

Harnessing the richness of the linguistic signal to predict pragmatic inferences

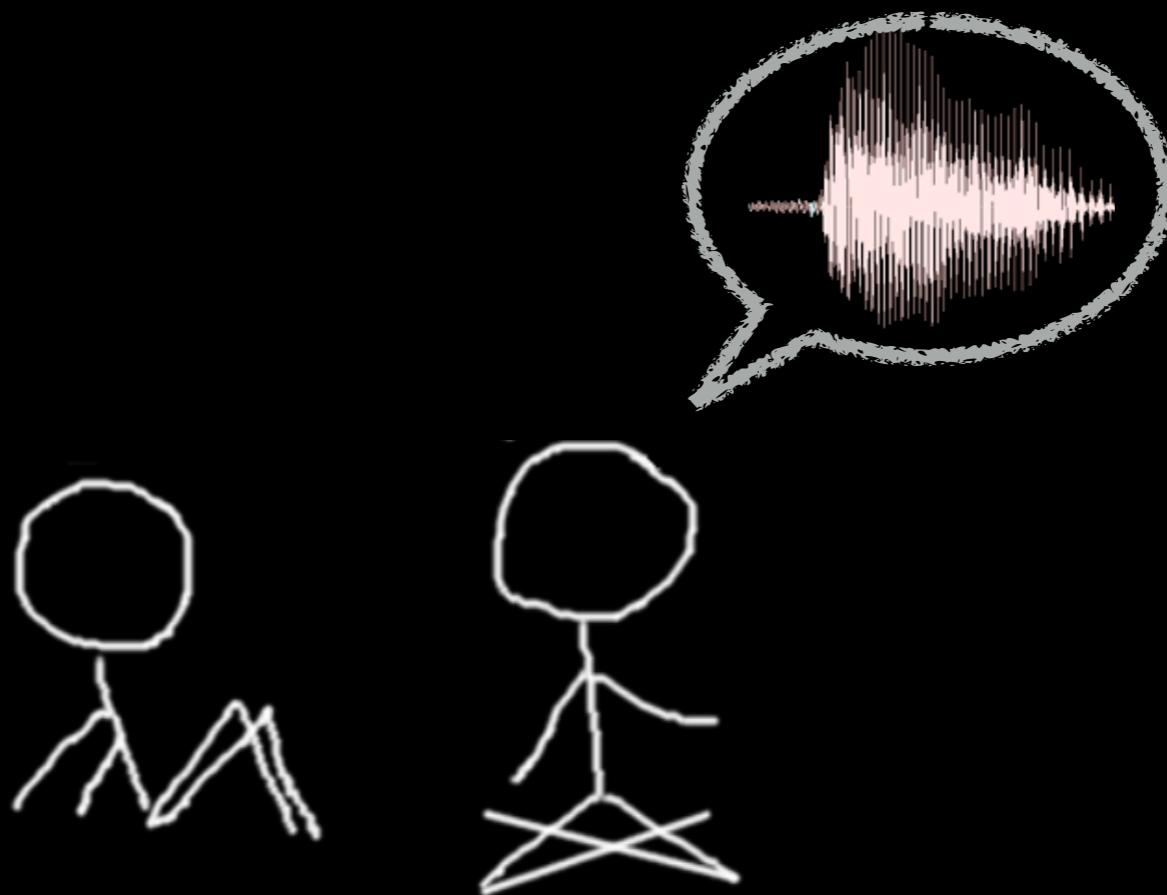
Judith Degen

Oct 31, 2019

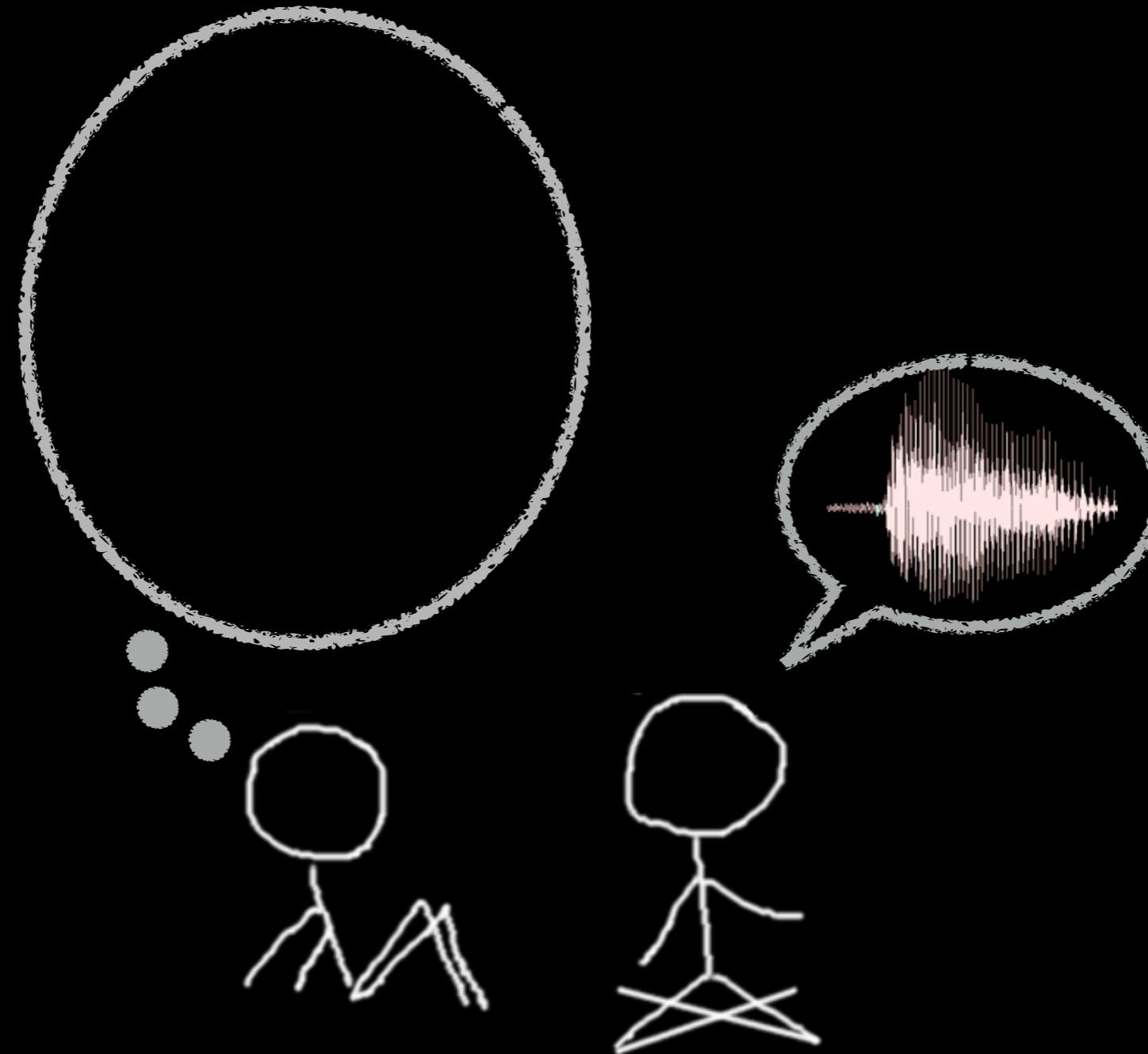
Cornell University





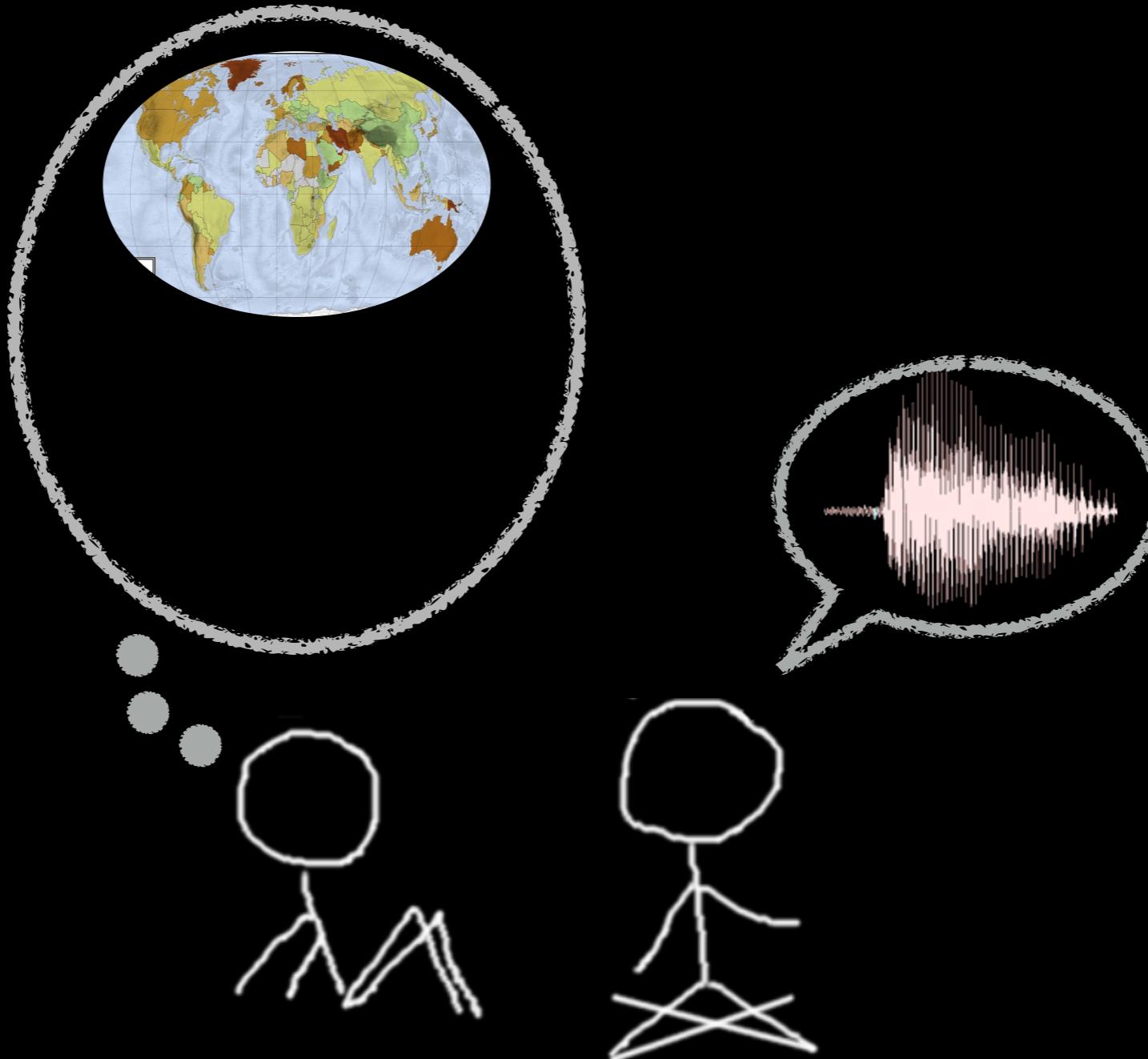


linguistic
signal



linguistic
signal

world
knowledge

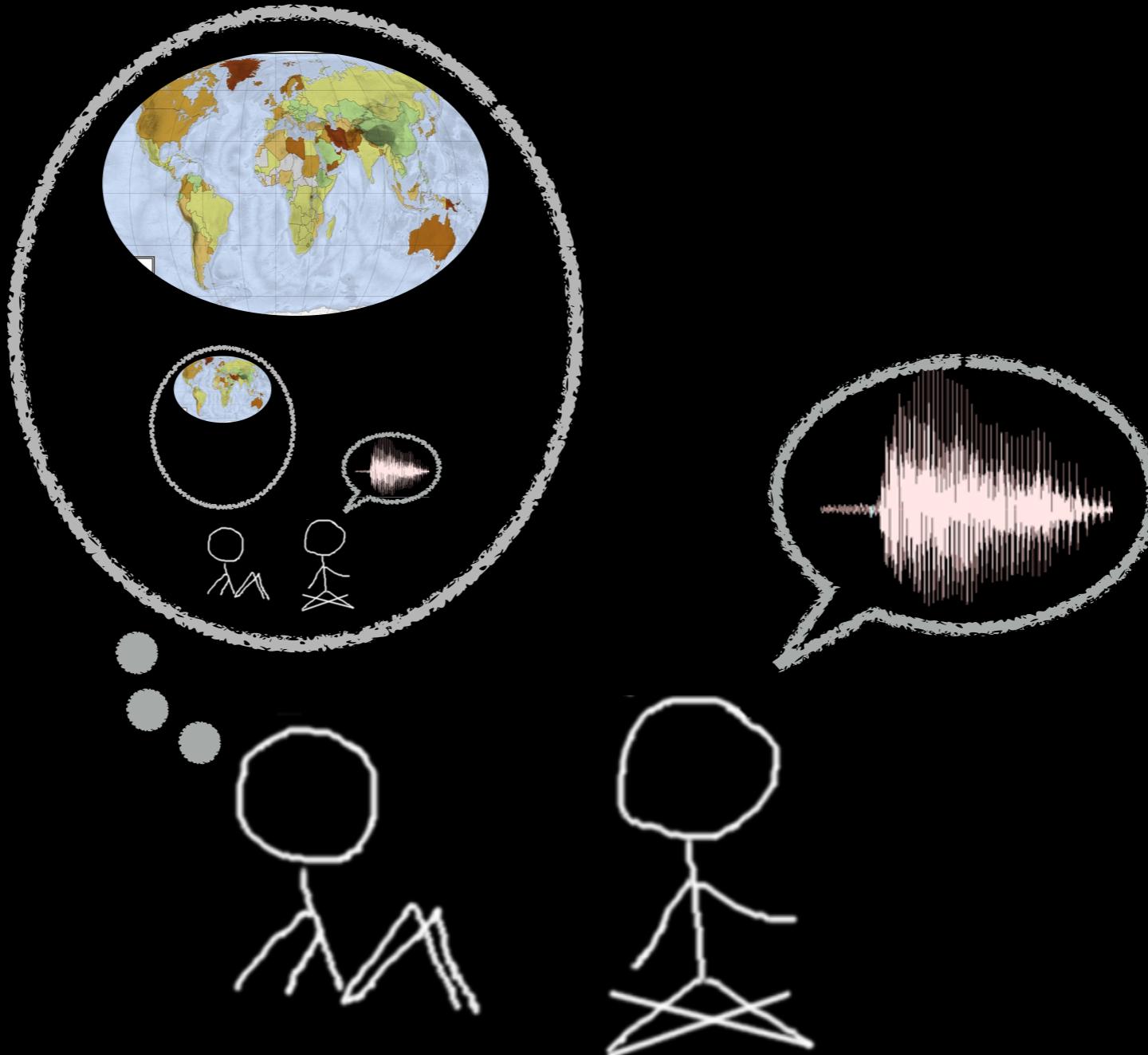


linguistic
signal

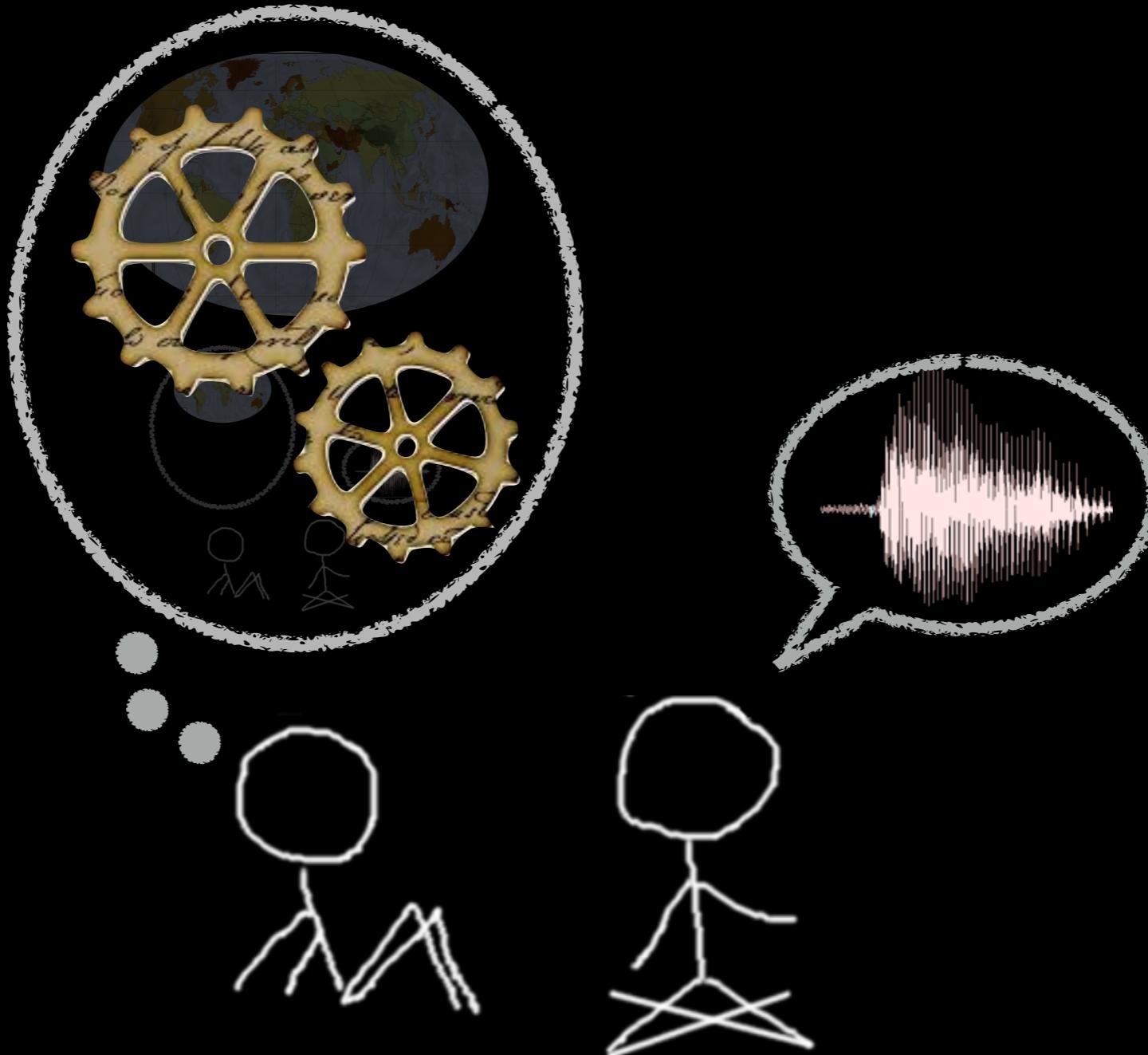
world
knowledge

context

linguistic
signal

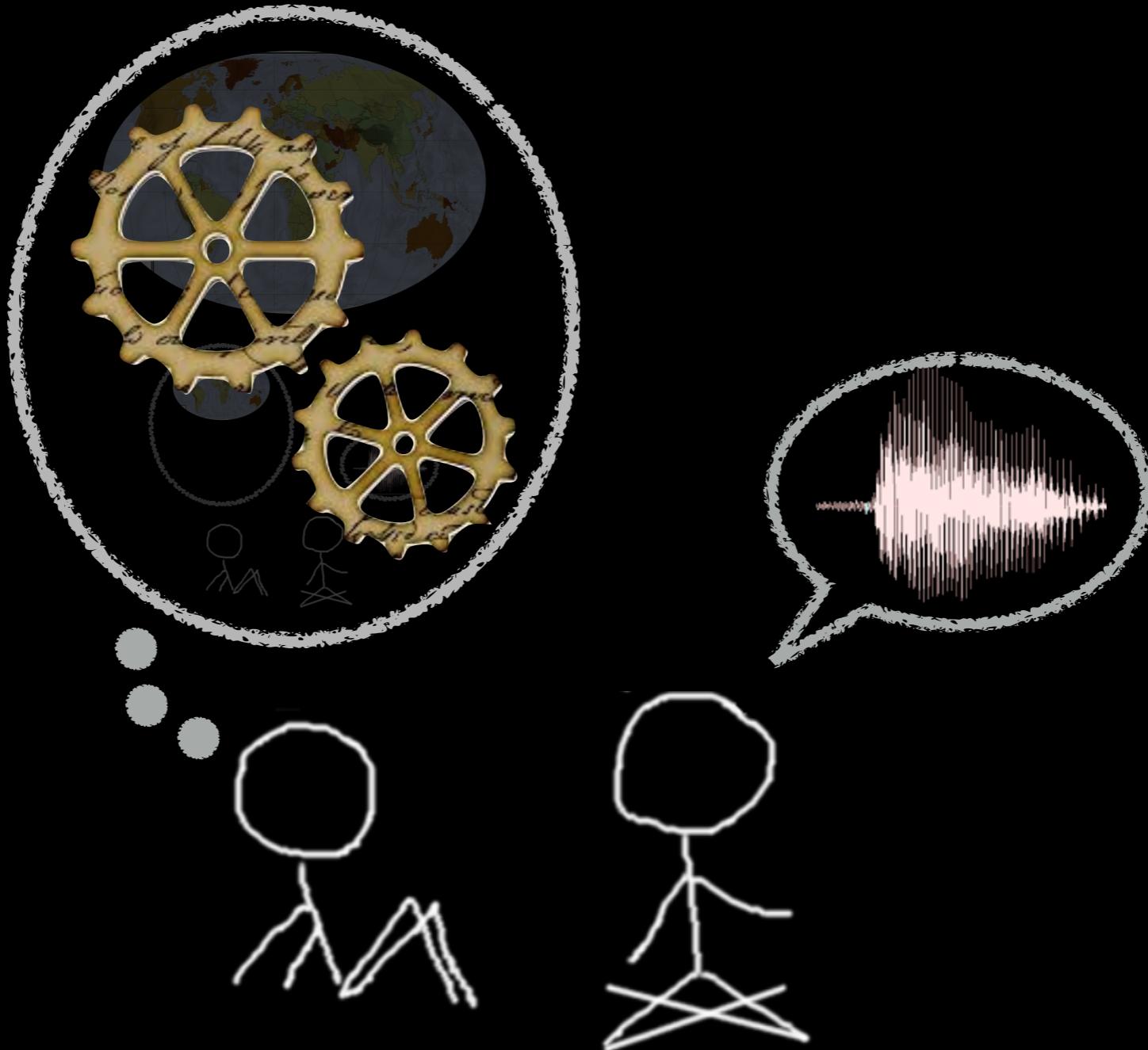


world
knowledge
reasoning
context



linguistic
signal

world
knowledge
reasoning
context



PRAGMATICS

Scalar implicature

(1) John: Was the exam easy?

Mary: Some of the students failed.

Inference: Some, but not **all** of the students failed.

Scalar implicature

(1) John: Was the exam easy?

Mary: Some of the students failed.

Inference: Some, but not **all** of the students failed.

(2) John: Who came to the party?

Mary: Ann or Greg.

Inference: Either Ann **or** Greg came, but not **both**.

Scalar implicature

(1) John: Was the exam easy?

Mary: Some of the students failed.

Inference: Some, but not **all** of the students failed.

(2) John: Who came to the party?

Mary: Ann or Greg.

Inference: Either Ann **or** Greg came, but not **both**.

(3) John: How was your date?

Mary: It was OK.

Inference: The date was **OK**, but not **great**.

Why study scalar implicature?

(1) John: Was the exam easy?

Mary: Some of the students failed.

Inference: Some, but not **all** of the students failed.

Inference: The exam was not easy.

Why study scalar implicature?

(1) John: Was the exam easy?

Mary: Some of the students failed.

Inference: Some, but not **all** of the students failed.

Inference: The exam was not easy.

(1) John: Is the teacher doing a good job?

Mary: Some of the students failed.

Inference: Some, but not **all** of the students failed.

Inference: ~~The exam was hard.~~

Inference: The teacher isn't doing a good job.

Why study scalar implicature?

(1) John: Was the exam easy?

Mary: Some of the students failed.

Inference: Some, but not **all** of the students failed.

Inference: The exam was not easy.

ROBUST

(1) John: Is the teacher doing a good job?

Mary: Some of the students failed.

Inference: Some, but not **all** of the students failed.

Inference: ~~The exam was hard.~~

Inference: The teacher isn't doing a good job.

