

# A gentle intro to item response theory

Alvin Tan & George Kachergis  
CogSci pre-workshop tutorial, 2024-07-01



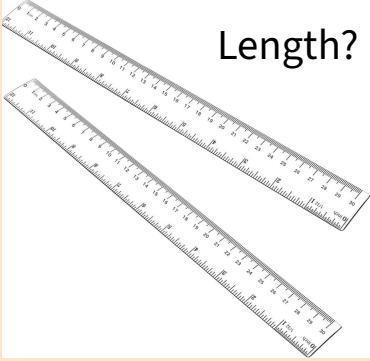
# Outline

- Why IRT?
- What is IRT?
- How to use IRT?
- Other IRT variants
- Cool stuff you can do with IRT
- Practical

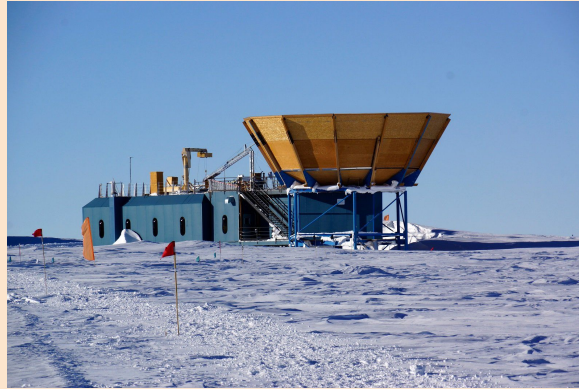
# **IRT motivation**

# How do we measure stuff?

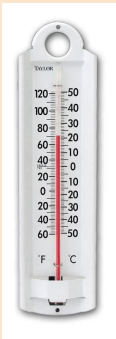
Length?



Cosmic microwave background radiation?



Temperature?



Personality?

Affect?

***Lack of direct  
observational  
access***

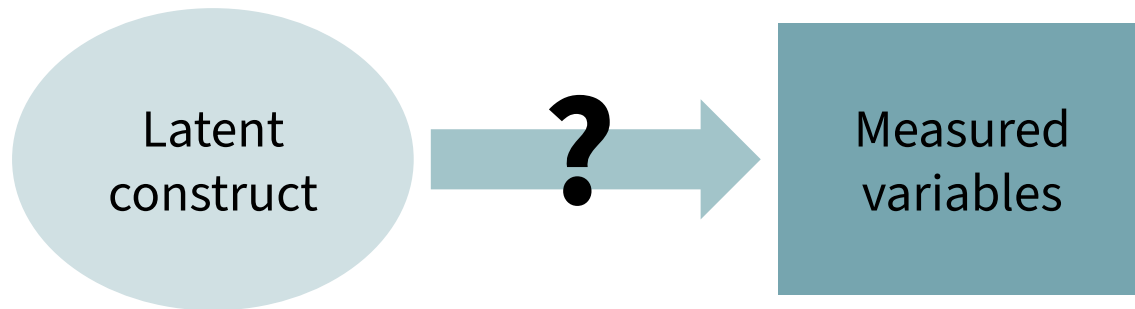
Memory?

Cognitive  
ability?

# The place for psychometrics

*What are the sources of variance?*

*Is my measurement consistent?*



*How are latent constructs organised?*

*How much error is there?*

# What is psychometrics?

*Validity*

***How do measurements relate to the underlying latent constructs?***

*Instrumentation*

- Scaling
- Reliability
- Bias
- Measurement invariance
- Differential item functioning

*Measurement model*

- **Classical test theory**
- **Item response theory**
- Generalisability theory

*Latent structure*

- Factor analysis
- Latent variable models
- Network theory
- Structural equation modelling

# Classical test theory

Suppose you had a test/questionnaire...

## Beck Depression Inventory

Name Justine Locke

Date 8 October 2021

Score 26

## Big Five Personality Test

### Instructions:

- For each statement, select the response that best reflects how you typically feel or behave.
- There are no right or wrong answers. Be honest and choose the response that feels most accurate for you. Please answer all questions.
- 1 = Strongly Agree, 2 = Agree, 3 = Neutral, 4 = Disagree, 5 = Strongly Disagree

### Openness to Experience

1. I enjoy trying new things and experiences.

1 2 3 4

☒ ☐ ☐ ☐

2. I am curious about the world and different cultures.

☐ ☒ ☐ ☐



## MacArthur-Bates CDI Words and Sentences

Copyright © 2007 The CDI Advisory Board.  
All rights reserved.  
Distributed by Paul W. Brookes Publishing Co.  
1-800-638-3775; 410-337-9580  
www.brookespublishing.com

Proper Mark

USE NO. 2 PENCIL ONLY

Improper Marks

### PART I WORDS CHILDREN USE

#### ADDITIONAL CHECKLIST

Understand many more words than they say. We are particularly interested in the words your child SAYS. Please go to the list and mark the words you have heard your child use. If your child uses a different pronunciation of a word (for "ruff" instead of "giraffe" or "sketti" for "spaghetti"), mark the word anyway. Remember that this is a "catalogue" words that are used by many different children. Don't worry if your child knows only a few of these right now.

#### EFFECTS AND ANIMAL SOUNDS (12)

<input type="checkbox"/> meow	<input type="checkbox"/> uh oh
<input type="checkbox"/> moo	<input type="checkbox"/> vroom
<input type="checkbox"/> ooh	<input type="checkbox"/> woof woof
<input type="checkbox"/> quack quack	<input type="checkbox"/> yum yum

...how would you determine a test-taker's level?

# Classical test theory

One obvious answer: just add it up (= sum scores/true scores)!

$$Y_i = T_i + e_i$$

$Y_i$ : observed score

$T_i$ : true score

$e_i$ : error

With  $k$  items,

$$(\sum Y_i)/k = (\sum T_i)/k + \underbrace{(\sum e_i)/k}$$

Expected value = 0

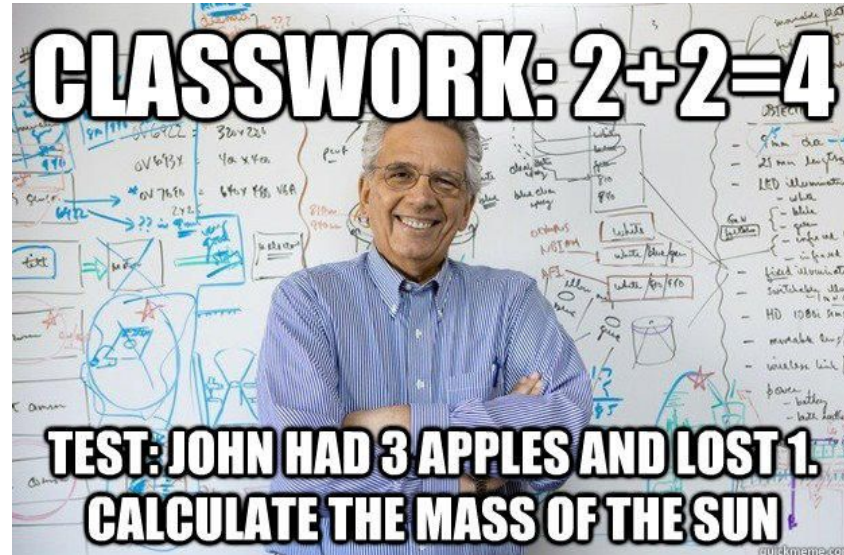
Variance decreases as  $k$  increases



# Classical test theory

One obvious answer: just add it up (= sum scores/true scores)!

