

Beyond Factual QA: Mentorship-Oriented Question Answering over Long-Form Multilingual Content

Anonymous ACL submission

Abstract

Question answering systems are typically evaluated on factual correctness, yet many real-world applications—such as education and career guidance—require *mentorship*: responses that provide reflection and guidance. Existing QA benchmarks rarely capture this distinction, particularly in multilingual and long-form settings. We introduce MENTORQA, the first multilingual dataset and evaluation framework for mentorship-focused question answering from long-form videos, comprising nearly 9,000 QA pairs from 180 hours of content across four languages. We define mentorship-focused evaluation dimensions that go beyond factual accuracy, capturing clarity, alignment, and learning value. Using MENTORQA, we compare Single-Agent, Dual-Agent, RAG, and Multi-Agent QA architectures under controlled conditions. Multi-Agent pipelines consistently produce higher-quality mentorship responses, with especially strong gains for complex topics and lower-resource languages. We further analyze the reliability of automated LLM-based evaluation, observing substantial variation in alignment with human judgments. Overall, this work establishes mentorship-focused QA as a distinct research problem and provides a multilingual benchmark for studying agentic architectures and evaluation design in educational AI. The dataset and evaluation framework are released at <https://anonymous.4open.science/r/MentorQA/>.

1 Introduction

Question answering (QA) systems are typically evaluated on factual correctness (Mulla and Gharpure, 2023; Fu et al., 2024). However, in many real-world settings—such as education, career guidance, and personal development—users also seek *mentorship*: answers that provide context, reflection, and guidance rather than isolated

English | Career/Education | Multi-Agent

Q: What should a student consider when choosing a research topic that aligns with their interests?

A: When selecting a research topic that aligns with their interests, a student should consider several key factors: 1. ****Interest and Passion**:** The student must genuinely...

Scores: QFlu 5, AFlu 5, QClar 4, AClar 5, QA-Algn 5, QMen 5, AMen 5

Chinese | Education | Single-Agent

Q: 千禧世代的出生年份是什麼? (*What are the birth years of the Millennial generation?*)

A: 千禧世代的出生年份是 1980 年至 1994 年或 1995 年。 (*The birth years of the Millennial generation are 1980 to 1994 or 1995.*)

Scores: QFlu 3, AFlu 4, QClar 5, AClar 2, QA-Algn 4, QMen 1, AMen 1

Hindi | Business | Dual-Agent

Q: भारत की ऐतिहासिक रूप से धनवान था तो उसके व्यापार का क्या रूप था? (*If India was historically wealthy... how was the nature of trade?*)

A: भारत ने चानक्य के समय से बीटा को बोर्ड ट्रेड के एक अच्छे देश के रूप में जाना था।..... डाट बिक था जो तब विश्व की दुनिया का नक्षा था। (*From Chanakya's time, India was known as.. and carried out trade with..*)

Scores: QFlu 3, AFlu 3, QClar 2, AClar 3, QA-Algn 5, QMen 1, AMen 3

Romanian | Mental Health | RAG

Q: Ce implică muncă prea puțin stimulantă și ce efecte are asupra sistemului nervos? (*What does work that is too little stimulating imply..... and affects nervous system?*)

A: Muncă prea puțin stimulantă implică simțirea de plictisire și blazare, ceea ce poate duce la o solicitare excesivă... (*Work that is not stimulating enough involves feelings of boredom.... lead to excessive strain on the nervous system, causing.....*)

Scores: QFlu 2, AFlu 4, QClar 2, AClar 4, QA-Algn 1, QMen 4, AMen 4

Table 1: MENTORQA multilingual examples across languages, topics, and generation models. Scores denote Fluency (Q/AFlu), Clarity (Q/AClar), QA Alignment (QA-Algn), Mentorship (Q/AMen) on 1-5 Likert scale.

facts (Ma and Ma, 2019; Kovács, 2016). Existing QA benchmarks and models rarely capture this distinction, particularly in multilingual and long-form scenarios. This gap is especially evident for long-form mentorship content, where relevant guidance is sparsely distributed across ex-

042
043
044
045
046
047

tended narratives. Answering mentorship-focused questions requires identifying and synthesizing insights that support learning and decision-making, going beyond surface-level correctness. While recent language models have advanced open-domain QA (Guo et al., 2024), it remains unclear how well they support mentorship-focused question answering across languages with different discourse norms and resource availability.

Prior QA research has largely focused on short-form inputs, factual accuracy, and monolingual evaluation (Faraby et al., 2024). Multilingual QA benchmarks emphasize linguistic fluency and correctness, but seldom assess whether responses provide meaningful guidance or learning value. As a result, datasets, metrics, and systematic evaluations for *mentorship-focused multilingual QA from long-form content* are largely missing. Moreover, the impact of architectural choices—such as retrieval augmentation or agentic coordination—on mentorship quality remains underexplored.

Contributions. **First, we introduce MENTORQA, the first multilingual dataset and evaluation framework designed for mentorship-focused question answering from long-form videos.** The dataset contains nearly 9,000 QA pairs derived from 180 hours of mentorship content across four languages. **Second, we define mentorship-focused evaluation dimensions that capture clarity, alignment, and guidance, and compare four QA architectures under controlled settings.** Our results show that Multi-Agent pipelines consistently achieve higher mentorship quality, with particularly strong gains for complex topics and lower-resource languages. **Finally, we examine the reliability of automated LLM-based evaluation across languages and metrics, finding substantial variation in alignment with human judgments. Together, our findings establish mentorship-focused QA as a distinct research problem and provide a multilingual benchmark for studying agentic architectures and evaluation methods in educational AI.**

2 Related Work

Educational and Mentorship QA from Long Transcripts. Recent work has explored extracting questions and answers from long educational recordings to support learning. Existing systems generate questions from lecture summaries (Chen and Yen, 2024), select or rewrite context for video-

based educational QA (Yu et al., 2025), or apply segment-level retrieval and extractive QA to podcasts and meeting transcripts (Elaryan, 2022; Prasad et al., 2023). Other efforts move toward personalization, including QA for language learning assistants (Sammoudi et al., 2025) and multi-agent debate frameworks for online education (Du et al., 2025). While these approaches address educational content, they primarily optimize for factual correctness or content coverage and are often monolingual or tied to a single modeling paradigm. In contrast, mentorship-focused QA emphasizes guidance, reflection, and learning value—objectives that remain underexplored in long-form, multilingual settings. Our work targets this gap by extracting mentorship-focused QA from long-form videos across multiple languages and evaluating architectural choices under a unified framework.

Agentic and RAG-Based QA Pipelines. Agentic models have been proposed to support complex, multi-step workflows for long-document understanding and QA. Prior work demonstrates how specialized agents can coordinate to decompose reading, reasoning, and generation tasks (Wang et al., 2025; Saadaoui et al., 2025), and recent surveys analyze common collaboration mechanisms in such systems (Tran et al., 2025). Retrieval-augmented generation (RAG) has also been widely adopted for domain-specific QA, including applications in health and risk assessment (Meng et al., 2025), as well as multi-stage pipelines for podcast and meeting analysis (Aquilina et al., 2023; Zhu et al., 2025). However, these approaches are typically evaluated on generic QA objectives or task efficiency, and rarely compare architectural paradigms under controlled conditions. Our work differs in two key ways: we systematically compare Single-Agent, Dual-Agent, RAG, and Multi-Agent pipelines using a fixed base model, and we evaluate them with respect to mentorship-focused quality rather than factual accuracy alone, revealing when agentic coordination is most beneficial.

Multilingual QA and Metric-Based Evaluation. A parallel body of work has examined evaluation methods for question generation and the extension of QA beyond English. Early studies investigated automatic QG metrics (Nema and Khapra, 2018), followed by reference-free and answer-based approaches such as RQUGE (Mohammadshahi et al., 2023) and QG-Eval (Fu et al., 2024), as well as LLM-as-a-judge frameworks for rubric-driven

evaluation (Wolfe, 2024). Multilingual QA benchmarks further explore QA and question generation across languages and domains (Ushio et al., 2023; Moreno-Cediel et al., 2024; Ruder and Sil, 2021; Asai et al., 2021). However, recent analyses show that many automatic metrics correlate weakly with human judgments, particularly for task-oriented dimensions such as usefulness or guidance (Fu et al., 2024). Motivated by these limitations, we introduce a streamlined evaluation framework with mentorship-focused metrics designed to better capture learning value and to enable systematic comparison between human and automated evaluation in multilingual QA.

3 The MENTORQA Dataset

We introduce MENTORQA, to our knowledge the first QA dataset focused on mentorship. Unlike traditional QA benchmarks that emphasize factual recall, MENTORQA centers on guidance, reflection, and practical insight. We define mentorship value as information that helps learners grow through advice, perspective, or actionable guidance. The dataset consists of 8,990 QA pairs drawn from 120 long-form mentorship videos, covering four languages and six mentorship topics. We next describe our data collection pipeline, including video selection and QA generation.

Video Collection and Processing. We collect long-form mentorship videos from YouTube, focusing on podcast-style conversations and panel discussions that provide reflective, experience-driven guidance. Each video averages over 1.5 hours, resulting in more than 100 hours of multilingual mentorship content. We process only the audio track, transcribing it with Whisper¹, which automatically detects the source language and produces transcripts used as input to the QA models.

Languages. The dataset covers four languages—English, Chinese, Hindi, and Romanian—selected to span typologically and culturally diverse settings while supporting reliable, high-quality annotation. All videos and QA pairs are balanced across languages and collected and evaluated by native or expert speakers, ensuring culturally informed analysis and evaluation.

Topic Selection across Languages. Videos are selected to cover six mentorship-related topics²: *Entrepreneurship, Education, Finance, Mental*

¹<https://github.com/openai/whisper>

²All the mentorship videos are in Appendix B.

Dataset Scale and Coverage	
Videos (long-form)	120
Total duration (hours)	180
Languages	4
Mentorship topics	6
QA-generation models	4
Total QA pairs	8,990
Evaluation and Quality Control	
Evaluation dimensions	4
LLM judges	9
Human annotators	12
Mentorship Topics (# QA pairs)	
Personal Growth	4,450
Career & Education	2,689
Mental Health	741
Business & Entrepreneurship	392
Physical Health	431
Finance	287

Table 2: MENTORQA statistics.

Health, Personal Growth, and Physical Health. Topic labels are automatically assigned using Qwen2.5-7B (Yang et al., 2025) and subsequently verified by expert annotators with proficiency in the respective languages.

