

Multi-Agent Reflection Loop: A Multimodal Framework for Cross-Cultural Narrative Adaptation

Anonymous ACL submission

Abstract

Cross-cultural narrative understanding demands language models to not only excel at text generation, but also perceive users' implicit cultural cognitive states and dynamically align with them. However, existing interactive narrative systems predominantly rely on static preset scripts or monolithic large language models (LLMs). These approaches fail to resolve semantic conflicts arising from multimodal user inputs (text, behavior, emotion), leading to feedback lacking cultural adaptability and even the generation of contextual illusions. To address these limitations, this paper proposes a Multimodal Culture-Aware Multi-Agent System (MC-MAS), which achieves high-precision cross-cultural narrative adaptation through a collaborative agent mechanism. Specifically, we design three specialized functional agents (behavioral, linguistic, cultural) to process heterogeneous user signals, and introduce a core coordinator agent. The coordinator employs a novel "reflection-reconstruction" loop mechanism, which can automatically detect cross-modal consistency conflicts and iteratively optimize narrative generation strategies. We validate the MC-MAS framework in a narrative scenario rooted in the exile literature of the Weimar Republic. Experimental results demonstrate that, compared with static models and single LLMs, our method significantly enhances the accuracy of cultural context alignment while preserving narrative coherence, and effectively alleviates users' cognitive load.

1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities in role-playing and interactive storytelling, enabling the construction of complex fictional worlds guided by user instructions (Bartle, 2004; Dietz et al., 2023). These advances have made LLMs promising engines for interactive narrative systems and serious games, particularly in culturally situated learning and heritage-

related contexts (Mortara et al., 2014; Titus and Ng'ambi, 2023). However, when such narratives are situated in cross-cultural contexts, the challenge extends beyond linguistic fluency or plot coherence. Effective interaction requires the system to accurately interpret users' culturally grounded cognitive states and dynamically adapt narrative strategies to achieve cultural alignment, an issue widely recognized in cross-cultural game-based learning research (Jossan et al., 2021; O'Hagan, 2009).

Crucially, cultural misalignment does not always manifest as explicit misunderstanding in textual input (Shen et al., 2025; Li and Li, 2025). Instead, it often emerges through inconsistencies across multiple modalities (Pescuma et al., 2023; Yang et al., 2025; Wafa et al., 2025). A user's text may appear logically coherent, while their in-game behavior, interaction patterns, or facial expressions (Zhang et al., 2024; Fan et al., 2023; Cloude et al., 2022) reveal confusion toward culturally specific metaphors, such as religious allusions or historical analogies. Such phenomena are consistent with findings in multimodal cultural meaning-making and perception research, where cultural understanding is distributed across linguistic, behavioral, and visual channels rather than confined to text alone (Cabezas-García and Reimerink, 2022; Xie, 2025). These discrepancies reflect deeper cognitive conflicts shaped by cultural background rather than surface-level language errors. As a result, cross-cultural narrative adaptation inherently becomes a multimodal inference problem, aligning with recent work on emotion-aware and multimodal interpretation in interactive systems (Vistorte et al., 2024).

Existing interactive narrative systems face several limitations in cross-cultural contexts:

- **Limited multimodal understanding:** Most systems rely primarily on textual input, ignoring behavioral and affective signals.

084	• Inability to reconcile conflicting cues: When	localization experts) and proposes the Monolin-	130
085	signals from different modalities contradict	gual Human Preference (MHP) metric, which tran-	131
086	each other, monolithic models cannot resolve	scends literal matching.	132
087	the conflicts effectively.		
088	• Insufficient evaluation of user understand-	Multimodal Analytics and Cultural Semantic	133
089	ing: Current approaches focus on narrative	Alignment To achieve deeper cognitive and se-	134
090	fluency rather than whether the system cor-	semantic understanding, research attention has turned	135
091	rectly interprets users' cognitive and cultural	to the deep fusion of heterogeneous data and the	136
092	states.	calibration of model values. In learning analyt-	137
093	To address these issues, we propose the Multi-	tics, MSC-Trans (Xie et al., 2025) and Yan et al.	138
094	modal Culture-Aware Multi-Agent System (MC-	(Yan et al., 2025) integrate data streams, such as fa-	139
095	MAS), which provides the following contributions:	cial expressions, postures, and logs, to resolve data	140
096	• Multimodal multi-agent framework: MC-	asynchrony issues and accurately detect student en-	141
097	MAS integrates specialized expert agents to	gagement, while Yan et al. (Yan et al., 2025) further	142
098	analyze behavioral, linguistic, and cultural sig-	demonstrate the synergy of Generative AI in imput-	143
099	nals, overcoming the limitation of text-only	ing missing longitudinal data. Addressing cultural	144
100	models.	consistency in language models, CultureLLM (Li	145
101	• Coordinator with reflection mechanism: A	et al., 2024a) effectively augments data using the	146
102	central coordinator reconciles conflicting out-	World Values Survey. Meanwhile, MENAValues	147
103	puts from different agents, enabling robust	(Zahraei and Asgari, 2025) reveals value inconsis-	148
104	interpretation of inconsistent cues.	tencies and “log-odds leakage” challenges during	149
105	• Human-level evaluation: We validate MC-	“cross-lingual value transfer” through benchmark-	150
106	MAS on a human-annotated benchmark,	ing, highlighting the critical importance of captur-	151
107	showing near-human performance in user pro-	ing profound cultural nuances.	152
108	filings and cross-cultural intent recognition.		
109	2 Related Work	3 Method	153
110	Multi-Agent Collaboration and Simulation En-	This section presents a multimodal analysis frame-	154
111	vironments Current education and professional	work based on a Multi-Expert Committee. Culture,	155
112	training are shifting from single-agent to multi-	Behavior, and Language experts analyze students'	156
113	agent collaborative environments, leveraging Gen-	expressions, actions, and text in parallel, while a	157
114	erative AI to simulate complex interaction dynam-	Referee Module fuses their outputs and resolves	158
115	ics. SimClass (Zhang et al., 2025) constructs a	conflicts, enabling accurate modeling of cultural	159
116	comprehensive classroom environment with di-	profiles and gaming personas for adaptive teaching.	160
117	verse peer roles, utilizing the Flanders Interac-		
118	tion Analysis System (FIAS) to reproduce authen-	3.1 Data Collection	161
119	tic teaching interactions. In the medical field,	Gardner's Theory of Multiple Intelligences argues	162
120	MEDCO (Wei et al., 2024) and <i>Agent Hospital</i> (Li	that traditional standardized testing, largely based	163
121	et al., 2024b) utilize expert-patient triads and the	on written assessments, offers a reductive view of	164
122	“Simulacrum-based Evolutionary Agent Learning”	cognition and overlooks the diversity of student	165
123	(SEAL) paradigm, respectively, to facilitate the evo-	potential. Although seven intelligence dimensions	166
124	lution of interdisciplinary reasoning and diagnostic	are defined, practical system constraints require a	167
125	capabilities. This collaborative model has also been	more focused implementation. Therefore, we col-	168
126	applied to the cultural domain; TRANSAGENTS	lect data from three primary modalities—facial ex-	169
127	(Wu et al., 2025) produces literary translations that	pressions, linguistic articulations, and operational	170
128	emphasize cultural resonance through a multi-agent	behaviors—to construct students' Cultural Profiles	171
129	virtual company workflow (including editors and	and Bartle gaming personas.	172
		Data collection is organized into three phases	173
		corresponding to three curriculum sessions. In each	174
		phase, players complete four tasks designed accord-	175
		ing to the Bartle Taxonomy to elicit behaviors of	176
		Achievers, Explorers, Socializers, and Killers. Data	177

