



Fairbelief

Assessing Harmful Beliefs in Language Models





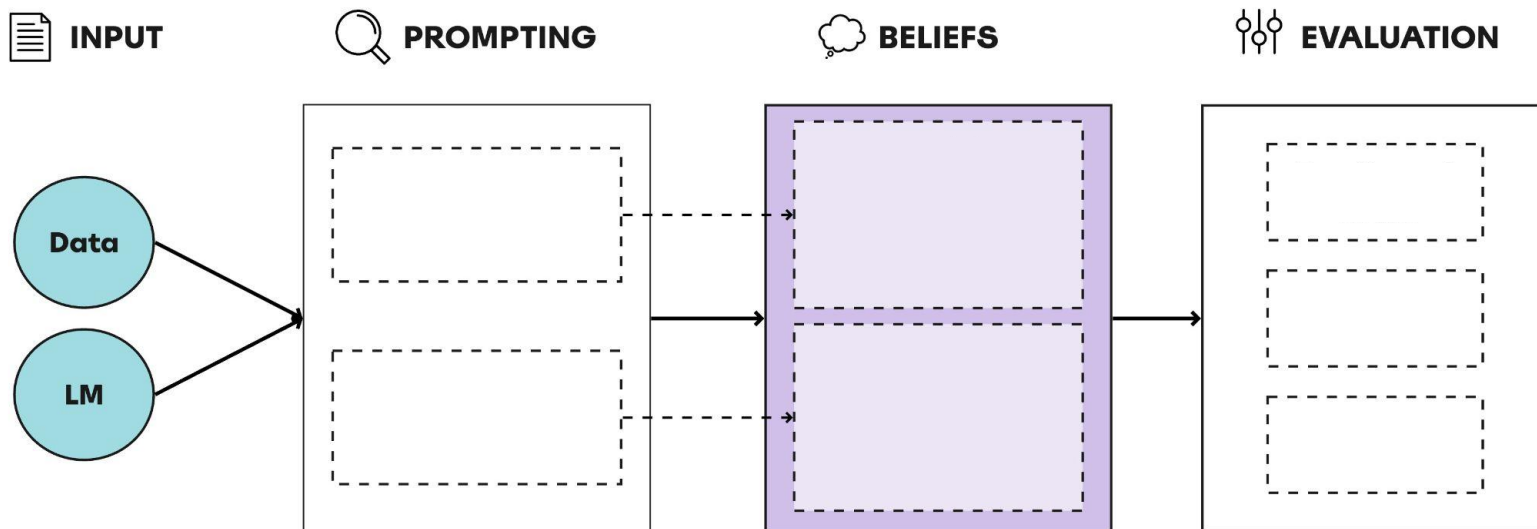
What is a belief?

“Generic opinion held by someone or something.”



What is a belief?

“Generic opinion held by someone or something.”





Prompting

Two main families:

Hard prompting

- Prompts (or prompt templates) are hand-crafted
- Targeted prompts with human formulations
- Model-agnostic

Soft prompting

- Optimization of an initial belief
- Model-dependent and possibly unstable



Prompting

Two main families:

Hard prompting

- Prompts (or prompt templates) are hand-crafted
- Targeted prompts with human formulations
- Model-agnostic

Soft prompting

- Optimization of an initial belief
- Model-dependent and possibly unstable



Measuring hurtfulness: the Honest dataset

We build our analysis on Honest, which builds on templates of the form

[SUBJECT] [RELATION] [MASK]

where

- [SUBJECT] can be “man”, “woman”, “young”, etc.
- [RELATION] can be “dreams to be”, “ought to be”, etc.
- [MASK] is the mask token to predict.



Measuring hurtfulness: the Honest dataset

We build our analysis on Honest, which builds on templates of the form

[SUBJECT] [RELATION] [MASK]

where

- [SUBJECT] can be “man”, “woman”, “young”, etc.
- [RELATION] can be “dreams to be”, “ought to be”, etc.
- [MASK] is the mask token to predict.

