Assessing students' DRIVE:

An evidence-based framework to evaluate learning through students' interactions with generative AI

Manuel Oliveira, Carlos Zednik, Gunter Bombaerts, Bert Sadowski, and Rianne Conijn

Department of Industrial Engineering and Innovation Sciences,

Eindhoven University of Technology

Corresponding author: Manuel Oliveira (<u>m.j.barbosa.de.oliveira@tue.nl</u>), Human-Technology Interaction, Eindhoven University of Technology, Den Dolech, Eindhoven, 5200MB, The Netherlands.

Preprint Notice

This manuscript is a preprint and has not yet undergone peer review. As such, the content, findings, and conclusions presented herein **should be cited and interpreted with caution.** This version of the manuscript may be subject to substantial revisions following the peer-review process. This preprint will be submitted to a specialty journal for formal peer review.

Funding

This research was co-funded by the 4TU.Centre for Engineering Education and BOOST! (Eindhoven University of Technology).

Acknowledgements

We are immensely grateful to all the dedicated research assistants, teachers, and teaching assistants [anonymized] who collected, scored, and annotated the data for this project.

CRediT statement

The following contributions are defined according to the <u>CRediT</u> (Contributor Roles Taxonomy):

- Conceptualization: MO, CZ
- Data Curation: MO
- Formal Analysis: MO
- Funding Acquisition: MO, RC, CZ, GB, BS
- Investigation: MO
- Methodology: MO
- Project Administration: MO
- Resources: MO, CZ, RC
- Supervision: MO, RC, CZ
- Validation: MO, RC, CZ
- Visualization: MO
- Writing Original Draft: MO
- Writing Review & Editing: MO, RC, CZ, GB

Abstract

As generative AI (GenAI) transforms how students learn and work, higher education must rethink its assessment strategies. This paper presents a taxonomy and conceptual framework (DRIVE) to evaluate student learning from GenAl interactions (prompting strategies), focusing on cognitive engagement (Directive Reasoning Interaction) and knowledge infusion (Visible Expertise). Despite extensive research mapping student GenAI writing behaviors, practical tools for assessing domain-specific learning remain underexplored. This paper shows how GenAl interactions inform such learning in authentic classroom contexts, moving beyond technical skills or low-stakes assignments. We conducted multi-methods analysis of GenAI interaction annotations (n=1450) from graded essays (n=70) in STEM writing courses. A strong positive correlation was found between high-quality GenAl interactions and final essay scores, validating the feasibility of this assessment approach. Furthermore, our taxonomy revealed distinct interaction profiles: High essay scores correlated with a "Targeted Improvement Partnership" focused on text refinement, while high interaction scores were linked to a "Collaborative Intellectual Partnership" centered on idea development. Conversely, below-average performance was associated with "Basic Information Retrieval" or "Passive Task Delegation". These findings demonstrate how the assessment method (output vs. process focus) may shape students' GenAl usage and learning depth. These findings demonstrate that the assessment method (output vs. process) shapes student AI use. Traditional assessment can reinforce text optimization, while process-focused evaluation may reward the exploratory partnership crucial for deeper learning. The DRIVE framework and related taxonomy offers educators a practical tool to design assessments that capture authentic learning in Al-integrated classrooms.

Keywords: Learning, Assessment, Academic Writing , Generative AI

1.Introduction

The emergence of generative artificial intelligence (GenAI) in higher education has fundamentally disrupted traditional methods of assessing student learning and raised questions about whether learning objectives should change (e.g., Bower et al., 2024; Xia et al., 2024). This is especially true in contexts where assessment focuses on text production given the ability of current GenAI applications to easily produce academic texts increasingly indistinguishable from human-generated work (Casal & Kessler, 2023; Clark et al., 2021; Fleckenstein et al., 2024; Porter & Machery, 2024). As students increasingly engage in dialogues with GenAI applications like ChatGPT (OpenAI, 2022) to develop their academic work (e.g., Ansari et al., 2024), conventional assessment approaches that evaluate only written outputs cannot effectively measure the acquisition of writing skills (e.g., Swiecki et al., 2022; Yan, 2023) or domain-specific knowledge. This fundamental shift in how academic work is produced demands a reimagining of assessment practices, as GenAI's rapid integration into education has rendered traditional evaluation methods ineffective. Consequently, scholarly attention is shifting from evaluating the final product to analyzing the interaction process, which offers a more transparent record of students' engagement and reasoning process (e.g., Swiecki et al., 2022).

Despite increasing research on GenAl's educational impact, a gap remains in understanding how these tools mediate learning. Early research began by documenting and interpreting behavioral patterns of how students interact with GenAI chatbots in contexts often removed from classrooms (e.g., Cheng et al., 2024; Pigg, 2024). These foundational studies provided initial frameworks for categorizing interactions, such as requesting, refining, and evaluating content (Pigg, 2024), or distinguishing between knowledge telling and knowledge transformation based on how students modify GenAl suggestions (Cheng et al., 2024). More recently, studies have begun to move beyond mere description to investigate how these interaction patterns might serve as cues to student learning, increasingly utilizing data from experimental academic tasks. For instance, by analyzing screen recordings of doctoral students in a controlled writing task, Nguyen et al. (2024) found that an iterative, highly interactive collaboration pattern was associated with higher writing performance, while a more linear, supplementary use was linked to lower scores. In a similar vein, Kim and collaborators (2025) conducted an experimental study where they analyzed logs of interactions between students and GenAI during an academic writing task devised specifically for the study (i.e., outside the classroom context). Next, they classified the verbs in these prompts using Bloom's Taxonomy (Anderson & Krathwohl, 2001), a well-known educational framework that categorizes cognitive skills along a spectrum ranging from lower-order (Remembering, Understanding) to higher-order ones (Applying, Analyzing, Evaluating, Creating). Based on these classifications, Kim et al (2025) identified distinct interaction patterns associated with different levels of AI literacy: students with high AI literacy used descriptive, context-rich prompts across various taxonomy levels, engaged collaboratively with AI, and viewed it as a tool for idea development. In contrast, those with low AI literacy relied on general prompts focused primarily on lower-order thinking skills, had briefer interactions, and used the AI mainly for content generation rather than engaging in complex cognitive processes. These differences in interaction patterns were linked to writing performance, with high AI literacy students achieving significantly higher scores in content, structure, and expression, and demonstrating more effective modification and synthesis of Al-generated content.

Although existing studies provide valuable insights into the relationship between student-Al interactions, cognitive processes, and learning outcomes within controlled settings, it remains unclear if the identified patterns of interaction with GenAl also emerge in authentic classroom contexts, where the decision to use AI is directly tied to student's grades. Moreover, the focus of this type of research has been on describing and finding structure in interaction patterns. However, there is a pressing need for studies focused on translating these findings into practical, evidence-based tools that higher education teachers can directly use to analyze student interaction logs with AI (i.e., prompts and respective AI outputs) and assess the evidence of writing skill acquisition. Developing such tools is crucial for leveraging the insights from interaction analysis to inform pedagogical practice and assessment in the age of GenAl.Responding to the need for practical assessment tools, this paper introduces and validates a taxonomy for analyzing student-GenAI interactions in authentic classroom settings. Our taxonomy is guided by a new conceptual framework we developed to evaluate the interaction process itself. As detailed later in this introduction, this framework is built on two core principles: first, assessing the degree to which the student actively and purposefully steers the dialogue with the AI, and second, evaluating the extent to which the student makes their own unique knowledge and ideas observable within that dialogue. To validate this approach, we test whether these interaction patterns correlate with traditional learning outcomes, namely essay scores, providing initial evidence for their use as learning proxies. Our methodology focuses on academic writing in general, with an emphasis on argumentative writing. This form of writing requires students to develop a debatable thesis, support it with logical evidence, and anticipate counterarguments (Toulmin, 1958), which in turn encompasses both skills that GenAI can readily replicate (e.g., text generation, basic argumentation) and struggles with (e.g., critical evaluation of self-generated content, integration of personal understanding of generated content), considering how GenAI systems exhibit significant limitations in comprehending their own outputs (West et al., 2023). By systematically categorizing and analyzing how students engage with GenAl throughout their writing and argumentative process, from initially prompting the system to critical evaluation and revision, we can identify the types of interactions that are associated with evidence of learning. Our proposed taxonomic approach aims to enable educators to further develop evidence-based assessment methods that remain relevant in an Al-integrated environment.

2. Background

2.1. The skill of argumentative writing

To contextualize the development of the taxonomic framework, we must first consider the nature of the academic skill it aims to evaluate: argumentative writing. Argumentative writing represents a foundational academic skill that extends beyond mere text composition to also involving critical thinking, evaluation of evidence, and logical reasoning (Andrews, 2015; Newell et al., 2011). Traditional assessment of argumentative writing has focused on evaluating the final product of a student's assignment (i.e., an essay) often according to a grading rubric designed by the teacher, which typically focuses on examining structural elements, coherence, use of evidence, and logical progression of arguments (Ferretti & Graham, 2019). However, the integration of GenAl into the writing process calls for innovative approaches to both instruction and assessment that consider how students leverage these tools in developing their argumentative competencies.

The literature on argumentative writing assessment has identified several key dimensions worth revisiting in the current discussion. Toulmin's (1958) model of argumentation, which identifies claims (i.e., statement the writer wants to improve), warrants (i.e. logical/persuasive connection between claim and evidence), backing (i.e., evidence supporting claim), and rebuttals(e.g., acknowledging alternative viewpoints) as essential components, has

informed numerous assessment frameworks (Erduran et al., 2004; Sampson & Clark, 2008). More recent approaches have expanded these frameworks to incorporate evaluations of source integration (Wingate, 2012), and the acknowledgement and integration of different perspectives in the argumentative process (Nussbaum & Schraw, 2007; Wolfe et al., 2009). These established assessment criteria provide a theoretical foundation for understanding the quality of argumentative writing, but are not yet able to account for the collaborative process that emerges when students engage with GenAl tools.