Accounts of scalar implicature

Accounts of scalar implicature

The default account

Levinson 2000

Basic assumptions:

- context is hard to integrate

Solution: two types of inferences

- fast, automatic, context-independent inferences

Generalized Conversational Implicature

- slow, effortful, context-dependent inferences

Particularized Conversational Implicature

Accounts of scalar implicature

The default account

Levinson 2000

Basic assumptions:

- context is hard to integrate

Solution: two types of inferences

- fast, automatic, context-independent inferences

Generalized Conversational Implicature

- slow, effortful, context-dependent inferences

Particularized Conversational Implicature

A contextualist account

Degen & Tanenhaus 2015

Basic assumptions:

- context is easy to integrate

Solution: efficient use of context

- listeners acquire a context-dependent speaker model:
 $P(\text{utterance} \mid \text{context}, \text{meaning})$
- listeners use available contextual cues to infer speaker meaning:
 $P(\text{meaning} \mid \text{utterance}, \text{context})$

Accounts of scalar implicature

The default account

Levinson 2000

Basic assumptions:

- context is hard to integrate

Solution: two types of inferences

- fast, automatic, context-independent inferences

Generalized Conversational Implicature

- slow, effortful, context-dependent inferences

Particularized Conversational Implicature

A contextualist account

Degen & Tanenhaus 2015

Basic assumptions:

- context is easy to integrate

Solution: efficient use of context

- listeners acquire a context-dependent speaker model:
 $P(\text{utterance} \mid \text{context}, \text{meaning})$
- listeners use available contextual cues to infer speaker meaning:
 $P(\text{meaning} \mid \text{utterance}, \text{context})$

Accounts of scalar implicature

??CONTEXT??

The default account

Levinson 2000

Basic assumptions:

- context is hard to integrate

Solution: two types of inferences

- fast, automatic, context-independent inferences

Generalized Conversational Implicature

- slow, effortful, context-dependent inferences

Particularized Conversational Implicature

A contextualist account

Degen & Tanenhaus 2015

Basic assumptions:

- context is easy to integrate

Solution: efficient use of context

- listeners acquire a context-dependent speaker model:
 $P(\text{utterance} \mid \text{context}, \text{meaning})$
- listeners use available contextual cues to infer speaker meaning:
 $P(\text{meaning} \mid \text{utterance}, \text{context})$

Sources of data in experimental pragmatics

- historically: introspective judgments
- judgment data from controlled experiments
- processing data from controlled experiments

Variability in scalar implicature

attributed to

- properties of the scale van Tiel et al 2016
- stress on cognitive system de Neys & Schaeken 2007
- idiosyncratic properties of participants
- context for a review: Degen & Tanenhaus 2019

What's lacking

- a clear picture of the naturalistic contexts that speakers produce scalar expressions in
- a clear picture of whether listeners make use of the contextual information available to them in naturalistic contexts

Overview

1. A study combining corpus analysis & web-based experiments on “some”
2. Using distributed meaning representations to predict human inference ratings

There is much more variability in scalar inferences than commonly assumed — but it's systematically context-dependent, and we can capture a lot of it by investigating the naturalistic signal

Case study: “some”

Scalar implicatures in the wild

Degen 2015

1. I like **some country music**.
2. It would certainly help them to appreciate **some of the things we have here**.
3. You sound like you have **some small ones** in the background.

Scalar implicatures in the wild

Degen 2015

1. I like **some country music**.

Inference? I like some, but not all, country music

2. It would certainly help them to appreciate **some of the things we have here**.

Inference? ...to appreciate some, but not all...

3. You sound like you have **some small ones** in the background.

Inference? ... some, but not all small ones...

Combining corpora & the web

1. extracted all 1390 utterances containing *some* from the Switchboard corpus of spoken American English
2. collected inference strength ratings for each item on Mechanical Turk (10 judgments per item)

Speaker A: i mean, they just have beautiful, beautiful homes and they have everything. the kids only wear name brand things to school and it's one of these things,

Speaker B: oh me. well that makes it hard for you, doesn't it.

Speaker A: well it does, you know. it really does because i'm a single mom and i have a thirteen year old now and uh, you know, it does.

Speaker B: oh, me.

Speaker A: i mean, we do it to a point but uh, not to where she feels different ,

Speaker B: yeah.

Speaker A:

but some of them are very rich

but **some, but not all** of them are very rich

How similar is the statement with 'some, but not all' (green) to the statement with 'some' (red)?

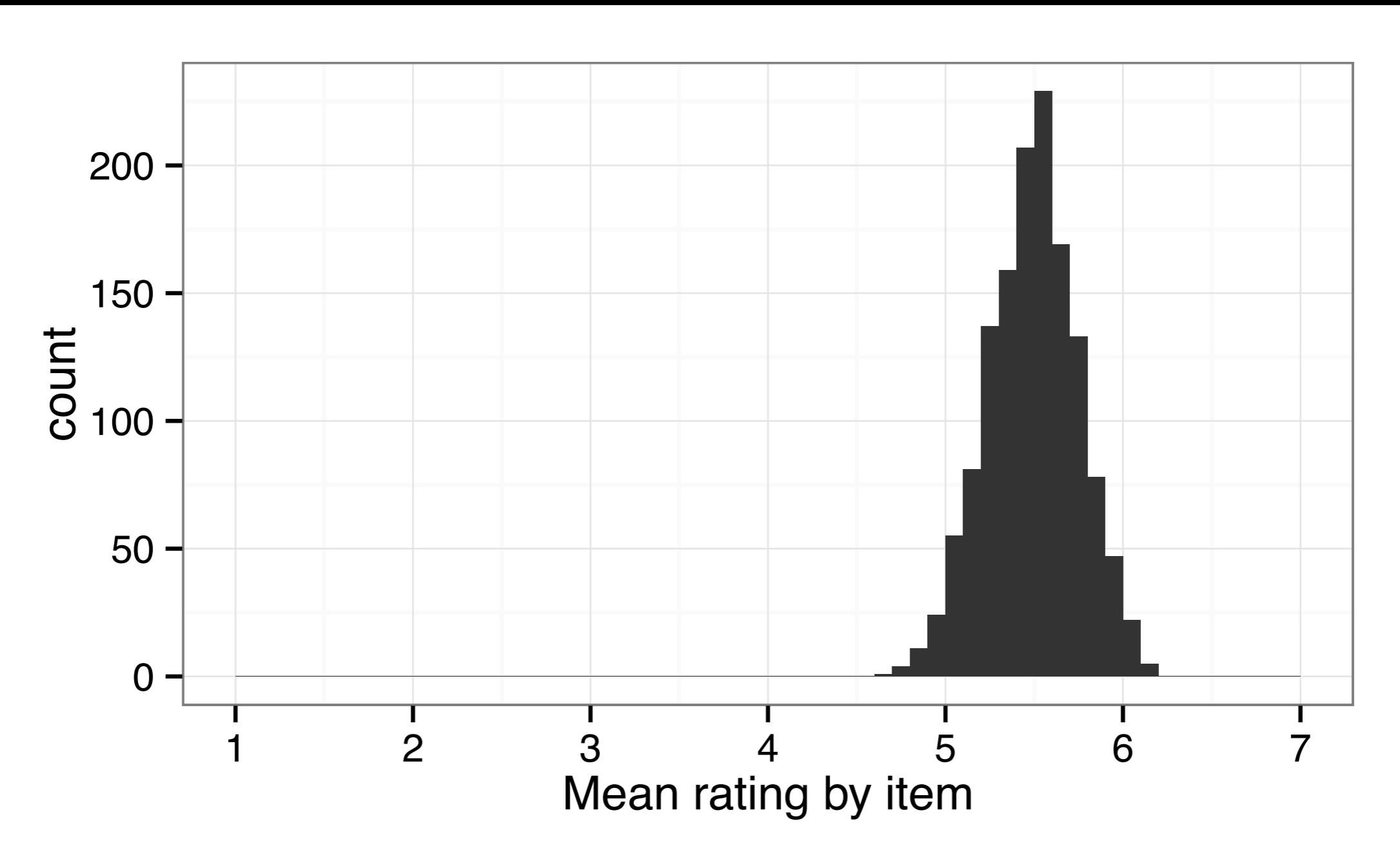
Very different meaning

Same meaning

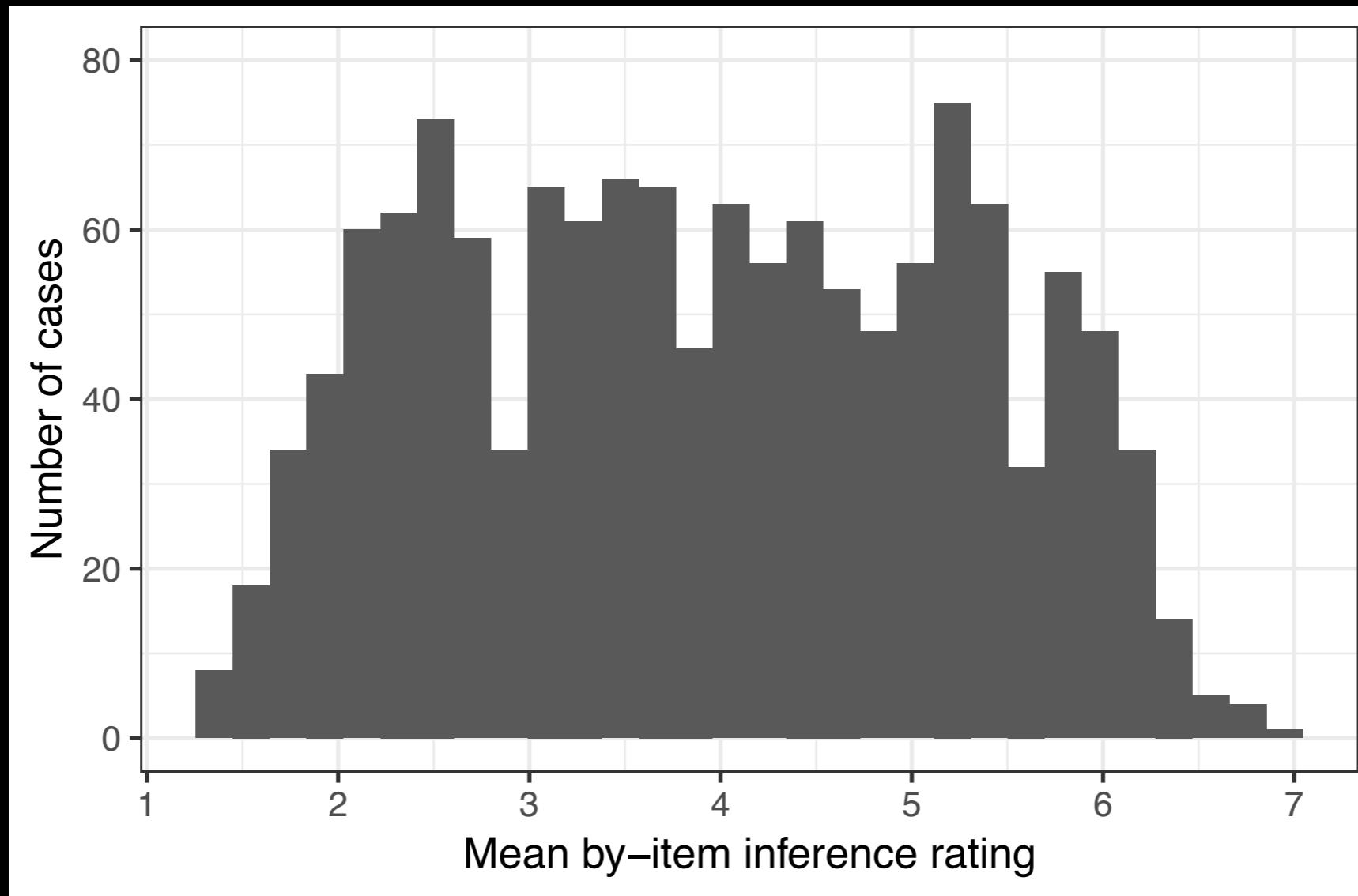
1 2 3 4 5 6 7

Continue

Default prediction



Variability in inference strength



large amount of variability in inference strength

Just noise?

Qualitative investigation

1. I like **some country music**.

6.9

2. It would certainly help them to appreciate **some of the things we have here**.

4

3. You sound like you have **some small ones** in the background.

1.5

Qualitative investigation

1. I like **some country music**.

6.9 Inference? I like some, but not all, country music

2. It would certainly help them to appreciate **some of the things we have here**.

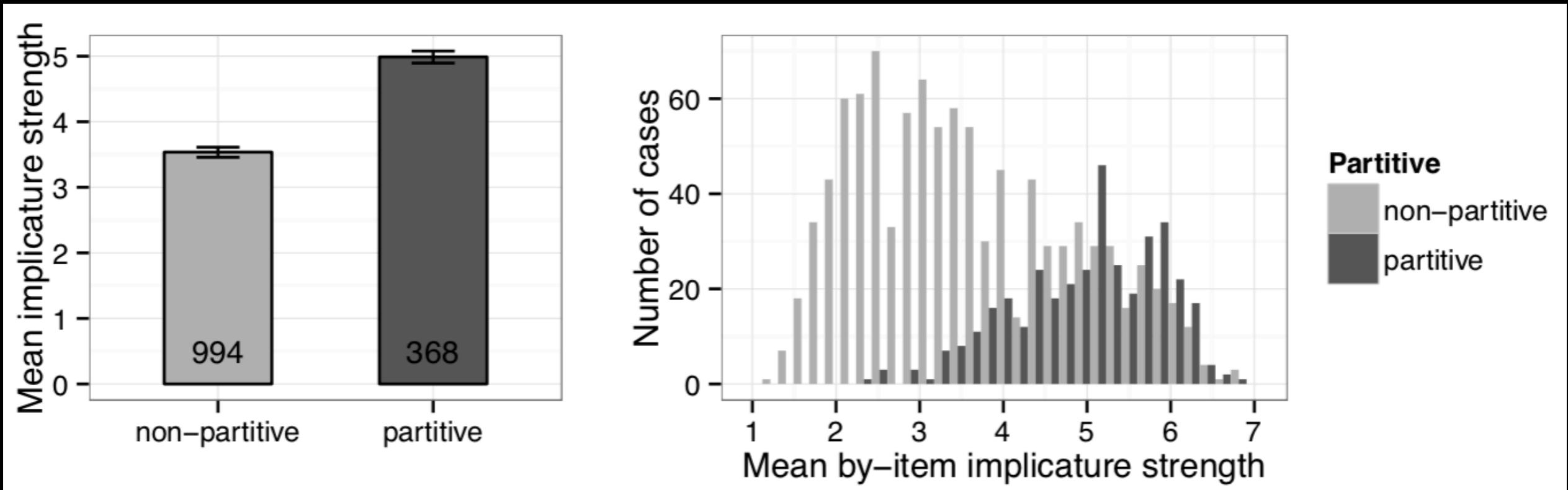
4 Inference? ...to appreciate some, but not all...

3. You sound like you have **some small ones** in the background.

1.5 Inference? ... some, but not all small ones...

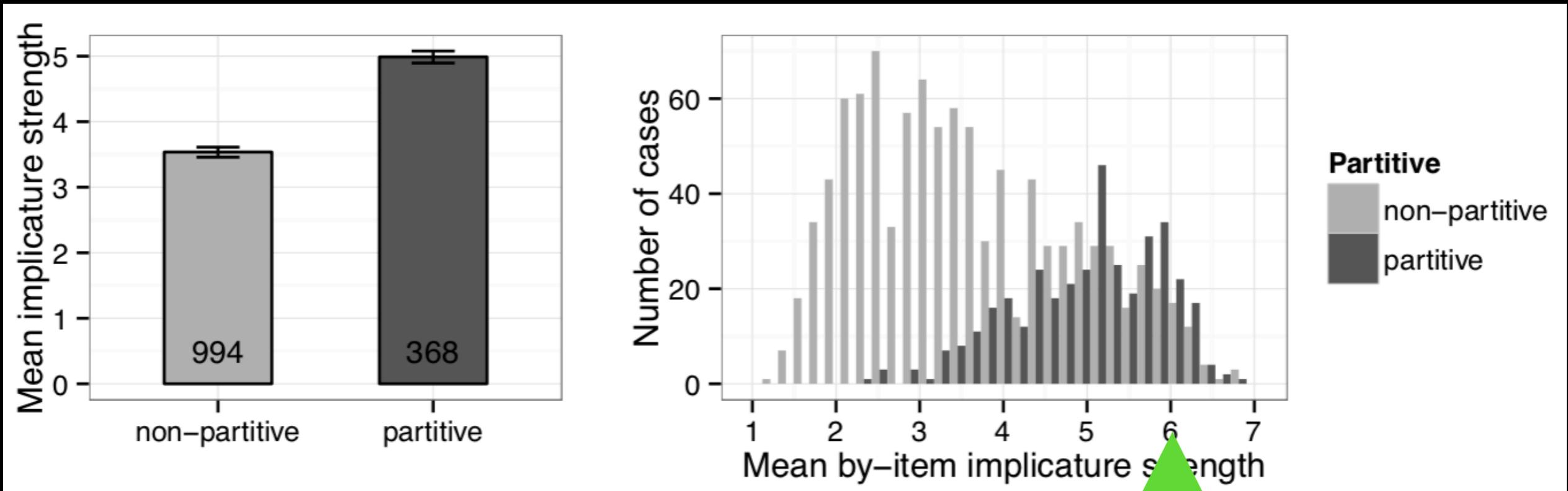
Stronger inferences...

...with **partitive** *some-NPs*.



Stronger inferences...

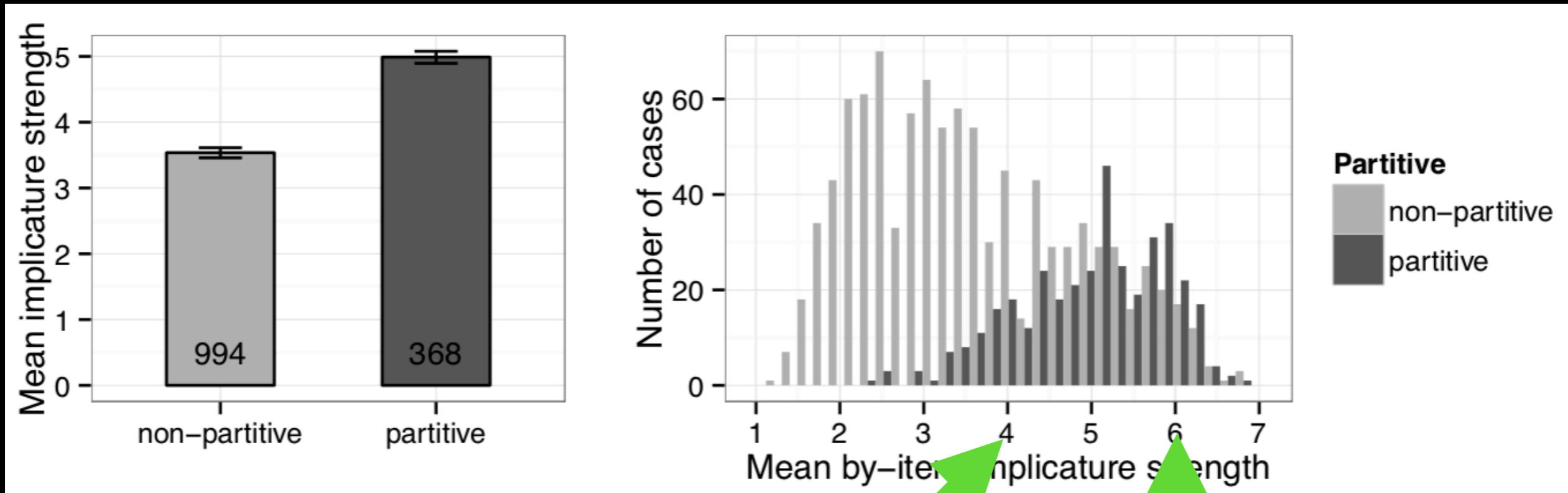
...with **partitive** *some-NPs*.



I've seen ***some of them*** on repeats

Stronger inferences...

...with **partitive** *some-NPs*.

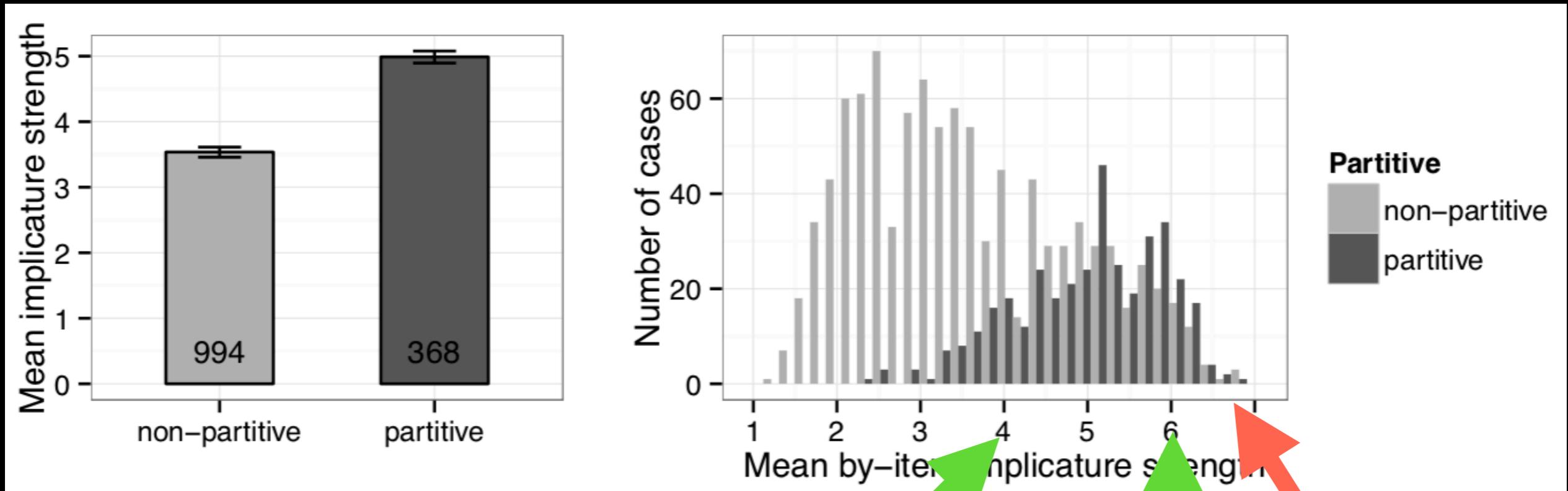


I've seen **some of them** on repeats

*It would certainly help them to appreciate
some of the things we have here.*

Stronger inferences...

...with **partitive** *some-NPs*.



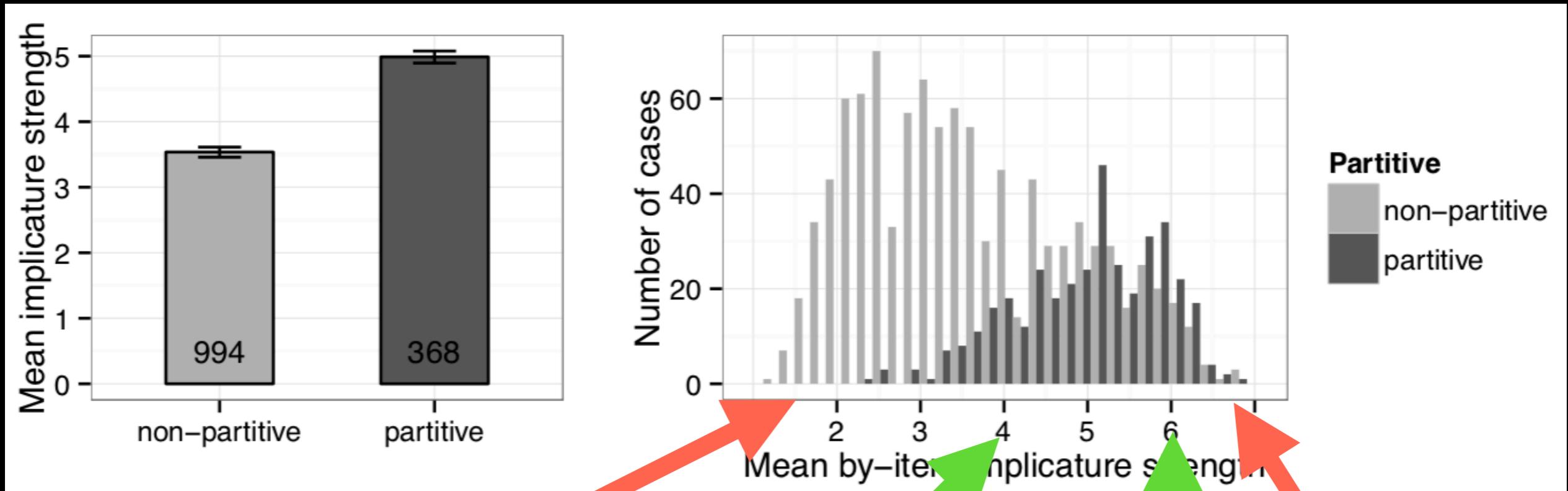
I've seen **some of them** on repeats

*It would certainly help them to appreciate
some of the things we have here.*

*I like **some country music**.*

Stronger inferences...

...with **partitive** *some-NPs*.



