

# Rubrik's Cube: Testing a New Rubric for Evaluating Explanations on the CUBE dataset



Diana Galvan-Sosa<sup>1,★</sup>, Gabrielle Gaudeau<sup>1,★</sup>, Pride Kavumba<sup>2,†</sup>, Yunmeng Li<sup>3</sup>,  
Hongyi Gu<sup>5</sup>, Zheng Yuan<sup>6</sup>, Keisuke Sakaguchi<sup>3,4</sup>, Paula Buttery<sup>1</sup>

<sup>1</sup>ALTA Institute, Computer Laboratory, University of Cambridge, <sup>2</sup>SB Intuitions, <sup>3</sup>Tohoku University,

<sup>4</sup>RIKEN, <sup>5</sup>NetMind.AI, <sup>6</sup>The University of Sheffield

## Abstract

The performance and usability of Large-Language Models (LLMs) are driving their use in explanation generation tasks. However, despite their widespread adoption, LLM explanations have been found to be unreliable, making it difficult for users to distinguish good from bad explanations. To address this issue, we present Rubrik's CUBE—an education-inspired rubric and a dataset of 26k explanations, written and later quality-annotated using the rubric by both humans and six open- and closed-source LLMs. The CUBE dataset focuses on two reasoning and two language tasks, providing the necessary diversity for us to effectively test our proposed rubric. Using Rubrik, we find that explanations are influenced by both task and perceived difficulty. Low quality stems primarily from a lack of conciseness in LLM-generated explanations, rather than cohesion and word choice. The full dataset, rubric, and code are available at [https://github.com/RubriksCube/rubriks\\_cube](https://github.com/RubriksCube/rubriks_cube).

## 1 Introduction

Explanations play a crucial role in the process of understanding why a decision was made. But, as illustrated in Figure 1, there exist many ways of expressing the rationale behind a choice. Large-Language Models (LLMs), with their inherent capacity for generating very different outputs given the same query, provide a compelling example of this phenomenon. In fact, these models are increasingly being used in applications which expect a detailed breakdown explaining why a decision was made (e.g., automated scoring, question generation, problem resolution; García-Méndez et al., 2024).

Unfortunately, LLM-generated explanations generally fall short of user expectations due to their unreliability (Kim et al., 2024). Indeed, they are

<sup>★</sup>Equal contribution, contact: {dg693,gjg34}@cam.ac.uk.

<sup>†</sup>Most of the author's contribution was performed while at LegalOn Technologies.

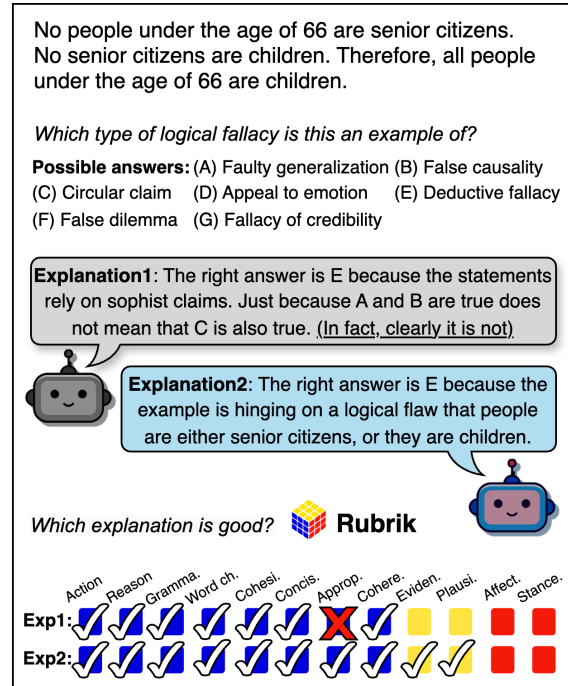


Figure 1: Two explanations of varying quality (in terms of APPROPRIATENESS and the provision of EVIDENCE) which present the logic behind an answer choice.

known to occasionally hallucinate, produce incorrect or misleading information, and struggle to back up their responses to queries, highlighting an overall deficiency in their reasoning capabilities (Huang and Chang, 2023; Saxena et al., 2024). As noted by Zhang et al. (2023a), these issues remain unaddressed, even by prompting strategies like “Let’s think step by step.” As a result, LLM-generated explanations lack transparency, and are a source of misinformation and limited knowledge (Sallam, 2023; Kabir et al., 2024). Consequently, the challenge has shifted from generating text to assessing its quality, a difficulty that has led some sites to temporarily ban the use of any generative AI (GenAI)<sup>1</sup>.

<sup>1</sup>See StackOverflow’s policy on the use of ChatGPT and other LLMs: <https://stackoverflow.com/questions/7712261/chatgpt-policy>

The most common practice in GenAI to determine the quality of a text is to rely on human evaluators. However, because such evaluators typically lack specific training, the exact evaluation criteria are left to their discretion (Clark et al., 2021). Inspired by the use of rubrics in education for the qualitative evaluation of complex and subjective tasks like essay writing (e.g., the IELTS writing rubric; Arnold, 2023), we design our very own rubric following Dawson (2017)’s best practices. In doing so, we align ourselves with the human-grounded evaluation proposed by Doshi-Velez and Kim (2017), which identifies and evaluates the “general notions” of the quality of an explanation without having a specific end goal.

We thus introduce Rubrik’s CUBE<sup>2</sup>, a task-independent rubric and a dataset to help evaluate the quality of LLM-generated explanations. Rubrik identifies the core components and features of a *good* explanation, differentiated by explanation type; CUBE contains 26k explanations drawn from instances of four distinct tasks, generated by both humans and a set of open- (Command R+, Gemma 2, Llama 3.1, Mixtral) and closed-source (GPT-4o, Claude Sonnet 3.5) models. We additionally include two custom agreement metrics that account for the hierarchical and nested nature of our rubric. Rubrik enables valuable insights on output quality, allowing us to identify distinct patterns in the explanations of all annotators. We observe that the explanation type depends on the task and its perceived difficulty. Specifically, our rubric revealed that low-quality LLM explanations are primarily due to not being concise and only rarely because of word choice or cohesion.

## 2 Background

We summarise different bodies of literature on the nature and qualities of explanations which, alongside insights from the education assessment literature, informed the design of our rubric.

### 2.1 Cognitive Science and Social Sciences

There is an open discussion in philosophy and other social sciences like psychology about what an explanation is and what makes the best explanation (Doshi-Velez and Kim, 2017; Gilpin et al., 2018; Miller, 2019a). From the psychology and

meta.stackoverflow.com/questions/421831/policy-generative-ai-e-g-chatgpt-is-banned

<sup>2</sup>Short for Commonsense reasoning, Usual logical fallacies, Basic reading comprehension, and Essay scoring.

cognitive science perspective, an explanation is something ubiquitous, diverse, and fundamental to humans’ sense of understanding. They come in a variety of forms and formats and are used for a variety of purposes (Keil, 2006), including: (1) understanding a *decision process* (2) understanding and predicting an *unexpected event*, and (3) filling a gap in knowledge (i.e., *learning*). It follows that a *good* explanation is inherently related to its purpose, which some suggest is shaped by what is being asked (Bromberger, 1992). In particular, authors like Lombrozo (2006) and Miller (2019a) argue that an explanation’s relation to cognition comes from an attempt to answer a *why-question*. Miller investigated the criteria that people use to evaluate explanations, finding that the most important are: PROBABILITY, SIMPLICITY, GENERALISE, and COHERENCE with prior beliefs. The truth of LIKELIHOOD is also identified as an important criterion. However, Miller notes that an explanation that includes this attribute is not always the *best* explanation.

### 2.2 Explainable AI

In the context of Explainable AI (XAI) and Machine Learning (ML) interpretability, an explanation should be able to reflect the internal decision process of a system. *Introspective* systems output this kind of explanation, while *justification* systems output evidence supporting a decision (Park et al., 2018). The most studied properties of explanation systems include FIDELITY, STABILITY, COMPREHENSIBILITY, GENERALISABILITY and CONSISTENCY (Fel et al., 2022). According to Wiegreffe and Marasović (2021), explanations are implicitly or explicitly designed to answer the *why-question* “*why is <input> assigned <label>*”. They identified HIGHLIGHTS (subsets of the input elements that explain a prediction) as one type of explanation in the Explainable NLP (EXNLP) literature, where COMPACTNESS, SUFFICIENCY and COMPREHENSIVENESS are the main attributes.

### 2.3 Natural Language Generation

In an attempt to find a consensus about how human evaluations of generated text should be designed and reported, Howcroft et al. (2020) examined twenty years of NLG papers that reported some form of human evaluation. Some of the most common criteria used to assess quality include FLUENCY, APPROPRIATENESS and CLARITY.

The Multidimensional Quality Metrics (MQM)





 <b>Typology of Explanations</b>	<b>COMPONENTS</b>	<b>DIMENSIONS</b>	
	necessary parts of an explanation	necessary qualities of a <i>good</i> explanation	
		<b>Language</b>	<b>Content</b>
Typ1.  <b>COMMENTARY</b>	1.a) Action 1.b) Reason	Grammaticality Word Choice Cohesion	Conciseness Appropriateness Coherence
Typ2.  <b>JUSTIFICATION</b>	2.a) Evidence		Plausibility
Typ3.  <b>ARGUMENT</b>	3.a) Affective appeal(s) and Qualifier(s)		Stance Clarity

Table 1: Overview of our evaluation rubric which identifies three hierarchical types of explanations, their necessary parts (COMPONENTS), and the features that distinguish the *good* from the *bad* ones (DIMENSIONS).

framework (Burchardt, 2013; Mariana, 2014; Freitag et al., 2021) has been widely applied to machine translation studies in recent years.<sup>3</sup> This hierarchical typology of quality issues provides a detailed and flexible approach for evaluating translation tasks. It has been applied to different domains of machine translation, such as literary translation (Karpinska and Iyyer, 2023) and chat translation (Li et al., 2025). It could be used as a metric for human evaluators to evaluate translation models, and could also be used to prompt models as evaluators (Park and Padó, 2024; Li et al., 2025). It identifies seven high-level error types (namely TERMINOLOGY, ACCURACY, LINGUISTIC CONVENTIONS, STYLE, LOCALE CONVENTIONS, AUDIENCE APPROPRIATENESS, and DESIGN AND MARKUP), which can be broken into multiple subtypes to enable fine-grained assessment. In the design of our rubric, we similarly arranged the significant features of explanations hierarchically to allow for both coarse and granular evaluations.

## 2.4 Education

Education, and specifically science education, has long focused on teaching students how to construct explanations, and assessing them (e.g., Sandoval, 2003; McNeill et al., 2006; McNeill and Krajcik, 2008; Zangori et al., 2013). For them, explanations “make sense of a phenomenon based on other scientific facts” (Ohlsson, 2002). They should begin with a statement of the *explanandum* (i.e., the phenomenon to be explained). Then, what makes a *good* explanation differs is “explanatory adequacy” (Brigandt, 2016) which consists in providing an understanding of how or why a phenomenon occurs (Chin and Brown, 2000).

<sup>3</sup>An updated version (MQM 2.0) is available from <https://themqm.org/the-mqm-full-typology/>.

In practice, assessing explanations is difficult (Berland and McNeill, 2012), so teachers generally rely on rubrics, like the one proposed by McNeill and Krajcik (2007), which provide clear, consistent, and objective sets of criteria for evaluation. More generally, rubrics are firmly established evaluation tools in written assessment and widely advocated in books by Walvoord and Anderson (1998); Huba and Freed (2000); Dunn et al. (2003); Stevens and Levi (2004); Freeman et al. (2016). Unfortunately, these practices are not currently being used beyond education, and no equivalent rubric exists for evaluating LLM explanations on a variety of tasks (beyond scientific explaining). To address this gap, we propose to draw on this literature to come up with our very own rubric.

## 3 A Systematic Quality Assessment Framework

This section introduces our proposed assessment framework in three parts. First, we detail the design decisions taken to develop the rubric, drawing upon the key principles outlined by Dawson (2017). Second, we provide a comprehensive overview of the rubric itself, outlining its key elements and their hierarchical relationships. Finally, we provide practical guidance on how to effectively use the rubric for the task of explanation quality assessment.

### 3.1 Designing an Assessment Rubric

Recognising that the foundation of an effective evaluation lies in its instrument, we carefully considered the design elements suggested by Dawson (2017). A key advantage of adhering to their framework is the streamlined design process and the enhanced transparency of the resulting rubric, facilitating easier comparisons with other instruments.

Design element	Decision
<i>Specificity</i> : the particular object of assessment	Assess the quality of explanations.
<i>Secrecy</i> : who the rubric is shared with, and when it is shared	It should be secret to the annotators. It is only shared with the evaluators.
<i>Exemplars</i> : work samples provided to illustrate quality	Examples of acceptable and not acceptable instances.
<i>Scoring strategy</i> : procedures used to arrive at marks and grades	A series of binary judgments (yes/no) all amounting to a binary decision (good/bad).
<i>Evaluative criteria</i> : overall attributes required of the explanation	Components and dimensions.
<i>Quality levels</i> : the number and type of levels of quality	Two quality levels (☺ good or ☹ bad).
<i>Quality definition</i> : explanations of attributes of different levels of quality	Motivated by different bodies of literature (social sciences, XAI, and NLG).
<i>Judgment complexity</i> : the evaluative expertise required of users of the rubric	Should be simple enough for <b>anyone</b> to use.
<i>Users and uses</i> : who makes use of the rubric, and to what end	Evaluators use for summative assessment.
<i>Creators</i> : the designers of the rubric	NLP researchers.

Table 2: Summary of the design decisions taken to develop our proposed rubric. The design elements are those suggested by Dawson (2017). The “annotators” are the humans or LLMs who write the explanations.

### 3.2 A Task-Agnostic Quality Rubric

A fundamental assumption underlying this work is that it is possible to account for the diverse nature of explanations (which can serve a wide range of goals as highlighted in Section 2.1) whilst also being able to recognise common features that generally characterise them. Through Section 2, we showed that different bodies of literature identify shared attributes of a *good* explanation. Using these attributes, our proposed rubric (henceforth Rubrik) classifies explanations into three goal-driven types. Each type is defined by the presence of specific COMPONENTS. The typology is hierarchical and nested, with subsequent types inheriting the COMPONENTS of preceding types and adding to them. Each explanation type also comes with its own set of attributes called DIMENSIONS: together, these capture the quality of an explanation of that type (*good* or *bad*). Much like COMPONENTS, DIMENSIONS are inherited and accumulate across the type hierarchy. This “building block”-like structure provides a robust framework for understanding how the form and features of explanations evolve alongside the distinct goals of each type. Table 1 presents an overview of our proposed rubric (see Table 6 in Appendix A for the full-sized, illustrated rubric) and Table 2 shows the design considerations and choices we made in developing it.

#### 3.2.1 Components

The three hierarchical and nested explanation types in Rubrik are: COMMENTARY, JUSTIFICATION, and ARGUMENT. The COMMENTARY is the foundational level and consists of two COMPONENTS: an ACTION and a REASON. The JUSTIFICATION

extends this base by incorporating an additional COMPONENT: an EVIDENCE. Finally, the ARGUMENT includes the elements of both the COMMENTARY and the JUSTIFICATION, as well as an additional unique element: the AFFECTIVE APPEAL(S) AND QUALIFIER(S). This progression, where each higher-level type nests the elements of the lower-level ones, results in an increasing richness of information. Providing an *understanding* of a decision process is the central goal of a COMMENTARY and a JUSTIFICATION. An ARGUMENT, while also considering the same goal, is more focused on *persuasion*. Formally,  $\text{COMMENTARY} \subseteq \text{JUSTIFICATION} \subseteq \text{ARGUMENT}$ . See Appendix A.1 for a more in-depth understanding of the reasoning that led to the definition of types and components.

#### 3.2.2 Dimensions

COMPONENTS provide the necessary structural elements of different types of explanations; DIMENSIONS are their requisite qualities. This distinction ensures that our rubric accounts for both what is being said (through the COMPONENTS) and how well it is communicated (through the DIMENSIONS).

The eight DIMENSIONS shown in Table 1 were chosen from a wider set of explanation qualities (see Table 5 in Appendix A.2) that have been studied, annotated or evaluated in the bodies of literature introduced in Section 2. We filtered out those that were too task-specific for our goal of creating a general-purpose rubric (e.g., FIDELITY, CONSISTENCY, TRANSPARENCY and INTERPRETABILITY specifically focus on the internal workings of AI models) or too vague (for e.g., CLARITY; see Section A.8). The eight remaining DIMENSIONS



were then put in one of two categories. **Language** assesses whether the explanation is well-formed; **Content** evaluates the ideas expressed by the explanation. This design choice was motivated by the fact that LLMs sometimes produce text that is only *good* on the surface but factually incorrect, inappropriate, or misleading (Huang et al., 2025). We describe our process in more detail in Appendix A.2.

These DIMENSIONS were then related to the COMPONENTS and explanation types introduced in the previous section. ACTION and REASON are pre-requisites for a COMMENTARY to be considered complete; but for it to be *good*, we must enforce certain linguistic requirements: it needs to be grammatical, cohesive, and use *context*-appropriate language. On the other hand, its content should be coherent and concise and match the expectations imposed by the defined *context*. Further, a JUSTIFICATION is contingent on the presence of EVIDENCE. Ensuring it is plausible and consistent with human reasoning is a further requirement for a *good* JUSTIFICATION. Finally, the presence of argumentative markers generally betrays the explainer’s intent to persuade the audience of their *stance* (i.e., their personal feelings towards the task). Whether this stance is clearly and unambiguously conveyed distinguishes a *good* from a *bad* ARGUMENT.

### 3.3 Scoring Strategy

To use Rubrik, evaluators must first establish the *context* of the explanations:

- What is the task? In our case, we will be looking at two reasoning and two language tasks (Section 4.1).
- Who is the target audience? In our case, NLP researchers (i.e., formal academic setting)
- What is their intended goal?