BUT this assumes that **all items behave exactly the same**

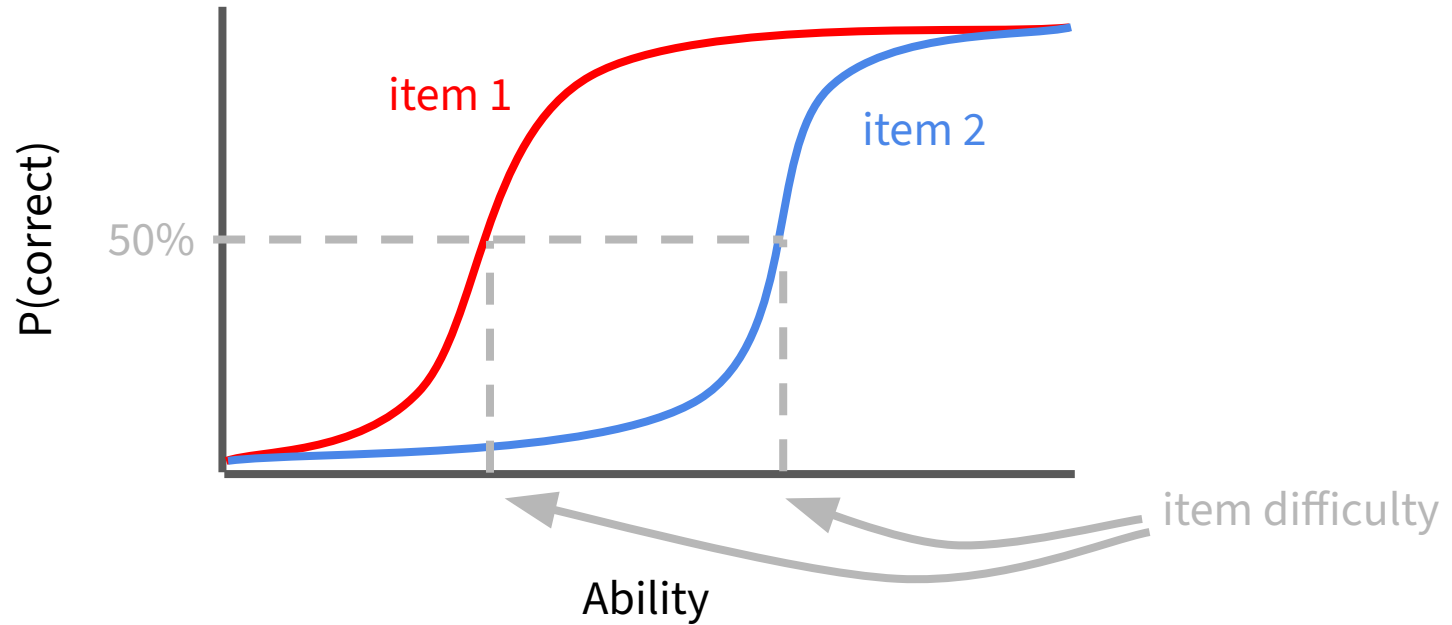


# Classical test theory

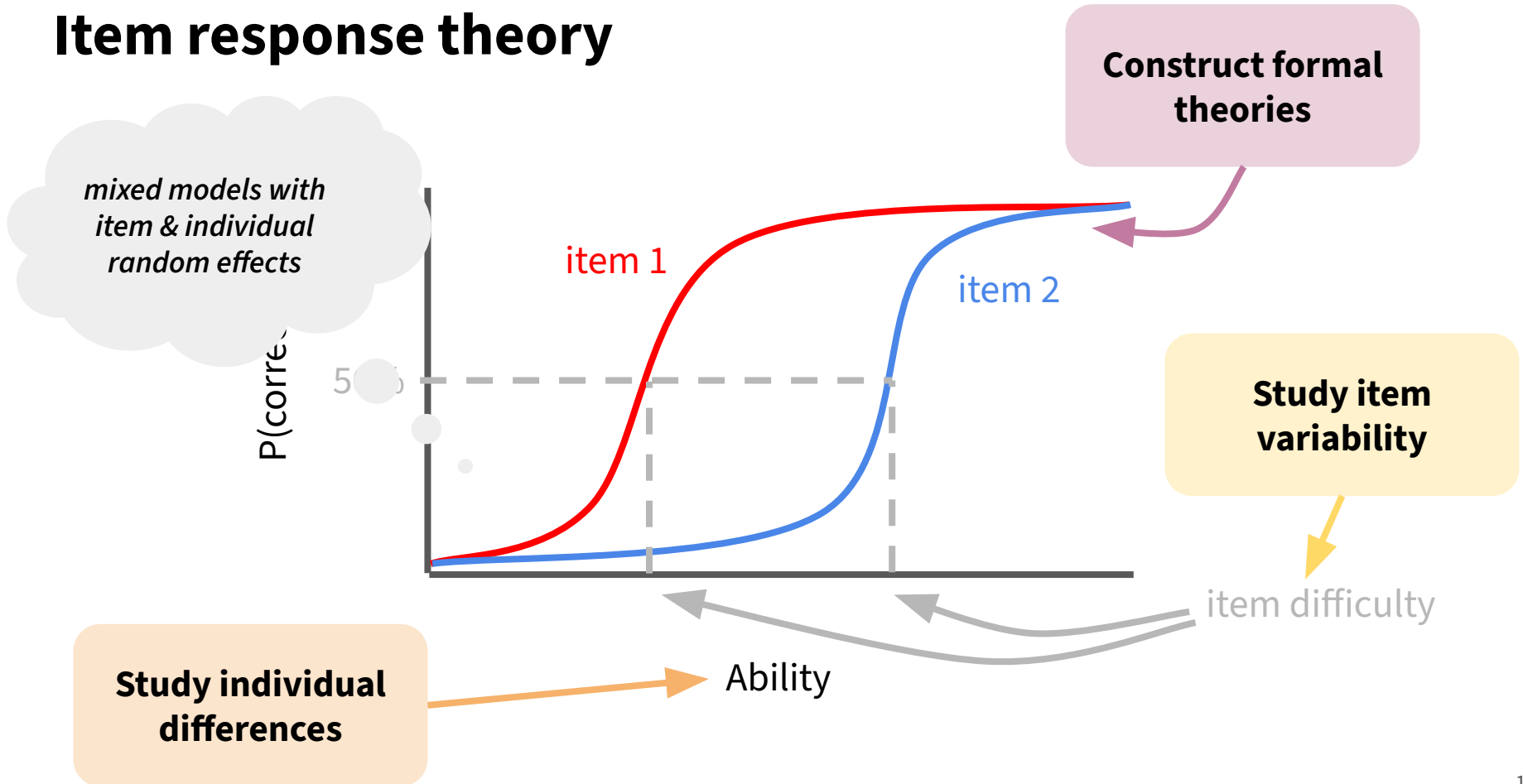
Consequences of this assumption:

- Only test-level information; no item-level information
- Reliability assumes parallel forms, but no way of testing
- Assumes measure is equally good over all levels of ability
- No generalisability to new tests → incentive to use the same ones

# Item response theory

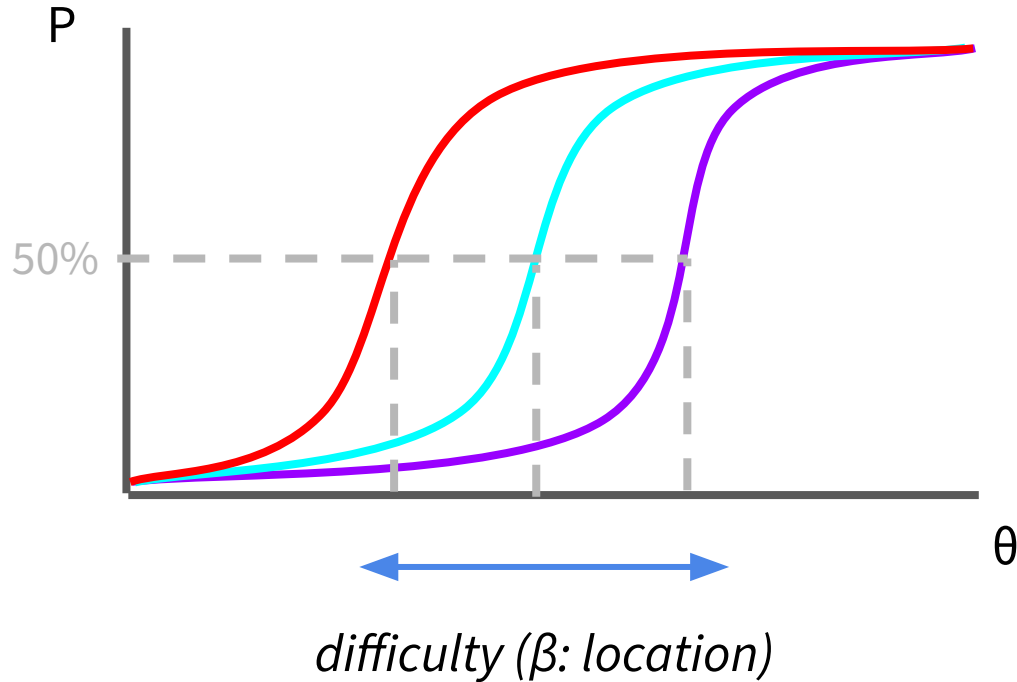


# Item response theory

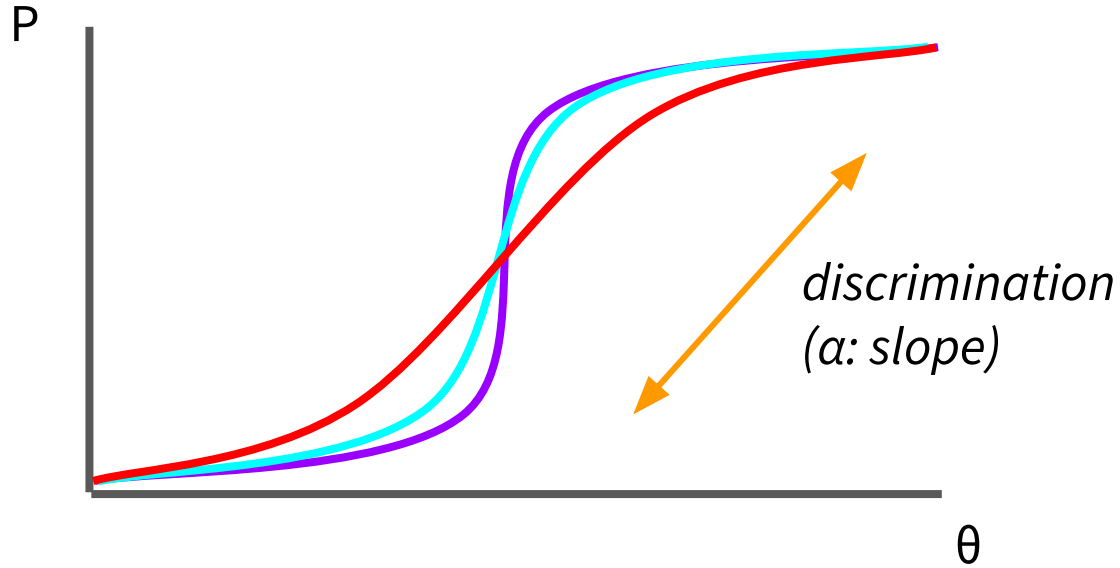


# **IRT basics**

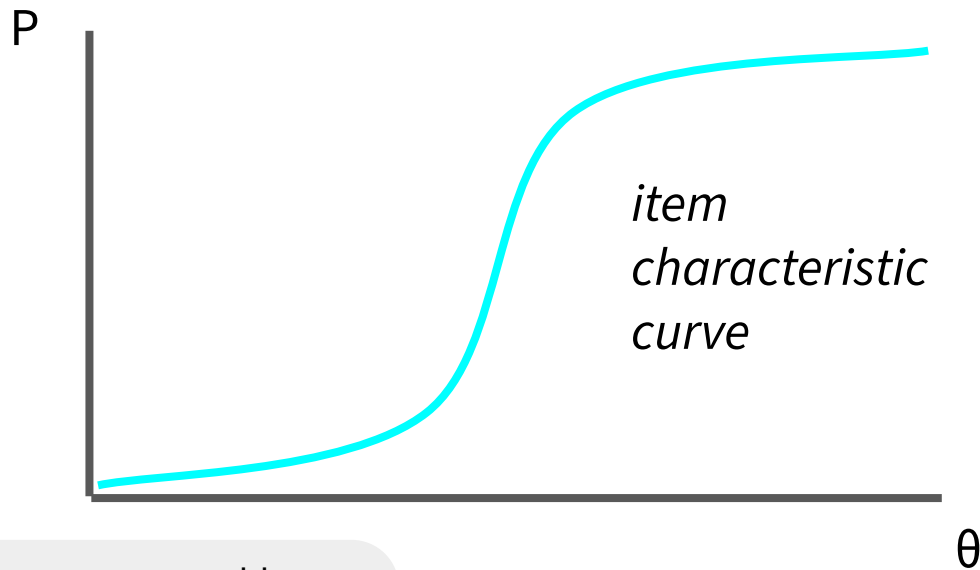
# Item response theoretic model (2PL)



# Item response theoretic model (2PL)



# Putting the L in 2PL



$$\begin{aligned} \text{logit}(P) &= \ln\left(\overbrace{P / (1 - P)}^{\text{odds}}\right) \\ &= \alpha(\theta - \beta) \end{aligned}$$

$$P(X_{ip} = 1 | \theta_p, \beta_i) = \frac{\exp(\alpha_i(\theta_p - \beta_i))}{1 + \exp(\alpha_i(\theta_p - \beta_i))}$$