English videos focus on academic and professional mentorship, including research directions, publishing, PhD pathways through panel discussion with experts in academia and industry, such as ACL Mentorship.³

Romanian videos center on education, personal development, and mental health through conversations with professionals from diverse backgrounds, such as authors, coaches or athletes.⁴

Chinese videos center on experience-driven mentorship in areas such as mental health, career growth, finance, and physical well-being.⁵

Hindi videos focus on practical guidance related to career advancement, business, and mental or physical health.⁶

QA Generation. We generate question-answer pairs from each video transcript using four complementary QA-generation models: *Single-Agent, Dual-Agent, RAG, and Multi-Agent*. These models are designed to address challenges of long-form, conversational mentorship content, including long-context reasoning and topic drift. For each video, each model generates approximately 20 mentorship-focused QA pairs, yielding about

³<https://www.youtube.com/@aclmentorship>

⁴<https://www.youtube.com/@MindArchitect>

⁵<https://www.youtube.com/@DrinkingLibrary>

⁶<https://www.youtube.com/@ranveerallahbadia>

80 QA pairs per video. Across 120 videos, this process results in 8,990 QA pairs (Table 2). Multilingual examples are shown in Table 1, and model architectures are described in Section 4.

Quality Assurance. Each QA pair is evaluated along four dimensions: *Fluency* and *Clarity* (linguistic quality), and *QA Alignment* and *Mentorship Value* (task-oriented quality). Human evaluation is conducted by three annotators per language, complemented by nine LLM-based judges to support scalable analysis and study alignment with human judgments. Evaluation is described in Section 5.

PII Anonymization. Although the source data are public, we proactively reduce the presence and downstream propagation of personally identifiable information (PII). We apply a layered anonymization strategy that combines automatic detection—using named entity recognition based on XLM-RoBERTa⁷, structured PII detection and anonymization via Presidio⁸, and LLM-based prompts to flag implicit identifiers—with human review, which is especially important for lower-resource languages and culturally specific entities.

Scalability. MENTORQA is released as an extensible, open-source resource. While the current version covers four languages, a limited number of channels, and a single video platform, the codebase and documentation support direct expansion to additional languages and content sources.

4 Mentorship-focused QA Generation

Mentorship-focused question answering from long-form conversational videos poses challenges that go beyond standard QA generation. In addition to long-context reasoning and topic shifts, mentorship content exhibits a critical property: *mentorship value is unevenly distributed*. Short segments rich in reflection, advice, or experience-driven insight often carry more educational value than longer descriptive or narrative portions. Effective mentorship-focused QA must therefore identify, prioritize, and allocate questions based on the *quality* of guidance rather than transcript length alone.

To study this problem systematically, we design four QA-generation pipelines of increasing structure: *Single-Agent*, *Dual-Agent*, *Retrieval-Augmented Generation (RAG)*, and *Multi-Agent*.

⁷<https://huggingface.co/Davlan/xlm-roberta-base-ner-hrl>

⁸<https://github.com/microsoft/presidio>

These pipelines progressively introduce (1) topic awareness, (2) answer grounding, and (3) explicit control over mentorship quality and question allocation.

Across all pipelines, we use the same base language model, *Qwen-2.5 7B-Instruct-1M* (Yang et al., 2025), configured in long-context mode. This controlled setup ensures that performance differences arise from architectural choices rather than model capacity. Full prompts are provided in subsection C.1.

4.1 Single-Agent

The Single-Agent baseline processes the full video transcript using a single long-context LLM. The model is instructed via a system prompt to act as an expert educational analyst and to generate mentorship-focused questions and answers. The prompt enforces broad topical coverage, discourages trivial or localized questions, and requires answers that provide educational or mentorship value. For each transcript, the model generates a fixed set of 20 question–answer pairs spanning the entire conversation. This baseline relies solely on the long-context capabilities of the underlying LLM, without explicit segmentation, retrieval, or coordination mechanisms.

4.2 Dual-Agent

While the Single-Agent baseline can generate high-quality questions, it often over-concentrates QA pairs in dense or early transcript regions. To address this, the Dual-Agent pipeline separates topic segmentation from QA generation.

A *Chunking Agent* first partitions the transcript into coherent, topic-based segments. The agent is instructed to cover the full conversation while avoiding overly small or fragmented segments. Segment lengths adapt dynamically to topic duration, and the content itself remains unchanged.

Each segment is then passed to a *QA Agent*, which generates a fixed number of mentorship-focused QA pairs per segment using the same prompt as the Single-Agent baseline. This topic-aware decomposition yields more balanced coverage across the conversation and improves alignment with underlying mentorship themes.

4.3 Retrieval Augmented Generation (RAG)

RAG explicitly separates question generation from answer grounding. In this pipeline, the model

321 first generates a fixed set of diverse, mentorship-
 322 focused questions from the full transcript. Each
 323 question is then answered using only relevant tran-
 324 script segments retrieved via similarity search. We
 325 construct a multilingual retrieval index using the
 326 BGE-M3 embedding model (Chen et al., 2024),
 327 which automatically chunks long transcripts and
 328 supports all four dataset languages. For each ques-
 329 tion, the most relevant transcript chunks are re-
 330 trieved and provided as context to a second LLM
 331 query, which generates the final answer. This
 332 approach enables global question coverage while
 333 keeping answers explicitly tied to relevant tran-
 334 script segments.

335 4.4 Multi-Agent

336 While the Dual-Agent model improves cover-
 337 age through topic-aware segmentation and RAG
 338 grounds answers via retrieval, both approaches
 339 rely on a single-pass question generation strategy.
 340 *In practice, we observe substantial variation in*
 341 *mentorship quality across topics: uniformly al-*
 342 *locating questions can over-represent weak seg-*
 343 *ments and under-represent parts of the conversa-*
 344 *tion rich in guidance.* To address this, we introduce
 345 a Multi-Agent framework that explicitly assesses
 346 mentorship quality and adaptively allocates ques-
 347 tions across topics based on their estimated value.
 348 As shown in Figure 1, the pipeline decomposes QA
 349 generation into specialized agents responsible for
 350 topic discovery, question exploration, quality as-
 351 sessment, allocation, and answer generation.

352 The process begins with the same *Chunking*
 353 *Agent* used in the Dual-Agent setting, which per-
 354 forms dynamic topic-based segmentation and as-
 355 signs each segment a concise title. For each seg-
 356 ment, a *Question Brainstorming Agent* deliberately
 357 over-generates candidate questions, emphasizing
 358 diversity and mentorship relevance. A *Scorer*
 359 *Agent* then assigns each candidate question q an
 360 integer mentorship score $score(q) \in \{1, \dots, 10\}$.
 361 To capture segment-level importance, we compute
 362 a *segment strength* score:

$$363 \quad S(i) = \mu(q) - \sigma(q), \quad (1)$$

364 where $\mu(q)$ and $\sigma(q)$ denote the mean and standard
 365 deviation of question scores within segment i . This
 366 formulation favors segments with consistently high
 367 mentorship value while penalizing segments with
 368 noisy or uneven question quality. Given a target
 369 of T questions per video, each segment receives a

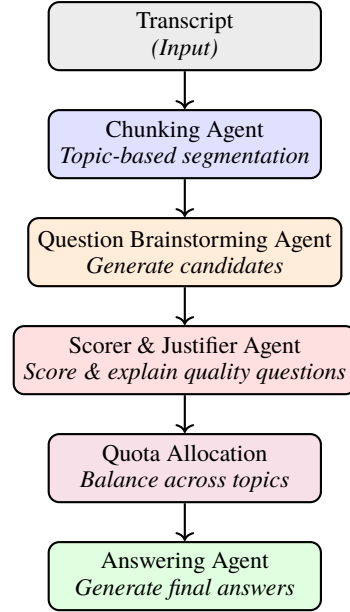


Figure 1: Multi-Agent QA generation pipeline. The framework decomposes QA generation from long-form transcripts into specialized agents for topic segmentation, question exploration, quality scoring, quota-based allocation, and grounded answer generation. This design enables explicit control over mentorship quality and balanced question allocation across topics.

quota proportional to its strength:

$$370 \quad \text{quota}(i) = \text{round}\left(\frac{S(i)}{\sum_{j \in \mathcal{S}} S(j)} \times T\right), \quad (2) \quad 371$$

372 with a minimum of one question per non-empty
 373 segment. Within each segment, only the top-
 374 ranked questions are retained. To improve trans-
 375 parency and reduce redundancy, a *Justifier Agent*
 376 generates short explanations for question selec-
 377 tion or rejection. Finally, an *Answering Agent*
 378 produces the final answers, grounding each re-
 379 sponse in its corresponding segment and empha-
 380 sizing mentorship-oriented explanations.

381 Overall, this design enables broader exploration
 382 of the question space, explicit assessment of men-
 383 torship quality, and adaptive allocation of ques-
 384 tions to the most valuable parts of long-form men-
 385 torship conversations.

386 5 Evaluation Methodology

387 We evaluate the QA-generation models using both
 388 human and automatic evaluation. Beyond stan-
 389 dard linguistic and alignment measures, we intro-
 390 duce mentorship-focused evaluation dimensions
 391 that capture guidance, reflection, and learning
 392 value—properties not addressed by current QA
 393 evaluation frameworks.

5.1 Evaluation Metrics

Evaluating mentorship-focused QA requires going beyond factual correctness. We build on QG-Eval (Fu et al., 2024), which distinguishes linguistic and task-oriented dimensions, and extend it to explicitly capture alignment and mentorship value. While Fluency and Clarity are adopted from prior work, we evaluate them separately for questions and answers to enable more fine-grained analysis.

Linguistic Metrics. We assess surface-level language quality using the following metrics, applied independently to questions and answers:

- **Fluency:** *The question or answer is grammatically correct and free of language errors.*
- **Clarity:** *The question or answer is easy to understand, specific, and unambiguous.*

Task-Oriented Metrics. While linguistic quality is a prerequisite for usable QA, our primary focus is on task-oriented metrics that capture alignment and mentorship effectiveness.

- **QA Alignment:** *The answer directly corresponds to what the question asks. The question–answer pair demonstrates proper alignment, where the answer adequately satisfies the question.*
- **Question Mentorship:** *The question provides learning, guidance, advice, or insights that would benefit from mentor expertise.*
- **Answer Mentorship:** *The answer provides guidance, wisdom, practical advice, or insights that help the reader learn or grow.*

Together, these metrics capture linguistic quality, semantic alignment, and mentorship effectiveness—dimensions not addressed by existing QA evaluation frameworks.

