

Exposía: Academic Writing Assessment of Exposés and Peer Feedback

Anonymous ACL submission

Abstract

We present Exposía, the first public dataset¹ that connects writing and feedback assessment in higher education, enabling research on educationally grounded approaches to academic writing evaluation. Exposía includes student research project proposals and peer and instructor feedback consisting of comments and free-text reviews. The dataset was collected in the “Introduction to Scientific Work” course of the Computer Science undergraduate program that focuses on teaching *academic writing skills* and *providing peer feedback on academic writing*. Exposía reflects the multi-stage nature of the academic writing process that includes drafting, providing and receiving feedback, and revising the writing based on the feedback received. Both the project proposals and peer feedback are accompanied by human assessment scores based on a fine-grained, pedagogically-grounded schema for writing and feedback assessment that we develop.

We use Exposía to benchmark state-of-the-art open-source large language models (LLMs) for two tasks: automated scoring of (1) the proposals and (2) the student reviews. The strongest LLMs attain high agreement on scoring aspects that require little domain knowledge but degrade on dimensions evaluating content, in line with human agreement values. We find that LLMs align better with the human instructors giving high scores. Finally, we establish that a prompting strategy that scores multiple aspects of the writing together is the most effective, an important finding for classroom deployment.

1 Introduction

Writing a research project proposal² is the number one challenge in academic education and represents a pressing concern across the academic community (Muneer et al., 2020). It is also an essential skill required for conducting research and for

¹The dataset will be released under CC BY-NC 4.0 license.

²We will use *research project proposal*, *research project exposé*, or *exposé* interchangeably in the paper.

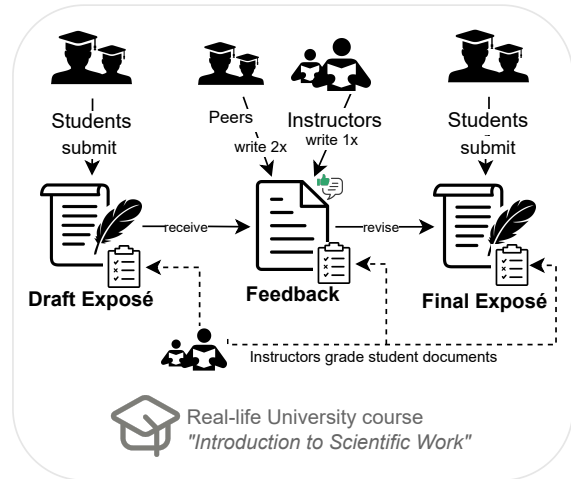


Figure 1: Overview of Exposía. In the university course “Introduction to Scientific Work”, students submit a *draft exposé*, receive *feedback* in the form of *comments* and a *free-text review*, revise and submit a *final exposé*. Feedback is produced by peers and instructors. Instructors grade draft and final exposés and student feedback.

achieving academic success at all levels (Kivunja, 2016; Grech, 2018). However, a lot of students lack the necessary competencies to develop a research proposal and its core components, including identifying a research question, formulating appropriate research methodology, and defining hypotheses and methods to solve the problem (Rastri et al., 2023).

At the same time, higher education and research in general face another pressing challenge: the peer review crisis (Lu et al., 2025). High-quality feedback is an essential pedagogical tool that has been shown to vastly improve learning outcomes of both authors and reviewers (Wu and Schunn, 2021). Yet, students are rarely taught this skill in a classroom setting. Together, these challenges highlight the need for instructional approaches that simultaneously support research project proposal writing and developing effective peer review skills.

Teaching students how to write an effective project exposé and how to provide high-quality feedback is crucial in university education for sev-

063 eral reasons. First, assessing both project exposés
064 and peer feedback requires substantial human ef-
065 fort and does not scale well (Grabau and Wilson,
066 1995). Second, developing automated assessment
067 tools requires modeling complex cross-document
068 relationships, connecting student writing, feedback,
069 and revisions. Moreover, peer feedback is often of
070 low quality and highly inconsistent (Purkayastha
071 et al., 2025). Finally, collecting data in an educa-
072 tional setting is extremely challenging (Section 3).
073 As a result, existing resources for computational
074 writing assessment are fragmented across subtasks
075 such as automated essay scoring (Ke and Ng, 2019),
076 peer feedback analysis (Staudinger et al., 2024), or
077 revision modeling (Ruan et al., 2024).

078 In this work, we present Exposía, the first corpus
079 that integrates drafting, feedback, revision, and fine-
080 grained exposé and feedback assessment scores
081 within a single dataset (Figure 1). Exposía is *con-*
082 *tinuously* collected in a newly developed univer-
083 sity Computer Science course. The course focuses
084 on teaching students (1) *how to write a research*
085 *project exposé* and (2) *how to write peer feedback*.
086 While this is an undergraduate course, the skills
087 it develops are equally applicable to graduate stu-
088 dents and other early-career researchers.

089 **Exposía** captures a complete research-proposal
090 writing workflow. As illustrated in Figure 1, each
091 student first produces an exposé **draft**. The draft
092 receives (i) **peer feedback** from fellow students
093 and (ii) **instructor feedback**, both provided as (a)
094 **inline comments** anchored to draft passages and
095 (b) a **free-text review**. The student revises the draft
096 based on the feedback to produce the **final** exposé
097 version. Figure 2 (top) details each component.

098 We further design fine-grained grading criteria
099 that target the key aspects of exposé and feed-
100 back writing (Figure 2, bottom). These criteria are
101 grounded in established research on exposé writ-
102 ing and feedback assessment and are tailored to
103 the academic disciplinary context of Exposía. The
104 exposé drafts, the final exposés, and the student
105 feedback are graded by course instructors, who
106 provide scores on each of the criteria.

107 To the best of our knowledge, Exposía is the first
108 dataset, collected in an educational setting, that
109 features this unique class of cross-document rela-
110 tions, enabling new applications for writing and
111 grading assistance. Because the corpus contains
112 paired drafts and final exposés, it enables quantify-
113 ing draft-to-final improvements on specific criteria
114 and relating these changes to the feedback received.
115 Overall, we find that exposés improve in response

116 to feedback, particularly on criteria related to the in-
117 domain content of the writing rather than form. We
118 further establish baseline models with open-source
119 LLMs for two tasks: **exposé scoring** and **feed-**
120 **back scoring**, based on the fine-grained grading
121 schema for writing and feedback assessment.

122 **Contributions** The main contribution of this
123 work is the introduction of a novel problem def-
124 inition of exposé and review scoring in the educa-
125 tional context. The result is the first public dataset
126 of student research exposés and feedback, accom-
127 panied by human assessment scores on exposés
128 and feedback, based on the fine-grained grading
129 schema that we develop. We present in-depth analy-
130 ses of the dataset, and models for the corresponding
131 tasks of exposé and feedback scoring, establishing
132 baselines that the community can build on.

133 2 Related Work

134 **Automated Essay Scoring** The task of Auto-
135 mated Essay Scoring (AES) (Page, 1966), aims
136 to reduce human grading effort and improve consis-
137 tency in writing assessment. Most widely used
138 AES benchmarks consist of essays written to
139 standardized prompts in general-domain language.
140 These essays are typically manually annotated with
141 a holistic score, complemented by scores for a
142 small set of assessment criteria (e.g., content, or-
143 ganization, language) (Mathias and Bhattacharyya,
144 2018; Gaudeau, 2025; Li and Ng, 2024; Xue et al.,
145 2021; Crossley et al., 2023; Ishikawa, 2023). Im-
146 portantly, these benchmarks do not include human
147 feedback, revisions, or feedback assessment scores.
148 **Feedback datasets.** A small number of resources
149 connect student writing with feedback and/or revi-
150 sion. *ArgRewrite* links multiple drafts of student
151 argumentative essays, annotates revision purposes,
152 and also includes instructor feedback (Zhang et al.,
153 2017; Kashefi et al., 2022). Pilan et al. (2020) re-
154 lease a teacher-feedback dataset pairing original
155 and revised student sentences with comments and
156 tags for a targeted error type, enabling analyses of
157 how feedback relates to revision outcomes. For
158 peer feedback, Rietsche et al. (2022) collect peer
159 reviews written by Master’s students and annotate
160 for specificity and perceived helpfulness. However,
161 existing datasets lack (i) passage-anchored inline
162 comments tied to draft spans, (ii) role-separated
163 *peer and instructor* feedback, and (iii) criterion-
164 level assessment scores on the text to be assessed
165 and the feedback. Exposía addresses these gaps by
166 capturing a full draft→feedback→revision work-

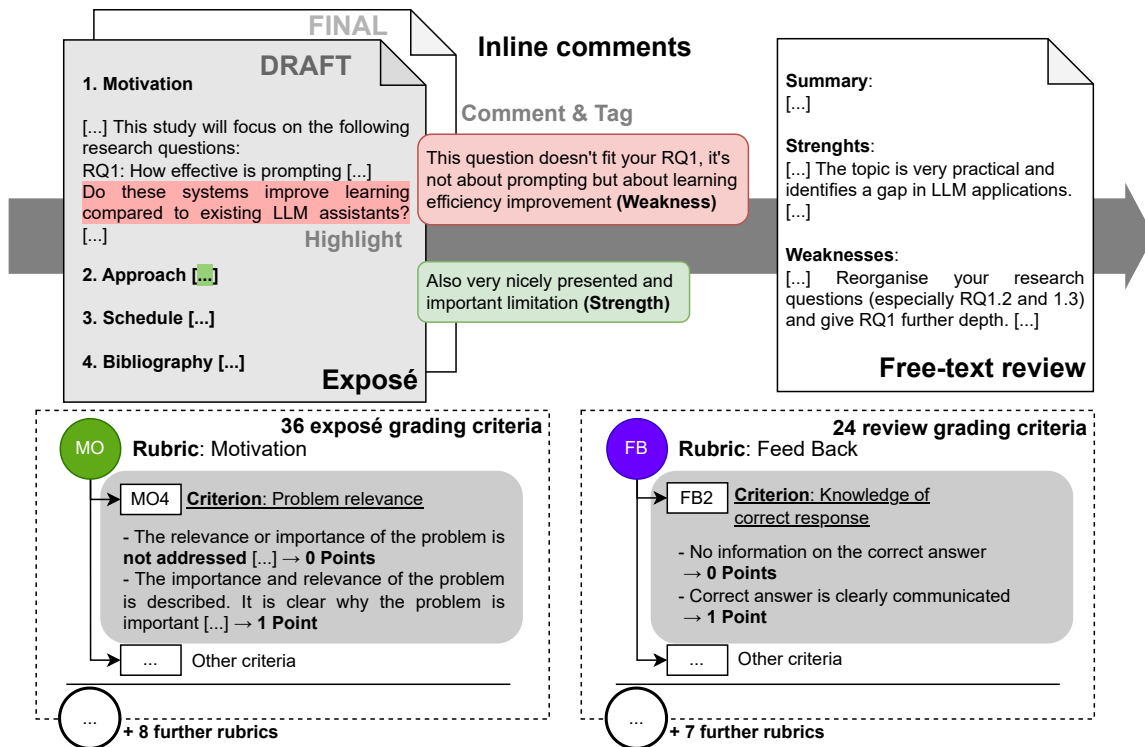


Figure 2: Example instance from Exposía. Top: A student **draft exposé** (top-left) is reviewed by an instructor or a student peer and receives **feedback**. The feedback consists of (1) passage-anchored **inline comments** (shown in the top middle in red and green), each associated with a **tag** (e.g., Weakness, Strength); and (2) a free-text **review** (top-right). Bottom: Exposés and reviews are graded by instructors and receive scores on multiple aspects (**criteria**) organized into high-level topics (**rubrics**). The **exposé grading criteria** and the **review grading criteria** are shown in the bottom left and bottom right of the figure, respectively.

flow with peer *and* instructor feedback.

LLM-based scoring and transfer to peer review

Recent work evaluates LLMs as scorers and emphasizes that educational use requires evidence beyond accuracy (Pack et al., 2024; Huang and Wilson, 2025). Exposía complements this line of work, by enabling fine-grained exposé and feedback assessment with LLMs and interpreting LLM-human alignment relative to human agreement (Section 6). Because our corpus links each exposé draft with feedback text and assessment scores, it is related to scholarly peer-review datasets (Kang et al., 2018; Dycke et al., 2023). However, these datasets provide only coarse-grained annotations. Importantly, they serve fundamentally different purposes, as they focus on scientific communication in a professional decision-making context, whereas Exposía is situated in an educational environment with a focus on learning and assessment. This creates a testbed for transfer experiments between educational feedback assessment and scholarly reviewing.

3 Course Setting and Study Protocol

Data collection protocol The dataset was collected in the newly-introduced Bachelor-level Com-

puter Science course *Introduction to Scientific Work* offered for the first time in winter term 2024. The corpus is the result of a targeted data collection.³ Collecting data in an educational setting is challenging due to the high sensitivity of the data. In particular, the collection, storage, and sharing of educational data often raise serious privacy and ethical issues (Guan et al., 2023).

Course structure The students wrote an exposé on one of three topics, covering various areas of Computer Science and Artificial Intelligence (AI):⁴ *Holographic AI Assistant* (Human-Computer Interaction), *Interactive LLM-based Teaching Agents* (Natural Language Processing), or *Generative AI in Healthcare* (Data Mining).

The course was broken down into three parts.⁵ In the beginning, alongside the lectures, the students received guidelines for writing an exposé (Appendix B) and the exposé grading criteria (Section 4.2). Each student produced an exposé draft. This was followed by a peer-review phase, where each student was assigned to write double-blind

³The institutional IRB has approved the related proposal.

⁴Topics were evenly assigned based on preference voting.

⁵The syllabus is provided in Appendix A.

reviews for two exposés written by their peers, selected uniformly at random. The students were provided with review grading criteria (Section 4.3). The instructional staff also provided a review for each exposé. This was followed by a revision phase to produce the final exposé, based on the feedback. Both the drafts and final versions, and the student reviews were graded by the instructional staff. The final grade was calculated as follows: draft exposé 30%, peer reviews 30%, and final exposé 40%.⁶ Reviewing was facilitated through [hidden].⁷

Continued large-scale data collection 193 students enrolled in the course,⁸ of whom 159 completed all requirements. The Exposía dataset refers to the *open data* that participants agreed to contribute for public release. All 12 instructors and about 30% of students consented to share their materials. Exposía currently comprises 55 exposés, 306 reviews, and 2,253 inline comments.

We emphasize that the dataset will continue to grow, as the course will be offered annually. We will continue to expand and update the dataset with additional data. In the winter term 2025, 600 students enrolled in the course. We expect at least 30% of the students to consent to the release of their materials, consistent with the previous semester. In general, 600 students are expected to enroll each year, with a higher proportion of students willing to contribute their content, as they understand the significance of data collection in the course that teaches important skills for *writing scientific papers* and for *providing effective feedback on scientific writing*. This is consistent with similar important collection efforts (Dycke et al. (2023)).

4 The Exposía Dataset

4.1 Exposía Dataset Overview

Figure 2 illustrates two main components of the Exposía: *exposés* and *feedback*.

Exposés Each exposé is available in two versions: a *draft* and a revised *final version* produced after receiving feedback. The exposés follow a standard Computer Science research proposal template with the following sections: (1) *Motivation*; (2) *Approach* divided into three parts – *Related Work*, *Theoretical Framework*, and *Methodology*; (3) a *Schedule* outlining the week-by-week plan; and (4) a *Bibliography* (Appendix B).

⁶Grades for each component are derived from criterion-based scores (see Section 4 and Appendix F.3).

⁷Reference hidden for anonymity.

⁸Enrollment was limited because the course was offered for the first time as a pilot.

Feedback: inline comments and reviews Each draft exposé received feedback from instructors and other students. Each feedback consists of *inline comments* and a free-text *review*. An *inline comment* consists of a span-level *highlight* on the draft text and a free-text *comment* connected to it (Figure 2). Each inline comment is also associated with a *tag* from a fixed taxonomy $\mathcal{Y} = \{Strength, Weakness, Highlight, Other\}$. A reviewer first adds inline comments. Then the reviewer uses them as the basis for composing a *review*, a document-level free-text feedback on the draft. Appendix C provides more detail.

4.2 Exposé Assessment Criteria and Rubrics

Both the draft and the final exposé are graded independently by an instructional staff member on 36 criteria and include 36 per-criterion *assessment scores* (Appendix Table 6). Each criterion is rated on an ordinal scale of [0,1], [0,1,2], or [0,-1,-2].

Designing exposé assessment criteria and rubrics

The extensive criterion-based grading scheme is designed to support exposé assessment across multiple aspects. These criteria were derived using 18 guidance documents for writing thesis exposés from major universities in Germany and Austria (Appendix E). These were supplemented with over 100 exposés previously submitted in the Computer Science Department. The 36 criteria center on the different aspects of writing a scientific proposal, with a focus on *content* assessment.

Guided by research in rubric construction and fairness-oriented assessment (Jonsson and Svingby, 2007; Mayring, 2000; Kuckartz and Rädiker, 2023), we organized the criteria into nine high-level categories referred to as *rubrics* (Table 1) that decompose evaluation into orthogonal dimensions, by grouping criteria that evaluate related aspects of the exposé writing. Rubrics are an effective assessment tool and are constructed to ensure *consistency in evaluation across a diverse student body and multiple evaluators* (Northern Illinois University Center for Innovative Teaching and Learning, 2012). See Appendix D.1 for more details.

4.3 Review Assessment Criteria and Rubrics

Each *student review* was independently graded by an instructional staff member and received 24 assessment scores, in line with the 24-criterion grading scheme (Appendix Table 7) organized into eight rubrics (Table 1), with each criterion rated on an ordinal scale of [0,1], [0,1,2], or [-1,1]. The rubrics are shown in Table 1.

Code	Rubric	Description	#Crit.
Exposé			
LA	Lang./Quality	Assessment of language quality and structural coherence of the text.	2
ME	Metadata	Evaluation of title appropriateness and creativity.	2
ST	Form/Structure	Evaluation of document formatting and structural requirements.	2
MO	Motivation	Assessment of problem identification, context establishment, and research question formulation.	7
AP	Approach	Assessment of state of the art analysis, theoretical framework.	8
MD	Methodology	Assessment of the methodology.	7
SC	Schedule	Assessment of project timeline structure, completeness, and realism.	4
BI	Bibliography	Assessment of bibliography consistency and literature relevance.	2
AD	Add'l points	Add points for great submissions or deduct points for major mistakes.	2
Review			
SC	Structure/Clarity	Assessment of overall structure and clarity.	2
LG	Language	Assessment of language quality and tone.	2
FU	Feed Up	Assessment of learning goals and success criteria.	2
FB	Feed Back	Assessment of feedback quality and knowledge transfer.	6
FF	Feed Forward	Assessment of future learning guidance.	2
CQ	Content Quality	Assessment of content accuracy and usefulness.	3
ER	Errors	Assessment of potential errors and issues in feedback.	5
AD	Add'l points	Add points for great reviews or deduct points for major mistakes.	2

Table 1: Exposé and review grading rubrics. The last column shows the number of criteria in the rubric.

