

Differential item function and adaptive tests of early language

Or How I Started to Worry about and Love Measurement

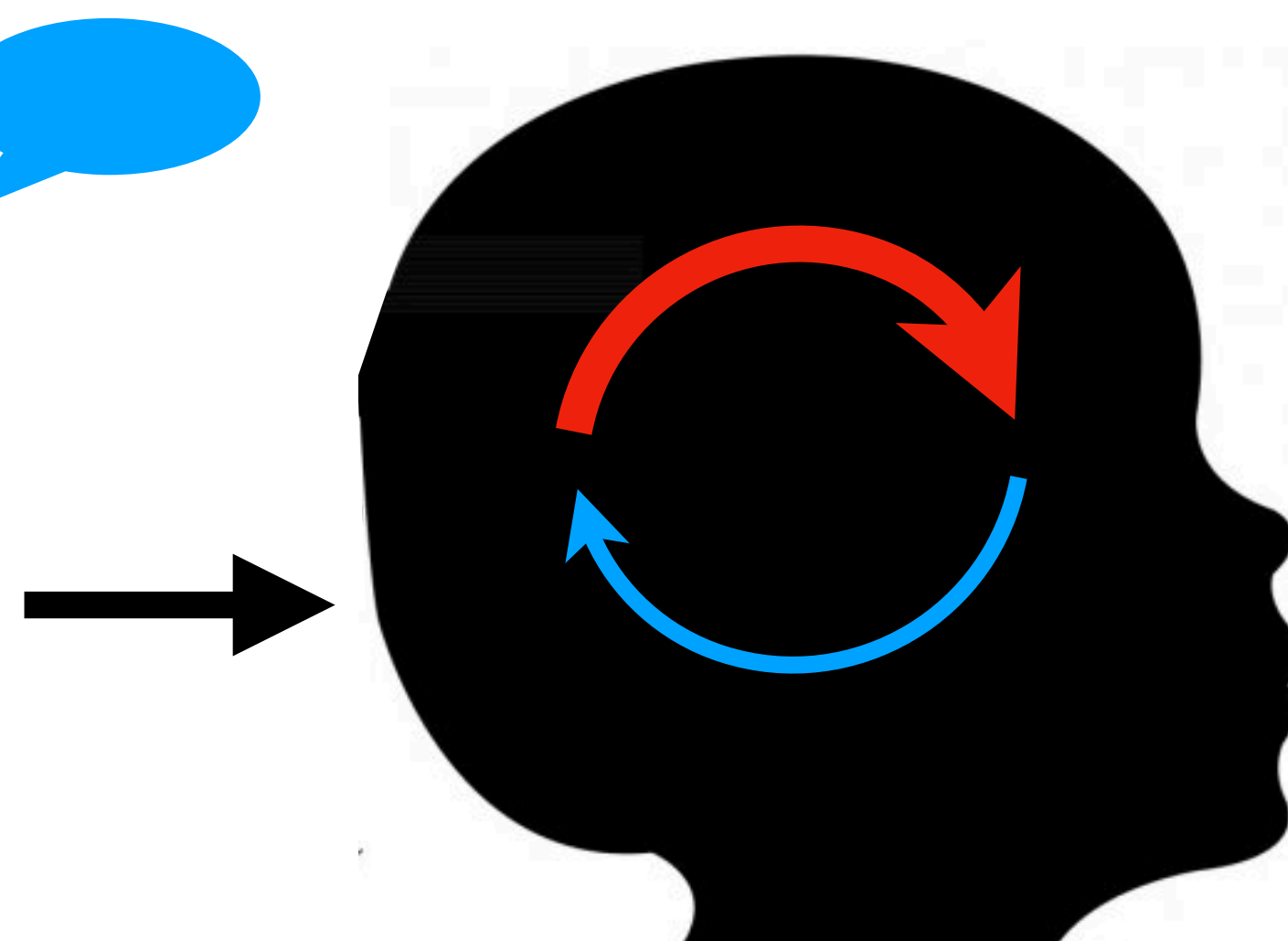
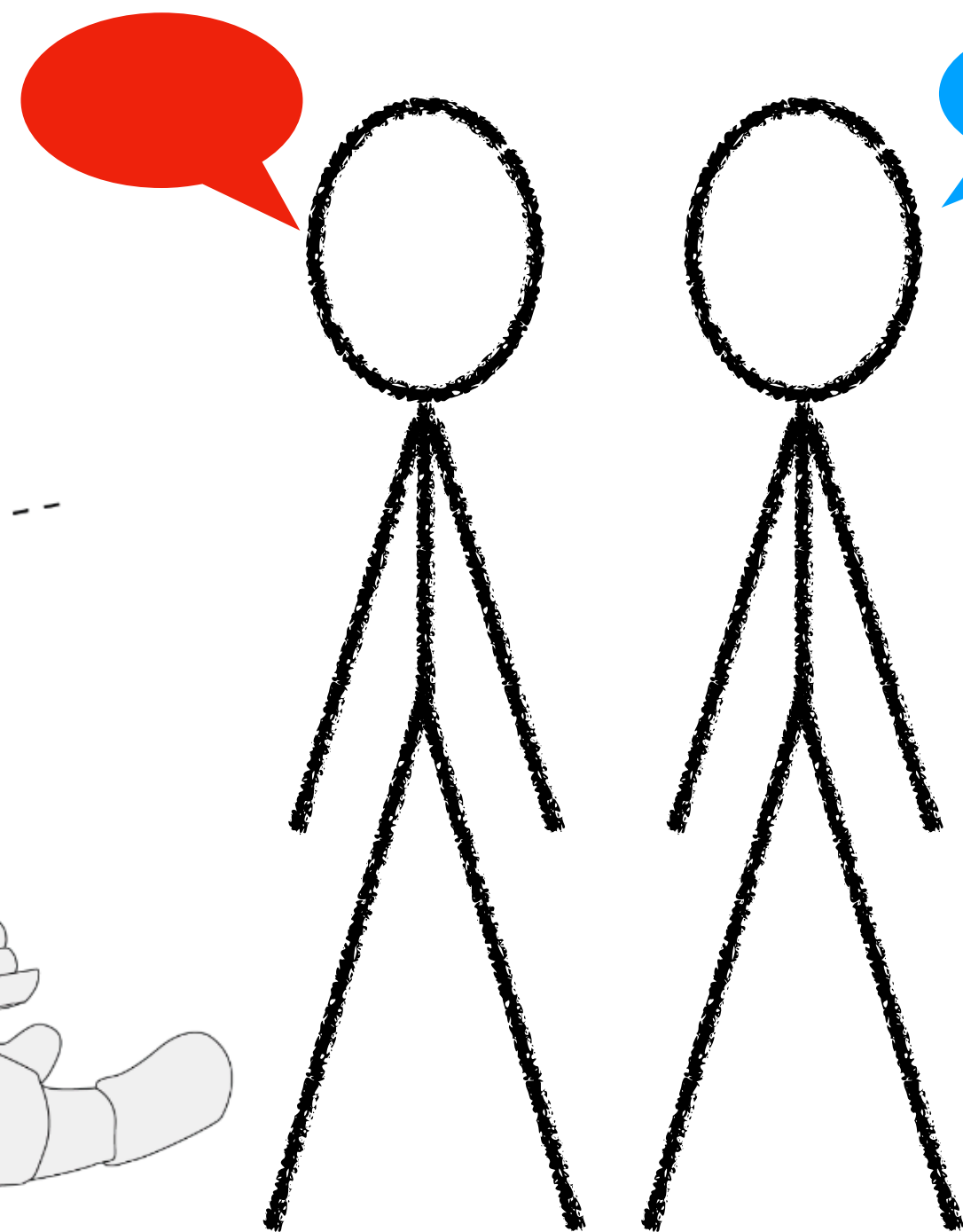
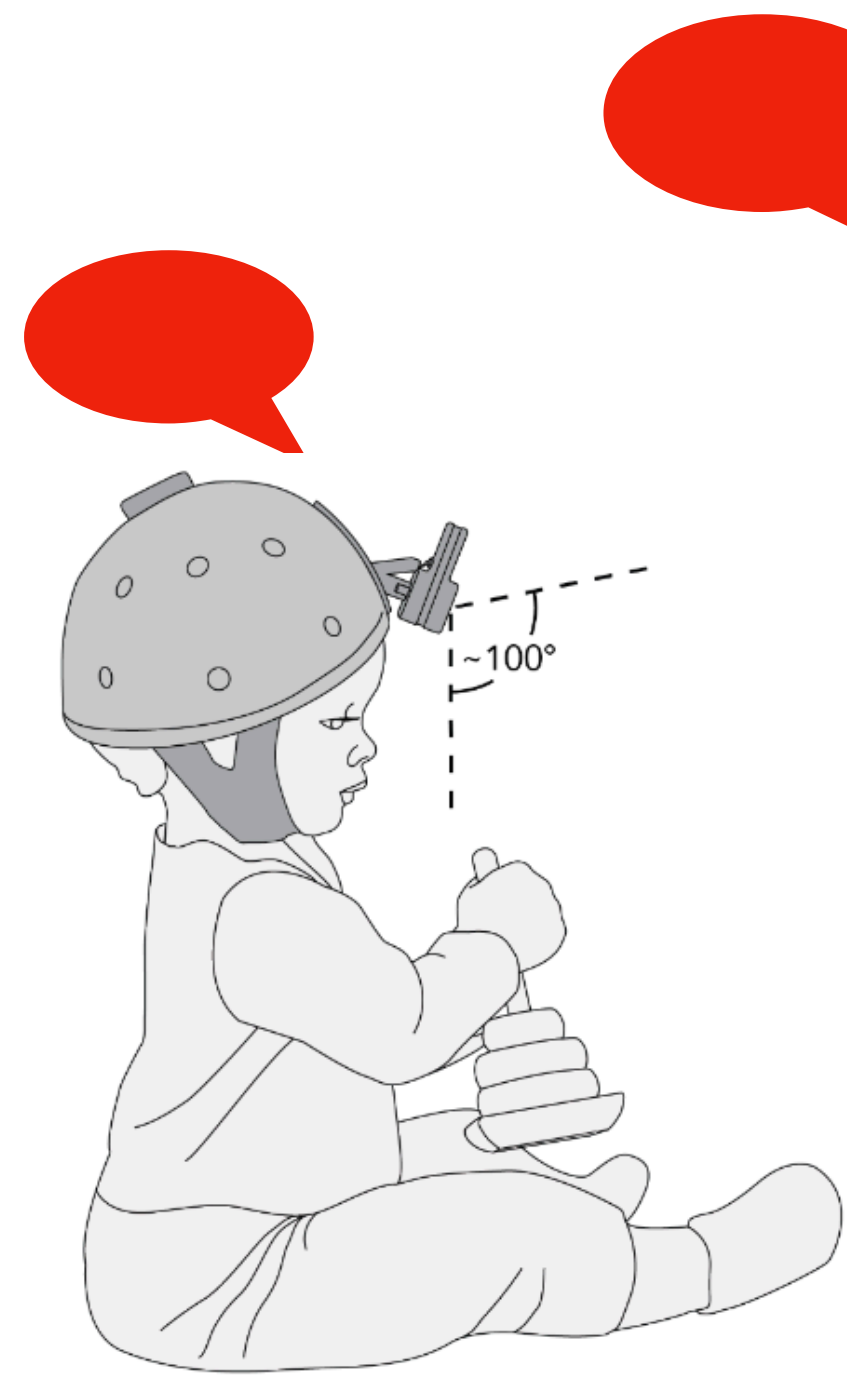
George Kachergis
July 24, 2024

“The history of science is the history of measurement.”
—James M. Cattell (1893), founder of *Psychological Review*

Input:
language environment
Measure early learning environments

Mechanisms:
of processing and learning
Test models (theories) of individual learning