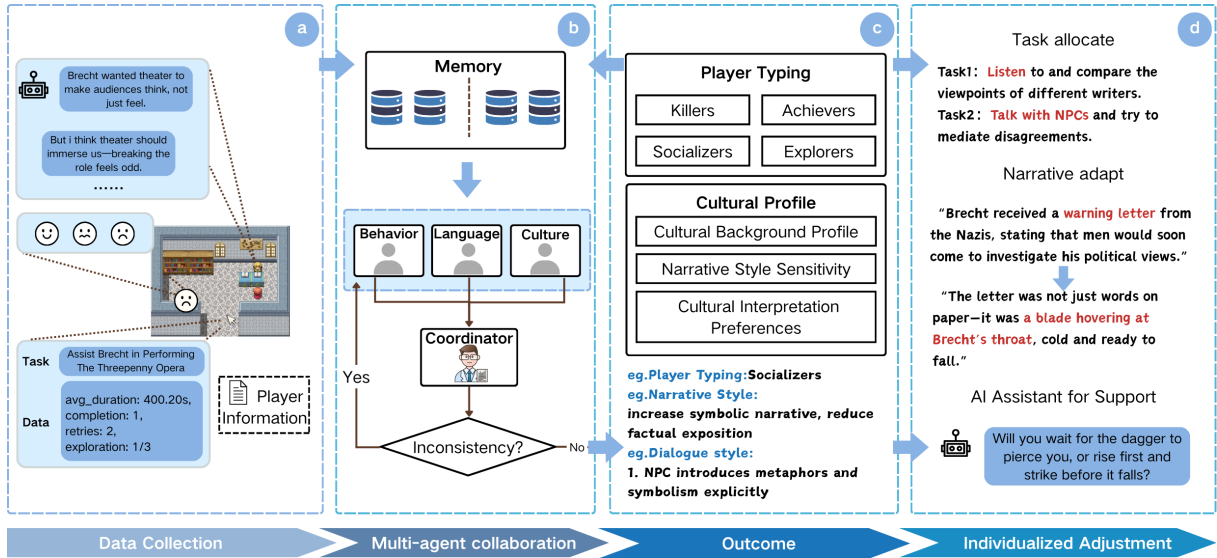


Figure 1: System workflow of the culturally adaptive game-based learning framework. (a) Data collection from player interactions, behavior logs, dialogues, and emotional cues; (b) multi-agent collaboration (behavior, language, culture agents with a coordinator for inconsistency resolution); (c) outcome generation including player typing and cultural profile; and (d) individualized adjustment through task allocation, narrative adaptation, and AI assistant support.

are analyzed sequentially, with outputs from each phase informing modeling in the next.

To capture facial expressions while preserving privacy, we adopt a discrete sequence sampling strategy. Instead of storing raw video, the system records emotion labels (e.g., Happy, Confused, Sad) at 5-second intervals, ensuring that no biometric facial data are retained. In parallel, we log fine-grained operational information, including task timestamps, failure counts, scene indices, AI dialogue history, and textual submissions. The collected data fields are detailed in Table 1.

Through this approach, we facilitate the unified temporal alignment of heterogeneous multimodal data streams, thereby laying a solid foundation for the subsequent analysis conducted by the Specialist Agents.

3.2 Workflow of Multi-Agent Reflective Cycles

Cultural adaptation in this system is implemented through a multi-agent collaborative mechanism. Learners’ multimodal in-game inputs (Figure 1-a) are processed by specialized agents, and their outputs are integrated and reconciled by a central master agent. This expert-like deliberative design improves interpretability and consistency while enabling structured reflection and iterative optimization. The overall workflow is shown in Figure 1-b.

Behavior Analysis Agent The Behavior Agent (A_{beh}) specializes in transforming discrete game logs into high-level representations of player motivational profiles. We formally define the player’s behavioral trajectory within a time window t as $\mathcal{B}_t = \{r_c, t_m, e_{ratio}, i_{freq}\}$. To bridge low-level metrics with psychological traits, A_{beh} functions as a probabilistic reasoner conditioned on behavioral psychology principles. Instead of simple rule-based mapping, the agent utilizes a Chain-of-Thought (CoT) reasoning process defined by our system instructions. It first interprets the raw metrics in the context of game mechanics (e.g., distinguishing "exploration" from "getting lost") and then maps these interpretations to the four Bartle player types (Killer, Achiever, Socializer, Explorer). The agent maximizes the conditional probability $P(\mathbf{P}_{bartle} | \mathcal{B}_t, \mathcal{K}_{psych})$, where \mathcal{K}_{psych} represents the injected domain knowledge regarding player typologies. The output is a normalized distribution vector \mathbf{P}_{bartle} , quantifying the player’s dominant and auxiliary motivations.

Language Analysis Agent The agent is tasked with analyzing the player’s dialogue content and facial expressions during gameplay. By extracting key features that reflect the player’s psychological state, it provides a basis for the dynamic adjustment of narrative content. We quantify the player’s expressive patterns through three dimen-

Table 1: Selected Interaction Log: Brecht’s Study Scenario

Time	Facial Expr.	Scene	User / AI Interaction	Input Buffer
0–5s	Neutral	Brecht’s living room	None	None
...				
60–65s	Anxious	Brecht’s study	User: “Is it... immersive experience?”	“Immersion”
70–75s	Happy	Brecht’s study	User: “Ah! The alienation effect? Breaking the fourth wall?”	“Alienation”
...				
115–120s	Happy	Brecht’s study	AI: “Task completed. Your understanding is correct.”	None

sions: Metaphorical Density, Syntactic Complexity, and Linguistic Profile (High/Low Context).

To implement this, the agent aligns the player’s text sequence (S_{text}) with facial expression labels (S_{face}). The feature extraction process is governed by a structured analysis protocol inspired by sociolinguistic frameworks. Specifically, the model is instructed to: Detect semantic alignment between player utterances and in-game lore to evaluate immersion; Parse syntactic structures to infer cognitive load (e.g., identifying fragmentation as a sign of confusion); Classify communication styles based on information density. Formally, the extraction can be denoted as $\mathcal{P}_{ling} = \text{Model}(S_{text}, S_{face} | I_{instruct})$, where $I_{instruct}$ encapsulates the definitions of the aforementioned dimensions. The resulting vector \mathcal{P}_{ling} serves as a real-time “probe,” distinguishing between states of *high immersion* and *cognitive overload*.

Cultural Context Analysis Agent The Cultural Agent, denoted as A_{cult} , functions as a bridge between the player’s identity background and their in-game actions, effectively acting as a cross-cultural communication and historical consultant. It integrates the player’s static profile $I_{profile}$ (e.g., nationality, native language) with a dynamic behavior log $\mathcal{B} = \{r_c, t_m, e_{ratio}, i_{freq}\}$ to construct a dynamic cultural model.

The core mechanism of A_{cult} involves Cultural Incongruity Detection. The agent integrates the static profile $I_{profile}$ with dynamic behavioral data to identify friction points. The reasoning logic follows a cause-effect derivation: it first scans for specific triggers (e.g., religious metaphors, historical allusions), checks the player’s reaction (e.g.,

hesitation, negative sentiment), and then evaluates whether the disconnect stems from a cultural knowledge gap.