# Item response theory

Consequences of this model:

- Item are on a meaningful scale (theta measures difficulty)
- Individuals are on a meaningful scale (latent trait scores)
- These two dimensions lie on the same scale
- Can calculate probability of an individual getting an item correct directly

$$P(X_{ip} = 1 | \theta_p, \beta_i) = \frac{\exp(\alpha_i(\theta_p - \beta_i))}{1 + \exp(\alpha_i(\theta_p - \beta_i))}$$

# Ex: Communicative Development Inventories (CDIs)

1. SOUND EFFECTS AND ANIMAL SOUNDS (12)					
baa baa	<input type="radio"/>	meow	<input type="radio"/>	uh oh	<input type="radio"/>
choo choo	<input type="radio"/>	moo	<input type="radio"/>	vroom	<input type="radio"/>
cockadoodledoo	<input type="radio"/>	ouch	<input type="radio"/>	woof woof	<input type="radio"/>
grrr	<input type="radio"/>	quack quack	<input type="radio"/>	yum yum	<input type="radio"/>

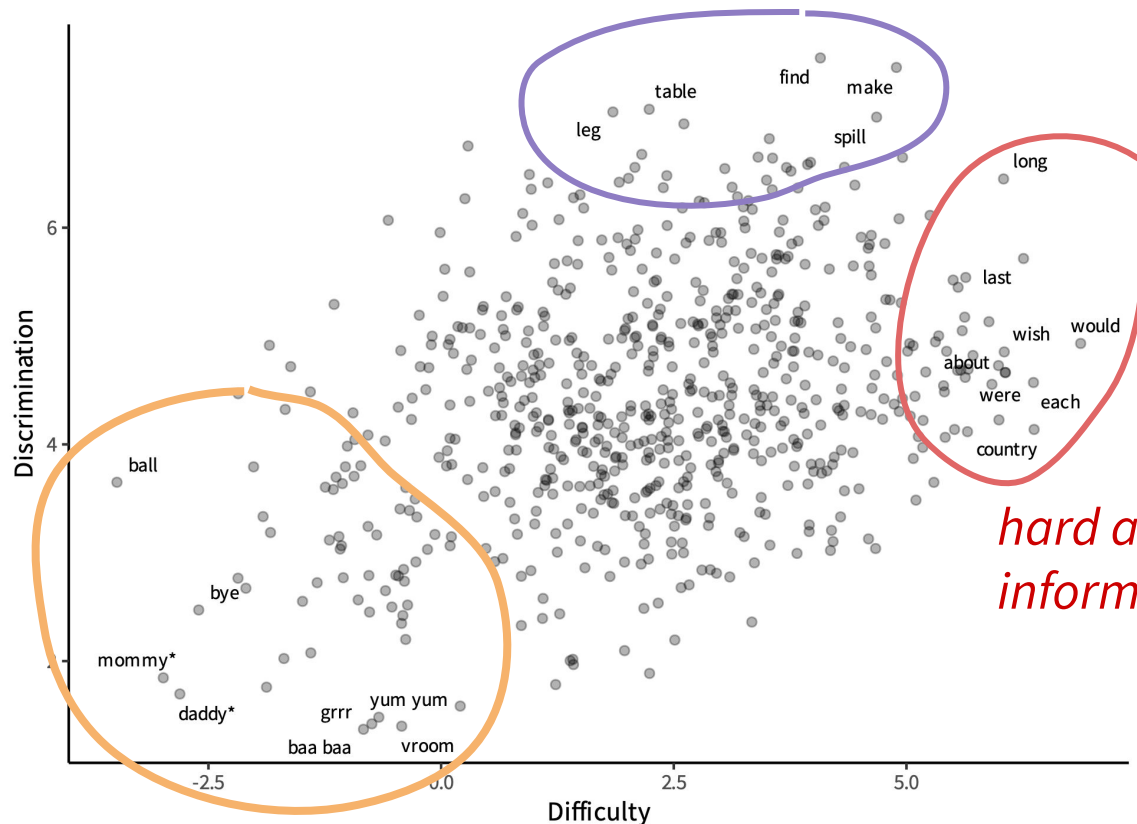
  

2. ANIMALS (Real or Toy) (43)					
alligator	<input type="radio"/>	duck	<input type="radio"/>	penguin	<input type="radio"/>
animal	<input type="radio"/>	elephant	<input type="radio"/>	pig	<input type="radio"/>
ant	<input type="radio"/>	fish	<input type="radio"/>	pony	<input type="radio"/>
bear	<input type="radio"/>	frog	<input type="radio"/>	puppy	<input type="radio"/>
bee	<input type="radio"/>	giraffe	<input type="radio"/>	rooster	<input type="radio"/>
bird	<input type="radio"/>	goose	<input type="radio"/>	sheep	<input type="radio"/>
bug	<input type="radio"/>	hen	<input type="radio"/>	squirrel	<input type="radio"/>
bunny	<input type="radio"/>	horse	<input type="radio"/>	teddybear	<input type="radio"/>
butterfly	<input type="radio"/>	kitty	<input type="radio"/>	tiger	<input type="radio"/>
cat	<input type="radio"/>	lamb	<input type="radio"/>	turkey	<input type="radio"/>
chicken	<input type="radio"/>	lion	<input type="radio"/>	turtle	<input type="radio"/>
cow	<input type="radio"/>	monkey	<input type="radio"/>	wolf	<input type="radio"/>
deer	<input type="radio"/>	moose	<input type="radio"/>	zebra	<input type="radio"/>
dog	<input type="radio"/>	mouse	<input type="radio"/>		
donkey	<input type="radio"/>	owl	<input type="radio"/>		