Research on technology-enhanced writing instruction has demonstrated that digital tools can support different phases of the writing process (Little et al., 2018; Zhang & and Zou, 2022). However, studies examining the specific impact of GenAl on argumentative writing remain limited. Initial investigations have documented students' utilization of GenAl for writing assignments (e.g., Kim et al., 2025) but, to the best of our knowledge, few studies have systematically analyzed how different patterns of GenAl interaction correlate with learning outcomes in the specific domain of argumentative writing.

2.2. Evidence of learning in the age of GenAl

Several theoretical educational frameworks have been guiding educators' understanding of teaching and learning over the past decades. A widely cited view by Marton and Säljö (1976) distinguishes between surface learning, focused on rote memorization, and deep learning, which involves actively seeking meaning, integrating new knowledge, and transforming understanding. Complementing this, the widely adopted Bloom's Taxonomy (Anderson & Krathwohl, 2001), known for its shift from noun categories to verb forms representing cognitive processes, provides a hierarchical structure for categorizing cognitive skills. This hierarchy ascends from lower-order thinking skills such as Remembering (recalling facts and basic concepts) and 'Understanding (explaining ideas or concepts), to higher-order thinking skills like Applying (using information in new situations), Analyzing (drawing connections among ideas, breaking material into constituent parts), Evaluating (justifying a stand or decision, critiquing), and Creating (producing new or original work). Educators often use these levels to design learning objectives and assessments (e.g., Britto & Usman, 2015). Evidence of learning is often inferred from a student's ability to demonstrate skills at the higher end of the taxonomy. For instance, an essay that not only recalled information but also analyzed different perspectives and created a new synthesis would be seen as indicative of "deeper" learning and more sophisticated cognitive processing. These frameworks have historically guided the assessment of student work, often focusing on the final product as the primary evidence of these cognitive processes. However, the advent of GenAI, which can generate sophisticated outputs that mimic human-like understanding and skill, requires a shift in focus. When students collaborate with GenAI, the final product alone offers an increasingly ambiguous signal of their learning, as it becomes challenging to disentangle the student's contribution from the Al's. One potential approach to circumvent this challenge might involve searching for learning evidence in the interaction process between a writer and AI systems, through the examination of how students steer these systems, how they evaluate their output, or decide to incorporate it in their writing. Existing frameworks, primarily designed to evaluate individual student output, may not adequately capture these nuanced interaction strategies or reveal the depth of student agency and critical engagement within the AI-assisted writing process.

3. DRIVE framework

As grounding for our current taxonomic framework, we provide a more nuanced lens for understanding learning in this new paradigm, we propose a conceptual framework: Directive Reasoning and Interaction, and Visible Expertise, or DRIVE. The primary purpose of this framework is to provide guidance on how to assess evidence of academic writing skill acquisition through the systematic examination of the interaction process between a student and a GenAI system during AI-assisted writing. It specifically seeks to identify behaviors indicative of what the student knows and how their actions produce visible evidence of skill acquisition, such as understanding domain-relevant theories, engaging in critical thinking, and introducing original, user-generated ideas into the dialogue with the system. At its core, the DRIVE framework posits that a crucial step in assessing AI-assisted writing processes is to observe the extent to which students actively and purposefully steer the interaction with GenAI, thereby making their learning, knowledge, and critical thinking visible. The DRIVE framework groups these observable markers of learning under two distinct categories, introduced below.

3.1.Directive Reasoning Interaction (DRI)

This component evaluates how actively and purposefully the student steers the interaction with the AI. It echoes the ideas of heutagogy, a framework of self-determined learning (Hase & Kenyon, 2007). Heutagogy is concerned with "learner-centred learning that sees the learner as the major agent in their own learning, which occurs as a result of personal experiences" (Hase & Kenyon, 2007, p. 112). In this model, the teacher (or AI) facilitates learning by providing scaffolding throughout the process, while the learner maintains ownership of their learning path. Framing the student-AI interaction through the lens of heutagogy allows us to conceptualize GenAl as a powerful resource that a self-determined learner can direct. This perspective also aligns with the Interactive, Constructive, Active, and Passive (ICAP) framework by Chi and Wylie (2014) which categorizes learning activities from shallow to deep. While passive engagement involves receiving information without active processing, and active engagement applies existing knowledge for retention, deeper learning arises from constructive and interactive engagement. Constructive engagement involves generating new ideas and outputs beyond learned material, enhancing problem-solving and transversal skills. Interactive engagement, the deepest form, entails collaborative idea generation, leading to novel inferences while fostering communication and collaboration skills. Crucially, in the context of Al-assisted writing, a student's high Directive Reasoning Interaction (DRI), characterized by taking a leading role, critically questioning AI outputs, and using their own reasoning to guide the dialogue, serves as tangible evidence of these deeper, more purposeful forms of engagement and self-determined agency. Essentially, high DRI means the student is more in command of the collaboration. DRI also aligns well with the principle of "active human agency", or the empowered capacity for a user to critically assess AI output and take steps to adjust it (Fanni et al., 2023; see also Lyons et al., 2021). This directive stance is not only crucial for maintaining a "human-in-command" approach but also serves as a crucial cognitive safeguard. Through the engagement in reasoning and intentional steering of the interaction, students can counter the negative effects of automation bias (i.e., tendency to uncritically accept Al-generated information) and mitigate the risks of skill atrophy associated with cognitive offloading, through which a person reduces cognitive effort by delegating a task to AI (e.g., see Gerlich, 2025; Wahn et al., 2023). A strong DRI profile can thus be understood as an observable proxy for a student's ability to maintain cognitive and ethical control in the collaborative process.

3.2.Visible Expertise (VE)

This component focuses on the extent to which the student makes their own knowledge, original ideas, and understanding visible within the interaction log. This concept resonates with earlier research-based pedagogical frameworks such as "Making Thinking Visible" from Harvard's Project Zero, which argues that for thinking to be truly understood, directed, and assessed, it must first be made observable to others (Ritchhart, 2011). In GenAl-assisted writing, visible expertise encompasses the demonstration of declarative and procedural knowledge and skills. This includes the application of domain-specific knowledge and crucial transversal skills, such as critical thinking, problem-solving, and adaptability. Furthermore, with the rise of GenAI, AI literacy, specifically the skills required to effectively and critically evaluate AI system outputs, is an increasingly important aspect of visible expertise that informs interaction patterns. When student prompts introduce specific course concepts, apply unique insights, or build upon pre-existing ideas with AI, they make their intellectual contribution and authorial voice evident. This demonstration of expertise is important because, as GenAl transforms learning, the ability to discern, critically engage, and contribute original thinking retains its essential value. Given GenAI's known limitations in reasoning ability and comprehending context, and its potential to produce unverified or biased content (Amirizaniani et al., 2024; e.g., Bender et al., 2021; Maleki et al., 2024; Shojaee et al., 2025), visible expertise also involves the capacity to critically assess and refine AI outputs, thereby countering risks like automation bias and cognitive offloading. VE directly addresses the fundamental challenge of evaluating student learning in GenAI-assisted assignments. For fair and effective educational assessment, teachers must clearly discern students' unique intellectual contributions within the interaction. This visibility offers a window into the student's learning process, allowing for an assessment of skill development that would otherwise be obscured in a final product (e.g., essay). In the classroom context, transparency is crucial for accountability and trust. Observing how students shape their interaction with GenAI over time allows teachers to more effectively evaluate their growth in light of the intended learning objectives (e.g., see Swiecki et al., 2022), especially when these take the technology into account.

This framework thus suggests that interaction patterns aligning with high "DRI" and "VE" are indicative of desirable profiles for using GenAI when writing argumentative essays, and potentially in other complex academic tasks. The more these aspects are visible in the interaction logs, the richer the evidence of learning available for assessment in the AI-assisted co-writing process. Conversely, interactions exhibiting limited DRI and VE would offer less tangible evidence of the student's active learning and skill acquisition through their engagement with the AI. Ultimately, the DRIVE framework is intended as a useful conceptual and analytical tool for educators. It aims to support the assessment of writing in a manner compatible with GenAI use in the classroom by shifting attention to the quality of the student's interaction and display their learning and knowledge throughout the process. This framework, therefore, underpins the development of our taxonomy tool, which aims to capture the desirable aspects encapsulated by the DRI and VE constructs, through the identification of observable interaction patterns that cue at these learning indicators in the context of argumentative essay writing.

4. Overview

4.1. Research Aims

This paper presents the development of a practice-oriented taxonomy for analyzing student-GenAl interactions, which is grounded in the DRIVE framework. Our taxonomy aims to identify strategies of engagement with GenAl technology and explore whether they can provide a meaningful window into student learning during academic writing. The present research is primarily exploratory and descriptive, and is guided by two central questions detailed below.

RQ1: How does a process-focused assessment of GenAl interaction quality relate to a traditional, output-focused assessment of essay quality?

This question seeks to validate our process-focused measure against traditional essay scores. A significant positive association would provide initial evidence that analyzing the interaction process is a valid approach for assessing student learning.

RQ2: What student-GenAl interaction patterns are associated with different levels of mastery, and do these patterns diverge depending on how mastery is measured?

This question uses our taxonomy to investigate the specific interaction types associated with mastery indicators. It is divided into two parts:

• RQ2a: How do GenAl interaction strategies connect with different levels of mastery based on traditional essay evaluations and GenAl interaction evaluations?

Here, we aim to identify which taxonomy classifications are associated with above-average versus below-average mastery on each measure. We expect that interaction types associated with higher mastery on both measures will reflect greater student agency over the technology and more visible integration of their own knowledge (core principles of DRIVE).