The output mechanism is two-fold:

- Risk Quantification ($O_{RiskAlert}$): A discrete classification (High/Medium/Low) predicting the likelihood of misunderstanding.
- Intervention Strategy ($O_{BridgingSuggestion}$): A generative task where the agent proposes specific remediation strategies, such as adding explanatory asides or simplifying cultural references, tailored to the specific type of detected incongruity.

Coordinator Agent: The Committee Chair Multimodal modeling inevitably engenders conflicts among agents. Sole reliance on initial inferences may lead to strategic distortion. To address this, we introduce a Reflection Round Mechanism that resolves multimodal conflicts through iterative feedback and model reconstruction. To address this, we introduce a Reflection-Arbitration Mechanism that strictly follows an “Identify-Remodel-Fuse-Update” loop:

Conflict Identification: The Coordinator aggregates outputs from sub-agents and calculates semantic discrepancies. A reflection round is triggered solely when conflicts exceed a pre-defined severity threshold.

Context-Aware Prompt Construction: Instead of generic re-prompting, the Coordinator dynamically synthesizes a meta-prompt. This prompt includes the conflicting conclusions, original evidence, and specific directives for re-evaluation, effectively contextualizing the reflection task.

Iterative Remodeling: Specialized agents re-process their inputs under the guidance of the meta-prompt. For instance, the Behavioral Agent may incorporate cultural context provided by the Coordinator to re-interpret a "retry" action not as failure, but as curiosity.

Revision Fusion & Decision: The Coordinator performs a final weighted aggregation of the revised outputs to generate the adaptive decision D_{final} .

This mechanism ensures that the final decision is not merely a statistical average but a result of deliberative consensus, balancing interpretability with robustness against single-modal biases. The detailed prompt specifications for each agent are provided in Appendix B.

3.3 Task and Expression Intelligent Modification

To ensure a tiered learning objective structure, we referenced Bloom’s Taxonomy (Blyth et al., 1966; Ormell, 1974) to divide each character’s learning tasks into three stages:

Stage One (Remember and Understand): Players complete initial tasks to familiarize themselves with the game’s fundamentals. This stage features four task types aligned with Bartle Player Types (Bartle, 2004) (Achiever, Explorer, Socializer, Killer) (Table 2) to cater to diverse learner preferences. During this process, the system collects multimodal player data, including text inputs, facial expressions, and operational behaviors.

Stages Two and Three (Apply and Analyze) (Evaluate and Create): The system employs a hierarchical dynamic adaptation mechanism (Figure 2). Vertically, data collected in prior stages feeds into subsequent stages to holistically adjust task allocation, narrative style, and difficulty design (inter-stage adaptation). Horizontally, the AI assistant receives real-time analysis results from the current stage, enabling instant adjustments to dialogue strategies and task prompts (intra-stage adaptation). For instance, the system increases annotation points and background extension tasks for exploration-oriented players, while introducing culturally contextualized explanations in narratives for players from diverse backgrounds.

This phased and personalized design not only deepens learners’ comprehension of literary knowledge but also enhances learning adaptability across cultural contexts through dynamic regulation. Figure 1 illustrates the dynamic adaptation of in-game

content. Specifically, Figure 1-a demonstrates the personalized task allocation tailored to player characteristics; Figure 1-b presents the script modifications informed by player profiles; and Figure 1-c highlights the individualized support offered by the AI system.

4 Experiment Setup

The following experimental setup is designed to assess both the inferential accuracy of the proposed multimodal and multi-agent reflective approach and its effects on players’ cognitive load in the adaptive literary learning system.

4.1 Accuracy Evaluation

To validate the effectiveness of the Multi-Agent Reflection-Reconstruction Loop within the MC-MAS framework, specifically in the tasks of User Profiling and Cultural Intent Recognition, we designed a comparative experiment benchmarked against human expert annotations. The data collection protocols were reviewed and approved by the Institutional Review Board (IRB) of the authors’ institution.

4.1.1 Dataset and Human Benchmark

Data Collection. Experimental data was derived from authentic gameplay logs. We recruited a total of 20 volunteer participants, consisting of undergraduate students majoring in German literature and enthusiasts with foundational proficiency in German (10 males, 10 females). From the raw dataset, we screened and selected 11 high-quality, complete interaction records (5 males, 6 females) to form the core test set. These records cover all three pedagogical phases of the game and contain unedited *Action Logs*, *Dialogue Text*, and *Facial Expression Sequences*. This multimodal dataset authentically reflects learners’ cognitive states and operational habits when confronting cross-cultural narratives, constituting the foundational material for system evaluation. All interaction logs and behavioral data were anonymized prior to analysis and do not include any information that could be used to identify individual participants.

Ground Truth Annotation. To establish a reliable evaluation benchmark (Ground Truth), we assembled an “Expert Annotation Committee” comprising 11 domain experts: one professor of German literature and ten postgraduate students in the same field. We provided the committee with the

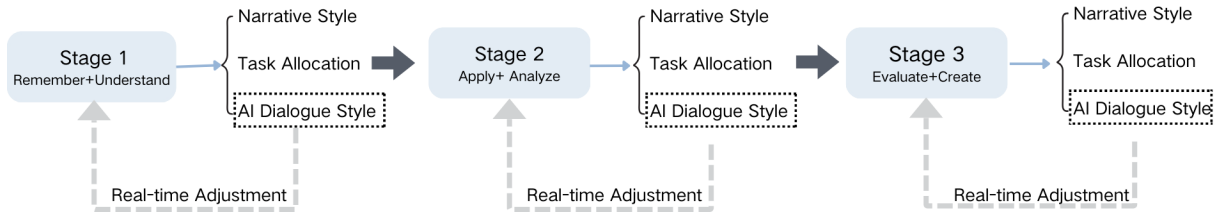


Figure 2: Cultural Adaptation for Segmented Operation. We divided the learning content of a writer into three stages according to Bloom’s taxonomy of educational objectives: Remember and Understand, Apply and Analyze, Evaluate and Create. The analysis of Narrative Style, Task Allocation, and AI Dialogue Style in current stage will influence the next stage, while AI Dialogue Style also exerts effects on the current stage.

Table 2: Player types and staged tasks used in our game-based study. Tasks are organized from Task 1 to Task 5.

Player Type	Task 1	Task 2	Task 3	Task 4	Task 5
Achiever	Level–challenge task	Achievement–unlock task	Time–limited objective	Statistics–oriented task	Completionist collection task
Explorer	Text exploration map	Culture–unlock task	Multi–path text navigation	Style restoration task	Easter-egg deep guidance
Socializer	AI critique & collaboration	Role-play dialog simulation	AI co-writing task	Scenario-based dialog	Virtual small-group discussion
Killer	AI-versus quiz	Literature challenge match	Beat-the-expert AI	Strategy analysis challenge	Point-capture tournament

11 complete gameplay datasets (including behavior logs, dialogue, and expression sequences) and conducted the annotation under a double-blind protocol. The annotation task encompassed two dimensions: Bartle Player Type: Experts synthesized the player’s behavior to categorize them into a single dominant type (Killer, Achiever, Socializer, Explorer). Cultural Cognitive Risk: Experts identified cultural misunderstandings within the interaction and determined the necessity of system intervention. The final benchmark labels were generated via Majority Voting. In cases of significant divergence, the professor adjudicated the final decision, ensuring the authority and reliability of the ground truth.

