

# Utility is in the Eye of the User: A Critique of NLP Leaderboards

EMNLP 2020







**Kawin Ethayarajh**



**Dan Jurafsky**



# Benchmark-based leaderboards have helped drive the creation of more accurate models.

Rank	Name	Model	URL	Score
1	HFL iFLYTEK	MacALBERT + DKM		90.7
2	Alibaba DAMO NLP	StructBERT + TAPT		90.6
3	PING-AN Omni-Sinitic	ALBERT + DAAF + NAS		90.6
4	ERNIE Team - Baidu	ERNIE		90.4
5	T5 Team - Google	T5		90.3
14	GLUE Human Baselines	GLUE Human Baselines		87.1

**Benchmark-based leaderboards have helped drive the creation of more accurate models.**

Rank	Name	Model	URL	Score
	 <b>GLUE</b>	acAl ructl	 The Natural Language Decathlon	
3	PING-AN Omni-Sinitic	ALBERT + DART + NAS		90.0
	 <b>SuperGLUE</b>			
14	GLUE Human Baselines	GLUE Human Baselines		87.1

**But this has been at the expense of other qualities that the NLP community cares about.**

**size?**

**inference latency?**

**fairness?**

**energy efficiency?**

**training time?**

**ease of use?**

# How to frame the divergence between what's incentivized by leaderboards and what's valued by practitioners?

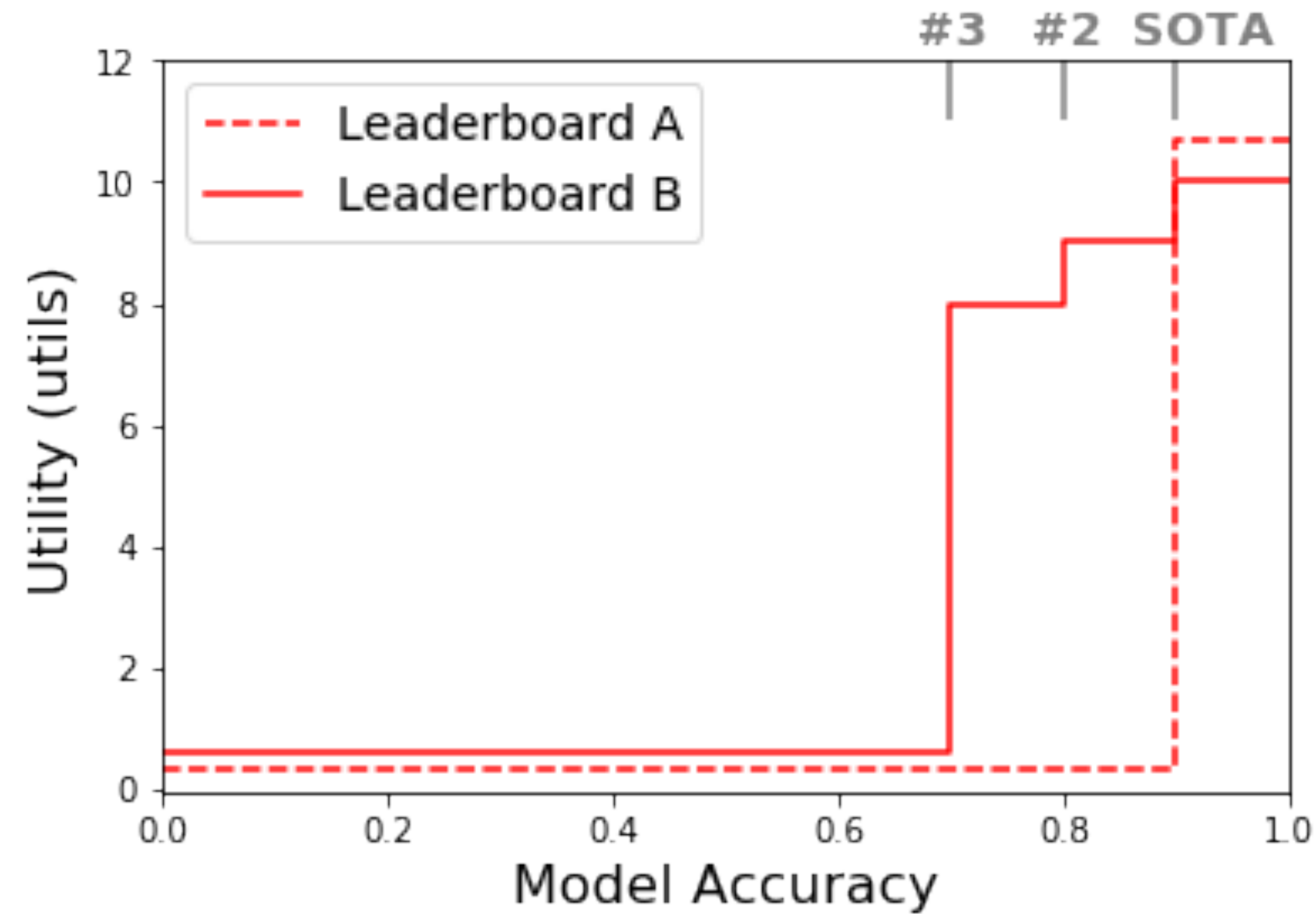
- Microeconomics!
- The *utility* of a good is the satisfaction that a *consumer* receives from it.
- Each consumer has a unique *utility function*.
- Both leaderboards and practitioners can be framed as consumers of models.