# Ex: CDI parameter space (2PL)

*fairly hard and  
informative*



*easy and not  
informative*

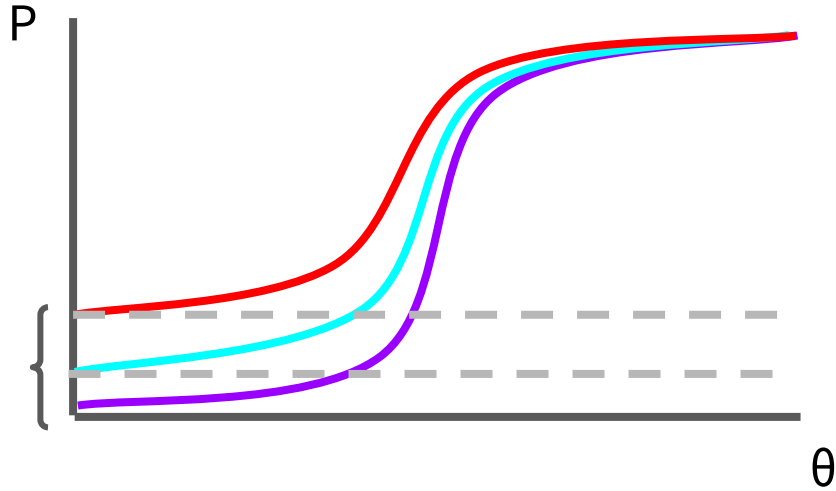
*hard and somewhat  
informative*

# **IRT variants**

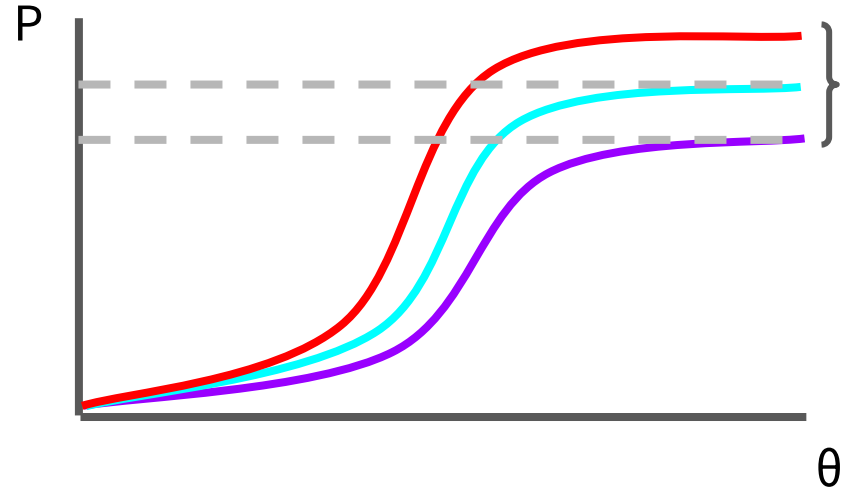
# IRT variants: Number of parameters

- 1PL: location (difficulty)
- 2PL: slope (discrimination)
- 3PL: lower bound (guessing)
- 4PL: upper bound (errors)

