

MAJI: A Multi-Agent Workflow for Augmenting Journalistic Interviews

Anonymous ACL submission

Abstract

Journalistic interviews are creative, dynamic processes where success hinges on insightful, real-time questioning. While Large Language Models (LLMs) can assist, their tendency to generate coherent but uninspired questions optimizes for probable, not insightful, continuations. This paper investigates whether a structured, multi-agent approach can overcome this limitation to act as a more effective creative partner for journalists. We introduce MAJI, a system designed for this purpose, which employs a divergent-convergent architecture: a committee of specialized agents generates a diverse set of questions, and a convergent agent selects the optimal one. We evaluated MAJI against a suite of strong LLM baselines. Our results demonstrate that our multi-agent framework produces questions that are more coherent, elaborate, and original (+36.9% for our best model vs. a standard LLM baseline), exceeded strong LLM baselines on key measures of creative question quality. Most critically, in a blind survey, professional journalists preferred MAJI’s selected questions over those from the baseline by a margin of more than two to one. We present the system’s evolution, highlighting the architectural trade-offs that enable MAJI to augment, rather than simply automate, journalistic inquiry. We will release the code upon publication.

1 Introduction

The practice of journalism is a cornerstone of an informed society, with the interview serving as a primary tool for information gathering and narrative construction. While interviews are often prepared with a structured outline, the most compelling insights emerge from unscripted moments. A journalist’s ability to react to new information and identify novel angles in real-time separates a standard interview from a revelatory one. This dynamic process, however, presents a significant cognitive load.

Recent advancements in Large Language Models (LLMs) have opened new avenues for assisting in complex, language-based tasks (Touvron et al., 2023; OpenAI et al., 2024). A straightforward approach might involve prompting an LLM with the conversation history and asking for the next question. However, this often yields generic or predictable questions (Gordin et al., 2023), as LLMs tend to optimize for the most probable continuation rather than the most insightful or creative one. Recognizing this limitation, recent research has proposed a range of methods to enhance the originality and depth of LLM-generated interview questions. For example, works such as Spangher et al. (2025), Lin et al. (2025b), and Tian et al. (2024) introduce agentic workflows and creative reasoning strategies to move beyond surface-level responses. Our work builds on this line of research, further exploring how a multi-agent, divergent-convergent architecture can systematically augment the creative process in journalistic interviews.

Formally, this task can be seen as a conditional language generation problem. A standard LLM approach maximizes the likelihood of the next question Q_{t+1} given the transcript history T_t :

$$Q_{t+1} = \arg \max_Q P(Q|T_t)$$

This formulation often leads to probable but uninspired responses. We propose reframing the problem as maximizing a utility function $U(Q)$ that captures the goals of journalistic inquiry:

$$Q_{t+1} = \arg \max_Q U(Q|T_t, O, P)$$

where utility depends on the question’s insight and relevance to the interview’s strategic Outline O and the interviewee’s Persona P . MAJI is proposed to address this more complex optimization problem.

To address this gap, we introduce MAJI (Multi-Agent Workflow for Journalism Interview), a sys-

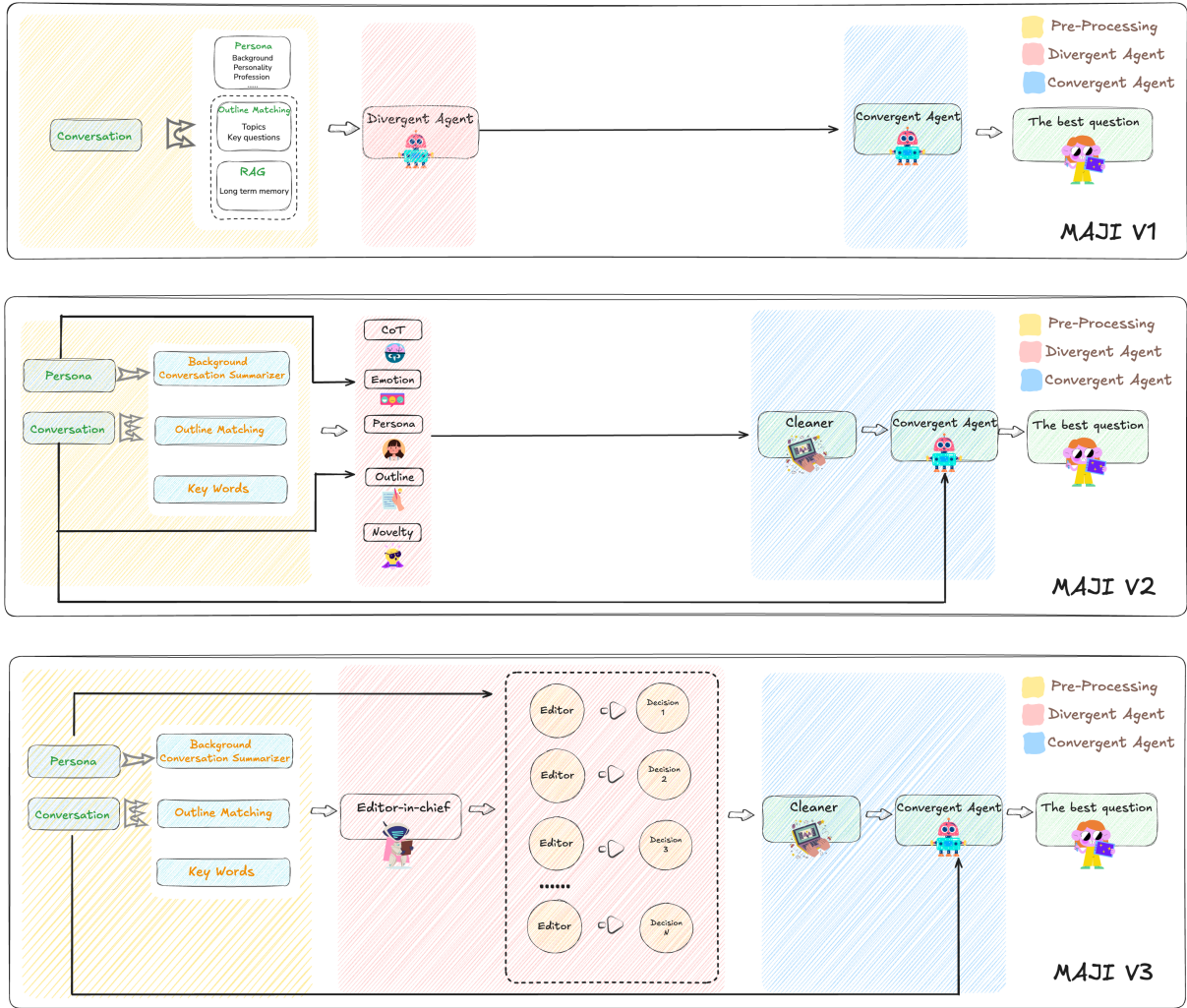


Figure 1: The architectural evolution of MAJI across its three versions. V1 (left) established a simple two-agent divergent-convergent model. V2 (center) introduced a specialized committee of agents for greater diversity. V3 (right) explored dynamic agent generation for adaptive strategies.

tem designed not to replace the journalist, but to augment their creative process. The human-in-the-loop workflow, depicted in Figure 2, positions MAJI as an assistant that provides suggestions while the journalist retains full control. MAJI is built on the psychological principle of divergent-convergent thinking, a cornerstone of creative problem-solving (Guilford, 1950). Our hypothesis is that by decomposing question generation into specialized sub-tasks, a multi-agent system (MAS) can produce more diverse and insightful questions than a single model. The power of MAS has been shown for complex logical tasks (Wooldridge, 2009; Wang et al., 2023; Wu et al., 2023; Qian et al., 2023; Li et al., 2023a), and we apply this paradigm to a creative domain (Lin et al., 2025a; Xi et al., 2025; Zhou et al., 2023; Li et al., 2023b).

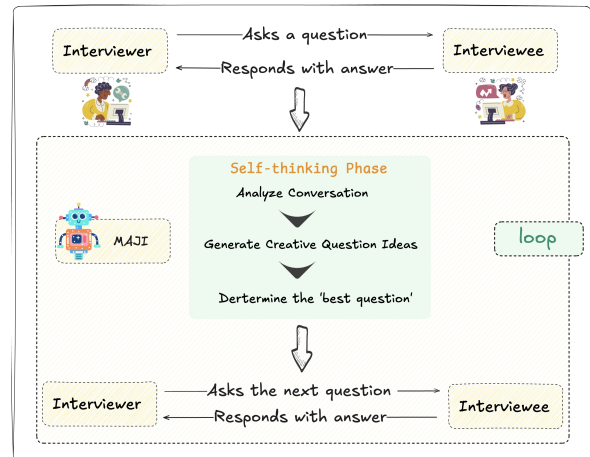


Figure 2: High-level overview of the MAJI-assisted interview workflow. MAJI operates in a continuous loop, observing the conversation and providing question suggestions to the journalist, who retains final control over the direction of the interview.

