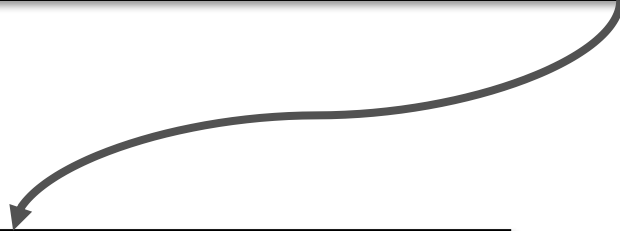# Commonsense benchmarks

Or how to measure that your model is actually doing some commonsense reasoning

# How do you know that a model is doing commonsense reasoning?
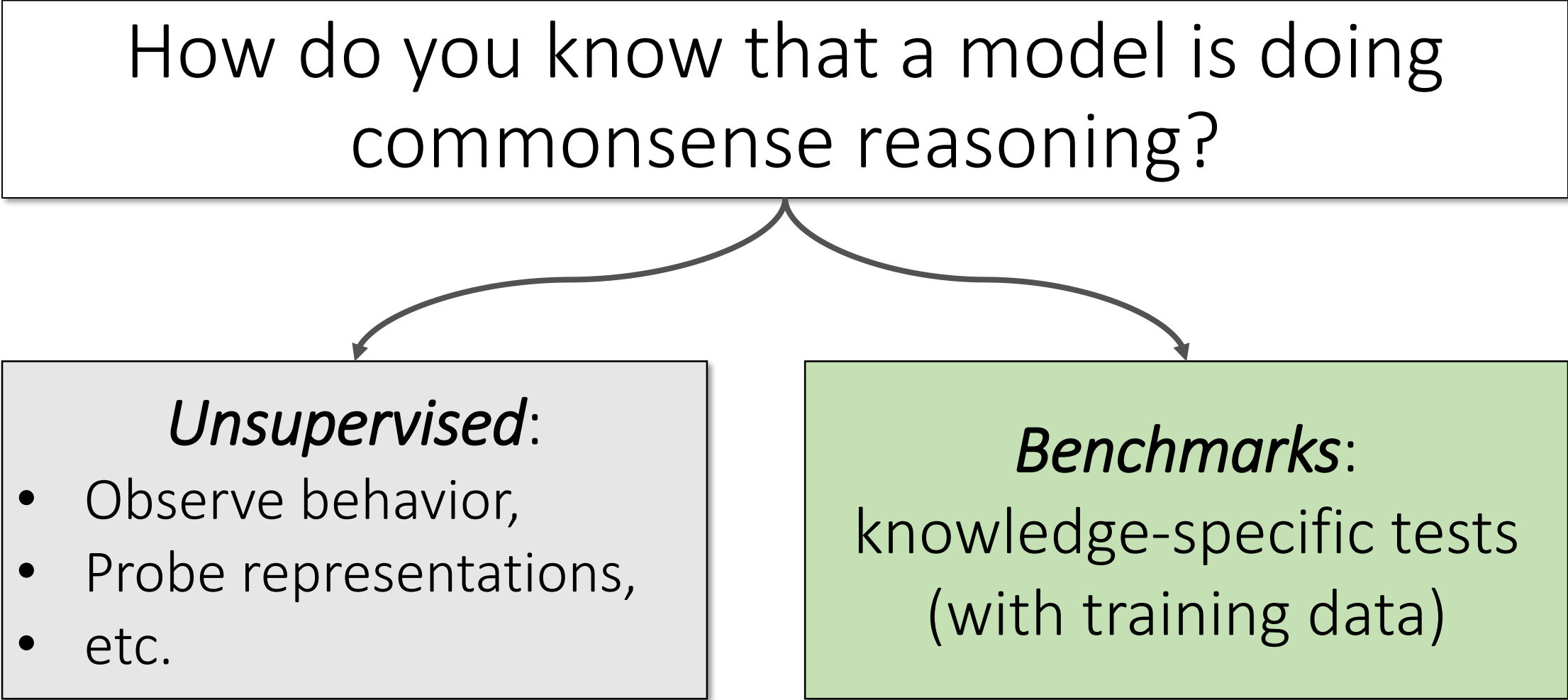
How do you know that a model is doing commonsense reasoning?

*Unsupervised*:
- Observe behavior,
- Probe representations,
- etc.

*Step 1*: Determine type of reasoning

# *Step 1*: Determine type of reasoning

Abductive reasoning

# *Step 1*: Determine type of reasoning

Abductive reasoning

Visual commonsense reasoning

Abductive NLI

Physical IQa

Social IQa

VCR

HELLA SWAG

https://leaderboard.allenai.org/

# *Step 1*: Determine type of reasoning

Abductive reasoning

Visual commonsense reasoning

Abductive NLI

Physical IQa

Social IQa

VCR

HELLA SWAG

# Reasoning about Social Situations

# Reasoning about Social Situations

Social IQa

Alex spilt food all over the floor and it made a huge mess.

What will Alex want to do next?

# Knowledge tested in SOCIAL IQA: ATOMIC



causes

drink too much

fall over

no intent

X needed to

X wanted to

stative

clumsy

careless

X is seen as

PersonX spills ___ all over the floor

has effect on X

gets dirty

slip on the spill

X will feel

embarrassed

upset

X will want

clean it up

get a broom

effects

# *Step 2*: Choosing a benchmark size

| | **Small scale** | **Large scale** |
|---|---|---|
| Creation | Expert-curated | Crowdsourced/automatic |
| Coverage | Limited coverage | Large coverage |
| Training | Dev/test only | Training/dev/test |
| Budget | Expert time costs | Crowdsourcing costs |

# *Step 2*: Choosing a benchmark size

|  | **Small scale** | **Large scale** |
|---|---|---|
| **Creation** | Expert-curated | Crowdsourced/automatic |
| **Coverage** | Limited coverage | Large coverage |
| **Training** | Dev/test only | Training/dev/test |
| **Budget** | Expert time costs | Crowdsourcing costs |

Winograd Schema Challenge (WSC),
Choice of Plausible Alternatives (COPA)

# Small commonsense benchmarks

**Winograd Schema Challenge (WSC)**
273 examples

**Choice of Plausible Alternatives (COPA)**
500 dev, 500 test

The city councilmen refused the demonstrators a permit because *they* **advocated** violence. Who is "*they*"?