I've seen **some of them** on repeats

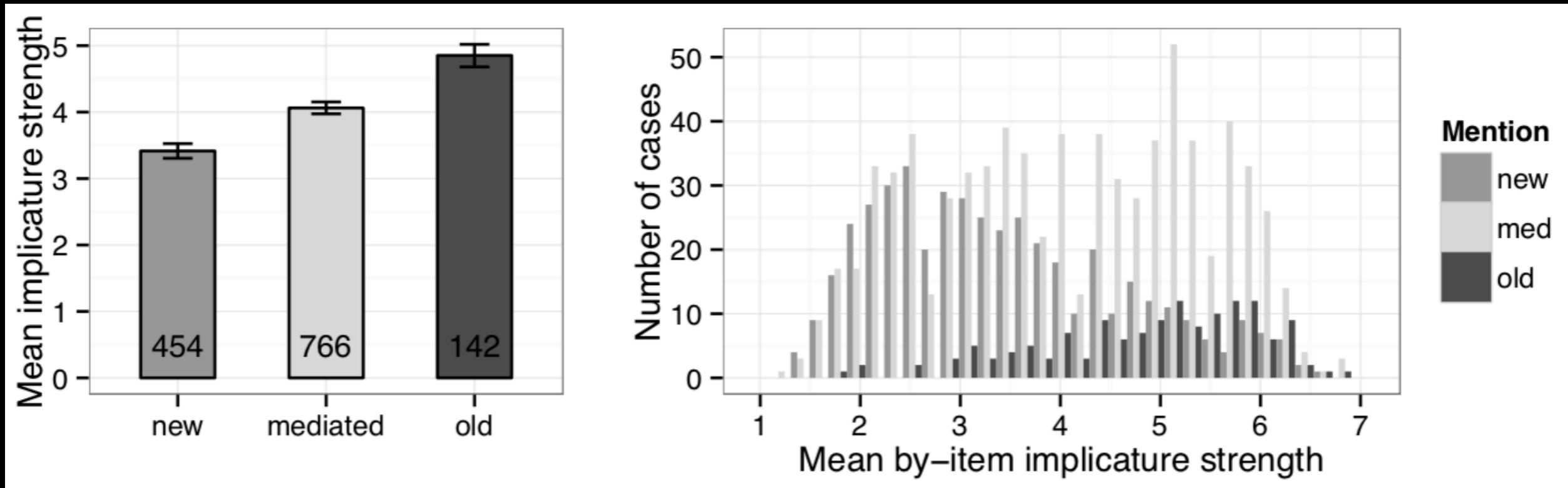
*It would certainly help them to appreciate
some of the things we have here.*

*You sound like you have **some small ones** in the background.*

*I like **some country music**.*

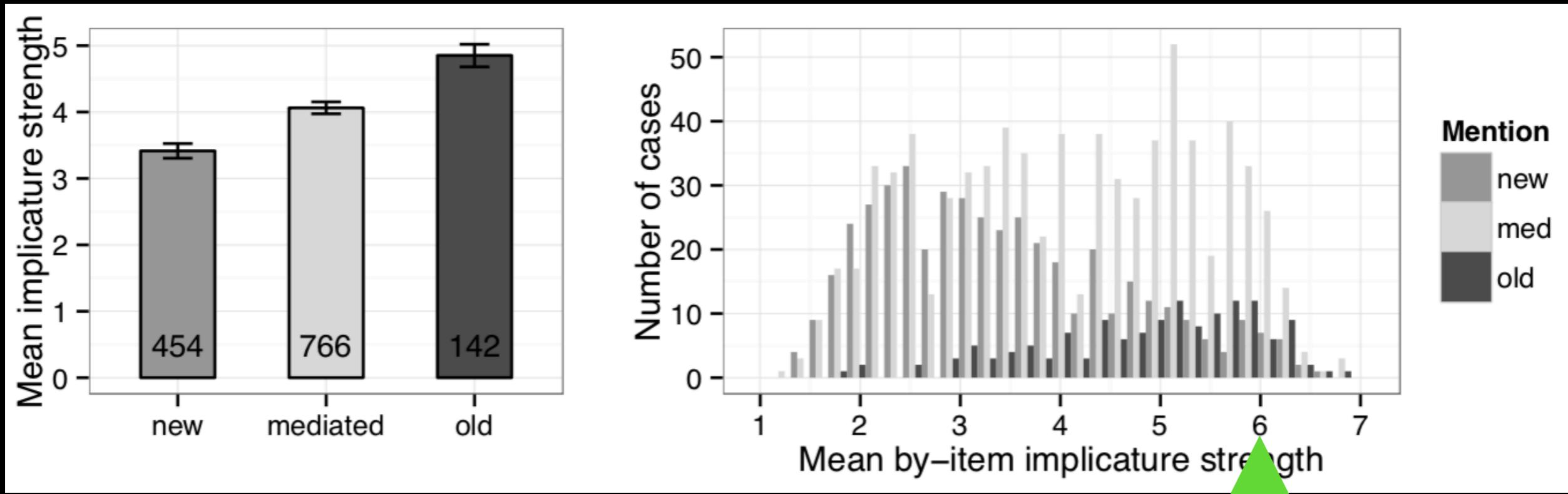
Stronger inferences...

...with **previously mentioned** NP referents.



Stronger inferences...

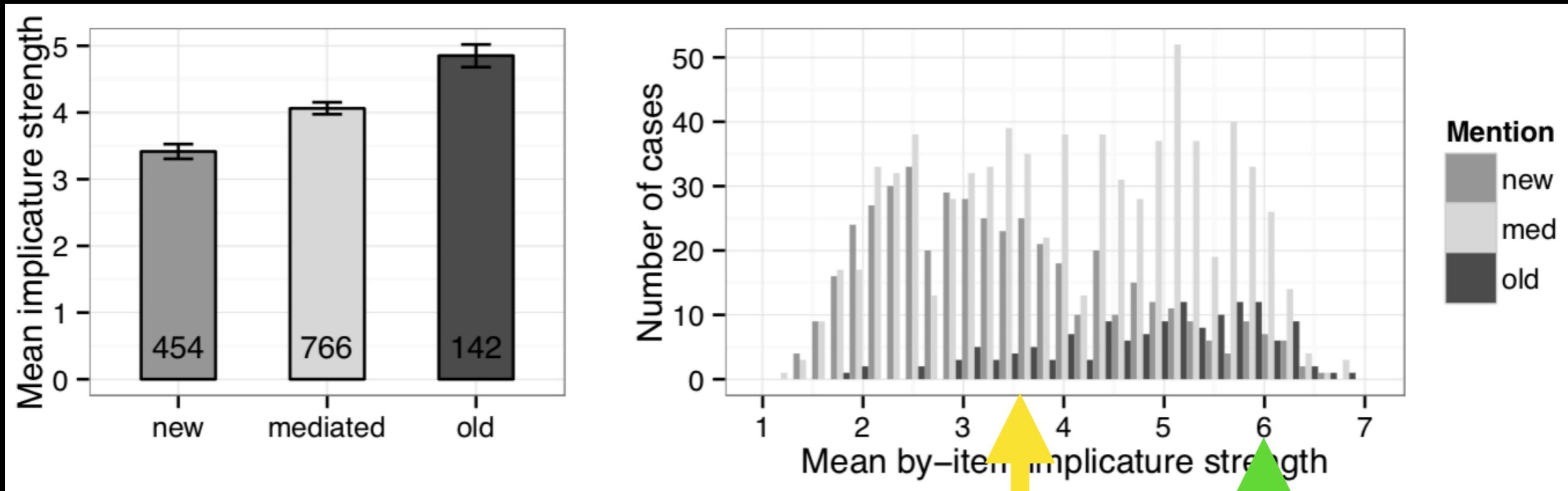
...with **previously mentioned** NP referents.



I've seen **some of them** on repeats

Stronger inferences...

...with **previously mentioned** NP referents.



I've seen **some of them** on repeats

We've got **some beets**.

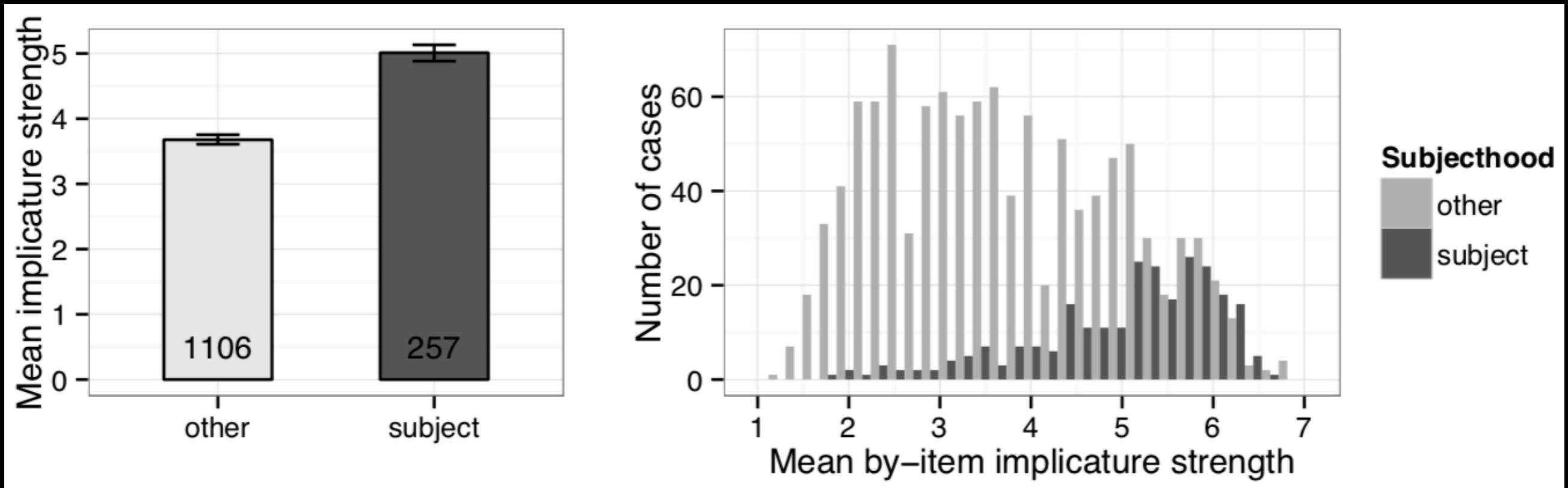
Stronger inferences...

...with **previously mentioned** NP referents.



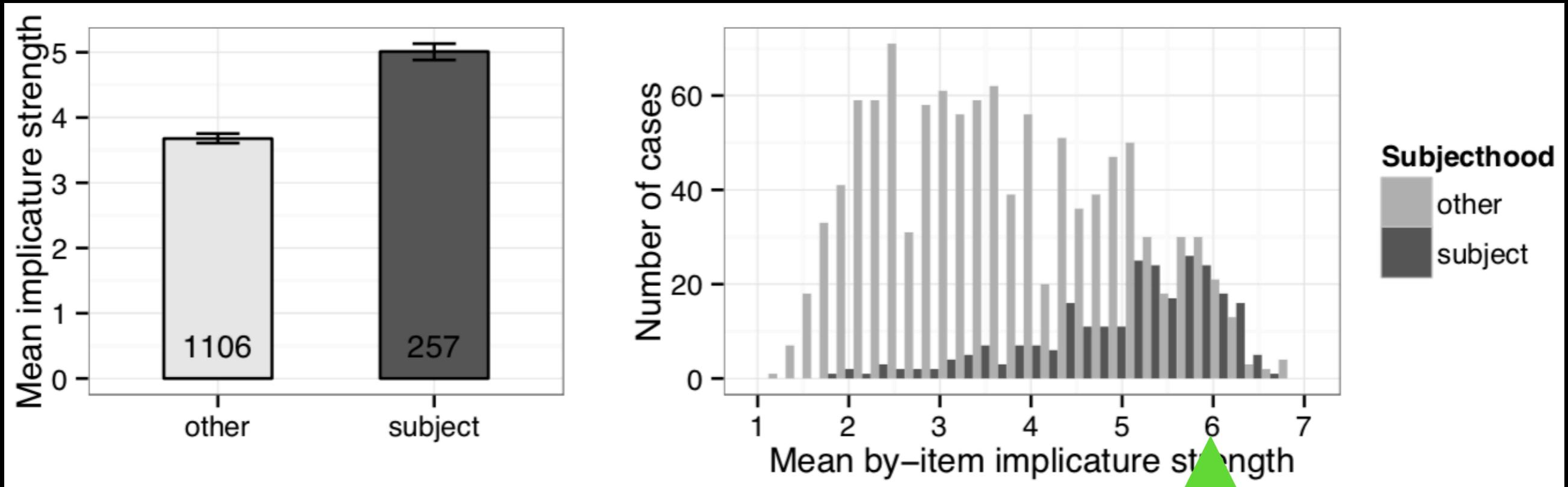
Stronger inferences...

...with *some*-NPs in **subject** position.



Stronger inferences...

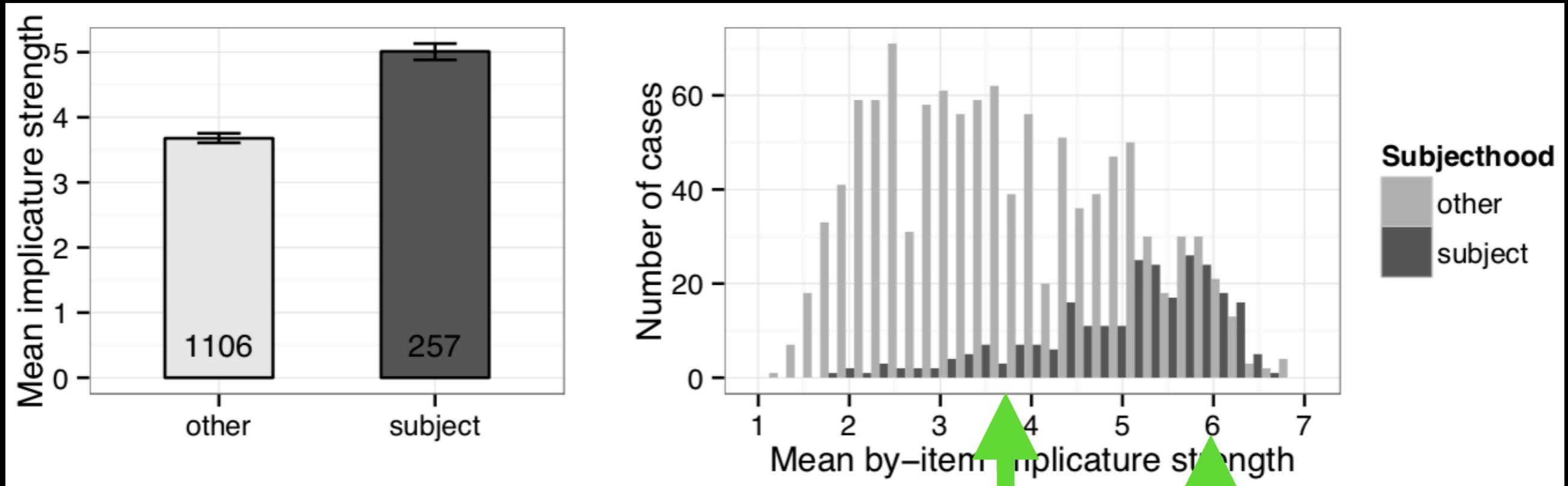
...with *some*-NPs in **subject** position.



Some kids are really having it.

Stronger inferences...

...with *some*-NPs in **subject** position.

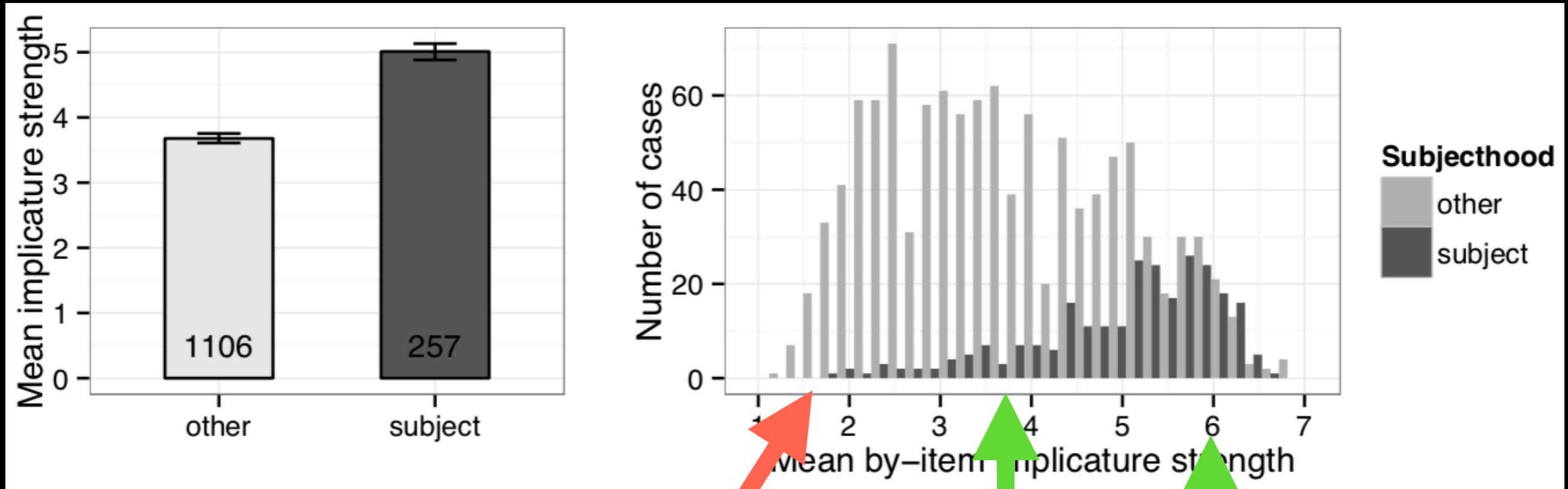


Some kids are really having it.

Occasionally, *some ice skating* will come on.

Stronger inferences...

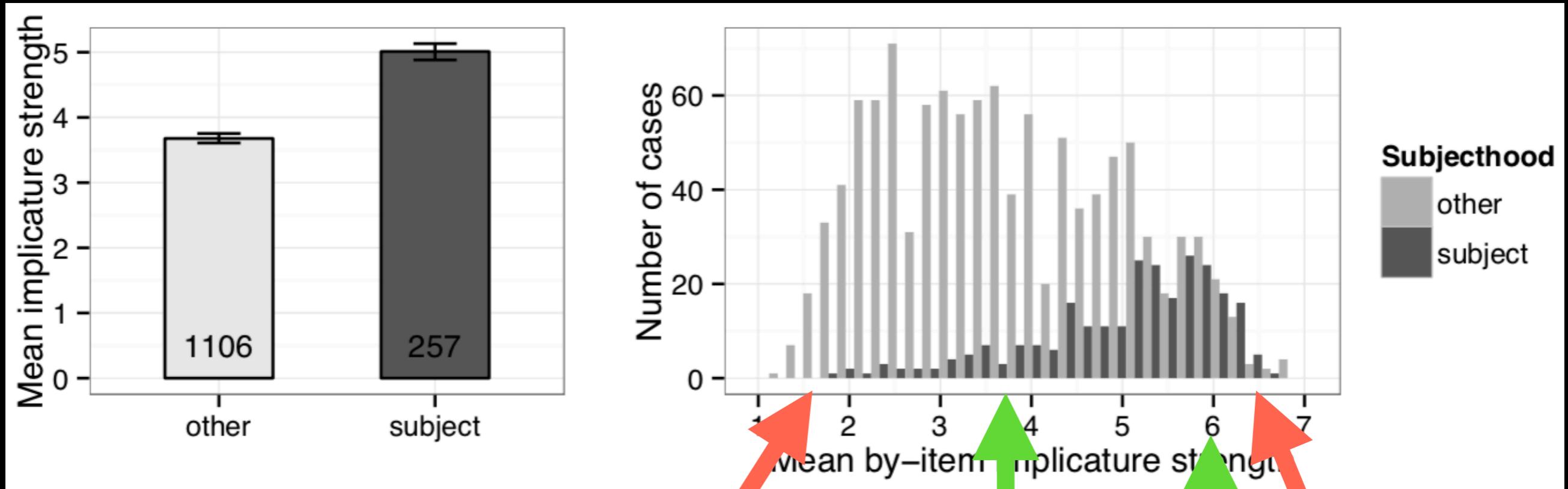
...with *some*-NPs in **subject** position.



Some kids are really having it.
Occasionally, *some ice skating* will come on.
That would take *some planning*.

Stronger inferences...

...with *some*-NPs in **subject** position.



Some kids are really having it.

Occasionally, *some ice skating* will come on.

That would take some planning.

I like some country music.

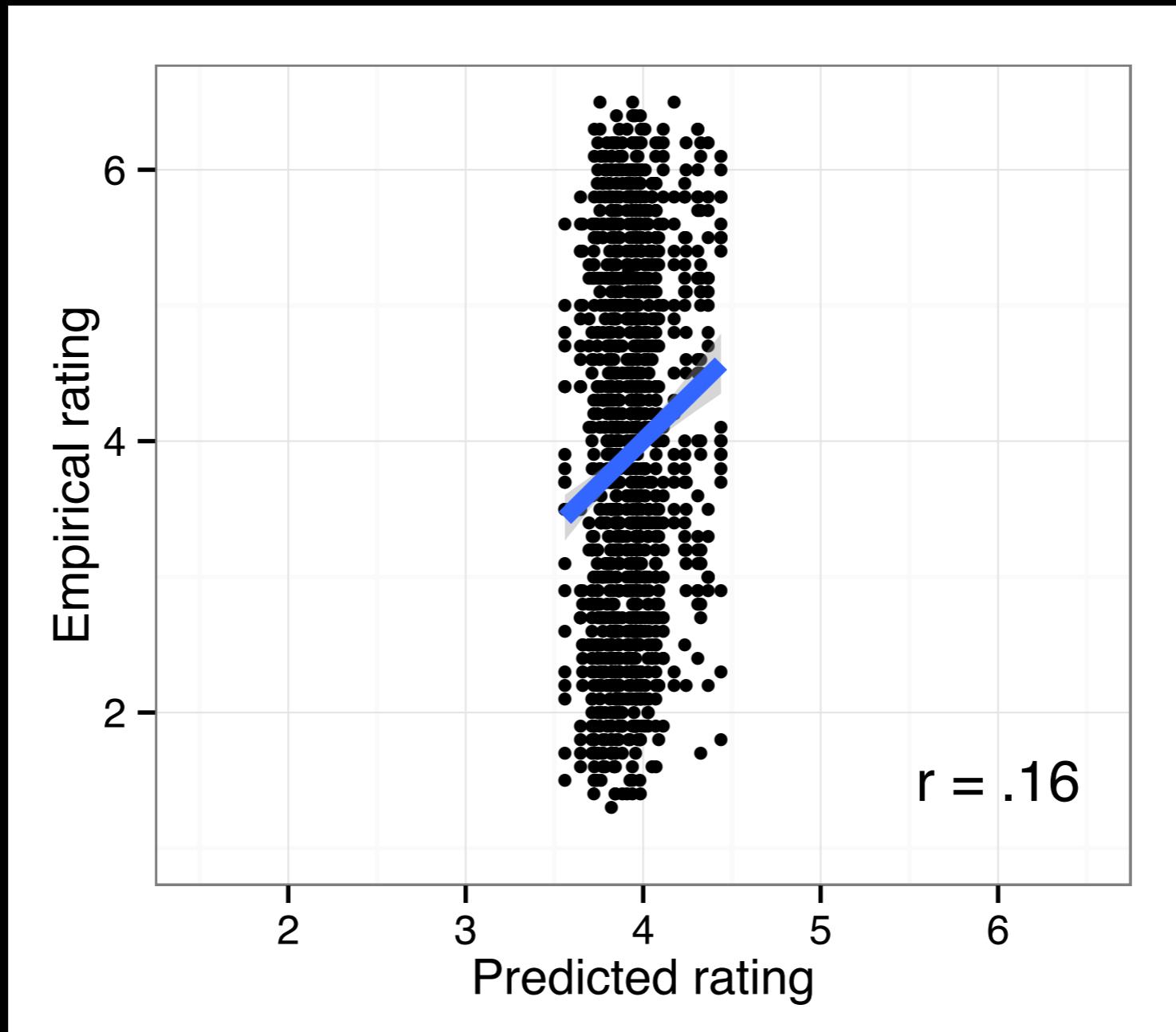
	Coef β	SE(β)	t	p
Intercept	4.01	0.06	68.7	<.0001
Partitive	0.91	0.09	9.6	<.0001
Strength	-0.50	0.05	-9.5	<.0001
Linguistic mention	0.31	0.07	4.4	<.0001
Subjecthood	0.41	0.10	4.2	<.0001
Modification	0.12	0.06	2.0	<.05
Sentence length	0.15	0.05	3.2	<.01
Partitive:Strength	0.39	0.10	4.1	<.0001
Linguistic mention:Subjecthood	0.17	0.21	0.8	<.44
Linguistic mention:Modification	0.34	0.13	2.6	<.01
Subjecthood:Modification	0.27	0.17	1.6	<.12
Linguistic mention:Subjecthood:Modification	0.61	0.42	1.4	<.16

Table 5 Model coefficients for the full model.

