# Introduction

# Why AI? Save time and effort.



SMEs & Test Developers



Instructors/Teachers



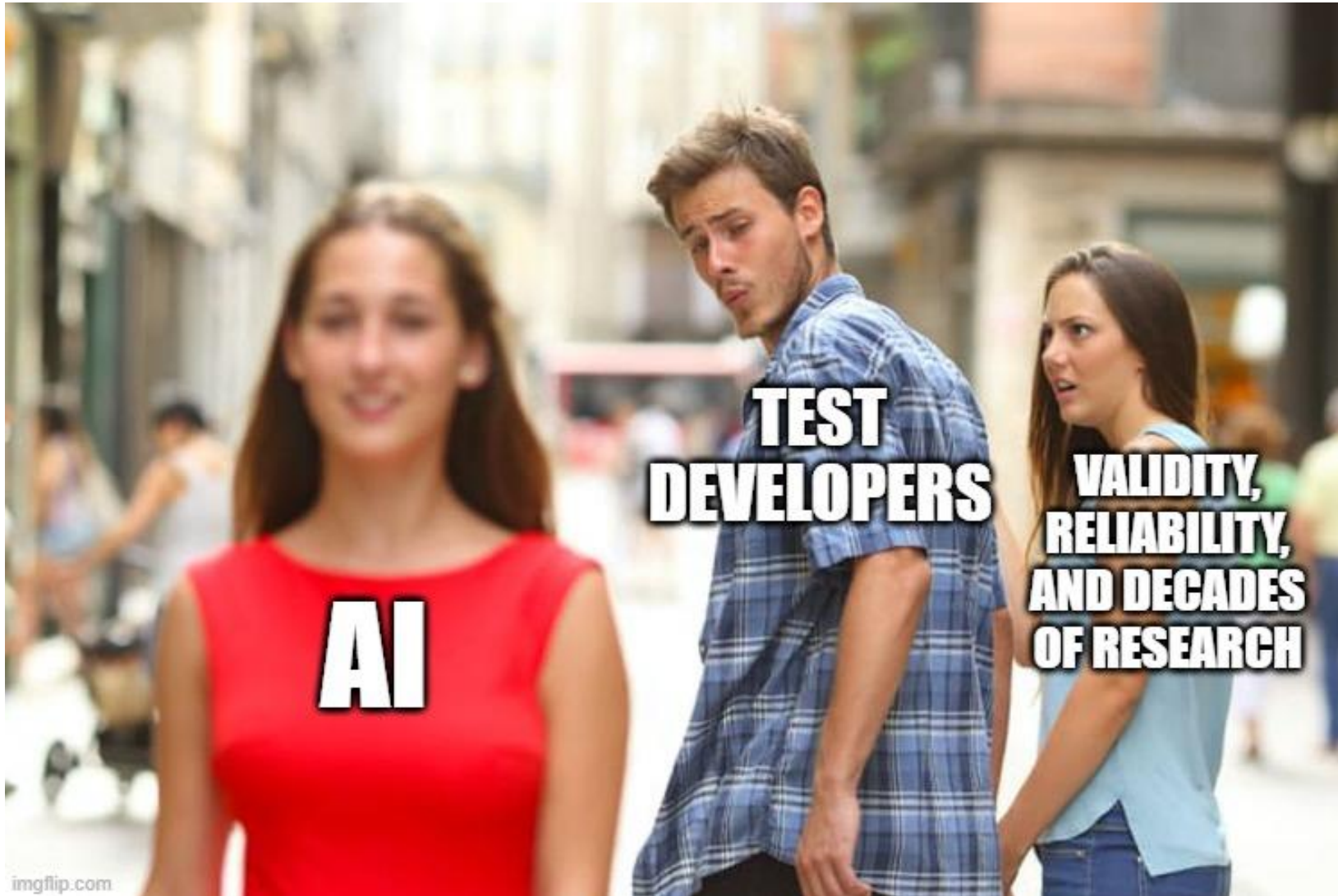Students/Candidates



Psychometricians

ASC
Assessment Systems

# Why AI?

- Pain points in assessment
  - Item writing
  - Review items
  - Assemble linear forms
  - Make the tests smarter and faster for examinees
  - Reduce essay marking time/costs
  - Maintain security in delivery
  - Review results
  - Provide feedback to examinees

