

# Easy2Hard-Bench: Standardized Difficulty Labels for Profiling LLM Performance and Generalization

Mucong Ding<sup>\*1</sup>, Chenghao Deng<sup>\*1</sup>, Jocelyn Choo<sup>1</sup>, Zichu Wu<sup>2</sup>, Aakriti Agrawal<sup>1</sup>, Avi Schwarzschild<sup>3</sup>, Tianyi Zhou<sup>1</sup>, Tom Goldstein<sup>1</sup>, John Langford<sup>4</sup>, Anima Anandkumar<sup>5</sup>, Furong Huang<sup>1</sup>

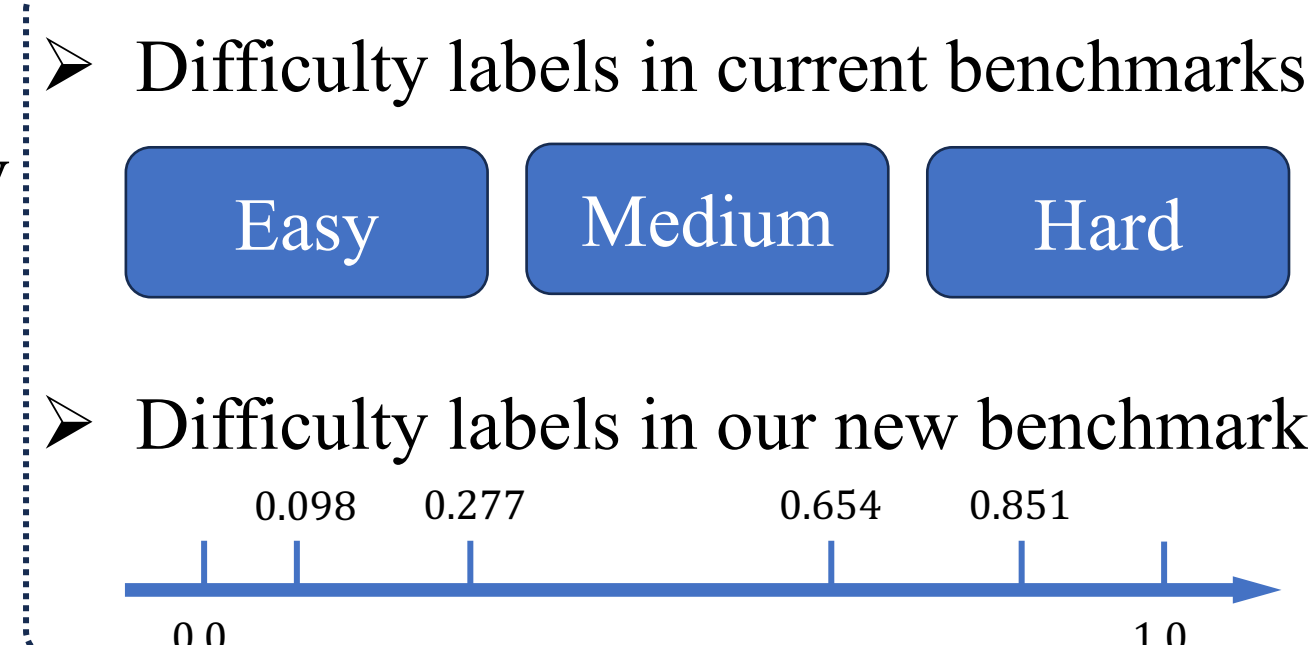
<sup>1</sup>University of Maryland, College Park <sup>2</sup>University of Waterloo <sup>3</sup>Carnegie Mellon University <sup>4</sup>Microsoft <sup>5</sup>California Institute of Technology

<sup>\*</sup>Equal contribution Correspondence: {mcding, dengch16, furongh}@umd.edu



## Introduction

- ❖ Motivation: difficulty labels of problems
  - Current benchmarks: categorical, low granularity
  - Ideal benchmarks: numerical, high granularity
- ❖ Why difficulty labels matters?
  - Model behavior within a wide difficulty range.
  - Generalization ability from easy to hard.



