

NEAS 2024

MANAGEMENT CONFERENCE

QUALITY REVOLUTION:

RETHINKING, RESHAPING, AND
REDEFINING EXCELLENCE IN ELT



NEAS

QUALITY ASSURANCE
IN EDUCATION AND
TRAINING

neas.org.au

Assessment in the Digital Age: Innovative Approaches to Testing and Evaluation

Duolingo

Brett Blacker & James Holden

Session Goals

Participants will understand:

- Different ways in humans are collaborating with AI to advance assessment, in particular related to building item banks and security
- How peers are approaching and collaborating with AI tools in the ELT context
- The latest research and development updates from Duolingo

Conference Themes

- *How can we redefine and measure excellence in ELT beyond traditional quality assurance methods, taking into account innovative approaches and revolutionary strategies?*
- *What transformative frameworks and best practices can be explored to reshape the future of ELT quality assurance, ensuring it remains relevant in an evolving educational landscape?*

Overview

1. Introduction
2. Collaborating with AI - Views on the present and future
3. Duolingo's approach to AI
 - Building Assessment using AI
 - Securing Assessment using AI
 - Latest from Duolingo Research / Development



About Duolingo

Duolingo's mission is to develop the best education in the world and make it universally available.

World's Most Downloaded Education App



Languages
Music



Literacy



Numeracy



World's Leading High-Stakes Digital Language Test

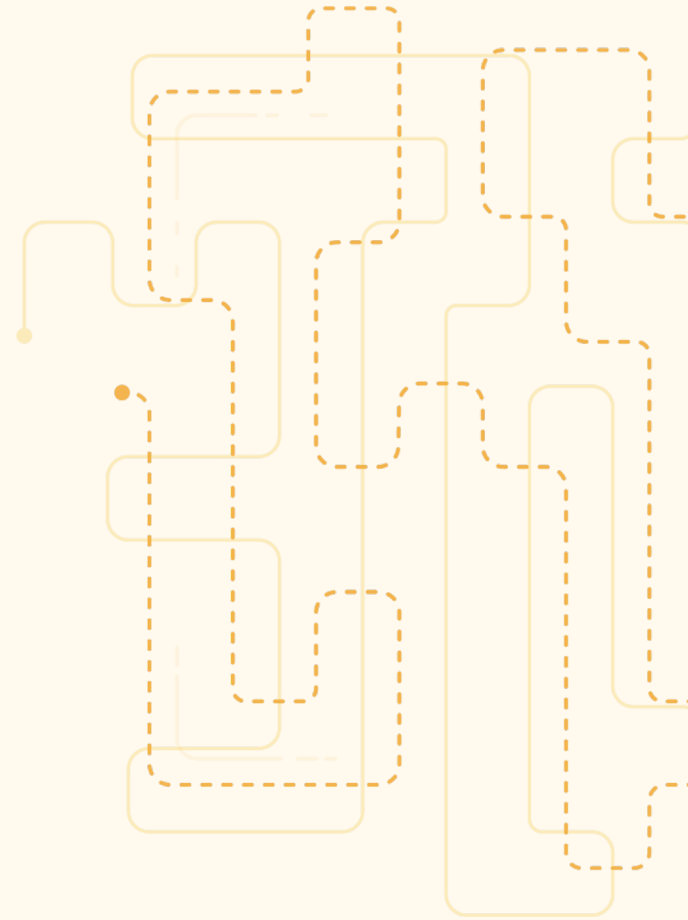


duolingo english test



Professor Luis von Ahn, Duolingo's co-founder and CEO

***Collaborating with AI -
Views on the present
and future***



Would you prefer to have a human doctor or an AI algorithm diagnose a medical condition?

Would you feel comfortable taking a plane flown by only AI?

*Would you trust an AI tool to
organise your finances?*

If a human was also involved in confirming and validating the AI's decisions and actions, would you feel more comfortable?

Collaborating with AI is about leveraging the best of machines and the best of humans.

We always want a 'human-in-the-loop'

Best of Machines

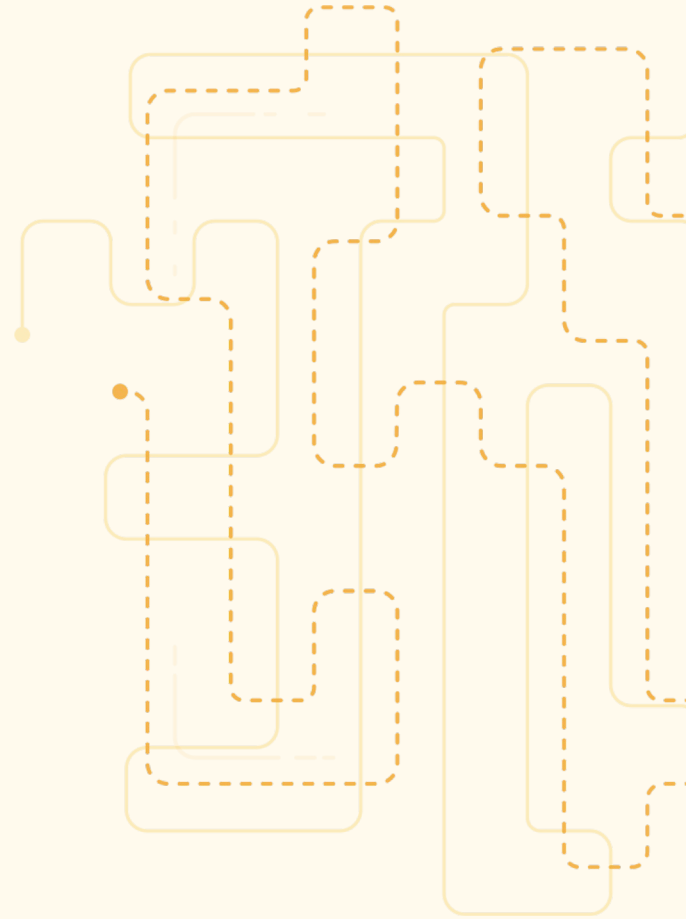
- **Data Processing and Analysis**
- **Consistency and Scalability**
- **Predictive Capabilities**

