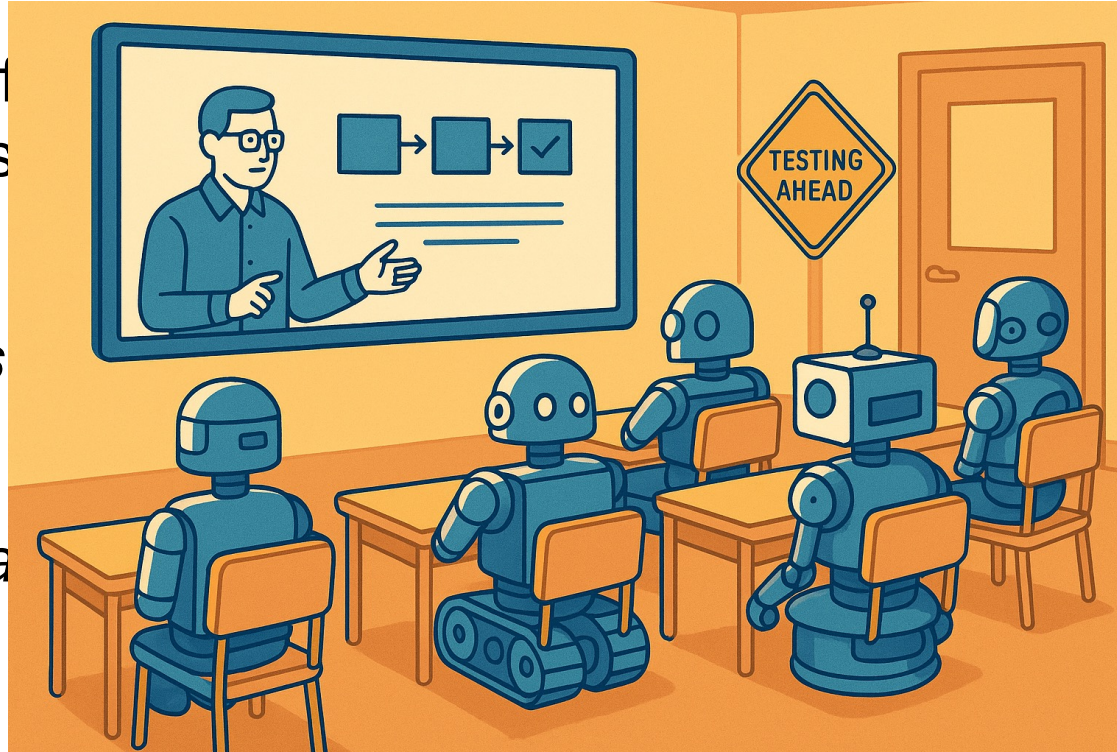As released, diff ... of preparation for any given tes ... ardt (2025)]

*"Some models ... en't"*

This is called *tra...*



***Train-before-test:*** Give each model the same benchmark specific fine-tuning before evaluation.

# Train-before-test harmonizes model rankings

# Tasks from the same category still disagree, unless ...



(a) Direct evaluation.

(b) Train-before-test.

Figure 3: Cross-category ranking agreement for direct evaluation (left) and train-before-test (right). We consider language understanding (LU), commonsense reasoning (CR), question answering (QA), physics/biology/chemistry (PBC), math (Math), and medicine (Med) categories. Kendall's $\tau$ is averaged across all pairs of benchmarks that belong to two given categories. The diagonal represents the intra-category agreement and the others represent the inter-category agreement. train-before-test improves both intra- and inter-category ranking agreement in all instances.

# Downstream agrees with perplexity under TbT



Direct evaluation.

