

INTERNATIONAL JOURNAL OF
EDUCATION, PSYCHOLOGY
AND COUNSELLING
(IJEPC)

<https://gaexcellence.com/ijepc>



TRUST IN STUDENTS' VOLUNTARY USE OF CHATGPT AND GENERATIVE AI IN HIGHER EDUCATION: A SYSTEMATIC LITERATURE REVIEW

Mohd Badli Ramli^{1*}, Suhaizal Hashim², Mohd Zulfadli Rozali², Aida Aisyah Moktar Abdullah²

¹Department of Mechanical Engineering, Politeknik Ibrahim Sultan, 81700 Pasir Gudang Johor, Malaysia

 mohdbadliramli@gmail.com

 <https://orcid.org/0009-0009-9257-9366>

²Faculty of Technical and Vocational Education, Universiti Tun Hussein Onn, 86400 Parit Raja, Johor, Malaysia

 suhaizal@uthm.edu.my
mzulfadli@uthm.edu.my
aydaaisyah94@yahoo.com

 <https://orcid.org/0000-0002-0236-1892>
<https://orcid.org/0000-0002-5163-7437>
<https://orcid.org/0009-0004-6764-2275>

*Corresponding Author

Article Info:

Article history:

Received date: 31.12.2025

Revised date: 12.01.2026

Accepted date: 26.02.2026

Published date: 11.03.2026

To cite this document:

Ramli, M. B., Hashim, S., Rozali, M. Z., Abdullah, A. A. M. (2026). Trust In Students' Voluntary Use of ChatGPT and Generative Ai in Higher Education: A Systematic Literature Review. *International Journal of Education, Psychology and Counselling*, 11(62), 680-701.

Abstract:

The rapid diffusion of generative artificial intelligence (GenAI), particularly tools such as ChatGPT, has intensified students' voluntary reliance on systems that operate under epistemic uncertainty. While existing research on GenAI adoption predominantly emphasises behavioural intention and usage outcomes, the psychological role of trust in enabling reliance remains conceptually fragmented. This systematic literature review synthesises empirical evidence on trust and techno-trust in students' voluntary use of GenAI within higher education contexts. Guided by PRISMA 2020, a systematic search of Scopus and Web of Science identified 11 empirical studies published between 2023 and 2025 for descriptive synthesis. The findings reveal substantial heterogeneity in how trust is defined, operationalised, and positioned within empirical models. Trust is predominantly conceptualised as a cognitively oriented evaluation of system reliability or output credibility. Across studies, it is variably positioned as an antecedent, mediator, moderator, outcome, or remains implicitly embedded within adoption constructs. Behavioural intention and use-related outcomes dominate the literature, whereas reliance, acceptance, and calibrated trust receive comparatively limited empirical attention. By consolidating fragmented evidence through an integrative typology and analytical mapping of trust positioning, this review clarifies the inconsistent analytical roles assigned to trust in voluntary GenAI use. The findings reconceptualise trust as a psychological reliance

mechanism operating under epistemic risk rather than a peripheral adoption variable. Educationally responsible engagement with GenAI therefore depends not on maximising trust, but on cultivating calibrated and warranted reliance. This synthesis provides a clearer foundation for future research on trust calibration, epistemic judgement, and decision-making in the educational use of generative artificial intelligence.

DOI: 10.35631/IJEPC.1162041 **Keyword:**

Academic Integrity; ChatGPT; Generative Artificial Intelligence; Higher Education Students; Techno-Trust; Technology Adoption; Trust In AI; Voluntary Use



© The authors (2026). This is an Open Access article distributed under the terms of the Creative Commons Attribution (CC BY NC) (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact ijepec@gaexcellence.com.

Introduction

The rapid diffusion of generative artificial intelligence (GenAI) in higher education has created a distinctive paradox in students' learning practices. On the one hand, GenAI systems such as ChatGPT are widely recognised as opaque, probabilistic, and prone to producing hallucinated or unverifiable outputs. On the other hand, students increasingly choose to rely on these systems in their academic work, even when such reliance is entirely optional. Unlike institutionally mandated learning technologies, the use of GenAI by students is largely voluntary, unregulated, and self-initiated. This situation raises a fundamental question for educational research: why do students decide to rely on epistemically risky systems when reliance is not required? Addressing this paradox requires moving beyond general accounts of technology adoption to examine the psychological mechanisms that enable reliance under conditions of uncertainty.

A growing body of empirical research has examined students' engagement with GenAI through the lens of technology adoption, typically focusing on constructs such as perceived usefulness, ease of use, intention to use, or frequency of use. While these studies provide valuable insights into behavioural uptake, they often conflate adoption with reliance and treat trust as either implicit or secondary. Adoption reflects a decision to use a technology, whereas reliance involves a deeper judgement about whether outputs can be depended upon in cognitively and academically consequential tasks. Trust, in turn, represents a distinct psychological state that enables individuals to accept vulnerability when outcomes cannot be fully predicted or monitored. Despite this conceptual distinction, much of the existing GenAI literature continues to operationalise usage intention or self-reported use as proxies for trust, leaving the psychological basis of reliance insufficiently theorised.

This conceptual gap becomes particularly salient in voluntary use contexts. When technology use is optional, individuals are not compelled by institutional requirements or formal incentives but must decide for themselves whether a system is sufficiently reliable to warrant dependence. In such situations, perceived utility alone is insufficient to explain sustained engagement. Instead, students must negotiate uncertainty associated with algorithmic opacity, potential misinformation, and the epistemic consequences of delegating cognitive tasks to AI systems. Trust functions here as a critical mechanism that mediates the transition from initial experimentation to sustained reliance, shaping not only whether GenAI is used, but also how and to what extent it is integrated into academic decision-making.

The importance of trust in GenAI use is further amplified by concerns surrounding academic integrity and ethical responsibility. Decisions to rely on AI-generated content are embedded within normative expectations about authorship, accountability, and legitimate knowledge production. Students' judgements about whether GenAI outputs are acceptable for academic purposes therefore extend beyond considerations of convenience or efficiency. Such judgements involve assessments of credibility, appropriateness, and academic risk. In this sense, trust in GenAI intersects with epistemic vigilance and moral reasoning, rather than functioning solely as a utilitarian belief about system performance.

Despite its centrality, trust remains conceptually fragmented within the GenAI adoption literature. Across empirical studies, trust is variably positioned as an antecedent of use, a mediating mechanism, an outcome of experience, or an implicit assumption embedded within other constructs. These inconsistencies are compounded by the frequent integration of trust into broader adoption models without explicit theoretical justification. As a result, findings related to trust are dispersed across studies and lack cumulative coherence, making it difficult to draw clear conclusions about its role in students' voluntary use of GenAI.