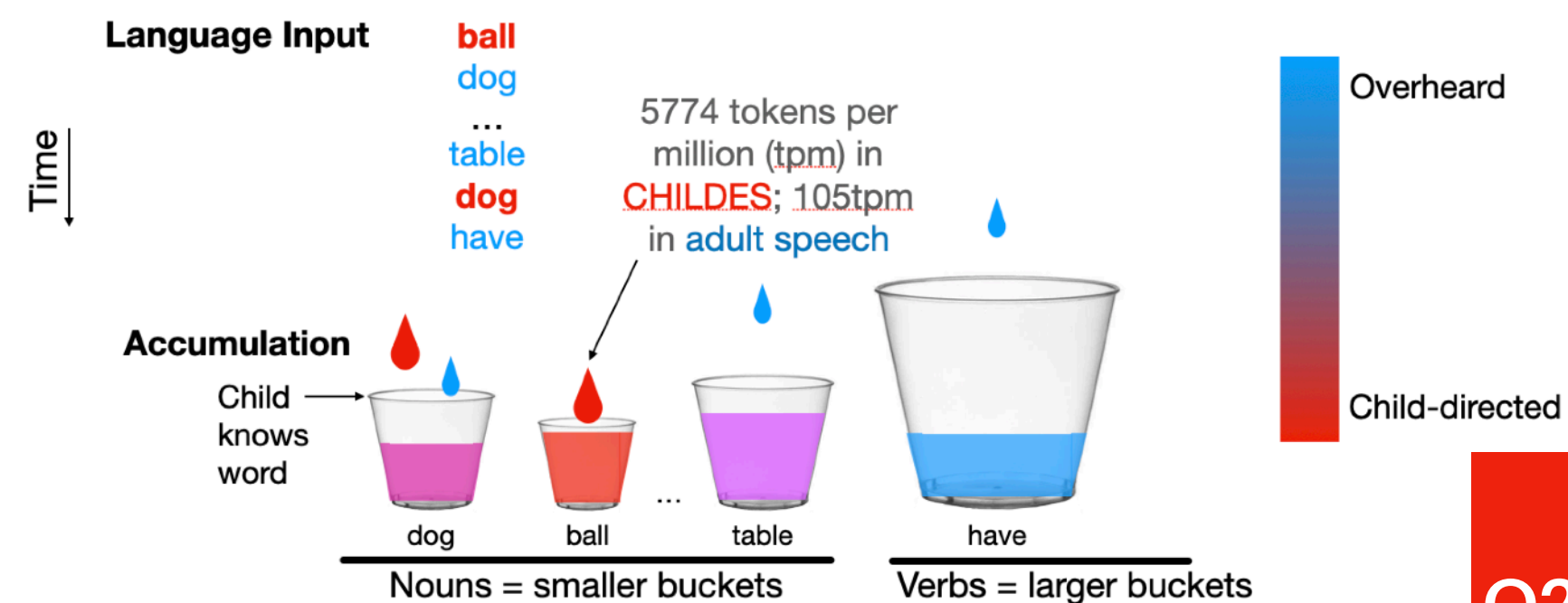
Uptake:
learning outcomes
Create short, valid, and fair tests of early language



Comprehension

Production

Ball!



Q1: How short can tests be?
Q2: How can we assess (& improve) fairness?

The MacArthur-Bates Communicative Development Inventory (CDI)

Parent-report measure of children's early language comprehension and production.

5. Food and Drink (30)

	Understands	Understands and says		Understands	Understands and says
apple	<input checked="" type="radio"/>	<input type="radio"/>	cheerios	<input type="radio"/>	<input type="radio"/>
banana	<input type="radio"/>	<input type="radio"/>	cheese	<input type="radio"/>	<input type="radio"/>
bread	<input type="radio"/>	<input type="radio"/>	chicken	<input type="radio"/>	<input type="radio"/>
butter	<input type="radio"/>	<input type="radio"/>	coffee	<input type="radio"/>	<input type="radio"/>
cake	<input type="radio"/>	<input checked="" type="radio"/>	cookie	<input type="radio"/>	<input type="radio"/>
candy	<input type="radio"/>	<input checked="" type="radio"/>	cracker	<input type="radio"/>	<input type="radio"/>
carrots	<input type="radio"/>	<input type="radio"/>	drink	<input type="radio"/>	<input type="radio"/>
cereal	<input type="radio"/>	<input type="radio"/>	egg	<input type="radio"/>	<input type="radio"/>




Impressive reliability and predictive validity

Fenson et al. (1994, 2007)


The MacArthur-Bates Communicative Development Inventory (CDI)

Wordbank Data Contributors Announcements About Us FAQ



Wordbank

An open database of children's vocabulary development




Questionnaire Mots et Gestes

USE NO. 2 PENCIL ONLY

I WORDS C

they say. We are particula
re heard your child use. If v

Technical Manual of the Japanese
MacArthur Communicative
Development Inventory:
Words and Grammar

日本語著者 綿巻 徹・小松たみ子

Japanese

Palabras y Enunciados
(Inventario II)

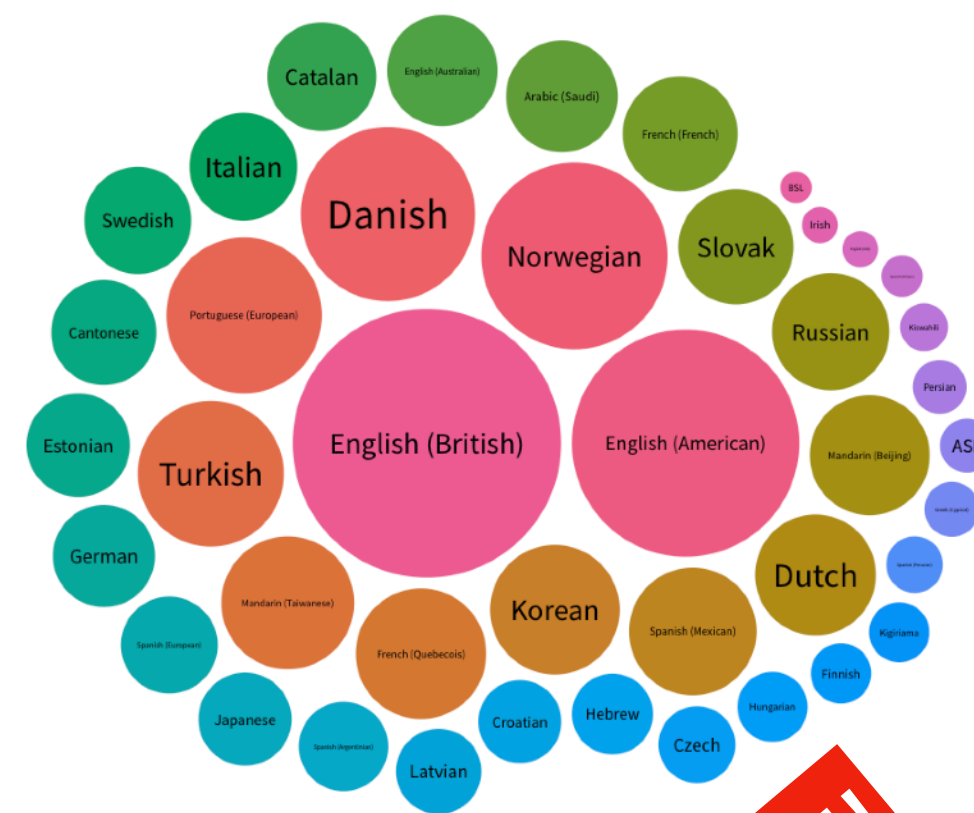
Donna Jackson-Maldonado, Ph.D.,
Elizabeth Bates, Ph.D., y Donna J. Thal, Ph.D.

Spanish

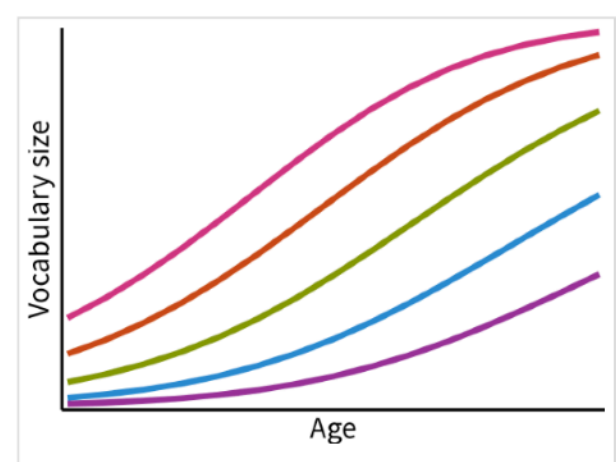
Polish

Slovak

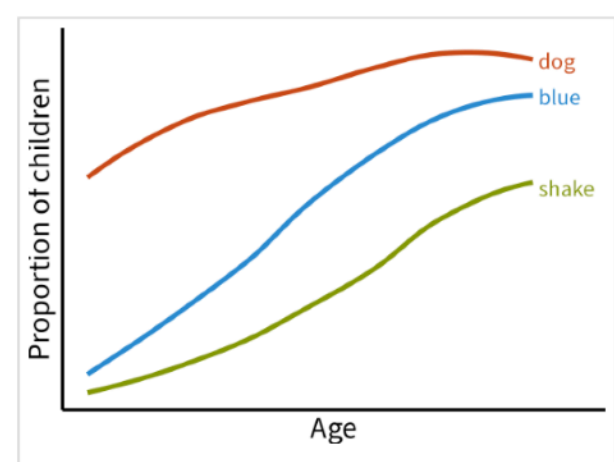
Wordbank contains data from 92,771 children and 105,290 CDI administrations, across 42 languages and 89 instruments:



2024!



Vocabulary Norms
Explore vocabulary size growth curves for various languages and demographic groups.



Item Trajectories
Explore developmental trajectories of individual words in various languages.

<https://wordbank.stanford.edu>

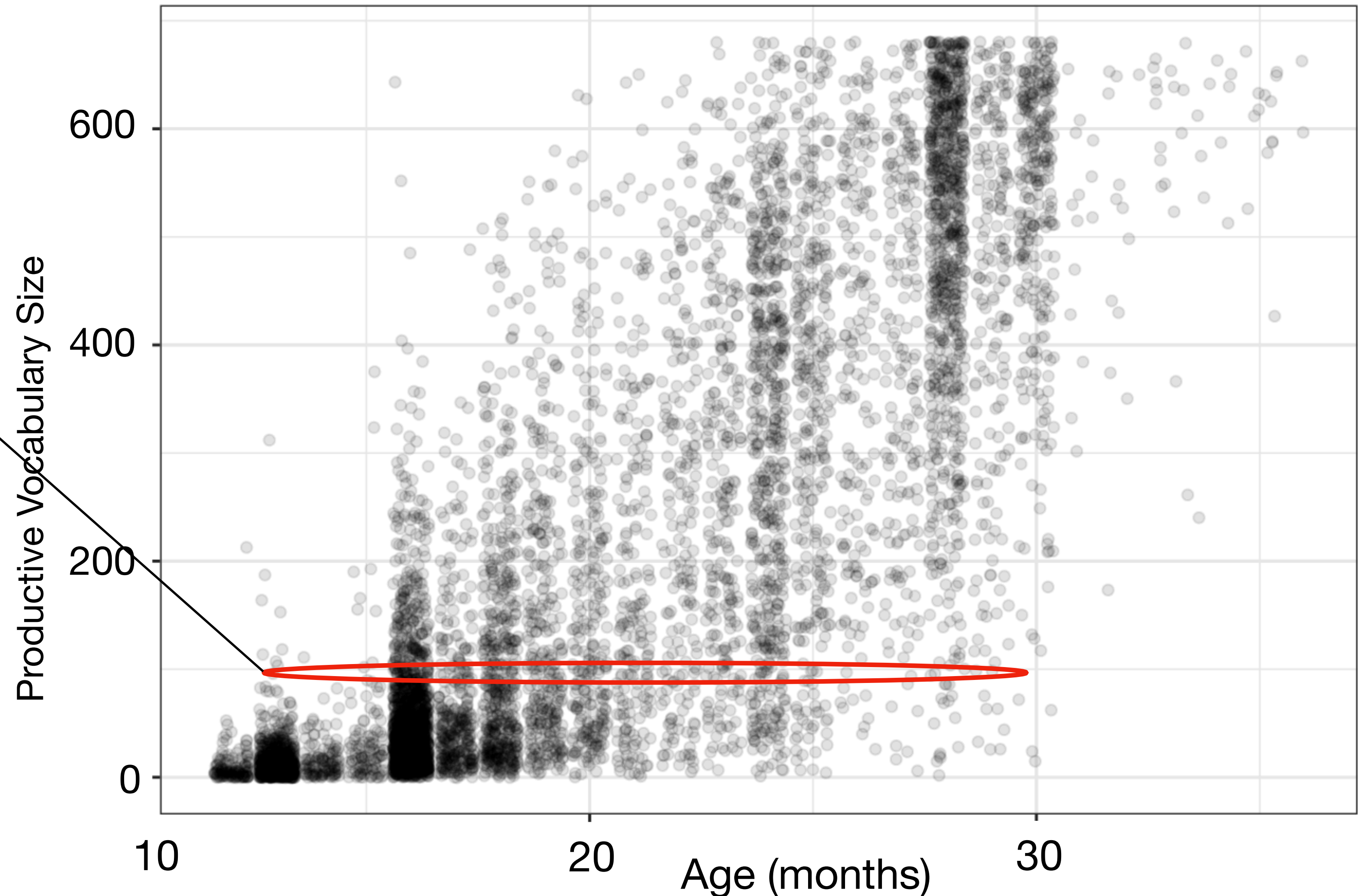
Frank, Braginsky, Yurovsky, & Marchman (2016)