### Our contributions

- Presenting the innovative datasets with fine-grained difficulty labels
- Broadening the difficulty range and diversifying the task domains
- Thoroughly evaluating the SoTA LLMs on our new benchmark
- Detailed Profiling LLMs' easy-to-hard generalization ability

## Easy2Hard-Bench Overview

**E2H-AMC**

Find the number of pairs of integers  $(a, b)$  with  $1 \leq a < b \leq 57$  such that  $a^2$  has a smaller remainder than  $b^2$  when divided by 57.

Difficulty:  $0.587 \pm 0.085$  Percentile: 70.3%

**E2H-GSM8K**

Jerry is twice as old as he was 5 years ago. How old will Jerry be in 3 years?

Difficulty:  $0.301 \pm 0.157$  Percentile: 30.3%

**E2H-ARC**

Which of the following takes place during fertilization in animals?  
 (A) production of sperm and egg (B) joining of sperm and egg  
 (C) division of egg (D) development of embryo

Difficulty:  $0.425 \pm 0.066$  Percentile: 30.4%

**E2H-Winogrande**

A very active Randy injured their elbow playing basketball against Lawrence, so \_ said sorry.  
 (A) Randy (B) Lawrence

Difficulty:  $0.124 \pm 0.164$  Percentile: 10.3%

**CODEFORCES**

You're given a string of lower-case Latin letters. Your task is to find the length of its longest substring that can be met in the string at least twice. These occurrences can overlap.

