



Development of a Complex Reasoning (CR) Test

Presentation to the Defense Advisory Committee on Military Personnel Testing (DACMPT)

Presenter: Mike Ingerick, Human Resources Research Organization (HumRRO)

December 16, 2022

BACKGROUND

• What is complex reasoning?

 Non-verbal reasoning; ability to analyze visual information and to solve problems using visual reasoning

• Why a complex reasoning test?

- Fluid intelligence has been found to be a strong predictor of training and job success
 - Complex (non-verbal) reasoning is one element of fluid intelligence
 - ASVAB Review Panel (2006) recommended that DoD consider adding tests of fluid intelligence to balance the ASVAB's composition (between fluid and crystalized intelligence)
- Potential benefits to the ASVAB testing program
 - Improved prediction of training and job success in military jobs
 - Lower susceptibility to test compromise
 - Less adverse impact; increased qualification rates for non-native and non-heritage English speakers

INITIAL DEVELOPMENT EFFORT: ABSTRACT REASONING TEST (ART)

- Developed by Susan Embretson
- Item format similar to Raven's Progressive Matrices (RPM)
 - Multiple choice, 6 or 8 response options per item
- DTAC commissioned the development of one form (30 items)
- Administered (for research purposes) to language training applicants (2017)
- Items were found to be relatively easy, time-consuming



Sample ART Item

CURRENT DEVELOPMENT EFFORT

Objective: Develop a complex (non-verbal) reasoning testing system to generate items for potential inclusion on ASVAB

- Employ non-proprietary Automated Item Generation (AIG) capability
 - Improve item development efficiency
 - Reduce or eliminate field-testing requirements
- Generate items with targeted properties
 - Items similar to Raven's Progressive Matrices (RPM) items
 - Items at appropriate difficulty for qualifying military applicants into jobs of varying complexity

CURRENT DEVELOPMENT EFFORT (CONTINUED)



PHASE 1 PILOT: DESIGN

Pilot intended to answer the following questions:

- Does performance differ by item type (transformation vs. logic)?
 Does performance align with item difficulty specs?
- Does the number of response options impact performance? Does including a "None of these are correct" (NOTAC) option impact performance?
- Does performance differ by gender, race, or ethnicity?
- How much time is needed to complete the items? Does completion time differ by item type or response option set? Does completion time differ by gender, race, or ethnicity?

PHASE 1 PILOT: DESIGN (CONTINUED)

Look at the 3 x 3 grid below. Identify the pattern(s).



Which of the following images best completes the pattern(s) in the grid?



Sample Transformation Item

Transformation item features:

- Number of transformations
 - Types of shapes
 - Orientation of shape(s)
 - Size of shape(s)
 - Number of shape(s)
 - Line weighting on shape(s)
- Direction(s) of transformations
 - Vertical
 - Horizontal
 - Diagonal

PHASE 1 PILOT: DESIGN (CONTINUED)

Logic item feature:

- Nature of logic rule
 - AND
 - OR
 - XOR (Exclusive OR)

Look at the 3 x 3 grid below. Identify the pattern(s).



Which of the following images best completes the pattern(s) in the grid?



Sample Logic Item ("AND")

PHASE 1 PILOT: DESIGN (CONTINUED)



Transformation Only



Condition 6: Logic First, Transformation Second

Condition 7: Scrambled

Transformation + Logic

PHASE 1 PILOT: DATA COLLECTION

Sample

- Non-military sample, ages 18–35, U.S. citizen
- Targeted *N* = 3,500 participants

Measures

- One complex reasoning form (24 items)
- Post-test questionnaire (demographics, perceived difficulty of items, testtaking experience)
- Two attention check items

Method

- Administered on Qualtrics platform
- No fixed time limit; recorded time to completion
- Desktop or laptop only

Dates of Data Collection:

19 July – 4 August (~2 weeks)

	N	% of Total Invited
Invited to participate in pilot	7,039	
Accessed pilot	4,459	63.3%
Completed pilot	3,778	53.7%
Completed pilot w/ valid data	3,491	49.6%

Final Analysis Sample for Pilot

PHASE 1 PILOT: DEMOGRAPHICS AND EDUCATIONAL BACKGROUND OF FINAL ANALYSIS SAMPLE All Participants

							% HS
Condition							Degree/GED/
Description	N	Mean Age	% Female	% Asian	% Black	% Hispanic	< 1 yr College
3 + NOTAC	501	27.2	59.3	6.6	10.8	12.0	36.3
4	503	26.7	62.2	6.2	11.3	9.5	37.9
4 + NOTAC	490	26.7	59.8	5.9	15.5	10.2	35.1
8	493	27.0	61.1	5.5	10.8	11.6	36.5
Transform 1 st , Logic 2 nd	507	26.4	59.0	6.5	11.4	10.1	32.3
Scrambled	492	26.6	57.3	7.7	11.6	10.0	33.3
Logic 1 st , Transform 2 nd	505	26.9	61.8	6.3	13.1	8.9	30.0
Overall	3,491	26.8	60.1	6.4	12.1	10.3	34.5