(a) The city councilmen
(b) The demonstrators

The city councilmen refused the demonstrators a permit because *they* **feared** violence. Who is "*they*"?

(a) The city councilmen
(b) The demonstrators

# Small commonsense benchmarks

**Winograd Schema Challenge (WSC)**
273 examples

**Choice of Plausible Alternatives (COPA)**
500 dev, 500 test

I hung up the phone.
What was the **cause** of this?

(a) The caller said goodbye to me.
(b) The caller identified himself to me.

The toddler became cranky.
What happened as a **result**?

(a) Her mother put her down for a nap.
(b) Her mother fixed her hair into pigtails.

# *Step 2*: Choosing a QA benchmark size

|  | **Small scale** | **Large scale** |
|---|---|---|
| Creation | Expert-curated | Crowdsourced/automatic |
| Coverage | Limited coverage | Large coverage |
| Training | Dev/test only | Training/dev/test |
| Budget | Expert time costs | Crowdsourcing costs |

*Challenge*: do to collect positive/negative answers?

# Challenge of collecting unlikely answers

**Goal**: negative answers have to be *plausible but unlikely*

- Automatic matching?
    - Random negative sampling won't work, too topically different
    - "smart" negative sampling isn't effective either
- Need better solution… maybe we can ask crowd workers?

# Collecting answers from crowdworders

**Context and Question**

Alex spilt food all over the floor and it made a huge mess.

WHAT HAPPENS NEXT

What will Alex want to do next?

# Collecting answers from crowdworders

**Context and Question**

Alex spilt food all over the floor and it made a huge mess.

WHAT HAPPENS NEXT

What will Alex want to do next?

# Collecting answers from crowdworders

**Context and Question**

Alex spilt food all over the floor and it made a huge mess.

WHAT HAPPENS NEXT

What will Alex want to do next?

**Free Text Response**

Handwritten ✔ and ✘ Answers

✔ mop up
✔ give up and order take out
✘ leave the mess
✘ run around in the mess

# Collecting answers from crowdworders

**Context and Question**

Alex spilt food all over the floor and it made a huge mess.

WHAT HAPPENS NEXT

What will Alex want to do next?

**Free Text Response**

Handwritten ✓ and ✗ Answers

✓ mop up
✓ give up and order take out
✗ leave the mess
✗ run around in the mess

amazon mturk
Requester

Problem: handwritten unlikely answers are **too easy to detect**

# *Problem*: annotation artifacts

# *Problem*: annotation artifacts

- Models can exploit artifacts in handwritten incorrect answers
  - Exaggerations, off-topic, overly emotional, etc.
  - See Schwartz et al. 2017, Gururangan et al. 2018, Zellers et al. 2018, etc.
- Seemingly "super-human" performance by large pretrained LMs (BERT, GPT, etc.)

# *Problem*: annotation artifacts

- Models can exploit artifacts in handwritten incorrect answers
  - Exaggerations, off-topic, overly emotional, etc.
  - See Schwartz et al. 2017, Gururangan et al. 2018, Zellers et al. 2018, etc.
- Seemingly "super-human" performance by large pretrained LMs (BERT, GPT, etc.)

# *Problem*: annotation artifacts

- Models can exploit artifacts in handwritten incorrect answers
  - Exaggerations, off-topic, overly emotional, etc.
  - See Schwartz et al. 2017, Gururangan et al. 2018, Zellers et al. 2018, etc.
- Seemingly "super-human" performance by large pretrained LMs (BERT, GPT, etc.)

# *Problem*: annotation artifacts

- Models can exploit artifacts in handwritten incorrect answers
  - Exaggerations, off-topic, overly emotional, etc.
  - See Schwartz et al. 2017, Gururangan et al. 2018, Zellers et al. 2018, etc.
- Seemingly "super-human" performance by large pretrained LMs (BERT, GPT, etc.)

# *Problem*: annotation artifacts

- Models can exploit artifacts in handwritten incorrect answers
  - Exaggerations, off-topic, overly emotional, etc.
  - See Schwartz et al. 2017, Gururangan et al. 2018, Zellers et al. 2018, etc.
- Seemingly "super-human" performance by large pretrained LMs (BERT, GPT, etc.)

How to make unlikely answers **robust to annotation artifacts**?

How to make unlikely answers **robust to annotation artifacts**?

SOCIAL IQA:
switch questions in annotation

# Question-Switching Answers (SOCIAL IQA)

Original Question

Alex spilt food all over the floor and it made a huge mess.

WHAT HAPPENS NEXT

What will Alex want to do next?

✓ mop up
✓ give up and order take out
✗
✗

# Question-Switching Answers (SOCIAL IQA)

## Original Question

Alex spilt food all over the floor and it made a huge mess.

### WHAT HAPPENS NEXT

What will Alex want to do next?

✔ mop up
✔ give up and order take out
✘
✘

## Question-Switching Answer

### WHAT HAPPENED BEFORE

What did Alex need to do before this?

# Question-Switching Answers (SOCIAL IQA)

## Original Question

Alex spilt food all over the floor and it made a huge mess.

### WHAT HAPPENS NEXT

What will Alex want to do next?

✔ mop up
✔ give up and order take out
✘
✘

## Question-Switching Answer

### WHAT HAPPENED BEFORE

What did Alex need to do before this?

✔ have slippery hands
✔ get ready to eat

# Question-Switching Answers (SOCIAL IQA)

## Original Question

Alex spilt food all over the floor and it made a huge mess.