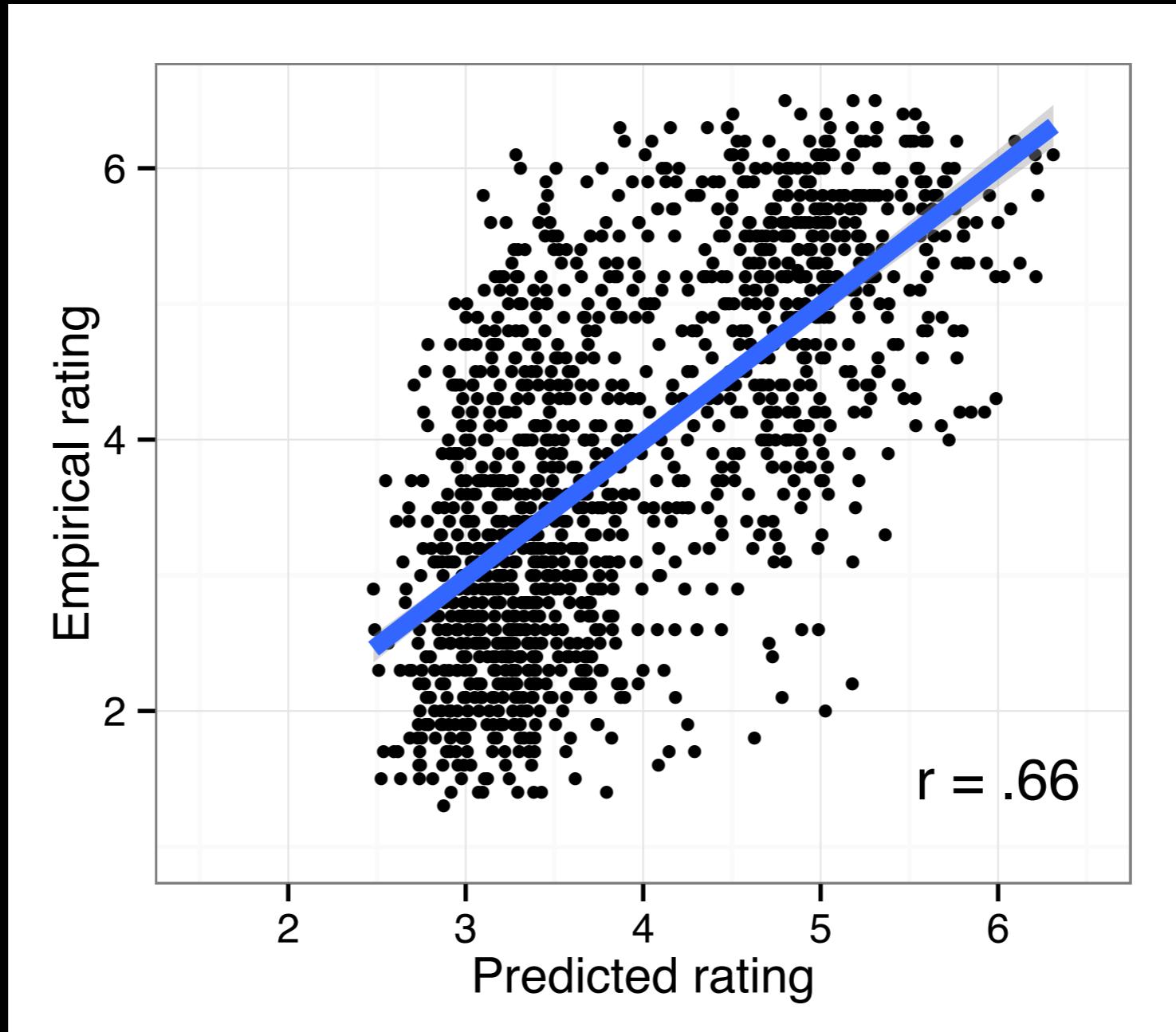
	Coef β	SE(β)	t	p
Intercept	4.01	0.06	68.7	<.0001
Partitive	0.91	0.09	9.6	<.0001
Strength	-0.50	0.05	-9.5	<.0001
Linguistic mention	0.31	0.07	4.4	<.0001
Subjecthood	0.41	0.10	4.2	<.0001
Modification	0.12	0.06	2.0	<.05
Sentence length	0.15	0.05	3.2	<.01
Partitive:Strength	0.39	0.10	4.1	<.0001
Linguistic mention:Subjecthood	0.17	0.21	0.8	<.44
Linguistic mention:Modification	0.34	0.13	2.6	<.01
Subjecthood:Modification	0.27	0.17	1.6	<.12
Linguistic mention:Subjecthood:Modification	0.61	0.42	1.4	<.16

Table 5 Model coefficients for the full model.

Model fit



Model fit



after adding fixed effects of context

Just noise?

Just noise?

No. Variability in ratings is systematically predicted by syntactic, semantic, and pragmatic features of context.

But in some of these cases,
“all” isn’t even an alternative!

You sound like you have **some small ones** in the background.

We've got **some beets**.

That would take **some planning**.

I like **some country music**.

I sold **some of them**.

I think **some parents** go a little bit overboard.

You sound like you have **all small ones** in the background.

We've got **all beets**.

That would take **all planning**.

I like **all country music**.

I sold **all of them**.

I think **all parents** go a little bit overboard.

1.5

You sound like you have **all small ones** in the background.

2.7

We've got **all beets**.

1.4

That would take **all planning**.

6.9

I like **all country music**.

6.8

I sold **all of them**.

6.4

I think **all parents** go a little bit overboard.

1.5

You sound like you have **all small ones** in the background.

2.7

We've got **all beets**.

1.4

That would take **all planning**.

6.9

I like **all country music**.

6.8

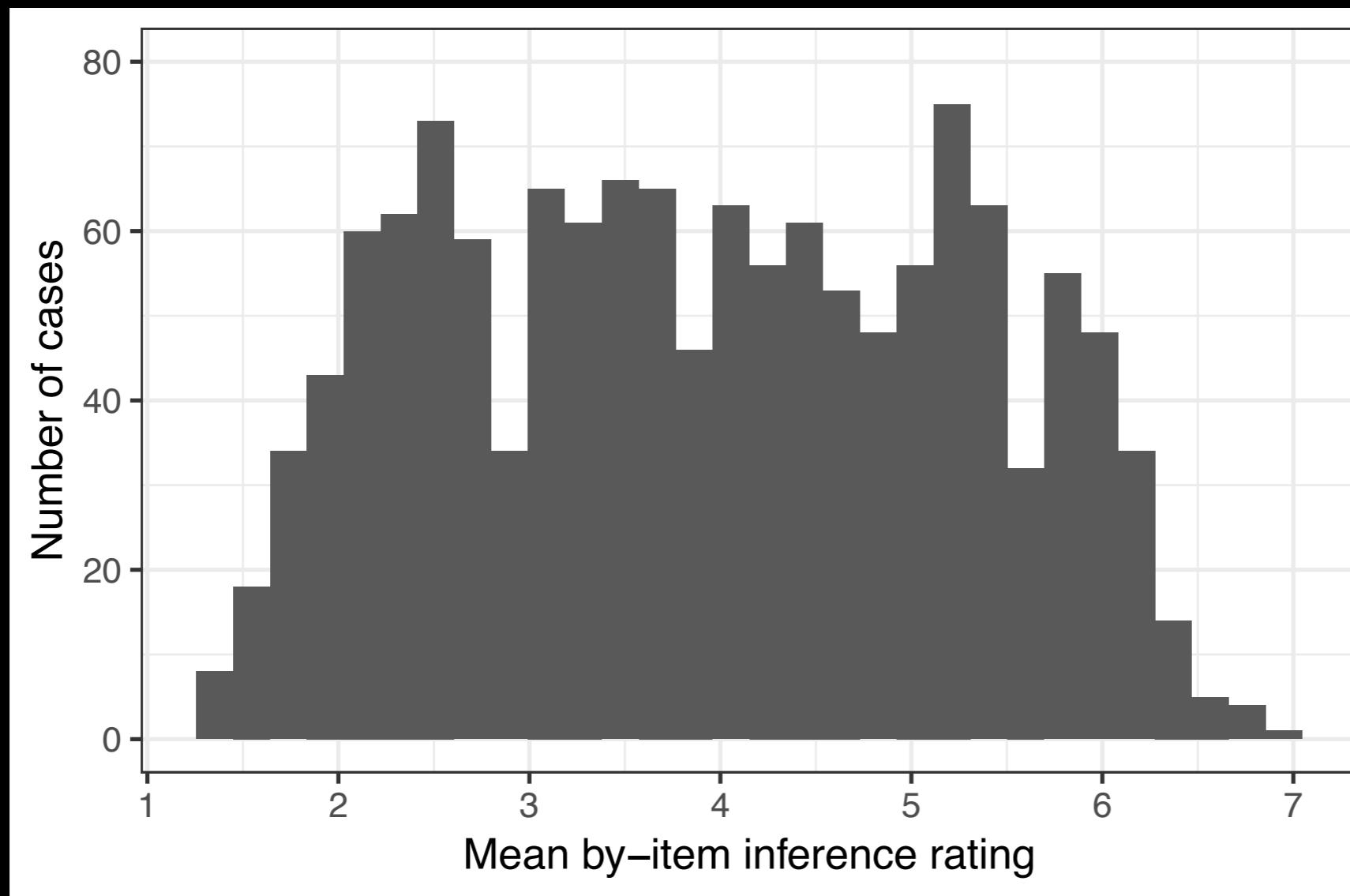
I sold **all of them**.

6.4

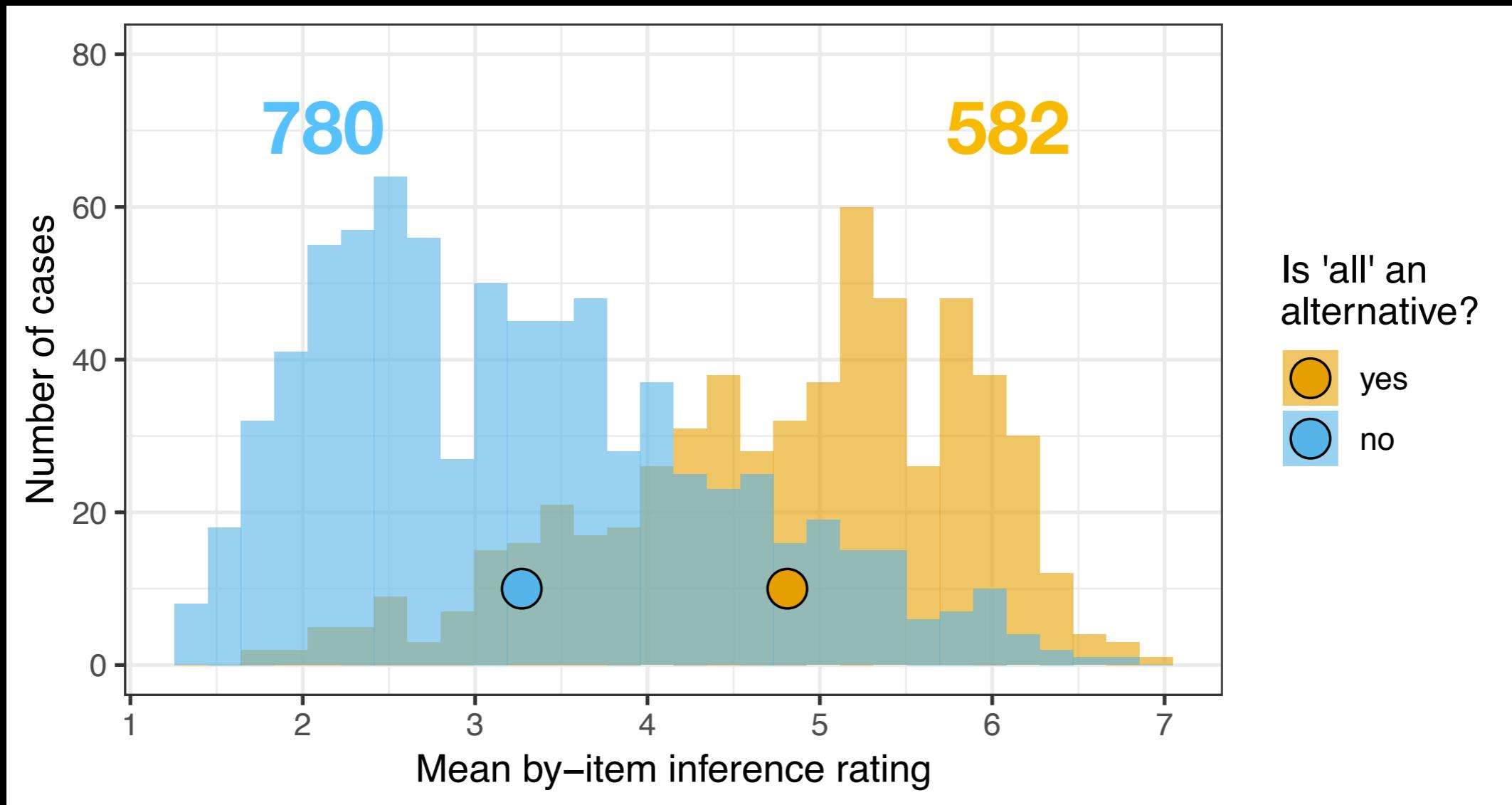
I think **all parents** go a little bit overboard.

All cases hand-annotated by 2 RAs for whether “some” can be replaced by “all” or only by “a lot (of)”

Variability in inference strength



Variability in inference strength



rating ~
 partitive + linguistic mention + subjecthood + ...
 + random effects

original model

	Coef β	SE(β)	t	p
Intercept	4.01	0.06	68.7	<.0001
Partitive	0.91	0.09	9.6	<.0001
Strength	-0.50	0.05	-9.5	<.0001
Linguistic mention	0.31	0.07	4.4	<.0001
Subjecthood	0.41	0.10	4.2	<.0001
Modification	0.12	0.06	2.0	<.05
Sentence length	0.15	0.05	3.2	<.01
Partitive:Strength	0.39	0.10	4.1	<.0001
Linguistic mention:Subjecthood	0.17	0.21	0.8	<.44
Linguistic mention:Modification	0.34	0.13	2.6	<.01
Subjecthood:Modification	0.27	0.17	1.6	<.12
Linguistic mention:Subjecthood:Modification	0.61	0.42	1.4	<.16

rating ~

(partitive + linguistic mention + subjecthood + ...)

* alternative

+ random effects

original model

	Coef β	SE(β)	t	p
Intercept	4.01	0.06	68.7	<.0001
Partitive	0.91	0.09	9.6	<.0001
Strength	-0.50	0.05	-9.5	<.0001
Linguistic mention	0.31	0.07	4.4	<.0001
Subjecthood	0.41	0.10	4.2	<.0001
Modification	0.12	0.06	2.0	<.05
Sentence length	0.15	0.05	3.2	<.01
Partitive:Strength	0.39	0.10	4.1	<.0001
Linguistic mention:Subjecthood	0.17	0.21	0.8	<.44
Linguistic mention:Modification	0.34	0.13	2.6	<.01
Subjecthood:Modification	0.27	0.17	1.6	<.12
Linguistic mention:Subjecthood:Modification	0.61	0.42	1.4	<.16

rating ~

(partitive + linguistic mention + subjecthood + ...)

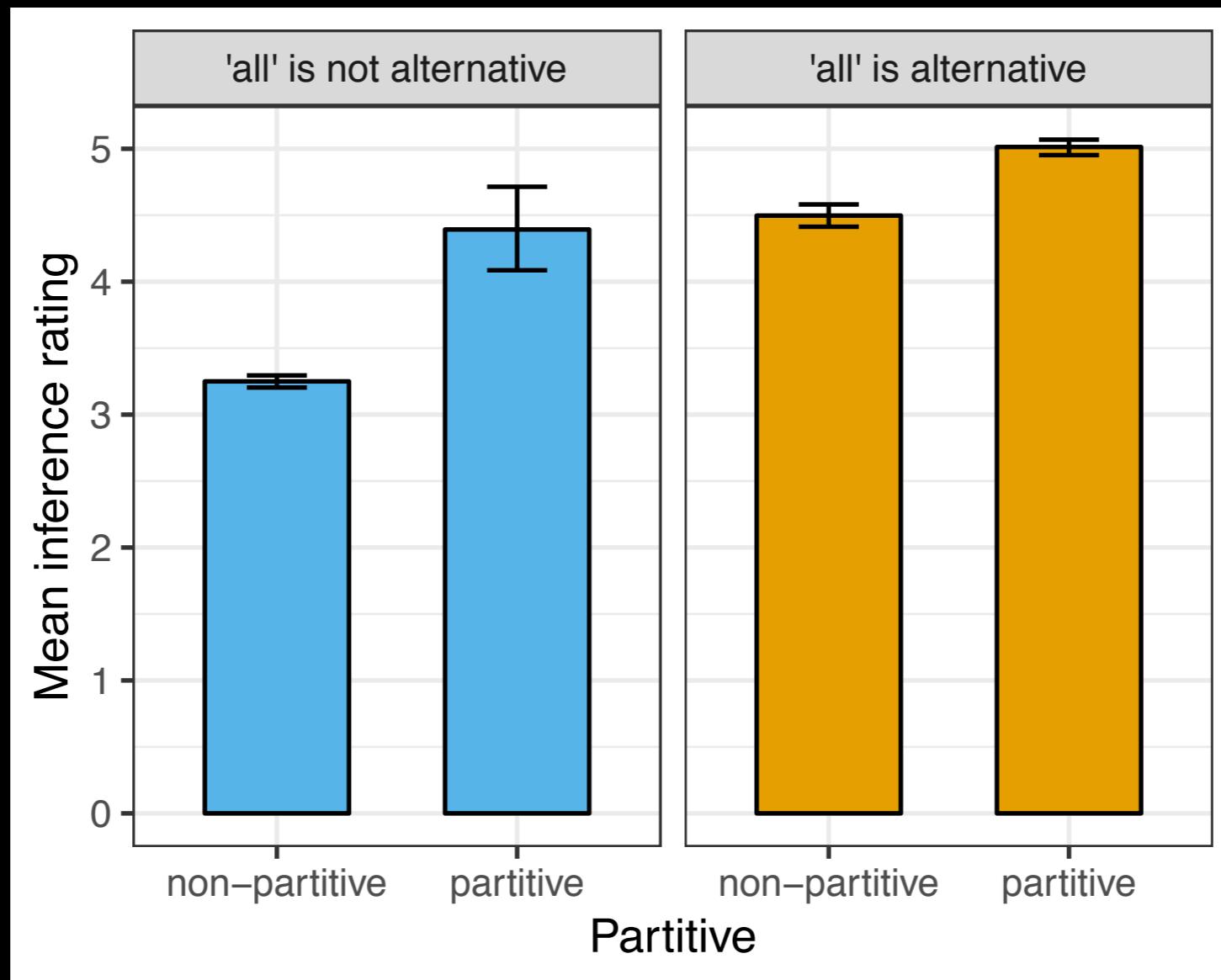
* alternative

+ random effects

	original model	original with alternative			
	Coef β	SE(β)	t	p	p
Intercept	4.01	0.06	68.7	<.0001	<.0001
Partitive	0.91	0.09	9.6	<.0001	<.0001
Strength	-0.50	0.05	-9.5	<.0001	<.0001
Linguistic mention	0.31	0.07	4.4	<.0001	<.0001
Subjecthood	0.41	0.10	4.2	<.0001	<.0001
Modification	0.12	0.06	2.0	<.05	<.0001
Sentence length	0.15	0.05	3.2	<.01	<.0001
Partitive:Strength	0.39	0.10	4.1	<.0001	<.0001
Linguistic mention:Subjecthood	0.17	0.21	0.8	<.44	>0.52
Linguistic mention:Modification	0.34	0.13	2.6	<.01	<.01
Subjecthood:Modification	0.27	0.17	1.6	<.12	<.01
Linguistic mention:Subjecthood:Modification	0.61	0.42	1.4	<.16	<.01

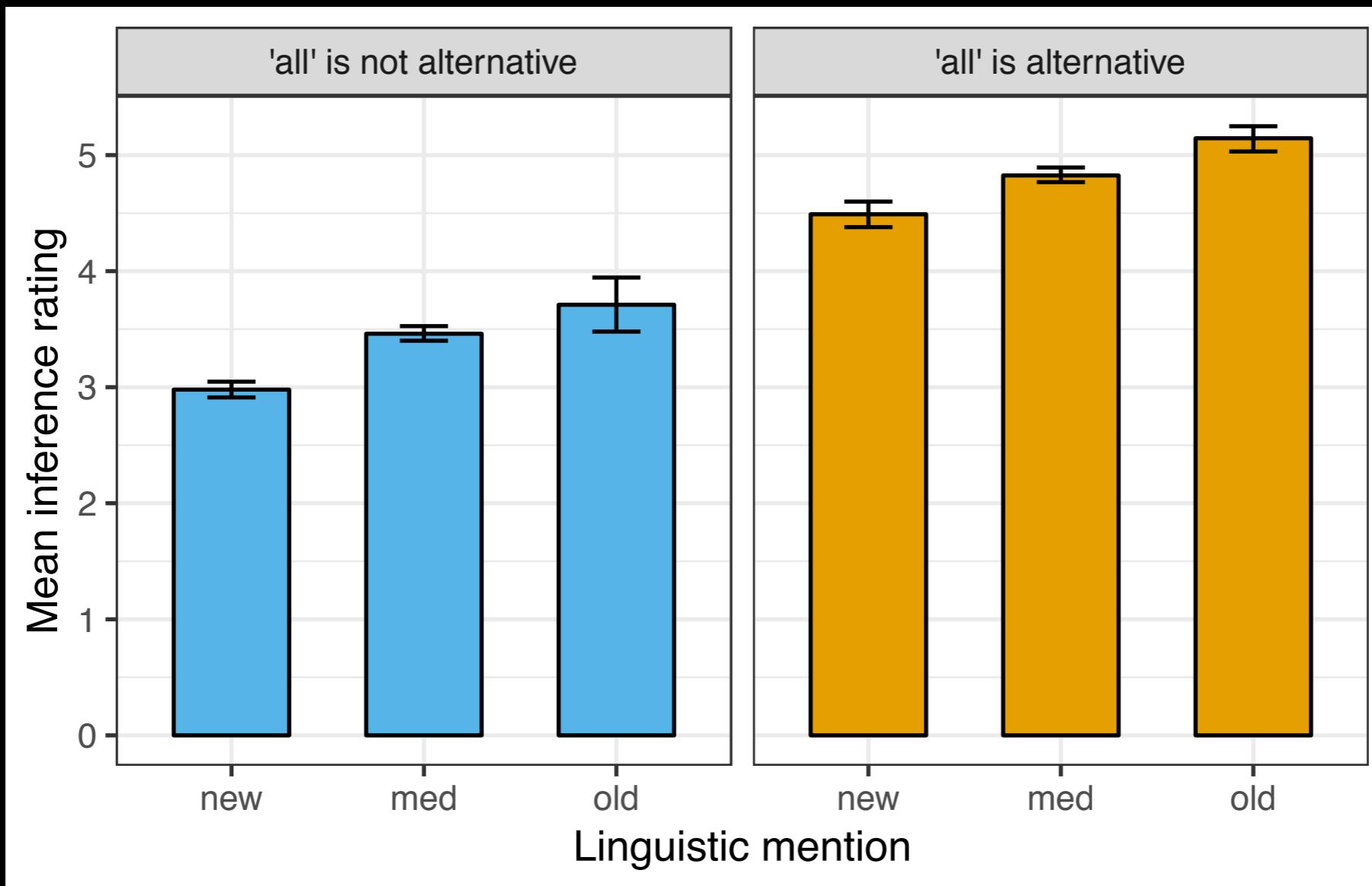
Stronger inferences...

...with **partitive** *some-NPs*.