Participants Completed in 30 Min or Less with HS Degree, GED, or < 1 yr of College

Condition							% HS		% < 1 yr
Description	N	Mean Age	% Female	% Asian	% Black	% Hispanic	Degree	% GED	College
3 + NOTAC	182	25.2	52.7	5.5	16.5	23.1	65.9	9.3	24.7
4	188	25.7	59.2	5.3	14.9	22.3	62.8	14.7	22.5
4 + NOTAC	172	25.2	60.4	4.6	22.6	16.2	70.3	11.6	18.0
8	179	24.9	61.6	3.4	13.4	25.7	64.4	11.7	23.9
Transform 1 st , Logic 2 nd	164	24.7	53.7	4.9	19.5	19.5	65.2	9.8	25.0
Scrambled	164	25.4	52.6	7.9	13.8	21.0	67.8	7.9	24.3
Logic 1 st , Transform 2 nd	151	25.0	64.0	3.6	12.2	18.3	65.2	9.7	25.0
Overall	1,200	25.2	57.7	5.0	16.1	20.9	65.9	10.7	23.3

PHASE 1 PILOT: SUMMING UP (ALL PARTICIPANTS, N = 3,491)

	Transform Only	Transform Only	Transform Only	Transform Only	Transform + Logic	Transform + Logic
Metric	8	3 + NOTAC	4	4 + NOTAC	Grouped	Scrambled
Unidimensionality	Yes	No	Yes	No	No	No
	α = .87	α = .78	α = .85	α = .83	α = .75	α = .75
Reliability	SEM = 1.98	SEM = 1.97	SEM = 1.99	SEM = 1.96	SEM = 2.10	SEM = 2.04
	<i>avg</i> CITC = .46	<i>avg</i> CITC = .31	<i>avg</i> CITC = .42	<i>avg</i> CITC = .37	<i>avg</i> CITC = .32	<i>avg</i> CITC = .30
	<i>M</i> = 12.29	<i>M</i> = 13.29	<i>M</i> = 15.00	<i>M</i> = 11.96	<i>M</i> = 10.32	<i>M</i> = 9.15
Observed Difficulty	SD = 5.48	SD = 4.21	SD = 5.15	<i>SD</i> = 4.76	SD = 4.19	<i>SD</i> = 4.08
	avg <i>p</i> = .51	avg <i>p</i> = .54	avg <i>p</i> = .63	avg <i>p</i> = .50	avg <i>p</i> = .43	avg <i>p</i> = .38
	F-M <i>d</i> = .21	F-M <i>d</i> = .31	F-M <i>d</i> = .22	F-M <i>d</i> = .16	F-M <i>d</i> = .10	F-M <i>d</i> = .22
Group Score	B-W <i>d</i> =39	B-W <i>d</i> =31	B-W <i>d</i> =58	B-W <i>d</i> =20	B-W <i>d</i> =34	B-W <i>d</i> =08
Differences	H-W <i>d</i> =17	H-W <i>d</i> =15	H-W <i>d</i> =21	H-W <i>d</i> =19	H-W <i>d</i> =22	H-W <i>d</i> = .04
	A-W <i>d</i> = .36	A-W <i>d</i> = .17	A-W <i>d</i> = .00	A-W <i>d</i> = .12	A-W <i>d</i> = .59	A-W <i>d</i> = .24
Completion Time	<i>M</i> = 12.74	<i>M</i> = 10.88	<i>M</i> = 11.36	<i>M</i> = 12.39	<i>M</i> = 13.54	<i>M</i> = 13.24
(30 minutes or less)	SD = 5.86	SD = 4.94	<i>SD</i> = 5.14	SD = 5.79	<i>SD</i> = 6.18	<i>SD</i> = 6.36
Derecived Difficulty	<i>M</i> = 3.92	<i>M</i> = 3.89	<i>M</i> = 3.98	<i>M</i> = 3.90	<i>M</i> = 3.50	<i>M</i> = 3.37
	SD = .95	SD = .92	SD = .89	SD = .95	SD = .96	<i>SD</i> = 1.00

PHASE 1 PILOT: SUMMING UP (ALL PARTICIPANTS, N = 3,491)