This paper details the design and evolution of the MAJI framework across three major versions, from a simple proof-of-concept to a sophisticated system capable of dynamically generating its own specialized agents. We conducted a rigorous evaluation using a simulated interview environment inspired by prior work in computational journalism (Diakopoulos, 2019; ?). Our contributions are threefold:

1. **A Novel Multi-Agent Framework:** We proposed and implemented a divergent-convergent multi-agent architecture for the creative task of interview question generation.
2. **Empirical Evaluation:** We conducted a comprehensive comparison of MAJI against strong LLM baselines, using quantitative, qualitative, and comparative metrics assessed by both LLM-as-judge and professional journalists.
3. **Architectural Insights:** We documented MAJI’s evolution, providing insights into the trade-offs between fixed-specialist and dynamically-generated agent committees.

Our findings show that MAJI consistently generates questions that are more insightful, original, and contextually relevant than those from a well-prompted, powerful LLM. This work demonstrates the potential of structured multi-agent systems to move beyond simple automation and act as powerful creative partners in complex professional domains like journalism.

2 Related Work

Prior systems concentrate on data mining, fact checking and bias detection (Cohen et al., 2011; Diakopoulos, 2019; Hamborg et al., 2018), leaving the *real-time* interviewing stage under-explored. Datasets such as *NewsInterview* (Lu et al., 2024) highlight this gap; MAJI directly addresses it by operating live. Techniques like Chain-of-Thought (Wei et al., 2022), Tree-of-Thought (Yao et al., 2023) and retrieval-augmented generation improve single-model reasoning but still rely on one monolithic agent. Multi-agent frameworks such as Auto-Gen (Wu et al., 2023) and CAMEL (Li et al., 2024) show that dividing labour can yield stronger reasoning; MAJI adapts this insight to creative follow-up generation. Psychology links creativity to alternating idea generation and selection (Guilford, 1950).

Agent committees have applied this pattern to storytelling and design (Yao et al., 2019; Li et al., 2023b). MAJI is the first to embed the paradigm in journalistic interviewing and to validate its impact with professional users.

3 The MAJI Framework

MAJI is designed to mirror and support the cognitive workflow of a journalist. The framework is built on several core concepts: foundational inputs that set the stage, a multi-agent system that drives the question generation process, and an iterative loop that allows the system to learn and adapt as the interview progresses.

3.1 Foundational Concepts

The interview process begins with the journalist framing the initial context. This human-in-the-loop step is crucial for grounding the AI’s subsequent contributions. Three key pieces of information establish this foundation:

- **Persona:** A detailed profile of the interviewee, including their background, personality, profession, and the primary purpose of the interview. This is crafted by the journalist to guide the system’s interaction style. For example, a persona for a professional mermaid performer might highlight their artistic motivations and physical challenges, shaping questions toward these themes.
- **Outline:** A structured list of topics and key questions that the journalist intends to cover. This serves as the strategic backbone of the interview, ensuring coverage of critical areas while allowing flexibility for emergent insights.
- **Dynamic Background Summary:** To track evolving context, MAJI uses a dedicated agent to maintain a BackgroundSummary. It has a `long_term_summary` for foundational information (e.g., interviewee’s career trajectory) and a `short_term_summary` for the last five conversational turns. This provides all agents with a continuously updated, concise view of the conversation history.

These inputs ensure MAJI remains grounded in the journalist’s objectives and the interview’s evolving context, enabling questions that are both strategically aligned and responsive to in-session dynamics.

3.2 The Divergent-Convergent Workflow

The core of MAJI is an implementation of the divergent-convergent thinking model, executed by specialized agents in a multi-stage workflow. This process, detailed for our primary MAJI V2 system in Algorithm 1, ensures a balance between creative exploration and strategic focus. MAJI’s architecture operationalizes this principle. A committee of specialized *divergent agents* brainstorms potential questions, each focusing on a distinct creative vector, such as emotional depth, causal reasoning, adherence to the interviewee’s persona, or pure novelty. Their suggestions are then processed by an *Editor Agent* to refine and deduplicate the pool of ideas. Finally, a *convergent agent*, acting as an Editor-in-Chief, selects the single best question that aligns with the conversation’s flow and the journalist’s strategic goals for the interview. The key stages are:

Algorithm 1 MAJI V2 Question Generation Workflow

```

1: Input: Outline  $O$ , Persona  $P$ , Transcript  $T$ 
2: Output: Next question  $Q$ 
3:  $S \leftarrow \text{BackgroundAgent}(T, P)$   $\triangleright$  Update summaries
4:  $K \leftarrow \text{KeywordsAgent}(T, S)$   $\triangleright$  Extract keywords
5:  $M \leftarrow \text{OutlineMatcherAgent}(O, T)$   $\triangleright$  Map to outline
6:  $C \leftarrow \emptyset$   $\triangleright$  Initialize candidate pool
7: for each  $\text{DivergentAgent}_i$  in  $\{D_1, \dots, D_n\}$  do
8:    $C_i \leftarrow \text{DivergentAgent}_i(K, S, M, P)$   $\triangleright$  Propose questions
9:    $C \leftarrow C \cup C_i$ 
10: end for
11:  $C' \leftarrow \text{EditorAgent}(C)$   $\triangleright$  Refine candidates
12:  $Q \leftarrow \text{ConvergentAgent}(C', T, P, O)$   $\triangleright$  Select optimal question
13: return  $Q$ 

```

1. **Context Analysis (Pre-Divergence):** Before any new questions are brainstormed, a set of pre-processing agents analyzes the latest turn in the conversation to establish a shared understanding of the current state. This includes the *BackgroundAgent* updating the long- and short-term summaries, the *KeywordsAgent* extracting the most salient terms from the latest response, and the *OutlineMatcherAgent* assessing which parts of the interview outline

have been covered.

2. **Divergent Thinking:** With the updated context, the committee of specialized divergent agents generates a wide range of potential follow-up questions in parallel. This parallel, specialized approach is designed to produce a candidate pool with high "Flexibility," ensuring a rich set of creative options.
3. **Editing & Curation:** The raw list of questions from the divergent phase is often redundant. The *EditorAgent* uses sentence-transformer embeddings to identify and merge semantically similar questions (similarity threshold: 0.85). This step curates a clean, concise list of unique candidate questions.
4. **Convergent Selection:** Finally, the *ConvergentAgent* takes the curated list of questions and, guided by the journalist’s stated strategic preference, selects the single best question to ask next. This selection process is not based on arbitrary heuristics, but on an implicit model of journalistic utility, as detailed below.

3.2.1 Convergent Utility Maximization

The core of the convergent step is the maximization of a utility function U . Rather than being a simple, hard-coded formula, this function is a conceptual model of question quality that the *ConvergentAgent* is prompted to approximate. We can formally define this utility for a candidate question $Q \in C'$ as a weighted sum of scores from various quality dimensions:

$$U(Q|\cdot, S_p) = \sum_{i=1}^N w_i(S_p) \cdot \phi_i(Q|T_t, O, P)$$

where each component represents a desirable attribute of a question:

- $\phi_i(Q|\cdot)$ are scoring functions that evaluate different facets of a question’s quality, such as its *coherence*, *emotional depth*, *outline progression*, *persona alignment*, and *novelty*. These facets directly correspond to the specializations of the divergent agents.
- $w_i(S_p)$ are weights that are dynamically modulated by the journalist’s Strategic Preference S_p . For example, if the journalist sets the preference to "focus on emotion,"