# How to frame the divergence between what's incentivized by leaderboards and what's valued by practitioners?

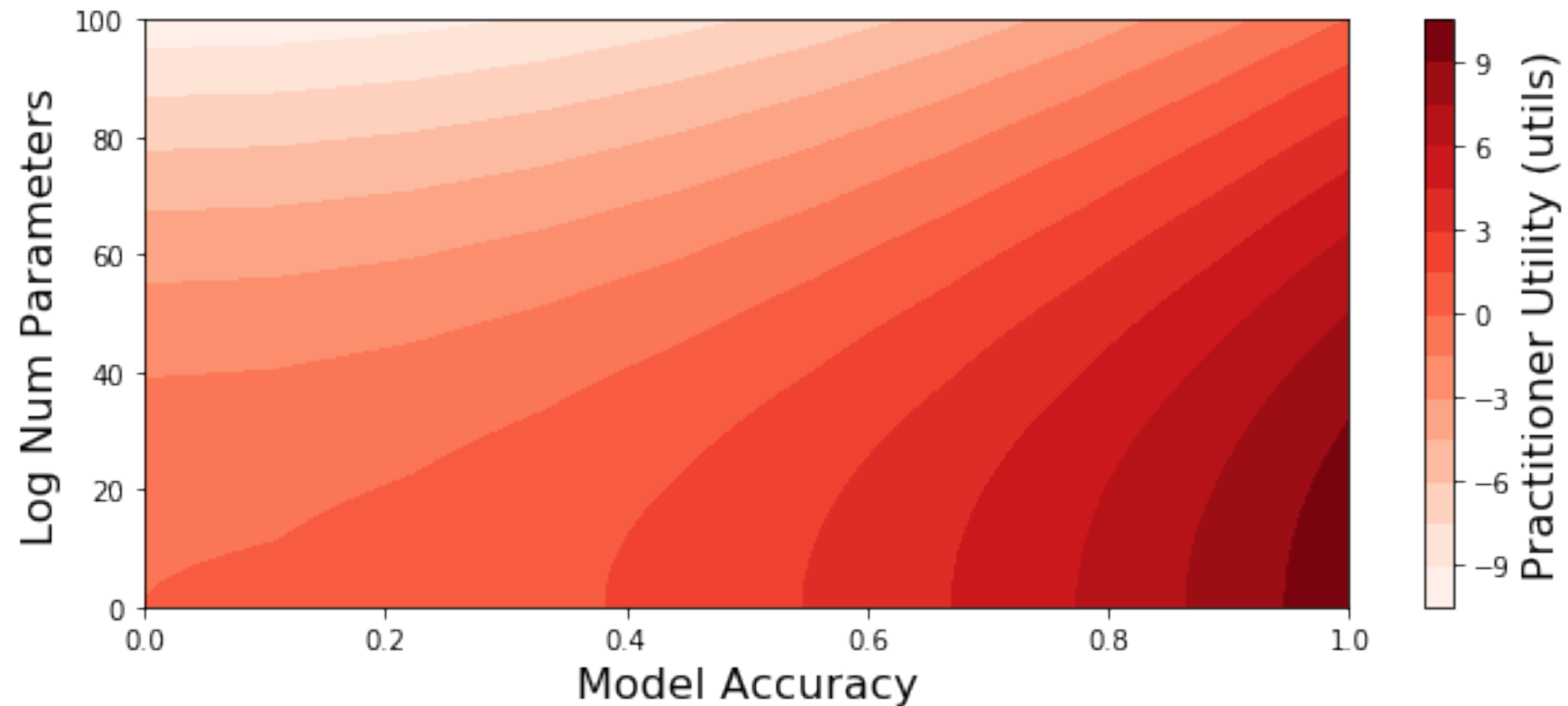
- Microeconomics!
- The *utility* of a good is the satisfaction that a *consumer* receives from it.
- Each consumer has a unique *utility function*.
- Both leaderboards and practitioners can be framed as consumers of models.
- **IDEA:** Compare leaderboards and practitioners using their utility functions.



