

---

# The Duolingo English Test Responsible AI Standards



Duolingo Research Report DRR-25-05  
August 13, 2025 (16 pages)  
<https://englishtest.duolingo.com/research>

**Jill Burstein**

## Abstract

As AI has become central to digital assessments, ensuring its responsible use is critical—especially in high-stakes contexts. The Duolingo English Test (DET) Responsible AI Standards was the **first** published, comprehensive framework that addressed responsible AI (RAI) for an educational assessment program. The standards are grounded in four ethical principles—Validity and Reliability, Fairness, Privacy and Security, and Accountability and Transparency. They guide AI use across the DET’s test design, measurement, and security frameworks, aiming to uphold test quality and equity. The standards further shape the integrity and fairness of the test by directly supporting test-taker experience and test validity. Informed by cross-disciplinary discussion, the DET RAI standards support risk mitigation, transparency, and test auditing. First published in 2023 as a living document, the DET RAI Standards remains open for public comment to foster ongoing conversation around ethical AI use in education and assessment.

## Keywords

AI for assessment, Duolingo English Test research, language assessment, responsible AI

---

## Corresponding author:

Jill Burstein  
Duolingo, Inc. 5900 Penn Ave, Pittsburgh, PA 15206, USA  
Email: [jill@duolingo.com](mailto:jill@duolingo.com)

## Contents

Contributors . . . . .	4
Invitation for Public Comment . . . . .	4
Version 3 Notes . . . . .	4
Legal Disclaimer . . . . .	4
Introduction . . . . .	5
The DET Responsible AI Standards . . . . .	6
The Duolingo English Test . . . . .	6
The DET Responsible AI Standards Development . . . . .	6
The Standards . . . . .	7
1. Validity & Reliability . . . . .	7
Goal 1.1 . . . . .	
Specify processes required to build a validity argument. . . . .	7
Goal 1.2 . . . . .	
Evaluate AI used in test item creation, item calibration, and scoring. . . . .	8
2. Fairness . . . . .	8
Goal 2.1 . . . . .	
Specify how the use of AI facilitates test-taker access, accessibility, accommodations, inclusion, and appeals. . . . .	8
Goal 2.2 . . . . .	
Specify test-taker and human expert demographic representation, algorithms known to contain or generate bias, and potential human bias. . . . .	9
3. Privacy & Security . . . . .	9
Goal 3.1 . . . . .	
Specify methods to ensure privacy and security associated with data origin, data collection and processing, and data management. . . . .	9
Goal 3.2 . . . . .	
Specify how to maintain test-taker privacy, item security, and test-taker security during test administration. . . . .	10
Goal 3.3 . . . . .	
Specify fair and reliable test security proctoring protocols, item pool development, and psychometric procedures for test security. . . . .	10
4. Accountability & Transparency . . . . .	11
Goal 4.1 . . . . .	
Assess how AI processes impact stakeholders. . . . .	11
Goal 4.2 . . . . .	
Document AI used for building the validity argument, test item creation, test item calibration, and scoring. . . . .	11

Goal 4.3	
Document processes for human-in-the-loop interactions with AI. . . . .	12
Goal 4.4	
Document human expert qualifications required for human-in-the-loop activities that support AI for the DET. . . . .	12
Goal 4.5	
Disseminate research about use of AI to various stakeholder communities. . . . .	12
Goal 4.6	
Publish information about how AI is used on the DET, and usage of test-taker data . . . . .	13
Where Theory Meets Practice	13
References	14

## Contributors

A number of Duolingo colleagues contributed ideas and expertise to the original or updated versions of these standards. Areas of expertise included: applied linguistics, computational psychometrics, language assessment, law, machine learning, statistics, and test security. Contributors were Alina von Davier, Kevin Yancey, Will Belzak, Klinton Bicknell, Carl Gottlieb, Rose Hastings, Ian Riggins, and Mark Zheng. Reviews were conducted by Ramsey Cardwell and Sophie Wodzak.

Pascale Fung, Chair Professor at the Department of Electronic and Computer Engineering, Hong Kong University of Science & Technology, played a key role in reviewing, organizing, and articulating the original version of the standards upon which all updates are based.

## Invitation for Public Comment

We developed these standards to advance thinking in the field of assessment with regard to the ethical use of AI for testing. The standards were written by leveraging industry guidelines and AI ethics, and – as noted in the Contributors section – by engaging with stakeholders from multiple, relevant disciplines. The Duolingo English Test (DET) Responsible AI (RAI) Standards were also informed by the American Educational Research Association, the American Psychological Association, and the National Council on Measurement in Education Standards (AERA & APA & NCME, 2014), the International Test Commission & Association of Test Publishers (ITC-ATP) guidelines for technology-based assessment (International Test Commission & Association of Test Publishers, 2022), and numerous academic and policy publications in AI ethics cited throughout this document. Continued updates of the DET RAI Standards also integrate public comments.