Once the *context* is defined, we can proceed with the evaluation. Given an explanation, the outcome of an evaluation with Rubrik is a **Type** for that explanation (NONE, COMMENTARY, JUSTIFICATION, ARGUMENT) and a related **Quality label** (☺ *good* or ☹ *bad*). The evaluation process follows our hierarchical typology: starting from the foundational level—the COMMENTARY—and going all the way to the ARGUMENT. We describe this process in detail below:

- First, we start by checking whether the two COMPONENTS of the COMMENTARY (namely ACTION and REASON) are present (✓) or absent (✗) in the explanation. If either COMPONENT is missing (✗), then the explanation is incomplete and classified as NONE, and the evaluation ends there. If, on the other hand, both are present (✓), then the explanation’s **Type** is at least a COMMENTARY.
- Next, we check whether the explanation satisfies (✓) each of the COMMENTARY’s six DIMENSIONS or not (✗). If the explanation fails to meet any of these (✗), then the explanation is a ☹ *bad* COMMENTARY and the evaluation ends there. If however, all six DIMENSIONS are satisfied (✓), then the explanation is at least a ☺ *good* COMMENTARY.
- Continue this procedure with the COMPONENTS and DIMENSIONS of the JUSTIFICATION. Specifically, if the explanation does not have EVIDENCE (✗), then the explanation is only a ☺ *good* COMMENTARY and the evaluation ends there. If it does (✓), then it is at least a JUSTIFICATION. Whether it is a ☺ *good* or ☹ *bad* JUSTIFICATION will depend on whether the EVIDENCE is judged as PLAUSIBLE (✓) or not (✗). If it is the latter, then the evaluation ends there; otherwise, the explanation is at least a ☺ *good* JUSTIFICATION.
- Repeat this process with the ARGUMENT’s COMPONENT and DIMENSION.

Notice that for each explanation type, we performed two validation steps: (1) Structure validation (determined by the COMPONENTS) and (2) Attribute validation (determined by the DIMENSIONS). At each step, the evaluator makes a series of binary judgements based on the presence (✓) or absence (✗) of COMPONENTS, and whether DIMENSIONS are satisfied (✓) or not (✗), using the definitions and examples included in the full rubric (Table 6) as reference.

## 4 Rubric Validation

The main motivation behind our proposed rubric is to allow for a more systematic evaluation of an explanation’s quality. In order to determine the effectiveness of our proposal, we designed a validation process aimed at addressing the following question: *Does the rubric effectively discriminate*

	Single annotations				Joint annotations						Single evaluations			Joint evaluations					Total
	Inst.	LLM	Total		Inst.	H	LLM	Total	Total	Total	Inst.	E	LLM	Inst.	E	H	LLM		
T1	<u>1000</u>	<u>890</u>	6	5340	<u>110<sup>‡</sup></u>	4	6	10	1100	6440	90 <sup>‡</sup>	900	1	20 <sup>‡</sup>	200	2	1		
T2	<u>1000</u>	<u>890</u>	6	5340	<u>110<sup>‡</sup></u>	4	6	10	1100	6440	90 <sup>‡</sup>	900	1	20 <sup>‡</sup>	200	2	1		
T3	<u>1000</u>	<u>890</u>	6	5340	<u>110<sup>‡</sup></u>	7	6	13	1430	6770	90 <sup>‡</sup>	1170	1	20 <sup>‡</sup>	260	2	1		
T4	<u>1000</u>	<u>890</u>	6	5340	<u>110<sup>‡</sup></u>	7	6	13	1430	6770	90 <sup>‡</sup>	1170	1	20 <sup>‡</sup>	260	2	1		
Total	4000	3560		21360	440				5060	26420	360	4140		80	920			5060	

Table 3: Instances and explanations (E) in CUBE. Double-underlined numbers represent the initial pool, divided into subsets (single-underlined) based on the annotators assigned. A (‡) denotes variations in evaluator assignment.

between high-quality and low-quality explanations, while simultaneously providing clear and concise guidance for evaluators? Given the absence of existing datasets for explanation assessment, the validation of this rubric required a tailored approach. This began with identifying an appropriate source of data, followed by gathering explanations, evaluating them using the rubric with three raters, and finally, measuring the inter-rater reliability to determine the consistency of the rubric’s application. The effectiveness of our rubric was evaluated by measuring the level of inter-rater agreement for each explanation.

## 4.1 Data Collection

We assume a decision-making scenario involving a set of choices, where one is selected. Thus, our data collection process required instances from tasks that could be framed as a series of multiple-choice questions (MCQ) with a single correct answer. To ensure a diverse set of explanations, we chose four different tasks, drawn from reasoning and language assessment. The reasoning tasks are: (T1) commonsense reasoning and (T2) fallacy detection. The language tasks are: (T3) reading comprehension and (T4) essay scoring. From an initial pool of 1000 instances from each task, we curated an *annotation set* of 440 total instances for annotation (110 from each dataset). A brief description of the datasets follows. Detailed selection criteria are described in Appendix B.

**Reasoning tasks.** For T1 and T2, we selected instances from the HellaSwag (Zellers et al., 2019) and Logic (Jin et al., 2022) datasets, respectively. Each instance in HellaSwag has a **context** and a set of four ENDINGS; the task is to select the most likely follow-up sentence. Logic consists of common logical fallacy examples collected from various online educational materials.

**Language tasks.** For T3 and T4, we selected instances from RACE (Lai et al., 2017) and the Write&Improve (W&I) (Bryant et al., 2019) cor-

pus, respectively. RACE consists of a series of passages and questions taken from English exams that evaluate a student’s ability in understanding and reasoning. Write&Improve<sup>4</sup> is an online web platform that assists English Language Learners with their writing (Yannakoudakis et al., 2018). The dataset contains submissions (defined as “essays”) that were annotated with a coarse CEFR<sup>5</sup> level (A, B or C) by trained raters.

### 4.1.1 Annotation

Two key decisions shaped the annotation process. First, we retained all annotations, regardless of the correctness of the chosen answer. This decision was driven by the need to explore the explanations associated with correct and incorrect answers, allowing for a more nuanced understanding of the explanatory quality. Second, human explanations were not treated as the gold standard. This allowed for a more objective comparison of human and LLM explanations, avoiding potential bias towards human responses. Below, we give a brief overview of the annotation process, but we refer the reader to Appendix C for more information.

**Human.** We recruited seven annotators: four general annotators (contractors) and three professionals with experience in language assessment. They were asked to answer a series of multiple-choice questions and explain their choices. While contractors covered all four tasks, experts focused on the language tasks. This process resulted in 880 explanations for T1 and T2, 1, 540 for T3 and T4.

**LLM-based.** We worked with six LLMs, including four open-source: Llama 3.1 (Dubey et al., 2024), Gemma 2 (Team et al., 2024), Mixtral (Jiang et al., 2024) Command R+, (Cohere for AI, 2024) and two closed-source models: GPT-4o (Ope-

<sup>4</sup><https://writeandimprove.com/>.

<sup>5</sup> Common European Framework of Reference for Languages (North and Piccardo, 2020) levels correspond to language proficiency levels ranging from A1 (elementary) to C2 (complete proficiency) from a second-language learner’s perspective.

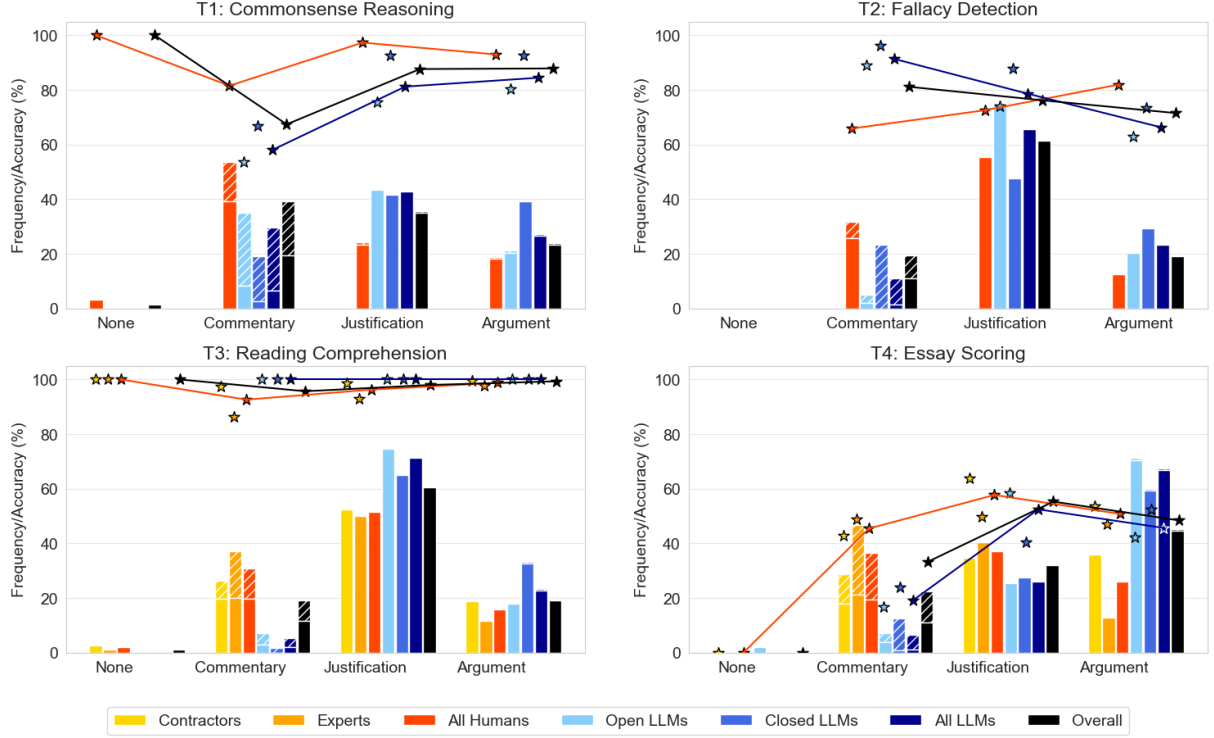


Figure 2: The bar plots show the frequencies (%) of the different explanation types in each group of annotators as judged by and averaged across the three evaluators (two humans and GPT-4o). The patterned fill indicates the proportion of *bad* explanations of each type; the solid fill shows the proportion of *good* explanations of each type. The scattered stars represent the accuracy (%) of each group of annotators (i.e., did they select the correct answer out of the possible multiple choices to a question) related to the type of explanation they produced as judged by and averaged across the three evaluators. We plot the accuracy lines for the following three groups: all human annotators, all LLMs, and all annotators (“Overall”).

nAI, 2024) and Claude 3.5 Sonnet (Anthropic, 2024). See Appendix C.2 for model versions. Models were prompted using a few-shot setting (see Appendix C.2.1). Explanations were generated for all instances, yielding a total of 24,000 explanations. Table 3 shows a more detailed breakdown of the number of annotations and evaluations.

#### 4.1.2 Evaluation

Data evaluation was performed by two expert evaluators and the same six LLMs on a subset of the *annotation set*: namely, 20 instances for each task. Thus, our *evaluation set* has a total of 920 explanations derived from 80 instances. Using two custom agreement metrics, we identified that out of the LLMs, GPT-4o most closely matched our human evaluators. As was previously done by Brassard et al. (2024) and Sottana et al. (2023), we took GPT-4o to act as our third evaluator to enhance the robustness of our analysis, and used it to automatically evaluate the 4,140 explanations from the remaining 360 instances of the *annotation set*. For details on the preliminary experiment and metrics,

see the Appendix D.

The raters followed the scoring strategy specified in Section 3.3. Unlike the human raters, GPT-4o limited its role to validating only the structure and attributes of the explanations. In other words, it did not render a final judgment on an explanation’s quality. This approach mitigated the risk of the model’s self-bias (as reported in Panickssery et al., 2024); further details on this potential source of bias are provided in Appendix E.

## 5 Discussion

**Agreement.** A key indicator of the utility of Rubrik is the level of agreement observed between the human evaluators who used it. Standard inter-rater agreement metrics are often inadequate for nested hierarchical data. Therefore, we designed a custom metric that accounts for both *superlabels* (explanation types) and *sublabels* (COMPONENTS and DIMENSIONS) in Rubrik, penalising discrepancies based on the difference in hierarchical level. Using this novel metric, we found an average inter-rater

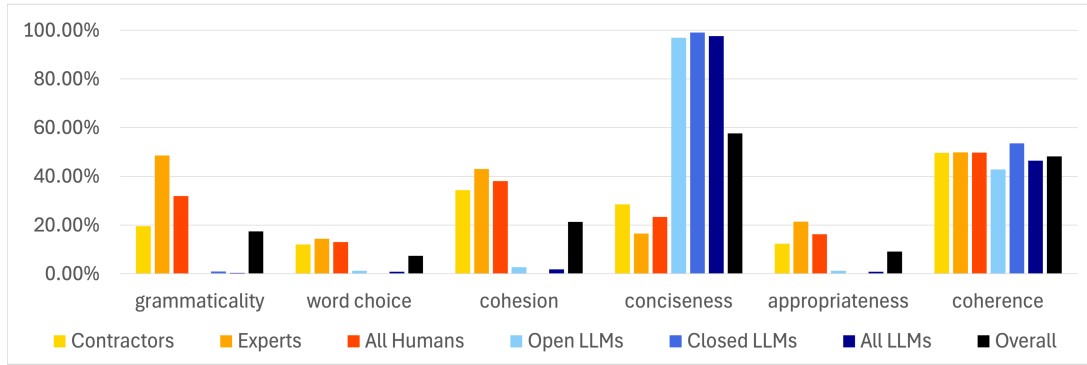


Figure 3: Plot showing the source of  $\oplus$  bad COMMENTARIES (i.e., which of the COMMENTARY’s DIMENSIONS was not met  $\times$ ) in the *evaluation set*. We average the frequencies across all three evaluators (two humans and GPT-4o).

agreement of 0.86 and 0.878 for superlabels and sublabels, respectively, among humans. In selecting the third evaluator, our preliminary experiments revealed that LLMs tended to favour JUSTIFICATIONS, potentially inflating agreement scores on this first metric. To address this, we designed a second metric that weights the evaluations based on a comparison with both human and LLM judgments, providing a more accurate measure of performance. Using both custom metrics, we obtained scores of 0.841 (superlabel) and 0.86 (sublabel) for metric one, and 0.476 for the second. The second metric led to the selection of GPT-4o as the third evaluator.

**Task Performance.** As mentioned in Section 4.1.1, we decided to keep explanations, even if they are associated with an incorrect answer. Just as explanations are inherently tied to their goal, we hypothesised that they might vary depending on the task, and how successful the annotators were. To explore this, we started by looking at the average performance of each annotator across tasks. Humans showed an average accuracy of T1: 70.46%, T2: 69.09%, T3: 80.78%, T4: 55.06%; LLMs showed T1: 78.94%, T2: 69.24%, T3: 87.42%, T4: 47.58% (as reported in Figure 7b). Overall, closed-source LLMs outperformed humans and open-source models. Interestingly, not only did T2 and T4 have the lowest accuracies, annotators also reported lower confidence on these tasks in comparison to T1 and T3 (see Appendix C.1.3). In fact, T4 proved to be the most challenging task for all annotators, while T3 was the least challenging. See Figure 7a in Appendix F.2 for a breakdown of these accuracies per annotator.

**Frequency of Explanation Types.** The bar plots in Figure 2 show the frequencies of each explanation type as judged by and averaged across the three evaluators (GPT-4o and two humans) in the

*evaluation set*. Overall, the evaluators judged explanations to be mostly JUSTIFICATIONS. A notable observation is the low frequency of negative types (i.e., NONE). A closer look at the data revealed that these assignments were predominantly made by human evaluators. Furthermore, we found that T4 had a much higher proportion of ARGUMENTS than other tasks, whereas T3, the easiest task, had comparatively fewer. These results reveal insights into the tendencies of humans and LLMs to generate JUSTIFICATIONS, whilst also highlighting the influence of task characteristics on the nature of generated explanations. T4 is a notoriously complex task that requires evaluators to go beyond simply recognising correct language use. They must also assess the effectiveness of the writing in achieving its intended purpose, which involves subjective judgments about argumentation, organisation, and style. While some interpretation might be involved in understanding the context in T1, T2 and T3 the range of acceptable interpretations is much narrower. Thus, our results suggest that the presence of ARGUMENTS is correlated with the subjectivity of the task. The relationship between ARGUMENTS and task subjectivity is reinforced by the findings of our follow-up survey, where human annotators expressed lower confidence in T4. Upon further inspection of the frequency of ARGUMENTS across tasks, we found that Sonnet 3.5, while similar in terms of accuracy to GPT-4o, is more likely to produce this type of explanation. Figure 8 in Appendix F provides a more granular view of these findings.

**Accuracy Across Types.** The scatter plot in Figure 2 relates the types of explanations produced by the annotators and their accuracy in each task (**Task Performance**). We observe an interesting trend in T1, T3 and T4: the “Overall” line shows that lower accuracy in a task is associated with the lowest type









Type	Example
 COMMENTARY	The right answer is D because the reason given to encourage Luke to eat is subjective and has nothing to do with his taste in food or any potential benefits. It is not a valid reason.
 COMMENTARY	The right answer is D because this statement is trying to make Luke <b>eat the sheep's brains with chopped liver and brussel sprouts</b> by making him feel guilty <b>about the poor, starving children in a third world country</b> . It's an appeal to his emotions, rather than presenting a logical argument.
 JUSTIFICATION	The right answer is A because the woman in the video is demonstrating how to make or destroy lipsticks, which is a process that involves using one's mouth and lips. This is the only option that describes an action that would require the use of the mouth and lips.
 JUSTIFICATION	The right answer is D because it is the most reasonable answer of the options. A is incorrect because <b>demonstrating how to make and destroy lipsticks is an unusual subject to demonstrate: the making of lipsticks in particular is complicated and challenging</b> . B is incorrect because one's tongue does not fall out when speaking. C is incorrect because you do not describe ordinarily describe women as walking into a classroom with "both [their] boobs" as this implies their boobs are separable. D is correct because sucking from a hookah is a fairly ordinary activity.
 ARGUMENT	The right answer is B because the essay is written in a somewhat coherent and understandable manner, but it lacks clarity, coherence, and proper sentence structure. The writer's emotions and thoughts are expressed, but the writing is not sophisticated or polished. The essay does not demonstrate a clear understanding of the topic or the ability to express complex ideas. The writer's use of language is simple, and the essay lacks depth and analysis. Therefore, it is best graded as Intermediate (grade B).
 ARGUMENT	The right answer is B because the essay demonstrates a basic understanding of grammar and sentence structure, <b>but there are some errors in word choice and sentence construction. The vocabulary used is relatively simple, and the ideas are expressed in a straightforward manner. However, the essay lacks a clear focus and organization, and the conclusion is abrupt.</b>

Table 4: Pairs of *good* and *bad* explanations by type. From top to bottom, the source of low-quality is CONCISENESS, PLAUSIBILITY, and STANCE CLARITY.

in Rubrik’s hierarchy. In other words, annotators tended to generate a COMMENTARY when their answers were incorrect whereas a JUSTIFICATION was primarily associated with correct answers and corresponded to the highest accuracy. T2, however, shows the opposite trend. Specifically, LLMs tend to generate an ARGUMENT (highest type in our hierarchy) whenever they answered incorrectly while humans generated a COMMENTARY. We hypothesise that the uneven behaviour on this task is due to the multi-label nature of T2. A similar variation was observed when we looked at the frequencies of the answer choices picked by the annotators (see Appendix F.1).

