What Neural Networks Can Teach Us About Language Acquisition



Alex Warstadt ETH Zurich

Presented at Aix-Marseille University 1 December 2022

In NLP, there are bigger and more powerful language models (LMs) all the time.

| Language models are few-shot learners | | | | |
|--|--|---|---|---|
| T Brown, B Mann, N Ryder Advances in neural 2020 - proceedings.neurips.cc | | | | |
| GPT-3 achieves strong perform | Scaling language models: Meth JW Rae, S Borgeaud, T Cai, K Millican | ods, analysis & insights from arXiv preprint arXiv, 2021 - arxiv. | n training gopher ^{org} | |
| the contract of the contract | We conclude that Gopher lifts the basel Gopher nears human rater performance o | Palm: Scaling language m <u>A Chowdhery, S Narang</u> , J Devlin, <u>N</u> | odeling with pathw M Bosma arXiv preprin | ays t arXiv, 2022 - arxiv.org |
| | ☆ Save 50 Cite Cited by 162 Related | We present the model FLOPs ut | tilization of PaLM 540B m | odel and PaLM model to several |
| different language models for code. First, we compare to the LaMDA 137B parameter | | | | |
| | | ☆ Save 59 Cite Cited by 251 | Related articles All 5 ver | sions 📎 |

In NLP, there are bigger and more powerful language models (LMs) all the time.

| Language models are few | -shot learners |
|--|--|
| T Brown, B Mann, N Ryder Ad | vances in neural 2020 - proceedings.neurips.cc |
| GPT-3 achieves strong perform | Scaling language models : Methods, analysis & insights from training gopher JW Rae, S Borgeaud, <u>T Cai, K Millican</u> - arXiv preprint arXiv, 2021 - arxiv.org |
| where GPT3 's few-shot learning st ☆ Save ⑰ Cite Cited by 6426 | We conclude that Gopher lifts the base Palm : Scaling language modeling with pathways Gopher nears human rater performance of <u>A Chowdhery, S Narang</u> , J Devlin, <u>M Bosma</u> arXiv preprint arXiv, 2022 - arxiv.org |
| | ☆ Save 57 Cite Cited by 162 Related We present the model FLOPs utilization of PaLM 540B model and PaLM model to several |
| | different language models for code. First, we compare to the LaMDA 137B parameter model (|

And researchers are racing to evaluate their strengths and limitations.

| Targeted syntactic evaluation | of language models | | | |
|---|------------------------------------|--|----------------------|--|
| R Marvin, T Linzen - arXiv preprint arXiv | v:1808.09031, 2018 - arxiv.org | | | |
| We present a dataset for evaluating | A primer in bertology: Wh | at we know about how bert works | | |
| model. We automatically construct a la | A Rogers, O Kovaleva, A Rumshis | ky - Transactions of the Association 2020 - direct mit edu | | |
| ☆ Save 579 Cite Cited by 292 Rel | Transformer-based models have p | u Emergent abilities of large language models | | |
| | understanding of what is behind th | e J Wei, Y Tay, R Bommasani, C Raffel, B Zoph arXiv preprint arX | iv, 2022 - arxiv.org | |
| | ☆ Save ፵ Cite Cited by 805 | when evaluated on a sufficiently large language model. Hence, their emergence cannot | | |
| L | | predicted by simply extrapolating performance on smaller-scale models. Emergent few-shot | | |
| | | ☆ Save 572 Cite Cited by 36 All 4 versions ≫ | | |

So what currently holds the state-of-the-art in grammar learning?

So what currently holds the state-of-the-a grammar learning?

Humans are data-efficient language learners.



Roadmap

1. The Argument

Why study language models as models of human learners? What kinds of questions can they address?

2. The Learner

Training LMs in more plausible learning environments, without major advantages over humans.

3. A Controlled Experiment

A proof-of-concept: Language models unlock new modes of testing learnability hypotheses

Part 1: The Argument

1. The Argument

Why study language models as models of human learners? What kinds of questions can they address? 2. The Learner

Training LMs in more plausible learning environments, without major advantages over humans. 3. A Controlled Experiment

A proof-of-concept: Language models unlock new modes of testing learnability hypotheses How do we make LMs into plausible cognitive models?

Why study LMs as models of human learners?

What kinds of questions can they address?

What Artificial Neural Networks Can Tell Us About Human Language Acquisition*

Alex Warstadt, Samuel R. Bowman

August 18, 2022



Language Deprivation Experiments



Pharaoh Psamtik (664 – 610 BCE)



Frederick II (1194-1250)



James IV (1473-1513)

Three reasons to study neural networks instead of humans:

1. Ethics

Ethics Expense

1. Ethics 2. Expense 3. Experimental paradigms

As with any scientific model, there are obvious limitations with LMs.



Debates in language acquisition often center around the sufficient conditions for human-learnability.

Suppose the model **SUCCEEDS** given some experimental manipulation. How likely are humans also to succeed?