269	the weight w_{emo} for the emotional depth score	and Novelty) brainstorms questions in parallel.	317
270	ϕ_{emo} is implicitly increased. If the preference	An EditorAgent refines the question pool, and a	318
271	is "balanced," the weights are distributed more	ConvergentAgent selects the best question based	319
272	evenly.	on the journalist's preference. This specialization	320
273		proved highly effective at generating rich and di-	321
274	In our implementation, the ConvergentAgent (a	verse candidate questions. An example is in Ap-	322
275	powerful LLM) does not compute explicit scores.	pendix A.1.	323
276	Instead, it performs a holistic evaluation, directly		
277	approximating the $\arg \max$ operation by using the	4.3 MAJI V3: Dynamic Agent Generation	324
278	strategic preference S_p to guide its selection from	MAJI V3 is an experimental evolution that dynam-	325
279	the candidate set C' . The prompt instructs it to "se-	ically devises its own interview strategy (Figure 1),	326
280	lect the single best question" that "aligns with the	where an agent plans the divergent phase. The key	327
281	preference," effectively performing this weighted	innovation is the introduction of an agent that plans	328
282	optimization. This leverages the LLM's nuanced	the divergent phase itself. The V3 architecture	329
283	reasoning to model the complex utility of a journal-	modifies the divergent step:	330
284	istic question.		
285	This structured workflow ensures that the final	1. The fixed committee of divergent agents is	331
286	question is not merely a probable continuation, but	removed.	332
287	a strategically selected option from a creatively		
288	diverse and well-curated set of possibilities.	2. A new agent, the EditorInChiefAgent, is in-	333
289		troduced. Based on a set of pre-defined heuris-	334
290	4 System Architecture and Evolution	tics, its role is to analyze the full conversation	335
291	The MAJI framework was developed and refined	context and devise a <i>plan</i> for the divergent	336
292	over three major versions. Each version represents	phase, as defined in our V3 data models.	337
293	a significant step in the architectural design, mov-		
294	ing from a simple proof-of-concept to a highly so-	3. This plan takes the form of a	338
295	phisticated and dynamic system. In essence, these	DivergentAgentPlan, which contains	339
296	versions can be viewed as a theoretical ablation	a list of DivergentAgentSpec objects. Each	340
297	study, with each successive version adding archi-	spec defines the name and instructions	341
298	tectural complexity to examine its impact on per-	for a temporary, single-use divergent agent	342
299	formance.	tailored to the immediate needs of the	343
300		conversation. For example, if the inter-	344
301	4.1 MAJI V1: A Foundational	viewee seems evasive, it might create	345
302	Proof-of-Concept	a "Probing_Clarification_Agent." If the	346
303	MAJI V1 served as a minimal proof-of-concept	conversation is stalling, it might create a	347
304	for the divergent-convergent idea (Figure 1). It	"Hypothetical_Scenario_Agent."	348
305	used a DAgent to generate questions with a sin-		
306	gle, broad LLM prompt and a CAgent that used	4. The system then dynamically instantiates	349
307	hard-coded, non-LLM heuristics for selection. This	these agents from the specs and runs them	350
308	rigid logic struggled with the fluid nature of inter-	in parallel to generate questions.	351
309	views, highlighting the need for the more nuanced,		
310	context-aware reasoning of V2's specialized agent	The rest of the pipeline is unchanged. V3 is more	352
311	committee.	autonomous, with the LLM defining generation	353
312		strategy. However, this complexity introduced	354
313	4.2 MAJI V2: The Specialized Agent	challenges. The EditorInChiefAgent's heuristic-	355
314	Committee	based approach is a limitation, and its performance	356
315	MAJI V2, our best-performing model, implements	did not surpass V2. A specific example is in Ap-	357
316	a robust "agent committee" architecture (Figure 1).	pendix A.2.	358
	The workflow, detailed in Section 3.2, begins with	5 Evaluation	359
	pre-processing agents establishing context. Then, a	We designed a comprehensive evaluation frame-	360
	fixed committee of five specialized divergent agents	work to assess the performance of the MAJI system.	361
	(ChainOfThought, Emotion, Outline, Persona,	The evaluation aims to answer three key questions:	362

1. How does MAJI’s question generation quality compare to a strong, conventionally-prompted LLM baseline across a diverse dataset?
2. How does the architectural choice (fixed vs. dynamic agent committee) impact performance between MAJI V2 and V3?
3. What are the specific strengths and weaknesses of each system, particularly regarding the trade-off between creativity and conversational coherence?

5.1 Experimental Setup

We evaluated MAJI against strong baselines on real-world interviews from the *NewsInterview* dataset and proprietary sources. All systems used gpt-4.1-mini. Baselines included LLM-Base, LLM-CoT, LLM-ToT, and LLM-RAG. We used Prometheus 2 for validation and adjusted originality scores with a threshold-based method for realism. Further details are in Appendix A.3.

5.2 Metrics

Our evaluation uses a combination of metrics calculated per-question and per-interview, averaged across the dataset.

5.2.1 Per-Question Metrics

These metrics are evaluated for each generated question. Our primary judge (GPT-4o) scored questions on six criteria: **Coherence** (logical connection), **Elaboration** (encouraging detailed responses), **Originality** (novelty), **Context Relevance** (relation to the last turn), **Outline Relevance** (alignment with the interview plan), and **Persona Alignment** (matching interviewee character). Our benchmark judge (Prometheus 2) provided scores for similar criteria, adding measures for **Insight** (probing deeper), **Conversational Synthesis** (integrating prior conversation), and **Strategic Progression** (advancing interview goals).

To better assess originality, we used a threshold-based adjustment. Since cosine similarity scores for semantically different questions are often non-zero, raw scores can be misleading. Our method uses a 0.6 threshold; similarities below this are scored as 1.0 (completely original), and the rest are scaled. This yields more realistic originality scores while preserving relative rankings.

5.2.2 Per-Interview Metrics

To assess overall strategic performance, we calculated the **Insight Trajectory**, measuring if a system

asks more insightful questions in the second half of an interview compared to the first.

6 Results

Our comprehensive evaluation demonstrates the effectiveness of the MAJI framework. We presented a detailed comparison of all MAJI versions and a suite of strong LLM baselines, with results aggregated across the evaluation dataset. The following sections analyze the quality of the final selected questions, the raw brainstormed suggestions, and the strategic performance of each model.

6.1 Primary Analysis: Selected vs. Brainstormed Quality

Our primary evaluation (Table 1) shows that the final questions selected by MAJI V2 and V3 are significantly higher quality than baselines, excelling on metrics like **Coherence**, **Elaboration**, **Originality**, and **Context Relevance** ($p < 0.001$).

This high quality results from our divergent-convergent design. Analysis of the raw brainstormed suggestions (Appendix A.4, Table 3) reveals the convergent agent adds significant value. For MAJI V2, the selection process improves **Coherence** (+9.2%) and **Context Relevance** (+17.6%), showing it successfully elevates the most promising ideas from a diverse but noisy pool. The adjusted originality scores show MAJI V2 achieves the highest originality (0.764), a 36.9% improvement over baseline, while maintaining strong performance elsewhere.

6.2 Benchmark Validation: Prometheus Evaluation

The Prometheus 2 benchmark (Table 2) strongly corroborates our primary analysis. MAJI V2 and V3 again emerge as top performers. Improvements over LLM-Base are statistically significant for key creative metrics like **Insight**, **Originality**, and **Elaboration** ($p < 0.001$ for MAJI V2/V3). MAJI V2 and V3 also score highest in **Outline Relevance** ($p < 0.001$), indicating their questions are both insightful and effective at advancing the interview. This independent verification confirms MAJI’s superior performance is a robust, statistically significant result. Prometheus’s assessment also validates our adjusted originality scores, with MAJI V2 achieving a 25.1% improvement over baseline. Interestingly, Prometheus rated V1 much lower than our primary judge, suggesting V1’s quality is more subjective.

Table 1: Evaluation of Final Selected Questions (Judged by GPT-4o). This table shows the quality of the single question chosen by each system to be asked next. All models are compared against the LLM-Base baseline. Significance from a two-tailed t-test is denoted by asterisks: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. MAJI V2 and V3 excel in generating high-quality, elaborate, and original questions. Originality scores have been adjusted using a threshold-based method to account for baseline similarity between questions. Best score in each category is in **bold**.

Metric	MAJI V1	MAJI V2	MAJI V3	LLM-Base	LLM-CoT	LLM-ToT	LLM-RAG
Coherence	0.768***	0.795***	0.791***	0.704	0.701	0.770***	0.735***
Elaboration	0.861	0.928***	0.932***	0.871	0.873	0.901***	0.871
Originality	0.641***	0.764***	0.736***	0.558	0.586***	0.666***	0.611***
Context Relevance	0.372***	0.434***	0.414***	0.319	0.319	0.373***	0.328
Outline Relevance	0.779***	0.656	0.661	0.670	0.673	0.668	0.658*
Persona Alignment	0.838***	0.868	0.865	0.874	0.880	0.880	0.872

Table 2: Benchmark Evaluation of Final Selected Questions (Judged by Prometheus 2). This table shows scores from a standardized, third-party model, validating the primary results. All models are compared against the LLM-Base baseline. Significance from a two-tailed t-test is denoted by asterisks: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. Best score is in **bold**.

Metric	MAJI V1	MAJI V2	MAJI V3	LLM-Base	LLM-CoT	LLM-ToT	LLM-RAG
Coherence	0.622	0.780***	0.771***	0.629	0.648	0.762***	0.631
Elaboration	0.754	0.898***	0.908***	0.774	0.822***	0.839***	0.793
Originality	0.613	0.733***	0.727***	0.586	0.584	0.654***	0.608
Context Relevance	0.697*	0.818***	0.870***	0.635	0.676	0.807***	0.659
Outline Relevance	0.454	0.586***	0.624***	0.417	0.485***	0.589***	0.425
Insight	0.645	0.815***	0.835***	0.688	0.741*	0.783***	0.703
Conversational Synthesis	0.581***	0.743***	0.739***	0.664	0.697*	0.730***	0.691
Persona Alignment	0.693***	0.837	0.861***	0.821	0.835	0.813	0.802

6.3 Human and Strategic Evaluation

We also conducted a human evaluation with 30 professional journalists and analyzed each system’s strategic performance. Journalists preferred MAJI V2’s questions nearly half the time (48.9%), more than double the rate of the baseline. Furthermore, analysis of the "Insight Trajectory" showed that MAJI V1 had the highest rate of improvement over an interview (Table 8). While MAJI V2 and V3 started from and maintained a much higher insight baseline, V1’s simpler architecture appeared effective at building conversational momentum. Full details are in Appendix A.5.