Difficulty:  $0.204 \pm 0.025$  Percentile: 29.1%

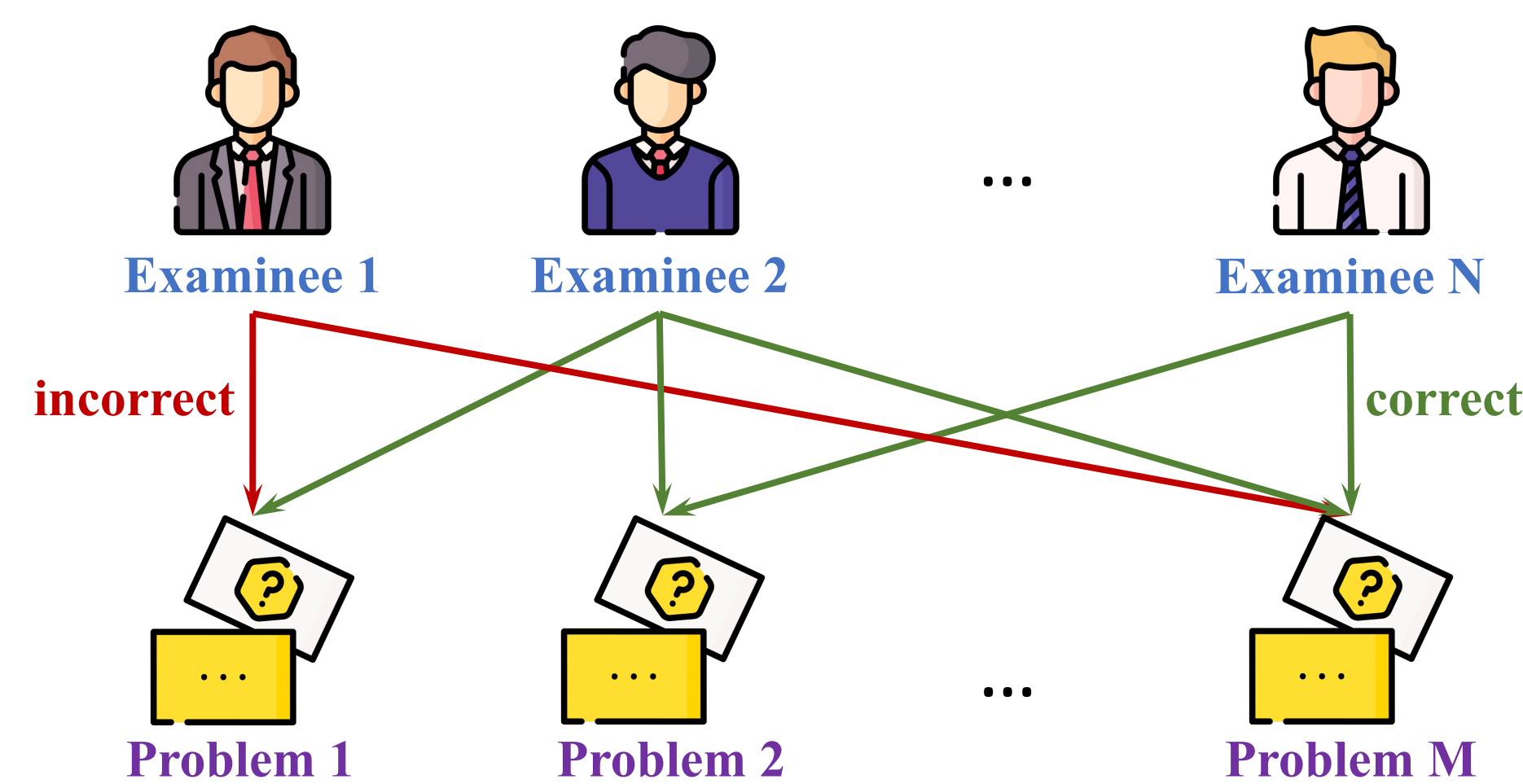