4.1.2 Ablation Study

To investigate the impact of different modalities and the “Reflection Loop” mechanism on system performance, we established the following five experimental settings for comparative analysis:

MC-MAS (Full): The proposed full model, incorporating multimodal inputs (Expression, Behavior Logs, Text) and the Coordinator’s Reflection Loop mechanism with specialized agents (Behavior, Language, and Culture).

w/o Expression: An ablated variant that removes facial expression data, relying solely on text input and behavior logs.

w/o Action: An ablated variant that removes

behavior log data, relying solely on text input and facial expressions.

w/o Dialogue: An ablated variant that removes text input, relying solely on behavior logs and facial expressions.

w/o Workflow (No-Reflection): This baseline retains all multimodal inputs but removes the *Coordinator Agent* and its “Reflection–Reconstruction” mechanism. It directly employs a single Large Language Model (qwen-plus) to process the multimodal data in a single pass.

Single-agent Self-Reflection: A single-agent baseline where one Large Language Model processes all multimodal inputs and performs one to two rounds of self-reflection using a generic reflection prompt, without agent decomposition or inter-agent coordination.

MC-MAS (No-specialization): A multi-agent baseline that preserves the Coordinator and reflection workflow but removes expert specialization. All agents share identical roles and prompts, differing only in random initialization, to isolate the effect of agent specialization.

4.1.3 Evaluation Metrics

To quantitatively assess the performance on the two core tasks, we define the following metrics:

Bartle Consistency Accuracy (ACC_{Bartle}). Given that the MC-MAS outputs a probability distribution vector P_{agent} over player types,

whereas the human annotation constitutes a single dominant category L_{human} , we employ Top-1 Accuracy for evaluation. A prediction is deemed correct if the player type with the highest probability assigned by the agent matches the human label. The formula is defined as:

$$ACC_{\text{Bartle}} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(\text{argmax}(P_{\text{agent}}^{(i)}) = L_{\text{human}}^{(i)}) \quad (1)$$

where N denotes the total number of samples, and $\mathbb{I}(\cdot)$ represents the indicator function.

Cultural Attribution Precision (P_{Culture}). This metric evaluates the system’s precision in identifying specific “cultural misunderstanding risks.” Let S_{model} denote the set of risk points identified by the model, and S_{human} denote the set of valid risk points verified by human experts. Precision is defined as the proportion of model-identified risk points that are validated as “necessary and correct” by the experts:

$$P_{\text{Culture}} = \frac{|S_{\text{model}} \cap S_{\text{human}}|}{|S_{\text{model}}|} \quad (2)$$

4.2 Cognitive Load Evaluation

To evaluate the practical teaching experience of the MC-MAS system in authentic instructional settings, we conducted a controlled user study with 20 participants. The participants were evenly divided into two groups (10 per group). Group A interacted with the Chapter 2 game using the full MC-MAS system with complete multimodal inputs and coordination mechanisms, while Group B completed the same game content using a baseline configuration that removed the MC-MAS framework and relied solely on an end-to-end large language model.

After completing the predefined in-game interaction tasks, participants immediately undertook two subjective evaluation measures to quantify the system’s instructional effectiveness and interaction experience.

NASA Task Load Index (NASA-TLX) (Hart and Staveland, 1988). To examine whether the proposed system effectively reduces cognitive barriers in cross-cultural learning scenarios, participants were asked to complete the NASA-TLX questionnaire, which comprises six dimensions, each rated

on a 1–100 scale. We focused on three core dimensions: Mental Demand, Frustration, and Perceived Performance. These dimensions were selected to investigate whether the multi-agent collaborative mechanism can maintain an appropriate level of learning challenge while alleviating frustration induced by cultural differences. For the detailed questionnaire, please refer to the Appendix C.

5 Main Results

On the basis of the experimental setup described in the previous section, this section reports the evaluation results of MC-MAS, with a particular focus on the framework’s modeling accuracy and user-centered cognitive load.

5.1 Accuracy Evaluation

Table 3 reports the performance of MC-MAS under different ablation settings, evaluated using Bartle player type prediction accuracy ($\text{Acc}_{\text{Bartle}}$) and cultural profile prediction performance (P_{Culture}). Overall, the full MC-MAS model achieves the best results on both metrics, validating the effectiveness of the proposed multi-modal, multi-agent architecture with a coordinated reflection mechanism.

Impact of Multimodal Inputs. The ablation results show that removing any modality leads to a substantial performance degradation, highlighting the complementary nature of multimodal behavioral signals. Among the three modalities, action logs play a dominant role in Bartle type prediction: removing action data causes $\text{Acc}_{\text{Bartle}}$ to drop sharply from 0.912 to 0.455. This indicates that player type inference relies heavily on observable in-game decision patterns rather than solely on expressive or linguistic cues.

In contrast, cultural profile prediction is more sensitive to dialogue data. When the dialogue modality is removed, P_{Culture} decreases significantly from 0.885 to 0.518, suggesting that cultural traits are more strongly reflected in language use and discourse patterns.

Removing facial expression features results in a moderate but consistent decline across both metrics, indicating that affective cues provide auxiliary information that enhances model performance without being the primary determinant. These findings suggest that MC-MAS benefits from modality specialization, where different input channels contribute unevenly but synergistically to distinct modeling objectives.

Table 3: Performance comparison across different model settings.

Model Setting	Missing Capability	ACC_{Bartle}	Δ	P_{Culture}	Δ
MC-MAS (Full)	–	0.912	–	0.885	–
<i>Ablation: Modalities</i>					
w/o Expression	Affective cues	0.818	↓0.094	0.754	↓0.131
w/o Action	Behavioral signals	0.455	↓0.457	0.812	↓0.073
w/o Dialogue	Linguistic interaction	0.727	↓0.185	0.518	↓0.367
<i>Ablation: Architecture</i>					
w/o Workflow (No-Reflection)	Iterative reflection	0.758	↓0.154	0.647	↓0.238
Single-agent Self-Reflection	Role diversity	0.805	↓0.107	0.723	↓0.162
MC-MAS (No-specialization)	Expert priors	0.836	↓0.076	0.779	↓0.106

Impact of Architectural Design. Architectural ablations further demonstrate the necessity of multi-agent coordination and structured reflection. Disabling the reflection workflow (*w/o Workflow*) leads to notable drops in both Acc_Bartle and $P_Culture$, indicating that simple feed-forward aggregation of agent outputs is insufficient for robust inference.

While introducing single-agent self-reflection partially recovers performance, it remains inferior to the full MC-MAS setting, suggesting that self-reflection alone cannot replace inter-agent deliberation. Moreover, the no-specialization variant, in which agents share identical roles, performs better than reflection-free configurations but still falls short of the full model. This highlights that performance gains stem from functional specialization rather than merely increasing the number of agents. For detailed performance comparison, see Appendix Figure 4.

5.2 Cognitive Load Evaluation

As shown in Figure 3, Full MC-MAS consistently outperforms the end-to-end LLM baseline across key NASA-TLX dimensions (measured on a scale of 0 to 100). Participants reported substantially higher perceived performance under Full MC-MAS (71.8 vs. 49.6), indicating improved task effectiveness and confidence. Meanwhile, mental demand was significantly reduced (55.6 vs. 68.9), suggesting that explicit multi-agent coordination alleviates cognitive effort during cross-cultural interpretation. Frustration was also notably lower (36.9 vs. 58.9), reflecting smoother interaction and reduced emotional strain. Overall, these results demonstrate that MC-MAS improves user experience by jointly enhancing perceived performance while reducing cognitive and affective load.