The cross-disciplinary approach helped to formulate the DET RAI Standards and related goals which contribute to the DET's validity, reliability, fairness, and security. We believe that through public engagement with stakeholders across communities impacted by AI in testing, the standards will promote the goal of using AI for good.

We invite public comment on the Duolingo English Test (DET) Responsible AI (RAI) Standards.

Comments can be sent to [englishtest-research@duolingo.com](mailto:englishtest-research@duolingo.com); please specify "Responsible AI Standards" in the email subject line.

## Version 3 Notes

This **third** version of the Duolingo English Test Responsible AI Standards has been revised to include updated references and editorial changes. The first version of these standards was published in April 2023.

This updated version replaces Burstein (2023).

## Legal Disclaimer

This document consists of general principles and goals for ethical AI usage in the context of the Duolingo English Test. While Duolingo makes good faith efforts to align our use of AI on the Duolingo English Test with this document, the assessment, the testing populations, and the technologies are in continuous flux, and at any given point in time there could be discrepancies between these standards and our internal processes. This document should not be relied upon or considered legally binding, and may be revised by Duolingo at any time. This document does not contain legal advice.

## Introduction

Increasingly more sophisticated artificial intelligence (AI) and AI-adjacent methods\* such as computational psychometrics † are embedded in digital learning and assessment platforms. While AI offers many benefits in educational contexts, it also introduces risks. In high-stakes settings, such as standardized assessments, responsible use of AI is essential. The Duolingo English Test (DET) Responsible AI Standards presented here illustrate how to design and apply AI responsibly in a high-stakes assessment context.

Frameworks, guidelines and standards are essential to explicitly address responsible AI in educational assessment. Classical assessment validity (Chapelle et al., 2008; Kane, 1992, 2013) and fairness (Kunnan, 2000) frameworks were developed for paper-and-pencil and first-generation, computer-based assessments. While those frameworks do capture ethical principles (such as validity and reliability, and fairness), they do not directly address the use of technology for assessment. As such, those frameworks are limited with regard to addressing issues related to responsible AI in modern, digital assessments. More recent assessment research leverages those earlier frameworks to discuss AI in terms of validity (such as, Burstein et al. (2025), Huggins–Manley et al. (2022), Williamson et al. (2012), and Xi (2010)). Ethical principles, including validity and reliability, and fairness, are embedded in the AERA & APA & NCME (2014) Standards, and are tied to the standards and their suggested practice. Ethical practices discussed in the those standards—such as those related to test-taker privacy, test security, and test documentation—inform current responsible AI standards for assessment. However, with regard to responsible AI, those standards address only automated essay scoring, highlighting the limited scope of AI for assessment at the time. Since the DET RAI Standards were first published in 2023, other testing organizations have published RAI guidelines (e.g., Johnson (2025)). To our knowledge, there are currently no detailed standards specifically addressing responsible AI across a modern, digital assessment ecosystem (Burstein et al., 2025)

There is a relatively small, but growing body of research focused on responsible AI for educational assessment (such as, Association of Test Publishers (2024), Burstein et al. (2024), Johnson et al. (2022), and LaFlair et al. (2022)). In earlier work, Aiken and Epstein (2000) discuss ethical considerations for AI in education. Earlier, Dignum (2021) proposed a high-level vision for responsible AI for education. Dieterle et al. (2022) and OECD (2023) discuss guidelines and issues associated with AI in assessment. By contrast, there is a substantial literature addressing broader responsible AI **guidelines, regulations, and governance** (e.g., Council of the European Union (2023) and United Nations (2024); **frameworks** (e.g., National Institute of Standards and Technology (2023)(NIST); and, **theory** (e.g., Bentley et al. (2023), Fjeld et al. (2020), Gianni et al. (2022), and Jobin et al. (2019)). Further, given the potential risk of AI use across different sectors, the computer science community has recommended systematic audits of AI applications to prevent potential harm (Mökander & Axente, 2023; Mökander & Floridi, 2021; Raji & Buolamwini, 2019; Raji et al., 2020). This broader literature can be leveraged as we think about responsible AI for education.