Best of Humans

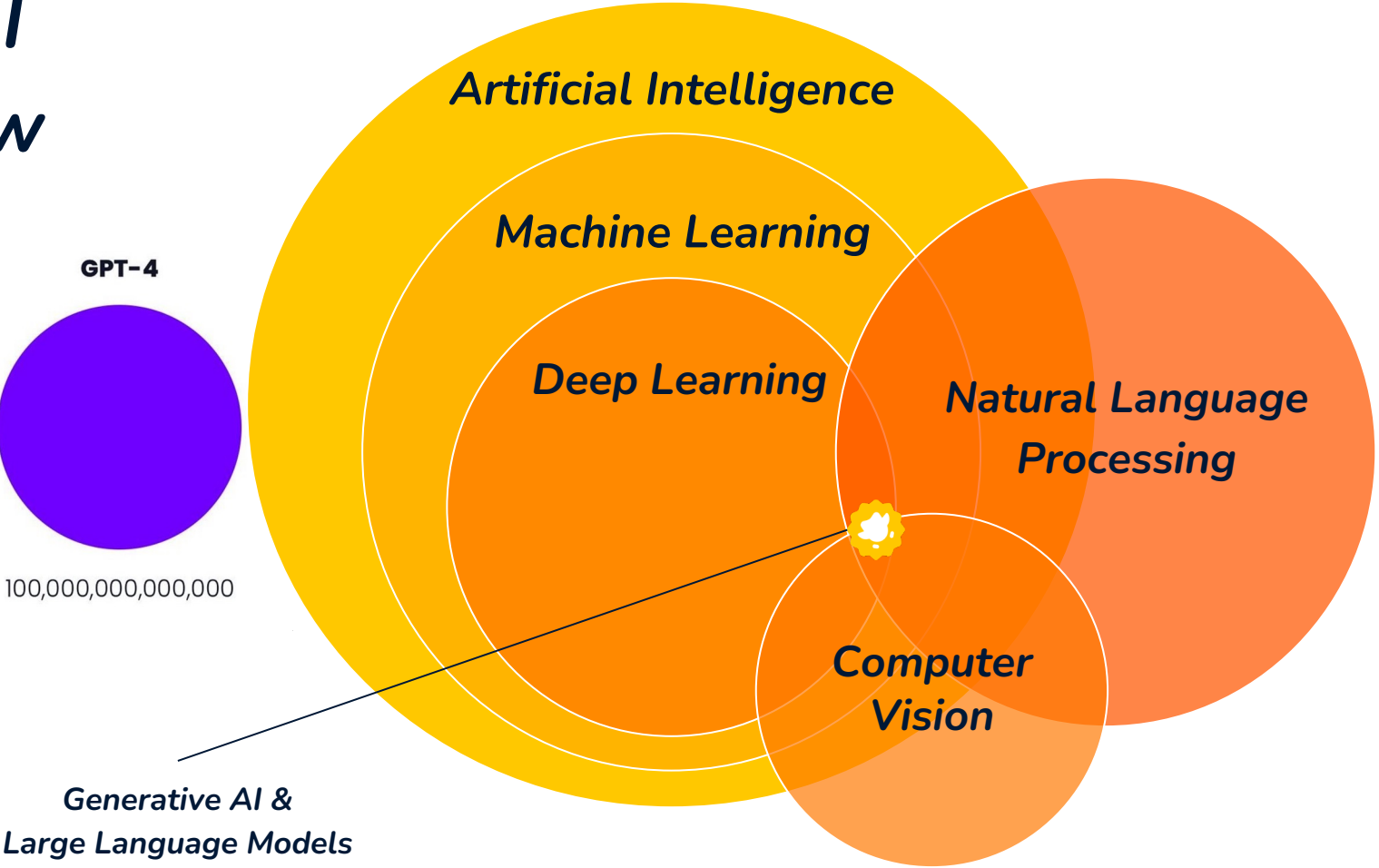
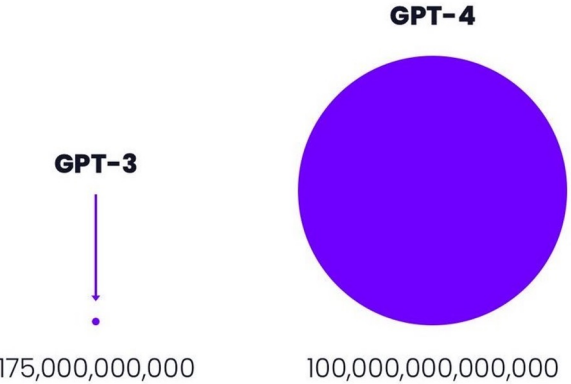
- **Creativity and Innovation**
- **Emotional Intelligence**
- **Adaptability and Intuition**

Humans program the AI - AI crunches the data - Humans validate and confirm

Human-in-the-loop AI Duolingo's approach

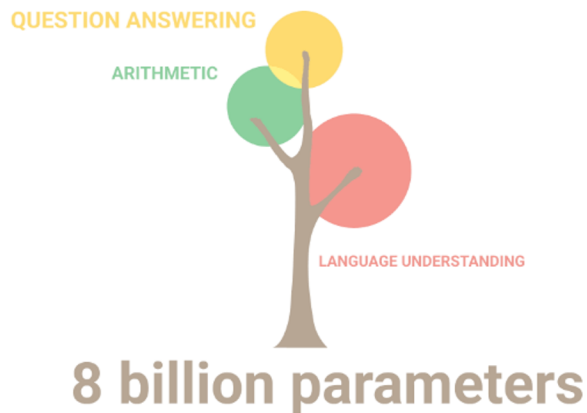


Quick AI overview



Generative AI & Large Language Models

More data = more capabilities



Med-PaLM (Google) was the first LLM to obtain a passing score on U.S. medical licensing questions, and in addition to answering both the multiple choice and open-ended questions accurately, it also asks questions, provides reasoning and is able to evaluate its own responses

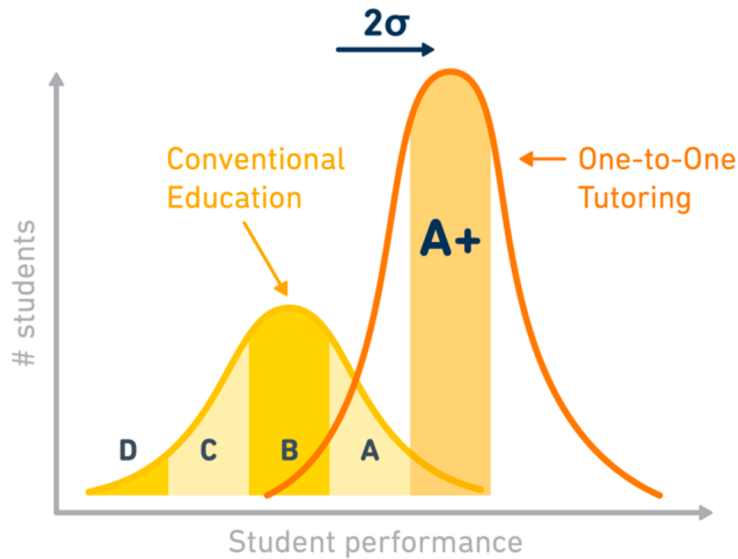
Why is Duolingo excited about AI?

Duolingo's mission is to develop the best education in the world and make it universally available.

- AI offers the best opportunity to achieve this mission and meet students 'where they are'
- Humans must be in the loop at all times to maximise the strengths of both humans and machines



Blooms 2 Sigma Problem (1984)

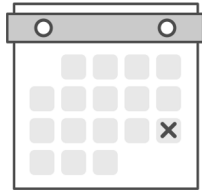


Bloom's 2-sigma problem

- The average student achievement with a 1-1 tutor is better than 98% of students from a traditional class
- How can we deliver 1-1 tutoring against the limitations of time, resources, and the varying needs of students?

AI = More inclusive, personalised learning which is no longer constrained by a learners resources