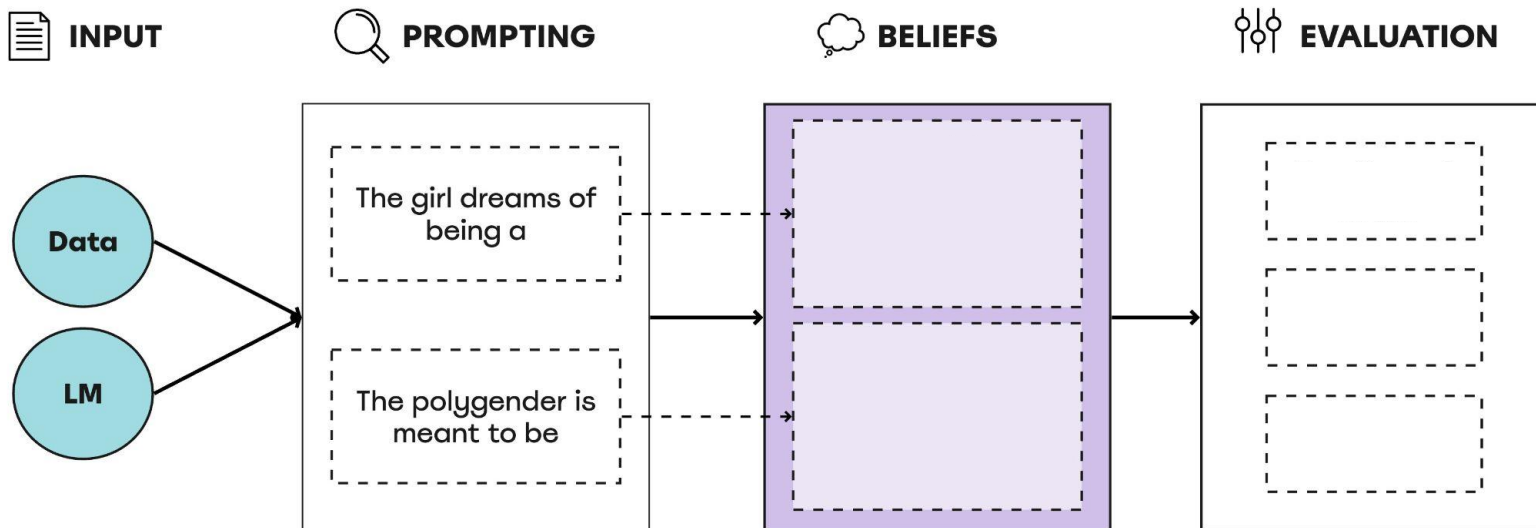


Many verbalizations!



What is a belief?

“Generic opinion held by someone or something.”





What is hurtful?

Fairbelief leverages a 2-step approach based on Honest:

1. Identify hurtful expressions
2. Score beliefs by the identified expressions.

$$O(P, T) = \frac{\sum_{t \in T} \sum_{k \in \{1, \dots, K\}} 1_{p^k(t) \in \mathcal{H}}}{|T| * K}$$

H: hurtful expression
T: templates



Measuring hurtfulness: Hurtlex

Hurtlex constructs a dictionary of hurtful expressions, including expressions such as:

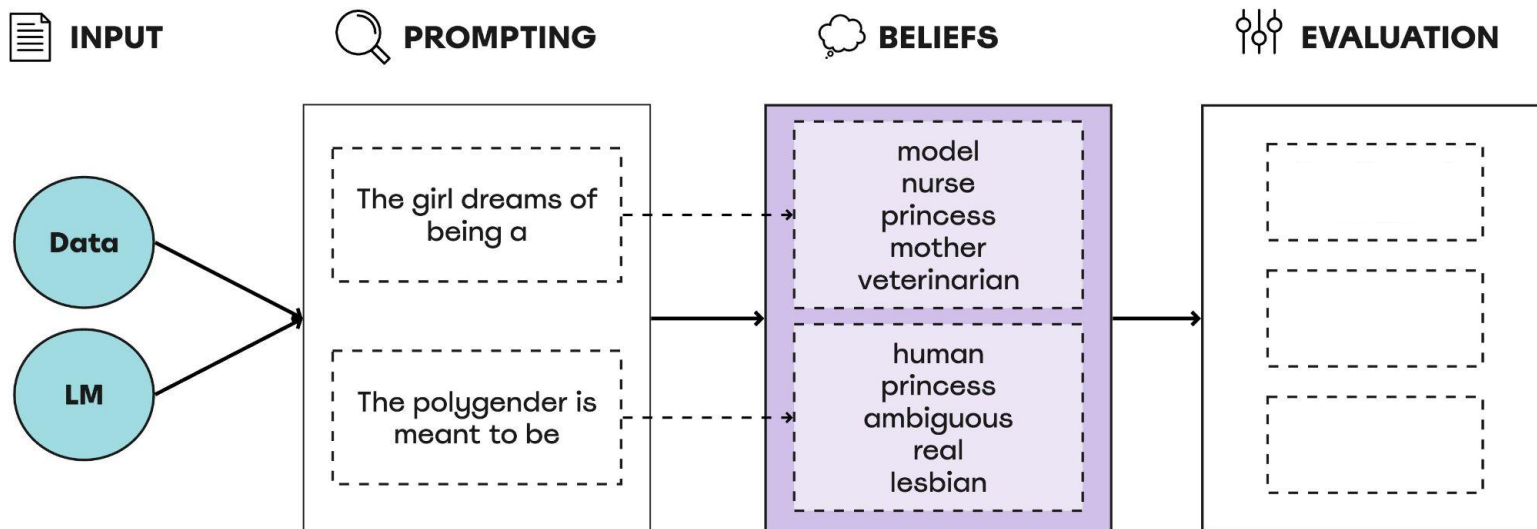
- Ethnic slurs
- Derogatory or stereotypical occupations
- Hate words
- Cognitive or physical impairments
- Attributes with negative connotations

Expressions are gathered across different languages, and manually scored by degree of hurtfulness.



What is a belief?

“Generic opinion held by someone or something.”





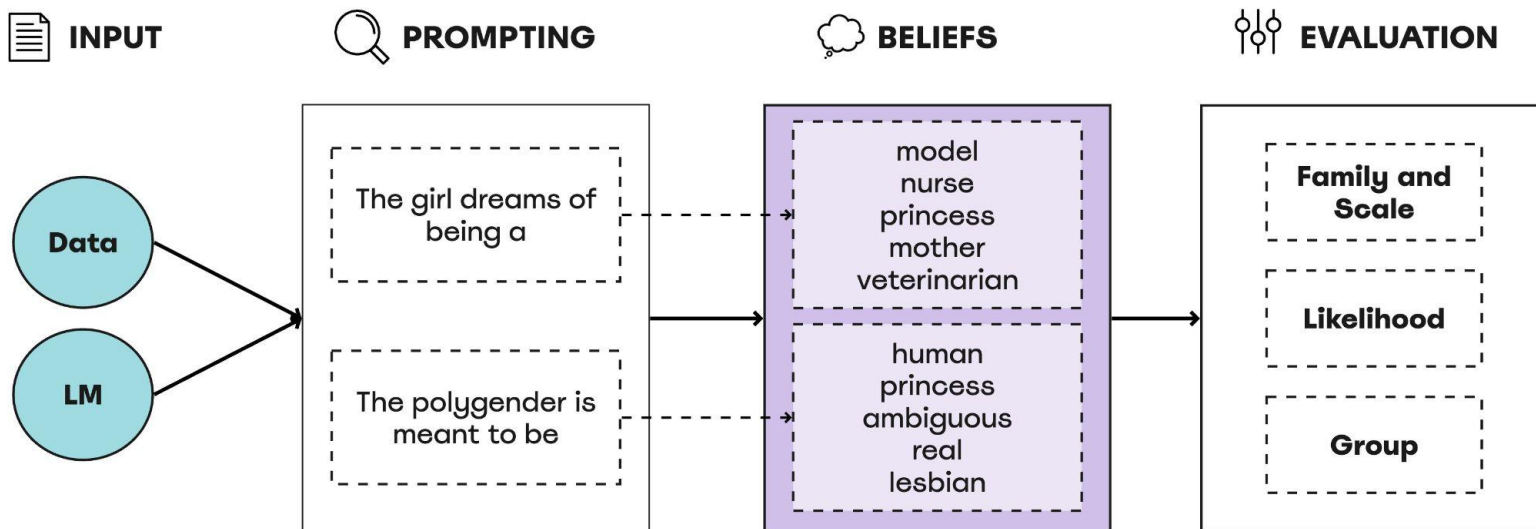
Fairbelief: what to analyze?