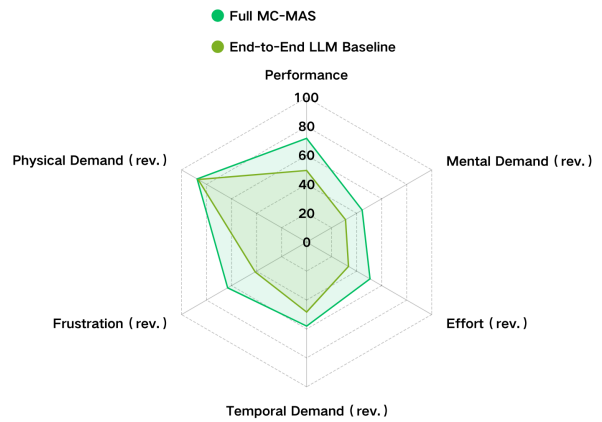


Figure 3: NASA-TLX comparison between Full MC-MAS and End-to-End LLM baseline. Larger areas indicate better performance after direction-aligned normalization.

6 Conclusion

This work examines adaptive narrative generation with LLMs in cross-cultural contexts, showing that multi-agent collaboration effectively enhances semantic and cultural alignment. By coordinating agents modeling linguistic, behavioral, and cultural signals, the framework achieves more accurate user profiling and cultural intention recognition than single-LLM baselines. A reflection-based coordination mechanism further resolves multimodal semantic conflicts and reduces hallucinations, especially in narratives with implicit cultural references. Experiments on a human-labeled Weimar-era literary benchmark validate the effectiveness of the approach. Future work will extend the framework to broader cultural domains and longitudinal user modeling.

606 Limitations

607 While the proposed multimodal and multi-agent
608 reflective framework significantly enhances large
609 language models' ability to infer players' cultural
610 backgrounds and gameplay profiles, several lim-
611 itations remain. One key constraint arises from
612 the deliberate abstraction of input data at the task
613 level. To mitigate hallucination risks associated
614 with long-context processing in large language
615 models, the current system prioritizes semantically
616 summarized interaction states over fine-grained be-
617 havioral traces. As a result, detailed low-level sig-
618 nals—such as continuous player movement trajec-
619 tories, cursor dynamics, or micro-level interaction
620 patterns—are not incorporated into the analytical
621 pipeline. Although this design choice improves
622 model stability and interpretability, it inevitably
623 limits the granularity of player modeling. Incorpor-
624 ating richer behavioral streams in future iterations
625 may enable more comprehensive representations of
626 player intent, engagement, and uncertainty, thereby
627 further strengthening adaptive interventions.

628 Another limitation concerns the scope of empiri-
629 cal validation. The evaluation system was designed
630 for Chinese learners engaging with German liter-
631 ary content, and the majority of participants shared
632 a relatively homogeneous cultural and linguistic
633 background. While the proposed multi-agent ar-
634 chitecture is not intrinsically tied to a specific dis-
635 cipline, language, or cultural setting, its effective-
636 ness in other cross-cultural learning contexts re-
637 mains empirically unverified. Consequently, the
638 current findings should be interpreted as evidence
639 of methodological feasibility rather than universal
640 generalizability. Future studies involving diverse
641 learner populations and cultural environments are
642 necessary to assess the robustness and transferabil-
643 ity of the framework across languages and educa-
644 tional domains.

645 Ethical Considerations

646 The use of multimodal data for adaptive player
647 modeling raises important ethical and privacy con-
648 siderations. In the present system, facial ex-
649 pression information is stored exclusively in a
650 labeled form (e.g., confused, engaged), and no
651 raw facial images are retained. Nevertheless,
652 the collection and analysis of multimodal sig-
653 nals—including language, speech, visual cues, and
654 behavioral patterns—inevitably involve potentially
655 sensitive personal information. Without rigorous

656 data anonymization, access control, and gover-
657 nance mechanisms, such systems may pose risks
658 to user privacy, particularly in educational contexts
659 involving minors or vulnerable populations.

660 A further ethical concern relates to cultural bias
661 embedded in large language models trained on
662 large-scale corpora. These models may inherit
663 and amplify cultural stereotypes or dominant cul-
664 tural narratives, which can undermine the goal of
665 fostering genuine cross-cultural understanding. In
666 adaptive educational systems, such biases risk over-
667 simplifying cultural representations or promoting
668 cultural homogenization under the guise of person-
669 alization. Balancing cultural authenticity with indi-
670 vidualized adaptation therefore remains a critical
671 design challenge.

672 Looking forward, advances in bias detection and
673 mitigation techniques, combined with participatory
674 design processes involving cross-cultural education
675 experts and target learner communities, are essen-
676 tial for developing more inclusive and culturally
677 responsive systems. Similarly, integrating privacy-
678 preserving data governance frameworks aligned
679 with educational ethics can help reduce risks re-
680 lated to data security and cultural inequity, while
681 maintaining the benefits of multimodal adaptation.

682 References

- 683 Richard A Bartle. 2004. *Designing virtual worlds*. New
684 Riders.
- 685 WAL Blyth, BS Bloom, and DR Krathwohl. 1966. Tax-
686 onomy of educational objectives. handbook i: Cog-
687 nitive domain. *Handbook 2: Affective domain. British journal of*
688 *educational studies*, 14(3). 689
- 690 Melania Cabezas-García and Arianne Reimerink. 2022.
691 Cultural context and multimodal knowledge repre-
692 sentation: Seeing the forest for the trees. *Frontiers in*
693 *Psychology*, 13:824932.
- 694 Elizabeth B Cloude, Daryn A Dever, Debbie L Hahs-
695 Vaughn, Andrew J Emerson, Roger Azevedo, and
696 James Lester. 2022. Affective dynamics and cogni-
697 tion during game-based learning. *IEEE Transactions*
698 *on Affective Computing*, 13(4):1705–1717.
- 699 Griffin Dietz, Nadin Tamer, Carina Ly, Jimmy K Le,
700 and James A Landay. 2023. Visual storycoder: A
701 multimodal programming environment for children's
702 creation of stories. In *Proceedings of the 2023 CHI*
703 *Conference on Human Factors in Computing Systems*,
704 pages 1–16.
- 705 Liyan Fan, Jinbo He, Yang Zheng, Yufeng Nie, Taolin
706 Chen, and Hongmei Zhang. 2023. Facial micro-