Designing review assessment criteria and rubrics

The main goal of assessing reviews was to capture the extent to which peer feedback supports learning, reflection, and revision. Accordingly, the rubrics operationalize key educational principles of effective feedback—clarity, constructiveness, specificity, and usefulness—through an analytic, criterion-based framework that is transparent, fair, and aligned with established feedback theory (Hattie and Timperley, 2007). As such, our rubric design is primarily informed by the feedback model in Hattie and Timperley (2007), which conceptualizes feedback as answering three questions: *Where am I going?* (Feed Up), *How am I going?* (Feed Back), and *Where to next?* (Feed Forward).

Drawing on Narciss (2008), we integrated the components of feedback as criteria within the 8 rubrics. We further incorporated insights from D’Arcy et al. (2024), which analyzes revisions in response to peer reviews in scientific writing, highlighting that *actionable comments*, those offering concrete, suggestions, are more likely to lead to substantive text revisions. As such, our rubrics emphasize argumentation, accuracy, and clear suggestions for improvement as indicators of pedagogically effective feedback (Shute, 2008).

Finally, while actionable feedback promotes learning, ineffective or misleading reviews can

have the opposite effect. To this end, inspired by Du et al. (2024), the criteria include a penalty for unhelpful feedback patterns (e.g., neglect of key elements, unsupported claims, or contradictory advice). Appendix D.2 provides more detail.

5 Exposía Dataset Analysis

5.1 Exposía Dataset Statistics

The dataset statistics are in Table 2. As the dataset comprises only materials for which consent was granted, some exposés lack the full set of three reviews (Table 18), and there exist 184 *orphaned* reviews because the corresponding exposés could not be made public. The total number of student reviews is 110, including 34 reviews with an existing exposé and 76 that are orphaned. Each draft exposé includes on average 41 inline comments. *Weakness* is the most frequent tag (59%) of all comments, followed by *Strength* (18%), *Other* (14%) and *Highlight* (9%).

Exposé revision analysis We analyze the impact of *feedback* by comparing how the final exposé changes relative to the draft: final versions increase in length from 528 to 649 words, on average. We also analyze revisions with respect to exposé criterion-based scores, averaged by rubric. The largest gains occur in *Approach*, *Motivation*, *Lan-*

Topic	Exposés	IC	Reviews
T1	19	798	43
T2	18	785	44
T3	18	670	35
<i>Orphaned reviews</i> *			184
Total Σ	55	2,253	306

Table 2: Statistics on the Exposía corpus. Topics: T1=*The Holographic AI Assistant*, T2=*Interactive LLM-based teaching agents*, T3=*Generative AI in Healthcare*. IC=Inline Comments. *Orphaned reviews counts reviews for which no corresponding exposé exists.

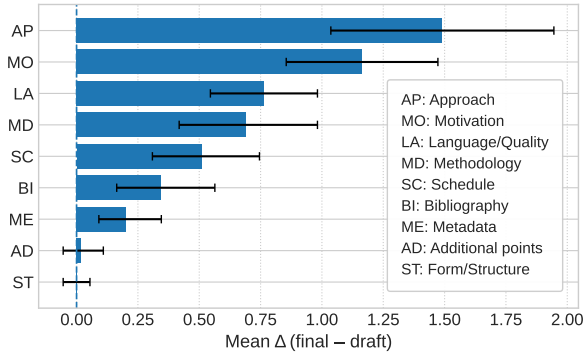


Figure 3: Per-rubric score improvements from draft to final exposé. Each bar shows the mean change in score. Error bars indicate nonparametric bootstrap 95% confidence intervals for each item.

365 *guage/Quality* and *Methodology* (Figure 3). The
366 most substantial revisions focus on conceptual and
367 argumentative elements (see Appendix F.1).

368 **Review analysis** We compare reviews written by
369 instructors vs. students (Appendix F.2): instructor
370 reviews are almost twice as short as student reviews,
371 but have higher lexical density and higher lexical
372 diversity. Appendix F.2 provides more analyses.

373 5.2 Inter-Annotator Agreement

374 We compute inter-annotator agreement (IAA) on
375 exposé and review criterion-based scores. Each
376 draft exposé and student review was independently
377 graded by two instructors on 36 exposé and 24
378 review scoring criteria, respectively.

379 **Rater training** During the semester, each draft
380 exposé and student review was graded by an instructor
381 only once; after the course, each item was graded
382 again for calculating IAA. We refer to instructors
383 who graded during the semester as *Group 1*. They
384 attended weekly sessions and discussed edge cases.
385 Group 1 raters were instructed to resolve uncer-
386 tainty in favor of the student, because grades were
387 at stake. *Group 2* includes a different subset of in-

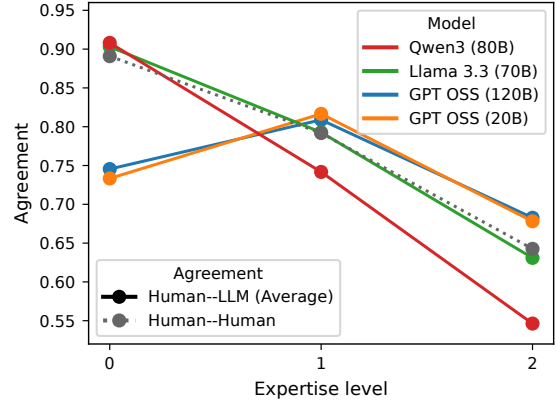


Figure 4: **Human-LLM agreement (QWA) for exposé scoring by expertise level.** The dotted line shows human-human QWA.

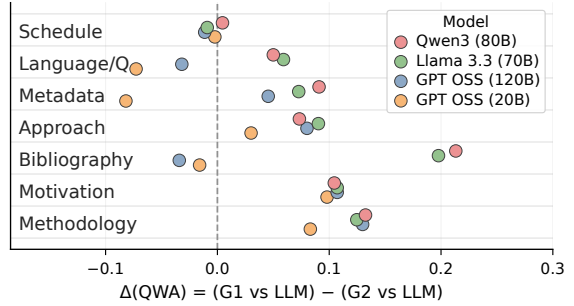


Figure 5: **Group asymmetry in human-LLM agreement by exposé rubric.** Each point shows the difference in agreement between an LLM and human raters (Group 1 (G1) vs. Group 2 (G2)). The vertical line at $\Delta = 0$ indicates equal agreement of the LLM with human G1 and G2 raters. Points to the right indicate the model agrees more with G1 raters than with G2 raters.

388 structors (with one who also participated in Group
389 1). They received the same guidelines, but they
390 were not instructed to resolve uncertainty in favor
391 of the student. See Appendix G.1 for more detail.

392 **IAA metrics** We report Krippendorff’s α (Krippen-
393 dorff, 2004),⁹ Pearson correlation coefficient r , and
394 raw quadratic weighted percent agreement (QWA).
395 We have several criteria with skewed distributions
396 where the majority label accounts for over 90%
397 cases (Tables 19 and 20). In this case, both α and
398 r can be misleading (Jeni et al., 2013; Zhao et al.,
399 2022), and QWA is the most informative metric.
400 Appendix G.2 and G.3 provide more detail.

401 **IAA on exposé assessment scores** Table 3 (top)
402 summarizes IAA results by exposé rubric. Overall
403 agreement yields an average QWA of 0.787, with
404 $\alpha = 0.198$, $r = 0.270$. The observed correlation

⁹Krippendorff’s α can be used in a setup such as ours, where we are comparing multiple raters where raters also belong to one of the two groups.

Rubric	QWA	α	r
Exposé rubrics			
Language/Quality (#2)	0.85	0.10	0.19
Metadata (#2)	0.83	0.00	0.31
Form/Structure (#2)	0.97	0.33	0.34
Motivation (#7)	0.76	0.11	0.17
Approach (#8)	0.77	0.22	0.26
Methodology (#7)	0.62	0.18	0.22
Schedule (#4)	0.93	0.36	0.44
Bibliography (#2)	0.78	0.19	0.42
Add'l points (#2)	0.98	0.32	0.34
All (avg.)	0.79	0.20	0.27
Review rubrics			
Structure/Clarity (#2)	0.86	0.27	0.38
Language (#2)	0.93	0.09	0.17
Feed Up (#2)	0.81	0.43	0.43
Feed Back (#6)	0.86	0.04	0.08
Feed Forward (#2)	0.78	0.02	0.04
Content Quality (#3)	0.94	0.22	0.24
Errors (#5)	0.93	0.10	0.12
Add'l points (#2)	0.98	-0.01	—
All (avg.)	0.89	0.13	0.18

Table 3: IAA on exposé and review scoring. Quadratic weighted agreement (QWA), Krippendorff’s (α), and Pearson r . r for Add'l points (review) is undefined due to zero variance; we exclude it from the average. Result for each rubric is averaged over all criteria in the rubric. The number of criteria per rubric is next to rubric name.

levels are comparable to those reported on automated essay scoring (Persing and Ng, 2014, 2013) and reflect the subjectivity of the task. The rubrics that attain the lowest agreement levels are those related to content (Motivation, Approach, Methodology). Table 19 reports IAA for all 36 criteria.

IAA on review assessment scores Table 3 (bottom) summarizes IAA results by review rubric, computed over student reviews. While QWA is higher for review rubrics than for exposé, α and Pearson coefficient are lower, due to the skewed distributions of many criteria (Appendix G.3). For criterion-level agreement and discussion, see Appendix G.2 and Appendix Table 20.

6 Experiments

6.1 Experimental Setup

Tasks We use Exposita for benchmarking performance on two tasks: (i) criterion-based exposé assessment and (ii) criterion-based review assess-

ment. Models predict scores for individual criteria. **LLM prompting** We assess LLM performance in a zero-shot setting. We evaluate four state-of-the-art open-weight instruction-tuned LLMs with different architectures and sizes (Appendix H). For both tasks, we use a unified prompting strategy and compare two prompting variants: (i) *single-criterion prompts*, which present and score one criterion at a time, and (ii) *combined prompts*, which list all criteria for the respective task in a single call. Cross-criterion synergy could potentially benefit the model, with a lower execution cost. Appendix I provides more detail about the prompts.

Evaluation Following studies on criterion-based scoring (Pack et al., 2024; Huang and Wilson, 2025), we evaluate performance, by comparing scores predicted with an LLM with a score assigned by a human rater, using QWA (Section 5.2). Results are reported on 55 exposés drafts and 110 student reviews. Unless otherwise specified, we compute QWA against Group 1 and Group 2 raters and average the results. For each exposé and review, a corresponding model predicts 31 and 22 scores, respectively,¹⁰ for a total of 1705 exposé criteria scores and 2420 review criteria scores.

6.2 Key Results

Table 4 shows main results on exposé and review assessment (averaged over all criteria – rows *All*). **Single-criterion vs. combined prompting** Across tasks and models, combined prompting yields *higher* agreement, while being more computationally efficient. The only exception is Qwen3 model on exposé scoring. Tables 22 and 23 tabulate criterion- and rubric-level results.

Performance by expertise level We further group criteria by *expertise level*: three levels for exposé, and two for review, based on the degree of domain knowledge required for assessment (Appendix D). Figure 4 demonstrates that exposé criteria requiring more in-domain knowledge display more variability, and, importantly, this aligns with human raters. Appendix Figure 7 reveals that for review scoring, the difference is not as dramatic, suggesting that review assessment does not require as much domain knowledge. See Appendix J for further discussion.

6.3 Additional Analyses

Systematic group bias We observe a tendency for Group 1 raters to assign higher scores (Appendix G). This matters for interpreting LLM-

¹⁰5 exposé and 2 review criteria are excluded (Appendix H).

Exp.	#Crit.	Human	Qwen3 (80B)		Llama 3.3 (70B)		GPT OSS (120B)		GPT OSS (20B)	
			\bar{A}_{all}	\bar{A}_{ind}	\bar{A}_{all}	\bar{A}_{ind}	\bar{A}_{all}	\bar{A}_{ind}	\bar{A}_{all}	\bar{A}_{ind}
Exposé criterion-based scoring (over 31 criteria)										
0	3	0.89	0.91	0.85	0.90	0.90	0.75	0.72	0.73	0.81
1	20	0.79	0.76	0.79	0.79	0.78	0.81	0.78	0.82	0.78
2	8	0.64	0.57	0.62	0.63	0.62	0.68	0.69	0.68	0.69
All	31	0.76	0.72	0.75	0.76	0.75	0.77	0.75	0.77	0.76
Review criterion-based scoring (over 22 criteria)										
1	17	0.89	0.89	0.85	0.91	0.84	0.82	0.75	0.77	0.71
2	5	0.86	0.87	0.82	0.86	0.77	0.67	0.63	0.64	0.63
All	22	0.88	0.89	0.84	0.90	0.82	0.78	0.72	0.74	0.69

Table 4: Quadratic weighted agreement (QWA) for the two tasks, grouped by expertise level (3 and 2 expertise levels for exposé and review criteria, respectively). Exp.=expertise level; #Crit.=number of criteria in the subset. Prompting variants: *all* (combined) and *ind.* (individual). Each LLM column shows human-LLM QWA averaged over both human raters and all criteria. Human column shows agreement between the human raters.

human agreement: LLM agreement with Group 1 is often higher than agreement with Group 2, consistent with LLMs being more generous (Figure 5). We do not observe similar behavior for review scoring (Appendix Figure 8). This is consistent with the more dramatic differences on expertise level for exposé scoring and suggests that review assessment does not require as much in-domain expertise.

Scalability and classroom deployment Combined prompting is both more effective and more efficient computationally (Appendix J and K). This is important because student activity is strongly deadline-driven (Figure 6). Thus, *combined prompting is recommended for classroom deployment.*

7 Discussion

Interpreting model performance relative to human raters Human agreement varies substantially across criteria and expertise levels, indicating that some aspects of academic writing and feedback are inherently more difficult to assess consistently. Within this context, LLMs achieve agreement levels that approach the consensus of human agreement. Notably, *performance degrades on expert-level criteria, mirroring the same dimensions where human agreement is weakest.*

LLM alignment with human groups The tendency of LLMs to align more closely with Group 1 raters suggests that LLMs adopt a more generous grading regime (Flodén, 2024), likely to originate from the positive bias in LLMs (Jain et al., 2025).

Potential tasks enabled by Exposita Beyond the

tasks studied here, Exposita enables a range of research directions. These include modeling draft-feedback-revision links across documents, with numerous applications, such as feedback text generation. More broadly, the dataset supports studying how feedback mediates learning and writing development in higher-education settings.

8 Conclusion

We introduce a novel task definition for exposé and review scoring in an educational context and present the first public dataset, Exposita, of student research exposés and feedback, accompanied by human assessment scores that are based on the fine-grained grading schema we develop.

By integrating drafting, feedback, and revision into a single resource, Exposita enables pedagogically grounded research on writing and feedback in higher education and provides a unified framework that allows for a joint study of feedback and exposé writing quality. Although the dataset is collected in the Computer Science domain, the task formulation and the assessment framework are general enough to inform the development of assessment models in other areas.

We establish baseline performance with state-of-the-art LLMs: models achieve high agreement with human raters on criteria that do not require in-domain expertise, while higher expert-level criteria remain challenging. Our findings provide insight and recommendations for classroom deployment.

534 Limitations

535 This study has several limitations. First, Exposía
536 is shaped by the course consent protocol: since
537 only consented materials are included, not all ex-
538 posés retain the complete set of three reviews. For
539 the same reason, some of the reviews in Exposía
540 are *orphaned*, as the author of the corresponding
541 exposé did not provide consent (Section 5.1). As
542 such, the consent structure reduces the number of
543 complete draft–review–revision instances in the
544 released data. We expect these constraints to di-
545 minish as the course is offered repeatedly and the
546 dataset continues to grow (Section 3). Although
547 the current size of Exposía is modest compared to
548 some other NLP datasets, it is substantial for the
549 educational domain, particularly given its complex-
550 ity and the richness of its annotations. We discuss
551 the time investment to create the various Exposía
552 components (feedback and criterion-based scores
553 on exposés and reviews) in Appendix F.3.

554 Second, although the course covers three areas –
555 Human–Computer Interaction, Natural Language
556 Processing, and Data Mining – all topics remain
557 within Computer Science and AI. However, we
558 believe that the task formulation and the assessment
559 framework are general enough and can inform the
560 development of assessment models in other areas.

561 Third, our model comparison prioritizes repro-
562 ducibility and transparency in an educational set-
563 ting and therefore focuses on openly available mod-
564 els. As a result, we do not evaluate proprietary
565 closed-source LLMs, which may behave differently
566 and limit the scope of our model comparison.

567 Finally, and most importantly, although LLM-
568 based scoring can be attractive for reducing instruc-
569 tional workload, our results should not be inter-
570 preted as support for fully automated grading: the
571 models exhibit systematic scoring tendencies and
572 may miss substantive issues that are crucial for
573 guiding revision, and automated grading can in-
574 centivize strategic “gaming” rather than genuine
575 improvement (Alqahtani et al., 2023); we therefore
576 view LLMs primarily as decision-support tools,
577 and as common in AES, recommend a human-
578 in-the-loop workflow when using them for assess-
579 ment.

580 Ethical Considerations

581 The data collection was approved by the Univer-
582 sity’s institutional review board and data protection
583 office [Details upon acceptance]. Data was col-
584 lected based on informed consent (all instructors

and 55 students), following a three-tier structure:
585 (1) general consent for data processing within the
586 review platform, (2) optional consent for contribut-
587 ing textual and highlight data, and (3) optional con-
588 sent for contributing interaction data. Participation
589 in the study was strictly voluntary and had no effect
590 on grading. No monetary compensation was pro-
591 vided to students, as participation was part of the
592 course setting; importantly, contributing research
593 data was optional. Instructors were compensated
594 via regular University employment contracts, fol-
595 lowing institutional pay scales and policies. The
596 compensation was not contingent on participation
597 in the study.

598 We release Exposía to enable research on aca-
599 demic writing, peer feedback, and assessment in
600 educational settings. We stress that our findings are
601 NOT to be interpreted as justification for fully auto-
602 matic grading: all assessment-oriented applications
603 of AI should be restricted to decision-support roles
604 with human oversight, in line with the limitations
605 outlined above.

607 References

- 608 Tariq Alqahtani, Hisham A. Badreldin, Mohammed Al-
609 rashed, Abdulrahman I. Alshaya, Sahar S. Alghamdi,
610 Khalid bin Saleh, Shuroug A. Alowais, Omar A. Al-
611 shaya, Ishrat Rahman, Majed S. Al Yami, and Ab-
612 dulkareem M. Albekairy. 2023. [The emergent role
613 of artificial intelligence, natural learning processing,
614 and large language models in higher education and re-
615 search](#). *Research in Social and Administrative Phar-
616 macy*, 19(8):1236–1242.
- 617 Scott Crossley, Yu Tian, Perpetual Baffour, Alex
618 Franklin, Youngmeen Kim, Wesley Morris, Meg Ben-
619 ner, Aigner Picou, and Ulrich Boser. 2023. [The en-
620 glish language learner insight, proficiency and skills
621 evaluation \(ellipse\) corpus](#). *International Journal of
622 Learner Corpus Research*, 9(2):248–269.
- 623 Mike D’Arcy, Alexis Ross, Erin Bransom, Bailey Kuehl,
624 Jonathan Bragg, Tom Hope, and Doug Downey. 2024.
625 [ARIES: A corpus of scientific paper edits made in
626 response to peer reviews](#). In *Proceedings of the 62nd
627 Annual Meeting of the Association for Computational
628 Linguistics (Volume 1: Long Papers)*, pages 6985–
629 7001, Bangkok, Thailand. Association for Computa-
630 tional Linguistics.
- 631 Jiangshu Du, Yibo Wang, Wenting Zhao, Zhongfen
632 Deng, Shuaiqi Liu, Renze Lou, Henry Peng Zou,
633 Pranav Narayanan Venkit, Nan Zhang, Mukund Sri-
634 nath, Haoran Ranran Zhang, Vipul Gupta, Yinghui Li,
635 Tao Li, Fei Wang, Qin Liu, Tianlin Liu, Pengzhi Gao,
636 Congying Xia, and 21 others. 2024. [LLMs assist
637 NLP researchers: Critique paper \(meta\)-reviewing](#).

