

# Ethical Thinking: Integration and Measurement in an AI Curriculum for Middle-High School Students

Glenda S. STUMP<sup>1</sup>, Shinyi KANG<sup>2</sup>, John MASLA<sup>3</sup>,  
Christina A. BOSCH<sup>3</sup>, Hal ABELSON<sup>4</sup>, Eric KLOPFER<sup>5</sup>,  
Cynthia BREAZEAL<sup>6</sup>

<sup>1</sup>*StumpWorks LLC, Melrose, Massachusetts, USA*

<sup>2</sup>*Department of Communication and Media Research, University of Zurich, Switzerland*

<sup>3</sup>*Scheller Teacher Education Program, Massachusetts Institute of Technology  
Cambridge, Massachusetts, USA*

<sup>4</sup>*Department of Electrical Engineering and Computer Science,  
Massachusetts Institute of Technology, Cambridge, Massachusetts, USA*

<sup>5</sup>*Department of Comparative Media Studies/Writing, Massachusetts Institute of Technology  
Cambridge, Massachusetts, USA*

<sup>6</sup>*Office of the Provost, Massachusetts Institute of Technology, Cambridge, Massachusetts, USA  
e-mail: gsstump1016@gmail.com, s.kang@ikmz.uzh.ch, j\_masla@education.mit.edu,  
cabosch@education.mit.edu, hal@mit.edu, klopfere@mit.edu, cynthiab@media.mit.edu*

Received: February 2025

**Abstract.** Ethical thinking and reasoning is considered a core component of artificial intelligence (AI) literacy. However, there is a lack of strategies and assessments to promote and measure students' ethical thinking in AI, particularly in K-12 education. In this paper, we discuss how RAICA, a project-based AI literacy curriculum integrates ethics into its framework and instructional resources. We employ a mixed-methods convergent design to obtain different yet complementary data on ethical thinking as an outcome mediated by diverse RAICA materials. Data analysis revealed that teachers utilized a variety of instructional strategies to foster students' ethical thinking and that students actively engaged in ethical thinking activities, resulting in a developing understanding of stakeholders and potential benefits/harms of AI. Our work makes a key contribution to AI education by providing empirical evidence to support mechanisms for integration and assessment of ethical thinking within AI literacy curricula.

**Keywords:** ethical thinking, AI curriculum, responsible design, middle school, high school.

## 1. Introduction

This paper presents an empirical study of K-12 artificial intelligence (AI) literacy education by focusing on one critical aspect: that of ethical thinking and reasoning. In current discussions around the increasing power and impact of AI, the call for ethics and responsibility in the design, deployment, and use of AI tools and technologies is undisputed (Barabas *et al.*, 2018; Green, 2019; Hagendorff, 2020; Torresen, 2018; Trotta *et al.*, 2023). The question is not whether ethics should be addressed, but how to make it a core part of the design and decision-making process before harm occurs.

This concern is particularly germane to AI literacy education at the K-12 level. As AI users and future designers, students should not only learn how AI works or develop their own projects with it but also critically engage with its ethical and societal implications (Ali *et al.*, 2019; DiPaola *et al.*, 2020; Holmes *et al.*, 2022; Long and Magerko, 2020; Ng *et al.*, 2021a; UNESCO, 2022; Williams *et al.*, 2023). Although existing AI literacy frameworks and literature emphasize this need, gaps remain – particularly in the availability of strategies and assessments that promote and measure students' ethical thinking related to AI (Williams, 2021). Our work proposes to fill this gap.

Our Responsible AI for Computational Action (RAICA) curriculum is designed for computer science classrooms, filling a need for computer science curricula that integrate ethical reasoning with technical skill building. Computer science curricula often prioritize the teaching of technical skills at the expense of emphasizing ethical reasoning, as ethics are often considered to be the domain of disciplines such as philosophy or social science and not of primary importance to computer science (Skirpan *et al.*, 2018). In RAICA, middle and high school students engage with AI through project-based learning, where ethical thinking is emphasized in both content, practice, and assessment.

During RAICA development, we continually refine and explore different ways to promote and assess students' ethical reasoning, adapting our approach based on emerging insights gained through teachers' real-world, school-based implementation of the curriculum. In this paper, we outline the integration of ethical thinking into the RAICA curriculum as part of responsible design, and we review findings emerging from our varied approaches to measuring students' thinking in this regard. Our research questions are:

- RQ1:** How do students engage with curricular activities designed to promote ethical thinking?
- RQ2:** How do students describe stakeholders associated with AI tools and applications in formative assessments?
- RQ3:** Does engagement with the RAICA curriculum enable students to identify stakeholders in AI tools and applications as evidenced through a pre-post measure?
- RQ4:** Does engagement with the RAICA curriculum enable students to identify potential benefits and harms of AI tools and applications as evidenced through a pre-post measure?
- RQ5:** What assessment approach can provide reliable and valid information regarding students' ethical thinking skills?

## 2. Background

### 2.1. Ethical Thinking in Adolescence

Schrier and colleagues define ethical thinking and reasoning as “the ability to analyze, assess, and reflect on our decisions and actions, and understand the consequences and complexities of social issues” (Schrier *et al.*, 2010, p. 255). Of the ethical frameworks commonly applied to ethical situations, Card and Smith (2020) argue that consequentialism is particularly relevant when considering fairness in machine learning, and by extension to an AI context. Consequentialism theories place outcomes as the priority and lead to decisions or actions that would produce favorable consequences according to an impartial view as to what is best (Card and Smith, 2020).

One factor that may underlie ethical thinking or reasoning about AI is an individual’s circle of moral concern, defined as the range of individuals one considers when evaluating the potential impact of their actions (Singer, 2011). Bloom (2004) extends this idea to child development, stating that young children begin with a very small circle; they are attached primarily to their family and those around them. Their circle, thus their moral perspective, expands as they interact with other individuals, engage in cooperative situations for mutual gain, hear stories that encourage perspective-taking, and are exposed to moral beliefs of previous generations. A crucial aspect of ethical thinking about AI is broadening one’s perspective, acknowledging it as a socio-technical system rather than merely a technical system (DiPaola *et al.*, 2020; Zhang *et al.*, 2023). AI tools and applications are easily utilized and dispersed electronically; those who may be affected by them, such as data workers in low-income countries, are unknown to developers or other users and thus invisible (see Grohmann and Araújo, 2021 as an example). This possibility of widespread impact on unknown others reinforces the need for individuals to develop the capacity for identification of ethical situations that may arise from AI technology and consideration of AI’s potential benefits and harms to those who may be affected. When in a position to create new technologies, they must be cognizant of design choices that reflect or conflict with their values and broader societal values.

An ideal time to guide students in honing their ethical thinking and reasoning skills is in middle or high school, as they are in the process of exploring and solidifying their values during adolescence (Mawasi *et al.*, 2022; Kuhn and Park, 2005; Sandoval, 2014). From a motivational perspective, adolescents may show increased interest and engagement in exploring AI ethical issues due to its popularization in movies and other forms of media (Mawasi *et al.*, 2022). AI education programs that include an ethics component can improve adolescents’ ability to identify and empathize with stakeholders, consider the far-reaching impacts of their design, and make design choices that reflect design values (DiPaola *et al.*, 2020, Zhang *et al.*, 2023). Presenting students with complex ethical problems exposes them to new, innovative ways of thinking that can have implications for interactions they may have with AI in the future. Kohlberg

and Hersh (1977) claim that students' general moral development can be aided by challenging them with complex issues, such as community or school-related problems to be solved, rather than ones that merely require mechanical application of rules. This would suggest a link between considering ethical complexities of AI tools and applications and moral growth.

## 2.2. *Embedding Ethical Thinking into Curriculum*

The United Nations Education and Scientific Organization (UNESCO) recommends that AI ethics curriculum is needed for all levels of learners and that curricula “promote cross-collaboration between AI technical skills education and humanistic, ethical, and social aspects of AI education” (UNESCO, 2022, p. 35). AI short stories (Forsyth *et al.*, 2021), interactive case studies, problem solving, and role-playing scenarios (Schuitema *et al.*, 2008) have been successful strategies utilized to engage students in ethical thinking. Interactive case studies and scenarios are also suggested pedagogical methods in UNESCO's AI Competency Framework for Students (Miao *et al.*, 2024). Contextualized, authentic tasks are effective methods for promoting student engagement in ethical thinking (Forsyth *et al.*, 2021; Schuitema *et al.*, 2008), whereas practice with ethical reasoning during collaborative class activities helps students to see it as important to the design process (Narayanan and Vallor, 2014). When placed in realistic situations that invite identification, caring, and concern, students can practice empathy, which is a critical aspect of ethical thinking (Schrier *et al.*, 2010). This aligns with Gilligan's theory of moral reasoning which places relationships and care as key components of ethical response (Gilligan, 2018).

Other scholars argue that traditional modes of ethics education such as case studies or ethical dilemma discussions do not lead to student understanding of the complex ethical and technical issues inherent in AI, but rather lead to students' shallow understanding of those concepts (Li *et al.*, 2024). They suggest a constructionist approach that includes tinkering with AI tools, building with AI projects, and engaging in critical reflection to foster deeper AI ethical understanding (Lin and Dai, 2025).

On a more granular level, learning science research provides significant support for the use of design strategies that improve student learning, specifically spaced and interleaved practice (Roediger and Pyc, 2012; Taylor and Rohrer, 2010). Spaced learning is achieved by returning to already-learned content at intervals of time. This practice provides practice at retrieval, shown to potentiate learning as well as long term retention (Roediger and Butler, 2011). Interleaving of practice, or studying relevant but different content across time, helps students to discern between types of problems or situations and decide the type of strategy they should use. Teachers can implement spaced practice and interleaved examples by reviewing topics covered in previous lessons and giving assignments that require use of content learned in previous lessons (Roediger and Pyc, 2012).

### 2.3. The RAICA Curriculum

The RAICA curriculum integrates ethical thinking into a process of responsible design, a novel adaptation of the design thinking process specifically geared toward computational action with AI (Wharton *et al.*, 2024). Early testing showed that students struggled with design thinking's empathy and problem definition phases; they engaged more during hands-on tool use than during explanations. Unlike empathy-driven design approaches where students identify and solve others' problems, RAICA's responsible design for computational action centers on effective tool engagement for solving problems students themselves care about. As a result, our model begins with playing (with the tool), followed by brainstorming (the project), then defining (project goals), and pausing/planning (to clarify stakeholders and impact) before prototyping and trying it out.

