

Do Language Models Accommodate Their Users?

A Study of Linguistic Convergence

Terra Blevins, Susanne Schmalwieser, Benjamin Roth



universität
wien

LM-generated text is extremely fluent

Jones and Bergen. [Large Language Models Pass the Turing Test](#).
arxiv, 2025.

C

Witness A	Witness B
are you human?	are you human?
yes	yeah
what do you want your future career to look like?	what do you want your future career to look like?
i would like to be a part of a company where I could do creative direction	thinking about grad school for psych, maybe research
what kind of creative direction?	what subject of research would you be interested in
managing the creative direction or certain projects and marketing	cognitive psych, specifically memory & decision making
what do you study in school?	are you in any labs now?
i study business psychology	yeah, research assistant in the Dobkins lab on campus
what made you decide to study that?	what do they research?
i was in interested in both business and psychology and I saw that they offered a major here that had both in 1	behavioral neuroscience, mostly sensory development & perception

LM-generated text is extremely fluent

C

Human

are you human?

yes

what do you want your future career to look like?

i would like to be a part of a company where I could do creative direction

what kind of creative direction?

managing the creative direction or certain projects and marketing

what do you study in school?

i study business psychology

what made you decide to study that?

i was in interested in both business and psychology and I saw that they offered a major here that had both in 1

LLaMa3-405B

are you human?

yeah

what do you want your future career to look like?

thinking about grad school for psych, maybe research

what subject of research would you be interested in

cognitive psych, specifically memory & decision making

are you in any labs now?

yeah, research assistant in the Dobkins lab on campus

what do they research?

behavioral neuroscience, mostly sensory development & perception

Jones and Bergen. [Large Language Models Pass the Turing Test](#).
arxiv, 2025.

LM fluency opens many questions about how models use language compared to humans, and about human-model interactions



Convergence: process where a speaker adapts their language use (vocabulary, dialect, speech patterns) to resemble that of their conversation partner

We test whether LMs exhibit **stylistic convergence** in dialogue settings

Measuring **Linguistic Convergence** in Language Models



Did you catch the game last night?

Speaker 1

Yeah, it was incredible! What a finish.



Speaker 2



I can't believe they scored in the last minute.

I know! I almost turned it off early.



That would've been a big mistake.

Definitely. Best game of the season so far.



Measuring **Linguistic Convergence** in Language Models



Did you catch the game last night?

Speaker 1

Yeah, it was incredible! What a finish.



Speaker 2



I can't believe they scored in the last minute.

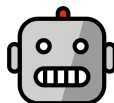
I know! I almost turned it off early.



That would've been a big mistake.

No kidding. I almost missed the best part!

LM



Measuring **Linguistic Convergence** in Language Models



Did you catch the game last night?

Speaker 1

Yeah, it was incredible! What a finish.



Speaker 2



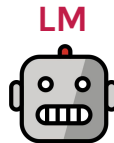
I can't believe they scored in the last minute.

I know! I almost turned it off early.



That would've been a big mistake.

No kidding. I almost missed the best part!



Prompt **LM** to take on role of **Speaker 2** after 5 turns in existing dialogues:

- Direct comparison of **LM** and **human** “generations”
- **Synthetic evaluation**

Measuring **Linguistic Convergence** in Language Models

Utterances:

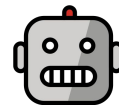
reference



That would've been a big mistake.

test (LM)

No kidding. I almost missed the best part!



Measuring **Linguistic Convergence** in Language Models

Utterances:

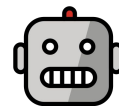
reference



That would've been a big mistake.

test (LM)

No kidding. I almost missed the best part!



control_{human}

Definitely. Best game of the season so far.



control_{random}

I think I left my jacket there.



baselines

Convergence Metrics

Given an utterance **a** and reference utterance **b**:

→ **LIWC** (*Linguistic Inquiry and Word Count*) Agreement

How much does the distribution of function word classes agree in **a** and **b**?