Against this backdrop, the present systematic literature review adopts a diagnostic and early evidence synthesis approach to clarify how trust and techno-trust are conceptualised and positioned in empirical studies of students' voluntary GenAI use in higher education. Focusing on studies published between 2023 and 2025, this review synthesises evidence from student-only samples within higher education and TVET contexts to examine how trust is defined and operationalised, to map the positioning of trust within empirical models of GenAI use, and to identify conceptual and methodological gaps that constrain theoretical integration. Rather than proposing a new adoption framework, this review organises existing evidence through an integrative trust typology and a trust positioning map, offering a clearer account of trust as a psychological reliance mechanism under conditions of uncertainty. By consolidating fragmented findings, the review provides a more coherent foundation for understanding trust in students' voluntary engagement with GenAI.

Conceptual Background

Trust has long been recognised as a foundational concept in human interaction with complex systems, particularly in contexts characterised by uncertainty, risk, and dependence. Within the technology literature, trust is commonly defined as a willingness to rely on a system despite the possibility of negative outcomes that cannot be fully controlled or predicted. This definition foregrounds vulnerability and reliance, rather than favourable attitudes or behavioural intentions alone. Classical accounts conceptualise trust as a psychological state that enables individuals to accept uncertainty based on expectations of system competence, reliability, and

benevolence. Importantly, trust is analytically distinct from related constructs such as satisfaction, attitudes, or intention to use, as it specifically concerns the decision to rely on a system when verification is incomplete and outcomes remain uncertain (Mayer et al., 1995; McKnight et al., 2002; Lee & See, 2004).

As technological systems evolved from passive tools to semi-autonomous and autonomous agents, the salience of trust increased correspondingly. In automated and artificial intelligence systems, users frequently delegate tasks, judgements, or decisions to technologies whose internal processes are opaque and difficult to observe. This opacity intensifies uncertainty and heightens dependence on system outputs, rendering trust a necessary psychological condition for effective interaction. Prior research demonstrates that trust in automation and AI shapes reliance behaviour, influences users' ability to calibrate trust appropriately, and affects continued use over time. Both over-reliance and under-reliance can undermine system effectiveness, underscoring trust as a dynamic and context-sensitive mechanism rather than a static belief (Lee & See, 2004; Hoff & Bashir, 2015; Lankton et al., 2015; Glikson & Woolley, 2020). The emergence of generative artificial intelligence represents a further shift in the nature of human–technology interaction, with distinctive implications for trust. Unlike traditional automation or rule-based systems, GenAI tools such as ChatGPT generate novel outputs that are probabilistic, non-deterministic, and often difficult to verify. These systems may produce fluent and persuasive responses that nonetheless contain factual inaccuracies, biased reasoning, or fabricated information. Consequently, trust in GenAI extends beyond assessments of technical reliability to encompass epistemic judgements concerning the credibility and trustworthiness of generated content. In this context, trust relates not only to whether a system functions as intended, but also to whether its outputs can be relied upon as a basis for knowledge construction and decision-making. This introduces a distinctive epistemic dimension of trust that is particularly salient in educational settings (Dwivedi et al., 2023; Kasneci et al., 2023; Amaro et al., 2024; Song, 2025).

Within higher education, the role of trust in GenAI use is further complicated by concerns surrounding academic integrity and ethical norms. Students increasingly employ GenAI tools for academic writing, coursework preparation, and problem solving, which blurs established boundaries between legitimate academic support and misconduct. Decisions to rely on AI-generated content are therefore embedded within normative expectations regarding originality, honesty, and accountability. In such contexts, trust is not merely a functional judgement about system performance, but also a moral and social consideration regarding whether the use of GenAI aligns with academic values. Students must assess not only the accuracy of GenAI outputs, but also their appropriateness and acceptability within institutional and disciplinary norms (Cotton et al., 2023; Subhani et al., 2025; Jo, 2024).

Despite its apparent importance, trust remains conceptually underdeveloped within dominant technology acceptance and adoption frameworks. In models such as the Technology Acceptance Model and the Unified Theory of Acceptance and Use of Technology, trust is frequently excluded or incorporated as an auxiliary construct without explicit theoretical justification. Empirical studies examining GenAI use reflect this inconsistency, with trust variably positioned as an antecedent, a mediating mechanism, an outcome of experience, or omitted altogether. As a result, findings related to trust are fragmented across studies and lack cumulative coherence, making it difficult to draw systematic conclusions about its role in students' voluntary use of GenAI (Al Amin et al., 2025; Shahzad et al., 2024; Rahman et al., 2023; Abdalla, 2025; Guo & Erdenebold, 2025).

Taken together, the literature points to a clear gap in the systematic understanding of trust and techno-trust in the context of GenAI use in higher education. While trust is widely acknowledged as influential, there remains no consolidated synthesis that clarifies how trust is defined, operationalised, and empirically positioned across studies focusing on students' voluntary engagement with GenAI. This lack of conceptual alignment constrains theoretical clarity and limits the integration of empirical findings. Addressing this gap requires a focused synthesis that maps how trust is conceptualised and positioned across studies, and that identifies conceptual and methodological inconsistencies that hinder cumulative knowledge development. Such clarification is necessary to support a more coherent understanding of trust as a psychological reliance mechanism in voluntary GenAI use.

Research Question

Existing literature indicates that trust plays an important role in students' use of generative artificial intelligence (GenAI) in educational settings, yet its treatment remains inconsistent across empirical studies. While prior research has examined trust in AI systems, substantial variation is evident in how trust is defined, conceptualised, and operationalised, as well as how it is positioned within empirical models. In many cases, trust appears as an auxiliary construct or is addressed only implicitly through related notions such as credibility or reliance. This fragmentation constrains the development of a coherent understanding of trust and techno-trust as psychological mechanisms underpinning students' voluntary use of GenAI. Consequently, a systematic synthesis is required to consolidate existing empirical evidence, identify dominant patterns, and evaluate conceptual and methodological gaps in the literature.

To ensure methodological clarity and alignment between the scope of this review and its guiding questions, this systematic literature review was framed using the PICOS elements. The population comprised higher education students, including university and TVET students. The phenomenon of interest focused on trust-related mechanisms in voluntary GenAI use, including trust in AI, techno-trust, perceived credibility, epistemic trust, and reliance. A formal comparator was not specified, as the review aimed to descriptively map and synthesise existing evidence rather than evaluate interventions. Outcomes included behavioural outcomes such as behavioural intention, actual use, continuance, and reliance-related outcomes, as well as psychological outcomes associated with trust. The study design was restricted to empirical journal articles published between 2023 and 2025 in the post-ChatGPT period.