Rating Scale. All metrics are rated on a 5-point Likert scale (1: strongly disagree to 5: strongly agree). Pilot studies showed that the 3-point scale used in prior work (Fu et al., 2024) was insufficient to capture nuanced distinctions in metric quality and resulted in lower inter-annotator agreement.

5.2 Human Evaluation

We recruit 12 annotators, organized into four language-specific groups corresponding to the dataset languages, with three native or fluent speakers per language. Annotators evaluate a balanced

sample of QA pairs drawn from five videos per language and evenly distributed across all four QA-generation models. This results in 240 unique QA pairs, each receiving three independent ratings. An example of the annotation interface is provided in Appendix D.

Inter-Annotator Agreement. We measure annotation agreement using Gwet’s AC2, computed separately for each evaluation metric across languages and models. We choose Gwet’s AC2 over Krippendorff’s alpha, which has been shown to produce misleading agreement estimates in skewed rating distributions common in AI system evaluation (Battisti and Ebling, 2024; Pradhan et al., 2025).

5.3 Automatic Evaluation

To study the reliability of automatic evaluation for mentorship-focused QA, we evaluate all QA pairs using nine multilingual LLM judges. Recent work shows that LLM-based evaluators outperform traditional NLG metrics and provide a scalable alternative to human evaluation (Fu et al., 2024; Kocmi and Federmann, 2023).

We select nine open-source or publicly documented multilingual models to ensure reproducibility: Qwen2.5-7B (Yang et al., 2025), Command-R-7B (Cohere and Aakanksha, 2025), Mistral-8B-Instruct (Mistral AI, 2025), Llama3.2-11B (Meta, 2024), Qwen3-8B (Team, 2025), Gemma3-12B-IT (Google, 2025), Pangea-7B (Yue et al., 2024), AyaExpand-8B (Dang, 2024), and Phi-4-15B (Microsoft, 2025).

LLM-Human Agreement. We assess LLM judge effectiveness by computing Gwet’s AC2 between LLM-based scores and human annotations for each evaluation metric, across languages and QA-generation models. This analysis allows us to quantify when automatic evaluation aligns with human judgments in mentorship-focused QA.

6 Evaluation Results

We evaluate mentorship-focused QA generation along four complementary axes: (1) the reliability of mentorship-focused evaluation, (2) the scalability of mentorship-focused evaluation, (3) the effectiveness of different QA-generation models, and (4) the consistency of these effects across languages and topics.

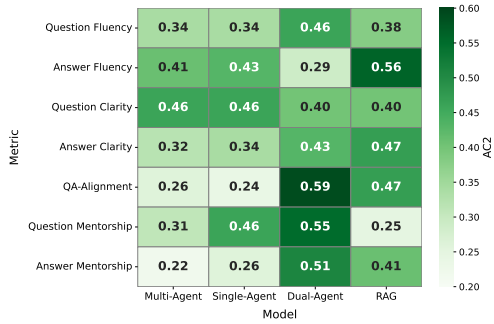


Figure 2: Inter-Annotator Agreement Scores.

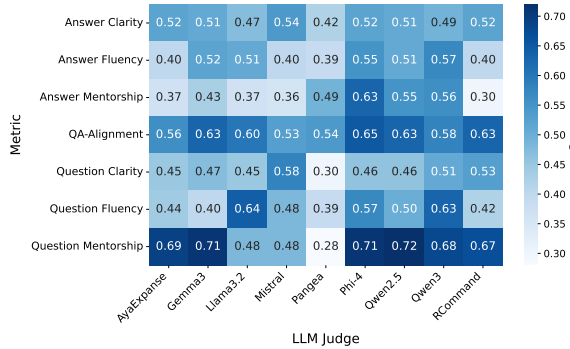


Figure 3: LLM-Human agreement scores. Qwen models and Phi-4 correlate the most with human results.

Reliability of Mentorship-Focused Evaluation.

We assess the reliability of the proposed evaluation dimensions using human annotations. As expected, inter-annotator agreement is highest for linguistic metrics—*Fluency* (Q: 0.38, A: 0.42) and *Clarity* (Q: 0.43, A: 0.39)—which capture surface-level language properties. Importantly, task-oriented metrics—*QA Alignment* (0.39) and *Mentorship* (Q: 0.39, A: 0.35)—achieve stable, moderate agreement, reflecting the inherently interpretive nature of mentorship quality.

Agreement varies across models and languages. Among models, *Dual-Agent* shows the highest agreement (0.46), followed by RAG (0.42), Single-Agent (0.36), and Multi-Agent (0.33). Across languages, agreement is highest for Romanian (0.67) and lower for Chinese (0.37), English (0.28), and Hindi (0.26) (more in Appendix D).

Overall, these patterns indicate that mentorship-focused evaluation is sensitive to both model behavior and linguistic context, mirroring real-world educational assessment and motivating future work on culturally aware and mentorship-specific evaluation frameworks.

Human-LLM Agreement in Automatic Evaluation. To assess the scalability of mentorship-focused evaluation, we compare LLM-based

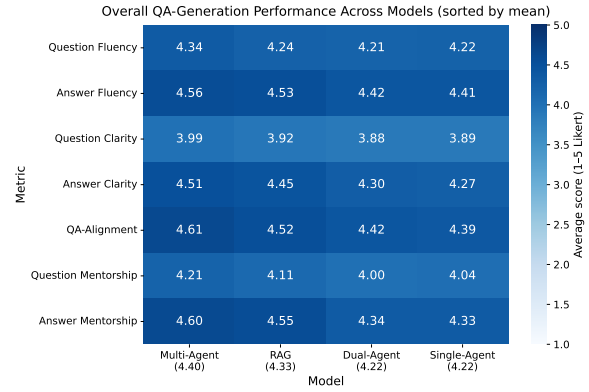


Figure 4: QA-generation performance across models and metrics (1–5 Likert). Models are sorted by mean.

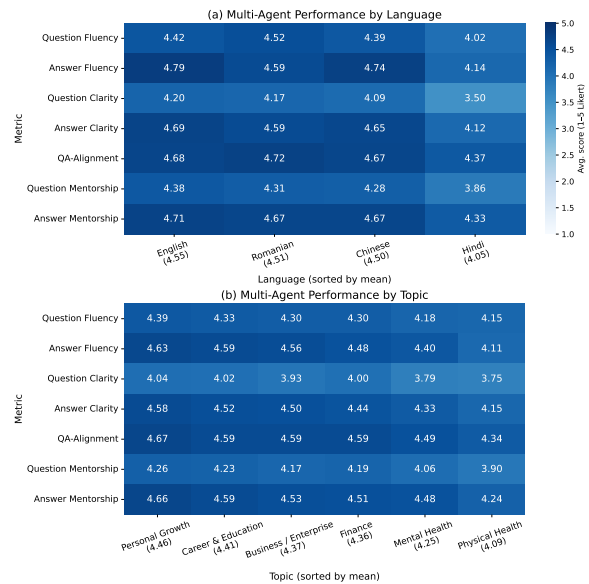


Figure 5: Multi-Agent performance across (a) language and (b) topic. Topics and languages are sorted by mean.

judges with human annotations (Figure 3). Agreement is highest for task-oriented metrics, with several judges (Qwen2.5, Qwen3, and Phi-4) reaching substantial agreement with human ratings (0.55–0.72), while linguistic metrics show moderate but stable alignment (0.40–0.57).

Qwen-family models (Qwen2.5, Qwen3) and Phi-4 consistently achieve the strongest agreement across metrics and languages, including lower-resource settings (results per language in Appendix D). At the same time, agreement varies substantially across judges, highlighting the challenge of automatically evaluating mentorship quality and motivating language- and metric-aware evaluation with multiple LLM judges rather than reliance on a single evaluator.

Comparison of QA-Generation Architectures. The Multi-Agent model achieves the strongest over-

all performance, ranking highest across all evaluation dimensions and obtaining the highest mean score (4.40), as shown in Figure 4. RAG follows closely (4.33), offering consistent improvements over the Single- and Dual-Agent baselines, particularly on *QA-Alignment*, *Question Mentorship*, and *Answer Mentorship*. In contrast, Single-Agent and Dual-Agent models perform nearly identically (4.22), suggesting that topic-aware chunking alone yields limited gains. Performance differences are most pronounced for *Mentorship* and *QA-Alignment* metrics, while *Fluency* remains relatively similar across models, highlighting the importance of agentic models for preserving mentorship quality. For detailed cross model and LLM-Judge results see Appendix D.2.

Qualitative comparisons of model outputs across different *Mentorship* scores are presented in Appendix D.2.

Analysis by Language and Topic. We analyze performance across languages and topics using the best-performing Multi-Agent model, with scores averaged over all LLM judges; cross model comparisons are reported in Appendix D.

Language Evaluation. As shown in Figure 5(a), *English achieves the highest overall performance (4.55)*, with *Romanian (4.51)* and *Chinese (4.50)* performing comparably across key metrics such as *QA-Alignment* and *Answer Mentorship*. In contrast, *Hindi scores lower overall (4.05)*, with the largest gaps in *Question Clarity (3.50)* and *Question Mentorship (3.86)*, consistent with known tokenization and training limitations for Devanagari script. This aligns with MMLU results, where Hindi scores substantially lower than English (large-traversaal, 2025). Importantly, the Multi-Agent architecture substantially narrows this gap: compared to Single-Agent performance (*QA-Alignment 3.96*, *Answer Fluency 3.72*), *Multi-Agent improves Hindi scores to 4.37 and 4.14*, respectively. While disparities remain, these gains suggest that *agentic coordination can partially mitigate language-specific errors*.

Topic Evaluation. As shown in Figure 5(b), QA-generation performance varies substantially by topic. *Personal Growth achieves the strongest overall performance (4.46)*, leading across metrics such as *QA-Alignment (4.67)* and *Answer Mentorship (4.66)*. *Career & Education (4.41)* and *Business/Enterprise (4.37)* follow closely, while *Physical Health exhibits the lowest overall scores (4.09)*,

with notably lower *Question Clarity (3.75)* and *Question Mentorship (3.90)*, reflecting the challenges posed by clinically dense language. In contrast, *Finance performs robustly*, achieving high *QA-Alignment (4.47)*, likely due to its structured discourse. Overall, these results suggest that *agentic coordination is most beneficial for technically complex mentorship topics*, while offering limited gains for domains well represented in pre-training data.

