

Domain Experts' Interpretations of Assessment Bias in a Scaled, Online Computer Science Curriculum

Benjamin Xie, Matt J. Davidson, Baker Franke, Emily McLeod,

Min Li Amy J. Ko

University of Washington

bxie@uw.edu



@benjixie
bxie@uw.edu

Hi

name + pronouns

I'm a PhD candidate at the University of Washington Seattle,

a university which acknowledges the Coast Salish peoples of this land, the land which touches the shared waters of all tribes and bands within the Duwamish (doo-amish), Puyallup (pee-all-up), Suquamish, Tulalip (too-lay-lip) and Muckleshoot nations.

And I'm EXCITED to share work that I've done with colleagues from the UW Information School, College of Education, DUB group, as well as with nonprofit [code.org](#)

Tests are not perfect measurement instruments

2

@benjixie
bxie@uw.edu

We use test scores for a lot of things.

University use scores for test such as the AP CS exams to determine if a student should be accepted into a university or major.

Teachers use tests for summative purposes such as grading.

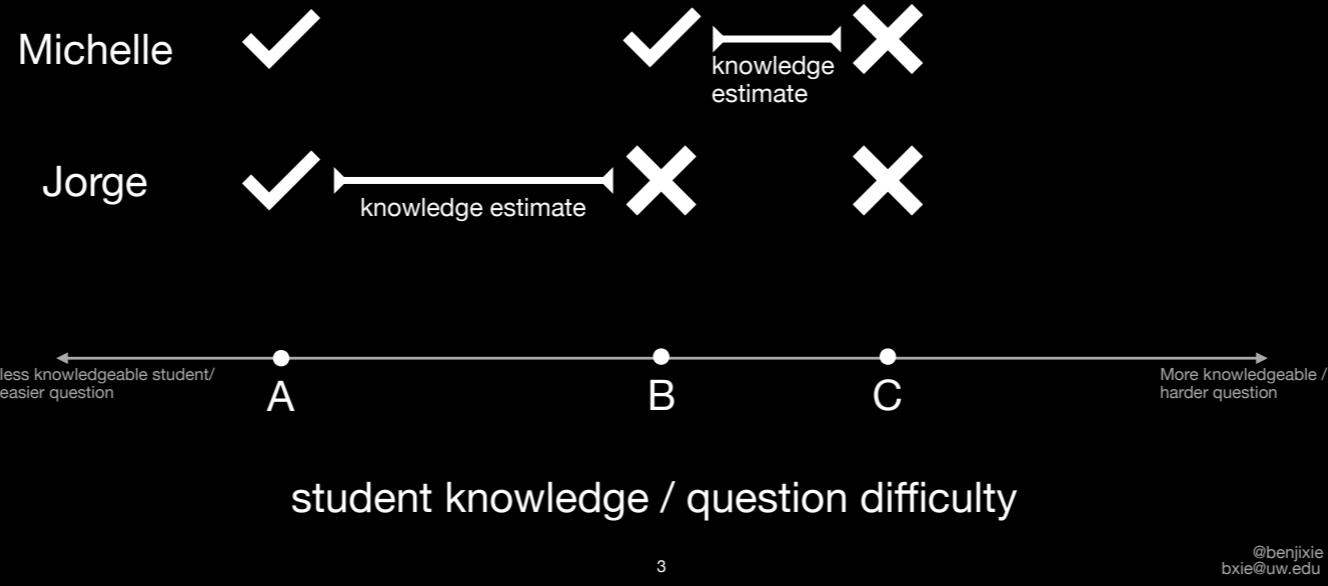
Students use tests to self-assess what they know, and the results can affect their self-efficacy and sense-of belongingness.

How do we know how good our tests are?

How different people interpret and use test scores is important, but tests are imperfect measures of knowledge.

Modeling students & questions (w/ IRT)

Student knowledge & question difficulty share continuous dimension



3

@benjixie
bxie@uw.edu

To understand how good a test is, we have to make a few assumptions.

Following Item Response Theory (IRT), we can assume that student knowledge and question difficulty are on the same continuous dimension.

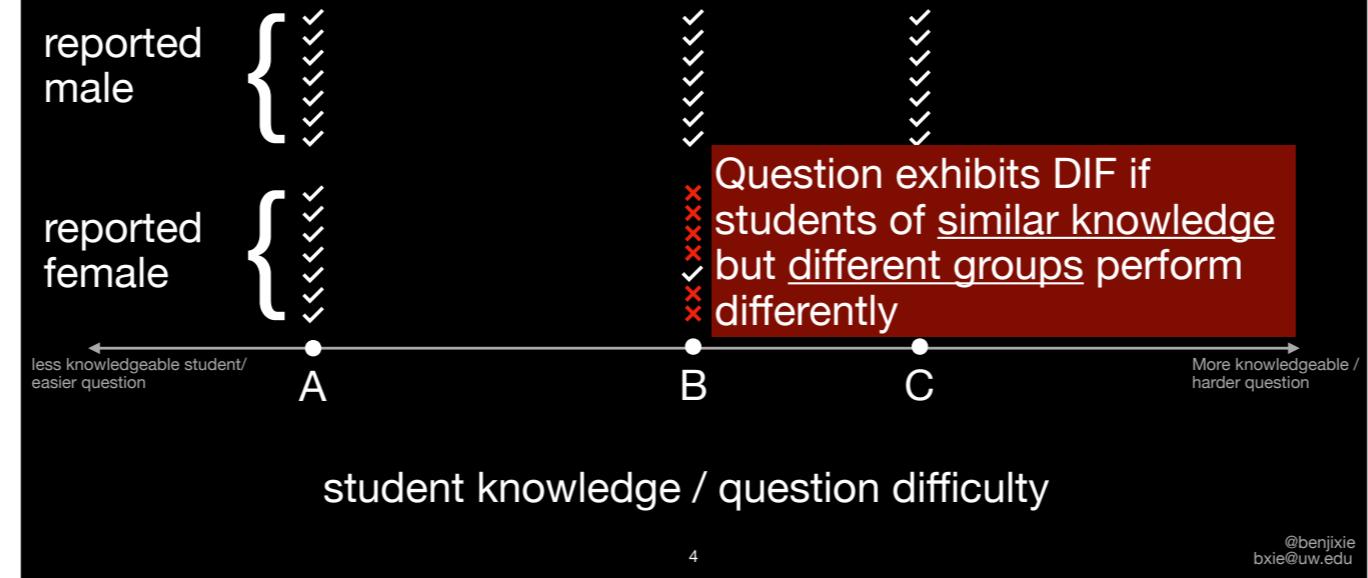
Say we have questions A, B, C, where A is the least difficult question and C is the most difficult.