# Why AI?

But…

# What is validity?

Evidence to support our intended interpretations of test scores

# Validity needs are a continuum

Quiz          Final exam          Admissions          Cert/Licensure

# Why AI?

- So, will it help us solve our pain points?
- Or, will it lead to new problems?
- What constraints are needed?
- The answer, in many cases: it has been providing solutions, for decades.
- LLMs are new… how do they compare to traditional AI in our use cases, and regarding validity?

ASC
Assessment Systems

# Part 2
# AI and ML:
# An Introduction

# New terms, not new ideas!

# Definitions

- Machine Learning
  - The use and development of computer systems that are able to learn and adapt without following explicit instructions, by using algorithms and **statistical models to analyze and draw inferences from patterns in data**. (*Oxford*)

- Artificial Intelligence
  - The theory and development of computer systems able to perform tasks that normally require human intelligence, such as visual perception, speech recognition, **decision-making**, and translation between languages. (*Oxford*)

ASC
Assessment Systems

# Definitions

- Machine Learning: **Unsupervised**
  - Looking for patterns with no particular goal
- Machine Learning: **Supervised**
  - Some sort of goal in mind (label/criterion/prediction)

# Machine learning & AI

- Common example is logistic model to predict binary outcome...
- "Natural language processing"

# Part 3
# Traditional AI and ML in Assessment and Psychometrics

# ML: Factor Analysis

- Evaluate dimensionality, find patterns (unsupervised ML)
- 3 levels: Items, subscores, tests/battery
- Invented in 1904!!!!!!!!!

Schmukle, Back, & Eglogg (2008)

# ML: Item Response Theory (IRT)

- Fits logistical models to find patterns and make sense of the data (unsupervised ML)

# ML: Cognitive Diagnostic Models (CDMs)

- Make a "Q matrix" for items and curriculum/blueprint
- Can be supervised or unsupervised

| Item | Excel Formulas | Quant reasoning | Tax knowledge | Reports |
|------|----------------|-----------------|---------------|---------|
| 1 | 1 | 0 | 0 | 1 |
| 2 | 0 | 0 | 1 | 1 |
| 3 | 0 | 0 | 1 | 0 |
| 4 | 1 | 0 | 0 | 0 |
| 5 | 0 | 1 | 1 | 0 |

ASC
Assessment Systems

# ML: Process Data

- Let's use ALL the data we have from assessment and learning
- We need complex ML to make sense of it



Provasnik (2021)

# AI: Automated Test Assembly

- Automate how a psychometrician would manually do complex assembly

- I used to do this!

# AI: Automated Item Generation

- Templates
- Can still be 6x as effective (Medical Council of Canada)

# AI: Enemy Item detection

- NLP + Text similarity/distance indices

| Item 1 | Content 1 | Item 2 | Content 2 | Similarity |
|---|---|---|---|---|
| I.C.077 | Glenohumeral ligaments are thickened folds | II.A.A.002 | The glenohumeral ligaments are thickened folds | 1 |
| I.C.075 | The three muscles which are used for | I.C.092 | The three muscles which are used for extension | 1 |
| DACBSP 0( | Appropriate rehab exercises for FAI can | N.Unspec | Appropriate rehab exercises for FAI can include: | 1 |
| I.C.070 | The tissue that is most susceptible to overuse | II.A.A.009 | The tissue injury that is most susceptible to | 0.979366 |
| II.E.068 | Infraspinatus trigger points can be stretched | II.E.069 | Supraspinatus trigger points can be stretched | 0.970059 |
| II.E.069 | Supraspinatus trigger points can be stretched | II.E.070 | Subscapularis trigger points can be stretched by: | 0.966984 |
| II.E.068 | Infraspinatus trigger points can be stretched | II.E.070 | Subscapularis trigger points can be stretched by: | 0.965954 |
| H.II.D.001 | It is essential that the team physician for | N.II.D.001 | It is essential that the team physician for female | 0.965832 |
| I.C.001 | Thrower's elbow symptoms are MOST often | I.C.027 | Tennis elbow symptoms are MOST often located | 0.916058 |
| H.II.E.027 | Self stretch of the pectineus muscle involve: | II.E.072 | Self-stretch of a piriformis muscle can be | 0.865478 |
| II.E.074 | The stretch position of the adductor longus | II.E.082 | The stretch position for the sartorius muscle is: | 0.865155 |
| II.A.A.125 | The carpal bone that is most common carpal | II.B.103 | The MOST common carpal bone fracture is the: | 0.816633 |