Overcoming traditional testing



Book a date
Pay >\$400



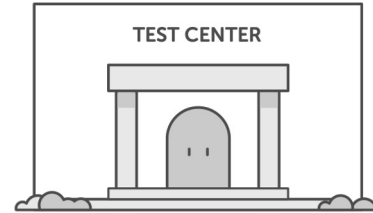
Generative AI item creation



Go to test-center or
even travel to a
different city



*Online on-demand test
Asynchronous AI assisted security*

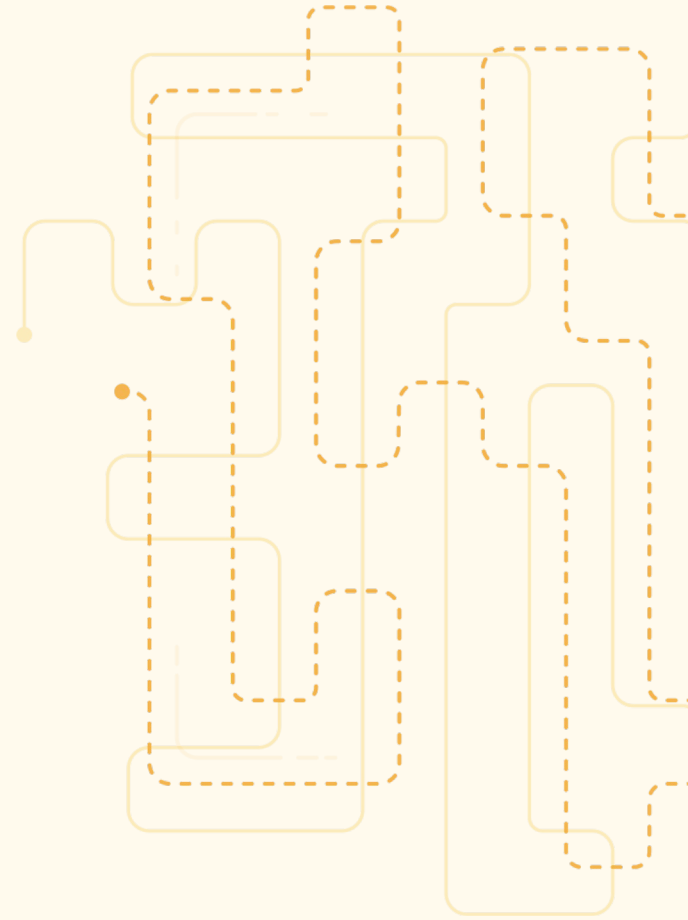


+3 hours to take
a test

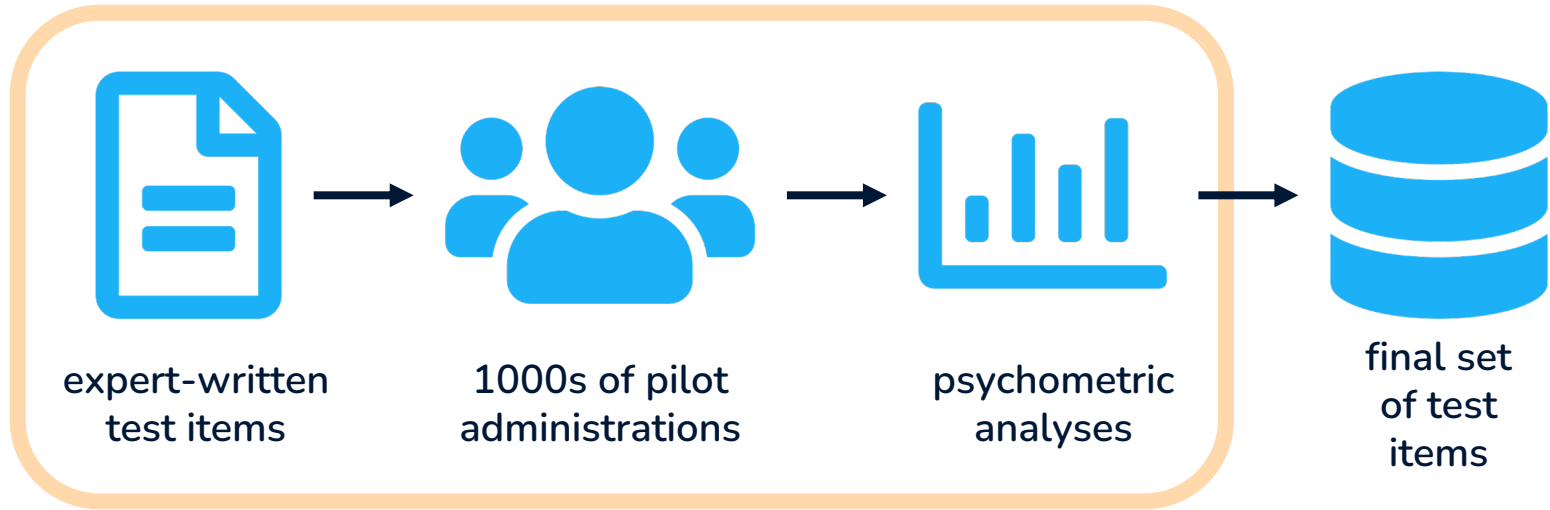


*Computer Adaptive
Testing*

Building assessment with AI



Conventional test development



Time-consuming, vulnerable &
expensive

The Promise of Large Language Models

- Large language models (LLMs) are trained to generate text from billions or trillions of words.
- LLM's excel at few-shot learning - *can mirror style, format, and content given extremely few examples*
- Allows rapid prototyping of new content and item types

 Gemini

 OpenAI

 Meta AI

ANTHROPIC

GPT-4 + DET

Item Generation

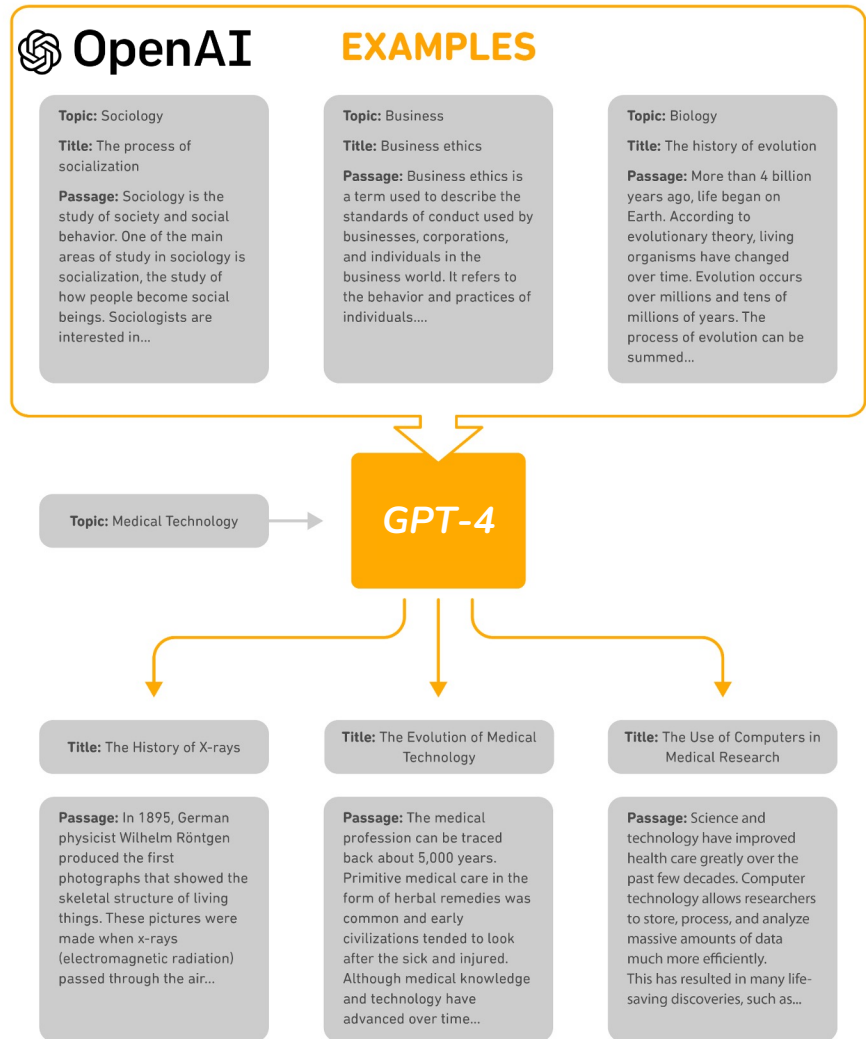
Goal

Create texts and associated materials, such as:

- Reading passages
- Conversations
- Comprehension tasks, correct answers and distractors
- Other information necessary for automated scoring

Method

Generate passages and titles conditioned by other academic content and a specified topic using prototypical examples



Example: Interactive Listening

Goal

Generate short conversations oriented around academic scenarios.

Challenges

1. Scenario selection - what types of situations to cover?
→ Human expert
1. Scenario expansion - adding details to create interesting & diverse conversations
→ LLMs



You will participate in a conversation about the scenario below.