**Explanation Quality Breakdown.** Regarding the quality of the explanations, the number of *bad* explanations was low and concentrated in COMMENTARIES across tasks. The analysis of sublabel frequencies (plotted in Figure 3) showed that the main source of a bad explanation was the lack of CONCISENESS, with open-source LLMs averaging 96.89% and closed-source LLMs averaging 99.06% on this sublabel. An example is shown in Table 4; the COMMENTARY is redundant, due to the repetition of details given in the question’s

context. This contrasts with the low frequency of WORD CHOICE, COHESION, APPROPRIATENESS and GRAMMATICALITY. On the other hand, CONCISENESS is less of a problem to humans, whose explanations are mostly judged as bad due to poor COHERENCE. Human explanations were different between contractors and experts. Bad explanations produced by experts were due to GRAMMATICALITY, while contractors struggled with COHERENCE. Figure 9 in Appendix F provides a more granular view of these findings.

## 6 Conclusion

This work introduces Rubrik, a novel evaluation rubric for assessing the quality of explanations, and a dataset. CUBE, which includes diverse explanations across four tasks, served as the testbed for evaluating Rubrik’s effectiveness. Rubrik’s design, rooted in educational principles, applies insights from education, XAI, and NLG literature. Our work contributes to the responsible integration of GenAI into critical decision-making processes, providing a foundation for future advancements in explanation quality assessment.

## Limitations

**Scoring strategy.** Given the scope of this work, we opted for a binary evaluation strategy, categorising explanations as either *good* or *bad*. The task of establishing criteria for a *good* explanation presented a significant challenge, necessitating the identification and definition of relevant attributes. A more nuanced scoring system that reflects varying degrees of quality would be desirable. However, while a Likert scale might be a convenient choice, developing a valid and reliable graded scale specifically for explanations requires considerably more research. Our primary goal in this initial study was to assess the viability of our proposed rubric in its simplest form, laying the groundwork for more nuanced evaluations in future work. Furthermore, our approach does not explicitly assess the quality of reasoning itself. While a *good* explanation is generally an indicator of a good reasoning, a poor explanation could stem from how the reasoning is communicated rather than from the reasoning process itself. Although this is a complex problem, the development of methods for directly assessing reasoning quality is an interesting direction for future research.

**Monolingual Data.** The different attributes (DIMENSIONS) of a *good* explanation were taken from studies that exclusively considered English data. In turn, our work only includes datasets in English as well. In principle, the DIMENSIONS and definitions presented here should extend to other languages. However, it is possible that some will change depending on the cultural heritage, literature, and history. Indeed, the very concept of explanations may differ depending on the linguistic community, which may influence how explanation types, COMPONENTS or DIMENSIONS are prioritised or understood.

**Annotators' Confidence Assessment.** After completing the annotation tasks, human annotators were surveyed about their experience, including a self-assessment of their performance. These responses provided valuable context for interpreting the data analysis results. As for LLM annotators, they were prompted to assign probabilities reflecting their confidence in each answer option's correctness. While logit analysis would have been ideal, we hypothesised that requesting that information in the prompt would be sufficiently accurate, especially given that logit access was not available across all models (due to some being closed-

source). However, the resulting probabilities often failed to sum to 100%, indicating a lack of consistent or meaningful probability assignment. Consequently, these assigned probabilities were not considered in the data analysis. Thus, we lack the means to make meaningful comparisons between human and LLM annotator confidence levels.

## Ethical Considerations

Prior to commencing the study, ethical approval was obtained from a relevant Ethics Committee. Informed consent was obtained from all participants, and their anonymity/confidentiality was ensured throughout the research process.

In light of [Baur \(2020\)](#)'s critique of the current "AI hype", we acknowledge the potential for misinterpretation of GenAI capabilities, particularly the risk of users over-relying on automatic explanations in tasks where human oversight is crucial. Our work aims to mitigate this risk by providing an objective evaluation framework for model outputs. This framework enables informed decision-making regarding the selection of the most appropriate resource—whether human or automated—for a given task. For instance, Rubrik can identify instances where a less complex model is sufficient, or conversely, when human expertise is required.

Finally, we also recognise the potential for misuse of our framework. Indeed, Rubrik could be exploited to deliberately generate misleading or poor-quality explanations. This could contribute to the spread of misinformation which poses a serious threat to informed decision-making. This risk highlights the importance of ensuring that the tool is used responsibly.

## Acknowledgments

We thank Øistein Andersen and Andrew Caines for their help recruiting annotators and their constructive suggestions and advice throughout the project. We also thank Camélia Guerraoui for her help conducting the preliminary experiments. Many thanks to the labmates at the NLIP Lab at the University of Cambridge, especially Marie Bexte and Iman Jundi for taking the time to review an earlier draft of this paper. We are deeply grateful to the annotators whose meticulous work was crucial for building our dataset. Our thanks also go to Diane Nicholls and her skilled team of annotators at Cambridge University Press & Assessment. Finally, we greatly appreciate the anonymous reviewers for their in-

sightful feedback, which significantly strengthened this manuscript.

This paper reports on research supported by Cambridge University Press & Assessment, by the JSPS KAKENHI Grant Numbers 25K03175, by JST Moonshot R&D Program Grant Number JP-MJMS2236, and by The Nakajima Foundation.

The third author's contributions were primarily completed while employed at LegalOn Technologies. This paper has not undergone internal review or approval process of LegalOn Technologies.

## Contributions of the Authors

This project was a large collaboration that would not have happened without dedicated effort from every co-author.

The *idea of the project* originated in discussions among Pride Kavumba, Diana Galvan-Sosa and Keisuke Sakaguchi. However, Gabrielle Gaudeau's entry as co-first author was essential in leading the project and *designing Rubrik* with Diana Galvan-Sosa. Paula Buttery, as an advisor, provided valuable input on its design.

The *data selection* was primarily carried out by Diana Galvan-Sosa (T2), Gabrielle Gaudeau (T4), Pride Kavumba (T1) and Yunmeng Li (T3). For the *collection of explanations*, Diana Galvan-Sosa and Gabrielle Gaudeau led the collection of human-generated explanations. Pride Kavumba and Hongyi Gu led the collection of LLM-generated explanations.

Pride Kavumba and Hongyi Gu led the *experimental implementation*, with Diana Galvan-Sosa, Gabrielle Gaudeau and Yunmeng Li actively participating in the experimental design. Zheng Yuan provided crucial expert advice and suggestions that shaped the final design.

*Analysis* of the experimental results were first done by Yunmeng Li and Gabrielle Gaudeau, and later updated by Diana Galvan-Sosa.

All co-authors contributed to *writing the paper*, especially Diana Galvan-Sosa, Gabrielle Gaudeau, Pride Kavumba, Yunmeng Li and Hongyi Gu.

## References

Chirag Agarwal, Sree Harsha Tanneru, and Himabindu Lakkaraju. 2024. [Faithfulness vs. Plausibility: On the \(Un\)Reliability of Explanations from Large Language Models](#). *arXiv preprint*. ArXiv:2402.04614 [cs].

Michael Alley. 1996. [Language: Being Concise](#). In Michael Alley, editor, *The Craft of Scientific Writing*, pages 119–127. Springer, New York, NY.

Mohammad Amiryousefi and Hossein Barati. 2011. Metadiscourse: exploring interaction in writing, ken hyland. *Continuum, London. Elixir Literature*, 40:5245–5250.

Anthropic. 2024. [Introducing computer use, a new claude 3.5 sonnet, and claude 3.5 haiku](#). Accessed: February 2025.

Kenneth C. Arnold, Krysta Chauncey, and Krzysztof Z. Gajos. 2020. [Predictive text encourages predictable writing](#). In *Proceedings of the 25th International Conference on Intelligent User Interfaces, IUI '20*, pages 128–138, New York, NY, USA. Association for Computing Machinery.

Paris Arnold. 2023. IELTS Writing Band Descriptors.

Kaan Y. Balta, Arshia P. Javidan, Eric Walser, Robert Arntfield, and Ross Prager. 2025. [Evaluating the Appropriateness, Consistency, and Readability of ChatGPT in Critical Care Recommendations](#). *Journal of Intensive Care Medicine*, 40(2):184–190. Publisher: SAGE Publications Inc STM.

Siu Wing Yee Barbara, Muhammad Afzaal, and Heshah Saleh Aldayel. 2024. A corpus-based comparison of linguistic markers of stance and genre in the academic writing of novice and advanced engineering learners. *Humanities and Social Sciences Communications*, 11(1):1–10.

Dorothea Baur. 2020. [Four reasons why hyping AI is an ethical problem](#). Accessed: February 14, 2025.

Robert De Beaugrande and Wolfgang U. Dressler. 1981. *Introduction to Text Linguistics*. Longman. Google-Books-ID: TmhiAAAAMAAJ.

Leema Berland and Katherine McNeill. 2012. [For whom is argument and explanation a necessary distinction? A response to Osborne and Patterson](#). *Science Education*, 96:808–813.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, and Marcos Zampieri. 2016. [Findings of the 2016 Conference on Machine Translation](#). In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 131–198, Berlin, Germany. Association for Computational Linguistics.

Antoine C Braet. 1992. Ethos, pathos and logos in aristotle's rhetoric: A re-examination. *Argumentation*, 6:307–320.

- Ana Brassard, Benjamin Heinzerling, Keito Kudo, Keisuke Sakaguchi, and Kentaro Inui. 2024. [ACORN: Aspect-wise Commonsense Reasoning Explanation Evaluation](#). *arXiv preprint ArXiv:2405.04818* [cs].
- Ingo Brigandt. 2016. [Why the Difference Between Explanation and Argument Matters to Science Education](#). *Science & Education*, 25.
- Sylvain Bromberger. 1992. *On what we know we don't know: Explanation, theory, linguistics, and how questions shape them*. University of Chicago Press.
- Gavin Brown. 2010. [The Validity of Examination Essays in Higher Education: Issues and Responses](#). *Higher Education Quarterly*, 64:276–291.
- Priscila G Brust-Renck, Rebecca B Weldon, and Valerie F Reyna. 2021. Judgment and decision making.
- Christopher Bryant, Mariano Felice, Øistein E. Andersen, and Ted Briscoe. 2019. [The BEA-2019 Shared Task on Grammatical Error Correction](#). In *Proceedings of the Fourteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 52–75, Florence, Italy. Association for Computational Linguistics.
- Christopher Bryant, Zheng Yuan, Muhammad Reza Qorib, Hannan Cao, Hwee Tou Ng, and Ted Briscoe. 2023. [Grammatical Error Correction: A Survey of the State of the Art](#). *Computational Linguistics*, pages 643–701. Place: Cambridge, MA Publisher: MIT Press.
- Aljoscha Burchardt. 2013. [Multidimensional quality metrics: a flexible system for assessing translation quality](#). In *Proceedings of Translating and the Computer 35*, London, UK. Aslib.
- Chris Callison-Burch, Cameron Fordyce, Philipp Koehn, Christof Monz, and Josh Schroeder. 2007. [\(Meta-\) Evaluation of Machine Translation](#). In *Proceedings of the Second Workshop on Statistical Machine Translation*, pages 136–158, Prague, Czech Republic. Association for Computational Linguistics.
- Michael Canale. 1983. From communicative competence to communicative language pedagogy 1. In *Language and Communication*. Routledge. Num Pages: 26.
- Mengyun Cao and Hai Zhuge. 2022. [Automatic evaluation of summary on fidelity, conciseness and coherence for text summarization based on semantic link network](#). *Expert Systems with Applications*, 206:117777.
- Christine Chin and David E. Brown. 2000. [Learning in Science: A Comparison of Deep and Surface Approaches](#). *Journal of Research in Science Teaching*, 37(2):109–138.
- Noam Chomsky. 1965. *Aspects of the Theory of Syntax*, 50 edition. The MIT Press.
- Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A Smith. 2021. All that's 'human' is not gold: Evaluating human evaluation of generated text. *arXiv preprint arXiv:2107.00061*.
- Jacob Cohen. 1960. [A coefficient of agreement for nominal scales](#). *Educational and Psychological Measurement*, 20:37 – 46.
- Cohere for AI. 2024. Introducing command r plus on microsoft azure. <https://cohere.com/blog/command-r-plus-microsoft-azure>. Accessed: 2025-02-14.
- James Collins. 1998. [Strategies for Struggling Writers](#). *College Composition and Communication*, 49:298.
- Ulla Connor. 1990. [Linguistic/Rhetorical Measures for International Persuasive Student Writing](#). *Research in the Teaching of English*, 24(1):67–87. Publisher: ncte.org.
- Scott Crossley, Yu Tian, Perpetual Baffour, Alex Franklin, Youngmeen Kim, Wesley Morris, Meg Benner, Aigner Picou, and Ulrich Boser. 2024. The English Language Learner Insight, Proficiency and Skills Evaluation (ELLIPSE) Corpus. *International Journal of Learner Corpus Research*. Status: forthcoming.
- Phillip Dawson. 2017. Assessment rubrics: towards clearer and more replicable design, research and practice. *Assessment & Evaluation in Higher Education*, 42(3):347–360.
- Randy Devillez. 2003. *Writing: Step by Step*. Kendall Hunt Publishing Company. Google-Books-ID: 79oAePQ7Of0C.
- Jean-Marc Dewaele. 2008. [“Appropriateness” in foreign language acquisition and use: Some theoretical, methodological and ethical considerations](#). 46(3):245–265. Publisher: De Gruyter Mouton Section: International Review of Applied Linguistics in Language Teaching.
- Finale Doshi-Velez and Been Kim. 2017. [Towards A Rigorous Science of Interpretable Machine Learning](#). *arXiv preprint*. ArXiv:1702.08608.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Lee Dunn, Chris Morgan, Meg O'Reilly, and Sharon Parry. 2003. *The Student Assessment Handbook: New Directions in Traditional and Online Assessment*. Routledge, London.
- M. Expósito-Ruiz, S. Pérez-Vicente, and F. Rivas-Ruiz. 2010. [Statistical inference: Hypothesis testing](#). *Allergologia et Immunopathologia*, 38(5):266–277. Publisher: Elsevier.