Accordingly, this systematic literature review was guided by the following research questions:

RQ1. What definitions and conceptualisations of trust are used in empirical studies examining the use of generative artificial intelligence in educational contexts?

RQ2. What types and dimensions of trust, such as system trust, information trust, techno-trust, or epistemic trust, have been identified in studies of students' voluntary GenAI use?

RQ3. How is trust positioned within existing empirical models, for example as an antecedent, mediator, moderator, or outcome?

RQ4. What behavioural and psychological outcomes, such as behavioural intention, actual use, reliance, or continuance, are associated with trust in studies of voluntary GenAI use?

RQ5. What conceptual and methodological gaps are evident in the current empirical literature on trust in students' use of generative artificial intelligence in education?

Material and Methods

This systematic literature review was conducted in accordance with the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) 2020 guidelines to ensure transparency, replicability, and auditability throughout the review process. The methodological procedures were designed to support a diagnostic synthesis of empirical evidence on trust and techno-trust in students' voluntary use of generative artificial intelligence within higher education contexts.

Identification

The identification phase involved a systematic search of academic literature to capture empirical studies examining trust-related mechanisms in students' voluntary use of generative artificial intelligence within higher education and TVET contexts. Two major bibliographic databases, Scopus and Web of Science, were selected due to their extensive coverage of peer-reviewed journals in education, social sciences, and technology-related research.

Search strategies were developed to target studies focusing on generative artificial intelligence tools such as ChatGPT, student populations, and educational contexts, with specific emphasis on trust and techno-trust related constructs. In Scopus, searches were conducted using the TITLE-ABS-KEY field, while an equivalent topic-based search strategy was applied in Web of Science. Boolean operators were used to combine GenAI-related terms, trust-related keywords, and education-specific descriptors to ensure comprehensive yet focused retrieval of relevant studies.

Database filters were applied to restrict the results to English-language journal articles published between 2023 and 2025, reflecting the post-ChatGPT period and ensuring the contemporary relevance of the evidence base. Only records retrieved from the final refined search iterations were retained for subsequent screening and eligibility assessment in line with PRISMA 2020 recommendations. The complete search strategies and database filters are summarised in Table 1.

Table 1: Search Strategy and Database Filters

Database	Search Field	Search String	Filters Applied	Date of Access
Scopus	TITLE-ABS-KEY	("trust" OR "technology trust" OR "techno-trust") AND ("ChatGPT" OR "generative AI" OR "generative artificial intelligence" OR "large language model*") AND (student* OR "higher education" OR university OR TVET OR education)	Year: 2023–2025; Document type: Article; Language: English	1 December 2025
Web of Science	Topic (TS)	("trust" OR "technology trust" OR "techno-trust") AND ("ChatGPT" OR "generative AI" OR "generative artificial intelligence" OR "large language model*") AND (student*	Year: 2023–2025; Document type: Article;	1 December 2025

OR “higher education” OR
university OR TVET OR
education)
Language:
English

Source: Author’s Own Work

Screening

All records identified from the Scopus and Web of Science searches were exported and consolidated into a single dataset. Duplicate records were removed prior to screening. The remaining records were screened based on their titles and abstracts to assess their relevance to the objectives of this review.

Predefined inclusion and exclusion criteria were applied consistently during the screening phase. Studies were included if they were empirical journal articles published between 2023 and 2025, written in English, situated in higher education or TVET contexts, and involved student populations. In addition, studies were required to explicitly or implicitly examine trust or techno-trust related mechanisms in the context of generative artificial intelligence use. Articles were excluded if they were conceptual in nature, conducted outside educational settings, focused exclusively on educators or organisations, or examined AI technologies without a trust-related focus. The inclusion and exclusion criteria applied during screening are summarised in Table 2.

Table 2: Inclusion and Exclusion Criteria

Criteria	Inclusion	Exclusion
Publication year	2023–2025	Before 2023
Document type	Journal articles	Conference papers, reviews, books
Language	English	Non English
Context	Higher education or TVET	Non educational contexts
Population	Students	Educators only or organizations
Focus	Trust or techno trust in GenAI use	AI studies without trust focus

Source: Author’s Own Work

Eligibility

Articles that passed the screening stage were subjected to full-text assessment to determine their eligibility for inclusion in the final synthesis. Full-text evaluation focused on the alignment of each study with the review objectives, the clarity of trust-related constructs or proxies, and the adequacy of empirical reporting.

Studies were excluded at this stage if trust or related mechanisms could not be identified or operationalised in relation to generative artificial intelligence use, if the empirical evidence was insufficient, or if the study did not support the analytical focus of this review. All exclusion decisions at the eligibility stage were documented using the predefined exclusion codes to ensure transparency and traceability within the PRISMA flow.

Data Abstraction and Analysis

Data from the eligible studies were extracted using a structured data abstraction form to ensure consistency across studies. Extracted information included bibliographic details, study context, research design, definitions and types of trust examined, the positioning of trust within empirical models, and reported behavioural or psychological outcomes associated with trust.

Beyond descriptive extraction, the analysis explicitly coded the analytical positioning of trust within each empirical study. Trust-related constructs were categorised according to their functional role, namely trust positioned as an antecedent, trust positioned as a mediating mechanism, or trust positioned as an outcome. This coding approach enabled systematic mapping of how trust has been integrated across empirical models of students' voluntary use of generative artificial intelligence, without imposing assumptions regarding its theoretical primacy.

Data analysis followed a descriptive and mapping-based synthesis aligned with the research questions and the diagnostic objectives of this review. Extracted data were synthesised to identify patterns in how trust was conceptualised, operationalised, and positioned across studies. No meta-analysis was conducted due to heterogeneity in study designs, measurement approaches, and trust operationalisation across the included studies. Quality appraisal of the included articles was conducted as a subsequent analytical step and is reported in the following section.

Quality Appraisal

Quality appraisal was conducted to enhance the transparency and interpretability of the evidence synthesised in this systematic literature review. Consistent with established practices in diagnostic and mapping-oriented SLRs, quality assessment was applied as an audit layer rather than as an exclusion criterion. This approach was adopted because the primary objective of the review was to map and synthesise existing empirical evidence on trust and techno-trust in students' voluntary use of generative artificial intelligence within educational contexts, rather than to evaluate intervention effectiveness or test causal relationships. Accordingly, all studies that met the inclusion criteria were retained for subsequent synthesis, while quality appraisal outcomes were used to contextualise the interpretation of findings.

Each included study was assessed using six quality appraisal criteria (QA1–QA6). For each criterion, studies were rated using a three-point scheme: Yes (Y) when the criterion was clearly and adequately addressed, partly (P) when the criterion was addressed with limited detail or clarity, and No (N) when the criterion was not addressed or not reported. In addition to categorical ratings, a total score and percentage were calculated to provide a descriptive indication of reporting quality and methodological clarity. Importantly, no minimum quality threshold was imposed for study exclusion. The quality appraisal results for the eleven empirical studies included in this review are presented in Table 3.