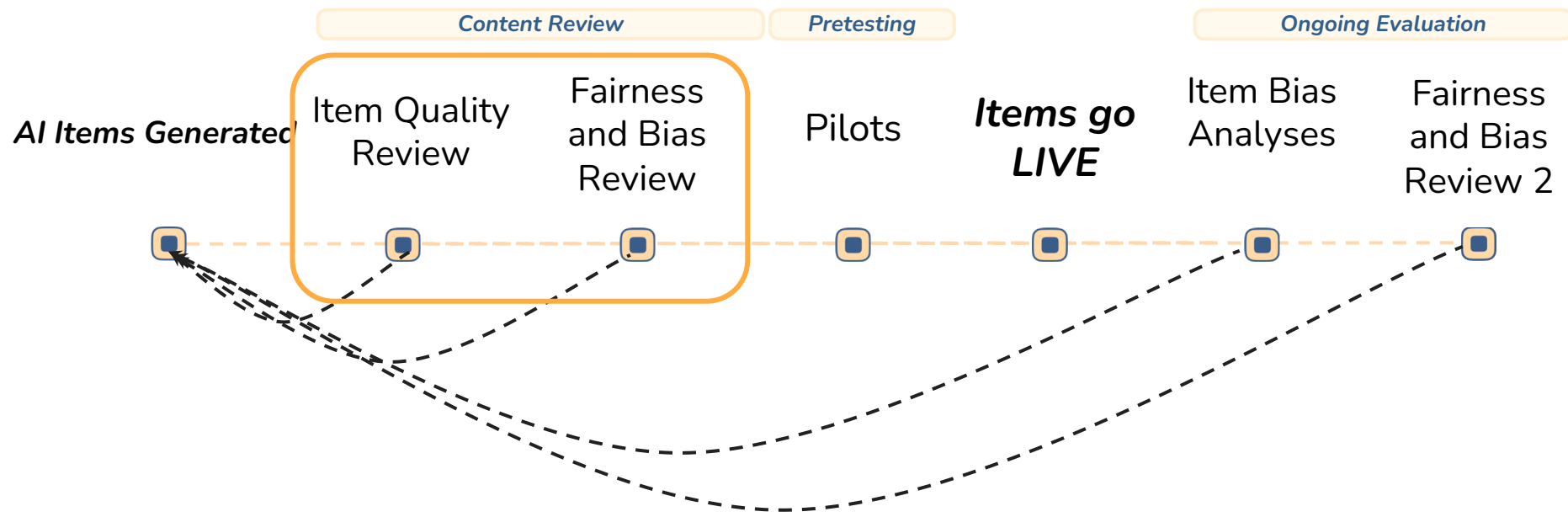
You are a student in a journalism class. After today's lecture on reporting, you approach your professor to ask them more about what it means to cultivate sources.

START

*AI item banks require
extensive human review*

Human-In-The-Loop AI

Goal: Use collaboration of humans and AI to perform a task that is difficult for both



FAB review process

The most important thing you can do to protect your family is to have a will. A will lets you appoint someone as the legal guardian of your children if something should happen to you. Many states require that the guardian has a family relationship with the child. It's also a great idea to set up life insurance for yourself and your children. Remember, you can never be too prepared for the unexpected.

DECISION Feedback

FAB Decision

Pass Fail

Flags ?

Content Warning

Global Relevance

Stereotype

Comments

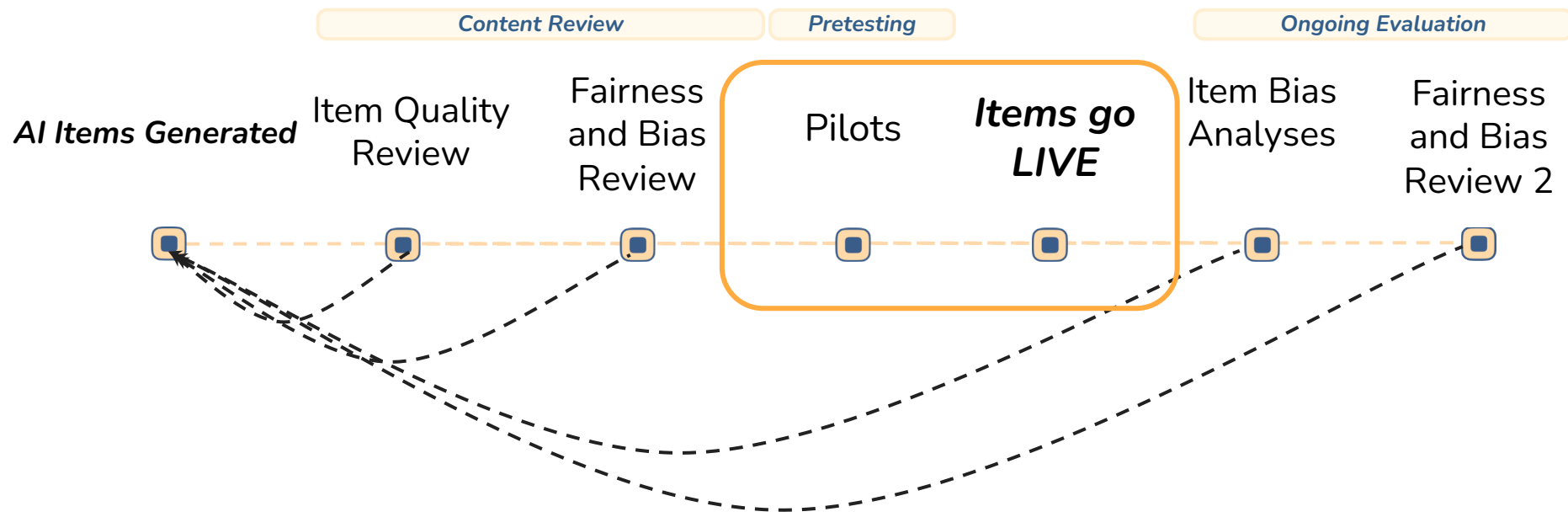
[BACK TO LIST](#) [SAVE](#) [SAVE AND REVISIT](#) [SUBMIT](#)

Calibration items

Text	Your Label	Our Label	Your Codes	Our Codes	Your Comment	Explanation
<p>As the general American public becomes informed about the cost of a university education, they are becoming more reluctant to pay for it. The rising costs of health care, rent, food, and commodities overall have made people take a closer look at their finances. Some people can't afford to pay to go to college. The fear of student loan debt has driven young people away from a college education.</p>	Fail	Fail	Content Warning	Content Warning	this topic might be sensitive to some, particularly those struggling with financial challenges related to education, although it is described in a general way without targeting or blaming specific groups. Since many of the TT demographic are planning to go to college it's best to fail this item	The passage focuses on income disparity and poverty.

Human-In-The-Loop AI

Goal: Use collaboration of humans and AI to perform a task that is difficult for both



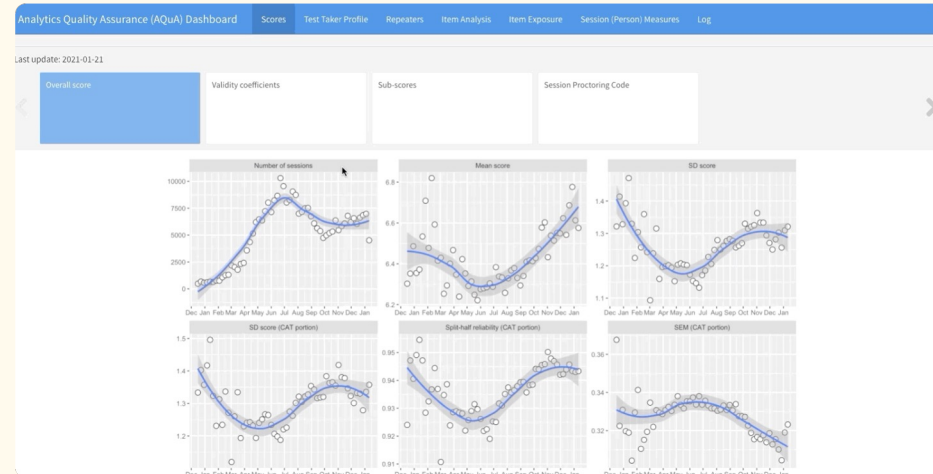
AQuAA

Analytics for Quality Assurance in Assessment

We utilise an internally developed dashboard that tracks and reports on all validity-related metrics of our test daily across the globe