**A leaderboard is a consumer whose preferences are perfectly revealed through its rankings: SOTA > #2 > ...**



**Practitioners derive utility from multiple properties of the model being consumed (e.g., accuracy, efficiency, latency).**



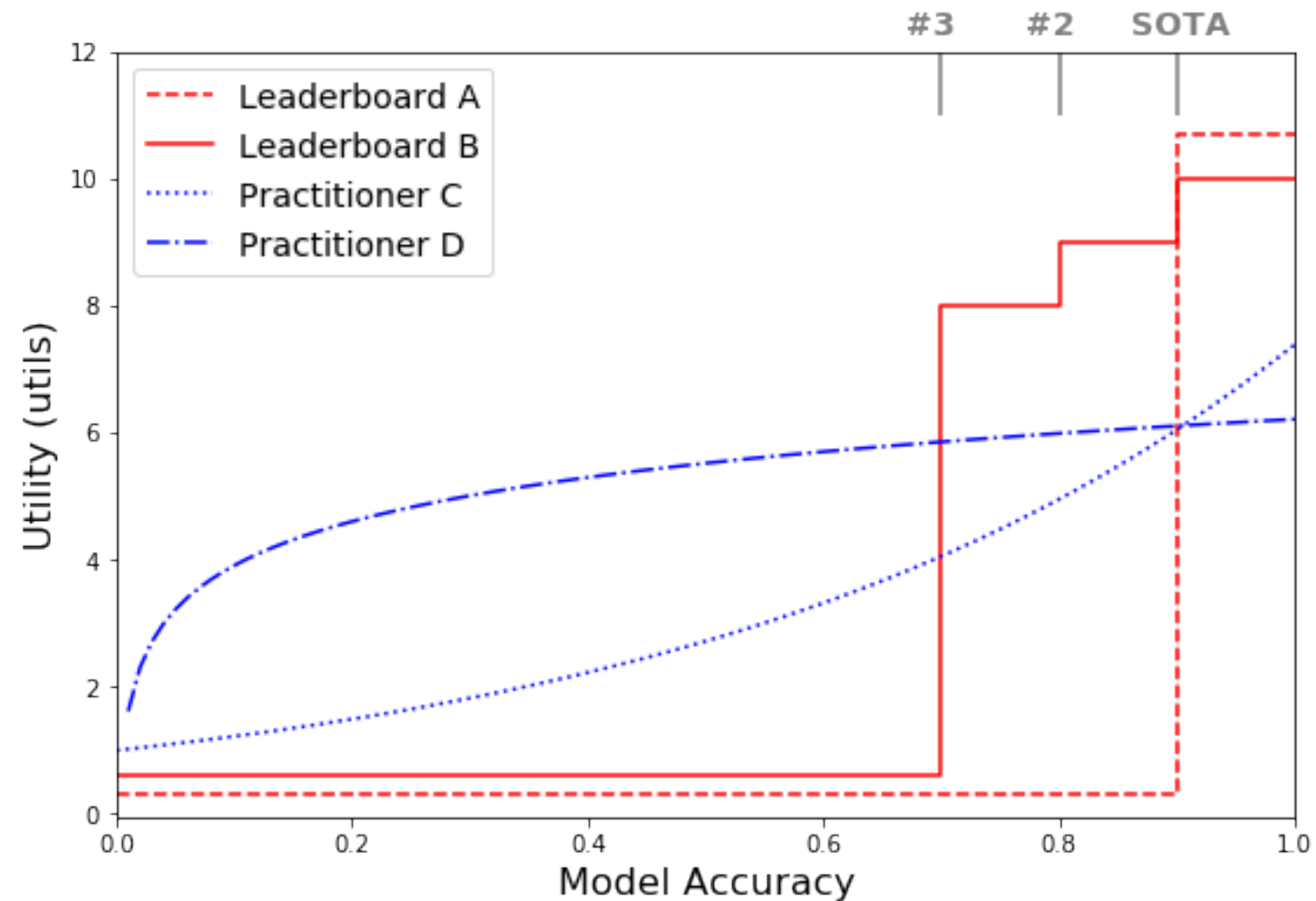


# **We can formally critique leaderboards by contrasting their utility functions with practitioners’.**

- We don’t know the exact shapes of utility functions, but we *do* know their properties: monotonicity, (in)sensitivity to certain attributes, etc.
- *Most* critiques apply to *most* leaderboards, but not all.