ASC
Assessment Systems

# AI: Computerized Adaptive Testing (CAT)

- Uses IRT to personalize the test for all examinees
- Reduces time by 50%, more secure, can measure off-grade, and more
- Automates Alfred Binet's 1908 IQ test with ML

# AI: Automated essay scoring (AES)

- NLP does not mean we are naturally processing
- We train a machine learning model for each prompt/rubric
- AI can now also do *evaluation and feedback*

| Score | school_board | three_reasons | leadership | ludacris |
|-------|--------------|---------------|------------|----------|
| 5 | 3 | 1 | 2 | 0 |
| 4 | 2 | 1 | 0 | 0 |
| 3 | 1 | 1 | 1 | 0 |
| 2 | 2 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 2 |

ASC
Assessment Systems

# AI: Remote Proctoring

- Break it down into solvable problems
- Looks for checklist of issues by evaluating still images based on ML models

# ML/AI outside the Assessment

- Prediction models:
  - Job performance, counterproductive work behavior, tenure
  - Same for Universities too!
- Resume/CV evaluation
- Adaptive learning
- CrystalKnows, Cambridge Analytica, etc
- So much more!

# Part 3
# LLMs

# Intro to LLMs

- LLMs are like essay scoring in that they break down texts into a regression-type data set
- Predict the next word or words based on prompts and previous words
- But waaaaay more complex

# AIG

- We've all heard about ChatGPT, but what does it mean for assessment?

I had such a

great    **great time**    lovely

Training Optimus Prime, M.D.: Generating Medical Certification Items by Fine-Tuning OpenAI's gpt2 Transformer Model

Matthias von Davier*

August 21st, 2019

**Abstract**

**Objective:** Showcasing Artificial Intelligence, in particular deep neural networks, for language modeling aimed at automated generation of medical education test items.

**Materials and Methods:** OpenAI's gpt2 transformer language model was retrained using PubMed's open access text mining database. The retraining was done using toolkits based on tensorflow-gpu available on GitHub, using a workstation equipped with two GPUs.

**Results:** In comparison to a study that used character based recurrent neural networks trained on open access items, the retrained transformer architecture allows generating higher quality text that can be used as draft input for medical education assessment material. In addition, prompted text generation can be used for production of distractors suitable for multiple choice items used in certification exams.

**Discussion:** The current state of neural network based language models can be used to develop tools in supprt of authoring medical education exams using retrained models on the basis of corpora consisting of general medical text collections.

**Conclusion:** Future experiments with more recent transformer models (such as Grover, TransformerXL) using existing medical certification exam item pools is expected to further improve results and facilitate the development of assessment materials.

**Objective**

The aim of this article is to provide evidence on the current state of automated item generation (AIG) usi deep neural networks (DNNs). Based on earlier work, a first paper that tackled this issue used character-bas

*Address for correspondence: mvondavier@nbme.org: Matthias von Davier, Ph.D., Distinguished Research Scientist, National Bo of Medical Examiners, Philadelphia, PA, 19105

