

# CS388: Natural Language Processing

## Lecture 10: Evaluation Principles, Dataset Artifacts

Greg Durrett



TEXAS

The University of Texas at Austin





# Announcements

---

- ▶ Final project proposals due next Tuesday
- ▶ P3 released next week



# Recap

---

- ▶ **Pretraining (BERT):**
  - ▶ Train a big model to fill in masked-out words, then adapt it to other tasks. Led to big gains in **question answering** and **NLI** performance. BART/T5, GPT-3, etc. push this further and extend it to other tasks
- ▶ **Decoding methods:** nucleus sampling > greedy for open-ended tasks
- ▶ **Two tasks we'll focus on today: Question answering (QA)...**
  - ▶ "What was Marie Curie the first female recipient of?"  
-> "The Nobel Prize" (find this span in a document)
- ▶ ...and **NLI**
  - ▶ "But I thought you'd sworn off coffee."  
*contradicts* "I thought that you vowed to drink more coffee."



# Today

---

- ▶ Evaluation in NLP: benchmarks and generalization
- ▶ Spurious correlations / dataset artifacts
- ▶ Debiasing

# Cross-Dataset Evaluation



# Principles of Evaluation Suites

---

- ▶ Training and testing on i.i.d. data with big neural models often yields very high performance
- ▶ “Solving” a task (getting human-level performance) may be useful, but often can’t tell us about our models more broadly
- ▶ A parable of single-task evaluation: SWAG



# SWAG: Weaknesses of single tasks

On stage, a woman takes a seat at the piano. She

- a) sits on a bench as her sister plays with the doll.
- b) smiles with someone as the music plays.
- c) is in the crowd, watching the dancers.
- d) **nervously sets her fingers on the keys.**

Rowan Zellers et al., 2019

- ▶ Text-only data comes from video captions in ActivityNet
- ▶ **Adversarial filtering** to produce the negative multiple choice answers
- ▶ They said models get ~60% and humans get ~85%, but BERT immediately solved this dataset when it was released

**while** convergence not reached **do**

- Split the dataset  $\mathcal{D}$  randomly up into training and testing portions  $\mathcal{D}^{tr}$  and  $\mathcal{D}^{te}$ .
- Optimize a model  $f_\theta$  on  $\mathcal{D}^{tr}$ .

**for** index  $i$  in  $\mathcal{D}^{te}$  **do**

- Identify easy indices:

$$\mathcal{A}_i^{easy} = \{j \in \mathcal{A}_i : f_\theta(x_i^+) > f_\theta(x_{i,j}^-)\}$$

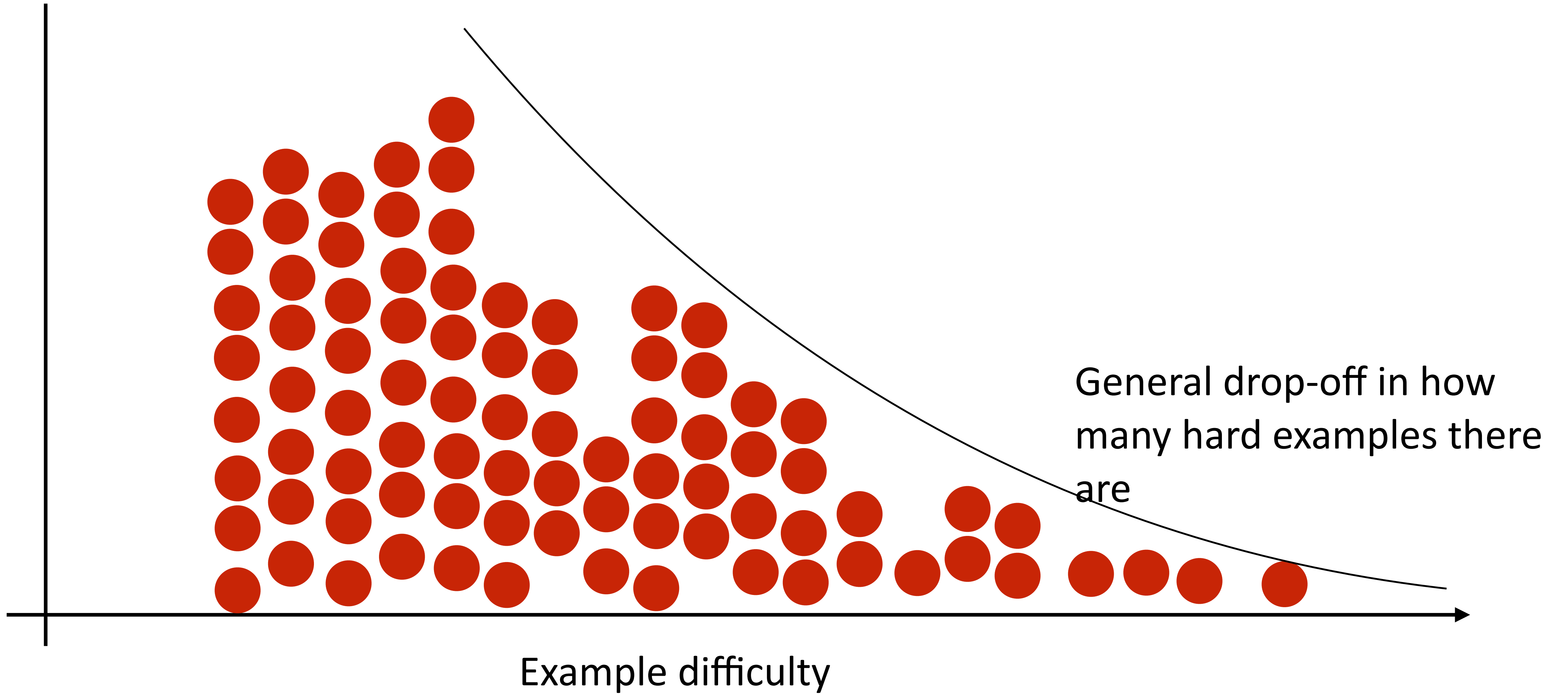
- Replace  $N^{easy}$  easy indices  $j \in \mathcal{A}_i^{easy}$  with adversarial indices  $k \notin \mathcal{A}_i$  satisfying  $f_\theta(x_{i,k}^-) > f_\theta(x_{i,j}^-)$ .