### WHAT HAPPENS NEXT
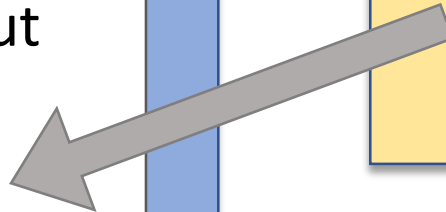
What will Alex want to do next?

✔ mop up
✔ give up and order take out

✘ have slippery hands
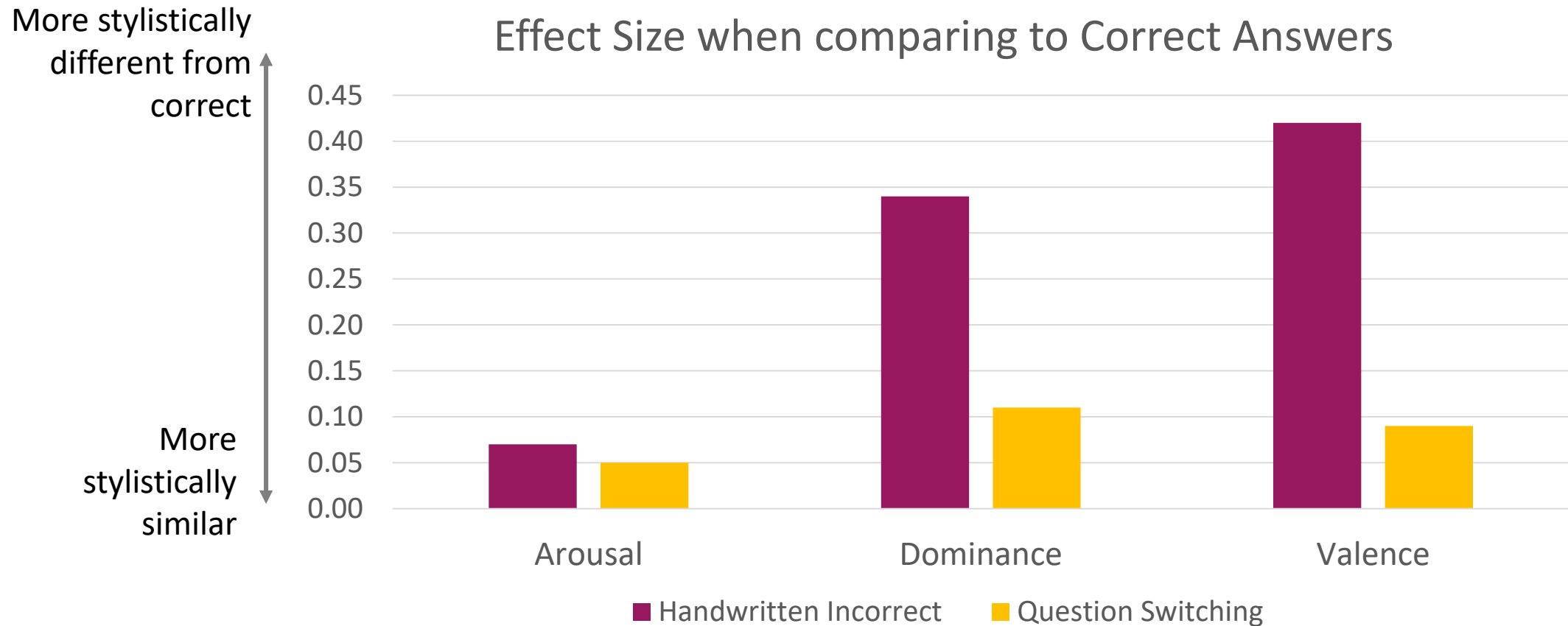✘ get ready to eat

## Question-Switching Answer

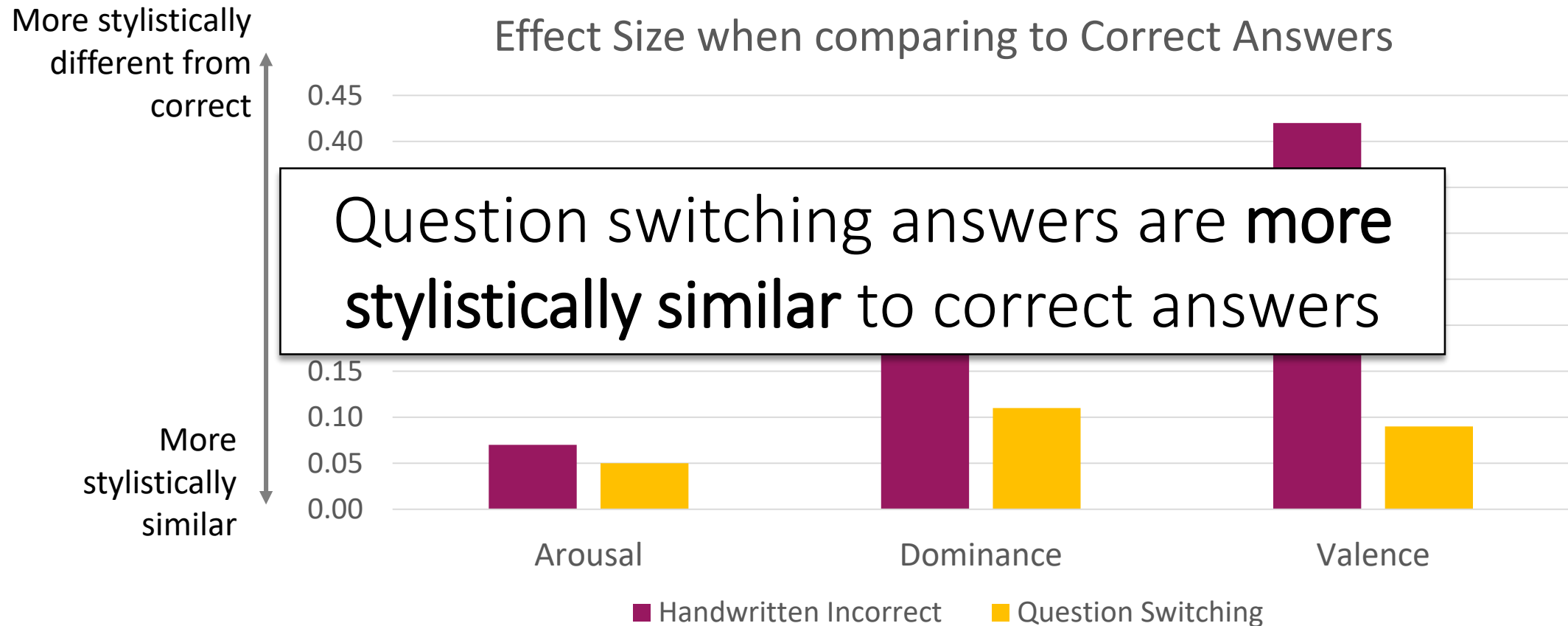### WHAT HAPPENED BEFORE

What did Alex need to do before this?