Another distinguishing aspect of RAICA compared to other curricula that teach about AI and ethical reasoning (e.g., DiPaola *et al.*, 2020; Lin and Dai, 2025; Williams *et al.*, 2023), is its ongoing development and refinement informed by classroom-based, teacher-implemented instruction. Few AI literacy interventions are taught by teachers as a curricular requirement within schools (as opposed to out-of-school or electives), through project-based learning, and as consecutive course content (over a dozen or more hours, rather than a workshop less than three hours) (Yue *et al.*, 2022).

Aligned with UNESCO's recommendation (2022), the RAICA curriculum weaves ethical thinking activities together with use of AI tools and applications in a constructionist approach to learning by making. The curriculum identifies three dimensions of ethical thinking: 1) consideration of stakeholders, 2) recognition of intended and unintended impacts of a proposed design, and 3) adoption of design values appropriate to the context.

RAICA module activities engage students in ethical thinking while learning pivotal AI concepts throughout each module and engages them in this thinking process at multiple points, thereby providing spaced and interleaved learning. For example, in a module entitled *Social Robots*, students participate in discussions related to potential benefits and harms as they learn about facial recognition technology; discussions about stakeholders and design values as they learn about affective computing; and discussions about stakeholders, potential benefits/harms, as well as design values as they apply the concepts of behavior, communication, and embodiment through building their own social robot projects. In three different contexts at repeated intervals, students consider one or more dimensions of ethical thought, thus providing authentic opportunities for them to consider the ethical implications of new information they are learning. Teachers are not guided to endorse a particular ethical stance for any of the situations introduced in the modules; instead, they focus on student understanding of the three dimensions of ethical thinking.

During the project build phase of each module, students use a structured graphic organizer to support discussion and diagramming as they brainstorm project ideas, clarify how their projects will be built, identify stakeholders, consider potential social impacts in terms of benefits/harms, and articulate design values. Design values are defined as

“the things we care about that guide us when making our project.” Students choose from a list of key characteristics—fairness, accessibility, security, safety, inclusivity, sustainability, transparency, justice, privacy, and accountability—and strive to incorporate one or more of them into their project.

RAICA curriculum materials provide teachers with means to contextualize student learning. The facilitators’ guide encourages teachers to align instruction with students’ needs by providing alternate approaches to activities and examples of ways to discuss AI topics. Teacher agency, or their sense of autonomy and voice in their professional practice is essential in this process. Molla and Nolan (2020) argue that teacher agency is relational; when others recognize it, teachers’ conscious and robust use of agency is positively influenced. At multiple points in the curriculum, RAICA’s Facilitators Guide supports teacher adaptation of materials, thereby supporting teacher agency for delivery of AI content that is tailored to students’ local context and culture.

In the current study, the RAICA content consisted of one or more of the following modules: *Picture This*, *Social Robots*, and *AppInventor+LLMs (Large Language Models)*. Ethical thinking activities in the modules are woven throughout the modules’ AI content instruction (see Table 1 for module descriptions).

The content of responsible design activities in the whole-class lessons, the individual exit tickets, and the graphic organizers align with specific supports in the facilitation guide and student-facing slides. The facilitation guide also links to additional resources designed to enhance teachers’ knowledge and pedagogy related to AI concepts, tools, and ethical concerns (e.g., What are neural networks? What is data bias? How can we talk about data bias in the classroom? and Using Teachable Machine).

Table 1  
RAICA modules implemented for this study

Module Name	Picture This	Social Robots	App Inventor+LLMs
Key Question	How can AI see?	How can humans interact and work creatively with AI?	How does AI make sense of human language?
AI Content	image classification, computer perception, neural networks	natural social interaction, affective computing, face detection	generative AI, large language model, tokens, attention mechanisms
Ethical Considerations	data bias, pros/cons of neural network layers	benefits/harms, rules, embodiment, personality	transparency, bias, environmental impact, human impact
Example Ethical Thinking Activity	Whole-class conversation about pros/cons of computer processes to interpret images compared with humans.	Data privacy scenarios. Students express their ethical reasoning first by physically positioning themselves along a “WORTH IT”/ “NOT WORTH IT” continuum in the classroom and then articulating their thinking to a partner.	Debate-style activity assigns small groups to “pro” and “anti-LLM” positions. Using real-world example cases of LLM use, students construct arguments for their assigned stance, later voting according to their real opinions.
Classroom Assessment	Checks for Understanding, Exit tickets, Completion of structured graphic organizer		

The materials introduce and reinforce core responsible design topics (e.g., stakeholders, design values, impact) as learning objectives. Together, the facilitation guide and additional resources are designed to build teachers' technological, pedagogical, and content knowledge (TPACK; Koehler *et al.*, 2014) related to AI. They are also intended to facilitate discussions that align directly with the scenario-based measure of ethical thinking.

#### 2.4. Measures of AI Literacy and Ethical Thinking

The development of assessments for AI literacy education has gained increasing attention, spanning attitudinal and conceptual dimensions (Almatrafi *et al.*, 2024; Lintner, 2024; Miao *et al.*, 2024; Zhang *et al.*, 2025). Efforts to evaluate K-12 students' ethical reasoning abilities have also emerged, taking the form of ethical matrix activity worksheets, interviews, and open-ended questionnaires (DiPaola *et al.*, 2020; Ali *et al.*, 2019). One commonly used method, the ethical matrix activity, adopts a consequentialist approach, where participants identify stakeholders and evaluate potential impacts from multiple value-based perspectives (O'Neil and Gunn, 2020; Mepham, 2000). However, currently-used methods face limitations, as none are scalable measures, nor have they been examined for content validity or psychometric robustness.

Beyond testing students' conceptual understanding or technical skills, a major challenge in assessing ethical thinking is determining which ethical principles and models should be prioritized, given that many are microsocial, or context-dependent (Goldsmith *et al.*, 2020; Donaldson and Dunfee, 1994). Ethical reasoning in AI-related contexts further requires foundational skills beyond ethical analysis itself, such as an understanding of technology and familiarity with discipline-specific vocabulary, complicating the development of adequate and reliable measures of students' ethical reasoning.

Existing ethical reasoning measures provide some insights but remain limited in scope. The Defining Issues Test (DIT; Rest, 1979) and its revised version, DIT-2 (Rest *et al.*, 1999), have been widely used to measure short-term ethical development through scenario-based items. Technology-specific assessments such as the Engineering Ethical Reasoning Instrument (Zhu *et al.*, 2014) and the Engineering and Science Issues Test (Borenstein *et al.*, 2010) have been adapted from the DIT-2. However, these assessments are primarily designed for non-K-12 contexts and have likewise been criticized for favoring a singular ethical model over others (Flanagan, 2009). Additionally, they rely on ranking predefined ethical issues without incorporating open-ended or constructed response elements (Goldsmith *et al.*, 2020).

Our work aims to address the above-mentioned gaps by working towards a scalable measure that captures students' ethical reasoning in AI-related contexts. We also seek to conduct validity testing of the new assessment as a low-stakes research instrument to ensure that it is an appropriate, psychometrically robust indicator of students' ethical thinking.

### 3. Method

This descriptive study draws on the data collected as part of design-based research (DBR) and design-based implementation research (DBIR) with RAICA, which seeks broadly to engage teachers and learners across contexts in computational action with AI. DBIR builds on DBR by seeking systems change (Fishman and Penuel, 2018) and has guided the overarching approach on RAICA, with room for smaller design experiments (e.g., the current study; Cobb *et al.*, 2003) within the pursuit of scalable innovations (Svihla, 2014; Fishman *et al.*, 2004).

#### 3.1. *Study Design*

A learning sciences methodology that values equity alongside excellence, DBR is the process of developing theory-based educational interventions, iterating those designs through dialectical collaborations with real-world stakeholders (e.g., practitioners, administrators, learners), and using ensuing results to contribute to design principles and learning theories (Fishman *et al.*, 2013; Prediger, 2024; Donaldson *et al.*, 2024). With a pragmatic approach to DBR as the overarching methodology, we employ a mixed-methods convergent design in the current study to obtain different yet complementary data to gain a more complete picture of students' ethical thinking as an outcome mediated by diverse RAICA materials (Creswell and Plano Clarke, 2018). After the initial design phase, we collected qualitative data on teacher and student experiences with the modules through four distinct methods (described below). We conducted qualitative analyses on each measure, transformed scenario codes into quantitative data to conduct statistical analyses, and then integrated the respective findings under our research questions to generate inferences and draw conclusions to inform subsequent iterations of the curriculum.

#### 3.2. *Participants and Sampling*

Five learning environments were recruited via convenience sampling for the current study – four schools in the northeastern U.S. and one out-of-school educational organization in southeastern Africa. We also worked with one additional out-of-school group, with an exclusive focus on middle school students with vision impairments in the U.S. Only one student's data was included from that site, in the scenario-based measure.

Three of the participating schools were private. Private School A was a day and boarding school outside of a major city. Private School B was a parochial school in a suburban setting. The schools served majority white student populations. Class size averaged ranged from 9 to 17 students, with ages between 11 and 13 years. Male and female students comprised approximately equal proportions of the total enrollment, though distributions varied by class. Slightly more than 10% of the enrolled students in

each class had a specialized education plan. Private School C was a day and boarding school in the suburbs that exclusively enrolled students with language-based disabilities. The school served a majority white student population. The three female teacher participants, one from each school, had significant experience in the field, participating in leadership and professional development opportunities in computer science education. In addition to extensive teaching experience, they had a high level of domain knowledge – for example, one was previously a professional software engineer, and at least one had a post-graduate degree.

One participating school was public. Public School D was a middle school, where 45% of the students were considered “high needs,” 28% were considered students with disabilities, and 20% were considered low-income.

The fifth participating site was an out-of-school technology education organization (Organization E). Based in a large refugee camp in southeastern Africa, youth participants were both male and female, ranged widely in ages, and all spoke three or more languages. Instruction at this site occurred using English materials that the teachers translated orally or by using software tools.