I didn't think it would work.

She never told me the truth.

Convergence Metrics

Given an utterance **a** and reference utterance **b**:

→ **LIWC** (*Linguistic Inquiry and Word Count*) Agreement

How much does the distribution of function word classes agree in **a** and **b**?

I didn't think it would work.

She never told me the truth.

High agreement on:

- Personal pronouns
- Negations
- Quantifiers (e.g., none!)

Low agreement on:

- Auxiliary verbs

Convergence Metrics

Given an utterance **a** and reference utterance **b**:

→ **LIWC** (*Linguistic Inquiry and Word Count*) Agreement

How much does the distribution of function word classes agree in **a** and **b**?

$LSM_x = 1 - |a - b| / (a + b)$ for each word class **X** = {personal and impersonal pronouns, articles, conjunctions, prepositions, auxiliary verbs, frequently used adverbs, negations, and quantifiers} → report **per-class** and **average** LSM

Common measure of linguistic convergence in computational research (*Ireland et al., 2011; Danescu-Niculescu-Mizil and Lee, 2011; Bhatt and Rios 2021; i.a.*)

Convergence Metrics

Given an utterance **a** and reference utterance **b**:

→ **LIWC Agreement**:
$$1 - \frac{(|a| - |b|)}{(|a| + |b|)}$$

→ **Utterance Length**:
$$1 - \frac{(|a| - |b|)}{(|a| + |b|)}$$

→ **Token Novelty**:
$$\frac{|\{w \in a\} \cap \{w \in b\}|}{|\{w \in a\}|}$$

→ **PROPN Overlap**:
$$|\text{PROPN}(a) \cap \text{PROPN}(b)|$$

Experiments

Comprehensive analysis with:

3 Datasets

- **DailyDialog**¹ general-topic spoken conversations by EN language learners
- **NPR**² radio interview transcripts
- **Movie Corpus**³ fictional conversations drawn from movie scripts

¹ (Li et al., 2017)

² (Majumder et al., 2020)

³ (Danescu-Niculescu-Mizil and Lee, 2011)

	Dataset Statistics		
	DailyDialog	Movie	NPR
Conversations	707	1,000	1,000
Avg. Turns	9.79	8.98	17.57
Avg. Turn Length	13.44	10.87	48.43
Replaced Turns	1,918	2,280	6,568

Experiments

Comprehensive analysis with:

3 Datasets

- **DailyDialog**¹ general-topic spoken conversations by EN language learners
- **NPR**² radio interview transcripts
- **Movie Corpus**³ fictional conversations drawn from movie scripts

¹ (Li et al., 2017)

² (Majumder et al., 2020)

³ (Danescu-Niculescu-Mizil and Lee, 2011)

16 Models

- **Gemma 3** (1B, 4B, 12B, 27B)
- **LLaMa 3** (1B, 3B, 8B, 70B)
- *Pretrained and instruction-tuned variants*

	Dataset Statistics		
	DailyDialog	Movie	NPR
Conversations	707	1,000	1,000
Avg. Turns	9.79	8.98	17.57
Avg. Turn Length	13.44	10.87	48.43
Replaced Turns	1,918	2,280	6,568

Experiments

Two control **baselines**:

- **Random baseline**: drawn at random from different conversation in dataset
- **Human baseline**: the original, human-authored utterance at timestep t

Reference

“I will take ten roses.”

Random

“Oh how nice. They’re bright rooms and the house is very quiet.”

Human

“Do you want to add some baby’s breath for that?”

Model-Gen.