• RQ2b: To what extent do the GenAl interaction patterns associated with different mastery levels overlap between the two assessment methods (traditional essay evaluation vs. GenAl interaction evaluation)?

This is a purely exploratory follow-up question. We have no specific hypothesis about the outcome. The goal is to investigate the degree to which the two assessment types (grading the final essay vs. grading the interaction process) are sensitive to the same, or different, types of student-GenAI engagement.

To address these questions, we analyze student-GenAl interaction logs and essay mastery data (i.e., grading scores) from university courses where Al-assisted writing was a graded component. By examining how students use GenAl for real coursework, we aim to provide initial evidence for the utility of the DRIVE framework and its associated taxonomy in understanding learning in Al-integrated settings.

5. Methodology

5.1. Overview

This study employs a multi-faceted approach to investigate whether student-GenAl interactions can serve as a meaningful proxy for learning in argumentative writing. Central to our methodology is the development and application of a new taxonomy designed to systematically classify student prompts from real-world classroom settings. We collected both the final written outputs (essays) and the process data (GenAl interaction logs). Student performance was then assessed using two distinct measures: a traditional, output-focused essay score and a novel, process-focused GenAl interaction quality score. Our analysis proceeded in two phases. For RQ1, we correlated the process- and output-focused performance scores to validate the former. For RQ2, we identified the interaction patterns characteristic of different mastery tiers on each measure (RQ2a) and then conducted an exploratory comparison to see if both assessment types prioritize the same patterns of GenAl engagement (RQ2b).

Table 1

Sample Descriptives

Course/Year	AI	Non-Al	Unknown	Total	Annotated	Total
(Academic	Users	Users	Al Use	Students	Essays	Annotations
Degree					(Al Users)	(Al Users)
Data Science Ethics 2023-2024 (BSc)	32 (21.2%)	119 (78.8%)	0 (0%)	151	21	369
Philosophy & Ethics AI 2023-2024 (MSc)	17 (12.9%)	106 (80.3%)	9 (6.8%)	132	16	309
Philosophy & Ethics AI 2024-2025 (MSc)	54 (33.3%)	107 (66.0%)	1 (0.6%)	162	33	772
Total	103 (23%)	332 (74.7%)	10 (2.3%)	445	70	1450

Prompt Statistics

(Annotated Essays Only, N = 70)

Measure	Prompts per Student	Prompt Length (characters)
Mean (SD)	20.71 (18.41)	505 (1026)
Median (IQR)	14.5 (16.75)	168 (390)
Min - Max	2 - 103	2 – 9828

Note. Percentages represent proportion within each course. Annotated essays represent the subset of AI user essays that underwent detailed interaction analysis.

5.2. Context and Participants

This research was conducted at a STEM university within three Bachelor or Master's level courses on philosophy and ethics, covering topics from human-technology interaction to the societal impact of artificial intelligence. In all courses, students were required to individually write a graded argumentative essay. Data were collected across these courses during the 2023-2024 and 2024-2025 academic years. As detailed in Table 1, a total of 445 students were enrolled across these courses. Of these, 103 students (23.2%) chose to use GenAI for their assignments under the condition that they would submit their interaction logs for assessment. The shared interaction log was formally graded using a marking rubric (see Table 2) and contributed to their final course grade¹. A subset of 70 Al-user essays, along with their corresponding GenAl interaction logs, were annotated using the proposed taxonomy (Appendix). A total of 1450 student-GenAl interactions (i.e. prompts) were annotated. The discrepancy between the 103 students who opted to use GenAI and the 70 essays that were annotated results from many interaction logs from AI users being unusable due to issues encountered during data collection and processing. Examples include broken hyperlinks to ChatGPT interaction logs shared by students, or messy screenshots of chat interactions that were difficult to incorporate into the dataset and were ultimately excluded. The "Unknown AI Use" category in Table 1 refers to cases with insufficient information regarding AI tool engagement. Among the 70 annotated AI user essays, ChatGPT was the most predominantly used GenAI tool (n = 48, 68.6%). One student (1.4%) used the chatbot Claude, and 21 students (30.0%) did not report their GenAl tool. The prompt-related statistics for the 70 annotated essays, including number of prompts per student (as derived from their interaction logs) and prompt length (as derived from number of characters in prompts), are summarized in Table 1.

¹ It should be noted that, according to a teacher (also a co-author), students were initially more willing to experiment with GenAI before this graded evaluation component was formally implemented in the course. The teacher noticed a visible decrease in the number of students deciding to use GenAI for their assignment after the decision to formally assess GenAI use.

Table 2

Evaluation Criteria for GenAl Interaction Logs in Argumentative Essay Writing

Criterion	Excellent (10-9)	Good (8-7)	Sufficient (6)	Insufficient (5-0)
AI for Writing	Prompts are clearly formatted and go far beyond the basic parameters of the assignment description, revealing expert-level mastery of using AI as a writing aid.	Prompts are clearly formatted and go considerably beyond the basic parameters of the assignment description, revealing considerable technical ability of using AI as a writing aid.	Prompts are clearly formatted and go beyond the basic parameters of the assignment description, revealing the basic ability of using Al as a writing aid.	No prompts provided, or prompts unclearly formatted. No visible effort to engineer prompts that go beyond the basic parameters of the assignment description.
Al for Argument ation	Extensive critical engagement of Al-generated content. Prompts reveal expert-level use of Al to improve argumentative structure.	Critical engagement of Al-generated content. Prompts reveal considerable efforts to use Al to improve argumentative structure.	Limited critical engagement of Al-generated content. Prompts reveal some effort to use Al to improve argumentative structure.	No critical engagement with Al-generated content. No meaningful effort to use Al to improve argumentative structure.
Al for Course Content	Prompts used to perform extensive content-related research. Prompts reveal deep and broad understanding of, and engagement with, the course material, at times going beyond that material.	Prompts used to perform considerable content-related research. Prompts reveal understanding of and engagement with the course material without going beyond that material.	Prompts used to perform some content-related research. Prompts reveal limited understanding of, or engagement with, the course material.	Prompts used insufficiently for content-related research. Prompts reveal no meaningful understanding of, or engagement with, the course material.

5.3. Data Collection Procedure

Over a 10-week period in each course, students completed a graded argumentative essay assignment. Students were informed that the use of GenAl tools (e.g., ChatGPT) was optional for their essay writing process, encompassing stages such as planning, researching, drafting, or

refining arguments. A condition for using GenAI was the submission of complete interaction logs (sequences of input prompts and AI outputs). To mitigate potential disparities in GenAI proficiency, all participating courses included at least one lecture on argumentative writing and basic techniques for using GenAI chatbots effectively, commonly referred to as prompt engineering. Scores reflecting traditional essay grades and experimental overall evaluations of student-GenAI interactions were collected. All data, including interaction logs, essays, and evaluation scores, were collected following informed consent from participating students and ethical approval granted by the Ethical Review Board of [anonymized]. Data were anonymized and stored securely for research purposes.

5.4. Course Learning Objectives

Across the courses included in this research, students are expected to develop the ability to critically engage with ethical, societal, and philosophical questions related to data science and artificial intelligence. A central learning objective is the capacity to construct well-reasoned, evidence-based arguments in written form. Students learn to identify and evaluate ethical and philosophical arguments, apply major ethical theories to contemporary technological contexts, and analyze value-laden concepts relevant to data-driven practices. They are also trained to read and critically interpret scholarly texts and to use research tools to investigate ongoing societal debates. Argumentative essay writing serves, thus, as a core integrative task through which students demonstrate their ability to synthesize conceptual understanding, ethical reasoning, and domain-specific analysis.

5.5. Measures

Two primary types of measures were used to assess student performance: traditional essay scores and GenAI interaction quality scores.

5.5.1. Traditional essay scores

Student essays were evaluated by course instructors using grading rubrics tailored to argumentative writing within the specific course contexts (Data Science Ethics or Philosophy & Ethics of AI). These scores represent an output-focused measure of performance, reflecting the quality of the final written product. Core assessment criteria that were common across these rubrics, independent of any AI tool usage, included: ability to define and contextualize an ethical problem or case relevant to the course; depth of ethical analysis and the construction of well-structured, coherent, and persuasive arguments; demonstration of critical thinking and reflection on complexities and diverse perspectives; effective integration of course concepts and adherence to academic writing style.

5.5.2. GenAl interaction quality scores

The quality of students' interactions with GenAI was assessed experimentally by course teachers and teaching assistants as part of the final course grade for students who chose to use GenAI in their essay assignments, using a set of criteria designed to evaluate GenAI use during argumentative essay writing, with an emphasis on the identification of learning indicators. These criteria are detailed in Table 2 and cover aspects such as AI for Writing, AI for Argumentation, and AI for Course Content. The criteria are aligned with the proposed taxonomy (see Appendix). They integrate course learning objectives and teachers' views of interaction quality. Although

these views can be subjective, the criteria link to our DRIVE framework by focusing on agentic cognitive engagement (DRI), seen in students steering prompts and critically revising AI output, and visible knowledge integration (VE), seen in students drawing on and developing their own disciplinary ideas during interaction with the AI. Overall, this score represents a more process-focused measure of performance, compared to the more final output-focused essay scores.

6. Development of the taxonomy

The taxonomy for classifying student-GenAl interactions was developed through an iterative, dual-approach process, carried out by two university teachers (one of whom is a co-author) and teaching assistants, all with expertise in argumentative essay writing. The top-down component of this process was firmly grounded in the intended learning outcomes of the courses under examination, which emphasize the ability to construct well-reasoned arguments, apply ethical reasoning, and critically engage with complex issues. To reflect these goals, the taxonomy was designed to identify how argumentative writing skills manifest in GenAl-supported processes. It organizes interaction patterns into three main categories aligned with core academic competencies: Writing, Content, and Argument. These dimensions were selected because they resonate with established components of argumentative writing.