As part of professional and institutional responsibility, the DET Responsible AI (RAI) Standards were originally published in 2023—and were the first for an assessment program. They aim to embrace ethical principles, and inform responsible AI practice intended to mitigate risk and minimize ethical debt associated with the use of AI for assessment across the DET’s assessment ecosystem – i.e, test design, measurement, and security (Burstein et al., 2025). Ethical debt is the notion that harmful consequences for people can

---

\*In this document, the term AI refers to both AI and AI-adjacent methods.

† See von Davier (2017) and von Davier et al. (2022).

accumulate when AI tools are developed without attention to responsible AI (Fiesler & Garrett, 2020). An example is facial recognition systems and fairness. In this case, Buolamwini and Gebru (2018) noted harms in that system accuracy varied based on skin color. This was an unanticipated, negative consequence that was not considered when building the facial recognition technology. Ethical debt then occurred, bringing a negative impact to people.

The DET RAI Standards were developed to inform responsible AI practices that mitigate risk and minimize ethical debt to ensure test validity.

## The DET Responsible AI Standards

### The Duolingo English Test

The Duolingo English Test (DET) is a measure of English language proficiency for communication and use in English-medium settings (Naismith et al., 2025). It is primarily used in a high-stakes context for higher education admissions decisions.

In high-stakes testing, outcomes can have a significant impact on test takers' educational opportunities. Therefore, human oversight is essential. The DET uses human-in-the-loop AI practices, in which human experts at critical stages of assessment development and evaluation. For example, human experts play a critical role in reviewing automatically generated content for fairness and bias (Burstein et al., 2024; Church et al., 2025).

The DET uses AI extensively for test design (e.g., automated generation of test items); measurement (e.g., item parameter estimation and automated scoring of written and spoken responses); and, test security (e.g., supporting decision-making in remote proctoring scenarios (Belzak et al., 2024). In addition, the DET's AI affordances contribute to a positive test-taker experience. Examples include: automated scoring of practice tests offering test takers an immediate estimate of their test performance; and, computer adaptive testing, allowing for a shorter test-taking experience.

The DET's use of AI aligns with Duolingo's mission to develop high-quality education and provide universal access. The mission reflects the principles of AI for Good. Specifically, the DET's use of AI supports a valid, reliable, fair, and secure test, while also promoting increased access and a positive test-taker experience.

### The DET Responsible AI Standards Development

Developed as part of professional responsibility, the DET RAI Standards align with four ethical principles: Validity and Reliability, Fairness, Privacy and Security, and Accountability and Transparency. The standards and their goals are designed to support: (a) auditing of AI-powered test processes across test design, measurement and security; (b) validity and reliability studies; and, (c) documentation for theoretical, qualitative, and quantitative research relevant to AI use on the DET, and responsible AI practices. To select the ethical principles for the standards, the DET research team completed five key activities.

**First**, to better understand the set of principles that were applicable to the DET, the team conducted a literature review of ethical principles used for AI (including, Fjeld et al. (2020), Floridi and Cowls (2022), Jobin et al. (2019), and Memarian and Doleck (2023)

**Second**, to examine alignment between domain-agnostic (e.g., NIST, 2023) and assessment-specific principles, the team reviewed assessment-specific standards (AERA & APA & NCME, 2014) and guidelines (including International Test Commission and Association of Test Publishers (2022) and OECD (2023)).

**Third**, the team engaged in multi-stakeholder collaboration. They worked with experts in applied linguistics, computational psychometrics, language assessment, law, machine learning, statistics, and security within Duolingo, and an independent, external RAI expert from computer science. The resulting principles are therefore the outcome of collaboration with a cross-disciplinary set of experts within and outside of the DET.

**Fourth**, after establishing the four ethical principles, the team worked with the external RAI expert collaborator to articulate the rationale and overall goal of each standard, and the more detailed subgoals (i.e., practical implementation of each standard).

**Finally**, the team published the standards as a living document that is freely available, and remains open for public comment to maintain multi-stakeholder conversation as we continue to update the standards.‡

## The Standards

### 1. Validity & Reliability

**Rationale:** Validity and Reliability standards are crucial to ensure that the test is suitable for its intended purpose. Validity standards involve evaluating construct relevance and accuracy (Kane, 1992, 2013), while Reliability standards focus on test score consistency.

**Goals summary:** To specify processes required to build a validity argument, and to evaluate AI used in test item creation, item calibration, and scoring.

#### Goal 1.1

Specify processes required to build a validity argument.