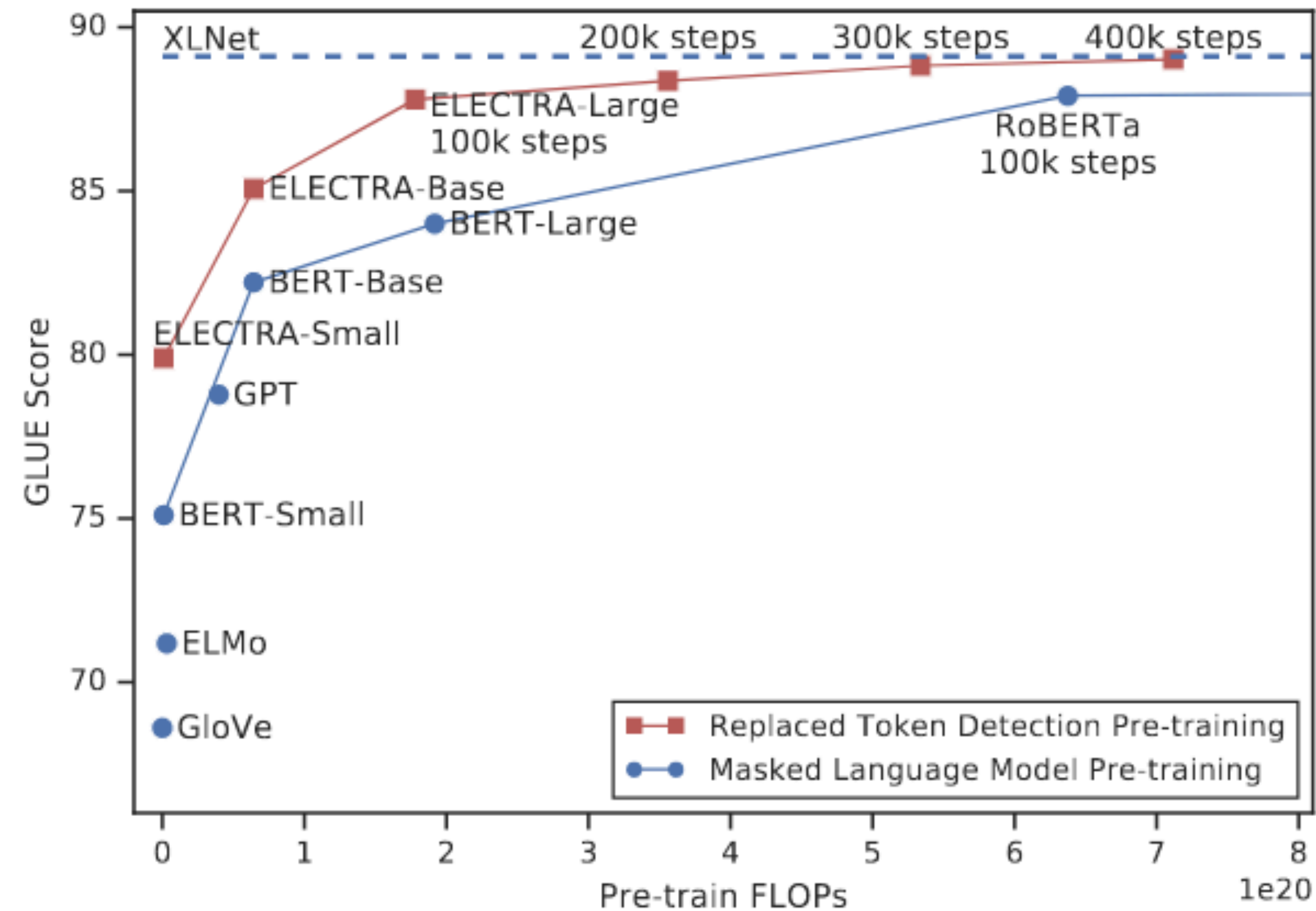
# Critique #1: Non-Smoothness of Utility

- Leaderboards only gain utility from increased accuracy when it improves rank.
- The utility of practitioners is smooth with respect to accuracy.