- "Writing" encapsulates interactions focusing on the mechanical and structural aspects of essay composition, including task-oriented actions like providing instruction, requesting content formatting, or requesting assistance to improve and organize specific sections (e.g., introduction, conclusion).
- "Content" captures interactions that center on knowledge construction and understanding, including actions such as requesting definitions, examples, or theoretical explanations, with a particular emphasis on course-specific material and critical engagement with Al-generated output.
- "Argument" encompasses interactions that specifically target logical and analytical aspects of writing such as interactions that develop and further refine argumentative elements (e.g., identifying different perspectives involved in a given discussion, improving articulation of arguments, strengthening one's thesis).

This top-down, pedagogically-informed structure was complemented and refined by a bottom-up analysis. This involved a hands-on review of actual student GenAl interaction logs, allowing the development team to gain insights into common, real-world user actions and patterns. This iterative process (illustrated in Figure 1) resulted in the final version of the taxonomy, which contains a total of 35 subcategories within three main categories: 13 subcategories under Writing, 10 under Content, and 12 under Argument. The full taxonomy can be found in the Appendix.



Figure 1. Illustration of the cycle of taxonomy development and refinement.

The current taxonomy is substantively informed by an analytical approach to human-AI collaboration patterns focused on student agency and knowledge construction, an approach that is further detailed in the present paper. Crucially, the taxonomy was also firmly grounded in the intended learning outcomes of the courses under examination. These outcomes emphasize students' ability to construct well-reasoned arguments, apply ethical and philosophical reasoning, and critically engage with complex societal issues. To reflect these goals, the taxonomy was designed to identify potential indicators of learning by rethinking how argumentative writing skills might manifest in GenAI-supported writing processes. It organizes interaction patterns into three main categories aligned with core academic competencies: Writing (mechanical and structural elements), Content (knowledge construction and information management), and Argument (logical reasoning and analytical thinking). These dimensions were selected because they resonate with established components of argumentative competence discussed in the literature. For example, the Writing category captures interactions related to textual coherence, structure, and clarity, which are fundamental to conveying an argument effectively, akin to the structural elements often evaluated in traditional rubrics. The Content category addresses how students engage with the substance of their arguments, including the

sourcing, evaluation, and integration of information and evidence (echoing aspects like Toulmin's 'backing' or the critical use of sources), a process that takes on new dimensions when information is co-constructed with GenAI. Finally, the Argument category directly targets interactions indicative of logical reasoning, the formulation and support of claims, the consideration of counterarguments or rebuttals (central to Toulmin's model and subsequent frameworks), and the overall analytical thinking involved in building a persuasive case.

7. Annotation of student-GenAI interactions

Across all courses, 103 out of 445 (23%) essays were (reported to have been) co-written with GenAI. A total of 70GenAI interaction logs, associated with the respective amount of graded essays, were annotated. For the remaining 33 cases the interaction logs were sometimes missing (e.g., broken hyperlinks to ChatGPT interaction logs, or missing files), or included a negligible amount of interactions focusing mainly on a few rephrasing requests (e.g., less than five minimally informative interactions). A total of four different annotators annotated the interaction logs. Annotators were instructed to classify each student prompt (i.e., their input) with the best fitting taxonomy item(s). To accommodate interactions that could be described by more than one item, annotators were free to decide whether to classify an interaction with a single or multiple category-subcategory items (e.g., Writing_Instructions and Content_Research). Interactions classified with multiple taxonomy items are referred to as 'Mixed' in our results.

To assess the reliability of the annotation, a subset of interaction logs (n=33, 772 annotations) of one of the three courses (Philosophy & Ethics 2024-2025) was annotated by three additional raters. Because there was a common second rater to three different raters, we computed the Cohen's Kappa metric of inter-rater agreement for each pair of raters. The average Cohen's Kappa was 0.44 (SD = 0.06), which according to the interpretation guidelines proposed by Landis and Koch (Landis & Koch, 1977), reflect moderate agreement (note this value is at the boundary between the "moderate" and "fair" levels of agreement proposed by these authors). It should be noted that the inter-rater reliability differed between the categories within the taxonomy. At the main category level, agreement was consistently moderate with an indication of higher agreement for Content classifications (ranging from 0.65 for Content and 0.57 for Writing, to 0.46 for Argument, all Kappa values with *ps* < .001). The agreement at the taxonomy subcategory levels (see Figure 2) was more heterogeneous with some classifications achieving very low agreement (e.g., writing_autoimprove, argument_improve, content_concept) and other very high (e.g., writing_introduction, content_bibliography, argument_objection). In general, however, these data suggest fair and higher agreement for most classifications.

Inter-rater agreement per taxonomy classification



Course: Philosophy & Ethics of Al 2024-2025

Figure 2. Inter-rater Agreement per Taxonomy Classifications.

8. Data Analysis

Data processing and analyses were conducted using R v.4.3.3 (R Core Team, 2024). Scripts of the analyses are available at the project's repository at Open Science Framework (<u>https://osf.io/32jq7/?view_only=feed50a5bad04bfab4f5bd60531510e7</u>).

To allow for the comparability across different course cohorts and grading scales, both traditional essay scores and GenAI interaction evaluation scores were standardized into z-scores within each course subset. A z-score of zero thus corresponded to the average score within the context of a specific course, while negative or positive z-scores quantified how much an individual score was below or above that course average, respectively. This normalization process accounted for the existing heterogeneity in scoring scale ranges across the courses.

To investigate RQ1, which is concerned with the relationship between traditional essay assessment and the experimental assessment of GenAI interaction quality, we calculated the

correlation between these measures using the z-scores associated with all the annotated essays (N=70). Specifically, we calculated both a Pearson product-moment correlation (*r*) and a Spearman rank correlation (*rho*). Additionally, reporting Spearman's rho is particularly useful for classroom-based data as it is a non-parametric measure that is less sensitive to outliers or non-normally distributed data, which are common characteristics of real-world educational datasets.

To address RQ2, which investigates whether the developed taxonomy can uncover patterns of student AI interaction associated with different levels of mastery, we analyzed both essay performance and GenAI interaction quality. In this study, we define "mastery" as a construct representing skill proficiency in two distinct dimensions:

- Essay mastery refers to the demonstrated proficiency in academic writing as reflected in the final essay scores, evaluated based on traditional essay writing quality criteria. They indirectly capture how successfully students incorporated content from GenAI interactions into a coherent academic argument.
- GenAl interaction mastery refers to demonstrated proficiency in productive engagement with generative AI tools, as assessed by expert graders using interaction quality criteria (Table 2). These criteria derive from the proposed DRIVE framework's concepts of Directive Reasoning Interaction and Visible Expertise, which emphasize strategic questioning, critical evaluation of AI outputs, and effective guidance of the AI system toward writing assignment-related goals.

The taxonomy descriptives were calculated to gain a sense of the most prevalent classifications in our sample of annotated interaction logs. Classifications with a prevalence below 1% were deemed practically irrelevant and were excluded from the analyses of RQ2a and 2b, as their interpretation within the context of the present RQs is less relevant, and these have negligible impact over the results (see online data materials for more details).

For RQ2a, we calculated the mean z-score and 95% confidence interval (CI) for each taxonomy classification across all annotated interactions with an overall prevalence above 1%. This allows us to identify which interaction types were associated with different mastery levels based on whether the 95% CI around the mean z-score was entirely above zero (Above Average mastery), entirely below zero (Below Average mastery), or included zero (Average mastery). This approach accounts for the uncertainty in our estimates and ensures that mastery level classifications are supported by sufficient statistical evidence. A z-score of zero represents the average mastery within each course context (as it was calculated within each classroom's sample), thus providing a meaningful reference point for interpreting mastery associations. We then developed qualitative profiles of GenAI interaction patterns by interpreting the taxonomy classifications most strongly associated with each mastery level through analysis of their mean z-scores, 95% CIs, and theoretical connections to the DRIVE framework.

For RQ2b, we examined whether both assessment methods were sensitive to the same interaction patterns or prioritized different GenAl usage strategies. We employed a dual analytical approach: first examining the degree of overlap between 95% CIs of mean z-scores for each taxonomy classification as an initial proxy for agreement between assessment approaches. Non-overlapping confidence intervals indicate potential disagreement between methods, while overlapping intervals suggest agreement but do not definitively rule out statistically significant differences. To address this limitation, we conducted exploratory paired t-tests comparing essay and GenAl interaction z-scores for each taxonomy classification, as both measures derive from identical classification observations. We applied a false discovery

rate (FDR; Benjamini & Hochberg, 1995) correction across all 22 comparisons to control for multiple testing. Additionally, we calculated Cohen's *d* effect sizes with 95% CIs to assess the practical significance of any detected differences. This approach allowed us to distinguish between cases where assessment methods truly converge versus those where subtle but meaningful systematic differences exist despite overlapping confidence intervals.

9. Results

9.1.RQ1: Relationship between traditional essay evaluations and GenAl interaction evaluations

For the 70 annotated essays, a Pearson product-moment correlation indicated a statistically significant, strong positive linear relationship between traditional essay assessment scores (output-focused) and GenAl interaction quality scores (process-focused) (r = 0.54, 95% CI [0.34, 0.68], t(68) = 5.24, p < .001). This suggests that students who demonstrated higher quality interactions with GenAl also tended to achieve higher traditional essay scores. A scatterplot illustrating this relationship is provided in Figure 3. This alignment between the two types of learning indicators (output-focused essay scores and process-focused GenAl interaction evaluations) lends support to the potential of GenAl interaction evaluations to provide insights into student learning, at least in the same capacity as essay scores allow for.