Fenson et al. (1994, 2007)



Beyond Vocabulary Size

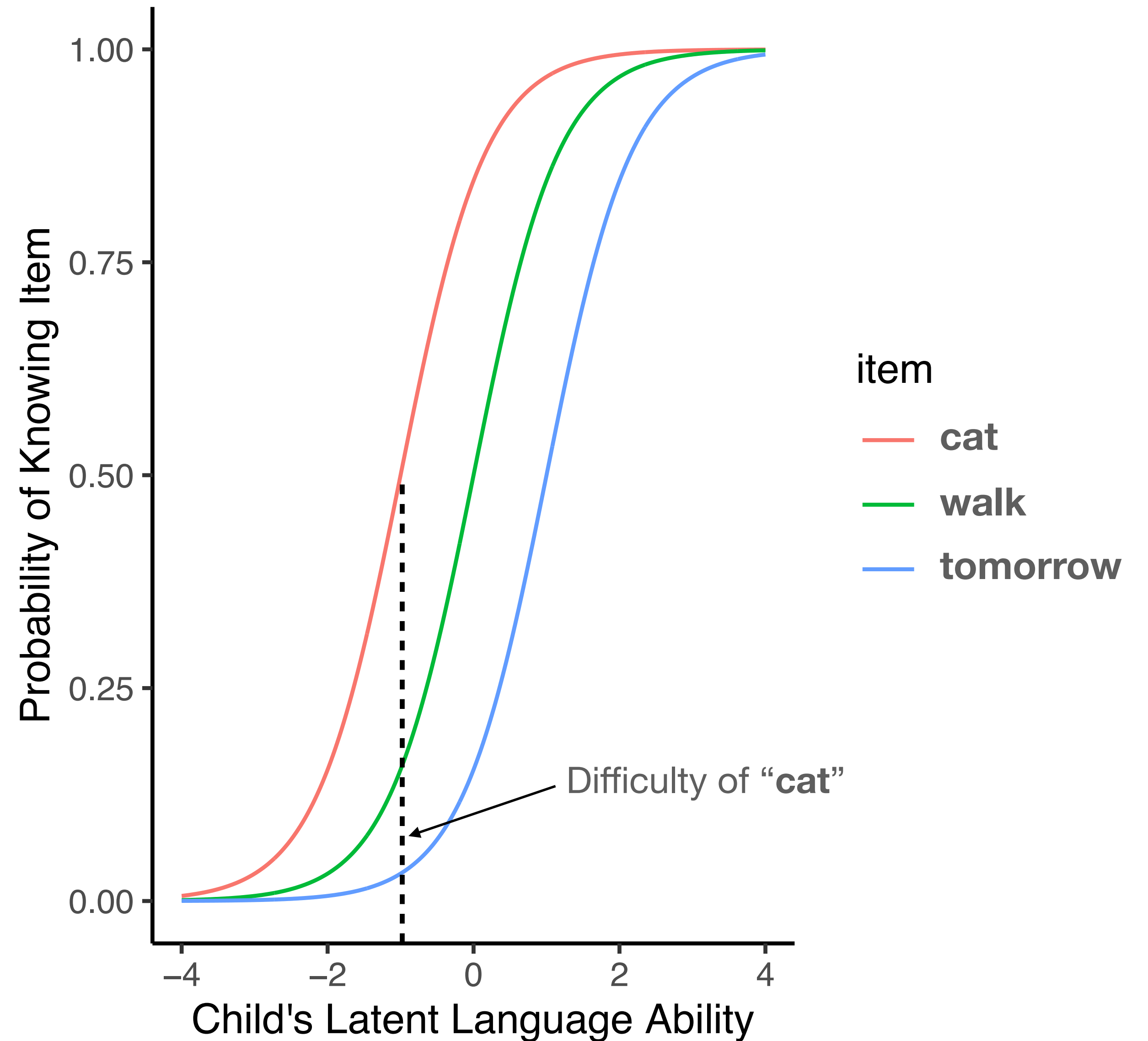
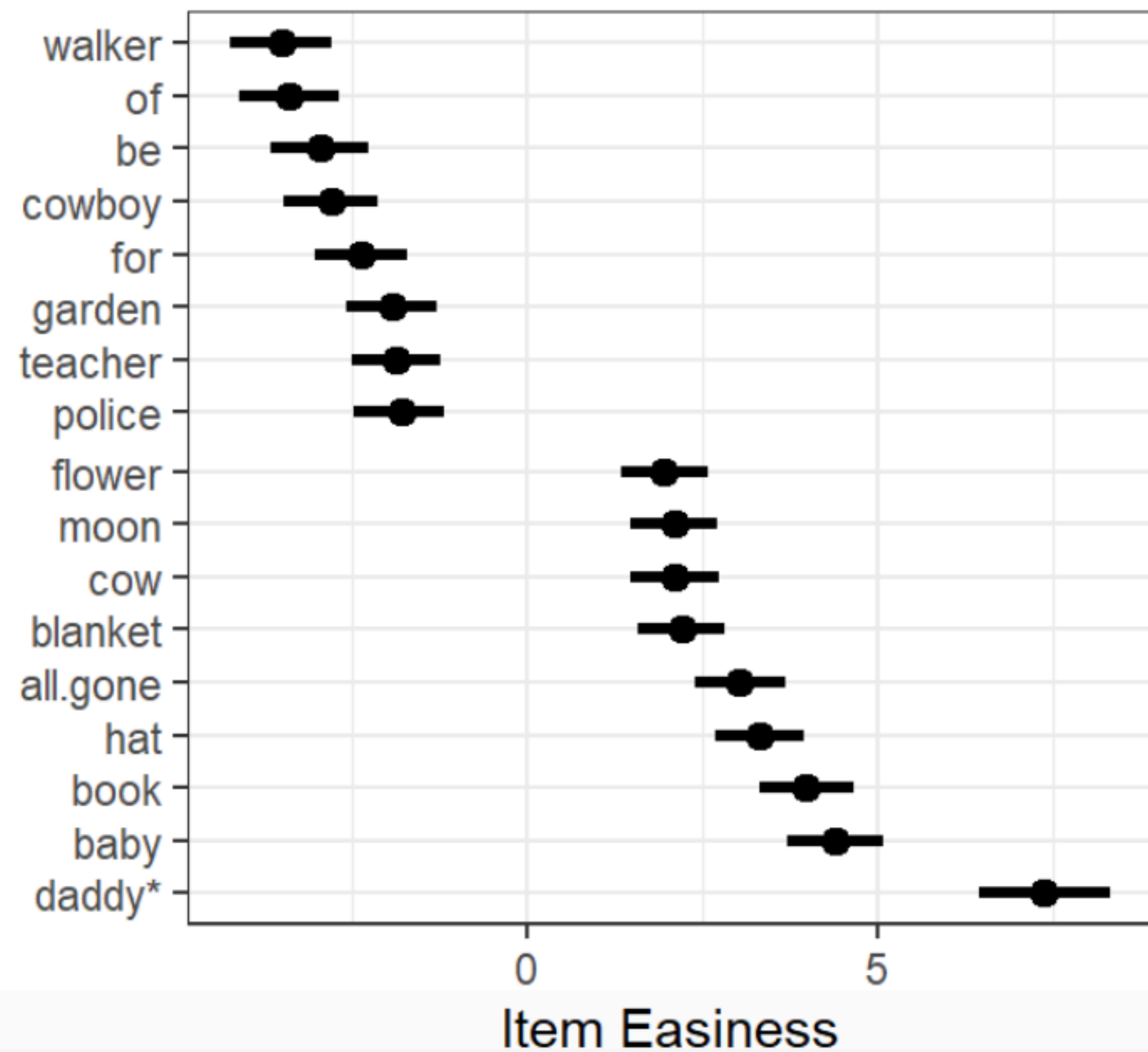
- Vocabulary size of 7,000 English-speaking children
- Do children with equal vocab size (e.g. 100) all have the same language ability?
- Words have different underlying difficulties (e.g., conceptual, phonological)
- Go beyond classical test theory ('sumscore') with psychometric models



Item-Response Theory (IRT)

- Jointly estimate 1) a latent ability θ , for each child j , and 2) a difficulty b_i for each item (word) i

Example Easiness Parameters



Computerized Adaptive Test: CDI-CAT



Traditional CDI

680 items
>20 minutes

Q1 Q2 Q3 Q4 Q5
Q6 Q7 Q8 Q9 Q10
... Q680

Sumscore

Can we create short, valid tests of early language?