englishtest.duolingo.com/research



Content Review

Item Development Timeline

Content Review

Pretesting

Yearly Evaluation

Item Quality
Review

Fairness and
Bias Review

Pilots

Items go
LIVE

Item
Analyses

Fairness and
Bias Review 2

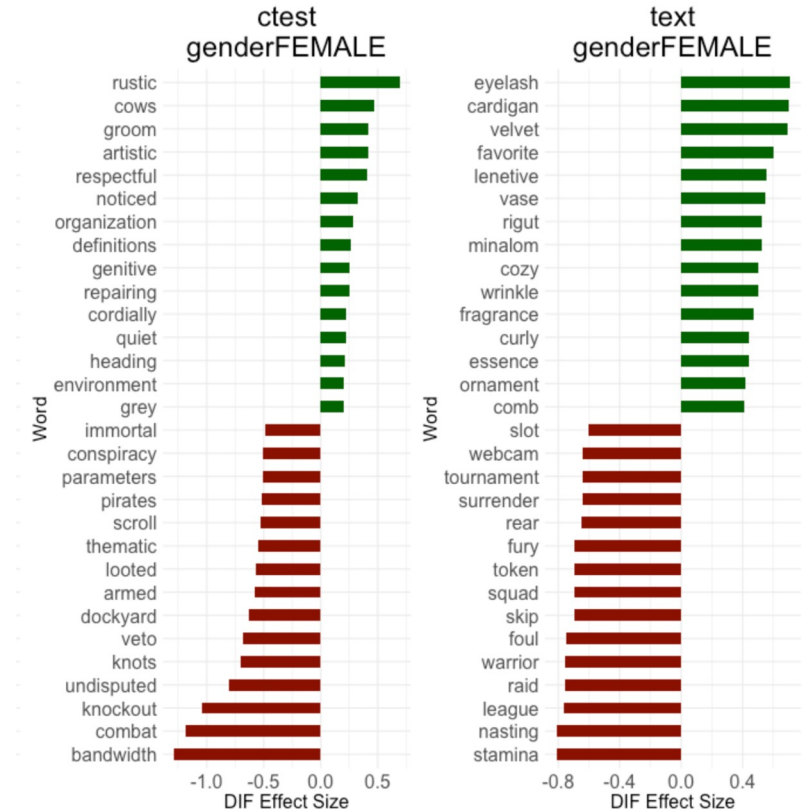


Differential Item Functioning (DIF)

A statistical method for testing whether items are “biased”

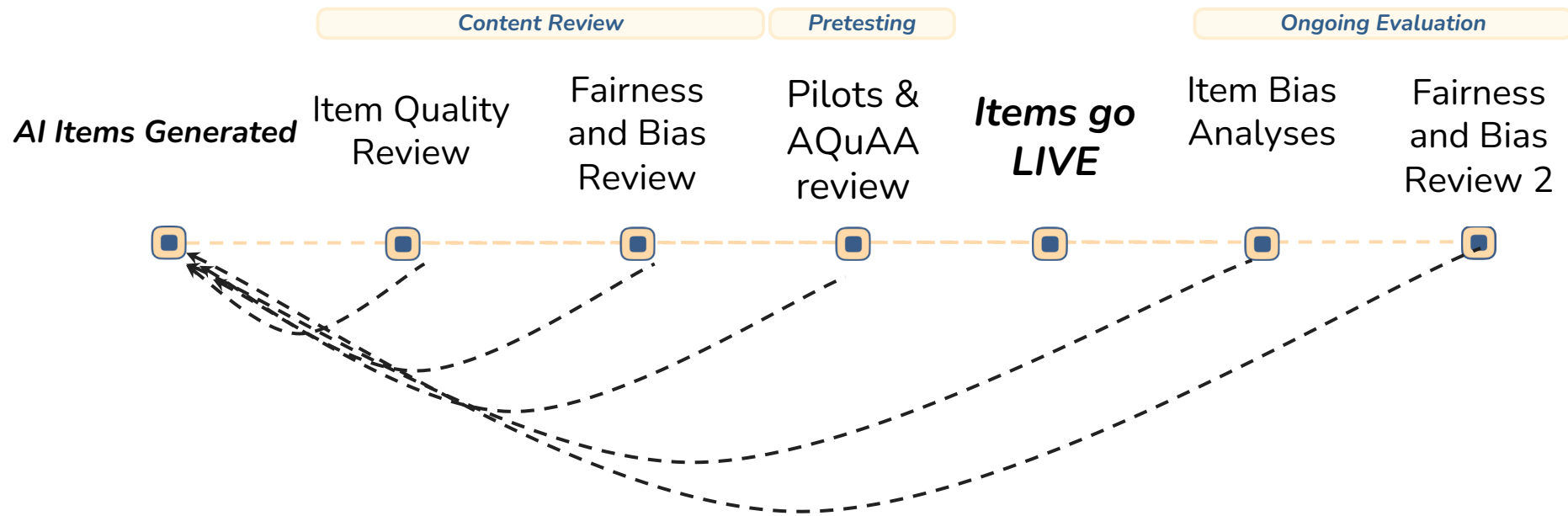
Item “bias”:

Individuals with the **same level of English proficiency** but **different backgrounds** (e.g., male and female) have **different probabilities** of answering an item correctly.



Human-In-The-Loop AI

Goal: Use collaboration of humans and AI to perform a task that is difficult for both



Benefits of AI item creation

Item banks which are more:

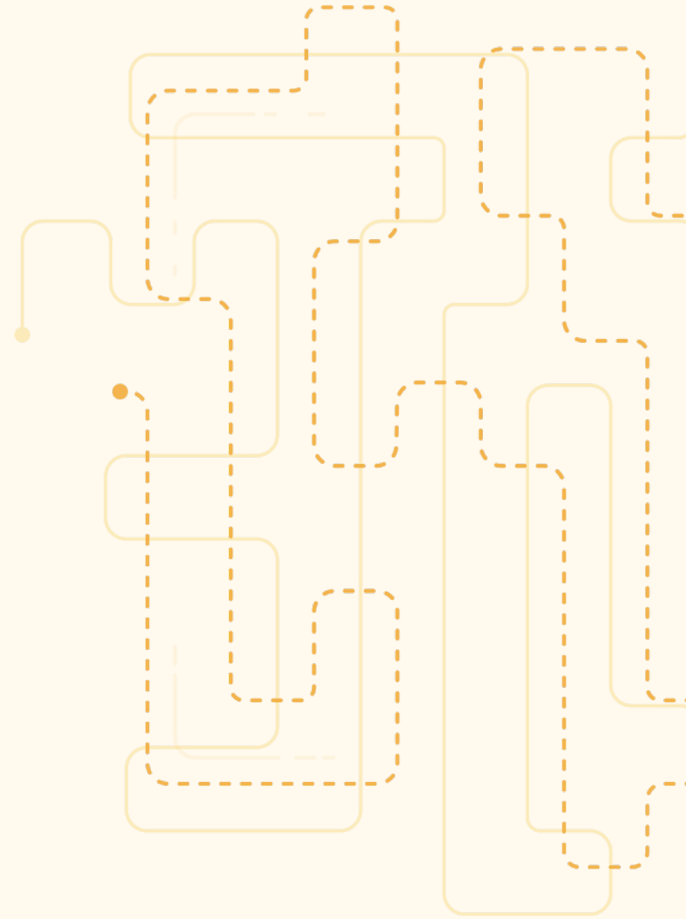
- Secure
- Aware and adaptable to test-taker population trends
- Efficient and affordable to create
- Scalable
- Customised

Our Own Contexts - Assessment Creation

Have you attempted to create your own assessments / item banks using LLMs?

What review process or human involvement do you use in your own assessment development?