✔ have slippery hands
✔ get ready to eat

# Comparing incorrect/correct answers' styles

More stylistically different from correct

More stylistically similar

## Effect Size when comparing to Correct Answers



- Handwritten Incorrect
- Question Switching

# Comparing incorrect/correct answers' styles

More stylistically different from correct

Effect Size when comparing to Correct Answers

More stylistically similar

| | | | | |
|---|---|---|---|---|
| 0.45 | | | | |
| 0.40 | | | | |

Question switching answers are **more stylistically similar** to correct answers

| | | | | |
|---|---|---|---|---|
| 0.15 | | | | |
| 0.10 | | | | |
| 0.05 | | | | |
| 0.00 | | | | |

Arousal          Dominance          Valence

■ Handwritten Incorrect     ■ Question Switching

# Adversarial Filtering (lite)



Unfiltered examples

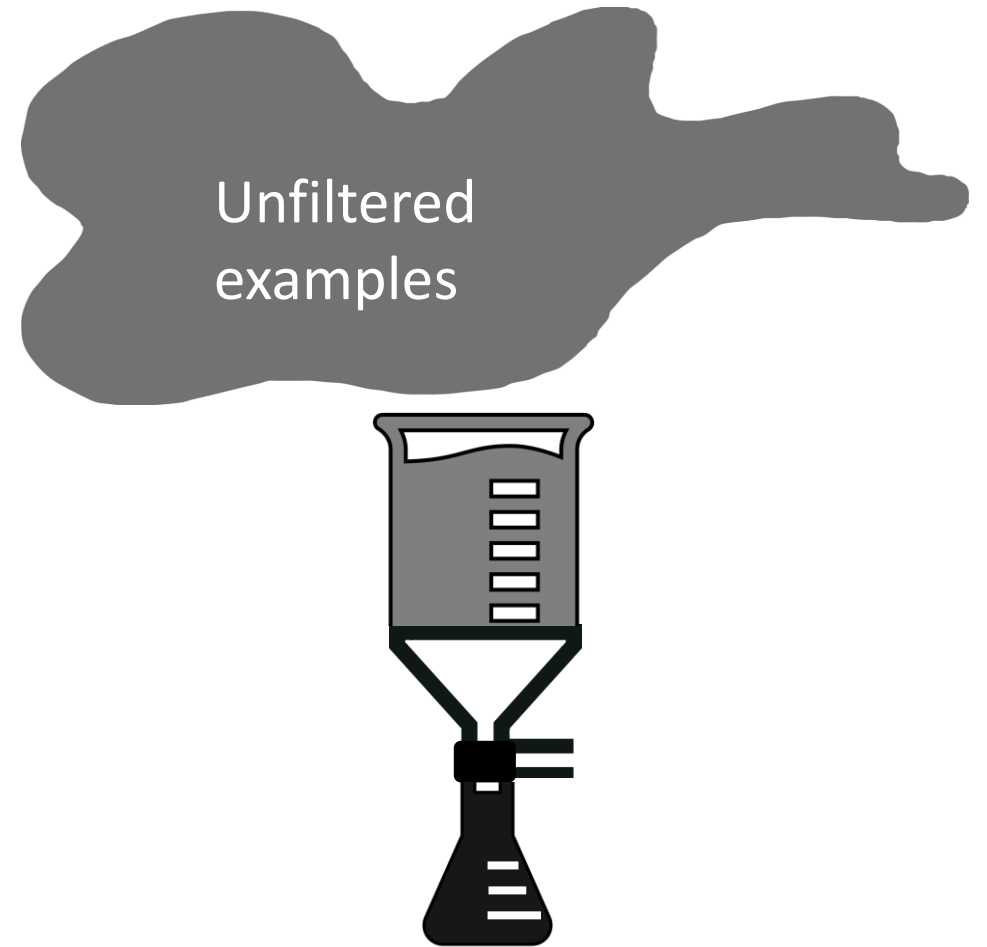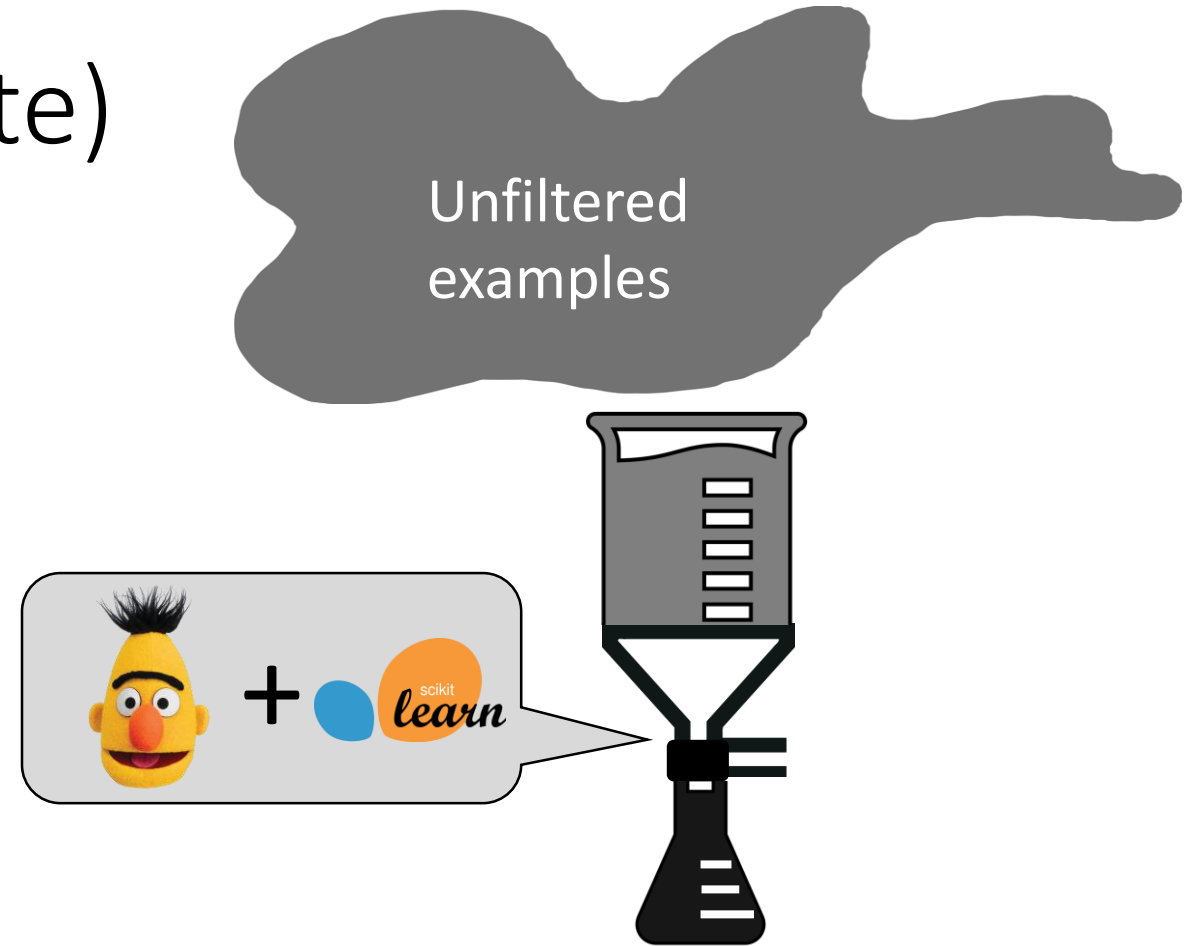*Goal*: remove examples with exploitable artifacts or spurious correlations