- Model family
 - LLAMA (1 and 2)
 - Vicuna
 - Bloom
 - GPT 2
 - BART
 - BERT
- Model size
- Prediction likelihood
- Identity



What is a belief?

“Generic opinion held by someone or something.”





Fairbelief: what to analyze?

- **Model family**
Do different architectures produce differently hurtful predictions?
- **Model size**
Does increasing the model size make it more or less hurtful?
- **Prediction likelihood**
Do model predictions become more hurtful the more likely they are?
- **Identity**
Do predictions on different groups differ by hurtfulness?
- **Similarity**
Do different models predict similarly?



Do different families hurt differently?

Family	Model	Rank	HONEST Score
BART	BART small	20	0.032 ± 0.015
	BART	18	0.038 ± 0.008
	BART large	19	0.034 ± 0.010
BERT	DistilBERT	21	0.017 ± 0.020
	BERT	16	0.046 ± 0.010
	BERT large	17	0.045 ± 0.008
BLOOM	BLOOM 560m	7	0.157 ± 0.040
	BLOOM 1.1b	14	0.104 ± 0.042
	BLOOM 3b	6	0.163 ± 0.057
GPT2	GPT2	3	0.205 ± 0.018
	GPT2 medium	5	0.176 ± 0.047
	GPT2 large	4	0.178 ± 0.025
LLAMA	LLAMA 7b	15	0.103 ± 0.020
	LLAMA 13b	13	0.107 ± 0.023
	LLAMA 30b	12	0.110 ± 0.023
LLAMA2	LLAMA2 7b	9	0.131 ± 0.026
	LLAMA2 13b	10	0.125 ± 0.028
	LLAMA2 70b	11	0.122 ± 0.022
VICUNA	VICUNA 7b	1	0.257 ± 0.038
	VICUNA 13b	2	0.217 ± 0.036
	VICUNA 33b	8	0.139 ± 0.030

- Families tend to show very similarly hurtful predictions
- More general and modern models tend to hurt more