1

# AIG

<table>
<tr><td colspan="11" align="center"><b>CD Ratings: Item Stems</b></td></tr>
<tr><td></td><td colspan="5" align="center"><b>HU Items</b></td><td colspan="5" align="center"><b>AI Items</b></td></tr>
<tr><td><b>Ord</b></td><td><b>Promote</b></td><td><b>Minor Edit & Promote</b></td><td><b>Major Edit Before Promote</b></td><td><b>Delete & Rewrite</b></td><td><b>Score</b></td><td><b>Ord</b></td><td><b>Promote</b></td><td><b>Minor Edit & Promote</b></td><td><b>Major Edit Before Promote</b></td><td><b>Delete & Rewrite</b></td><td><b>Score</b></td></tr>
<tr><td>1</td><td>11</td><td></td><td></td><td></td><td>33</td><td>1</td><td>11</td><td></td><td></td><td></td><td>33</td></tr>
<tr><td>2</td><td>11</td><td></td><td></td><td></td><td>33</td><td>2</td><td>11</td><td></td><td></td><td></td><td>33</td></tr>
<tr><td>3</td><td>11</td><td></td><td></td><td></td><td>33</td><td>3</td><td>11</td><td></td><td></td><td></td><td>33</td></tr>
<tr><td>4</td><td>10</td><td>1</td><td></td><td></td><td>32</td><td>4</td><td>11</td><td></td><td></td><td></td><td>33</td></tr>
<tr><td>5</td><td>10</td><td>1</td><td></td><td></td><td>32</td><td>5</td><td>11</td><td></td><td></td><td></td><td>33</td></tr>
<tr><td>6</td><td>10</td><td>1</td><td></td><td></td><td>32</td><td>6</td><td>11</td><td></td><td></td><td></td><td>33</td></tr>
<tr><td>7</td><td>10</td><td>1</td><td></td><td></td><td>32</td><td>7</td><td>11</td><td></td><td></td><td></td><td>33</td></tr>
<tr><td>8</td><td>10</td><td>1</td><td></td><td></td><td>32</td><td>8</td><td>11</td><td></td><td></td><td></td><td>33</td></tr>
<tr><td>9</td><td>10</td><td>1</td><td></td><td></td><td>32</td><td>9</td><td>10</td><td>1</td><td></td><td></td><td>32</td></tr>
<tr><td>10</td><td>10</td><td></td><td>1</td><td></td><td>31</td><td>10</td><td>10</td><td>1</td><td></td><td></td><td>32</td></tr>
<tr><td>11</td><td>9</td><td>2</td><td></td><td></td><td>31</td><td>11</td><td>10</td><td>1</td><td></td><td></td><td>32</td></tr>
<tr><td>12</td><td>9</td><td>2</td><td></td><td></td><td>31</td><td>12</td><td>10</td><td>1</td><td></td><td></td><td>32</td></tr>
<tr><td>13</td><td>9</td><td>2</td><td></td><td></td><td>31</td><td>13</td><td>10</td><td>1</td><td></td><td></td><td>32</td></tr>
<tr><td>14</td><td>9</td><td>2</td><td></td><td></td><td>31</td><td>14</td><td>10</td><td>1</td><td></td><td></td><td>32</td></tr>
<tr><td>15</td><td>10</td><td></td><td></td><td>1</td><td>30</td><td>15</td><td>10</td><td>1</td><td></td><td></td><td>32</td></tr>
<tr><td>16</td><td>9</td><td>1</td><td>1</td><td></td><td>30</td><td>16</td><td>10</td><td>1</td><td></td><td></td><td>32</td></tr>
<tr><td>17</td><td>9</td><td>1</td><td>1</td><td></td><td>30</td><td>17</td><td>10</td><td></td><td>1</td><td></td><td>31</td></tr>
<tr><td>18</td><td>8</td><td>3</td><td></td><td></td><td>30</td><td>18</td><td>10</td><td></td><td>1</td><td></td><td>31</td></tr>
<tr><td>19</td><td>8</td><td>3</td><td></td><td></td><td>30</td><td>19</td><td>10</td><td></td><td>1</td><td></td><td>31</td></tr>
<tr><td>20</td><td>8</td><td>3</td><td></td><td></td><td>30</td><td>20</td><td>10</td><td></td><td>1</td><td></td><td>31</td></tr>
<tr><td>21</td><td>9</td><td>1</td><td></td><td>1</td><td>29</td><td>21</td><td>9</td><td>2</td><td></td><td></td><td>31</td></tr>
<tr><td>22</td><td>9</td><td></td><td>2</td><td></td><td>29</td><td>22</td><td>9</td><td>2</td><td></td><td></td><td>31</td></tr>
<tr><td>23</td><td>8</td><td>2</td><td>1</td><td></td><td>29</td><td>23</td><td>9</td><td>2</td><td></td><td></td><td>31</td></tr>
<tr><td>24</td><td>7</td><td>4</td><td></td><td></td><td>29</td><td>24</td><td>9</td><td>2</td><td></td><td></td><td>31</td></tr>
