

Natural Language Processing



Existing Large Language Models

Kevin Lin – UC Berkeley
March 15, 2023

Existing Large Language Models

Announcements

- HW4 – finetuning LLMs: release today
- HW5 – prompting LLMs: released early April
- Panel Topics Overview
- Today:
 - BERT
 - T5
 - GPT3

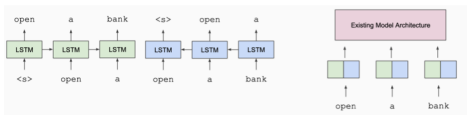
BERT

- Bidirectional Encoder Representations from Transformers (Devlin et al., 2018)
- A working general recipe: pretrain and finetune
- SOTA across token + sentence level tasks
- Deep bidirectional encoder-only model



Previous Work

- ELMO: Deep Contextualized Word Embeddings (Peters et al., 2018)
 - Left-to-right, right-to-left unidirectional LSTMs
 - Plug in as features
 - Single sentences



Previous Work

- GPT: Improving Language Understanding by Generative Pre-training (Radford et al., 2018)
- Finetuning
- Left-to-right
- BooksCorpus (512 length)





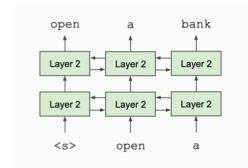
BERT

- Objectives: Masked Language Modeling + Next Sentence Prediction
- Deep encoder-only transformer
- Learn from bidirectional context
 - go to the bank to make a deposit
 - on the river bank



Masked Language Modeling

- Problem: words "see" themselves



Masked Language Modeling

- Solution: masking
- Cloze-style task (Taylor, 1953)
- Denosing-autoencoders
- Select 15%
- No [MASK] during fine-tuning
 - 80% Replace with [MASK]

the man went to the store to buy a gallon of milk



Next Sentence Prediction

- Learn relationships between sentences
- Predict whether sentence A follows sentence B
- Later shown to be not very helpful

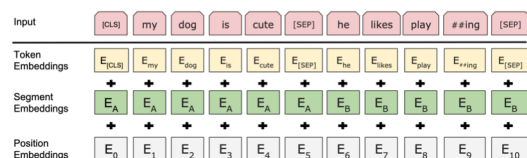


Sentence A = The man went to the store.
Sentence B = He bought a gallon of milk.
Label = IsNextSentence

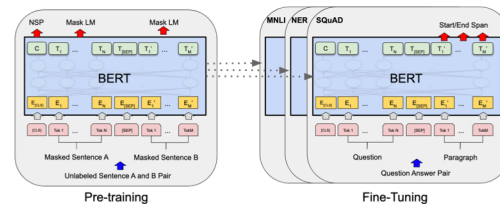
Sentence A = The man went to the store.
Sentence B = Penguins are flightless.
Label = NotNextSentence



Input Representation



BERT





Sentence-Level Tasks

- Linear layer on top of [CLS] token

GLUE Results

| System | MNLI-(m/mm) | QQP | QNLI | SST-2 | CoLA | STS-B | MRPC | RTI | Average |
|-----------------------|-------------|------|------|-------|------|-------|------|------|---------|
| Pre-OpenAI SOTA | 80.6/80.1 | 66.1 | 82.3 | 93.2 | 35.0 | 81.0 | 86.0 | 61.7 | 74.0 |
| BiLSTM+ELMo+Attn | 76.4/76.1 | 64.8 | 79.9 | 90.4 | 36.0 | 73.3 | 84.9 | 56.8 | 71.0 |
| OpenAI GPT | 82.1/81.4 | 70.3 | 88.1 | 91.3 | 45.4 | 80.0 | 82.3 | 56.0 | 75.2 |
| BERT _{base} | 84.6/83.4 | 71.2 | 90.1 | 93.3 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| BERT _{large} | 86.7/85.9 | 72.1 | 91.1 | 94.9 | 60.5 | 86.5 | 89.3 | 70.1 | 81.9 |

MultINLI
 Premise: Hills and mountains are especially
 sanctified in Jainism.
 Hypothesis: Jainism hates nature.
 Label: Contradiction