- Thomas Fel, David Vigouroux, Rémi Cadène, and Thomas Serre. 2022. How good is your explanation? algorithmic stability measures to assess the quality of explanations for deep neural networks. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 720–730.
- Yang Feng, Wanying Xie, Shuhao Gu, Chenze Shao, Wen Zhang, Zhengxin Yang, and Dong Yu. 2020. [Modeling Fluency and Faithfulness for Diverse Neural Machine Translation](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 59–66. ISSN: 2374-3468, 2159-5399 Issue: 01 Journal Abbreviation: AAAI.
- Anita Fetzer. 2012. [Textual coherence as a pragmatic phenomenon](#). In Kasia M. Jaszczolt and Keith Allan, editors, *The Cambridge Handbook of Pragmatics*, Cambridge Handbooks in Language and Linguistics, pages 447–468. Cambridge University Press, Cambridge.
- Anita Fetzer. 2018. Appropriateness in context.
- Freeman, Richard, Lewis, Roger (BP Professor of Learning Development, and University of Humberside). 2016. *Planning and Implementing Assessment*. Routledge, London.
- Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. [Experts, errors, and context: A large-scale study of human evaluation for machine translation](#). *Transactions of the Association for Computational Linguistics*, 9:1460–1474.
- Milton Friedman. 1940. [A Comparison of Alternative Tests of Significance for the Problem of m Rankings](#). *The Annals of Mathematical Statistics*, 11(1):86–92. Publisher: Institute of Mathematical Statistics.
- Silvia García-Méndez, Francisco de Arriba-Pérez, and María del Carmen Somoza-López. 2024. A review on the use of large language models as virtual tutors. *Science & Education*, pages 1–16.
- Leilani H Gilpin, David Bau, Ben Z Yuan, Ayesha Bajwa, Michael Specter, and Lalana Kagal. 2018. Explaining explanations: An overview of interpretability of machine learning. In *2018 IEEE 5th International Conference on data science and advanced analytics (DSAA)*, pages 80–89. IEEE.
- Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2013. [Continuous Measurement Scales in Human Evaluation of Machine Translation](#). In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 33–41, Sofia, Bulgaria. Association for Computational Linguistics.
- Sylviane Granger, Estelle Dagneaux, Fanny Meunier, and Magali Paquot. 2009. *International Corpus of Learner English. Version 2. Handbook and CD-ROM*.
- M. A. K. Halliday and Ruqaiya Hasan. 2014. [Cohesion in English](#). Routledge, London.
- Sandra G Hart. 2006. Nasa-task load index (nasa-tlx); 20 years later. In *Proceedings of the human factors and ergonomics society annual meeting*, volume 50, pages 904–908. Sage publications Sage CA: Los Angeles, CA.
- SG Hart. 1988. Development of nasa-tlx (task load index): Results of empirical and theoretical research. *Human mental workload/Elsevier*.
- Derrick Higgins, Jill Burstein, Daniel Marcu, and Claudia Gentile. 2004. [Evaluating Multiple Aspects of Coherence in Student Essays](#). In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 185–192, Boston, Massachusetts, USA. Association for Computational Linguistics.
- David M. Howcroft, Anya Belz, Miruna-Adriana Clinciu, Dimitra Gkatzia, Sadid A. Hasan, Saad Mahamood, Simon Mille, Emiel van Miltenburg, Sashank Santhanam, and Verena Rieser. 2020. [Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions](#). In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 169–182, Dublin, Ireland. Association for Computational Linguistics.
- Xinyu Hu, Mingqi Gao, Sen Hu, Yang Zhang, Yicheng Chen, Teng Xu, and Xiaojun Wan. 2024. [Are LLM-based Evaluators Confusing NLG Quality Criteria?](#) *arXiv preprint*. ArXiv:2402.12055 [cs].
- Jie Huang and Kevin Chen-Chuan Chang. 2023. [Towards reasoning in large language models: A survey](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1049–1065, Toronto, Canada. Association for Computational Linguistics.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2025. [A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions](#). *ACM Trans. Inf. Syst.*, 43(2):42:1–42:55.
- Mary Huba and Jann Freed. 2000. Learner-Centered Assessment on College Campuses: Sifting the Focus from Teaching to Learning. *Community College Journal of Research and Practice*, 24.
- Dell Hymes. 1972. On Communicative Competence. In *Sociolinguistics*, pages 269–293. Harmondsworth: Penguin.
- Alon Jacovi and Yoav Goldberg. 2021. [Aligning Faithful Interpretations with their Social Attribution](#). *Transactions of the Association for Computational Linguistics*, 9:294–310. Place: Cambridge, MA Publisher: MIT Press.
- Arshia P. Javidan, Tiam Feridooni, Lauren Gordon, and Sean A. Crawford. 2024. [Evaluating the progression](#)

- of artificial intelligence and large language models in medicine through comparative analysis of ChatGPT-3.5 and ChatGPT-4 in generating vascular surgery recommendations. *JVS-Vascular Insights*, 2:100049.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.
- Zhijing Jin, Abhinav Lalwani, Tejas Vaidhya, Xiaoyu Shen, Yiwen Ding, Zhiheng Lyu, Mrinmaya Sachan, Rada Mihalcea, and Bernhard Schoelkopf. 2022. [Logical Fallacy Detection](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 7180–7198, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Samia Kabir, David N. Udo-Imeh, Bonan Kou, and Tianyi Zhang. 2024. [Is Stack Overflow Obsolete? An Empirical Study of the Characteristics of ChatGPT Answers to Stack Overflow Questions](#). In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI '24, pages 1–17, New York, NY, USA. Association for Computing Machinery.
- Marzena Karpinska and Mohit Iyyer. 2023. [Large language models effectively leverage document-level context for literary translation, but critical errors persist](#). In *Proceedings of the Eighth Conference on Machine Translation*, pages 419–451, Singapore. Association for Computational Linguistics.
- Zixuan Ke and Vincent Ng. 2019. [Automated Essay Scoring: A Survey of the State of the Art](#). pages 6300–6308.
- Frank C Keil. 2006. Explanation and understanding. *Annu. Rev. Psychol.*, 57(1):227–254.
- Yoonsu Kim, Jueon Lee, Seoyoung Kim, Jaehyuk Park, and Juho Kim. 2024. Understanding users' dissatisfaction with chatgpt responses: Types, resolving tactics, and the effect of knowledge level. In *Proceedings of the 29th International Conference on Intelligent User Interfaces*, pages 385–404.
- Klaus Krippendorff. 2011. [Computing krippendorff's alpha-reliability](#).
- Cherise Kristoffersen. 2019. [Where do my words come from? Towards methods for analyzing word choice in primary level writing](#). *Apples - Journal of Applied Language Studies*, 13(3):59–75. Number: 3.
- Kristopher Kyle, Scott Crossley, and Cynthia Berger. 2018. [The tool for the automatic analysis of lexical sophistication \(TAALES\): version 2.0](#). *Behavior Research Methods*, 50(3):1030–1046.
- Kristopher Kyle and Scott A. Crossley. 2015. [Automatically Assessing Lexical Sophistication: Indices, Tools, Findings, and Application](#). *TESOL Quarterly*, 49(4):757–786.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. [RACE: Large-scale ReAding comprehension dataset from examinations](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 785–794, Copenhagen, Denmark. Association for Computational Linguistics.
- Shengjie Li and Vincent Ng. 2024. [ICLE++: Modeling Fine-Grained Traits for Holistic Essay Scoring](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 8465–8486, Mexico City, Mexico. Association for Computational Linguistics.
- Yunmeng Li, Jun Suzuki, Makoto Morishita, Kaori Abe, and Kentaro Inui. 2025. [MQM-chat: Multi-dimensional quality metrics for chat translation](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 3283–3299, Abu Dhabi, UAE. Association for Computational Linguistics.
- Tania Lombrozo. 2006. The structure and function of explanations. *Trends in cognitive sciences*, 10(10):464–470.
- Elizabeth Cloninger Long. 2007. [College writing resources with readings](#). New York : Pearson/Longman.
- Andrea A Lunsford, Kirt H Wilson, and Rosa A Eberly. 2008. *The SAGE handbook of rhetorical studies*. Sage Publications.
- Valerie R Mariana. 2014. *The Multidimensional Quality Metric (MQM) framework: A new framework for translation quality assessment*. Brigham Young University.
- Marianna Martindale, Marine Carpuat, Kevin Duh, and Paul McNamee. 2019. [Identifying Fluently Inadequate Output in Neural and Statistical Machine Translation](#). In *Proceedings of Machine Translation Summit XVII: Research Track*, pages 233–243, Dublin, Ireland. European Association for Machine Translation.
- Sandeep Mathias and Pushpak Bhattacharyya. 2018. [ASAP++: Enriching the ASAP Automated Essay Grading Dataset with Essay Attribute Scores](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Philip M. McCarthy and Scott Jarvis. 2007. [vocd: A theoretical and empirical evaluation](#). *Language Testing*, 24(4):459–488.
- Danielle McNamara and Com. 2010. *Cohesion, coherence, and expert evaluations of writing proficiency*. Journal Abbreviation: Proceedings of the 32nd Annual Conference of the Cognitive Science Society

- Publication Title: Proceedings of the 32nd Annual Conference of the Cognitive Science Society.
- Danielle S. McNamara, Arthur C. Graesser, Philip M. McCarthy, and Zhiqiang Cai. 2014. *Automated Evaluation of Text and Discourse with Coh-Metrix*. Cambridge University Press. Google-Books-ID: xSPeAgAAQBAJ.
- Katharine L. McNeill and Joseph Krajcik. 2007. Middle school students' use of appropriate and inappropriate evidence in writing scientific explanations. In *Thinking with data*, Carnegie Mellon symposia on cognition, pages 233–265. Lawrence Erlbaum Associates Publishers, Mahwah, NJ, US.
- Katherine McNeill, David Lizotte, Joseph Krajcik, and Ronald Marx. 2006. [Supporting Students' Construction of Scientific Explanations by Fading Scaffolds in Instructional Materials](#). *Journal of the Learning Sciences*, 15:153–191.
- Katherine L. McNeill and Joseph Krajcik. 2008. [Scientific explanations: Characterizing and evaluating the effects of teachers' instructional practices on student learning](#). *Journal of Research in Science Teaching*, 45(1):53–78.
- Tim Miller. 2019a. Explanation in artificial intelligence: Insights from the social sciences. *Artificial intelligence*, 267:1–38.
- Tim Miller. 2019b. [Explanation in artificial intelligence: Insights from the social sciences](#). *Artificial Intelligence*, 267:1–38.
- Eleni Miltsakaki. 2004. [Evaluation of text coherence for electronic essay scoring systems](#). *Natural Language Engineering*, 10:25–55.
- Brian North and Enrica Piccardo. 2020. *Common European Framework of Reference for Languages: Learning, Teaching, Assessment Common European Framework of Reference for Languages: Learning, Teaching, Assessment. Companion volume Language Policy Programme Education Policy Division Education Department Council of Europe*.
- Stellan Ohlsson. 2002. Generating and understanding qualitative explanations. In *The psychology of science text comprehension*, pages 91–128. Lawrence Erlbaum Associates Publishers, Mahwah, NJ, US.
- OpenAI. 2024. [Hello gpt-4o](#). Accessed: February 2025.
- Arjun Panickssery, Samuel Bowman, and Shi Feng. 2024. Llm evaluators recognize and favor their own generations. *Advances in Neural Information Processing Systems*, 37:68772–68802.
- Dojun Park and Sebastian Padó. 2024. [Multi-dimensional machine translation evaluation: Model evaluation and resource for Korean](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 11723–11744, Torino, Italia. ELRA and ICCL.
- Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. 2018. Multimodal explanations: Justifying decisions and pointing to the evidence. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8779–8788.
- Isaac Persing and Vincent Ng. 2013. [Modeling Thesis Clarity in Student Essays](#). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 260–269, Sofia, Bulgaria. Association for Computational Linguistics.
- Isaac Persing and Vincent Ng. 2015. [Modeling Argument Strength in Student Essays](#). In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 543–552, Beijing, China. Association for Computational Linguistics.
- Maxime Peyrard. 2019. [A Simple Theoretical Model of Importance for Summarization](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1059–1073, Florence, Italy. Association for Computational Linguistics.
- J. B. Pride. 1972. *Sociolinguistics : selected readings*. Harmondsworth, Penguin.
- Philip Quinn and Shumin Zhai. 2016. [A Cost-Benefit Study of Text Entry Suggestion Interaction](#). In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, CHI '16, pages 83–88, New York, NY, USA. Association for Computing Machinery.
- Malik Sallam. 2023. Chatgpt utility in healthcare education, research, and practice: systematic review on the promising perspectives and valid concerns. In *Healthcare*, volume 11, page 887. MDPI.
- William A. Sandoval. 2003. [Conceptual and Epistemic Aspects of Students' Scientific Explanations](#). *The Journal of the Learning Sciences*, 12(1):5–51. Publisher: Taylor & Francis, Ltd.
- Yash Saxena, Sarthak Chopra, and Arunendra Mani Tripathi. 2024. Evaluating consistency and reasoning capabilities of large language models. *arXiv preprint arXiv:2404.16478*.
- Yi Song, Michael Heilman, Beata Beigman Klebanov, and Paul Deane. 2014. [Applying Argumentation Schemes for Essay Scoring](#). In *Proceedings of the First Workshop on Argumentation Mining*, pages 69–78, Baltimore, Maryland. Association for Computational Linguistics.
- Andrea Sottana, Bin Liang, Kai Zou, and Zheng Yuan. 2023. [Evaluation metrics in the era of GPT-4: Reliably evaluating large language models on sequence to sequence tasks](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8776–8788, Singapore. Association for Computational Linguistics.



- Micol Spitale, Minja Axelsson, and Hatice Gunes. 2024. [Appropriateness of LLM-equipped Robotic Well-being Coach Language in the Workplace: A Qualitative Evaluation](#). *arXiv preprint*. ArXiv:2401.14935 [cs].
- Christian Stab and Iryna Gurevych. 2014. [Annotating Argument Components and Relations in Persuasive Essays](#). In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 1501–1510, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.
- Manfred Stede. 2002. [Lexical Choice Criteria in Language Generation](#).
- Dannelle D. Stevens and Antonia J. Levi. 2004. *Introduction to Rubrics: An Assessment Tool to Save Grading Time, Convey Effective Feedback and Promote Student Learning*. Stylus Publishing, LLC. Publication Title: Stylus Publishing, LLC ERIC Number: ED515062.
- Susan Strauss and Parastou Feiz. 2013. [Discourse analysis: Putting our worlds into words](#). *Discourse Analysis: Putting our Worlds into Words*, pages 1–411.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. 2024. [Gemma 2: Improving open language models at a practical size](#). *arXiv preprint arXiv:2408.00118*.
- Antonio Toral and Víctor M. Sánchez-Cartagena. 2017. [A Multifaceted Evaluation of Neural versus Phrase-Based Machine Translation for 9 Language Directions](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 1063–1073, Valencia, Spain. Association for Computational Linguistics.
- Stephen Toulmin. 1958. *The Uses of Arguments*, 1 edition. Cambridge University Press.
- Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C. J. Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, and Paul van Mulbregt. 2020. [SciPy 1.0: fundamental algorithms for scientific computing in Python](#). *Nature Methods*, 17(3):261–272. Publisher: Nature Publishing Group.
- Barbara E. Walvoord and Virginia Johnson Anderson. 1998. *Effective Grading: A Tool for Learning and Assessment*. Jossey-Bass Publishers, 350 Sansome St. ERIC Number: ED416810.
- Xiaojun Wan, Jianwu Yang, and Jianguo Xiao. 2007. [Manifold-ranking based topic-focused multi-document summarization](#). In *Proceedings of the 20th international joint conference on Artificial intelligence, IJCAI’07*, pages 2903–2908, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Johnny Tian-Zheng Wei, Khiem Pham, Brian Dillon, and Brendan O’Connor. 2018. [Evaluating Syntactic Properties of Seq2seq Output with a Broad Coverage HPSG: A Case Study on Machine Translation](#). *arXiv preprint*. ArXiv:1809.02035 [cs].
- Sarah Wiegrefe and Ana Marasović. 2021. [Teach me to explain: A review of datasets for explainable natural language processing](#). In *NeurIPS Datasets and Benchmarks*.
- Yuxiang Wu and Baotian Hu. 2018. [Learning to Extract Coherent Summary via Deep Reinforcement Learning](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1). Number: 1.
- Helen Yannakoudakis, Øistein Andersen, Ardeshir Geranpayeh, Ted Briscoe, and Diane Nicholls. 2018. [Developing an automated writing placement system for ESL learners](#). *Applied Measurement in Education*, 31.
- Jin-ge Yao, Xiaojun Wan, and Jianguo Xiao. 2017. [Recent advances in document summarization](#). *Knowledge and Information Systems*, 53(2):297–336.
- Laura Zangori, Cory T. Forbes, and Mandy Biggers. 2013. [Fostering student sense making in elementary science learning environments: Elementary teachers’ use of science curriculum materials to promote explanation construction](#). *Journal of Research in Science Teaching*, 50(8):989–1017.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. [HellaSwag: Can a machine really finish your sentence?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.
- Jiegen Zhang. 2006. [A Text-based Approach to Cohesion and Coherence](#). Ph.D. thesis.
- Muru Zhang, Ofir Press, William Merrill, Alisa Liu, and Noah A Smith. 2023a. [How language model hallucinations can snowball](#). *arXiv preprint arXiv:2305.13534*.
- Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2023b. [Siren’s Song in the AI Ocean: A Survey on Hallucination in Large Language Models](#). *arXiv preprint*. ArXiv:2309.01219 [cs].



## A Rubric Creation

### A.1 Components

As the foundational type in Rubrik, a COMMENTARY embodies the most basic type of explanation, with its primary objective being to provide an understanding of a decision-making process. Throughout this work, we assume a situation where there is an explicit set of choices, and one choice is selected over the others. Then, a decision is the behavioural ACTION of choosing among alternative options (Brust-Renck et al., 2021) and it is complemented by the REASON that guided that choice. If there is EVIDENCE to support the decision, a COMMENTARY then transitions to a JUSTIFICATION. Note that in either case, the underlying principle of objectivity remains consistent across both types. A subjective approach to presenting a decision process shifts the main goal of *understanding* the underlying rationale to *persuading* the audience. This idea aligns with the definition of an ARGUMENT, which is the result of an activity aimed at convincing a reasonable critic of the acceptability of a standpoint (Lunsford et al., 2008).

When considering the nature of argumentation, it is common to refer to the seminal work of Toulmin (1958), who provided a framework for constructing, analysing, and evaluating arguments. However, we adopt a different perspective, drawing upon the principles of rhetoric. Although there are some similarities between WARRANT–REASON and BACKING–EVIDENCE, this does not hold for the relationship between CLAIM–ACTION. In Toulmin’s framework, a warrant supports the claim and the backing further supports the warrant, but a claim is always assumed to be linked to a *standpoint*. Rhetorical argumentation, on the other hand, commonly refers to Aristotle’s trio *ethos-logos-pathos* (Braet, 1992), where *ethos* refers to the credibility of the speaker, *pathos* refers to the emotional state of the audience and *logos* refers to what is true. We can identify a relationship between LOGOS–COMMENTARY through the REASON component and ETHOS–JUSTIFICATION through EVIDENCE. It is then left to PATHOS to introduce the elements of *persuasion*. Considering that a stance is usually implicit in discourse, we focus on linguistic markers: metadiscourse features used by writers to express stance (Barbara et al., 2024). Thus, we merge into one component the essence of *pathos*, usually expressed in discourse through AFFECTIVE APPEAL(S), and features from Hyland’s Interper-

sonal Model of Metadiscourse (Amiryousefi and Barati, 2011): hedges, boosters, attitude and engagement markers (i.e., QUALIFIERS).

### A.2 Dimensions

We conducted an extensive review of NLP literature including work in Natural Language Generation (NLG) such as Machine Translation (MT) and Educational NLP (including Grammatical Error Correction and Automated Essay Scoring), but also in Linguistics and Cognitive Science. In doing so, we recorded the names of qualities (or DIMENSIONS) that people have looked for in explanations or argumentative writing more generally, and, when present, their definitions. We also kept note of how these qualities have been evaluated in a target text, using either human annotators or automated methods. See Table 5 for the exhaustive list.

DIMENSION NAME
APPROPRIATENESS
Adequacy
Clarity
COHERENCE
COHESION
Completeness
CONCISENESS
Consistency
Comprehensibility
Comprehensiveness
Correctness
Factuality
Faithfulness
Fidelity
Fluency
GRAMMATICALITY
Interpretability
Organisation
Persuasiveness
PLAUSIBILITY
Readability
Reasonableness
Transparency
Truth of likelihood
Usefulness
WORD CHOICE

Table 5: Exhaustive list of the quality DIMENSIONS of explanation we found when surveying the literature. We highlight in CAPITAL LETTERS the names of the DIMENSIONS we included in our rubric *verbatim*.