Correlation GenAl Interaction and Essay Z-Scores

Figure 3. Relationship Between Traditional Essay Scores and GenAl Interaction Evaluation Scores.

It should be noted that while the essay z-score distribution met the normality assumption, the GenAl interaction z-score distribution marginally failed the Shapiro-Wilk normality test (W = 0.965, p = .047). As an additional check, a Spearman's rank correlation was calculated to confirm the relationship using a non-parametric test. This yielded an identical result (rho = 0.54, p < .001).

9.2. RQ2: What student-GenAl interaction patterns are prevalent across different levels of mastery, and do these patterns diverge depending on how mastery is measured?

We first describe the overall pattern of taxonomy classifications before focusing on the descriptives per mastery level.

9.2.1. Taxonomy descriptives

The frequencies at which taxonomy classifications were observed during the annotation of student-GenAI interaction logs collected from the three courses are shown in Figure 4, both at the main taxonomy category level (Figure 4-A) and at the subcategory level (Figure 4-B, all above 1% frequency). The overall pattern for the main categories indicates that the most prevalent category of interactions relate to Writing aspects (41.3%), followed by Content (28.7%) and Argument (22.3%). A total of 7.7% of the interactions were annotated with more than one category, categorized as "Mixed".

Within the Writing category, Writing_Improve (improving spelling, style or grammar of input text) was the most prominent subcategory, accounting for 13.4% of the total interactions, followed by Writing_Evaluate (requesting evaluation of essay section; 7%) and Writing_Miscellaneous (prompting system in a non-specific technical way, 4.8%). It should be noted that the subcategory Writing_Miscellaneous is a "catch-all" classification, and in that sense, its underrepresentation (or overrepresentation) in the results may be interpreted as desirable (or undesirable), as it hints at interactions hard to classify with the current taxonomy content.

For Content, the most common subcategory was Content_Research (asking AI to define ideas or find related ideas to user's input; 5.6%), Content_Bibliography (asking for references, 5.2%), followed by Content_Elaboration (requesting additional detail incorporating course content, 4.6%) and Content_Idea (elaborating on existing well-formulated ideas, 4.1%).

Finally, for Argument, Argument_Improve (improving the structure given argument, 5.9%) was most common, followed by Argument_Objection (providing an objection for a given argument 4.5%) and Argument_Justify (requesting AI to provide reasons for an input claim, 3.4%).





Figure 4. Overall Descriptives of Taxonomy Annotations for All Courses.

9.2.2. RQ2a: How do GenAl interaction strategies connect with different levels of mastery based on traditional essay evaluations and GenAl interaction evaluations?

The following analyses examine how interaction types connect with mastery levels across both traditional essay quality and GenAI interaction quality assessments, revealing how different GenAI usage patterns relate to performance under output-focused versus process-focused evaluation approaches. Figure 5 shows the mean z-scores (+ 95% CIs) by taxonomy classification for both essay scores (in blue) and GenAI interaction scores (in red). Confidence intervals including zero (z-score) reflect average mastery levels, while intervals entirely below or above zero reflect below-average or above-average mastery levels, respectively. This confidence interval approach provides statistically rigorous classification by ensuring that mastery level designations are supported by sufficient evidence rather than chance variation.



Mean Z-Scores with 95% Cls by Taxonomy Classification

Figure 5. GenAl Interaction Classifications And Mastery Level: Essay and GenAl Interaction Mean Z-Scores + 95% Confidence Intervals Per Taxonomy Classification

9.2.2.1. Essay z-scores and taxonomy classifications

Above-average essay mastery was associated with a "Targeted Improvement Partnership" approach, characterized by three distinct but complementary student-GenAI interaction strategies. Writing_Improve dominated this profile (n = 194 or 13.4% of annotations, mean z = 0.13, 95% CI [0.03, 0.24]), reflecting actions such as the systematic refinement of spelling, style, and grammar in existing text. This was complemented by sophisticated analytical engagement through Content_Critical interactions (n = 23 or 1.6% of annotations, mean z = 0.25, 95% CI [0.03, 0.47]), where students critically engaged with AI-generated content by asking for clarifications or corrections. This profile was further defined by Argument_Relate interactions (n = 27 or 1.9% of annotations, mean z = 0.36, 95% CI [0.13, 0.59]), which involved requests to connect or relate two concepts or ideas. This set of strategies suggests that students who achieved higher essay scores engaged GenAI as a targeted text improvement tool, by systematically improving their input work (essay sections) through (inferred) critical evaluation and conceptual integration rather than by seeking comprehensive assistance from GenAI.

Below-average essay mastery was characterized by a "Basic Information Retrieval" prompting strategy, including only two interaction types with z-score confidence intervals entirely below zero (average). Content_Research showed the strongest negative relationship (n = or 5.6% of annotations, mean z = -0.41, 95% CI [-0.72, -0.11]), involving requests for AI to define ideas or identify related concepts. Content_Example interactions also demonstrated negative associations (2.6%, mean z = -0.18, 95% CI [-0.35, -0.01]), where students asked for specific examples of general cases or issues. This constrained profile suggests that students with lower essay performance primarily used AI for foundational information gathering rather than sophisticated content development or critical engagement.

The predominance of interactions categorized as average (77% of classifications) suggests that most GenAl usage patterns neither significantly enhanced nor detracted from essay writing quality as traditionally assessed (i.e., output focus). This pattern emphasizes the specificity of interaction types that correlate with essay performance and suggests that only a few types of prompting strategies (as identified by the current taxonomy) appear to be connected with very high and very low writing quality as assessed traditionally/

Relating back to the DRIVE framework, the above-average profile demonstrates a moderate display of Directive Reasoning Interaction (DRI) through the apparent targeted steering of the AI toward specific essay improvement tasks. The pattern also suggests an emerging Visible Expertise (VE) as inferred from critical evaluation of AI output, or the requests for assisting with conceptual integration within the essay's narrative. By contrast, the below-average profile shows less evidence of DRI, with interactions focused primarily on information extraction (vs. a more collaborative development of the essay), and minimal VE, as these prompts sought more basic or foundational definitional support (vs. demonstrating original thinking or knowledge synthesis through the usage of GenAI).

9.2.2.2.GenAl interaction z-scores and taxonomy classifications

Above-average GenAl interaction mastery was associated with a "Collaborative Intellectual Partnership" approach, characterized by four interaction strategies that demonstrate an engagement with (Gen)Al as a thinking partner/assistant. Argument_ConceptualClarity emerged as the strongest positive indicator (n = 19 or 1.3% of annotations, mean z = 0.76, 95% CI [0.52, 1.00]), involving requests to simplify or improve the definition of concepts. This was

complemented by Argument_Relate interactions (n = 27 or 1.9% of annotations, mean z = 0.49, 95% CI [0.27, 0.72]), where students asked AI to connect or relate two concepts or ideas in the course of the essay writing process. Content_Idea interactions formed a substantial component of this profile (n = 60 or 4.1% of annotations, mean z = 0.39, 95% CI [0.21, 0.56]), where students brought their own well-motivated original ideas or questions to the AI and requested confirmation, elaboration, or discussion of these concepts (assumedly generated outside of the dialogue, likely by the student themself). This profile is further characterized by Content_Critical interactions (n = 23 or 1.6% of annotations, mean z = 0.35, 95% CI [0.01, 0.68]), where students critically engage with AI-generated content by asking for clarifications or corrections of the target content (e.g., AI output, student input, or a hybrid content). This combination of strategies suggests that students with higher GenAI interaction scores engaged AI as an intellectual collaborator, leveraging the technology for conceptual refinement, knowledge synthesis, and critical dialogue.

Below-average GenAl interaction mastery was characterized by a "Passive Task Delegation" approach, which included only one interaction type. Writing_Instructions demonstrated the sole negative association (n = 44 or 3.0% of annotations, mean z = -0.37, 95% CI [-0.60, -0.14]), involving specifications of tasks in terms of course assignment descriptions, typically through copy-pasting or uploading assignment instructions. This singular profile suggests that students with lower GenAl interaction scores primarily used Al as a direct recipient of student input rather than engaging in collaborative knowledge construction or strategic dialogue. This may be hinting at lower levels of confidence or trust in the capabilities of the Al system, although that remains an open question that cannot be addressed by the current data.

The overwhelming prevalence of average-classified interactions (82% of classifications) indicates that most GenAI usage patterns demonstrated neither exceptional mastery nor deficiency when evaluated against the DRIVE framework's process-focused criteria. This finding highlights the distinctiveness of interaction types that correlate with high or low quality GenAI engagement and suggests that effective collaborative partnership with GenAI requires specific strategic approaches rather than simply general usage competency.

Through the lenses of the DRIVE framework, the above-average profile shows a strong Directive Reasoning Interaction (DRI) demonstrated through the (inferred) strategic steering toward conceptual development and knowledge integration. This was coupled with a clearer display of Visible Expertise (VE) through actions demonstrating original idea contribution and critical evaluation of AI outputs. This pattern suggests a behavioral profile where students engage with GenAI as an intellectual collaboration rather than treating it as a mere tool. In contrast, the below-average profile demonstrates minimal DRI, with interactions focused on task specification rather than strategic guidance, and negligible VE, as these actions only show the ability to provide instructions to the system without any signs of user knowledge incorporation, knowledge synthesis, or critical engagement with AI throughout the collaborative process.

9.2.3.RQ2b: To what extent do the GenAl interaction patterns associated with different mastery levels overlap between the two assessment methods (traditional essay evaluation vs. GenAl interaction evaluation)?