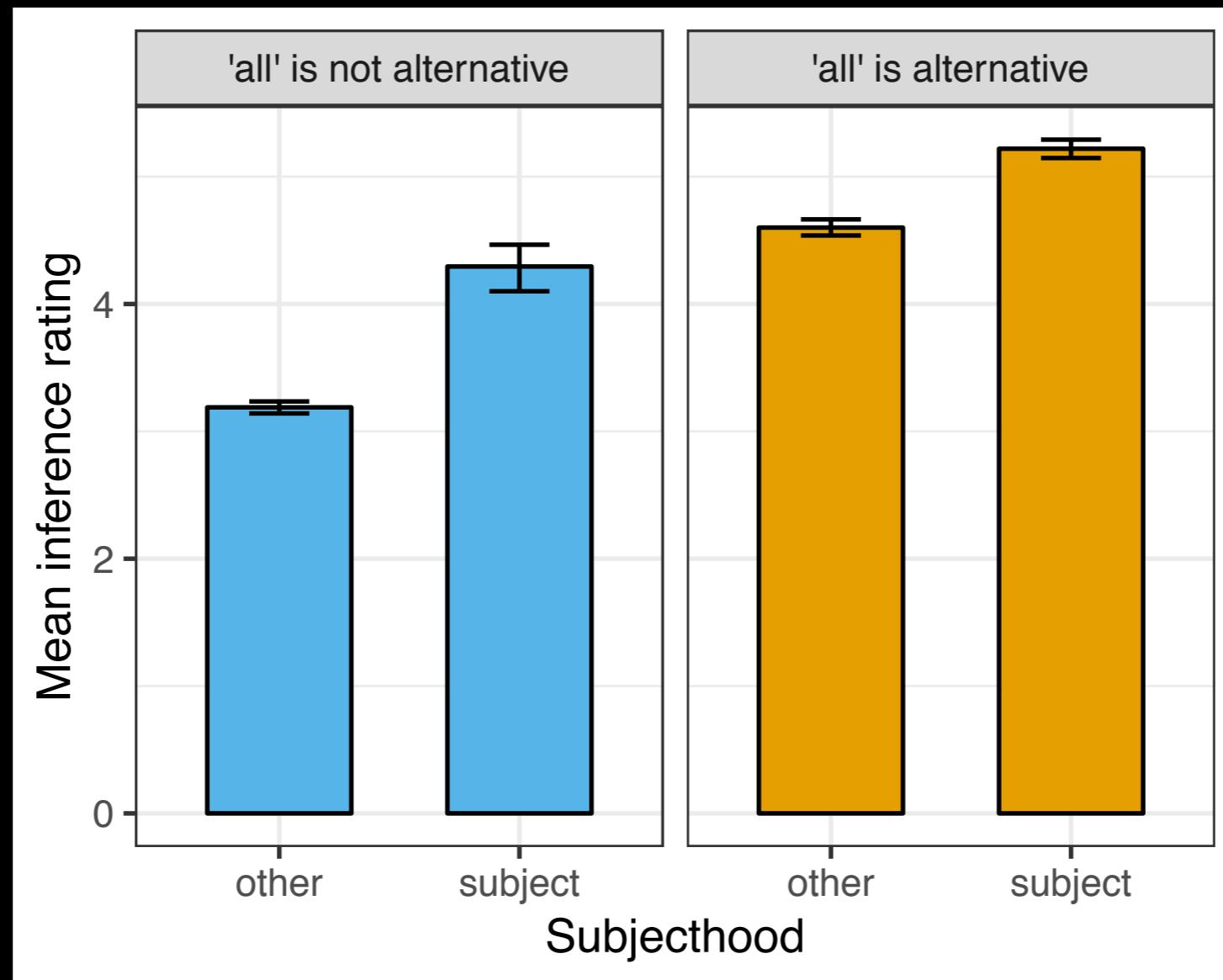
Stronger inferences...

...with **previously mentioned** NP referents.



Stronger inferences...

...with *some*-NPs in **subject** position.



Just noise?

Just noise?

No. Variability in ratings is systematically predicted by syntactic, semantic, and pragmatic features of context.

No. Replication by Eiteljoerge et al 2019 in child-directed speech

Just noise?

No. Variability in ratings is systematically predicted by syntactic, semantic, and pragmatic features of context.

Implications for theories of
pragmatic inference

Implications for theories of pragmatic inference

The status of scalar implicatures
as GCIs is highly questionable.

**How many features? Do they
need to be hand-mined?**

Predicting inference strength from distributed meaning representations

Schuster, Chen, & Degen, on arXiv tomorrow!



Sebastian
Schuster

Yuxing
Chen

Ultimate goal:

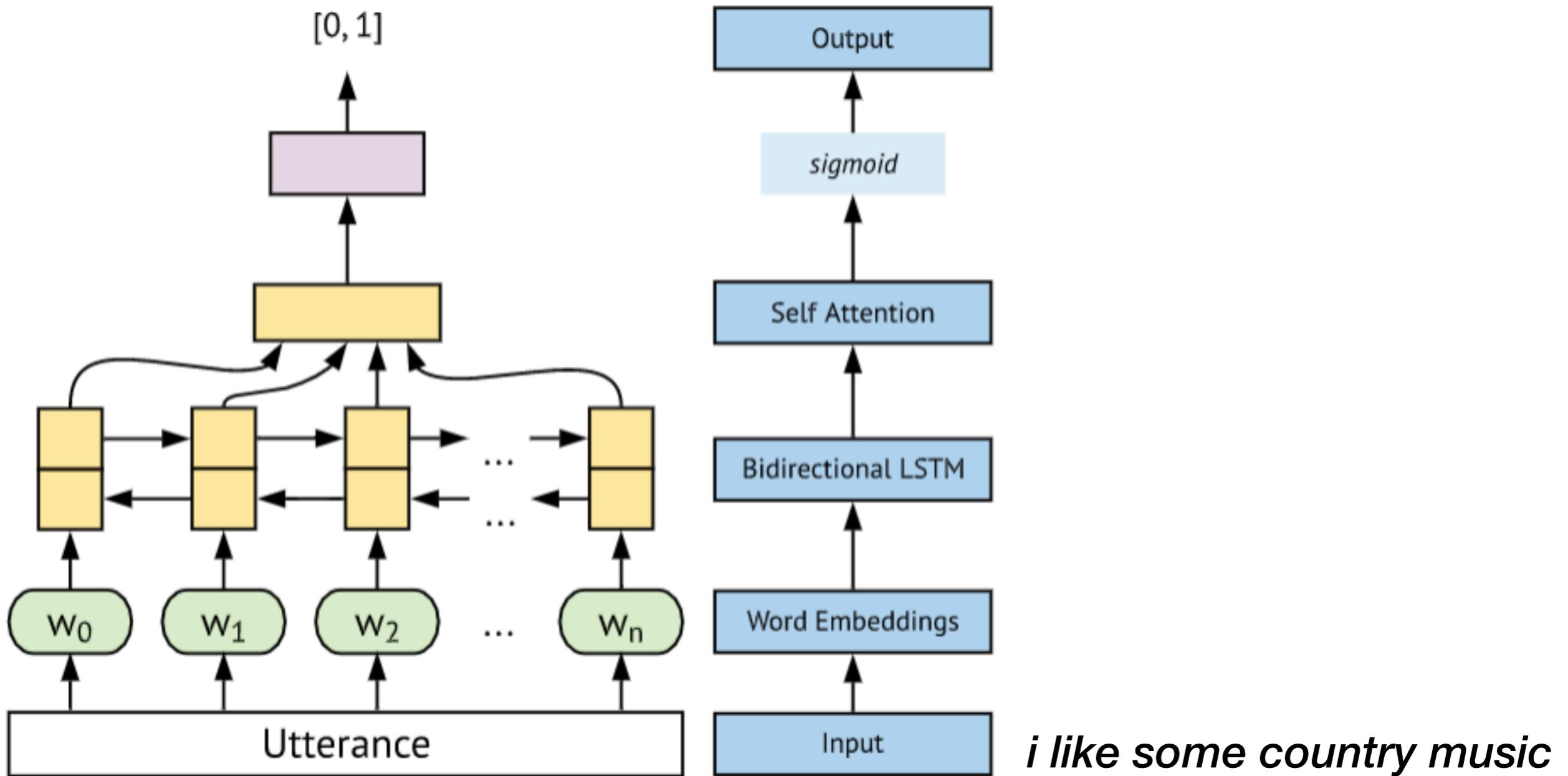
Use distributed vector-based meaning representation methods from NLP to infer which, if any, linguistically encoded features of context listeners use in drawing inferences, to help inform pragmatic theory.

More proximate goal:

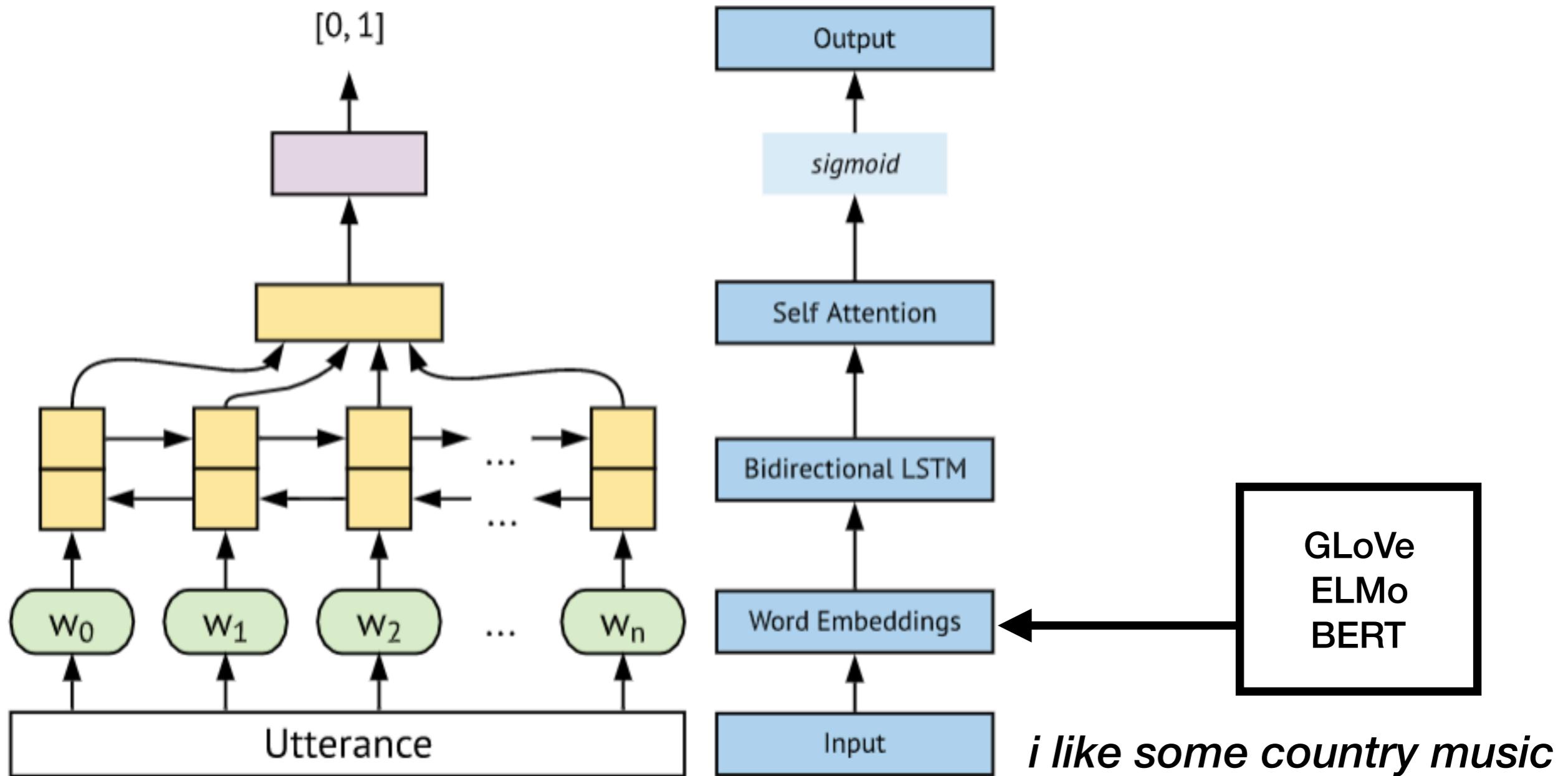
Use distributed vector-based meaning representation methods from NLP to test whether any of these methods

- reliably predict inference ratings
- capture the identified context effects

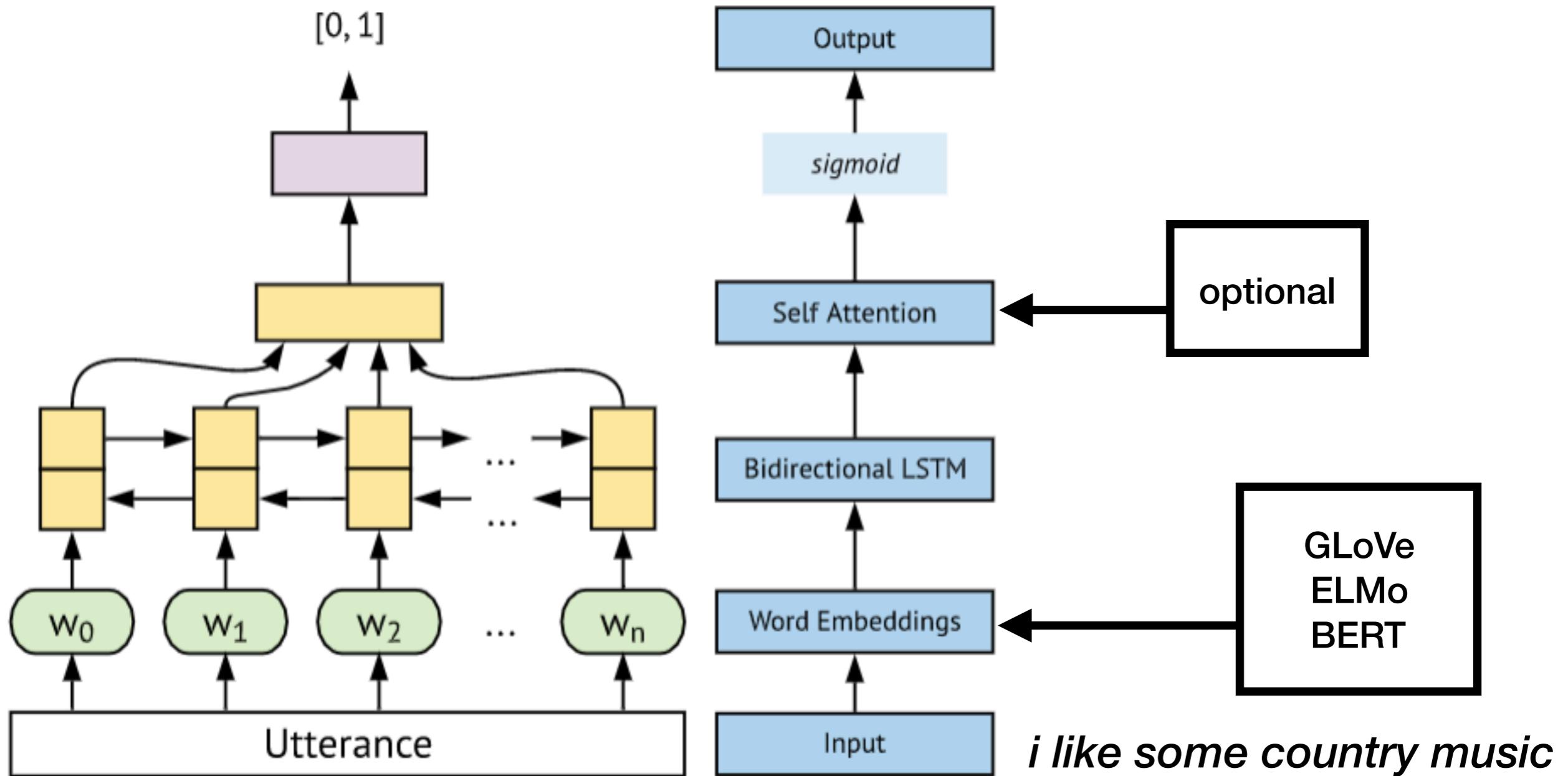
Model architecture



Model architecture



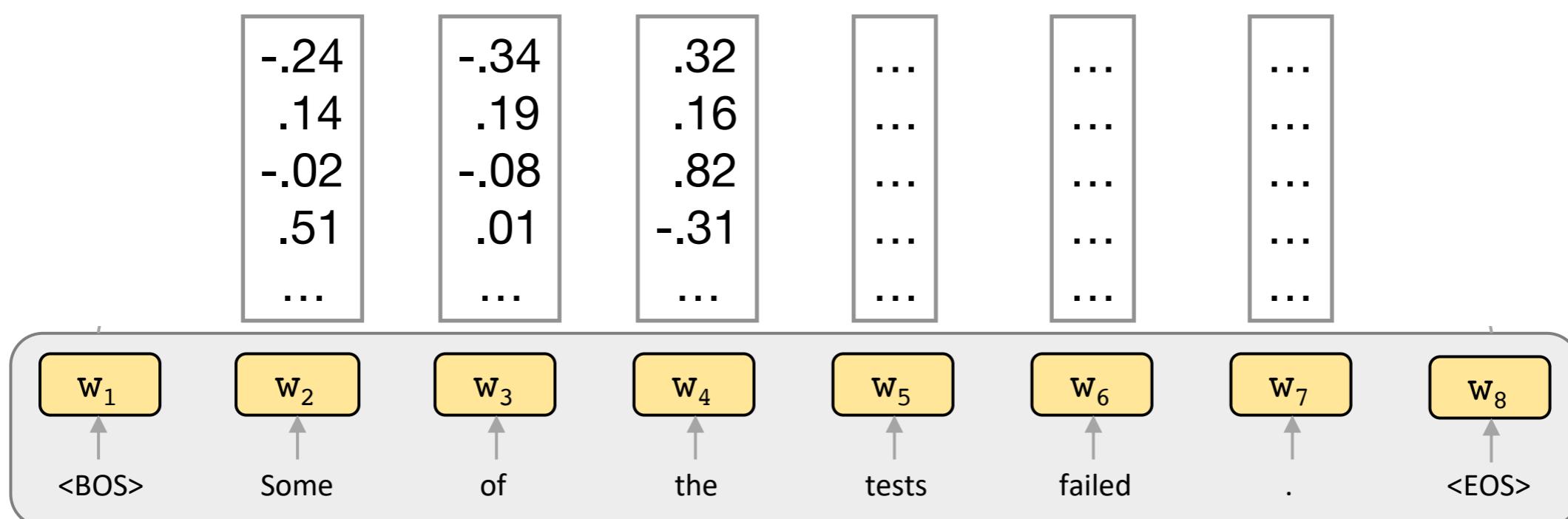
Model architecture

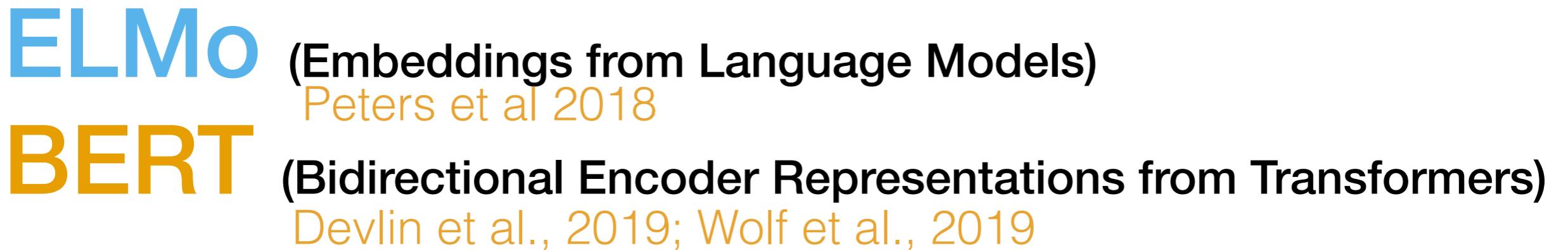


GloVe (Global Vectors for word representation)

- captures meaning in vector space
- based on co-occurrence statistics of words
- 100-dimensional vector for each word, pre-trained on 6 billion tokens from Wikipedia 2014 and Gigaword 5
- words around “some” encoded in pretrained 100-dimensional GloVe vectors

Pennington et al 2014





- **contextual** word embeddings (considers entire sentence before assigning a word in it an embedding)
- captures that the same word can have different meanings in different sentences
 1. **Apple** announced the new iPhone today.
 2. **Google** announced a new browser last week.
 3. I ate an **apple** for breakfast.
 4. I ate an **orange** after dinner.
- ELMo: based on word sequence modeling (bi-directional LSTM)
- BERT: based on transformers (also bi-directional)
- pre-trained

Context preceding sentence

Speaker A: i mean, they just have beautiful, beautiful homes and they have everything. the kids only wear name brand things to school and it's one of these things,

Speaker B: oh me. well that makes it hard for you, doesn't it.

Speaker A: well it does, you know. it really does because i'm a single mom and i have a thirteen year old now and uh, you know, it does.

Speaker B: oh, me.

Speaker A: i mean, we do it to a point but uh, not to where she feels different ,

Speaker B: yeah.

Speaker A:
but some of them are very rich

either did or didn't include context in generating the sentence embedding
(context may be important for capturing factors like linguistic mention)

Context preceding sentence

Speaker A: i mean, they just have beautiful, beautiful homes and they have everything. the kids only wear name brand things to school and it's one of these things,

Speaker B: oh me. well that makes it hard for you, doesn't it.

Speaker A: well it does, you know. it really does because i'm a single mom and i have a thirteen year old now and uh, you know, it does.

Speaker B: oh, me.

Speaker A: i mean, we do it to a point but uh, not to where she feels different ,

Speaker B: yeah.

Speaker A:

but some of them are very rich

either did or didn't include context in generating the sentence embedding
(context may be important for capturing factors like linguistic mention)

Context preceding sentence

Speaker A: i mean, they just have beautiful, beautiful homes and they have everything. the kids only wear name brand things to school and it's one of these things,

Speaker B: oh me. well that makes it hard for you, doesn't it.

Speaker A: well it does, you know. it really does because i'm a single mom and i have a thirteen year old now and uh, you know, it does.

Speaker B: oh, me.

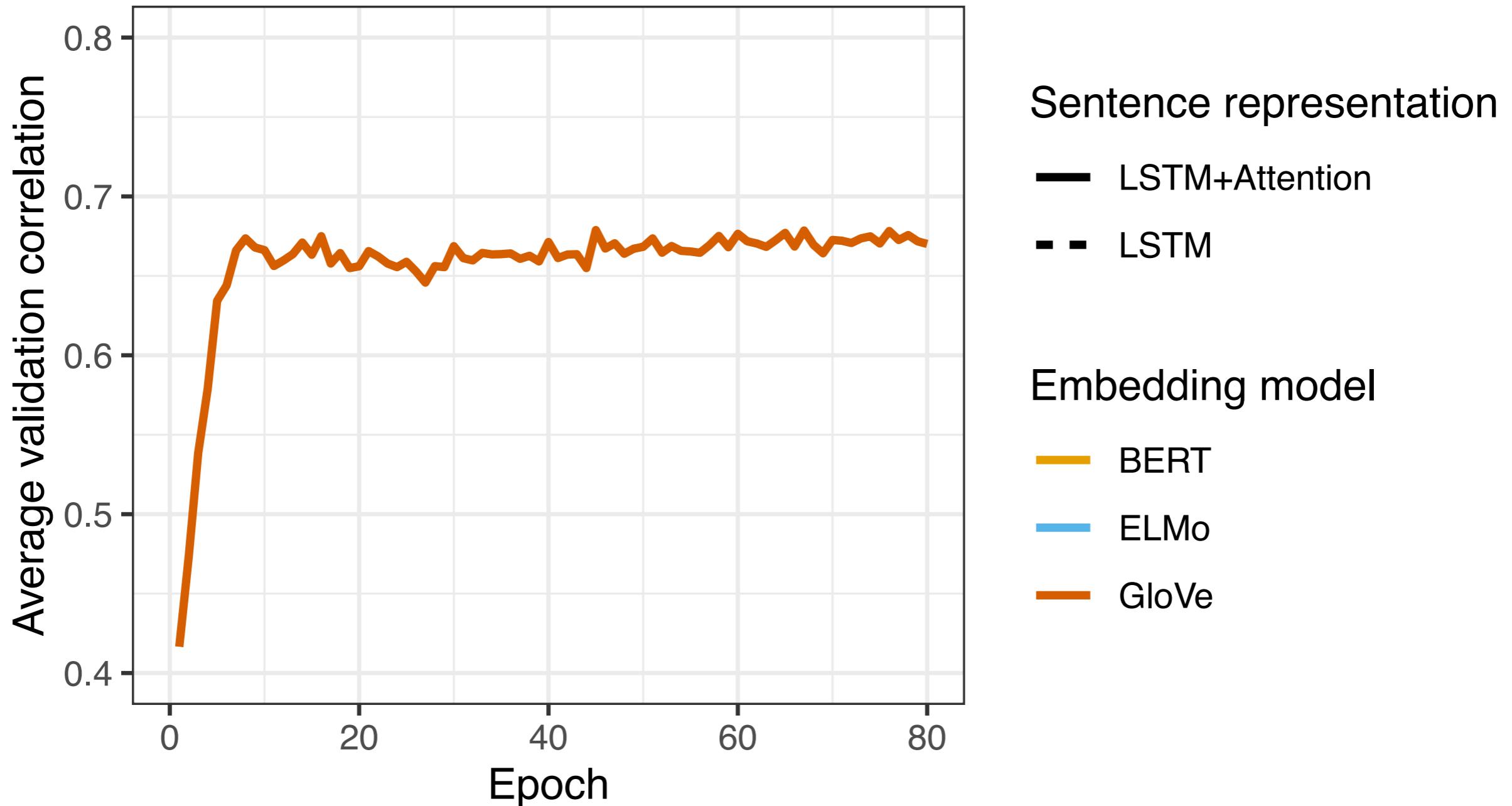
Speaker A: i mean, we do it to a point but uh, not to where she feels different ,

Speaker B: yeah.