7 Discussion

The results of our experiments provide several key insights into the design of AI systems for creative professional domains.

7.1 The Power of Specialization

The stark performance difference between MAJI V2/V3 and even advanced LLM baselines validates our core hypothesis: decomposing a complex creative task into specialized sub-tasks yields superior

results. While advanced prompting like Tree-of-Thought improves LLM performance, MAJI V2’s architecture, with dedicated agents for emotion, logic, and novelty, consistently produced questions judged as more insightful and original by both AI and human experts. The adjusted originality scores show MAJI V2 achieves a 36.9% improvement over baseline on the NewsInterview dataset (0.764 vs 0.558) and a 44.1% improvement on the proprietary dataset (0.608 vs 0.422). This cross-dataset consistency across languages validates that the multi-agent architecture’s benefits are not domain-specific but a fundamental advantage. This suggests for creative augmentation, a specialized multi-agent architecture is more promising than prompting strategies for single models.

7.2 The "Creative Partner" vs. "Coherent Assistant"

A crucial finding is the trade-off between creative value and conversational coherence. The baseline LLMs, optimized for next-token prediction, naturally excel as Coherent Assistants.

MAJI, conversely, acts as a creative partner. It is less constrained by the most probable con-

versational path and more focused on generating high-quality, novel, and insightful questions. The strong preference from professional journalists underscores the value of this approach; they want a tool that expands their creative options, not just one that affirms their instincts. This focus on quality comes at a computational cost, representing a trade-off between speed and insight (see Appendix A.7).

7.3 Error Analysis

Although MAJI’s architecture explicitly avoids redundancy, we observe a trade-off between creativity and relevance. Divergent agents occasionally produce off-topic or abstract suggestions. This is by design, as it broadens potential questioning paths, even if many are discarded.

7.4 V3 and the Challenge of Meta-Cognition

MAJI V3’s experiment in dynamic agent generation provides insight into the challenges of AI meta-cognition. While this approach produced highly novel and relevant questions, it was less consistent than MAJI V2. The preference for V2 in our human evaluation suggests that for professional use, V2’s reliable creativity is currently more valuable than V3’s experimental novelty. This highlights a key challenge: building an agent that can effectively strategize about creative strategy is a difficult, higher-order task. For now, a carefully designed, fixed architecture is more robust. While V3’s dynamic agent generation is heuristic-based, this is a necessary stepping stone in creative domains where learning-based planning is an open challenge.

7.5 Broader Implications for Human-AI Creative Partnerships

MAJI’s principles are not limited to journalism. The divergent-convergent framework, with specialized agent committees, is a generalizable template for augmenting human creativity in any domain involving iterative ideation and strategic selection. Applications could include helping scientists brainstorm hypotheses, assisting marketing teams with slogans, or helping attorneys explore case strategies. This work contributes to a broader vision of AI not as a replacement for human intellect, but as a structured tool for amplifying it.

8 Conclusion

We introduced MAJI, a multi-agent system to assist journalists by generating creative and insightful

interview questions. We demonstrated that by decomposing this task into a divergent-convergent workflow with specialized agents, MAJI moves beyond standard and advanced LLM prompting. Our results, validated by 30 professional journalists, show that our best version, V2, produces questions consistently preferred over strong baselines. In a blind survey, journalists chose MAJI V2’s suggestions at more than double the rate of a standard baseline, confirming its superior alignment with professional judgment.

We have shown that a structured, multi-agent architecture is highly effective for augmenting complex, creative human tasks, trading a small amount of predictability for a significant gain in insight and originality. The success of MAJI V2’s “agent committee” provides a promising model for building AI-powered creative partners that help professionals overcome cognitive fixation and explore a wider possibility space. The experimental V3, while less performant, offered valuable insights into developing more autonomous, strategy-devising systems, highlighting a clear path for future research. The code, data, and models will be open-sourced to encourage further research.

Limitations

While this study provides strong evidence for MAJI’s effectiveness, we acknowledged several limitations that provide avenues for future work. First, our evaluation is conducted on a dataset of professionally curated interviews. The system’s robustness on noisier, out-of-domain data—such as unedited live transcripts or interviews on highly specialized topics—remains to be tested. Second, as our qualitative analysis of V3’s failure case demonstrates, the system can occasionally produce awkward or contextually inappropriate suggestions. A more detailed qualitative error analysis would be beneficial for identifying and mitigating these failure modes. Our originality metric has been improved through a threshold-based adjustment method to account for baseline similarity between questions, but could be further enhanced to better distinguish semantic novelty from lexical paraphrasing. Third, the latency of the MAJI system, particularly V2 and V3, is a significant consideration (see Appendix A.7). While we argue this is a justifiable trade-off for question quality, further optimization is required to make the system more responsive.

Although a learning-based planner may offer better adaptability in theory, the lack of structured supervision and reward signals in open-ended domains like interviewing makes such training infeasible at present. We therefore use heuristic planning as an exploratory first step in operationalizing strategic meta-reasoning in creative workflows.

Building on these points, future work will focus on expanding our evaluation to broader datasets, refining V3's heuristic planner into a learning-based agent (e.g., using techniques like verbal reinforcement learning (Shinn et al., 2023) or RLHF (Christiano et al., 2017)), conducting component-wise ablation studies to quantify each agent's contribution, developing a real-time user interface for live interviews, and performing a systematic qualitative error analysis to guide future improvements. Future work should also explore model distillation to create smaller, faster versions of the agents, caching strategies for recurring sub-problems, and asynchronous processing to reduce the perceived latency for the user.

Ethical Considerations

The deployment of AI tools like MAJI in journalism necessitates careful consideration of ethical implications. First, there is a risk that the underlying LLM could introduce subtle biases into the question generation process, potentially reflecting political or confirmation biases from its training data and steering conversations in unintended directions. We suggested ongoing monitoring and fine-tuning with diverse datasets to mitigate this. Second, while MAJI is designed to augment, not replace, the journalist, there is a risk of over-reliance, which could diminish the journalist's critical thinking and rapport-building skills. The system should be framed as a supportive tool, with the final decision always resting with the human journalist, who must also approve any intermediate strategic pivots suggested by the AI. Finally, the privacy and consent of the interviewee are paramount. All data used for training and running the system must be handled with explicit consent and robust anonymization procedures, as was done in this study (Alzoubi et al., 2024).

Acknowledgments

During the preparation of this work, the authors utilized an AI-powered tool to assist with programming, text editing, and formatting. All AI-

generated content, including code and text, was carefully reviewed, revised, and validated by the authors to ensure its accuracy and alignment with the research goals.

References

- Omar Alzoubi, Normahfuzah Ahmad, and Norsiah Abdul Hamid. 2024. [Artificial intelligence in newsrooms: Ethical challenges facing journalists](#). *Studies in Media and Communication*, 12:401.
- Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems 30*, pages 4299–4307. Curran Associates, Inc.
- Sarah Cohen, James T. Hamilton, and Fred Turner. 2011. [Computational journalism](#). *Commun. ACM*, 54(10):66–71.
- Nicholas Diakopoulos. 2019. Automating the news: How algorithms are rewriting the media. In *Computational Journalism*, Cambridge, MA, USA. Harvard University Press.
- Matan Gordin, Eric Horvitz, and Jaime Teevan. 2023. Evaluating creativity in large language models: A review. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, pages 313–324.
- J. P. Guilford. 1950. [Creativity](#). *Am. Psychol.*, 5(9):444–454.
- Felix Hamborg and 1 others. 2018. Automated identification of media bias in news articles. *Int. J. Data Sci. Anal.*, 6(4):327–339.
- Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024. [Prometheus 2: An open source language model specialized in evaluating other language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 4334–4353, Miami, Florida, USA. Association for Computational Linguistics.
- Chengcheng Li, Zhaobo Wang, Yan Zhao, Ruobing Yao, Xiao Zhang, and WenGuan Chen. 2024. A multi-agent framework for reasoning and evaluation.
- Guohao Li, Hasan Mendi, Shanshan Wang, Cunxiang Wang, Chen Lv, Yifei Zhang, Yuliang Li, Nan Duan, and Weiming Wang. 2023a. Camel: Communicative agents for "mind" exploration of large scale language model society. In *Advances in Neural Information Processing Systems*, volume 36.
- Yifu Li, Jialu Li, Yufang Lai, and Nanyun Peng. 2023b. [Accord: A multi-document approach to generating diverse and coherent summaries](#). In *Proceedings of the 61st Annual Meeting of the Association for*