Confidence interval overlap analysis revealed substantial convergence between assessment methods, with 21 of 22 taxonomy classifications (95.5%) demonstrating overlapping CIs. Only Content_Idea showed clear disagreement, with essay evaluation classifying it as below average (95% CI [-0.32, 0.19]) while GenAI interaction evaluation rated it as above average (95% CI [0.21, 0.56]). This high level of agreement aligned closely with the strong positive correlation (r = 0.54) between assessment methods identified in RQ1. However, this overlap analysis provides

a conservative test that may miss statistically meaningful differences when intervals overlap but distributions differ significantly. To explore this possibility, we conducted paired t-tests comparing essay and GenAI interaction z-scores for each taxonomy classification, applying FDR correction across all 22 comparisons to control for multiple testing.

This exploratory statistical analysis uncovered a more nuanced picture, suggesting additional classifications with significant differences (FDR-corrected). Beyond the already-identified Content Idea (p < .001, d = -0.50, 95% CI [-0.72, -0.28]), four additional disagreements emerged. Argument ConceptualClarity demonstrated the largest effect (p = .004, d = -0.71, 95% CI [-1.09, -0.33]), followed by Writing Miscellaneous (p = .001, d = -0.40, 95% CI [-0.61, -0.20]), Writing Instructions (p = .028, d = 0.46, 95% CI [0.13, 0.80]), and Content Research (p = .013, d = -0.31, 95% CI [-0.50, -0.11]). The pattern of disagreements suggests a slight degree of systematic assessment differences in terms of what they may indirectly incentivize through their evaluation focus. Process-focused GenAI interaction evaluation assigned substantially higher scores to conceptualization-related work (Argument ConceptualClarity, Content Idea) and flexible AI engagement or diversity of prompts (Writing Miscellaneous). By contrast, output-focused essay evaluation showed a relative preference for structured task specification (Writing Instructions) and compensatory information-seeking (Content Research). Of note, 17 of 22 classifications (77.3%) demonstrated negligible effect sizes, indicating that most interaction patterns receive fundamentally similar evaluations across both methods.

This divergence pattern, despite small, suggests that traditional essay assessment may undervalue exploratory behaviors in GenAI interactions that process-focused evaluation rewards as cues to effective student-GenAI collaboration, while simultaneously undervaluing certain foundational interaction patterns that contribute to final product quality. The statistically significant disagreements suggest a small tension between optimizing output quality versus rewarding a more sophisticated engagement with GenAI, which may eventually translate into practical implications for how assessment design shapes student AI usage patterns in educational contexts.

10. Discussion

The primary aim of this research was to address the evolving challenge of assessing student learning in AI-integrated writing environments by investigating the utility of analyzing student interactions with GenAI, particularly within the context of argumentative essay writing. We sought to determine if a novel taxonomy, grounded in the proposed DRIVE (Directive Reasoning Interaction & Visible Expertise) framework, could reveal patterns of student-GenAI engagement associated with mastery levels on both traditional output-focused essay assessments and process-focused evaluations of GenAI interaction quality (RQ2). A central aim in this exploration focused on determining whether we could identify meaningful learning indicators by analyzing student interactions with GenAI tools ((RQ1). This was based on the idea that learning can be visible through how students engage with these systems. Such an approach could either replace traditional product-only assessments or complement them for a more complete view of student learning within AI-infused classrooms.

Our findings lend initial support to this assessment approach. A significant positive relationship was found between traditional essay scores and GenAl interaction quality evaluations (RQ1), suggesting that the assessment of the interaction process itself can provide valuable insights into student performance that align, at least moderately, with established measures. Furthermore, addressing the need for a more nuanced understanding of the quality

tied to how students engage with these tools, the application of our taxonomy (RQ2) allowed us to map distinct profiles of GenAl interaction associated with different mastery tiers on both measures (i.e. essay scores and GenAl interaction evaluations). High-performing students, as judged by either essay quality or GenAl interaction quality, generally exhibited more sophisticated engagement patterns, though, as will be discussed, the specific characteristics of these high-performing interactions varied depending on the assessment focus.

10.1. How GenAI interaction strategies connect with learning indicators

Our results revealed meaningful patterns in how GenAl interaction strategies connect with mastery indicators, with differences that emerged depending on assessment approach. While the strong positive correlation between methods suggests substantial agreement, these systematic differences draw attention to distinct GenAl interaction profiles with important implications for teachers navigating the challenges of writing and domain-specific knowledge assessment in GenAl-infused classrooms.

Traditional essay evaluations favored a "Targeted Improvement Partnership" prompting approach characterized by systematic text refinement, analytical evaluation of AI outputs, and strategic conceptual integration. Interactions connected with below-average essay scores demonstrated a "Basic Information Retrieval" pattern, predominantly seeking definitional support and basic examples, which suggests a compensatory rather than collaborative AI use. For teachers accustomed to evaluating final products, this pattern may feel familiar, yet it raises questions about whether such assessments capture the full scope of student learning in the new GenAl-infused classroom. Process-focused GenAl interaction evaluations revealed a strikingly different pattern. Above-average performance was associated with a "Collaborative Intellectual Partnership" involving sophisticated conceptual refinement, original idea development, strategic connection between concepts, and critical evaluation of AI outputs. Below-average interactions reflected a "Passive Task Delegation" pattern, characterized by basic task specification to the Al, without any signs of engaging with the technology as an intellectual partner. This divergence may present a practical challenge: teachers may have to decide whether to prioritize polished written outputs or evidence of intellectual collaboration with AI, or both. The distinction between using GenAl as a support tool versus an intellectual partner, and between text optimization-focused versus exploration-focused strategies, crucially shapes what cognitive processes we recognize and reward in AI-integrated learning environments. This decision directly influences how students approach AI as a learning tool.

Despite the potentially context-specific nature of these interpretations, our findings align with broader observations in the field. For instance, our sophisticated, co-creative interaction patterns linked to higher mastery mirror Kim et al.'s (2025) findings that higher AI literacy correlates with more diverse and complex interaction strategies. Kim et al. (2025) specifically found that high AI literacy users employed descriptive, context-rich prompts across various cognitive levels and engaged GenAI for idea development, leading to significantly higher writing performance. Conversely, low AI literacy was associated with general, lower-order prompts and primary use for content generation. This suggests that our process-focused GenAI interaction patterns, particularly those reflecting high Directive Reasoning Interaction (DRI) and Visible Expertise (VE), may closely approximate AI literacy distinctions, with our "collaborative intellectual partnership" profile resonating with high AI literacy behaviors. This is further supported by Nguyen et al.'s (2024) work, which found that doctoral students employing iterative, highly interactive strategies with GenAI achieved superior writing performance compared to those using it merely as an information source.

Our work also builds upon existing literature by addressing the assessment of domain-specific learning within prompting interactions, moving beyond a sole focus on technical prompting skills or general interaction descriptions. While existing research has taken important steps in detailing prompt construction (e.g., Chen et al., 2023; Giray, 2023; Heston & Khun, 2023; White et al., 2023) and characterizing interaction patterns with connections to writing outcomes (Cheng et al., 2024; Nguyen et al., 2024; Pigg, 2024), they often treat prompting as a means to achieve a desired output or a reflection of general AI literacy. Our work builds upon this foundation by shifting focus to the assessment of domain-specific learning directly from these interactions. We propose a conceptual framework and a practical tool (taxonomy) designed to guide educators in evaluating evidence of students' conceptual understanding and mastery, particularly within academic argumentative writing. This approach, along with the empirical validation of its utility, aims to bridge the gap between observing how students use GenAI and identifying what substantive learning occurs about the subject matter itself in a GenAI-infused classroom setting.

10.2. Assessment focus shapes GenAl interaction strategies

The systematic differences between assessment methods, though modest in magnitude, illuminate a fundamental tension in AI-integrated education: what cognitive processes do we value and reward? This divergence intensifies a long-standing pedagogical debate about assessing the learning process versus evaluating only the final product (e.g., see Swiecki et al., 2022). Assessment design may inadvertently shape student AI engagement in ways that conflict with educational goals.

The pattern reveals a potential washback effect (Alderson & Wall, 1993) which refers to the influence that tests or assessments exert on teaching and learning. In our context, traditional output-focused assessment may encourage students to treat AI as an efficiency tool for text optimization, while process-focused evaluation promotes intellectual collaboration and exploratory engagement. This divergence suggests that traditional assessments might inadvertently undervalue interaction patterns indicative of deeper learning and knowledge co-construction, which are central to the DRIVE framework we are proposing. Our results suggest that traditional output-focused essay assessments have the potential to induce negative washback in AI-assisted writing, incentivizing students to use GenAI for superficial refinement and task efficiency rather than profound knowledge co-construction or critical engagement. By contrast, process-focused GenAl interaction evaluations may foster positive washback, promoting higher-order skill development by rewarding creative ideation and collaborative intellectual partnership. While these differences are small, they may accumulate over time to fundamentally shape how students develop AI collaboration skills, a capability increasingly critical for lifelong learning. Students who learn to use AI for strategic text refinement develop one set of valuable skills, while those who engage AI as an intellectual partner develop another. However, the latter approach may better prepare students for a future where AI collaboration requires sophisticated reasoning and critical evaluation skills rather than mere task delegation. This resonates with Messick's (1989, 1996) argument for consequential validity: a test's legitimacy must account for its educational consequences, particularly its capacity to drive meaningful learning with GenAl in our Al-infused society.

For educators, assessment choices carry consequences beyond immediate grading decisions, as even modest systematic differences in what interactions are rewarded may signal to students which cognitive approaches are valued, potentially influencing their long-term relationship with AI tools.