# Critique #1: Non-Smoothness of Utility

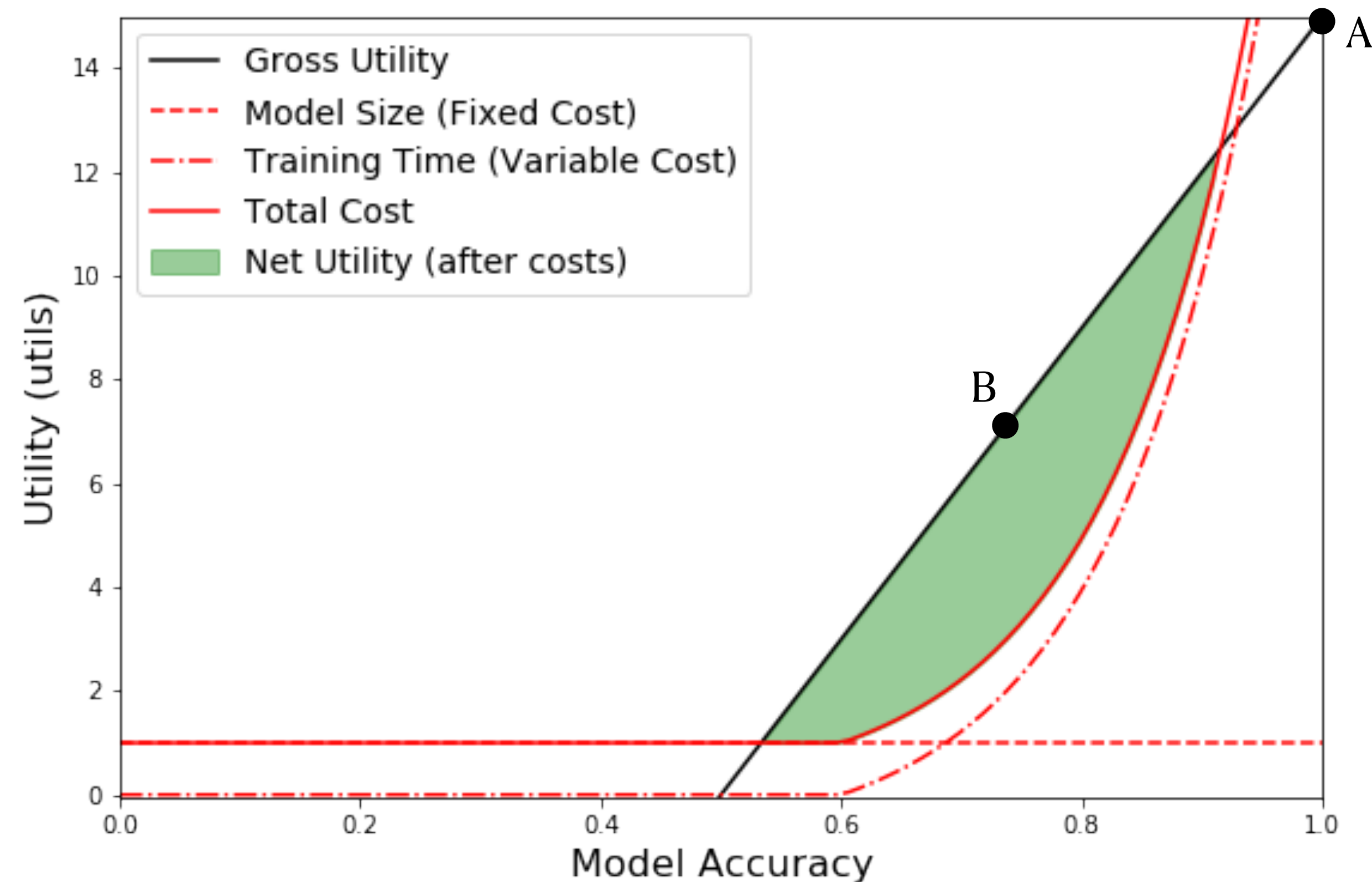
- Practitioners who are content with less-than-SOTA — e.g., for low latency or Green AI — have few options; those who want competitive-with-SOTA have many.



[ Clark et al., 2020 ]

# Critique #2: Cost-Ignorance

- Leaderboards rank by prediction *value*: accuracy, F1 score, exact match rate, etc.
- They ignore prediction *costs*: size, latency, energy efficiency, training time, etc.



# Critique #2: Cost-Ignorance

- Practitioners can't afford to be cost-ignorant (especially the poorly-resourced)!
- Cost-sensitive rankings would
  - incentivize the creation of more low-cost models like ELECTRA
  - allow practitioners to better estimate (net) model utility

# Critique #3: Robustness

- Over-fitting via resubmission is possible, even on private test sets.
- Most practitioners — but not most leaderboards — would gain utility from
  - robustness to adversarial examples
  - generalization to OOD data



# Critique #3: Robustness

- Over-fitting via resubmission is possible, even on private test sets.
- Most practitioners — but not most leaderboards — would gain utility from
  - robustness to adversarial examples
  - generalization to OOD data



