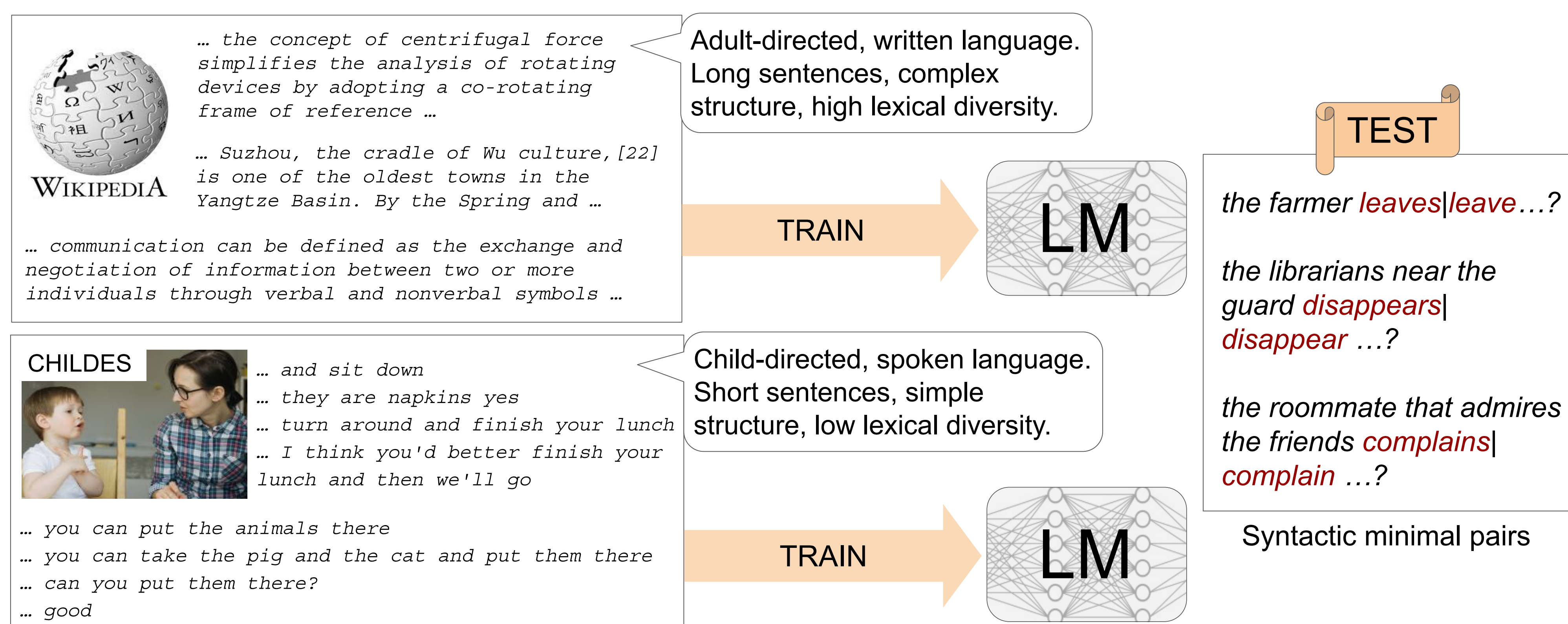


## Motivation

LMs have human-like linguistic skills, but are trained in non-human ways. What if we trained LMs in more **developmentally plausible** ways?

Training on **Child-Directed Language (CDL)** has been argued to **benefit syntax learning** in English LMs (Huebner et al., 2021; Salhan et al. 2024), **but ...**

**Are CDL benefits consistent across models, languages, and syntactic benchmarks?**



We find no consistent benefit of training LMs on Child-Directed Language for syntactic learning; effects vary by model architecture, language, and benchmark, and are partly driven by lexical frequency and corpus composition.