**Processes include theoretical and empirical evaluations that directly inform or address AI used to build a validity argument for test score use.**

1. Develop a description for the test target domain – i.e., English language proficiency – to ensure that test items are aligned with the domain being measured.
2. Evaluate AI scoring system accuracy and fairness, leveraging human expertise. Examples include agreement with human raters, accuracy of system features used for scoring constructed responses, and evaluations of scoring bias.
3. Develop (a) explainable scoring methods, and (b) interpretable AI features used for scoring that have clear alignment with domain constructs.
4. Conduct empirical investigations of item reliability, ensuring reliability for AI-generated items.
5. Evaluate extrapolation through empirical investigations to illustrate relationships between automatically-generated items, test-taker scores, and relevant external measures that suggest proficiency in English skills. Examples of external measures include relationships to other tests, relationships between test-taker's linguistic input and the target domain.

---

‡See Burstein et al. (2024) for discussion and illustrations about how the standards are practically applied as part of the DET's test development, measurement, and security, as well as the standard's impact on the test's validity.

## Goal 1.2

### Evaluate AI used in test item creation, item calibration, and scoring.

1. Identify AI methods for **item creation**, leveraging human expertise to efficiently create valid and reliable test items. An example of human expertise is human review of items from automated item generation.
2. Conduct human evaluations of the quality of items created using AI, such as reviewing outputs from automated item generation.
3. Identify AI methods that can be efficiently used for valid and reliable test **item calibration**.
4. Conduct evaluations that confirm the accuracy of AI for predicting item parameters (such as item difficulty), leveraging human expertise for quality assurance.
5. Identify AI methods that efficiently produce valid and reliable **scores** for test-taker responses.
6. Conduct evaluations to confirm the accuracy of AI for scoring test-taker responses, leveraging human expertise for quality assurance.
7. System changes and re-evaluations are conducted in the case of suboptimal evaluation outcomes for item creation, item calibration, and scoring.

## 2. Fairness

**Rationale:** Fairness standards are required to promote democratization and social justice through increased access, accommodations, and inclusion (Burstein et al., 2025; International Test Commission & Association of Test Publishers, 2022; Naismith et al., 2025) represent test-taker demographics, and avoid algorithms known to contain or generate bias (Belzak, 2022; Johnson et al., 2022).

**Goals Summary:** To specify how the use of AI facilitates test-taker access, accessibility and inclusion; and to specify test-taker demographic representation, and algorithms known to contain or generate bias.

### Goal 2.1

#### Specify how the use of AI facilitates test-taker access, accessibility, accommodations, inclusion, and appeals.

1. Identify AI methods to increase **test-taker access** globally, as part of the DET test-taker experience mission. For example, the DET is available remotely, online, and 24/7. Other access considerations include, but are not limited to, test costs, access to devices, and testing time.
2. Adopt **design principles** in compliance with **accessibility standards** for test takers that offer a generally accessible user interface and allow for test accommodations, due to factors such as low vision or physical limitations.
3. Ensure that AI or AI-adjacent capabilities do not impact design, such that accessibility compliance might be violated.
4. Ensure that test accommodations are not adversely affected by AI.
5. Develop and apply fairness and bias item review principles for **inclusion** that eliminate construct-irrelevant barriers, and ensure that cultural and linguistic factors do not impede accessibility and inclusion for the DET test-taker population.
6. Maintain an **appeals process**, allowing test takers to appeal test results where there is suspicion of cheating.

## Goal 2.2

Specify test-taker and human expert demographic representation, algorithms known to contain or generate bias, and potential human bias.

1. Evaluate and document **demographic representation** in data sets used to build AI. Documentation should describe how representative (inclusive) the data are with regard to DET test takers. For example, the selection of data for AI system development, such as human rater scoring of written responses, should consider the underrepresentation of test-taker language groups which could lead to bias in test-taker outcomes.
2. Evaluate and document known **algorithmic bias** in AI used in DET ecosystem processes (i.e., test security, design, and measurement).
3. Evaluate and document **bias** associated with automatically-generated item content (e.g., fairness and bias review guidelines), and proficiency measurement.
4. Identify and document **human expert demographic representation**.
5. Evaluate and document **bias** associated with human expert ratings of test-taker production tasks.

## 3. Privacy & Security

**Rationale:** The Privacy and Security standards are needed to ensure that we (a) comply with relevant laws and regulations governing the collection and use of test taker data; (b) ensure test taker privacy and (c) to ensure secure test administration. (See Duolingo English Test (2021), Liao, Attali, Lockwood, and von Davier (2022), Liao, Attali, von Davier, and Lockwood (2022), and Wodzak (2021)).

**Goals Summary:** To specify methods to ensure privacy and security associated with data origin, data collection and processing, and data management; to specify how to maintain test-taker privacy, item security, and test-taker security during test administration; and to specify fair and reliable test security proctoring protocols, item pool development and psychometric procedures for test security.