Takeaways. Our results yield four key findings: (1) Task-oriented metrics are reliable yet inherently more subjective than linguistic metrics. (2) LLM-evaluation must be language- and metric-aware, motivating multi-judge strategies. (3) Multi-Agent models consistently outperforms simpler architectures on mentorship and alignment dimensions. (4) Agentic coordination is particularly beneficial for complex topics and lower-resource languages.

7 Conclusion

We introduced MENTORQA, the first multilingual dataset and evaluation framework for mentorship-focused question answering from long-form videos. By treating mentorship as a distinct QA objective—prioritizing guidance, clarity, and learning value beyond factual correctness—we find limitations in existing benchmarks, particularly in multilingual and long-context settings. Through a controlled comparison of Single-Agent, Dual-Agent, RAG, and Multi-Agent pipelines, we show that agentic coordination consistently improves mentorship quality, with the largest gains for complex topics and lower-resource languages, where it acts as a partial performance equalizer. Our analysis of automated evaluation further reveals substantial variation in alignment with human judgments across languages and metrics, highlighting the need for task-aware and language-sensitive evaluation. Together, these findings establish mentorship-focused QA as a rich and novel research direction and provide a multilingual benchmark for studying model architectures and evaluation methods in educational AI. We hope our work encourages the community to move beyond factual-driven evaluation toward models and metrics that better support learning, reflection, and guidance across cultures and languages.

632 **Limitations**

633 Our work represents an initial step toward
634 mentorship-focused question answering. While it
635 provides a comprehensive evaluation of QA archi-
636 tectures, several limitations remain.

637 **Generality of Findings.** Our study focuses on
638 mentorship-focused QA derived from long-form
639 videos across four languages and a limited set of
640 domains. This design enables controlled compar-
641 ison of architectures and evaluation dimensions,
642 but it does not capture the full diversity of mentor-
643 ship contexts, interaction styles, or cultural norms
644 found in real-world settings. In particular, our find-
645 ings may not transfer to highly interactive mentor-
646 ship scenarios or to languages with substantially
647 different resource profiles. Rather than aiming
648 for universal coverage, our goal is to establish
649 mentorship-focused QA as a distinct problem and
650 to provide a benchmark that supports further eval-
651 uation and analysis. We encourage future work
652 to extend our open-source framework to additional
653 languages and domains.

654 **A Fixed Base Model.** To isolate the effects of archi-
655 tectural and agentic design, we fix the under-
656 lying base language model (Qwen 2.5) across all
657 experiments. This choice enables fair comparison
658 but limits insight into how multi-agent pipelines in-
659 teract with other foundation models. Newer mul-
660 tilingual or instruction-tuned models may exhibit
661 different trade-offs or amplify agentic gains, which
662 we invite future work to explore.

663 **Evaluation and Scalability.** Our evaluation com-
664 bines human judgments with LLM-based judges to
665 assess mentorship quality across multiple dimen-
666 sions. Although LLM judges provide scalable sig-
667 nal, we observe notable variation in their alignment
668 with human judgments across languages and met-
669 rics. Accordingly, automated evaluation should
670 be viewed as complementary rather than a substi-
671 tute for human assessment. Multi-agent pipelines
672 also incur higher computational cost than single-
673 agent or RAG-based approaches, which may limit
674 deployment in resource-constrained settings. Our
675 analysis prioritizes comparative effectiveness, and
676 we see efficiency-aware agentic designs as an im-
677 portant direction for future work.

678 **Subjectivity of Mentorship Evaluation.** Mentor-
679 ship quality is inherently subjective, as reflected in
680 the moderate inter-annotator agreement observed
681 for mentorship-focused metrics. Annotators may
682 reasonably prioritize different aspects of guidance,

683 particularly across cultural and linguistic contexts.
684 While our evaluation captures this variability, it
685 does not explicitly model disagreement or multiple
686 valid interpretations of mentorship quality. Devel-
687 oping evaluation protocols that better account for
688 such diversity remains an open challenge.

689 **Downstream Educational Impact.** Our evalua-
690 tion focuses on mentorship QA quality rather than
691 downstream educational outcomes. Future work
692 could examine how mentorship-focused QA sup-
693 ports comprehension, reflection, and skill develop-
694 ment in real educational settings.

695 **Ethical Considerations**

696 Mentorship content often involves personal experi-
697 ences, sensitive advice, and implicit power dynam-
698 ics. Although all videos used in MENTORQA are
699 publicly available, we emphasize that mentorship
700 quality is not equivalent to factual authority, and
701 generated QA pairs should not be treated as pre-
702 scriptive guidance without human oversight. Au-
703 tomated systems may overgeneralize advice, miss
704 contextual nuance, or reflect cultural biases present
705 in training data.

706 We also note ethical risks in automated evalua-
707 tion. LLM-based judges may encode language-
708 or culture-specific preferences, which can disad-
709 vantage certain communication styles or commu-
710 nities. Our findings underscore the continued im-
711 portance of human evaluation, transparent metric
712 design, and cautious interpretation of automatic
713 scores, particularly in multilingual and educational
714 settings.

715 We release MENTORQA as a research resource
716 rather than a deployable mentorship system, and
717 encourage future work to incorporate participatory
718 evaluation, educator-in-the-loop validation, and
719 safeguards that prioritize learner well-being.

720 **Anonymization and Personal Information**

721 MENTORQA is constructed from publicly available
722 mentorship videos. While the source content is
723 public, we take steps to minimize the presence and
724 propagation of personally identifiable information
725 (PII) in the released dataset. We process only
726 automatically generated transcripts and do not
727 release raw audio or video. The released dataset
728 contains only de-identified text and focuses on
729 general mentorship themes, guidance, and reflect-
730 ive insights rather than personal narratives tied
731 to specific individuals. We follow established
732 best practices for ethical data use from publicly

733	available sources and release the dataset for	11783–11803, Miami, Florida, USA. Association for	786
734	research purposes only.	Computational Linguistics.	787
735	References		
736	Andrew Aquilina, Sean Diacono, Panagiotis Papa-	Google. 2025. Gemma-3-12b-it. https://huggingface.co/google/gemma-3-12b-it .	788
737	petrou, and Maria Movin. 2023. An end-to-end	Multilingual large	789
738	workflow using topic segmentation and text sum-	vision-language model (Italian emphasis), hosted on	790
739	marisation methods for improved podcast compre-	Hugging Face.	791
740	hension . <i>arXiv preprint</i> .		
741	Akari Asai, Xinyan Yu, Jungo Kasai, and Hannaneh	Shasha Guo, Lizi Liao, Cuiping Li, and Tat-Seng Chua.	792
742	Hajishirzi. 2021. One question answering model for	2024. A survey on neural question generation: Meth-	793
743	many languages with cross-lingual dense passage re-	ods, applications, and prospects . In <i>Proceedings</i>	794
744	trieval . <i>arXiv preprint</i> .	of the Thirty-Third International Joint Conference	795
745		on Artificial Intelligence (IJCAI 2024) , Survey Track,	796
746	Alessia Battisti and Sarah Ebling. 2024. Automatic	pages 8038–8046, Macau, China.	797
747	annotation elaboration as feedback to sign language		
748	learners . In <i>Proceedings of the 18th Linguistic Anno-</i>	Tom Kocmi and Christian Federmann. 2023. Large lan-	798
749	tation Workshop (LAW-XVIII) , pages 46–60, St. Ju-	guage models are state-of-the-art evaluators of trans-	799
750	lians, Malta. Association for Computational Linguis-	lation quality . In <i>Proceedings of the 24th Annual</i>	800
751		Conference of the European Association for Machine	801
752	Jianlyu Chen, Shitao Xiao, Peitian Zhang, Kun	Translation , pages 193–203, Tampere, Finland. Euro-	802
753	Luo, Defu Lian, and Zheng Liu. 2024. M3-	pean Association for Machine Translation.	803
754	embedding: Multi-linguality, multi-functionality,		
755	multi-granularity text embeddings through self-	Géza Kovács. 2016. Effects of in-video quizzes on	804
756	knowledge distillation . In <i>Findings of the Asso-</i>	mooC lecture viewing . In <i>Proceedings of the Third</i>	805
757	ciation for Computational Linguistics: ACL 2024 ,	<i>(2016) ACM Conference on Learning@Scale</i> , pages	806
758	pages 2318–2335, Bangkok, Thailand. Association	31–40, Edinburgh, UK. Association for Computing	807
759	for Computational Linguistics.	Machinery.	808
760		large-traversaal. 2025. Qwen-2.5-14b-hindi.	809
761	Mao-Siang Chen and An-Zi Yen. 2024. E-qgen: Ed-	https://huggingface.co/large-traversaal/Qwen-2.	810
762	ucational lecture abstract-based question generation	5-14B-Hindi. Large language model fine-tuned for	811
763	system . In <i>Proceedings of the Thirty-Third Interna-</i>	Hindi, hosted on Hugging Face.	812
764	tional Joint Conference on Artificial Intelligence (IJ-		
765	CAI 24) , Demonstrations Track, pages 8631–8634,	Lin Ma and Yuchun Ma. 2019. Automatic question gen-	813
766	Macau, China. IJCAI/AAAI.	eration based on mooc video subtitles and knowledge	814
767		graph . In <i>Proceedings of the 2019 7th International</i>	815
768	Team Cohere and Aakanksha. 2025. Command a: An	Conference on Information and Education Technol-	816
769	enterprise-ready large language model .	ogy (ICIET 2019) , pages 49–53, Aizu-Wakamatsu,	817
770		Japan. Association for Computing Machinery.	818
771	Team Cohere & John Dang. 2024. Aya expande: Com-	Wenjun Meng, Yuzhe Li, Lili Chen, and Zhaomin Dong.	819
772	bining research breakthroughs for a new multilingual	2025. Using the retrieval-augmented generation to	820
773	frontier .	improve the question-answering system in human	821
774		health risk assessment: The development and appli-	822
775	Jing Du, Guangtao Xu, Wenhao Liu, Dibin Zhou, and	cation . <i>Electronics</i> , 14(2):386.	823
776	Fuchang Liu. 2025. Enhancing online learning		
777	through multi-agent debates for cs university stu-	Meta. 2024. Llama-3.2-11b-vision. https://	824
778	dents . <i>Applied Sciences</i> , 15(11):5877.	huggingface.co/meta-llama/Llama-3.2-11B-Vision .	825
779		Multimodal vision language model with text+image	826
780	A. Elaryan. 2022. Question-answering for segment	inputs, released under the Llama 3.2 Community	827
781	retrieval on podcast recordings . Master’s thesis, Uni-	License; optimized for visual recognition, reasoning,	828
782	versity of Oklahoma. Master’s thesis.	and captioning.	829
783		Microsoft. 2025. Phi-4. https://huggingface.co/	830
784	Said Al Faraby, Adiwijaya Adiwijaya, and Ade Romad-	microsoft/phi-4 . Large language model (Phi-4)	831
785	hony. 2024. Review on neural question generation	hosted on Hugging Face.	832
	for education purposes . <i>International Journal of Ar-</i>	Mistral AI. 2025. Ministral-8b-instruct-	833
	tificial Intelligence in Education , 34:1008–1045.	2410. https://huggingface.co/mistralai/	834
		Ministral-8B-Instruct-2410. Language model,	835
		instruct-tuned, released under the Mistral Research	836
		License.	837
		Alireza Mohammadshahi, Thomas Scialom, Majid Yaz-	838
		dani, Pouya Yanki, Angela Fan, James Henderson,	839
		and Marzieh Saeidi. 2023. Rqge: Reference-free	840