Likelihood that humans show same result

Human is at a great advantage

Model is at a great advantage

Suppose the model FAILS given some experimental manipulation. How likely are humans also to succeed? Likelihood that humans show same result



Human is at a great advantage

Model is at a great advantage







• Increase the relevance of positive by results by reducing the advantages that models have over humans.

- Increase the relevance of positive by results by reducing the advantages that models have over humans.
- Increase the chance of positive results by providing models with advantages that humans have.

- Increase the relevance of positive by results by reducing the advantages that models have over humans.
- Increase the chance of positive results by providing models with advantages that humans have.

In other words, the best model learners will be ones whose environments and innate abilities are as rich as possible without being richer than those of humans.

Advantages ANNs Have

Data quantity



Data domain



Orthography



Advantages ANNs Have

Data quantity



Data domain



Orthography



Human vs. model linguistic input (# of word tokens)



Part 2: The Learner

1. The Argument

Why study language models as models of human learners? What kinds of questions can they address?

2. The Learner

Training LMs in more plausible learning environments, without major advantages over humans. 3. A Controlled Experiment

A proof-of-concept: Language models unlock new modes of testing learnability hypotheses

MiniBERTas





Five Sets of Probing Methods



Five Sets of Probing Methods

None



BLiMP: The Benchmark of Linguistic Minimal Pairs for English

Alex Warstadt¹, Alicia Parrish¹, Haokun Liu², Anhad Mohananey², Wei Peng², Sheng-Fu Wang¹, Samuel R. Bowman^{1,2,3}

30B

Pretraining Dataset Size

Minimal Pairs

A pair of two nearly identical sentences which differ in grammatical acceptability.

Betsy is <u>eager</u> to sleep.

Betsy is <u>easy</u> to sleep.



BLiMP Categories

Morphology

- Anaphor agr.
- Determiner-noun agr.
- Irregular forms
- Subj-verb agr

Syntax

- Argument structure
- Binding
- Control/raising
- Ellipsis
- Filler-gap (wh) deps.
- Island constraints

Semantics

- NPIs
- Quantifiers
Data Sample

| Phenomenon | Acceptable example | Unacceptable example |
|-------------------|---|---|
| Anaphor agr. | Many girls insulted <u>themselves</u> . | Many girls insulted <u>herself</u> . |
| Detnoun agr. | Rachelle had bought that <u>chair</u> . | Rachelle had bought that chairs. |
| Subject-verb agr. | These casseroles <u>disgust</u> Kayla. | These casseroles <u>disgusts</u> Kayla. |
| Filler-gap | Brett knew <u>what</u> many waiters find. | Brett knew <u>that</u> many waiters find. |
| Island effects | Which <u>bikes</u> is John fixing? | Which is John fixing <u>bikes</u> ? |

BLiMP Learning Curves







BLiMP: Human vs. Human-scale model



UPCOMING Shared task @ CMCL/CoNLL 2023



BabyLM Challenge

Sample-efficient pretraining on a developmentally plausible corpus

Approximate Timeline:

December 2022: Training data released February 2023: Shared evaluation pipeline published June 2023: Submissions for presentation at CMCL due September 2023: CMCL, Submissions for presentation at CoNLL November 2023: CoNLL

Part 3: A Controlled Experiment

1. The Argument

Why study language models as models of human learners? What kinds of questions can they address? 2. The Learner

Training LMs in more plausible learning environments, without major advantages over humans. 3. A Controlled Experiment

A proof-of-concept: Language models unlock new modes of testing learnability hypotheses

Testing the Poverty of the Stimulus: Controlled experiments at the scale of human language learning

with Yian Zhang, Haau-Sing Li, and Samuel R. Bowman





There's a long history of debate about the role of innate bias in the acquisition of hierarchical syntactic rules in favor of alternative linear rules.

Chomsky (1965, 1971); Crain & Nakayama (1987); Lewis & Elman (2001); Pullum & Scholz (2002); Legate & Yang (2002); Reali & Christiansen (2005); Perfors, Tenenbaum, & Regier (2011); Berwick et al. (2011); Hsu, Chater, & Vitanyi (2013); McCoy, Frank, & Linzen (2018, 2020); Warstadt & Bowman (2020)



A theory that attributes possession of certain linguistic universals to a languageacquisition system [...] implies that only certain kinds of symbolic systems can be acquired [....] Specifically, grammatical transformations are necessarily "structure-dependent" [i.e., hierarchical] [....] It is impossible, however, [for the language-acquisition system to learn] a transformation such a simple operation as reflection of an arbitrary string.

(Chomsky, 1965)

Subject Auxiliary Inversion

The zebra **does** chuckle.

Does the zebra chuckle?

Example: McCoy et al. (2020)

MOVE-FIRST: Move the <u>linearly first</u> auxiliary to the front of the sentence.

does the zebra does chuckle

MOVE-FIRST: Move the <u>linearly first</u> auxiliary to the front of the sentence.

does the zebra does chuckle

MOVE-MAIN: Move the <u>main verb's</u> auxiliary to the front of the sentence.



MOVE-FIRST: Move the <u>linearly first</u> auxiliary to the front of the sentence.