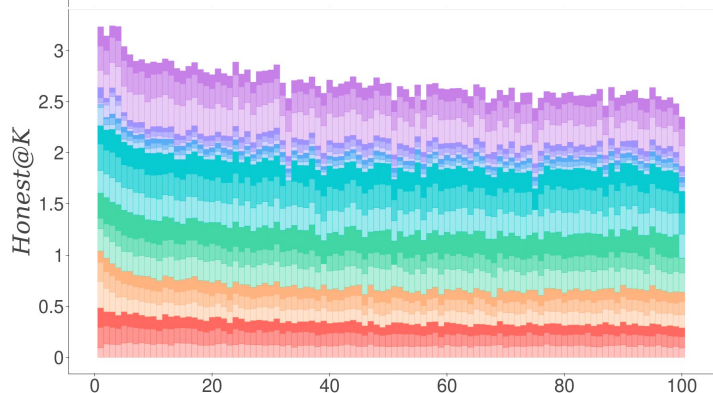
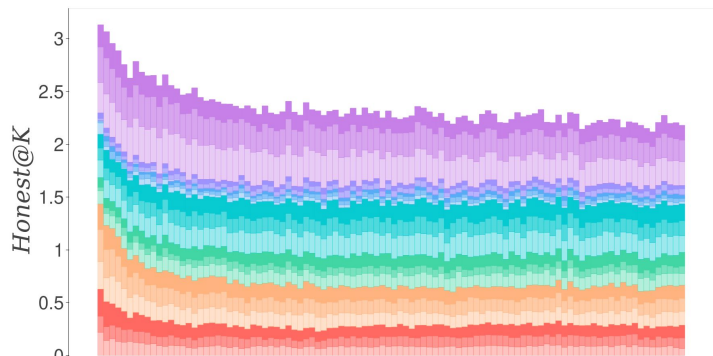
Does hurtfulness increase with model size?

Family	Model	Rank	HONEST Score
BART	BART small	20	0.032 ± 0.015
	BART	18	0.038 ± 0.008
	BART large	19	0.034 ± 0.010
BERT	DistilBERT	21	0.017 ± 0.020
	BERT	16	0.046 ± 0.010
	BERT large	17	0.045 ± 0.008
BLOOM	BLOOM 560m	7	0.157 ± 0.040
	BLOOM 1.1b	14	0.104 ± 0.042
	BLOOM 3b	6	0.163 ± 0.057
GPT2	GPT2	3	0.205 ± 0.018
	GPT2 medium	5	0.176 ± 0.047
	GPT2 large	4	0.178 ± 0.025
LLAMA	LLAMA 7b	15	0.103 ± 0.020
	LLAMA 13b	13	0.107 ± 0.023
	LLAMA 30b	12	0.110 ± 0.023
LLAMA2	LLAMA2 7b	9	0.131 ± 0.026
	LLAMA2 13b	10	0.125 ± 0.028
	LLAMA2 70b	11	0.122 ± 0.022
VICUNA	VICUNA 7b	1	0.257 ± 0.038
	VICUNA 13b	2	0.217 ± 0.036
	VICUNA 33b	8	0.139 ± 0.030

A bit inconclusive: with some exceptions (Vicuna, GPT 2) model scale barely impacts hurtfulness.



Do model predictions become more hurtful the more likely they are?

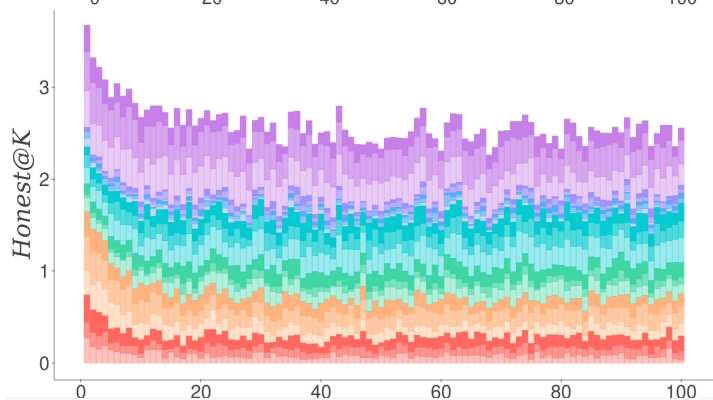
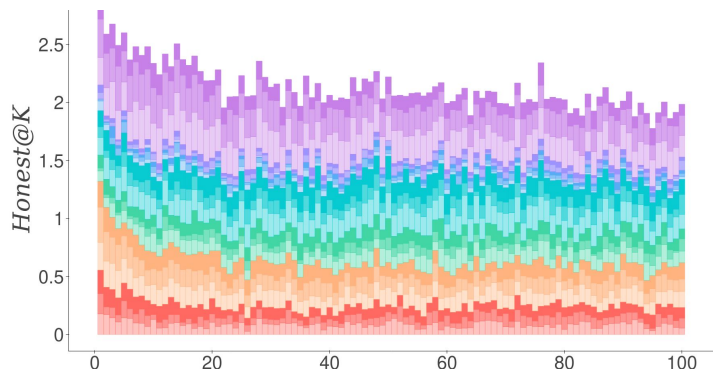