## Evaluation setup

- 2 models architectures: GPT2 (Causal LM) & RoBERTa (Masked LM)
- 3 language datasets: English, French, German
- 3 syntactic benchmarks: BLiMP, Zorro, CLAMS

MIXED RESULTS! → CHILDES  $\approx$  Wiki CHILDES  $>$  Wiki CHILDES  $<$  Wiki

	Model	Training Data	BLiMP	Zorro	CLAMS		
					English	French	German
GPT-2	CLM	CHILDES	0.61 $\pm$ 0.02	<b>0.76 <math>\pm</math> 0.04</b>	0.60 $\pm$ 0.01	0.64 $\pm$ 0.01	0.69 $\pm$ 0.03
		Wiki	0.61 $\pm$ 0.02	0.69 $\pm$ 0.04	<b>0.71 <math>\pm</math> 0.01</b>	<b>0.80 <math>\pm</math> 0.01</b>	<b>0.81 <math>\pm</math> 0.01</b>
RoBERTa	MLM	CHILDES	0.59 $\pm$ 0.03	0.66 $\pm$ 0.05	0.57 $\pm$ 0.02	0.59 $\pm$ 0.02	0.70 $\pm$ 0.01
		Wiki	0.59 $\pm$ 0.03	0.67 $\pm$ 0.03	<b>0.63 <math>\pm</math> 0.01</b>	<b>0.69 <math>\pm</math> 0.01</b>	<b>0.75 <math>\pm</math> 0.01</b>

## Better evaluation: Frequency Informed Testing (FIT)

Problem: Comparing syntactic accuracy of models trained on different corpora is tricky because of different vocabularies!

Our solution: Generate minimal pairs balanced by corpus-specific lexical frequency, ensuring a fair spread of high- and low-frequency items *for each corpus*.