Table 3: Quality Appraisal Results for Included Studies

Study	QA1	QA2	QA3	QA4	QA5	QA6	Total	%
Song (2025)	Y	Y	Y	Y	Y	Y	6.0	100
Guo & Erdenebold (2025)	Y	Y	Y	Y	Y	Y	6.0	100
Al-Dmour et al. (2025)	Y	Y	Y	Y	Y	Y	6.0	100
Al Amin et al. (2025)	Y	Y	Y	Y	Y	Y	6.0	100
Abdalla (2025)	Y	Y	Y	Y	Y	Y	6.0	100
Shahzad et al. (2024)	Y	Y	Y	Y	Y	Y	6.0	100
Rahman et al. (2023)	Y	Y	Y	Y	Y	Y	6.0	100
Amaro et al. (2024)	Y	Y	Y	Y	Y	Y	6.0	100
Kıyak et al. (2025)	Y	Y	Y	Y	Y	Y	6.0	100
Subhani et al. (2025)	Y	Y	Y	N	Y	Y	5.0	83
Jo (2024)	Y	Y	Y	N	Y	Y	5.0	83

Source: Author's Own Work

Quality Appraisal Criteria

QA1. Clarity of research objectives.

The study clearly states its objectives or research questions and maintains consistency with its stated focus.

QA2. Appropriateness of research design.

The research design is appropriate for addressing the stated objectives and the study context.

QA3. Clarity of context and sample.

The educational context and student sample characteristics are clearly described.

QA4. Operationalisation of the trust construct.

Trust or techno-trust is explicitly defined and measured or clearly operationalised through theoretically grounded proxies.

QA5. Transparency of analysis and reporting.

Data analysis procedures and reporting of findings are clearly described and traceable.

QA6. Alignment of discussion with objectives.

The discussion of findings is aligned with the stated objectives, with specific reference to trust-related constructs or mechanisms.

For Kıyak et al. (2025), behavioural measures such as advice taking and weight of advice were treated as theoretically grounded proxies for perceived credibility and trust within an experimental context. For Subhani et al. (2025) and Jo (2024), QA4 was rated as No because trust was not measured as an explicit construct, although these studies were retained to inform comparative adoption-related factors and to highlight trust-related conceptual gaps within the literature.

Overall, the quality appraisal indicates that the included studies demonstrate acceptable levels of methodological clarity and reporting transparency for the purposes of descriptive synthesis. Rather than functioning as a ranking mechanism, the appraisal outcomes provide contextual cues for interpreting variations in how trust is defined, operationalised, and positioned across empirical models of students' voluntary GenAI use.

PRISMA Flow Diagram

The study selection process for this systematic literature review was conducted in accordance with the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) 2020 guidelines to ensure transparency and traceability across all stages of study identification, screening, and inclusion. A systematic search of the Scopus and Web of Science databases resulted in the identification of 53 records, comprising 13 records from Scopus and 40 records from Web of Science.

Following the consolidation of records from both databases, 11 duplicate articles were removed. The remaining 42 unique records were subsequently screened based on their titles and abstracts to assess their relevance to the objectives of this review. At this screening stage, 8 records were excluded as they did not meet the predefined inclusion criteria, primarily due to the absence of an explicit or implicit focus on trust-related mechanisms in the context of generative artificial intelligence use in education.

The full texts of the remaining 34 articles were then assessed for eligibility. During this stage, 23 articles were excluded based on predefined exclusion criteria. These exclusions were mainly attributable to studies that did not operationalise trust or techno-trust in relation to generative artificial intelligence use, were not situated within higher education or TVET contexts, or did not involve student populations. All exclusion decisions were documented using established exclusion codes (E1–E3) to maintain consistency and auditability.

As a result of this selection process, 11 empirical studies met all inclusion criteria and were included in the final synthesis of this systematic literature review. The overall study selection process is illustrated in Figure 1 using the PRISMA 2020 flow diagram.