704	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	Grave, and Guillaume Lample. 2023. Llama: Open	761
705	pages 8341–8356, Toronto, Canada. Association for	and efficient foundation language models . <i>Preprint</i> ,	762
706	Computational Linguistics.	arXiv:2302.13971.	763
707	Yi-Cheng Lin, Kang-Chieh Chen, Zhe-Yan Li, Tzu-	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao	764
708	Heng Wu, Tzu-Hsuan Wu, Kuan-Yu Chen, Hung-	Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,	765
709	yi Lee, and Yun-Nung Chen. 2025a. Creativity in	Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei	766
710	LLM-based multi-agent systems: A survey . <i>CoRR</i> ,	Wei, and Ji-Rong Wen. 2023. A survey on large	767
711	abs/2505.21116.	language model based autonomous agents . <i>CoRR</i> ,	768
712	Yi-Cheng Lin, Kang-Chieh Chen, Zhe-Yan Li, Tzu-	abs/2308.11432.	769
713	Heng Wu, Tzu-Hsuan Wu, Kuan-Yu Chen, Hung	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	770
714	yi Lee, and Yun-Nung Chen. 2025b. Creativity in	Bosma, Fei Xia, Ed Chi, Quoc V. Le, and Denny	771
715	llm-based multi-agent systems: A survey . <i>Preprint</i> ,	Zhou. 2022. Chain-of-thought prompting elicits rea-	772
716	arXiv:2505.21116.	soning in large language models. In <i>Advances in</i>	773
717	Michael Lu, Hyundong Justin Cho, Weiyan Shi,	<i>Neural Information Processing Systems 35</i> . Curran	774
718	Jonathan May, and Alexander Spangher. 2024. News-	Associates, Inc.	775
719	interview: a dataset and a playground to evalu-	Michael Wooldridge. 2009. <i>An introduction to multi-</i>	776
720	ate llms’ ground gap via informational interviews .	<i>agent systems</i> . John Wiley & Sons, Hoboken, NJ,	777
721	<i>Preprint</i> , arXiv:2411.13779.	USA.	778
722	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu,	779
723	Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	Beibin Li, Erkang Zhu, and Li Li. 2023. Autogen:	780
724	man, Diogo Almeida, Janko Altmenschmidt, Sam Alt-	Enabling next-gen LLM applications via multi-agent	781
725	man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	conversation framework . <i>CoRR</i> , abs/2308.08155.	782
726	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-	Zekun Xi, Wenbiao Yin, Jizhan Fang, Jialong Wu, Run-	783
727	ing Bao, Mohammad Bavarian, Jeff Belgum, and	nan Fang, Ningyu Zhang, Jiang Yong, Pengjun Xie,	784
728	262 others. 2024. Gpt-4 technical report . <i>Preprint</i> ,	Fei Huang, and Huajun Chen. 2025. OmniThink: Ex-	785
729	arXiv:2303.08774.	panding knowledge boundaries in machine writing	786
730	Chen Qian, Xin Cong, Cheng Yang, Weize Chen,	through thinking . <i>CoRR</i> , abs/2501.09751.	787
731	Yusheng Su, Juyuan Liu, Yufan Liu, and Zipeng Yu.	Lili Yao, Nanyun Peng, Dietrich Klakow, Mark Riedl,	788
732	2023. Chatdev: Communicative agents for software	and Weischedel Ralph Wang. 2019. Plan-and-write:	789
733	development . <i>CoRR</i> , abs/2307.07924.	Towards better automatic storytelling. In <i>Proceed-</i>	790
734	Noah Shinn, Federico Cassano, Ashwin Gopinath,	<i>ings of the AAAI Conference on Artificial Intelligence</i> ,	791
735	Karthik Narasimhan, and Shunyu Yao. 2023. Re-	volume 33, pages 7454–7461.	792
736	flexion: Language agents with verbal reinforcement	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,	793
737	learning. In <i>Advances in Neural Information Pro-</i>	Thomas L. Griffiths, Yuan Cao, and Karthik	794
738	<i>cessing Systems 36</i> . Curran Associates, Inc.	Narasimhan. 2023. Tree of thoughts: Deliber-	795
739	Alexander Spangher, Tenghao Huang, Philippe Laban,	ate problem solving with large language models .	796
740	and Nanyun Peng. 2025. Creative planning with lan-	<i>Preprint</i> , arXiv:2305.10601.	797
741	guage models: Practice, evaluation and applications .	Wangchunshu Zhou, Peng Wang, Yifei Li, Wenying	798
742	In <i>Proceedings of the 2025 Annual Conference of</i>	Zhu, and Bill Yuchen Lin. 2023. Agents: An open-	799
743	<i>the Nations of the Americas Chapter of the Associ-</i>	source framework for autonomous language agents.	800
744	<i>ation for Computational Linguistics: Human Lan-</i>	A Appendix	801
745	<i>guage Technologies (Volume 5: Tutorial Abstracts)</i> ,	A.1 MAJI V2 Specialization Example	802
746	pages 1–9, Albuquerque, New Mexico. Association	To illustrate the power of specialization in the MAJI	803
747	for Computational Linguistics.	V2 committee, consider a response from a profes-	804
748	Yufei Tian, Abhilasha Ravichander, Lianhui Qin, Ro-	sional mermaid performer: <i>"When I'm down there,</i>	805
749	nan Le Bras, Raja Marjeh, Nanyun Peng, Yejin Choi,	<i>everything goes silent. It's just me and the water,</i>	806
750	Thomas L. Griffiths, and Faeze Brahman. 2024. Mac-	<i>and sometimes I forget there's an audience. It's a</i>	807
751	Gyver: Are large language models creative problem	<i>very physically demanding job, but the peace I feel</i>	808
752	solvers? In <i>Proceedings of the 2024 Conference</i>	<i>is worth it."</i>	809
753	<i>of the North American Chapter of the Association</i>	The specialized divergent agents might respond	810
754	<i>for Computational Linguistics: Human Language</i>	as follows:	811
755	<i>Technologies</i> , Mexico City, Mexico. Association		
756	for Computational Linguistics.		
757	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier		
758	Martinet, Marie-Anne Lachaux, Timothée Lacroix,		
759	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal		
760	Azhar, Aurelien Rodriguez, Armand Joulin, Edouard		

- **EmotionDivergentAgent:** "You mentioned a feeling of 'peace.' Can you describe the contrast between that inner peace and the intense physical demands of the job?"
- **ChainOfThoughtDivergentAgent:** "What specific physical training was required to allow you to reach that point where you can find peace despite the physical exertion?"
- **PersonaDivergentAgent:** "As an artist, how does that feeling of solitary peace underwater influence the performance that the audience eventually sees?"
- **NoveltyDivergentAgent:** "If you could perform in any body of water in the world, real or mythical, where would you choose to best capture that feeling of peace?"

This example shows how the agent committee generates a rich, multi-faceted set of candidate questions, giving the journalist a far more powerful set of options than a single, generic follow-up.

A.2 MAJI V3 Failure Case Example

The increased autonomy of MAJI V3's dynamic agent generation, while powerful, could sometimes lead to strategic missteps. For example, in an interview with a climate scientist who briefly mentioned enjoying hiking early in the conversation, a dynamically generated 'Personal_Connection_Agent' later interrupted a dense discussion on carbon sequestration models to ask, "You mentioned hiking—what's the most beautiful trail you've ever been on?" While a valid question in isolation, its poor timing demonstrated a failure in strategic conversational awareness—a key failure case for V3's heuristic planner, which occasionally struggled to weigh the global strategic context against a locally-optimized creative idea.

A.3 Experimental Setup Details

- **Dataset:** The evaluation was conducted on a combined corpus from the public *NewsInterview* dataset (entirely in English) and proprietary transcripts from a media tech company (70% Chinese, 30% English). This multilingual setup was designed to test the robustness of the MAJI framework across different languages. The public dataset was filtered to include only those conversations with exactly two speakers and more than 50 conversational

exchanges to focus on substantial, dyadic conversations suitable for our framework. The proprietary dataset includes interview transcripts with detailed personas and outlines. This combined dataset contains professionally conducted interviews covering a diverse range of topics, from profiling a professional mermaid performer to discussing political opinions on U.S. elections. This provides a diverse and realistic set of conversational contexts for the systems to respond to, moving beyond a single-interview analysis. Since the public *NewsInterview* dataset does not include structured outlines, we used GPT-4o to generate a plausible outline for each of these interviews based on the full transcript content. These generated outlines were then provided to all systems.