**end for**

**end while**



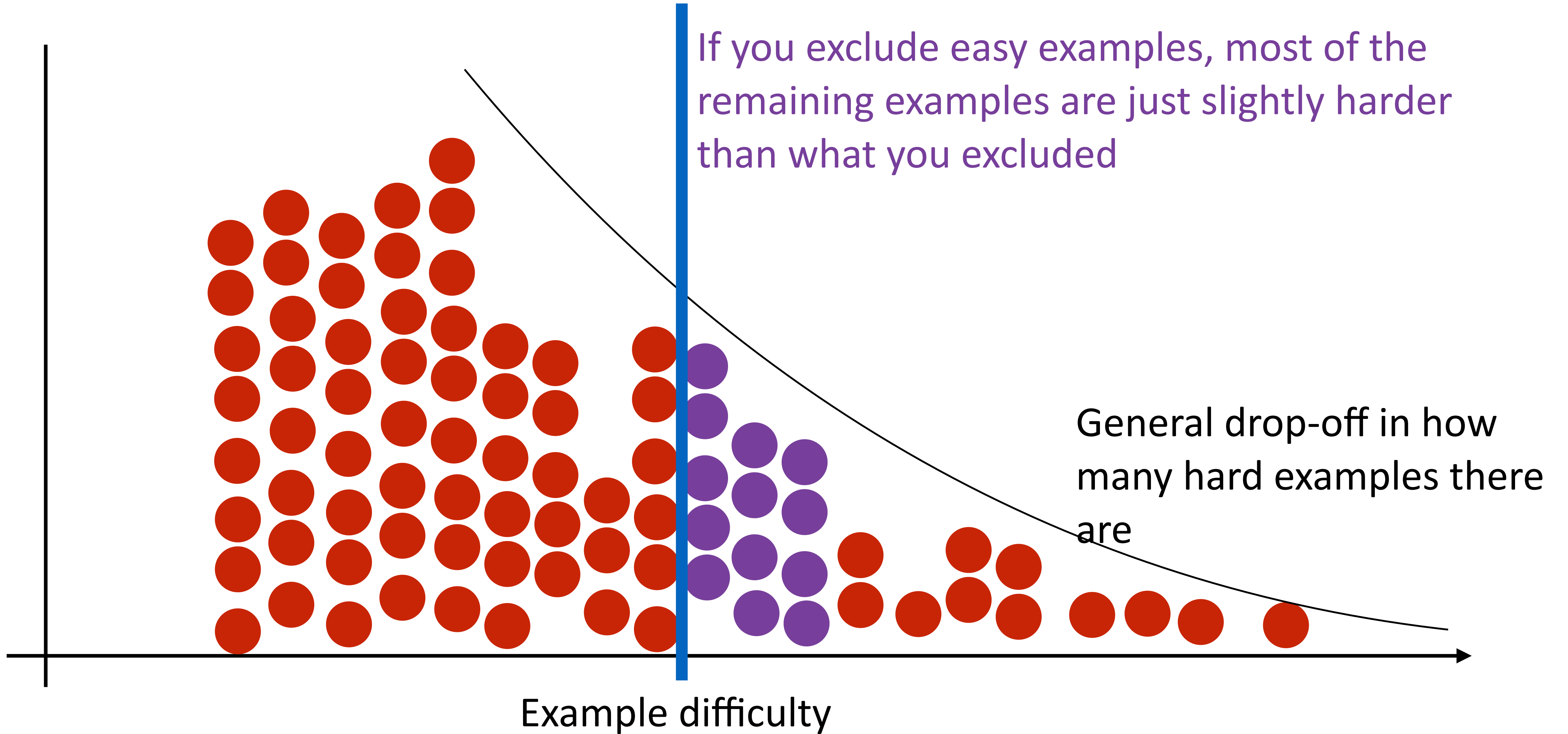
# Intuition







# Intuition





# Principles of Evaluation Suites

---

- ▶ Training and testing on i.i.d. data with big neural models often yields very high performance
- ▶ “Solving” a task (getting human-level performance) may be useful, but often can’t tell us about our models more broadly
- ▶ **Designing a single, difficult task is really challenging!**
- ▶ To assess big models, we need **evaluation suites (benchmarks)** like GLUE
- ▶ What makes a good evaluation suite of tasks?



# Principles of Evaluation Suites

---

- ▶ Difficulty: even if some task can be solved by hand-engineering, it should be hard to solve all  $N$  tasks
- ▶ Diverse: doing well on it should say something useful
- ▶ Good “yardstick”: should understand where human performance is and what good performance on the task would mean
- ▶ GLUE was the first of these, but it wasn’t really diverse and it was too easy. Next step: SuperGLUE



# SuperGLUE: Task Requirements

---

- ▶ Task substance: *“Tasks should test a system’s ability to understand and reason about texts in English.”*
- ▶ Task difficulty: *“Tasks should be beyond the scope of current state-of-the-art systems, but solvable by most college-educated English speakers.”* (notably they excluded domain-specific tasks, which have become more popular these days, e.g., the bar exam)
- ▶ Evaluatable: this is challenging to find!
- ▶ Public dataset, good license, etc.



# SuperGLUE: Performance

Model	Avg	BoolQ	CB	COPA	MultiRC	ReCoRD	RTE	WiC	WSC	AX <sub>b</sub>	AX <sub>g</sub>
Metrics		Acc.	F1/Acc.	Acc.	F1 <sub>a</sub> /EM	F1/EM	Acc.	Acc.	Acc.	MCC	GPS Acc.
Most Frequent	47.1	62.3	21.7/48.4	50.0	61.1 / 0.3	33.4/32.5	50.3	50.0	65.1	0.0	100.0/ 50.0
CBoW	44.3	62.1	49.0/71.2	51.6	0.0 / 0.4	14.0/13.6	49.7	53.0	65.1	-0.4	100.0/ 50.0
BERT	69.0	77.4	75.7/83.6	70.6	70.0 / 24.0	72.0/71.3	71.6	<b>69.5</b>	<b>64.3</b>	23.0	97.8 / 51.7
BERT++	<b>71.5</b>	79.0	<b>84.7/90.4</b>	73.8	70.0 / 24.1	72.0/71.3	79.0	<b>69.5</b>	<b>64.3</b>	38.0	99.4 / 51.4
Outside Best	-	<b>80.4</b>	- / -	<b>84.4</b>	<b>70.4*/24.5*</b>	<b>74.8/73.0</b>	<b>82.7</b>	-	-	-	- / -
Human (est.)	89.8	89.0	95.8/98.9	100.0	81.8*/51.9*	91.7/91.3	93.6	80.0	100.0	77.0	99.3 / 99.7