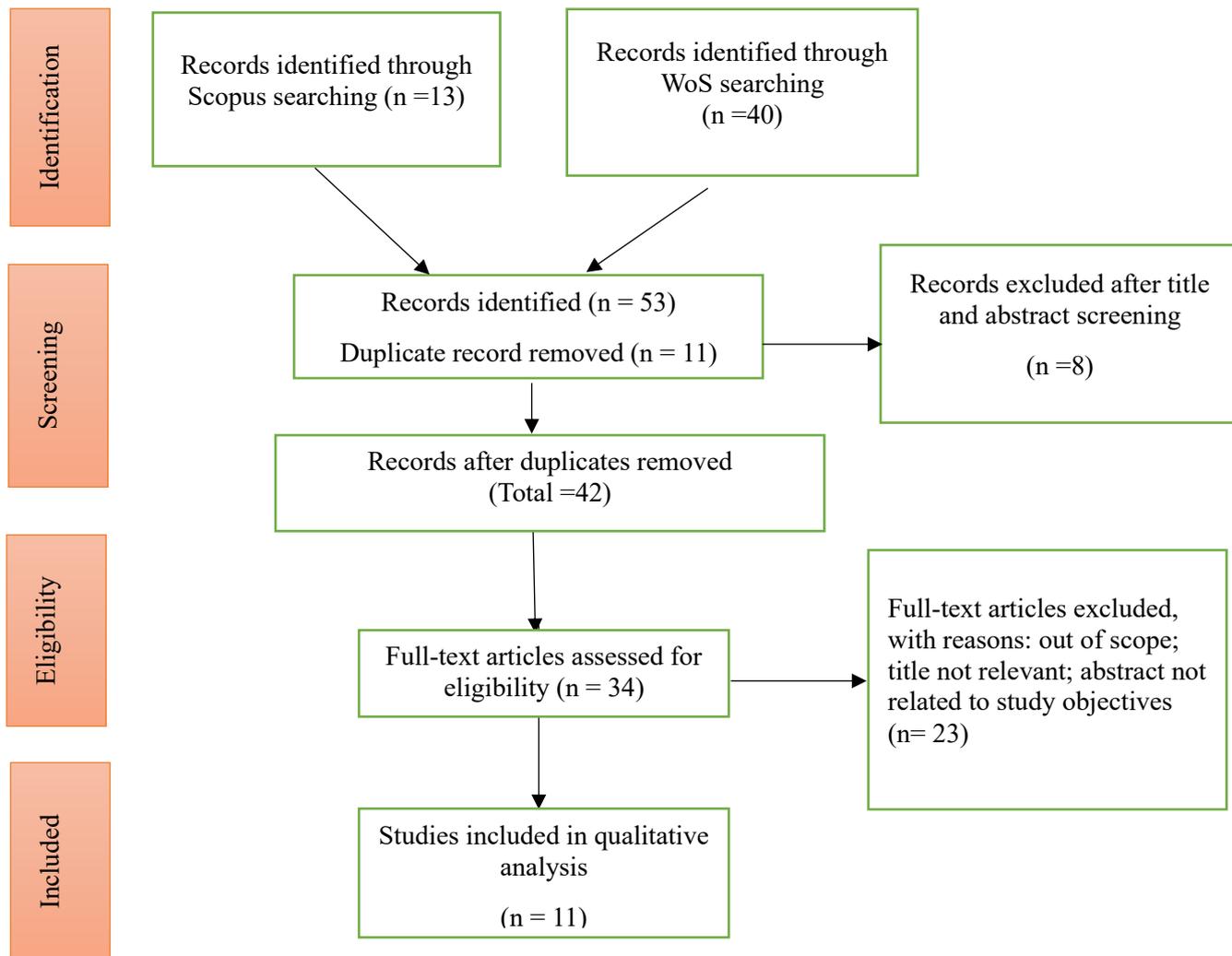


Figure 1: PRISMA 2020 Flow Diagram Illustrating the Study Selection Process

Source: Author's Own Work

Results

The results of this systematic literature review are based on 11 empirical studies focusing exclusively on student populations in higher education and TVET contexts. All findings reported in this section are descriptive and derived directly from the locked dataset. No additional interpretation, theoretical inference, or evaluative judgement is introduced beyond the synthesis of reported evidence.

Characteristics of Included Studies

Across the included studies, trust was conceptualised in varied and non-uniform ways, reflecting the absence of a shared definitional framework in the empirical literature on students' voluntary use of generative artificial intelligence. A majority of studies provided explicit definitions of trust, typically drawing on established trust literature or adapting existing definitions to the context of AI systems. These definitions commonly framed trust as a belief in the reliability, dependability, or credibility of the GenAI system or its outputs.

Table 4: Integrative Typology of Trust in Students' Voluntary Use of Generative Artificial Intelligence**

Trust orientation	Conceptual basis	How trust is inferred or operationalised	Typical outcomes examined	Epistemic risk profile
System-oriented trust	Cognitive evaluation of system reliability and dependability	Self-reported trust scales focusing on system performance	Behavioural intention, adoption, acceptance	Overgeneralisation of system competence
Output-oriented trust	Credibility and perceived quality of AI-generated content	Survey items measuring output accuracy, usefulness, or credibility	Acceptance, use intention	Hallucination and misinformation risk
Behavioural trust (reliance-based)	Willingness to rely inferred from behaviour under uncertainty	Behavioural proxies such as advice-taking or weighting of AI recommendations	Reliance behaviour, compliance	Over-delegation of judgement
Implicit or unmodelled trust	Trust assumed or embedded within adoption factors	Trust not explicitly measured or theorised	Use behaviour, adoption factors	Epistemic blind spots and unexamined reliance

Source: Author's Own Work

Table 4 presents an integrative typology that synthesises how trust has been conceptualised and operationalised across empirical studies on students' voluntary use of generative artificial intelligence. Rather than cataloguing individual studies, the typology organises trust orientations based on their underlying conceptual basis, modes of inference, and associated outcomes. The table highlights that trust is most commonly treated as a cognitively oriented evaluation of system reliability or output credibility, while behavioural and implicit forms of trust receive comparatively less explicit attention. Importantly, each trust orientation is associated with distinct epistemic risks, underscoring that trust in generative artificial intelligence is not uniformly beneficial and must be understood in relation to how reliance is enacted.

Several studies emphasised trust as a willingness to rely on AI systems under conditions of uncertainty, aligning with classical conceptualisations of trust as a psychological state that enables vulnerability. In contrast, a smaller number of studies did not articulate a formal definition of trust. In these cases, trust was operationalised implicitly through related constructs such as perceived credibility, confidence in system outputs, or reliance on AI generated information, without an accompanying conceptual explanation.

Overall, trust was predominantly conceptualised as a cognitively oriented construct focusing on system reliability and information quality. Affective, relational, or multidimensional conceptualisations of trust were largely absent from the reviewed studies. Table 5 summarises how trust was defined and conceptualised across the included studies.

Table 5: Characteristics of Included Studies

Study	Year	Context	Population	Research Design	Focus of Study
Rahman et al.	2023	Higher education	University students	Quantitative survey	Adoption of ChatGPT
Jo	2024	Higher education	University students	Quantitative survey	Use of GenAI in learning
Amaro et al.	2024	Higher education	University students	Quantitative survey	Trust and GenAI use
Shahzad et al.	2024	Higher education	University students	Quantitative survey	Adoption intention
Song	2025	Higher education	University students	Quantitative survey	Trust in ChatGPT
Abdalla	2025	Higher education	University students	Quantitative survey	GenAI acceptance
Al Amin et al.	2025	Higher education	University students	Quantitative survey	Trust as mediator
Guo & Erdenebold	2025	Higher education	University students	Quantitative survey	Trust and usage behaviour
Al-Dmour et al.	2025	Higher education	University students	Quantitative survey	AI adoption factors
Kiyak et al.	2025	Higher education	University students	Experimental study	Trust and advice-taking
Subhani et al.	2025	Higher education	University students	Quantitative survey	Adoption and price value

Source: Author's Own Work

Definitions and Conceptualisations of Trust

The positioning of trust within empirical models varied considerably across the reviewed studies, indicating a lack of consensus regarding its analytical role in explaining students' voluntary use of generative artificial intelligence. Trust appeared in multiple positions within empirical models, including as an antecedent, a mediating mechanism, a moderating variable, an outcome, or remaining unmodelled altogether.

In several studies, trust was positioned as an antecedent predicting adoption related outcomes such as behavioural intention, acceptance, or use behaviour. Other studies treated trust as a mediating mechanism that transmitted the effects of antecedent factors to behavioural outcomes. In a smaller number of cases, trust was incorporated as a moderator influencing the strength of relationships between predictors and outcomes. Some studies did not explicitly model trust but were retained because they provided indirect evidence of trust related gaps through their focus on adoption factors without corresponding trust constructs.

This descriptive variation in trust positioning highlights the fragmented manner in which trust has been integrated into empirical models of GenAI use. The distribution of trust positioning across studies is summarised in Table 6 and visually represented in the Trust Positioning Map.

