

How Do Professors Format Exams?

An Analysis of Question Variety at Scale

(c) The polynomial $x - 1 \pmod{15}$ has at most 3 zeros.

True ☒ False

(ii) [2 pts] Continuous

☒ Self-Driving Car ☐ Pacman

(a) (4 points) ☒ $A_0 = 0$ ☐ $A_1 = 1$ ☐ $A_k = k$
☐ $A_0 \neq 0$ ☒ $A_1 \neq 1$ ☒ $A_k \neq k$

2) What protocol is used to find the hardware address?

- a) RARP
b) ICMP
c) None of the above

d. What are the two species present at high concentration?

NaCl , NaNO_2

3. Complete the following:

```
typedef struct blo {  
    item_t *head;  
    pthread_mutex_t blo_lock;  
} blo_t;
```

a) Binary

b) Multiple Choice

c) Short Writing

Paul Laskowski, Sergey Karayev, and Marti Hearst

Learning at Scale 2018

WHO



Paul Laskowski

UC Berkeley

Professor at School of Information



Sergey Karayev

Gradescope, Inc.

Co-founder, Director of Research



Marti Hearst

UC Berkeley

Professor at School of Information & CS

WHAT

- We analyzed questions on nearly 1,800 paper exams graded in real STEM courses
- We annotated the type of each question (multiple choice, short/medium/long writing, drawing)
- We found interesting differences by type in position, use across subjects, student performance, and reliability
 - e.g. 5+ binary choice questions are required to match reliability of 1 long writing question

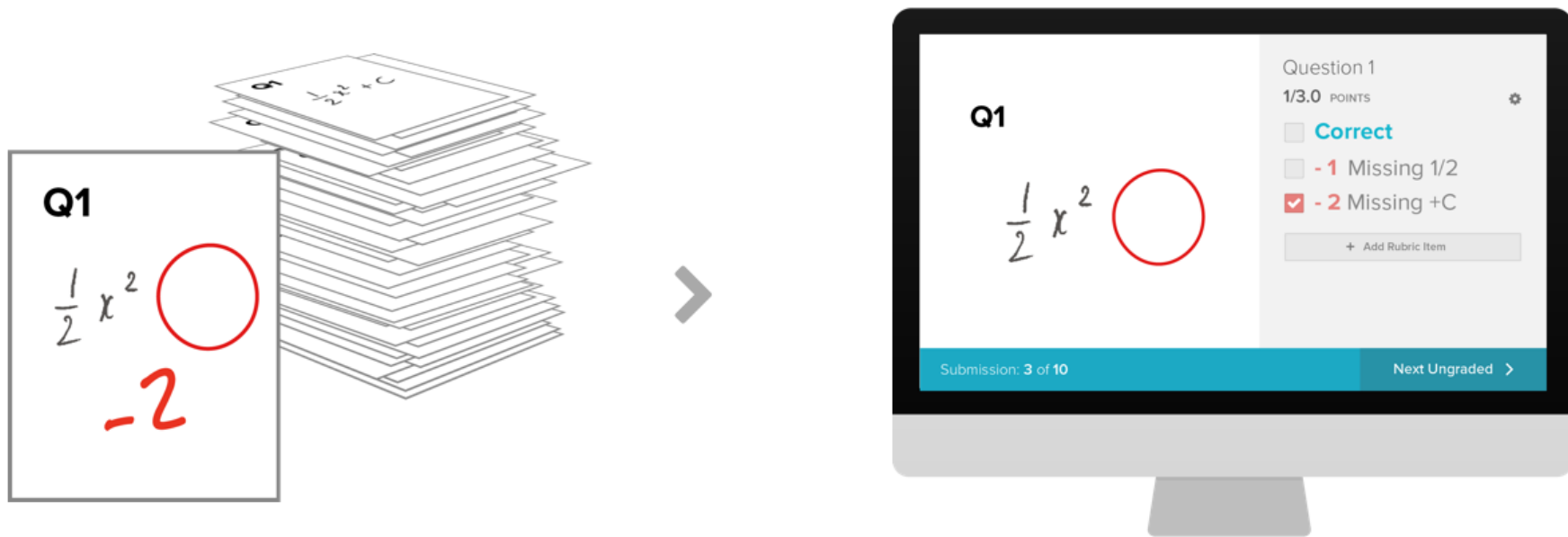
WHY

- Unique vantage point on exams “in the wild”
- Near term: get an **understanding** of what is happening today
- Longer term: **improve** exam writing for both for instructors and students

RELATED WORK

- The trend: oral open-response -> written open-response -> multiple-choice
- Good assessment is valid and reliable.
 - Ideally, summative assessment is also discriminative and easy to grade without bias.
- Many studies on open-response vs. multiple-choice
 - Well-written MCs *can* be as good as essays
 - MCs are hard to write, easy to grade. Essays are the opposite.
- Ordering effects are unclear

DATA SOURCE




Gradescope: a Fast, Flexible, and Fair System for Scalable Assessment of Handwritten Work

Arjun Singh, [Sergey Karayev](#), Kevin Gutowski, Pieter Abbeel

Learning at Scale 2017

Q1

$$\frac{1}{2}x^2$$


Question 1

1/3.0 POINTS

☐

Correct

☐

-1 Missing 1/2

☒

-2 Missing +C

+ Add Rubric Item

Submission: 3 of 10

Next Ungraded >

- No constraints on type of questions
- Consistent scoring due to shared, modifiable rubric

DATA ANNOTATION

- 40M+ student answers
 - graded by 6,000+ instructors (mostly STEM)
 - at 600+ schools (mostly US higher-ed)
- Built special annotation interface, hired educators
- Annotated over 50% of 50,000+ questions that had at 100+ student answers at the time