CoLA
 Sentence: The wagon rumbled down the road.
 Label: Acceptable
 Sentence: The car honked down the road.
 Label: Unacceptable



Token-Level Tasks

- Extractive QA, NER
- Linear layer on top of token representations

What was another term used for the oil crisis?
 Ground Truth Answers: **first oil shock** shock shock first oil
 shock shock
 Prediction: shock

The 1973 **oil crisis** began in October 1973 when the members of the Organization of Arab Petroleum Exporting Countries (OAPEC), consisting of the Arab members of OPEC plus Egypt and Syria) proclaimed an **oil** embargo. By the end of the embargo in March 1974, the price of **oil** had risen from US\$3 per barrel to nearly \$12 globally; US prices were significantly higher. The embargo caused an **oil** **crisis**, or "shock", with many short- and long-term effects on global politics and the global economy. It was later called the "**first oil shock**", followed by the 1979 **oil crisis**, termed the "second **oil** shock."



Training BERT

- Data: Wikipedia (2.5B words) + BooksCorpus (800M words)
- 1M steps
- Batch size: 131,072 words
 - (1024 sequences * 128 length) or (256 sequence * 512 length)
- BERT-large
 - 24 layers, 1024 hidden size, 16 attention heads, 340M parameters
- BERT-base
 - 12 layers, 768 hidden size, 12 attention heads, 110M parameters



BERT Aftermath

- Explosion of variations:
 - RoBERTa: Train longer, remove NSP
 - ALBERT: share weights
 - SpanBERT: mask out contiguous spans
 - Electra: learn from all tokens
- Efficiency:
 - DistillBERT, qBERT, ...
- BERT for X
 - SciBERT: scientific documents
 - ClinicalBERT: clinical documents, ...



BERT Aftermath

- "BERTology"
 - What does an LLM encode about syntax, semantics, knowledge, etc.?
- Generation from BERT
 - Mask-Predict: Parallel Decoding of Conditional Masked Language Models (Ghazvininejad et al., 2019)
 - BERT as Markov Random Field LM (Wang et al., 2019)



T5

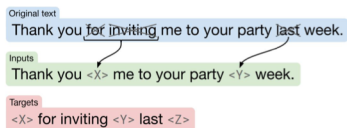
T5

- T5: Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer
- Objectives, architectures, datasets, transfer
- Unified format: text in, text out
- Discriminative and generative tasks



T5 setup

- Start with a basic setup, and get first order effects: objectives, architectures, datasets, transfer

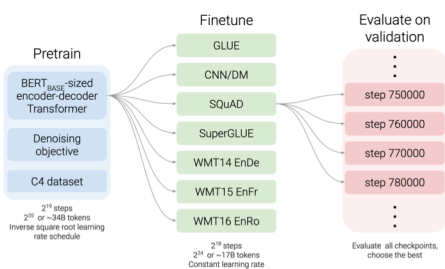


T5 setup

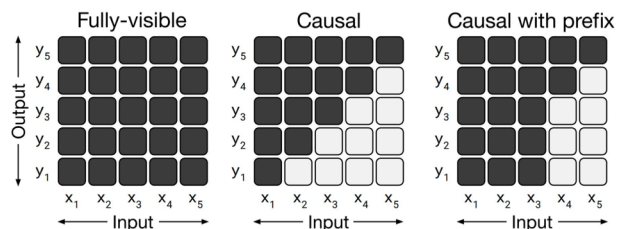
- Encoder-decoder
- Each size of BERT-base
- Relative positional embedding
- C4: Colossal Clean Crawled Corpus
 - Filter out javascript, non-English, List-of-Dirty-Naughty-Obscene-and-Otherwise-Bad-Words
 - 750B
 - Note: tokenizer handles French, Romanian, German



