**GloVe**  **ELMo**  **GPT**  **BERT**  **RoBERTa**  · · · · · ·
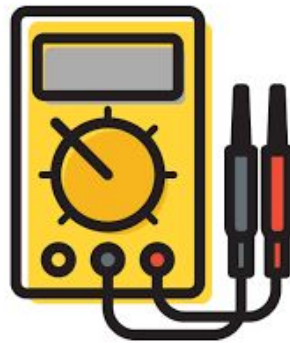
## Leaderboard

SQuAD2.0 tests the ability of a system to not only answer reading comprehension questions, but also abstain when presented with a question that cannot be answered based on the provided paragraph. How will your system compare to humans on this task?

| Rank | Model | EM | F1 |
|------|-------|-----|-----|
| | Human Performance<br>*Stanford University*<br>(Rajpurkar & Jia et al. '18) | 86.831 | 89.452 |
| 1<br>Mar 20, 2019 | BERT + DAE + AoA (ensemble)<br>*Joint Laboratory of HIT and iFLYTEK Research* | **87.147** | **89.474** |
| 2<br>Mar 15, 2019 | BERT + ConvLSTM + MTL + Verifier (ensemble)<br>*Layer 6 AI* | 86.730 | 89.286 |
| 3<br>Mar 05, 2019 | BERT + N-Gram Masking + Synthetic Self-Training (ensemble)<br>*Google AI Language*<br>https://github.com/google-research/bert | 86.673 | 89.147 |

# Why?

# Probes

Linguistic, factual, commonsense, etc.

- Well motivated tests that encode and measure correspondence to human knowledge/intelligence (e.g. linguistic annotation, factual query, etc.)
- Better test score
  - → better learned ability
  - → better explain the **"why?"**

# Current Probes
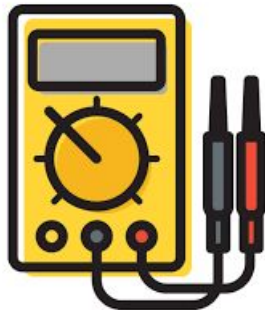


GloVe  ELMo  **GPT**  BERT  RoBERTa  ......

Linguistic, factual, commonsense, etc.

😄

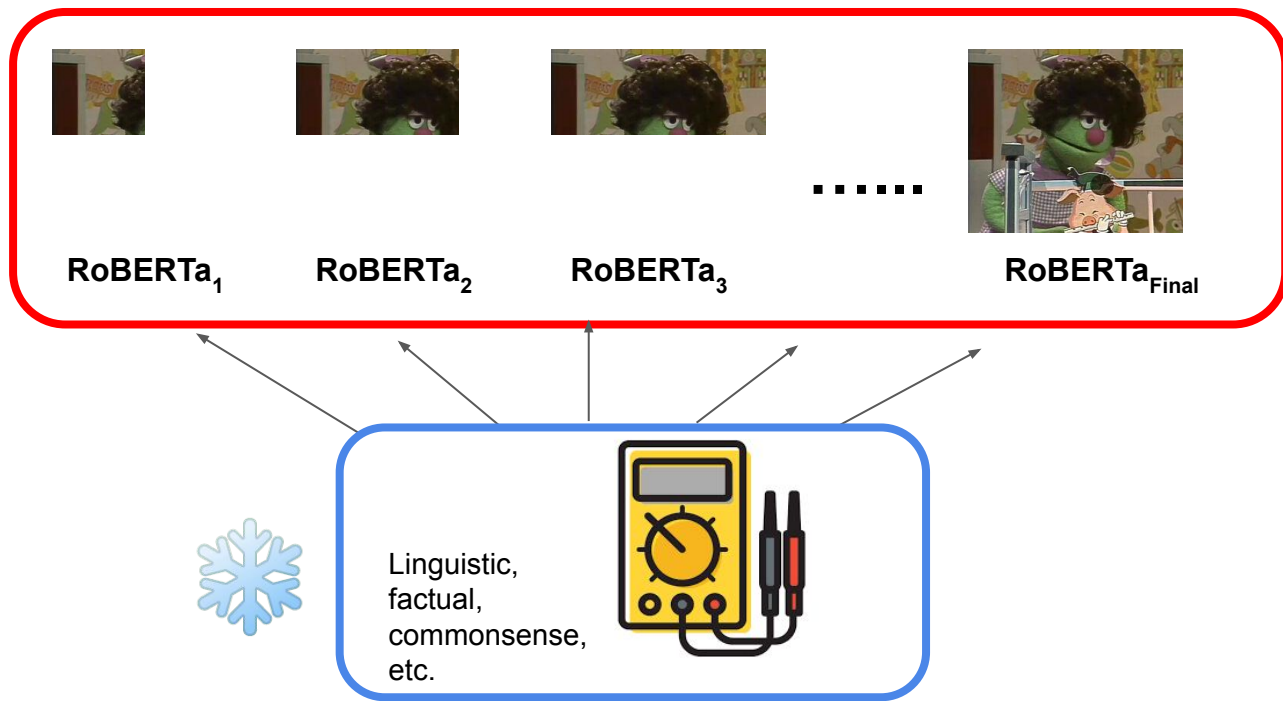- Compare models "shoulder-to-shoulder" by an interpretable metric