Say two students, Michelle and Jorge, answer these questions.

Based on their responses to these questions to estimate their knowledge levels. Michelle's is between B and C because she got B correct but C incorrect.

By a similar logic, Jorge's knowledge is between A and B.

Differential Item Functioning (DIF) in test questions



Now say we wanted to look at a group of students who reported as male and a group who reported as female. All students got questions A and C correct. We would expect them to get question B correct as well. And say all the students who reported as male do get B correct.

CLICK

But say we observe that most students who report as female get question B wrong.

CLICK

This is a toy example that demonstrates Differential Item Functioning (DIF)

CLICK

where students of similar knowledge levels but different groups (genders in this case) perform differently on an item, question B in this case.

DIF is a technique to identify potential bias in test questions.

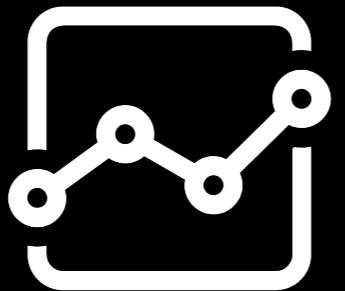
What if DIF could signal opportunities for better pedagogy?

5

@benjixie
bxie@uw.edu

Traditionally, DIF has been used as a filter...
what if we could use it to improve equitably we teach?

Connecting data on DIF w/ domain experts



data on DIF (test bias)



teachers



students



curriculum designers

@benjixie
bxie@uw.edu

Data on DIF can help identify or substantiate nuanced patterns of disparities or bias. But we need the domain expertise of stakeholders such as students, teachers, and curriculum designers to interpret and use these findings to address inequities.

How do domain-experts use data on test bias by gender and race for equity?

7

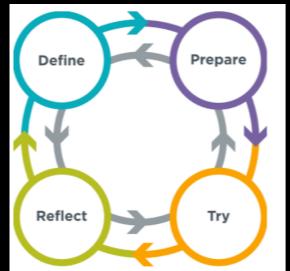
@benjixie
bxie@uw.edu

So this paper explores how domain-experts (curriculum designers) might be able to interpret and use gender and race-based DIF for equity related goals.

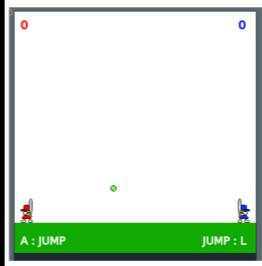
Code.org CS Discoveries (CSD) '19-'20

- Computing as creative form of expression
- 19,617 students (most 11 - 16 yrs old)
- 17 questions for formative use

Unit 2: web dev



Unit 1: problem solving



Unit 3: interactive games

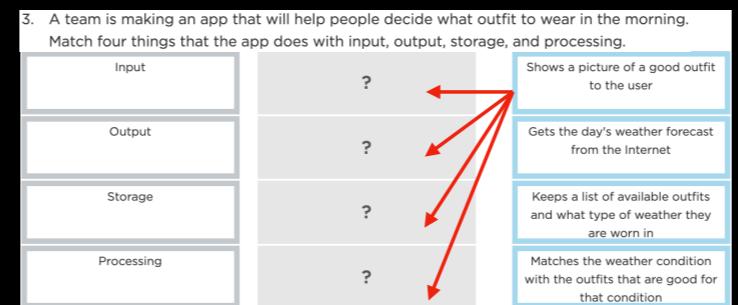
8

@benjxie
bxie@uw.edu

17 CSP items: multiple choice, matching

Matching
(Unit 1, Q4)

Multiple-
choice
(U2, Q3)



3. A group of artists is making a website. When should they use classes?
- A. They want to make their h3 tags bigger than their h1 tags.
 - B. They want headings to be one color, but paragraphs to be a different color.
 - C. They want some images to float left, and other images to float right.
 - D. They want all the pages on their site to have the same style.
 - E. They want make a special color using RGB codes.

9

@benjixie
bxie@uw.edu

The 17 questions I analyzed were either matching questions or multiple choice questions.

Matching questions required students to place options in their correct locations.

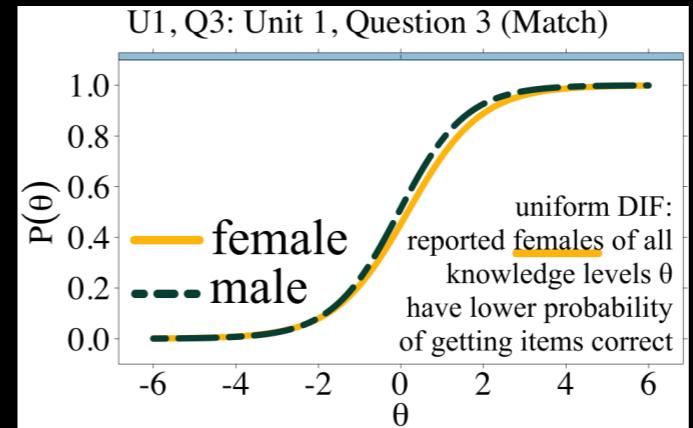
Multiple choice questions required students to choose one or two options.