T5 setup

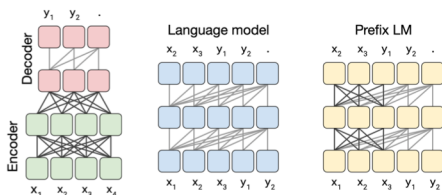


T5 architecture



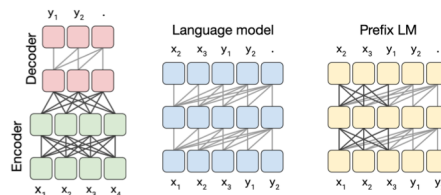
T5 Architecture

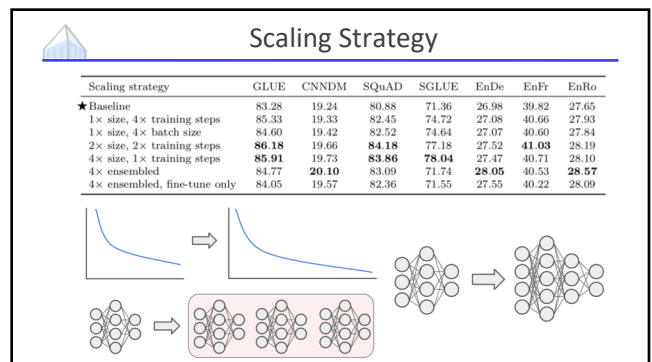
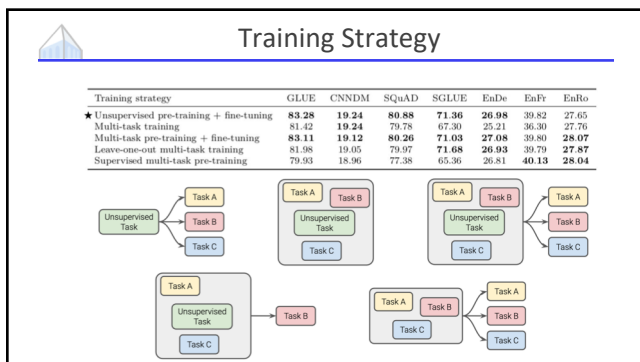
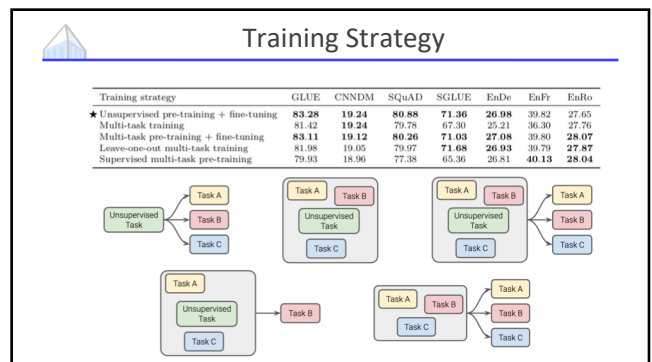
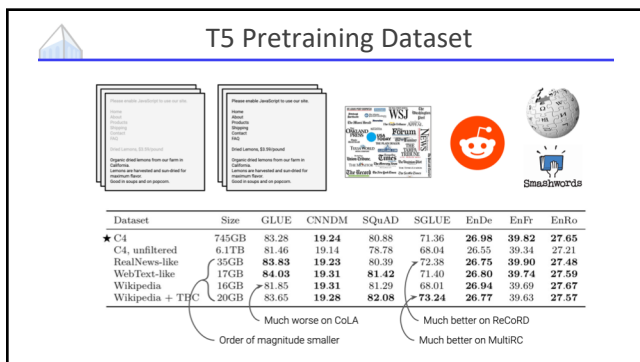
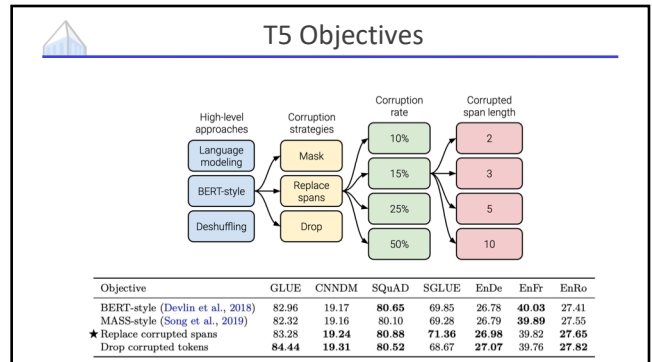
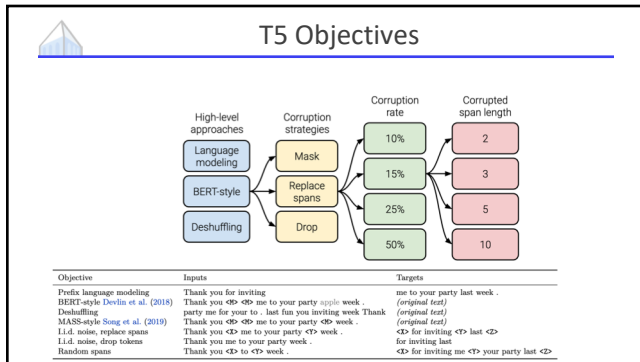
| Architecture | Params | Cost | GLUE | CNN/DM | SQuAD | SGlUE | EnDe | EnFr | EnRo |
|-------------------|--------|------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| ★ Encoder-decoder | 2P | M | 83.28 | 19.24 | 80.88 | 71.36 | 26.98 | 39.82 | 27.65 |
| Enc-dec, shared | P | M | 82.81 | 18.78 | 80.63 | 70.73 | 26.72 | 39.03 | 27.46 |
| Enc-dec, 6 layers | P | M/2 | 80.88 | 18.97 | 77.59 | 68.42 | 26.38 | 38.40 | 26.95 |
| Language model | P | M | 74.70 | 17.93 | 61.14 | 55.02 | 25.09 | 35.28 | 25.86 |
| Prefix LM | P | M | 81.82 | 18.61 | 78.94 | 68.11 | 26.43 | 37.98 | 27.39 |