• Use pre-trained representations

• Iteratively remove data that's easiest to predict by a linear classifier (e.g., logistic)

• Robust examples remain

HellaSwag (Zellers et al., 2019)
AF-lite (Le Bras et al., 2019)

# Adversarial Filtering (lite)

*Goal*: remove examples with exploitable artifacts or spurious correlations
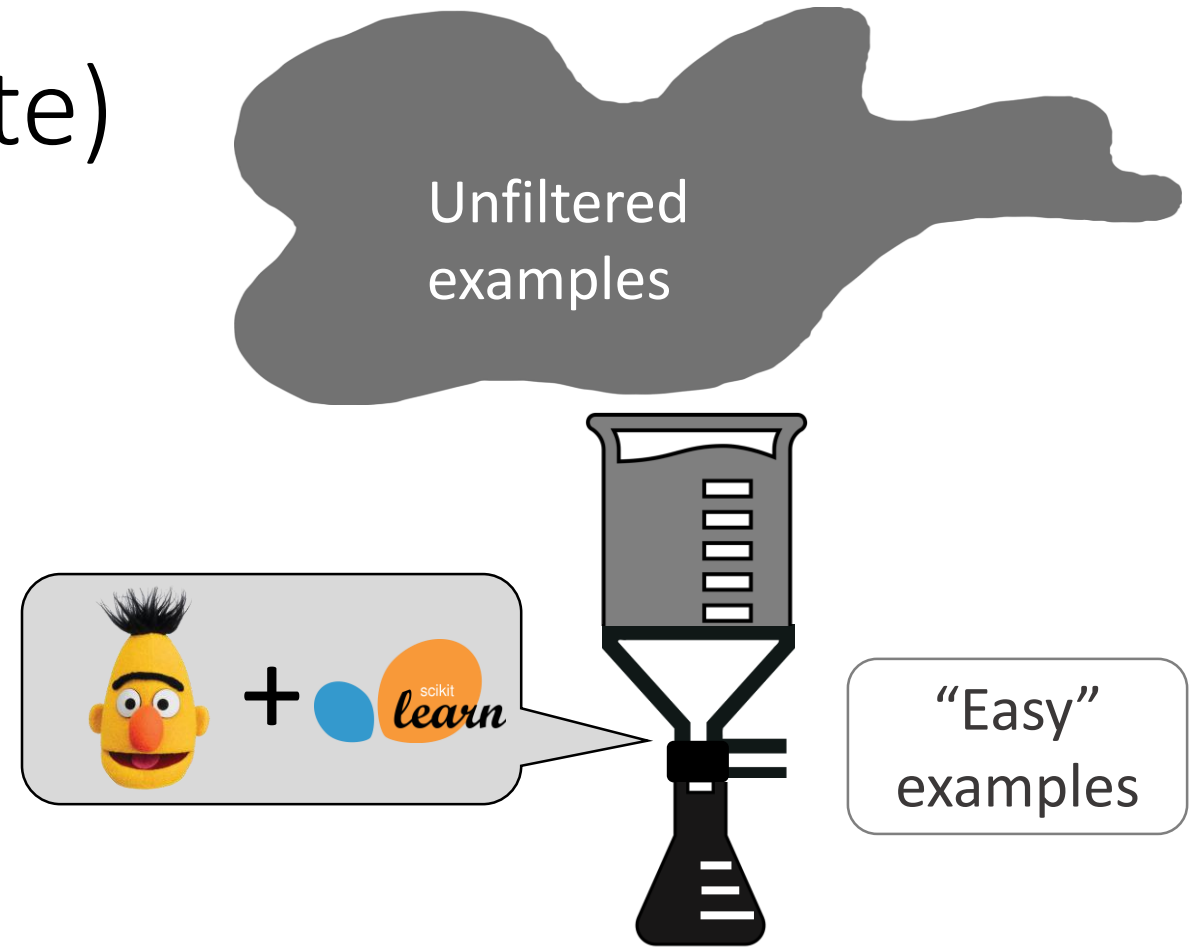
- Use pre-trained representations

- Iteratively remove data that's easiest to predict by a linear classifier (e.g., logistic)

- Robust examples remain

Unfiltered examples

HellaSwag (Zellers et al., 2019)
AF-lite (Le Bras et al., 2019)