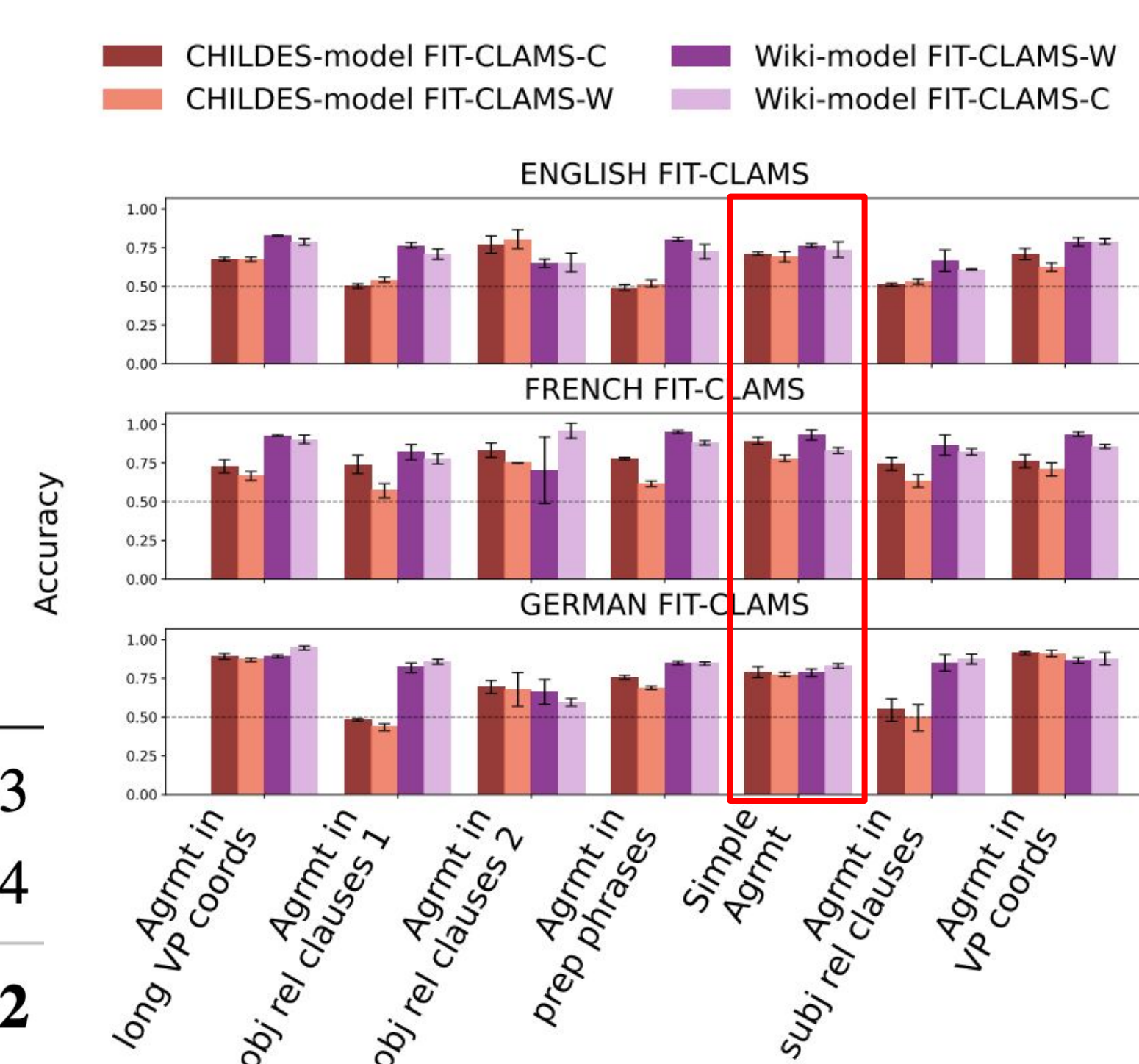
EN Nouns	Bin	Freq	Df	EN Verbs	Bin	Freq	Long VP	Df
roommate, roommates	0	2	CHI	awaits, await	0	2	<i>the guests</i>	CHI
resident, residents	1	6	CHI	complains, complain	1	8	<i>about the noise</i>	CHI
librarian, librarians	2	13	CHI	arrives, arrive	2	17	<i>at the station</i>	CHI
officer, officers	3	36	CHI	disappears, disappear	2	42	<i>from the scene</i>	CHI
toddler, toddlers	4	90	CHI	bows, bow	4	243	<i>to the king</i>	CHI
farmer, farmers	5	264	CHI	hides, hide	4	391	<i>from the chicken</i>	CHI
policeman, policemen	6	380	CHI	leaves, leave	6	1793	<i>the room</i>	CHI
doctor, doctors	7	656	CHI	sits, sit	7	4219	<i>in the car</i>	CHI
man, men	8	2156	CHI	thinks, think	8	14710	<i>about the trip</i>	CHI
daddy, daddies	9	7027	CHI	goes, go	9	27620	<i>to the new store</i>	CHI

## Results of CLM models on FIT-CLAMS

As expected, LMs perform better on minimal pairs created with in-distribution lexical items.