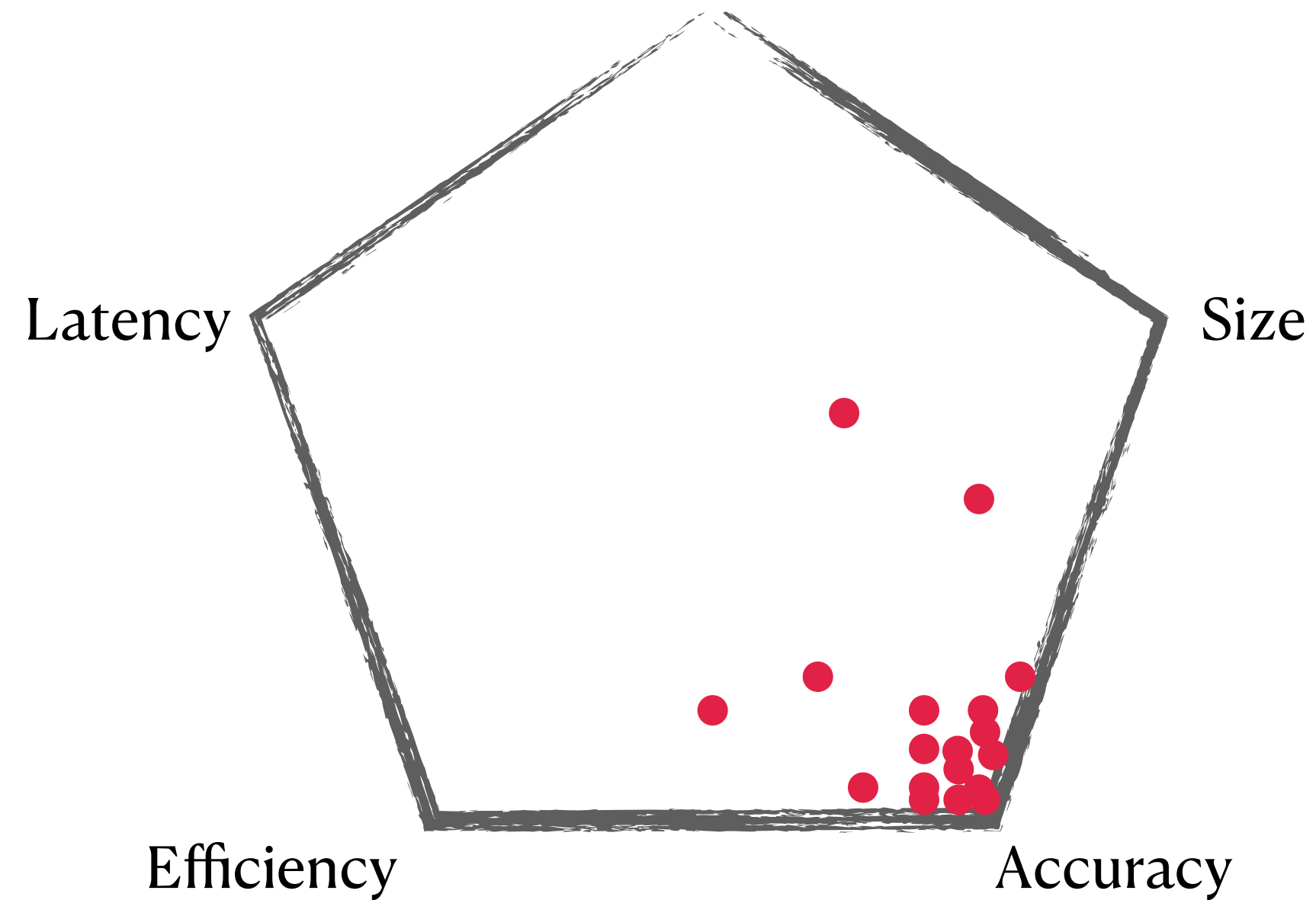
# The Future of Leaderboards: One for Every User

- Every practitioner has a unique utility function — no one-size-fits-all leaderboard.
- Leaderboards should demand transparency: require the reporting of metrics that are of practical concern (e.g., costs, adversarial performance, etc).
- Then allow users to dynamically re-rank models based on their priorities over these statistics (i.e., align leaderboard's utility with their own).