750	Sheng Lu, Iliia Kuznetsov, and Iryna Gurevych. 2025. Identifying aspects in peer reviews . In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 6145–6167, Suzhou, China. Association for Computational Linguistics.	805
751		806
752		807
753		
754		
755	Sandeep Mathias and Pushpak Bhattacharyya. 2018. ASAP++: Enriching the ASAP automated essay grading dataset with essay attribute scores . In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , Miyazaki, Japan. European Language Resources Association (ELRA).	808
756		809
757		810
758		811
759		812
760		813
761		814
762	Philipp Mayring. 2000. Qualitative content analysis . <i>Forum Qualitative Sozialforschung / Forum: Qualitative Social Research</i> , 1(2).	815
763		816
764		817
765	Dr. Rizwana Muneer, Rida Batool, and Saima Zehra. 2020. A qualitative study on challenges that post-graduate students face in research proposal writing at university level . <i>International Journal of Social Sciences: Current and Future Research Trends</i> , 5(01):1–6.	818
766		819
767		820
768		821
769		822
770		823
771	S. Narciss. 2008. Feedback strategies for interactive learning tasks. In J. M. Spector, M. D. Merrill, J. van Merriënboer, and D. M. Driscoll, editors, <i>Handbook of research on educational communications and technology</i> , pages 125–144. Routledge, New York.	824
772		825
773		826
774		
775		
776	Northern Illinois University Center for Innovative Teaching and Learning. 2012. Rubrics for assessment. https://www.niu.edu/citl/resources/guides/instructional-guide . Retrieved from Northern Illinois University Center for Innovative Teaching and Learning website.	827
777		828
778		829
779		830
780		831
781		832
782	OpenAI. 2025. gpt-oss-120b & gpt-oss-20b model card .	833
783	Austin Pack, Alex Barrett, and Juan Escalante. 2024. Large language models and automated essay scoring of english language learner writing: Insights into validity and reliability . <i>Computers and Education: Artificial Intelligence</i> , 6:100234.	834
784		835
785		
786		
787		
788	Ellis B. Page. 1966. The imminence of grading essays by computer. <i>Phi Delta Kappan</i> , 47:238–243.	836
789		837
790	Isaac Persing and Vincent Ng. 2013. Modeling thesis clarity in student essays . In <i>Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 260–269, Sofia, Bulgaria. Association for Computational Linguistics.	838
791		839
792		840
793		841
794		
795		
796	Isaac Persing and Vincent Ng. 2014. Modeling prompt adherence in student essays . In <i>Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1534–1543, Baltimore, Maryland. Association for Computational Linguistics.	842
797		843
798		844
799		845
800		846
801		847
802	Ildiko Pilan, John Lee, Chak Yan Yeung, and Jonathan Webster. 2020. A dataset for investigating the impact of feedback on student revision outcome . In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 332–339, Marseille, France. European Language Resources Association.	848
803		849
804		850
		851
		852
		853
		854
		855
		856
		857
		858
	Sukannya Purkayastha, Zhuang Li, Anne Lauscher, Lizhen Qu, and Iryna Gurevych. 2025. LazyReview: A dataset for uncovering lazy thinking in NLP peer reviews . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3280–3308, Vienna, Austria. Association for Computational Linguistics.	
	Afifa Rastri, Yanti Sri Rezeki, Urai Salam, Dwi Riyanti, and Surmiyati Surmiyati. 2023. An analysis of students’ problems in writing a research proposal . <i>Actiya: Journal of Teaching and Education</i> , 5:57–71.	
	Roman Rietsche, Andrew Caines, Cornelius Schramm, Dominik Pfütze, and Paula Buttery. 2022. The specificity and helpfulness of peer-to-peer feedback in higher education . In <i>Proceedings of the 17th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2022)</i> , pages 107–117, Seattle, Washington. Association for Computational Linguistics.	
	Qian Ruan, Iliia Kuznetsov, and Iryna Gurevych. 2024. Are large language models good classifiers? a study on edit intent classification in scientific document revisions . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 15049–15067, Miami, Florida, USA. Association for Computational Linguistics.	
	Valerie J. Shute. 2008. Focus on formative feedback . <i>Review of Educational Research</i> , 78(1):153–189.	
	Moritz Staudinger, Wojciech Kusa, Florina Piroi, and Allan Hanbury. 2024. An analysis of tasks and datasets in peer reviewing . In <i>Proceedings of the Fourth Workshop on Scholarly Document Processing (SDP 2024)</i> , pages 257–268, Bangkok, Thailand. Association for Computational Linguistics.	
	Qwen Team. 2025. Qwen3 technical report .	
	Yong Wu and Christian D. Schunn. 2021. The effects of providing and receiving peer feedback on writing performance and learning of secondary school students . <i>American Educational Research Journal</i> , 58(3):492–526.	
	Jin Xue, Xiaoyi Tang, and Liyan Zheng. 2021. A hierarchical bert-based transfer learning approach for multi-dimensional essay scoring . <i>IEEE Access</i> , 9:125403–125415.	
	Fan Zhang, Homa B. Hashemi, Rebecca Hwa, and Diane Litman. 2017. A corpus of annotated revisions for studying argumentative writing . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1568–1578, Vancouver, Canada. Association for Computational Linguistics.	

859	Xinshu Zhao, Guangchao Charles Feng, Song Harris	C Additional Details on the Exposita	903
860	Ao, and Piper Liping Liu. 2022. Interrater reliability	Dataset	904
861	estimators tested against true interrater reliabilities.		
862	<i>BMC Med Res Methodol</i> , 22(1):232.	Exposés Each exposé is available in two versions:	905
863		a <i>draft</i> and a revised <i>final version</i> produced after	906
864	A Course Syllabus	receiving feedback.	907
865	This appendix documents the syllabus of the under-	Each version is provided as (i) a compiled PDF	908
866	graduate computer science course <i>Introduction to</i>	file and (ii) the original L ^A T _E X source, including	909
867	<i>Scientific Work</i> , which provided the instructional	the separate bibliography file (.bib) and all textual	910
868	setting for the data collection described in Sec-	content.	911
869	tion 3.	Reviews Every exposé draft is accompanied by	912
870	The course covers foundational concepts and	two <i>student reviews</i> written by fellow students tak-	913
871	practical skills required for conducting scientific	ing the course and at least one <i>instructor review</i> .	914
872	research. Here we present the weekly lectures:	Instructional staff included nine instructors and con-	915
873	1. What is science? Course overview and intro-	sisted of PhD students, postdoctoral researchers,	916
874	duction	and teaching assistants holding a Master’s degree	917
875	2. Empirical methods and research design	(Appendix G.1 (Group 1)). The feedback was	918
876	3. Literature review and academic search strate-	double-blind among students, but not with respect	919
877	gies	to the instructional staff. During the revision pro-	920
878	4. Project proposals (exposés); <i>start writing ex-</i>	cess, students could identify whether a specific	921
879	<i>posé draft</i>	feedback came from an expert or a student.	922
880	5. Scientific writing	Each review is released as (i) an HTML docu-	923
881	6. L ^A T _E X for academic writing	ment that preserves the original formatting used	924
882	7. Research Questions and practical exercise on	on the review platform and (ii) a plain-text render-	925
883	academic literature search	ing. For all reviews, we provide basic metadata	926
884	8. Introduction to peer review	including a synthetic reviewer name, reviewer role,	927
885	9. Feedback; <i>start writing peer feedback</i>	document size, the associated exposé (if available),	928
886	10. Generative AI in scientific work	and an anonymized identifier.	929
887	11. Research ethics I	For all reviews for which the author provided ad-	930
888	12. Research ethics II	ditional consent for behavioral logging, we release	931
889	13. Workshop: Software development	the full edit history of the review text as stored	932
890	14. Workshop: Research and project management	by the platform, capturing the successive revisions	933
891	B Exposé Template and Writing	made by the reviewer while composing their feed-	934
892	Guidelines	back. This results in sequences of review drafts that	935
893	Panel <i>Exposé template</i> shows the content of the	can be aligned chronologically via timestamps.	936
894	L ^A T _E X exposé template that students received as	Inline comments and voting Beyond document-	937
895	part of the course and used to create their exposés.	level feedback, the platform supports fine-grained,	938
896	While the complete instructions were provided in	span-based comments on the exposé draft. For	939
897	an extensive slide-based lecture (Appendix A), the	each exposé, we collect a set of inline comments	940
898	excerpt reproduced here is the concise reference	created by students and instructors while reading	941
899	sheet embedded directly in the template. Because	the draft. Each comment is attached to a character-	942
900	this template guided how students organized and	level span in the PDF-rendered text of the exposé	943
901	formatted their drafts, we include it to ensure trans-	and is associated with (i) optional free-text con-	944
902	parency and to document the pedagogical context	tent, (ii) a categorical tag from a fixed taxonomy	945
	in which the dataset was created.	$\mathcal{Y} = \{Strength, Weakness, Highlight, Other\}$ and	946
		(iii) metadata about the author role (student or in-	947
		structor) and the underlying review document it	948
		belongs to.	949
		The platform further allows other participants,	950
		but not the comment author, to cast binary votes	951
		(<i>thumbs up / thumbs down</i>) on each comment.	952
		These votes were intended as judgments of com-	953

Exposé template

1 Motivation

Hook – Introduce the topic and catch the reader’s interest.

- *Concrete*: Quote, example, real-world event.
- *Helpful*: Directly useful or valuable information.
- *Convincing*: Strong argument or empirical fact.
- *Shocking*: Unexpected or counterintuitive finding.
- *Personal*: Research journey or applied experience.
- *Understandable*: Link complex ideas to everyday life.

Anchor – Situate the problem concisely.

- Discipline, field, or domain.
- How the problem has been addressed so far.
- Current academic relevance.
- Deficits or criticism in existing literature.

Teaser – Identify the gap and position your research.

Research Questions

- Break the main RQ into meaningful sub-questions.
- Clarify terminology; define if needed.
- Set clear boundaries and scope.

2 Approach

State of the Art (SOTA)

- Refer to relevant existing work.
- Highlight research deficits and criticism.
- Combine different traditions or perspectives.

Theoretical Framework

- Introduce theories used in the study.
- Explain key concepts and definitions.
- Describe how the theory informs the RQs.

Note: Theories need not be SOTA but should be connected to it.

Methodology

- Specify method (quantitative, qualitative, mixed, source research).
- Define target population and sampling.

- State data sources and availability.
- Clarify agreements or permissions.
- Anticipate potential difficulties.
- Discuss methodological limitations.
- Explain relevance and scope.
- Describe variables and their relationships.
- Outline practical implementation (software, documents, helpers).
- Clarify epistemic goals (knowledge to be gained).
- If building on prior work, state progress so far.

3 Schedule

Provide a chronological work plan with realistic milestones.

Tips:

- Break tasks into concrete steps.
- Include checkpoints and buffer time.
- Use weekly or monthly units depending on project size.
- Seek advice from experienced researchers.

Example Schedule

Month	Week	Task	Description
1	1	Planning	Work plan, coordination
1	2–3	Literature Review	Search and initial analysis
1	4	Literature Evaluation	Critical review and selection
2	5–6	Theoretical Framework	Develop theory and methodology
2	7	Concept Development	Draft interview/ analysis guide
2	8	Scenario Planning	Design simulations or experiments
3	9–10	Recruitment	Acquire participants
3	11	Data Collection	Conduct interviews/ experiments
3	12	Data Preparation	Clean and pre-analyze data
4	13–14	Qualitative Analysis	In-depth qualitative analysis
4	15–16	Quantitative Analysis	Statistical analysis
5	17	Writing Results	Compose results section
5	18	Discussion	Interpret findings
5	19	Methodology Review	Finalize methodology
5	20	Full Revision	Edit and review all chapters
6	21	Final Editing	Final proofreading, formatting
6	22–24	Finalization	Buffer, printing, submission

Reference

[1] Helmut Balzert, Christian Schäfer, Marion Schröder, Uwe Kern, Roman Bendisch, and Klaus Zeppenfeld. *Wissenschaftliches Arbeiten: Ethik, Inhalt und Form wissenschaftlicher Arbeiten, Handwerkszeug, Quellen, Projektmanagement, Präsentation*. W3L-Verlag, Herdecke/Witten, 2011. ISBN: 978-3-937137-59-9.

ment helpfulness. However, interviews with participants indicate that student authors sometimes used the voting mechanism as a personal *checklist* during revision (marking comments as *done* rather than as helpful), despite explicit instructions from the instructors to use it purely as a helpfulness signal. We therefore provide the raw voting data, together with information about voter role, but caution that votes do not form a clean ground truth for perceived usefulness.

Behavioral interaction logs For users who provided explicit additional consent, we release interaction logs from the review platform, covering 45 students and 9 instructors. Each record contains an action type, a JSON-encoded data payload, a timestamp, and an anonymized user identifier. Logged events are timestamped UI interactions (Table 5); an overview of activity over the semester is shown in Figure 6 and discussed in Appendix L. Actions including navigation and visibility changes (e.g., Navigation step, PDF page show/hide, Browser Tab show/hide), scrolling and resizing (e.g., Annotation Page scroll, PDF scroll, Sidebar scroll, PDF (Browser) resize), interface controls and modal events (e.g., Sidebar click, Topbar click, Modal show/Modal hide, Modal Button click), and lightweight editing interactions (e.g., Text select, Copy text, Text paste).

Authorship metadata and anonymization Each exposé, review, inline comment, vote, and interaction event is associated with an anonymized author identifier and a role label (*student, junior, senior*). To support more interpretable examples and qualitative analysis, we replace real names with consistently applied synthetic names generated with Faker.¹¹ No raw platform identifiers, email addresses, or other personally identifying information are released. This setup enables the analysis of role-specific behavior and bias while preserving participant privacy.

D Criteria and Rubrics for Exposé and Review Assessment

This appendix provides an overview of exposé and review criteria used to grade the exposé and student reviews, respectively.

Action	Students (n=45)	Instructors (n=9)	Total
Mouse move	73,811	131,272	205,083
PDF page show/hide	53,378	102,120	155,498
Annotation Page scroll	18,944	31,080	50,024
Browser Tab show/hide	16,810	31,657	48,467
PDF (Browser) resize	13,377	25,475	38,852
Text select	7,403	8,043	15,446
Sidebar scroll	2,445	7,248	9,693
Navigation step	2,705	3,742	6,447
Sidebar click	2,882	3,382	6,264
Topbar click	1,826	3,512	5,338
PDF scroll	2,206	2,179	4,385
Copy text	1,403	2,371	3,774
Table click	1,811	1,164	2,975
Text paste	1,208	1,435	2,643
Browser load	715	1,364	2,079
Modal show	889	324	1,213
Modal hide	396	304	700
Search used	102	369	471
Modal Button click	259	176	435
Total	202,570	357,217	559,787

Table 5: Number of behaviour data rows per action, split by user group.

Expertise levels Each criterion is assigned an *expertise level*¹², indicating the type of judgment needed to assess the quality of the exposé or review with respect to this criterion. We introduce a three-level expertise scheme (expertise levels by criteria are shown in Tables 6 and 7 for exposé and review criteria, respectively):

- **Level 0 (formal).** No academic or expert knowledge required. Largely objective. Sample exposé criteria: *Metadata available; Template used; Number of pages*. Review criteria: *None*.
- **Level 1 (academic).** Comprehension-based judgments of clarity, coherence, and basic research sense-making, but no specific in-domain expertise is required. Sample exposé criteria: *Language quality; Hook existing*. Sample review criteria: *Summary available; Language*.
- **Level 2 (expert).** Expert judgment of the soundness and adequacy of the reasoning and its alignment with disciplinary standards. Sample exposé criteria: *Teaser (RQ); Relevance of theoretical framework*. Sample re-

¹²Two co-authors inspected all criteria and labeled each for the expertise level required for scoring an exposé for a given criterion.

¹¹<https://faker.readthedocs.io/en/master>

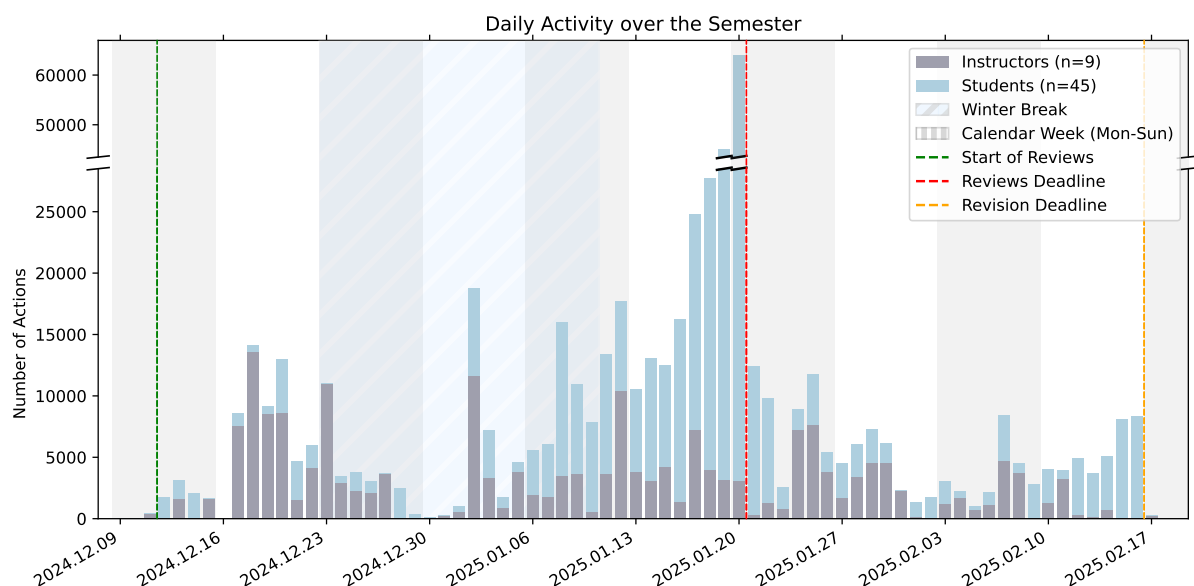


Figure 6: Daily number of actions across the semester for instructors ($n=9$) and students ($n=45$). Bars show per-day activity (instructors in gray, students in blue). Alternating vertical shading indicates calendar weeks (Mon–Sun), and the hatched region marks the winter break. Dashed lines mark the start of reviews (green), the review deadline (red), and the revision deadline (orange). The y-axis is broken to accommodate the peak activity before the review deadline.

view criteria: *Knowledge of correct response*;
Learning goals.