### Goal 3.1

Specify methods to ensure privacy and security associated with data origin, data collection and processing, and data management.

1. Ensure that **data provenance, governance, and management** comply with the Duolingo [privacy policy](#), external privacy policies (where appropriate), and applicable laws such as the European Union's General Data Protection Regulation (GDPR).
2. Define and document **data requirements** with regard to DET intended uses, stakeholders, and the geographic areas where the DET is administered that do not violate privacy terms or security (such as including personally-identifying information without consent).
3. Define, document, and implement methods to ensure that **data provenance** complies with Duolingo [privacy policy](#), and security policies with regard to the origin of the data (e.g., open-access corpora, test taker), how it was obtained (e.g., test-taker consent), and changes applied.
4. Define, document, and implement methods to ensure that **data governance** complies with Duolingo [privacy policy](#) and, where appropriate, external privacy or security policies during data collection and processing, including data cleaning, annotation, enrichment, and aggregation, sharing, and use.

Document how stakeholder data is used, including but not limited to biometric and personally-identifying data (e.g., IDs for security), process (such as keystroke profiles), and product response data and test scores.

5. Define, document, and implement **data management** procedures to ensure compliance with Duolingo and, where appropriate, external privacy or security policies.

### Goal 3.2

Specify how to maintain test-taker privacy, item security, and test-taker security during test administration.

1. Define, document, and implement methods for **test-taker verification** in the context of test onboarding to ensure that test-taker identity can be authenticated.
2. Define, document, and implement methods to ensure that verified **test-taker identity** (i.e., personally-identifiable information) is secure.
3. Define and document **test-taking rules**, such as prohibiting headphones, except in cases of accommodations, and mitigate actual or perceived cheating behaviors to support **test-taker integrity** (see [Test Rules in FAQs](#); [Test Security Rules](#)).
4. Define, document, and implement **test administration** processes that mitigate cheating through use of external resources – i.e., test administration through a desktop application.

### Goal 3.3

Specify fair and reliable test security proctoring protocols, item pool development, and psychometric procedures for test security.

1. Define, document, and implement human-in-the-loop AI **proctoring protocols** that fairly and reliably identify novel and known cheating behaviors.
2. Define and document the algorithm used for proctoring support, proctor training for use of AI, and bias management – i.e., how proctors identify and report perceived AI bias.
3. Define, document, and implement methods (such as, automated item generation) to support scaling of a large, and continuously refreshed test **item pool**. Methods include human expert monitoring protocols for tracking item exposure and test overlap. Larger item pools mitigate the risk that a single test taker, or multiple test takers are likely to see the same test item, or set of ordered test items during repeated sessions (i.e., a test taker registers for and takes the test multiple times).
4. Define, document, and implement security protocols to prevent item breach from external attackers.
5. Define, document, and implement **psychometric procedures**, such as test-retest reliability that can track anomalies in test-taker performance that can reveal test-taker cheating behavior.

## 4. Accountability & Transparency

**Rationale:** To gain trust from stakeholders, it is essential that the DET have Accountability & Transparency standards for proper governance of AI used on the test. Through documentation and explanations, we are holding ourselves accountable.

**Goals Summary:** To assess how AI processes impact stakeholders; to document AI used for building the validity argument, test item creation, test item calibration, and scoring; to document processes for human-in-the-loop interactions with AI; document human expert qualifications required for human-in-the-loop activities that support AI for the DET; to disseminate research about use of AI to various stakeholder communities; and to publish information about how AI is used on the DET, and usage of test-taker data.

### Goal 4.1

Assess how AI processes impact stakeholders.

**Stakeholders include test takers and organizations who use the DET, such as universities.**

1. Document how **ML algorithms** (a) are used on the test and DET support resources§, (b) are used for test design, measurement and security, and (c) impact stakeholders (i.e., high-stakes decisions based on DET outcomes).
2. Document **unintended risks** resulting from AI, such as biased scores or construct-irrelevant variance related to design, such as a test-taker's unfamiliarity with "drag-and-drop" options. Unintended risks may result in negative consequences for stakeholders, such as unfair or inappropriate admissions decisions.
3. Document **external factors** that result in a need to modify AI. Examples might include: a) institutional policy changes, such as, new language proficiency requirements require new item types; b) modifications in institutional use cases, such as, different rating practices; and, c) demographic changes, such as, increases in particular language groups that might impact differential item functioning (DIF).
4. Document how AI is used for DET **stakeholder support**. Examples of support include test readiness resources for test takers, and score interpretation guidance for organizations.