	Transform Only	Transform Only	Transform Only	Transform Only	Transform + Logic	Transform + Logic
Metric	8	3 + NOTAC	4	4 + NOTAC	Grouped	Scrambled
Unidimensionality	Yes	No	Yes	No	No	No
	α = .87	α = .78	α = .85	α = .83	α = .75	α = .75
Reliability	SEM = 1.98	SEM = 1.97	SEM = 1.99	SEM = 1.96	SEM = 2.10	SEM = 2.04
	<i>avg</i> CITC = .46	<i>avg</i> CITC = .31	<i>avg</i> CITC = .42	<i>avg</i> CITC = .37	<i>avg</i> CITC = .32	<i>avg</i> CITC = .30
	<i>M</i> = 12.29	<i>M</i> = 13.29	<i>M</i> = 15.00	<i>M</i> = 11.96	<i>M</i> = 10.32	<i>M</i> = 9.15
Observed Difficulty	<i>SD</i> = 5.48	SD = 4.21	SD = 5.15	SD = 4.76	SD = 4.19	SD = 4.08
	avg <i>p</i> = .51	avg <i>p</i> = .54	avg <i>p</i> = .63	avg <i>p</i> = .50	avg <i>p</i> = .43	avg <i>p</i> = .38
	F-M <i>d</i> = .21	F-M <i>d</i> = .31	F-M <i>d</i> = .22	F-M <i>d</i> = .16	F-M <i>d</i> = .10	F-M <i>d</i> = .22
Group Score	B-W <i>d</i> =39	B-W <i>d</i> =31	B-W <i>d</i> =58	B-W <i>d</i> =20	B-W <i>d</i> =34	B-W <i>d</i> =08
Differences	H-W <i>d</i> =17	H-W <i>d</i> =15	H-W <i>d</i> =21	H-W <i>d</i> =19	H-W <i>d</i> =22	H-W <i>d</i> = .04
	A-W <i>d</i> = .36	A-W <i>d</i> = .17	A-W <i>d</i> = .00	A-W <i>d</i> = .12	A-W <i>d</i> = .59	A-W <i>d</i> = .24
Completion Time	<i>M</i> = 12.74	<i>M</i> = 10.88	<i>M</i> = 11.36	<i>M</i> = 12.39	<i>M</i> = 13.54	<i>M</i> = 13.24
(30 minutes or less)	<i>SD</i> = 5.86	SD = 4.94	SD = 5.14	SD = 5.79	<i>SD</i> = 6.18	SD = 6.36
Derectived Difficulty	<i>M</i> = 3.92	<i>M</i> = 3.89	<i>M</i> = 3.98	<i>M</i> = 3.90	<i>M</i> = 3.50	<i>M</i> = 3.37
Perceived Difficulty	SD = .95	SD = .92	SD = .89	SD = .95	SD = .96	<i>SD</i> = 1.00

PHASE 1 PILOT: SUMMING UP (COMPLETED \leq 30 MINUTES WITH HS DEGREE/GED/< 1 YR OF COLLEGE, N = 1,200)

Transform Only		Transform Only	Transform Only	Transform Only	Transform + Logic	Transform + Logic
Metric	8	3 + NOTAC	4	4 + NOTAC	Grouped	Scrambled
Unidimensionality						
	α = .88	α = .79	α = .86	α = .84	α = .67	α = .69
Reliability	SEM = 1.78	SEM = 1.78	SEM = 1.92	SEM = 1.76	SEM = 2.11	SEM = 1.98
	<i>avg</i> CITC = .39	<i>avg</i> CITC = .28	<i>avg</i> CITC = .40	avg CITC = .37	<i>avg</i> CITC = .27	<i>avg</i> CITC = .26
	<i>M</i> = 11.60	M = 12.56	<i>M</i> = 13.89	<i>M</i> = 11.66	<i>M</i> = 9.71	<i>M</i> = 8.27
Observed Difficulty	SD = 5.15	SD = 3.89	SD = 5.12	SD = 4.39	SD = 3.68	SD = 3.56
	avg <i>p</i> = .48	avg <i>p</i> = .52	avg <i>p</i> = .58	avg <i>p</i> = .49	avg <i>p</i> = .40	avg <i>p</i> = .35
Group Score Differences	F-M <i>d</i> = .20	F-M <i>d</i> = .26	F-M <i>d</i> = .26	F-M <i>d</i> = .18	F-M <i>d</i> =02	F-M <i>d</i> = .34
Completion Time	<i>M</i> = 12.32	<i>M</i> = 10.59	<i>M</i> = 11.34	<i>M</i> = 12.04	<i>M</i> = 13.24	<i>M</i> = 12.79
(30 minutes or less)	<i>SD</i> = 5.70	SD = 4.91	<i>SD</i> = 5.41	<i>SD</i> = 5.51	<i>SD</i> = 6.16	<i>SD</i> = 6.21
Derectived Difficulty	<i>M</i> = 3.83	M = 3.79	<i>M</i> = 3.86	<i>M</i> = 3.83	<i>M</i> = 3.41	<i>M</i> = 3.26
Perceived Difficulty	SD = .97	SD = .93	<i>SD</i> = .91	SD = .96	SD = .97	<i>SD</i> = 1.00