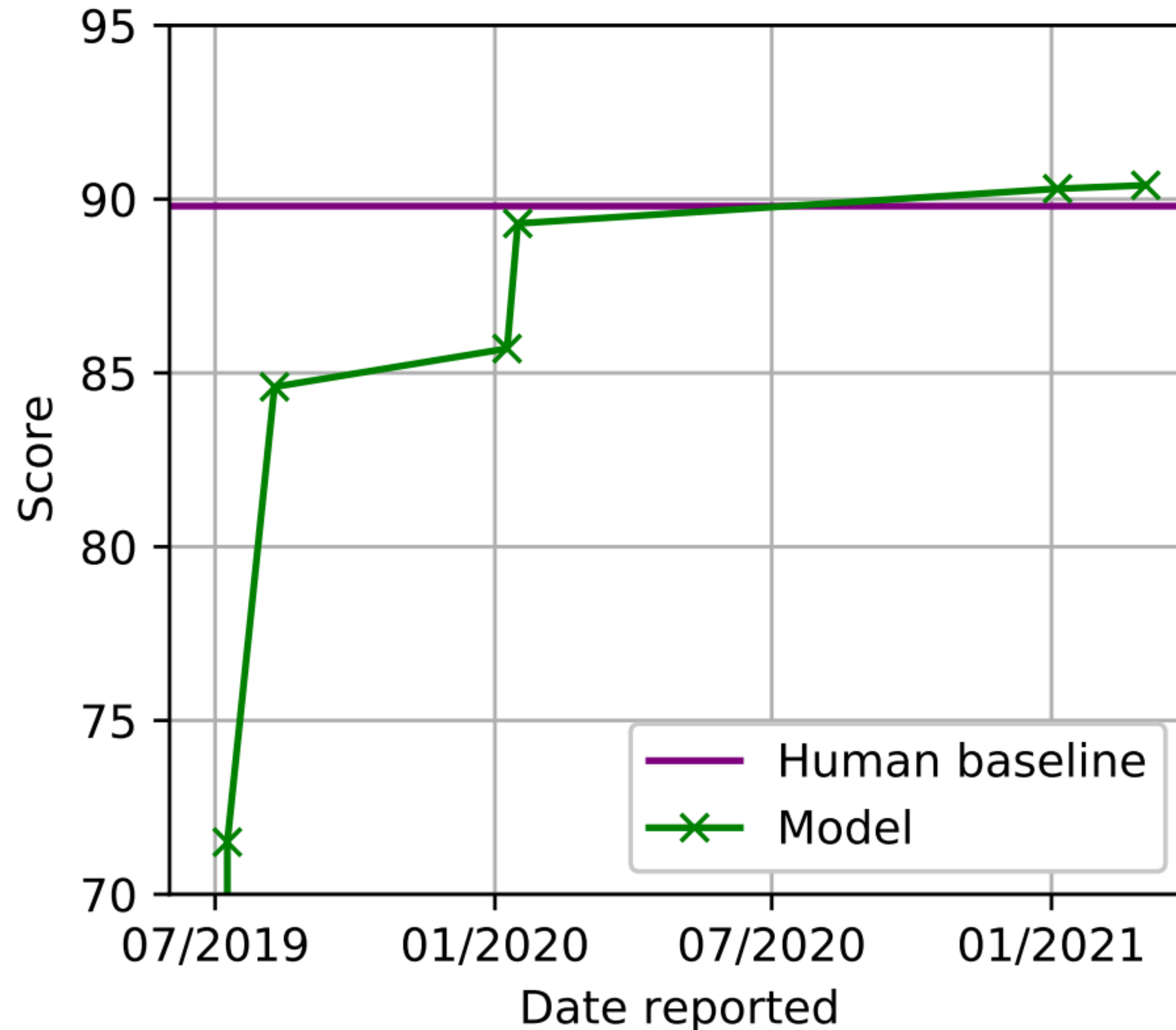
- ▶ RoBERTa in 2019: 84.6
- ▶ DeBERTa in 2020: 90.3. **Even SuperGLUE was solved quickly!**



# SuperGLUE: Performance

As reported in BIGBench:

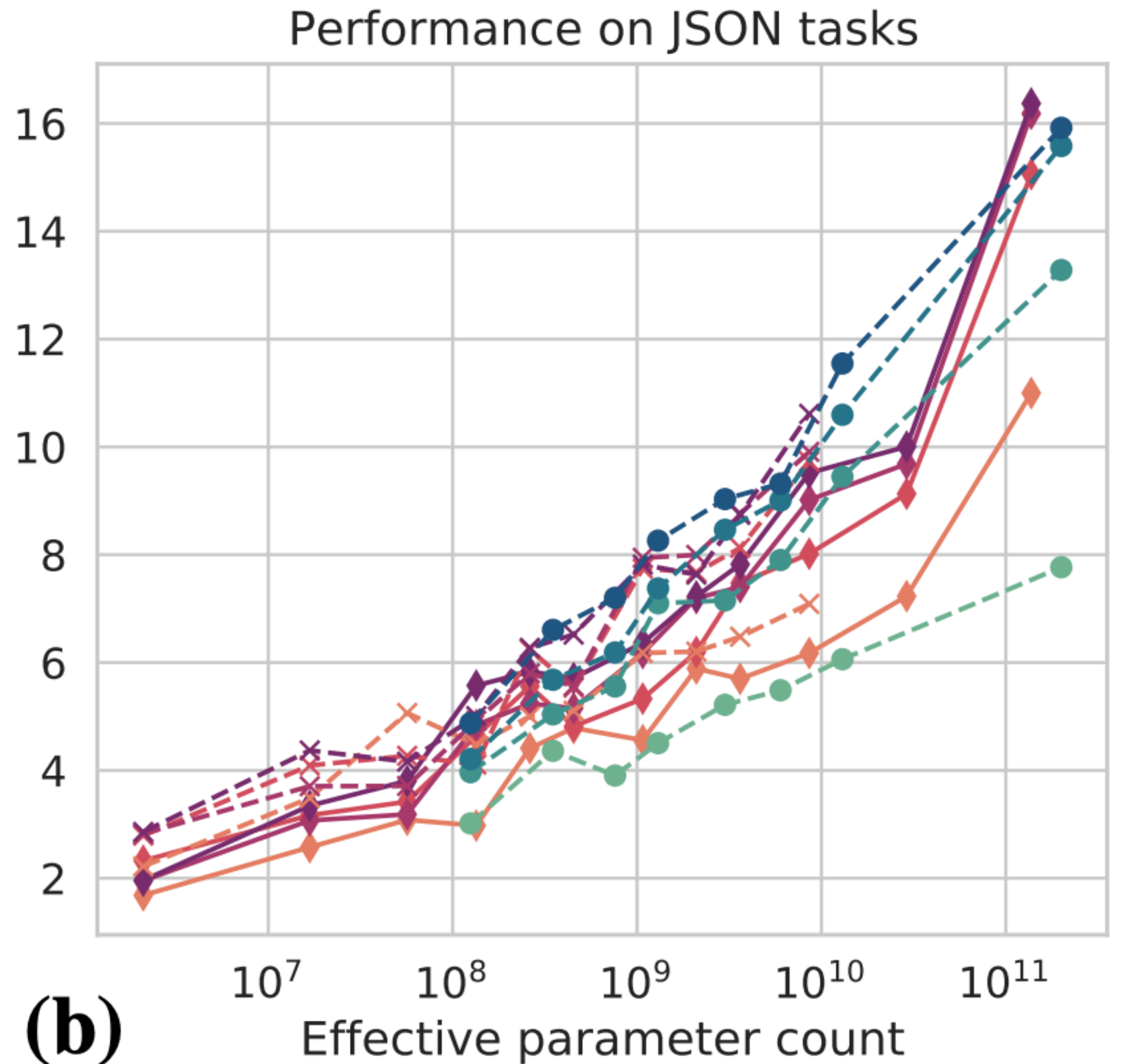
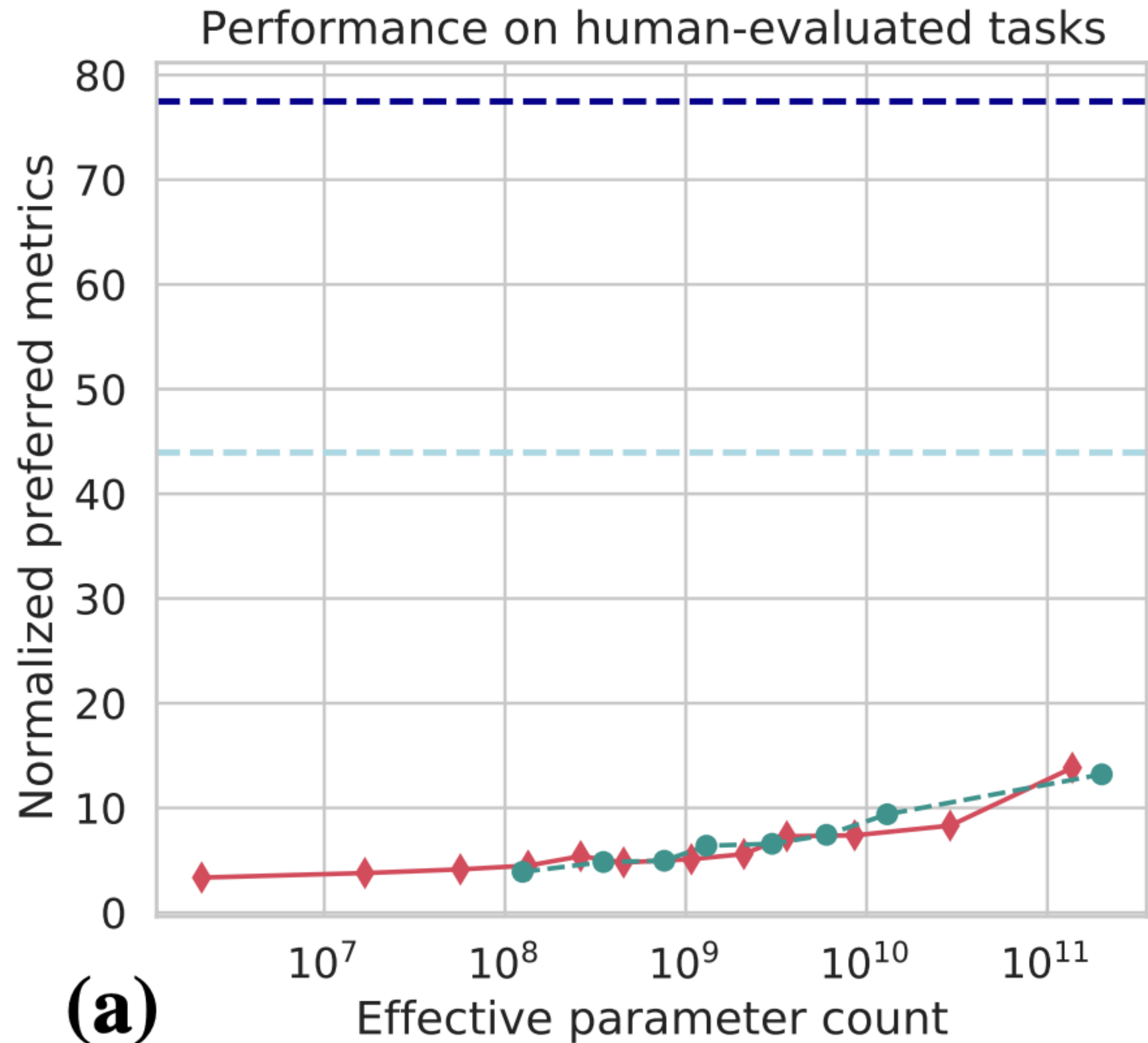
SuperGLUE state of the art over time