Below we describe how we defined and chose the eight DIMENSIONS that are represented in Rubrik. We also introduce a few of the many qualities that were considered and explain why they were excluded, as a demonstration of our overall process. Though we cannot be exhaustive at this time, we rigorously researched each and every one of the

dimensions mentioned in Table 5. The final definitions we used in the automated evaluation prompts are provided in Appendix E. The full rubric with examples is shown in Section A.13.

### A.3 Grammaticality

GRAMMATICALITY, though essential, was surprisingly hard to define. This was largely due to the fact that grammar has a long-standing tradition in a variety of fields—including Linguistics, Psychology, Education, and Cognitive Science—which have each contributed different perspectives and theories over time. As a result there is no single, universally accepted definition. Definitions which originate from the field of Linguistics tend to be highly theoretical, and as a result, quite impractical. A classic example is Chomsky (1965, Chapter 1, p.2) for whom the “grammar of a language purports to be a description of the ideal speaker-hearer’s intrinsic competence”, which has been criticised for being too abstract and disconnected from actual language use (Pride, 1972, Chapter 18). On the other hand, most NLP studies assume that the definition of GRAMMATICALITY is common knowledge and avoid going through the trouble of formally defining it in the context of their work (e.g., Wei et al., 2018). In fact, it is openly admitted that “Grammatical Error Correction” is something of a misnomer as it is now commonly understood to encompass errors that are not always strictly grammatical in nature” (Bryant et al., 2023).

However, to avoid relying on our intuition of what a grammatical explanation is, we needed to bridge the gap between theory and practice, and find a definition that could be both pragmatic and grounded in the literature. We did find one in a paper by Hu et al. (2024, Table 10), similarly focused on the evaluation of LLM outputs, which defines GRAMMATICALITY as measuring “whether the target text is grammatically correct without any lexical or syntax errors, regardless of its content and meaning. Consider whether the target text itself complies with the English standard usage and rules of grammar, such as tense errors, misspellings, incorrect prepositions, collocation misusages, and so on.” In using this definition, it is quite straightforward to classify GRAMMATICALITY as a **Language** DIMENSION as it in no way attends to the content of the text.

### A.4 Conciseness

In contrast, we found CONCISENESS to be well-documented across many literatures and much less controversial. In Education, “concise writing gets to the point quickly and does not introduce unnecessary information” (Long, 2007, p.25) and requires you to “cut fat” into your writing by “eliminating redundancies, eliminating writing zeroes, reducing sentences to simplest form, and cutting bureaucratic waste” (Alley, 1996, Chapter 8). Similarly, in NLP, Cao and Zhuge (2022) define it as a measure of “non-redundancy” in text, sometimes through the number of repeated words (Peyrard, 2019) or through computing sentence similarities (Wan et al., 2007).

We finally opted for Kabir et al. (2024)’s comprehensive taxonomy of three conciseness issues:

*Redundant* sentences reiterate information stated in the question or in other parts of the answer. *Irrelevant* sentences talk about concepts that are out of the scope of the question being asked. And lastly, *Excess* sentences provide information that is not required to understand the answer.

Not only were these issues identified when evaluating ChatGPT answers, a task closely related to ours, we additionally felt that they encompassed all the elements that were individually picked out in previous definitions. Note that since this definition is concerned with redundant, irrelevant or excess information, not just language, we decided to classify CONCISENESS as a **Content** dimension.

### A.5 Fluency

For a while, we considered fluency, an important notion in Machine Translation, which is generally evaluated by humans (e.g., Callison-Burch et al., 2007; Graham et al., 2013; Bojar et al., 2016), or using automated metrics (e.g., Toral and Sánchez-Cartagena, 2017; Martindale et al., 2019; Feng et al., 2020). In the first case, we found that human annotators were almost never provided with a proper definition of fluency and expected to use their intuition of what the word meant via prompts like “how do you judge the fluency of this translation?” in Callison-Burch et al. (2007) or “read the text below and rate it by how much you agree that: the text is fluent English” in (Graham et al., 2013). In the latter case, the metrics used were

only considered to be proxies for fluency which was never actually defined.

As with GRAMMATICALITY, [Hu et al. \(2024, Table 9\)](#) provided the following definition: “[fluency] measures the quality of individual sentences, are they grammatically correct, non-repetitive, and in accord with common English usage, with clear meanings”, which seemed to overlap both our definitions for CONCISENESS and GRAMMATICALITY. Since our goal was to reach a set of well-delineated, atomic dimensions, we chose to discard it.

## A.6 Cohesion

COHESION is a very important notion in Linguistics and is classically defined by [Halliday and Hasan \(2014, p.4\)](#) as:

occur[ring] where the INTERPRETATION of some element in the discourse is dependent on that of another. The one PRESUPPOSES the other, in the sense that it cannot be effectively decoded except by recourse to it. When this happens, a relation of cohesion is set up, and the two elements, the presupposing and the presupposed, are thereby at least potentially integrated into the text.

Unfortunately, as with GRAMMATICALITY, this definition is not accessible to most people and is far too theoretical.

However, COHESION is also widely present in Education, particularly in writing assessment and teaching literature, due to the common idea that a written text’s quality is highly related to its level of COHESION ([McNamara and Com, 2010](#)). This belief is reflected in the literature about writing (e.g., [Collins, 1998](#), [Devillez, 2003](#)) and the rubrics that teachers use to assess writing (e.g., [Arnold, 2023](#); [Crossley et al., 2024](#)). It is notably defined by [McNamara and Com \(2010\)](#) as follows:

Cohesion refers to the presence or absence of explicit cues in the text that allow the reader to make connections between the ideas in the text. For example, overlapping words and concepts between sentences indicate that the same ideas are being referred to across sentences. Likewise, connectives such as ‘because’, ‘therefore’, and ‘consequently’, inform the reader that there are relationships between ideas and the nature of those relationships.

Or more simply as the “appropriate use of transition phrases” by [Ke and Ng \(2019, Table 1\)](#). For our purposes, we prefer these pragmatic definitions to those offered by Linguistics.

From these definitions, it seems that COHESION is only concerned with **Language** not the content of a text. In fact, the dimension has also been examined through automated tools like Coh-Metrix ([McNamara et al., 2014](#)) or TAACO ([Kyle and Crossley, 2015](#)), which use a compound of linguistic metrics like the Type Token Ratio (TTR; [McCarthy and Jarvis, 2007](#)) as proxies for COHESION.

## A.7 Coherence

A related notion to COHESION is COHERENCE. It has been defined in Linguistics as a “continuity of sense” by [Beaugrande and Dressler \(1981, p.84\)](#), or more concretely as “the state of being logically consistent and connected” ([Fetzer, 2012](#)). It is also an important notion in Document Summarisation, where COHERENCE is similarly defined as “what makes multiple sentences semantically, logically and syntactically coherent” ([Yao et al., 2017](#)). It is also frequently evaluated writing assessment either by humans (e.g., [Higgins et al., 2004](#)) or via automated methods (e.g., [Higgins et al., 2004](#); [Mitsakaki, 2004](#); [Wu and Hu, 2018](#)).

Where COHESION is an “overt (or explicit) linguistic-surface phenomenon, [...] coherence is a covert (or implicit) deep-structure phenomenon”. But while COHERENCE is more concerned with meaning (i.e., **Content**) than form ([Fetzer, 2012](#)), it also “depends on a number of factors, including explicit cohesion cues, implicit cohesion cues (which are more closely linked to text coherence than are explicit cues), and nonlinguistic factors such as prior knowledge and reading skill” ([Kyle and Crossley, 2015](#)). They are thus “interdependent” notions ([Zhang, 2006](#)). To portray this in our rubric, we chose to similarly relate both DIMENSIONS: an explanation should thus not be labelled as coherent without first being judged as cohesive.

## A.8 Clarity

We first encountered this quality while looking at writing education papers, where clarity generally “refers to how clearly an author explains the thesis of her essay, i.e., the position she argues for with respect to the topic on which the essay is written” ([Persing and Ng, 2013](#)). It also appears in the ICLE++ corpus of persuasive student essays ([Granger et al., 2009](#); [Li and Ng, 2024](#)), an

important dataset in the field of Automated Written Assessment. However, the definitions we found were far too vague and we struggled to find more formal or practical descriptions of the term which seemed to support Beaugrande and Dressler (1981, Chapter 2)’s claim that clarity is “too vague and subjective to be reliably defined and quantified”. We ultimately decided to drop this DIMENSION.

## A.9 Word Choice

The WORD CHOICE DIMENSION is broadly defined as “the choice and aptness of the vocabulary used” (Mathias and Bhattacharyya, 2018). It is frequently included in written assessment rubrics (e.g., see the very detailed 6-point rubric for this dimension in the ASAP<sup>6</sup> corpus) and the focus of automated assessment research (e.g., Kyle and Crossley, 2015; Kyle et al., 2018; Kristoffersen, 2019).

We also came across Stede (2002)’s work on lexical choice for NLG:

Generally speaking, the point of “interesting” language generation (that is, more than merely mapping semantic elements one-to-one onto words) is **to tailor the output to the situation at hand**, where “situation” is to be taken in the widest sense, including the regional setting, the topic of the discourse, the social relationships between discourse participants, etc.

Though not explicitly defining WORD CHOICE, the above citation introduces the idea that every “interesting” or *good* utterance (or in our case, explanation) is made within a given “situation” and thus evaluating the language of that utterance should be context-dependent. It is this **context** that dictates what is “apt” (Mathias and Bhattacharyya, 2018). Realising that it is necessary to define an evaluation *context* before starting any kind of evaluation (see Section 3.3) was a turning point for our rubric.

Now, *context*-appropriateness relies on both form and content. However, due to the strong emphasis on evaluating WORD CHOICE as a surface-level feature, not a content one, in automated assessment research, we chose to classify it as a **Language** DIMENSION.

<sup>6</sup> The original dataset and annotation guidelines can be downloaded from <https://www.kaggle.com/c/asap-aes/data>.

## A.10 Appropriateness

APPROPRIATENESS defined in Linguistics by Canale (1983) as “the extent to which particular communicative functions [...] and **ideas** are judged to be proper in a given situation” or as “an optimal mapping between context and speech, or as ‘natural speech,’ is also connected intrinsically with the sociocultural notions of politeness and impoliteness” by Fetzer (2018). This term also occasionally appears in AI literature as something we must ensure in the systems we develop, and thus, evaluate (e.g., Spitale et al., 2024; Javidan et al., 2024; Balta et al., 2025;). There, it is more often related to other qualities such as safety, consistency, and readability. Hence, APPROPRIATENESS is a complex, multi-faceted dimension which also relies on *context*.

For our purpose, we needed to relate this DIMENSION to WORD CHOICE. For this, we turned to the prominent sociolinguist, Dell Hymes who “pointed out that appropriateness [depend] both on linguistic and sociocultural competence” (Dewaele, 2008), and defined it as “what to say to whom in what circumstances and how to say it” in Hymes (1972, p.277). We deem that this last part, “how to say it” is already encompassed by our definition of WORD CHOICE. Further, “to whom in what circumstances” refers to our very own definition of the *context*, which leaves us with the “what to say” for APPROPRIATENESS, that is, the **Content**.

## A.11 Plausibility

In reading around the topic of explanations in AI, we came across the following trait: “the **truth of likelihood** of an explanation is considered an important criterion of a good explanation” in a paper by Miller (2019b). The term was used to refer to facts that were judged as “either true or likely to be true by the explainee.” We note that in no way is our rubric intended to evaluate the truth condition of explanations. However, we felt that it was important that our rubric allows for JUSTIFICATION to be evaluated as *bad* or of *bad* quality if their EVIDENCE was deemed implausible by the evaluator. After some research, we could not find any other mention of the “truth of likelihood” and sought a more general name for our DIMENSION.

A related notion was PLAUSIBILITY which was present in similar literature and already being used to evaluate explanations. For instance, Agarwal et al. (2024) who define plausible explanations as



being “seemingly logical and coherent to human users” or as “being convincing towards the model prediction, regardless of whether the model was correct or whether the interpretation is faithful” by [Jacovi and Goldberg \(2021\)](#). Though not exactly similar, the latter introduces the idea that using PLAUSIBILITY as criteria for a *good* explanation might encourage deception. As a result, the authors advise against pursuing this DIMENSION.

Taking this warning into consideration, it was important to us to centre our definition of PLAUSIBILITY around the EVIDENCE component (2.a), and we modified [Agarwal et al. \(2024\)](#)’s Definition 1, substituting the word “explanation” with “evidence”:

An evidence\* is considered plausible if it is coherent with human reasoning and understanding.

### A.12 Stance Clarity

Whenever we found a mention of ARGUMENTS in the literature, the concept of persuasiveness was almost always mentioned. It thus seemed natural that it would be included in our rubric. We first looked at the notion of “argument strength” in persuasive writing which is defined, in an admittedly very circular fashion, as “the strength of the argument an essay makes for its thesis” and evaluated by [Persing and Ng \(2015\)](#). In a similar vein, we discovered work by [Song et al. \(2014\)](#) and [Stab and Gurevych \(2014\)](#) which designed argument schemes for annotating arguments manually in student essays. Yet, none of the definitions we found seemed right.

We then turned to persuasiveness in rhetoric, and found [Connor \(1990, Table 5\)](#)’s Persuasive Appeals Scale. Though very useful, we struggled to see whether these were in fact COMPONENTS or indeed a DIMENSION, and where to fit them in our rubric. After some iterations, we arrived at the fact that the presence of AFFECTIVE APPEALS and QUALIFIERS in an argument help us understand what the explainer’s “stance” is, that is, their personal “feeling, attitude, perspective, or position as enacted in discourse” ([Strauss and Feiz, 2013](#)). By that point, it felt like persuasiveness was too vague and we coined the term “Stance Clarity” for our last DIMENSION.

### A.13 Full Rubric

A concise overview of Rubrik is presented in Section 3.2, Table 1. This appendix provides the com-

plete details of the full-sized, illustrated rubric in Table 6 and Table 6 (contd.).

## B Data Selection

Considering the fact that the four datasets we chose to work with were all of different sizes, we chose to only work with a subset of each dataset: namely  $n = 1000$  instances for each task. Thus, our *base set* has a total of 4000 instances.

We collected a set of human-written (see Section C.1) and LLM-generated explanations (see Section C.2). Due to limitations in time and resources, only a subset of the 1000 instances was shown to the annotators: namely  $n = 110$  instances for each task. Thus, our *annotation set* has 440 instances. The following subsections detail the subset selection criteria.

### B.1 Commonsense Reasoning

**Base set.** Each CONTEXT in the HELLASWAG dataset is taken either from ActivityNet’s video captions or WikiHow’s how-to-articles. During the annotator’s training (see Section C.1.1), questions whose context made reference to a video were constantly flagged as “*not clear or ambiguous*”. Thus, we filtered instances that include the word “*camera*”, “*video*” or “*clip*”. After that, instances were selected randomly, making sure that the correct answers were distributed as evenly as possible across the four options (A-D), with roughly 25% assigned to each.

**Annotation set.** Since the *base set* already had an even distribution of the four answer choices, we selected a proportionally representative subset of 110 instances. See Table 7 for a summary of this selection process.

### B.2 Fallacy Detection

**Base set.** [Jin et al. \(2022\)](#) classified fallacies in the LOGIC dataset into 13 fallacy types. Due to potential overlap between some of the initial types and dataset imbalance, we focused on a subset of 7 types.

Selecting instances within the 30-300 character range effectively eliminated instances requiring specialised political or religious knowledge, ensuring consistent annotation based on general knowledge. After manual inspection, we removed some duplicated instances and statements that were not exactly fallacies, but rather someone’s opinion on a topic. We also identified a few instances that were

	COMPONENTS	DIMENSIONS	
		Language	Content
COMMENTARY	<p><b>ACTION (1.a):</b> does the explanation clearly indicate the decision or choice being made (e.g., specifying the selected answer)? For e.g.,</p> <ul style="list-style-type: none"> <li>Acceptable: “<b>The correct answer is A.</b>”</li> <li>Not acceptable: “Because it is the final part of the sequence.”</li> </ul>	<p><b>GRAMMATICALITY:</b> is the explanation grammatically correct, free of lexical or syntax errors? <i>Small typos are acceptable, but the errors should <b>not</b> impede comprehension in any way.</i> For e.g.,</p> <ul style="list-style-type: none"> <li>Acceptable: “The correct answer is A because nowadays our <b>society</b> is based on consumerism and the way in which we are producing is contaminating the <b>word.</b>”</li> <li>Not acceptable: “The correct answer is A because <b>now a day</b> our <b>society it is bassed in consumer, so that become the word more contaminate</b> to produce the products <b>that we demanding.</b>”</li> </ul>	<p><b>CONCISENESS:</b> is the explanation free of any redundant, irrelevant, or excess sentences (that is, not required to understand the answer)? For e.g., given that the answer choice D is “next she explains how to use the lawnmower and other tools and then she cuts the grass,”</p> <ul style="list-style-type: none"> <li>Acceptable: “The correct answer is D because it accurately reflects the sequence of events.”</li> <li>Not acceptable: “The correct answer is D because <b>she explains how to use the lawnmower and other tools, and then she cuts the grass.</b>”</li> </ul>
	<p><b>REASON (1.b):</b> does the explanation provide reasoning or insight into why the decision or choice was made, explaining the underlying logic or rationale for the Action? For e.g.,</p> <ul style="list-style-type: none"> <li>Acceptable: “The right answer is C, <b>because it is the final part of the sequence.</b>”</li> <li>Not acceptable: “The correct answer is A.”</li> </ul>	<p><b>WORD CHOICE:</b> is the language used in the explanation tailored to the given <i>context</i> (task, audience, purpose)? And are the sentences in the explanation well-formed? For e.g.,</p> <ul style="list-style-type: none"> <li>Acceptable: “The correct answer is A because the essay lacks fluency, has many incorrect clauses and missing words. And while the overall meaning can be deduced, the essay does not demonstrate an accurate grasp of language.”</li> <li>Not acceptable: “<b>Answer A.</b> lack of fluency, incorrect clauses and missing words, <b>meaning</b> can be found but does not demonstrate an accurate grasp of language.”</li> </ul> <p><b>COHESION:</b> does the explanation make appropriate use of transition phrases (e.g., connectives like “because”, “therefore”, and “consequently”, overlapping words across sentences, etc.)? For e.g.,</p> <ul style="list-style-type: none"> <li>Acceptable: “The correct answer is C because the man is on roller blades, not on a skateboard. Further, he is not talking to anyone and therefore cannot possibly ‘continue speaking’.”</li> <li>Not acceptable: “The correct answer is C, <b>because</b> the man is on roller blades, not a skateboard, <b>and</b> is not talking to anyone in the example <b>so cannot</b> ‘continue speaking’.”</li> </ul>	<p><b>APPROPRIATENESS:</b> is the explanation culturally appropriate, matching expectations for the given <i>context</i>? For e.g.,</p> <ul style="list-style-type: none"> <li>Acceptable: “The right answer is B because the tenses are properly used and the story makes sense.”</li> <li>Not acceptable: “The right answer is B because the tenses are properly accorded and <b>(within the slightly odd context)</b> the story makes sense.”</li> </ul> <p><b>COHERENCE:</b> does the explanation appropriately transition between ideas, i.e., does it make sense as a whole (e.g., good context-relatedness, semantic consistency, and inter-sentence causal and temporal dependencies, etc.)? For e.g., given the start of explanation “The correct answer is D, because no information about Liu’s relationship to science subjects specifically is given in the passage,”</p> <ul style="list-style-type: none"> <li>Acceptable: “therefore the fact that they like chemistry is implied and ambiguous.”</li> <li>Not acceptable: “therefore the fact that they like <b>cheese</b> is implied and ambiguous.”</li> </ul>

Table 6: Extended rubric with definitions and illustrative examples for each of the COMPONENTS and DIMENSIONS (continued on next page).