Table 6: Positioning of Trust within Empirical Models

Study	Explicit Definition of Trust	Source or Basis of Definition	Conceptual Focus
Rahman et al. (2023)	Yes	Adapted from prior trust literature	System reliability
Jo (2024)	No	Implicit (credibility-related items)	Output credibility
Amaro et al. (2024)	Yes	Trust in AI literature	Reliability of AI output
Shahzad et al. (2024)	Yes	Adapted trust definition	Trust in system
Song (2025)	Yes	Single-item trust definition	Confidence in ChatGPT
Abdalla (2025)	Yes	Prior technology trust studies	System dependability
Al Amin et al. (2025)	Yes	Trust as psychological mechanism	Willingness to rely
Guo & Erdenebold (2025)	Yes	Trust in AI framework	Trust in AI output
Al-Dmour et al. (2025)	Yes	Adapted technology trust	System trust
Kiyak et al. (2025)	No	Behavioural proxy (advice-taking)	Reliance behaviour
Subhani et al. (2025)	No	Trust not explicitly defined	Adoption-related factors

Source: Author's Own Work

Types and Dimensions of Trust Identified

The reviewed studies reported a range of behavioural and psychological outcomes associated with trust in generative artificial intelligence. Behavioural outcomes were the most frequently examined and included behavioural intention, actual use behaviour, adoption behaviour, and continued use of GenAI tools. In studies where trust was positioned as an antecedent or mediator, trust was associated with these outcomes either directly or indirectly within empirical models.

Psychological outcomes related to trust were reported less frequently. These outcomes included acceptance of GenAI tools, confidence in AI generated outputs, and reliance related behaviour. In experimental contexts, trust was sometimes inferred through behavioural proxies such as advice taking behaviour rather than measured as an explicit self-reported construct.

Overall, the outcomes associated with trust primarily reflected students' behavioural engagement with GenAI technologies, while psychological dimensions such as reliance and acceptance received comparatively less empirical attention. Studies that did not explicitly model trust nonetheless contributed to identifying areas where trust related mechanisms were

absent or underrepresented in adoption focused research. Table 7 summarises the behavioural and psychological outcomes associated with trust across the included studies.

Table 7: Outcomes Associated with Trust in Included Studies

Study	Trust Type or Label Used	Operationalisation	Measurement Approach	Notes
Rahman et al. (2023)	System trust	Single construct	Self-reported scale	Focus on reliability
Jo (2024)	Output credibility	Single construct	Self-reported items	Trust implicit
Amaro et al. (2024)	Trust in AI	Single construct	Survey-based	Output reliability
Shahzad et al. (2024)	Technology trust	Single construct	Self-reported scale	System-level trust
Song (2025)	Trust in ChatGPT	Single item	Self-reported item	Confidence in output
Abdalla (2025)	System trust	Single construct	Survey-based	Dependability focus
Al Amin et al. (2025)	Techno-trust	Single construct	Survey-based	Psychological mechanism
Guo & Erdenebold (2025)	Trust in AI output	Single construct	Survey-based	Information trust
Al-Dmour et al. (2025)	Technology trust	Single construct	Survey-based	System trust
Kiyak et al. (2025)	Behavioural trust proxy	Not applicable	Behavioural measure	Advice-taking
Subhani et al. (2025)	Not specified	Not specified	Trust not measured	Trust-gap evidence

Source: Author's Own Work

Summary of Descriptive Patterns

Several descriptive patterns emerged from the synthesis of the included studies. First, trust was conceptualised inconsistently, with variation in whether it was explicitly defined, implicitly operationalised, or inferred through behavioural indicators. Second, trust was predominantly treated as a cognitively oriented construct focused on system reliability and output credibility, with limited attention to affective or multidimensional aspects.

Third, the positioning of trust within empirical models varied widely, appearing as an antecedent, mediator, moderator, outcome, or remaining unmodelled. No dominant or standardised positioning of trust was observed across studies. Fourth, trust was most frequently associated with behavioural outcomes related to adoption and use, while psychological outcomes such as reliance and acceptance were less frequently examined and sometimes inferred rather than directly measured.

All findings reported in this section are descriptive and aim to provide a structured overview of how trust has been conceptualised, positioned, and associated with outcomes in empirical studies of students' voluntary use of generative artificial intelligence.

Discussion

This systematic literature review provides a diagnostic synthesis of how trust has been conceptualised, positioned, and operationalised in empirical studies examining students' voluntary use of generative artificial intelligence in educational contexts. Rather than reiterating descriptive findings, this discussion advances a theoretical interpretation of trust as a psychological mechanism that enables reliance under conditions of uncertainty, opacity, and epistemic risk. The discussion is organised around three interrelated themes, namely trust as a psychological reliance mechanism, the interaction between trust, voluntariness, and academic integrity, and directions for future research.

Trust as a Psychological Reliance Mechanism

The findings of this review indicate that trust in the context of generative artificial intelligence is best understood as a psychological mechanism that enables students to rely on system outputs when verification is incomplete, and outcomes are uncertain. Across the reviewed studies, trust was predominantly conceptualised in cognitive terms, focusing on perceptions of system reliability, output credibility, or dependability. While such conceptualisations are consistent with established trust literature, the way trust operates in voluntary GenAI use suggests important qualitative distinctions in reliance behaviour.

First, some forms of reliance can be characterised as blind reliance, where students accept AI generated outputs with minimal scrutiny or critical evaluation. In these cases, reliance is not grounded in calibrated judgement but in overconfidence, convenience, or the persuasive fluency of AI responses. Second, convenience driven use reflects a more instrumental orientation, where students rely on GenAI primarily to reduce effort, save time, or simplify academic tasks, without necessarily forming a stable trust judgement about the system itself. Third, warranted or calibrated trust involves a more reflective reliance, where students engage in epistemic checking, contextual evaluation, and selective use of AI outputs based on task demands and perceived risk.

These distinctions are important because they challenge the implicit assumption that higher levels of trust are inherently desirable. In the context of generative artificial intelligence, trust that is poorly calibrated may lead to inappropriate reliance, over delegation of cognitive tasks, or uncritical acceptance of inaccurate information. The reviewed literature, however, rarely differentiates between these forms of reliance, often treating trust as a unidimensional construct. This lack of differentiation obscures the psychological function of trust and limits the ability of empirical models to capture how students actually decide when, how, and to what extent to rely on GenAI systems.

Trust, Voluntariness, and Academic Integrity

The role of trust becomes particularly salient in educational contexts characterised by voluntary use. Unlike institutionally mandated learning technologies, generative artificial intelligence tools are adopted through self-initiated decisions, often outside formal curricular structures or

explicit guidance. In such contexts, trust functions as a prerequisite for use rather than an outcome of prolonged interaction. Students must decide whether to rely on GenAI outputs despite awareness of algorithmic opacity, hallucination risks, and variable output quality.