Securing assessment with AI



Our Own Contexts - Security

How are digital innovations being used in security (personal, home, finance)?

Is digital security more or less trusted?

The live proctoring disadvantage

Live proctoring relies on a chain of trust consisting of dozens of humans.

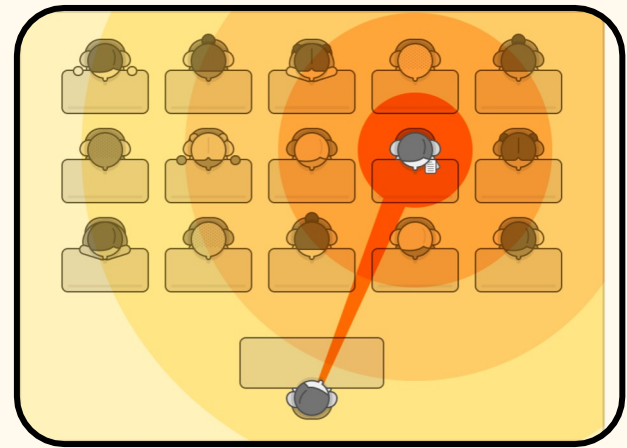
- A single point of failure can disrupt the entire chain.

Live proctoring doesn't allow for multiple rounds of review.

- No evidence of test taker behavior other than the proctor's narrative, based on memory.

Since proctors and test takers are in the same room, they cannot be anonymous.

A single instance of cheating compromises the entire testing cohort.



Los Angeles MAGAZINE

CRIME EDUCATION

A Chinese Cheating Ring at UCLA Reveals an Industry Devoted to Helping International Students Scam Grades

Operation [redacted] was overshadowed by Operation Varsity Blues, but it's just as scandalous

By Christopher Beam - April 26, 2019



On paper, Liu Cai was a model student. After moving to the United States from Beijing, he majored in biology at UCLA and volunteered at the Boys & Girls Club. A former teacher, Jess Eshewer, remembers

Defining a digital threat model

<i>Digital Threats</i>	<i>Attacker</i>	<i>Mitigations</i>
Test theft	Scrapers	Large item pool, adaptive engine, Screenshot prevention etc.
Getting third party aid during test	Cheaters, Cheating rings	Monitoring system softwares, preventing window switching, multi-layered evidence collection, no human<>test taker contact
Identity and account abuse	Imposters	ID verification, previous session and account matching, ID history, shared video response etc.

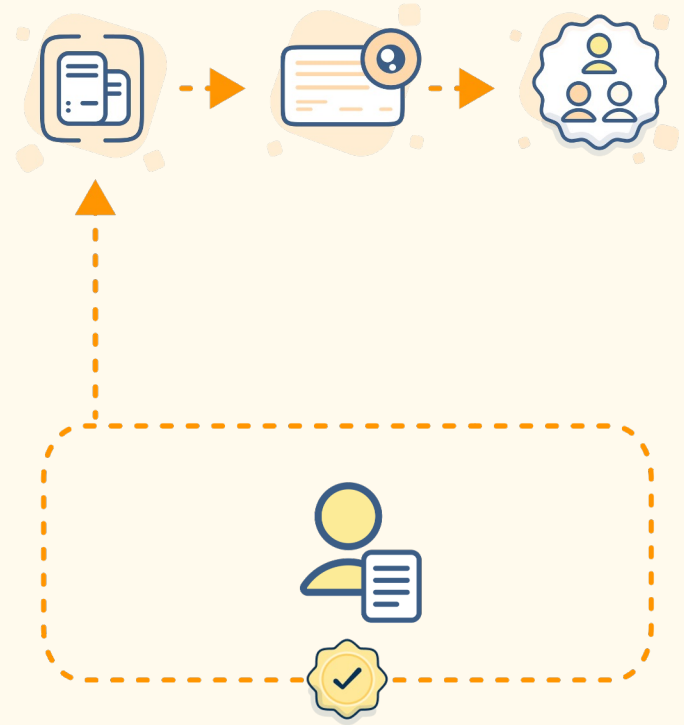
“Human-in-the-loop” AI

Remote

- Testers & proctors located all over the world who have no interaction with each other

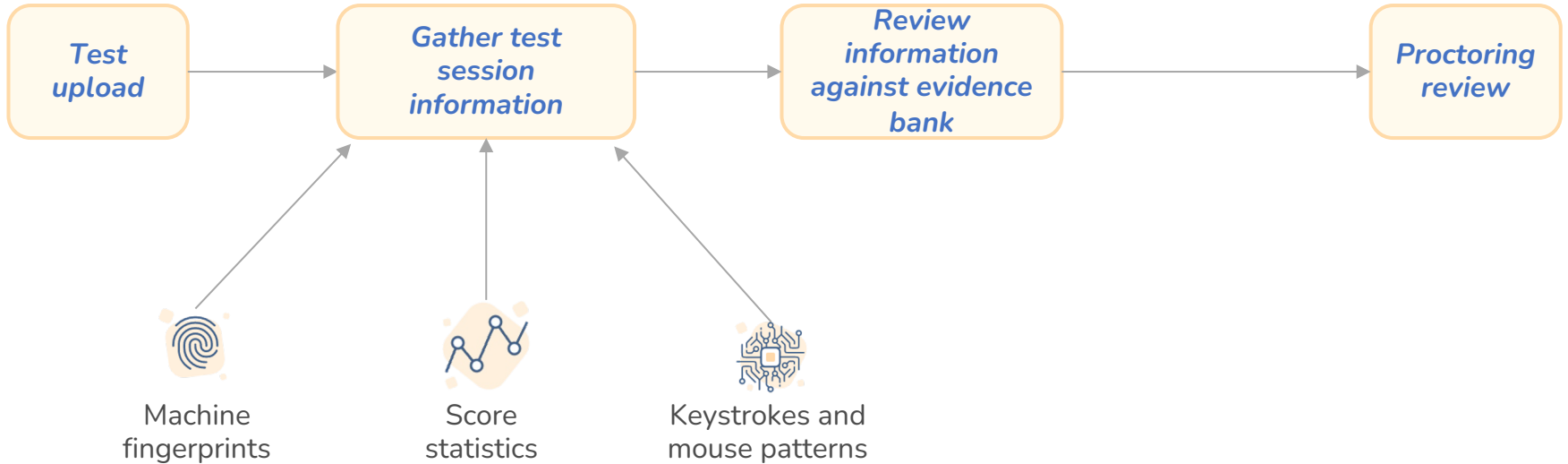
Record-and-review

- Test video is recorded and uploaded
- AI scans 150+ evidence categories
- AI flags and test video is reviewed by human proctors