All questions scored dichotomously (entirely correct or entirely wrong)

Found evidence of bias by gender, race

Question **biased** if >5% difference in probability that students w/ same knowledge, different groups get question correct ($p < 0.001$, medium or large effect size)

- 2 questions disadvantaged reported **females** vs reported **male**
- 13 questions disadvantaged **AHNP** (African/Black, Hispanic/Latinx, Native American/Alaskan Native, and Pacific Islander) vs **WA** (white, Asian)



10

@benjixie
bxie@uw.edu

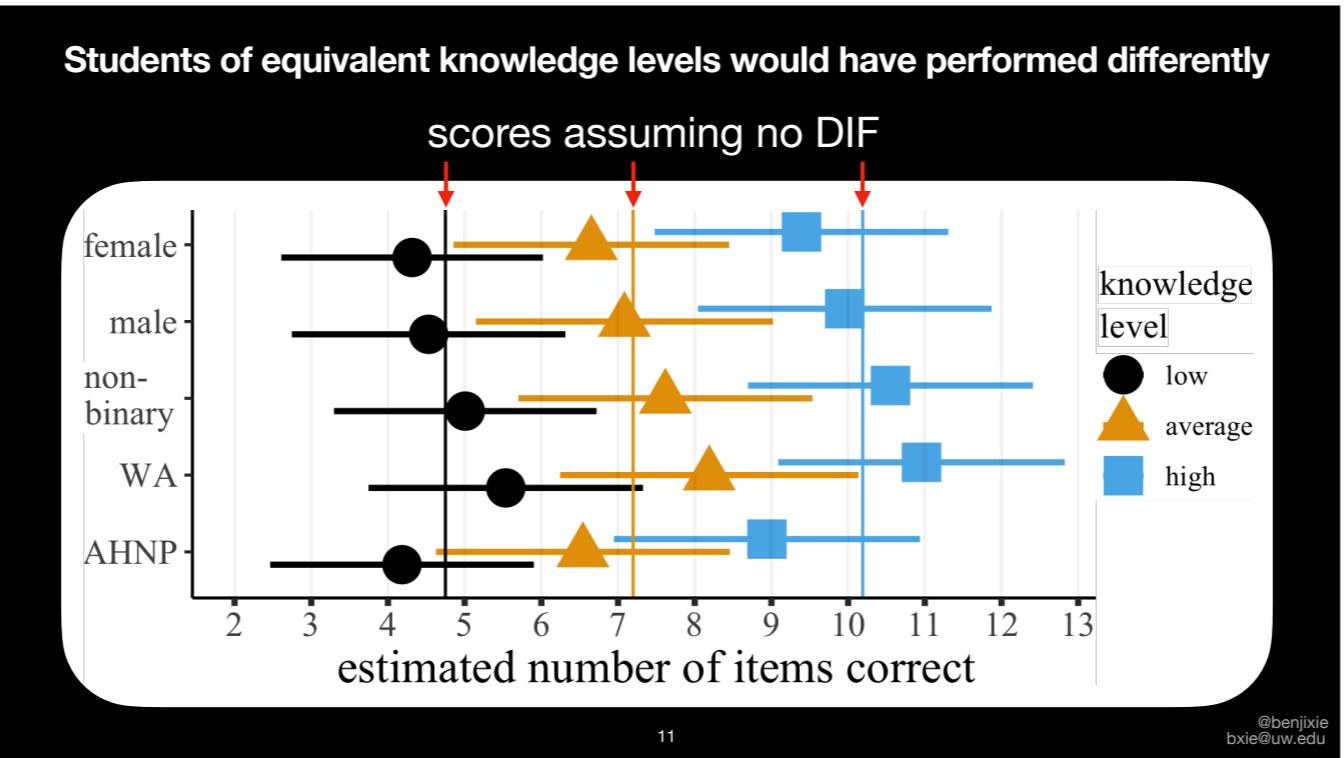
Our quantitative analysis focused on checking for potential test bias by gender and race.

We say a question is biased if on average, a student from a disadvantaged group is at least 5% less likely to get that question correct compared to a student of similar knowledge from the other group. This is equivalent to checking for statistical significance with a medium or large effect size.

We found two questions disadvantaged students who reported as female compared to students who reported as male.

CLICK* and describe trace plot

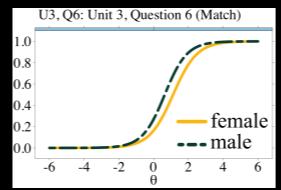
Most test questions disadvantaged AHNP students (African/Black, Hispanic/Latinx, Native/Alaskan Native, and Pacific Islander) compared to WA students (white, Asian)



Put together, we can say that students of equivalent knowledge but different genders or races would score differently on the CSD assessments.

So as a whole, this test disadvantages AHNP and reported female students the most, and advantages WA students the most.

How curriculum designers interpreted DIF



@benjixie
bxie@uw.edu

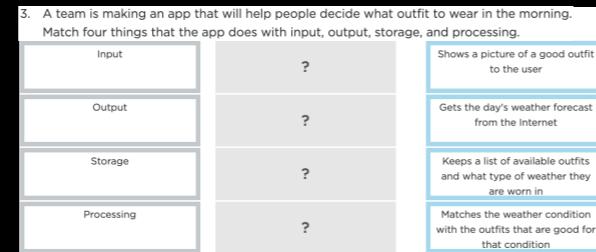
DIF does not tell us the cause of this bias or what to do about it.

So to understand that, we conducted a remote workshop where 7 Code.org curriculum designers interpreted DIF data.

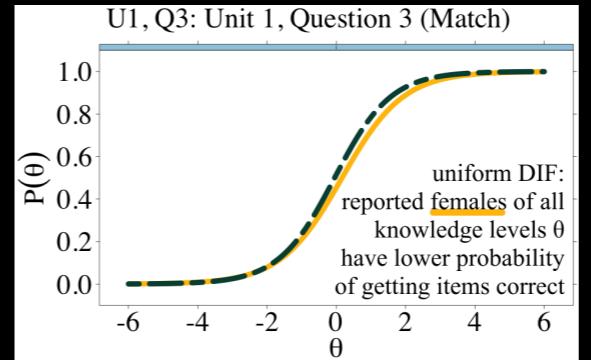
All this was in an effort to understand a new use for DIF: improving equity in learning by informing domain experts of potential issues.

