The Unreasonable Effectiveness of Easy Training Data for Hard Tasks



Peter Hase^{1,2} Mohit Bansal² Peter Clark¹ Sarah Wiegreffe¹ ¹Allen Institute for AI ²UNC Chapel Hill peter@cs.unc.edu

Main Finding



Accuracy on College STEM Questions

Model fit to 3rd grade questions *almost as good* as model fit to college questions

Mixtral-8x7b model, prompted with 10 examples

- How will models generalize from easy train data to hard test data?
 - Easy = easy to label
 - Hard = hard to label
- Why does this matter?

Effectively supervising models is challenging for many problems of interest

- Is easy-to-hard generalization possible?
 - Pretrained LMs have a lot of latent knowledge and skills

We want to

- Understand how well models generalize based on easy data
 → maybe we only need easy data
- 2. Understand how difficult the *scalable oversight* problem is

We want to

- Understand how well models generalize based on easy data
 → maybe we only need easy data
- 2. Understand how difficult the *scalable oversight* problem is

We want to

- Understand how well models generalize based on easy data
 → maybe we only need easy data
- 2. Understand how difficult the *scalable oversight* problem is

Challenging to train models when outputs are difficult to evaluate (Amodei et al., 2016)

We want to

- Understand how well models generalize based on easy data
 → maybe we only need easy data
- Understand how difficult the scalable oversight problem is
 → maybe scalable oversight is not always difficult

Hase et al.

Measuring Easy-to-Hard Generalization

We introduce the Supervision Gap Recovered

- 89.7 Easy Unsupervised 83.1
- 89.9 Hard Unsupervised 83.1

SGR = 97%



Research Questions

- 1. How Can We Measure Data Hardness? Do Different Approaches Agree?
- 2. Can We Do Well on Hard Data by Training on Easy Data?
- 3. What Are the Cost-Benefit Tradeoffs of Collecting Easy vs. Hard Training Data?
- 4. Is Easy-To-Hard Generalization Consistent Across Model Scale and Train-Test Hardness Gap Size?

What can we measure?

- 1. Education / grade level
- 2. Expert rating
- 3. Required cognitive skills
- 4. Question length
- 5. Answer length
- 6. Compositional reasoning steps
- Model-based hardness (datapoint loss w/ weaker LM)

What can we measure?

- 1. Education / grade level
- 2. Expert rating
- 3. Required cognitive skills
- 4. Question length
- 5. Answer length
- 6. Compositional reasoning steps
- Model-based hardness (datapoint loss w/ weaker LM)

Data we use...

- 3rd grade to college STEM
- Compositional reasoning in math and general-knowledge trivia

| ARC | MMLU-STEM-5 | StrategyQA | GSM8k |
|---------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------|
| n = 4521 | n = 1746 | n = 2290 | n = 8792 |
| Grade Level (3-8) Difficulty Score (1-3) Bloom Skill (1-5) Question Num. Words Answer Num. Chars Num. Reasoning Steps MDL | Grade Level (HS vs. College) Difficulty Score Bloom Skill Question Num. Words Answer Num. Chars Num. Reasoning Steps MDL | Grade Level Difficulty Score Bloom Skill Question Num. Words Answer Num. Chars Num. Reasoning Steps MDL | Grade Level Difficulty Score Bloom Skill Question Num. Words Answer Num. Chars Num. Reasoning Steps MDL |

4 datasets6 human hardness measures1 model-based measure

Hase et al.

RQ1: How Can We Measure Hardness?

We use human and model-based hardness

Diverse measures, all seem to capture something about labeling difficulty

| | Hardness Measure | Easy | Medium | Hard |
|--------------------------------------------------|----------------------------------------|-----------------|--------|-----------------|
| We need to define <i>easy</i> and <i>hard</i> | ARC Grade | 3-5 | 6-7 | 8 |
| | ARC Expert Difficulty | 1 | 2 | 3 |
| | ARC Bloom Skill | 1-2 | 3 | 4-5 |
| | MMLU Grade | High School | | College |
| | StrategyQA Reasoning | 1-2 | 3 | 4-5 |
| | GSM8k Reasoning | 2-3 | 4-5 | 6-11 |
| | Question Length, Answer Length, MDL | 30th percentile | | 70th percentile |