REGION 1

+ ADD REGION

☐ Unsuitable

Save

RESPONSE AREA

- ☐ Other
- ☐ Bubble or Square MC
- ☐ Other kinds of MC
- ☐ Small fill-in-the-blank
- ☐ Line-sized
- ☐ Paragraph-sized
- ☐ Page-sized
- ☐ Matrix or grid
- ☐ Diagram, plot, or chem structure
- ☒ Box for drawing

TYPICAL RESPONSE LENGTH

- ☒ N/A
- ☐ Character, word, or words
- ☐ Sentence or a couple of lines
- ☐ Paragraph or more

RESPONSE TYPE

Binary MC

Multiple MC

Text

Code

Math

Science Symbols

Diagram

Diagram Annotation

Chemical Structure

Self-boxed

Justification/Extra Work

Final Answer

Not English

Other




Remove

☒ Set Bounding Box

Rubric 6.0 points, floor, ceiling

Sample Answers 2 answers

Mask pages

- 0  Correct
- 3  Incorrect phasing
- 6  Fully incorrect

3. Draw a 7p orbital. Show all nodes and phasing for full credit (6 pts)

1

3. Draw a 7p orbital. Show all nodes and phasing for full credit (6 pts)

1



Review Groups

Grade

Next Annotated

Next Random

Save & Next Random >

REGION 1

+ ADD REGION

☐ Unsuitable

Save

RESPONSE AREA

- ☐ Other
- ☐ Bubble or Square MC
- ☐ Other kinds of MC
- ☐ Small fill-in-the-blank
- ☐ Line-sized
- ☐ Paragraph-sized
- ☐ Page-sized
- ☐ Matrix or grid
- ☐ Diagram, plot, or chem structure
- ☒ Box for drawing

TYPICAL RESPONSE LENGTH

- ☒ N/A
- ☐ Character, word, or words
- ☐ Sentence or a couple of lines
- ☐ Paragraph or more

RESPONSE TYPE

Binary MC

Multiple MC

Text

Code

Math

Chemical Structure

Self-boxed

Justification/Ext

Remove

☒ Set Bounding Box

FINAL DATASET

- Almost 23,000 questions
- From almost 1,800 exams
- Corresponding to graded answers of 120,000+ students

QUESTION TYPE SAMPLES

(c) The polynomial $x - 1 \pmod{15}$ has at most 3 zeros.

True ☐ False ☒

(ii) [2 pts] Continuous

☒ Self-Driving Car

☐ Pacman

a) Binary

(e) Pre-image resistance (first pre-image resistance)

Given a hash value h , it should be computationally infeasible to find x such that $h = h(x)$

(c) (2 pts) What is the distance to default?

Distance to default is defined as $\ln\left(\frac{V}{B}\right)$
 Since the debt level is 90% of the value from value $\frac{1}{0.9} = \frac{V}{B} = 1.1111 \dots$
 $\ln\left(\frac{1}{0.9}\right) = 0.105$

d) Medium Writing

(a) (4 points) ☒ $A_0 = 0$ ☐ $A_1 = 1$ ☐ $A_k = k$
☐ $A_0 \neq 0$ ☒ $A_1 \neq 1$ ☒ $A_k \neq k$

2) What protocol is used to find the hardware address?

a) RARP

☒ b) ICMP

☒ c) None of the above

b) Multiple Choice

d. What are the two species present at high concentration?

$\text{NaCl}, \text{NaNO}_2$

3. Complete the following:

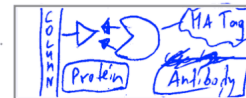
```
typedef struct blo {
    item_t *head;
    pthread_mutex_t blo_lock;
} blo_t;
```

c) Short Writing

3. Draw a 7p orbital. Show all nodes and phasing for full credit (6 pts)



A). (4 pts) Use the HA-tag to purify your favorite protein.
 Draw a diagram showing the purification setup.



e) Long Writing

f) Drawing

(c) The polynomial $x - 1 \pmod{15}$ has at most 3 zeros.

True

False

(ii) [2 pts] Continuous

☒ Self-Driving Car

☐ Pacman

a) Binary

(e) Pre-image resistance (first pre-image resistance)

(a) (4 pts)

2) What

zeros.

(a) (4 points) $\textcircled{\bullet} A_0 = 0$ $\textcircled{} A_1 = 1$ $\textcircled{} A_k = k$
 $\textcircled{} A_0 \neq 0$ $\textcircled{\bullet} A_1 \neq 1$ $\textcircled{\bullet} A_k \neq k$

2) What protocol is used to find the hardware address?

- a) RARP
- ☒ b) ICMP
- ☒ c) None of the above

b) Multiple Choice

$\equiv k$
 $\neq k$

d. What are the two species present at high concentration?

NaCl , NaNO_2

re address?

3. Complete the following:

```
typedef struct blo {  
    item_t *head;  
    pthread_mutex_t blo_lock;  
} blo_t;
```

c) Short Writing

(e) Pre-image resistance (first pre-image resistance)

Given a hash value h , it should be computationally infeasible to find x such that $h = h(x)$

(c) (2 pts) What is the distance to default?