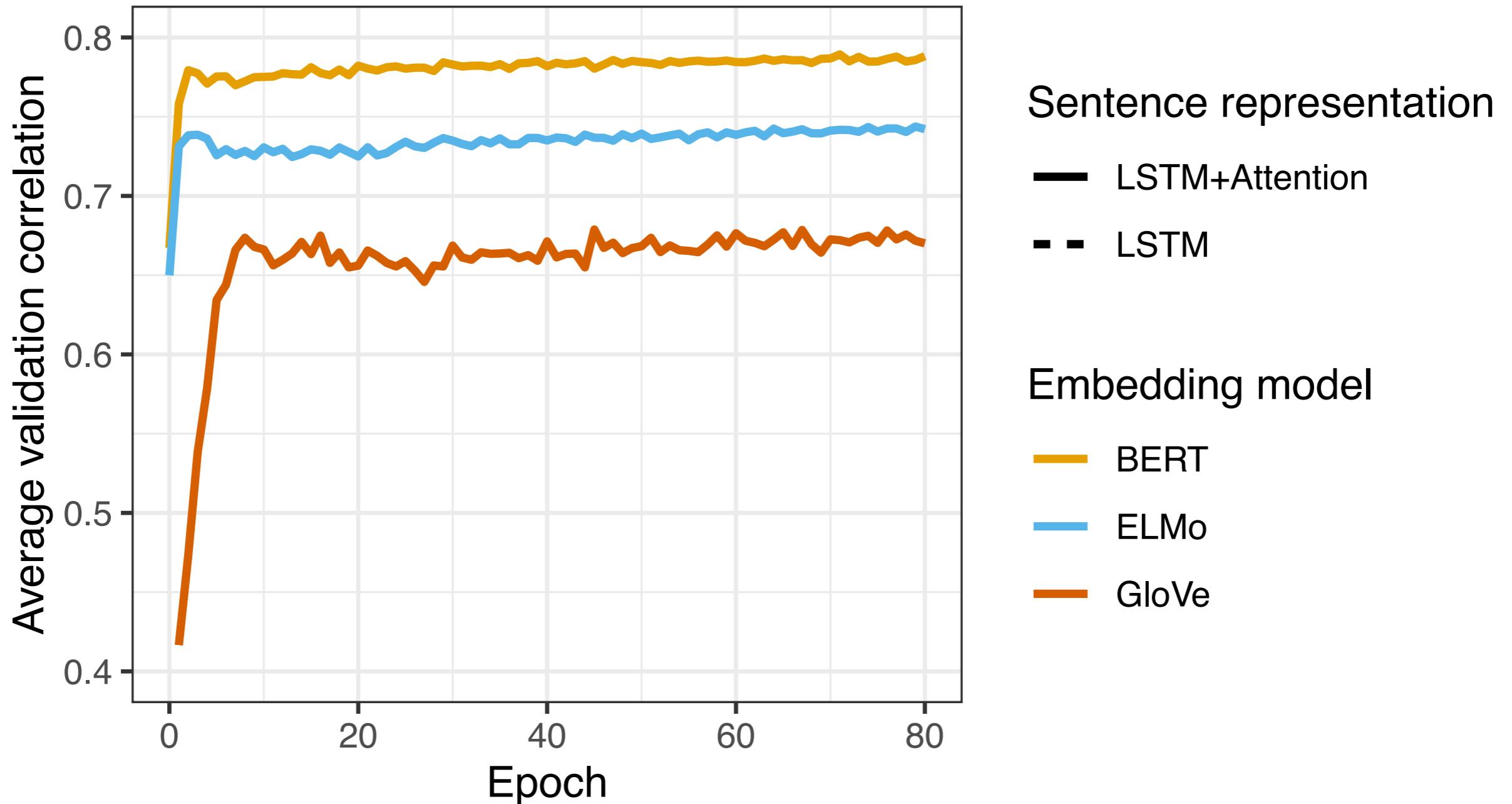
Speaker A:
but some of them are very rich

either did or didn't include context in generating the sentence embedding
(context may be important for capturing factors like linguistic mention)

Results on validation set (30%)

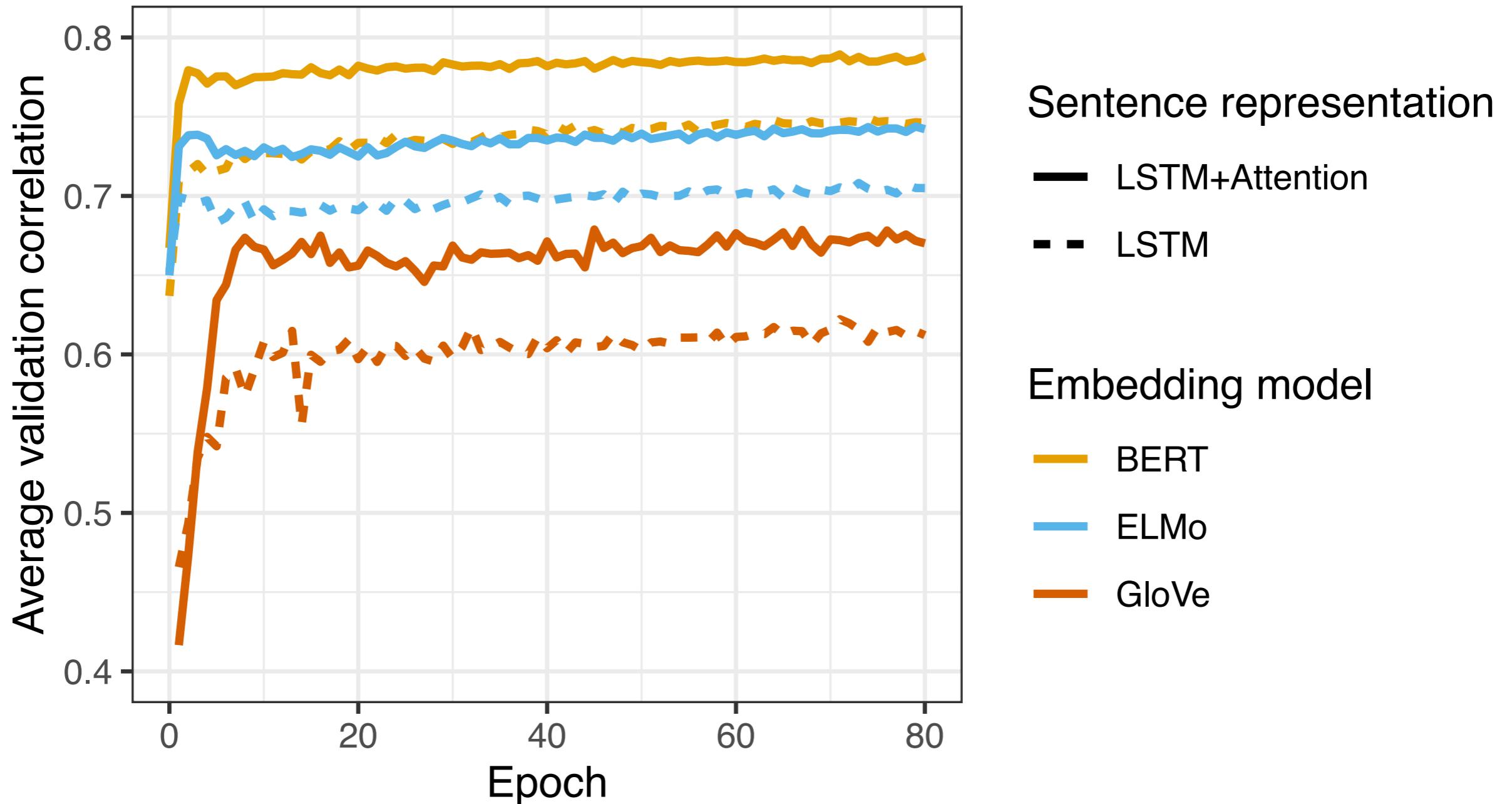


Results on validation set (30%)



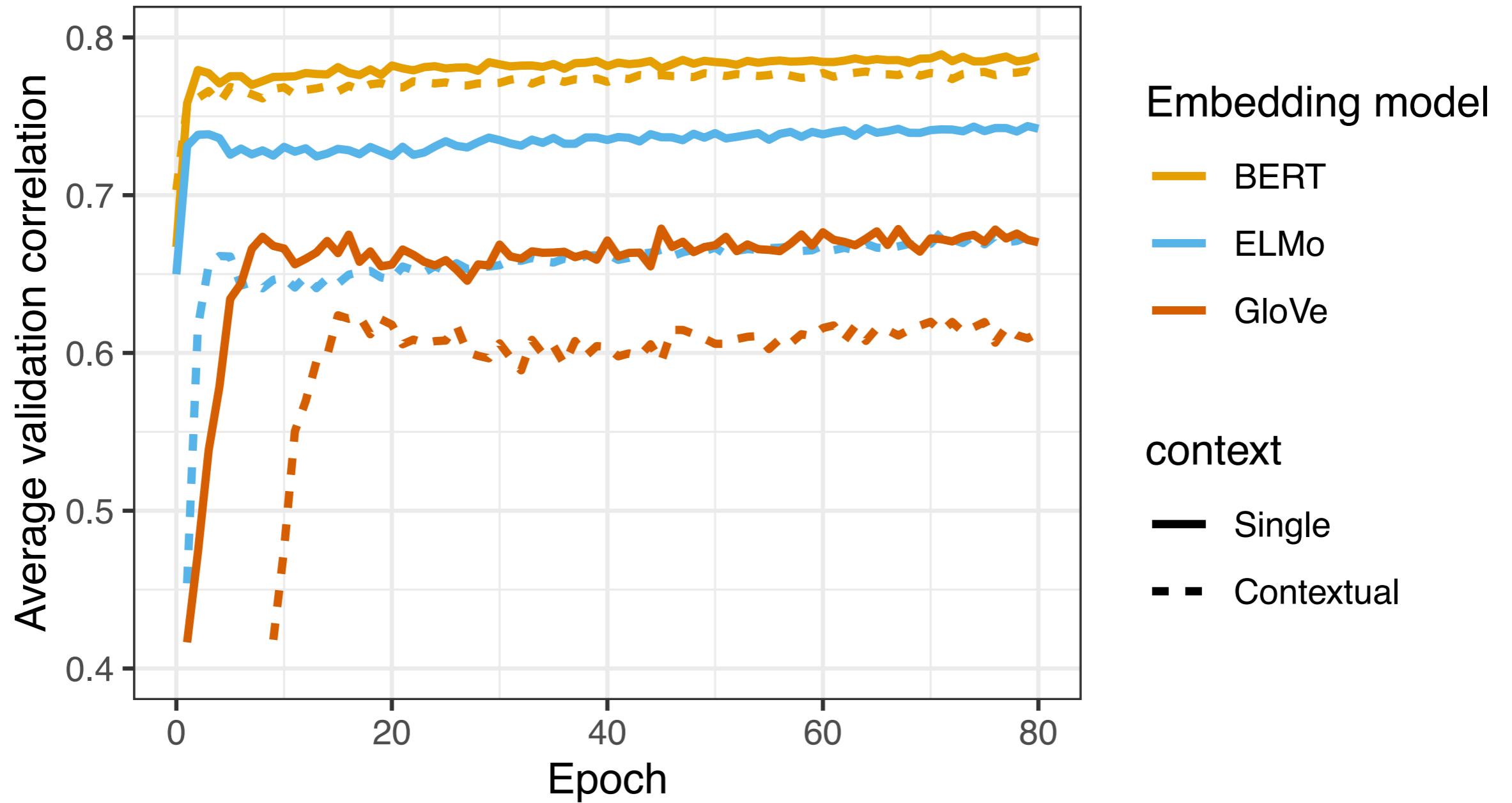
contextual embeddings do better than static ones

Results on validation set (30%)



removing attention hurts performance

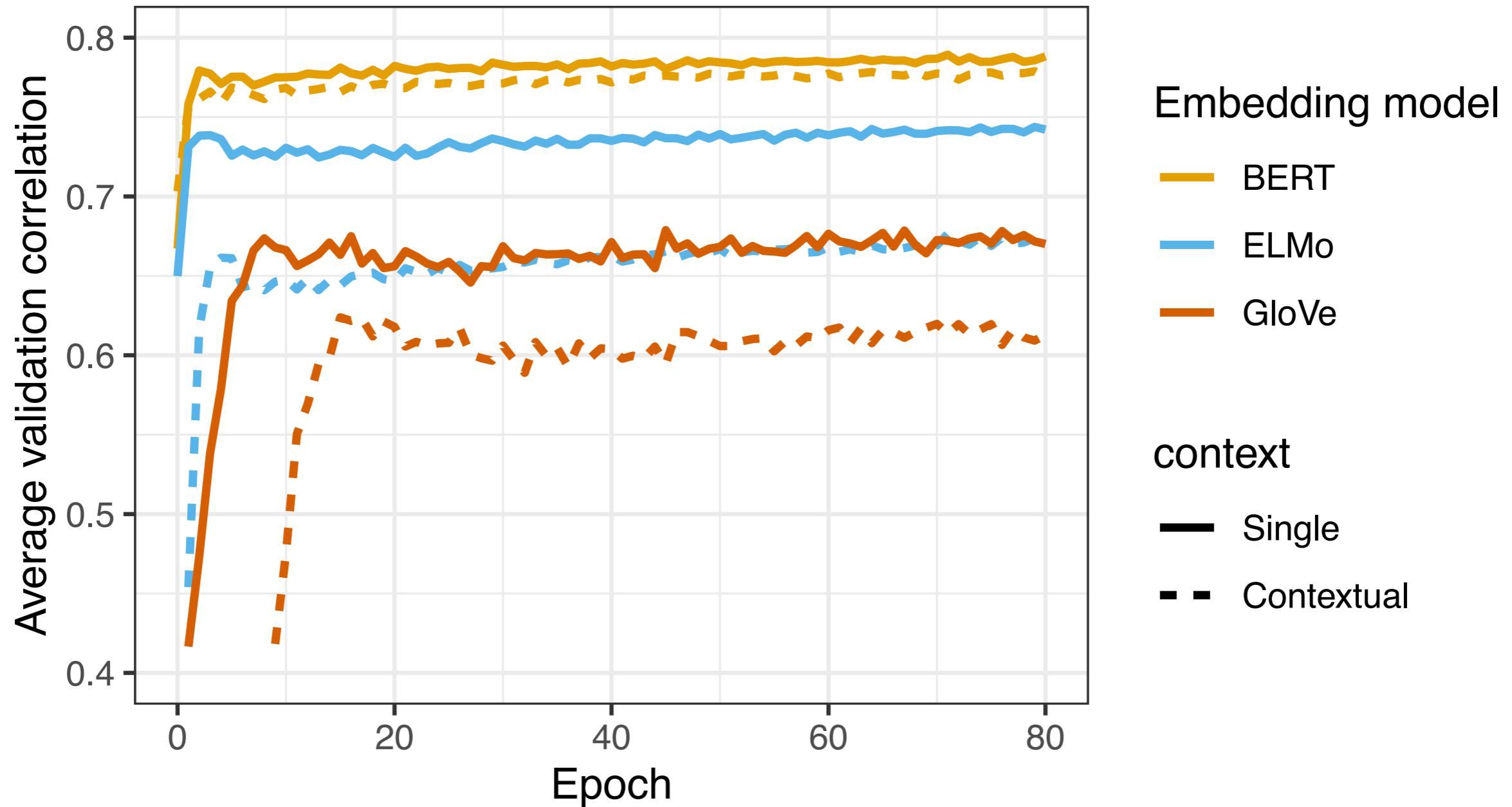
Results on validation set (30%)



adding context hurts performance

Best model ($r=.78$):
BERT
LSTM+attention
no context

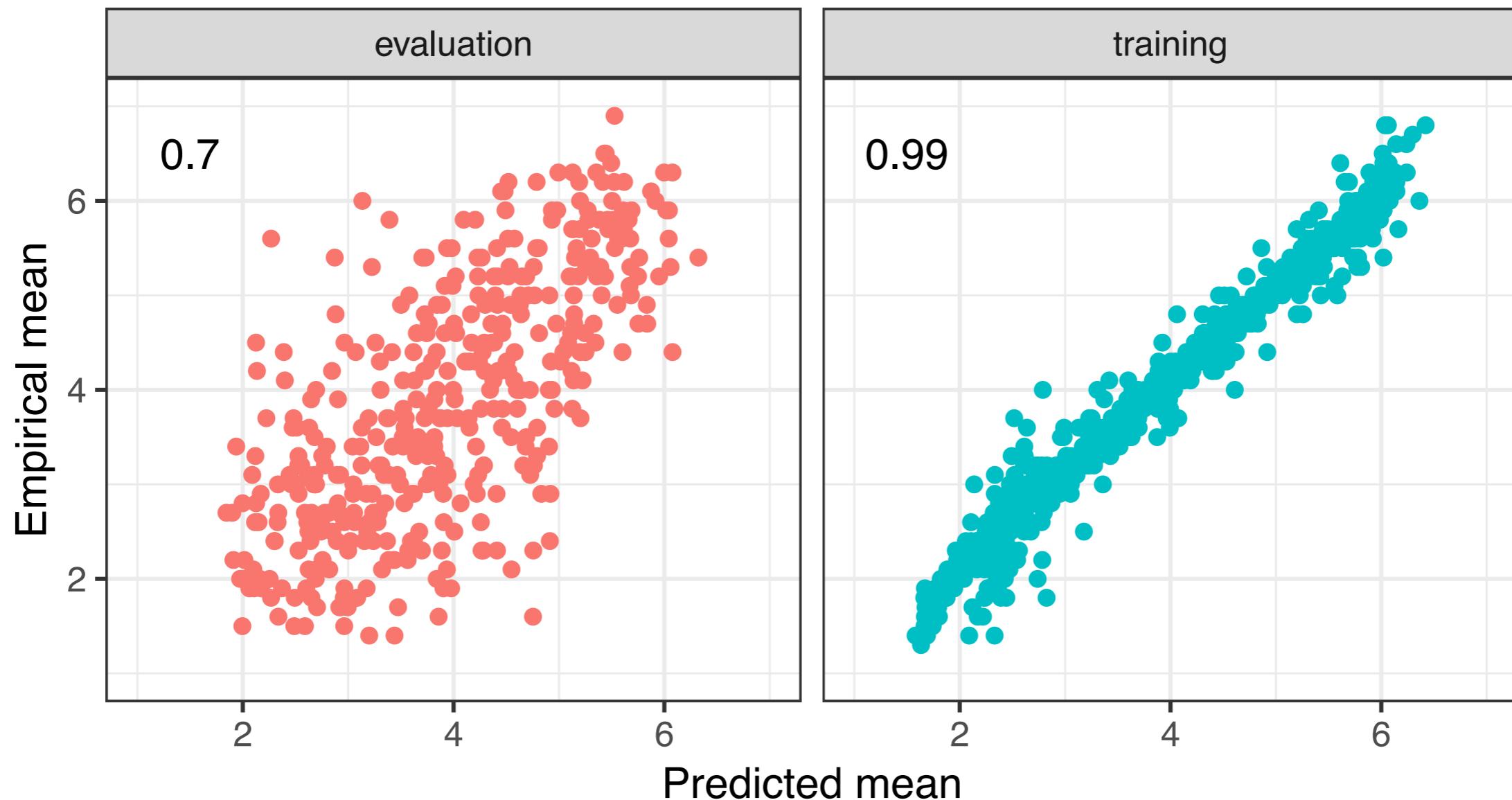
Results on validation set (30%)



adding context hurts performance

Model predictions

Best model: BERT — LSTM + attention — no-context



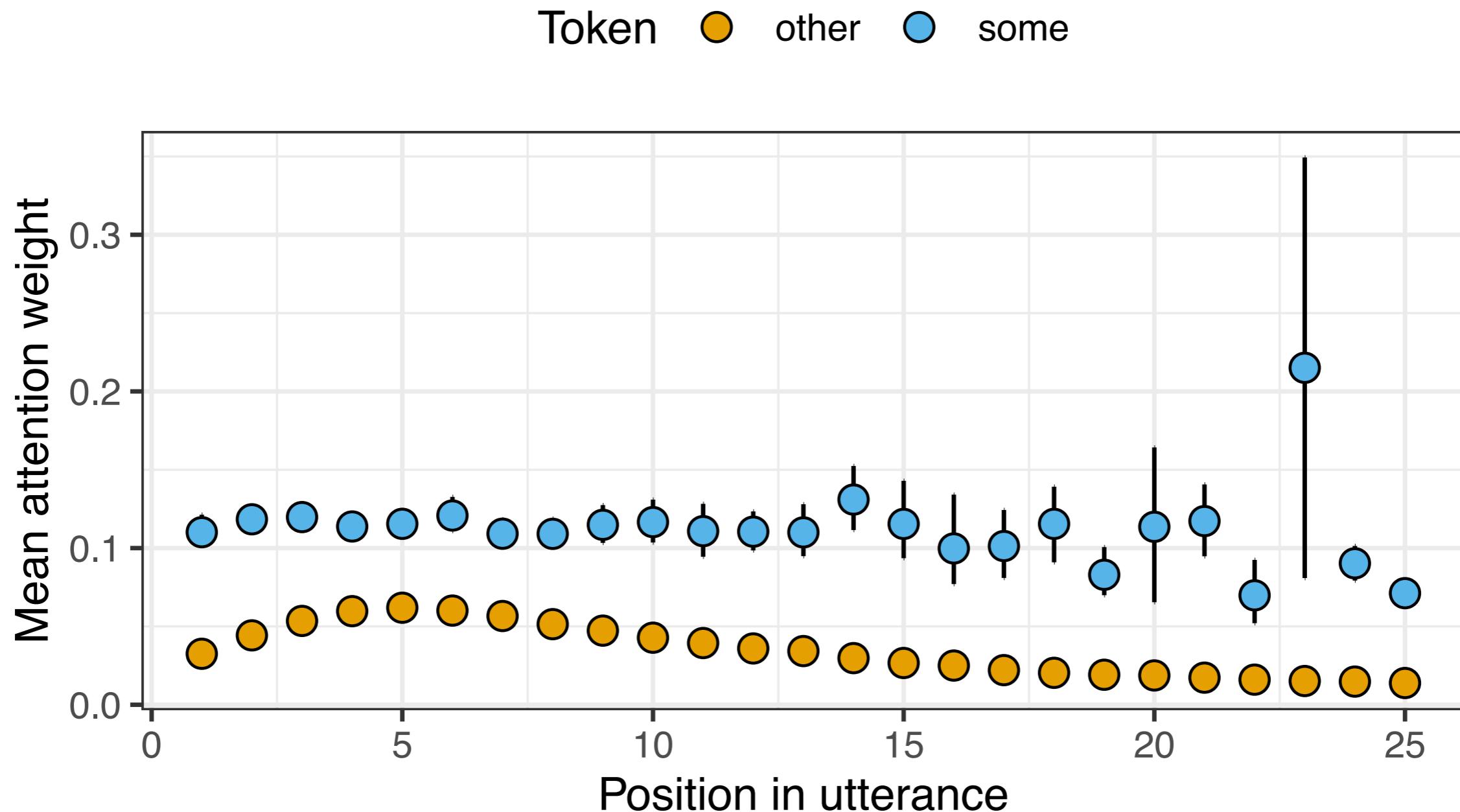
Attention weight analysis

Lee et al., 2017; Ding et al., 2017; Wiegreffe and Pinter, 2019

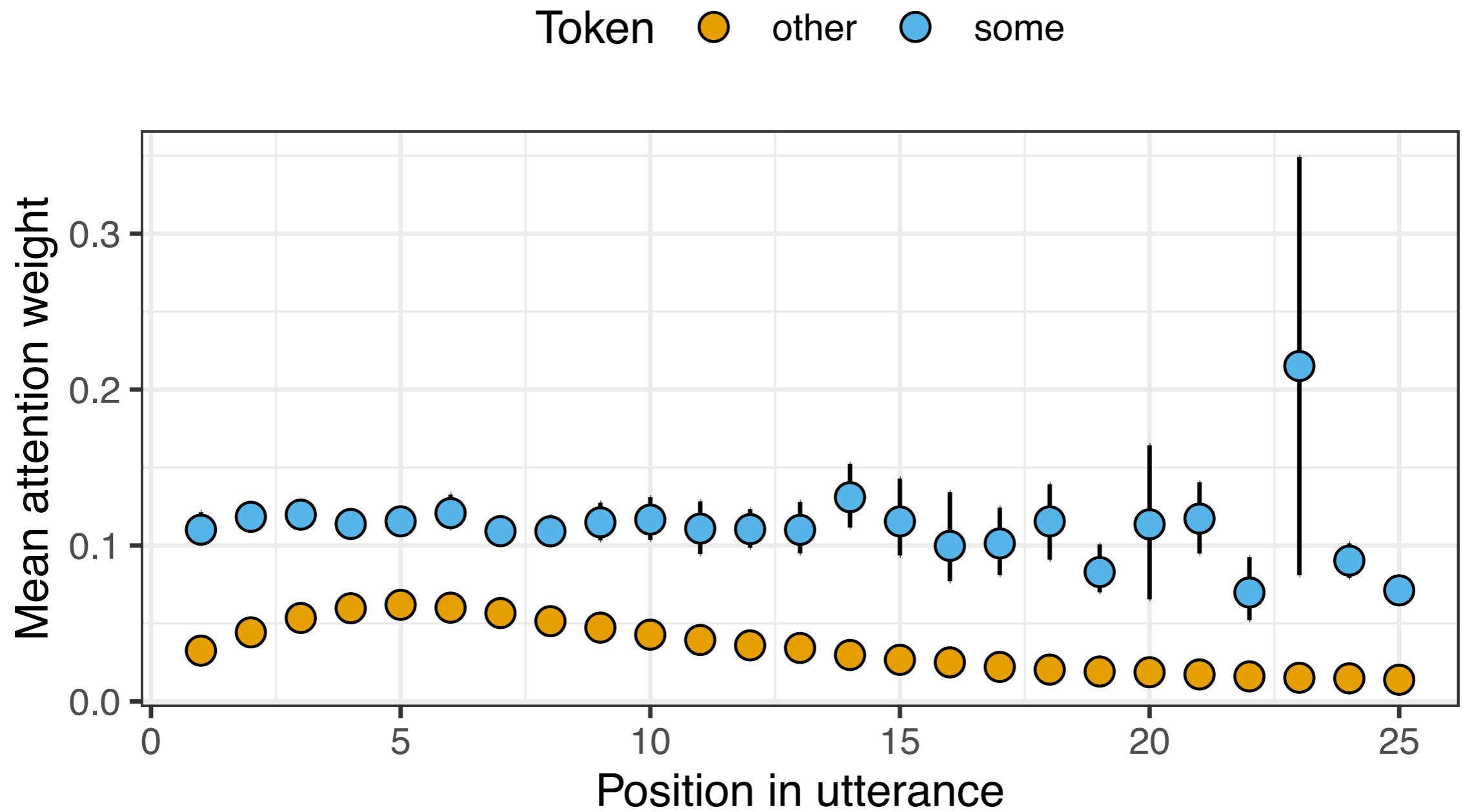
Is there any evidence that the model learned to pay attention to a priori relevant utterance tokens?