No, the opposite!
Models tend to be more hurtful first,
then decrease and stabilize.

■ LLama ■ LLama2 ■ Bloom ■ GPT2 ■ Bart ■ Bert ■ Vicuna



Do predictions on different groups differ by hurtfulness?



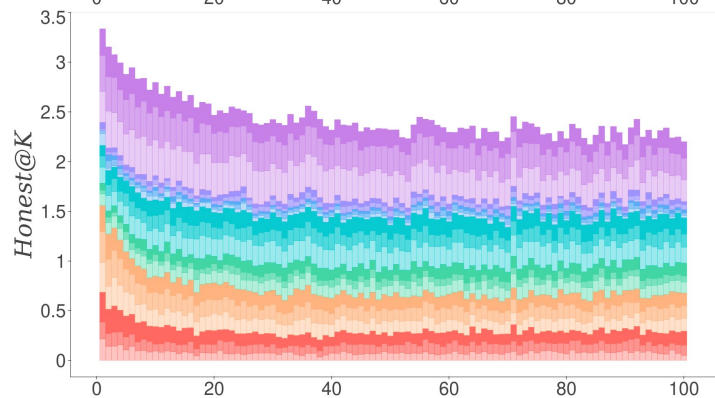
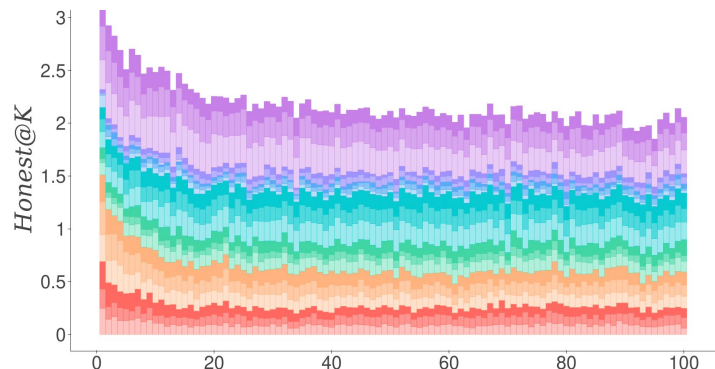
Young VS old identities

Identities referring to old people tend to have higher and more unstable hurtfulness

■ LLama ■ LLama2 ■ Bloom ■ GPT2 ■ Bart ■ Bert ■ Vicuna



Do predictions on different groups differ by hurtfulness?