Here are a few high level takeaways, but I point you to the paper to read more about our findings.

Considering question properties relative to student identities



“Female students are performing lower on matching [type] questions...”



Curriculum designers considered how test design may have introduced bias, with some identifying how matching type questions disadvantaged students who reported as female.

Alignment between assessment and curriculum

“Comments are not very well emphasized in CS Discoveries... this **may be the very first [time]** that students are seeing this idea of putting a comment to a block of code.”

6. Look at the following code, and match the comments to where they should go in the program.

```
1 var player = createSprite(200, 200, -1);
2 player.setAnimation("alien");
3 var bubblelet = 400;
4 var bubbleX = 50;
5 // comment
6 stroke("white");
7
8 function draw() {
9   background("airplaneBlue");
10  if (keyWentDown("space")) {
11    // comment
12    player.velocityY = -3;
13    // comment
14    player.velocityY = player.velocityY + 0.2;
15
16  drawSprites();
17  // comment
18  bubbleX = bubbleX + randomNumber(-1, 1);
19
20  ellipse(bubbleX, bubbleY, 10, r, -1);
21
22 }
```

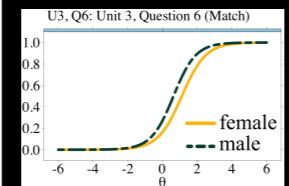
Line 5	?	Jump up
Line 11	?	Fall down
Line 14	?	Shake back and forth
Line 17	?	Give the bubble an outline
Line 19	?	Float up

@benjixie
bxie@uw.edu

Curriculum designers also considered how the curriculum may or may not have prepared students for the test questions.

So in one case, curriculum designers acknowledged that commenting code was a skill worth learning, but may not have been well taught prior to this test question. Considering how specific aspects test and curriculum design may contribute to bias is a potential first step to making changes that support more equitable learning experiences.

Implications: Equitable action by contextualizing data w/ domain expertise



data on DIF
(test bias)

Identify nuanced
patterns, bias



curriculum
designers

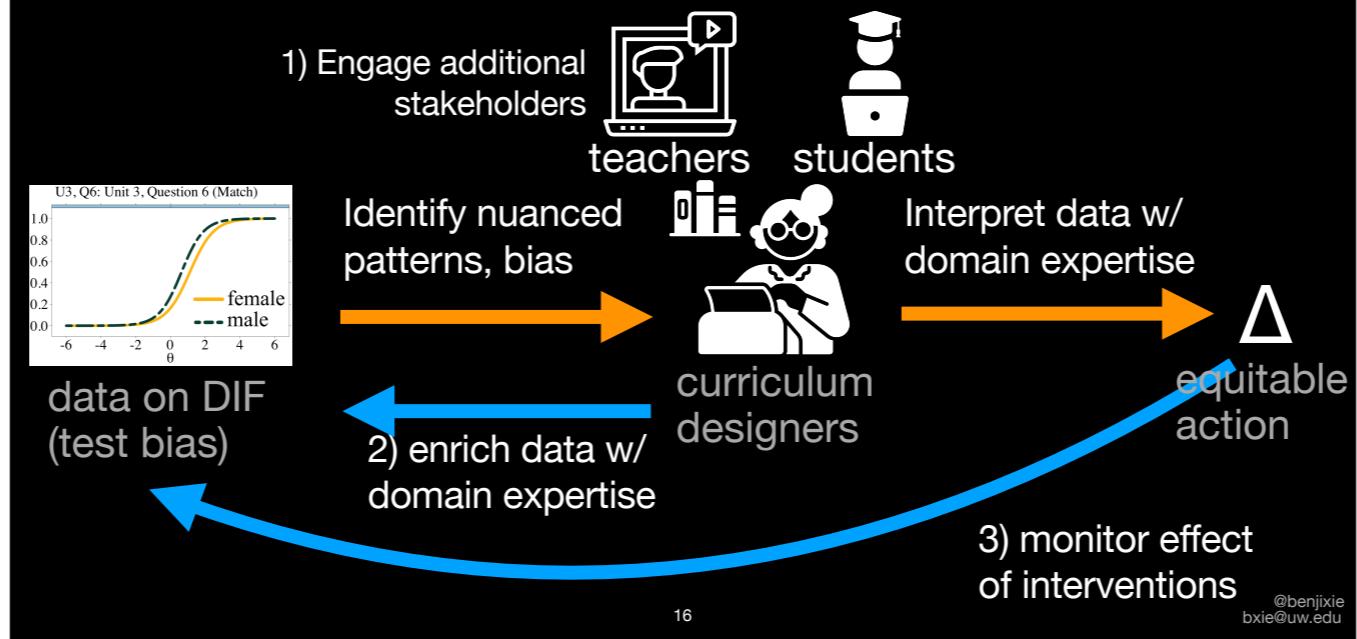
Interpret data w/
domain expertise



equitable
action

Iterating towards more equitable learning experiences requires measuring factors we cannot easily intuit, and using domain expertise to contextualize these findings with understanding we cannot easily measure

Future work: More stakeholders, feedback loops



Some potential future work includes engaging additional stakeholders,
Using human-centered AI techniques to enrich data with domain expertise,
And monitoring the effect of interventions