**E2H-Lichess**

Chess

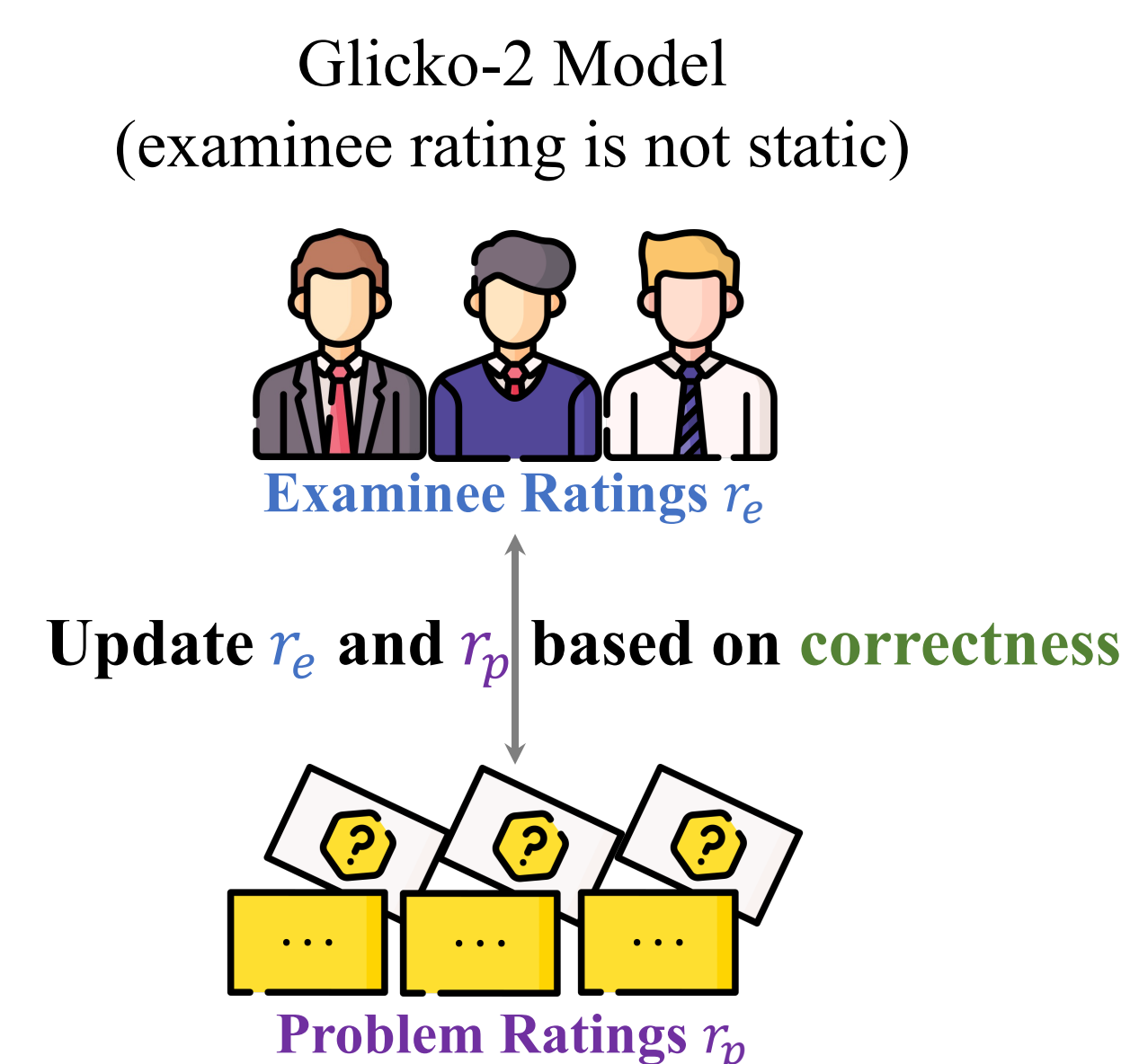