Within the participating sites, computer science class teachers and students aged 11 and older were eligible to participate; students in the current study ranged from 11–15 years. We obtained administrator permission to conduct the study, and adult participants (teachers) completed informed consent. All students in the classes implementing RAICA participated in the curriculum as part of the regularly occurring instruction, but only those who confirmed assent (through the online pre-post questionnaire) and had parental consent were considered participants in the research. All participants had limited experience with AI and AI tools. As research participants, teachers participated in optional interviews. They were compensated for this activity as it was conducted outside of their normal instructional duties.

### 3.3. *Data Collection*

In this section we discuss the measures, procedures, and contexts/participants associated with each data collection activity.

#### 3.3.1. *Classroom Observation Tool*

The structured classroom observation protocol was adapted from Lee, Anderson and Hsiao’s (2019) work and utilized to record class engagement and teacher implementation. The protocol consisted of four parts. Part A asked for basic information like date, school, class size, and part of the curriculum being taught. Part B provided space to record teacher and student behaviors and observers’ inferences and impressions. Part C included questions for a debrief interview with the teacher. Part D, to be completed after the observation, included broader questions for observer reflection, such as “How would you summarize the teacher’s facilitation style in supporting students’ collaborative learning?”

A total of three trained observers employed the protocol at Private Schools A and B during single-day site visits in the fall and spring semesters of 2024. Each classroom observation had two trained observers. Each observer recorded teacher and student behaviors at regular five-minute intervals during the 50-minute lesson. During whole group instruction, the recorded observations were typically of the same phenomena. During small group work time, observers' observations differed because they were checking in with separate groups. During the debrief interview, observers reported on teachers' perspectives on the session, including things that went well and things they would do differently. Across five different class periods, 57 middle school students were observed during the RAICA module's project build phase.

### 3.3.2. Teacher Interview Tool

The teacher semi-structured interview protocol was designed to gather information regarding teachers' impression of module implementation including the responsible design activities that included ethical thinking. The interview questions asked teachers about their growth in Technological, Pedagogical Content Knowledge (TPACK; Koehler *et al.*, 2014), the context they worked in, and feedback regarding adaptations they made to the curriculum to support student learning.

After obtaining teacher consents for interview and recording, two researchers trained on the protocol conducted each recorded interview which lasted approximately 60 minutes. All three private school teachers, including the two observed, agreed to be interviewed about their experiences implementing RAICA. Private School A teacher participated in a pair of pre/post interviews around *Picture This* implementation in Spring 2024. Private School B teacher participated in a focus-group interview with Private School A teacher following the implementation of *Social Robots* in Fall 2024. Private School C teacher participated in an exit interview after implementing *App Inventor + LLMs* in Fall 2024.

### 3.3.3. Student Exit Tickets

Each lesson in the RAICA curriculum contains an "Exit Ticket" – a formative assessment teachers could deliver at the end of a lesson to re-orient students toward the learning goals, check for understanding, and provide feedback for both students and teachers. Teachers distributed them via a paper copy or a link to a digital version, depending on their preference and teaching style. Students worked individually, and had approximately five minutes to complete the assessment. Usage of exit tickets varied by teacher and lesson, resulting in a range of student completion data. The exit ticket analyzed for the current study was explicitly tied to ethical thinking through stakeholder engagement, asking students to define the term ("What does stakeholder mean?") and to plan methods of stakeholder input ("How do you plan to get feedback from stakeholders about your project?").

Ninety-seven middle school students from Private Schools A and B as well as Public School D completed the exit ticket used for data analysis in this study. The exit ticket data were from the *Picture This* module implementations in 2024.

### 3.3.4. Pre-post Measure

A low-stakes, open-ended ethical scenario measure was part of a questionnaire designed to provide a broader pre-post assessment of students' AI understanding. The broader assessment, intended to last 20 minutes, also included items measuring students' attitudes toward AI and their conceptual grasp of AI. Positioned at the end of the assessment, the scenario concerned a school setting where an AI system was being considered as a judge for an annual drawing contest. The scenario was followed by three open-ended questions, each designed to prompt reflection on key ethical dimensions: stakeholders, potential benefits, and potential harms (see Fig. 1).

The measure was grounded in the consequentialist premise that ethical thinking involves recognizing and articulating thoughts about a decision's societal effects. However, it did not assume a more specific, value-driven format like the ethical matrix (O'Neil and Gunn, 2020; Mepham, 2000) to ensure understandability across diverse student populations regardless of their prior exposure to ethics education.

Teachers followed a script to administer the questionnaire via Qualtrics during class time before and after they taught a RAICA module. All participating students completed the pre-module scenario measure except those from Public School D. The only students who completed both pre- and post- scenario measures were from Private Schools A and B, and did so around the *Social Robots* module. On average, students took about 19 minutes to complete the full questionnaire. A total of 54 middle and high school students completed the scenario-based pre-measure, and 29 of them also completed the post-measure.

The last set of questions ask you to evaluate the use of AI in the scenario below.

Imagine your school hosts an annual drawing contest where all students can submit drawings on any topic of their choice. This year, the school is considering using an AI system to judge the contest. The AI will score and rank the drawings based on their technical skill and creativity. In an upcoming assembly, the school will discuss the potential benefits and harms of using the AI judge.

Who might be impacted or involved in the use of the AI judge in the drawing contest?

What are three benefits people might mention about using the AI judge?

What are three harms people might mention about using the AI judge?

Fig. 1. Ethical Thinking Scenario Question.

### 3.4. Data Analysis

In this section, we describe our processes to examine and interpret the data. To establish trustworthiness of our qualitative coding, we provide transparent reporting of our analytic procedures and samples from the raw data (O'Connor and Joffe, 2020). As part of our team approach, we utilized multiple group discussions and group consensus to ensure consistency among multiple coders (Saldaña, 2021).

#### 3.4.1. Class Observations

One research team member reviewed the classroom observations and selected segments of class in which students were engaged in responsible design, which was primarily in the project plan, build, and test phases. The unit of analysis for this data was the observation that occurred every five minutes. This yielded 52 segments, some of which were repetitive, as both researchers were observing the same moment. In cases in which both researchers were observing the whole class, the most detailed observation was retained for analysis. In cases in which researchers were observing different groups during the same time period, both observations were retained. This process left 44 observations, which were then coded using Chi's Interactive, Constructive, Active, Passive (ICAP) framework (Chi, 2009).

The ICAP framework categorizes overt, observable student behaviors according to modes of engagement (interactive, constructive, active, or passive), and links each mode to underlying cognitive processes that influence the depth of student learning. The modes are hierarchical, with the interactive mode referring to behaviors that lead to deeper learning, as opposed to the passive mode which refers to behaviors that lead to more shallow learning (Menekse *et al.*, 2013; Chi *et al.*, 2018). Interactive, constructive, and active modes of engagement are all considered active learning, but an interactive mode leads to deeper learning than a constructive mode, both interactive and constructive modes lead to deeper learning than an active mode, and all three lead to deeper learning than a passive mode. For example, in our work if students were watching a video, the initial code applied was 'watching video' and the category according to the ICAP framework was 'passive' as that behavior aligns with a passive mode of engagement. If students were taking notes, the initial code was simply 'notetaking' and the ICAP category was 'active' as this demonstrated an active mode of engagement according to ICAP. In this study, we expected that students' constructive or interactive mode of engagement in responsible design activities would support better student understanding of the three dimensions of ethical thinking (stakeholders, potential benefits/harms, design values).

Two researchers coded one fourth of the data using the four ICAP modes of engagement as apriori codes. The initial comparison of results revealed poor intercoder agreement. Subsequent refining of the coding frame provided greater agreement. Any discrepancies were resolved through discussion, and the remaining data were then coded by one researcher.

### 3.4.2. *Teacher Interviews*

The interview recordings were transcribed and segmented for analysis by three researchers. The unit of analysis was at the level of teacher response, meaning that we considered an entire response to the interviewer's question as one unit. For analysis, we chose segments in which teachers addressed their perception of responsible design within the lesson. This process resulted in 13 segments for coding. One researcher used thematic analysis on the 13 segments, using an inductive approach to identify recurring patterns in teachers' responses to arrive at categories and common themes. As analysis can be an interpretive process subject to a single researcher's position and past experiences (Saldaña, 2021), two additional researchers reviewed the applied codes to mitigate any impact of personal bias. All three researchers then discussed and agreed upon the final results.

### 3.4.3. *Student Exit Tickets*

Student answers to exit tickets were compiled in a central spreadsheet. The unit of analysis was at the level of individual student responses to each question on the exit ticket. A single researcher inductively coded answers to open response questions to summarize general patterns in student thought. These included themes in student definitions of the word "stakeholder" and key words used when generating ideas to engage stakeholders. These summaries were reviewed by the entire team, who discussed insights the summaries revealed about student understanding along with proposed curriculum revisions.

### 3.4.4. *Pre-post Measure*

The open-ended responses to the three questions on potential stakeholders, benefits, and harms of the AI scenario were retrieved from Qualtrics and compiled on a spreadsheet prior to qualitative analysis. The unit of analysis was at the level of individual student responses. Key phrases were first extracted as in-vivo codes, which were then grouped into broader codes, and eventually categories following the approaches outlined by Mayring (2021) and Saldaña (2021). For example, a category for stakeholder responses was "Students," which included codes of "Students" and "Kids." A category under benefits, "Operational Benefits of AI," encompassed remarks like "AI reduces the workload for judges," "AI removes the need for judges," and "AI reduces cost for organizers." For harms, an example category was "Lack of Human Traits in AI," which included statements such as "AI cannot value art," "AI cannot feel," and "AI does not acknowledge the different aspects of art" (see Appendix A for full list of categories, corresponding codes, and examples).

Responses that indicated a lack of knowledge, such as "I don't know," or those left blank were coded and categorized separately. The total number of stakeholders, benefits, and harms mentioned by each respondent were also recorded. In counting, we excluded any miscellaneous responses, which were undecipherable, ambiguous, or irrelevant.

Three researchers independently coded subsets of responses to ensure consistency in the analysis. Researchers 1 and 2 conducted open coding of 100% of the stakeholder

data from the pre-measure to establish a preliminary coding frame. Researcher 1 then developed a coding frame for student responses related to AI benefits and harms. Researchers 1 and 3 then coded 50% of the pre-measure data related to stakeholders, benefits, and potential harms using the originally established coding frames and compared their results. After refining the original coding frames and discussing numerous examples, Researcher 1 then coded the remaining data.