🙁

- Would model perform better in the next optimization update?

- How does the model learn?

# Probe Across Time

# Probing Across Time

😄

- Understand the underlying learning curriculum

- Longer observation increases our confidence in concluding how well model learns tested knowledge

# Choice of probes

Diverse **probe formulation** **probed knowledge**

- Linguistic
  (Contextual) embedding $\rightarrow$ $\mathbf{f}_{linear}$ $\rightarrow$ annotation

- Factual
  Score(✅) > Score(❌), e.g. Score=perplexity of a sentence

- Commonsense
  Pr(✅ | slot-filling query) > Pr(❌ | slot-filling query)

- Reasoning
  A ferry and a floatplane are both a type of *[MASK]*.

  ✅ vehicle ❌ airplane ❌ boat

# Baseline*

- **Random Guess**:  1 / (# labels)

- **{Random, GloVe} Vector** $\rightarrow \mathbf{f}_{\text{linear}} \rightarrow$ annotation

- **Original RoBERTa** probes the officially released checkpoint of RoBERTa base to see if our checkpoints are pretrained properly and can achieve reasonable performance

**\*** applicable to different types of probes

Legend: Random Guess; GloVe + Linear Clf.; Our Checkpoints; Learning Progress-90%; Learning Progress-97%; Random Vector + Linear Clf.; Original RoBERTa$_{BASE}$; exp. moving average curve; Learning Progress-95%

Rows (left labels): LKT (**Linguistic**); BLiMP (**Linguistic**); LAMA (**Factual& Commonsense**); CAT (**Commonsense**); oLMpics (**Reasoning**)

Panel titles:
Row LKT: Ave. Performance; POS Tagging; Syntactic Chunking; NER; Syntactic Arc Pred.; Syntactic Arc Class.
Row BLiMP: Ave. Performance; Irregular Forms; Determiner-Noun Agree.; Subject-Verb Agree.; Filler-Gap; Island Effects
Row LAMA: Ave. Performance; Google RE; SQuAD; T-REx; ConceptNet
Row CAT: Ave. Performance; Conjunction Acceptability; Winograd; Sense Making; SWAG; Argument
Row oLMpics: Ave. Performance; Taxonomy Conjunction; Antonym Negation; Objects Comparison; Always Never; Multi-Hop Composition

X-axis: Number of Pretraining Steps

**Learning curriculum**

**TL;DR**

Linguistic (knowledge)