841	metric for evaluating question generation by answering the question. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 6845–6867, Toronto, Canada. Association for Computational Linguistics.	
842		
843		
844		
845		
846	Antonio Moreno-Cediel, Jesús Ángel del Hoyo-Gabaldón, Eva García-López, Antonio García-Cabot, and David de Fitero-Domínguez. 2024. Evaluating the performance of multilingual models in answer extraction and question generation. <i>Scientific Reports</i> , 14:15477.	
847		
848		
849		
850		
851		
852	Nikahat Mulla and Prachi Gharpure. 2023. Automatic question generation: a review of methodologies, datasets, evaluation metrics, and applications. <i>Progress in Artificial Intelligence</i> , 12(1):1–32.	
853		
854		
855		
856	Preksha Nema and Mitesh M. Khapra. 2018. Towards a better metric for evaluating question generation systems. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3950–3959, Brussels, Belgium. Association for Computational Linguistics.	
857		
858		
859		
860		
861		
862	Anu Pradhan, Alexandra Ortan, Apurv Verma, and Madhavan Seshadri. 2025. Llm-as-a-judge: Rapid evaluation of legal document recommendation for retrieval-augmented generation.	
863		
864		
865		
866	Archiki Prasad, Trung Bui, Seunghyun Yoon, Hanieh Deilamsalehy, Franck Dérnoncourt, and Mohit Bansal. 2023. MEETINGQA: Extractive question-answering on meeting transcripts. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15000–15025, Toronto, Canada. Association for Computational Linguistics.	
867		
868		
869		
870		
871		
872		
873		
874	QwenLM Community. 2025. Improving qwen3 tokenizer efficiency for multilingual use case (issue #1400). https://github.com/QwenLM/Qwen3/issues/1400 . GitHub issue discussing tokenizer inefficiency in Qwen3 for multilingual scenarios including Hindi, Italian, and German.	
875		
876		
877		
878		
879		
880	Sebastian Ruder and Avi Sil. 2021. Multi-domain multilingual question answering. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: Tutorial Abstracts</i> , pages 17–21, Punta Cana, Dominican Republic & Online. Association for Computational Linguistics.	
881		
882		
883		
884		
885		
886	S. Saadaoui, Y. Li, M. Zimmer, and L. Qin et al. 2025. Coordinated llm multi-agent systems for collaborative question answering. <i>Knowledge-Based Systems</i> , 302:109827.	
887		
888		
889		
890	Mohammad Sammoudi, Ahmad Habaybeh, Huthaifa I. Ashqar, and Mohammed Elhenawy. 2025. Question-answering (qa) model for a personalized learning assistant for arabic language. In <i>Intelligent Systems, Blockchain, and Communication Technologies (ISBCom 2024), Lecture Notes in Networks and Systems, vol. 1268</i> , pages 356–367. Springer, Cham.	
891		
892		
893		
894		
895		
896		
	Qwen Team. 2025. Qwen3 technical report.	897
	Khanh-Tung Tran, Dung Dao, Minh-Duong Nguyen, Quoc-Viet Pham, Barry O’ Sullivan, and Hoang D. Nguyen. 2025. Multi-agent collaboration mechanisms: A survey of llms. <i>arXiv preprint</i> .	898 899 900 901
	Asahi Ushio, Fernando Alva-Manchego, and Jose Camacho-Collados. 2023. A practical toolkit for multilingual question and answer generation. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)</i> , pages 86–94, Toronto, Canada. Association for Computational Linguistics.	902 903 904 905 906 907 908
	Kesen Wang, Daulet Toibazar, Abdulrahman Alfulayt, Abdulaziz S. Albadawi, Ranya A. Alkahtani, Asma A. Ibrahim, Haneen A. Alhomoud, Sherif Mohamed, and Pedro J. Moreno. 2025. Multi-agent interactive question generation framework for long document understanding. <i>arXiv Preprint</i> .	909 910 911 912 913 914
	Cameron R. Wolfe. 2024. Llm-as-a-judge. Substack article.	915 916
	An Yang, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoyan Huang, Jiandong Jiang, Jianhong Tu, Jianwei Zhang, Jingren Zhou, Junyang Lin, Kai Dang, Kexin Yang, Le Yu, Mei Li, Minmin Sun, Qin Zhu, Rui Men, Tao He, Weijia Xu, Wenbiao Yin, Wenyuan Yu, Xiafei Qiu, Xingzhang Ren, Xinlong Yang, Yong Li, Zhiying Xu, and Zipeng Zhang. 2025. Qwen2.5-1m technical report. <i>arXiv preprint</i> .	917 918 919 920 921 922 923 924
	Mengxia Yu, Bang Nguyen, Olivia Zino, and Meng Jiang. 2025. Context selection and rewriting for video-based educational question generation. <i>arXiv Preprint</i> .	925 926 927 928
	Xiang Yue, Yueqi Song, Akari Asai, Seungone Kim, Jean de Dieu Nyandwi, Simran Khanuja, Anjali Kantharuban, Lintang Sutawika, Sathyanarayanan Ramamoorthy, and Graham Neubig. 2024. Pangea: A fully open multilingual multimodal llm for 39 languages. <i>arXiv preprint arXiv:2410.16153</i> .	929 930 931 932 933 934
	Jie Zhu, Junhui Li, Yalong Wen, Xiandong Li, Lifan Guo, and Feng Chen. 2025. M ³ finmeeting: A multilingual, multi-sector, and multi-task financial meeting understanding evaluation dataset. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 244–266, Vienna, Austria. Association for Computational Linguistics.	935 936 937 938 939 940 941
	A Data Collection	942
	B Video Sources	943
	To see the full detailed source of all videos, see Table 3 and Table 4.	944 945
	Topic Selection across Languages.	946

947	Topic Annotation. To ensure accurate topic categorization, each QA pair in the dataset was manually annotated and verified by human annotators, establishing ground truth topic labels for all evaluations.	
948		
949		
950		
951		
952	Videos are selected to cover six mentorship-related topics: <i>Entrepreneurship, Education, Finance, Mental Health, Personal Growth, and Physical Health</i> . Topic labels are automatically assigned using <i>Qwen2.5-7B</i> (Yang et al., 2025) and subsequently verified by expert annotators with proficiency in the respective languages.	
953		
954		
955		
956		
957		
958		
959	English. English videos are collected from the ACL Mentorship channel ⁹ . ACL Mentorship is a Year-Round Mentorship initiative of the Association for Computational Linguistics (ACL), aimed at supporting students and early-career researchers entering the field of NLP. The channel hosts virtual mentorship sessions featuring mentors from academia and industry, covering topics such as choosing a research direction, writing papers, pursuing a PhD, and maintaining work–life balance. The initiative promotes equal access to career guidance and currently supports a global community of over 2,500 members.	
960		
961		
962		
963		
964		
965		
966		
967		
968		
969		
970		
971		
972	Romanian. Romanian videos are collected from the MindArchitect channel ¹⁰ , which features conversations with coaches, educators, athletes, and other specialists. Topics span education and career development (e.g., “ <i>How do we prepare young people for the labor market?</i> ”), as well as personal growth and mental health (e.g., “ <i>From perfectionism and self-criticism to the joy of play and mindful creation</i> ”).	
973		
974		
975		
976		
977		
978		
979		
980		
981	Chinese. The Chinese videos are collected from multiple public channels rather than a single mentorship program. We curate videos that exhibit strong mentorship characteristics across mental health, career growth, finance, and physical health. The selected content primarily consists of podcast-style conversations, discussions, and interview-based dialogues involving two or three participants, typically led by domain professionals. These discussions emphasize experience-based guidance, reflective reasoning, and practical advice drawn from real world practice. Although the videos originate from diverse creators, they are unified by their focus on mentorship-focused discourse and sustained explanatory depth.	
982		
983		
984		
985		
986		
987		
988		
989		
990		
991		
992		
993		
994		
995		
		Hindi. Hindi videos are collected from long-form interview and podcast recordings that feature in-depth conversations with domain experts, industry practitioners, entrepreneurs, and high-performing competitive exam candidates. These interactions emphasize experiential knowledge, preparation methodologies, career decision-making, and reflective discussions on challenges and learnings. The conversational format allows speakers to elaborate on nuanced perspectives and real-world insights, making the content particularly rich for extracting mentorship-focused question–answer pairs.
		996
		997
		998
		999
		1000
		1001
		1002
		1003
		1004
		1005
		1006
		1007
		1008