Following the coding process, response patterns were descriptively summarized through frequency counts. Pre-post differences in the average number of stakeholders, benefits, and harms mentioned were analyzed using paired samples t-tests. To explore potential shifts in the categorical distribution of responses, Fisher's Exact test was used in place of Chi-Square tests due to violations of the expected cell count assumption. Visual comparisons were also created through bar charts.

After the qualitative and quantitative analyses, we integrated the data collected from these different perspectives (Creswell and Plano Clarke, 2018) to gain a better understanding of students' ethical thinking when engaged with the RAICA curriculum and how ethical thinking can be effectively measured.

#### 4. Results

In this section, we present the results of our analyses, organized by the research questions they address.

**RQ1:** How do students engage with curricular activities designed to promote ethical thinking?

Qualitative analysis of class observation data and teacher interviews revealed a positive pattern of student engagement in module activities along with teachers' varied instructional strategies to elicit more desired modes of engagement. During moments related to responsible design, our results showed that students were almost always actively engaged according to the ICAP framework. An active mode included behaviors such as physically demonstrating willingness to participate through a raised hand in response to a teacher question or writing on graphic organizers to plan projects. In the majority of observations (33 or 75%), students engaged with the material in constructive or interactive modes. Constructive modes of engagement included combining knowledge of AI concepts from class with prior background knowledge to brainstorm a project relevant to their interests or seeking out additional resources for their projects during lunch time. The Interactive mode included behaviors such as student-teacher or peer-peer dialogue and building projects in groups. Only three of the 44 (6.8%) instances were coded as "passive." This mode of engagement included behaviors such as listening to teacher directions or watching an informational video.

Analysis of teacher interviews revealed ways in which they engaged students with ethical thinking during RAICA modules as well as students' reaction to those activities. Two primary themes related to our research question. The first theme related to pedagogies teachers used to increase students' comprehension of AI's ethical implica-

tions. All interviewees described using active learning strategies to engage students, e.g., using real life examples or video to facilitate discussion of bias and recognition of potential benefits and harms. One teacher elaborated on the use of discussion, articulating their belief that discussion of bias was a good segue to thinking critically about AI for students. Another teacher stated that they used constructionism, or hands-on making, followed by a reflective discussion about ethical implications of students' creations.

One teacher described their own confidence in probing students' thinking about ethical considerations after module implementation as they described how they engaged their students.

*“So we feel much more knowledgeable, much more confident in how to describe it [AI], especially from an ethical implication. And as we send kids out to the world to create, you know, what are your values as [a] designer? What are your [project's] harms? What are your [project's] benefits? Are you thinking about that? Do you have a diverse team that are you getting feedback from? So it really helps us prepare these kids to think more than it [AI] being something on our phones, you know, in Google Maps and talking to Siri or Alexa.”*

The second theme was that engagement in the RAICA curriculum created a positive impact on students. One teacher gave three examples of students demonstrating combinations of creativity, joy, rigorous effort, accountability, and leadership skills during the build and testing phases of the curriculum. These are the points in which ethical thinking is inherent in students' activities.

**RQ2:** How do students describe stakeholders associated with AI tools and applications in formative assessments?

Analysis of students' responses to class exit tickets revealed that immediately following instruction, some students were not clear about the concept of stakeholder and ways to meaningfully include them in projects. When asked to define “stakeholder,” only 59% of the 97 students correctly identified it as anyone affected by the project. A significant portion (28%) responded with the partially correct answer of someone who used the project or benefited from it. The remaining 13% of students defined it with a misconception, including the overly general “everyone” and association with the financial definition of the term, “the owner of the organization.”

When asked to plan methods for gathering stakeholder feedback, 50% of the students mentioned a specific method for gathering feedback, including seeking out an expert and creating a survey. Thirty-eight percent of the students responded with general statements like “ask them for feedback,” and “contact them.” Some students (12%) misunderstood the question to be asking them to predict feedback they might get from stakeholders.

**RQ3:** Does engagement with the RAICA curriculum enable students to identify stakeholders in AI tools and applications as evidenced through a pre-post measure?

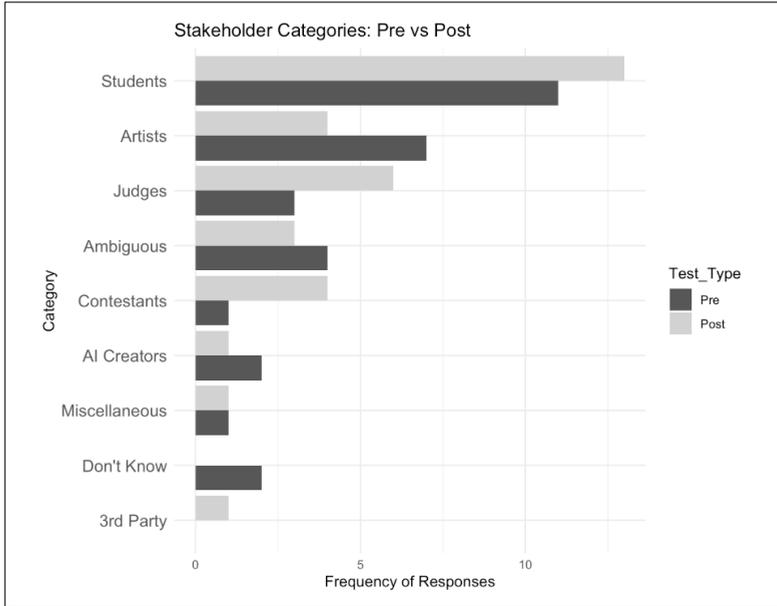


Fig. 2. Stakeholder Categories: Pre vs. Post.

Our analysis of students' pre-post scenario measures suggests that students may have deepened their understanding of stakeholders after completing a RAICA module, although not evident in the number of responses they provided. After engaging with the *Social Robots* module, students showed minimal change in the number and kinds of stakeholders they identified when presented with an AI scenario. From pre- to post-measure ( $n = 29$ ) the average number of stakeholders identified rose from 1.00 to 1.14, but the difference was not statistically significant,  $t(27) = 1.00$ ,  $p = 0.33$ . A Fisher's Exact test also showed no significant change in the distribution of stakeholder categories,  $p = 0.67$ .

Despite the non-significant changes observed, some trends are noteworthy (see Fig. 2). One is the increased mention of judges as the potential stakeholder from 5% in pre to 9% post. This increase may be hinting at an expansion in students' circle of concern, particularly since judges do not directly interact with the AI, unlike the other often-cited stakeholders like contestants. Students may have been considering the potential for the workload of judges to be reduced, or their roles to be displaced, by AI. Additionally, whereas references to students and contestants as stakeholders increased, from 17% to 20% and 2% to 6% respectively, mentions of artists declined from 11% to 6%. However, students and artists were often explicitly described as contestants in many of the responses, so these references may have referred to overlapping or identical groups.

**RQ4:** Does engagement with the RAICA curriculum enable students to identify potential benefits and harms of AI tools and applications as evidenced through a pre-post measure?

Our analysis of students' pre-post scenario measures suggests that although students did not name a greater number of benefits or harms, their responses showed an increased appreciation for those consequences. When comparing pre- and post- measure responses on the potential benefits and harms of AI application ( $n = 29$ ), there was a small, non-significant increase in the average number of benefits mentioned from 1.93 to 1.96,  $t(27) = 0.15$ ,  $p = 0.88$ . For harms, there was a non-significant decrease from 2.00 to 1.86,  $t(27) = -0.67$ ,  $p = 0.50$ . The Fisher's Exact test revealed no significant differences in the distribution of the benefit categories,  $p = 0.89$ , nor in the distribution of the harm categories,  $p = 0.75$ .

The benefit categories from pre to post show minimal overall change, in accordance with the lack of statistically significant shifts (see Fig. 3). However, several category-level patterns are noteworthy, especially when examined alongside corresponding categories on the bar chart of harm categories, which display more marked changes in student perceptions (see Fig. 4).

First, there was an increase in responses that cite improved accuracy and consistency as a benefit of AI usage, from 10% to 16%. This upward shift suggests a growing perception among students that AI systems, when rule-based, may outperform humans in making consistent, error-free judgement. Such a trend is mirrored by a sharp decline in mentions of accuracy-related harms, or concerns that AI might produce incorrect or inconsistent judgments, from 27% to 14%. Taken together, these patterns imply a growth in student trust in AI's procedural reliability, which could stem from a tendency to equate automation with accuracy and consistent execution of rules. However, this trust may be overlooking the fact that AI systems and their workings depend heavily on

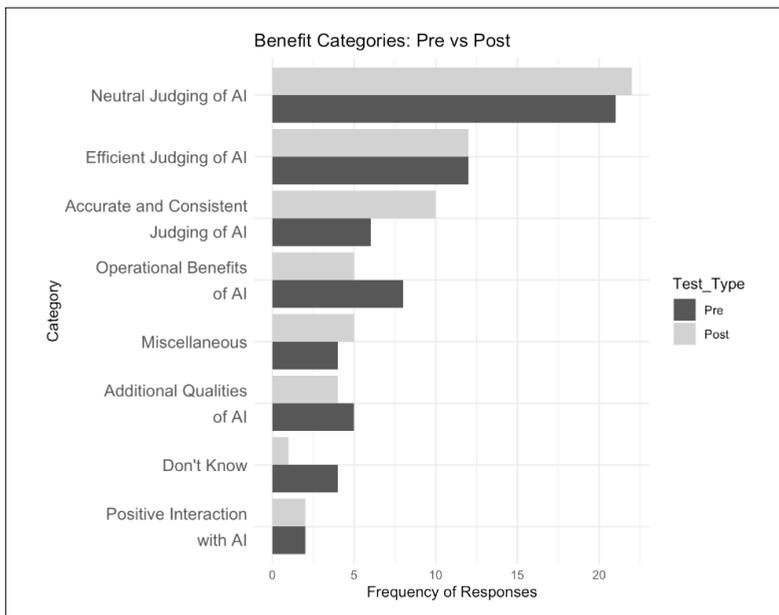


Fig. 3. Benefit Categories: Pre vs. Post.