Distance to default is defined as $\ln\left(\frac{V}{B}\right)$

Since the debt level is 90% of the value
from value $\frac{1}{0.9} = \frac{V}{B} = 1.1111 \dots$

$$\ln\left(\frac{1}{0.9}\right) = 0.105$$

d) Medium Writing

1. (10 points) An unknown compound contains carbon, hydrogen, and oxygen. A 15g sample was combusted to produce 36.62g of carbon dioxide and 14.99 of water. What is the formula of the unknown compound?

$$36.62 \text{ g } \text{CO}_2 \times \frac{12.01 \text{ g}}{44.01 \text{ g}} = 9.997 \text{ g C} \quad \begin{matrix} 14.99 \\ 11.67127 \\ \hline 3.321790 \end{matrix}$$

$$14.99 \text{ g } \text{H}_2\text{O} \times \frac{2.02 \text{ g H}_2}{18.019} = 1.6812 \text{ g H}$$

$$9.997 \text{ g C} \times \frac{1 \text{ mol C}}{12.01 \text{ g C}} = 0.8323 \text{ mol C}$$

$$1.6812 \text{ g H} \times \frac{1 \text{ mol}}{1.01 \text{ g H}} = 1.6645 \text{ mol H}$$

$$3.3217 \text{ g O} \times \frac{1 \text{ mol O}}{15.999 \text{ g}} = 0.20762 \text{ mol O}$$

$$\frac{0.8323}{0.20762} = 4 \text{ C} \quad \frac{1.6645}{0.20762} = 8 \text{ H} \quad \frac{0.20762}{0.20762} = 1$$



4. (20 pts) Validation is a key concept in Minsky's financial instability hypothesis. Write an essay to explain this concept. Pay attention to how validation influences lending standards. Include at least one example of validation from the subprime lending crisis. What are the implications of validation leading up to the Great Recession for future business cycles?

The hypothesis centers around the concept that bubbles exist in the marketplace, because of the inherent fluctuation in the expectations and actions of people. The changes in the validation for lending standards, longer or shorter prospects change.

When the economy is doing well, lending standards are loose and banks and lenders validate easily. The economy continues to do well & grow, until credit and lending increases due to the expansion of credit. (because validation is easy). What started as hedge finance transforms into speculative finance, where too much money has been lent out to cover all the spending and investment.

Lastly, this transitions into a ponzi finance environment where debts are paid through increasing debt. This is like the subprime lending crisis, where validation and lending went over the edge because of the expanded credit.

e) Long Writing

esis. Write an essay to explain this
 de at least one example of validation from
 up to the Great Recession for future

upt the cycles
 the inherent fluctuation
 angle. The changes
 change or interest

debt on large and
 my culture to do
 only of the
 is easy). What
 into speculative
 been left out
 investment.

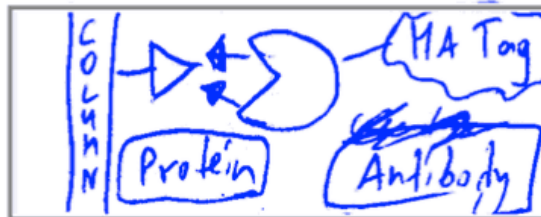
Finance environment
 debt. This is
 validation and
 of the expected

3. Draw a 7p orbital. Show all nodes and phasing for full credit (6 pts)



A). (4 pts) Use the HA-tag
 to purify your favorite protein.

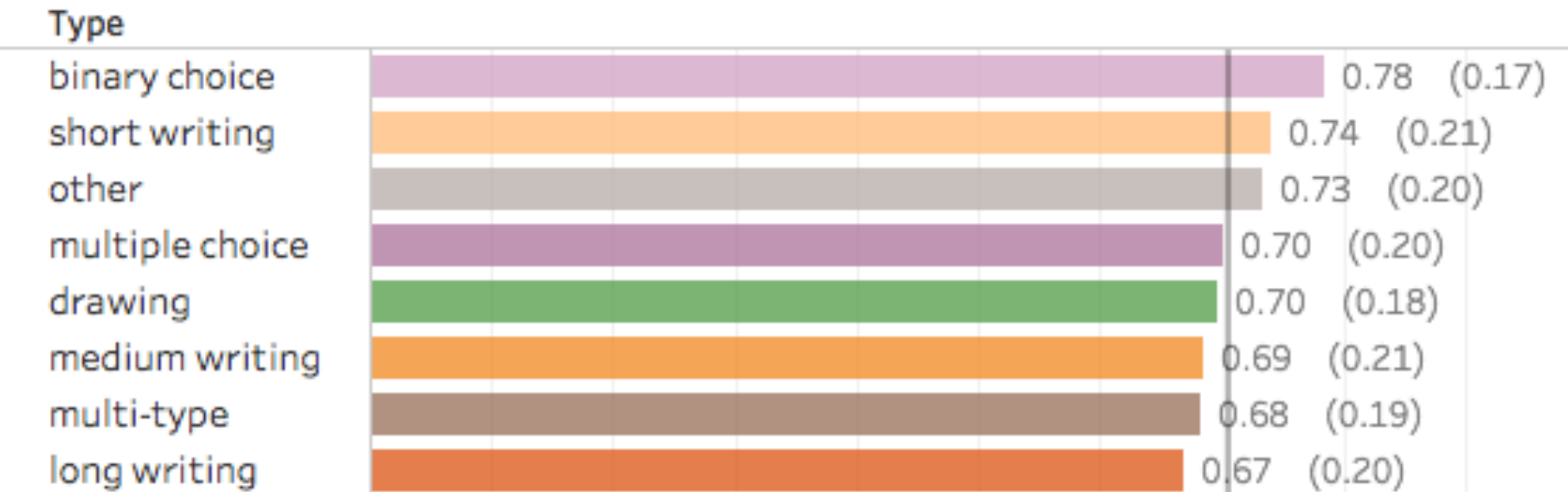
Draw a diagram showing
 the purification setup.