⁹<https://www.youtube.com/@aclmentorship>

¹⁰<https://www.youtube.com/@MindArchitect>

1009	C Methodology		
1010	Implementation Details.	All models use the Py-	
1011		Torch and Hugging Face Transformers stacks, and	
1012		are run in inference-only mode without any fine-	
1013		tuning. Inference is executed on a GPU node	
1014		equipped with four NVIDIA L40S GPUs, each	
1015		with 48 GB of memory, which allows us to host	
1016		the model once and parallelize calls from different	
1017		videos, segments, or agents across devices.	
1018	C.1 Prompt Design		
	D Evaluation		1019
		Figure 10 shows the interface used in the human	1020
		annotation process.	1021
	D.1 Can the metrics accurately assess the		1022
	generated QA pairs?		1023
	D.1.1 Inter-Annotator Agreement		1024
		Figure 11 reveals substantial variation in inter-	1025
		annotator agreement across languages. Roma-	1026
		nian exhibits the highest and most consistent	1027
		agreement across metrics (0.49–0.81), demon-	1028
		strating strong to near-perfect reliability in evaluat-	1029
		ing mentorship-focused quality. Chinese displays	1030
		moderate to fair agreement (0.27–0.47), with	1031
		metrics like <i>Answer Clarity</i> and <i>QA-Alignment</i>	1032
		achieving reasonable inter-annotator consensus.	1033
		English and Hindi show lower but still fair	1034
		agreement (0.13–0.37 and 0.20–0.34 respec-	1035
		tively), with Hindi’s challenges potentially stem-	1036
		ming from tokenization issues and weaker model	1037
		support for Devanagari script, making consistent	1038
		evaluation more difficult even for human annota-	1039
		tors. Overall, these results suggest that annota-	1040
		tion reliability is highly language-dependent , re-	1041
		inforcing the need for culturally aware evaluation	1042
		frameworks when designing mentorship-focused	1043
		QA assessment protocols.	1044
	D.1.2 LLM-Human Agreement		1045
		LLM-Human Cohen kappa agreement scores per	1046
		language are shown in Figure 12	1047
	D.2 How do QA-generation models perform		1048
	across metrics?		1049
		All model performance across LLM-Judges is	1050
		shown in Figure 13.	1051
		Qualitative results per model are shown below:	1052
	Single-Agent Q:	What is the most important cri-	1053
		terion when reviewing PhD applications?	1054
	A:	Passion is the most important criterion when re-	1055
		viewing PhD applications.	1056
	Scores:	QFlu 4, AFlu 5, QClar 5, AClar 1, QA-	1057
		Algn 5, QMen 4, AMen 1	1058
	Dual-Agent Q:	Why is it important to interact	1059
		with students according to the speaker in the tran-	1060
		script?	1061
	A:	Interacting with students is one of the most re-	1062
		warding aspects of the academic path, as it allows	1063
		the speaker to work with brilliant PhD students and	1064

Model Type	Prompt
Single-Agent	<p>System: You are an expert at educational content analysis who can identify key questions and provide clear answers.</p> <p>Prompt: Read the following transcript carefully and identify the 20 most important questions discussed in this content. For each question:</p> <ol style="list-style-type: none"> 1. Ensure the question captures a key concept or important information from the transcript 2. Provide a clear, accurate answer to the question based only on information in the transcript 3. Make sure questions and answers cover different aspects of the content the whole transcript and don't overlap significantly 4. Select questions and answers in a balanced way from throughout the entire content, not concentrating too heavily on any single section or part 5. Answers should be in proper detail length and include only the relevant information answering the question properly with educational/mentorship value. 6. Use the same Language as of the original content. <p>Strictly Format your response as a list of question-answer pairs, with each pair clearly marked (e.g., "Question 1:", "Answer 1:"). Strictly to make your response as structured as possible so it can be easily parsed. Also avoid any other extra words in the start and beginning and only the strict structured response.</p> <p>Transcript: {transcript}</p>

Figure 6: Single-Agent prompt.

1065 other students across various levels, including un- 1098
1066 dergraduates, masters, and postdoctoral fellows. 1099
1067 **Scores:** QFlu 5, AFlu 5, QClar 4, AClar 4, QA- 1100
1068 Algn 5, QMen 4, AMen 3

1069 **RAG Q:** How do PhD admission processes dif- 1101
1070 fer internationally? 1102

1071 **A:** PhD admission processes vary internationally 1103
1072 in terms of interview conduct and committee struc- 1104
1073 tures. For instance, MIT has a stringent committee 1105
1074 system where individual professors' admissions 1106
1075 decisions are heavily scrutinized, whereas Stanford 1107
1076 uses a rolling admissions process. Universities like 1108
1077 UIUC or Hong Kong allow faculty members signif- 1109
1078 icant autonomy in admitting PhD students. This 1110
1079 diversity highlights the importance of considering 1111
1080 a broader range of institutions beyond prestigious 1112
1081 American universities. 1113

1082 **Scores:** QFlu 5, AFlu 5, QClar 4, AClar 3, QA- 1114
1083 Algn 5, QMen 5, AMen 3 1115

1084 **Multi-Agent Q:** What differences exist between 1116
1085 PhD admission processes in various academic in- 1117
1086 stitutions? 1118

1087 **A:** "The differences between PhD admission pro- 1119
1088 cesses in various academic institutions primarily 1120
1089 revolve around the structure of the selection com- 1121
1090 mittees, the role of individual faculty members, 1122
1091 and the criteria used for evaluating applications. 1123
1092 For instance: 1124

1093 1. ****Committee Structure**:** At some institu- 1125
1094 tions, like MIT, there is a strict committee process 1126
1095 where individual professors cannot independently 1127
1096 admit students without consensus from the com- 1128
1097 mittee. This contrasts with universities like UIUC

1098 or Hong Kong, where faculty members have more 1099
1100 autonomy to select students based on their own 1101
1102 preferences. 1103

2. ****Faculty Autonomy**:** In some systems, 1104
1105 such as those in the U.S., faculty members play a 1106
1107 significant role in admitting PhD students. They 1108
1109 may have the authority to interview and accept can- 1110
1111 didates directly, which can vary widely depending 1111
1112 on the institution and department. 1113

3. ****Application Review Process**:** The num- 1114
1115 ber of applications reviewed by each committee 1116
1117 member differs. Some may skim through a large 1117
1118 volume of applications, while others may focus 1118
1119 on specific candidates they have been contacted 1119
1120 by. For example, one panelist mentioned receiv- 1120
1121 ing numerous unsolicited emails, indicating that 1121
1122 some professors may not read every application 1122
1123 thoroughly. 1123

4. ****Criteria for Selection**:** While the core 1124
1125 process involves evaluating academic records, let- 1125
1126 ters of recommendation, and research proposals, 1126
1127 the emphasis on these factors can differ. For ex- 1127
1128 ample, a close match in research interests might 1128
1129 be more beneficial for some institutions, especially 1129
1130 smaller programs, compared to larger committees 1130
1131 that consider a broader range of criteria. 1131

1132 Understanding these differences can help PhD 1132
1133 aspirants tailor their applications and outreach 1133
1134 strategies to align better with the specific require- 1134
1135 ments and preferences of the institutions they are 1135
1136 targeting." 1136

1137 **Scores:** QFlu 5, AFlu 5, QClar 5, AClar 5, QA- 1129
1138 Algn 5, QMen 5, AMen 5 1130

Model Type	Prompt
Dual-Agent	<p>System: You are an expert at analyzing transcripts and segmenting them by topic. The transcript has been split into numbered lines where each line represents a complete thought. Identify topic boundaries and assign concise topic titles.</p> <p>Prompt: Output a JSON list of dictionaries with these keys: - "topic": Concise descriptive title (3-7 words) - "start_line": First line number of this topic - "end_line": Last line number of this topic</p> <p>Rules: 1. Topics must cover consecutive line numbers 2. Entire transcript must be covered without gaps or overlaps 3. The first topic must start at line 1 4. The last topic must end at line {total_lines} 5. Use line numbers exactly as provided 6. Output ONLY the JSON with no additional text</p> <p>Numbered Transcript: {numbered_transcript}</p>
	<p>System: You are an expert at educational content analysis who can identify key questions and provide clear answers.</p> <p>Prompt: Read the following transcript carefully and identify the {num_questions} most important questions discussed in this content. For each question:</p> <ol style="list-style-type: none"> 1. Ensure the question captures a key concept or important information from the transcript chunk 2. Provide a clear, accurate answer to the question based only on information in this chunk 3. Make sure questions and answers cover different aspects of the content the whole chunk and don't overlap significantly 4. Select questions and answers in a balanced way from throughout the entire content, not concentrating too heavily on any single section. 5. Answers should be in proper detail length and include only the relevant information answering the question properly with educational/mentorship value. 6. Use the same Language as of the original content. <p>Strictly Format your response as a list of question-answer pairs, with each pair clearly marked (e.g., "Question 1:", "Answer 1:"). Strictly to make your response as structured as possible so it can be easily parsed. Also avoid any other extra words in the start and beginning and only the strict structured response.</p> <p>Transcript: {transcript}</p>

Figure 7: Dual-Agent prompt.

D.3 How does QA-generation performance vary across languages and mentorship topics?

Language Evaluation. As shown in Figure 14, English demonstrates the strongest performance across all metrics, which is expected given its prevalence in training data. Romanian and Chinese follow closely, performing comparably to English on critical dimensions like QA-Alignment and Mentorship metrics. However, Hindi exhibits a notable performance drop across all evaluated metrics. This disparity can be attributed to architectural and data constraints in the Qwen2.5 model. The model employs a Byte-Level Byte Pair Encoding (BBPE) tokenizer that is specifically optimized for English and Chinese (Yang et al., 2025). When processing Hindi text in Devanagari script, the tokenizer fails to recognize com-

plete word units and instead fragments individual words into numerous byte-level tokens. This sub-optimal tokenization, combined with limited Hindi representation in the training corpus, directly impacts the model's ability (QwenLM Community, 2025) to generate and evaluate Hindi content effectively. This interpretation is further supported by the model's reported MMLU benchmark scores: English achieves 74.37 while Hindi scores only 52.16 (large-traversaal, 2025), reflecting a substantial performance gap. The combination of tokenization inefficiency and insufficient training exposure explains the observed quality degradation in Hindi-generated QA pairs.