D.1 Exposé Criteria

Table 6 lists all the 36 exposé assessment criteria, grouped into nine high-level rubrics. These criteria define the scoring rules applied to each submitted exposé.

Refinement of criteria and rubrics for exposé assessment Before the course was run, the designed criteria and rubrics were piloted on exposés previously submitted within the department: two members of the instructional staff assessed together a small sample of old student exposés, and other instructional staff members went over these to give feedback. This feedback was used to verify that every criterion is assessable and the criteria definitions are clear. Criteria definitions were adjusted accordingly for clarity and consistency.

D.2 Review Criteria

Table 7 lists all the 24 review assessment criteria grouped into eight high-level rubrics. These criteria define the scoring rules applied to each submitted student review.

Refinement of criteria and rubric for review assessment As with the exposé rubrics, all criteria and rubrics were reviewed by instructors and piloted on sample reviews before course deployment.

Code	Criterion	Expert level	Point descriptors
Rubric: Language/Quality (LA) <i>Assessment of language quality and structural coherence of the text.</i>			
LA1	Language quality	1	0: many linguistic errors (e.g. wrong tense, wrong personal pronouns) 1: few linguistic errors, clear and comprehensible sentences 2: no linguistic errors, clear and comprehensible sentences with consistent terminology
LA2	Common Thread	1	0: No recognizable common thread, the structure is unclear and jumpy 1: Partly recognizable, but the structure is not continuous 2: Text has a continuous structure that is easy to follow, common thread is clear and recognizable from beginning to end
Rubric: Metadata (ME) <i>Evaluation of title appropriateness and creativity.</i>			
ME1	Preliminary title	1	0: Title does not match the topic 1: Title fits the topic; similar title to the one in the topic suggestions 2: Title fits the topic; creative title changed in a novel way and/or fits into the story of the exposé
ME2	Metadata available	0	0: Name, date and matriculation number are not entered correctly/changed 1: Name, date, matriculation number are correct and updated
Rubric: Form/Structure (ST) <i>Evaluation of document formatting and structural requirements.</i>			
ST1	Template used	0	0: Latex template is not used 1: Latex template is used
ST2	Number of pages	0	0: > three pages (motivation + approach only) 1: <= three sides (motivation + approach only)
Rubric: Motivation (MO) <i>Assessment of problem identification, context establishment, and research question formulation.</i>			
MO1	Hook existing	1	0: Hook not available 1: Hook present and suitable for the topic
MO2	Anker	1	0: The problem is not mentioned at all. 1: The problem is stated in general terms without specific details. It is recognizable which problem is being addressed, but important information is omitted. 2: The problem is described in detail. All relevant aspects of the problem are mentioned and clarify its scope and possible effects.
MO3	Domain	1	0: The domain or context in which the problem occurs is not mentioned.

(continued on next page)

Code	Criterion	Expert level	Point descriptors
MO4	Problem relevance	1	<p>1: The domain or the relevant expert area of the problem is mentioned. It is clear in which environment the problem exists.</p> <p>0: The relevance or importance of the problem is not addressed. It remains unclear why this problem is relevant or worth discussing.</p>
MO5	Problem handling	1	<p>1: The importance and relevance of the problem is described. It is made clear why the problem is important and what impact or consequences it could have (e.g. for specific individuals, organisations or society as a whole).</p> <p>0: current handling of the problem or current criticism is not given</p> <p>1: current handling of the problem or current criticism is given</p>
MO6	Teaser (RQ)	2	<p>0: There is no research question available</p> <p>1: A research question exists, but it has weaknesses, e.g. due to unclear terminology or vague formulation. It is difficult to recognise what is to be investigated and how the question fits the problem description.</p> <p>2: The research question is formulated clearly and precisely. It clearly conveys what is to be investigated, and fits thematically with the problem description.</p>
MO7	RQ Limitation	2	<p>0: The research question is either too general or too broad, so that it does not seem realistic to answer it within the scope of a Bachelor's thesis, to answer it in the context of a Bachelor's thesis.</p> <p>1: The research question is precise and clearly defined so that it can realistically be addressed in the context of a Bachelor's thesis. The question is formulated in concrete terms and describes specifically which aspects are to be investigated without deviating from the topic.</p>

Rubric: Approach (AP)

Assessment of state of the art analysis, theoretical framework.

AP1	SOTA	1	<p>0: SOTA not available</p> <p>1: SOTA present and correctly cited</p>
AP2	SOTA Relevance	1	<p>0: The relevance of SOTA is not present or unclear.</p> <p>1: The relevance of SOTA is mentioned, but is only partially comprehensible or weakly presented.</p> <p>2: It is clearly shown how SOTA is relevant to the topic of the exposé and why it is relevant to the research.</p>
AP3	SOTA Weaknesses	1	<p>0: SOTA's weak points are not mentioned</p> <p>1: SOTA's weaknesses are mentioned, but only superficially or unclearly described</p> <p>2: SOTA's weaknesses are clearly identified and precisely explained; it is clear which weaknesses are relevant for own research</p>

(continued on next page)

Code	Criterion	Expert level	Point descriptors
AP4	SOTA Delimitation	1	0: No differentiation of own performance from the current state of research 1: The own achievement is differentiated from SOTA; it is explicitly emphasized, what is new or unique about your own research
AP5	SOTA Combination	2	0: No meaningful combination of existing material (even if reference is made to existing material, but its combination or consolidation is unconvincing or unclear) 1: Existing material is clearly and meaningfully combined; the added value of the combination is clear.
AP6	Theoretical Framework	1	0: The theoretical framework does not exist or is completely misleading 1: The theoretical framework is present, but unclear, poorly structured or difficult to understand 2: The theoretical framework is clearly and logically structured. The theoretical approaches are presented in a understandable and comprehensible
AP7	Relevance of theoretical framework	2	0: The theoretical framework shows no connection to the research topic, or the theory is irrelevant 1: The theoretical framework clearly demonstrates how the selected theories directly relate to the research topic and support the research question
AP8	Methodology Availability	1	0: Methodology not available 1: Methodology available and suitable for answering the RQ(s)

Rubric: Methodology (MD)
Assessment of the methodology.

MD1	Methodology Completeness	1	0: Methodology is available, but not all RQ(s) are addressed 1: Methodology is available and addresses all points of the research questions
MD2	Methodology Relevance	1	0: The relevance of the methodology for the research project is not clear or is unclear 1: It is clearly and precisely demonstrated why the chosen methodology is relevant to the research project and how it contributes to the research questions
MD3	Methodology Target group	1	0: The target group is not specified or there is no precise differentiation 1: The target group is clearly and precisely defined and it is clear why it is relevant to the methodology
MD4	Methodology Existing Material	1	0: No reference is made to the existing material or the material is inappropriate 1: The existing material is appropriately and accurately related to the topic; it is clear that this material is relevant to the methodology

(continued on next page)

Code	Criterion	Expert level	Point descriptors
MD5	Methodology Difficulties	2	<p>0: Potential difficulties in the practical implementation of the methods are not addressed. Potential challenges that could hinder the research process are overlooked.</p> <p>1: Potential difficulties in the practical implementation of the methods are described clearly and precisely. Concrete solutions or strategies to overcome these difficulties are also presented, in order to organize the research process as efficiently and target-oriented as possible.</p>
MD6	Methodology Possibilities/restrictions	2	<p>0: The possibilities or limitations of the methods used are not addressed, or they are only mentioned superficially, without concrete details. It remains unclear how these methods contribute to answering the research question and which limitations could possibly influence the validity of the results.</p> <p>1: The possibilities and limitations of the methods are described clearly and precisely. It is clear what strengths the methods bring to the study and what weaknesses or limitations could possibly influence the results. Limitations could possibly influence the results. The choice of methods is critically reflected upon and related to the research project.</p>
MD7	Methodology Details	2	<p>0: No additional explanations are given on the methodology or implementation</p> <p>1: The methodology is explained comprehensively and precisely with additional explanations (e.g. details on implementation), so that the procedure is clearly comprehensible</p>

Rubric: Schedule (SC)

Assessment of project timeline structure, completeness, and realism.

SC1	Schedule Availability	0	<p>0: Tabular schedule not available</p> <p>1: Tabular schedule in chronological order available</p>
SC2	Schedule Completeness	1	<p>0: Tabular schedule contains gaps</p> <p>1: Tabular schedule available from start to finish and reasonably divided into blocks</p>
SC3	Schedule Block Description	0	<p>0: not available or only headlines available</p> <p>1: individual blocks described (not only headlines)</p>
SC4	Schedule Realistic Relevance	1	<p>0: The time distribution is obviously unrealistic and important work steps are missing. The research questions are not or only insufficiently addressed</p> <p>1: The timeline contains the most important steps, but some of the timelines are unrealistic or unclear. It is partially unclear how the steps contribute to answering the research questions. The timeline does not appear to be fully aligned with the Bachelor's thesis context.</p>

(continued on next page)

Code	Criterion	Expert level	Point descriptors
			2: The timeline is clearly structured, realistic and aligned with the scope of the Bachelor's thesis. All important work steps are included and organised in a sensible time frame. Each step contributes to answering the research questions, and the entire plan fits coherently into the exposé and the context of the thesis.
Rubric: Bibliography (BI) <i>Assessment of bibliography consistency and literature relevance.</i>			
BI1	Bibliography Consistency	0	0: Bibliography with different citation styles 1: Bibliography is standardized (consistent citation style) and essentially all information on the literature is available
BI2	Key literature	2	0: Literature does not fit the topic 1: Literature fits the topic 2: Literature fits the topic and most of it is highly relevant
Rubric: Additional points (AD) <i>Allow additional points for very good submissions and major mistakes in the submission.</i>			
AD1	Additional points	1	0: No additional strengths beyond existing criteria. 1: Some additional strengths beyond existing criteria. 2: Clear exceptional strengths beyond existing criteria.
AD2	Negative points	1	0: No major issues beyond existing criteria. -1: Notable issues beyond existing criteria. -2: Serious issues that strongly reduce overall quality.

Table 6: 36 criteria used to grade student exposés. We organize the criteria into nine rubrics, by grouping together criteria that evaluate related aspects of exposé writing.

Code	Criterion	Expert	Point descriptors
Rubric: Structure and Clarity (SC) <i>Assessment of overall structure and clarity.</i>			
SC1	Summary available	1	<p>0: No summary available or, the summary does not reflect understanding of the performance</p> <p>1: Simple summary that demonstrates basic understanding of the performance</p> <p>2: Concise and accurate summary that clearly reflects understanding of performance and builds confidence in feedback</p>
SC2	Clarity	2	<p>0: Feedback needs considerable improvement in clarity and comprehensibility, comments are disjointed or difficult to follow</p> <p>1: Feedback is largely clear, but sometimes lacks specific references to particular lines, pages, sections or figures/tables</p> <p>2: Feedback is clear, well-structured and easy to understand, with coherent comments and precise references to specific points in the text</p>
Rubric: Language (LG) <i>Assessment of language quality and tone.</i>			
LG1	Language	1	<p>0: Many linguistic errors (e.g. wrong tense, wrong personal pronouns)</p> <p>1: Almost no linguistic errors, clear and comprehensible sentences with standardised terminology</p>
LG2	Tone	1	<p>0: The tone is inappropriate, sounds judgemental or reproachful, criticism comes across as destructive or personal</p> <p>1: The tone is generally respectful, avoids personal attacks, but remains unclear or less friendly in places</p> <p>2: The tone is consistently calm, respectful and polite; criticism is constructive, factual and only directed at the text, not the person</p>
Rubric: Feed Up (FU) <i>Assessment of learning goals and success criteria.</i>			
FU1	Learning Goals	2	<p>0: No learning objectives, unclear what is to be achieved, no orientation</p> <p>1: Partly unclear learning objectives, generally formulated, offer only limited orientation</p> <p>2: Clear learning objectives, specific, challenging, comparable and provide clear direction</p>
FU2	Success Criteria	1	<p>0: No success criteria, success unclear or vaguely formulated, objectives not defined</p> <p>1: Partly unclear success criteria, available but imprecise or general, orientation partly missing</p> <p>2: Clear success criteria, concrete, measurable goals and clear definition of success</p>
Rubric: Feed Back (FB) <i>Assessment of feedback quality and knowledge transfer.</i>			
<i>(continued on next page)</i>			

Code	Criterion	Expert	Point descriptors
FB1	Knowledge of result	1	0: No information as to whether the task has been solved or not 1: Information given as to whether the task has been solved or not
FB2	Knowledge of correct response	2	0: No information on the correct answer 1: Correct answer is clearly communicated
FB3	Knowledge about task constraints	1	0: No information on rules, requirements or restrictions of the task 1: Clear information on rules, requirements or restrictions of the task (e.g. specifications)
FB4	Knowledge about concepts	1	0: No information on conceptual knowledge 1: Basic information on relevant concepts is provided
FB5	Knowledge about mistakes	1	0: No information about errors 1: Information on the number, location, source or type of errors
FB6	Self-feedback	1	0: Self-feedback is available that has not been combined with task, process or self-regulation (including feedback on the learner's abilities (i.e. intelligence or talent)) 1: No self-feedback is available or, if available, in combination with task, process or self-regulation

Rubric: Feed Forward (FF)

Assessment of future learning guidance.

FF1	Self-regulation	1	0: No indication of how to approach the task, no recognisable approaches to self-monitoring or self-assessment 1: Either indications of strategies for working on the task and recognising errors (e.g. suggestions for improvement, targeted search for information) or approaches to self-monitoring and self-assessment (e.g. reflection on own understanding, strategies or willingness to seek help)
FF2	Learning Skills	1	0: No information on how to develop or improve learning skills 1: Information on how to improve the learning objective and how to deal with challenges

Rubric: Content Quality (CQ)

Assessment of content accuracy and usefulness.

CQ1	Correctness	2	0: The content has significant deficiencies in terms of correctness 1: The content is largely correct 2: The content is completely correct
CQ2	Actionability	1	0: No specific instructions available, feedback is difficult to implement 1: Some instructions available, but only partially clear or specific

(continued on next page)

Code	Criterion	Expert	Point descriptors
CQ3	Argumentation	1	<p>2: Clear, specific feedback with actionable instructions, e.g. comparisons with existing methods, application to other tasks, definitions, explanations, discussions, additional analyses or specific suggestions to remove confusion</p> <p>0: Statements or conclusions are not supported by evidence or comprehensible explanations. The argumentation remains superficial or unsystematic, making the content unconvincing.</p> <p>1: Statements are supported by comprehensible evidence, examples or explanations. The argumentation shows a clear link between evidence and conclusions, which strengthens the quality of the content.</p>
Rubric: Errors (ER)			
<i>Assessment of potential errors and issues in feedback.</i>			
ER1	Neglect	2	<p>0: All important details have been considered and no unnecessary questions are asked</p> <p>-1: Important details have been overlooked, leading to unnecessary questions or criticism</p>
ER2	Vague Critic	1	<p>0: Criticism is specific and clearly states what is missing or can be improved</p> <p>-1: Unspecific criticism, mentions missing components without clearly stating what is missing</p>
ER3	Out-of-Scope	1	<p>0: Proposals remain within the scope of the task and are relevant</p> <p>-1: Suggestions refer to methods or analyses that go beyond the intended scope</p>
ER4	Missing Reference	1	<p>0: Alternative methods, etc. (if available) are supported with justification or references</p> <p>-1: Alternative methods, etc. (if available) are suggested without justification or references</p>
ER5	Contradiction	1	<p>0: Feedback is consistent in itself, no contradictory statements</p> <p>-1: Contradictory feedback, e.g. criticism and praise for the same methods</p>
Rubric: Additional points (AD)			
<i>Allow additional points for very good reviews or deduct points for major mistakes.</i>			
AD1	Additional points	1	<p>0: No additional strengths identified.</p> <p>1: Notable additional strengths beyond other criteria.</p>
AD2	Negative points	1	<p>0: No additional major issues identified.</p> <p>-1: Significant issues beyond other criteria.</p>

Table 7: 24 criteria used to grade student reviews. We organize the criteria into eight rubrics, by grouping together criteria that evaluate related aspects of review.

E Summary of Collected Sources for Exposé Criteria

Across the full corpus of 110 exposé guidelines we reviewed, we identified materials from universities and higher-education institutions in Germany, Austria, and Switzerland, complemented by student-facing writing portals (e.g., Scribbr, Finito24, Studieren.de), which we further investigated. In terms of disciplinary scope, our primary sources are published by faculties and units in humanities and education; the social sciences (including sociology and political science); medicine; economics and business; and applied/STEM fields (notably architecture, computer science, and civil engineering). Roughly one quarter of the sources are published in English, primarily from internationalized programs (e.g., TU Berlin, Universität Wien, TU Dresden, Uni Magdeburg), while the remainder appear in German and serve national undergraduate or graduate audiences. About one-fifth of the documents explicitly target Bachelor- or Master-level theses, and a smaller subset (around ten) focus on doctoral or dissertation exposés, such as those from the University of Leipzig Graduate Centre, Uni Magdeburg Medical Faculty, and Uni Freiburg’s Graduate School. Only a few pages, most notably TU Berlin’s and Kassel’s departmental guides, explicitly include project or grant proposal contexts.

Thematically, nearly all guidelines converge on the same structural expectations: a clear *Motivation* or problem statement, a theoretically informed *Approach* or methodology, a feasible *Schedule* or timeline, and a complete *Bibliography*. Additional sections such as *Research Questions*, *Related Work*, and *Formal Aspects* (format, citation, or layout) occur frequently, especially in English-language handbooks and institutional templates.

Across this heterogeneity, the documents show remarkable consistency in pedagogical intent: *to ensure students articulate a focused question, justify methodological choices, and plan realistic execution*. These recurrent patterns directly informed the rubric domains of *Motivation*, *Approach*, *Schedule*, and *Bibliography* used for annotation and model benchmarking in our dataset.