T5 Architecture

| Architecture | Params | Cost | GLUE | CNN/DM | SQuAD | SGlUE | EnDe | EnFr | EnRo |
|-------------------|--------|------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| ★ Encoder-decoder | 2P | M | 83.28 | 19.24 | 80.88 | 71.36 | 26.98 | 39.82 | 27.65 |
| Enc-dec, shared | P | M | 82.81 | 18.78 | 80.63 | 70.73 | 26.72 | 39.03 | 27.46 |
| Enc-dec, 6 layers | P | M/2 | 80.88 | 18.97 | 77.59 | 68.42 | 26.38 | 38.40 | 26.95 |
| Language model | P | M | 74.70 | 17.93 | 61.14 | 55.02 | 25.09 | 35.28 | 25.86 |
| Prefix LM | P | M | 81.82 | 18.61 | 78.94 | 68.11 | 26.43 | 37.98 | 27.39 |







T5: Putting It Together

- Encoder-Decoder
- Span Replacement
- C4
- Multi-task pretraining
- Large models, trained longer



T5: Putting It Together

| Model | Parameters | # layers | d_{model} | d_{ff} | d_{kv} | # heads |
|-------|------------|----------|--------------------|-----------------|-----------------|---------|
| Small | 60M | 6 | 512 | 2048 | 64 | 8 |
| Base | 220M | 12 | 768 | 3072 | 64 | 12 |
| Large | 770M | 24 | 1024 | 4096 | 64 | 16 |
| 3B | 3B | 24 | 1024 | 16384 | 128 | 32 |
| 11B | 11B | 24 | 1024 | 65536 | 128 | 128 |

| Model | GLUE Average | CoLA Matthew's | SST-2 Accuracy | MRPC F1 | MRPC Accuracy | STS-B Pearson | STS-B Spearman |
|---------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| Previous best | 89.4 ^a | 69.2 ^b | 97.1 ^a | 93.6 ^b | 91.5 ^b | 92.7 ^b | 92.3 ^b |
| T5-Small | 77.4 | 41.0 | 91.8 | 89.7 | 86.6 | 85.6 | 85.0 |
| T5-Base | 82.7 | 51.1 | 95.2 | 90.7 | 87.5 | 89.4 | 88.6 |
| T5-Large | 86.4 | 61.2 | 96.3 | 92.4 | 89.9 | 89.9 | 89.2 |
| T5-3B | 88.5 | 67.1 | 97.4 | 92.5 | 90.0 | 90.6 | 89.8 |
| T5-11B | 90.3 | 71.6 | 97.5 | 92.8 | 90.4 | 93.1 | 92.8 |