The Multi-Agent model demonstrates a consistent performance advantage, effectively raising the ceiling for every language compared to the other three architectures. This improvement is

Model Type	Prompt
RAG	<p>System: You are an expert content analyst. Your task is to read a long transcript and identify potential questions that cover the most important educational and mentorship-related topics discussed.</p> <p>Prompt: Based on the following transcript, generate a list of exactly 20 diverse and high-value questions.</p> <p>Guidelines for questions:</p> <ul style="list-style-type: none"> - Focus on key concepts, advice, and actionable insights. - Ensure questions span the entire transcript, from beginning to end. - Avoid trivial or overly specific questions. - Phrase them as clear, standalone questions. <p>Format your output STRICTLY as a numbered list. Do not add any other text before or after the list.</p> <p>Transcript: {transcript}</p>
	<p>System: You are an expert Q&A agent. You will be given a question and a set of context paragraphs. Your task is to synthesize a clear and accurate answer based ONLY on the retrieved context.</p> <p>Prompt: Please answer the following question using ONLY the information from the context provided below.</p> <p>Question: {question}</p> <p>Context: {context}</p> <p>If the context does not contain the answer, state that the information is not available in the provided context.</p>

Figure 8: RAG prompt.

evident across both high- and low-resource languages. A particularly striking improvement is observed for Hindi, which reveals the Multi-Agent model’s role as a performance equalizer. In the Single-Agent architecture, Hindi substantially underperforms compared to Chinese, English, and Romanian, with QA-Alignment dropping to 3.96 and Answer Fluency falling to 3.72—both failing to reach the 4.0 threshold. However, the Multi-Agent model significantly boosts these scores, elevating QA-Alignment to 4.37 and Answer Fluency to 4.14. This finding suggests that even when a base LLM exhibits inherent weaknesses in a specific language, implementing a Multi-Agent architecture can effectively compensate for these limitations. The approach acts as a performance equalizer, enabling the model to deliver quality that approaches or matches that of high-resource languages, thereby bypassing the model’s native capability constraints.

Topic Evaluation. As illustrated in Figure 15, topic complexity substantially impacts LLM judge performance across content domains. To ensure accurate topic categorization, each QA pair in the dataset was manually annotated and verified by human annotators, establishing ground truth topic labels for all evaluations. **Personal Growth.** This domain achieves the highest scores in both Mentorship (4.55) and Alignment (4.56). The preva-

lence of self-help content and motivational literature in pre-training corpora enables models to generate well-structured guidance with relative ease, as the domain relies on broadly applicable principles rather than specialized knowledge. **Physical Health.** This domain presents the greatest challenge, registering the lowest scores in Mentorship (4.04) and Alignment (4.12). Qualitative analysis reveals that speakers frequently provided medical-style guidance incorporating complex anatomical terminology (e.g., hippocampus, amygdala, prefrontal cortex), biochemical nomenclature (e.g., cortisol, dopamine, serotonin), bodily fluids, and pharmaceutical references. Models demonstrated insufficient precision for assessing such nuanced clinical content, frequently producing vague or overly generalized evaluations. **Finance.** Despite having the smallest dataset (287 samples), Finance achieved a robust Alignment score of 4.47, comparable to more straightforward topics. This success stems from the structured, rule-based nature of financial discourse. Speakers typically provided logically structured guidance adhering to consistent frameworks (e.g., profit-loss analysis, investment strategies), enabling accurate LLM evaluation even with limited training examples.

As demonstrated in Figure 15, architectural complexity significantly impacts performance on challenging domains while showing minimal differen-

Model Type	Prompt
Multi-Agent	<p>System: You are an expert at analyzing transcripts and segmenting them by topic. The transcript has been split into numbered lines where each line represents a complete thought. Identify topic boundaries and assign concise topic titles.</p> <p>Prompt: Output a JSON list of dictionaries with these keys: - "topic": Concise descriptive title (3-7 words) - "start_line": First line number of this topic - "end_line": Last line number of this topic</p> <p>Rules: 1. Topics must cover consecutive line numbers 2. Entire transcript must be covered without gaps or overlaps 3. The first topic must start at line 1 4. The last topic must end at line {total_lines} 5. Use line numbers exactly as provided 6. Output ONLY the JSON with no additional text</p> <p>Numbered Transcript: {numbered_transcript}</p>
	<p>System: You are a curious and insightful analyst.</p> <p>Prompt: Based on the following text, generate a list of potential questions with high educational or mentorship value. Format your output as a simple numbered list. Avoid generating duplicate questions with similar meanings.</p> <p>Text Segment: {segment_content}</p>
	<p>System: You are an expert content evaluator. Return ONLY a number 1-10.</p> <p>Prompt: Rate the following question for educational/mentorship value (1=poor, 10=excellent). Question: {question} Hint: {context_hint}. Return just a number.</p> <p>System: You are a content analyst. Your task is to provide a clear and concise justification for a question's selection status.</p> <p>Prompt: You will be given a question, its source topic, and its "Selected" or "Rejected" status. Your job is to explain *why* that status makes sense.</p> <p>**Question:** {question_item['question']} **Source Topic:** {question_item['source_segment']['topic']} **Selection Status:** {question_item['status']}</p> <p>Based on the information above, provide a concise reason. Use the following as **sample reasons for style and inspiration, but you are not limited to them**:</p> <ul style="list-style-type: none"> - **Good 'Selected' reasons:** <ul style="list-style-type: none"> - "Addresses the key concept of 'model gluing', providing high educational value." - "Asks about career transition challenges, which has direct mentorship value." - **Good 'Rejected' reasons:** <ul style="list-style-type: none"> - "The question is too basic and offers little insight." - "This topic is already covered by another, more specific selected question." - "This question has a low score based on the result of agent 3."
	<p>System: You are a helpful Q&A assistant. Your task is to answer the given question using ONLY the information from the provided context.</p> <p>Prompt: **Context:** {context} — **Question:** {question}</p> <p>Based on the context above, what is the answer of the following questions ?</p>

Figure 9: Multi-Agent prompt.

1225 tiation on straightforward topics. **Physical Health.**
1226 This domain exhibits the most pronounced per-
1227 formance gap. Single Agent achieves Answer
1228 Mentorship of 3.78 and QA Alignment of 3.95,
1229 with marginal improvement in Dual Agent (3.94
1230 and 4.01). Multi Agent substantially elevates
1231 these to 4.24 and 4.34, demonstrating clear ad-
1232 vantages in handling clinically complex content.
1233 **Finance.** This logic intensive domain reveals
1234 similar architectural differentiation. Single Agent
1235 produces Answer Mentorship of 4.21, Q Mentor-
1236 ship of 3.86, QA Alignment of 4.23, and Answer
1237 Clarity of 4.23. Multi Agent markedly improves
1238 these to 4.51, 4.20, 4.60, and 4.50, respectively.
1239 **Personal Growth and Career.** These domains
1240 show minimal performance variation across archi-

tectures, with Single Agent achieving scores com-
parable to Multi Agent. This reflects the abun-
dance of motivational content in training data, en-
abling simpler architectures to generate effective
responses. These results reveal domain depen-
dent performance differentiation: common top-
ics (Career, Personal Growth) show architectural
equivalence, while challenging domains (Physical
Health, Finance) demonstrate significant separa-
tion—Single and Dual Agent models degrade sub-
stantially, whereas Multi Agent maintains robust
performance.

1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252

⚠ Please do not rate neutral (3) unless necessary or required.

Question Fluency
(Evaluates grammar, syntax, and sentence structure.) The question is grammatically correct and free from language errors.

Question:
What strategies can researchers employ to ensure they are contributing meaningfully to the field of large language models without simply following popular trends?

Rating: 1 - Strongly Disagree
 2 - Disagree
 3 - Neutral
 4 - Agree
 5 - Strongly Agree

Answer Fluency
(Evaluates grammar, syntax, and sentence structure.) The answer is grammatically correct and free from language errors.

Answer:
Researchers aiming to contribute meaningfully to the field of large language models (LLMs) while avoiding the pitfalls of chasing popular trends should consider several strategic approaches. 1. **Identify Niche Areas:** Recognize that the field is vast, and not every researcher needs to focus on the latest, most hyped models. By identifying specific niches or underexplored areas within LLMs, researchers can carve out unique contributions. 2. **Understand Limitations:** Given the industry dominance and resource constraints, researchers should critically examine the limitations of existing LLMs. This includes understanding biases, ethical considerations, and practical applications. 3. **Focus on Analysis and Application:** Since much of the current research is driven by the ability to analyze and apply LLMs, researchers can focus on these aspects to offer valuable insights. Analyzing the biases and errors in LLMs can lead to more robust and ethically sound models. 4. **Stay Informed but Not Focused Solely on Trends:** While staying informed about the latest developments is crucial, researchers should not let themselves be swayed by fleeting trends. Keeping an eye on foundational and long-term impacts rather than short-term hype can lead to more enduring contributions. 5. **Engage with Academic and Industry Perspectives:** Balancing academic research with industry insights can provide a comprehensive view of the field. Academics can offer theoretical depth, while industry experts bring real-world applicability and constraints. By employing these strategies, researchers can ensure their contributions are meaningful and impactful, fostering advancements that go beyond mere trend-following.

Figure 10: GUI of the Human Evaluation Process.