# BIG-bench

- ▶ 204 tasks, 444 authors





# BIG-bench

- ▶ “Beyond the Imitation Game” — aim to learn more than what’s possible from model vs. human performance
- ▶ Particular emphasis on scaling
- ▶ Primarily for pre-trained models without fine-tuning. Therefore, not all tasks have large training (or even test!) sets







# MMLU

- ▶ MMLU task (Hendrycks et al., 2020): 57 high school/college/professional exams:

<b>Conceptual Physics</b>	When you drop a ball from rest it accelerates downward at $9.8 \text{ m/s}^2$ . If you instead throw it downward assuming no air resistance its acceleration immediately after leaving your hand is	
	(A) $9.8 \text{ m/s}^2$	✓
	(B) more than $9.8 \text{ m/s}^2$	✗
	(C) less than $9.8 \text{ m/s}^2$	✗
	(D) Cannot say unless the speed of throw is given.	✗
<b>College Mathematics</b>	In the complex $z$ -plane, the set of points satisfying the equation $z^2 =  z ^2$ is a	
	(A) pair of points	✗
	(B) circle	✗
	(C) half-line	✗
	(D) line	✓

Figure 4: Examples from the Conceptual Physics and College Mathematics STEM tasks.



# MMLU

- ▶ MMLU task (Hendrycks et al., 2020): 57 high school/college/professional exams:

---

-	Random	25.0
-	Average human rater	34.5
May 2020	GPT-3 5-shot	43.9
Mar. 2022	Chinchilla 5-shot	67.6
Apr. 2022	PaLM 5-shot	69.3
Oct. 2022	<b>Flan-PaLM 5-shot</b>	<b>72.2</b>
	<b>Flan-PaLM 5-shot: CoT + SC</b>	<b>75.2</b>
-	Average human expert	89.8



# MMLU

Model	Finetuning Mixtures	Tasks	Norm. avg.	MMLU		BBH	
				Direct	CoT	Direct	CoT
540B	None (no finetuning)	0	49.1	71.3	62.9	49.1	63.7
	CoT	9	52.6 (+3.5)	68.8	64.8	50.5	61.1
	CoT, Muffin	89	57.0 (+7.9)	71.8	66.7	56.7	64.0
	CoT, Muffin, T0-SF	282	57.5 (+8.4)	72.9	<u>68.2</u>	57.3	64.0
	CoT, Muffin, T0-SF, NIV2	1,836	<u>58.5</u> (+9.4)	<u>73.2</u>	68.1	<u>58.8</u>	<u>65.6</u>

- ▶ Human performance estimates are ~90 on MMLU, ~80 on Big-Bench (BBH). Even getting close on **these** tasks!