<tr><td>25</td><td>7</td><td>3</td><td>1</td><td></td><td>28</td><td>25</td><td>9</td><td>2</td><td></td><td></td><td>31</td></tr>
<tr><td>26</td><td>7</td><td>3</td><td>1</td><td></td><td>28</td><td>26</td><td>9</td><td>2</td><td></td><td></td><td>31</td></tr>
<tr><td>27</td><td>7</td><td>3</td><td>1</td><td></td><td>28</td><td>27</td><td>9</td><td>2</td><td></td><td></td><td>31</td></tr>
<tr><td>28</td><td>7</td><td>2</td><td>2</td><td></td><td>27</td><td>28</td><td>9</td><td>2</td><td></td><td></td><td>31</td></tr>
<tr><td>29</td><td>6</td><td>4</td><td>1</td><td></td><td>27</td><td>29</td><td>9</td><td>2</td><td></td><td></td><td>31</td></tr>
<tr><td>30</td><td>5</td><td>6</td><td></td><td></td><td>27</td><td>30</td><td>10</td><td></td><td></td><td>1</td><td>30</td></tr>
<tr><td>31</td><td>5</td><td>6</td><td></td><td></td><td>27</td><td>31</td><td>9</td><td>1</td><td>1</td><td></td><td>30</td></tr>
<tr><td>32</td><td>5</td><td>6</td><td></td><td></td><td>27</td><td>32</td><td>9</td><td>1</td><td>1</td><td></td><td>30</td></tr>
<tr><td>33</td><td>5</td><td>6</td><td></td><td></td><td>27</td><td>33</td><td>9</td><td>1</td><td>1</td><td></td><td>30</td></tr>
<tr><td>34</td><td>6</td><td>3</td><td>2</td><td></td><td>26</td><td>34</td><td>9</td><td>1</td><td>1</td><td></td><td>30</td></tr>
<tr><td>35</td><td>5</td><td>5</td><td>1</td><td></td><td>26</td><td>35</td><td>9</td><td>1</td><td>1</td><td></td><td>30</td></tr>
<tr><td>36</td><td>4</td><td>7</td><td></td><td></td><td>26</td><td>36</td><td>9</td><td>1</td><td>1</td><td></td><td>30</td></tr>
<tr><td>37</td><td>5</td><td>4</td><td>2</td><td></td><td>25</td><td>37</td><td>9</td><td>1</td><td>1</td><td></td><td>30</td></tr>
<tr><td>38</td><td>5</td><td>4</td><td>2</td><td></td><td>25</td><td>38</td><td>8</td><td>3</td><td></td><td></td><td>30</td></tr>
<tr><td>39</td><td>5</td><td>4</td><td>2</td><td></td><td>25</td><td>39</td><td>8</td><td>3</td><td></td><td></td><td>30</td></tr>
<tr><td>40</td><td>4</td><td>6</td><td>1</td><td></td><td>25</td><td>40</td><td>8</td><td>3</td><td></td><td></td><td>30</td></tr>
<tr><td>41</td><td>6</td><td>2</td><td>2</td><td>1</td><td>24</td><td>41</td><td>8</td><td>3</td><td></td><td></td><td>30</td></tr>
<tr><td>42</td><td>5</td><td>3</td><td>3</td><td></td><td>24</td><td>42</td><td>9</td><td>1</td><td></td><td>1</td><td>29</td></tr>
<tr><td>43</td><td>3</td><td>7</td><td>1</td><td></td><td>24</td><td>43</td><td>8</td><td>2</td><td>1</td><td></td><td>29</td></tr>
<tr><td>44</td><td>3</td><td>7</td><td>1</td><td></td><td>24</td><td>44</td><td>7</td><td>4</td><td></td><td></td><td>29</td></tr>
<tr><td>45</td><td>3</td><td>7</td><td>1</td><td></td><td>24</td><td>45</td><td>7</td><td>3</td><td>1</td><td></td><td>28</td></tr>
<tr><td>46</td><td>4</td><td>4</td><td>3</td><td></td><td>23</td><td>46</td><td>7</td><td>2</td><td>2</td><td></td><td>27</td></tr>
<tr><td>47</td><td>1</td><td>7</td><td>3</td><td></td><td>20</td><td>47</td><td>6</td><td>3</td><td>2</td><td></td><td>26</td></tr>
<tr><td>48</td><td>1</td><td>7</td><td>3</td><td></td><td>20</td><td>48</td><td>5</td><td>5</td><td>1</td><td></td><td>26</td></tr>
<tr><td>49</td><td>2</td><td>5</td><td>3</td><td>1</td><td>19</td><td>49</td><td>4</td><td>6</td><td>1</td><td></td><td>25</td></tr>
<tr><td>50</td><td>1</td><td>5</td><td>5</td><td></td><td>18</td><td>50</td><td>3</td><td>4</td><td>4</td><td></td><td>21</td></tr>
<tr><td><b>Total</b></td><td>341</td><td>158</td><td>47</td><td>4</td><td></td><td><b>Total</b></td><td>450</td><td>75</td><td>23</td><td>2</td><td></td></tr>
</table>