This reliance decision is not merely technical or instrumental, but also moral and cognitive in nature. Academic integrity concerns, including plagiarism, over delegation of learning tasks, and inappropriate substitution of student effort, introduce an additional layer of judgement. Trust in GenAI therefore intersects with epistemic vigilance, which refers to users' capacity to critically evaluate information sources and assess their legitimacy. When epistemic vigilance is weak, trust may slide into blind reliance, increasing the risk of misconduct or superficial learning. When epistemic vigilance is strong, trust becomes conditional and selective, supporting responsible and context appropriate use.

The reviewed studies suggest that trust in GenAI is often embedded within broader judgements about legitimacy and appropriateness, rather than being reducible to perceived usefulness or ease of use. However, empirical models rarely make this moral cognitive dimension explicit. As a result, trust is frequently treated as a functional belief about system performance, rather than as a judgement that mediates the relationship between technological capability, personal responsibility, and academic norms. Recognising this dimension is essential for understanding why students may continue to rely on GenAI even when they are aware of its limitations and risks.

Future Research Agenda

The synthesis of existing evidence highlights several directions for future research that can advance theoretical clarity without prematurely proposing new models or frameworks. First, future studies should move beyond static measures of trust and examine trust calibration, that is, how students adjust their reliance on GenAI in response to experience, feedback, and task complexity. Understanding whether trust becomes more warranted or more blind over time is critical in educational settings.

Second, there is a need for greater use of experimental and behavioural research designs. While survey-based studies dominate the current literature, behavioural measures such as advice taking, verification behaviour, or error detection can provide more direct insight into how trust shapes actual reliance. Such designs are particularly well suited to capturing the dynamic and situational nature of trust in generative systems.

Third, longitudinal approaches are needed to examine how patterns of reliance evolve over extended periods of GenAI use. Most existing studies capture trust at a single point in time, which limits understanding of how repeated interaction, institutional norms, and assessment practices shape students' trust judgements. Longitudinal evidence would allow researchers to distinguish between temporary convenience driven use and more stable forms of calibrated reliance.

Taken together, these directions point towards a research agenda that treats trust not as a static attitudinal variable, but as a dynamic psychological mechanism embedded within voluntary, uncertain, and normatively charged educational contexts.

Limitation

Several limitations should be considered when interpreting the findings of this systematic literature review. First, the core synthesis was deliberately restricted to empirical studies involving student-only populations within higher education and TVET contexts. This boundary was established to maintain conceptual clarity and population validity in examining trust-related mechanisms in students' voluntary use of generative artificial intelligence. While this focus strengthens the internal coherence of the review, it limits the transferability of the findings to other stakeholder groups such as educators, institutional leaders, or policymakers. Studies involving mixed populations were therefore not included in the core synthesis and were considered only for contextual reference.

Second, the reviewed literature exhibited substantial heterogeneity in how trust was defined, operationalised, and measured. Although all included studies addressed trust-related mechanisms, several relied on implicit conceptualisations, proxy indicators, or narrowly focused measures centred on system reliability or output credibility. This variation constrained direct comparability across studies and precluded quantitative synthesis. As a result, the review necessarily adopted a descriptive and mapping-oriented approach rather than effect size comparison or meta-analytic integration.

Third, the available evidence was unevenly distributed across outcome domains. Most empirical studies focused on adoption-related outcomes such as behavioural intention or initial use, whereas relatively fewer examined post-adoption phenomena including sustained reliance, calibration over time, or behavioural adjustment. Consequently, the findings offer stronger insights into the role of trust in early-stage engagement with generative artificial intelligence than in longer-term or repeated use contexts. This imbalance reflects prevailing patterns in the primary literature rather than limitations of the review design.

Fourth, the synthesis was constrained by reporting limitations in the primary studies. Several articles did not fully specify contextual conditions, details of the generative artificial intelligence tools examined, or the precise analytical role of trust within their models. In line with the review protocol, unclear or missing information was treated conservatively to avoid inferential overreach. Nevertheless, incomplete reporting restricted deeper examination of contextual variation and model-level differentiation in trust mechanisms.

Finally, caution is warranted when generalising the findings across diverse educational settings. Although the included studies spanned multiple regions, most relied on cross-sectional designs and self-reported measures. Trust in generative artificial intelligence is likely to be sensitive to contextual, cultural, and temporal factors, including institutional norms, assessment practices, and evolving exposure to AI tools. These dynamics may not be fully captured within the current evidence base. Taken together, these limitations underscore the need to interpret the findings within the scope of a diagnostic synthesis aimed at mapping early empirical patterns rather than establishing definitive conclusions about trust dynamics in generative artificial intelligence.

Conclusion and Implications

This systematic literature review examined how trust has been conceptualised, positioned, and operationalised in empirical studies of students' voluntary use of generative artificial intelligence in educational contexts. Taken together, the evidence indicates that trust plays a

pivotal role in enabling students to rely on generative AI systems under conditions of uncertainty, opacity, and epistemic risk. However, the review also demonstrates that trust is frequently treated as a simplified or auxiliary construct, often reduced to perceptions of reliability or credibility, and inconsistently integrated within empirical models. These patterns suggest that the central challenge in educational contexts is not to maximise trust in generative artificial intelligence, but to cultivate forms of reliance that are contextually appropriate, reflective, and epistemically warranted.

From an educational perspective, indiscriminate or uncalibrated trust in generative artificial intelligence carries significant risks, including over-delegation of cognitive tasks, diminished epistemic vigilance, and potential erosion of academic integrity. The findings of this review therefore point towards a normative shift in how trust is understood in education. Rather than being framed as an outcome to be increased, trust should be regarded as a psychological judgement that governs when, how, and to what extent students choose to rely on AI-generated outputs. The educational objective, in this sense, is warranted reliance, where students engage with generative artificial intelligence in a manner that is selective, critical, and aligned with academic values, rather than blind acceptance or convenience-driven use.

Looking ahead, future research on trust in generative artificial intelligence should prioritise the examination of trust calibration processes. This includes understanding how students form, adjust, and regulate trust judgements across different tasks, contexts, and levels of uncertainty. Greater attention should be given to decision-making under uncertainty, particularly how students balance efficiency, accuracy, and ethical considerations when engaging with AI-generated information. Such work would benefit from moving beyond static attitudinal measures towards designs that capture judgement, verification behaviour, and reliance decisions in situ.

In addition, there is a need for more behavioural and longitudinal research approaches that can illuminate how patterns of reliance evolve over time. Experimental designs, behavioural indicators, and repeated-measures studies offer promising avenues for examining how trust develops through interaction, feedback, and experience with generative artificial intelligence. By focusing on calibration, judgement, and reliance rather than simple adoption outcomes, future research can deepen understanding of trust as a dynamic psychological mechanism embedded within voluntary and normatively charged educational environments.