# Evaluation Under Distribution Shift



# Model Performance

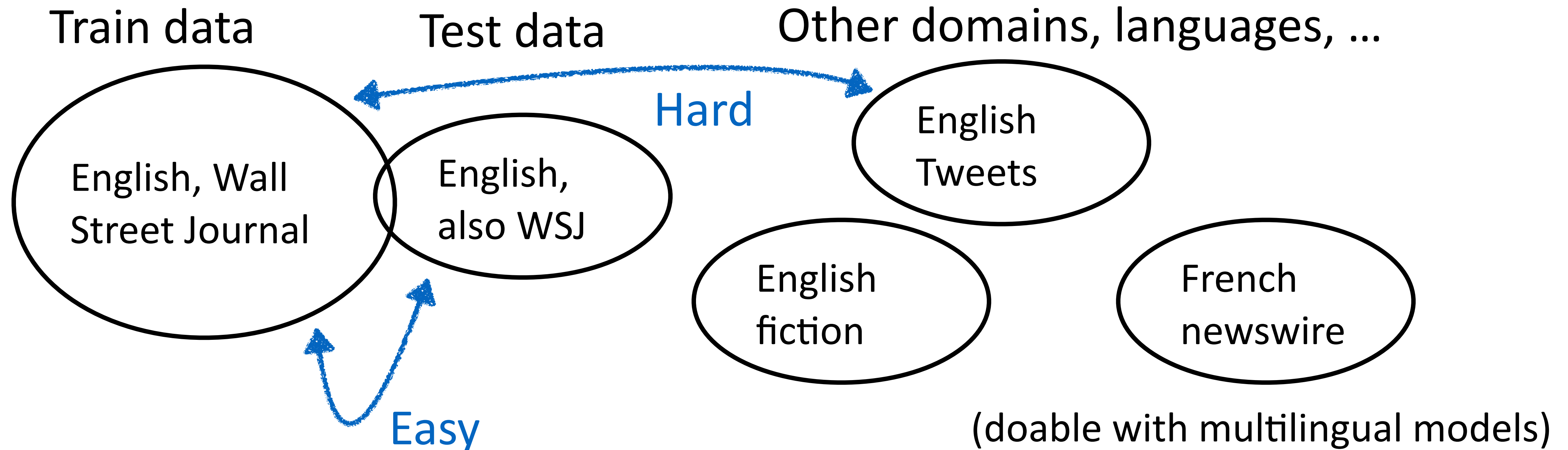
---

- ▶ If models can be fine-tuned on each of  $n$  tasks in an evaluation suite and perform very well on the held-out test dataset, have we solved everything we want?
- ▶ What can go wrong?



# Generalization

- ▶ If a model does well on train but poorly on test data, it *doesn't generalize*
- ▶ A model can do well on its test data and still fail to generalize *out of distribution* — arguably an even more important notion
- ▶ Many notions of generalization. Example: POS tagging





# Generalization: QA

Train data

SQuAD: factoid questions with answers on Wikipedia

Test data

SQuAD

Other *domains*

Science questions

Unanswerable questions

French questions

Other types of reasoning, such as *multi-hop questions*

*Who won the Nobel in Chemistry the year Marie Curie won the Nobel in Physics?*



# Generalization

---

- ▶ Just doing well on a single test set **is not that useful**
- ▶ We want POS taggers, QA systems, and more that can generalize to new settings so we can deploy them in practice
- ▶ Sometimes, you can get **very good test performance** but the model **generalizes very poorly**. How does this happen?



# Annotation Artifacts, Reasoning Shortcuts: QA



# Annotation Artifacts

---

- ▶ Some datasets might be easy because of how they're constructed, especially in QA and NLI

*What becomes of Macbeth?*

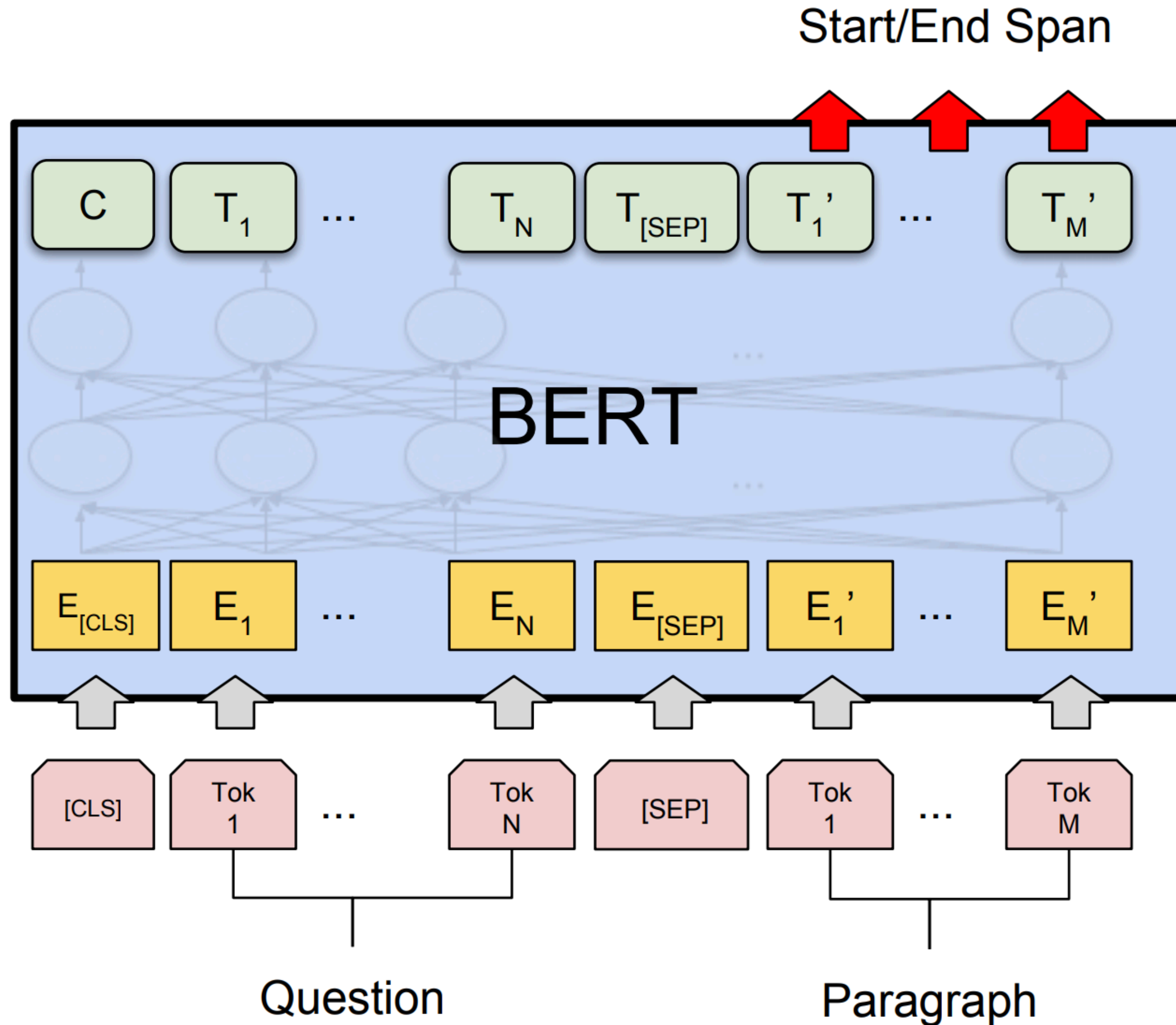
*What does Macduff do to Macbeth?*

*What violent act does Macduff perform upon Macbeth?*

- ▶ All questions have the same answer. But some are more easily guessable



# Reminder: QA with BERT





# QA: Answer Type Heuristics

---

*What degree did Martin Luther receive on October 19, 1512?*

*On October 19, 1512, Luther was awarded his doctorate of theology and, on October 21, 1512, was received into the senate of the theological faculty of the University of Wittenberg. He spent the rest of his career in this position at the University of Wittenberg.*

- ▶ What should the model be doing? Corresponding Martin Luther with Luther, matching October 19, 1512 between question and passage



# QA: Answer Type Heuristics

---

*What degree did Martin Luther receive?*

*What degree \_\_\_\_?*

*On October 19, 1512, Luther was awarded his doctorate of theology and, on October 21, 1512, was received into the senate of the theological faculty of the University of Wittenberg. He spent the rest of his career in this position at the University of Wittenberg.*

- ▶ Only one possible degree here! Model only needs to see “what degree” and will not learn to use the rest of the context!