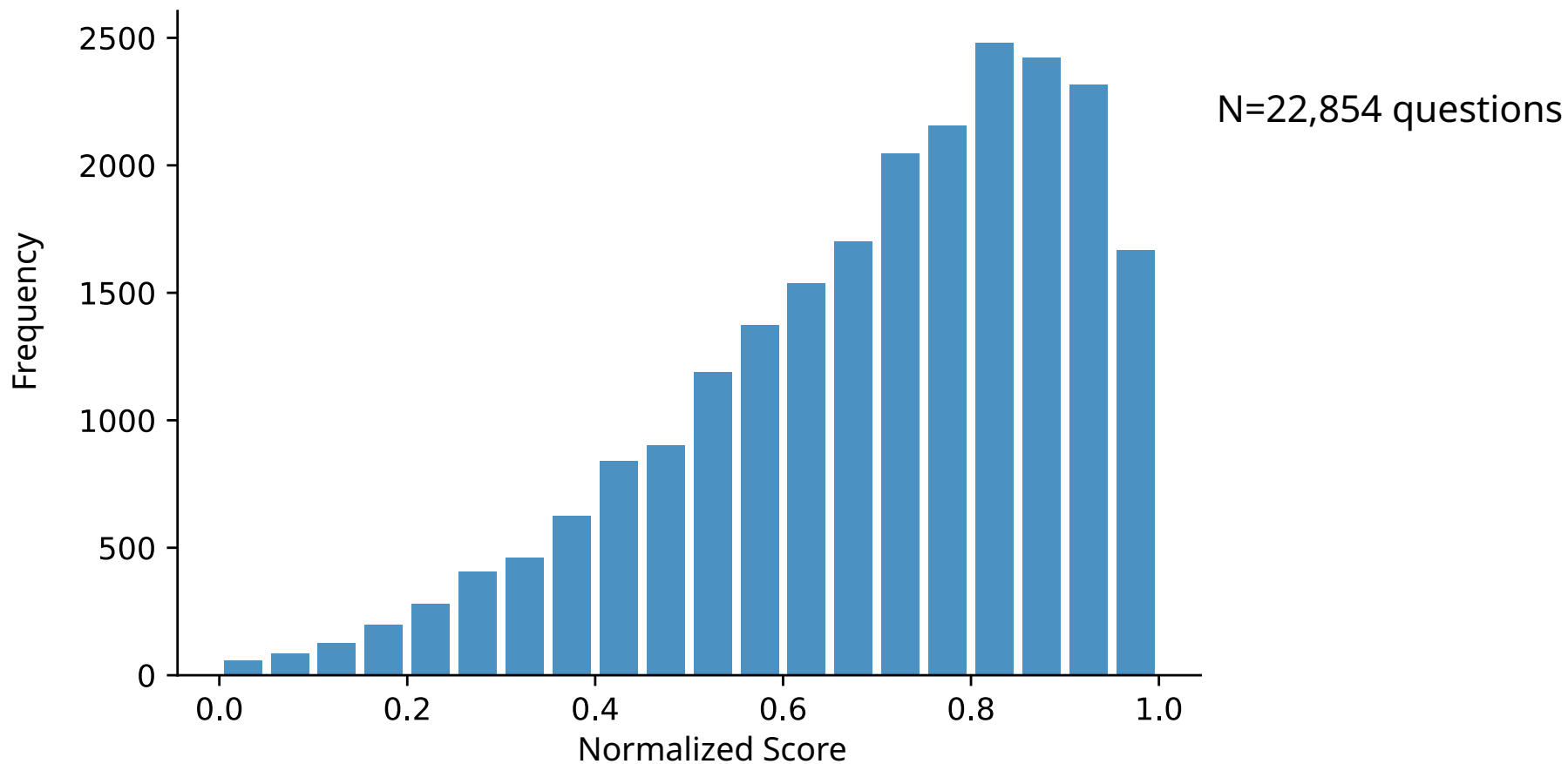
f) Drawing

HOW DO SCORES DIFFER BY QUESTION TYPE?