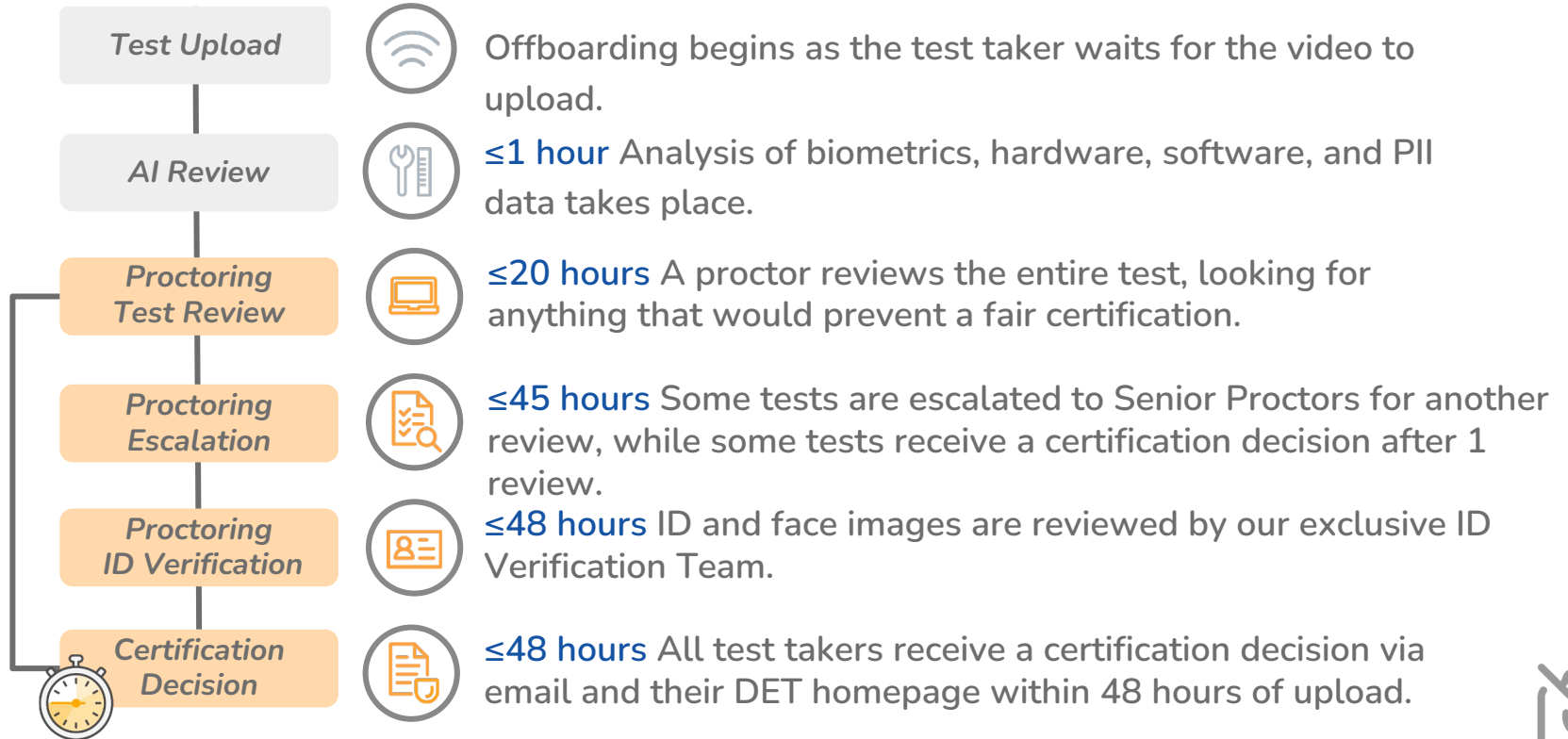


*Asynchronous
proctoring approach*

How Do We Detect Cheating?



The proctoring process



Example - Plagiarism Detection

Why is plagiarism detection necessary?

- Some test takers could cheat by searching the Internet for answers or use external resources
- Identifying and preventing such malicious behaviors is necessary to maintain the integrity of the test
- Plagiarism detection can:
 - Catch these malicious behaviors
 - Have a deterrent effect on these malicious behaviors

Plagiarism detection

The similarity of text responses are evaluated against external resources AND historical DET responses.

The screenshot displays a plagiarism detection tool interface. At the top, it shows 'Similar Texts Found!' with a 'Max similarity: 60.00%' and 'Current similarity: 60.00%'. The interface is divided into several sections:

- Current Response:** Contains a student's text about technology and screen time. Several phrases are highlighted in yellow, such as 'ho consistently spend more than hours per day watching tv are more likely to be overweight' and 'teens who play violent video games'. A callout box labeled 'Similar sentences are highlighted' points to these yellow highlights.
- Past exam:** Shows a similar text from a previous exam session, with overlapping phrases highlighted in yellow. A callout box labeled 'Similar historical test sessions' points to this section.
- Similar Results:** Lists external sources with similarity scores. For example, 'Past exam' has a similarity of 60.00%, and 'Internet: h.org' has a similarity of 47.89%. A callout box labeled 'Similar Internet content' points to these results.

At the bottom, there are two buttons: 'Not Plagiarism' and 'Plagiarism'. A callout box labeled 'Shortcuts for proctor decisions' points to these buttons.

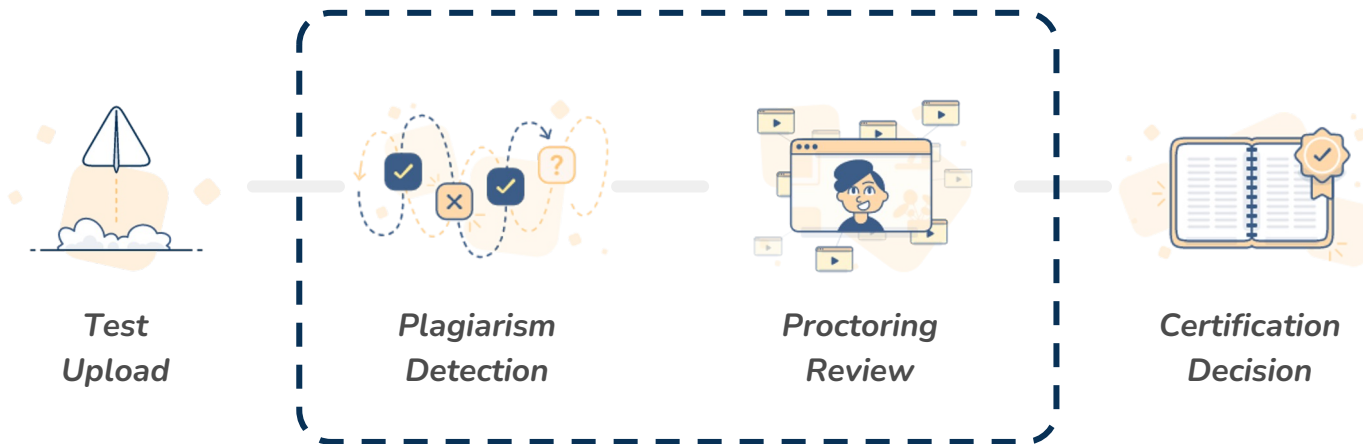


Plagiarism detection

- Expanded to spoken responses in 2023
- Allows our human proctors to focus on observing behavior



Plagiarism detection



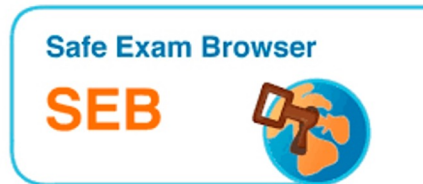
Among all the sessions that are flagged by the plagiarism detection algorithm, 93.6% are determined as having malicious behaviors by proctors.

Humans always make the final certification decision.

Digital security options for the ELT sector

Safe Exam Browser

- Accessible (free and open source - Works on computers & tablets)
- Security Strengths (Website blocking, Virtual Machine Detection, Windows Switching prevention)
- No registration system, ID monitoring, nor comprehensive cheating detection



Exam.net

- Secure browsers and screen lockdown, as well as 'Background cheat detection'
- 3 security modes, tailored to the testing environment
- Comprehensive ability to create and mark assessments, however \$\$\$