# QA: Answer Type Heuristics

---

- ▶ Question type is powerful indicator. Only a couple of locations in this context!

*Where \_\_\_\_\_?*

*On October 19, 1512, Luther was awarded his doctorate of theology and, on October 21, 1512, was received into the senate of the theological faculty of the University of Wittenberg. He spent the rest of his career in this position at the University of Wittenberg.*

*Who \_\_\_\_\_?*

*When \_\_\_\_\_?*



# QA: Answer Type Heuristics

---

- ▶ Question type is powerful indicator. Only a couple of locations in this context!

*Where \_\_\_\_\_?      Who \_\_\_\_\_?      When \_\_\_\_\_?*

*On October 19, 1512, Luther was awarded his doctorate of theology and, on October 21, 1512, was received into the senate of the theological faculty of the University of Wittenberg. He spent the rest of his career in this position at the University of Wittenberg.*

- ▶ What will happen if we train on this data?
  - ▶ Will loss decrease?
  - ▶ How will the model learn to “behave”?

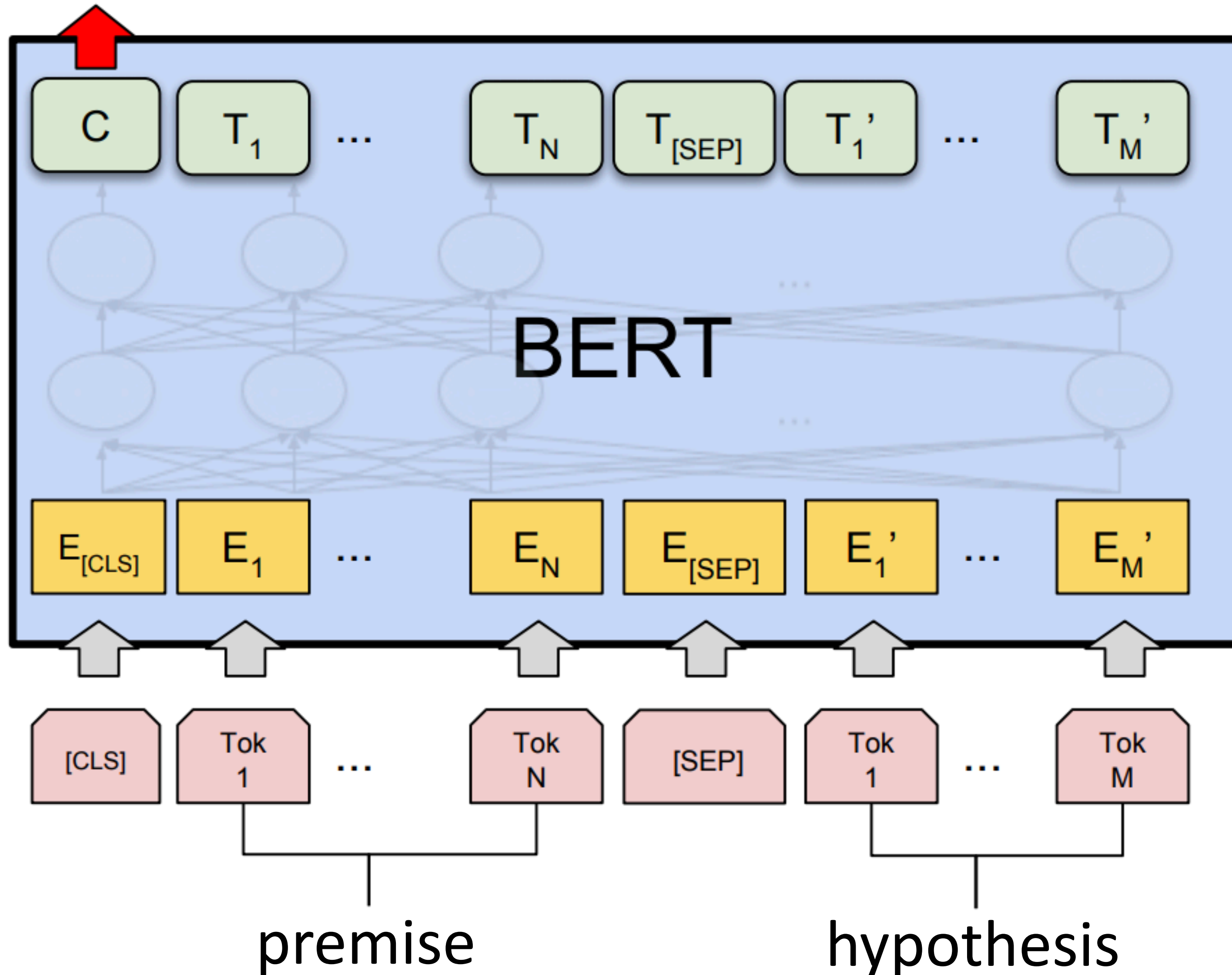
# Annotation Artifacts, Reasoning Shortcuts: NLI





# Reminder: NLI with BERT

entailed/neutral/contradiction





# NLI: Hypothesis-only Baselines

---

Premise: *A woman on a deck is selling bamboo sticks.*

Label?

Hypothesis: *A man is selling bamboo sticks*

Hypothesis: *A man is juggling flaming chainsaws*

Hypothesis: *Eighteen flying monkeys are in outer space*

- ▶ Not all of these things have the same likelihood of being true a priori
- ▶ What might the model learn to do in this case?



# NLI: Hypothesis-only Baselines

---

<b>Premise</b>	A woman selling bamboo sticks talking to two men on a loading dock.
<b>Entailment</b>	There are <b>at least three people</b> on a loading dock.
<b>Neutral</b>	A woman is selling bamboo sticks <b>to help provide for her family.</b>
<b>Contradiction</b>	A woman is <b>not</b> taking money for any of her sticks.

---

- ▶ What's different about this neutral sentence?
  - ▶ To create neutral sentences: annotators *add information*
- ▶ What's different about this contradictory sentence?
  - ▶ To create contradictions: annotators *add negation*
- ▶ These are not broadly representative of what can happen in other settings. There is no “natural” distribution of NLI, but this is still very restrictive



# NLI: Hypothesis-only Baselines

---

<b>Premise</b>	A woman selling bamboo sticks talking to two men on a loading dock.
<b>Entailment</b>	There are <b>at least three people</b> on a loading dock.
<b>Neutral</b>	A woman is selling bamboo sticks <b>to help provide for her family</b> .
<b>Contradiction</b>	A woman is <b>not</b> taking money for any of her sticks.