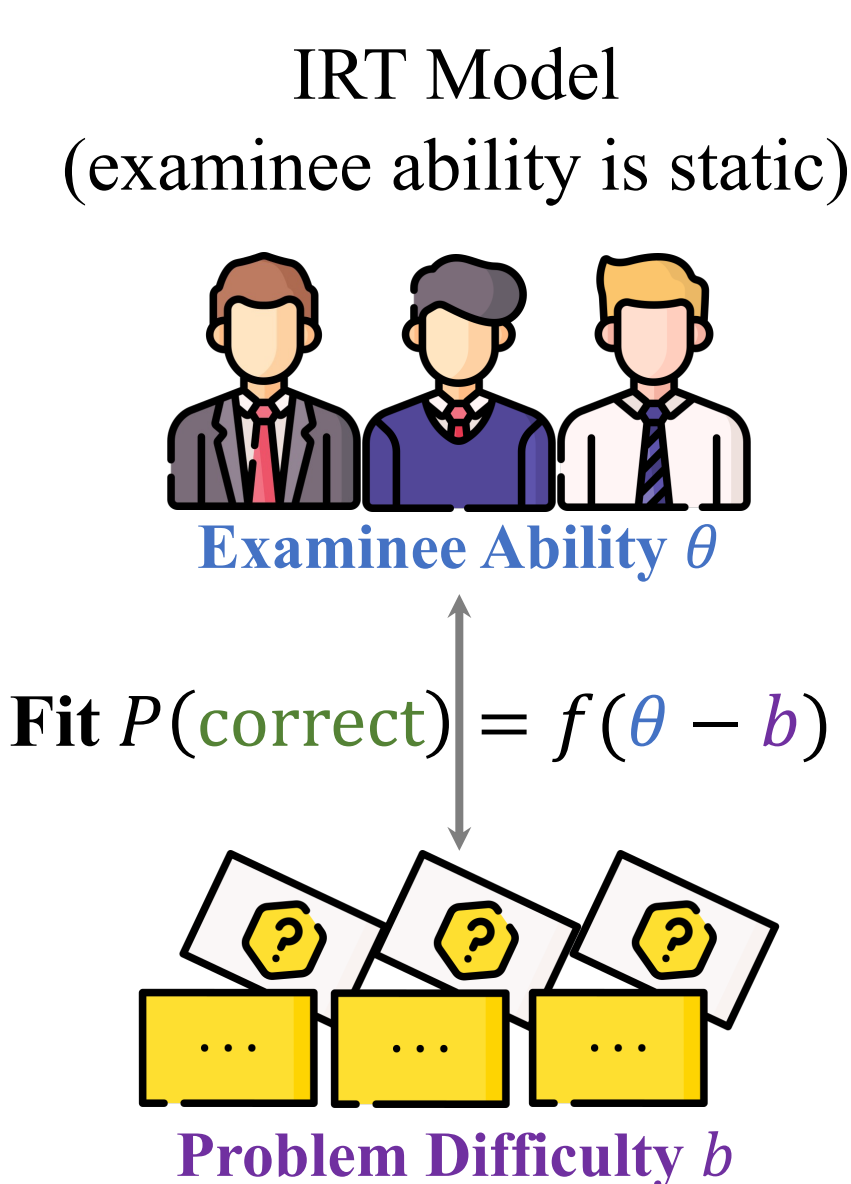
Difficulty:  $0.163 \pm 0.038$  Percentile: 30.9%