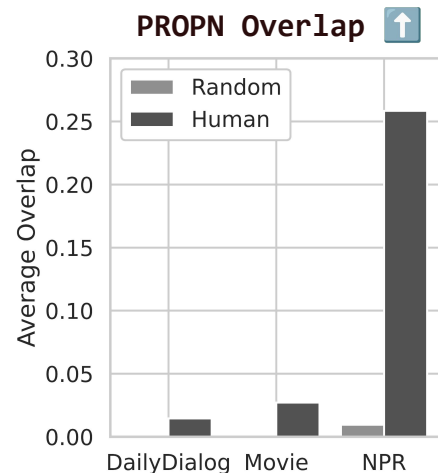
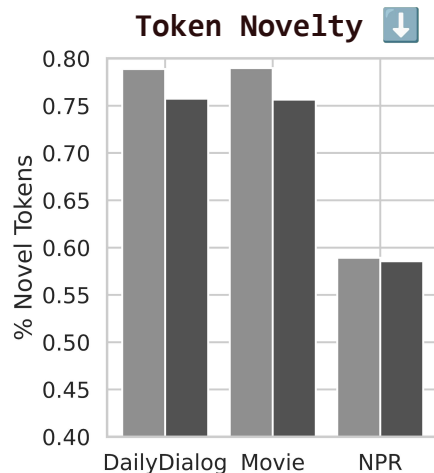
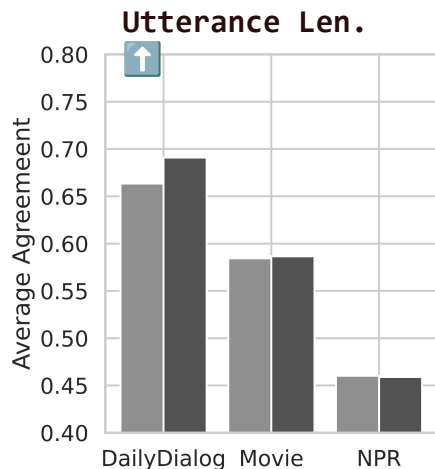
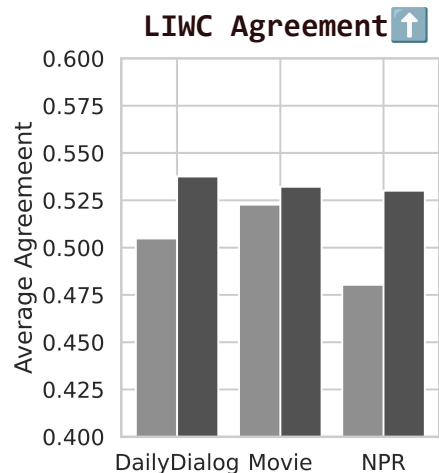
“Do you want them delivered?”

* From DailyDialog

Experiments

Two control **baselines**:

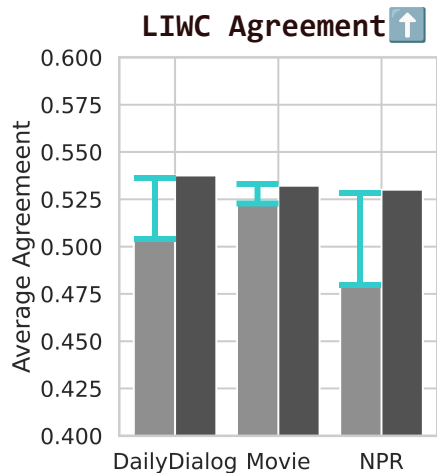
- **Random baseline**: drawn at random from different conversation in dataset
- **Human baseline**: the original, human-authored utterance at timestep t



Experiments

Two control **baselines**:

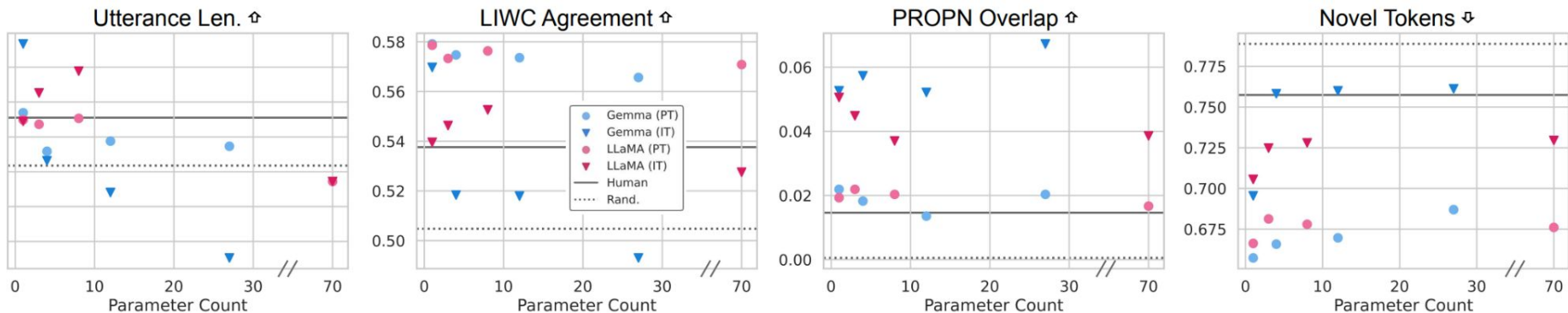
- **Random baseline**: drawn at random from different conversation in dataset
- **Human baseline**: the original, human-authored utterance at timestep t



If **model generations** exhibit agreement...

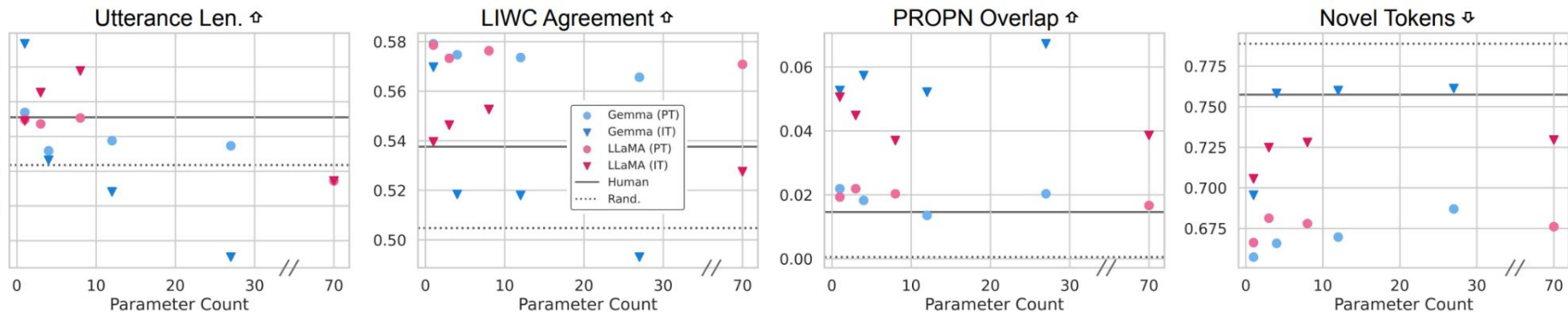
- \leq than **random**: converges at or worse than chance in dataset's distribution
- Between **random** and **human**: converges, but less than human interlocutors
- $>$ than **human**: model over-converges relative to humans