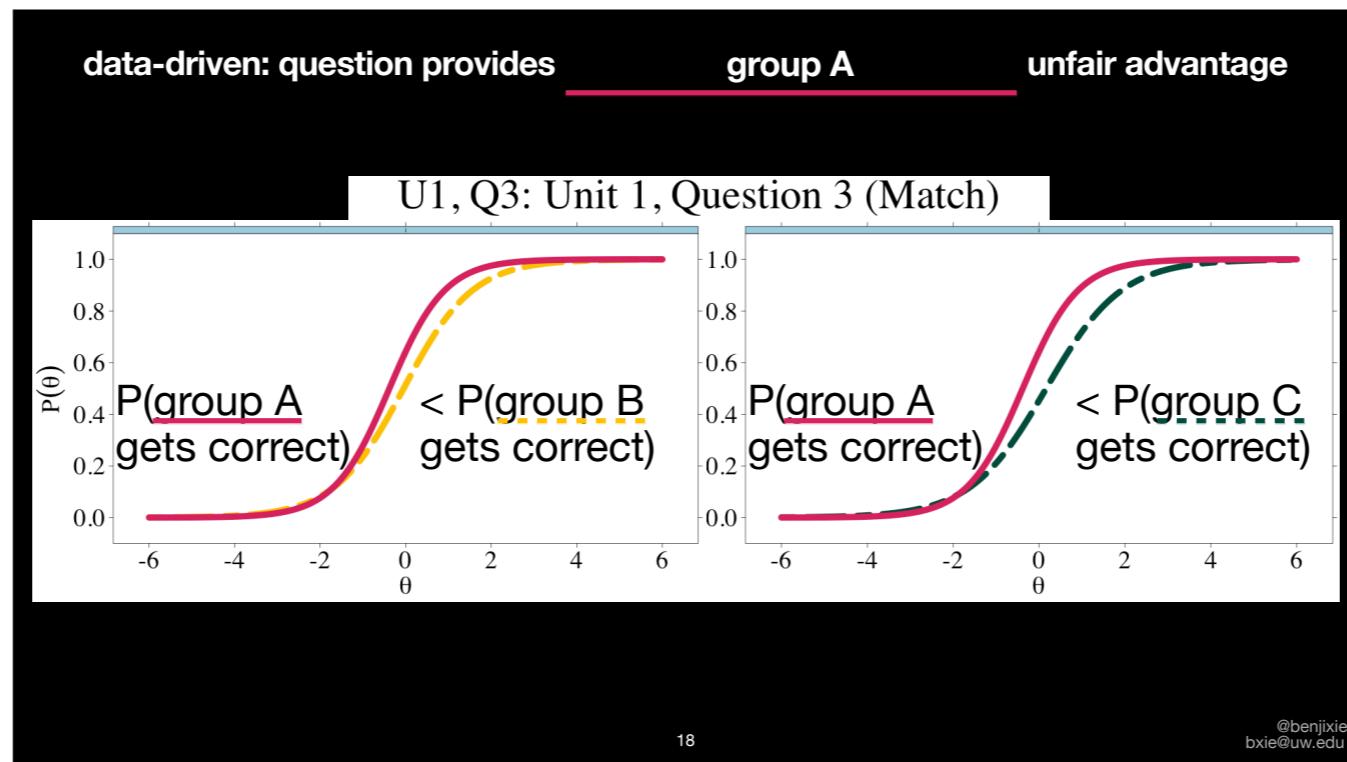
**friendly reminder: data w/o
context can be misleading**

17

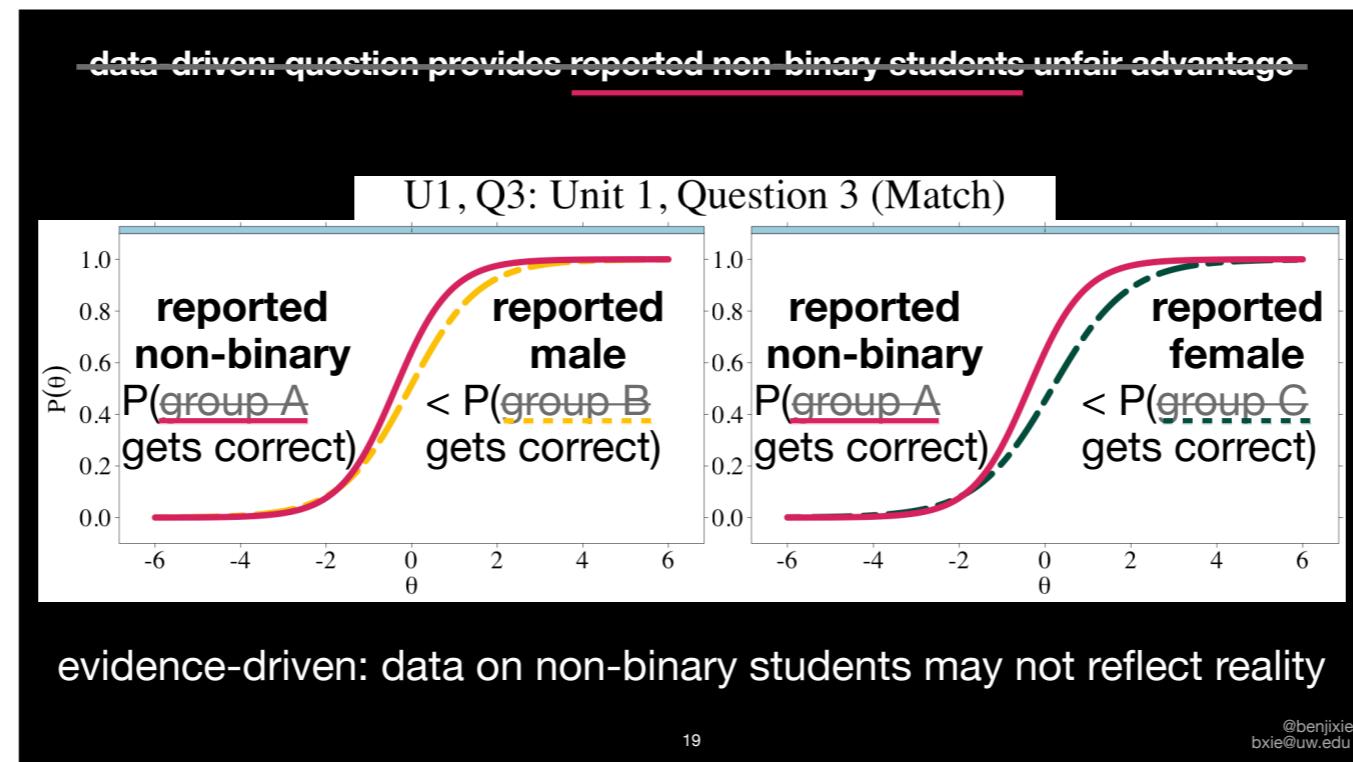
@benjixie
bxie@uw.edu

(time check)

