

Sean F. Reardon, Demetra Kalogrides, Erin M. Fahle, Anne Podolsky, Rosalia Zarate

## Background

Standardized test construction takes measures to ensure that test items are not biased against any group of students. However, the methods used to remove biased items may fail to identify bias due to test format – the mode of student response. Several studies suggest that males perform better relative to females on multiple-choice items, and females perform better relative to males on short answer and extended-response questions, though the evidence is not conclusive (Lindberg, Hyde, Petersen, & Linn 2010; Beller & Gafni, 2000; Gamer & Engelhard, 1999; DeMars, 1998; Bolger & Kellaghan, 1990).

State accountability test formats vary significantly across the U.S. (Figure 1). Given the evidence that test format may lead to gender bias, it is unclear to what extent variable item formats affect the measurement of gender achievement gaps and “distort” comparisons across states.

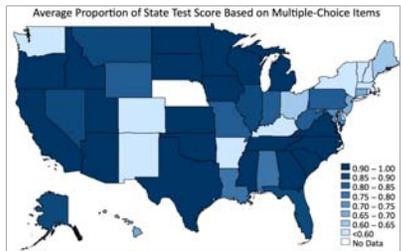


Figure 1

## Current Study

Using the test scores of roughly eight million students tested in the fourth and eighth grades in English Language Arts (ELA) and math during the 2008-09 school year, we estimate state- and district-level math and ELA gender achievement gaps on each state’s accountability tests, as well as on “audit” tests that do not vary in format across states and districts. We then characterize the extent to which the variation of gender achievement gaps on standardized tests across the U.S. can be explained by differing state accountability test formats.

## Data

We use 4<sup>th</sup> and 8<sup>th</sup> grade math and reading achievement data in the 2008-09 school year.

EDFacts State Accountability Test Data: district-level achievement data for all students in all states on their respective state accountability tests, which vary in format across states.

National Assessment for Educational Progress (NAEP) State Data: student-level achievement data in math and reading from a representative sample of students in every state. The format of the NAEP test is the same across states.

Northwest Evaluation Association Measures of Academic Progress (MAP) Data: student-level achievement data for all students in districts that purchase the MAP assessment. The MAP test format is 100% multiple-choice.

## Methods

We compare gender gaps on tests with different variable item structures that are administered to the same population of students. Specifically, we compare gaps on the state assessments with gaps on a test with a constant format across states/districts (NAEP or MAP). In this framework, a systematic association between measured gender gaps and the item structure suggests that item format may differentially affect male and female students’ scores.

We estimate gender achievement gaps with the V-statistic (Ho & Reardon, 2012; Reardon & Ho, 2015). We then model the association as:

$$G_{st} = \gamma_s + \delta p_{st} + u_{st}^*$$

$$= \gamma_s + (\delta p_{sa})T_{st} + (\delta p_n)(1 - T_{st}) + u_{st}^*$$

where  $G_{st}$  is the male-female achievement gap on test  $t$  in state/district  $s$ ,  $p_{st}$  is the proportion of score based on constructed response items on test  $t$  in state  $s$ ;  $\gamma_s$  is the state gap that would be observed if a test had only multiple-choice items; and  $u_{st}^*$  is any other gender bias. The state test is denoted as  $a$ , and the national test as  $n$ .  $T_{st}$  is an indicator for the state test.

Figures 2 (top) & 3 (bottom)

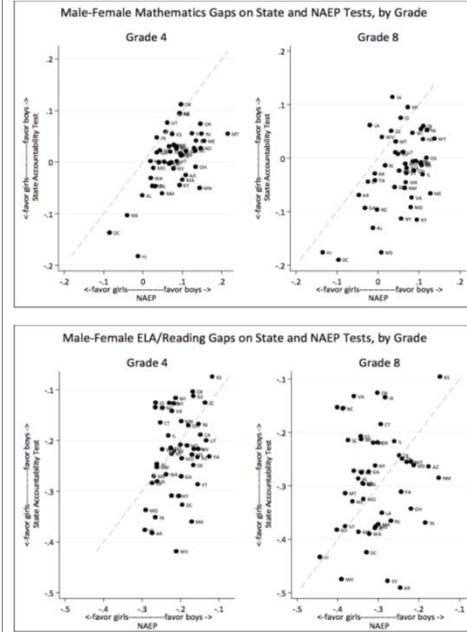


Table 1: Relationship between Proportion of Score from Multiple Choice Items on State Tests and the Size of Gender Gaps

	State-Level					District Level				
	Pooled Across Grades and Subjects	Math	ELA	Math	ELA	Pooled Across Grades and Subjects	Math	ELA	Math	ELA
	Grade 4	Grade 8	Grade 4	Grade 8	Grade 4	Grade 8	Grade 4	Grade 8	Grade 4	Grade 8
<b>Model 1</b>										
Proportion Short Response+ Extended Response	-0.208*** (0.050)	-0.135** (0.041)	-0.109 (0.075)	-0.223* (0.084)	-0.376** (0.113)	-0.241** (0.069)	-0.126 (0.122)	-0.151* (0.068)	-0.296** (0.098)	-0.389*** (0.101)
P(Grade/Subject Interactions Are Equal)	.118					.060				
<b>Model 2</b>										
Proportion Short Response	-0.146* (0.064)	-0.078 (0.061)	-0.067 (0.137)	-0.177+ (0.094)	-0.264* (0.122)	-0.248** (0.087)	-0.159 (0.103)	-0.160 (0.134)	-0.358 (0.213)	-0.326** (0.107)
Proportion Extended Response	-0.264*** (0.050)	-0.178*** (0.044)	-0.141 (0.135)	-0.267* (0.102)	-0.5*** (0.115)	-0.234* (0.099)	-0.088 (0.178)	-0.144 (0.104)	-0.246+ (0.128)	-0.437*** (0.109)
P(Short Response=Extended Response)	.127	.189	.749	.391	.090	.913	.625	.936	.684	.259
N	382	96	94	96	96	5746	1530	1380	1534	1302

Beller, M., & Gafni, N. (2000). Can item format (multiple-choice vs. open-ended) account for gender differences in mathematics achievement? *Sex Roles*, 42(1-2), 1-21.  
 Bolger, N., & Kellaghan, T. (1990). Method of measurement and gender differences in scholastic achievement. *Journal of Educational Measurement*, 27(2), 165-174.  
 DeMars, C. E. (1998). Gender differences in mathematics and science on a high school proficiency exam: The role of response format. *Applied Measurement in Education*, 11(3), 279-299.  
 Gamer, M., & Engelhard Jr, G. (1999). Gender differences in performance on multiple-choice and constructed-response

mathematics items. *Applied Measurement in Education*, 12(1), 29-51.  
 Ho, A. D., & Reardon, S. F. (2012). Estimating achievement gaps from test scores reported in ordinal “proficiency” categories. *Journal of Educational and Behavioral Statistics*, 37(4), 489-517.  
 Lindberg, S. M., Hyde, J. S., Petersen, J. L., & Linn, M. C. (2010). New trends in gender and mathematics performance: a meta analysis. *Psychological bulletin*, 136(6), 1123.  
 Reardon, S. F., & Ho, A. D. (2015). Practical Issues in Estimating Achievement Gaps From Coarsened Data. *Journal of Educational and Behavioral Statistics*, 40(2), 158-189.