Our Own Contexts - Security in ELT

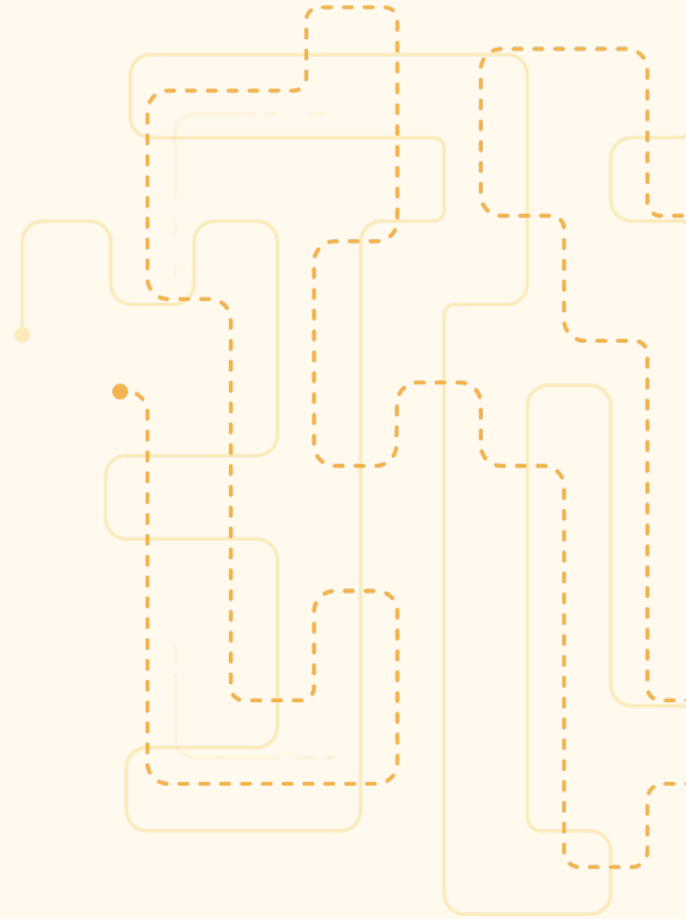
How are you currently managing security and proctoring in your own assessments?

Have you used any digital security tools or software?

Do you plan to / are under pressure to do this?

What are the potential risks?

Research & Development Update



24'V8 new items: Interactive Writing

- New task taps into **higher level writing skills (e.g., elaboration, revision)**, allowing test takers to demonstrate their academic writing proficiency better.

5:00 QUIT TEST

● **Write about the topic below for 5 minutes.** ● Write a follow-up response for 3 minutes.

Think about something that you have made with your own hands (such as a cake, a piece of furniture, or a painting). Describe the experience and the results of your work. Include specific details.

Your response

CONTINUE AFTER 3 MINUTES

3:00 QUIT TEST

HISTORY

Prompt

Think about something that you have made with your own hands (such as a cake, a piece of furniture, or a painting). Describe the experience and the results of your work. Include specific details.

Initial Response

I painted my room. I used a whole day. It wasn't easy but it was nice. So, I mixed the bucket of paint with water. I didn't know the right proportion but I did it and I started using it on the wall. It was very hectic and I did it alone. Since it was my first time, it wasn't all that nice but I took it like that. I loved the fact that I did it all by myself and it was also a good time for me.

STEP 2 OF 2

Write a follow-up response for 3 minutes.

Can you expand on your initial response by writing about the sense of accomplishment and pride in one's own abilities as a result of your work.

Your response

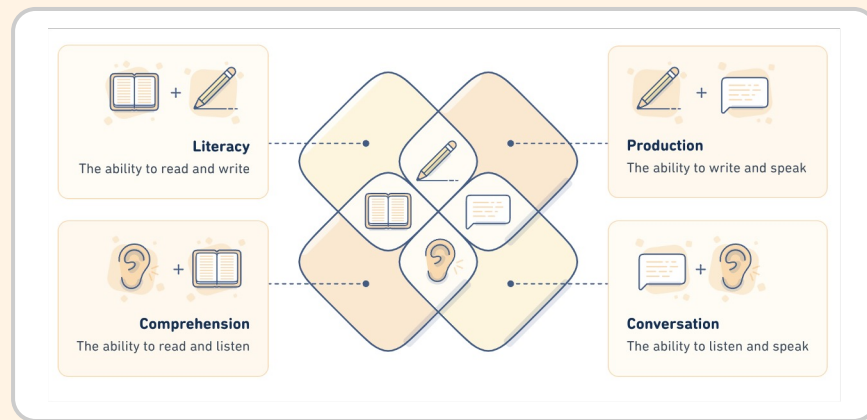
CONTINUE AFTER 1 MINUTE

DET 2024

New Scoring Coming in July

- All tests will now include new **Speaking, Writing, Reading, and Listening subscores** alongside our existing integrated modality subscores.

“DET requires a score of 115, with no subscore below 105” etc.



DET validity & reliability research

2023 - Predictive and Academic Validity: DET scores display a strong correlation with academic evaluations of students' English proficiency and subsequent academic performance ([Isbell et al. \(2023\)](#))

2023 - Concurrent Validity: In a 2023 study of over 5000 students, DET scores display significant, strong correlations with scores from TOEFL and IELTS English language tests taken by the same individuals. ([Cardwell et al. \(2023\)](#))

2023 - Test-Retest Reliability: Overall test-retest reliability of **0.93**, along with the reliability of subscores ranging from 0.90 to 0.92, underscores the DET's scoring consistency. In March 24 test-retest reliability has reached **0.95** (*Journal of applied psychological measurement - Will be published in 2024*)

2024 - Task Duration: The Impact of Task Duration on the Scoring of Independent Writing Responses. Equally high test reliability and criterion validity (writing) was demonstrated by a 5 minute task, compared to a 20 minute task. ([Naismith et al. 2023](#))

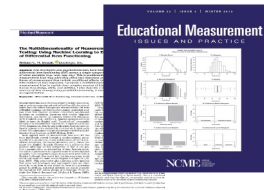
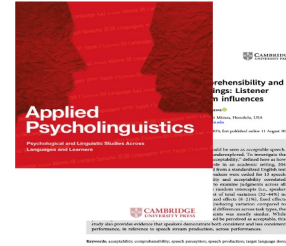
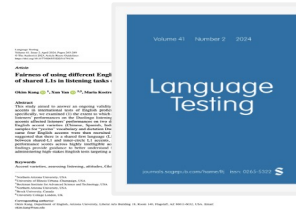
Ongoing - External Validity:

Endorsed a high-stakes English test by Australian English Quality Assurance organisation NEAS.

Endorsed with full CEFR validity alignment by the UK Chartered Institute of Linguists

Successful tracer studies from the world's most elite universities

- englishtest.duolingo.com/scores
englishtest.duolingo.com/research



Technical manual

- Introduction to the test
- Item type & construct coverage
- Test development and scoring
- Test administration & security
- Test taker demographics and performance statistics
- Accessibility, fairness & bias
- Quality assurance

Duolingo English Test: Technical Manual

 duolingo english test
Duolingo Research Report
August 8, 2022 (111 pages)
<https://english-test.duolingo.com/research>

Ramsey Cardwell*, Geoffrey T. LaFlair*, Ben Naismith*, and Burr Settles*

Abstract

The Duolingo English Test Technical Manual provides an overview of the design, development, administration, and scoring of the Duolingo English Test. Furthermore, the Technical Manual reports validity, reliability, and fairness evidence, as well as test-taker demographics and the statistical characteristics of the test. This is a living document whose purpose is to provide up-to-date information about the Duolingo English Test, and it is updated on a regular basis (last update: August 8, 2022).

Contents

1	Introduction	3
2	Purpose	3
3	Item Type Construct Descriptions	4
3.1	C-Test	5
3.2	Yes/No Vocabulary (Text)	5
3.3	Yes/No Vocabulary (Audio)	6
3.4	Dictation	6
3.5	Elicited Imitation (Read-aloud)	6
3.6	Interactive Reading	8
3.7	Extended Writing & Writing Sample	10

*Duolingo, Inc.

Corresponding author:
Geoffrey T. LaFlair, PhD
Duolingo, Inc. 5900 Penn Ave, Pittsburgh, PA 15206, USA
Email: english-test-research@duolingo.com

Collaborating with AI is about leveraging the best of machines and the best of humans.

We always want a 'human-in-the-loop'

Best of Machines

- **Data Processing and Analysis**
- **Consistency and Scalability**
- **Predictive Capabilities**

Best of Humans

- **Creativity and Innovation**
- **Emotional Intelligence**
- **Adaptability and Intuition**

Humans program the AI - AI crunches the data - Humans validate and confirm