707	expression recognition impairment and its relationship with social anxiety in internet gaming disorder. <i>Current Psychology</i> , 42(24):21021–21030.	763
708		764
709		765
710	Stanley G Hart and Lowell E Staveland. 1988. Development of nasa-tlx (task load index): Results of empirical and theoretical research. In <i>Advances in Psychology</i> , volume 52, pages 139–183. North-Holland.	766
711		767
712		768
713		
714	Kanwar Sunreep Jossan, Andrea Gauthier, and Jodie Jenkinson. 2021. Cultural implications in the acceptability of game-based learning. <i>Computers & Education</i> , 174:104305.	769
715		770
716		771
717		772
718	Cheng Li, Mengzhuo Chen, Jindong Wang, Sunayana Sitaram, and Xing Xie. 2024a. Culturellm: Incorporating cultural differences into large language models. <i>Advances in Neural Information Processing Systems</i> , 37:84799–84838.	773
719		774
720		775
721		776
722		
723	Junkai Li, Yunghwei Lai, Weitao Li, Jingyi Ren, Meng Zhang, Xinhui Kang, Siyu Wang, Peng Li, Ya-Qin Zhang, Weizhi Ma, and 1 others. 2024b. Agent hospital: A simulacrum of hospital with evolvable medical agents. <i>arXiv preprint arXiv:2405.02957</i> .	777
724		778
725		779
726		780
727		781
728	Yishu Li and Yongjian Li. 2025. The paradox of cultural similarity: Teacher adaptation through the lens of linguistic, ethical, cultural, and historical friction. <i>International Journal of Intercultural Relations</i> , 109:102295.	782
729		783
730		784
731		785
732		
733	Michela Mortara, Chiara Eva Catalano, Francesco Bellotti, Giusy Fiucci, Minica Houry-Panchetti, and Panagiotis Petridis. 2014. Learning cultural heritage by serious games. <i>Journal of Cultural Heritage</i> , 15(3):318–325.	786
734		787
735		788
736		789
737		790
738	Minako O’Hagan. 2009. Towards a cross-cultural game design: an explorative study in understanding the player experience of a localised japanese video game. <i>The Journal of Specialised Translation</i> , (11):211–233.	791
739		792
740		793
741		
742		
743	Christopher P Ormell. 1974. Bloom’s taxonomy and the objectives of education. <i>Educational Research</i> , 17(1):3–18.	794
744		795
745		796
746	Valentina N Pescuma, Dina Serova, Julia Lukassek, Antje Sauer mann, Roland Schäfer, Aria Adli, Felix Bildhauer, Markus Egg, Kristina Hülk, Aine Ito, and 1 others. 2023. Situating language register across the ages, languages, modalities, and cultural aspects: Evidence from complementary methods. <i>Frontiers in psychology</i> , 13:964658.	797
747		798
748		
749		
750		
751		
752		
753	Lingzhi Shen, Xiaohao Cai, Yunfei Long, Imran Razzak, Guanming Chen, and Shoaib Jameel. 2025. Calm: Culturally self-aware language models. In <i>The Thirtieth Annual Conference on Neural Information Processing Systems</i> .	799
754		800
755		801
756		802
757		
758	Simone Titus and Dick Ng’ambi. 2023. Digital gaming for cross-cultural learning: Development of a social constructivist game-based learning model at a south african university. <i>International Journal of Game-Based Learning (IJGBL)</i> , 13(1):1–20.	803
759		804
760		805
761		806
762		807
	Angel Olider Rojas Vistorte, Angel Deroncele-Acosta, Juan Luis Martín Ayala, Angel Barrasa, Caridad López-Granero, and Mariacarla Martí-González. 2024. Integrating artificial intelligence to assess emotions in learning environments: a systematic literature review. <i>Frontiers in psychology</i> , 15:1387089.	808
		809
		810
		811
		812
		813
		814
		815
		816
	Abeer A Wafa, Mai M Eldefrawi, and Marwa S Farhan. 2025. Advancing multimodal emotion recognition in big data through prompt engineering and deep adaptive learning. <i>Journal of Big Data</i> , 12(1):210.	
	Hao Wei, Jianing Qiu, Haibao Yu, and Wu Yuan. 2024. Medco: Medical education copilots based on a multi-agent framework. In <i>European Conference on Computer Vision</i> , pages 119–135. Springer.	
	Minghao Wu, Jiahao Xu, Yulin Yuan, Gholamreza Haf-fari, Longyue Wan, Weihua Luo, and Kaifu Zhang. 2025. (perhaps) beyond human translation: Harnessing multi-agent collaboration for translating ultra-long literary texts. <i>Transactions of the Association for Computational Linguistics</i> , 13:901–922.	
	Fang Xie. 2025. From words to multimodalities: Compliment perceptions across lingua cultures. <i>Journal of Pragmatics</i> , 239:94–116.	
	Nan Xie, Zhengxu Li, Haipeng Lu, Wei Pang, Jiayin Song, and Beier Lu. 2025. Msc-trans: A multi-feature-fusion network with encoding structure for student engagement detecting. <i>IEEE Transactions on Learning Technologies</i> .	
	Lijuan Yan, Xiaotao Wu, and Yi Wang. 2025. Student engagement assessment using multimodal deep learning. <i>PLoS One</i> , 20(6):e0325377.	
	Chengcheng Yang, Zhiyao Liang, Tonglai Liu, Zeng Hu, and Dashun Yan. 2025. Mgm-net: Mamba-guided multimodal reconstruction and fusion network for sentiment analysis with incomplete modalities. <i>Electronics</i> , 14(15):3088.	
	Pardis Sadat Zahraei and Ehsaneddin Asgari. 2025. I am aligned, but with whom? mena values benchmark for evaluating cultural alignment and multilingual bias in llms. <i>arXiv preprint arXiv:2510.13154</i> .	
	He Zhang, Lu Yin, and Hanling Zhang. 2024. Using subjective emotion, facial expression, and gaze direction to evaluate user affective experience and predict preference when playing single-player games. <i>Ergonomics</i> , 67(12):1863–1883.	
	Zheyuan Zhang, Daniel Zhang-Li, Jifan Yu, Linlu Gong, Jinchang Zhou, Zhanxin Hao, Jianxiao Jiang, Jie Cao, Huiqin Liu, Zhiyuan Liu, and 1 others. 2025. Simulating classroom education with llm-empowered agents. In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 10364–10379.	