Links for Table 8:

1. <https://www.tu.berlin/en/ak/study-and-teaching/final-theses/writing-an-expose>
2. <https://www.oei.fu-berlin.de/soziologie/>

studiumlehre/handreichungen_zum_studium/How-to-write-an-expose.pdf	1103
3. https://home.uni-leipzig.de/~gsgas/fileadmin/Recommendations/Expose_Recommendations.pdf	1105
4. https://expae.med.ovgu.de/unimagdeburg_mm/Downloads/Kliniken/KPAE/EXPAE/Guidelines+for+writing+an+expos%C3%A9+for+a+doctoral+project.pdf	1108
5. https://www.wvf.uni-freiburg.de/promhabil/merkblatt-erstellung-expose	1112
6. https://www.h2.de/fileadmin/user_upload/Promotionszentren/PromZ_SGW/4eng_Expose-Hinweise_PromZ_SGW.pdf	1114
7. https://www.archland.uni-hannover.de/fileadmin/archland/FORSCHUNG/promotionsausbildung/Requirements_admission_and_expose_.pdf	1117
8. https://tu-dresden.de/ing/informatik/smt/im/studium/thesen-and-research-projects/writing-an-expose	1121
9. https://www.hfp.tum.de/fileadmin/w00cjd/governance/pdf/Instructions_for_Preparing_an_Expose.pdf	1124
10. https://www.gsi.uni-muenchen.de/personen/wiss_mitarbeiter/kruck/general_guidelines.pdf	1127
11. https://sts.univie.ac.at/fileadmin/user_upload/i_sts/Studium/Master_STS/05_Services_for_current_students/Master_Thesis/Guidelines_for_Writing_a_Master_Thesis_Expose.pdf	1130
12. https://www.uni-kassel.de/fb11agrar/index.php?eID=dumpFile&t=f&f=572&token=eb5a0183632b6762477026bbd8f1694660f1b3f6	1135
13. https://www.bwl.uni-mannheim.de/media/Lehrstuehle/bwl/Schoen/Teaching/Master/Expose_Guide/Master_Thesis_Expose_Guide_v6.pdf	1138
14. https://www.fh-erfurt.de/fileadmin/Bilder/Seitenbereiche/Fakultaeten/BKR/B/Studium/MA_Sustainable_Engineering_of_Infrastructure/Expose_Recommendations_SEI.pdf	1142
15. https://www.scribbr.de/anfang-abschlussarbeit/expose-bachelorarbeit/	1146
16. https://www.wu.ac.at/fileadmin/wu/d/i/ign/IGN_BA-Expos%C3%A9_Requirements_De.pdf	1149
17. https://www.finito24.de/expose-schreiben/	1151
18. https://www.uni-bremen.de/fileadmin/user_upload/fachbereiche/fb7/gscm/Dokumente/Structure_of_an_Expose.pdf	1152

F Exposita Dataset Statistics

This section provides additional descriptive statistics for Exposita that complement the summary in Section 5 of the main paper. To give a more detailed view of the dataset’s composition, we report basic statistics on the exposés, reviews, and assessment scores.

Source	Lang	BA	MA	Diss.	Grant/Proj.	Link
TU Berlin	EN	✓	✓	✓	✓	[1]
FU Berlin	EN	✗	✓	✓	✗	[2]
Uni Leipzig	EN	✗	✗	✓	✗	[3]
Uni Magdeburg	EN	✗	✗	✓	✗	[4]
Uni Freiburg	DE/EN	✗	✗	✓	✗	[5]
Hochschule Magdeburg	EN	✗	✗	✓	✗	[6]
Uni Hannover	EN	✗	✗	✓	✗	[7]
TU Dresden (Informatik)	EN	✓	✓	✗	✗	[8]
Hochschule für Politik München	EN	✓	✓	✗	✗	[9]
LMU München	EN	✓	✓	✗	✗	[10]
Universität Wien (STS)	EN	✗	✓	✗	✗	[11]
Uni Kassel (Agrar & Food Marketing)	EN	✗	✓	✗	✗	[12]
Uni Mannheim (BWL)	EN	✗	✓	✗	✗	[13]
FH Erfurt	EN	✗	✓	✗	✗	[14]
Scribbr (DE)	DE	✓	✗	✗	✗	[15]
WU Wien	DE	✓	✗	✗	✗	[16]
Finito24	DE	✓	✗	✗	✗	[17]
Uni Bremen (Wirtschaftswiss.)	EN	✗	✗	✗	✗	[18]

Table 8: University and institutional exposé guidelines with explicit scope for Bachelor, Master, Dissertation, or Grant/Project proposals.

F.1 Exposé Revision Statistics

Table 9 reports section-wise text statistics for *Motivation*, *Approach*, and *Schedule* in draft vs. final versions (mean \pm SD). Raw \LaTeX is converted to plain text with `pylatexenc`¹³; sentence/word segmentation uses `spaCy`. We report **length/structure** (#Words, #Sentences and average word length), **lexical/readability** (type–token ratio (TTR) and Flesch Reading Ease (FRE); [Flesch, 1948](#) via `textstat`¹⁴), and **draft–final similarity** (cosine similarity of SentenceTransformer embeddings¹⁵ and normalized token-level Levenshtein similarity; [Levenshtein, 1965](#)).

Across sections, final exposés are longer in words and sentences, with the largest expansion in *Approach*; average word length changes only marginally. Lexical diversity (TTR) decreases slightly from draft to final in all three sections. The negative FRE values for *Schedule* likely reflect formatting artifacts (fragmentary, table-like content) interacting with plain-text conversion and

segmentation.

Similarity and topic structure High embedding cosine similarities in Table 9 indicate strong semantic continuity between draft and final, even where token-level edit similarity is lower (notably for *Motivation* and *Approach*). Table 10 reports pairwise embedding cosine similarities within and across topics: within-topic pairs are generally more similar than cross-topic pairs, with the largest within-topic variance for *Generative AI in Healthcare*.

Bibliography size and stability Table 11 reports topic-wise bibliography statistics for draft vs. final exposés (mean \pm SD), measured as the number of distinct Bib \TeX entry keys and their draft–final overlap (Jaccard similarity over key sets). Draft bibliography size differs substantially by topic ($T1 < T2 < T3$). Despite these initial gaps, revisions show a consistent pattern across topics: authors add about two references on average while removing fewer than one, suggesting modest expansion rather than wholesale replacement. Correspondingly, mean draft–final Jaccard similarity remains relatively high (overall 0.78), indicating that bibliographies are generally stable across revi-

¹³<https://github.com/phfaist/pylatexenc>

¹⁴<https://github.com/textstat/textstat>

¹⁵`sentence-transformers/all-mpnet-base-v2`

Metric	Motivation		Approach		Schedule	
	Draft	Final	Draft	Final	Draft	Final
Length and structure						
#Words	402 (± 133)	428 (± 125)	748 (± 258)	851 (± 221)	248 (± 83)	270 (± 131)
#Sentences	21.8 (± 8.5)	24.0 (± 7.8)	40.9 (± 18.4)	47.7 (± 15.3)	7.3 (± 8.3)	8.4 (± 8.7)
Avg. word length	5.38 (± 0.45)	5.43 (± 0.44)	5.41 (± 0.49)	5.46 (± 0.50)	5.55 (± 0.37)	5.50 (± 0.36)
Lexical and readability						
Lexical diversity	0.54 (± 0.06)	0.53 (± 0.06)	0.45 (± 0.07)	0.43 (± 0.06)	0.60 (± 0.06)	0.59 (± 0.06)
Flesch reading ease	43.0 (± 14.5)	41.6 (± 13.3)	42.1 (± 15.3)	41.0 (± 15.9)	-50.4 (± 90.4)	-34.2 (± 90.1)
Draft–Final similarity						
Cosine similarity	0.96 (± 0.04)		0.95 (± 0.06)		0.93 (± 0.14)	
Levenshtein similarity	0.64 (± 0.24)		0.62 (± 0.26)		0.81 (± 0.24)	

Table 9: Chapter-level textual statistics for exposés by version (draft vs. final). Entries show mean (\pm SD) for length, lexical diversity, readability, and draft–final similarity metrics.

Topic	Motivation		Approach		Schedule	
	Draft	Final	Draft	Final	Draft	Final
T1	0.59 (± 0.13)	0.58 (± 0.13)	0.55 (± 0.14)	0.56 (± 0.14)	0.52 (± 0.16)	0.52 (± 0.16)
T2	0.60 (± 0.12)	0.66 (± 0.09)	0.60 (± 0.11)	0.62 (± 0.12)	0.59 (± 0.13)	0.61 (± 0.12)
T3	0.48 (± 0.14)	0.47 (± 0.14)	0.40 (± 0.17)	0.41 (± 0.16)	0.54 (± 0.13)	0.50 (± 0.11)
Overall	0.37 (± 0.18)	0.38 (± 0.18)	0.35 (± 0.18)	0.35 (± 0.19)	0.48 (± 0.16)	0.46 (± 0.15)

Table 10: Pairwise cosine similarity between exposés (each exposé is compared with all others within the same topic and version). Entries show mean (\pm SD) over all unordered document pairs, computed separately per chapter. Topic mapping: T1=*The Holographic AI Assistant*, T2=*Interactive LLM-based teaching agents*, T3=*Generative AI in Healthcare*.

sions with targeted additions.

Draft–final improvements on rubric scores Table 12 reports per-rubric statistics for draft and final exposés. It complements Section 5.1 by additionally showing (i) the remaining gap between the mean final score and the rubric maximum and (ii) the observed score ranges (min–max) for each version. As in the main text (Section 5.1), the largest mean gains occur in conceptually central rubrics such as *Approach* and *Motivation*. In contrast, smaller changes in *Metadata*, *Bibliography*, and *Form/Structure* largely reflect limited headroom: these rubrics have low maxima, and draft scores are already close to the ceiling. Finally, *Additional points* changes little because it represents bounded instructor adjustments rather than points that students can systematically earn through revision (see Appendix F.3 for the corresponding scoring rules). Overall, the revision signal is therefore most visible in rubrics with substantial remaining headroom,

and is concentrated in substantive content and argumentation.

F.2 Review Statistics

Table 13 reports text statistics for review documents by author role (mean \pm SD), computed analogously to the exposé statistics. Student reviews are substantially longer than instructor reviews, while instructors show higher lexical diversity.

Instructor roles We distinguish between *junior* and *senior* instructors based on their academic stage and experience in reviewing and research. Junior instructors are recruited Master’s students who have completed a Bachelor’s thesis and thus have basic academic writing experience, whereas senior instructors are course staff (PhD researchers and a postdoc) with substantially more experience in academic peer review and supervision; see Appendix G.1 for details on rater composition and training.

Topic	Draft	\bar{n}_{final}	Final		Similarity
	\bar{n}_{draft}		\bar{n}_{added}	\bar{n}_{removed}	Jaccard
T1	6.74 ± 6.02	8.63 ± 5.63	2.16 ± 3.86	0.26 ± 0.56	0.73 ± 0.40
T2	9.22 ± 6.39	10.78 ± 6.31	2.06 ± 2.62	0.50 ± 1.25	0.76 ± 0.27
T3	11.83 ± 8.01	13.06 ± 8.74	1.89 ± 2.97	0.67 ± 2.83	0.84 ± 0.26
Overall	9.22 ± 7.04	10.78 ± 7.11	2.04 ± 3.15	0.47 ± 1.77	0.78 ± 0.32

Table 11: Bibliography size and draft–final changes by topic. We report mean (\pm SD) number of entries in draft and final versions, average additions/removals, and mean draft–final entry-set Jaccard similarity computed over Bib_T_X keys. Topic mapping: T1=*The Holographic AI Assistant*, T2=*Interactive LLM-based Teaching Agents*, T3=*Generative AI in Healthcare*.

Rubric	Mean (\pm SD)			Gaps		Observed bounds	
	Draft	Final	Δ	Gap	Max	Draft	Final
Approach	7.91 (\pm 2.25)	9.40 (\pm 2.06)	1.49 (\pm 1.74)	1.60	11.00	1.00–11.00	1.00–11.00
Motivation	6.93 (\pm 1.46)	8.09 (\pm 1.01)	1.16 (\pm 1.17)	0.91	9.00	4.00–9.00	6.00–9.00
Language/Quality	2.69 (\pm 1.02)	3.45 (\pm 0.78)	0.76 (\pm 0.81)	0.55	4.00	0.00–4.00	1.00–4.00
Methodology	4.07 (\pm 1.32)	4.76 (\pm 0.57)	0.69 (\pm 1.08)	0.24	5.00	0.00–5.00	3.00–5.00
Schedule	4.15 (\pm 0.84)	4.65 (\pm 0.55)	0.51 (\pm 0.83)	0.35	5.00	0.00–5.00	3.00–5.00
Bibliography	2.56 (\pm 0.85)	2.91 (\pm 0.29)	0.35 (\pm 0.77)	0.09	3.00	0.00–3.00	2.00–3.00
Metadata	2.75 (\pm 0.51)	2.95 (\pm 0.23)	0.20 (\pm 0.48)	0.05	3.00	1.00–3.00	2.00–3.00
Additional points	-0.02 (\pm 0.36)	0.00 (\pm 0.27)	0.02 (\pm 0.30)	2.00	2.00	-2.00–1.00	-1.00–1.00
Form/Structure	1.96 (\pm 0.19)	1.96 (\pm 0.19)	0.00 (\pm 0.19)	0.04	2.00	1.00–2.00	1.00–2.00

Table 12: Per-rubric statistics for draft and final exposé scores. We report mean (\pm SD), average remaining gap to the rubric maximum, and observed score bounds (min–max) for each version.

Inline comment distributions Table 14 summarizes comment volume and length by role: junior instructors contribute the most comments overall, while students and senior instructors write longer comments on average. Table 15 reports the distribution of span-level tags. Totals are similar but not identical across the two tables because span annotations are only recorded for comments that are attached to a specific text span; page-level comments have no associated span annotation. Across roles, *Weakness* is the most frequent tag, followed by *Strength*, *Other*, and *Highlight*.

F.3 Assessment scores coverage

Table 16 summarizes the number of assessments that include a complete set of criterion-based scores for exposé drafts and final submissions, by rater group. Draft exposés were graded by instructors during the semester (Group 1) and re-scored after the course by Group 2 to compute inter-annotator agreement (Section 5.2). Final exposés, in contrast, were not re-scored post-semester: due to the substantial time required for criterion-based grading, Group 2 focused on a draft subset rather than du-

plicating the full set of final-submission grades. A similar pattern holds for student review grading (Table 17): reviews were graded during the semester and selectively re-scored for agreement analyses, but criterion-based scores are only available for student-written reviews (not for instructor-written reviews).

Annotation effort We emphasize that constructing this dataset is unusually labor-intensive. In Group 1, producing a full review with inline comments and criterion-based scores typically required on the order of 2–4 hours early in the semester, decreasing to roughly 1–2 hours as raters gained experience. Criterion-based grading itself is also high effort: grading a single exposé typically takes about 1–2 hours, while grading a single review requires 20–40 minutes. Despite these restrictions, the dataset currently comprises 55 exposés and 306 reviews and includes unusually rich annotation artifacts: inline comments and comment tags, free-text reviews, draft and final exposé versions, and hundreds of instructor assessment scores across 36 exposé and 24 review criteria, which we

Metric	Review author role		
	Student	Junior Instructor	Senior Instructor
Length and structure			
#Words	680 (± 455)	285 (± 134)	347 (± 116)
#Sentences	33.7 (± 21.3)	16.7 (± 8.4)	18.7 (± 6.3)
Avg. word length	5.31 (± 0.44)	5.16 (± 0.28)	5.18 (± 0.49)
Lexical and readability			
Lexical diversity	0.46 (± 0.08)	0.56 (± 0.09)	0.55 (± 0.08)
Flesch reading ease	50.2 (± 13.3)	58.9 (± 9.9)	54.3 (± 14.2)

Table 13: Textual statistics of review texts by author role. Entries show mean (\pm SD).

Role	Comments	Avg chars	Avg words
St. (n=31)	640	77.0	12.6
Jr. (n=7)	1121	53.8	9.1
Sr. (n=5)	492	68.8	12.0
Σ (n=43)	2253	63.6	10.7

Table 14: Basic statistics of comments by role, including total counts and average length in characters and words. St.=Students, Jr.=Junior Instructors, Sr.=Senior Instructors.

use throughout our analyses.

Scoring and special rubric behavior Within each submission (draft exposé, final exposé, peer feedback), criterion-level scores are summed to yield rubric points, and rubric points are summed to obtain total submission points. Most rubrics use simple addition; we introduce a small number of rubric-specific rules to maintain consistency across topics and feedback situations. The exposé rubric *Methodology* is capped at 5 points, and the review rubric *Feedback* is capped at 4 points. For example, in the exposé rubric, *Existing Material* may be inapplicable when the proposed methodology relies on newly collected primary data rather than reusing an existing dataset; the cap ensures that full credit remains attainable by satisfying an appropriate subset of criteria. The review rubric *Errors* employs a deduction-only scheme, starting at 4 points and decreasing by 1 point for each identified error, with a minimum score of 0. Finally, the instructor-only rubric *Additional points* permits bounded adjustments in exceptional cases (Exposé: $[-2, +2]$, Review: $[-1, +1]$); these adjustments do not increase the maximum attainable points.

G Inter-annotator agreement

This section documents the computation of inter-annotator agreement (IAA), complementing the criterion-based scoring description in Section 5.2 for Exposita. It also provides the reference point for interpreting LLM–human agreement in the main experiments. Because our criterion scores are discrete and ordinal, we prioritize agreement measures that account for ordered categories and report distributional diagnostics for criteria affected by prevalence skew.

G.1 Raters, training, and scoring design

Raters and training Twelve raters participated and were organized into two groups. In total, 7 were *junior instructors* (student helpers) and 5 were *senior instructors* (4 PhD researchers and 1 postdoc). One PhD researcher participated in both groups.

Group 1 (in-semester): 6 student helpers, 2 PhD researchers, 1 postdoc; weekly calibration (jointly graded initial exposés, then addressed edge cases). Because grades were at stake, Group 1 was instructed to resolve uncertainty in favor of the student.

Group 2 (post-semester): 1 student helper, 3 PhD researchers; scored subsets of draft exposés and reviews using a recorded calibration video prepared by a Group-1 senior instructor.

Rater groups and coverage Table 16 indicates that Group 1 grades are predominantly assigned by junior instructors, whereas Group 2 grades are predominantly assigned by senior instructors. This shift mainly reflects post-semester availability: once the course ended, much of the teaching staff was no longer available, so re-grading was carried out primarily by senior instructors; final

Tag	Metric	St. (N=31)	Jr. (N=7)	Sr. (N=5)	Σ
Highlight	Raw count	90	28	58	176
	Avg. length	112.68 (± 188.62)	122.79 (± 149.17)	168.34 (± 227.29)	132.63 (± 198.48)
Other	Raw count	107	145	52	304
	Avg. length	120.66 (± 236.88)	41.77 (± 61.91)	90.83 (± 132.55)	77.93 (± 160.87)
Strength	Raw count	131	200	73	404
	Avg. length	254.63 (± 263.11)	178.40 (± 239.87)	196.52 (± 285.21)	206.39 (± 258.44)
Weakness	Raw count	308	691	290	1289
	Avg. length	79.20 (± 106.17)	59.17 (± 127.02)	97.88 (± 154.42)	72.66 (± 130.22)
Σ	Raw count	636	1064	473	2173

Table 15: Distribution of span-level annotation tags by role. We report raw counts and average span length (\pm SD). St.=Students, Jr.=Junior Instructors, Sr.=Senior Instructors.

Version / Topic	Group 1			Group 2		
	Jr.	Sr.	Σ	Jr.	Sr.	Σ
Draft (N = 55)						
T1 (N=19)	14	5	19	2	17	19
T2 (N=18)	14	4	18	4	14	18
T3 (N=18)	15	3	18	2	16	18
Σ Draft	43	12	55	8	47	55
Final (N = 55)						
T1 (N=19)	15	4	19	0	0	0
T2 (N=18)	14	4	18	0	0	0
T3 (N=18)	15	3	18	0	0	0
Σ Final	44	11	55	0	0	0

Table 16: Distribution of criterion-based exposé assessments across topics (T1–T3), rater roles (Jr. Junior Instructors, Sr. Senior Instructors), rater groups (Group 1 in-semester, Group 2 post-semester), and versions (draft vs. final). T1: “The Holographic AI Assistant”; T2: “Interactive LLM-based teaching agents”; T3: “Generative AI in Healthcare”.