# Para para parameter



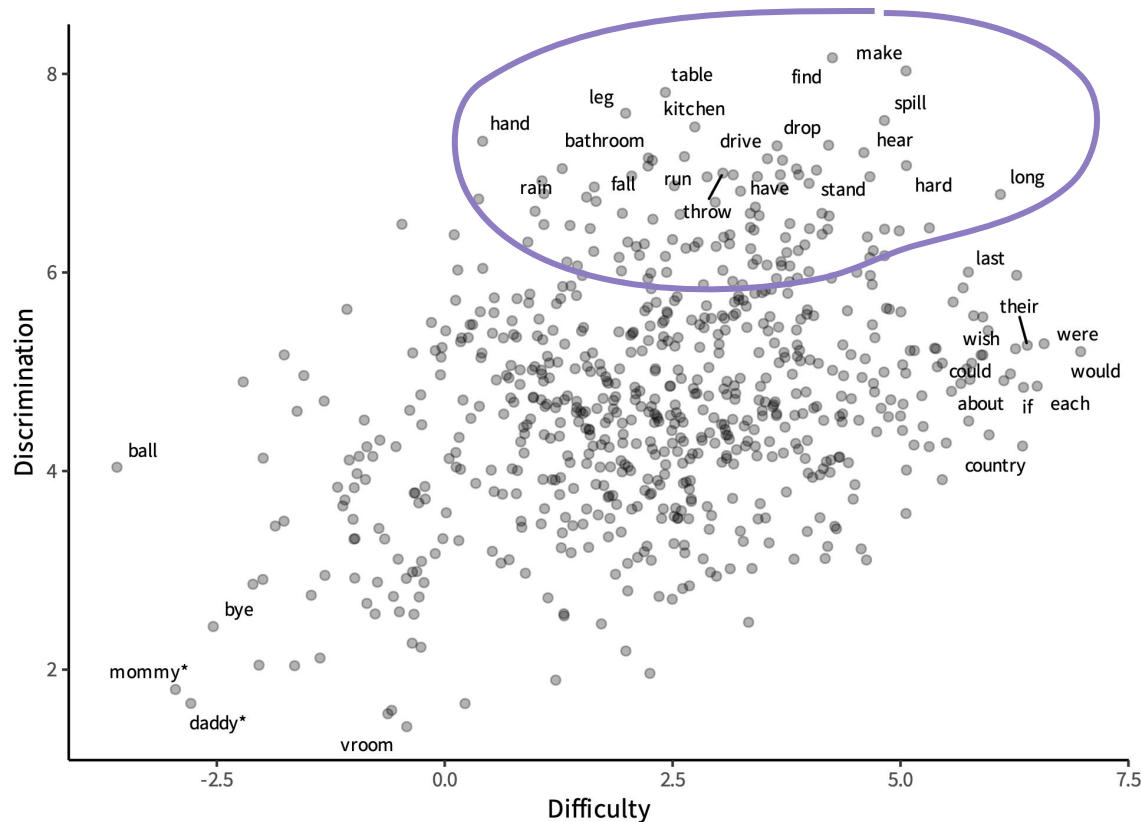
3PL: *lower asymptote*



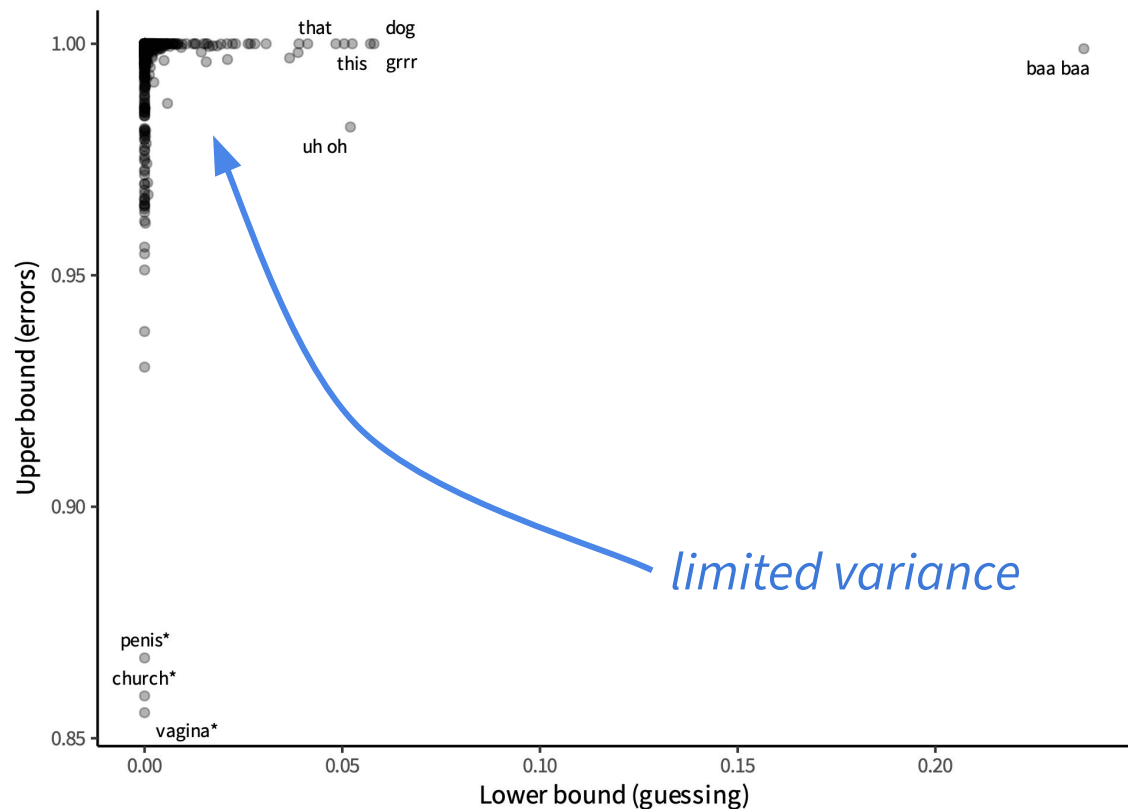
4PL: *upper asymptote*

# Ex: CDI parameter space (4PL)

*more informative items*



# Ex: CDI parameter space (4PL)





## Ex: CDI model fits

Model	BIC
1PL	4788209
<b>2PL</b>	<b>4668081</b>
4PL	4676954

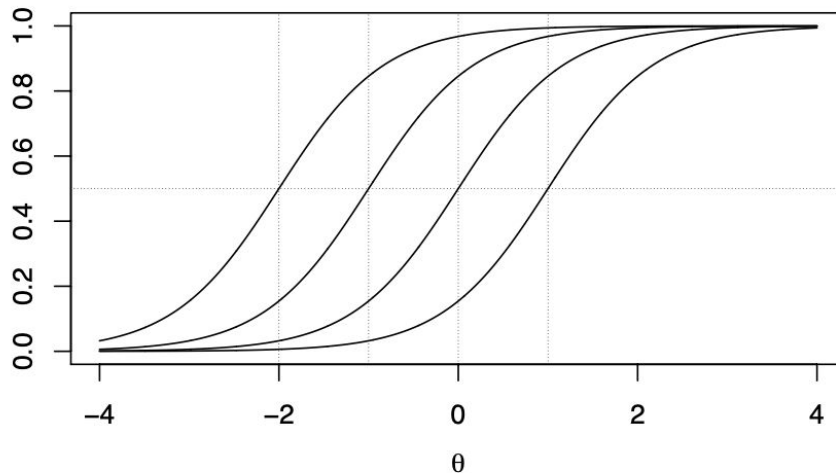
# IRT variants: Item type

- Vanilla IRT operates over binary observations
- How do we move from dichotomous to polytomous measures?
- Solution: decompose polytomous items into sets of dichotomous comparisons
- Two main classes: difference models (e.g., GRM) and divide-by-total models (e.g., PCM)

# Difference models: Graded response model

Dichotomise based on being above/below a threshold; e.g., with  $k = 5$

- $P_1 = \Pr(y_{ip} > 0)$
- $P_2 = \Pr(y_{ip} > 1)$
- $P_3 = \Pr(y_{ip} > 2)$
- $P_4 = \Pr(y_{ip} > 3)$



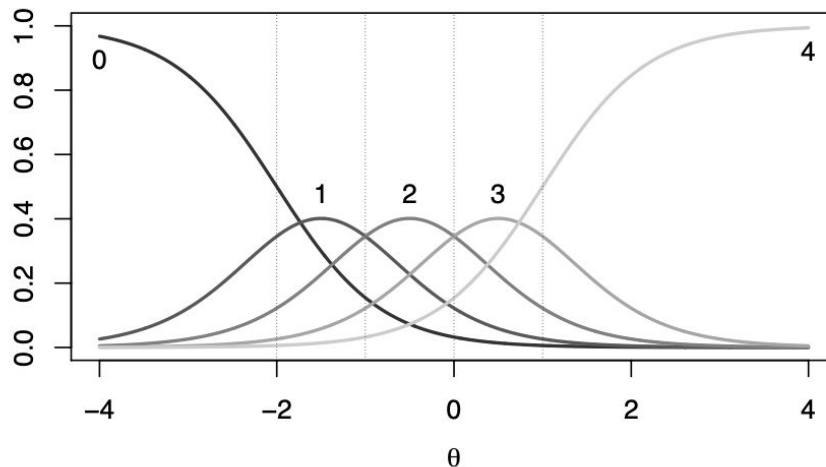
*cumulative response curves*

# Difference models: Graded response model

Dichotomise based on being above/below a threshold; e.g., with  $k = 5$

- $P_1 = \Pr(y_{ip} > 0)$
- $P_2 = \Pr(y_{ip} > 1)$
- $P_3 = \Pr(y_{ip} > 2)$
- $P_4 = \Pr(y_{ip} > 3)$

$$\Pr(y_{ip} = 2) = P_3 - P_2$$



*category response curves*

# **IRT techniques**

# IRT techniques: Equating

How do you handle multiple versions of an instrument?