14

PHASE 2: RECOMMENDED ITEM SPECS

- Transformation items only
- No NOTAC response option
- Four response options
 - Refine item difficulty model and item selection to ensure appropriate level of difficulty and minimize group score differences by raceethnicity, where feasible

PHASE 2: PROPOSAL FOR PILOT

Objective

- Collect data on refined pool of CR items representative of the population of CR items with a participant sample representative of military applicants
- Use results to refine CR item specs and select potential pool of CR items for follow-on research (Computational Thinking) and field testing

Design and Measures

- 125+ CR items, multiple forms (5 or more forms, 25 items each, with a subset of items common to all forms)
- General mental ability (GMA) test (e.g., retired form or items from APT/PiCAT)
- Post-test questionnaire (demographics, perceived difficulty of items, test-taking experience)
- Two attention check items

Sample

- Non-military sample, ages 18–35, U.S. citizen, HS Degree/GED/<1 year of college
- Targeted *N* = 3,000 participants (~ 600 participants per form)

Method

- Administer on Qualtrics platform
- Participants randomly assigned to one CR form
- No fixed time limit; record time to completion
- Desktop or laptop only

ACKNOWLEDGMENTS

- Scott Oppler, HumRRO
- Matthew Brown, HumRRO
- Katherine Klein, HumRRO

Backup Slides









Sample Transformation Item



Look at the 3 x 3 grid below. Identify the pattern(s).

Which of the following images best completes the pattern(s) in the grid?







Sample Transformation Item



Sample Logic Item

Which of the following images best completes the pattern(s) in the grid?









Look at the 3 x 3 grid below. Identify the pattern(s).

Sample Logic Item



PHASE 1 PILOT: PRACTICE-ORDER EFFECTS (TRANSFORMATION ONLY CONDITIONS)

Condition Label	Condition Description	n	М	SD	Min	Max
C1	3 + NOTAC	260	13.29	4.21	4	22
C2	3 + NOTAC	241	12.87	4.12	1	22
C3	4	249	15.00	5.15	0	24
C4	4	254	15.28	5.21	2	24
C5	4 + NOTAC	245	12.12	4.64	2	24
C6	4 + NOTAC	245	11.96	4.76	1	24
C7	8	239	12.50	5.69	0	23
C8	8	254	12.09	5.27	0	24
C9 & C10	Transform 1 st , Logic 2 nd	507	10.55	4.16	1	24
C11 & C14	Scrambled	492	9.15	4.08	1	23
C12 & C13	Logic 1 st , Transform 2 nd	505	10.09	4.23	0	24

• Observed no significant effects for:

- Block order, *F*(1, 1979) = 0.68, *p* = .41
- Item block, F(1, 1979) = 1.73, p = .19
- Block order * Item block, *F*(3, 1979) = 1.14, *p* = .65
- Response option condition * Block order * Item block, F(3, 1979) = 0.89, p = .44

Observed significant effects for:

- Response option condition, *F*(3, 1979) = 40.96, *p* < .001
- Response option condition * Item block, F(1, 1979) = 38.55, p < .001
 - Item block effects were only observed in two of four response option conditions
 - 3 + NOTAC (d = .31 favoring Block 1)
 - 8 option conditions (d = .22 favoring Block 2)

PHASE 1 PILOT: PRACTICE-ORDER EFFECTS (TRANSFORMATION + LOGIC CONDITIONS)