# Adversarial Filtering (lite)

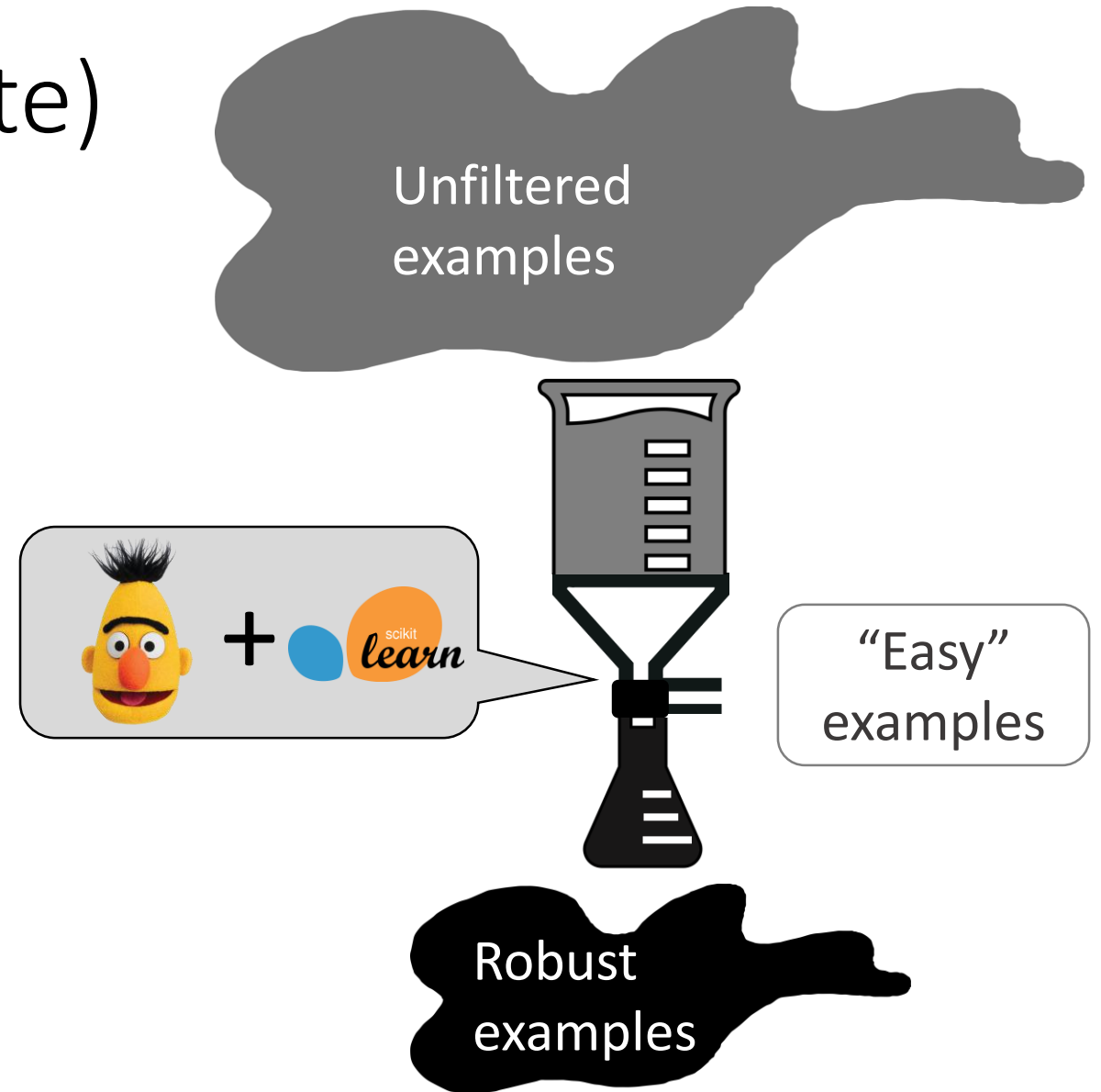*Goal*: remove examples with exploitable artifacts or spurious correlations

- Use pre-trained representations
- Iteratively remove data that's easiest to predict by a linear classifier (e.g., logistic)
- Robust examples remain

Unfiltered examples

HellaSwag (Zellers et al., 2019)
AF-lite (Le Bras et al., 2019)

# Adversarial Filtering (lite)

**Unfiltered examples**

*Goal*: remove examples with exploitable artifacts or spurious correlations

- Use pre-trained representations

- Iteratively remove data that's easiest to predict by a linear classifier (e.g., logistic)
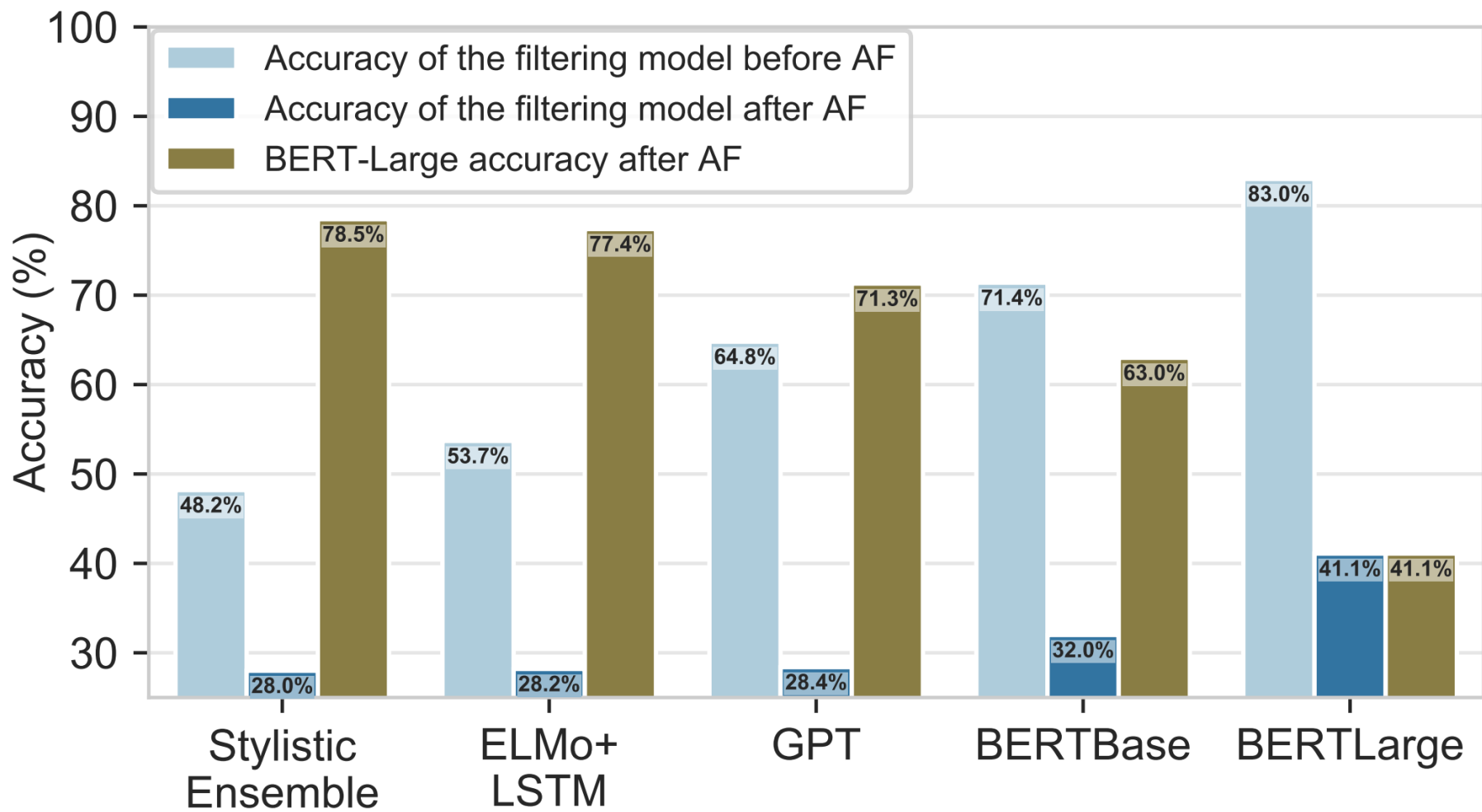
- Robust examples remain