∨

Factual ≅ Commonsense

Reasoning

# In fact,
# we didn't mention...

RoBERTa₁      RoBERTa₂      RoBERTa₃      RoBERTa_Fina

Linguistic,
factual,
commonsense
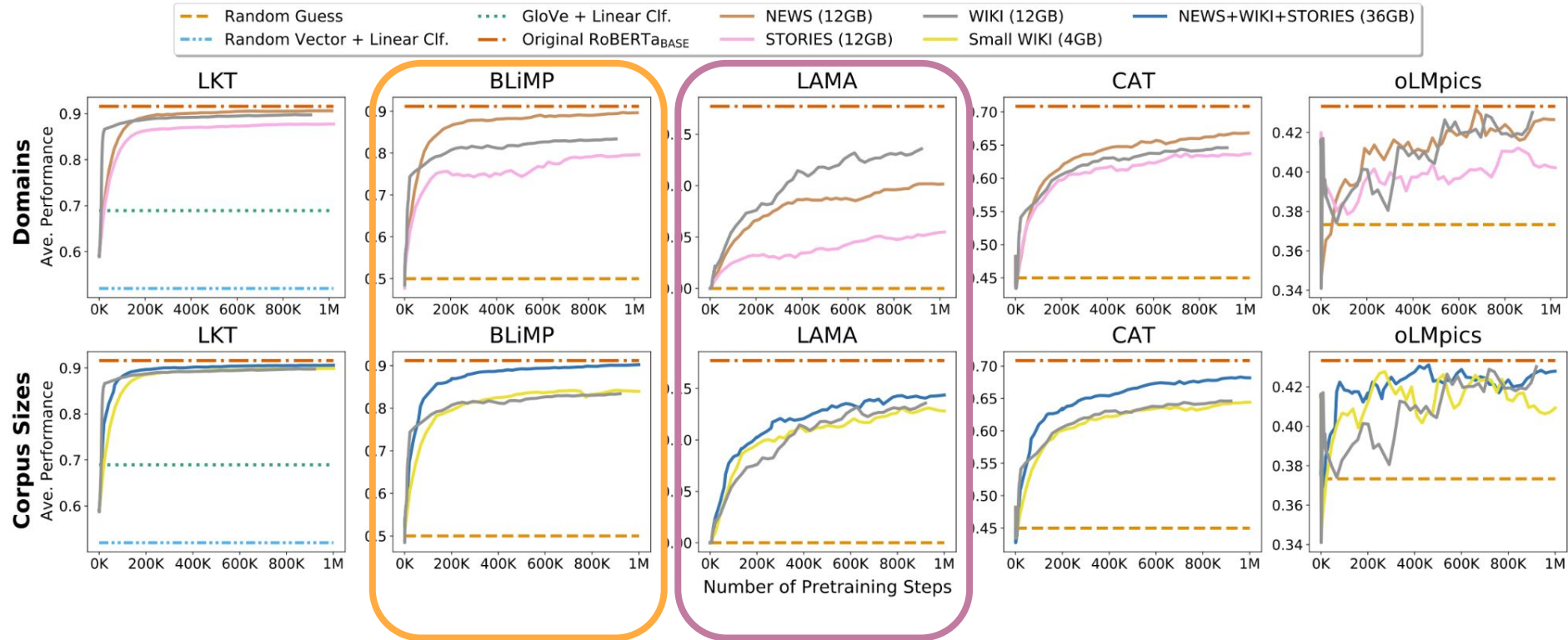, etc.

# Varying Pretraining Corpus

**Domains:**

- English WIKI (12 GB)

- NEWS (12 GB)

- STORIES (12 GB)

**Corpus Size:**

- Small English WIKI (4 GB)

- English WIKI (12 GB)
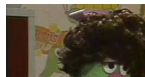
- English WIKI + NEWS + STORIES (36 GB)

Legend: Random Guess; Random Vector + Linear Clf.; GloVe + Linear Clf.; Original RoBERTa_BASE; NEWS (12GB); STORIES (12GB); WIKI (12GB); Small WIKI (4GB); NEWS+WIKI+STORIES (36GB)

**TL;DR:**

- Observed learning curriculum remains the same

- Domains affect learning more than corpus sizes

# Research Benchmarks

❄️

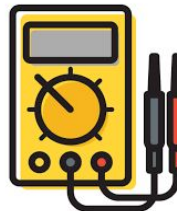Finetuning: CoLA, MNLI, SQuAD, etc.