10.3. Implications for Educational Practice

Building on our findings, this section provides practical recommendations for teachers integrating GenAl into their teaching practices. While 77.3% of interaction patterns showed negligible differences between assessment approaches, the systematic differences in five interaction types carry important implications for practice. Teachers should adopt a complementary evaluation approach, combining traditional essay assessment with interaction log evaluation to capture both final product quality and intellectual collaboration skills. Given that assessment focus subtly shapes student AI engagement, teachers must proactively scaffold both strategic text refinement and exploratory intellectual partnership as distinct but valuable collaboration approaches. This requires targeted instruction in Directive Reasoning Interaction (DRI) and Visible Expertise (VE) skills to enforce the constructive alignment between course learning objectives, pedagogical activities, and assessment methods (see Biggs, 1996). Importantly, grading rubrics must distinguish between using AI as a basic support tool (e.g., proofreading) and using it for a more dynamic intellectual partnership, so they can signal which cognitive engagement approaches are valued in the course context. Crucially, the grading procedure must ensure fairness regardless of GenAl adoption, so that no student is penalized for using or not using GenAI. Since even modest assessment choices may influence students' long-term AI usage skills, intentional design becomes increasingly critical as these tools become ubiquitous in education. While automated classification may eventually assist with interaction log evaluation (e.g., via prompt engineering and fine-tuning techniques), human oversight remains essential for accurately assessing sophisticated aspects of student-AI collaboration. Given GenAl's rapid evolution and uncertain impact on education, teachers should actively engage with educational research to adapt their practices thoughtfully (e.g., see Bauer et al., 2025).

10.4. Limitations and Future Directions

Despite the utility of the present findings, we can identify several limitations in our work that can inform future research directions. The descriptive and exploratory scope of this work is a primary limitation, resulting from sample and resource constraints that were insufficient for advanced statistical modeling (e.g., latent profile analysis or multilevel regression). The current version of our taxonomy also risks contextual overfitting from its development within specific philosophy courses, although its core categories (Writing, Content, Argument) are broadly applicable to academic writing and its reasoning elements across most disciplines. Future work requires larger datasets and taxonomy refinement by revising categories with low inter-rater agreement. decomposing "catch-all" subcategories, and integrating GenAl literacy competencies (e.g., see Jin et al., 2025). Furthermore, the rapid evolution of GenAI makes findings on transient tool functionalities guickly obsolete. Future studies should, therefore, target durable learning principles, such as the metacognitive skill of discerning when to use AI versus relying on unassisted cognition. Finally, the potential for "meta-prompting" (i.e., fabricating user engagement logs based on AI use evaluation rubrics) is a fundamental threat to GenAI interaction assessment validity. While the high effort involved may deter such malpractice, addressing this risk requires innovative strategies, like incorporating mandatory student reflections on their AI-assisted process (e.g., Nikolic et al., 2023).

11. Conclusion

The increasing integration of GenAI in higher education presents both opportunities and challenges for assessing student learning. Our work offers a novel perspective by moving

beyond evaluating just the final product or general prompting skills. We propose a conceptual framework (DRIVE) and a practical taxonomy that allows educators to discern evidence of domain-specific learning directly from students' interactions with GenAI, particularly within academic argumentative writing. For researchers, this contribution means advancing the understanding of human-GenAI interaction in learning contexts. It shifts the focus from merely observing tool use to identifying how students' evolving prompts and interactive strategies reflect their deepening conceptual understanding. This opens new avenues for studying the intricate cognitive processes involved when GenAI assists (or not) in knowledge construction. For teachers, this work provides a concrete approach to assessing learning in GenAI-compatible classrooms. It offers a way to look beyond concerns of GenAI misuse, instead guiding them to interpret student interactions with GenAI as rich indicators of authentic engagement and mastery, thereby promoting more effective and meaningful human-GenAI educational partnerships.

Declaration of generative AI and AI-assisted technologies in the writing process

During the preparation of this work the author(s) used several Generative AI chatbots (Microsoft Copilot, Gemini, Claude) in order to structure and improve the readability of the text throughout. After using this tool/service, the author(s) critically reviewed and edited the content as needed and take(s) full responsibility for the content of the publication.

References

Alderson, J. C., & Wall, D. (1993). Does washback exist? Applied Linguistics, 14(2), 115–129.

https://doi.org/10.1093/applin/14.2.115

Amirizaniani, M., Martin, E., Sivachenko, M., Mashhadi, A., & Shah, C. (2024). Can LLMs

Reason Like Humans? Assessing Theory of Mind Reasoning in LLMs for Open-Ended

Questions. Proceedings of the 33rd ACM International Conference on Information and

Knowledge Management, 34–44. https://doi.org/10.1145/3627673.3679832

Anderson, L. W., & Krathwohl, D. R. (2001). A taxonomy for Learning, teaching, and assessing:

A revision of Bloom's taxonomy of educational objectives. Addison Wesley Longman.

Andrews, R. (2015). Critical Thinking and/or Argumentation in Higher Education. In M. Davies &

R. Barnett (Eds.), The Palgrave Handbook of Critical Thinking in Higher Education (pp.

49-62). Palgrave Macmillan US. https://doi.org/10.1057/9781137378057_3

Bauer, E., Greiff, S., Graesser, A. C., Scheiter, K., & Sailer, M. (2025). Looking Beyond the

Hype: Understanding the Effects of AI on Learning. *Educational Psychology Review*, 37(2), 45. https://doi.org/10.1007/s10648-025-10020-8

- Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the Dangers of
 Stochastic Parrots: Can Language Models Be Too Big? . Proceedings of the 2021
 ACM Conference on Fairness, Accountability, and Transparency, 610–623.
 https://doi.org/10.1145/3442188.3445922
- Benjamini, Y., & Hochberg, Y. (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society*, *57*(1), 289–300. https://doi.org/10.2307/2346101
- Biggs, J. (1996). Enhancing teaching through constructive alignment. *Higher Education*, *32*(3), 347–364. https://doi.org/10.1007/BF00138871
- Bower, M., Torrington, J., Lai, J. W. M., Petocz, P., & Alfano, M. (2024). How should we change teaching and assessment in response to increasingly powerful generative Artificial Intelligence? Outcomes of the ChatGPT teacher survey. *Education and Information Technologies*, 29(12), 15403–15439. https://doi.org/10.1007/s10639-023-12405-0
- Britto, R., & Usman, M. (2015). Bloom's taxonomy in software engineering education: A systematic mapping study. 2015 IEEE Frontiers in Education Conference (FIE), 1–8. https://doi.org/10.1109/FIE.2015.7344084
- Casal, J. E., & Kessler, M. (2023). Can linguists distinguish between ChatGPT/AI and human writing?: A study of research ethics and academic publishing. *Research Methods in Applied Linguistics*, 2(3), 100068. https://doi.org/10.1016/j.rmal.2023.100068
- Chen, B., Zhang, Z., Langrené, N., & Zhu, S. (2023). Unleashing the potential of prompt engineering in Large Language Models: A comprehensive review (arXiv:2310.14735). arXiv. https://doi.org/10.48550/arXiv.2310.14735
- Cheng, Y., Lyons, K., Chen, G., Gašević, D., & Swiecki, Z. (2024). Evidence-centered assessment for writing with generative AI. *Proceedings of the 14th Learning Analytics*

and Knowledge Conference, 178–188. https://doi.org/10.1145/3636555.3636866

- Chi, M. T. H., & and Wylie, R. (2014). The ICAP Framework: Linking Cognitive Engagement to Active Learning Outcomes. *Educational Psychologist*, *49*(4), 219–243. https://doi.org/10.1080/00461520.2014.965823
- Clark, E., August, T., Serrano, S., Haduong, N., Gururangan, S., & Smith, N. A. (2021). All That's 'Human' Is Not Gold: Evaluating Human Evaluation of Generated Text. In C.
 Zong, F. Xia, W. Li, & R. Navigli (Eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)* (pp. 7282–7296). Association for Computational Linguistics. https://doi.org/10.18653/v1/2021.acl-long.565
- Erduran, S., Simon, S., & Osborne, J. (2004). TAPping into argumentation: Developments in the application of Toulmin's argument pattern for studying science discourse. *Science Education*, *88*(6), 915–933.
- Fanni, R., Steinkogler, V. E., Zampedri, G., & Pierson, J. (2023). Enhancing human agency through redress in Artificial Intelligence Systems. AI & SOCIETY, 38(2), 537–547. https://doi.org/10.1007/s00146-022-01454-7
- Ferretti, R. P., & Graham, S. (2019). Argumentative writing: Theory, assessment, and instruction. *Reading and Writing*, 32(6), 1345–1357. https://doi.org/10.1007/s11145-019-09950-x

Fleckenstein, J., Meyer, J., Jansen, T., Keller, S. D., Köller, O., & Möller, J. (2024). Do teachers spot AI? Evaluating the detectability of AI-generated texts among student essays. *Computers and Education: Artificial Intelligence*, *6*, 100209. https://doi.org/10.1016/j.caeai.2024.100209

- Gerlich, M. (2025). Al tools in society: Impacts on cognitive offloading and the future of critical thinking. *Societies*, *15*(1), Article 1. https://doi.org/10.3390/soc15010006
- Giray, L. (2023). Prompt engineering with ChatGPT: A guide for academic writers. Annals of

Biomedical Engineering, 51(12), 2629–2633.