Computerized Adaptive Test: CDI-CAT



Traditional CDI

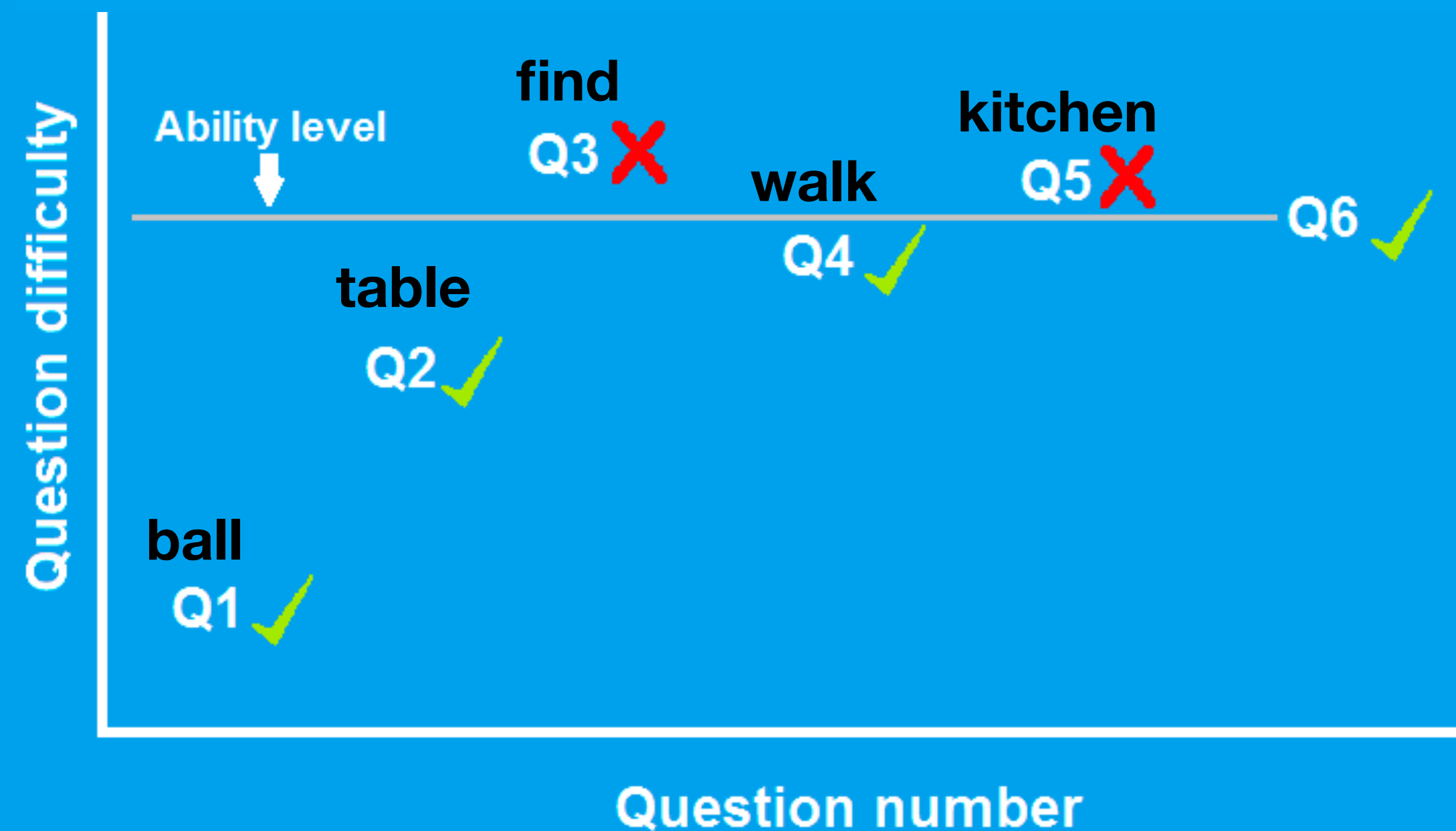
680 items
>20 minutes

Q1	Q2	Q3	Q4	Q5
Q6	Q7	Q8	Q9	Q10
				... Q680

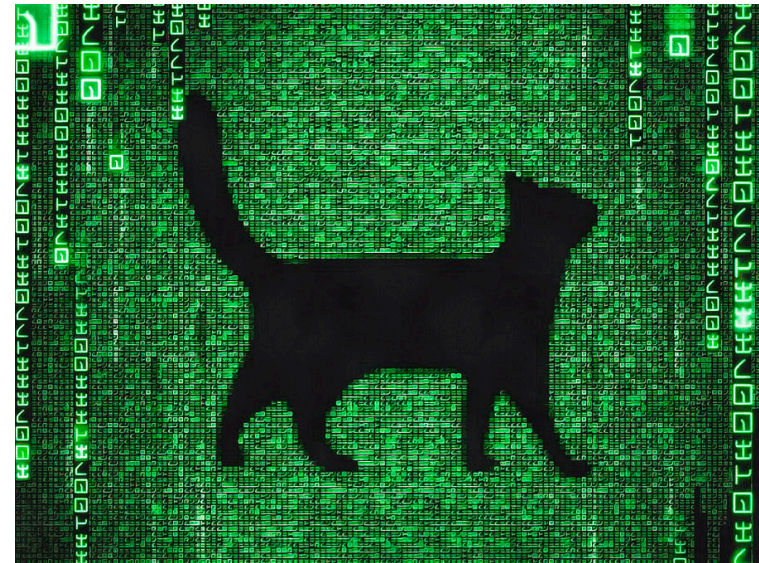


Sumscore

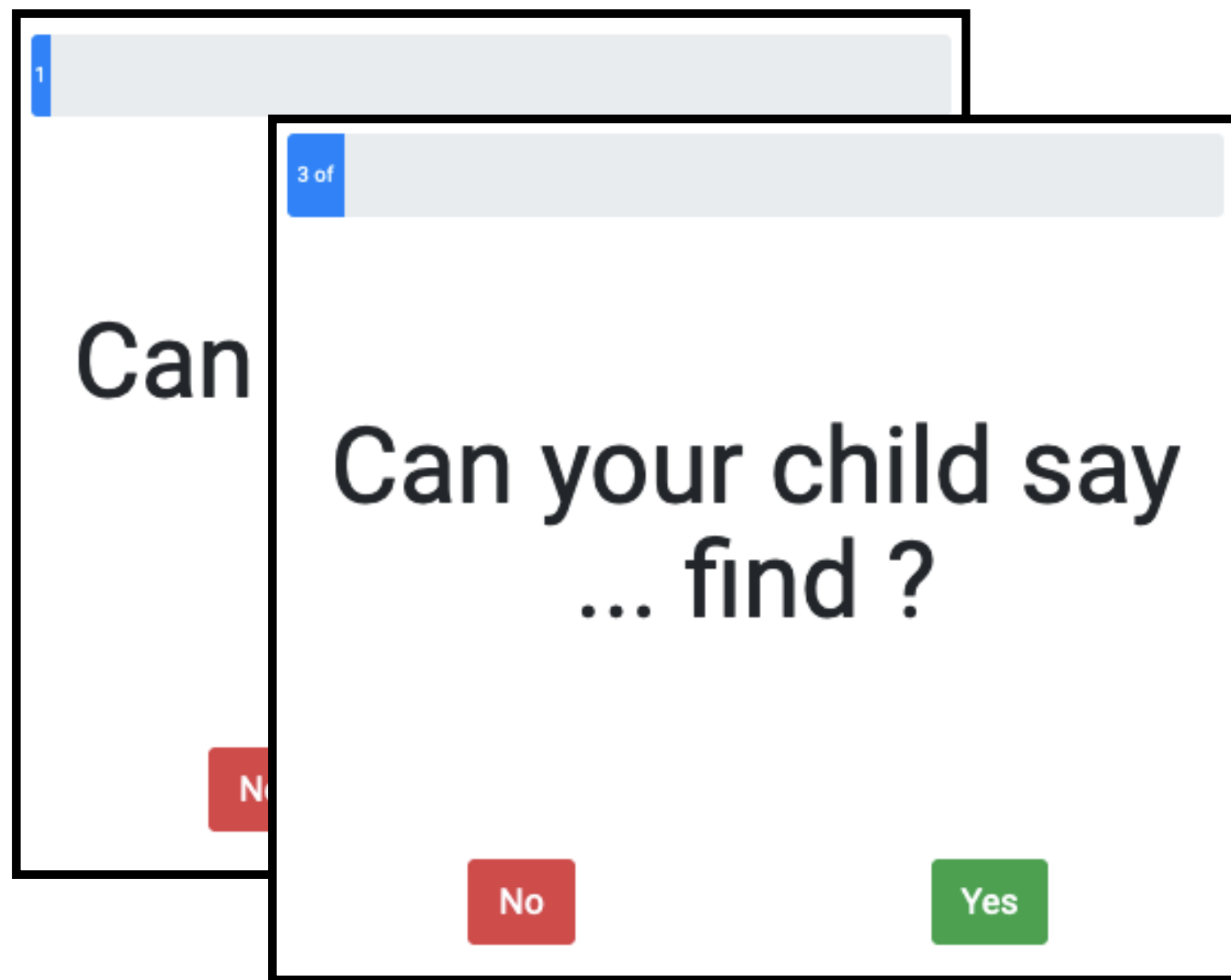
Computerized Adaptive Test (CAT)



Computerized Adaptive Test: CDI-CAT

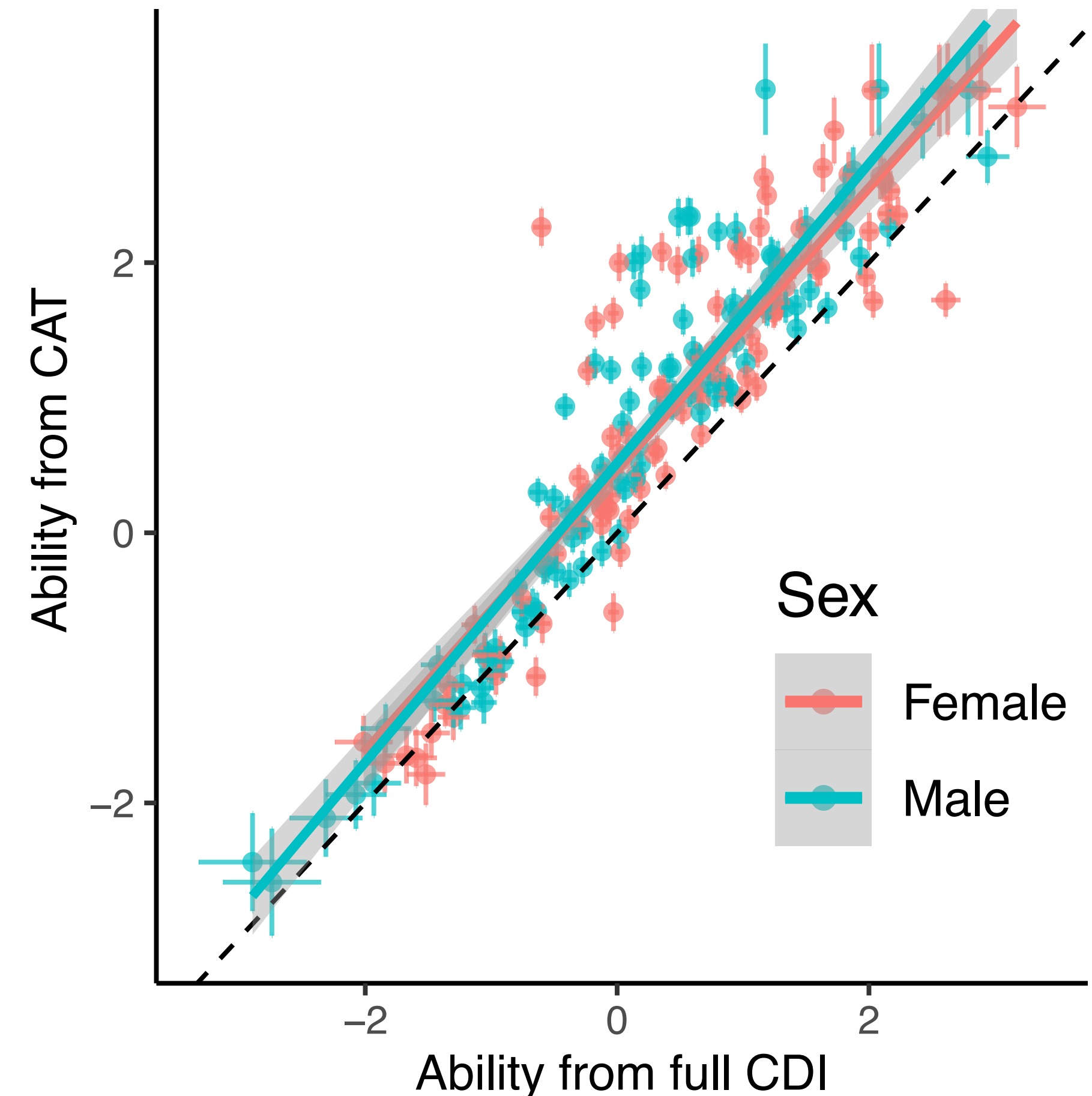


Computerized Adaptive Test (CAT) with 25-50 questions
(vs. >650) reliable across a wide age range (12-36 mos)



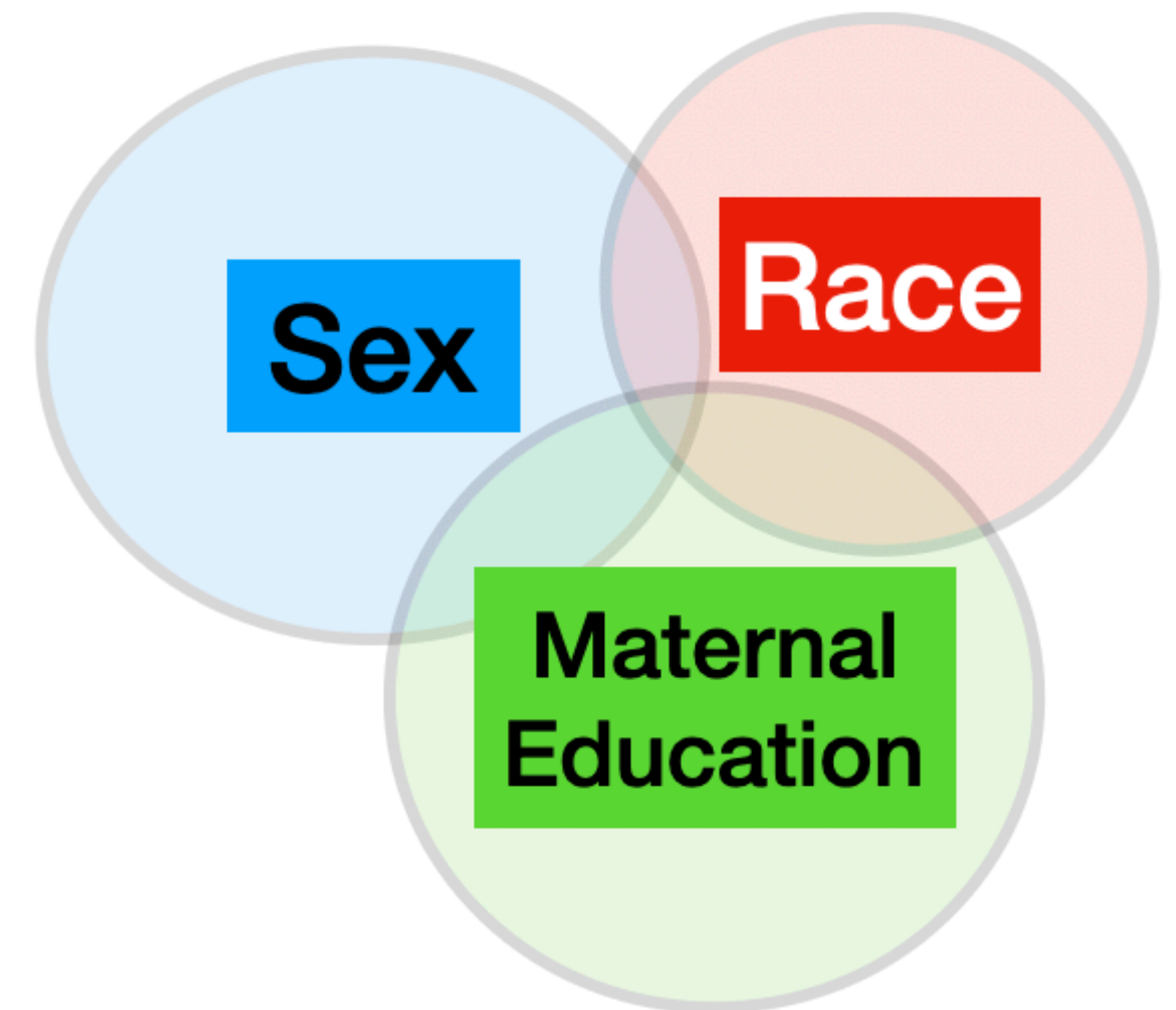
- Validation study (N=200) showed a strong association between full CDI & CAT ability ($r = 0.92$)
- Comprehension *and* production
- English & Spanish in NIH Baby Toolbox (& WebCDI)
- Added French (2023) & Japanese (2024)