Condition Label	Condition Description	n	М	SD	Min	Max
C1	3 + NOTAC	260	13.29	4.21	4	22
C2	3 + NOTAC	241	12.87	4.12	1	22
C3	4	249	15.00	5.15	0	24
C4	4	254	15.28	5.21	2	24
C5	4 + NOTAC	245	12.12	4.64	2	24
C6	4 + NOTAC	245	11.96	4.76	1	24
C7	8	239	12.50	5.69	0	23
~~~	0	054	10.00	E 07		04
C9 & C10	Transform 1 st , Logic 2 nd	507	10.55	4.16	1	24
C11 & C14	Scrambled	492	9.15	4.08	1	23
C12 & C13	Logic 1 st , Transform 2 nd	505	10.09	4.23	0	24

#### Observed no significant effects for:

- Logic block order, *F*(1, 1498) = 0.30, *p* = .86
- Response option conditions, F(3, 1498) = 2.14, p = .12
- Logic block order * Logic item block, F(1, 1498) = 0.08, p = .77
- Logic block order * response option condition, F(2, 1498) = 0.66, p = .52
- Logic item block * response option condition, F(2, 1498) = 0.09, p = .91
- Response option condition * Logic block order * Logic item block, F(2, 1498) = 0.40, p = .67
- Observed significant effect for:
  - Logic item block, *F*(1, 1498) = 7.08, *p* = .01

# PHASE 1 PILOT: DIMENSIONALITY

### **Modified Parallel Analyses**

- Compared observed, second-factor eigenvalues to simulated second-factor eigenvalues (from 100 randomly generated samples) to determine whether the items were unidimensional
- Transformation items met the criteria for unidimensionality in the Transformation only conditions, except for the NOTAC conditions (3 + NOTAC, 4 + NOTAC)
- Items in all three Transformation + Logic conditions were <u>not</u> unidimensional, presentation order made no difference

Condition			Observed	Simulated	
Label	<b>Condition Description</b>	n	Eigenvalue	Eigenvalue	р
C1 + C2	3 + NOTAC	501	3.32	1.78	.01
C3 + C4	4	503	1.80	1.37	.13
C5 + C6	4 + NOTAC	490	3.65	1.82	.01
C7 + C8	8	493	1.81	1.47	.06
C9 + C10	Transform 1 st , Logic 2 nd	507	2.31	1.29	.01
C11 + C14	Scrambled	492	2.07	1.20	.02
C12 + C13	Logic 1 st , Transform 2 nd	505	2.58	1.23	.01

# PHASE 1 PILOT: DIMENSIONALITY (CONTINUED)

#### **Confirmatory Factor Analysis**

• Transformation & Logic Items (responses from Transformation + Logic conditions, combined)

- Two-factor model provided significantly superior model fit compared to a single-factor model ( $\Delta \chi^2$  (1) = 406.555,  $\Delta$ CFI = .08,  $\Delta$ RMSEA = .01)

Model	Model χ ²	df	CFI	RMSEA	SRMR	∆ <b>X²</b>	∆CFI	
One-factor	1297.65	252	.79	.05	.05			
Two-factor	891.09	251	.87	.04	.05	406.55	.08	.01

#### • NOTAC Items (responses from 3 + NOTAC and 4+ NOTAC, combined)

- Two-factor model provided better overall model fit compared to the one-factor model ( $\Delta \chi^2$  (1) = 255.28,  $\Delta CFI$  = .08,  $\Delta RMSEA$  = .01)
  - Second factor contained all six items where the NOTAC options was the keyed (correct) answer
- Model fit was further improved from the two-factor model by cross-loading all six items to both latent factors but constraining the correlation between factors ( $\Delta \chi^2$  (5) = 95.16,  $\Delta CFI$  = .03,  $\Delta RMSEA$  = .00)

Model	Model x ²	df	CFI	RMSEA	SRMR	∆ <b>X²</b>	∆CFI	
One-factor	1065.98	252	.74	.06	.08			
Two-factor	810.70	251	.82	.05	.08	255.28	.08	.01
Alt Two-factor *	742.32	246	.85	.05	.07	68.38	.03	.00

# PHASE 1 PILOT: INTERNAL CONSISTENCY (ALL PARTICIPANTS)

The NOTAC and Transformation + Logic items were less internally consistent.

- The four-response option and eight-response option conditions showed slightly higher internal consistency ( $\alpha$  = .85 .87) than the NOTAC conditions ( $\alpha$  = .78 .83) [CITCs (avg *r* = .42 .46) compared to the NOTAC items (avg *r* = .32 .37)]
- Internal consistency estimates for the Transformation + Logic items were < .80, regardless of presentation order (α's = .74 - .76)
- Among the NOTAC conditions, performance on the six items where NOTAC was the keyed (correct) response were less correlated with performance on the other NOTAC items
  - CITCs (avg r = .47, ranging from .29 to .60) compared to the NOTAC items (avg r = .20, ranging between -.19 and .48)