817	A Experimental Setup and	
818	Hyperparameters	
819	We provide the detailed implementation settings for	
820	our multi-agent experiments in Table 4 to ensure	
821	reproducibility.	
822	B prompt	
823	This section illustrate the prompt design for the Be-	
824	havior Analysis Agent, Language Analysis Agent,	
825	Cultural Context Analysis Agent, and Coordinator	
826	Agent.	
827	B.1 Behavior Analysis Agent (A_{beh})	
828	Role. The agent acts as a senior game data analyst	
829	and behavioral psychology expert, specializing in	
830	the Bartle taxonomy of player types.	
831	Input. The agent receives raw player behavior	
832	logs $\mathcal{B} = \{r_c, t_m, e_{\text{ratio}}, i_{\text{freq}}\}$ along with the con-	
833	flict opinion provided by the Committee Leader	
834	(S_{Conflict}).	
835	Task. Based on the input data, the	
836	agent infers the player’s probability distribution	
837	over the four Bartle dimensions: $\mathbf{P}_{\text{bartle}} =$	
838	[Killer, Achiever, Socializer, Explorer].	
839	Theoretical Framework. The inference is	
840	grounded in the Bartle taxonomy, where Achiev-	
841	ers prioritize in-game accomplishments, Explorers	
842	focus on environmental interaction, Socializers em-	
843	phasize player-to-player interaction, and Killers	
844	seek competition and dominance.	
845	Reasoning Anchors.	
846	• Achiever: High repetition with increasing suc-	
847	cess rates, direct trajectories toward task objec-	
848	tives, and resource hoarding behavior.	
849	• Explorer: Redundant movement trajectories, fre-	
850	quent low-yield interactions, and unconventional	
851	operations.	
852	• Socializer: Extended dwell time in public areas,	
853	active chat or emoji usage, and a preference for	
854	team-based play.	
855	• Killer: High PvP participation, preference for	
856	burst or control-oriented equipment, and disrupt-	
857	ive behavior in non-competitive zones.	
858	Output Format. The agent outputs (1) a normal-	
859	ized probability distribution vector $\mathbf{P}_{\text{bartle}}$ (sum =	
860	100%), and (2) a step-by-step reasoning explana-	
861	tion grounded in observable behavioral evidence	
862	such as timestamps, movement trajectories, and	
863	interaction heatmaps.	
	B.2 Language Analysis Agent (A_{lang})	864
	Role. The agent functions as an expert in computa-	865
	tional linguistics and affective computing.	866
	Input. The agent simultaneously analyzes the	867
	player’s dialogue text (\mathcal{T}), facial expression data	868
	(\mathcal{E}), and the conflict opinion provided by the Com-	869
	mittee Leader (S_{Conflict}).	870
	Task. The agent performs cross-modal analysis	871
	to infer the player’s implicit cognitive and emo-	872
	tional states by aligning linguistic features with	873
	facial expression cues.	874
	Analysis Focus.	875
	• Metaphorical Depth: Examines the use of	876
	metaphors or domain-specific game terminology,	877
	reflecting the player’s level of immersion.	878
	• Syntactic Complexity: Assesses whether sen-	879
	tence structures are complex or fragmented, serv-	880
	ing as indicators of cognitive state and cognitive	881
	load.	882
	• Emotional Consistency: Evaluates the align-	883
	ment between verbal content and facial expres-	884
	sions, including potential discrepancies such as	885
	positive language accompanied by negative facial	886
	cues (e.g., forced smiling).	887
	Output Format. The agent generates the fol-	888
	lowing outputs, with approximately two sentences	889
	for each item:	890
	1. Player’s Linguistic Profile: A description of	891
	the player’s contextual depth (high or low) and	892
	salient syntactic characteristics.	893
	2. Player’s Emotional State: A synthesized in-	894
	terpretation of the player’s emotional condition	895
	derived from cross-modal alignment.	896
	B.3 Cultural Context Analysis Agent (A_{cult})	897
	Role. The agent acts as an expert in cross-cultural	898
	communication and a historian specializing in the	899
	Weimar Republic period.	900
	Input. The agent analyzes the player’s static	901
	profile I_{profile} (including nationality and native lan-	902
	guage), behavioral logs $\mathcal{B} = \{r_c, t_m, e_{\text{ratio}}, i_{\text{freq}}\}$,	903
	and the conflict opinion provided by the Committee	904
	Leader (S_{Conflict}).	905
	Task. The agent evaluates the cultural distance	906
	between the player’s background (source culture)	907
	and the current in-game context (target culture),	908
	with a particular focus on potential cultural misun-	909
	derstanding risks.	910
	Key Checks.	911

Table 4: Hyperparameter settings and implementation details for the multi-agent experiments.

Parameter	Value	Description
Model Architecture		
Base Model	Qwen-Plus	Accessed via Alibaba Cloud DashScope API.
Generation Configuration		
Temperature	0.7	Set to encourage diversity while maintaining logic.
Top-p	Default	Using the default setting of the API.
Agent Settings		
Max Dialogue Turns	5	Maximum interaction rounds per session.
System Prompt	See App. B	Role-specific instructions for each agent.

912 • **Cultural Symbol Identification:** Detects
 913 culture-specific references embedded in the cur-
 914 rent scenario (e.g., biblical allusions or political
 915 musical performances).

916 • **Misunderstanding Prediction:** Predicts poten-
 917 tial misinterpretations based on the player’s na-
 918 tive language and observed behavioral patterns.

919 **Output Format.** The agent outputs the cultural
 920 assessment result $\mathcal{P}_{\text{Culture}}$, consisting of:

- 921 1. **Risk Alert Level** ($O_{\text{RiskAlert}}$): High, Medium,
 922 or Low.
- 923 2. **Bridging Suggestions** ($O_{\text{BridgingSuggestion}}$):
 924 Concrete recommendations aimed at reducing
 925 cultural distance and mitigating misunderstand-
 926 ing.

927 **B.4 Coordinator Agent** (A_{coord})

928 **Role.** The agent acts as the chief architect of adap-
 929 tive learning, coordinating and supervising a panel
 930 of expert agents, including the Behavior, Language,
 931 and Cultural Context Agents.

932 **Input.** The agent receives analytical reports
 933 from all subordinate agents, including detected
 934 conflicts (S_{Conflict}) and the current iteration count
 935 (N_{iter}).

936 **Task.** The agent synthesizes cross-agent reports
 937 to determine whether a re-evaluation through a re-
 938 flection loop is required. If necessary, it specifies
 939 the conflict context, triggers re-analysis, and tracks
 940 iterative updates until consensus is reached or the
 941 maximum iteration limit (5) is exceeded. Upon
 942 consensus, the agent generates the final adaptation
 943 strategy.

944 **Adjudication Protocols.**

945 • **Conflict Detection:** Identifies semantic or func-
 946 tional inconsistencies across agents (e.g., effi-
 947 cient behavioral patterns conflicting with high
 948 cultural misunderstanding risk).

949 • **Reflection Trigger:** Instructs subordinate agents
 950 to re-evaluate their analyses based on the identi-
 951 fied conflict context.

952 • **Decision Priority:** Prioritizes pedagogical effi-
 953 cacy over gameplay progression.

954 **Output Format.**

955 • **If conflict exists:**

- 956 – **Re-evaluation Needed:** Yes.
- 957 – **Conflict Description:** A detailed explana-
 958 tion of the detected cross-modal or cross-
 959 agent inconsistency.

960 • **If consensus is reached:**

- 961 – A final JSON configuration specifying:
 962 * **Task_Allocation** (adjustment strat-
 963 egy),
 964 * **Narrative_Style** (adjustment strat-
 965 egy),
 966 * **AI_Assistant_Tone** (tone setting).

967 **C NASA-TLX Questionnaire**

968 This appendix reports the post-experiment ques-
 969 tionnaire used to assess participants’ perceived task
 970 load after gameplay. The questionnaire is based
 971 on the NASA Task Load Index (NASA-TLX) and
 972 was adapted to the interactive narrative and cross-
 973 cultural learning context of this study.