MEAN SCORES BY QUESTION TYPE

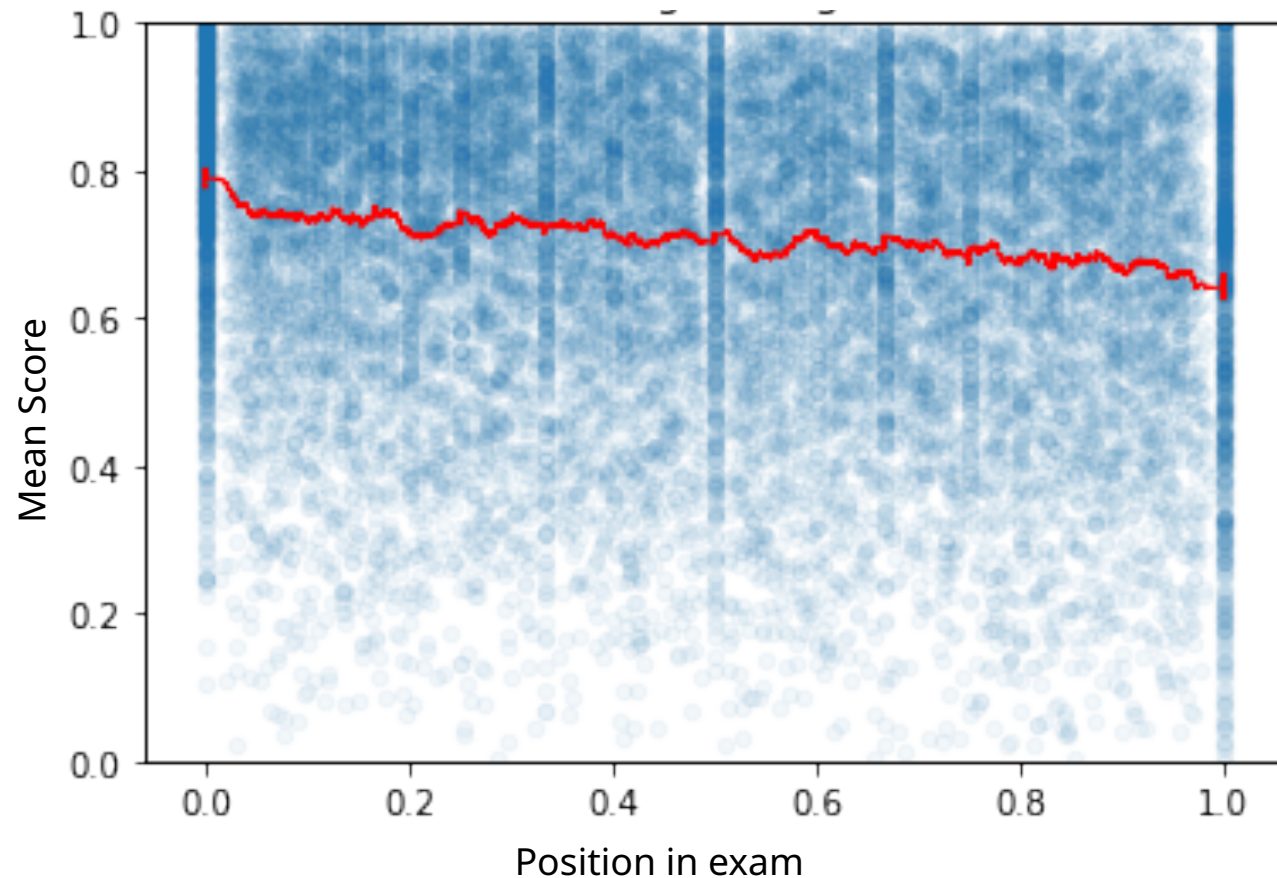


OVERALL DISTRIBUTION OF QUESTION MEAN SCORES



HOW DO INSTRUCTORS SEQUENCE
QUESTIONS?

QUESTION SCORE BY POSITION IN EXAM

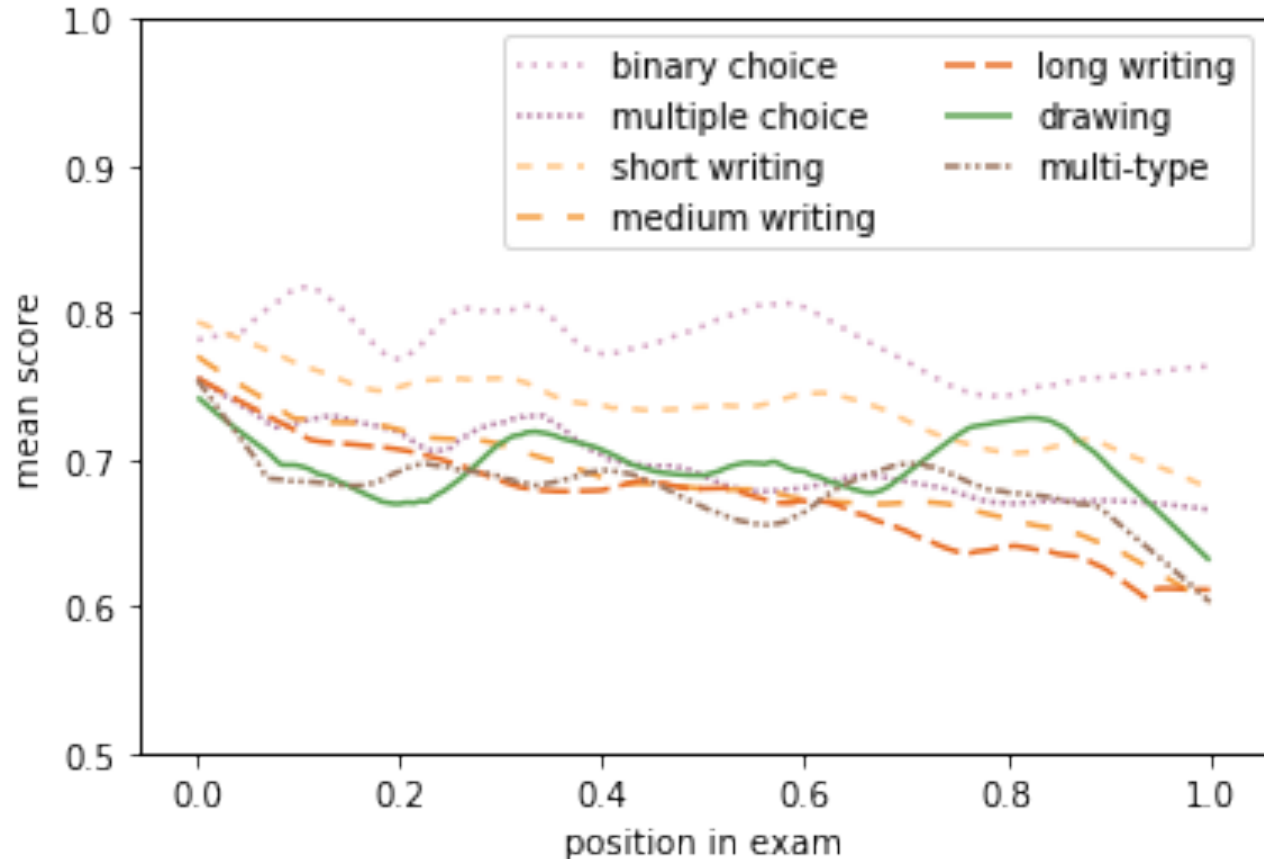


- Mean scores drop with position in exam.

FIRST HALF VS SECOND HALF

- More frequently occur in **first** half of exam:
 - Binary choice (70%)
 - Multiple choice (56%)
 - Short writing (53%)
- More frequently occur in **second** half of exam:
 - Long writing (60%)
 - Medium writing (53%)

QUESTION SCORES BY POSITION IN EXAM



- Mean scores drops with position in exam for all question types.

HOW RELIABLE ARE DIFFERENT QUESTION TYPES?