In the spirit of pride month, I did want to share one finding which also serves as yet another reminder that data without context invites misinterpretation



Here are two trace plots which show the probably that learners from different groups will get a test question correct given different latent knowledge levels (θ). These two plots show that this question unfairly advantages group A, the solid red line. So if you had three students of the same knowledge level, the student from group A would still perform better compared to students from groups B and C.



These three groups are actually reported genders students.

So the data-driven result is that students who report as non-binary have an advantage for this item and 4 others, and that we must make changes to better support students who report as male or female.

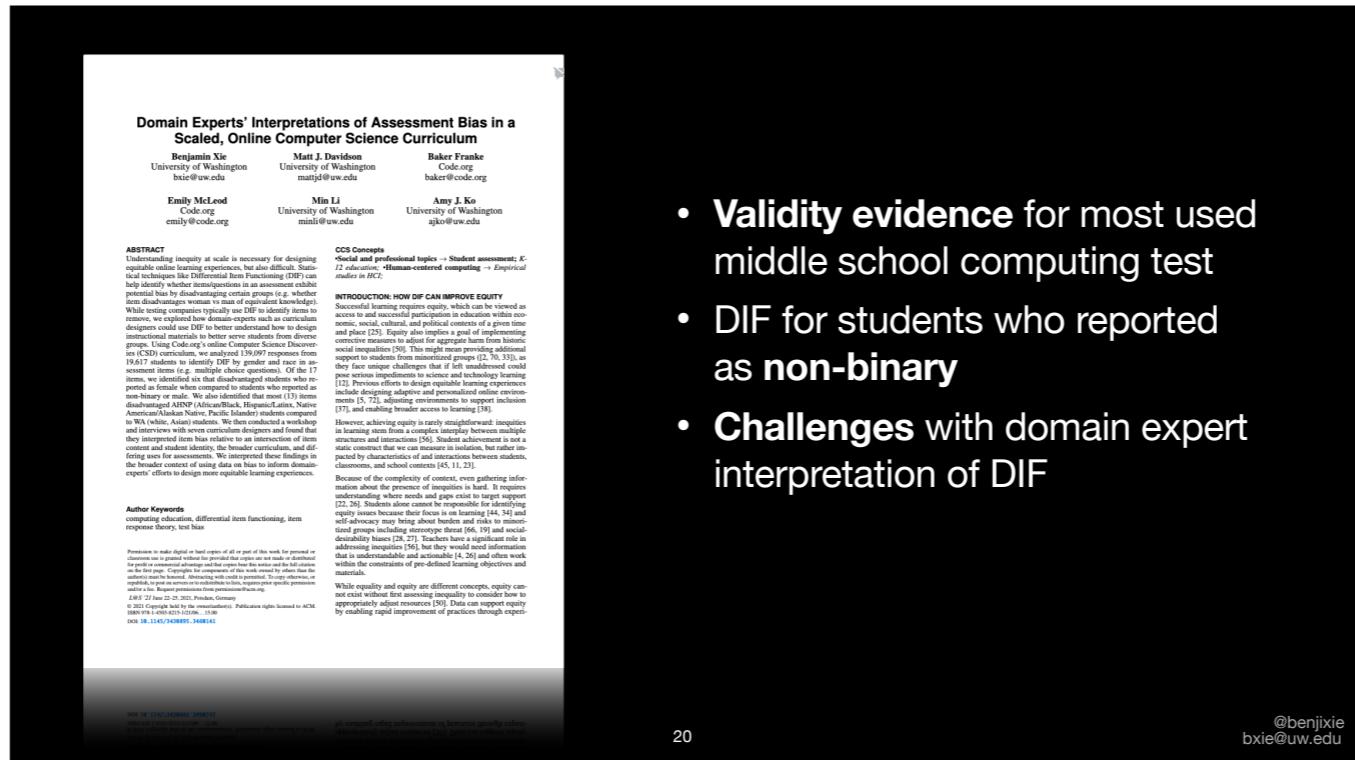
If you have a feeling that something isn't quite right, we did too. Because there is a mountain of computing and STEM education literature detailing the minoritization of non-binary students.

So we consulted with learning scientists and critical scholars and found out that coming out as non-binary in primary school is a constant and challenging process and that a majority of non-binary students may have been assigned female at birth. So the evidence driven result is that our data and labels may not reflect identity of students.

So the more accurate statement is that this data-driven result is actually misleading because of its lack of context. So let's replace that....

CLICK

with a more contextualized evidence-driven finding that suggests that data on non-binary students may not reflect student experience. So all this is to say that, analyzing data is dangerous, so we need to consult with domain-experts to contextualize and validate our supposed findings.



- **Validity evidence** for most used middle school computing test
- DIF for students who reported as **non-binary**
- **Challenges** with domain expert interpretation of DIF

20

@benjixie
bxie@uw.edu

I encourage you to read the paper to learn more about validity evidence, why our analysis of non-binary students' responses required contextualization, and challenges with interpreting DIF.

**Data helps us identify existence
and extent of biases.
Domain expertise helps us identify
causes, take equitable action**

But the main take away is this:

Domain Experts' Interpretations of Assessment Bias in a Scaled, Online Computer Science Curriculum

Benjamin Xie | bxie@uw.edu | [@benjixie](#) | benjixie.com/las21



On job market! Post-doc, research scientist.