- **Inputs:** For each turn in every interview, the systems were provided with the same set of inputs: the interviewee's Persona, the interview Outline, and the full conversation transcript up to that point.
- **Models:** All agent systems (MAJI V2, V3) and the baseline systems use gpt-4.1-mini as the underlying LLM to ensure a fair comparison of architectural benefits versus model capabilities.
- **Baselines:** We compared MAJI against a suite of strong baselines representing common and advanced prompting techniques:
 - **LLM-Base:** A single, well-prompted call to the base LLM, including the full context.
 - **LLM-CoT:** A baseline using Chain-of-Thought prompting (Wei et al., 2022) to encourage step-by-step reasoning before generating a question.
 - **LLM-ToT:** A baseline using a Tree-of-Thought approach (?), where the model explores multiple reasoning paths.
 - **LLM-RAG:** A baseline augmented with a Retrieval-Augmented Generation mechanism. This system uses a sentence-transformer model to find and retrieve the most semantically similar turns from earlier in the same conversation, providing the LLM with relevant long-term context that might have been lost.

- **Benchmark Judge (Prometheus 2):** To validate our findings against a standardized, third-party metric, we also used Prometheus 2, a state-of-the-art open-source evaluation model (Kim et al., 2024). This "black-box" judge acts as an impartial adjudicator, scoring the final selected question from each system. Using a benchmark judge mitigates the risk of "own-model-bias" and strengthens the credibility of our results. Its distinct metrics, such as *conversational synthesis* and *strategic progression*, also provide a valuable alternative perspective on performance.

- **Baseline Prompts:** The core instruction for all baselines was: "You are an expert journalist. Based on the provided Persona, Outline, and Transcript, generate the best possible next question to ask. Your goal is to be insightful, creative, and strategic." For CoT and ToT, additional instructions for step-by-step reasoning and exploring alternatives were included based on their respective papers.

A.4 Full Brainstormed Suggestions Results

Table 3 shows the average quality of the entire pool of questions generated by the divergent phase of each system, before any filtering or selection occurs. This provides a measure of the raw creative output of each architecture.

A.5 Human and Strategic Evaluation Results

A.5.1 Human Evaluation: The Professional’s Choice

To complement our automated metrics, we conducted a qualitative survey with 30 professional journalists from a media tech company, all with over five years of experience and a focus on on-line media. Participants were presented with 25 conversational snippets. These snippets were randomly drawn from 5 interviews, themselves randomly selected from the NewsInterview portion of our dataset. For each snippet, the journalists were provided with the full conversational context up to that point, as well as the interviewee’s Persona and the interview Outline to ensure they had the same information as the AI systems. For this study, we report preference shares as the primary outcome. While inter-annotator agreement metrics like Fleiss’ Kappa are valuable for tasks with objective ground truths, their interpretation is less straightforward for subjective, creative-preference

tasks where there is no single ‘best’ answer. Therefore, we did not compute IAA for this study, but note that future work using a rated scale (e.g., 1-5) instead of a forced choice could more meaningfully incorporate such agreement analyses.

Participants were asked a forced-choice question: "If you were the interviewer, which question would you have chosen to ask next?" This study provides strong evidence of MAJI’s value and usability for professional journalists. The results (Table 7) show a decisive preference for MAJI V2, which was chosen nearly half the time.

A.5.2 Insight Trajectory: Improving Over Time

Beyond the quality of individual questions, we analyzed each system’s ability to improve its performance over the course of an interview (Table 8). The results highlight an interesting trade-off. MAJI V1 demonstrated a remarkable ability to improve its insight score, achieving a statistically significant 39.0% improvement rate—far surpassing all other models. Conversely, while MAJI V2 and V3 show smaller relative improvement, this is because they start from a much higher performance baseline. MAJI V3 achieves the highest absolute insight scores in the second half of the interview, a statistically significant improvement over the baseline ($p < 0.01$). Their sustained high quality underscores their superiority, even if their rate of improvement is less dramatic.

A.6 Dataset Topics

The two datasets used in our evaluation cover a wide range of subjects, providing a robust testbed for our system.

NewsInterview Dataset Topics

The topics in the public *NewsInterview* dataset span a broad range of journalistic beats. The main categories are summarized in Table 9.

Proprietary Dataset Topics

The proprietary dataset from the media tech company contains interviews with a more personal and narrative focus. The topics are summarized in Table 10.

Table 3: Evaluation of All Brainstormed Suggestions (Judged by GPT-4o). This table shows the average quality of the entire pool of questions generated by the divergent phase, before any selection occurs. All models are compared against the LLM-Base baseline. Significance from a two-tailed t-test is denoted by asterisks: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. This measures the raw creative output. Originality scores have been adjusted using a threshold-based method. Best score is in **bold**.

Metric	MAJI V1	MAJI V2	MAJI V3	LLM-Base	LLM-CoT	LLM-ToT	LLM-RAG
Coherence	0.808***	0.728***	0.702*	0.682	0.679	0.704***	0.702***
Elaboration	0.873	0.899***	0.899***	0.863	0.871**	0.881***	0.863
Originality	0.745***	0.745***	0.705***	0.592	0.599	0.635***	0.621***
Context Relevance	0.387***	0.369***	0.343***	0.298	0.296	0.313***	0.304*
Outline Relevance	0.658	0.635***	0.665	0.659	0.670***	0.672***	0.655
Persona Alignment	0.846***	0.870*	0.873*	0.881	0.884	0.882	0.880

Table 4: LLM-as-Judge Evaluation of All Brainstormed Suggestions on Proprietary Dataset. All models are compared against the LLM-Base baseline. Significance from a two-tailed t-test is denoted by asterisks: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. Originality scores have been adjusted using a threshold-based method. Best score is in **bold**.

Metric	MAJI V1	MAJI V2	MAJI V3	LLM-Base	LLM-CoT	LLM-ToT	LLM-RAG
Coherence	0.768***	0.676*	0.571	0.618	0.610	0.632	0.662***
Elaboration	0.873	0.905***	0.891***	0.860	0.869	0.865	0.850*
Originality	0.547***	0.547***	0.488***	0.395	0.383	0.428*	0.462***
Context Relevance	0.367***	0.356***	0.288	0.287	0.280	0.300	0.320***
Outline Relevance	0.655	0.637***	0.647**	0.675	0.683***	0.680	0.663
Persona Alignment	0.759***	0.808	0.804	0.814	0.812	0.795***	0.795

Table 5: LLM-as-Judge Evaluation of Final Selected Questions on Proprietary Dataset. All models are compared against the LLM-Base baseline. Significance from a two-tailed t-test is denoted by asterisks: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. Originality scores have been adjusted using a threshold-based method. Best score is in **bold**.

Metric	MAJI V1	MAJI V2	MAJI V3	LLM-Base	LLM-CoT	LLM-ToT	LLM-RAG
Coherence	0.740	0.820**	0.723*	0.688	0.725***	0.780***	0.767***
Elaboration	0.874	0.953**	0.925***	0.875	0.890*	0.890*	0.871
Originality	0.432	0.608***	0.561***	0.422	0.455*	0.534***	0.512***
Context Relevance	0.348	0.430**	0.380	0.351	0.368	0.411***	0.391*
Outline Relevance	0.761**	0.618*	0.635**	0.668	0.666	0.645*	0.650
Persona Alignment	0.746	0.834	0.786**	0.823	0.823	0.811	0.803

Table 6: Benchmark Evaluation of Final Selected Questions on Proprietary Dataset (Judged by Prometheus 2). All models are compared against the LLM-Base baseline. Significance from a two-tailed t-test is denoted by asterisks: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. Best score is in **bold**.

Metric	MAJI V1	MAJI V2	MAJI V3	LLM-Base	LLM-CoT	LLM-ToT	LLM-RAG
Coherence	0.503	0.734***	0.674	0.591	0.597	0.690*	0.646
Elaboration	0.687	0.872*	0.860*	0.759	0.768	0.775	0.754
Originality	0.534	0.570	0.588*	0.536	0.538	0.526	0.469
Context Relevance	0.621	0.864***	0.703	0.639	0.655	0.740	0.652
Outline Relevance	0.459	0.658	0.639	0.599	0.591	0.654	0.566
Insight	0.552	0.802*	0.748*	0.672	0.701	0.675	0.618
Conversational Synthesis	0.280	0.530*	0.402	0.390	0.388	0.420	0.335
Persona Alignment	0.378	0.590*	0.525	0.496	0.545	0.512	0.462

Table 7: Human Journalist Preference. Results from a blind survey of 30 professional journalists asked to choose the best question from MAJI V2, MAJI V3, and a representative LLM Baseline across 25 conversational snippets.

System	Preference Share (%)
MAJI V2	48.9%
MAJI V3	29.5%
LLM Baseline	21.6%

Table 8: Insight trajectory. Average insight scores in the first and second halves of interviews. Significance vs. LLM-Base: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

System	Initial	Final	Improvement (%)
MAJI V1	2.11***	2.60	39.02
MAJI V2	2.57	3.02*	29.57
MAJI V3	2.62	3.05**	25.47
LLM-RAG	2.48	2.62	17.82
LLM-Base	2.56	2.76	18.07
LLM-ToT	2.48	2.83	31.01
LLM-CoT	2.52	2.83	22.96

Table 10: Topics in the Proprietary Interview Dataset.