<https://webcdi.stanford.edu/>



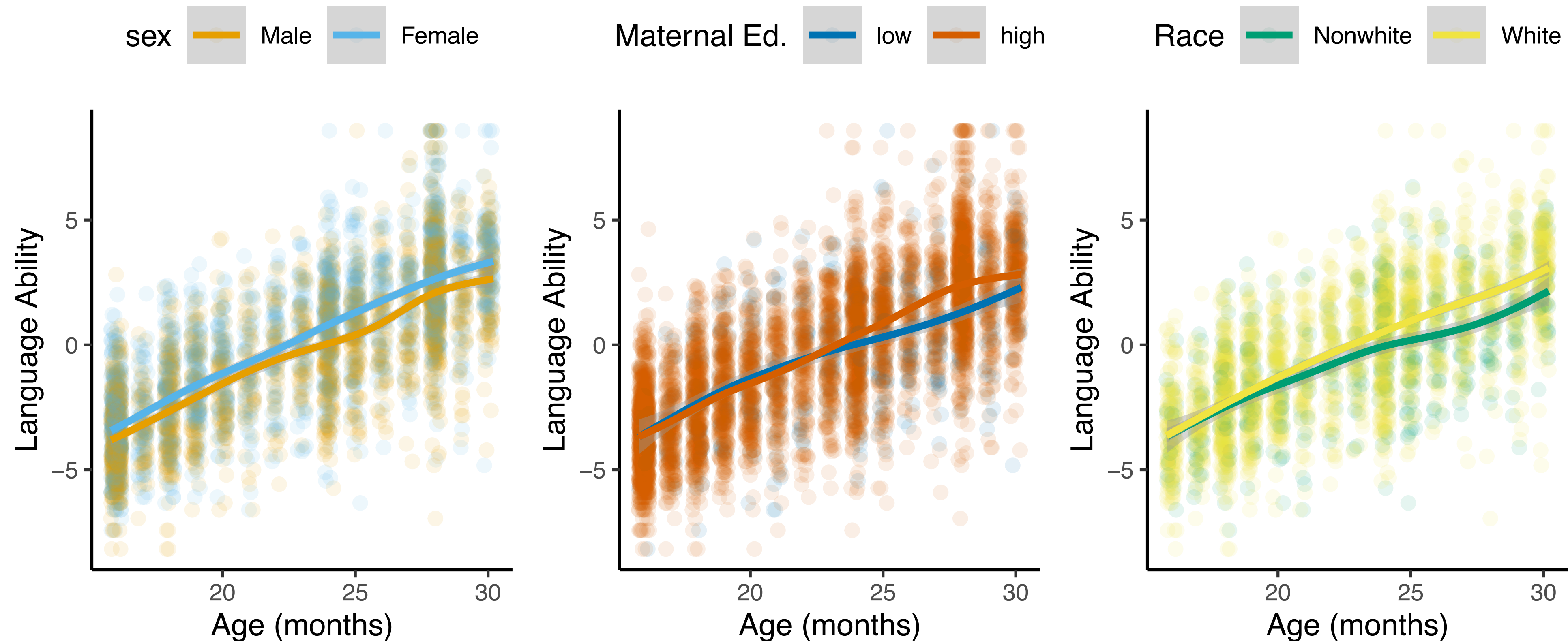
Identifying Measurement Bias

- Wordbank CDI data show demographic differences in vocabulary size advantaging 1) females, 2) white children, and 3) children of highly-educated mothers (a proxy for SES) (Eriksson et al., 2012 Frank et al., 2021)
- Sex-related differences in language skill that persist until high school (Peterson, 2018)
- Could the set of CDI items be biased?
- How to evaluate CDI items for potential bias?



Identifying Measurement Bias

- Ability vs. age by demographic group in a baseline Rasch (1PL) model



- Using a multigroup Rasch model, we estimate item difficulties that are allowed to vary by demographic group. DIF exists if the difference between group parameters is non-zero.
- The question becomes: How many items are (significantly) biased in favor of each group? Is the item bank representative of all possible items?

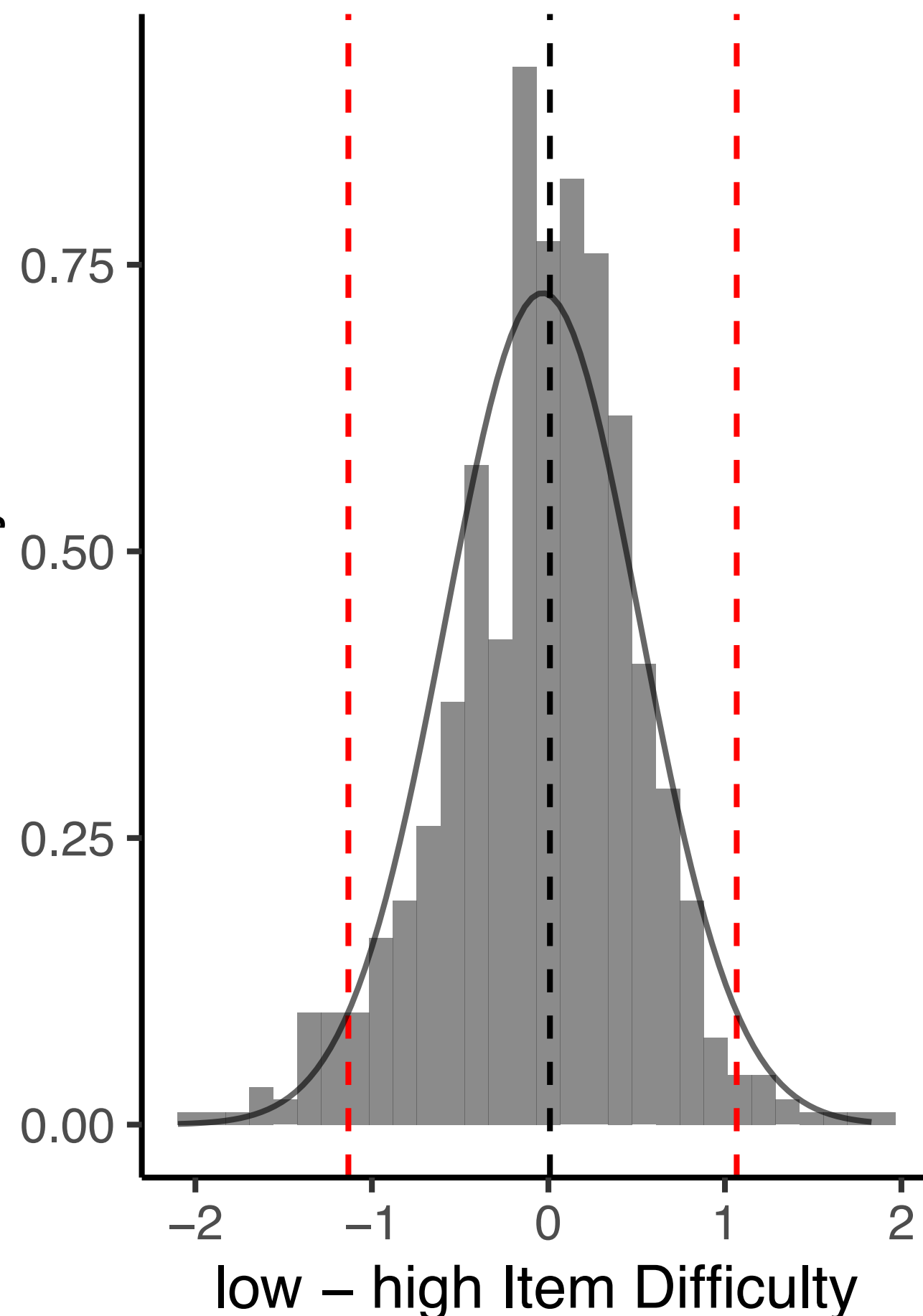
Differential Item Function

- DIF can decrease the validity of a test: imagine groups A & B have no mean difference in ability, but some items are easier for group A (e.g., farm equipment for rural children).
- If many of these items are selected to be on a test, the test will overestimate the ability of those in group A, and underestimate the ability of those in group B: the test is unfair.
- Of course, a true ability difference may exist between groups — regardless, the selection of items can either inflate or deflate the actual group difference.