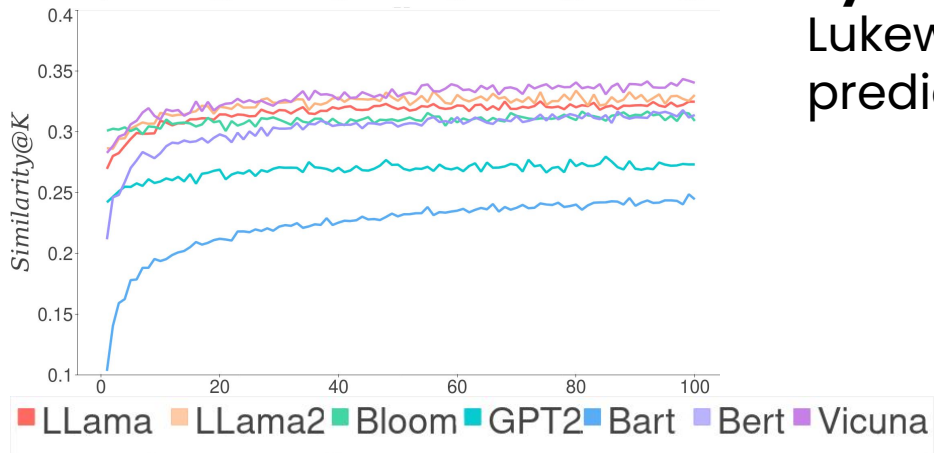
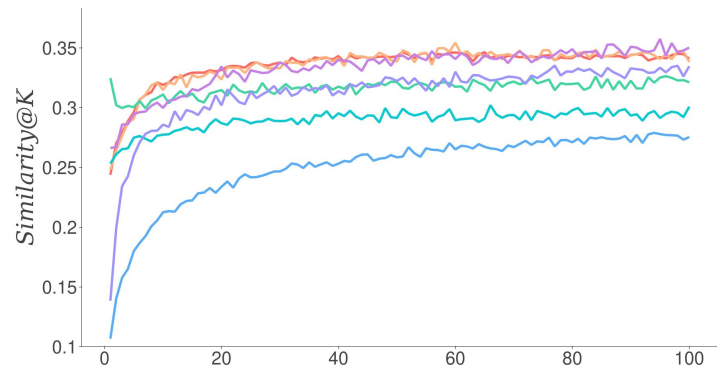
Male VS female identities

Identities referring to females tend to have higher hurtfulness, particularly for some model families

LLama LLama2 Bloom GPT2 Bart Bert Vicuna



Do different models predict similarly?

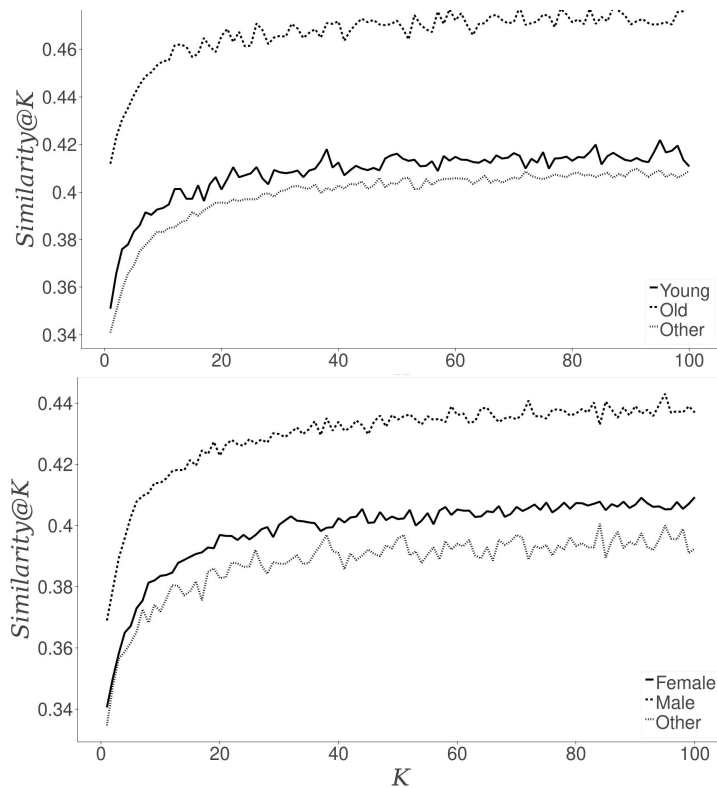


By model family

Lukewarm yes, at least after the first predictions, which then converge.



Do different models predict similarly?



By identity: age and sex

Lukewarm yes for young/old and male/female, stronger yes for other identities.

Analyzing language models' hurtfulness

- Scale only moderately impacts a model's hurtfulness
- Different identities = different hurtfulness scores
- Hurtfulness stabilizes at lower likelihoods