	COMPONENTS	DIMENSIONS	
		Language	Content
<b>JUSTIFICATION</b> 	<p><b>EVIDENCE (2.a):</b> does the explanation provide concrete evidence (can be both explicit or implicit) that supports the reasoning, such as information from the question's context or general knowledge? For e.g.,</p> <ul style="list-style-type: none"> <li>Acceptable: "The right answer is C, because it finishes the sequence, <b>describing the effect of bowling the ball and what happens as a result.</b>"</li> <li>Not acceptable: "The right answer is C, because is is the final part of the sequence."</li> </ul>		<p><b>PLAUSIBILITY:</b> is the provided <b>EVIDENCE</b> plausible and consistent with human reasoning, considering the context and general world knowledge? For e.g.,</p> <ul style="list-style-type: none"> <li>Acceptable: "The correct answer is A ('Jack picks the cheese') because <b>we are told that he enjoys eating 'mozzarella' in the morning.</b>"</li> <li>Not acceptable: "The correct answer is A ('Jack picks the cheese') because <b>my name is also Jack and I personally love cheese for breakfast.</b>"</li> </ul>
<b>ARGUMENT</b> 	<p><b>AFFECTIVE APPEAL(S) (3.a):</b> does the explanation use vivid, or emotionally charged language (e.g., metaphors) to evoke feelings in the audience? For e.g.,</p> <ul style="list-style-type: none"> <li>Acceptable: "The expression in the final section is very <b>heartfelt</b>; the tone is <b>excitable</b> and <b>keen</b> throughout."</li> <li>Not acceptable: "The final section reflects the writer's strong feelings on this issue."</li> </ul> <p><b>QUALIFIERS(S) (3.a):</b> does the explanation make use of hedges, boosters, attitude markers, self-mentions, or engagement markers? For e.g.,</p> <ul style="list-style-type: none"> <li>Acceptable: "The right answer is B, because the text is keeping with what is <b>presumably</b> a tour guide's voice: <b>intentionally</b> using clunky and <b>overly</b> expressive words."</li> <li>Not acceptable: "The right answer is B, because the text is keeping with the original tour guide's voice."</li> </ul>		<p><b>STANCE CLARITY:</b> is the explainer's stance (their personal feelings towards the task) clearly and unambiguously conveyed through affective appeals or qualifiers? Note that the stance can be implicit unlike the <b>ACTION</b>. For e.g.,</p> <ul style="list-style-type: none"> <li>Acceptable: "The correct answer is A (beginner) because this text is <b>undeniably</b> of a low English level."</li> <li>Not acceptable: "The correct answer is A (beginner) because this text is <b>clearly</b> of a low English level although the final section is <b>incredibly</b> well written."</li> </ul>

Table 6 (contd.): Extended rubric with definitions and illustrative examples for each of the COMPONENTS and DIMENSIONS.

Correct answer	Base set	Ann set
A	267	27
B	228	28
C	266	27
D	239	28
Total	1000	110

Table 7: Distribution of questions across each possible correct answer for T1’s *base set* and *annotation set*.

incorrectly labelled (i.e., were assigned the wrong fallacy type). Those were re-labelled and kept in the final subset. Table 8 shows the final distribution of our subset.

Logical Fallacy	Inc	Base set	Ann set
Faulty Generalisation	✓	289	17
Ad Hominem	✗		
Ad Populum	✗		
False Causality	✓	154	15
Circular Claim	✓	112	15
Appeal to Emotion	✓	109	15
Fallacy of Relevance	✗		
Deductive Fallacy	✓	120	15
Intentional Fallacy	✗		
Fallacy of Extension	✗		
False Dilemma	✓	118	17
Fallacy of Credibility	✓	95	16
Equivocation	✗		
Total		1000	110

Table 8: Distribution of instances across each fallacy type for T2’s *base set* and *annotation set*.

**Annotation set.** This task was originally framed as a classification task. For the purposes of this research, we adapted the task to follow an MCQ format, where the CONTEXT was the fallacy statement, and each of the fallacy types was listed as ANSWER CHOICES. We aimed for a balanced distribution of correct answers across the seven options (A-G). Instances were selected randomly from the *base set*. See Table 8 for a summary of this selection process.

### B.3 Reading Comprehension

**Base set.** RACE data is grouped by difficulty (RACE-M: middle school; RACE-H: high school). To better understand the dataset, authors subdivided questions into five reasoning categories. Since the *Passage Summarization* and *World Knowledge* do not fully require students to carefully read the passage to answer, we focused on the other three question types: *Detail Reasoning*, *Whole Picture Reasoning*, and *Attitude Analysis*. Specifically, answers to *Detail Reasoning* questions cannot simply

be found by matching the questions to the reading passages and require test-takers to provide reasons for their choices. For *Whole Picture Reasoning* questions test the students’ overall understanding of a story. *Attitude Analysis* questions ask about the opinions or attitudes of the author or characters of the reading passages.

Unfortunately, the questions have not been labelled with these reasoning categories in the published dataset; hence, we manually selected the data based on the description and examples given by Lai et al. (2017) and reviewed them to ensure quality.

Question type	Inc	Base set	Ann set
Detail reasoning	✓	400	36
Whole-picture reasoning	✓	400	37
Passage summarization	✗		
Attitude analysis	✓	200	37
World knowledge	✗		
Total		1000	110

Table 9: Distribution of text passages across each question type for T3’s *base set* and *annotation set*.

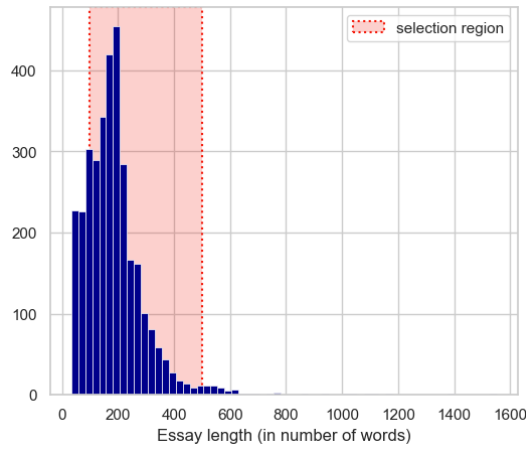
**Annotation set.** Each question in RACE has four answer choices (A-D). We aimed for a balanced distribution of instances of correct answers across options within each question type. Instances were randomly selected from the *base set*, targeting a proportion of approximately 25% per option. See Table 9 for a summary of this selection process.

### B.4 Essay Scoring

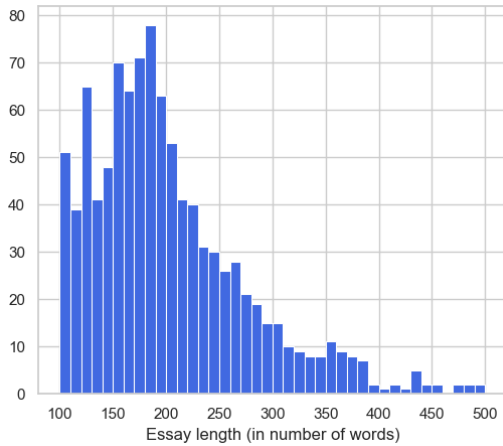
**Base set.** In the W&I corpus, essays range between 33 and 1,551 words in length. Figure 4a plots this distribution. We chose to exclude essays of less than 100 words, and more than 500 words, to avoid selecting essays sitting on either extreme of this distribution. Indeed, essays that are too short might contain too little information to be interesting to evaluate; essays that are long might exceed the limits of LLM contexts or prove too time-taking to annotate for humans. This step left us with a remaining total of 2,598 essays (833 A-scored essays, 1,039 B-scored essays, and 726 C-scored essays). Then, we randomly sampled 333 essays from each CEFR level group (334 for the B level) to obtain our *base set* of 1000 essays. We additionally randomly selected 3 essays (one of each CEFR level) from the remaining pool of essays to be used as examples in our experiments.

**Annotation set.** For our *annotation set*, we

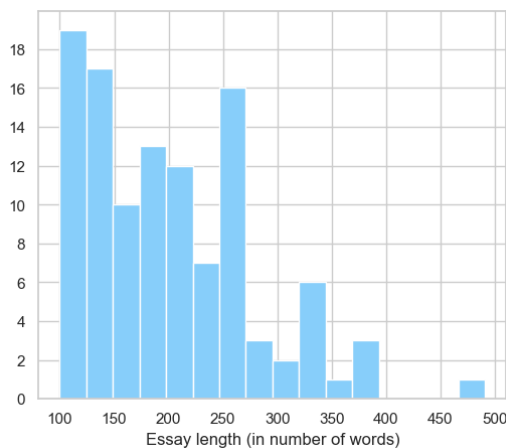




a ) W&I corpus word count distribution. We highlight in orange the region from which the *base set* essays were selected.



b ) *Base set* word count distribution.



c ) *Annotation set* word count distribution.

Figure 4: Plotting the word count distributions

again selected randomly from the *base set*, aiming for a balanced distribution of essays across the three CEFR levels. See Table 10 for a summary of this selection process.

Essay Grade	W&I	Base set	Ann set
A	1430	333	36
B	1100	334	37
C	770	333	37
Total	3300	1000	110

Table 10: Distribution of W&I essays across each CEFR level for T4’s *base set* and *annotation set*.

Essay Grade	W&I		Base set		Ann set	
	$\mu$	$\sigma$	$\mu$	$\sigma$	$\mu$	$\sigma$
A	125	70	163	56	150	51
B	211	100	207	73	205	71
C	262	132	235	71	245	77
Overall	186	113	201	73	201	78

Table 11: Mean ( $\mu$ ) and standard deviation ( $\sigma$ ) word count of the essays in the W&I corpus, the *base set*, and the *annotation set* (rounded to the nearest integer).

## C Data Collection

### C.1 Human Annotators

We recruited seven human annotators: four research assistants (contractors) and three professional annotators (experts). One of the main authors, along with a senior researcher, led the contractors’ recruiting efforts, which included conducting interviews with potential candidates. We selected individuals who appeared to have strong abilities in **attention to detail**, **assessment**, and **strong language skills**. These skills were essential for completing the assigned reasoning and language tasks. The PA’s were annotators who were specially trained EFL (English as a Foreign Language) teachers and examiners. The annotators were paid an hourly rate of £22.59 for their work. We anonymised the annotations by removing any personally identifiable information. Each annotator was identified with a randomly assigned ID (e.g., 000005FB, 000004E4)

#### C.1.1 Training

All annotators received a detailed annotation guide that introduced the four tasks and provided a number of annotated examples (question + answer choices + correct answer) for each task. The exam-

ples were intended to help them familiarise themselves with the tasks. Since T2 necessitates some familiarity with fallacious reasoning, this task was further supported by an appendix with definitions of all fallacy types.<sup>7</sup> We did not include explanations to avoid biasing the annotators as to what a *good* explanation should look like. The annotation guide also included a series of guidelines they should abide by during the annotation process.

Upon reading the annotation guide, the annotators were asked to write explanations for each of the annotated examples contained in the guide. Their explanations were then reviewed by two of the main authors to ensure they were acceptable in terms of format and length.<sup>8</sup> Unless absolutely necessary, annotators did not receive any feedback on their explanations.

Subsequently, each annotator received an invitation-only Google Spreadsheet with a set of 15 to 40 examples per task.<sup>9</sup> Before beginning their annotation work, the annotators were reminded that:

1. They were asked to dedicate exactly 20 minutes per task (for a total of 1h20min) and should not necessarily aim to complete all the questions provided in the allocated time.
2. At the end of each 20 minute set, the annotators were told to move onto to the next task without delay and asked not to go back to any previous task (even if they had time to spare).
3. They were asked to select only one single answer per question from the set of potential answers, and to not explain their decision process during the training phase.
4. Within one task, they were allowed to attempt the questions in any given order. However, they were asked not to spend more than 5 minutes on a single question. In order to manage their time more efficiently, it was also recommended that they (1) flag difficult questions as they found them, moving immediately to the

next one. In other words, they should first **focus on answering the questions where they felt confident** and only if they had time to spare, (2) go back to the flagged questions and try to solve them. Questions could be **flagged as either “too difficult” or “not clear or ambiguous”**.

5. Finally, they were allowed to consult the annotation guide at any time.

When the training was complete, their work was marked by two of the main authors of this paper and sent back to the annotators who were then asked to review their answers in order to learn from their mistakes.

### C.1.2 Annotation Process

As shown in Table 12, we followed a two-phase iterative approach. Phase 1 included a small batch from the T2, T3 and T4’s *annotation set*. Note that T1 data was excluded due to necessary revisions based on training feedback (see Section B.1). Once completed, explanations underwent the same review process as those used during the annotation training. Our training scheme proved to be effective, resulting in minimal necessary corrections to the annotations. Phase 2 included the remaining instances in the *annotation set*.

Phase	T1	T2	T3	T4
1	0	28	28	28
2	110	82	82	82
Total	110	110	110	110

Table 12: Distribution of task instances across each annotation phase.

Annotators generally adhered to the allocated time frame of 5 minutes per instance, which translated to approximately 7 hours of annotation in Phase 1 and 30 hours in Phase 2. Upon completion, their files were marked and formatted as a JSON file.

### C.1.3 Follow-up Survey

After completing the annotation, we asked the annotators to take a brief follow-up survey. We collected task load data for each of the four tasks using all six NASA-TLX items on a 9-point scale (1-10) (Hart, 1988, 2006). We considered the items individually, as well as their sum, as has been done in prior work (e.g., Quinn and Zhai, 2016; Arnold et al., 2020).

<sup>7</sup>Specifically, the information provided by Jin et al. (2022) in their Appendix D.

<sup>8</sup>Since the guide does not specify a minimum length for the explanations, we made sure annotators wrote complete sentences as opposed to disjointed notes.

<sup>9</sup>The number varied according to the difficulty of each task. For example, the questions in T2 were short but required more specific knowledge while T3 questions contained longer but easier-to-read texts.

Figure 5 shows box-plot representations of the responses from the NASA-TLX surveys, on which we performed Friedman tests (Friedman, 1940) using the `friedmanchisquare` function of the `scipy` Python library (Virtanen et al., 2020). Taking the accepted standard  $\alpha = 0.05$  as the significance threshold (Expósito-Ruiz et al., 2010), we found significant differences for performance ( $\chi^2 = 8.11$ ,  $p\text{-value} = 0.044$ ) only. Note that the performance item in the NASA-TLX survey is framed as follows: “How successful do you think you were in accomplishing the goals of the task set by the experimenter (or yourself)? How satisfied were you with your performance in accomplishing these goals?”. Hence, annotators generally reported a lower sense of accomplishment and satisfaction in T2 and T4, than in T1 and T3.

In the survey, we also included the two open-ended questions to learn more about the annotators’ individual approaches to writing the explanations: specifically, whether they had a particular audience in mind, and what they thought the purpose of the explanations was. We include the exact wording of the questions below:

**Q1:** *The intended recipient of our writing shapes our choice of language and style. Different audiences have different expectations, knowledge levels, and interests. When writing your explanations, did you have a specific audience in mind, or were you writing for a general audience?*

**Q2:** *Explanations can serve a range of purposes: (1) provide an understanding of why a choice was made, (2) justify how that choice was made by providing some evidence, (3) convince others that the choice was correct, and (4) other. When writing your explanations, what were you trying to achieve?*

In response to **Q1**, some annotators reported targeting a “specific” audience, such as researchers or students. On the other hand, one annotator explicitly aimed for a general audience. Others assumed an educated readership with basic linguistic knowledge of English without necessarily being specific about who they might be. Notably, one annotator expressed frustration towards the lack of clarity regarding the intended readership. The diversity in the annotators’ conceptual audiences is very much echoed in the variety of tones used and the level

of depth of the explanations we collected (refer to Table 4 for example).

In response to **Q2**, five out of the six annotators that completed the survey chose (1) as their intended purpose which roughly matches our idea of what a COMMENTARY should do. The remaining annotator sought to justify their choice with evidence (2). While annotators assumed similar strategies, it is interesting to see that they in fact often went well beyond simply providing an understanding of why a choice was made and provided a majority of JUSTIFICATIONS instead (see Figure 2).