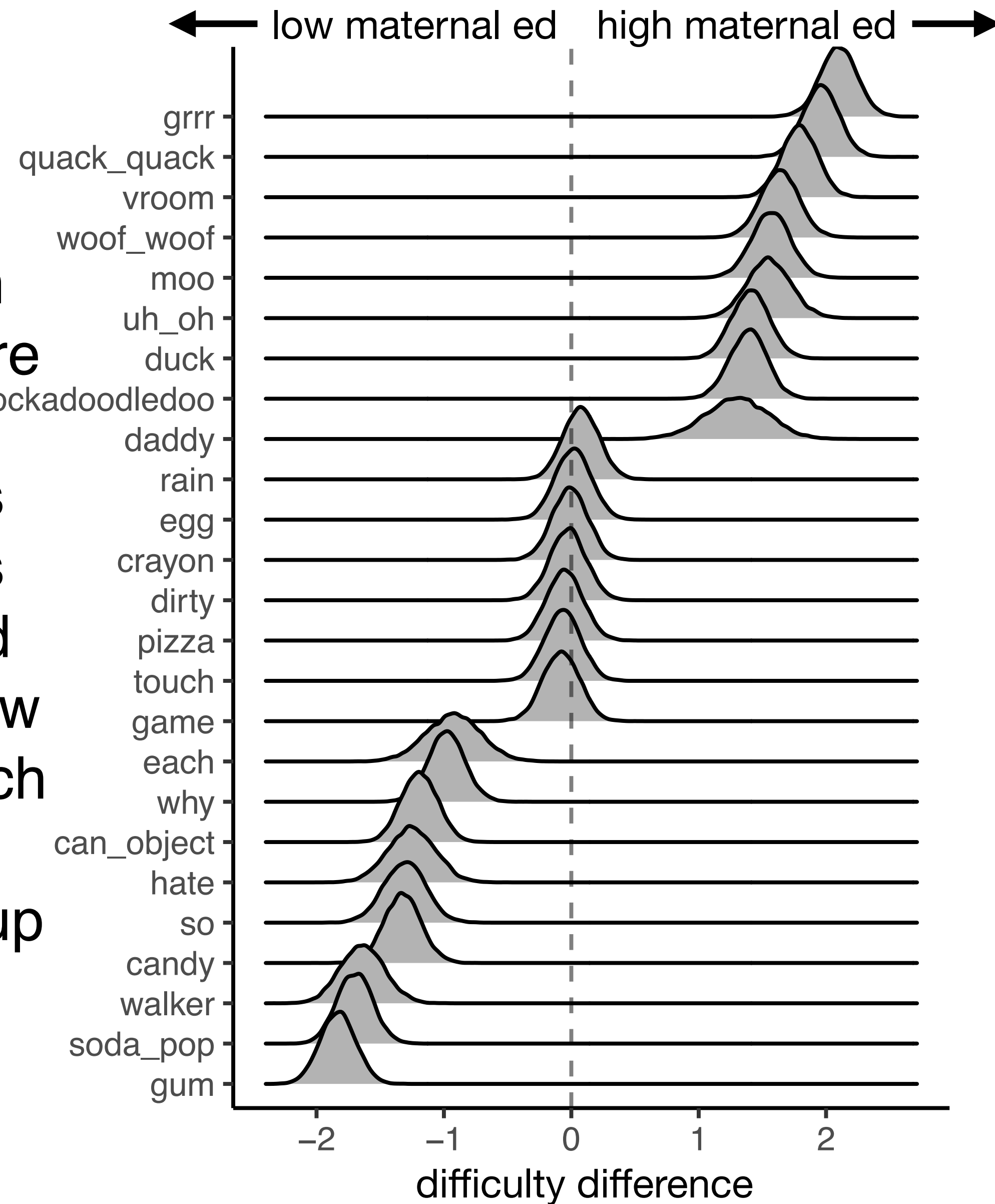


Identifying Measurement Bias

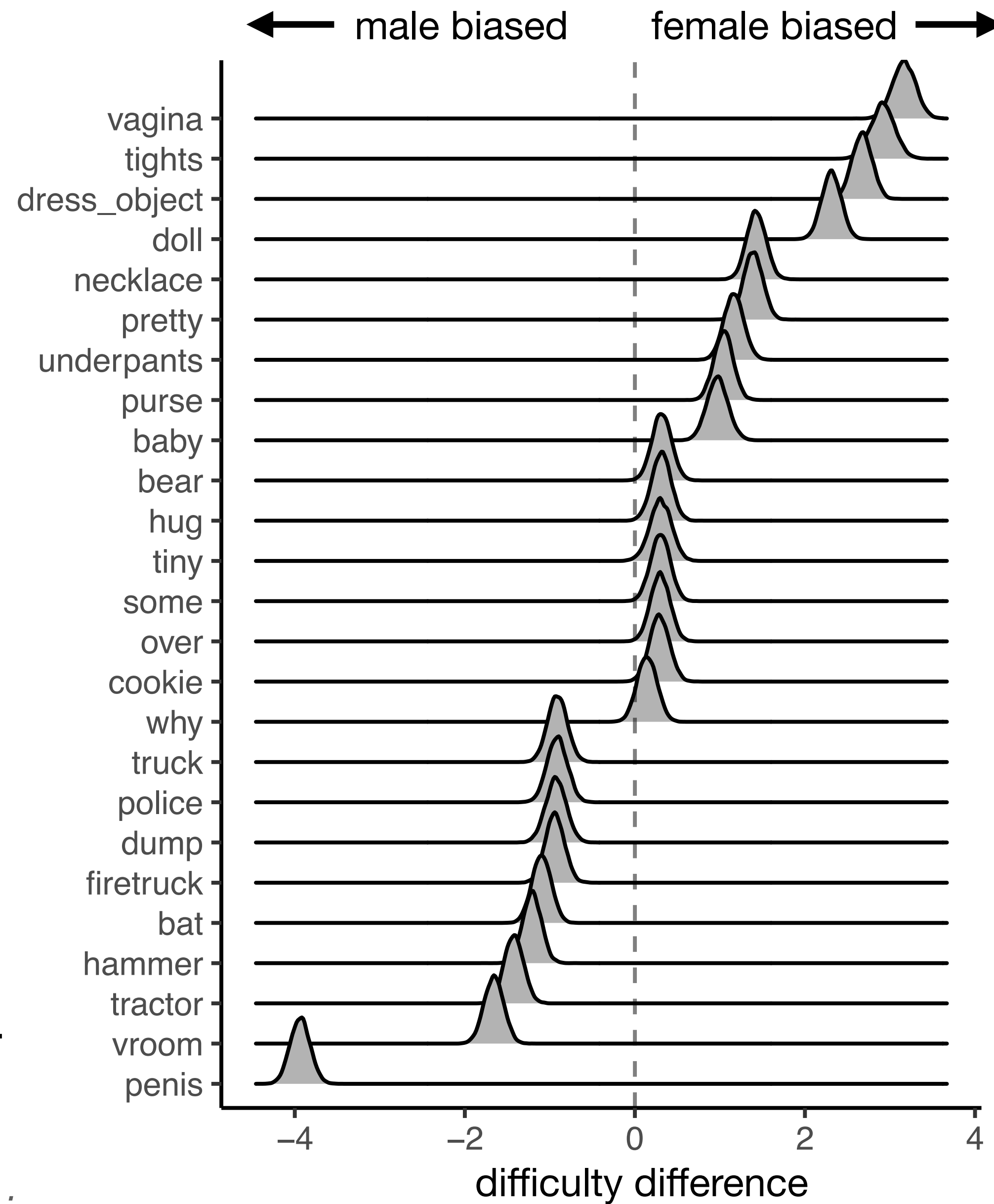
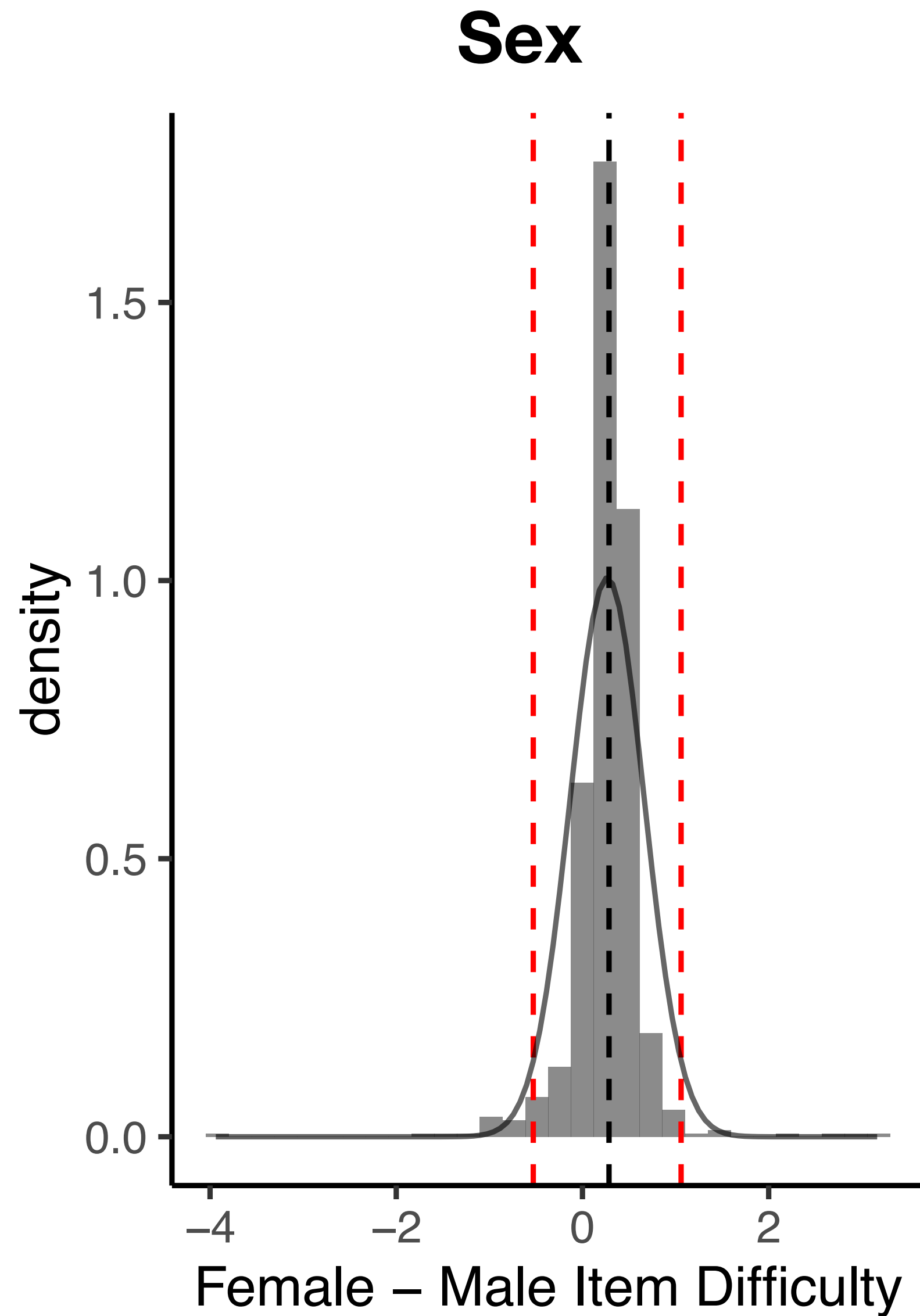
A distribution of item difficulty differences for low vs. high maternal education groups:



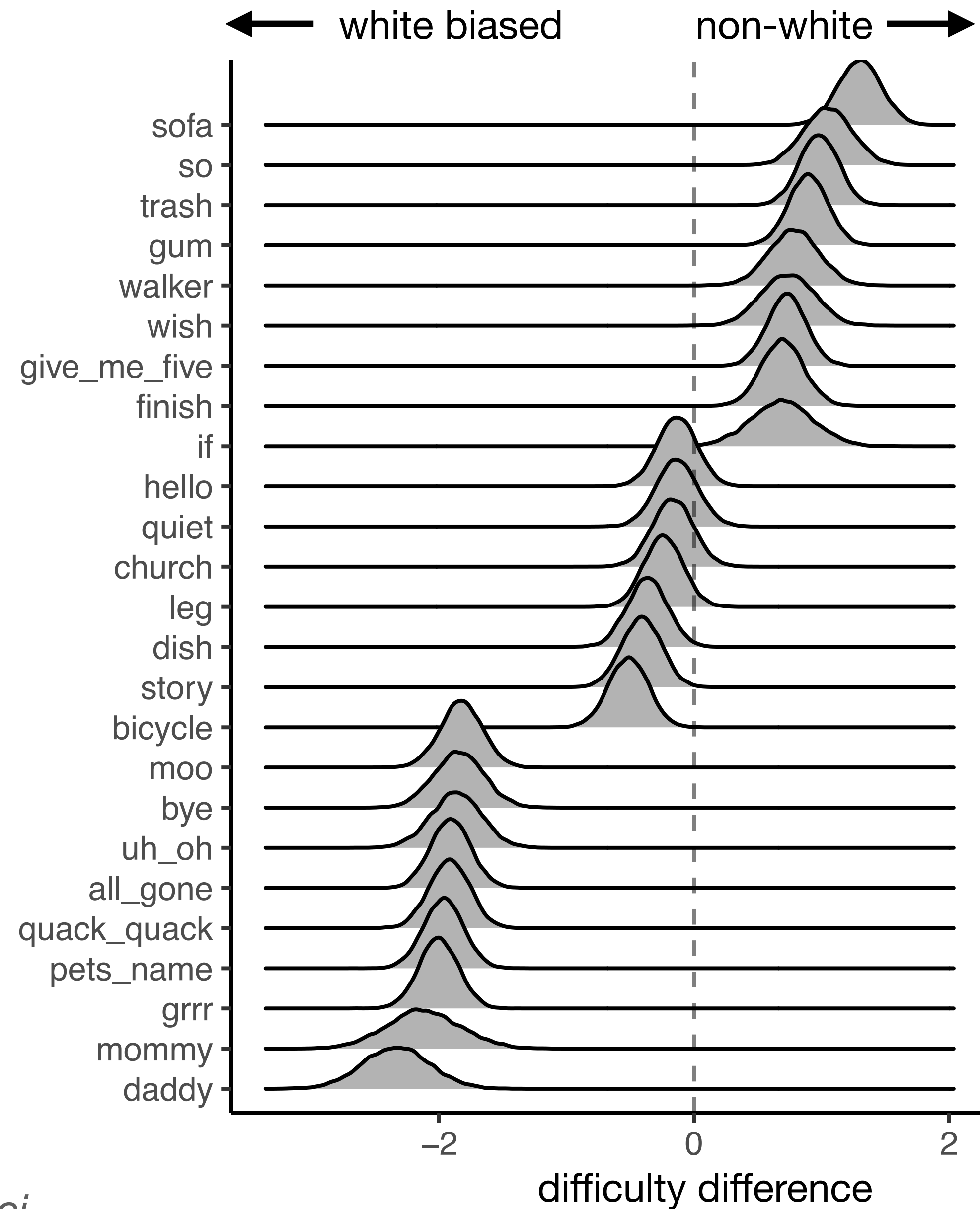
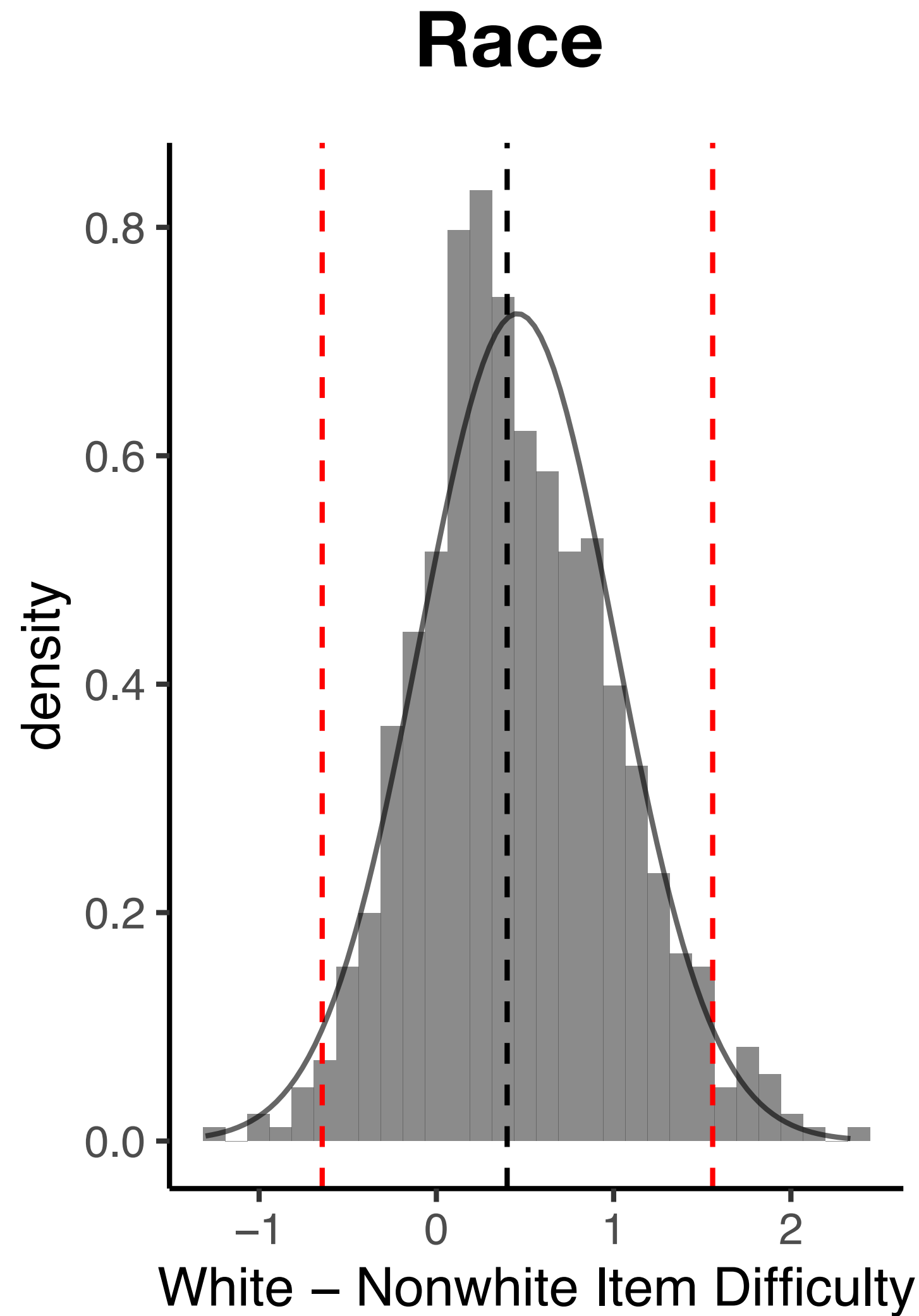
To understand between-item variation (i.e., which items are especially biased), we used GLIMMER (Graphs of Logits Imputed Multiply with Means Equal; Stenhaug, Frank, and Domingue, 2021), which draw from a fitted multigroup Rasch model that assumes mean language ability in each group is the same (pushing all variance into item difficulty).



Identifying Measurement Bias



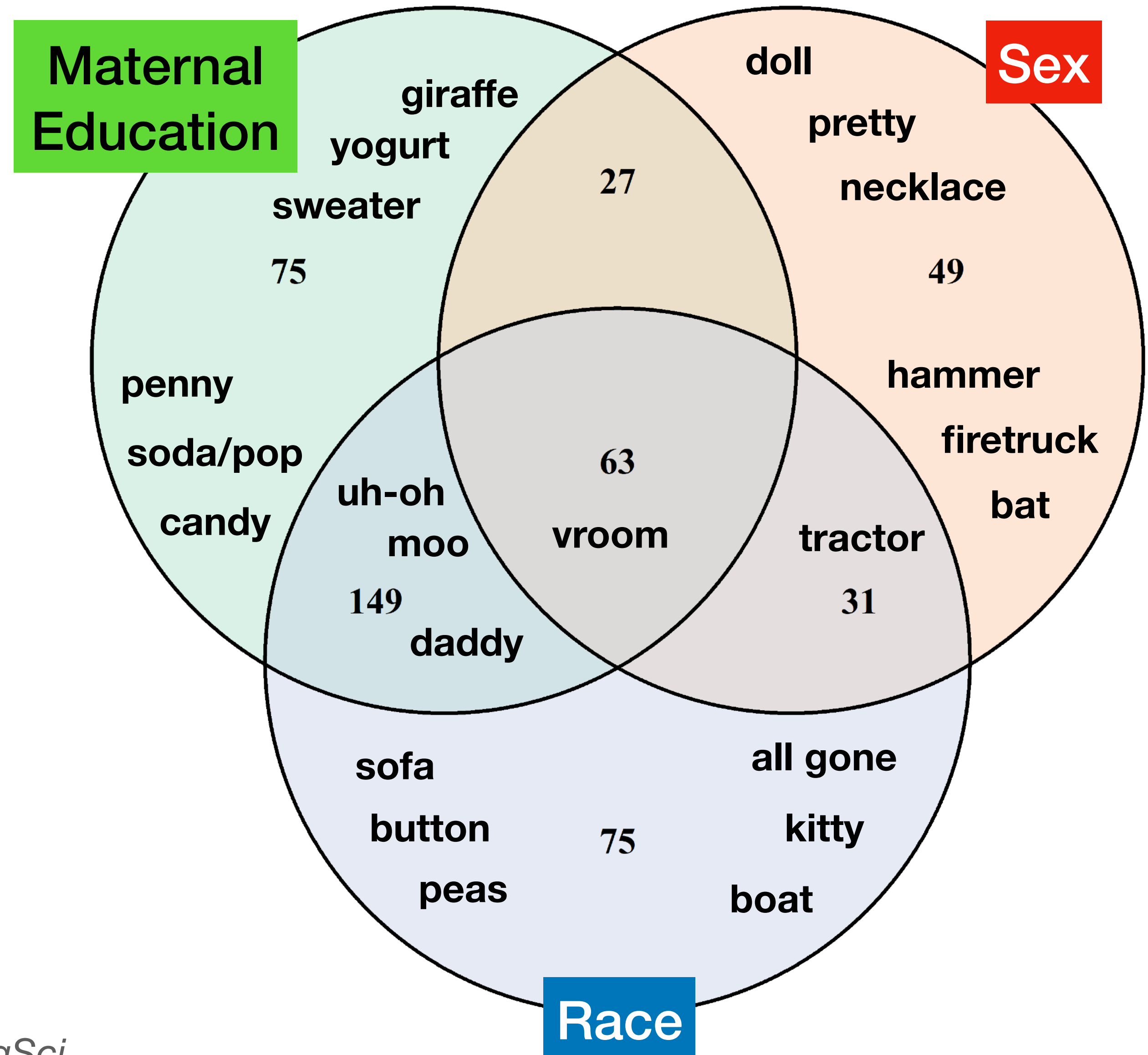
Identifying Measurement Bias



Reducing Measurement Bias

Many items showed significant bias favoring one or more demographic groups.

But we don't want to get rid of all/most of these items: How does eliminating the extrema change demographic-based differences?