Experiment Setup

- Models
 - Llama-2 models (7b, 13b, 70b)
 - Mixtral-8x7b, Llama-2 70b chat, Qwen-72b
- Training Methods
 - o ICL, *n*≤10
 - Linear probing, *n*=160
 - QLoRA, *n*=160
- Unsupervised Baseline
 - Zero-shot prompted model (better than fully supervised weaker model)
- Results averaged over 5 random seeds

The Supervision Gap Recovered is 70-100% across hardness measures

I lama-2-70b ICL with $k \le 10$











| Dataname | Hardness Measure | SGR Estimate | Test Hardness | n |
|------------|---------------------|--------------------------------|---------------|------|
| ARC | Grade Level | $0.96 \pm 0.10 \ (p < 1e-4)$ | Hard | 1588 |
| ARC | 1/2/3 Difficulty | $0.98 \pm 0.36 \ (p = 0.0033)$ | Hard | 1588 |
| ARC | Bloom Skill | $1.00 \pm 0.20 \ (p < 1e-4)$ | Hard | 1588 |
| MMLU | HS vs. College | $0.97 \pm 0.59 \ (p = 0.0158)$ | Hard | 603 |
| StrategyQA | Num Reasoning Steps | $0.72 \pm 0.93 \ (p = 0.0788)$ | Hard | 427 |
| GSM8k | Num Reasoning Steps | $0.79 \pm 0.60 \ (p = 0.0125)$ | Hard | 333 |

We just saw these SGR values

Hase et al.

Llama-2-70b ICL with k≤10

| Dataname | Hardness Measure | SGR Estimate | Test Hardness | n |
|------------|---------------------|--------------------------------|---------------|------|
| ARC | Grade Level | $0.96 \pm 0.10 \ (p < 1e-4)$ | Hard | 1588 |
| ARC | 1/2/3 Difficulty | $0.98 \pm 0.36 \ (p = 0.0033)$ | Hard | 1588 |
| ARC | Bloom Skill | $1.00 \pm 0.20 \ (p < 1e-4)$ | Hard | 1588 |
| MMLU | HS vs. College | $0.97 \pm 0.59 \ (p = 0.0158)$ | Hard | 603 |
| StrategyQA | Num Reasoning Steps | $0.72 \pm 0.93 \ (p = 0.0788)$ | Hard | 427 |
| GSM8k | Num Reasoning Steps | $0.79 \pm 0.60 \ (p = 0.0125)$ | Hard | 333 |
| ARC | Grade Level | $1.00 \pm 0.09 \ (p < 1e-4)$ | All | 3521 |
| ARC | 1/2/3 Difficulty | $0.96 \pm 0.08 \ (p < 1e-4)$ | All | 3521 |
| ARC | Bloom Skill | $0.98 \pm 0.08 \ (p < 1e-4)$ | All | 3521 |
| MMLU | HS vs. College | $1.00 \pm 0.27 \ (p = 0.0001)$ | All | 1746 |
| StrategyQA | Num Reasoning Steps | $0.87 \pm 0.32 \ (p < 1e-4)$ | All | 2290 |
| GSM8k | Num Reasoning Steps | $0.98 \pm 0.39 \ (p = 0.0003)$ | All | 2065 |

SGR values even higher when testing on "all" data

Hase et al.

Llama-2-70b ICL with k≤10

Hase et al.

RQ2: How Good Is Easy-to-Hard Generalization?

Easy supervision is 70-100% as good as hard supervision

- Previous experiments used equal amounts of cleanly labeled easy and hard data
- This is actually unrealistic
- Hard data is more expensive and labels are noisier
- What if hard data is 2x as costly to collect?
- What if hard data is 2x as noisy as easy data?
 - 2x as much high school data as college data in MMLU
 - Expert error rate in GPQA (grad questions) more than 2x expert error rate in MMLU (undergrad questions)

Easy training data can be better than hard data

Llama-2-70b with linear probe

Testing on MMLU-STEM-5



What if Hard Data Is 2x Costlier to Label?

Easy training data can be better than hard data

Llama-2-70b with linear probe

Testing on MMLU-STEM-5



Collecting easy data can be better than hard data due to data cost and label noise

- What happens as models get better?
- What happens as the train-test hardness gap grows?

The Supervision Gap Recovered Is Similar Across Model Size



ICL with k=10



When train-test gap is big enough...