1352 submissions were graded only during the semester.

1353 **Scope of IAA** Our agreement analyses are necessarily restricted to items (exposés and reviews) 1354 that have assessments from two instructors. For 1355 exposés, we report agreement only for *draft* sub- 1356 missions, as final submissions were not scored post- 1357 semester due to the substantial time required for 1358 criterion-based assessment (Appendix F.3). For 1359 reviews, we focus on the subset that has criterion- 1360 based grades from Group 1 (student reviews). Fi- 1361 nally, for a subset of reviews ($n = 60$), we col- 1362 lected duplicate retrospective ratings from Group 2, 1363

Reviews	Group 1			Group 2		
	Jr.	Sr.	Σ	Jr.	Sr.	Σ
Student reviews (N = 110)						
E (N = 34)	28	6	34	0	59	59
$\neg E$ (N = 76)	67	9	76	0	111	111
Σ	95	15	110	0	170	170

Table 17: Distribution of criterion-based student review assessments by instructor role (Jr. Junior Instructors, Sr. Senior Instructors) and rater group (Group 1 vs. Group 2), stratified by whether the reviewed author has an associated exposé (E) or not ($\neg E$).

enabling a within-group agreement comparison in 1364 addition to cross-group agreement (Table 21). Be- 1365 yond criterion-level results, we also report aggre- 1366 gate agreement (i) by rubric domain and (ii) by 1367 criterion expertise level (see Appendix D). 1368

1369 G.2 Human–human agreement

1370 **Agreement reporting** We report raw quadratic 1371 weighted agreement (QWA), Krippendorff’s α , 1372 Pearson correlation coefficient r , and simple mean 1373 difference (“mean bias”) between groups. **QWA** is 1374 the observed agreement between two raters on an 1375 ordinal scale, where disagreements are weighted 1376 by squared distance so that larger score differ- 1377 ences are penalized more strongly than adjacent- 1378 category mismatches. Let x_i, y_i be two raters’ 1379 scores for item i on an ordinal scale with dis- 1380 crete domain $\mathcal{D} = \{d_{\min}, \dots, d_{\max}\}$ and span 1381 $\Delta = d_{\max} - d_{\min} > 0$. We use the quadratic

Reviews per exposé	Group 1					Group 2			
	N	Jr.	Sr.	St.	Σ	N	Jr.	Sr.	Σ
Exposés (total N = 55)									
0 review(s)	0	0	0	0	0	30	0	0	0
1 review(s)	23	17	6	0	23	24	4	20	24
2 review(s)	23	21	7	18	46	1	0	2	2
3 review(s)	9	8	3	16	27	0	0	0	0
Σ totals	55	46	16	34	96	55	4	22	26
Orphaned reviews (total N = 184)									
Σ totals	184	88	20	76	184	0	0	0	0

Table 18: Distribution of reviews per exposé by review group (Group 1 vs. Group 2). In the Exposés block, N is the number of exposés in the bucket; role columns (Jr./Sr./St.) are counts of reviews written by that role (columns that are all-zero for a group are omitted); and Σ is the total number of reviews. Orphaned reviews are reported separately.

weight

$$w(x_i, y_i) = \max\left(0, 1 - \left(\frac{x_i - y_i}{\Delta}\right)^2\right),$$

and define raw quadratic weighted agreement as

$$\text{QWA} = \frac{1}{N} \sum_{i=1}^N w(x_i, y_i),$$

where N is the number of co-rated items.

We use quadratic weighting because it preserves the magnitude of disagreement and remains interpretable under strong prevalence skew, a setting in which chance-corrected coefficients (e.g., κ or α) can be misleading. Accordingly, for highly skewed criteria (e.g., *Number of pages*, *Template used*, *Negative points*), we emphasize QWA.

α values range between -1 and 1, with 1 being perfect agreement, 0 being agreement by chance, and negative values indicating systematic disagreement; r measures the correlation between a system’s predicted scores and the annotator-assigned scores; it ranges from -1 to 1. A positive (negative) value implies that the two sets of predictions are positively (negatively) correlated. Inter-annotator agreement results for exposé criteria are reported in Table 19, while agreement for review criteria is summarized in Tables 20 and 21.

G.3 Failure modes of chance-corrected agreement

This section interprets the agreement patterns observed in the exposé and review analyses (see Section 5.2 in the main paper). In particular, we focus on cases where chance-corrected reliability or correlation is negative, undefined, or unexpectedly low, despite high raw agreement, and relate these outcomes to prevalence effects, calibration differences between rater groups, and ambiguities in rubrics or training.

Prevalence effects under extreme class imbalance (“rare-event” criteria)

For highly skewed binary (or near-binary) criteria, the expected agreement under a chance model is already very high, so a handful of rare, non-overlapping deviations can yield $\alpha < 0$ even when raw agreement is high. This explains the slightly negative values for *Template used* (ST1) and *Negative points* (AD2), where almost all submissions receive the default label, but the few exceptional ratings occur on different submissions. In such settings, Pearson’s r is typically close to zero because the rating vectors are nearly constant, and mean bias $\Delta\mu$ helps diagnose whether the deviations are symmetric (ST1: $\Delta\mu \approx 0$) or slightly one-sided (AD2: small negative $\Delta\mu$ consistent with marginally more penalizing scores in one group).

Annotation guideline/interpretation issues (training or rubric ambiguity)

Some negative/undefined reliability values are attributable to concrete, identifiable rubric failures. For *Metadata available* (ME2), metadata was present for all submissions; all disagreements originated from a single rater in Group 2, indicating a training/interpretation error (i.e., the rater did not apply the intended check). This also leads to an undefined Pearson correlation because one group’s ratings are effectively constant (zero variance). For *Template used* (ST1), the few “0” assignments were triggered by superficial deviations (e.g., renamed chapter headings or bibliography compilation issues) despite the underlying template being used, revealing an unclear rubric description about what counts as a template violation. For *Methodology Availability* (AP8), we found edge cases where the methodology is present only as a very short list; these items were inconsistently judged as available vs. unavailable across groups, suggesting that the criterion definition is underspecified (i.e., availability vs. sufficiency/quality).

Systematic calibration shifts and low signal in subjective criteria

For *Anker* (MO2), the negative (near-zero) α is consistent with a systematic calibration shift between groups: Group 2 assigns lower scores on average (positive mean bias $\Delta\mu = 0.309$), and disagreements are predominantly adjacent-category mismatches (e.g., 2 vs. 1), which remain relatively benign under quadratic weighting but still depress chance-corrected reliability. For *Problem relevance* (MO4), negative coefficients are driven by sparse and subjective “0” decisions: the negative label is rare and tends to be applied to different submissions across groups.

Prevalence effects under extreme class imbalance in review scoring

A closely related prevalence-driven failure mode appears in review scoring for error and penalty criteria (Tables 20 and 21). These criteria are effectively *rare-event detection* tasks on a highly skewed label space: the non-penalizing default (typically 0) dominates, while penalizing assignments (e.g., -1 in *Errors (ER)*) occur only rarely. As a result, chance-corrected reliability and Pearson correlation become unstable or uninformative. A small number of non-overlapping penalizations can yield $\alpha \leq 0$ despite high raw agreement, and near-constant rating vectors drive r toward zero or render it undefined when one rater shows zero variance.

Overall, we emphasize quadratic-weighted raw agreement as our primary agreement signal and use α , r , $\Delta\mu$, and the pooled distributions to (i) flag prevalence-induced artifacts, (ii) diagnose calibration differences, and (iii) identify rubric/training issues that can be addressed in the future.

H Experimental Details

Models We evaluate a diverse set of recent open-weight instruction-tuned LLMs that cover different architectures and scales. Our selection includes (i) a large dense decoder-only Transformer from the Llama 3.3 family (meta-llama/Llama-3.3-70B-Instruct; 70B parameters) (Grattafiori et al., 2024), (ii) a high-sparsity Mixture-of-Experts (MoE) model with long-context support (Qwen/Qwen3-Next-80B-A3B-Instruct; 80B total parameters with 3B active) (Team, 2025) and an independent MoE reasoning model in two flavors, (iii) the high reasoning model (openai/gpt-oss-120b, 117B parameters with 5.1B active parameters) and (iv) the smaller low latency variant (openai/gpt-oss-20b, 21B parameters with 3.6B active parameters) (OpenAI, 2025).

All models are fully open-weight and can be run on two local NVIDIA A100 80GB GPUs, which is important for classroom deployment where data protection and limited consent often require on-premise inference instead of cloud-based APIs. Inference is performed using the Hugging Face TRANSFORMERS library¹⁶ with deterministic greedy decoding (`do_sample = False`).

Prompts For both tasks, we use a unified prompting strategy (Appendix I) with a system message that (i) assigns the model the role of an expert evaluator and (ii) enumerates the relevant criteria, including their names, short descriptions, allowed integer score ranges, and the verbal descriptions associated with each score level (Tables 6 and 7). The model is required to return a single JSON object with fields `criterion_name`, `assigned_score`, and a short `justification`. The additional justifications are requested to explain the decision-making process, as it is crucial in an environment where people are affected by the LLM outcome (Kolářová and Schmude, 2026).

Criteria excluded from the LLM experiments

The exposé and review rubrics contain structural or format-related checks (e.g., LaTeX template usage, page limits) and content-related criteria. We exclude the first one, as the LLM cannot successfully assess it because the generated PDF is needed for correct evaluation. We also exclude additional points that can be given by the human evaluator, as these occur only in isolated cases and are intended for exceptional circumstances. The following exposé criteria are excluded: Metadata available; Template used; Number of pages; Additional points; Negative points. The following review criteria are excluded: Additional points; Negative points.

I Prompt Templates

This appendix provides the exact prompt formats used in all baseline experiments. The rubric text included in the templates is automatically generated from the released JSON rubric files, ensuring that the prompts align with the annotation scheme used in Exposita.

For both tasks, we use a unified prompt structure. The system message (i) assigns the model the role of an expert evaluator and (ii) enumerates the relevant criteria, including their names, short descriptions, allowed integer score ranges, and the

¹⁶<https://huggingface.co/docs/transformers>

1556
1557
1558
1559
1560
1561
1562
1563
1564

verbal descriptions associated with each score level (Tables 6 and 7). The model is required to return a JSON *list* with one object per criterion, containing the fields `criterion_name`, `assigned_score`, and a short justification. The additional justifications are requested to explain the decision-making process, as it is crucial in an environment where people are affected by the LLM outcome (Kolářová and Schmude, 2026).

I.1 Expose Scoring Prompt

The user message provides the topic, the full exposé text, and the bibliography.

1565
1566
1567

System message

You are an expert evaluator specializing in assessing student exposes. Your task is to carefully evaluate the given expose based on the specified grading criteria.

The grading criteria are as follows:

<Rubric Name>
<Rubric description>

Criterion: <Criterion Name>
Description: <Criterion description>
Allowed scores:
- <Points> points: <Description>
...

<Next Rubric>
...

- You must:
1. Evaluate the expose **one** criterion by criterion.
 2. For **every** listed criterion, assign exactly **one** score using **only** the scoring options and point values provided.
 3. Justify each score with a clear, concise and **short** explanation.

Your response **MUST** be a list of dictionaries in **valid JSON** schema:

```
[
  {
    "criterion_name": "<the EXACT criterion name from the criteria list, without any other words>",
    "assigned_score": "<an INTEGER, one of the allowed point values for this criterion>",
    "justification": "<SHORT justification of why you decided to give this score>"
  },
  ...
]
```

Constraints:

1568

- Return **one** object per criterion; do not skip any criteria.
- Use only the criteria and scoring options that are explicitly provided.
- Return a response that consists solely of valid JSON, with no additional explanatory text before or after it.
- Make sure the JSON is syntactically valid (no trailing commas, correct quotation marks, etc.).

1569

User message

The topic is: "<TOPIC>".

Here is the expose you have to grade:

<EXPOSE TEXT>

Bibliography:
<BIBLIOGRAPHY>

1570

I.2 Review Scoring Prompt

The user message provides the full written review text.

1571

1572

1573

System message

You are an expert evaluator specializing in assessing the quality of student peer reviews.

Your task is to carefully evaluate the reviewer's performance based on the specified grading criteria.

Grading criteria:

<Rubric Name>
<Rubric description>

Criterion: <Criterion Name>
Description: <Criterion description>
Allowed scores:
- <Points> points: <Description>
...

<Next Rubric>
...

- You must:
1. Evaluate the peer review **one** criterion by criterion.
 2. For **every** listed criterion, assign exactly **one** score using **only** the scoring options and point values provided for that criterion.
 3. Base your evaluation strictly on the quality of the full review text.
 4. Justify each score with a clear, concise and **short** explanation.

Output format:

1574

Return a list of objects following this **valid JSON** schema:

```
[
  {
    "criterion_name": "<EXACT criterion name from the criteria list>",
    "assigned_score": <INTEGER, one of the allowed point values for this criterion>,
    "justification": "<SHORT justification of why you decided to give this score>"
  },
  ...
]
```

Constraints:

- Return **one object per criterion**; do not skip any criteria.
- Use only the criteria and scoring options that are explicitly provided.
- The response must consist solely of valid JSON, with no additional explanatory text.
- Make sure the JSON is syntactically valid (no trailing commas, correct quotation marks, etc.).

User message

Here are the materials of the peer review you have to grade.

=== Full written text of the review ===
<FULL_REVIEW_TEXT>

J Additional Experimental Results

Agreement by expertise level Table 4 reports human-LLM agreement (QWA) for criterion-level scoring, grouped by criterion expertise. Across both tasks, agreement decreases monotonically with expertise, mirroring the human-human reference. For exposé scoring, low-expertise criteria (Exp. 0; three criteria) yield high agreement for the strongest models under combined prompting (e.g., Llama 3.3: $\bar{A}_{\text{all}}=0.90$; Qwen3: 0.91), whereas Exp. 2 criteria are substantially harder (Llama 3.3: 0.63; Qwen3: 0.57); humans show the same drop in the Human column (0.89 at Exp. 0 vs. 0.64 at Exp. 2), consistent with higher-expertise criteria being more subjective. Aggregated over all 31 exposé criteria, model agreement is close to the human baseline (Human: 0.76), and combined prompting is typically comparable or slightly better than single-criterion prompting, with Qwen3 as the main exception ($\bar{A}_{\text{ind}}=0.75$ vs. $\bar{A}_{\text{all}}=0.72$). The comparatively lower Exp. 0 values for GPT-OSS in

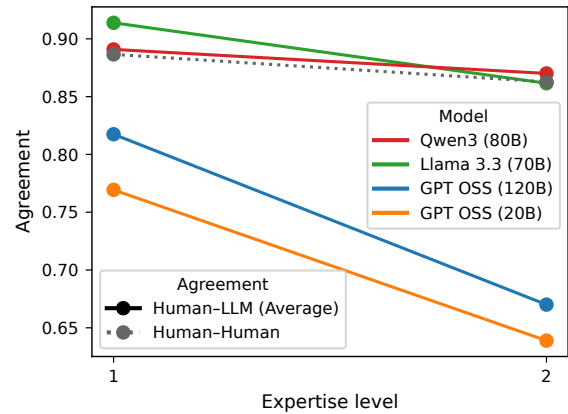


Figure 7: **Human-LLM agreement (QWA) for review scoring by expertise level.** The dotted line shows human-human QWA.

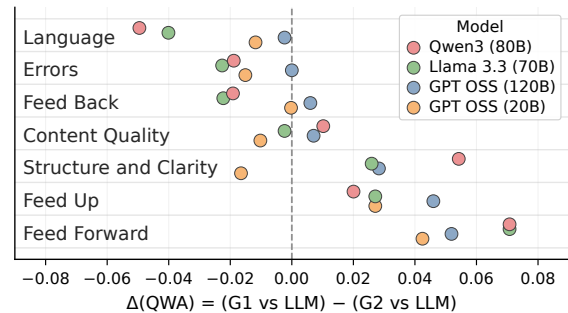


Figure 8: **Group-asymmetry in human-LLM agreement by review rubric.** Each point shows the difference in agreement between an LLM and human raters (Group 1 (G1) vs. Group 2 (G2)). The vertical line at $\Delta = 0$ indicates equal agreement of the LLM with human G1 and G2 raters. Points to the right indicate the model agrees more with G1 raters than with G2 raters.

the aggregate are explained by the small size of this bin and a single outlier criterion: Table 22 shows that GPT-OSS matches humans very well on the schedule checks (SC1/SC3; QWA \approx 0.91–0.98), but performs unusually poorly on BI1 (Bibliography Consistency; e.g., GPT-OSS 120B: $A_{G1}^{\text{all}}=0.27$, $A_{G2}^{\text{all}}=0.40$), which dominates the Exp. 0 average. For review scoring, overall agreement is higher (Human: 0.88 over 22 criteria), and Llama 3.3 and Qwen3 reach near-human alignment under combined prompting (Llama 3.3: $\bar{A}_{\text{all}}=0.90$; Qwen3: 0.89); the expertise effect persists, with Exp. 2 criteria yielding lower agreement than Exp. 1.

Scalability and classroom deployment Across all models and both tasks, combined prompting is more efficient in most cases than single-criterion prompting: it reduces mean runtime per submission by about $1.3\times$ – $5.3\times$ and avoids the linear cost

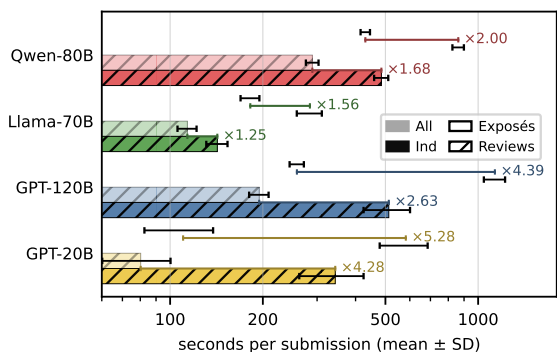


Figure 9: Runtime per submission for criterion-based assessments across models. Bars show mean seconds per submission with standard deviation bars. Within each task, **All** denotes combined prompting (light fill) and **Ind** denotes single-criterion prompting (dark fill). Review runs are hatched; Exposé runs are unhatched. Ratios show $\times(\text{Ind}/\text{All})$, highlighting the additional latency introduced by single-criterion prompting.

growth from issuing one model call per criterion (Figure 9). Single-criterion prompting also substantially increases total token processing because the full input is repeated across prompts, with particularly large overhead for long exposés and the review setting. These efficiency differences become especially consequential under the *observed course usage dynamics*. As shown by the behavioural interaction logs (Appendix Figure 6), student activity is strongly deadline-driven, with a pronounced spike in actions immediately before the review deadline, whereas instructor activity is relatively steady. This implies that in practice, LLM-based scoring must be provisioned for short periods of peak demand rather than average-day usage. As such, the overhead of single-criterion prompting is particularly problematic, while combined prompting better supports batching and throughput during deadline-induced bursts. Consequently, *combined prompting is the recommended operating point for classroom deployment under constrained on-premise resources, with model choice primarily determining the remaining latency-quality trade-off*. Runtime and token usage are reported in Appendix K (Tables 24–25).