## Fine-grained Difficulty Estimation

- ❖ Data source: examinees' attempts on problems (correct/incorrect)

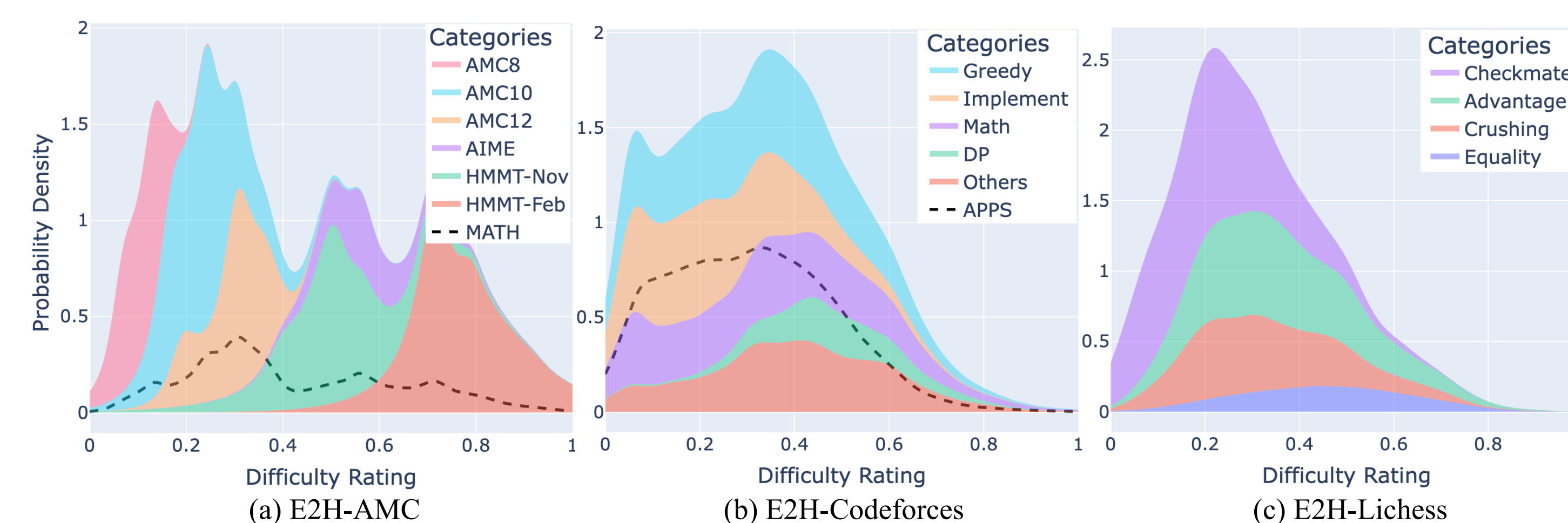


- ❖ The Method depends on whether examinee's ability is static or not

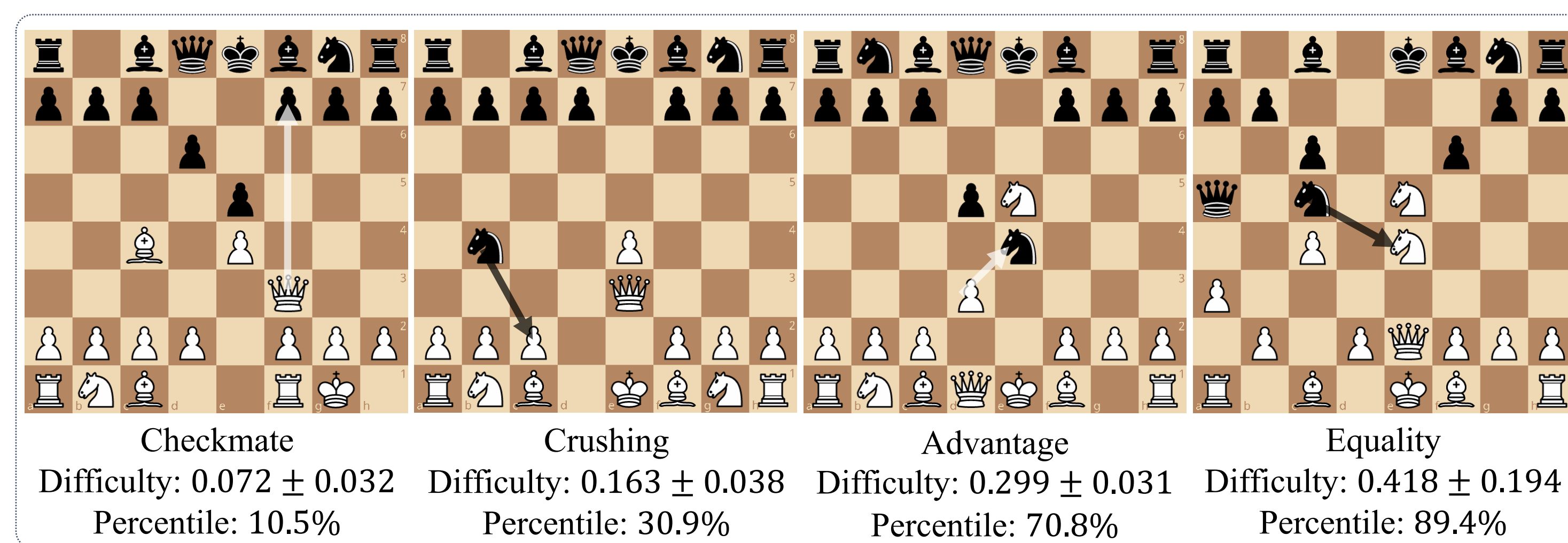


- ❖ The fine-grained standardized difficulty labels

- Shows that the problems within each domain cover a wide range of difficulties
- Verifies that our datasets include more challenging problems