The supervision gap recovered is robust across model scale Easy-to-hard generalization may decline with very large train-test gaps

Discussion

- Are our tasks hard enough to provide generalizable results?
 - We personally couldn't annotate MMLU
 - We consider 3rd grade to college generalization
- How are the LMs actually doing this?
 - Training elicits some latent knowledge/skill *that is hardness-invariant*
 - Not merely learning the task format
- Why not test for knowledge/skills not in the train data?
 - Wouldn't that be *true* generalization?
 - Our aim is to elicit knowledge we suspect the model may know, without knowing it ourselves – not teach something new

Conclusion

- 1. How Can We Measure Hardness? Diverse human and model-based measurements
- 2. How Good Is Easy-to-Hard Generalization? Easy supervision is 70-100% as good as hard supervision
- **3.** Cost-Benefit Tradeoffs of Easy vs. Hard Data Collecting easy data can be better than hard data
- Scaling Model Size & Train/Test Hardness
 Results robust across model size
 Huge train-test gaps could be an issue

Thank You!

PDFs + code: https://peterbhase.github.io/research/

Contact Info:

Peter Hase, UNC Chapel Hill peter@cs.unc.edu https://peterbhase.github.io

Other Work

2020

Evaluating Explainable AI: Which Algorithmic Explanations Help Users Predict Model Behavior?

Leakage-Adjusted Simulatability: Can Models Generate Non-Trivial Explanations of Their Behavior in Natural Language?

2021

When Can Models Learn From Explanations? A Formal Framework for Understanding the Roles of Explanation Data

The Out-of-Distribution Problem in Explainability and Search Methods for Feature Importance Explanations

Do Language Models Have Beliefs? Methods for Detecting, Updating, and Visualizing Model Beliefs

2022

VisFIS: Visual Feature Importance Supervision with Right-for-the-Right-Reason Objectives

2023

Does Localization Inform Editing? Surprising Differences in Causality-Based Localization vs. Knowledge Editing in Language Models

Can Sensitive Information Be Deleted From LLMs? Objectives for Defending Against Extraction Attacks

Interpretability + Model Control



Examples

MMLU College-level Computer Science Example



Examples



Correct on hard problem given only easy data in prompt

Hase et al.

RQ1: How Can We Measure Hardness?

ARC MDL (QLoRA) 0.09 0.2 0.09 0.09 0.19 0.82 0.65 **Rank** Correlation MDL (Linear Probe) 0.06 0.21 0.05 0.12 0.2 0.41 0.65 1 1.0 MDL (ZS Prompt) 0.08 0.14 0.08 0.06 0.14 0.41 0.82 0.5 Bloom Skill 0.31 0.24-0.030.26 0.14 0.2 0.19 Hardness measures do not 0.01/2/3 Difficulty 0.13 0.16-0.01 0.26 0.06 0.12 0.09 correlate strongly -0.5Grade Level 0.14 0.05 -0.01-0.030.08 0.05 0.09 -1.0Answer Num. Chars 0.11 0.05 0.16 0.24 0.14 0.21 0.2 Question Num. Words 0.11 0.14 0.13 0.31 0.08 0.06 0.09 Linear Probe Dake LS Prompt er BloomSt Answer hum. Grade Question Num.

Hase et al.

RQ1: How Can We Measure Hardness?



















5









10

15











Model-based hardness: Minimum description length (MDL)

- (Voita and Titov, 2020)
- How "long" does it take a model to learn the datapoint?
- Average loss
 - Avg across n = {5, 20, 80, 340, 900} training points
- Training
 - Linear classifier
 - QLoRA
 - Zero-shot "MDL" with n = {0}
- Avg over some "weaker" models
 - Falcon-7b, Mistral-7b, Persimmon-8b, Llama-1-7b

Hase et al.

RQ1: How Can We Measure Hardness?





39

Hase et al.



Hase et al.



41



Hard Test Performance As a Function of Training Hardness (Across Models)

Hard Test Accuracy vs. Train Data Source

80

60

40

20

Train Data Source









Hase et al.

ARC Question Num. Words 100 89.6 87.3 90 82.7

80

70

60 50





Unsupervised Easy



















Results robust across training methods



Easy is barely worse than Medium

Llama-2-70b ICL with k≤10



Test Data Leakage?



Train Data Source

Task Format Prompts - Hard Test Data



Hard Test Accuracy vs. Train Data Source

Train Data Source

Task Format Prompts - All Test Data



Effect of Reasoning



Differences with Weak-to-Strong Paper

- 1. The baseline in SGR vs. PGR
- 2. We train on easy or hard data, not both
- 3. Human hardness variables in addition to model-based
- 4. All experiments with publicly available data and models (up to 70b params)
- 5. No early stopping
- 6. No new methods in our paper