QUESTION RELIABILITY

$$\text{reliability} = \frac{\text{signal}}{\text{signal} + \text{noise}}$$

QUESTION RELIABILITY

*Ideal: give each student an
infinite number of questions*

$$\text{reliability} = \frac{\text{signal}}{\text{signal} + \text{noise}} = \frac{\text{variance in student ability}}{\text{total score variance}}$$

Observed data

QUESTION RELIABILITY

~~Ideal: give each student an infinite number of questions~~

Estimate with a linear mixed model
(Hoyt 1941, Cronbach 1951)

$$\text{reliability} = \frac{\text{signal}}{\text{signal} + \text{noise}} = \frac{\text{variance in student ability}}{\text{total score variance}}$$

Observed data

OUR MODEL

$$score_{s,q,a,t} = G_t + K_{s,t} + A_{s,a} + u_{s,q,a,t}$$



Global question
type average



Student-assignment
affinity



Observed score of student s
on question q in assignment a
of type t

Student aptitude
for the type

Student-question
affinity

OUR MODEL

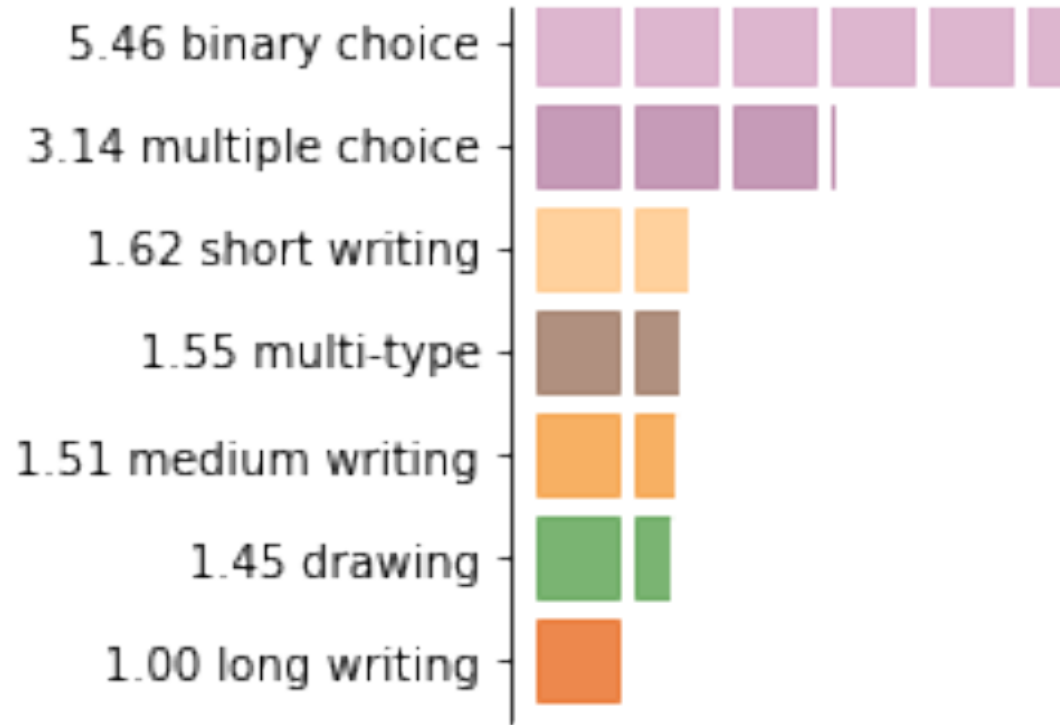
$$score_{s,q,a,t} = G_t + K_{s,t} + A_{s,a} + u_{s,q,a,t}$$

- Key assumption is that each effect is independent.
- Many factors not accounted for:
 - Students that work hard inspire others to work hard.
 - Instructors make up for hard questions/assignments with easy ones.
 - Consecutive questions cover similar material

RELIABILITY ESTIMATES

	Reliability	Error
Binary	0.036	0.002
Multiple Choice	0.061	0.001
Short Writing	0.112	0.001
Medium Writing	0.120	0.001
Long Writing	0.170	0.001
Drawing	0.144	0.002
Multi-type	0.117	0.001

NUMBER OF QUESTIONS GIVING EQUIVALENT RELIABILITY



EBEL, R. L. Can teachers write good true-false test items?
Journal of Educational Measurement 12, 1 (1975), 31–35.

TO WHAT DEGREE ARE DIFFERENT STUDENTS
SUITED FOR DIFFERENT QUESTION TYPES?