- Classic comparison also used for AES: how does ChatGPT compare to human quality?
- Jones, Becker, & Ferris (2023)

ASC
Assessment Systems

# Item Review

- First: code your rules into your item banker
- Second: LLMs?
  1. Ask LLM to review
  2. Provide your review criteria to LLM
  3. Fine-tune LLM based on past reviews

# Item Review

*Item Writing Rules*

*1. No colon at the end of the question stem.*

*2. Do not repeat text in the answers.*

Question 1

The capital of South Korea is:

A. Seoul

B. Tokyo

C. London

D. Paris

Question 2

What is 2 + 2?

A. The answer is 4.

B. The answer is 6.

C. The answer is 8.

D. The answer is 10.

# Item Review

Here's the feedback based on your rules:

Question 1

Rule 1: Violates the rule about not using a colon at the end of the question stem.

Rule 2: This question follows the rule, as the answers don't repeat text.

Question 2

Rule 1: This question adheres to the formatting rule (no colon).

Rule 2: Violates the rule, as all answers repeat "The answer is."

Suggested Revisions:

Question 1

The capital of South Korea is

A. Seoul

B. Tokyo

C. London

D. Paris

Question 2

What is 2 + 2?

A. 4

B. 6

C. 8

D. 10

Let me know if you need further assistance!

ASC
Assessment Systems

# Item Alignment

- Give it your curriculum and item bank, tell it to align all items to the curriculum as well as anything else like Bloom's Taxonomy



*AI: "Image of teacher using Bloom's Taxonomy"*

# AI: Automated essay scoring (AES)

- Feed your essays and rubrics to an LLM
- Can also fine-tune based on some human ratings
- Produces scores on the rubrics for each student

ASC
Assessment Systems

# AI: Automated essay scoring (AES)

- Here are some results from US high school equivalency exam
  - R or Python with basic NLP/ML: 0.65 – 0.80
  - LLAMA-3 without fine-tuning: 0.18
  - LLAMA-3 with fine-tuning (long run time with high costs): 0.81

ASC
Assessment Systems

# AI: Automated essay scoring (AES)

- Funny thing is, every website inquiry I have received thinks that #2 above is the only approach, and that it is a panacea
- They also think that:
  - No humans ever need to mark anything
  - We don't even need to establish rubrics, the AI will just "figure out what growth mindset means"