- **RQ:** How do domain-experts use data on test bias by gender and race for equity?
- **takeaway:** Data identifies inequalities, unfairness, bias. Equitable action requires engaging domain experts.
- **design implication:** Domain experts can connect bias to causes, identify equitable action

Questions I have:

1. If minoritized groups are small, what should we do with data about them?
2. Ideally, who should be interpreting and using our findings? How do we support them?
3. How should we balance privacy and protection with checking for bias?
4. How do people qualify our labels/representations of them?



Information School
UNIVERSITY of WASHINGTON



COLLEGE OF EDUCATION
UNIVERSITY of WASHINGTON



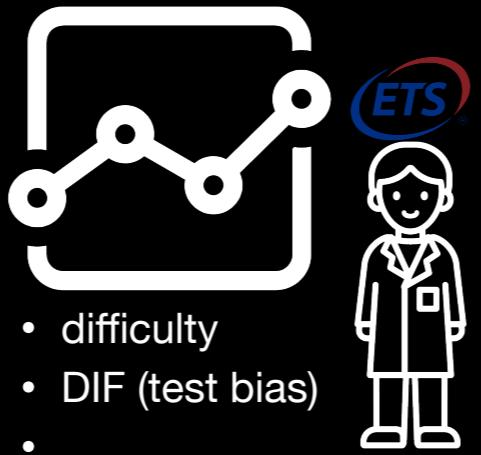
@benjixie
bxie@uw.edu

The link to read the full paper and access supplementary material is benjixie.com/las21.

I'm legally obligated to inform everyone I run into that I'm on the job market, seeking postdoc and research scientist opportunities and would love to chat. And I'll conclude by pointing you to a list of 4 questions that keep me up at night that can inspire some discussion and ideation.

Supplementary slides

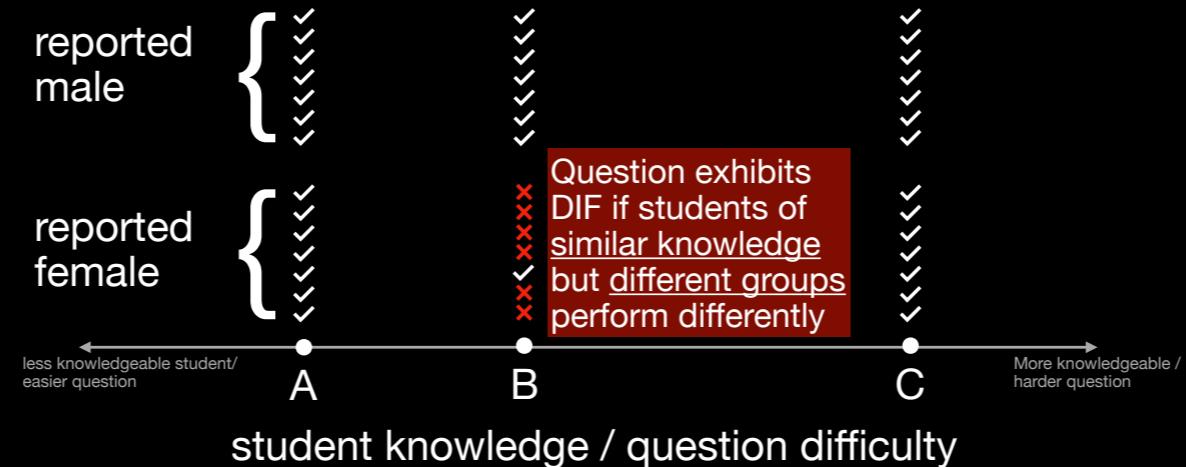
DIF normally used by psychometricians to remove “bad” questions



@benjixie
bxie@uw.edu

DIF is often used by educational testing companies to ensure high stakes exams are fair. For example, say someone at ETS was creating questions for this upcoming years' AP Computer Science Principles exam. They would likely use DIF techniques to identify questions that exhibited DIF and remove them because they may disadvantage certain groups (by gender or race for example).