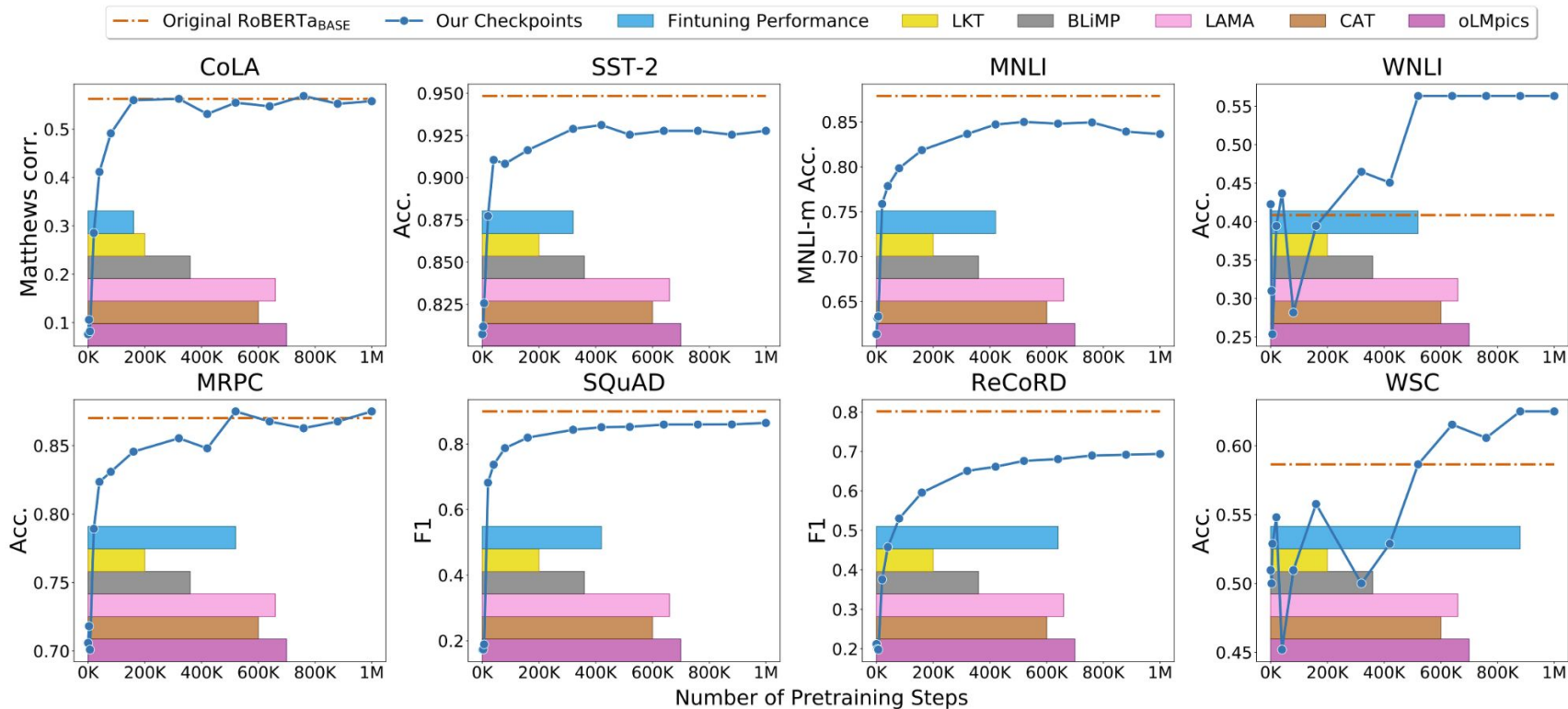**RoBERTa$_1$**     **RoBERTa$_2$**     **RoBERTa$_3$**          **RoBERTa$_{Fina}$**
l

❄️

Linguistic, factual, commonsense, etc.

**TL;DR:**

Among finetuning tasks, **ordering of difficulties exists** -- more knowledge required, more difficult

**Main Contribution:**

- Most systematic work of learning dynamics for pretraining yet

- Learning curriculum:

  Linguistic (knowledge) $>$ Factual $\cong$ Commonsense $\gg$ Reasoning

- Domain diversity matters more than just corpus size

- Ordering of difficulties among downstream tasks

- As models evolve and new probes emerge, ***probing across time*** framework can serve as a general framework to inform progress on both fronts

# Thanks!

Check our paper for more details and discussion!