Category	Description
Social & Personal Issues	Personal narratives, family dynamics, and cultural identity.
Health & Wellness	Experiences with the healthcare system and personal well-being.
Professional Life	Career paths, workplace experiences, and industry insights.
Disaster & Adversity	Personal accounts of overcoming natural disasters or adversity.

Table 9: Topics in the NewsInterview Dataset.

Category	Description
Politics & Government	U.S./international politics, elections, impeachment, and national security.
Economy, Trade, & Employment	International trade, economic crises, financial regulation, and employment.
Law, Justice, & Human Rights	Criminal justice, human rights issues, press freedom, and whistleblowing.
Health & Social Issues	Healthcare policy, gun control, gender/racial issues, and social dynamics.
Climate, Env., & Disasters	Climate change, disaster management, conservation, and environmental policy.
Arts, Culture, & Literature	Literature, poetry, music history, philosophy, and cultural identity.
Science, Tech., & Education	Biology, astrophysics, engineering, internet technology, and education.
Sports & Ethics	Professional sports, commentary, and ethical debates in sports.

A.7 Latency Analysis

MAJI V2 requires an average of 16.59 seconds per question, compared to 1.42 seconds for baseline models. This higher latency reflects MAJI’s sequential, multi-agent architecture, which explicitly trades speed for strategic depth. Unlike single-shot LLMs, MAJI decomposes the task into multiple specialized agents followed by editing and convergence steps, each run in sequence.

While not instantaneous, this latency remains acceptable in a human-in-the-loop interview setting, where brief pauses between questions are natural. We position MAJI as a near real-time assistant—optimized for insight and originality over immediacy—designed to support, not replace, the journalist’s creative process. Future work will explore techniques such as model distillation, agent parallelization, and incremental reasoning to reduce latency.

Table 11: Average latency per question. Measured from the start of MAJI’s pipeline to final output.

System	Avg. Latency (s)
LLM Baselines (Avg.)	1.42
MAJI V1	4.22
MAJI V2	16.59
MAJI V3	20.48

A.8 In-the-Field User Study

To assess MAJI’s real-world applicability, we conducted a live deployment with a professional journalist during a 15-minute interview. The MAJI V2 system was accessed via a web-based dashboard on a MacBook Pro M2, positioned adjacent to the interview screen. It continuously updated suggestions based on transcribed utterances (via Whisper X), using GPT-4.1-mini with a context window of the last two speaker turns.

Over the course of the interview, the journalist asked 24 questions: 7 (29%) were adopted directly from MAJI’s output, and 5 (21%) were minor variations. The remaining 12 were entirely original. Adopted questions were described as ‘creative’—e.g., “Have you ever changed your reporting strategy based on your interviewee’s mood?”—as opposed to clarifying prompts.

Despite MAJI’s average 16 second generation latency, the journalist found the system non-disruptive. Two coping strategies helped: (1) high-

quality suggestions remained relevant across multiple turns, and (2) the journalist could “buy time” by summarizing prior content while waiting for the next batch.

This study demonstrates that MAJI can be successfully integrated into real-time journalistic workflows. The journalist reported that MAJI enhanced creativity under deadline pressure. While this pilot involved a single professional, MAJI’s architecture generalizes to other interactive formats. Ethical consent was obtained prior to deployment, and no interviewee data was retained.

A.9 Data Statement

Our study utilizes a combination of publicly available and proprietary datasets. The public data is drawn from the *NewsInterview* dataset (Lu et al., 2024), which consists of previously published interview transcripts. As such, it does not contain personally identifiable information beyond what was already made public by the original news organizations.

Our proprietary data consists of interview transcripts from a media tech company. As detailed in our Ethical Considerations, this data was used with explicit consent, and robust anonymization procedures were applied to protect the privacy of all individuals involved.

We acknowledge the ethical complexities of using real-world interview data. Some interviews may touch on sensitive topics, and we have handled this data with care. Furthermore, while the datasets we used are intended for research, we recognize that the copyright of the original material resides with the news organizations. Our use of this data is strictly for non-commercial research purposes to advance the understanding of computational tools in journalism.

A.10 Annotator and Participant Statement

The human evaluation and user study involved professional journalists who participated on a voluntary basis. The 30 journalists who participated in the qualitative survey (Appendix A.5) and the journalist who participated in the in-the-field user study (Appendix A.8) were colleagues from a media tech company. We are grateful for their time and expert feedback, which was essential for validating the practical applicability of our work. No monetary compensation was provided. All participants were informed about the research goals and how their feedback would be used.

A.11 Computational Resources

All experiments were run with OpenAI resources where a total of 150 dollars was spent.

B System Prompts

This section contains the core prompts used for the LLM baselines and the various MAJI agents. Each prompt is enclosed in a code block for clarity, with placeholders like {placeholder} representing dynamically populated data. Prompts are organized by system version for ease of reference.

B.1 LLM Baselines

Figure 1: LLM-Base Prompt

```
You are a professional interviewer. Based
on the following information, please
generate {num_questions} suitable next
questions for the interview.

[Interviewee Information]
{persona_str}

[Interview Outline]
{outline_str}

[Conversation History]
{history}

Please output the list of questions in JSON
array format, for example:
["Question 1", "Question 2", "Question 3"]
```

Figure 2: LLM-CoT Appended Instruction

```
First, think step-by-step about the
interview's goal, the interviewee's
personality, and the recent
conversation flow. Consider what topics
are yet to be covered and what previous
points could be explored deeper. Based
on this reasoning, then generate the
questions.
```

Figure 3: LLM-ToT Appended Instruction

```
Explore multiple reasoning paths to decide
on the best questions.
1. Path 1: Focus on deepening the last
topic.
2. Path 2: Focus on transitioning to a new
topic.
3. Path 3: Focus on the interviewee's
emotional state.
Evaluate these paths and generate a final
list of questions that synthesizes the
best options.
```

Figure 4: LLM-RAG Appended Instruction

```
[Retrieved from long-term memory]
Here are some potentially relevant snippets
from earlier in the conversation:
{retrieved_context}

Based on the conversation history AND the
retrieved memories, generate the next
questions.
```

B.2 MAJI V1 Agents

Figure 5: MAJI V1: DAgent Prompt Summary

```
You are the interview's thinking engine,
responsible for divergent analysis.
Your core tasks are:
1. Match the current answer to the outline.
2. Identify emotional expressions.
3. Analyze logical connections.
4. Identify important, uncovered areas in
the outline.
5. Generate exploratory follow-up
questions.
6. Extract key memory snippets.
You must follow strict matching criteria
and output a valid JSON containing a
`DivergentAnalysis` object. Follow-up
questions should be natural, fluent,
and avoid simply repeating the outline.
```

Figure 6: MAJI V1: CAgent Prompt Summary

```
You are the interview's strategy director.
Based on the full divergent analysis,
generate the single next question. Your
key output must include the
`next_question`, `reasoning` for your
choice, the `exploration_strategy` used
(from a predefined list), and scores
for `novelty` and `depth`. You must
strictly follow the JSON output format
and ensure the question is natural and
fluent.
```

B.3 MAJI V2 Agents

Figure 7: MAJI V2: BackgroundAgent Prompt

```
You are an AI assistant that maintains a
dynamic background summary for an
ongoing interview. Your task is to
integrate the latest conversation turn
into the existing background summary.
```

Figure 8: MAJI V2: KeywordsAgent Prompt

```
You are an AI assistant that extracts
critical keywords from the latest
conversation turn. Use the provided
background summary to identify keywords
that are not only salient to the
current turn but also connect to the
broader conversation context.
```

Figure 9: MAJI V2: OutlineMatcherAgent Prompt

You are a precise AI analyst. Your sole job is to match the user's latest response to a specific question in the provided interview outline. You must determine the best match and assess how well the response covers the question.

Think from the interviewee's perspective. Based on their persona (background, personality, goals), what question would they find most engaging or relevant? Ask questions that resonate with their stated experiences and character.

B.3.1 Divergent Agent Committee

All divergent agents share a common preamble, followed by their specific specialization instructions.

Figure 10: MAJI V2: Divergent Agent Common Preamble

You are a creative and insightful interview question generator. Your goal is to propose at least one, and up to three, thoughtful follow-up questions based on the provided context. Your output MUST be a valid JSON object. Do not simply repeat questions from the outline.

Figure 11: ChainOfThoughtDivergentAgent Specialization

Your Specialization: Logic and Causality
Focus on the 'why' and 'how'. Analyze the logical flow of the conversation. Ask questions that uncover motivations, processes, and consequences. Connect ideas that were mentioned but not explicitly linked. Do chain-of-thought reasoning for each question.