### Goal 4.2

Document AI used for building the validity argument, test item creation, test item calibration, and scoring.

1. Document theoretical claims, and empirical studies to support the **DET validity argument** – i.e., evidence that the test is suitable for its intended purpose.
2. For **item creation**, document rationales for, and descriptions of item generation methods, including data and algorithms, human expert processes (such as fairness and bias review), and evaluation methods that provide a clear explanation of system performance metrics, and fairness evaluations, such as mitigating bias with regard to item content generation.

---

§DET support mechanisms include: test readiness resources for test takers, score interpretation guidance for organizations, and how associated AI contributes to impact (e.g., practice test automated scoring is out of sync with certified test scoring).

3. For **test item calibration**, document rationales for, and descriptions of AI for predicting item parameters, such as item difficulty.
4. For **test scoring**, document (a) rationales, descriptions, and alignment between item subconstructs and computationally-derived features used for scoring. This is relevant for constructed-response tasks involving spoken and written responses; and, (b) rationales for, and descriptions of measurement methods used to generate DET scores.

### Goal 4.3

#### Document processes for human-in-the-loop interactions with AI.

**These processes include, system development, evaluation, and data preparation and oversight.**

1. Document **sustainable processes requiring 100% human expertise** (e.g., data annotation), or human expert supervision, such as fairness & bias review, monitor language model hallucination for automatically-generated item content (e.g., Ji et al. (2023)). Sustainability is measured in terms of time, costs, and required resources (e.g., need for third parties, such as hiring contractors to support human fairness and bias review, and annotation supporting AI development).
2. Document **qualifications of human experts**, such as software engineers, AI researchers, and assessment researchers, who are responsible for supervision during system development and evaluation, and piloting and operational deployment phases.
3. Document **demographic composition of human experts** who participate in annotation (e.g., scoring written responses) and reviewing tasks (e.g., fairness and bias review).
4. Document **supports to help individuals understand and carry out their responsibilities** in relation to interacting with AI. Supports may include system UX, alert and reporting functions, and rubrics.
5. Document **content to help individuals' understand how AI is applied** on the DET. Document (1) AI systems' intended uses, such as text generation, scoring, test security, (2) empirical evaluations and interpretations of AI system behavior, and 3) acknowledgement of potential automation bias – specifically, favoring system outputs.

### Goal 4.4

#### Document human expert qualifications required for human-in-the-loop activities that support AI for the DET.

1. Document qualifications criteria for **human experts** that are specific to activities that support human-in-the-loop AI. Such documentation may be applied for hiring practices, such as job descriptions.

### Goal 4.5

#### Disseminate research about use of AI to various stakeholder communities.

1. Document research to illustrate how the **DET validity argument** was constructed with attention to AI.
2. Disseminate research about **theoretical, and quantitative, qualitative and mixed-methods research** through peer-reviewed, external publications and presentations See [DET Research page](#).
3. Prioritize **peer-reviewed, open-access venues and publications**, and provide **public access** to peer-reviewed presentations.

4. Document, disseminate, and update **white papers about internal DET research** through the DET website, such as [DET Research page](#).
5. Disseminate **external media publications**, such as blogs, that provide clear and plain language explanations of complex DET processes, such as [test security](#). These publications are intended to render complex concepts transparent and accessible for the broader stakeholder community.

## Goal 4.6

### Publish information about how AI is used on the DET, and usage of test-taker data

1. Display on the DET website documentation about **how AI is used on the DET** so that it is understandable to stakeholders. For example, include content about automated test item creation and scoring, and test security on the DET website. Information should be publicly-available, such as in [FAQs](#) on the DET website, or a section devoted to explanations about how AI are used on the test.
2. Display documentation on the DET website. about **how stakeholder data is used**, including but not limited to biometric and personally-identifying data, such as IDs for security; process, such as keystroke profiles; and, product response data and test scores.

## Where Theory Meets Practice

As AI becomes increasingly integrated into educational assessments, it is critical for test developers and assessment researchers to focus not only on creating RAI standards, but also on evaluating the maturity of those standards (Dotan et al., [2024](#)). A key part of that evaluation would involve developing and implementing concrete practices that reflect and uphold the values that the standards represent.

While this document centers on the DET RAI Standards, it should be noted that DET maintains an internal database of RAI practices and related documentation associated with the standards. Readers can explore how the test's practices align with DET RAI standards in Burstein et al. ([2024](#)) and Church et al. ([2025](#)).