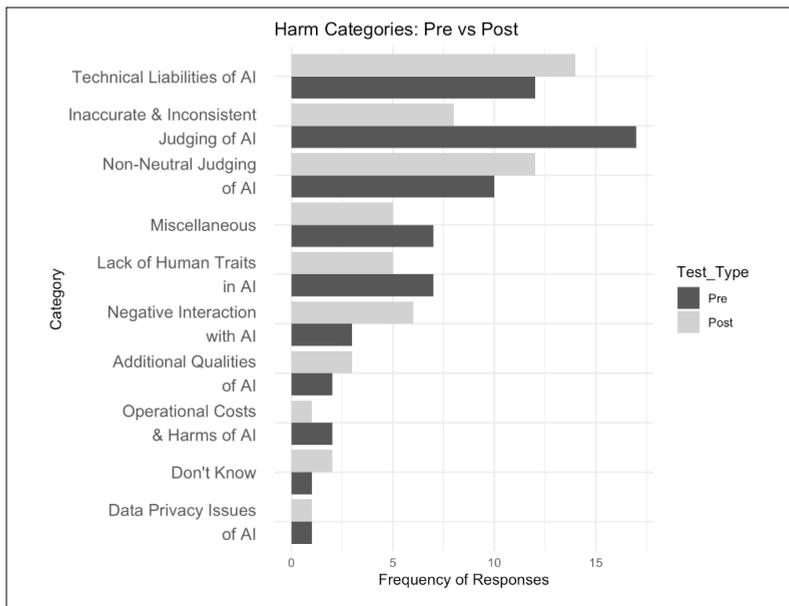


Fig. 4. Harm Categories: Pre vs. Post.

the humans who design, program, and train them. Therefore, making human influences on AI more visible to students to increase awareness of this potential fallacy and inaccuracy seems essential.

Second, we saw a decline in mentions of operational advantages, such as reduced cost and labor, from 13% to 8% among benefits. This decline, coupled with stakeholder responses showing increased mentions of judges as potential stakeholders, hints at students' growing awareness of how AI implementation will affect human roles. Such a shift, however, was not matched by a rise in mentions of operational harms, like job shortages or displacement, that remained relatively low at 3% pre and 2% post. This discrepant trend warrants further investigation.

Interestingly, AI's reduction, or lack, of bias was the most cited benefit across both pre- and post-measures, at 34% and 36% respectively. This may be because students are interpreting "bias" as human bias, such as favoritism or discrimination, rather than technological or algorithmic bias. Such an interpretation could manifest in their perception of AI as inherently neutral or objective. This aligns with prior literature that also noted students' perception of AI as an entirely objective tool (Mertala *et al.*, 2022). However, this perspective shows early signs of shifting. Following the module, mentions of AI bias as a harm increased from 16% to 21%, implying that some students started to question whether AI judges are impartial. The co-existence of both views underscores a tension between assuming AI as neutral and learning that it is a socially constructed, biased system.

The final trend worth highlighting is the doubling in mentions of harms related to human-AI interaction, such as the concern that contestants might feel hurt or offended by an AI's decision, from 5% to 11%. This suggests a growing awareness among stu-

dents of the social implications of AI use, particularly regarding the emotional impact of being judged by an AI.

**RQ5:** What assessment approach can provide reliable and valid information regarding students' ethical thinking skills?

In response to our final research question, we assert that our efforts to measure students' ethical thinking thus far have provided valuable insights with regards to the usefulness of scenario-based assessments and potential measurement biases to detect and mitigate.

First, we learned that the scenarios utilized in our pre-post measure establish a context that can help students to identify with the situation presented. Recent recommendations (Miao *et al.*, 2024) support use of scenarios for ethical thinking instruction, and accordingly we argue that they are appropriate for measurement of the construct as well. We learned that using open-ended items provides valuable information regarding students' understanding of our three ethical thinking dimensions, particularly about potential benefits and harms.

Second, analysis of only pre-test responses ( $n = 54$ ), which included students from four U.S. classes and one educational organization in southeastern Africa, revealed potential context-based differences in how students engaged with the scenario questions. Students in southeastern Africa identified, on average, fewer stakeholders, benefits, and harms than those in the U.S. The gap was especially notable for harms, where students in southeastern Africa mentioned an average of 0.85 harms versus 1.80 in the U.S., a statistically significant difference,  $U = 148$ ,  $p = .013$ . These differences may reflect a cultural variation in ethical thinking, as scholarship suggests that cultural context is closely linked with moral reasoning (Wilhelm and Gunawong, 2016; Tsui and Windsor, 2001). Prior research has also revealed a significant difference in moral reasoning depending on whether it is engaged in one's native language or foreign language (Costa *et al.*, 2014; Geipel *et al.*, 2015). Both abroad and in the U.S., a number of students in this study were not native English speakers, which may have impacted their experience with the ethical reasoning aspects of the curriculum.

Lastly, we recognize that the differences in students' pre-test responses may stem from potential cultural or contextual bias in the scenarios themselves. Being cognizant of these potential biases is essential to ensure the cultural validity and applicability of the questions across diverse learning contexts and student populations. To that end, we are currently conducting cognitive interviews at multiple implementation sites to examine students' thinking about the questions and their responses.

## **5. Discussion**

Taken together, data from our sample of class observations and teacher interviews suggest that teachers utilized a variety of approaches to engage students in considering ethical implications of AI, and that students responded to these efforts with behaviors indicating productive modes of engagement according to the ICAP framework (Chi, 2009). Class observations revealed that students were participating in class discussions, documenting

their work, planning, and executing their project designs, and discussing their work with peers. Classroom observers' impressions noted that the classroom was 'vibrant' during the project creation stage, suggesting that collaborative project-based learning is a feasible way to engage learners in ethical reasoning. AI curriculum developers should note our results considering previous findings that direct instruction was the most frequent pedagogical approach adopted in K-12 AI teaching units (Yue *et al.*, 2022).

The teacher's strategy of engaging students in making tangible projects to facilitate reflection on ethical implications aligns with the premise that constructionism is an optimal way to make abstract concepts like AI ethics more accessible (Dai, 2025; Papert and Harel, 1991). An important consideration for teacher professional development around AI curricula is the incorporation of TPACK (Koehler *et al.*, 2014) around AI tools, active learning pedagogies, and ethical considerations so that teachers have flexible knowledge regarding optimal ways to enhance student understanding of ethical thinking. Our data shows that teachers were able to effectively use the curriculum to support students in developing ethical reasoning about AI. While we do not imply causality, this finding provides support for our strategic inclusion of supports for TPACK in the RAICA curriculum. Almatrafi *et al.* (2024) found that teachers are the least frequently targeted population in AI literacy efforts, ranking lower than students, workforce, and family. These findings suggest that further research is needed to probe the methodology and significance of educator TPACK in the context of AI literacy interventions targeting students.

Students' in-class exit ticket responses showed that they had only a developing understanding of project stakeholders and appropriate ways that AI designers could or should interact with them. The extent of student misconceptions and partially correct responses suggests a need for direct instruction regarding vocabulary acquisition, clarifying the definition of stakeholder as anyone who is affected by the project, rather than someone who just uses it. Students would also benefit from explicit instruction about methods for obtaining stakeholder feedback. Although overreliance on the practice of direct instruction may be incompatible with project-based learning, some scholars assert that learners benefit from explicit instruction when dealing with novel information (Kirschner *et al.*, 2006). Future iterations of the curriculum would benefit from this idea, scaffolding students to identify stakeholders and feedback methods in a hypothetical project before doing so in their own work.

Results from the pre-post scenario measure likewise revealed a developing recognition of potential stakeholders, benefits, and harms when presented with an AI scenario within a familiar school context. Most students identified fewer benefits and harms than were prompted by the questions. Pre-post comparison results showed no statistically significant shifts in the number nor categories of responses but did reveal a growth in students' attention to the potential biases in AI systems, along with an increased perception that AI could be more computationally accurate than human decision-making. However, mentions of the potential environmental and operational harms of AI remained infrequent, and there were continued assumptions that AI is less biased than human decision-makers. These patterns point to aspects of AI ethics that may not naturally surface in learners' initial reasoning, suggesting the need for more intentional opportunities to engage with them in the curricula.

Our scenario measure results could have been influenced simply by survey fatigue and cognitive overload, as it was the final component of the pre-post measure. Students responded to two sets of attitude items and nine AI concept items prior to responding to the ethical thinking scenario. In future administrations, randomizing the order of items in the larger assessment may help reduce the potential for these issues.

Our work highlights the value of integrating multiple forms of data as part of the assessment approach. Formative assessment data, such as class observations, student exit tickets, and teacher interviews, as well as data from instruments designed specifically for research, are equally important assessment approaches, as both support the iterative process of design, assessment, and revision, necessary to produce curriculum that is effective, realistic, and achieves a primary goal of equitable AI education.

As with many education studies in naturalistic settings, this work has limitations. In addition to the impact of language and culture described earlier, our convenience sampling strategy resulted in study participants from only schools and teachers interested in implementing AI modules. Although we recruited and implemented our modules in different school contexts, providing data from diverse cultures, not all students in the participating schools gave assent for participation, resulting in a self-selected sample. These factors contributed to the small sample sizes for our different data types. A small sample size limited our ability to draw statistically robust conclusions from the pre-post scenario measure.

Another limitation is that we did not conduct repeated observations in the classrooms where students completed the pre-post measure and thus, teachers' facilitation practices over time may have varied from those observed, potentially affecting students' responses on the exit tickets and pre-post measure. However, we argue that our analyses across samples and over time still provide important feedback for the next iteration of the curriculum as well as significant insights for those developing AI education curricula. Lastly, we acknowledge that the importance of researchers' reflexive thought during qualitative research cannot be overstated. Our team approach with repeated discussion and iterations of our coding frames could still allow for potential bias.

## **6. Conclusion and Future Research**

Finding effective ways to incorporate ethical thinking into AI curricula and accurately measure the outcomes are critical components of improving AI literacy education. This work is one step in that direction. More broadly, at the intersection of technology, curriculum design, and ethical thinking, the pedagogical innovation embodied in the RAICA curriculum can improve not only students' AI literacy, but can effect epistemological change as well. Students as creators (of AI projects) and teachers as facilitators (of AI related content), embodied in RAICA module pedagogy, can shift students' understanding of the nature and acquisition of knowledge from a static phenomenon passively consumed to a dynamic one actively created or co-constructed. By hands-on making and building, students are not passively receiving information, but rather actively constructing their understanding as they acquire, evaluate, and use information (Papert and Harel,

1991). As facilitators of this AI curriculum, educators are invited to become co-learners with students rather than compelled to assume a position of authority about what is likely an unfamiliar technological and content domain for them. This paradigm shift may extend beyond the boundaries of AI literacy to learning in general and aligns with the goals of DBIR to create positive systems change.