Differential Item Functioning (DIF) in test questions



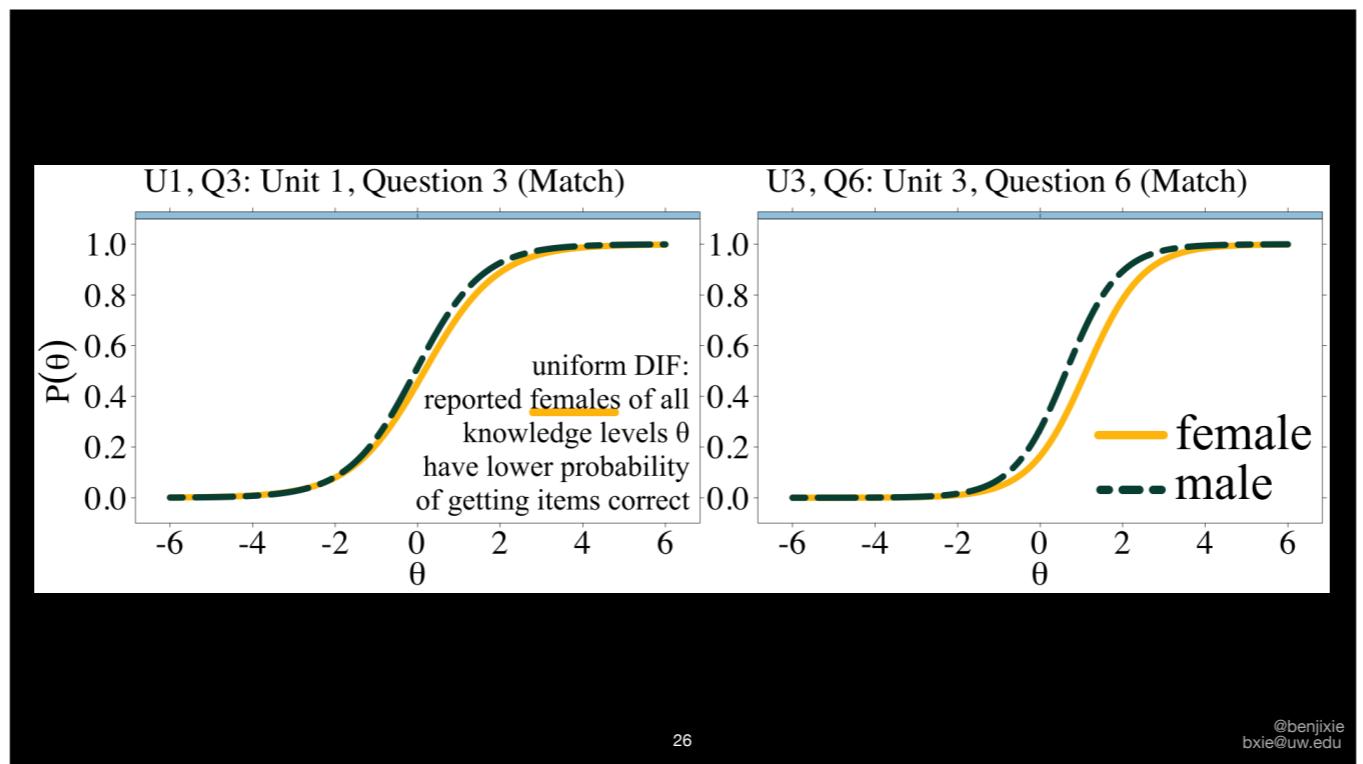
Now say we wanted to look at a group of students who reported as male and a group who reported as female. All students got questions A and C correct.

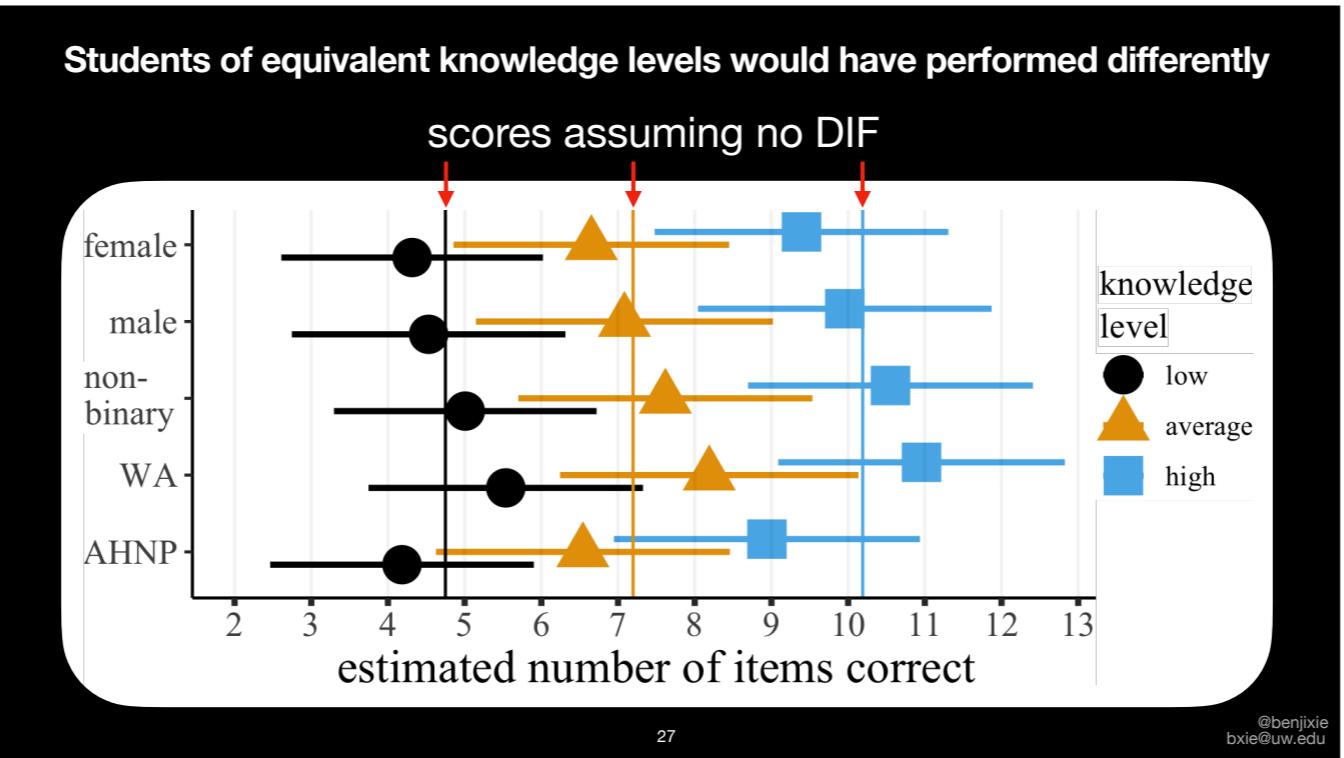
We would expect them to get question B correct as well. And say all the students who reported as male do get B correct.

But say we observe that most students who report as female get question B wrong.

This is a toy example that demonstrates Differential Item Functioning (DIF), where students of similar knowledge levels but different groups (genders in this case) perform differently on an item, question B in this case.

DIF is a technique to identify potential bias in test questions.





Put together, we can say that students of equivalent knowledge but different genders or races would score differently on the CSD assessments.

So as a whole, this test disadvantages AHNP and reported female students the most, and advantages WA students the most.

Questions that keep me up at night

1. If minoritized groups are small,
what should we do with data about them?
2. Ideally, who should be interpreting and using our findings?
How do we support them?
3. How should we balance privacy and protection with checking for
bias?
4. How do people qualify our labels/ representations of them?