---

- ▶ Models can detect new information or negation easily
- ▶ Models can do very well *without looking at the premise*

Performance of models that only look at the hypothesis:  
~70% on 3-class SNLI dataset

---

	Hyp-only model	Majority class	
SNLI	69.17	33.82	<b>+35.35</b>
MNLI-1	55.52	35.45	<b>+20.07</b>
MNLI-2	55.18	35.22	<b>+19.96</b>

---



# NLI: Heuristics (HANS)

Heuristic	Definition	Example
Lexical overlap	Assume that a premise entails all hypotheses constructed from words in the premise	<b>The doctor was paid by the actor.</b> ————→ The doctor paid the actor. WRONG
Subsequence	Assume that a premise entails all of its contiguous subsequences.	The doctor near <b>the actor danced.</b> ————→ The actor danced. WRONG
Constituent	Assume that a premise entails all complete subtrees in its parse tree.	If <b>the artist slept</b> , the actor ran. ————→ The artist slept. WRONG

- ▶ Word overlap supersedes actual reasoning in these cases
- ▶ They create a test set (HANS) consisting of cases where heuristics like word overlap are misleading. Very low performance



# Evidence of Spurious Correlations: Contrast Sets

---

- ▶ How do we control for annotation artifacts? Things like “premises and hypotheses overlap too much” aren’t easy to see!
- ▶ For any particular effect like lexical overlap, we could try to annotate data that “breaks” that effect
- ▶ Issue: breaking one correlation may just result in another one surfacing. How do we “break” them all at the same time?
- ▶ Solution: construct new examples through *minimal edits that change the label*.



# Evidence of Spurious Correlations: Contrast Sets

---

Hardly one to be faulted for his ambition or his vision, it is genuinely unexpected, then, to see all Park's effort add up to so very little. ... The premise is promising, gags are copious and offbeat humour abounds but it all fails miserably to create any meaningful connection with the audience.

*(Label: Negative)*

Hardly one to be faulted for his ambition or his vision, **here we see all Park's effort come to fruition.** ... The premise is **perfect**, gags are **hilarious** and offbeat humour abounds, **and it creates a deep** connection with the audience.

*(Label: Positive)*

- ▶ By minimally editing an example, we control for pretty much all of the possible shortcuts that apply to the original.
- ▶ E.g., [summary starts with "*Hardly*" -> negative] is a pattern that could not hold anymore



# Evidence of Spurious Correlations: Contrast Sets

<b>Dataset</b>	<b># Examples</b>	<b># Sets</b>	<b>Model</b>	<b>Original Test</b>		<b>Contrast</b>
NLVR2	994	479	LXMERT	76.4	61.1	(-15.3)
IMDb	488	488	BERT	93.8	84.2	(-9.6)
MATRES	401	239	CogCompTime2.0	73.2	63.3	(-9.9)
UD English	150	150	Biaffine + ELMo	64.7	46.0	(-18.7)
PERSPECTRUM	217	217	RoBERTa	90.3	85.7	(-4.6)
DROP	947	623	MTMSN	79.9	54.2	(-25.7)



# Solutions



# Broad Solutions

---

- ▶ Most solutions involve changing what data is trained on
  - ▶ Subset of data
  - ▶ Soft subset (i.e., reweight the existing examples)
  - ▶ Superset: add adversarially-constructed data, contrast sets, etc.
- ▶ For subsets: what do we train on?
  - ▶ Don't train on stuff that allows you to cheat
  - ▶ Train on examples that teach the real task rather than shortcuts



# Dataset Cartography

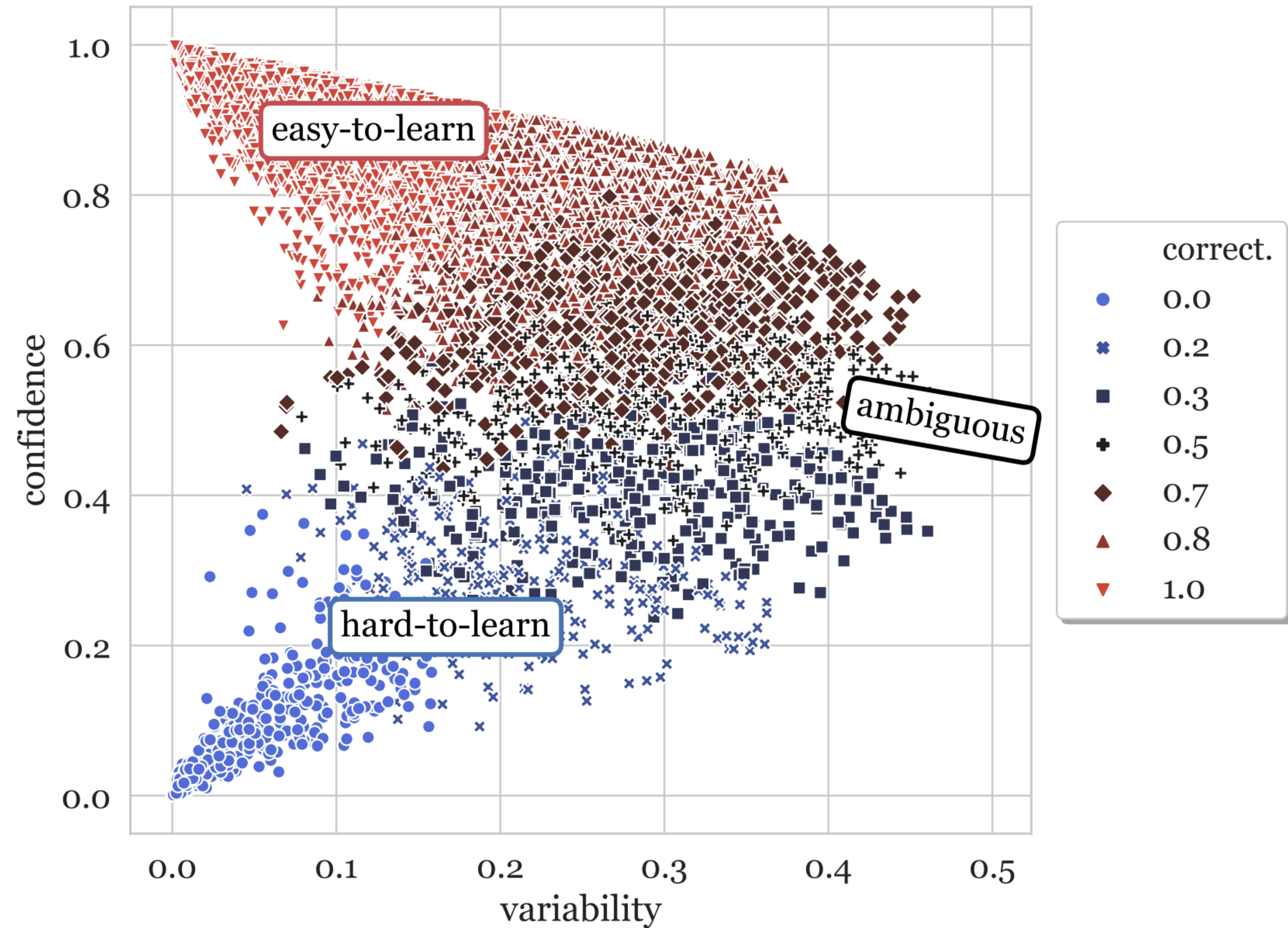
---

- ▶ What happens with each particular example during training?
- ▶ Spurious correlations are *easy to learn*: a model should learn these early and always get them right
- ▶ Imagine a very challenging example
  - ▶ Model prediction may change a lot as it learns this example, may be variable in its predictions
- ▶ Imagine a mislabeled example
  - ▶ Probably just always wrong unless it gets overfit