					Transform 1 st , Logic		Logic 1 st , Transform
Statistic	3 + NOTAC	4	4 + NOTAC	8	2nd	Scrambled	2nd
α	.78	.85	.83	.87	.74	.75	.76
Avg CITC	.32	.42	.37	.44	.31	.30	.33
Min CITC	19	.23	.01	.23	.10	.07	.14
5 th Pct	12	.25	.05	.27	.11	.07	.15
25 th Pct	.22	.34	.33	.41	.29	.22	.22
50 th Pct	.33	.42	.38	.45	.34	.30	.35
75 th Pct	.44	.50	.47	.48	.40	.41	.39
95 th Pct	.51	.60	.60	.60	.53	.50	.56
Max CITC	.51	.60	.60	.60	.54	.51	.56

# (COMPLETED IN <30 MINUTES AND HS DEGREE/GED/< 1 YR OF COLLEGE)

Same pattern of results as for All Participants, internal consistency estimates were lower for the Transformation + Logic items ( $\alpha$ 's < .70).

					Transform 1 st , Logic		Logic 1 st , Transform
Statistic	3 + NOTAC	4	4 + NOTAC	8	2nd	Scrambled	2nd
α	.79	.86	.84	.88	.64	.69	.70
Avg CITC	.28	.41	.37	.40	. 25	.26	.28
Min CITC	28	.21	05	.19	03	.06	.01
5 th Pct	28	.21	04	.19	02	.07	.02
25 th Pct	.13	.33	.30	.31	.13	.15	.18
50 th Pct	.33	.40	.38	.39	.25	.25	.26
75 th Pct	.40	.50	.46	.46	.35	.36	.40
95 th Pct	.54	.56	.65	.61	.55	.39	.58
Max CITC	.54	.56	.65	.61	.58	.47	.58

### PHASE 1 PILOT: TEST SCORE DISTRIBUTIONS BY CONDITION (ALL PARTICIPANTS)



#### 12 Transformation + 12 Logic Items



### PHASE 1 PILOT: TEST SCORE DISTRIBUTIONS BY CONDITION (COMPLETED IN < 30 MINUTES AND HS DEGREE/ GED/< 1 YR OF COLLEGE)



# PHASE 1 PILOT: TRANSFORMATION AND LOGIC SUBSCORE DISTRIBUTIONS BY CONDITION

### **All Participants**



### PHASE 1 PILOT: TRANSFORMATION AND LOGIC SUBSCORE DISTRIBUTIONS BY CONDITION (CONTINUED)

Completed in < 30 minutes and HS Degree/ GED/< 1 yr of College



## PHASE 1 PILOT: ITEM DIFFICULTY BY TRANSFORMATION ONLY CONDITION

#### All Participants

- Items in the four-response option items were less difficult (avg p = .63) compared to:
  - 3 + NOTAC conditions (avg p = .54)
  - 8-response option conditions (avg p = .51)
  - 4 + NOTAC conditions (avg p = .50)
- Among NOTAC conditions, the six items where NOTAC was keyed (correct) response were more difficult
  - Item difficulties (avg p = .53) compared to items where NOTAC was keyed response (avg p = .23)

#### Completed in 30 Mins or Less with HS Degree/GED/< 1 yr of College

The four-response option items were less difficult (avg p = .58) compared to:

- 3 + NOTAC conditions (avg p = .52)
- 4 + NOTAC conditions (avg p = .49)
- 8-response option conditions (avg p = .48)

Statistic	3 + NOTAC	4	4 + NOTAC	8
Average p	.54	.63	.50	.51
Min <i>p</i>	.07	.29	.08	.15
5 th Pct	.10	.32	.10	.17
25 th Pct	.37	.51	.33	.36
50 th Pct	.52	.60	.46	.46
75 th Pct	.75	.71	.67	.59
95 th Pct	.94	.97	.96	.94
Max p	.94	.97	.96	.94

Statistic	3 + NOTAC	4	4 + NOTAC	8
Average p	.52	.58	.49	.48
Min <i>p</i>	.05	.20	.05	.14
5 th Pct	.08	.23	.08	.15
25 th Pct	.30	.46	.29	.31
50 th Pct	.46	.52	.38	.40
75 th Pct	.62	.59	.57	.49
95 th Pct	.95	.92	.93	.92
Max p	.95	.96	.98	.97

*Note*. Results calculated after screening participants for completion times (equal to or less than 30 minutes) and highest educational attainment (HS Degree/GED/< 1 yr of College)

## PHASE 1 PILOT: ITEM DIFFICULTY BY TRANSFORMATION + LOGIC CONDITION (ALL PARTICIPANTS)

#### **Transformation Items**

- Items were more difficult in the Scrambled condition (avg p = .47) than in the grouped conditions:
  - Transformation  $1^{st}$  condition (avg p = .57)
  - Logic  $1^{st}$  condition (avg p = .55)

Transformation Items	Transform 1 st , Log 2 nd	Logic 1 st , Transform 2 nd	Scrambled
Average <i>p</i>	.57	.55	.47
Min <i>p</i>	.25	.24	.19
5 th Pct	.31	.29	.23
25 th Pct	.44	.42	.31
50 th Pct	.69	.68	.60
75 th Pct	.54	.50	.43
90th Pct	.91	.90	.89
Max p	.96	.96	.94

#### Logic Items

- Items were more difficult than the Transformation items, regardless of presentation order:
  - Transformation  $1^{st}$  condition (avg p = .31)
  - Logic  $1^{st}$  condition (avg p = .29)
  - Scrambled condition (avg *p* = .29)

Logic Items	Transform 1 st ,	Logic 1 st , Transform 2 nd	Scrambled
Average <i>p</i>	.31	0.29	0.29