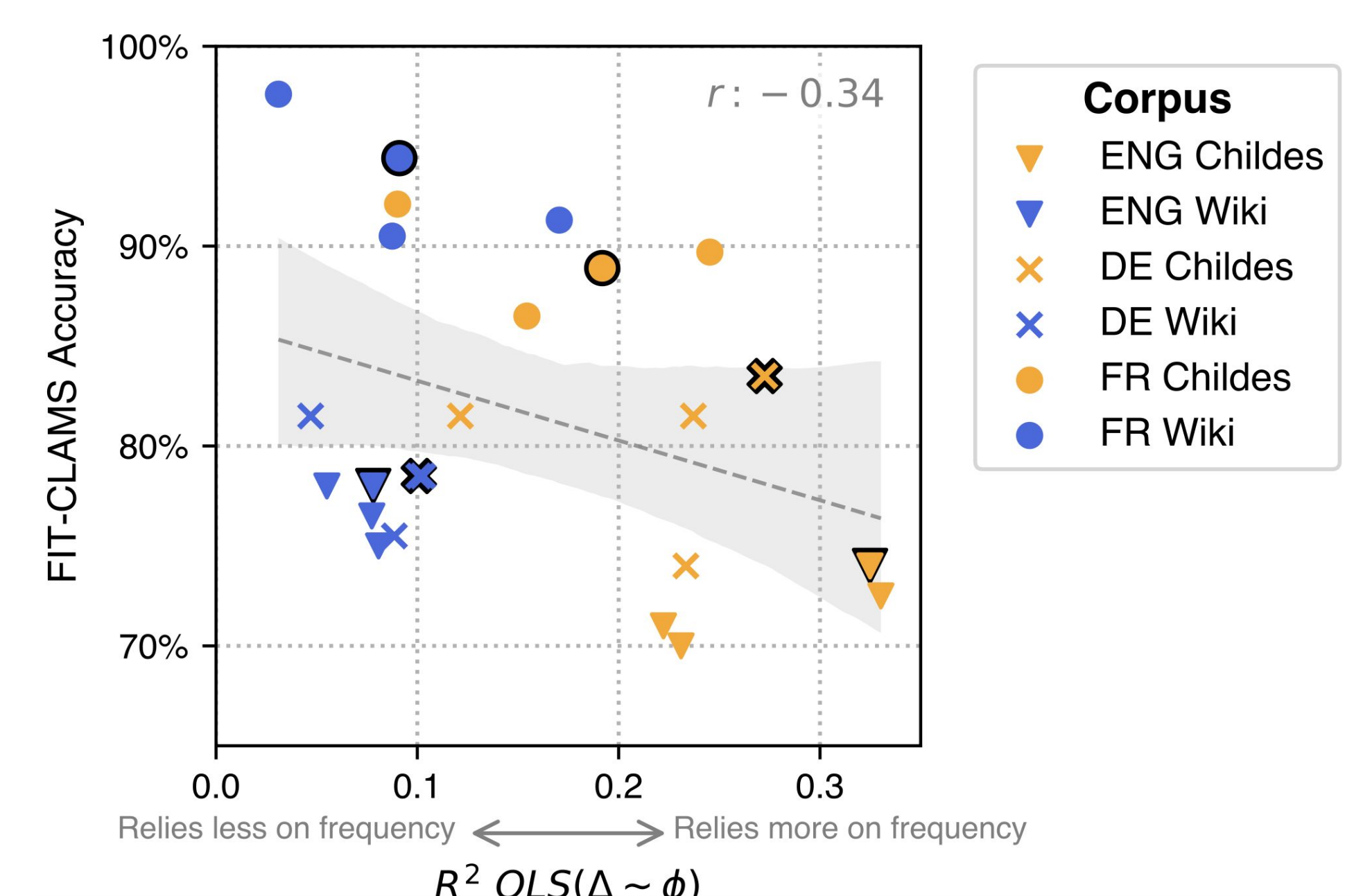
However, **LMs trained on Wiki still outperform those trained on CHILDES** in most cases including Simple Agreement, even when controlling for frequency effects

Training	Eval. lex.	EN	FR	DE
CHILDES	CHILDES	0.63 $\pm$ 0.02	0.78 $\pm$ 0.04	0.73 $\pm$ 0.03
	Wiki	0.63 $\pm$ 0.03	0.67 $\pm$ 0.03	0.69 $\pm$ 0.04
Wiki	CHILDES	<b>0.72 <math>\pm</math> 0.03</b>	<b>0.86 <math>\pm</math> 0.02</b>	<b>0.83 <math>\pm</math> 0.02</b>
	Wiki	<b>0.75 <math>\pm</math> 0.02</b>	<b>0.88 <math>\pm</math> 0.06</b>	<b>0.82 <math>\pm</math> 0.03</b>



## Regression Analysis

How does lexical frequency in minimal pairs affect model accuracy? ⇨ Models that rely more on lexical frequency (higher fit) tend to perform worse on FIT-CLAMS.



## Conclusions

- CDL shows no clear benefit for syntax learning
- This holds for current modeling approaches: LMs trained in **static, non-interactive** environments, without feedback, developmental grounding, or cognitive constraints, unlike human learners
- CDL may still hold promise in interactive, situated learning environments, shifting focus toward the **communicative and contextual factors essential to language acquisition**