Our analyses of data from class observations, teacher interviews, exit tickets, and a pre-post measure suggest that our sample of teachers are indeed adopting a facilitator role, and many students are engaging actively with module materials and with peers to co-construct knowledge as they build and present AI projects. As described earlier, our results also indicate areas of the curriculum where ethical thinking instruction needs to be strengthened. Aligned with our DBR methodology, our next steps include continued iteration of curriculum and assessment, triangulating data to improve instructional materials.

In tandem with curriculum development, we are in the early stages of developing a low-stakes, scalable pre-post measure of ethical thinking following our pre-post scenario administration. Our work begins to answer the question of how ethical thinking can be evaluated effectively and at scale. Our open-ended measure situating ethical reasoning within a real-world, age-appropriate scenario elicited short, easily coded responses that would support development of a rubric for evaluation across users and environments. Successful use of this question format informs others who are developing scalable assessment frameworks and in RAICA, it will provide a way for teachers to evaluate students' response to instruction.

From a research perspective, we realize that additional work needs to be conducted to measure the ethical thinking process itself—how it is manifested in student analysis of AI applications and in their planning or creation of an AI project. Probing students' thinking and reasoning processes as they engage with project creation is an ideal method to assess complex and multifaceted processes (Brennen and Resnik, 2012). This approach is often constrained by the time required to conduct a thorough assessment and a lack of rubrics or measurement scales that provide evidence for valid interpretation of students' knowledge and skills (Ng *et al.*, 2021b). Our challenge is to refine our assessment of ethical thinking to measure not only what students know about ethical thinking, but to measure how they actualize that understanding in the design of projects for authentic contexts.

Our next steps include refining our operationalization of ethical thinking as a construct to include the thinking processes to be measured. This includes recognition of potential ethical issues when evaluating AI tools or contexts that disregard design values such as fairness, accessibility, security, safety, inclusivity, sustainability, transparency, justice, privacy, or accountability. We plan to refine pre-post assessment scenarios and related questions to measure these latent variables. Increasing sample size in future administrations and randomizing the order of items in the larger assessments may help reduce the issues of survey fatigue and cognitive overload mentioned earlier. Another potential improvement to the measure design would be to develop and test a range of scenarios to assess how students transfer their understanding across different contexts. We also anticipate developing a rubric to evaluate responses to support the scalability and adaptability of our evaluation approach. Our goal is to produce an assessment tool that will produce reliable scores that can be interpreted as valid evidence of students' ethical knowledge and skills.

## Acknowledgments

This work was supported by a grant from DP World, United Arab Emirates. The authors also acknowledge Randi Williams, Angela Daniel, Sarah Wharton, Lydia Guterman, Mary Cate Gustafson-Quiett and Andy Stoiber for their support in integrating ethical thinking activities into the RAICA curriculum.

## References

- Ali, S., Payne, B.H., Williams, R., Park, H.W., Breazeal, C. (2019). Constructionism, Ethics, and Creativity: Developing Primary and Middle School Artificial Intelligence Education. <https://www.media.mit.edu/publications/constructionism-ethics-and-creativity/>
- Almatrafi, O., Johri, A., Lee, H. (2024). A systematic review of AI literacy conceptualization, constructs, and implementation and assessment efforts (2019–2023). *Computers and Education Open*, 6, 100173. <https://doi.org/10.1016/j.caeo.2024.100173>
- Barabas, C., Virza, M., Dinakar, K., Ito, J., Zittrain, J. (2018, January). Interventions over predictions: Reframing the ethical debate for actuarial risk assessment. In: *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, PMLR 81, 62–76.
- Bloom, P. (2004). *Descartes' Baby: How the Science of Child Development Explains what Makes us Human*. Random House.
- Brennan, K., Resnick, M. (2012). New frameworks for studying and assessing the development of computational thinking. In: *Proceedings of the 2012 Annual Meeting of the American Educational Research Association, Vancouver, Canada*, 1, 25.
- Borenstein, J., Drake, M.J., Kirkman, R., Swann, J.L. (2010). The Engineering and Science Issues Test (ESIT): A discipline-specific approach to assessing moral judgment. *Science and Engineering Ethics*, 16, 387–407. <https://doi-org.libproxy.mit.edu/10.1007/s11948-009-9148-z>
- Card, D., Smith, N.A. (2020). On consequentialism and fairness. *Frontiers in Artificial Intelligence*, 3(34). <https://doi.org/10.3389/frai.2020.00034>
- Chi, M.T.H. (2009). Active-Constructive-Interactive: A conceptual framework for differentiating learning activities. *Topics in Cognitive Science*, 1(1), 73–105. <https://doi.org/10.1111/j.1756-8765.2008.01005.x>
- Chi, M.T., Adams, J., Bogusch, E.B., Bruchok, C., Kang, S., Lancaster, M., Levy, R., Li, N., McEldoon, K.L., Stump, G.S., Wylie, R. (2018). Translating the ICAP theory of cognitive engagement into practice. *Cognitive Science*, 42(6), 1777–1832.
- Cobb, P., Confrey, J., diSessa, A., Lehrer, R., Schauble, L. (2003). Design experiments in educational research. *Educational Researcher*, 32(1), 9–13. <https://doi.org/10.3102/0013189X032001009>
- Creswell, J.W., Plano Clark, V.L. (2018). *Designing and Conducting Mixed Methods Research (3rd Ed.)*. Sage Publications, Los Angeles.
- Costa, A., Foucart, A., Hayakawa, S., Aparici, M., Apesteguia, J., Heafner, J., Keysar, B. (2014). Your morals depend on language. *PLoS ONE*, 9(4), e94842. <https://doi.org/10.1371/journal.pone.0094842>
- Dai, Y. (2025). Integrating unplugged and plugged activities for holistic AI education: An embodied constructionist pedagogical approach. *Education and Information Technologies*, 30, 6741–6764. <https://doi.org/10.1007/s10639-024-13043-w>
- DiPaola, D., Payne, B.H., Breazeal, C. (2020). Decoding design agendas: an ethical design activity for middle school students. In: *IDC'20: Proceedings of the Interaction Design and Children Conference*, 1–10. <https://doi.org/10.1145/3392063.3394396>
- Donaldson, T., Dunfee, T. W. (1994). Toward a unified conception of business ethics: Integrative social contracts theory. *The Academy of Management Review*, 19(2), 252–284. <https://doi.org/10.2307/258705>
- Donaldson, J.P., Han, A., Yan, S., Lee, S., Kao, S. (2024). Learning experience network analysis for design-based research. *Information and Learning Sciences*, 125(1/2), 22–43. <https://doi.org/10.1108/ILS-03-2023-0026>
- Flanagan, O. (2009). *Varieties of Moral Personality: Ethics and Psychological Realism*. Harvard University Press.

- Forsyth, S., Dalton, B., Foster, E.H., Walsh, B., Smilack, J., Yeh, T. (2021). Imagine a more ethical AI: Using stories to develop teens' awareness and understanding of artificial intelligence and its societal impacts. In: *2021 Conference on Research in Equitable and Sustained Participation in Engineering, Computing, and Technology (RESPECT)* IEEE, 1–2. <https://doi.org/10.1109/RESPECT51740.2021.9620549>
- Fishman, B., Marx, R.W., Blumenfeld, P., Krajcik, J., Soloway, E. (2004). Creating a framework for research on systemic technology innovations. *Journal of the Learning Sciences*, 13(1), 43–76. [https://doi.org/10.1207/s15327809jls1301\\_3](https://doi.org/10.1207/s15327809jls1301_3)
- Fishman, B.J., Penuel, W.R. (2018). Design-based implementation research. In: *International handbook of the learning sciences* (pp. 393–400). Routledge.
- Fishman, B.J., Penuel, W.R., Allen, A.-R., Cheng, B.H., Sabelli, N. (2013). Design-based implementation research: An emerging model for transforming the relationship of research and practice. *Teachers College Record: The Voice of Scholarship in Education*, 115(14), 136–156. <https://doi.org/10.1177/016146811311501415>
- Geipel, J., Hadjichristidis, C., Surian, L. (2015). The foreign language effect on moral judgment: The role of emotions and norms. *PLoS ONE* 10(7): e0131529. <https://doi.org/10.1371/journal.pone.0131529>
- Gilligan, C. (2018). Revisiting “In a Different Voice”. *LEARNing Landscapes*, 11(2), 25–30. <https://doi.org/10.36510/learnland.v11i2.942>
- Goldsmith, J., Burton, E., Dueber, D.M., Goldstein, B., Sampson, S., Toland, M.D. (2020, April). Assessing Ethical Thinking about AI. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(9), 13525–13528. <https://doi.org/10.1609/aaai.v34i09.7075>
- Green, B. (2019, December). “Good” isn’t good enough. In: *Proceedings of the AI for Social Good workshop at NeurIPS*, 17, 1–7. [https://aiforsocialgood.github.io/neurips2019/accepted/track3/pdfs/67\\_aisg\\_neurips2019.pdf](https://aiforsocialgood.github.io/neurips2019/accepted/track3/pdfs/67_aisg_neurips2019.pdf)
- Grohmann, R., Araújo, W.F. (2021). Beyond mechanical turk: The work of Brazilians on global AI platforms. In: Verdegem, P. (Ed.) *AI for Everyone? Critical Perspectives*. University of Westminster Press, London, 247–266. <https://doi.org/10.16997/book55.n>
- Hagendorff, T. (2020). The Ethics of AI ethics: An evaluation of guidelines. *Minds & Machines*, 30, 99–120. <https://doi.org/10.1007/s11023-020-09517-8>
- Holmes, W., Porayska-Pomsta, K., Holstein, K., Sutherland, E., Baker, T., Shum, S.B., Santos, O.C., Rodrigo, M.T., Cukurova, M., Bittencourt, I.I., Koedinger, K.R. (2022). Ethics of AI in education: Towards a community-wide framework. *International Journal of Artificial Intelligence in Education*, 32(3), 504–526. <https://doi.org/10.1007/s40593-021-00239-1>
- Koehler, M.J., Mishra, P., Kereluik, K., Shin, T.S., Graham, C.R. (2014). The technological pedagogical content knowledge framework. In: Spector, J.M., Merrill, M.D., Elen, J., Bishop, M.J. (Eds.), *Handbook of Research on Educational Communications and Technology*. Springer Science+Business Media, New York, 101–111. [https://doi.org/10.1007/978-1-4614-3185-5\\_9](https://doi.org/10.1007/978-1-4614-3185-5_9)
- Kirschner, P.A., Sweller, J., Clark, R.E. (2006). Why Minimal Guidance During Instruction Does Not Work: An Analysis of the Failure of Constructivist, Discovery, Problem-Based, Experimental, and Inquiry-Based Teaching. *Educational Psychologist*, 41(2), 75–86. [https://doi.org/10.1207/s15326985ep4102\\_1](https://doi.org/10.1207/s15326985ep4102_1)
- Kohlberg, L., Hersh, R.H. (1977). Moral development: A review of the theory. *Theory Into Practice*, 16(2), 53–59. <https://doi.org/10.1080/00405847709542675>
- Kuhn, D., Park, S.H. (2005). Epistemological understanding and the development of intellectual values. *International Journal of Educational Research*, 43(3), 111–124. <https://doi.org/10.1016/j.ijer.2006.05.003>
- Lee, I., Anderson, E., Hsiao, L. (2019). Teachers with GUTS: Developing teachers as computational thinkers through supported authentic experiences in computer modeling and Simulation [Conference presentation]. STEM+CI Summit, Alexandria, VA. <https://stemcsummit.edc.org/slides/Anderson.pdf>
- Long, D., Magerko, B. (2020). What is AI Literacy? Competencies and Design Considerations. In: *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, 1–16. <https://doi.org/10.1145/3313831.3376727>
- Li, L., Yu, F., Zhang, E. (2024). A Systematic Review of Learning Task Design for K-12 AI Education: Trends, Challenges, and Opportunities. *Computers and Education: Artificial Intelligence*, 6, 100217. <https://doi.org/10.1016/j.caeai.2024.100217>
- Lin, Z., Dai, Y. (2025, April). Fostering Epistemic Insights into AI Ethics through a Constructionist Pedagogy: An Interdisciplinary Approach to AI Literacy. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(28), 29171–29177. <https://doi.org/10.1609/aaai.v39i28.35190>
- Lintner, T. (2024). A systematic review of AI literacy scales. *npj Science of Learning*, 9(1), 50. <https://doi.org/10.1038/s41539-024-00264-4>