Attention to “some”

Attention to “some”



Attention to “some”

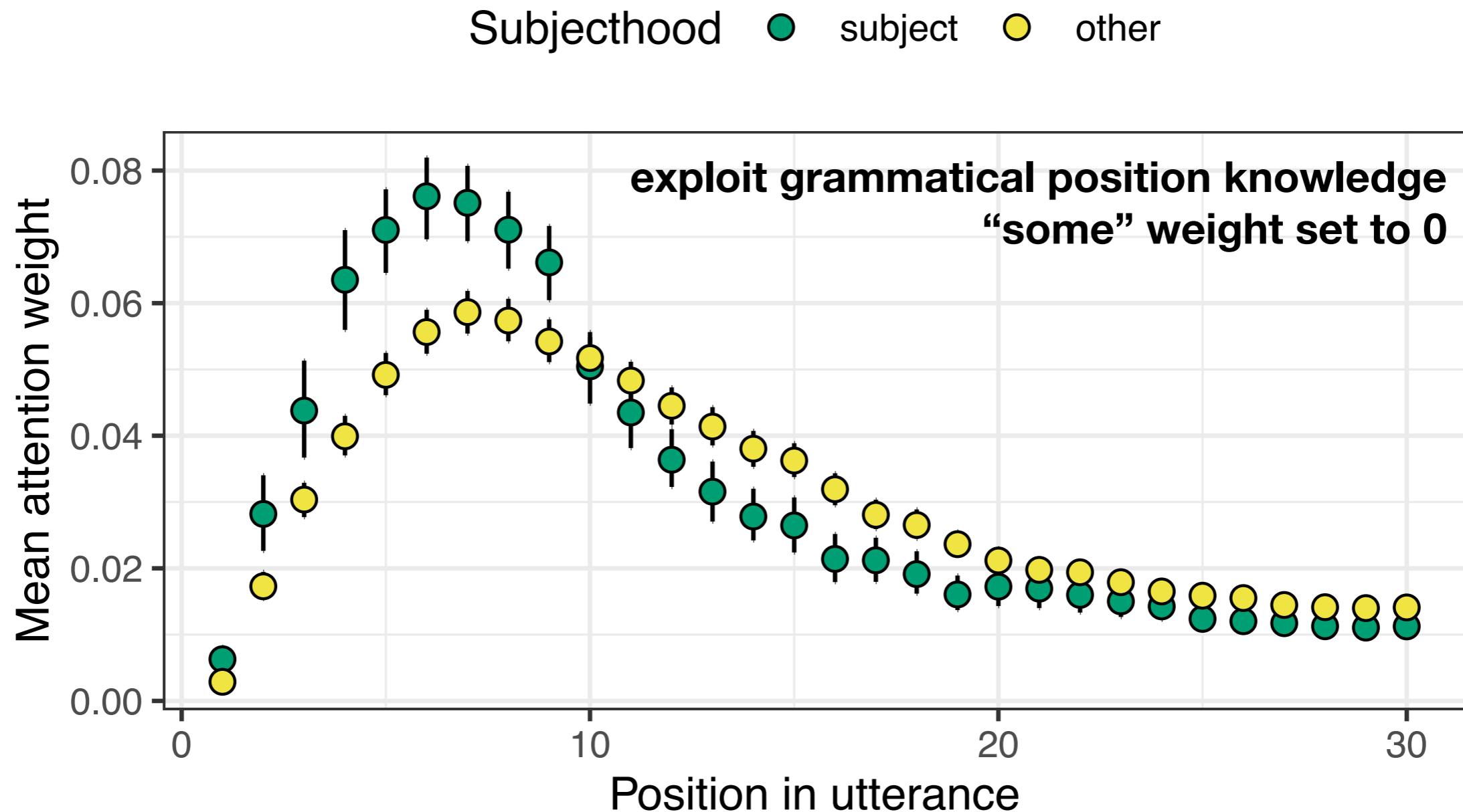


more attention to “some” than other tokens throughout sentence

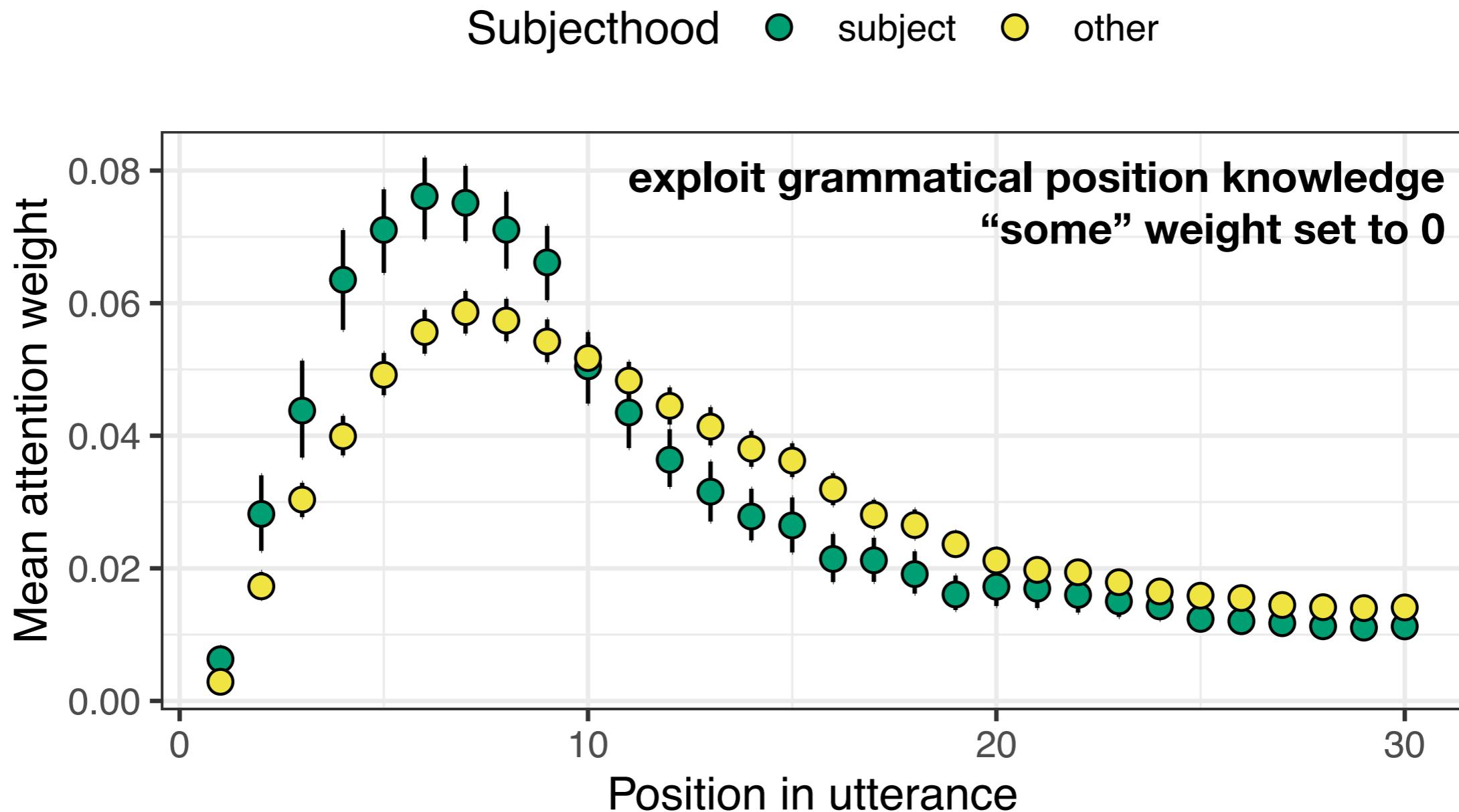
Attention to “some”-NP

**exploit grammatical position knowledge
“some” weight set to 0**

Attention to “some”-NP

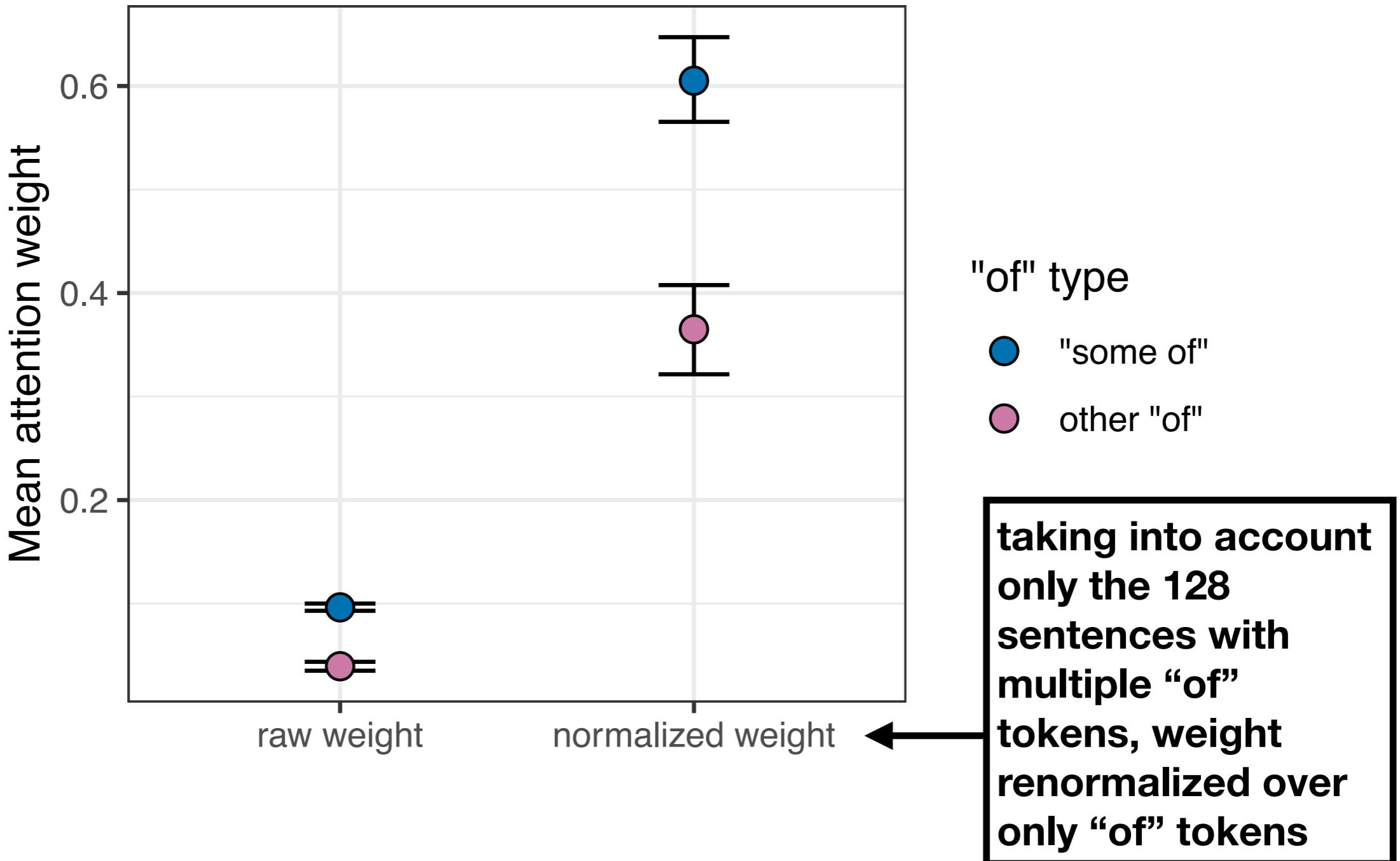


Attention to “some”-NP

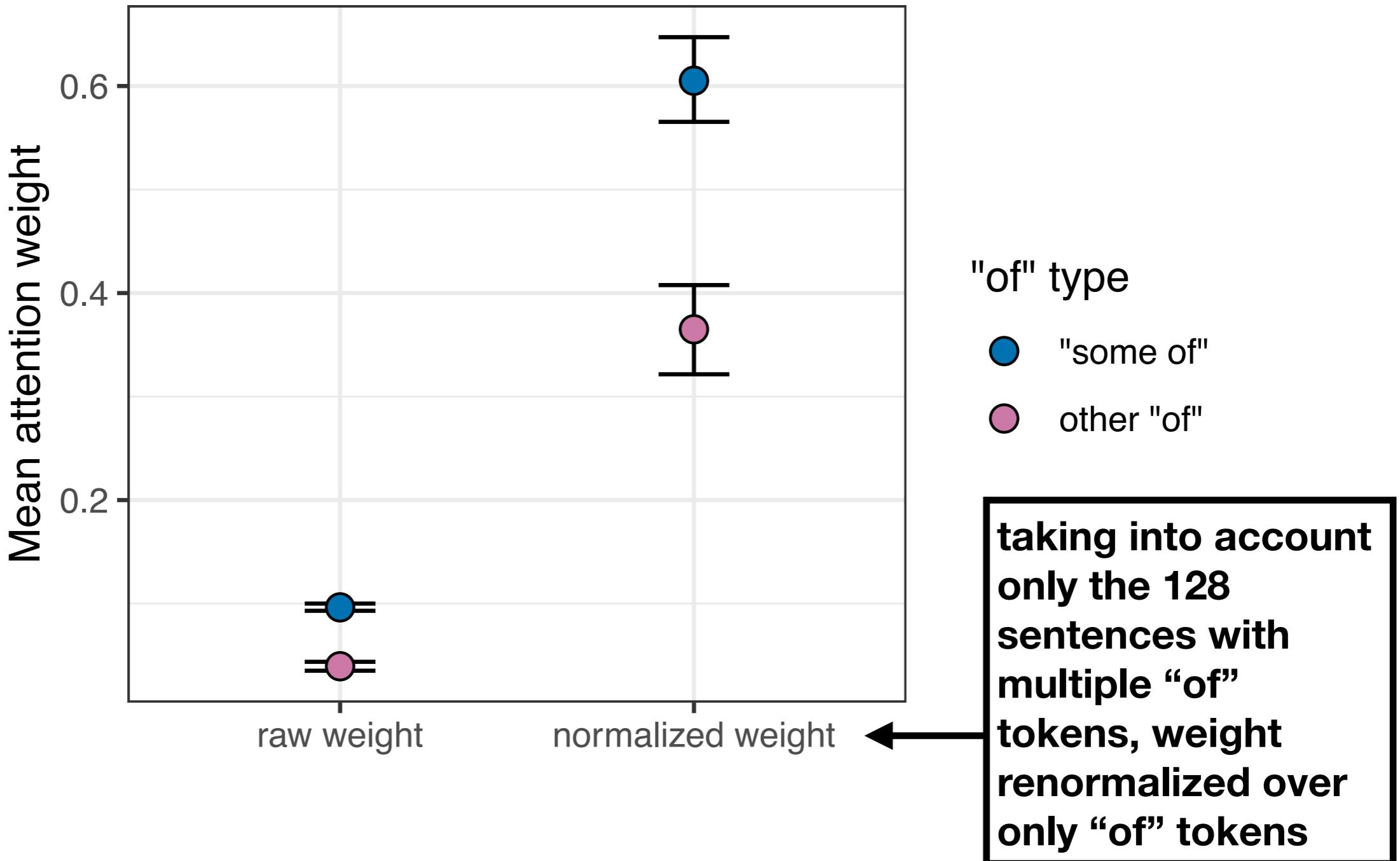


more attention to early positions when “some”-NP is subject
more attention to late positions when “some”-NP is not subject

Attention to “of”



Attention to “of”



more attention to “of” when it’s part of a “some”-NP

Attention weight analysis

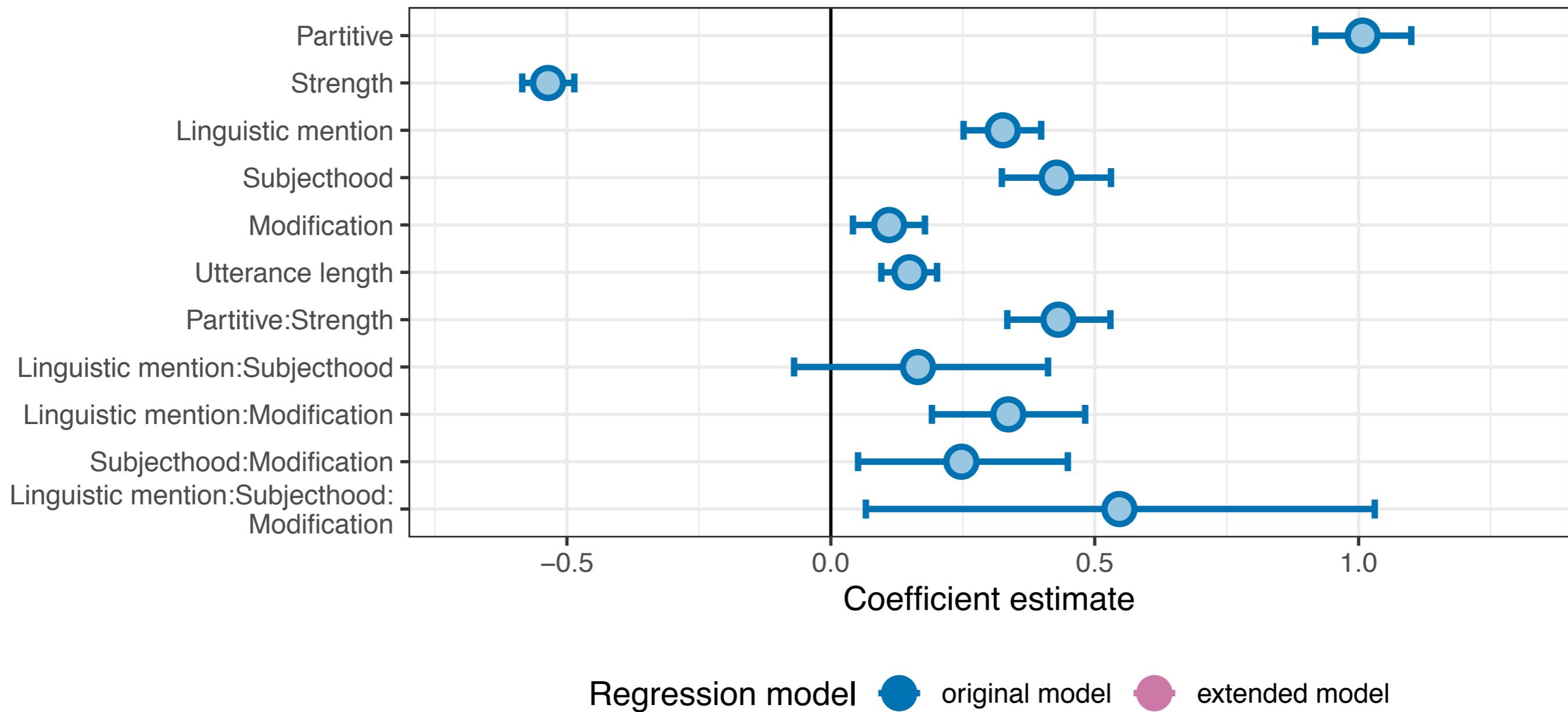
Is there any evidence that the model learned to pay attention to a priori relevant utterance tokens?

Yes!

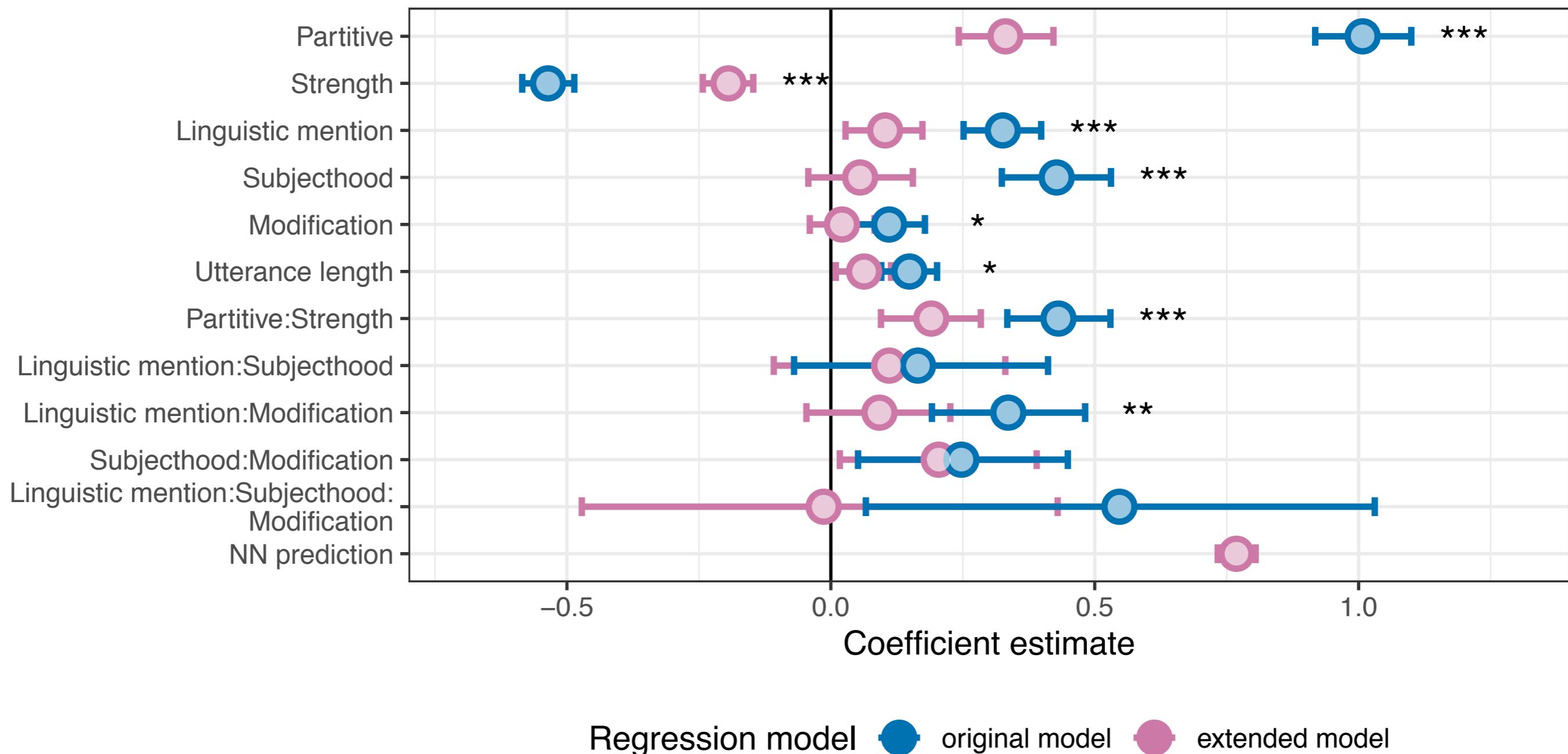
Quantitative analysis

Is there any evidence that the model captures the same effects that the hand-mined feature model did?