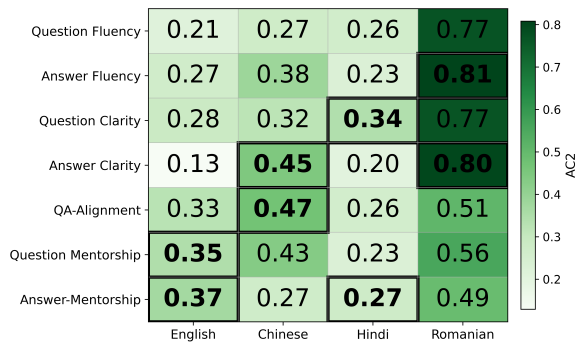


Figure 11: Inter-Annotator Agreement AC2 Scores by Language

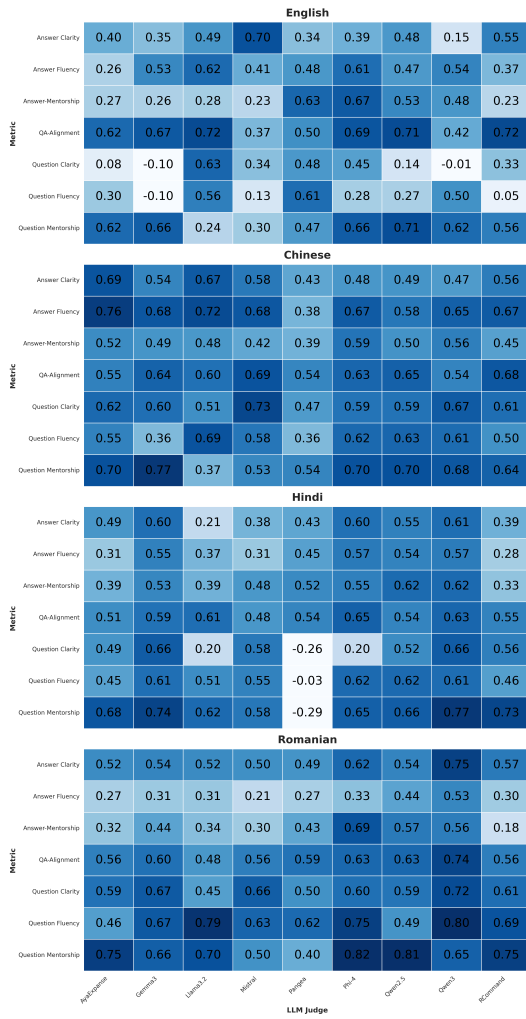


Figure 12: LLM-Human AC2 Agreement Scores

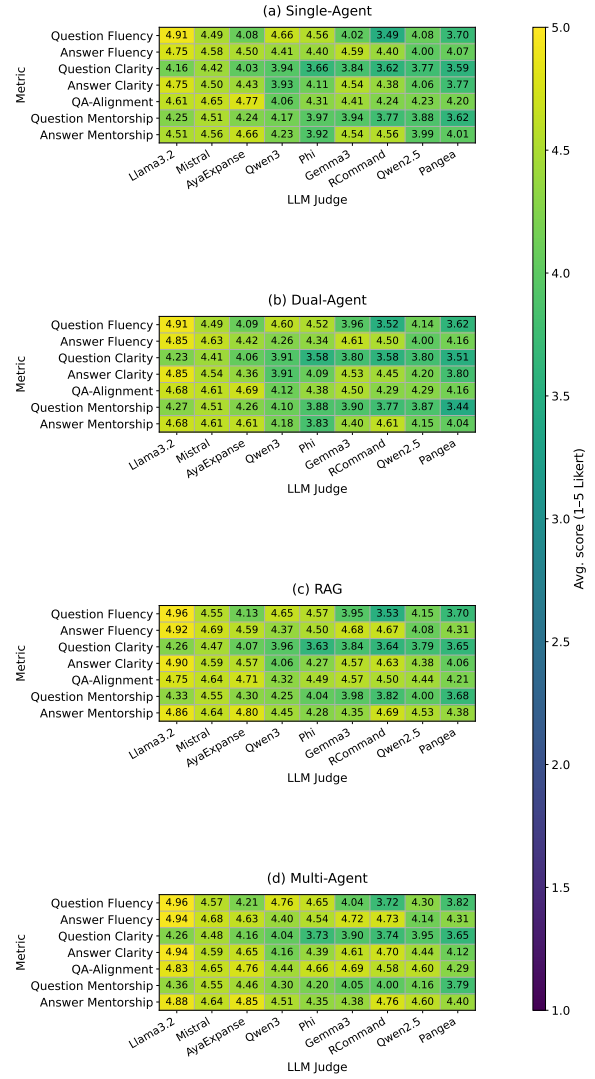


Figure 13: Judge-wise performance across QA-generation models: (a) Single-Agent, (b) Dual-Agent, (c) RAG, (d) Multi-Agent. Scores are averaged over metrics using a shared 1-5 Likert scale.

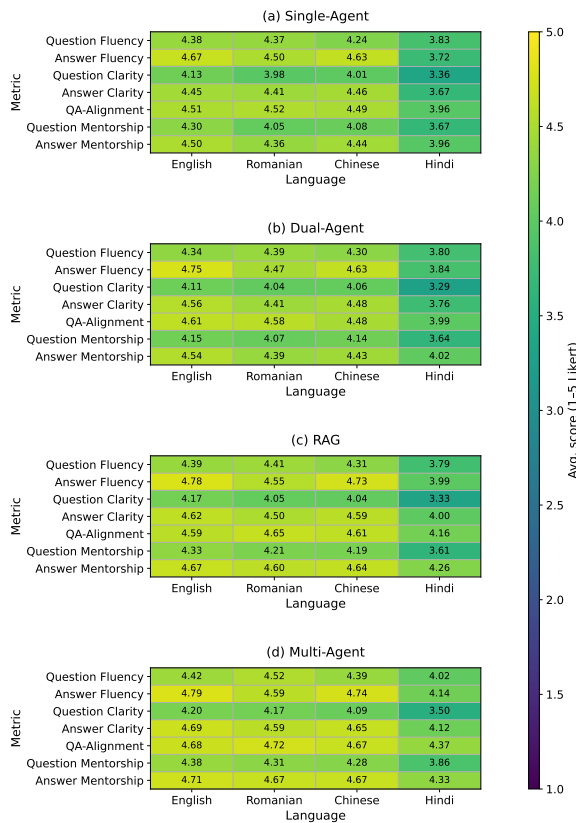


Figure 14: Language-wise performance across QA-generation models: (a) Single-Agent, (b) Dual-Agent, (c) RAG, (d) Multi-Agent. Scores are averaged over human and LLM judges (1–5 Likert).

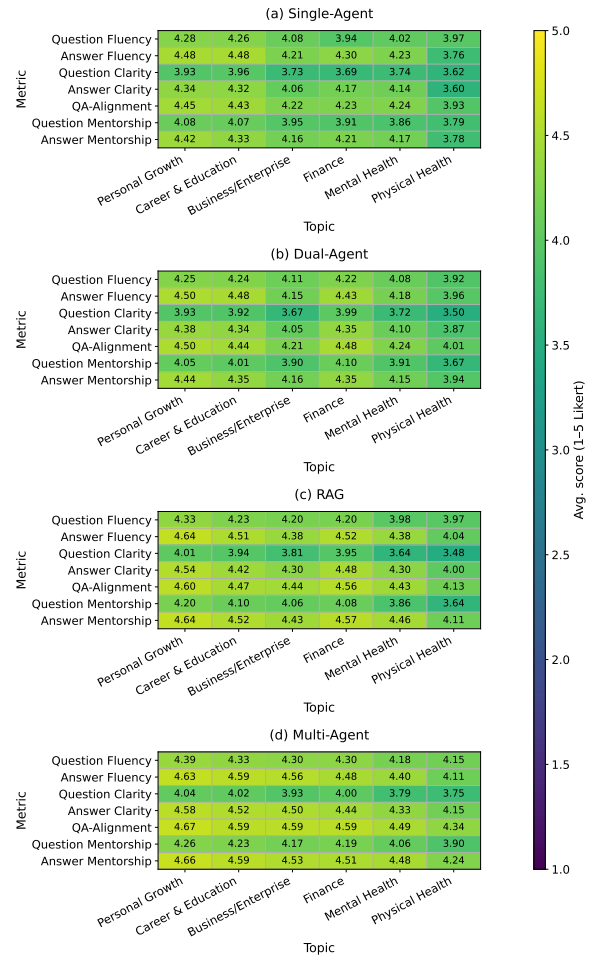


Figure 15: Topic-wise performance across QA-generation models: (a) Single-Agent, (b) Dual-Agent, (c) RAG, and (d) Multi-Agent. Scores are averaged over human and LLM judges (1–5 Likert scale).

Table 3: Mentorship Videos

Video ID	URL	Language	Video ID	URL	Language
1	Video Link	English	31	Video Link	Chinese
2	Video Link	English	32	Video Link	Chinese
3	Video Link	English	33	Video Link	Chinese
4	Video Link	English	34	Video Link	Chinese
5	Video Link	English	35	Video Link	Chinese
6	Video Link	English	36	Video Link	Chinese
7	Video Link	English	37	Video Link	Chinese
8	Video Link	English	38	Video Link	Chinese
9	Video Link	English	39	Video Link	Chinese
10	Video Link	English	40	Video Link	Chinese
11	Video Link	English	41	Video Link	Chinese
12	Video Link	English	42	Video Link	Chinese
13	Video Link	English	43	Video Link	Chinese
14	Video Link	English	44	Video Link	Chinese
15	Video Link	English	45	Video Link	Chinese
16	Video Link	English	46	Video Link	Chinese
17	Video Link	English	47	Video Link	Chinese
18	Video Link	English	48	Video Link	Chinese
19	Video Link	English	49	Video Link	Chinese
20	Video Link	English	50	Video Link	Chinese
21	Video Link	English	51	Video Link	Chinese
22	Video Link	English	52	Video Link	Chinese
23	Video Link	English	53	Video Link	Chinese
24	Video Link	English	54	Video Link	Chinese
25	Video Link	English	55	Video Link	Chinese
26	Video Link	English	56	Video Link	Chinese
27	Video Link	English	57	Video Link	Chinese
28	Video Link	English	58	Video Link	Chinese
29	Video Link	English	59	Video Link	Chinese
30	Video Link	English	60	Video Link	Chinese

Table 4: Mentorship Videos (Continue)

Video ID	URL	Language	Video ID	URL	Language
61	Video Link	Hindi	91	Video Link	Romanian
62	Video Link	Hindi	92	Video Link	Romanian
63	Video Link	Hindi	93	Video Link	Romanian
64	Video Link	Hindi	94	Video Link	Romanian
65	Video Link	Hindi	95	Video Link	Romanian
66	Video Link	Hindi	96	Video Link	Romanian
67	Video Link	Hindi	97	Video Link	Romanian
68	Video Link	Hindi	98	Video Link	Romanian
69	Video Link	Hindi	99	Video Link	Romanian
70	Video Link	Hindi	100	Video Link	Romanian
71	Video Link	Hindi	101	Video Link	Romanian
72	Video Link	Hindi	102	Video Link	Romanian
73	Video Link	Hindi	103	Video Link	Romanian
74	Video Link	Hindi	104	Video Link	Romanian
75	Video Link	Hindi	105	Video Link	Romanian
76	Video Link	Hindi	106	Video Link	Romanian
77	Video Link	Hindi	107	Video Link	Romanian
78	Video Link	Hindi	108	Video Link	Romanian
79	Video Link	Hindi	109	Video Link	Romanian
80	Video Link	Hindi	110	Video Link	Romanian
81	Video Link	Hindi	111	Video Link	Romanian
82	Video Link	Hindi	112	Video Link	Romanian
83	Video Link	Hindi	113	Video Link	Romanian
84	Video Link	Hindi	114	Video Link	Romanian
85	Video Link	Hindi	115	Video Link	Romanian
86	Video Link	Hindi	116	Video Link	Romanian
87	Video Link	Hindi	117	Video Link	Romanian
88	Video Link	Hindi	118	Video Link	Romanian
89	Video Link	Hindi	119	Video Link	Romanian
90	Video Link	Hindi	120	Video Link	Romanian