Operational latency under peak-demand grading windows In classroom deployment under constrained on-premise resources, using LLM scores as instructor support (e.g., batch pre-scoring or rubric-aligned decision support) introduces a practical scheduling choice: scores can either be *pre-computed* ahead of grading, or compute ca-

capacity must be provisioned to serve requests *on demand* throughout the grading period. Under deadline-driven bursts (Appendix Figure 6), it can translate directly into waiting times, because even combined prompting incurs minute-level latency per submission and shows substantial min-max spread (Tables 24–25). To illustrate, consider a cohort of $N=600$ exposés scored on the same two-GPU setup used in our experiments. For exposés, the fastest combined configuration in Table 24 (GPT OSS 20B, All) yields $600 \times 110.2s \approx 18.4$ hours on average; larger models further increase this budget (e.g., Qwen3 80B, All: ≈ 71.7 hours). For student reviews, the volume is larger ($2N$ submissions); in the worst-case combined configuration, runtimes of $\approx 289.6s$ per review (Table 25) translate into ≈ 96.5 hours of wall-clock time, making peak-time live scoring particularly susceptible to queuing delays.

Finally, beyond latency, we observed occasional out-of-memory failures for unusually long inputs (e.g., exposés with many and very long comments), reinforcing the need for robust deployment to plan around worst-case input lengths (e.g., advance computation, batching, and input-length safeguards) rather than relying on unconstrained live scoring near deadlines.

Is review text enough for review assessment?

Overall, we did not find an indication that review scoring requires context beyond the review text. In supplementary experiments, we augmented the review input jointly with additional context from the exposé and inline comments, but did not observe consistent improvements over the review-only setting.

K Runtime and Token Consumption

In this appendix, we focus exclusively on inference-time cost (runtime and token usage) under these fixed conditions. Prompting variants (combined vs. single-criterion) are defined in Section 6.1. For each model/condition, we measure wall-clock runtime both per full submission (s/assessment) and per model call (s/prompt). Token counts are computed with the respective model tokenizer on the fully formatted prompt (system + user messages) and the generated JSON response, and reported as tokens/prompt split into tokens/input and tokens/output. All values are aggregated across evaluation instances and reported as mean (SD) and [min-max], capturing variation in input length.

LLM scoring throughput: combined prompting is essential Tables 24 and 25 report runtime and token usage for exposé and review scoring, respectively. Across both tasks, combined prompting is essential for throughput: it requires a single forward pass per submission and therefore yields the lowest latency (110–430s on average across models for exposé scoring). In contrast, single-criterion prompting multiplies calls per submission, increasing total runtime to 284–1133s on average for exposés, with worst cases up to 23 minutes. The same multiplicative effect appears in token processing: although single-criterion outputs are shorter per call, repeating the full input across criteria yields an order-of-magnitude increase in total processed tokens per submission (approximately 13× in our setup). Review scoring exhibits the same trend (Table 25), with combined prompting consistently yielding lower latency and substantially fewer processed tokens per submission than single-criterion prompting.

Limitations These behavioural patterns require *particular* caution because they reflect only *partial observation*. Not all students consented to behavioural logging, so the curves describe only the consenting subset and may not generalize to the full cohort. The instructor group is also small ($n=9$), limiting statistical power and making aggregates sensitive to individual work styles. Moreover, instructor activity is not directly comparable to student activity because instructors carried a substantially higher workload: while students completed only two reviews without grading, instructors collectively reviewed across the cohort and additionally assigned criterion-based grades. Finally, logs capture only on-platform actions; substantial reviewing work may occur offline (e.g., reading drafts externally or drafting feedback elsewhere). Thus, peaks and troughs should be interpreted as *platform activity* rather than complete measures of effort or time-on-task.

L Behaviour Data: Interaction Logs and Semester Dynamics

To complement the textual artefacts in Exposita, we also record fine-grained *behavioural interaction logs* from the review platform used in the course.

Temporal dynamics across the semester Figure 6 shows daily action counts for students and instructors. Two qualitative patterns stand out.

First, student activity is strongly deadline-driven: actions rise sharply in the days immediately preceding the review deadline (red dashed line). In contrast, instructor activity is comparatively smoother and distributed across the semester, consistent with instructional work that is scheduled and paced over time rather than concentrated in the final days.

Second, activity drops to near-zero during the winter break (hatched region), including the Christmas/New Year period, and then resumes afterward.

Implications for LLM-assisted assessment workflows The pronounced pre-deadline student spike (Figure 6) suggests that any LLM-based support integrated into the reviewing/assessment pipeline (e.g., rubric guidance, feedback scaffolding, or automated checks) will face *highly non-uniform demand*. Compute budgets and latency constraints should be evaluated under peak load conditions rather than on average-day usage.

Code	Criterion	Distribution (%)	QWA	α	r	$\Delta\mu$
LA1	Language quality	0 (15), 1 (55), 2 (30)	0.82	0.15	0.17	0.13
LA2	Common Thread	0 (4), 1 (65), 2 (32)	0.87	0.05	0.20	0.38
Average Language			0.85	0.10	0.19	0.25
ME1	Preliminary title	0 (3), 1 (42), 2 (55)	0.86	0.10	0.31	0.44
ME2*	Metadata available	0 (10), 1 (90)	0.80	-0.10	—	0.20
Average Metadata			0.83	0.00	0.31	0.32
ST1*	Template used	0 (2), 1 (98)	0.96	-0.01	-0.02	0.00
ST2*	Number of pages	0 (3), 1 (97)	0.98	0.66	0.70	0.02
Average Form/Structure			0.97	0.33	0.34	0.01
MO1	Hook existing	0 (15), 1 (85)	0.76	0.10	0.10	0.05
MO2	Anker	0 (5), 1 (63), 2 (33)	0.85	-0.00	0.06	0.31
MO3*	Domain	0 (5), 1 (95)	0.95	0.38	0.49	0.05
MO4	Problem relevance	0 (17), 1 (83)	0.65	-0.20	-0.21	-0.02
MO5	Problem handling	0 (35), 1 (65)	0.64	0.20	0.24	0.18
MO6	Teaser (RQ)	0 (7), 1 (56), 2 (36)	0.85	0.15	0.21	0.22
MO7	RQ Limitation	0 (46), 1 (54)	0.58	0.17	0.25	0.27
Average Motivation			0.76	0.11	0.16	0.15
AP1	SOTA	0 (13), 1 (87)	0.85	0.35	0.37	0.07
AP2	SOTA Relevance	0 (9), 1 (56), 2 (35)	0.86	0.28	0.39	0.36
AP3	SOTA Weaknesses	0 (23), 1 (55), 2 (22)	0.83	0.25	0.30	0.27
AP4	SOTA Delimitation	0 (30), 1 (70)	0.73	0.36	0.40	0.16
AP5	SOTA Combination	0 (34), 1 (66)	0.65	0.23	0.29	0.20
AP6	Theoretical Framework	0 (20), 1 (47), 2 (33)	0.78	0.15	0.20	0.29
AP7	Relevance of theoretical framework	0 (35), 1 (65)	0.64	0.20	0.20	0.00
AP8*	Methodology Availability	0 (9), 1 (91)	0.82	-0.09	-0.08	0.11
Average Approach			0.77	0.22	0.26	0.18
MD1	Methodology Completeness	0 (45), 1 (55)	0.69	0.38	0.46	0.24
MD2	Methodology Relevance	0 (26), 1 (74)	0.73	0.30	0.30	0.05
MD3	Methodology Target group	0 (25), 1 (75)	0.73	0.27	0.30	0.13
MD4	Methodology Existing Material	0 (35), 1 (65)	0.55	0.02	0.03	-0.13
MD5	Methodology Difficulties	0 (44), 1 (56)	0.56	0.12	0.19	0.25
MD6	Methodology Possibilities/restrictions	0 (56), 1 (44)	0.60	0.19	0.27	0.25
MD7	Methodology Details	0 (53), 1 (47)	0.49	-0.01	0.01	0.18
Average Methodology			0.62	0.18	0.22	0.14
SC1*	Schedule Availability	0 (4), 1 (96)	0.96	0.49	0.57	0.04
SC2*	Schedule Completeness	0 (4), 1 (96)	0.96	0.49	0.57	0.04
SC3*	Schedule Block Description	0 (6), 1 (94)	0.91	0.24	0.28	-0.05
SC4	Schedule Realistic Relevance	0 (7), 1 (71), 2 (22)	0.90	0.21	0.36	0.25
Average Schedule			0.93	0.36	0.44	0.07
BI1	Bibliography Consistency	0 (19), 1 (81)	0.80	0.36	0.49	0.20
BI2	Key literature	0 (15), 1 (35), 2 (49)	0.76	0.03	0.35	0.64
Average Bibliography			0.78	0.19	0.42	0.42
AD1*	Additional points	0 (97), 1 (2), 2 (1)	0.99	0.67	0.70	0.00
AD2*	Negative points	-2 (2), -1 (1), 0 (97)	0.96	-0.02	-0.03	-0.02
Average Additional points			0.98	0.32	0.34	-0.01
Average all criteria			0.79	0.20	0.27	0.16

Table 19: Quadratic weighted agreement (QWA) between Group 1 and Group 2 raters for each exposé criterion, with rubric-wise and overall averages. Columns show agreement (Agr), reliability (α), correlation (r), and mean bias ($\Delta\mu$). Rows marked with * have a highly skewed class distribution ($\max_{cp}(c) \geq 90\%$).

Code	Criterion	Distribution (%)	QWA	α	r	$\Delta\mu$
SC1	Summary available	0 (14), 1 (30), 2 (56)	0.90	0.49	0.61	0.11
SC2	Clarity	0 (10), 1 (38), 2 (53)	0.82	0.04	0.15	0.38
Average Structure and Clarity			0.86	0.27	0.38	0.25
LG1*	Language	0 (6), 1 (94)	0.88	-0.06	-0.06	-0.07
LG2*	Tone	0 (1), 1 (7), 2 (92)	0.98	0.23	0.39	-0.06
Average Language			0.93	0.08	0.17	-0.06
FU1	Learning Goals	0 (29), 1 (28), 2 (42)	0.83	0.51	0.50	-0.08
FU2	Success Criteria	0 (32), 1 (31), 2 (37)	0.79	0.35	0.36	0.05
Average Feed Up			0.81	0.43	0.43	-0.01
FB1*	Knowledge of result	0 (5), 1 (95)	0.95	0.26	0.40	-0.05
FB2*	Knowledge of correct response	0 (8), 1 (92)	0.90	-0.05	-0.05	-0.05
FB3	Knowledge about task constraints	0 (33), 1 (67)	0.65	0.07	0.08	0.07
FB4	Knowledge about concepts	0 (13), 1 (87)	0.82	-0.00	0.05	-0.12
FB5*	Knowledge about mistakes	0 (8), 1 (92)	0.90	-0.05	-0.05	0.01
FB6*	Self-feedback	0 (3), 1 (97)	0.97	-0.01	—	0.03
Average Feed Back			0.86	0.04	0.08	-0.02
FF1*	Self-regulation	0 (6), 1 (94)	0.95	-0.02	-0.02	-0.03
FF2	Learning Skills	0 (39), 1 (61)	0.60	0.06	0.10	0.17
Average Feed Forward			0.78	0.02	0.04	0.07
CQ1	Correctness	0 (1), 1 (28), 2 (71)	0.91	0.14	0.16	0.12
CQ2	Actionability	0 (1), 1 (32), 2 (66)	0.92	0.22	0.21	0.08
CQ3*	Argumentation	0 (10), 1 (90), 2 (0)	0.98	0.30	0.36	-0.06
Average Content Quality			0.94	0.22	0.24	0.05
ER1	Neglect	-1 (16), 0 (84)	0.86	0.40	0.41	-0.05
ER2*	Vague Critic	-1 (7), 0 (93)	0.91	0.12	0.12	0.00
ER3*	Out-of-Scope	-1 (4), 0 (96)	0.93	-0.03	-0.03	-0.01
ER4*	Missing Reference	-1 (2), 0 (98)	0.96	-0.01	—	-0.04
ER5*	Contradiction	-1 (2), 0 (98)	0.98	-0.00	-0.01	0.00
Average Errors			0.93	0.10	0.12	-0.02
AD1*	Additional points	0 (100), 1 (0)	1.00	—	—	0.00
AD2*	Negative points	-1 (2), 0 (98)	0.96	-0.01	—	-0.04
Average Additional points			0.98	-0.01	—	-0.02
Average all criteria			0.89	0.13	0.18	0.02

Table 20: Quadratic weighted agreement (QWA) between review Group 1 and Group 2 raters for each review criterion, with rubric-wise and overall averages. Columns show agreement (Agr), reliability (α), correlation (r), and mean bias ($\Delta\mu$). Rows marked with * have a highly skewed class distribution ($\max_{cp}(c) \geq 90\%$).

Code	Criterion	Distribution (%)	QWA	α	r	$\Delta\mu$
SC1	Summary available	0 (13), 1 (43), 2 (43)	0.91	0.60	0.62	0.00
SC2	Clarity	0 (12), 1 (48), 2 (40)	0.88	0.47	0.45	0.17
Average Structure and Clarity			0.89	0.53	0.54	0.08
LG1*	Language	0 (4), 1 (96)	0.95	0.38	0.38	0.02
LG2	Tone	0 (1), 1 (12), 2 (88)	0.93	0.02	0.00	0.07
Average Language			0.94	0.20	0.19	0.04
FU1	Learning Goals	0 (28), 1 (28), 2 (44)	0.75	0.31	0.33	0.25
FU2	Success Criteria	0 (34), 1 (35), 2 (31)	0.72	0.13	0.15	0.23
Average Feed Up			0.74	0.22	0.24	0.24
FB1*	Knowledge of result	0 (6), 1 (94)	0.88	-0.05	-0.06	-0.02
FB2	Knowledge of correct response	0 (11), 1 (89)	0.85	0.23	0.22	-0.02
FB3	Knowledge about task constraints	0 (37), 1 (63)	0.53	0.00	0.03	0.17
FB4	Knowledge about concepts	0 (13), 1 (87)	0.77	-0.00	-0.01	0.03
FB5*	Knowledge about mistakes	0 (10), 1 (90)	0.87	0.27	0.27	0.03
FB6*	Self-feedback	0 (4), 1 (96)	0.92	-0.03	-0.04	0.02
Average Feed Back			0.80	0.07	0.07	0.04
FF1*	Self-regulation	0 (9), 1 (91)	0.82	-0.09	-0.10	-0.02
FF2	Learning Skills	0 (48), 1 (52)	0.57	0.14	0.14	0.10
Average Feed Forward			0.69	0.02	0.02	0.04
CQ1	Correctness	0 (0), 1 (39), 2 (61)	0.91	0.27	0.27	0.02
CQ2	Actionability	0 (2), 1 (34), 2 (63)	0.89	0.23	0.22	0.05
CQ3	Argumentation	0 (12), 1 (88), 2 (0)	0.94	-0.12	-0.13	0.00
Average Content Quality			0.91	0.13	0.12	0.02
ER1	Neglect	-1 (15), 0 (85)	0.80	0.22	0.22	-0.03
ER2*	Vague Critic	-1 (5), 0 (95)	0.90	-0.04	-0.05	0.00
ER3*	Out-of-Scope	-1 (2), 0 (98)	0.95	-0.02	-0.02	-0.02
ER4*	Missing Reference	-1 (1), 0 (99)	0.98	0.00	—	0.02
ER5*	Contradiction	-1 (2), 0 (98)	0.95	-0.02	-0.02	-0.02
Average Errors			0.92	0.03	0.03	-0.01
AD1*	Additional points	0 (100), 1 (0)	1.00	—	—	0.00
AD2*	Negative points	-1 (1), 0 (99)	0.98	0.00	—	0.02
Average Additional points			0.99	0.00	—	0.01
Average all criteria			0.86	0.13	0.14	0.04

Table 21: Quadratic weighted agreement (QWA) within Group 2 raters for each review criterion, with rubric-wise and overall averages. Columns show agreement (Agr), reliability (α), correlation (r), and mean bias ($\Delta\mu$). Rows marked with * have a highly skewed class distribution ($\max_{cp}(c) \geq 90\%$).