## C.2 LLM Annotators

Six different models were used to generate annotations. They were chosen based on coverage of different model sizes, architectures and diversity of sources:

- *Llama-3.1-8B-Instruct*<sup>10</sup> belongs to the family of Llama3.1 models published by Meta AI under the Llama3 community license. It incorporates a context window of 128k length and is pre-trained on a corpus of about 15 trillion tokens.
- *gemma-2-9b-it*,<sup>11</sup> a lightweight open-source model from Google that also supports a 128k length context window. It was trained on 8 trillion tokens of data covering web documents, code, mathematics and more.
- *Mixtral-8x7B-Instruct-v0.1*,<sup>12</sup> a pre-trained generative, sparse, mixture of experts model from Mistral AI. It has a context window of 32k tokens and is pre-trained on data extracted from open web.
- *c4ai-command-r-plus-08-2024*<sup>13</sup> is a 104B parameter multilingual model released from Cohere For AI. It supports a context length of 128K.
- *GPT-4o*,<sup>14</sup> a multimodal model from OpenAI capable of processing and generating text,

<sup>10</sup><https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct>

<sup>11</sup><https://huggingface.co/google/gemma-2-9b-it>

<sup>12</sup><https://huggingface.co/mistralai/Mixtral-8x7B-Instruct-v0.1>

<sup>13</sup><https://huggingface.co/CohereForAI/c4ai-command-r-plus-08-2024>

<sup>14</sup><https://openai.com/index/hello-gpt-4o/>

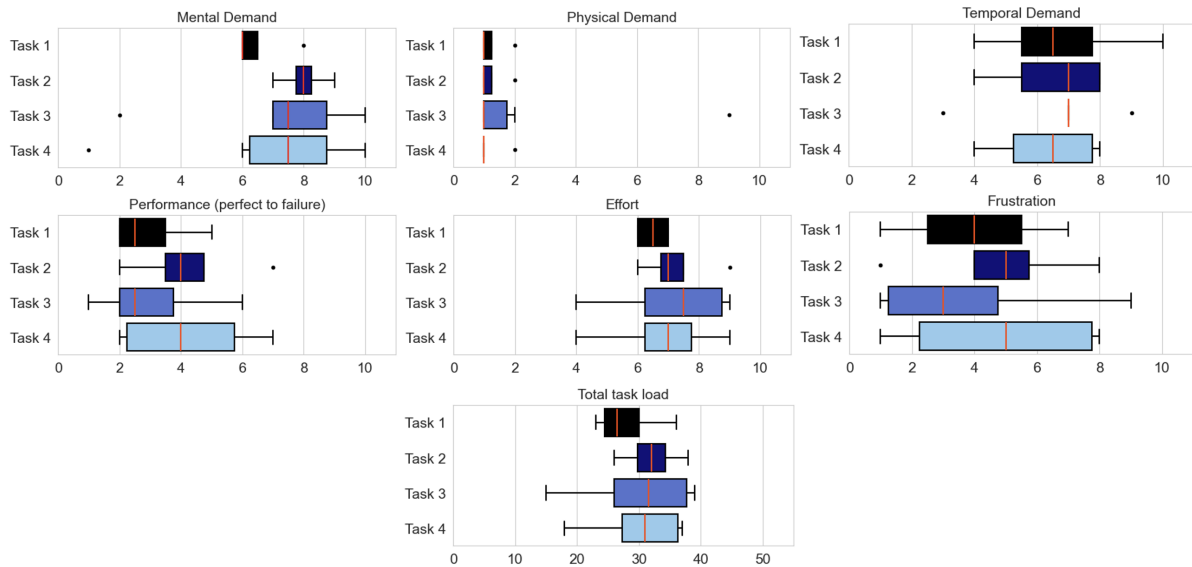


Figure 5: Box-plots of the six NASA-TLX items on a 9 point scale and their sum total. The median is shown in red.

images, and audio. The parameter count of GPT-4o has not been publicly disclosed.

- *Claude 3.5 Sonnet* (*claude-3-5-sonnet-20240620*),<sup>15</sup> an LLM model from Anthropic with improvements in reasoning, language understanding, and coding. The parameter count of Claude 3.5 Sonnet has not been publicly disclosed.

All open-source models were run on NVIDIA A100 GPUs using bf16 precision. We used the latest checkpoints of all open-weight models available at the time of the experiment, along with the default pretrained tokenizers provided for each model. A temperature of 0 was used for all models, including Sonnet 3.5 and GPT-4o, which we accessed via API (for some HuggingFace models, we used 0.01 or set `do_sample=False` due to implementation constraints).

### C.2.1 Prompts for Eliciting Explanations

To elicit explanations from the model, we use a structured prompting approach. Each dataset is associated with a specific prompt designed to guide the model in generating explanations. Additionally, all prompts are preceded by a common system prompt:

You are a helpful, pattern-following assistant. Use the following instructions to respond to user inputs. 1. Start your

answer with a prefix that says "The right answer is: ". 2. Explain the response given in Step 1, with a prefix that says "Because: ". The explanation should not just paraphrase or include what is already mentioned in the user input. 3. Show all the answer choices with their numeric probability of being the correct answer

Below, we present the prompts used for each dataset.

### C.3 HellaSwag Prompt

Each model was given 4 examples to guide its responses. For brevity, these examples are omitted from the prompt shown below.

```
## Examples

Please choose the most plausible ending
(event) for the given context. There is
only one correct answer. After
selecting a correct answer, explain why
you selected that option. The examples
do not include an explanation but you
will need to provide it when answering
the question.

For reference, we provide below four
examples that have already been solved
for you.

{% for example in examples %}
Example {{loop.index}}
{{example}}
{% endfor %}

## Exercise

Context: {ctx_a}
```

<sup>15</sup><https://www.anthropic.com/news/3-5-models-and-computer-use>



Question: Choose the option that best completes the above story.

Options:

```
{% for ending in endings %}
{{ 'ABCD'[loop.index0] }} {{ctx_b}} {{
ending }}
{% endfor %}
```

## C.4 RACE Prompt

We provided 4 examples per query to improve model performance. The prompt format is shown below, excluding the examples for CONCISENESS.

## Examples

In this task, you will be presented with a series of articles. Each is followed by a question which relates to the information provided in the text, and four possible answers. Select only **one** of these options as the correct answer, and explain your choice.

For reference, we provide below four examples that have already been solved for you.

```
{% for example in examples %}
**Example {{loop.index}}**
{{example}}
{% endfor %}
```

# Exercise

Article: {article}

Question: {question}

Options:

```
{% for option in options %}
{{ 'ABCD'[loop.index0] }} {{ option }}
{% endfor %}
```

## C.5 W&I Prompt

Models received 3 examples as part of the prompt structure. The displayed prompt excludes these examples for clarity.

# Task

In this task, you will be presented with a series of essays. Annotate each of these with exactly **one** of three grades: A (beginner), B (intermediate), C (advanced), and then explain your choice.

For reference, we provide below three examples that have already been solved for you.

## Examples

```
{% for example in examples %}
**Example {{loop.index}}**
{{example}}
{% endfor %}
```

## Exercise

Essay: {{full\_text}}

Question: If you were to assign a grade to this essay, what would it be?

Options:

1. Beginner (grade A)
2. Intermediate (grade B)
3. Advanced (grade C)

## C.6 Logic Prompt

Each model was given 7 examples to guide its responses. For brevity, these examples are omitted from the prompt shown below.

## Examples

Please identify the type of logical fallacy. There is only **one** correct answer. After selecting a correct answer, explain why you selected that option.

For reference, we provide below seven examples that have already been solved for you.

```
{% for example in examples %}
**Example {{loop.index}}**
{{example}}
{% endfor %}
```

## Exercise

Statement: {source\_article}

Question: Which type of logical fallacy is this an example of?

Options:

- A. Faulty generalisation
- B. False causality
- C. Circular claim
- D. Appeal to emotion
- E. Deductive fallacy
- F. False dilemma
- G. Fallacy of credibility

## D Custom Agreement Metric

**First metric.** Cohen's  $\kappa$  (Cohen, 1960) and Krippendorff's  $\alpha$  (Krippendorff, 2011) are among the most frequently used inter-rater reliability metrics. However, their direct application is best suited to nominal or categorical data. Even with adaptations like weighted kappa, these coefficients struggle to capture the full inter-relationship of hierarchical

nested data. To bridge this gap, we introduced a custom metric that specifically accounts for the nested dependencies in CUBE. Our custom metric accounts for the *superlabels* (NONE, COMMENTARY, JUSTIFICATION, ARGUMENT) and *sublabels* (i.e., all DIMENSIONS) in Rubrik. In both cases, the metric penalises discrepancies between ratings, with the penalty proportional to the difference in the hierarchical level. For example, consider the cases shown in Table 13 and Table 14.

Case	Rater 1	Rater 2	Diff.	Agree. (%)
1	COMMENTARY	JUSTIFICATION	1	67
2	COMMENTARY	ARGUMENT	2	50
3	NONE	ARGUMENT	3 to 4	0 to 25

Table 13: Superlabel agreement. NONE denotes the case where either of the COMMENTARY’s COMPONENTS are missing, namely Action (1.a) and Reason (1.b).

From the *superlabel* point of view, there is a partial agreement in Case 1 since a JUSTIFICATION has the two components (ACTION and REASON) of a COMMENTARY and an additional one (namely, EVIDENCE). Thus, the difference in the raters’ judgement is 1. From the *sublabel* point of view, the agreement range is higher as it takes into consideration all the elements of a COMMENTARY (8: 2 COMPONENTS, 6 DIMENSIONS) and a JUSTIFICATION (10: 3 COMPONENTS, 7 DIMENSIONS).

Case	Rater 1	Rater 2	Diff.	Agree. (%)
1	COMMENTARY	JUSTIFICATION	1-8 of 10	90-20
2	COMMENTARY	ARGUMENT	4-10 of 12	66-17
3	NONE	ARGUMENT	11-12 of 12	8-0

Table 14: Sublabel agreement. The difference (Diff.) column shows a range, taking both COMPONENTS and DIMENSIONS into consideration.

As explained in Section 3.3, a good COMMENTARY is the base of a good JUSTIFICATION. This means that Rater 2 judged with met (✓) all the elements of a COMMENTARY. The disagreement with Rater 1 comes from them judging with not met (✗) one or more of the six dimensions. The same logic applies to Cases 2 and 3.

**Second metric.** The first agreement metric accounts for partial agreement between LLMs and human annotators. We tested all LLMs as evaluators on the same subset judged by humans. However, we observe that LLMs often rate an explanation as JUSTIFICATION over the other options, compromising their ability to detect other types (see Table 16).

This highlighted the need for an additional custom metric, which we designed based on a weighted F1 score to penalise over-centralization on a single label. The class weights are derived from both human evaluations and LLM evaluations from all six models. In our approach, we first calculate the distribution percentage of each superlabel in human evaluation  $p_i^{human}$  for label  $i$ . We then calculate the average distribution percentage of each superlabel across all 6 LLM evaluations denoted as  $p_i^{LLM}$ . These two percentages are combined as the class weight:

$$w_i = \lambda p_i^{human} + (1 - \lambda) p_i^{LLM}$$

where  $\lambda$  is a hyperparameter representing the relative importance of human evaluations vs. LLM evaluations. The derived class weights are then incorporated into the calculation of the weighted F1 score.

As shown in Table 15, our first metric points to Command R+ as the model with higher agreement with human evaluators. However, a closer look at the distribution of the explanation types assigned show that the high agreement is due to identifying an explanation as JUSTIFICATION nearly always. Our second metric penalises this behaviour, ranking Command R+ as the least effective evaluator.

## E Rubric Evaluation Prompts

To evaluate explanations generated by the model, we use a structured prompting approach based on a rubric. Each dataset is associated with a specific prompt designed to guide the model in assessing explanations. Below is the prompt template that encodes the evaluation rubric.

Note that the prompt **does not** ask the model to judge whether an explanation is 😊 *good* or 😞 *bad*. This choice reflects the insights of Panickssery et al. (2024), who found that out-of-the-box LLMs, such as GPT-4 and Llama 2, have non-trivial (over 50%) accuracy at distinguishing themselves from other LLMs and humans. As a result, these models tend to recognise and favour their own generations. Thus, our prompt only specifies the evaluation criteria to decide whether a given COMPONENT or DIMENSION is met (✓) or not met (✗). This approach successfully mitigated self-preference; GPT-4o, our third evaluator, judged its own outputs as *bad* at a comparable low rate to other models’ outputs. Recall from Section 3.3 that an explanation is deemed good if, and only if, it meets all the criteria. While

Task	Agreement	Humans	Open models				Closed Models	
			Llama 3.1	Gemma 2	Command R+	Mixtral	GPT-4o	Sonnet 3.5
T1	Superlabel	0.814	0.693	0.799	0.797	<b>0.812</b>	0.794	<b>0.800</b>
	Sublabel	0.823	0.706	0.795	0.826	<b>0.829</b>	0.807	<b>0.811</b>
T2	Superlabel	0.910	0.832	0.862	<b>0.873</b>	0.869	0.878	<b>0.879</b>
	Sublabel	0.923	0.865	0.888	<b>0.903</b>	0.898	<b>0.902</b>	0.899
T3	Superlabel	0.830	0.830	0.838	0.843	<b>0.847</b>	0.844	<b>0.854</b>
	Sublabel	0.869	0.862	0.866	0.881	<b>0.887</b>	0.872	<b>0.881</b>
T4	Superlabel	0.887	0.797	<b>0.817</b>	0.810	0.774	<b>0.846</b>	0.833
	Sublabel	0.897	0.807	0.804	<b>0.853</b>	0.787	<b>0.860</b>	0.851
Overall	Superlabel	0.860	0.788	0.829	<b>0.831</b>	0.825	0.841	<b>0.842</b>
	Sublabel	0.878	0.810	0.838	<b>0.866</b>	0.850	<b>0.860</b>	<b>0.860</b>

Table 15: Overview of agreements scores, calculated with the first metric. In bold, the highest score by superlabel and sublabel, comparing the performance of open- vs. closed-source models.

Annotator	NONE	COMMENTARY	JUSTIFICATION	ARGUMENT	Second-metric-score	Second-metric-rank
Human_annotator 1	0	293	406	221	-	-
Human_annotator 2	5	264	229	422	-	-
LLama 3.1	87	47	450	336	0.405	5
Gemma 2	<b>9</b>	<b>222</b>	<b>561</b>	<b>128</b>	<b>0.464</b>	<b>2</b>
Command R+	4	20	894	2	0.346	6
Mixtral	5	240	654	21	0.427	4
GPT-4o	<b>14</b>	<b>107</b>	<b>685</b>	<b>114</b>	<b>0.476</b>	<b>1</b>
Sonnet 3.5	5	126	742	47	0.444	3

Table 16: Aggregated label counts for each annotator and metric score. In bold are the results from the two best-ranked LLM evaluators. In both cases, there is a better balance in the judgement of explanation types.

this condition establishes a rigorous baseline for assessment, this scoring strategy is flexible. The specific conditions can be adjusted to fit varying research objectives or contextual needs.

```
{# Base template for rubric scoring #}  
# Explanation Judging Task
```

Your task is to evaluate a set of explanations in a given context. We define the context (\*\*Task\*\*, \*\*Audience\*\*, and \*\*Purpose\*\*) in the following way:

**\*\*Task\*\***: you will be shown a series of multiple-choice questions relating to one of four tasks (commonsense reasoning, fallacy detection, reading comprehension and essay scoring) in the following format:

1. **\*\*Question\*\***: The question being answered.
2. **\*\*Answer Choices\*\***: The possible answer choices for that question.
3. **\*\*Correct Answer\*\***: The correct answer to the question.
4. **\*\*User Answer\*\***: The answer provided by the user.
5. **\*\*Explanation\*\***: The explanation provided by the user to support their answer.

**\*\*Audience\*\***: you should assume that the audience of the explanations is adult, English-proficient, and provided in a formal academic setting.

**\*\*Purpose\*\***: the explanations should provide an understanding of why a certain answer was chosen for a given multiple-choice question.

---

## ## Evaluation Criteria

For the given explanation, please answer the following questions with either **\*\*Yes\*\*** or **\*\*No\*\***. Note that you **\*\*should not consider the correctness of the user's answer\*\*** when evaluating the explanation. Focus solely on the quality of the explanation according to the criteria provided.

1. **\*\*Action\*\***: Does the explanation clearly indicate the decision or choice being made (e.g., specifying the selected answer)?
  - Answer **\*\*Yes\*\*** if it does. For example "The correct answer is A."
  - Answer **\*\*No\*\*** if it does not. For example "Because it is the final part of the sequence."
2. **\*\*Reason\*\***: Does the explanation provide reasoning or insight into why the decision or choice was made, explaining the underlying logic or rationale for the **\*\*Action\*\***?
  - Answer **\*\*Yes\*\*** if it does. For example "The right answer is C because it is the final part of the sequence."
  - Answer **\*\*No\*\*** if it does not. For example "The correct answer is A."
3. **\*\*Grammaticality\*\***: Is the explanation grammatically correct and



free of lexical or syntax errors? Small typos are acceptable, but the errors should not impede comprehension in any way.

- Answer **\*\*Yes\*\*** if it is. For example "The correct answer is A because nowadays our society is based on consumerism and the way in which we are producing is contaminating the world."

- Answer **\*\*No\*\*** if it is not. For example "The correct answer is A because now a day our socity it is bassed in consumer, so that become the word more contaminate to produce the products that we demanding."

4. **\*\*Word Choice\*\***: Is the language used in the explanation tailored to the given context (task, audience, purpose)? And are the sentences in the explanation well-formed?

- Answer **\*\*Yes\*\*** if they are. For example "The correct answer is A because the essay lacks fluency. There are many incorrect clauses and missing words. And while the overall meaning can be deduced, the essay does not demonstrate an accurate grasp of language (e.g., frequent spelling and punctuation errors)."