**ASC** Assessment Systems

# Part 3
# Recommendations and Future

# Using AI

- We have been using ML and AI in psychometrics for more than a century!
- Most of it is highly specific, developed to meet certain needs
- Not a panacea, but very powerful in the right hands
  - E.g., IRT revolutionized a lot of tricky problems

ASC
Assessment Systems

# Using LLMs

- They are being used more and more
- Very powerful, but not a panacea
- Huge time saver but still needs humans
- Consider the results from Jones et al. – better results than human writers, but tons of items still need editing

ASC
Assessment Systems

# Thoughts on AES

- Being able to explain and support the model is very important… that is <u>validity</u>!
- Using LLMs might have the same accuracy, but is a total black box; a lawsuit waiting to happen
  - Especially true without fine-tuning, which also means that no humans have rated any students
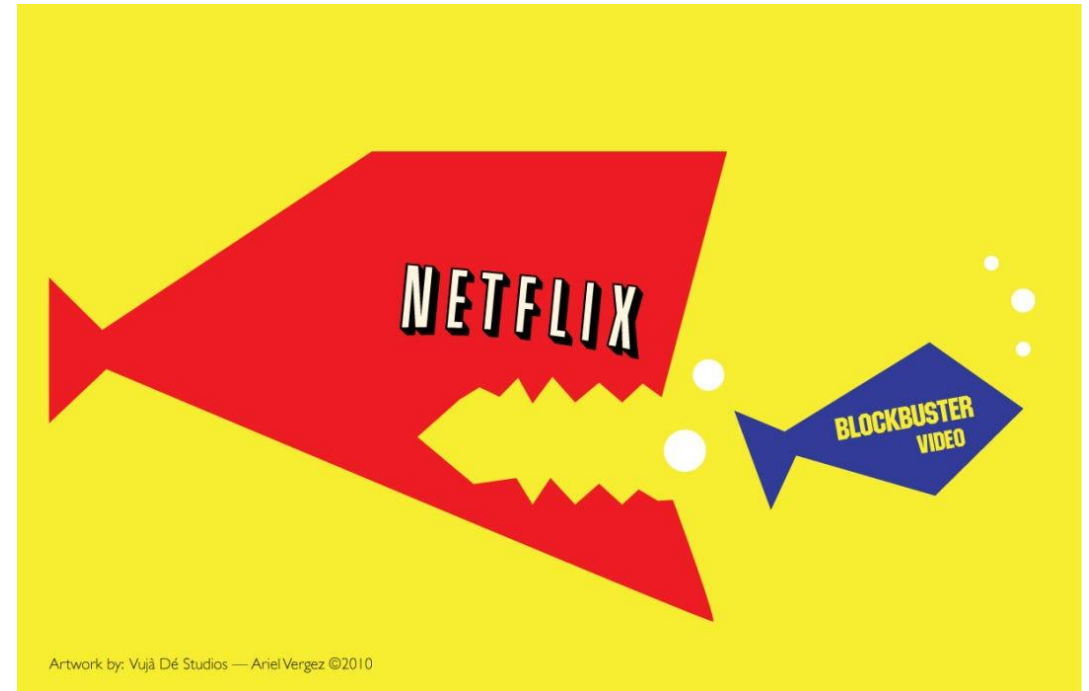
ASC
Assessment Systems

# Security?

- Do you want to feed your content to an LLM to make items?
  - Past items
  - Your textbooks etc.
- How about data from past item reviews or item statistics?
- For AES, how about all your student essays and human marks?
- Public vs self-hosted

# The future?

- As in many other areas, LLMs and other AI are making an impact
- Even if they only provide 10% improvement, you do not want to be left behind
- It is irresponsible to avoid AI altogether
  - Friend that was told to stop all AI, even spam filters!
- So you need to decide what to use, and with what guardrails and post hoc checks



Artwork by: Vujà Dé Studios — Ariel Vergez ©2010

# The future?



AI

# Questions?

nate@assess.com