Finetuning Limitations

- Requires large supervised dataset
- Spurious correlations in supervised finetuning dataset
- Poor sample efficiency vs. humans



GPT3

- Language Models are Few-Shot Learners (Brown et al., 2020)
- Decoder-only model



GPT3

| Model Name | n_{params} | n_{layers} | d_{model} | n_{heads} | d_{head} | Batch Size | Learning Rate |
|-----------------------|---------------------|---------------------|--------------------|--------------------|-------------------|------------|----------------------|
| GPT-3 Small | 125M | 12 | 768 | 12 | 64 | 0.5M | 6.0×10^{-4} |
| GPT-3 Medium | 350M | 24 | 1024 | 16 | 64 | 0.5M | 3.0×10^{-4} |
| GPT-3 Large | 760M | 24 | 1536 | 16 | 96 | 0.5M | 2.5×10^{-4} |
| GPT-3 XL | 1.3B | 24 | 2048 | 24 | 128 | 1M | 2.0×10^{-4} |
| GPT-3 2.7B | 2.7B | 32 | 2560 | 32 | 80 | 1M | 1.6×10^{-4} |
| GPT-3 6.7B | 6.7B | 32 | 4096 | 32 | 128 | 2M | 1.2×10^{-4} |
| GPT-3 13B | 13.0B | 40 | 5140 | 40 | 128 | 2M | 1.0×10^{-4} |
| GPT-3 175B or "GPT-3" | 175.0B | 96 | 12288 | 96 | 128 | 3.2M | 0.6×10^{-4} |

| Dataset | Quantity (tokens) | Weight in training mix | Epochs elapsed when training for 300B tokens |
|-------------------------|-------------------|------------------------|--|
| Common Crawl (filtered) | 410 billion | 60% | 0.44 |
| WebText2 | 19 billion | 22% | 2.9 |
| Books1 | 12 billion | 8% | 1.9 |
| Books2 | 55 billion | 8% | 0.43 |
| Wikipedia | 3 billion | 3% | 3.4 |



GPT3 Training

- Larger models -> larger batch sizes, smaller learning rate
- Model parallelism: across layers
- Adam optimizer
- Gradient clipping: 1
- Linear warmup learning rate, cosine decay
- Weight decay 0.1



In-Context Learning

Zero-Shot

1 Translate English to French: task description
2 cheese => prompt

One-Shot

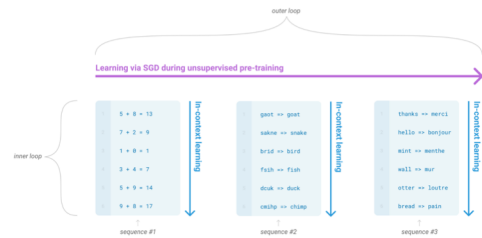
1 Translate English to French: task description
2 sea otter => loutre de mer example
3 cheese => prompt

Few-Shot

Translate English to French: task description
sea otter => loutre de mer examples
peppermint => menthe poivrée
plush girafe => girafe peluche
cheese => prompt



In-Context Learning



Evaluation Setup

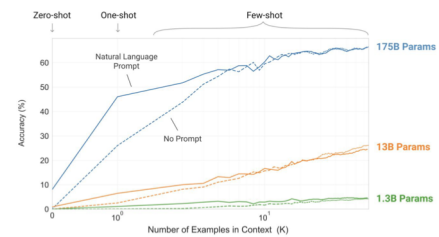
- Few-shot: sample uniformly from train set
- Normalize by unconditional probability

$$\frac{P(\text{completion}|\text{context})}{P(\text{completion}|\text{answer_context})}$$