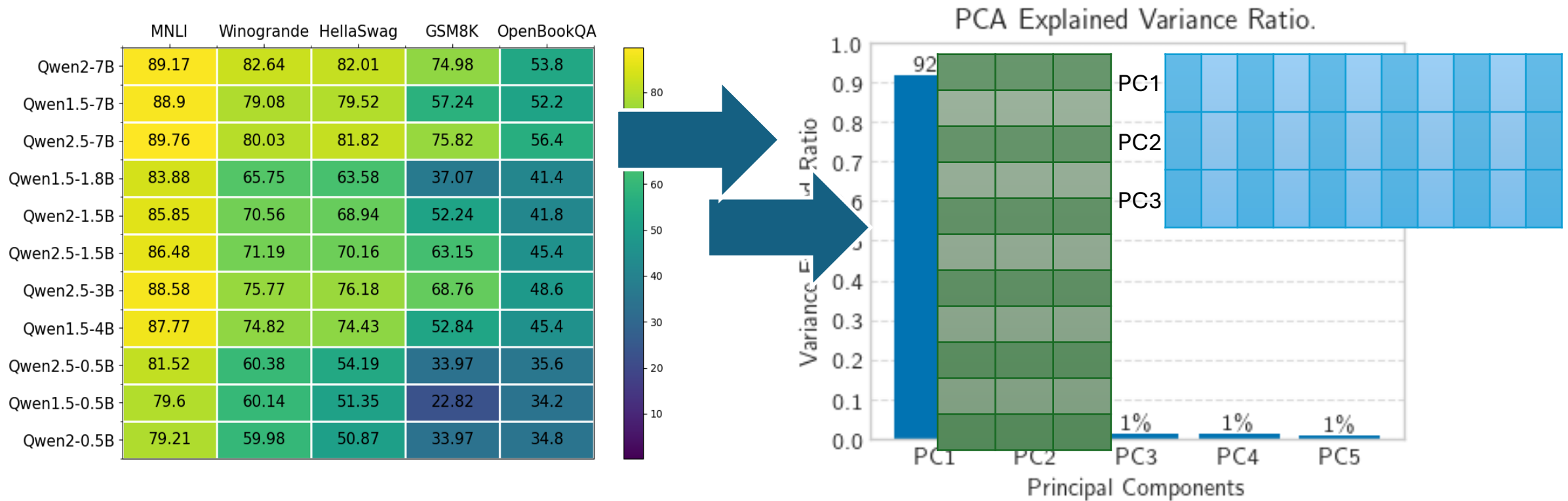|  | Wiki | Stack | Arxiv | MNLI | QQP | MedMCQA | QNLI | NQ-Open | SST-2 | Winogrande | HellaSwag | Social-IQA | MathQA | ANLI | PIQA | SciQ | CommonsenseQA | BoolQ | CoLA | GSM8K | WiC | OpenBookQA | MRPC | HeadQA | RTE | ARC-Easy | ARC-Challenge |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Wiki | 1 | 0.75 | 0.8 | 0.34 | 0.35 | 0.48 | 0.4 | 0.42 | 0.34 | 0.71 | 0.73 | 0.55 | 0.56 | 0.4 | 0.8 | 0.59 | 0.47 | 0.5 | 0.21 | 0.34 | 0.21 | 0.6 | 0.25 | 0.7 | 0.32 | 0.69 | 0.65 |
| Stack | 0.75 | 1 | 0.85 | 0.4 | 0.52 | 0.55 | 0.49 | 0.36 | 0.42 | 0.62 | 0.69 | 0.67 | 0.64 | 0.34 | 0.71 | 0.68 | 0.55 | 0.52 | 0.26 | 0.37 | 0.29 | 0.66 | 0.25 | 0.69 | 0.41 | 0.72 | 0.69 |
| Arxiv | 0.8 | 0.85 | 1 | 0.35 | 0.41 | 0.52 | 0.46 | 0.38 | 0.4 | 0.66 | 0.74 | 0.62 | 0.59 | 0.37 | 0.75 | 0.6 | 0.52 | 0.53 | 0.26 | 0.38 | 0.27 | 0.62 | 0.27 | 0.72 | 0.35 | 0.66 | 0.69 |

Train-before-test.

|  | Wiki | Stack | Arxiv | MNLI | QQP | MedMCQA | QNLI | NQ-Open | SST-2 | Winogrande | HellaSwag | Social-IQA | MathQA | ANLI | PIQA | SciQ | CommonsenseQA | BoolQ | CoLA | GSM8K | WiC | OpenBookQA | MRPC | HeadQA | RTE | ARC-Easy | ARC-Challenge |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Wiki | 1 | 0.78 | 0.82 | 0.78 | 0.71 | 0.67 | 0.65 | 0.79 | 0.79 | 0.85 | 0.87 | 0.79 | 0.71 | 0.64 | 0.89 | 0.8 | 0.61 | 0.81 | 0.68 | 0.64 | 0.67 | 0.8 | 0.54 | 0.82 | 0.57 | 0.83 | 0.77 |
| Stack | 0.78 | 1 | 0.84 | 0.81 | 0.76 | 0.74 | 0.61 | 0.7 | 0.63 | 0.77 | 0.8 | 0.64 | 0.78 | 0.68 | 0.74 | 0.81 | 0.67 | 0.81 | 0.68 | 0.71 | 0.74 | 0.77 | 0.66 | 0.79 | 0.65 | 0.86 | 0.8 |
| Arxiv | 0.82 | 0.84 | 1 | 0.81 | 0.78 | 0.68 | 0.65 | 0.68 | 0.66 | 0.8 | 0.83 | 0.69 | 0.71 | 0.68 | 0.79 | 0.83 | 0.64 | 0.82 | 0.72 | 0.63 | 0.7 | 0.79 | 0.61 | 0.86 | 0.61 | 0.82 | 0.75 |

# Train-before-test makes score matrix rank one

- Conduct principal component analysis (PCA) on the multi-task score matrix.