- Mawasi, A., Nagy, P., Finn, E., Wylie, R. (2022). Using Frankenstein-themed science activities for science ethics education: An exploratory study. *Journal of Moral Education*, 51(3), 353–369. <https://doi.org/10.1080/03057240.2020.1865140>
- Mayring, P. (2021). *Qualitative content analysis: A step-by-step guide*. SAGE Publications Ltd.
- Menekse, M., Stump, G.S., Krause, S., Chi, M.T. (2013). Differentiated overt learning activities for effective instruction in engineering classrooms. *Journal of Engineering Education*, 102(3), 346–374. <https://doi.org/10.1002/jee.20021>
- Mephram, B. (2000). A Framework for the Ethical Analysis of Novel Foods: The Ethical Matrix. *Journal of Agricultural and Environmental Ethics*, 12(2), 165–176. <https://doi.org/10.1023/A:1009542714497>
- Mertala, P., Fagerlund, J., Calderon, O. (2022). Finnish 5th and 6th grade students' pre-instructional conceptions of artificial intelligence (AI) and their implications for AI literacy education. *Computers and Education: Artificial Intelligence*, 3, 100095. <https://doi.org/10.1016/j.caeai.2022.100095>
- Miao, F., Shiohira, K., Lao, N. (2024). AI competency framework for students. UNESCO. <https://www.unesco.org/en/articles/ai-competency-framework-students>
- Molla, T., Nolan, A. (2020). Teacher agency and professional practice. *Teachers and Teaching*, 26(1), 67–87. <https://doi.org/10.1080/13540602.2020.1740196>
- Narayanan, A., Vallor, S. (2014). Why software engineering courses should include ethics coverage. *Communications of the ACM*, 57(3), 23–25. <https://doi.org/10.1145/2566966>
- Ng, D.T.K., Leung, J.K.L., Chu, S.K.W., Qiao, M.S. (2021a). AI Literacy: Definition, Teaching, Evaluation and Ethical Issues. *Proceedings of the Association for Information Science and Technology*, 58(1), 504–509. <https://doi.org/10.1002/pr2.487>
- Ng, D.T.K., Leung, J.K.L., Chu, S.K.W., Qiao, M.S. (2021b). Conceptualizing AI literacy: An exploratory review. *Computers and Education: Artificial Intelligence*, 2, 100041. <https://doi.org/10.1016/j.caeai.2021.100041>
- O'Connor, C., Joffe, H. (2020). Intercoder Reliability in Qualitative Research: Debates and Practical Guidelines. *International Journal of Qualitative Methods*, 19. <https://doi.org/10.1177/160940691989922>
- O'Neil, C., Gunn, H. (2020). Near-term artificial intelligence and the ethical matrix. In: Liao, M. S. (Ed.), *Ethics of Artificial Intelligence*, Oxford University Press, New York, 235–69.
- Papert, S., Harel, I. (1991). Situating Constructionism. In: Harel, I., Papert, S. (Eds.), *Constructionism*. Ablex Publishing, 1–11.
- Prediger, S. (2024). Conjecturing is not all: Theorizing in design research by refining and connecting categorical, descriptive, and explanatory theory element. *EDeR. Educational Design Research*, 8(1). <https://doi.org/10.15460/eder.8.1.2120>
- Rest, J.R. (1979). *Revised Manual for the Defining Issues Test: An Objective Test of Moral Judgment Development*. Minnesota Moral Research Projects.
- Rest, J.R., Narvaez, D., Thoma, S.J., Bebeau, M.J. (1999). DIT2: Devising and testing a revised instrument of moral judgment. *Journal of Educational Psychology*, 91(4), 644–659. <https://doi.org/10.1037/0022-0663.91.4.644>
- Roediger, H.L. III, Butler, A.C. (2011). The critical role of retrieval practice in long-term retention. *Trends in Cognitive Sciences*, 15(1), 20–27. <https://doi.org/10.1016/j.tics.2010.09.003>
- Roediger, H.L. III, Pyc, M.A. (2012). Inexpensive techniques to improve education: Applying cognitive psychology to enhance educational practice. *Journal of Applied Research in Memory and Cognition*, 1(4), 242–248. <https://doi.org/10.1016/j.jarmac.2012.09.002>
- Saldaña, J. (2021). *The Coding Manual for Qualitative Researchers (4th Ed.)*. Sage Publications Ltd.
- Sandoval, W. (2014). Science Education's Need for a Theory of Epistemological Development. *Science Education*, 98(3), 383–387. <https://doi.org/10.1002/sce.21107>
- Schrier, K., Diamond, J., Langendoen, D. (2010). Using Mission US: For Crown or Colony? To Develop Historical Empathy and Nurture Ethical Thinking. In: Schrier, K., Gibson, D. (Eds.), *Ethics and Game Design: Teaching Values Through Play*. Information Science Reference, Hershey, PA, 255–273.
- Schuitema, J., Dam, G. ten, Veugelers, W. (2008). Teaching strategies for moral education: a review. *Journal of Curriculum Studies*, 40(1), 69–89. <https://doi.org/10.1080/00220270701294210>
- Skirpan, M., Beard, N., Bhaduri, S., Fiesler, C., Yeh, T. (2018). Ethics Education in Context: A Case Study of Novel Ethics Activities for the CS Classroom. In: *Proceedings of SIGCSE '18: The 49th ACM Technical Symposium on Computing Science Education*, Baltimore, MD, USA, February 21–24, 2018 (SIGCSE '18) (pp. 940–945). <https://doi.org/10.1145/3159450.3159573>
- Singer, P. (2011). *The Expanding Circle: Ethics, Evolution, and Moral Progress*. Princeton University Press.
- Svihla, V. (2014). Advances in Design-Based research. *Frontline Learning Research*, 2(4), 35–45. <https://doi.org/10.14786/flr.v2i4.114>

- Taylor, K., Rohrer, D. (2010). The effects of interleaved practice. *Applied Cognitive Psychology*, 24(6), 837–848. <https://doi.org/10.1002/acp.1598>
- Torresen, J. (2018). A Review of Future and Ethical Perspectives of Robotics and AI. *Frontiers in Robotics and AI*, 4. <https://doi.org/10.3389/frobt.2017.00075>
- Trotta, A., Ziosi, M., Lomonaco, V. (2023). The future of ethics in AI: challenges and opportunities. *AI & Society*, 38, 439–441. <https://doi-org.libproxy.mit.edu/10.1007/s00146-023-01644-x>
- Tsui, J., Windsor, C. (2001). Some Cross-Cultural Evidence on Ethical Reasoning. *Journal of Business Ethics*, 31, 143–150. <https://doi-org.libproxy.mit.edu/10.1023/A:1010727320265>
- UNESCO. (2022). Recommendation on the Ethics of Artificial Intelligence. UNESCO. <https://unesdoc.unesco.org/ark:/48223/pf0000381137>
- Wharton, S., Gustafson-Quiett, M., Bosch, C., Davis, E., Breazeal, C., Abelson, H., Klopfer, E. (2024). Responsible design: A design thinking process for students creating with AI. In: *Play Make Learn, Annual Conference*. Madison, WI, United States. Poster session.
- Wilhelm, W.J., Gunawong, P. (2016). Cultural dimensions and moral reasoning: a comparative study. *International Journal of Sociology and Social Policy*, 36(5/6), 335–357. <https://doi.org/10.1108/IJSSP-05-2015-0047>
- Williams, R. (2021). A Review of Assessments in K-12 AI Literacy Curricula. Accessed (December, 2024) from [https://randi-c-dubs.github.io/K12-AI-ed/Constructionist\\_AI\\_Assessments.pdf](https://randi-c-dubs.github.io/K12-AI-ed/Constructionist_AI_Assessments.pdf)
- Williams, R., Ali, S., Devasia, N., DiPaola, D., Hong, J., Kaputsos, S.P., Jordan, B., Breazeal, C. (2023). AI + Ethics Curricula for Middle School Youth: Lessons Learned from Three Project-Based Curricula. *International Journal of Artificial Intelligence in Education*, 33(2), 325–383. <https://doi.org/10.1007/s40593-022-00298-y>
- Yue, M., Jong, M.S.-Y., Dai, Y. (2022). Pedagogical Design of K-12 Artificial Intelligence Education: A Systematic Review. *Sustainability*, 14, 15620. <https://doi.org/10.3390/su142315620>
- Zhang, H., Lee, I., Ali, S., DiPaola, D., Cheng, Y., Breazeal, C. (2023). Integrating Ethics and Career Futures with Technical Learning to Promote AI Literacy for Middle School Students: An Exploratory Study. *International Journal of Artificial Intelligence in Education*, 33(2), 290–324. <https://doi-org.libproxy.mit.edu/10.1007/s40593-022-00293-3>
- Zhang, H., Perry, A., Lee, I. (2025). Developing and Validating the Artificial Intelligence Literacy Concept Inventory: An Instrument to Assess Artificial Intelligence Literacy among Middle School Students. *International Journal of Artificial Intelligence in Education*, 35, 398–438. <https://doi.org/10.1007/s40593-024-00398-x>
- Zhu, Q., Zoltowski, C.B., Feister, M.K., Buzzanell, P.M., Oakes, W.C., Mead, A.D. (2014). The Development of an Instrument for Assessing Individual Ethical Decision-making in Project-based Design Teams: Integrating Quantitative and Qualitative Methods. In *ASEE Annual Conference and Exposition, Conference Proceedings*, 10060. <https://doi.org/10.18260/1-2--23130>