Code	Qwen3 (80B)				Llama 3.3 (70B)				GPT OSS (120B)				GPT OSS (20B)			
	A_{G1}^{all}	A_{G1}^{ind}	A_{G2}^{all}	A_{G2}^{ind}	A_{G1}^{all}	A_{G1}^{ind}	A_{G2}^{all}	A_{G2}^{ind}	A_{G1}^{all}	A_{G1}^{ind}	A_{G2}^{all}	A_{G2}^{ind}	A_{G1}^{all}	A_{G1}^{ind}	A_{G2}^{all}	A_{G2}^{ind}
LA1	0.83	0.85	0.85	0.91	0.78	0.84	0.77	0.83	0.68	<u>0.59</u>	0.73	0.67	0.72	<u>0.57</u>	0.77	0.64
LA2	0.87	0.88	0.76	0.85	0.87	0.87	0.75	0.75	0.88	0.90	0.89	0.86	0.82	0.87	0.91	0.81
Language	0.85	0.87	0.80	0.88	0.82	0.85	0.76	0.79	0.78	0.75	0.81	0.76	0.77	0.72	0.84	0.73
ME1	0.91	0.92	0.82	0.86	0.91	0.90	0.84	0.87	0.91	0.84	0.87	0.84	0.82	0.81	0.90	0.86
Metadata	0.91	0.92	0.82	0.86	0.91	0.90	0.84	0.87	0.91	0.84	0.87	0.84	0.82	0.81	0.90	0.86
MO1	0.87	0.87	0.82	0.82	0.87	0.84	0.82	0.82	0.84	0.73	0.78	0.60	0.84	0.71	0.78	0.69
MO2	0.86	0.88	0.74	0.76	0.86	0.86	0.74	0.74	0.92	0.87	0.82	0.76	0.90	0.86	0.85	0.74
MO3	0.96	0.98	0.91	0.93	0.98	0.98	0.93	0.93	0.98	0.98	0.93	0.93	0.98	0.98	0.93	0.93
MO4	<u>0.38</u>	0.82	0.42	0.84	0.75	0.82	0.76	0.84	0.82	0.82	0.84	0.84	0.82	0.82	0.84	0.84
MO5	0.75	0.75	0.56	0.56	0.75	0.75	0.56	0.56	0.76	0.75	0.55	<u>0.56</u>	0.73	0.75	0.55	0.56
MO6	0.87	0.85	0.79	0.78	0.86	0.87	0.78	0.79	0.92	0.89	0.85	0.85	0.90	0.87	0.81	0.81
MO7	0.67	0.69	0.40	<u>0.42</u>	0.67	0.67	<u>0.40</u>	<u>0.40</u>	0.75	0.65	0.47	0.71	0.67	0.76	<u>0.40</u>	<u>0.49</u>
Motivation	0.77	0.83	0.66	0.73	0.82	0.83	0.71	0.72	0.85	0.81	0.75	0.75	0.83	0.82	0.74	0.72
AP1	0.89	0.91	0.87	0.84	0.91	0.91	0.84	0.84	0.93	0.85	0.85	0.82	0.93	0.89	0.89	0.85
AP2	0.86	0.85	0.76	0.71	0.87	0.84	0.79	0.70	0.89	0.88	0.81	0.77	0.86	0.86	0.91	0.74
AP3	0.75	0.70	0.64	0.60	0.75	0.72	0.72	0.65	0.78	0.77	0.73	0.66	0.85	0.75	0.87	0.67
AP4	<u>0.15</u>	0.78	<u>0.11</u>	0.62	0.80	0.78	0.64	0.62	0.82	0.76	0.65	0.67	0.84	0.82	0.67	0.65
AP5	0.62	0.76	0.49	0.56	0.80	0.76	0.60	0.56	0.84	0.75	0.56	0.62	0.82	0.84	0.62	0.67
AP6	0.83	0.77	0.76	0.66	0.80	0.76	0.74	0.63	0.81	0.84	0.88	0.80	0.78	0.83	0.88	0.82
AP7	<u>0.16</u>	<u>0.65</u>	<u>0.15</u>	0.65	<u>0.64</u>	<u>0.65</u>	0.64	0.65	<u>0.67</u>	0.71	0.71	0.60	<u>0.56</u>	0.75	0.65	0.71
AP8	0.96	0.96	0.85	0.85	0.96	0.96	0.85	0.85	0.96	0.96	0.85	0.85	0.96	0.96	0.85	0.85
Approach	0.65	0.80	0.58	0.69	0.82	0.80	0.73	0.69	0.84	0.82	0.76	0.72	0.82	0.84	0.79	0.75
MD1	0.67	0.67	0.44	<u>0.44</u>	0.67	0.67	<u>0.44</u>	<u>0.44</u>	0.69	<u>0.64</u>	<u>0.45</u>	0.58	0.65	0.78	<u>0.49</u>	0.58
MD2	0.73	0.76	0.69	0.71	0.75	0.76	0.69	0.71	0.76	0.75	0.71	0.69	0.76	0.76	0.71	0.71
MD3	0.80	0.82	0.67	0.69	0.80	0.82	0.67	0.69	0.85	0.78	0.73	0.69	0.87	0.78	0.75	0.73
MD4	<u>0.56</u>	<u>0.58</u>	0.73	0.71	<u>0.56</u>	<u>0.58</u>	0.73	0.71	<u>0.62</u>	0.64	0.71	0.73	<u>0.60</u>	<u>0.58</u>	0.73	0.67
MD5	0.71	0.69	0.45	0.44	0.71	0.71	0.49	0.45	0.78	0.64	0.56	0.64	0.82	0.70	0.71	0.70
MD6	0.58	<u>0.56</u>	<u>0.33</u>	<u>0.31</u>	<u>0.58</u>	<u>0.56</u>	<u>0.36</u>	<u>0.31</u>	0.67	0.73	<u>0.45</u>	<u>0.55</u>	0.64	0.65	0.56	<u>0.55</u>
MD7	0.56	<u>0.56</u>	<u>0.38</u>	<u>0.38</u>	<u>0.56</u>	<u>0.56</u>	<u>0.38</u>	<u>0.38</u>	<u>0.58</u>	<u>0.60</u>	<u>0.44</u>	<u>0.45</u>	<u>0.60</u>	<u>0.60</u>	<u>0.42</u>	<u>0.42</u>
Methodology	0.66	0.66	0.53	0.52	0.66	0.67	0.54	0.53	0.71	0.68	0.58	0.62	0.71	0.70	0.62	0.62
SC1	0.98	0.98	0.95	0.95	0.98	0.98	0.95	0.95	0.96	0.96	0.96	0.96	0.98	0.96	0.95	0.96
SC2	0.98	0.98	0.95	0.95	0.98	0.98	0.95	0.95	0.98	0.95	0.95	0.91	0.98	0.98	0.95	0.95
SC3	0.91	0.91	0.96	0.96	0.91	0.91	0.96	0.96	0.91	0.93	0.96	0.98	0.91	0.91	0.96	0.96
SC4	0.85	0.82	0.85	0.78	0.90	0.78	0.95	0.72	0.85	0.83	0.88	0.81	0.85	0.77	0.88	0.74
Schedule	0.93	0.92	0.93	0.91	0.94	0.91	0.95	0.89	0.93	0.92	0.94	0.92	0.93	0.91	0.93	0.90
BI1	0.93	0.72	0.72	0.59	0.91	0.91	0.71	0.71	<u>0.27</u>	<u>0.13</u>	<u>0.40</u>	<u>0.33</u>	<u>0.20</u>	<u>0.49</u>	<u>0.40</u>	<u>0.55</u>
BI2	0.90	0.92	0.68	0.73	0.90	0.88	0.71	0.66	0.86	0.91	0.80	0.72	0.91	0.88	0.75	0.70
Bibliography	0.91	0.82	0.70	0.66	0.91	0.90	0.71	0.68	0.57	0.52	0.60	0.52	0.56	0.69	0.57	0.62
All criteria	0.75	0.80	0.66	0.70	0.81	0.80	0.71	0.70	0.81	0.77	0.74	0.72	0.79	0.79	0.75	0.73

Table 22: Quadratic-weighted agreement between LLM scores and human ratings for exposé criteria. Columns report A_{G1} (Group 1 vs. LLM) and A_{G2} (Group 2 vs. LLM) for the prompting variants shown (all, ind). Bold/underlined cells mark column-wise maxima/minima. Average rows summarize agreement within rubrics and across all criteria.

Code	Qwen3 (80B)				Llama 3.3 (70B)				GPT OSS (120B)				GPT OSS (20B)			
	A_{G1}^{all}	A_{G1}^{ind}	A_{G2}^{all}	A_{G2}^{ind}	A_{G1}^{all}	A_{G1}^{ind}	A_{G2}^{all}	A_{G2}^{ind}	A_{G1}^{all}	A_{G1}^{ind}	A_{G2}^{all}	A_{G2}^{ind}	A_{G1}^{all}	A_{G1}^{ind}	A_{G2}^{all}	A_{G2}^{ind}
SC1	0.83	0.83	0.85	0.84	0.84	0.85	0.88	0.89	0.85	0.83	0.87	0.87	0.86	0.88	0.88	0.91
SC2	0.92	0.92	0.80	0.81	0.90	0.86	0.81	0.80	0.91	0.86	0.83	0.85	0.84	0.85	0.85	0.83
Structure and Clarity	0.88	0.87	0.82	0.83	0.87	0.85	0.84	0.85	0.88	0.84	0.85	0.86	0.85	0.87	0.87	0.87
LG1	0.90	0.91	0.98	0.97	0.91	0.88	0.97	0.94	0.86	0.67	0.85	0.62	0.75	0.69	0.75	0.64
LG2	0.98	0.98	0.99	0.98	0.98	0.98	0.99	0.99	0.98	0.99	0.99	1.00	0.96	0.97	0.98	1.00
Language	0.94	0.94	0.99	0.97	0.94	0.93	0.98	0.97	0.92	0.83	0.92	0.81	0.85	0.83	0.87	0.82
FU1	<u>0.75</u>	<u>0.70</u>	<u>0.74</u>	<u>0.69</u>	<u>0.74</u>	<u>0.78</u>	<u>0.72</u>	<u>0.73</u>	<u>0.64</u>	<u>0.63</u>	<u>0.57</u>	<u>0.55</u>	0.66	0.55	0.60	0.48
FU2	<u>0.75</u>	<u>0.71</u>	<u>0.73</u>	<u>0.69</u>	<u>0.76</u>	<u>0.78</u>	<u>0.72</u>	<u>0.75</u>	<u>0.64</u>	<u>0.66</u>	<u>0.62</u>	<u>0.65</u>	0.58	0.61	0.59	0.60
Feed Up	0.75	0.71	0.73	0.69	0.75	0.78	0.72	0.74	0.64	0.64	0.59	0.60	0.62	0.58	0.59	0.54
FB1	0.94	0.94	0.99	0.99	0.94	0.92	0.99	0.96	<u>0.23</u>	<u>0.10</u>	<u>0.20</u>	<u>0.06</u>	<u>0.25</u>	<u>0.09</u>	<u>0.22</u>	<u>0.05</u>
FB2	0.92	0.93	0.97	0.96	0.91	0.93	0.95	0.96	<u>0.24</u>	<u>0.09</u>	<u>0.19</u>	<u>0.05</u>	<u>0.08</u>	<u>0.10</u>	<u>0.04</u>	<u>0.04</u>
FB3	0.79	0.79	<u>0.73</u>	<u>0.73</u>	<u>0.79</u>	0.78	<u>0.73</u>	0.74	0.77	0.62	0.73	0.63	0.75	0.59	0.69	0.58
FB4	0.84	0.84	0.96	0.96	0.84	0.84	0.96	0.96	0.83	0.74	0.93	0.78	0.84	0.84	0.96	0.91
FB5	0.95	0.95	0.94	0.94	0.95	0.94	0.94	0.92	0.95	0.92	0.94	0.92	0.95	0.76	0.94	0.79
FB6	1.00	0.99	0.97	0.96	0.99	0.99	0.98	0.98	0.96	0.83	0.95	0.82	<u>0.34</u>	0.41	<u>0.37</u>	<u>0.42</u>
Feed Back	0.91	0.91	0.93	0.92	0.90	0.90	0.93	0.92	0.66	0.55	0.66	0.54	0.54	0.47	0.54	0.46
FF1	0.96	0.96	0.99	0.99	0.96	0.96	0.99	0.99	0.96	0.90	0.99	0.92	0.91	0.71	0.90	0.70
FF2	<u>0.78</u>	<u>0.78</u>	<u>0.61</u>	<u>0.61</u>	<u>0.78</u>	<u>0.78</u>	<u>0.61</u>	<u>0.61</u>	0.75	<u>0.37</u>	0.62	<u>0.50</u>	0.78	<u>0.36</u>	0.71	0.51
Feed Forward	0.87	0.87	0.80	0.80	0.87	0.87	0.80	0.80	0.86	0.63	0.81	0.71	0.84	0.53	0.80	0.60
CQ1	0.96	0.94	0.92	0.92	0.94	0.88	0.93	0.86	0.92	0.95	0.90	0.92	0.84	0.90	0.84	0.90
CQ2	0.94	0.94	0.92	0.92	0.93	0.94	0.93	0.91	0.94	0.94	0.92	0.92	0.92	0.93	0.93	0.92
CQ3	<u>0.77</u>	0.98	0.80	0.99	0.98	0.98	0.99	0.99	0.98	0.96	0.99	0.96	0.98	0.96	0.99	0.97
Content Quality	0.89	0.95	0.88	0.94	0.95	0.93	0.95	0.92	0.95	0.95	0.94	0.93	0.91	0.93	0.92	0.93
ER1	0.83	<u>0.63</u>	0.88	<u>0.68</u>	0.83	<u>0.44</u>	0.88	<u>0.45</u>	0.75	0.65	0.77	0.70	0.80	0.82	0.83	0.87
ER2	0.92	0.74	0.92	0.75	0.93	<u>0.65</u>	0.95	<u>0.61</u>	0.89	0.91	0.89	0.91	0.86	0.94	0.88	0.94
ER3	0.96	0.91	0.97	0.92	0.96	<u>0.93</u>	0.97	<u>0.94</u>	0.96	0.92	0.97	0.93	0.96	0.92	0.97	0.93
ER4	0.80	<u>0.25</u>	0.84	<u>0.26</u>	0.94	<u>0.17</u>	0.98	<u>0.13</u>	<u>0.42</u>	<u>0.32</u>	<u>0.38</u>	<u>0.32</u>	<u>0.36</u>	<u>0.36</u>	<u>0.38</u>	<u>0.34</u>
ER5	0.99	0.95	0.99	0.95	0.99	0.94	0.99	0.94	0.99	0.97	0.99	0.97	0.99	0.95	0.99	0.95
Errors	0.90	0.69	0.92	0.71	0.93	0.63	0.95	0.62	0.80	0.75	0.80	0.77	0.79	0.80	0.81	0.81
All criteria	0.89	0.84	0.89	0.84	0.90	0.83	0.90	0.82	0.79	0.72	0.78	0.72	0.74	0.69	0.74	0.69

Table 23: Quadratic-weighted agreement between LLM scores and human ratings for review criteria. Columns report A_{G1} (Group 1 vs. LLM) and A_{G2} (Group 2 vs. LLM) for the prompting variants shown (all, ind). Bold/underlined cells mark column-wise maxima/minima. Average rows summarize agreement within rubrics and across all criteria.

Model	Mode	s/assessment	s/prompt	tokens/prompt	tokens/input	tokens/output
GPT OSS (20B)	All [min-max]	110.20 (\pm 27.69) [84.88-289.55]	110.20 (\pm 27.69) [84.88-289.55]	8,010.67 (\pm 1,018.78) [6,384.00-11,178.00]	5,587.85 (\pm 932.59) [4,092.00-8,765.00]	2,422.82 (\pm 629.58) [1,874.00-6,489.00]
Llama 3.3 (70B)	All [min-max]	182.10 (\pm 12.79) [163.18-212.11]	182.10 (\pm 12.79) [163.18-212.11]	7,181.78 (\pm 973.95) [5,594.00-10,311.00]	5,606.69 (\pm 947.16) [4,091.00-8,839.00]	1,575.09 (\pm 105.46) [1,399.00-1,858.00]
GPT OSS (120B)	All [min-max]	258.01 (\pm 13.74) [227.52-304.39]	258.01 (\pm 13.74) [227.52-304.39]	8,227.11 (\pm 915.79) [6,692.00-11,292.00]	5,679.85 (\pm 932.59) [4,184.00-8,857.00]	2,547.25 (\pm 148.73) [2,265.00-3,067.00]
Llama 3.3 (70B)	Ind [min-max]	284.30 (\pm 26.54) [237.56-370.23]	9.17 (\pm 1.37) [5.43-14.43]	3,228.31 (\pm 948.58) [1,656.00-6,555.00]	3,162.76 (\pm 947.61) [1,613.00-6,490.00]	65.55 (\pm 9.79) [42.00-117.00]
Qwen3 (80B)	All [min-max]	430.37 (\pm 15.13) [394.29-460.62]	430.37 (\pm 15.13) [394.29-460.62]	7,473.75 (\pm 1,046.99) [5,775.00-11,079.00]	5,719.84 (\pm 1,024.62) [4,126.00-9,301.00]	1,753.91 (\pm 64.05) [1,602.00-1,874.00]
GPT OSS (20B)	Ind [min-max]	582.36 (\pm 103.30) [452.01-878.16]	18.79 (\pm 20.85) [7.34-357.11]	3,553.87 (\pm 1,047.56) [1,821.00-11,982.00]	3,162.31 (\pm 933.03) [1,633.00-6,431.00]	391.56 (\pm 478.39) [123.00-8,091.00]
Qwen3 (80B)	Ind [min-max]	862.02 (\pm 36.97) [789.80-946.85]	27.81 (\pm 3.81) [19.10-41.08]	3,359.56 (\pm 1,027.13) [1,705.00-7,050.00]	3,275.90 (\pm 1,025.04) [1,648.00-6,952.00]	83.66 (\pm 15.20) [51.00-138.00]
GPT OSS (120B)	Ind [min-max]	1133.41 (\pm 89.76) [869.61-1398.68]	34.35 (\pm 6.90) [19.20-97.15]	3,383.58 (\pm 932.15) [1,771.00-6,902.00]	3,160.07 (\pm 933.04) [1,627.00-6,431.00]	223.51 (\pm 67.47) [100.00-854.00]

Table 24: Runtime and token usage for LLM scoring. Entries report mean (\pm SD) and a second row with [min-max]. Mode: All = combined prompting; Ind = single-criterion prompting. Token statistics are reported per prompt (total, input, output).

Model	Mode	s/assessment	s/prompt	tokens/prompt	tokens/input	tokens/output
GPT OSS (20B)	All [min-max]	80.25 (\pm 20.05) [60.56-345.62]	80.25 (\pm 20.05) [60.56-345.62]	4,465.78 (\pm 606.62) [3,641.00-10,957.00]	2,539.31 (\pm 402.22) [2,117.00-6,330.00]	1,926.47 (\pm 472.59) [1,449.00-8,091.00]
Llama 3.3 (70B)	All [min-max]	113.76 (\pm 8.10) [98.39-135.10]	113.76 (\pm 8.10) [98.39-135.10]	3,574.05 (\pm 434.77) [3,041.00-7,464.00]	2,513.88 (\pm 404.74) [2,088.00-6,329.00]	1,060.17 (\pm 73.19) [920.00-1,264.00]
Llama 3.3 (70B)	Ind [min-max]	142.31 (\pm 11.27) [121.03-210.90]	6.47 (\pm 0.95) [4.26-11.25]	982.85 (\pm 406.97) [518.00-4,845.00]	924.60 (\pm 405.24) [471.00-4,774.00]	58.25 (\pm 8.26) [40.00-93.00]
GPT OSS (120B)	All [min-max]	194.64 (\pm 14.13) [167.25-252.35]	194.64 (\pm 14.13) [167.25-252.35]	4,485.63 (\pm 447.10) [3,823.00-8,242.00]	2,540.12 (\pm 403.33) [2,117.00-6,330.00]	1,945.51 (\pm 146.74) [1,670.00-2,555.00]
Qwen3 (80B)	All [min-max]	289.58 (\pm 13.70) [254.40-359.69]	289.58 (\pm 13.70) [254.40-359.69]	3,705.31 (\pm 483.11) [3,087.00-7,550.00]	2,567.16 (\pm 460.99) [2,066.00-6,295.00]	1,138.15 (\pm 50.89) [1,004.00-1,361.00]
GPT OSS (20B)	Ind [min-max]	343.64 (\pm 81.03) [240.61-792.63]	15.62 (\pm 15.37) [4.91-344.67]	1,316.56 (\pm 561.36) [687.00-9,596.00]	955.08 (\pm 402.72) [505.00-4,779.00]	361.47 (\pm 363.39) [108.00-8,091.00]
Qwen3 (80B)	Ind [min-max]	485.12 (\pm 25.59) [421.38-576.16]	22.05 (\pm 2.33) [15.73-40.33]	970.92 (\pm 407.79) [516.00-4,822.00]	902.12 (\pm 405.52) [449.00-4,740.00]	68.80 (\pm 9.00) [45.00-145.00]
GPT OSS (120B)	Ind [min-max]	512.82 (\pm 88.18) [375.80-1048.64]	23.31 (\pm 9.87) [8.89-465.17]	1,171.90 (\pm 420.31) [634.00-5,883.00]	955.08 (\pm 402.72) [505.00-4,779.00]	216.81 (\pm 96.63) [85.00-4,827.00]

Table 25: Runtime and token usage for LLM scoring. Entries report mean (\pm SD) and a second row with [min-max]. Mode: All = combined prompting; Ind = single-criterion prompting. Token statistics are reported per prompt (total, input, output).