Min <i>p</i>	.21	.19	.14
5 th Pct	.22	.20	.14
25 th Pct	.23	.22	.24
50 th Pct	.35	.33	.37
75th Pct	.35	.33	.37
90th Pct	.47	.43	.47
Max p	.54	.46	.51

### PHASE 1 PILOT: ITEM DIFFICULTY BY TRANSFORMATION + LOGIC CONDITION (COMPLETED < 30 MINS WITH HS DEGREE/GED/< 1 YR OF COLLEGE)

#### **Transformation Items**

- Items were more difficult in the Scrambled condition (avg *p* = .42) than in the grouped conditions:
  - Transformation 1st condition (avg p = .54)
  - Logic  $1^{st}$  condition (avg p = .52)

Transformation Items	Transform 1 st , Log 2 nd	Logic 1 st , Transform 2 nd	Scrambled
Average <i>p</i>	.54	.52	.42
Min <i>p</i>	.23	.16	.16
5 th Pct	.27	.23	.16
25 th Pct	.39	.38	.20
50 th Pct	.52	.51	.29
75 th Pct	.64	.63	.41
95 th Pct	.89	.88	.88
Max p	.96	.97	.97

#### Logic Items

- Items were more difficult than the Transformation items, regardless of presentation order:
  - Transformation 1st condition (avg p = .29)
  - Logic  $1^{st}$  condition (avg p = .26)
  - Scrambled condition (avg *p* = .26)

Logic Items	Transform 1 st , Log 2 nd	Logic 1 st , Transform 2 nd	Scrambled
Average p	.29	.26	.26
Min <i>p</i>	.18	.15	.11
5 th Pct	.19	.16	.11
25 th Pct	.23	.18	.15
50 th Pct	.27	.25	.20
75 th Pct	.32	.32	.32
95 th Pct	.44	.40	.43
Max p	.54	.42	.47

## PHASE 1 PILOT: OBSERVED ITEM DIFFICULTY BY ESTIMATED DIFFICULTY (TRANSFORMATION ONLY)

Estimated item difficulty based on item features:

- Number of transformations
  - Vertical or horizontal
  - Diagonal (top-down or bottom-up)
- Rotation

Estimated item difficulty correlated with observed item difficulty at r = -0.64

• 95% CI = -0.84, -0.32



### PHASE 1 PILOT: GROUP SCORE DIFFERENCES BY CONDITION All Participants

#### Gender

 Female participants consistently scored higher than male participants across all conditions, on average (observed d's = .06 to .31).

#### **Race and Ethnicity**

- Asian participants scored equal to or higher than White, non-Hispanic participants, on average (observed d's = .00 to .80).
- Black participants scored lower than White, non-Hispanic participants, on average (observed d's = -.08 to -.58).
- White, Hispanic participants also scored lower than White, non-Hispanic participants in all but the Scrambled condition (observed *d's* = .04 to -.21).

#### Completed in 30 Mins or Less with HS Degree/GED/< 1 yr of College

- Female participants consistently scored higher than male participants across all Transformation Only conditions, on average same as for all participants (observed *d's* = .18 to .26).
- Results were mixed on the Transformation + Logic conditions (observed *d*'s = -.03 to .34).

		Asian,	Black,	10/1-14
Condition		non-	non-	white,
Description	Female	Hispanic	Hispanic	Hispanic
3 + NOTAC	.31 (.35)	.17 (.22)	31 (40)	15 (19)
4	.22 (.27)	.00 (.00)	58 (69)	21 (25)
4 + NOTAC	.16 (.18)	.12 (.14)	20 (.24)	19 (23)
8	.21 (.24)	.36 (.42)	39 (45)	17 (19)
Transform 1 st , Logic 2 nd	.14 (.18)	.38 (.52)	40 (54)	15 (20)
Scrambled	.22 (.29)	.24 (.32)	08 (11)	.04 (.05)
Logic 1 st , Transform 2 nd	.06 (.08)	.80 (1.12)	28 (37)	29 (38)

Note. d's corrected for measurement error in the CR items are reported in the parentheses.

<b>Condition Description</b>	Female
3 + NOTAC	.26 (.33)
4	.26 (.30)
4 + NOTAC	.18 (.21)
8	.20 (.23)
Transform 1 st , Logic 2 nd	03 (05)
Scrambled	.34 ( .50)
Logic 1 st , Transform 2 nd	01 (01)

*Note.* Results calculated after screening participants for completion times (equal to or less than 30 minutes) and highest educational attainment (HS Degree/GED/< 1 yr of College). *d*'s corrected for measurement error in the CR items are reported in the parentheses.

# PHASE 1 PILOT: COMPLETION TIME BY CONDITION (ALL PARTICIPANTS, COMPLETED IN $\leq$ 30 MINUTES)

- Among the Transformation Only conditions, four-response options (with or without NOTAC) resulted in shorter average completion times than five- or eight-response options