- Answer **\*\*No\*\*** if they are not. For example "Answer A. lack of fluency, incorrect clauses and missing words, meaning can be found but does not demonstrate an accurate grasp of language"

5. **\*\*Cohesion\*\***: Does the explanation make appropriate use of transition phrases (e.g., connectives like "because", "therefore", "consequently", overlapping words across sentences, etc.)?

- Answer **\*\*Yes\*\*** if it does. For example "The correct answer is C because the man is on roller blades, not on a skateboard. Further, he is not talking to anyone and therefore cannot possibly 'continue speaking.'"

- Answer **\*\*No\*\*** if it does not. For example "The correct answer is C, because the man is on roller blades, not a skateboard, and is not talking to anyone in the example so cannot 'continue speaking'".

6. **\*\*Conciseness\*\***: Is the explanation free of any redundant, irrelevant, or excess sentences (that is, not required to understand the answer)?

- Answer **\*\*Yes\*\*** if it is. For example "The correct answer is D because it accurately reflects the sequence of events."

- Answer **\*\*No\*\*** if it is not. For example, given that the option D was "next she explains how to use the lawnmower and other tools and then she cuts the grass", the following explanation is not concise: "The correct answer is D because the sentence mentions that she explains how to use the lawnmower and other tools, and then she cuts the grass. Option D accurately reflects the sequence of events."

7. **\*\*Appropriateness\*\***: Is the explanation culturally appropriate, matching expectations for the given context?

- Answer **\*\*Yes\*\*** if it is. For example "The right answer is B because the tenses are properly used and the story makes sense."

- Answer **\*\*No\*\*** if it is not. For example "The right answer is B

because the tenses are properly used and (within the slightly odd context) the story makes sense."

8. **Coherence**: Does the explanation appropriately transition between ideas? That is, does the explanation make sense as a whole (e.g., good context-relatedness, semantic consistency, and inter-sentence causal and temporal dependencies, etc.)?

- Answer **Yes** if it does. For example "The correct answer is D, because no information about Liu's relationship to science subjects specifically is given in the passage, therefore the fact that they like chemistry is implied and ambiguous."

- Answer **No** if it does not. For example "The correct answer is D, because no information about Liu's relationship to science subjects specifically is given in the passage, therefore the fact that they like cheese is implied and ambiguous."

9. **Evidence**: Does the explanation provide concrete evidence (can be both explicit or implicit) that supports the reasoning, such as information from the question's context or general knowledge?

- Answer **Yes** if it does. For example "The right answer is C, because it finishes the sequence, describing the effect of bowling the ball and what happens as a result."

- Answer **No** if it does not. For example "The right answer is C, because it is the final part of the sequence."

10. **Plausibility (of the evidence)**: Is the provided evidence plausible and consistent with human reasoning, considering the context and general world knowledge?

- Answer **Yes** if it is. For example "The correct answer is A ('Jack picks the cheese') because we are told that he enjoys eating 'mozzarella' in the morning."

- Answer **No** if it is not. For example "The correct answer is A ('Jack picks the cheese') because my name is also Jack and I personally love cheese for breakfast."

11. **Affective Appeals**: Does the explanation use vivid, or emotionally charged language (e.g., metaphors) to evoke feelings in the audience?

- Answer **Yes** if it does. For example "The expression in the final section is very heartfelt; the tone is excitable and keen throughout."

- Answer **No** if it does not. For example "The final section reflects the writer's strong feelings on this issue."

12. **Qualifiers**: Does the explanation make use of hedges, boosters, attitude markers, self-mentions, or engagement markers to clarify the writer's stance (i.e., the explainer's personal feelings towards the task)? Note that the stance can be implicit unlike the **Action**.

- Answer **Yes** if it does. For example "The right answer is B, because the text is keeping with what is presumably a tour guide's voice: intentionally using clunky and overly expressive words."

- Answer **\*\*No\*\*** **if** it does not. For example "The right answer is B, because the text is keeping with the original tour guide's voice."

13. **\*\*Stance Clarity\*\***: Is the explainer's stance (their personal feelings towards the task) clearly and unambiguously conveyed through affective appeals or qualifiers? Note that the stance can be implicit unlike the Action.

- Answer **\*\*Yes\*\*** if it is. For example "The correct answer is A (beginner) because this text is undeniably of a low English level."

- Answer **\*\*No\*\*** if it is not. For example "The correct answer is A (beginner) because this text is clearly of a low English level although the final section is incredibly well written."

---

## ## Expected Output

Your answers should be formatted as follows:

1. Action: **\*\*Yes\*\*** or **\*\*No\*\***
2. Reason: **\*\*Yes\*\*** or **\*\*No\*\***
3. Grammaticality: **\*\*Yes\*\*** or **\*\*No\*\***
4. Word Choice: **\*\*Yes\*\*** or **\*\*No\*\***
5. Cohesion: **\*\*Yes\*\*** or **\*\*No\*\***
6. Conciseness: **\*\*Yes\*\*** or **\*\*No\*\***
7. Appropriateness: **\*\*Yes\*\*** or **\*\*No\*\***
8. Coherence: **\*\*Yes\*\*** or **\*\*No\*\***
9. Evidence: **\*\*Yes\*\*** or **\*\*No\*\***
10. Plausibility: **\*\*Yes\*\*** or **\*\*No\*\***
11. Affective Appeals: **\*\*Yes\*\*** or **\*\*No\*\***
12. Qualifiers: **\*\*Yes\*\*** or **\*\*No\*\***
13. Stance Clarity: **\*\*Yes\*\*** or **\*\*No\*\***

---

## ## Question

{% block question -%}

{{ task\_question }}

{%- endblock -%}

## ## Answer Choices

{% block choices %}

{% for choice in choices %}

{{ 'ABCDEFGH'[loop.index0] }} {{ choice }}

{% endfor %}

```
{% endblock %}

## Correct Answer
{{correct_answer}}

## User Answer
{{user_answer}}

## Explanation
{{explanation}}
```



## Dataset-Specific Evaluation Prompts

In the above template, the main difference between datasets is the format of the question and the options. Below, we show how each dataset-specific question and option block is customised.

### E.1 HellaSwag

```
{% extends "rubric_prompt" %}

{% block question -%}

{{ ctx_a }}

{%- endblock %}

{% block choices %}

{% for ending in endings %}

{{ 'ABCD'[loop.index0] }}) {{
ctx_b}} {{ ending }}

{% endfor %}

{% endblock %}
```

### E.2 RACE

```
{% extends "rubric_prompt" %}

{% block question -%}

Article: {text}

Question: {question}

{%- endblock %}
```

### E.3 WANDI

```
{% extends "rubric_prompt" %}

{% block question %}
Essay: {text}
{% endblock %}

{% block choices -%}

1. Beginner (grade A)
2. Intermediate (grade B)
3. Advanced (grade C)
```

```
{%- endblock %}
```

### E.4 Logic

```
{% extends "rubric_prompt" %}

{% block question %}

Statement: {{text}}

Question: {{question}}

{% endblock %}

{%- block choices -%}

A. Faulty generalisation
B. False causality
C. Circular claim
D. Appeal to emotion
E. Deductive fallacy
F. False dilemma
G. Fallacy of credibility

{%- endblock -%}
```

## F Detailed Analysis Results

This section delves deeper into the data, offering additional insights to complement the summary provided in Section 5.

### F.1 Answer Frequencies

First, we report the frequencies of the answer choices picked by different groups of annotators during the annotation phase, and compare these to the actual distribution of correct answers in each task on the *annotation set* in Figure 6. Recall that we explicitly tried to get as uniform a distribution across the different answer choices as possible in the *annotation set* (as described in Appendix B).

Overall, we note that while human annotators sometimes refused to choose an answer between those provided (NONE), the LLMs almost never refused to answer. This may be because LLMs have a tendency to overestimate their ability to answer questions (Zhang et al., 2023b).

In T1 and T3, the answer frequencies of all annotators seem fairly balanced, with the only notable difference being that human annotators also responded NONE. In T2, however, we can see that

the grouped Open LLMs (Command R+, Mixtral, Llama 3.1 and Gemma 2) seem to significantly favour answers A, B and D at the expense of answers C and G, while the other groups of annotators remain relatively close to the actual frequency distribution. We should note that despite the fact that the *annotation set* is more or less balanced, in Jin et al. (2022) authors state that more than a single fallacy type may apply to a single instance. This may explain the variation observed. Specifically, they identified “common among incorrect but reasonable predictions” in their task, which “are debatable cases where multiple logical fallacy types seem to apply”.

In T4, we notice a stark difference between humans and LLMs annotators. On one hand, LLMs almost never assign C (advanced) scores to essays, and overwhelmingly assign B (intermediate) scores around 65% of the time. While human annotators use the whole range of the scale, though still showing signs of a strong central tendency or severity by only assigning around half the actual proportion of advanced scores. Interestingly, experts annotators, that are professionally trained to assess the work of language learners, did not distinguish themselves from the contractors we hired who had very similar frequency distributions in the two language tasks. Overall, evaluators failed to identify advanced essays, focusing most of their attention on the middle of the rating scale. Essay scoring is a notoriously complex and subjective task (Brown, 2010), and we intentionally did not provide any scoring rubric to the annotators. They thus lacked a proper point of reference for the scale, which seems to be the source of the frustration reported by one annotator (see Section C.1.3).

## F.2 Detailed Accuracy

Next, in Figure 7 we report the performance or accuracy (%) of the individual annotators and their groups, in each of the tasks, as well as their overall average performance across the four tasks.

Looking at the average performance across the four tasks, closed LLMs seem to perform the best, while open LLMs perform the worst, with humans (contractors and experts) performing just slightly better than the open models. The two closed models exhibited comparable average performance across the four tasks, but Sonnet 3.5 is more consistently good across the four tasks, whereas GPT-4o is very good at Reading Comprehension (T3) and less good at Essay Scoring (T4).

Overall, these graphs make it apparent that Essay Scoring (T4) was the hardest with an average accuracy of roughly 52% (across all annotators), while Reading Comprehension (T3) was by far the easiest with an average accuracy reaching almost 84%.

As in the previous section, we note that humans were overall quite consistent. The experts were ever so slightly better at Essay Scoring (T4) than the contractors, but this difference is very small. We had expected them to do much better due to being professionally trained to perform language assessment tasks. Further, while this background should have directly impacted their capacity to do well in T4, we also expected them to do better than the contractors in T3 given the language-related nature of their day-to-day work. However, contractors were in fact ever so slightly better at Reading Comprehension (T3). These findings suggest that we do not always necessarily need to hire professionals, and that professional expertise can be matched by a rigorous selection process and sufficient training of annotators.



Figure 6: Frequencies of the answers picked by the different groups of annotators during the annotation phase. We also show the **Actual** distribution of correct answers in black in the *annotation set*.

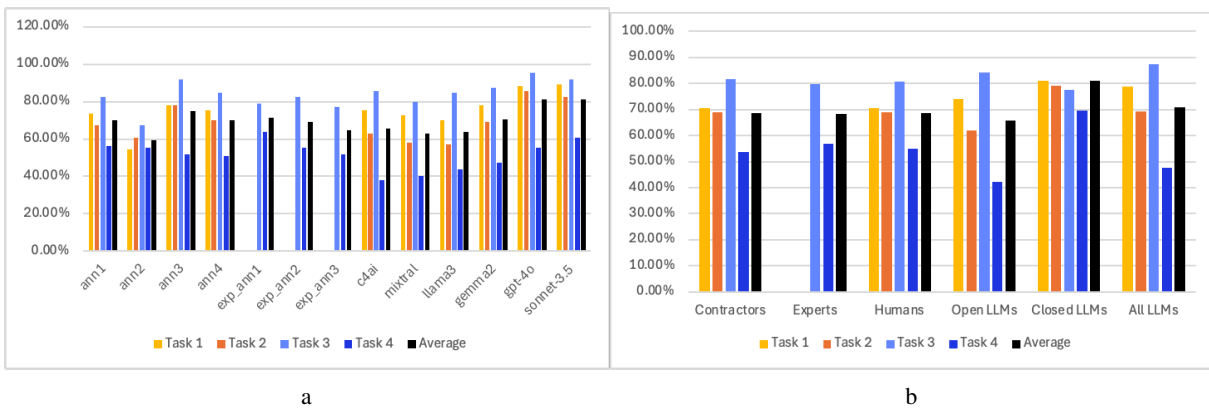


Figure 7: Accuracy results of the different annotators in each of the tasks. On the left, 7a shows the individual annotator performance, and on the left, 7b shows the performance by group of annotators. We also include the **Average** accuracy across the four tasks of each annotator or group in black.

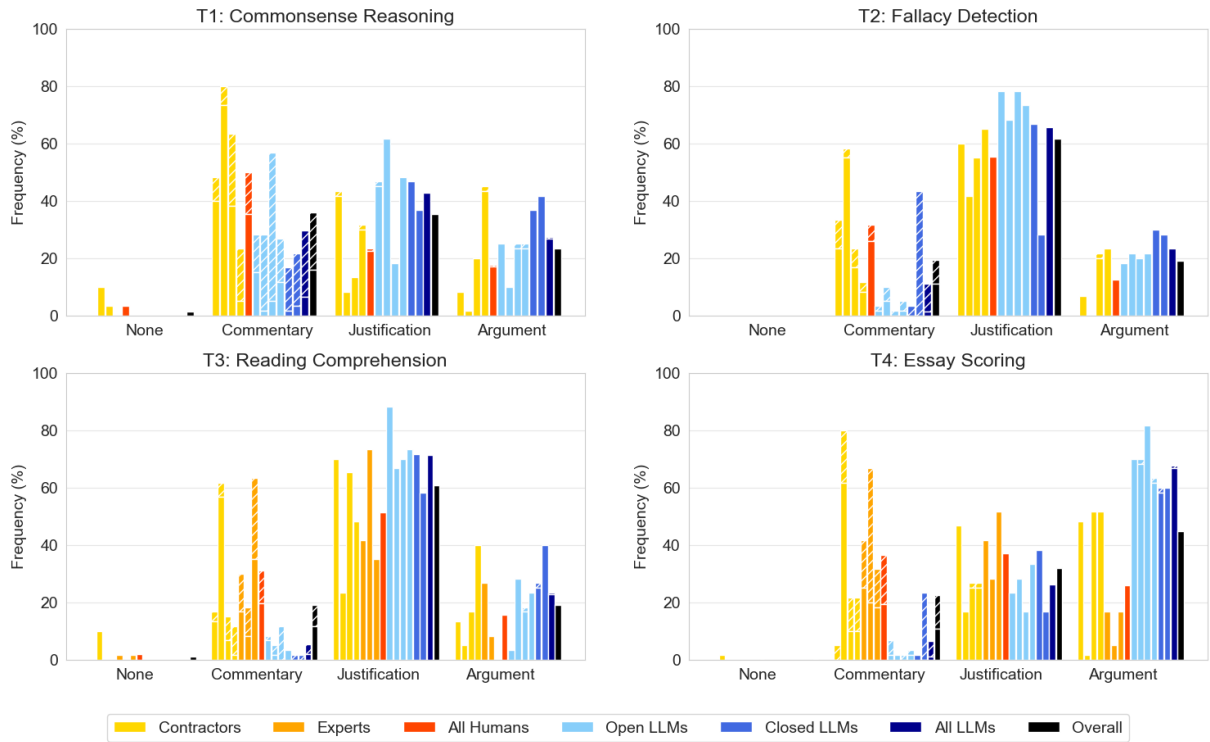


Figure 8: A breakdown of the bar plots in Figure 2 which shows the frequencies (%) of the different explanation types for each individual human and LLM annotator (4 contractors, 3 experts in **T3** and **T4**, 4 open-models—Command R+, Mixtral, LLama 3.1, Gemma 2—and 2 closed-models—GPT-4o and Sonnet 3.5—in this order) in the *evaluation set*. We also include the average frequencies across all annotators (in black). We average the frequencies across all three evaluators (two humans and GPT-4o).

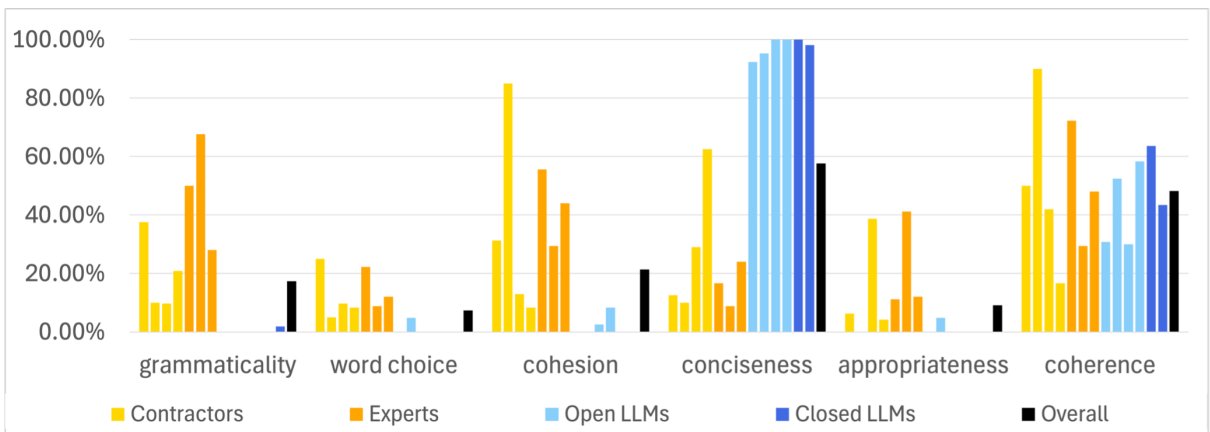


Figure 9: A breakdown of Figure 3 which shows the sources of the *bad* COMMENTARIES for each individual human and LLM annotator (4 contractors, 3 experts in **T3** and **T4**, 4 open-models—Command R+, Mixtral, LLama 3.1, Gemma 2—and 2 closed-models—GPT-4o and Sonnet 3.5—in this order) in the *evaluation set*. We also include the average frequencies across all annotators (in black). We average the frequencies across all three evaluators (two humans and GPT-4o).