-
- Acknowledgements:** Not applicable.
- Funding Statement:** This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.
- Conflict of Interest Statement:** The authors declare that there is no conflict of interest regarding the publication of this paper. All authors have contributed to this work and approved the final version of the manuscript for submission.
- Ethics Statement:** This study did not involve any human participants, animals, or sensitive data requiring ethical approval. The authors confirm that the research was conducted in accordance with accepted academic integrity and ethical publishing standards.
- Author Contribution Statement:** Mohd Badli Ramli conceived and designed the study, established the conceptual and analytical framework, conducted the systematic review procedures, performed the data screening and synthesis, and led the writing of the manuscript.
Suhaizal Hashim provided academic supervision, strengthened the theoretical positioning of the study, and offered strategic guidance to enhance the scholarly depth of the discussion.
Mohd Zulfadli contributed methodological expertise, reviewed the research design and analytical rigor, and provided critical input to ensure the robustness and validity of the review process.
Aida Aisyah supported manuscript preparation through language refinement, formatting alignment with journal requirements, reference verification, and final proofreading prior to submission.
All authors reviewed, revised, and approved the final version of the manuscript.
-

References

- Abdalla, A. M. (2025). Students' acceptance of generative artificial intelligence in higher education: Examining the role of trust. *Education and Information Technologies*, 30(2), 1–20.
- Al Amin, M., Hossain, M. A., & Islam, M. S. (2025). Techno-trust as a mediator in students' adoption of generative AI tools in higher education. *Computers & Education: Artificial Intelligence*, 6, 100201. <https://doi.org/10.1016/j.caeai.2024.100201>.
- Al-Dmour, H., Alshurideh, M., & Al Kurdi, B. (2025). Determinants of artificial intelligence adoption in higher education: The role of technology trust. *Education and Information Technologies*, 30(1), 1–24.
- Amaro, J., Costa, C., & Reis, L. P. (2024). Trust in generative artificial intelligence: Students' perceptions of ChatGPT in academic contexts. *IEEE Transactions on Learning Technologies*, 17(1), 45–57. <https://doi.org/10.1109/TLT.2023.3324567>.
- Cotton, D. R. E., Cotton, P. A., & Shipway, J. R. (2023). Chatting and cheating: Ensuring academic integrity in the era of ChatGPT. *Innovations in Education and Teaching International*, 60(6), 1–12. <https://doi.org/10.1080/14703297.2023.2190148>.
- Dwivedi, Y. K., Hughes, D. L., Ismagilova, E., Aarts, G., Coombs, C., Crick, T., ... Williams, M. D. (2023). So what if ChatGPT wrote it? Multidisciplinary perspectives on opportunities, challenges and implications of generative conversational AI for research, practice and policy. *International Journal of Information Management*, 71, 102642. <https://doi.org/10.1016/j.ijinfomgt.2023.102642>.
- Gefen, D., Karahanna, E., & Straub, D. W. (2003). Trust and TAM in online shopping: An integrated model. *MIS Quarterly*, 27(1), 51–90. <https://doi.org/10.2307/30036519>.
- Glikson, E., & Woolley, A. W. (2020). Human trust in artificial intelligence: Review of empirical research. *Academy of Management Annals*, 14(2), 627–660. <https://doi.org/10.5465/annals.2018.0057>.
- Guo, Y., & Erdenebold, T. (2025). Trust in AI output and students' continuance intention to use generative AI in higher education. *Computers & Education*, 198, 104750.
- Hoff, K. A., & Bashir, M. (2015). Trust in automation: Integrating empirical evidence on factors that influence trust. *Human Factors*, 57(3), 407–434. <https://doi.org/10.1177/0018720814547570>.
- Jo, H. (2024). University students' use of generative AI for learning: Perceptions of credibility and usage behaviour. *Education Sciences*, 14(2), 155. <https://doi.org/10.3390/educsci14020155>.
- Kasneci, E., Sessler, K., Küchemann, S., Bannert, M., Dementieva, D., Fischer, F., Kasneci, G. (2023). ChatGPT for good? On opportunities and challenges of large language models for education. *Learning and Individual Differences*, 103, 102274. <https://doi.org/10.1016/j.lindif.2023.102274>.
- Kıyak, M., Wiese, E., & Hancock, P. A. (2025). Advice-taking from artificial intelligence: Behavioural indicators of trust in generative systems. *Human-Computer Interaction*, 40(1), 1–32.
- Lankton, N. K., McKnight, D. H., & Tripp, J. F. (2015). Technology, humanness, and trust: Rethinking trust in technology. *Journal of the Association for Information Systems*, 16(10), 880–918. <https://doi.org/10.17705/1jais.00411>.
- Lee, J. D., & See, K. A. (2004). Trust in automation: Designing for appropriate reliance. *Human Factors*, 46(1), 50–80. <https://doi.org/10.1518/hfes.46.1.50.30392>

- Mayer, R. C., Davis, J. H., & Schoorman, F. D. (1995). An integrative model of organizational trust. *Academy of Management Review*, 20(3), 709–734. <https://doi.org/10.5465/amr.1995.9508080335>.
- McKnight, D. H., Choudhury, V., & Kacmar, C. (2002). Developing and validating trust measures for e-commerce: An integrative typology. *Information Systems Research*, 13(3), 334–359. <https://doi.org/10.1287/isre.13.3.334.81>.
- Page, M. J., McKenzie, J. E., Bossuyt, P. M., Boutron, I., Hoffmann, T. C., Mulrow, C. D., ... Moher, D. (2021). The PRISMA 2020 statement: An updated guideline for reporting systematic reviews. *BMJ*, 372, n71. <https://doi.org/10.1136/bmj.n71>.
- Rahman, M. M., Islam, M. A., & Rahman, M. S. (2023). Determinants of students' intention to use ChatGPT in higher education: The role of trust. *Education and Information Technologies*, 28(5), 1–23.
- Shahzad, F., Xiu, G., Wang, J., & Shahbaz, M. (2024). Exploring students' adoption of generative AI tools: The mediating role of technology trust. *Interactive Learning Environments*, 32(4), 1–17.
- Song, Y. (2025). Students' trust in ChatGPT and its influence on acceptance in higher education. *Computers & Education: Artificial Intelligence*, 6, 100198. <https://doi.org/10.1016/j.caeai.2024.100198>.
- Subhani, M. I., Hasan, S. M., & Mehmood, A. (2025). Factors influencing students' adoption of generative AI tools: Evidence from higher education. *Education and Information Technologies*, 30(2), 1–2.