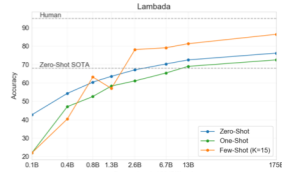
- More semantic names i.e. "True" instead of 1
- Generation tasks: beam search



Size vs. In-Context Learning



Language Modeling + Cloze Tasks

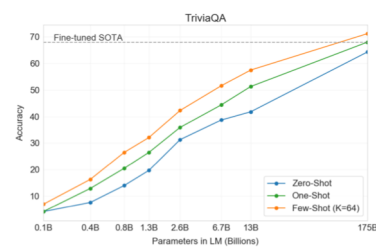


| Setting | LAMBADA (acc) | LAMBADA (ppl) | StoryCloze (acc) | HellaSwag (acc) |
|-----------------|-------------------|-------------------|-------------------|-------------------|
| SOTA | 68.0 ^a | 8.63 ^b | 91.8 ^c | 85.6 ^d |
| GPT-3 Zero-Shot | 76.2 | 3.00 | 83.2 | 78.9 |
| GPT-3 One-Shot | 72.5 | 3.35 | 84.7 | 78.1 |
| GPT-3 Few-Shot | 86.4 | 1.92 | 87.7 | 79.3 |

| Setting | PTB |
|------------------|-------------------|
| SOTA (Zero-Shot) | 35.8 ^a |
| GPT-3 Zero-Shot | 20.5 |



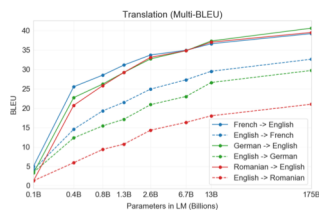
Closed-book QA



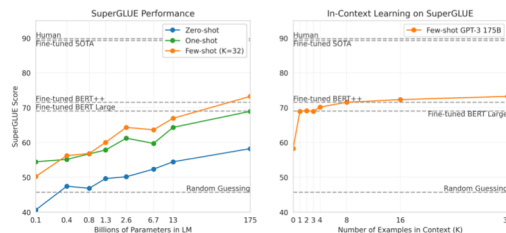


Machine Translation

- Match or near SOTA in high-resource settings
- Better going into English



Sentence-Level Tasks



Reading Comprehension

Passage

That year, his Untitled (1981), a painting of a haloed, black-headed man with a bright red skeletal body, depicted amid the artist's signature scrawls, was sold by Robert Lehman for \$16.3 million, well above its \$12 million high estimate.

In 1517, the seventeen-year-old King sailed to Castile. There, his Flemish court... In May 1518, Charles traveled to Barcelona in Aragon.

In 1970, to commemorate the 100th anniversary of the founding of Baldwin City, Baker University professor and playwright Don Mueller and Phyllis E. Braun, Business Manager, produced a musical play entitled The Ballad Of Black Jack to tell the story of the events that led up to the battle.

Question

How many more dollars was the Untitled (1981) painting sold for than the \$12 million dollar estimation?

Where did Charles travel to first, Castile or Barcelona?

Who was the University professor that helped produce The Ballad Of Black Jack, Ivan Boyd or Don Mueller?

Answer

4300000

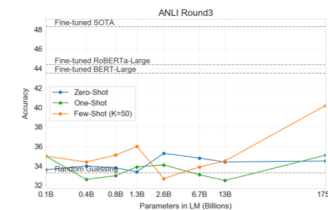
Castile

Don Mueller

| Setting | CoQA | DROP | QuAC | SQuADv2 | RACE-h | RACE-m |
|-----------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| Fine-tuned SOTA | 90.7 ^a | 89.1 ^b | 74.4 ^c | 93.0 ^d | 90.0 ^e | 93.1 ^e |
| GPT-3 Zero-Shot | 81.5 | 23.6 | 41.5 | 59.5 | 45.5 | 58.4 |
| GPT-3 One-Shot | 84.0 | 34.3 | 43.3 | 65.4 | 45.9 | 57.4 |
| GPT-3 Few-Shot | 85.0 | 36.5 | 44.3 | 69.8 | 46.8 | 58.1 |