MOVE-FIRST: Move the <u>linearly first</u> auxiliary to the front of the sentence.



MOVE-MAIN: Move the <u>main verb's</u> auxiliary to the front of the sentence.

MOVE-MAIN: Move the main verb's auxiliary to the front of the sentence. has the man seen the cat who has gone

The Poverty of the Stimulus Argument relies on quantifying "sufficient evidence":

"Surely, if children hear enough sentences like those [below], then they could reject the [move-first] hypothesis. But if such evidence is virtually absent from the linguistic data, one cannot but conclude that children do not entertain the [move-first] hypothesis, because the knowledge of structure dependency is innate."

(Legate & Yang, 2001)

e.g., Has the man who has gone seen the cat?

The INDIRECT Evidence Hypothesis:

Indirect evidence may be sufficient for a learner without hierarchical bias to eliminate move-first.

Proponents of Indirect Evidence

"While a child may not receive direct evidence about the correctness of a particular hierarchical phrase structure rule..., there is vast indirect evidence for the general superiority of syntax with that structure throughout language. A learner who adopts a hierarchical phrase structure framework for describing the syntax of English will arrive at a much simpler, more explanatory account of her observations than a learner who adopts a linear framework." (Perfors et al., 2011)

Syntactic filtering



Syntactic filtering

Training data: 1B words from books & Wikipedia

- Percent filtered: 1.7%
- Accuracy: 98%
- Recall (% of direct evidence removed): 99%



Distribution of direct evidence (by domain)



Models

Filtered Condition

Unfiltered Condition (control)

24 RoBERTa models pretrained from scratch.

- 2 main conditions
- 4 sizes
- 3 runs (failed runs discarded)





Masked Sentence A

Unlabeled Sentence A and B Pa

Results: General acceptability judgments on BLiMP

Question: Did the removal of direct evidence have effects on unrelated phenomena?

Results: General acceptability judgments on BLiMP

Question: Did the removal of direct evidence have effects on unrelated phenomena?



Answer: No

Results: General acceptability judgments on BLiMP

This result holds across all phenomena in BLiMP.



Results: Subject Aux Inversion

Question: Did the removal of direct evidence affect learning of the target phenomenon?

Results: Subject Aux Inversion

Question: Did the removal of direct evidence affect learning of the target phenomenon?

Answer: Yes, in the written domain.



Results: Subject Aux Inversion (fine-grained)

Question: Did the removal of direct evidence differentially affect Only Move-Main examples?

Results: Subject Aux Inversion (fine-grained)

Question: Did the removal of direct evidence differentially affect Only Move-Main examples?

Answer: Yes





Results: Subject Aux Inversion (BEST CASE)

Question: Is indirect evidence sufficient to acquire the hierarchical rule?

Results: Subject Aux Inversion (BEST CASE)

Question: Is indirect evidence sufficient to acquire the hierarchical rule?

Answer: Yes



Results: Subject Aux Inversion (BEST CASE)

This result holds across all test cases in the written domain.





The results support the indirect evidence hypothesis, but with important caveats.

- How reproducible is the best model's success?
- How important are small amounts of direct evidence that passed through the filter?
- Can models succeed with the same data-volume limitations as humans?

Discussion: What does indirect evidence for hierarchical structure look like?

1. Classic constituency tests

Fragment answers

<u>Who</u> has seen the cat? [The man who was here this afternoon]

Coordination

John and [the man who was here this afternoon] are friends.

Pronominalization [The man who was here this afternoon] left. <u>He</u> saw the cat.
Discussion: What does indirect evidence for hierarchical structure look like?

2. Other hierarchical rules

Subject Verb Agreement [The man who saw the cats] <u>is</u> here.

Passivization

I greeted [the man who saw the cat.] \rightarrow [The man who saw the cat] was greeted by me.

Conclusion

Computational model learners allow us to:

- Make causal inferences about the effects of the environment on language learning.
- Conduct controlled experiments on the scale of human language learning.

BUT we are still far from having developmentally plausible learners.

Why are humans more data efficient than LMs?

BUT we are still far from having developmentally plausible learners.

Why are humans more data efficient than LMs?

Multimodal input

Interactive learning



Open questions

- How do we quantify indirect evidence, and locate relevant indirect evidence in the input?
- Does indirect evidence drive typological tendencies across the world's languages?
- How do we compare human and model developmental trajectories?
- How do we optimize LM training at small scales?
- How do we incorporate interactive learning signals into LM training?

Thank you!

TALEP Marseille for hosting me, especially Mitja Nikolaus, Abdellah Fourtassi, Carlos Ramisch, and Sylvie Ros!

Collaborators: Sam Bowman, Amanpreet Singh, Alicia Parrish, Yian Zhang, Haokun Liu, Haau-Sing Li, Sheng-Fu Wang, Anhad Mohananey, Wei Peng.

Audiences who helped improve previous versions of this talk: NYU Text-as-Data, Stanford Language & Cognition Lab, CoLaLa (UC Berkeley), UT Austin, ETH Zurich, ENS Paris, IST Unbabel Lisbon,

This work was funded in part by the NSF.