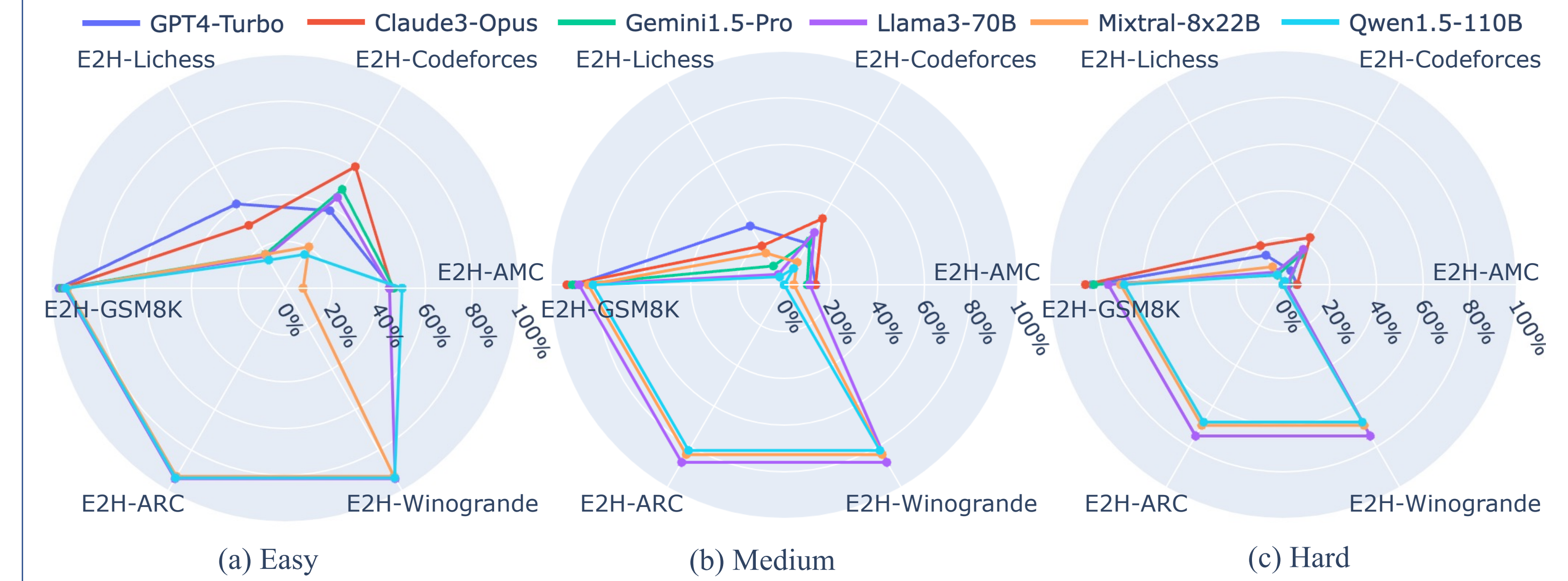


- Corresponds with human recognition

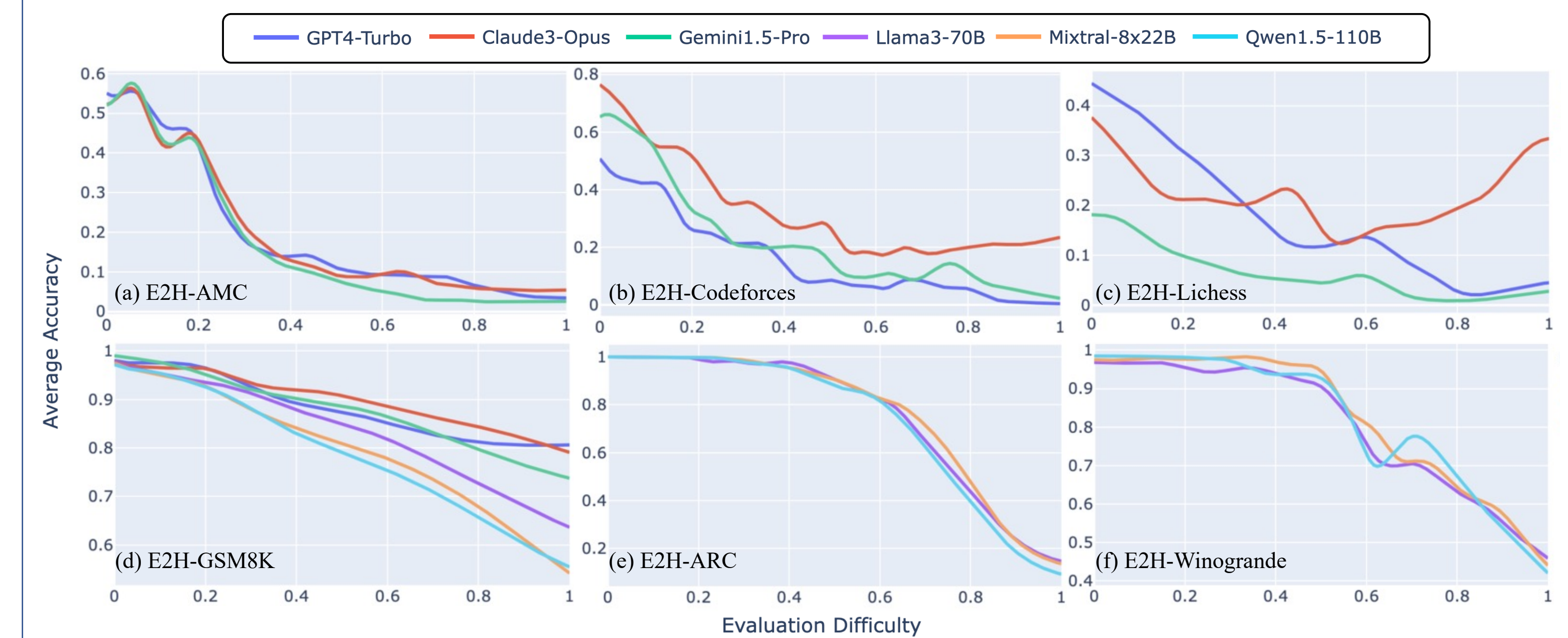


## Benchmarking SoTA LLMs

- ❖ 6 SoTA on all datasets, segmented into easy, medium, and hard levels
  - The performance notably decreases as difficulty increases
  - The newly curated are much more challenging than the pre-existing ones