## References

- AERA & APA & NCME. (2014). Standards for educational & psychological testing. American Educational Research Association.
- Aiken, R. M., & Epstein, R. G. (2000). Ethical guidelines for ai in education: Starting a conversation. *International Journal of Artificial Intelligence in Education*, 11, 163–176.
- Association of Test Publishers. (2024). Creating responsible and ethical ai policies for assessment organizations [White Paper].
- Belzak, W. (2022). The multidimensionality of measurement bias in high-stakes testing: Using machine learning to evaluate complex sources of differential item functioning. *Educational Measurement: Issues and Practice*, 1–10.
- Belzak, W., Lockwood, J. R., & Attali, Y. (2024). Measuring variability in proctor decision making on high-stakes assessments: Improving test security in the digital age. *Educational Measurement: Issues and Practice*.
- Bentley, C., Aicardi, C., Poveda, S., Magela Cunha, L., Kohan Marzagao, D., Glover, R., Rigley, E., Walker, S., Compton, M., & Acar, O. (2023). A framework for responsible ai education.
- Buolamwini, J., & Gebru, T. (2018). Gender shades: Intersectional accuracy disparities in commercial gender classification. *Conference on Fairness, Accountability and Transparency*, 77–91.
- Burstein, J., LaFlair, G., & von Davier, A. (2025). A theoretical assessment ecosystem for a digital-first assessment—the duolingo english test (Duolingo Research Report No. DRR-25-03). Duolingo.
- Burstein, J., LaFlair, G. T., Yancey, K., von Davier, A. A., & Dotan, R. (2024). Responsible ai for test equity and quality: The duolingo english test as a case study. *arXiv preprint arXiv:2409.07476*.
- Burstein, J. (2023). Duolingo english test responsible ai standards. Updated March, 29, 2024.
- Chapelle, C., Enright, M., & Jamieson, J. (2008). Building a validity argument for the test of english as a foreign language. Routledge.
- Church, J., Park, Y., & Burstein, J. (2025). Guidelines for fair test content: The duolingo english test example (tech. rep.). The Duolingo English Test.
- Council of the European Union. (2023). Artificial intelligence act: Council and parliament strike a deal on the first rules for ai in the world [Press release]. <https://www.consilium.europa.eu/en/press/press-releases/2023/12/09/artificial-intelligence-act-council-and-parliament-strike-a-deal-on-the-first-worldwide-rules-for-ai/>
- Dieterle, E., Dede, C., & Walker, M. (2022). The cyclical ethical effects of using artificial intelligence in education. *AI & Society*, 1–11.
- Dignum, V. (2021). The role and challenges of education for responsible ai. *London Review of Education*, 19(1), 1–11.
- Dotan, R., Blii-Hamelin, B., Madhavan, R., Matthews, J., Scarpino, J., & Anderson, C. (2024). A flexible maturity model for ai governance based on the nist ai risk management framework (Technical Report). IEEE. <https://ieeusa.org/product/a-flexible-maturity-model-for-ai-governance>
- Duolingo English Test. (2021). Duolingo english test: Security, proctoring, and accommodations (tech. rep.). Duolingo.
- Fiesler, C., & Garrett, N. (2020, September). Ethical tech starts with addressing ethical debt [Wired Ideas]. <https://www.wired.com/story/opinion-ethical-tech-starts-with-addressing-ethical-debt/>