- Avoid practice effects
- Prevent cheating
- Updating

Group 1

Form 1		
		Form 2

Group 2



*anchor items*

→ **concurrent calibration**

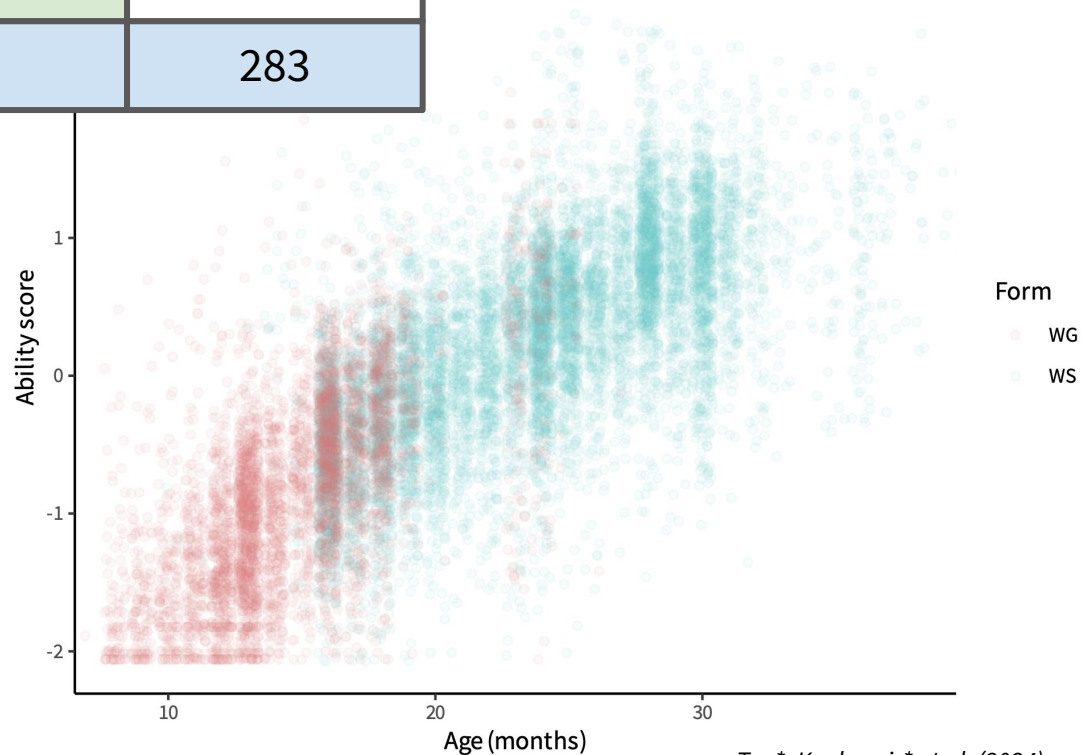


# Ex: Stitching different CDI forms

WG: 12–18mo

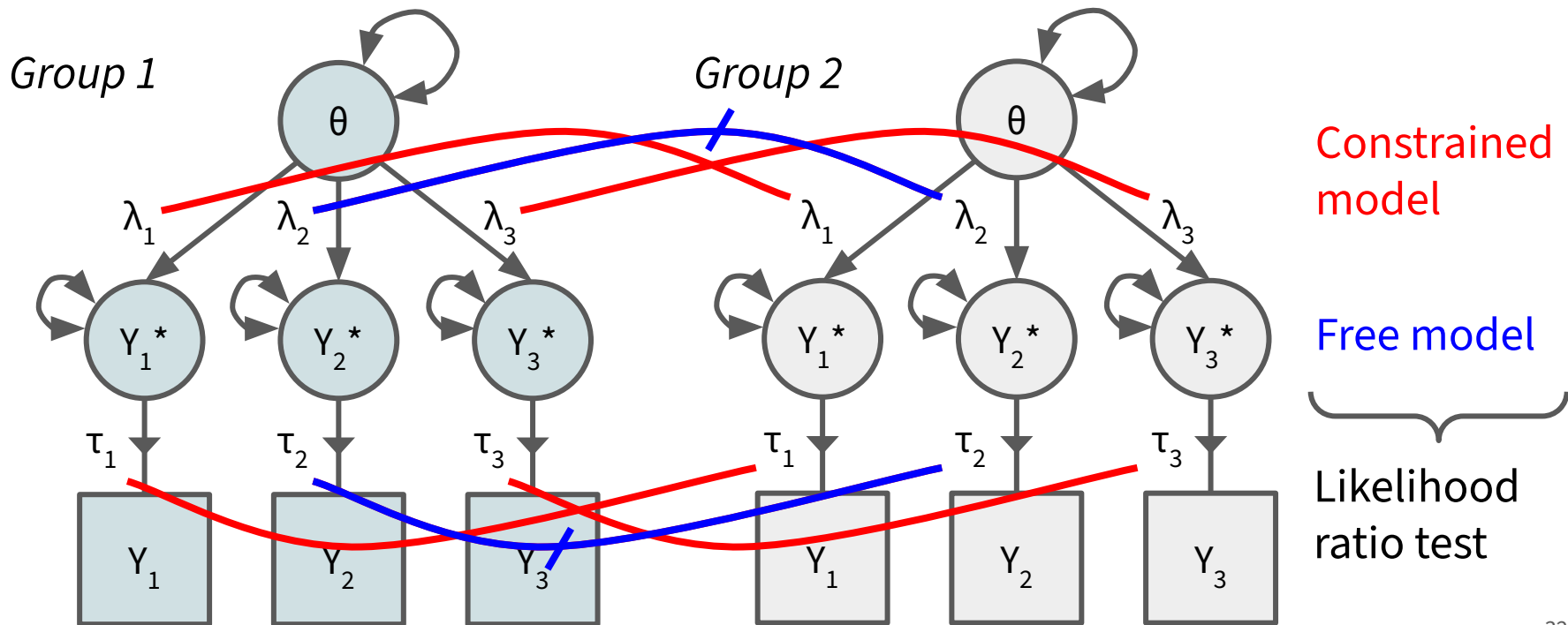
WS: 16–36mo

2	397	
	397	283



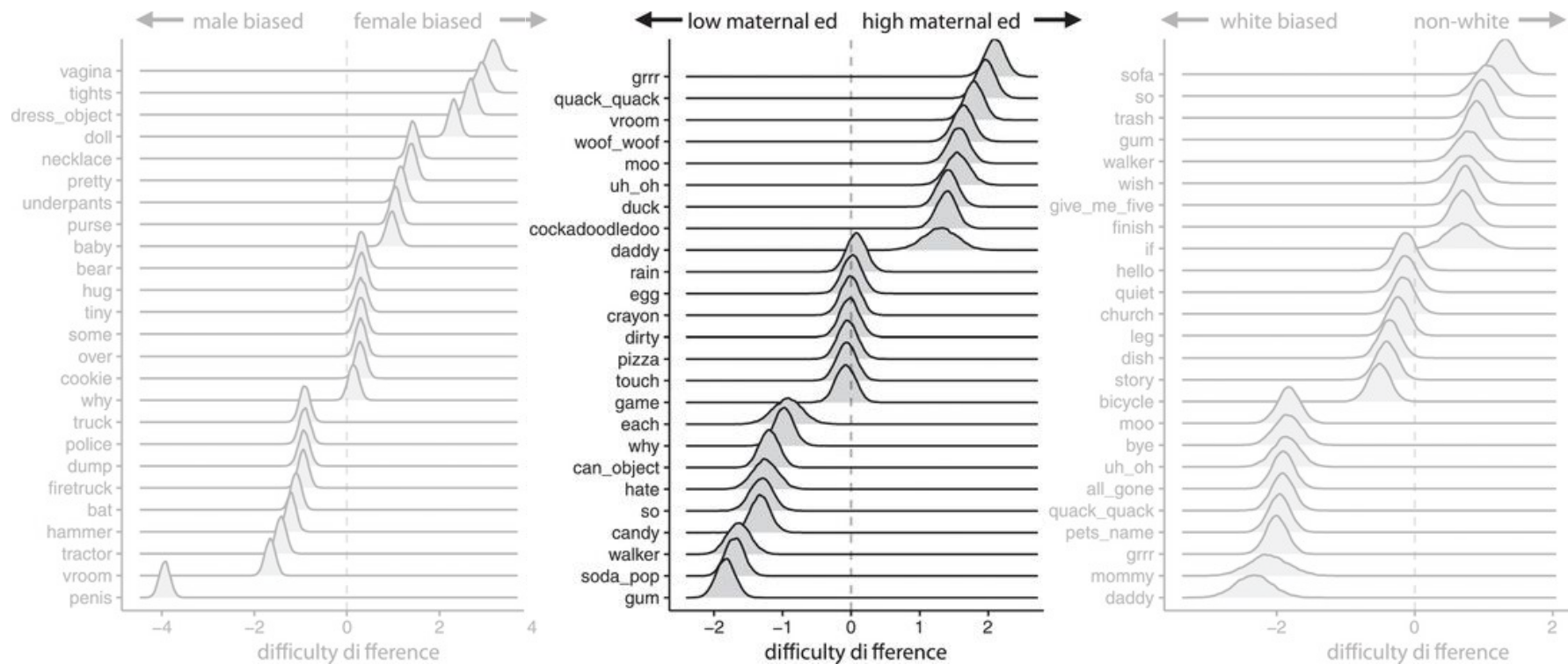
# IRT techniques: Differential item functioning

How do we know if items have measurement invariance across groups?