# Diverse Preferences, Diverse Models

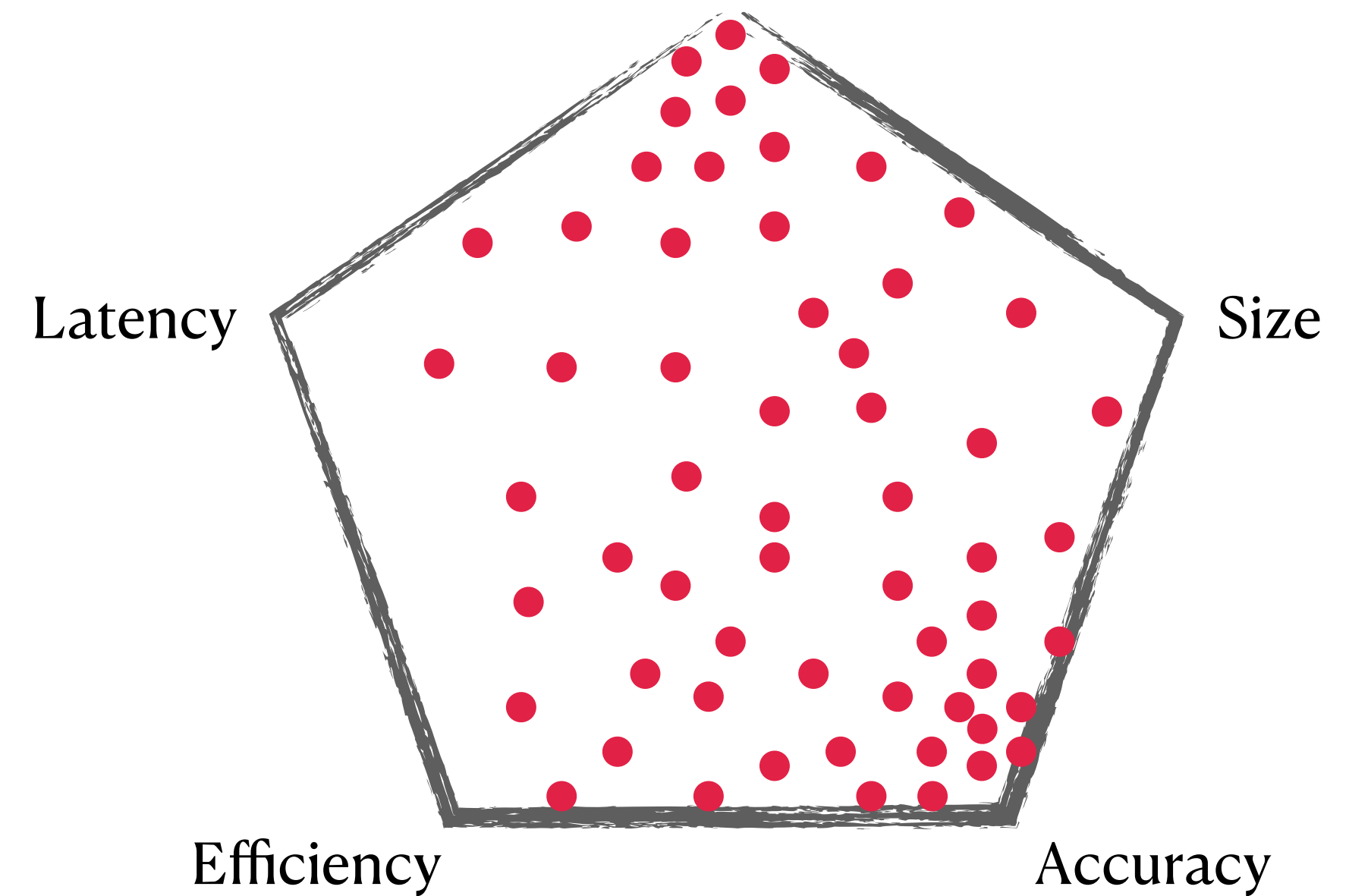
2020

Accuracy on Worst-off Group



A More Enlightened Age

Accuracy on Worst-off Group



Thank you!