# Data Maps

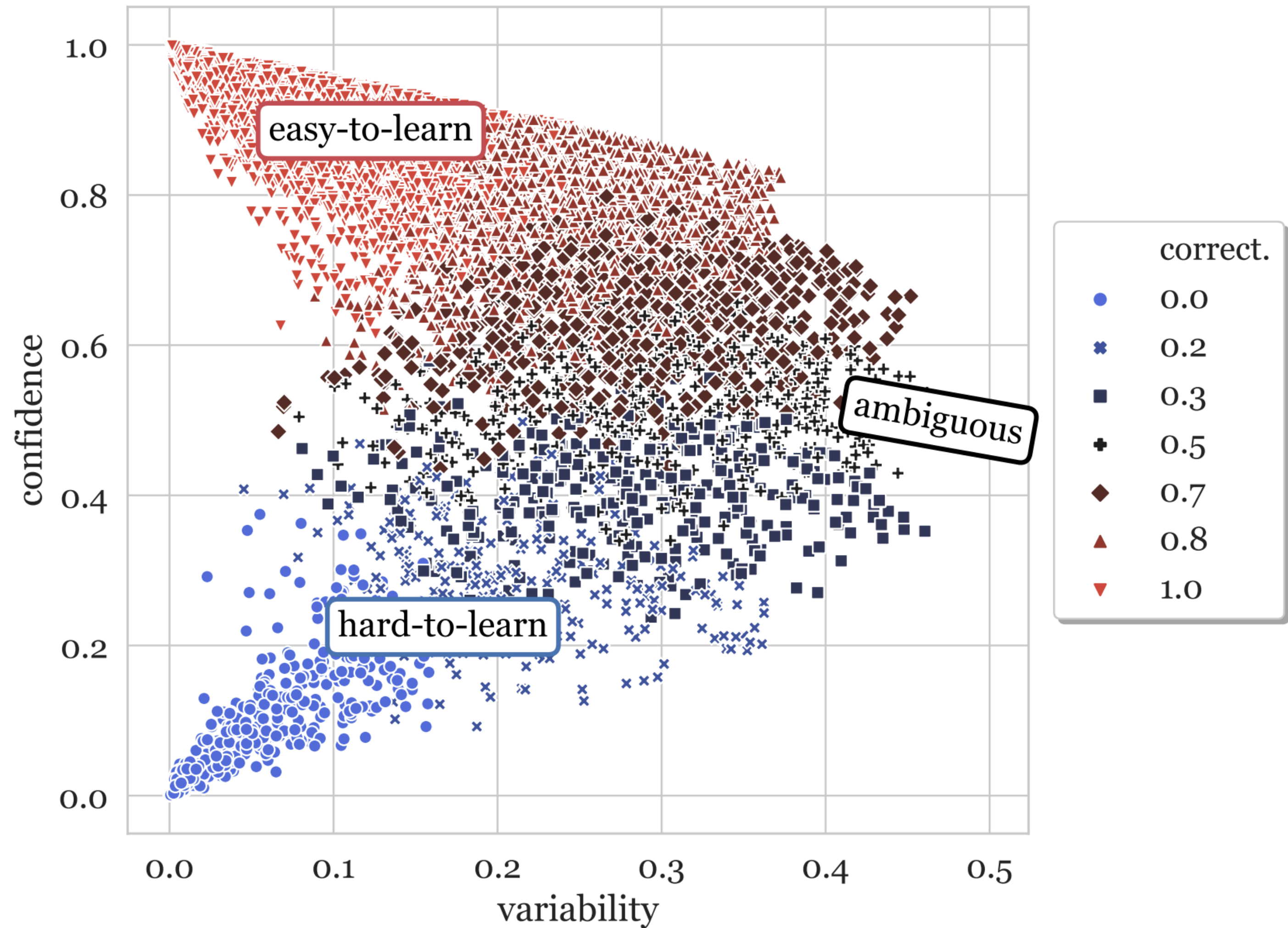
- ▶ Confidence: mean probability of correct label
- ▶ Variability: standard deviation in probability of the correct label
- ▶ Ambiguous examples: possible learnable (model knows it sometimes but not other times), but hard!





# Data Maps

- ▶ What to do with them?
- ▶ Training on hard-to-learn or ambiguous examples leads to better performance out-of-domain





# Debiasing

- ▶ Other ways to identify easy examples other than data maps
- ▶ Train some kind of a weak model and discount examples that it fits easily

$$\mathcal{L}(\theta_d) = - \left(1 - p_b^{(i,c)}\right) y^{(i)} \cdot \log p_d$$

one-hot label vector

log probability of each label

probability under a copy of the model trained for a few epochs on a small subset of data (bad model)



# Debiasing

Method	MNLI (Acc.)		
	dev	<i>HANS</i>	$\Delta$
BERT-base	84.5	61.5	-
Reweighting <small>known-bias</small>	83.5 <sup>‡</sup>	69.2 <sup>‡</sup>	+7.7
Reweighting <small>self-debias</small>	81.4	68.6	+7.1
Reweighting <small>♠ self-debias</small>	82.3	69.7	<b>+8.2</b>

- ▶ On the challenging HANS test set for NLI, this debiasing improves performance substantially
- ▶ In-domain MNLI performance goes down



# Debiasing

- ▶ Other work has explored similar approaches using a known bias model

$$\hat{p}_i = \textit{softmax}(\log(p_i) + \log(b_i))$$

probabilities from learned bias model — like the weak model from Utama et al. (prev. slides), but you define its structure

- ▶ *Ensembles* the weak model with the model you actually learn.
- ▶ Your actual model learns the *residuals* of the weak model: the difference between the weak model's output distribution and the target distribution.
- ▶ This lets it avoid learning the weak model's biases!





# Takeaways

---

- ▶ Strong neural models trained on “tough” datasets may fail to generalize because they learn annotation artifacts
- ▶ By reweighting data or changing the training paradigm, you can learn a model that generalizes better
- ▶ Most gains will show up **out-of-domain**. Very hard to get substantial improvements on the same dataset, unless you consider small subsets of examples (e.g., the toughest 1% of examples by some measure)
- ▶ As more “generalist” LLMs are learned, this problem goes away...but there’s always a tradeoff when you want to fine-tune them for certain tasks
- ▶ Next time: further understanding in-context learning