**G. Stump** is an Education Research Consultant at StumpWorks LLC since her retirement from MIT in 2023. She served as the Research Lead for RAICA during the initial years of development and implementation. Her work at MIT focused on development and evaluation of technology-enhanced STEM education and teacher/faculty professional development in international contexts. Her research also focuses on motivation and pedagogic theory, STEM teaching methods, and development of instruments to measure students' motivation as well as conceptual understanding of emergence and artificial intelligence. She holds a Ph.D. in Educational Psychology and a certificate in Educational Technology from Arizona State University.

**S. Kang** is a Research and Teaching Assistant in the Department of Communication and Media Research (IKMZ) at the University of Zurich. She holds a Master's Degree in Education from Harvard University and a Bachelor's Degree in Interactive Media and Business from New York University, Abu Dhabi. Her past work includes conducting AI literacy research at the Massachusetts Institute of Technology (MIT) and leading a blockchain education initiative in Abu Dhabi. Shinyi's current research explores how people practice and express creativity through digital media, and how these experiences are evolving with technological advancements.

**J. Masla** is an Assessment Specialist with MIT's Responsible AI for Social Empowerment and Education (MIT RAISE) initiative. He holds a Master's Degree in Education (M.Ed) from the University of Washington. His previous experience as a classroom teacher informs his work.

**C. Bosch** is a research scientist and currently serves as the Research Lead for RAICA implementation and evaluation, focusing on design-based research methods as a way of addressing issues around AI education in schools. Her work explores a range of experiences with/in education systems in various U.S. contexts and internationally, with particular interest in evidence-based instruction, inclusive curriculum design, teacher professional development, and research-practice partnerships that advance access, equity, and interest in life-long learning. She holds a Ph.D. in Special Education from the University of Massachusetts Amherst, a M.Ed. in Mind, Brain, and Education from Harvard University, and a M.A. in Special Education from American University.

**H. Abelson** is the Class of 1922 Professor of Computer Science and Engineering in the Department of Electrical Engineering and Computer Science at MIT. He has dedicated his work to making information technology more accessible to all and empowering everyone, especially young people, through computer science. Among other contributions, he was a founding director of both Creative Commons and the Free Software Foundation, creator of the MIT App Inventor platform, and co-author of the widely-used textbook *Structure and Interpretation of Computer Programs* (SICP).

**E. Klopfer** is Professor and Director of the Scheller Teacher Education Program and The Education Arcade at MIT. He is also co-PI of MIT's RAISE initiative in AI education. His research focuses on technology and pedagogy for building understanding of science, technology, engineering and mathematics (STEM) and systems. He has special interests in games, simulations and computing pathways to STEM learning.

**C. Breazeal** is a Professor of Media Arts and Sciences at MIT, where she founded and directs the Personal Robots group at the Media Lab. She is also the MIT Dean of Digital Learning and recognized as a pioneer of AI Literacy at MIT. She is the founding Director of the MIT-wide initiative on Responsible AI for Social Empowerment (RAISE), a research and outreach effort that advances access and inclusivity in AI education to people of all ages and backgrounds.

## Appendix A

### Coding Schemes for Analysis of Stakeholder, Benefit, and Harm Responses

Table A1  
Coding Scheme for Stakeholders

Category	Code	Examples
Students	Students	<i>"The students who submitted an art project."</i>
	Kids	<i>"Younger kids or kids that do not have artistic talent."</i>
Artists	Artists	<i>"The artists might be as the AI might not recognize the full length of their art."</i>
	Drawers	<i>"bad or good drawers"</i>
Judges	Judges	<i>"Human judges"</i>
	Teachers	<i>"Teachers might be impacted because they won't have to judge."</i>
Contestants	Participants / contestants / competitors	<i>"The people who enter the contest would be involved."</i>
AI Creators	Programmers / Coders	<i>"Someone who is really good at coding and AI."</i>
3rd Party	Viewers	<i>"The viewers"</i>
Ambiguous	People	<i>"People who are creative and think of something fun that humans understand, but not robots."</i>
	Everyone	<i>"Everyone because unless it gets hacked it cannot prefer a certain person."</i>
Miscellaneous	Irrelevant response	<i>"The contest won't be unfair."</i>
	Undecipherable response	<i>"judgeai"</i>
Don't Know	I don't know	<i>"I don't know"</i>
No Response	No response	

Table A2  
Coding Scheme for Potential Benefits

Category	Code	Examples
Efficient Judging of AI	AI makes judging easier	<i>"the AI judge would be easier than having multiple people decide on one thing"</i>
	AI makes judging quicker	<i>"It run things faster"</i>
Neutral Judging of AI	AI is less biased or not biased	<i>"They are not biased"</i>
	AI reduces or eliminates favoritism	<i>"There would be no judge that has a favorite"</i>

Continued on next page

Table A2 – continued from previous page

Category	Code	Examples
Accurate & Consistent Judging of AI	AI makes judging consistent	<i>“All are judged on the same scale/mind”</i>
	AI follows the rules	<i>“It will take all the contest as it was been said and the information will be taken as it was provided”</i>
	AI is not influenced by emotions, opinions, or thoughts	<i>“The A.I. has no feelings towards anyone”</i>
	AI is more accurate	<i>“It has no mistakes when judging”</i>
Positive Interaction with AI	AI provides feedback to contestants	<i>“It could give feedback [on] what could be better”</i>
	AI minimizes or reduces offense to contestants	<i>“They do not insult the contender”</i>
Operational Benefits of AI	AI reduces the workload for judges	<i>“Human judges get a break”</i>
	AI removes the need for judges	<i>“There doesn’t have to be a human judge”</i>
	AI reduces cost for organizers	<i>“They would not have to pay a real judge”</i>
Additional Qualities of AI	AI is technologically interesting	<i>“It would be cool to see how AI judges the art”</i>
	AI can do other things	<i>“[AI can] identify cheating”</i>
	AI is superior	<i>“[AI] is more modern and advanced”</i>
Miscellaneous	Irrelevant response	<i>“Development of a country”</i>
	Undecipherable response	<i>“it will help as.it will move as in abad way”</i>
Don’t Know	I don’t understand / I don’t know	<i>“I don’t know”</i>
No Response	No response	

Table A3  
Coding Scheme for Potential Harms

Category	Code	Examples
Non-Neutral Judging of AI	AI is biased or more biased	<i>“The AI might be biased”</i>
	AI is unfair	<i>“Might judge unfairly”</i>
	AI does not acknowledge the diversity of art	<i>“[the AI might...] not include all styles of drawing”</i>
Operational Costs and Harms of AI	AI incurs costs	<i>“It needs a lot of money”</i>
	AI takes away human jobs	<i>“The real judges might lose their jobs.”</i>
	AI requires infrastructure and/or resources	<i>“It requires work and maintenance”</i>

Continued on next page

Table A3 – continued from previous page

Category	Code	Examples
Technical Liabilities of AI	AI malfunctions or does not work	<i>“There might be a malfunction”</i>
	AI breaks	<i>“[The AI...] can fall and break”</i>
	AI is prone to hacking	<i>“it can be hacked”</i>
	AI is slower and/or less efficient	<i>“Not being fast”</i>
Negative Interaction with AI	AI can hurt contestants’ feelings	<i>“People might disagree with the AI and would probably get them aggravated”</i>
	AI cannot give feedback	<i>“The AI judge cannot give feedback”</i>
	AI feedback is less than humans	<i>“It’s always better to have a qualified judge.”</i>
Inaccurate and Inconsistent Judging of AI	AI is incorrect	<i>“It could get confused and mess up the judging”</i>
	There is not sufficient data-set	<i>“AI hasn’t been sufficiently trained”</i>
	AI does not recognize the drawing	<i>“[AI] judges bad when things looks similar”</i>
Lack of Human Traits in AI	AI is limited to the data and information it has been given	<i>“it only thinks of the information given to it”</i>
	AI cannot value art	<i>“It’s a computer so it has no value of art”</i>
	AI cannot feel	<i>“AI doesn’t have feelings”</i>
Data Privacy Issues of AI	AI does not acknowledge the different aspects of art	<i>“AI might pay attention to one aspect of art”</i>
	There are privacy issues around AI usage	<i>“AI collects people’s data”</i>
Additional Qualities of AI	AI is superior	<i>“AI is too accurate”</i>
	AI is less interesting	<i>“[AI is...] less entertaining”</i>
	AI has sentience	<i>“Some may believe AI has a mind of its own and it will take over our world”</i>
Miscellaneous	Undecipherable response	<i>“can short curit”</i>
	Irrelevant response	<i>“Causing death”</i>
	Ambiguous response	<i>“Programming”</i>
Don’t Know	I don’t understand / I don’t know	<i>“I don’t know”</i>
No Response	No response	