https://doi.org/10.1007/s10439-023-03272-4

- Hase, S., & Kenyon, C. (2007). Heutagogy: A child of complexity theory. Complicity: An International Journal of Complexity and Education, 4(1). https://doi.org/10.29173/cmplct8766
- Heston, T. F., & Khun, C. (2023). Prompt engineering in medical education. *International Medical Education*, 2(3), Article 3. https://doi.org/10.3390/ime2030019
- Jin, Y., Martinez-Maldonado, R., Gašević, D., & Yan, L. (2025). GLAT: The generative AI literacy assessment test. *Computers and Education: Artificial Intelligence*, 9, 100436. https://doi.org/10.1016/j.caeai.2025.100436
- Kim, J., Yu ,Seongryeong, Lee ,Sang-Soog, & and Detrick, R. (2025). Students' prompt patterns and its effects in AI-assisted academic writing: Focusing on students' level of AI literacy. *Journal of Research on Technology in Education*, 0(0), 1–18. https://doi.org/10.1080/15391523.2025.2456043
- Landis, J. R., & Koch, G. G. (1977). The measurement of observer agreement for categorical data. *Biometrics*, 33(1), 159. https://doi.org/10.2307/2529310
- Little, C. W., Clark, J. C., Tani, N. E., & Connor, C. M. (2018). Improving writing skills through technology-based instruction: A meta-analysis. *Review of Education*, 6(2), 183–201. https://doi.org/10.1002/rev3.3114
- Lyons, H., Velloso, E., & Miller, T. (2021). Conceptualising contestability: Perspectives on contesting algorithmic decisions. *Proc. ACM Hum.-Comput. Interact.*, 5(CSCW1), 106:1-106:25. https://doi.org/10.1145/3449180
- Maleki, N., Padmanabhan, B., & Dutta, K. (2024). AI Hallucinations: A Misnomer Worth Clarifying. 2024 IEEE Conference on Artificial Intelligence (CAI), 133–138. https://doi.org/10.1109/CAI59869.2024.00033

Messick, S. (1989). Meaning and Values in Test Validation: The Science and Ethics of

Assessment. Educational Researcher, 18(2), 5–11.

https://doi.org/10.3102/0013189X018002005

- Messick, S. (1996). Validity and washback in language testing. *Language Testing*, *13*(3), 241–256. https://doi.org/10.1177/026553229601300302
- Newell, G. E., Beach, R., Smith, J., & VanDerHeide, J. (2011). Teaching and learning argumentative reading and writing: A review of research. *Reading Research Quarterly*, 46(3), 273–304. https://doi.org/10.1598/RRQ.46.3.4
- Nguyen, A., Hong ,Yvonne, Dang ,Belle, & and Huang, X. (2024). Human-AI collaboration patterns in AI-assisted academic writing. *Studies in Higher Education*, *49*(5), 847–864. https://doi.org/10.1080/03075079.2024.2323593
- Nikolic, S., Daniel, S., Haque, R., Belkina, M., Hassan, G. M., Grundy, S., Lyden, S., Neal, P., & Sandison, C. (2023). ChatGPT versus engineering education assessment: A multidisciplinary and multi-institutional benchmarking and analysis of this generative artificial intelligence tool to investigate assessment integrity. *European Journal of Engineering Education*, 48(4), 559–614.

https://doi.org/10.1080/03043797.2023.2213169

- Nussbaum, E. M., & Schraw, G. (2007). Promoting argument-counterargument integration in students' writing. *The Journal of Experimental Education*. https://doi.org/10.3200/JEXE.76.1.59-92
- OpenAI. (2022). *ChatGPT* (Version December 15) [Large language model; Large language model]. https://chat.openai.com/chat
- Pigg, S. (2024). Research writing with ChatGPT: A descriptive embodied practice framework. *Computers and Composition*, *71*, 102830. https://doi.org/10.1016/j.compcom.2024.102830
- Porter, B., & Machery, E. (2024). Al-generated poetry is indistinguishable from human-written poetry and is rated more favorably. *Scientific Reports*, *14*(1), 26133.

https://doi.org/10.1038/s41598-024-76900-1

- R Core Team. (2024). *R: A language and environment for statistical computing* (Version 4.3.3) [Computer software]. R Foundation for Statistical Computing.
- Ritchhart, R. (2011). *Making thinking visible: How to promote engagement, understanding, and independence for all learners*. Jossey-Bass.
- Sampson, V., & Clark, D. B. (2008). Assessment of the ways students generate arguments in science education: Current perspectives and recommendations for future directions. *Science Education*, 92(3), 447–472.
- Shojaee, P., Mirzadeh, I., Alizadeh, K., Horton, M., Bengio, S., & Farajtabar, M. (2025). The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity. *arXiv Preprint arXiv:2506.06941*.
- Swiecki, Z., Khosravi, H., Chen, G., Martinez-Maldonado, R., Lodge, J. M., Milligan, S., Selwyn,
 N., & Gašević, D. (2022). Assessment in the age of artificial intelligence. *Computers and Education: Artificial Intelligence*, *3*, 100075. https://doi.org/10.1016/j.caeai.2022.100075

Toulmin, S. E. (1958). The uses of argument (Repr. of updated ed). Cambridge University Press.

- Wahn, B., Schmitz, L., Gerster, F. N., & Weiss, M. (2023). Offloading under cognitive load:
 Humans are willing to offload parts of an attentionally demanding task to an algorithm.
 PLOS ONE, *18*(5), e0286102. https://doi.org/10.1371/journal.pone.0286102
- West, P., Lu, X., Dziri, N., Brahman, F., Li, L., Hwang, J. D., Jiang, L., Fisher, J., Ravichander,
 A., Chandu, K., Newman, B., Koh, P. W., Ettinger, A., & Choi, Y. (2023, October 13). *The Generative AI Paradox: "What It Can Create, It May Not Understand."* The Twelfth
 International Conference on Learning Representations.
 https://openreview.net/forum?id=CF8H8MS5P8
- White, J., Fu, Q., Hays, S., Sandborn, M., Olea, C., Gilbert, H., Elnashar, A., Spencer-Smith, J.,
 & Schmidt, D. C. (2023). A Prompt Pattern Catalog to Enhance Prompt Engineering with ChatGPT (arXiv:2302.11382). arXiv. https://doi.org/10.48550/arXiv.2302.11382

- Wingate, U. (2012). 'Argument!'helping students understand what essay writing is about. *Journal of English for Academic Purposes*, *11*(2), 145–154.
- Wolfe, C. R., Britt, M. A., & Butler, J. A. (2009). Argumentation schema and the myside bias in written argumentation. *Written Communication*, *26*(2), 183–209. https://doi.org/10.1177/0741088309333019
- Xia, Q., Weng, X., Ouyang, F., Lin, T. J., & Chiu, T. K. F. (2024). A scoping review on how generative artificial intelligence transforms assessment in higher education. *International Journal of Educational Technology in Higher Education*, *21*(1), 40. https://doi.org/10.1186/s41239-024-00468-z
- Yan, D. (2023). Impact of ChatGPT on learners in a L2 writing practicum: An exploratory investigation. *Education and Information Technologies*, 28(11), 13943–13967. https://doi.org/10.1007/s10639-023-11742-4
- Zhang, R., & and Zou, D. (2022). Types, features, and effectiveness of technologies in collaborative writing for second language learning. *Computer Assisted Language Learning*, 35(9), 2391–2422. https://doi.org/10.1080/09588221.2021.1880441

Appendix

Taxonomy to evaluate student-GenAI interactions

Category	Туре	Meaning	
Writing	Instructions	User specifies the task, in terms of the course's assignment description (e.g. copy-paste or upload)	
	Criteria	User specifies the task in more detail, by providing the evaluation criteria for the assignment, from the assignment rubric (usually, copy-paste)	
	Evaluate	User asks the machine to evaluate a draft against the provided criteria (or without criteria).	
	Improve	User provides a phrase, paragraph, or essay to be improved by the machine for e.g. spelling, style or grammar.	
	Format	User asks for improved formatting (including e.g. bibliographical formatting)	
	Organization	User asks for feedback or improvement of essay structure.	
	Introduction	User asks the machine to provide an effective introduction.	
	Conclusion	User asks the machine to provide an effective conclusion.	
	Role	User specifies the role/character/expertise the language model should take.	
	AutoComplete	User asks machine to append or expand on text, without providing specific guidance about the content.	
	Summarize	User asks machine to summarize text (e.g. an uploaded article).	
	Content Removal	User ask machine to delete existing text (e.g., deleting a specific paragraph or sentence)	
	Miscellaneous	User prompting system in a non-specific technical way.	
Content	Bibliography	User asks for bibliographic references on a specific topic.	
	Example	User asks the machine to provide specific example for a general case or issue.	

	Research	User asks the machine to define an idea, or to identify related ideas to one, given by the user.
	Definitions	User provides the machine with definitions to/elaborations of key technical terms discussed in the course (e.g. "data activism").
	Case	User describes a relevant case from class/their own research.
	Idea	User provides well-motivated original idea or question and asks for confirmation/elaboration/discussion.
	Concept	User introduces a keyword concept from the course material and asks the machine to define it or apply it to a case.
	Elaboration	User provides a relevant sentence/paragraph and asks the machine to elaborate and provide additional detail, mentioning specific course-related content.
	Theory	User asks the machine to appeal to a philosophical or ethical theory (e.g. consequentialism), named or not.
	Critical	User critically engages with AI-generated content, asking for clarification or correction
Argument	Context	User asks the machine to describe or analyze the context of a real world case, technology, or news story. E.g. setting the case into a broader debate.
	Case Research	User asks the machine to describe or analyze the details of a given case.
	Stakeholders	User asks the machine to identify the stakeholders for a case or technology.
	Values	User asks the machine to specify the values of the stakeholders in a case.
	Moral Problem	User asks the machine to formulate a moral problem or identify an ethical issue with a particular case or technology
	Objection	User asks the machine to provide an objection and/or a response to a given claim.
	Justify	User asks the machine to provide reasons for a given claim

Structure	User asks the machine to impose a particular logical structure onto a text.
Improve	User asks the machine to improve the argumentative structure (according to given criteria).
Relate	User asks the machine to relate or connect two concepts or ideas.
Conceptual Clarity	User asks the machine to simplify or otherwise improve the definition of concepts.
Thesis	User asks the machine to make a thesis/conclusion more precise, concise, or clear.