Convergence Results: DailyDialog



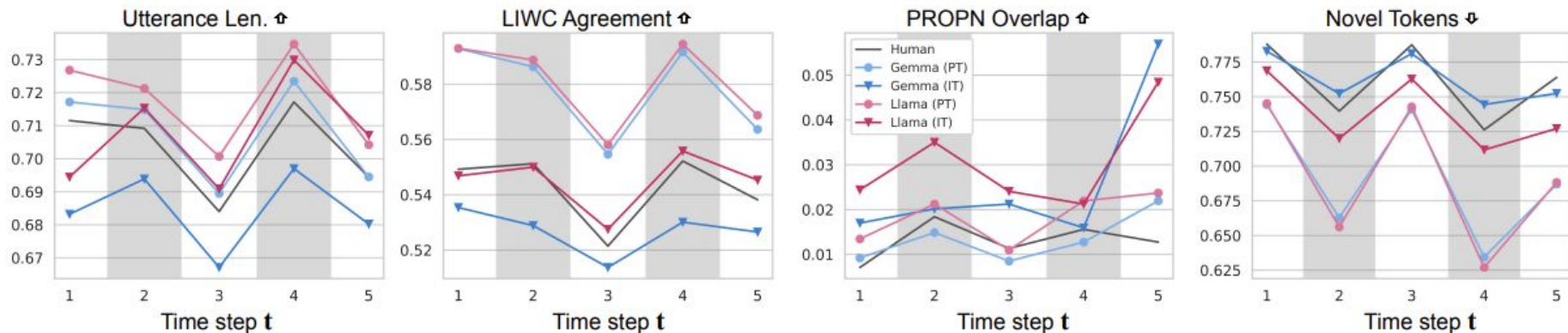
- LMs almost always converge more than the random baseline...
- ...and often **over-converge** relative to the human baseline

Convergence Results: DailyDialog



- LMs almost always converge more than the random baseline...
- ...and often **over-converge** relative to the human baseline
- Pretrained models adapt **more** than their **instruction-tuned** counterparts (except on PROPn overlap)
- While larger models trend toward human-level convergence, size effects are **not** statistically significant*

Convergence Results: Across Time

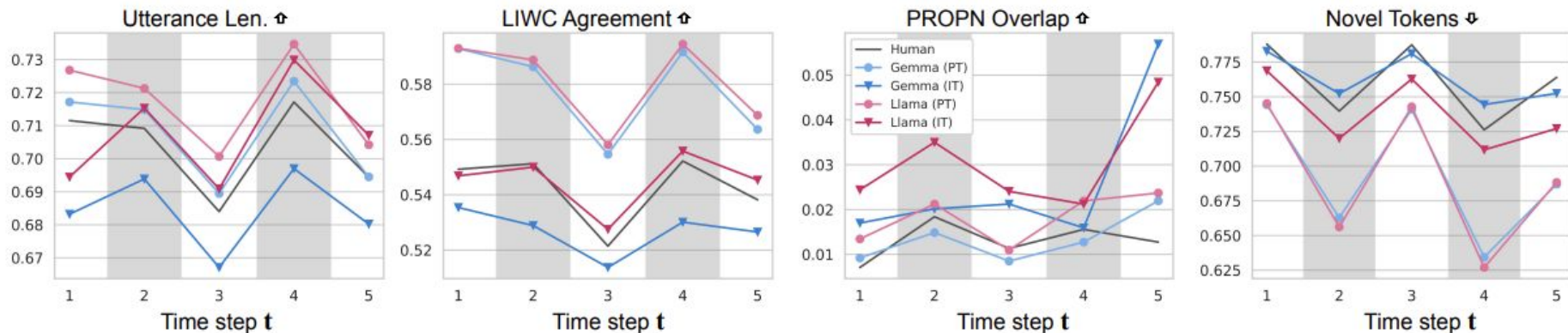


- Models mirror the human zigzag pattern across time, suggesting model sensitivity to **speaker roles**, with IT vs. PT trends holding at each timestep
- Exception: PROP N Overlap, where models demonstrate strong recency bias towards prior turn $t-1$

* on DailyDialog

Convergence Results: Across Time

More findings in the paper!



- Models mirror the human zigzag pattern across time, suggesting model sensitivity to **speaker roles**, with IT vs. PT trends holding at each timestep
- Exception: PROP N Overlap, where models demonstrate strong recency bias towards prior turn **$t-1$**

* on DailyDialog

Takeaways

LMs strongly mirror the style of their context, often **over-converging** relative to human utterances

Model convergence is driven by **training scheme** (PT > IT) more than model size

LMs show sensitivity to speaker roles and **recency** bias when considering multiple prior conversational turns

Check out
the paper for
more details!



Questions?



<http://blvns.github.io>



t.blevins@northeastern.edu



<https://www.benjaminroth.net/>



benjamin.roth@univie.ac.at