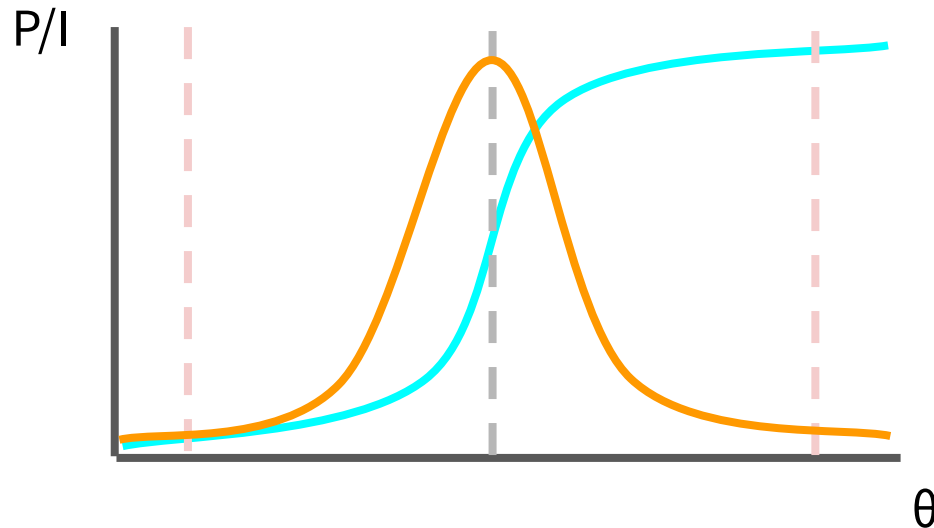


# Ex: DIF in CDIs



# IRT techniques: Item information

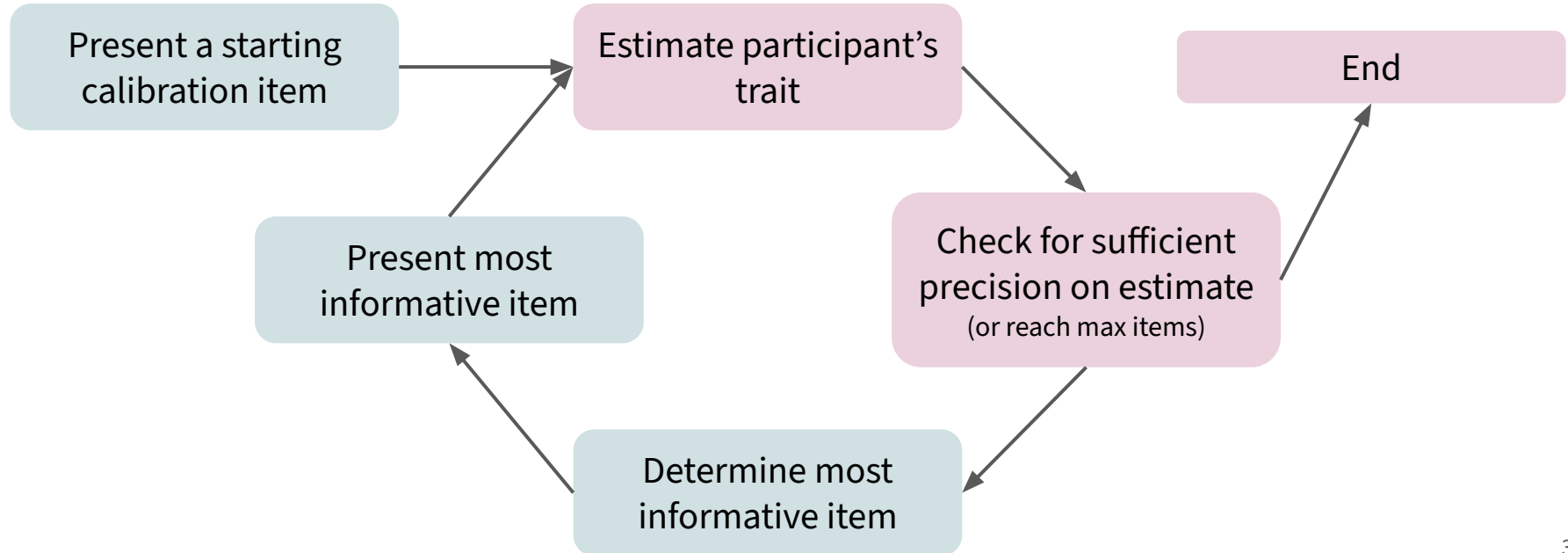
Because items are at different locations, they are differently informative about individuals of different trait levels



$$\begin{aligned} I(\theta) &= dP(\theta) / d\theta \\ &= P(\theta) (1 - P(\theta)) \end{aligned}$$

# IRT techniques: Computerised adaptive testing

Instead of giving a full instrument, choose items that are most informative given your current estimate of the participant's trait

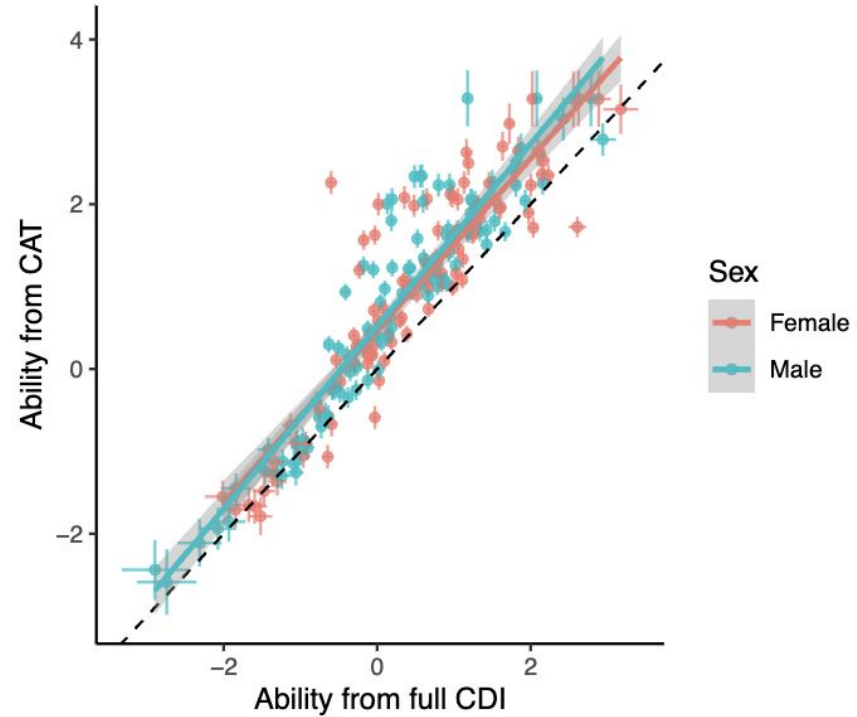


# Ex: CDI-CAT

CDI:  
~600–800 items



**CDI-CAT:**  
**~35 items**



# Summary

- IRT helps us better handle inter-item variability
- IRT is effectively logistic regression with item-level parameters
- IRT models can have various numbers of parameters, and can handle different types of categorical variables
- IRT enables us to do lots of statistical techniques that are not possible with CTT (e.g., equating, DIF, CAT)

# ty+q?

Slides, links, resources:

<https://psychometrics-workshop.github.io/>