"Easy" examples

HellaSwag (Zellers et al., 2019)
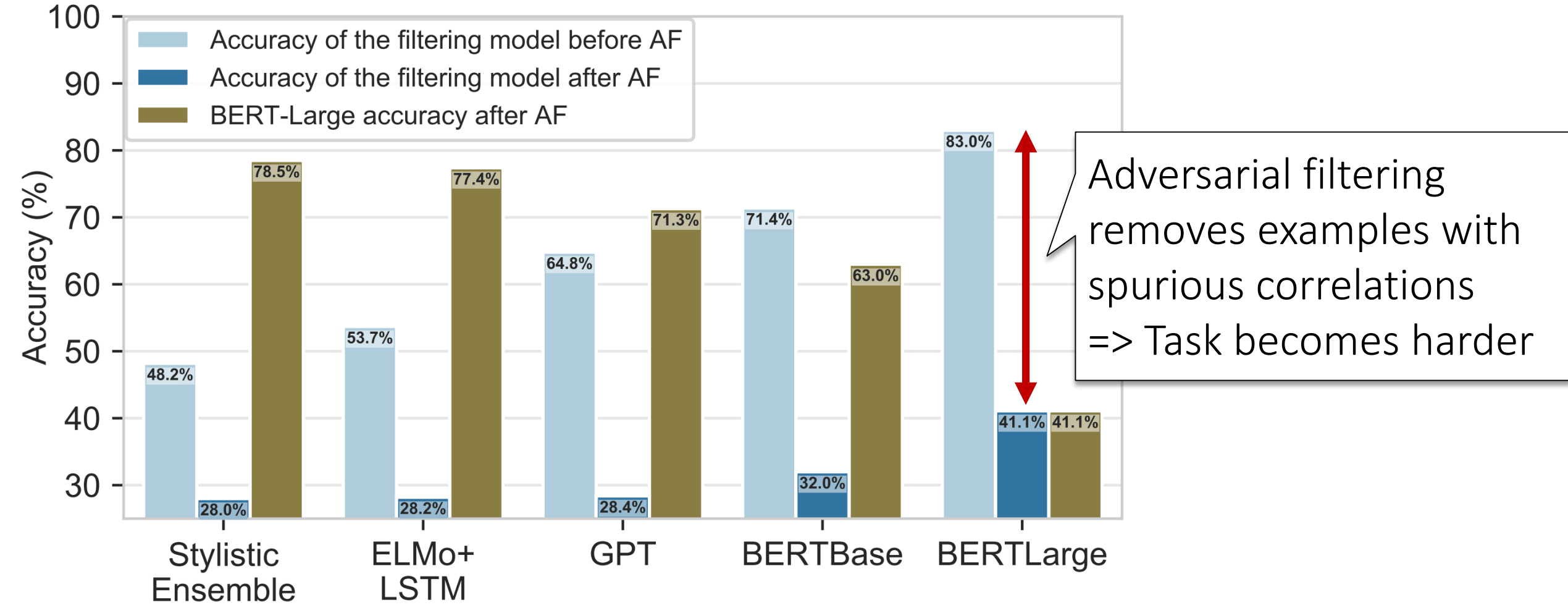AF-lite (Le Bras et al., 2019)

# Adversarial Filtering (lite)



*Goal*: remove examples with exploitable artifacts or spurious correlations

- Use pre-trained representations
- Iteratively remove data that's easiest to predict by a linear classifier (e.g., logistic)
- Robust examples remain

HellaSwag (Zellers et al., 2019)
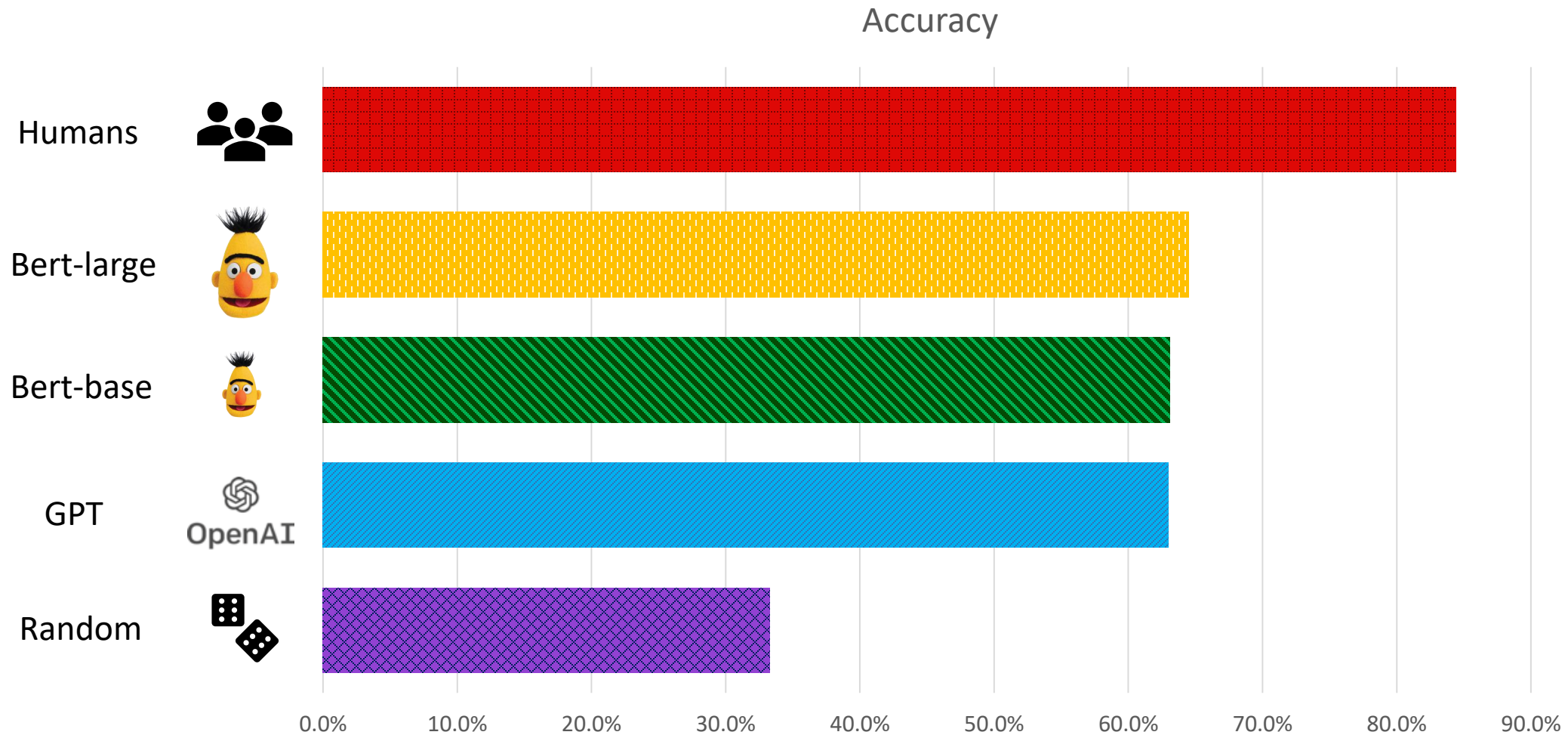AF-lite (Le Bras et al., 2019)