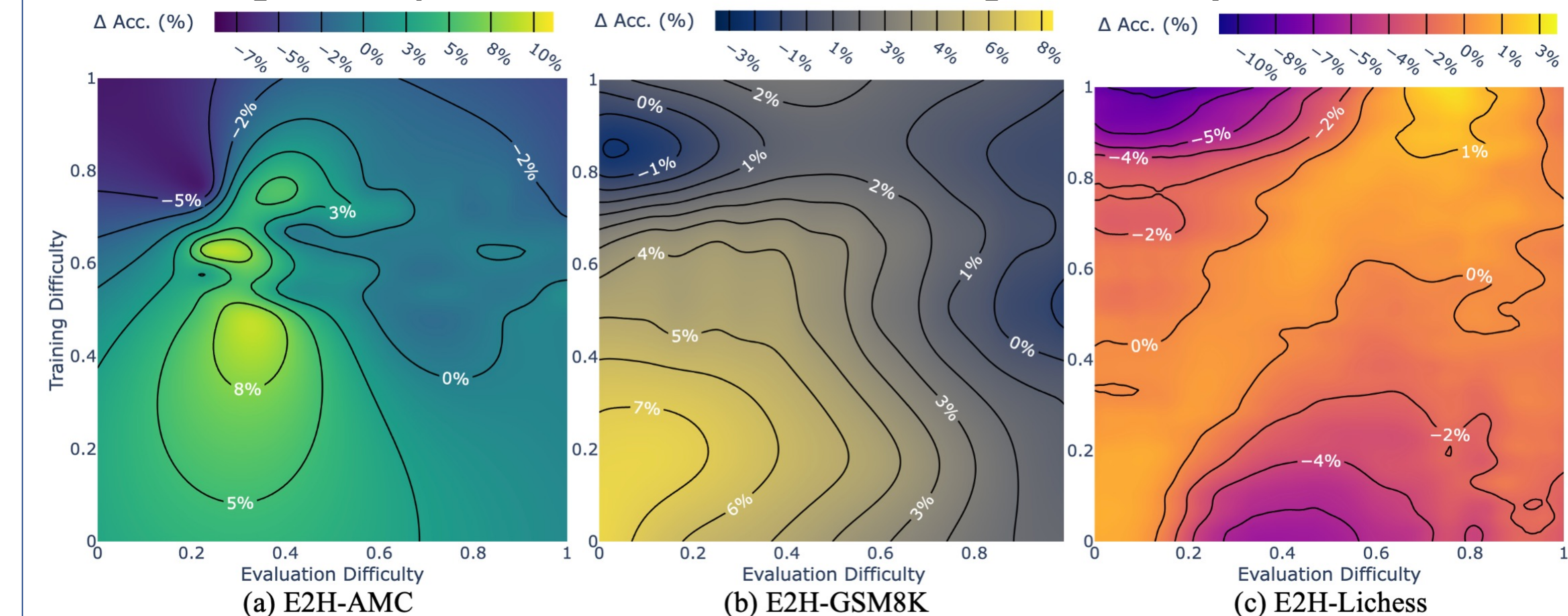
- ❖ Model behavior against increasing difficulty levels for each dataset
  - Most models show monotonic decreasing accuracies
  - The extent of this decline varies significantly among models and datasets



## Profiling Easy2Hard Generalizations

- ❖ The first to deliver detailed easy-to-hard generalization results across continuous, wide-range of difficulties on LLMs.

- Y-axis: trained on subsets of varying difficulty via Supervised Fine-Tuning (SFT)
- X-axis: evaluated across all evaluation difficulties
- Color gradient: performance difference to training on randomly difficulties



- ❖ Our preliminary observation

- Generalization benefits when training and evaluation difficulties are similar
- Training on more challenging samples poses increased generalization difficulties.