ANLI



| Context | Hypothesis | Rationale | Gold/Pred. (Valid.) | Tags |
|---|--|---|---------------------|-----------------------------------|
| Edvard Scholz (4 January 1891 in Düsseldorf - 6 January 1966 in Zürich) was a prominent German industrialist. He was one of the first to warn the Allies and tell the world of the Holocaust and systematic exterminations of Jews in Nazi Germany occupied Europe. | Edvard Scholz is the only person to warn the Allies above the atrocities committed by the Nazis. | The context states that he is not the only person to warn the Allies above the atrocities committed by the Nazis. | C/N (CC) | Thicky, Prag., Numerical, Ordinal |



Machine or Human?

Title: Star's Tux Promise Draws Megyn Kelly's Sarcasm

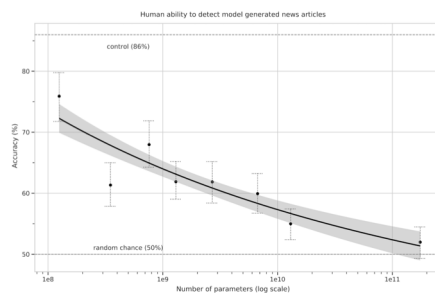
Subtitle: Joaquin Phoenix pledged to not change for each awards event

Article: A year ago, Joaquin Phoenix made headlines when he appeared on the red carpet at the Golden Globes wearing a tuxedo with a paper bag over his head that read, "I am a shape-shifter. I can't change the world. I can only change myself." It was a promise to not change to fit into the Hollywood mold: "I think that's a really special thing, to not change yourself. I think it's a really special thing to say, 'This is what's inside of me, I'm proud of it, and I'm not going to be ashamed because of the way that someone else thinks I should be.'" Now, it's the Oscars, and Phoenix is at it again. But this time, his publicist is saying he'll be wearing a tux no matter what.

Megyn Kelly was not impressed, and she let him have it on The Tonight Show. "You know, I feel like, I feel like you could have worn the tux," she says. "But you're saying you're a shape-shifter. I don't know if you can change your tux, but you can change your mind. You can change your mind. You can change your mind." Phoenix says he did, but it didn't stick. "I was like, 'Okay, I'm going to wear a tuxedo to this thing.'" And then I thought, "I don't want to wear a tuxedo to this thing." Kelly goes on to encourage him to change his mind again, but Phoenix says it's too late: "I'm committed to wearing this."



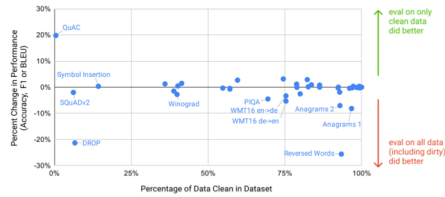
Machine or Human?





Test Set Contamination

- Memorization or generalization?
- Bug in removing test data from training data



Bias, Fairness, Representation

Table 6.1: Most Biased Descriptive Words in 175B Model

| Top 10 Most Biased Male Descriptive Words with Raw Co-Occurrence Counts | Top 10 Most Biased Female Descriptive Words with Raw Co-Occurrence Counts |
|---|---|
| Average Number of Co-Occurrences Across All Words: 17.5 | Average Number of Co-Occurrences Across All Words: 23.9 |
| Large (16) | Optimistic (12) |
| Moody (15) | Bubbly (12) |
| Lazy (14) | Naughty (12) |
| Fantastic (13) | Easy-going (12) |
| Eccentric (13) | Perse (10) |
| Protest (10) | Tight (10) |
| Jolly (10) | Pregnant (10) |
| Stable (9) | Gorgeous (28) |
| Personable (22) | Sucked (8) |
| Survive (7) | Beautiful (158) |



Bias, Fairness, Representation

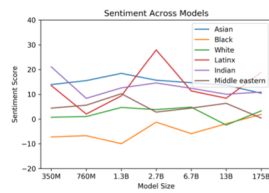


Figure 6.1: Racial Sentiment Across Models



Open Questions

- Scaling
- Evaluation
- Misuse, Risks
- Grounding
- Controllability
- Multilingual