Quantitative comparison with hand-mined model



Quantitative comparison with hand-mined model



Quantitative analysis

Is there any evidence that the model captures the same effects that the hand-mined feature model did?

Yes! In fact, most hand-mined feature effects barely survive, and some don't.

Minimal pair analysis

Linzen et al. 2016; Gulordava et al. 2018; Chowdhury and Zamparelli 2018; Marvin and Linzen 2018; Futrell et al. 2019; Wilcox et al. 2019

Is there any evidence that the model can generalize what it learned to entirely new, artificial sentences?

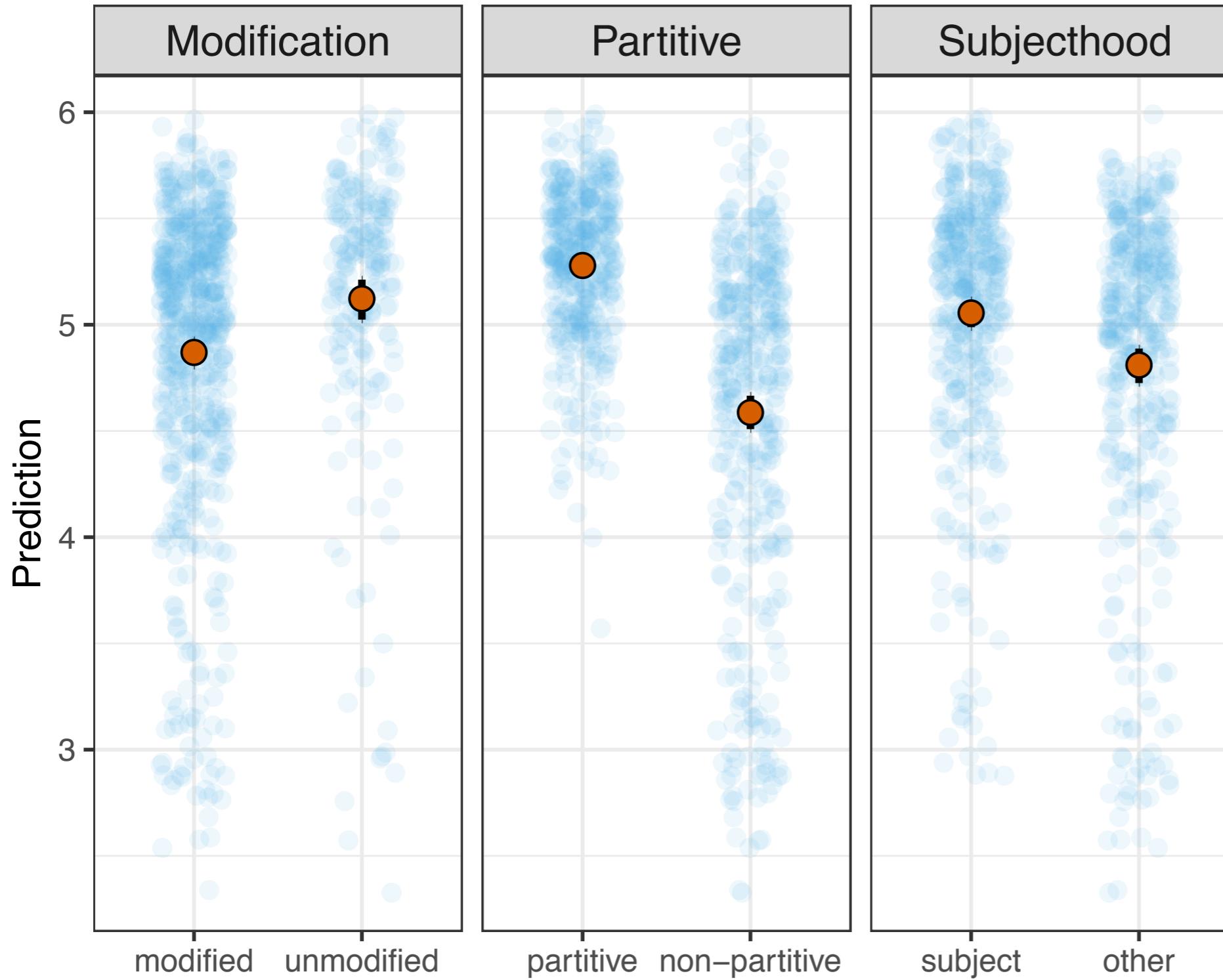
Artificial dataset

Generate sentences that cross factors of interest:
partitive, subjecthood, modification

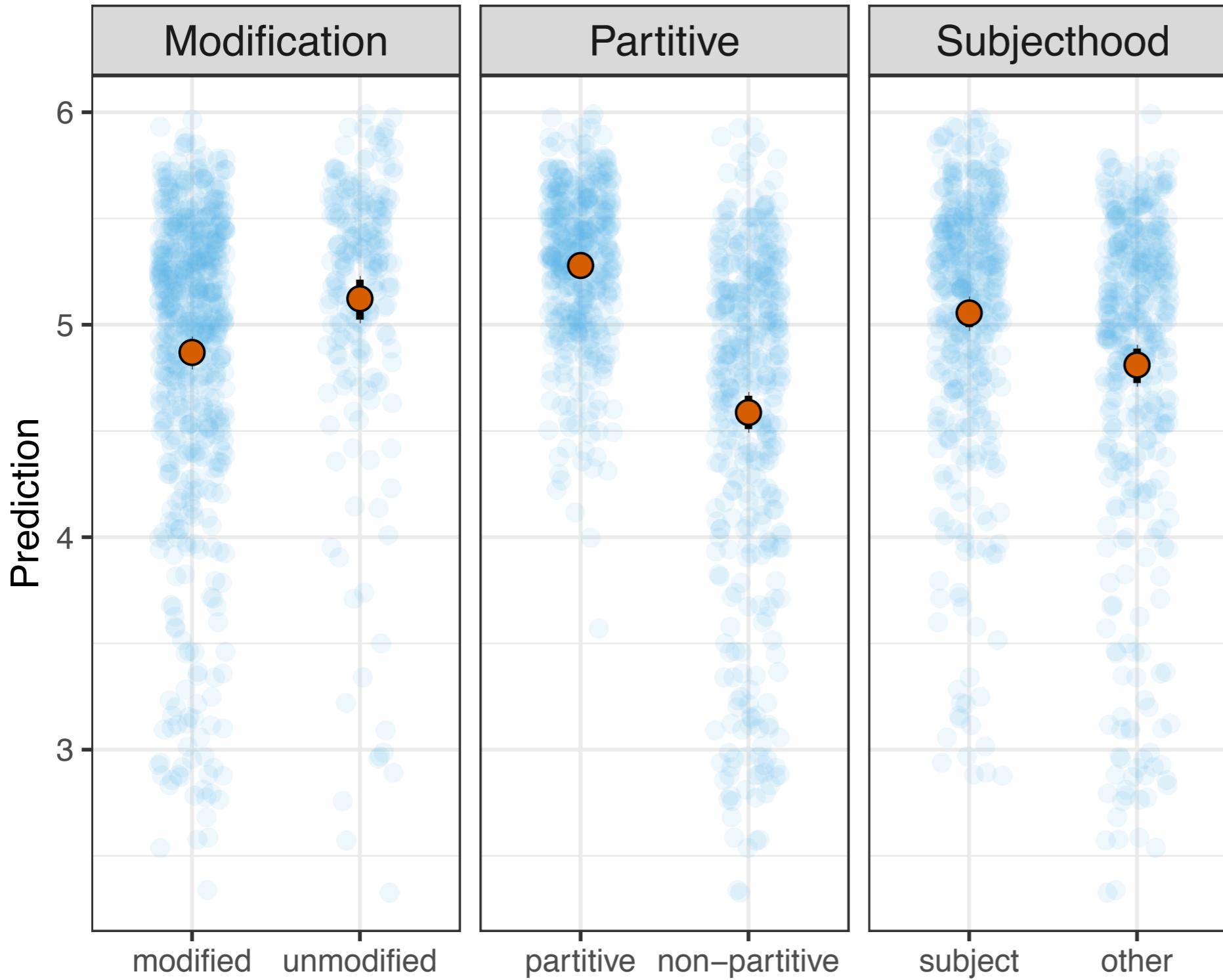
1. **Some (of the) waiters** poured the white wine that my friend really likes.
2. The white wine that my friend really likes was poured by **some (of the) waiters**.
3. The waiters poured **some (of the) white wine that my friend really likes**.
4. **Some (of the) white wine that my friend really likes** was poured by the waiters.
5. Some **attentive** waiters **at the gallery opening** poured the white wine that my friend really likes.
6. ...

25 items, 32 variants of each item = 800 sentences

Qualitative model results



Qualitative model results



the model qualitatively retrieves the partitive, subjecthood, and modification effects on an artificial dataset

Minimal pair analysis

Is there any evidence that the model can generalize what it learned to entirely new, artificial sentences?

Yes!

Context, revisited

Why does the model not learn to use the context beyond the target sentence?

2 possibilities:

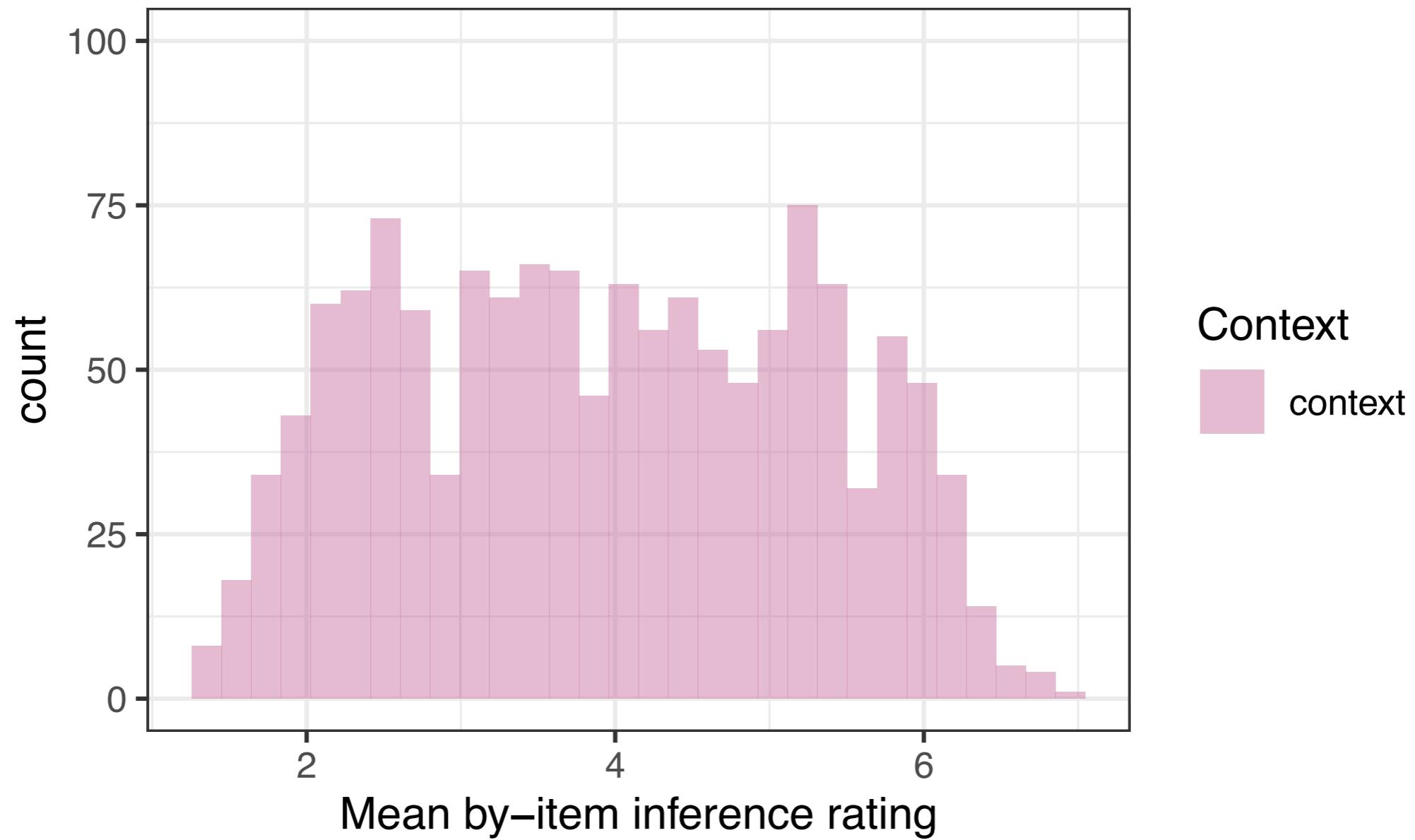
1. humans don't use context in their ratings
2. model has inadequate representation of context

To address: re-ran experiment without displaying context (680 participants, 10 judgments per item).

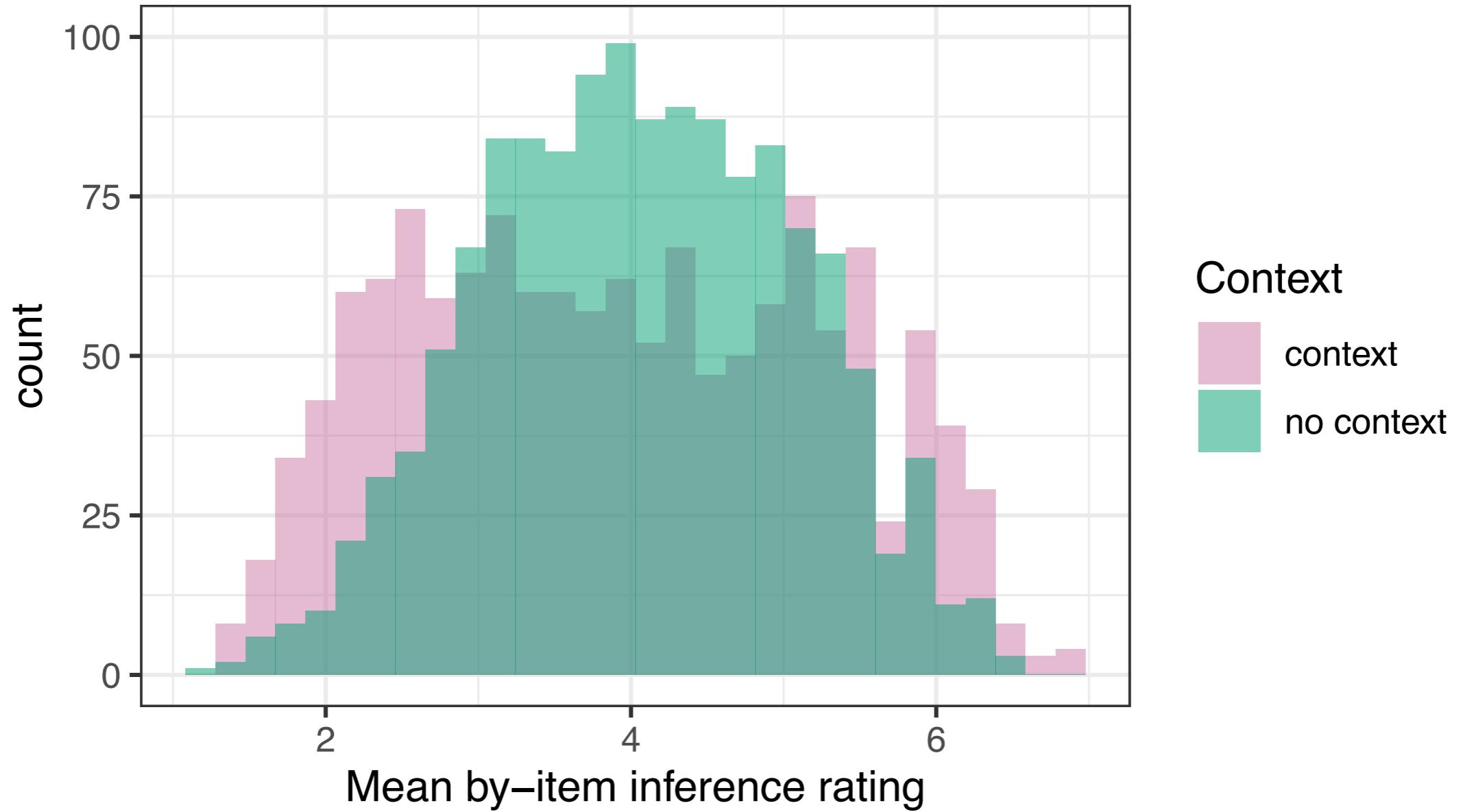
If ratings don't change, 1.

If ratings do change, 2.

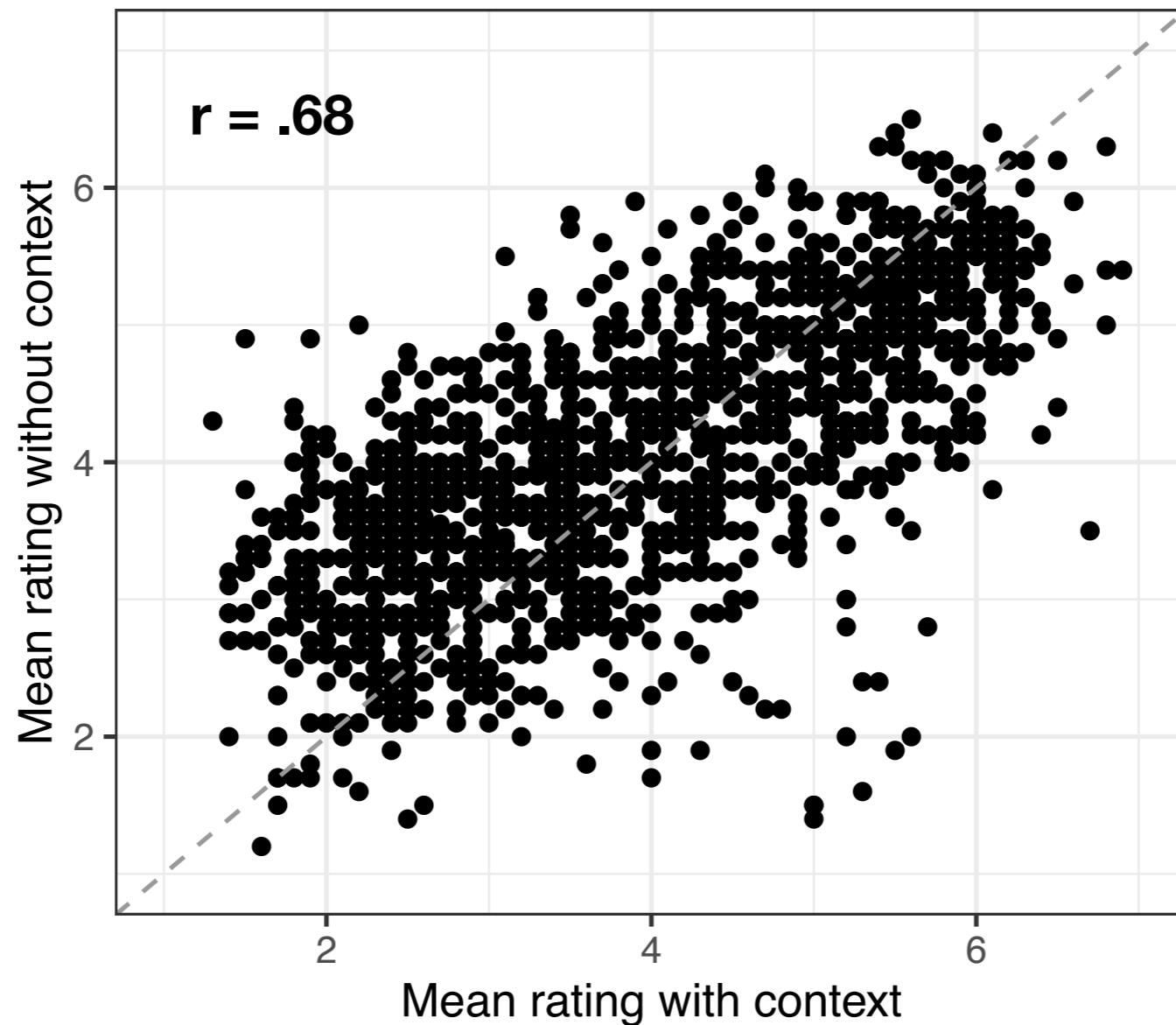
Mean ratings



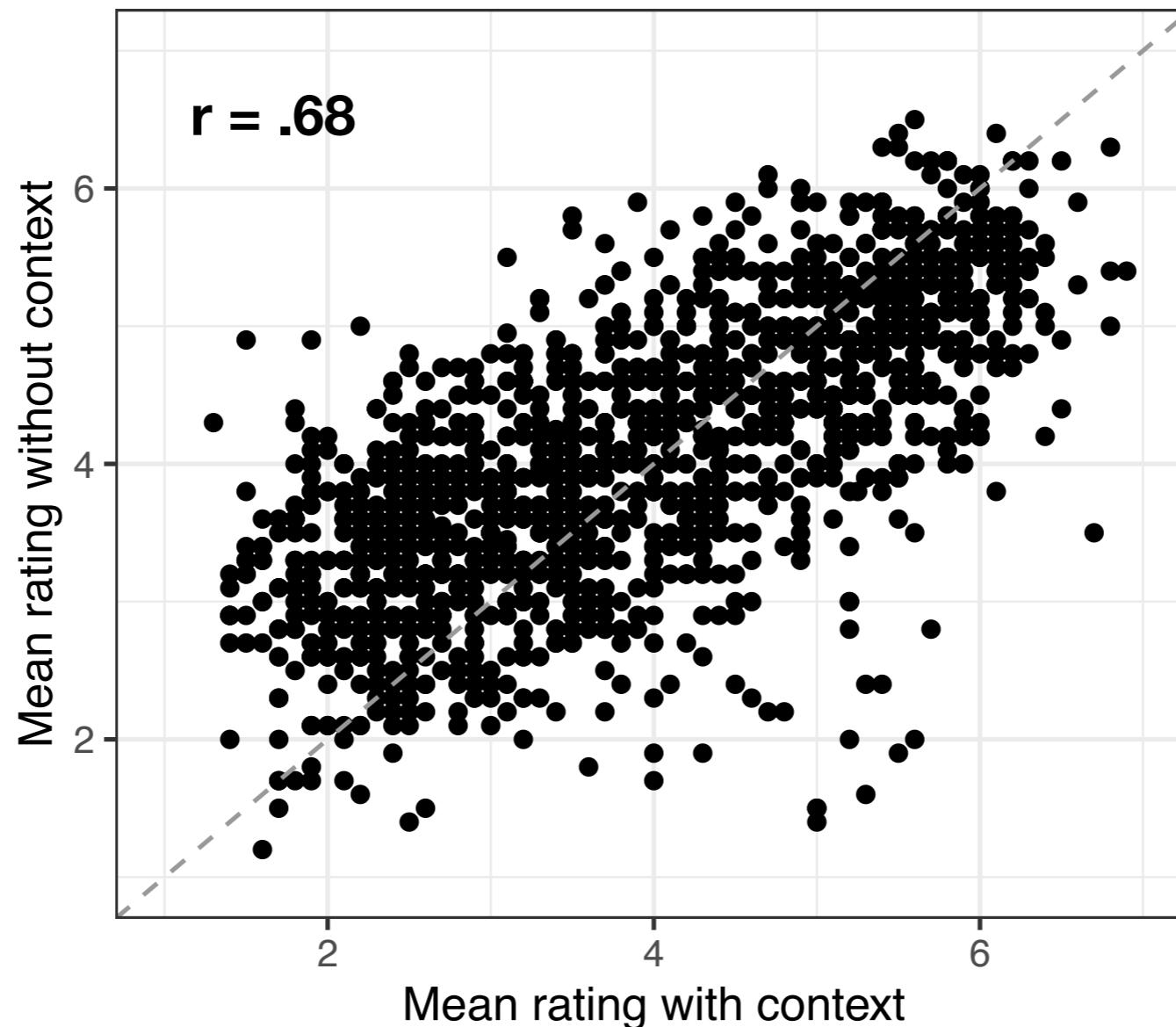
Mean ratings



Mean ratings



Mean ratings



some information about inference is in broader conversational context —> model has poor context representation

Conclusion

There is much more variability in scalar inferences than commonly assumed — but it's systematically context-dependent, and we can capture a lot of it by inspecting the naturalistic signal.

Recent advances in NLP offer a promising avenue for informing pragmatic theory if we can develop good methods for probing the black box neural representations.

Thank you!

Research assistants

Jane Boettcher

Leyla Kursat

Andrew Watts

Collaborators

Yuxing Chen

Sebastian Schuster