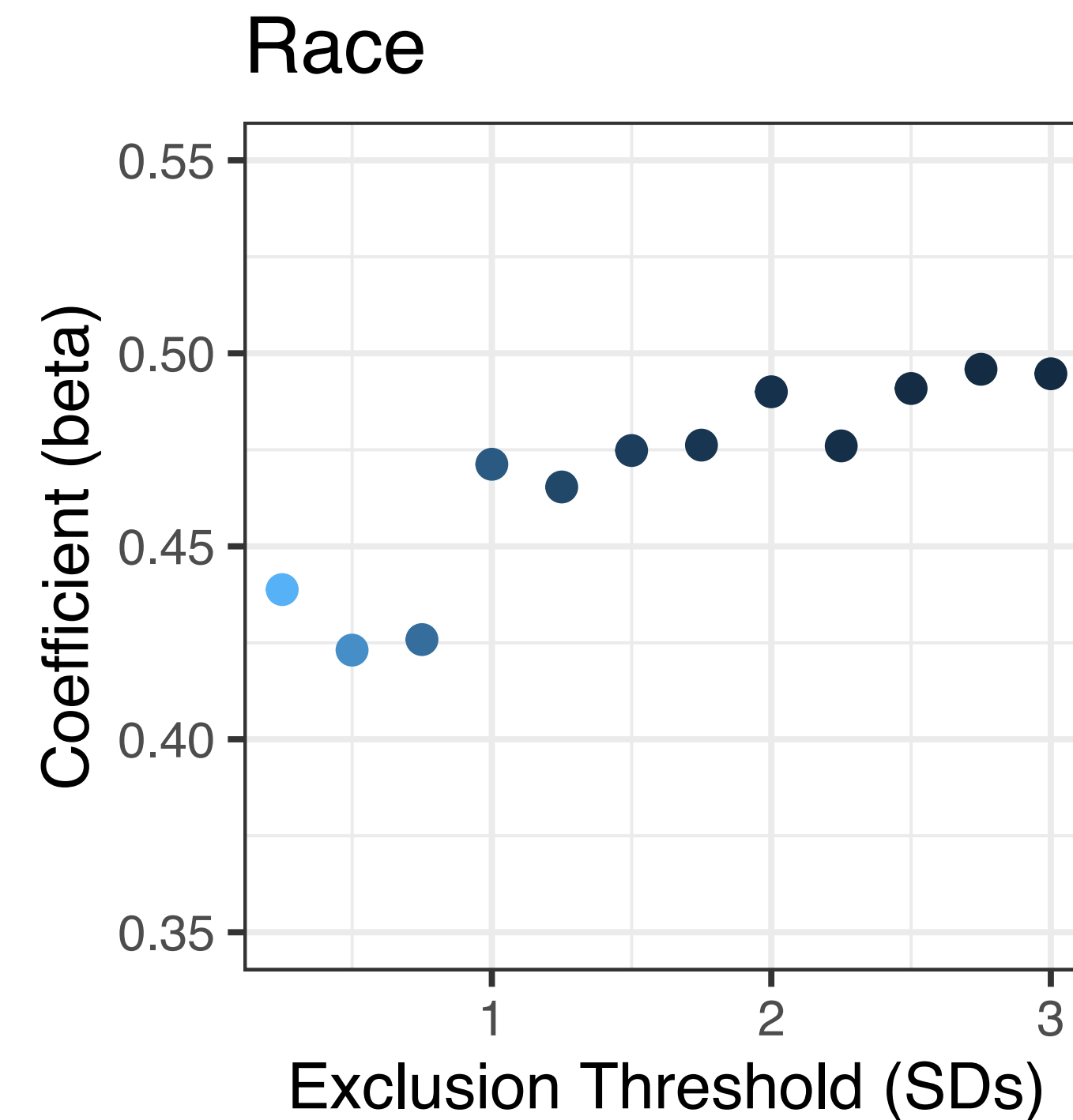
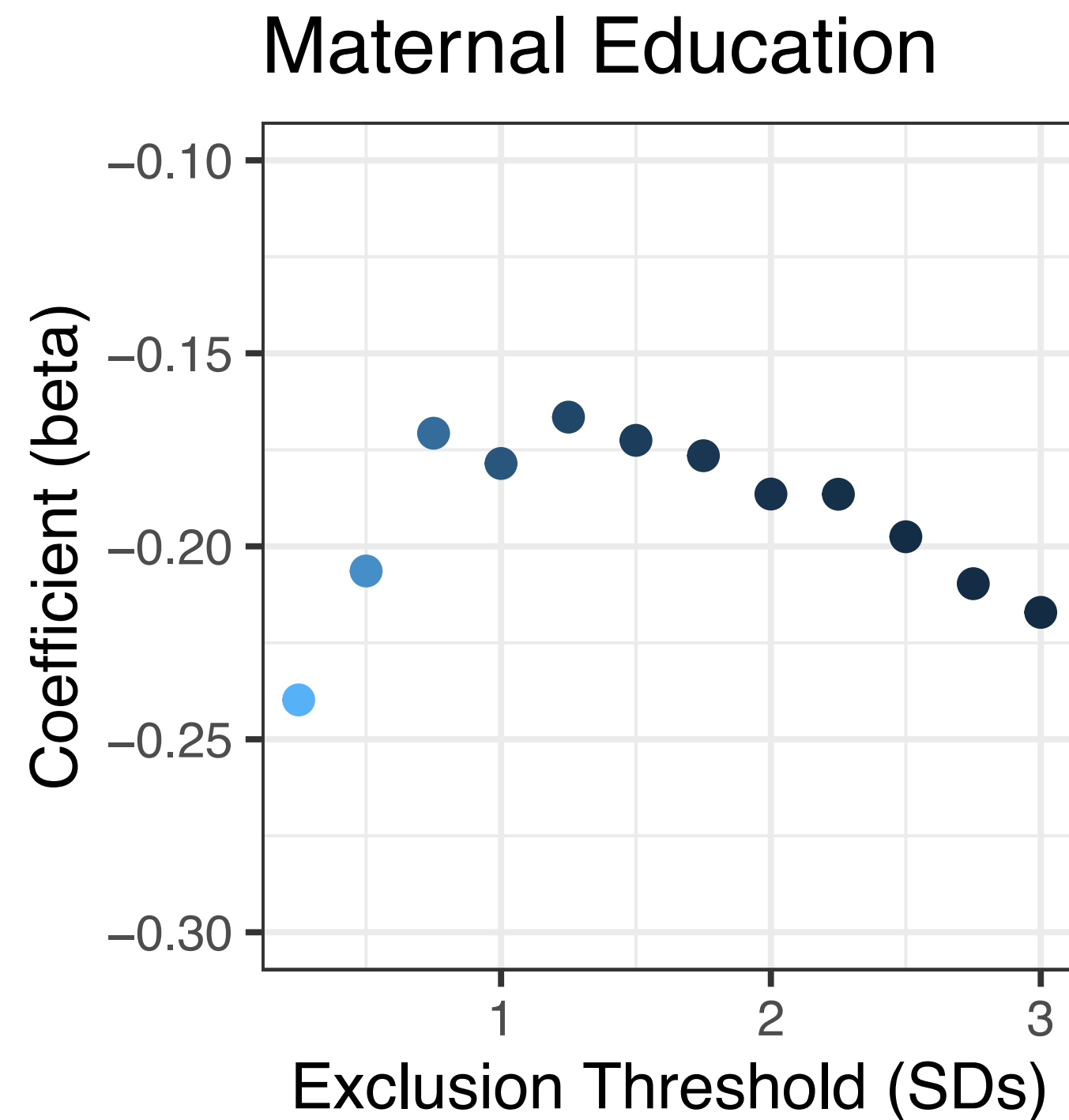
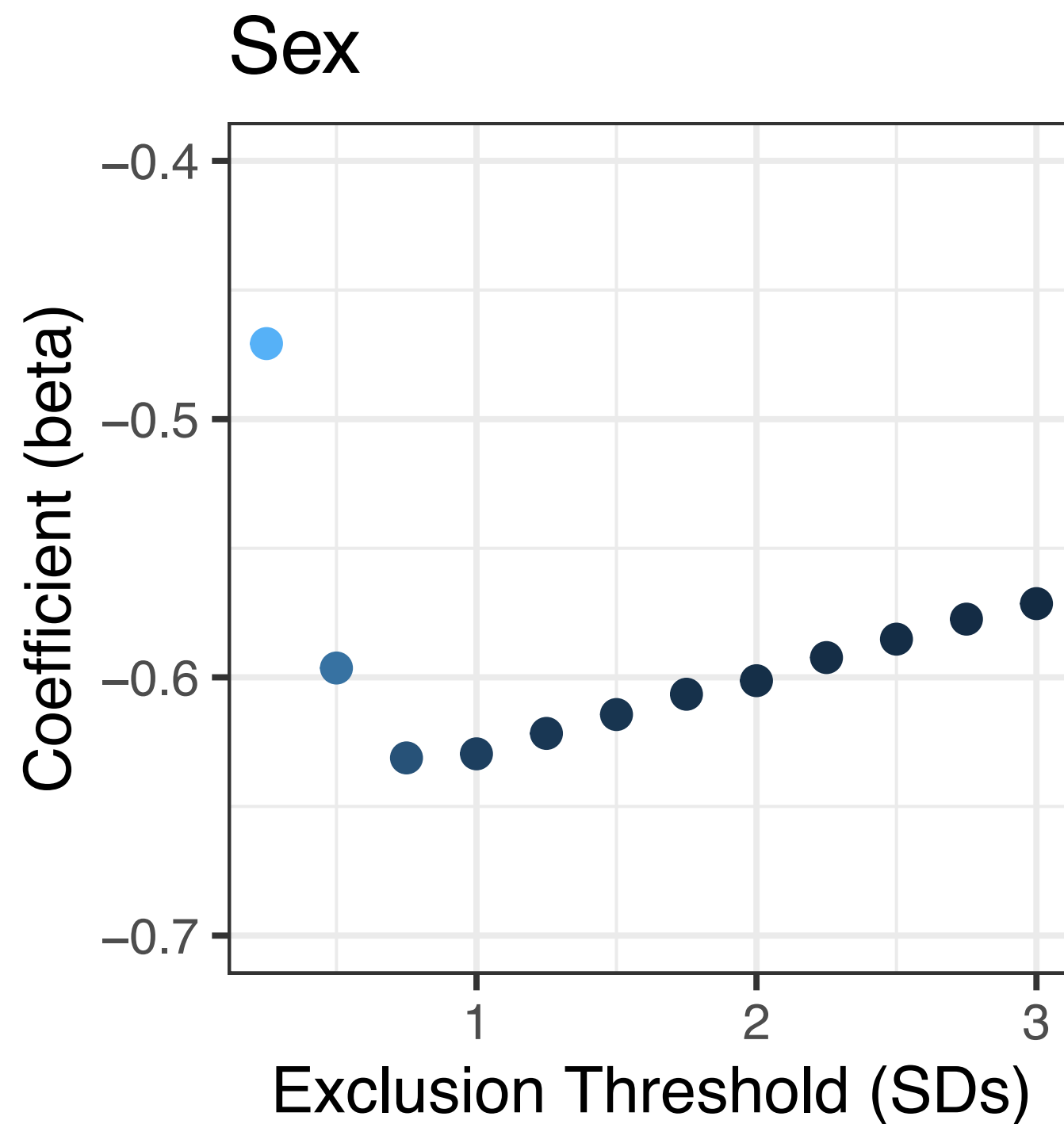


Reducing Measurement Bias

Suggestion: prune items showing extreme disadvantage for any demographic group



Pruning 59 extrema reduces the size of SES- and race-based demographic effects. Sex differences persisted.



Conclusion

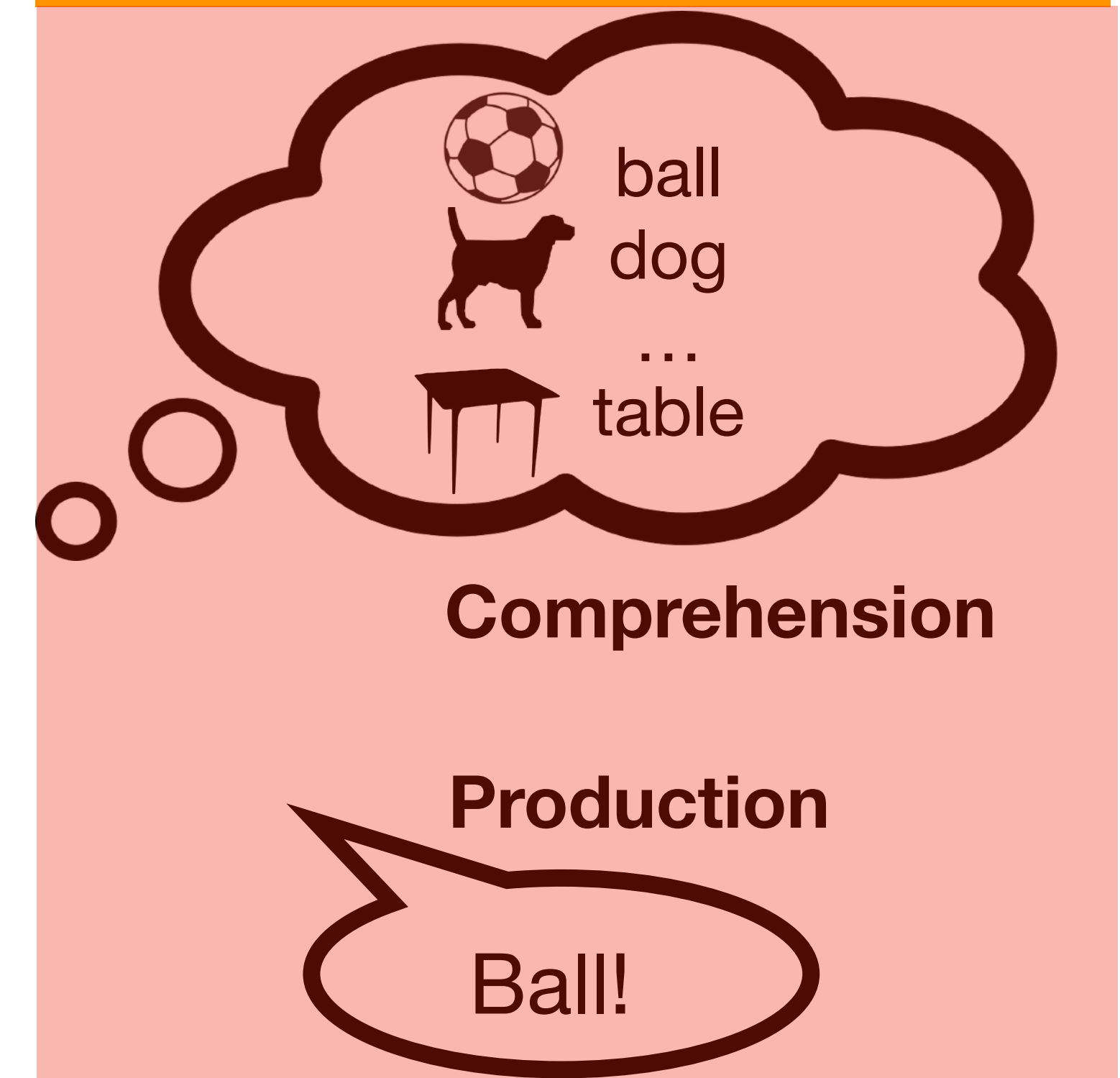
- Computerized Adaptive Tests (CATs) can evaluate children's language ability in a few minutes, with high validity
- How short can tests of early language be? ~25 (and up to 50) words
- Identified demographically-biased items on the English CDI, suggested removing a small number of extrema to mitigate SES and racial bias
- How can we assess & improve test fairness? *Evaluate Differential Item Functioning, eliminate outliers, and (ToDo) consider replacements*

“The history of science is the history of measurement.”

—James M. Cattell (1893), founder of *Psychological Review*

Uptake: learning outcomes

Create short, valid, and fair tests of early language



Thank you! — Questions?

kachergis@stanford.edu

<https://kachergis.com>

Thanks to my collaborators on these projects:



Michael Frank



Virginia Marchman



Bria Long



Alvin Wei Ming Tan



Nathan Francis

...and thanks to the **Language & Cognition** lab.