$$score_{s,q,a,t} = G_t + K_{s,t} + A_{s,a} + u_{s,q,a,t}$$



Observed score of student s
on question q in assignment a
of type t



Global question
type average



Student aptitude
for the type



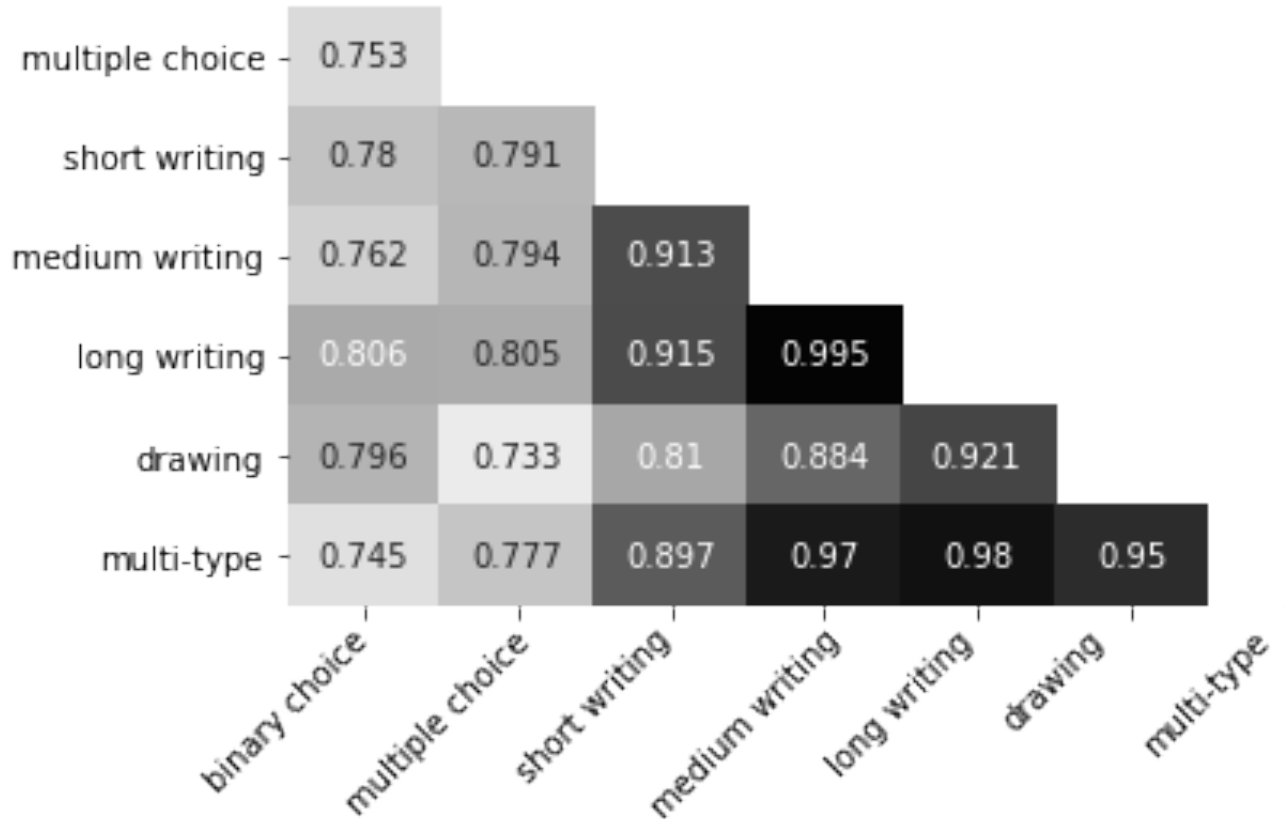
Student-assignment
affinity

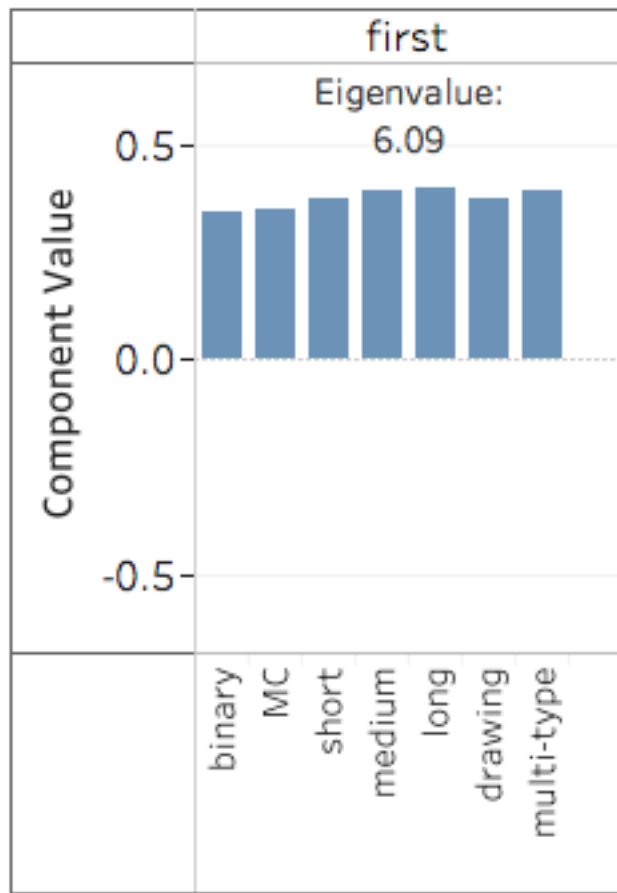


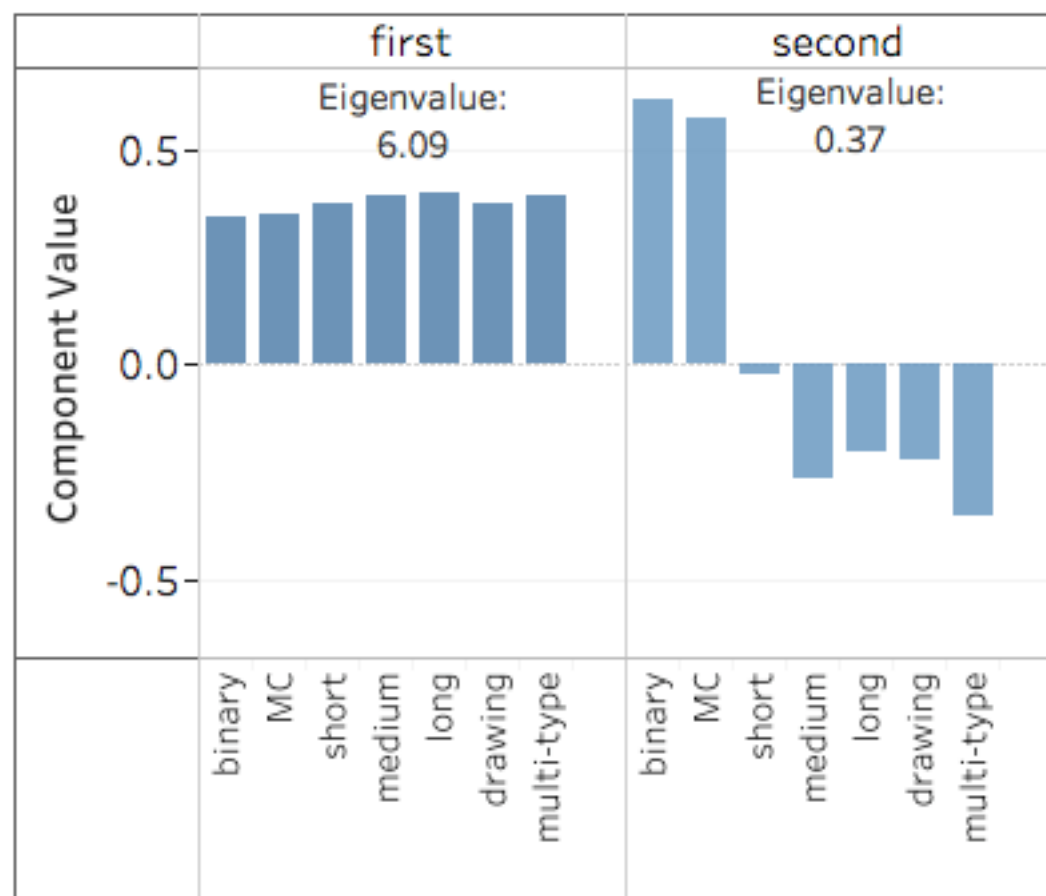
Student-question
affinity

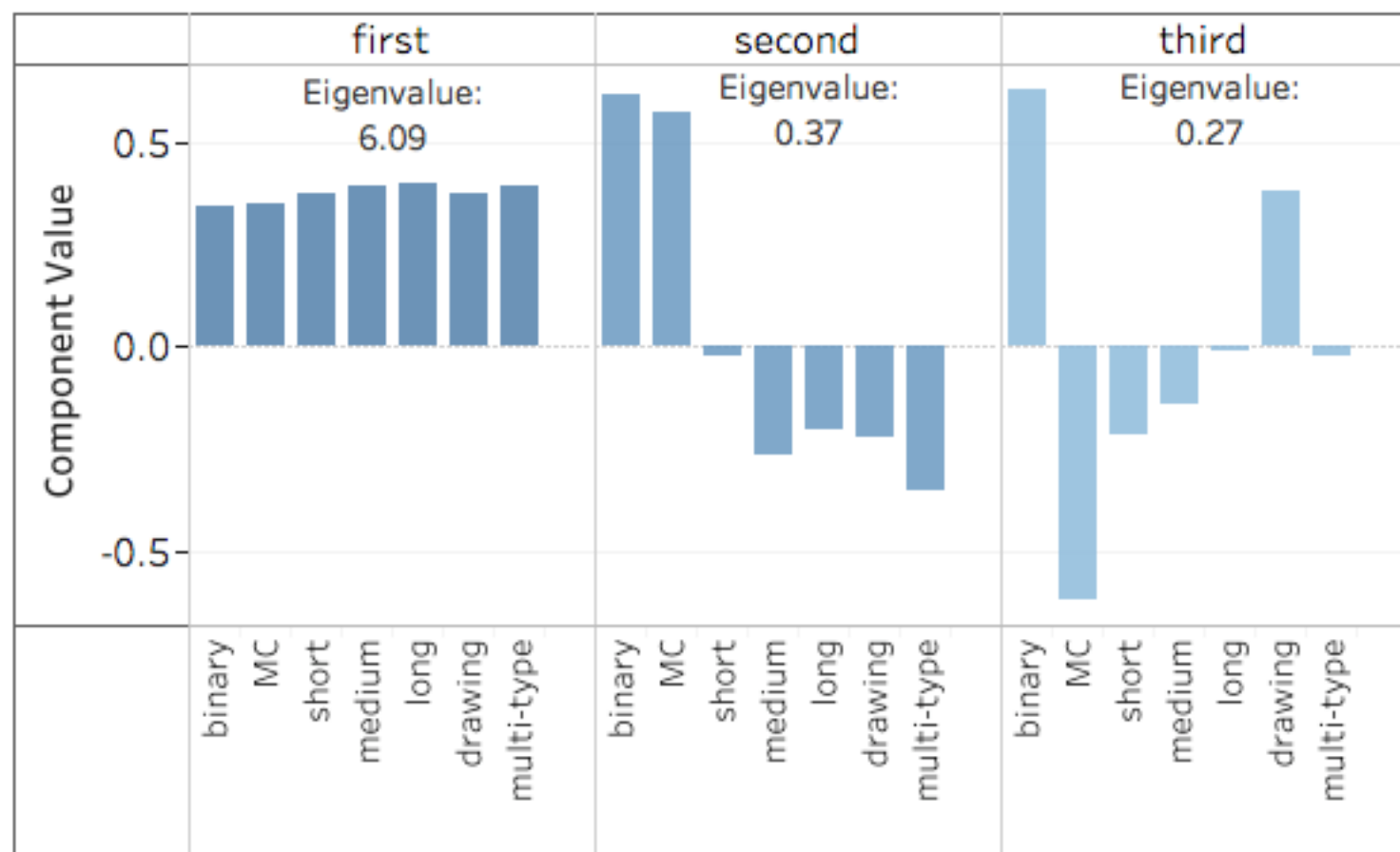
	Binary	MC	Short	$t = \text{Medium}$	Long	Drawing	Multi
Student A							
$s = \text{Student B}$				$K_{s,t}$			
Student C							
Student D							
Student E							
...							

QUESTION TYPE APTITUDE CORRELATIONS









SUMMARY

- Revealed some common patterns in university STEM exams today
 - Start has binary and multi-MC, ending has open response
 - Mean score drops with position in exam
- Found significant differences in question type reliability
- Found that student aptitude for question types is largely correlated
 - But: binary/multiple choice questions are less correlated with others
- More interesting findings in the paper!

THANK YOU!

LET'S COLLABORATE:
SERGEY@GRADESCOPE.COM