- Our result: A single factor (PC1) dominates model performances on 24 tasks
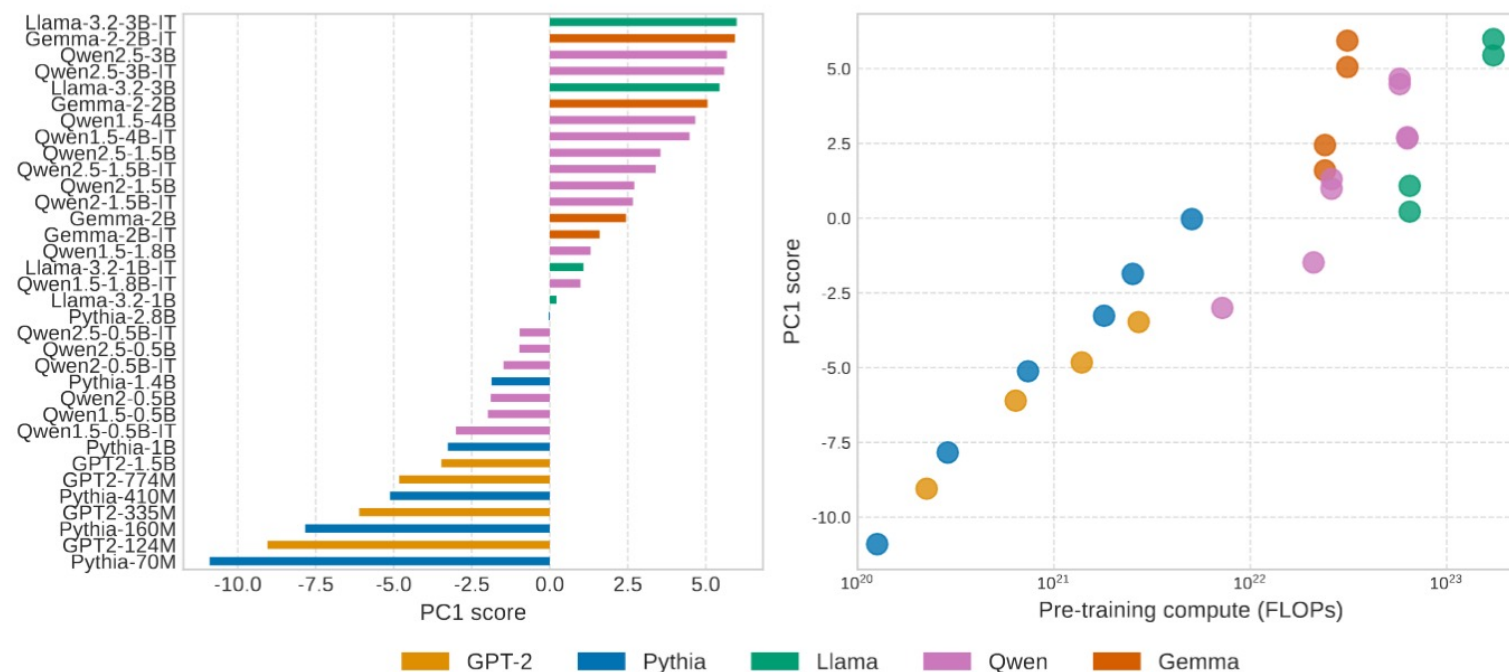
# PC1 correlates with model scale



Figure 7: PC1 scores under train-before-test correlates with scale and pre-training compute.

- PC1 score stands for something **useful for all tasks**.
  - *All dimensions of PC1 is positive.*

# Model potential is what really matters

- Train-before-Test measures **model potential** after development
  - Model potential rankings in any benchmark extend to others
  - Model potential correlates with perplexity of models
  - Model potential is of rank one

# Take-away

Ranking is all you need
Currently benchmarking is broken for LLMs
But there's a fix: Use train-before-test.

## Thanks!