Performance of models on the WikiHow portion of HellaSwag (Zellers et al., 2019)
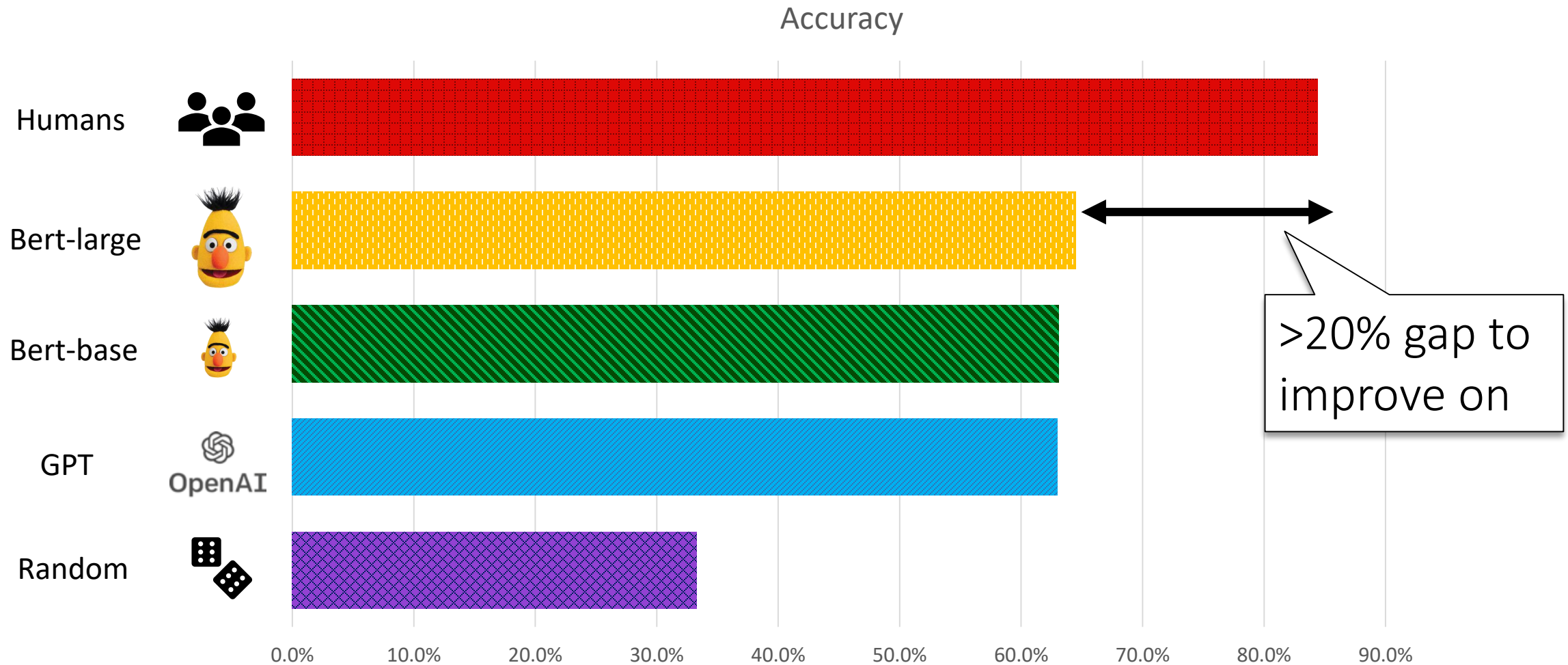with different AF settings and different training models

Performance of models on the WikiHow portion of HellaSwag (Zellers et al., 2019)
with different AF settings and different training models

# Model performance on SOCIAL IQA

# Model performance on SOCIAL IQA

Accuracy



>20% gap to improve on

# Challenging SOCIAL IQA examples for BERT-large

Although Aubrey was older and stronger, they lost to Alex in arm wrestling.

Remy gave Skylar, the concierge, her account so that she could check into the hotel.

**How would Alex feel as a result?**

**What will Remy want to do next?**

ashamed — how **Aubrey** would feel, not Alex

✓ boastful

they need to practice more

lose her credit card

arrive at a hotel — what Remy did **before**

✓ get the key from Skylar

Need more robust, person-centric reasoning

Need better notion of causes vs. effects

# Commonsense benchmarks

Social commonsense

Naïve Psychology

ROC story

Social IQa

WSC

COPA

VCR

WinoGrande
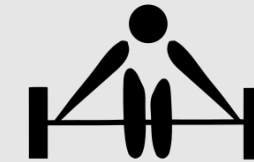
Physical commonsense

Physical IQa

HellaSwag

SWAG

Abductive NLI

CommonsenseQA

JHU Ordinal Commonsense

MCTaco

Temporal commonsense

ReCORD

CosmosQA

MultiRC

Commonsense reading comprehension