- Fjeld, J., Achten, N., Hilligoss, H., Nagy, A., & Srikumar, M. (2020). Principled artificial intelligence: Mapping consensus in ethical and rights-based approaches to principles for ai (tech. rep.). Berkman Klein Center for Internet & Society.
- Floridi, L., & Cows, J. (2022). A unified framework of five principles for ai in society. In S. Carta (Ed.), *Machine learning and the city: Applications in architecture and urban design* (pp. 535–545). John Wiley & Sons Ltd. <https://doi.org/10.1002/9781119815075.ch45>
- Gianni, R., Lehtinen, S., & Nieminen, M. (2022). Governance of responsible ai: From ethical guidelines to cooperative policies. *Frontiers in Computer Science*, 4.
- Huggins–Manley, A. C., Booth, B. M., & D'Mello, S. K. (2022). Toward argument-based fairness with an application to ai-enhanced educational assessments. *Journal of Educational Measurement*, 59(3), 362–388.
- International Test Commission & Association of Test Publishers. (2022). Guidelines for technology-based assessment (tech. rep.). International Test Commission.
- Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y. J., Madotto, A., & Fung, P. (2023). Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12), 1–38.
- Jobin, A., Ienca, M., & Vayena, E. (2019). The global landscape of ai ethics guidelines. *Nature Machine Intelligence*, 1, 389–399. <https://doi.org/10.1038/s42256-019-0088-2>
- Johnson, M. S., Liu, X., & McCaffrey, D. F. (2022). Psychometric methods to evaluate measurement and algorithmic bias in automated scoring. *Journal of Educational Measurement*, 59(3), 338–361.
- Johnson, M. S. (2025). Responsible ai for measurement and learning: Principles and practices (tech. rep.). ETS.(Research Report No. RR-25-03). [https://www.ets.org/Media/Research/pdf ...](https://www.ets.org/Media/Research/pdf...)
- Kane, M. T. (1992). An argument-based approach to validity. *Psychological Bulletin*, 112(3), 527.
- Kane, M. T. (2013). Validating the interpretations and uses of test scores. *Journal of Educational Measurement*, 50(1), 1–73.
- Kunnan, A. J. (2000). Fairness and justice for all. In A. J. Kunnan (Ed.), *Fairness and validation in language assessment* (pp. 1–14). Cambridge University Press.
- LaFlair, G. T., Langenfeld, T., Baig, B., Horie, A. K., Attali, Y., & von Davier, A. A. (2022). Digital-first assessments: A security framework. *Journal of Computer Assisted Learning*, 38(4), 1077–1086.
- Liao, M., Attali, Y., Lockwood, J. R., & von Davier, A. A. (2022). Maintaining and monitoring quality of a continuously administered digital assessment. *Frontiers in Education*, 7.
- Liao, M., Attali, Y., von Davier, A. A., & Lockwood, J. R. (2022). Quality assurance in digital-first assessments [Virtual, 2021]. In *Quantitative psychology: The 86th annual meeting of the psychometric society* (pp. 265–276). Springer.
- Memarian, B., & Doleck, T. (2023). Fairness, accountability, transparency, and ethics (fate) in artificial intelligence (ai), and higher education: A systematic review. *Computers and Education: Artificial Intelligence*, 5. <https://doi.org/10.1016/j.caeai.2023.100152>
- Mökander, J., & Axente, M. (2023). Ethics-based auditing of automated decision-making systems: Intervention points and policy implications. *AI & Society*, 38, 153–171. <https://doi.org/10.1007/s00146-021-01286-x>
- Mökander, J., & Floridi, L. (2021). Ethics-based auditing to develop trustworthy ai. *Minds and Machines*, 31(2), 323–327. <https://doi.org/10.1007/s11023-021-09557-8>
- Naismith, B., Cardwell, R., LaFlair, G., Nydick, S., & Kostromitina, M. (2025). Duolingo English Test: Technical manual (Duolingo Research Report). Duolingo. <https://go.duolingo.com/dettechnicalmanual>

- National Institute of Standards and Technology. (2023). Artificial intelligence risk management framework (ai rmf 1.0) (tech. rep.). U.S. Department of Commerce. <https://doi.org/10.6028/NIST.AI.100-1>
- OECD. (2023). Advancing accountability in ai: Governing and managing risks throughout the lifecycle for trustworthy ai (tech. rep. No. 349). Paris.
- Raji, I. D., & Buolamwini, J. (2019). Actionable auditing: Investigating the impact of publicly naming biased performance results of commercial ai products. Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society, 429–435.
- Raji, I. D., Smart, A., White, R. N., Mitchell, M., Gebru, T., Hutchinson, B., Smith-Loud, J., Theron, D., & Barnes, P. (2020). Closing the ai accountability gap: Defining an end-to-end framework for internal algorithmic auditing. Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (FAT\*20), 33–44. <https://doi.org/10.1145/3351095.3372873>
- United Nations. (2024, September). Governing ai for humanity. [https://www.un.org/sites/un2.un.org/files/governing\\_ai\\_for\\_humanity\\_final\\_report\\_en.pdf](https://www.un.org/sites/un2.un.org/files/governing_ai_for_humanity_final_report_en.pdf)
- von Davier, A. A. (2017). Computational psychometrics in support of collaborative educational assessments. *Journal of Educational Measurement*, 54(1), 3–11.
- von Davier, A. A., Mislevy, R. J., & Hao, J. (2022). Computational psychometrics: New methodologies for a new generation of digital learning and assessment: With examples in r and python. Springer Nature.
- Williamson, D. M., Xi, X., & Breyer, F. J. (2012). A framework for evaluation and use of automated scoring. *Educational Measurement: Issues and Practice*, 31(1), 2–13.
- Wodzak, S. (2021, September). What if tests were delightful? <https://blog.duolingo.com/what-if-tests-were-delightful/>
- Xi, X. (2010). Automated scoring and feedback systems: Where are we and where are we heading? *Language Testing*, 27(3), 291–300.