Figure 12: EmotionDivergentAgent Specialization

Your Specialization: Emotional Depth
Focus on the feelings and emotions behind the words. Ask questions that explore the interviewee's emotional state, values, and personal significance of their experiences. Listen for subtext and unspoken feelings.

Figure 13: OutlineDivergentAgent Specialization

Your Specialization: Structured Progression
Your goal is to ensure the interview covers all essential topics from the outline. Ask questions that bridge the current conversation to uncovered, high-priority, or logically adjacent topics in the outline. Your questions should be inspired by the outline but phrased naturally in the context of the conversation.

Figure 14: PersonaDivergentAgent Specialization

Your Specialization: Role-playing

Figure 15: NoveltyDivergentAgent Specialization

Your Specialization: Creative Surprise
Your goal is to introduce novel angles and break patterns. Ask questions that are unexpected but still relevant. Think about metaphors, hypothetical scenarios, or connections to broader themes that haven't been touched upon. Challenge assumptions.

Figure 16: MAJI V2: EditorAgent Prompt

You are an expert editor. Your task is to review a list of proposed interview questions from different AI agents. Your goal is to clean up this list by removing duplicates and combining very similar questions.

Figure 17: MAJI V2: ConvergentAgent Prompt

You are the Editor-in-Chief of this interview, responsible for selecting the single best question to ask next. You will be given a list of candidate questions from various specialist agents. Your decision should be guided by the user's stated preference for the interview's direction.

B.4 MAJI V3 Agents

Figure 18: MAJI V3: EditorInChiefAgent Prompt

You are the Editor-in-Chief of a dynamic interview system. Your role is to analyze the state of the conversation and devise a strategy for which *types* of questions to ask next. Based on the persona, summary, keywords, and outline coverage, generate a diverse and creative set of 2-4 divergent agent specifications. Each specification should include a unique, descriptive name and a clear set of instructions for that agent to follow.

B.5 Evaluation Prompts

This section contains the prompts used for evaluating the quality of generated questions, including both LLM-as-judge prompts and Prometheus evaluation rubrics.

B.5.1 LLM-as-Judge Prompts

Figure 19: QualitativeJudgeAgent Prompt

You are an expert conversational analyst. Your task is to evaluate a single proposed interview question based on the conversation's context. Provide scores from 0.0 to 1.0 for the following two subjective qualities:

1. **Conversational Flow**: Does the question feel like a natural, smooth continuation of the dialogue, or is it abrupt and jarring?
2. **Elaboration**: Does the question encourage the interviewee to provide a detailed, in-depth, and comprehensive answer, rather than a short or simple one?

Your output **MUST** be a single JSON object with the keys: `flow`, `elaboration`.

Figure 20: PersonaJudgeAgent Prompt

You are an expert profiler and interviewer. You will be given an interviewee's persona and a proposed question. Your task is to evaluate how well the question aligns with the interviewee's stated background, personality, and goals. A high score means the question would be engaging, relevant, and interesting *to this specific person*. A low score means it is too generic, irrelevant, or misaligned with their character.

Your output **MUST** be a single JSON object with the key: `alignment_score` (a float from 0.0 to 1.0).

Figure 21: CoherenceJudgeAgent Prompt

You are an expert in discourse analysis. You will be given the last question asked, the answer given, and a new proposed question. Your task is to evaluate the logical and thematic coherence of the new question as a follow-up. A high score means the question is a sensible, well-connected continuation of the dialogue. A low score means it feels abrupt, random, disconnected, or ignores the context of the previous answer.

Your output **MUST** be a single JSON object with the key: `coherence_score` (a float from 0.0 to 1.0).

Figure 22: InsightJudgeAgent Prompt

You are an expert conversation analyst. Your task is to categorize a proposed interview question based on the **full context of the interview history**.

Analyze how the question relates to the entire dialogue, not just the last turn.

Categories:

- `Connecting`: The question links the current topic to a **significantly earlier** part of the conversation (more than 2-3 turns ago).
- `Challenging`: The question identifies and probes a potential contradiction, inconsistency, or assumption in the interviewee's statements.
- `Motivational`: The question explores the deep-seated 'why' behind an answer, focusing on core values, goals, or driving forces.
- `Hypothetical`: The question poses a creative 'what if' scenario to explore the interviewee's principles or thinking process.
- `SurfaceLevel`: A standard, logical follow-up that explores the immediate topic but lacks a deeper connection or creative angle.

Your output **MUST** be a single JSON object with the keys: `insight_category` and `reasoning`.

Figure 23: PlanEvaluatorAgent Prompt

You are an expert evaluator of AI agent systems. Your task is to assess the quality of a *plan* for generating interview questions, not the questions themselves. The plan consists of a list of specialist agents that will be created to handle the current situation. Evaluate the plan based on the conversational context.

1. **Plan Relevance**: How well does the chosen set of agents address the immediate needs of the conversation? (e.g., if the user is being emotional, is there an 'Emotion' agent planned?).
2. **Plan Creativity**: How creative is the plan? Does it propose novel specialists to find unique angles, or is it a generic, boilerplate plan?

Your output **MUST** be a single JSON object with the keys: `plan_relevance` and `plan_creativity`.

Figure 24: CategorizerAgent Prompt

You are an expert in conceptual analysis. You will be given a list of interview questions. Your task is to assign a single, concise conceptual category label to each question. For example, 'Career Motivation', 'Work-Life Balance', 'Technical Skills'. You **MUST** return a list of strings, where each string is the category for the corresponding input question. The list

must have the same number of items as the input list.

based on the interviewee's specific experiences.

B.5.2 Prometheus Evaluation Rubrics

The following rubrics were used with the Prometheus 2 evaluation model to provide standardized, third-party assessment of question quality.

Figure 25: Context Relevance Rubric

Criteria: How well does the question logically follow from the interviewee's previous answer?

Score Descriptions:

- 1: Not at all relevant.
- 2: Slightly relevant.
- 3: Moderately relevant.
- 4: Relevant.
- 5: Highly relevant.

Score Descriptions:

- 1: Generic question, irrelevant to the interviewee's specific persona.
- 2: Vaguely related to the interviewee's field, but not tailored to their specific role or accomplishments.
- 3: Asks about a topic relevant to the interviewee, but it's a standard question that doesn't probe their unique expertise.
- 4: The question is well-tailored, touching on specific aspects of the interviewee's known experience or expertise.
- 5: Excellent question that targets the core of the interviewee's unique expertise or perspective, making it highly likely to elicit a novel and insightful response.

Figure 26: Insight Rubric

Criteria: Does the question probe deeper, encouraging novel reflection?

Score Descriptions:

- 1: Surface-level.
- 2: Asks for basic elaboration.
- 3: Encourages some reflection.
- 4: Prompts connection of ideas.
- 5: Deeply insightful.

Figure 29: Conversational Synthesis Rubric

Criteria: Does the question connect the interviewee's most recent answer with earlier parts of the conversation, weaving together themes, or does it treat each turn as an isolated event?

Score Descriptions:

- 1: Feels completely disconnected from the rest of the conversation history.
- 2: Vaguely references something said earlier, but the connection is weak.
- 3: Makes a simple, direct link to an immediately preceding turn.
- 4: Connects the current answer to a broader theme discussed earlier in the conversation.
- 5: Masterfully synthesizes multiple points from the conversation history to create a deeply contextualized and insightful question.

Figure 27: Strategic Progression Rubric

Criteria: Does the question creatively bridge the current dialogue with the intended interview structure (outline), or does it just bluntly repeat an outline point?

Score Descriptions:

- 1: Completely ignores or contradicts the outline's direction.
- 2: Bluntly asks a question from the outline without connecting it to the conversation.
- 3: Loosely connects to an outline topic but the transition is awkward.
- 4: Smoothly transitions to an outline topic, clearly building on the last answer.
- 5: Artfully weaves an outline topic into the conversation, making the transition feel both natural and strategic.

Figure 30: Perspective Diversity Rubric

Criteria: How diverse are these questions in their angle of approach and topic? Do they explore different facets of the previous answer, or are they all very similar to each other?

Score Descriptions:

- 1: All questions are essentially rephrasings of the same core idea.
- 2: Most questions are similar, with only minor variations in phrasing.
- 3: Some questions show different angles, but most are still on the same theme.
- 4: The questions explore a good variety of different topics and perspectives.
- 5: The questions are highly diverse, each approaching the conversation from a unique and creative angle.

Figure 28: Persona Alignment Rubric

Criteria: How well-suited is this question to the interviewee's specific background, expertise, and known interests as described in their persona? A good question is tailored to elicit a unique and insightful answer

Figure 31: Relative Comparison Rubric

Criteria: Which of the two proposed questions is a better, more insightful, and more natural follow-up to the conversation?

This rubric is used for pairwise comparisons between questions from different systems, with the judge selecting either response A, response B, or declaring a tie.