- Transformation + Logic conditions required more time to complete than the Transformation Only conditions, on average, magnitude of the differences varied

Condition Description	n	М	SD	d _{c-8}
3 options + NOTAC	468	10.88	4.94	34
4 options	474	11.36	5.14	25
4 options + NOTAC	453	12.39	5.79	06
8 options	455	12.74	5.86	
Transform 1 st , Logic 2 nd	447	14.10	6.49	.22
Scrambled	438	13.24	6.36	.08
Logic 1 st , Transform 2 nd	440	12.98	5.87	.04

## PHASE 1 PILOT: COMPLETION TIME (IN MINUTES) BY CONDITION (COMPLETED < 30 MINUTES)

24 Items – Transformation Only



#### 12 Transformation + 12 Logic Items



# PHASE 1 PILOT: COMPLETION TIME (IN MINUTES) BY CONDITION (HS DEGREE/GED/< 1 YR OF COLLEGE)

30 -30 -3+NOTAC 20-20-10-10-0 -30 -0-20-30-10-Count Count - 0 30 -10-4+NOTAC 20-0-10-30 -0-30 -20-20ω 10-10-0-30 10 20 Ó

24 Items – Transformation Only

**Completion Time** 

#### 12 Transformation + 12 Logic Items



# PHASE 1 PILOT: PERCEIVED DIFFICULTY BY CONDITION

#### All Participants

- All participants reported their perceived performance on the CR items using a 5-point scale (1 = "Less than 20% of items correct" to 5 = "More than 80% of items correct").
- Significant effect of experimental conditions on perceived CR performance, *F* (6, 3484) = 36.65, *p* < .001 (all participants, top table).</li>
- Participants in Transformation Only conditions (3+NOTAC, 4, 4+NOTAC, and 8) reported more correct responses compared to participants in all three Transformation + Logic conditions (Transform 1st, Logic 1st, and Scrambled), on average.

# Completed in 30 Mins or Less with HS Degree/GED/< 1 yr of College

 Observed similar effects among the subset of participants who completed in 30 minutes or less with HS Degree/GED/< 1 yr of College (bottom table).

Condition Description	n	М	SD	<b>d</b> _{C-8}
3 options + NOTAC	501	3.89	.92	03
4 options	503	3.98	.89	.07
4 options + NOTAC	490	3.90	.95	02
8 options	493	3.92	.95	N/A
Transform 1st , Logic 2nd	507	3.53	.96	41
Scrambled	492	3.37	1.00	56
Logic 1st , Transform 2nd	505	3.46	.95	48

<b>Condition Description</b>	n	М	SD	<b>d</b> _{C-8}
3 options + NOTAC	182	3.79	.93	04
4 options	188	3.86	.91	.03
4 options + NOTAC	172	3.83	.96	.00
8 options	179	3.83	.97	N/A
Transform 1st , Logic 2nd	164	3.59	.94	25
Scrambled	164	3.26	1.00	58
Logic 1st , Transform 2nd	151	3.21	1.00	63

Note. Results calculated after screening participants based on completion times (equal to or less than 30 minutes) and highest educational attainment (HS Degree/GED/< 1 yr of college)

# PHASE 1 PILOT: PERCEIVED DIFFICULTY BY CONDITION AND GENDER, RACE, AND ETHNICITY

#### **All Participants**

#### Gender

- Significant difference in perceived performance by gender, *F* (2, 3488) = 21.78, *p* < .001.</li>
- Male participants consistently reported scoring higher compared to female participants, avg d = .22.

#### **Race and Ethnicity**

 Differences in perceived performance among racial or ethnic groups were mixed and varied by condition.

#### Completed in 30 Mins or Less with HS Degree/GED/< 1 yr of College

#### Gender

 Males reported consistently higher performance than females (avg d = .25).

Condition Description	Female	Asian, non- Hispanic	Black, non- Hispanic	White, Hispanic
3 + NOTAC	3.85 (.94)	3.97 (.77)	3.69 (.82)	3.92 (.93)
4	3.91 (.93)	4.03 (.95)	3.89 (.90)	3.77 (.88)
4 + NOTAC	3.84 (.94)	4.00 (.80)	3.87 (1.06)	3.86 (1.07)
8	3.88 (.91)	4.15 (.53)	3.72 (1.03)	3.88 (.96)
Transform 1 st , Logic 2 nd	3.39 (.96)	3.55 (.97)	3.78 (.90)	3.61 (.85)
Scrambled	3.24 (.96)	3.24 (1.00)	3.12 (1.12)	3.47 (.92)
Logic 1 st , Transform 2 nd	3.34 (.96)	3.75 (.76)	3.64 (.92)	3.51 (.87)

Condition Description	Female	
3 + NOTAC	3.72 (.96)	
4	3.82 (.93)	
4 + NOTAC	3.75 (.93)	
8	3.78 (.99)	
Transform 1 st , Logic 2 nd	3.40 (.95)	
Scrambled	3.17 (.98)	
Logic 1 st , Transform 2 nd	3.07 (.99)	

*Note*: Results calculated after screening participants based on completion times (equal to or less than 30 minutes) and highest educational attainment (HS Degree/GED/< 1 yr of college)