974	C.1 Instructions		
975	Participants were asked to rate each item on a continuous scale from 0 (very low) to 100 (very high),	the most significant type based on the highest frequency of occurrence or behavior at key decision moments. The specific definitions are as follows:	1019
976	according to their subjective experience during the game session.		1020
977			1021
978			
979	C.2 NASA-TLX Dimensions		
980	Mental Demand How much mental activity was required to complete the game tasks (e.g., understanding narrative content, interpreting cultural metaphors, selecting dialogue options, and making decisions)?	Achievers Focus on obtaining points, levels, equipment, or other concrete measures of success. Their core motivation is derived from receiving immediate feedback and visible improvement within the game.	1022
981			1023
982			1024
983			1025
984			1026
985	Physical Demand How much physical effort was required to interact with the system (e.g., mouse clicking and keyboard input)?	Explorers Focus on discovering deep system mechanisms, hidden content, or unknown map areas. They enjoy the process of "figuring out how the game works" more than merely clearing the game.	1027
986			1028
987			1029
988	Temporal Demand How much time pressure did you experience during interactions with NPCs or while completing tasks in the game?	Socializers Focus on interactions with other players or NPCs. The game itself serves merely as a social backdrop; interpersonal communication and emotional connection are their core motivations.	1030
989			1031
990			1032
991	Performance How successful do you think you were in accomplishing the learning objectives and narrative tasks in the game?	Killers Focus on imposing influence on others (players or the game environment), typically manifested through competition, destruction, exerting control, or gaining a sense of superiority by dominating the environment.	1033
992			1034
993			1035
994	Effort How much mental or physical effort was required to achieve your level of performance?		1036
995			
996	Frustration How frustrated, stressed, or discouraged did you feel during the game (e.g., due to cultural unfamiliarity or unclear system feedback)?		
997			
998			
999	D Instructions for Expert Evaluators		
1000	D.1 Background and Objective	Task 2: Identification of Cultural Misunderstandings	1042
1001	This project aims to validate the accuracy of an automated system in analyzing player game behavior and cultural adaptability. As an expert evaluator, you will be provided with game interaction data collected from 11 anonymous players (including dialogue logs, key decision paths, and behavioral timestamps). Your task is to independently classify each player's characteristics and identify cultural misunderstandings based on this raw data. Your judgments will serve as the Ground Truth for the experiment, utilized for comparative validation against the output of the automated workflow.	Please combine the interactive context of the player and the game content to identify moments of cognitive dislocation caused by cultural background differences. The criterion for judgment is: mark instances where the player's behavior or textual feedback clearly indicates a failure to understand the German literary context (e.g., specific social norms of the Weimar Republic era, the intent of Brecht's "Verfremdungseffekt"), or attempts to misinterpret the game situation by applying their own cultural background. Please indicate the specific trigger point (Context) and the nature of the potential misunderstanding in your evaluation report.	1043
1002			1044
1003			1045
1004			1046
1005			1047
1006			1048
1007			1049
1008			1050
1009			1051
1010			1052
1011			1053
1012			1054
1013	D.2 Evaluation Tasks	D.3 Risk Control and Evaluation Principles	1055
1014	Task 1: Bartle Taxonomy Classification	To ensure the objectivity and ethical compliance of the evaluation results, the process must strictly adhere to the following guidelines:	1056
1015	Please carefully review the player's game behavior patterns to determine which category of the Bartle Taxonomy their dominant motivation best fits. When making a determination, please select	First, regarding Subjectivity Bias , given that players may exhibit multiple personality traits simultaneously, please ensure classification is based on Dominant evidence , avoiding projections based on personal gaming preferences.	1057
1016			1058
1017			1059
1018			1060
			1061
			1062
			1063
			1064
			1065

1066 Second, regarding **Over-interpretation**, please
1067 maintain a cautious attitude when identifying cul-
1068 tural misunderstandings. Only mark an instance
1069 as a cultural misunderstanding when supported by
1070 clear evidence (e.g., specific expressions of confu-
1071 sion, repeated ineffective attempts), avoiding the
1072 misclassification of common operational errors as
1073 cultural barriers.

1074 Finally, all evaluations must be conducted under
1075 strict **Independence & Confidentiality**. Please
1076 complete the annotation independently without ref-
1077 erencing the automated system’s results or dis-
1078 cussing with other experts to ensure the validity
1079 of the blind test. All data is for this research only
1080 and is strictly prohibited from external disclosure.

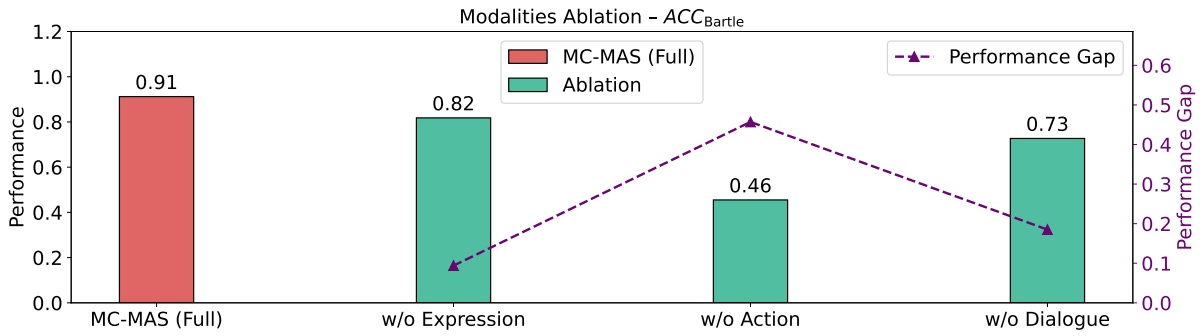
1081 **E Recruitment and Compensation Details**

1082 In this study, we recruited a total of 20 game par-
1083 ticipants to engage in the gameplay experiments.
1084 Additionally, we invited 11 annotation experts to
1085 perform data quality assessment and annotation.
1086 All participants were recruited on a voluntary ba-
1087 sis.

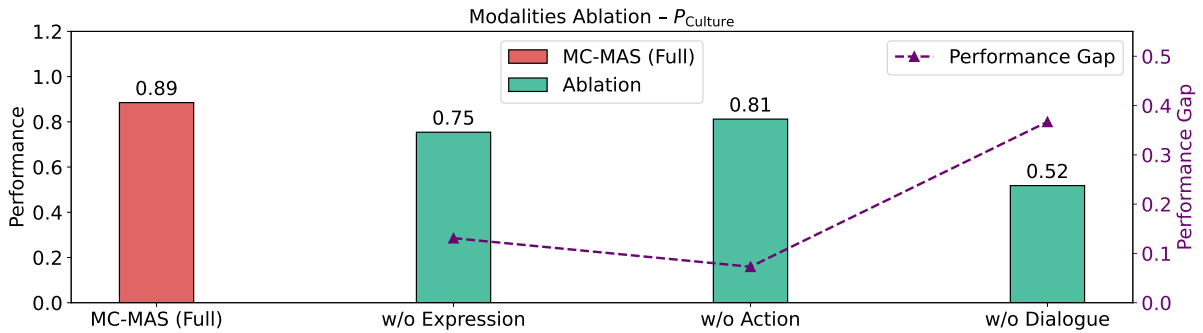
1088 Regarding compensation, each participant (in-
1089 cluding both game players and experts) received a
1090 payment of 20 CNY or a gift of equivalent value
1091 upon completion of the task. This compensation
1092 plan was designed in accordance with the estimated
1093 workload and local ethical standards, and was re-
1094 viewed and approved by the Institutional Review
1095 Board (IRB) of the authors’ institution.

1096 **F Figure**

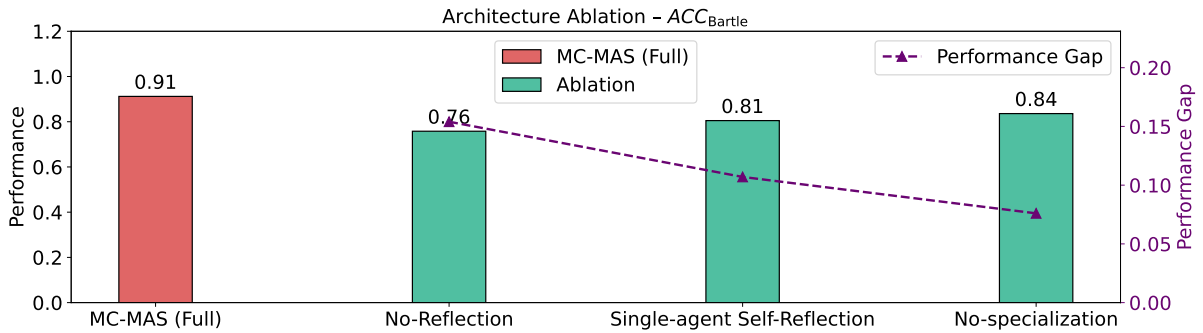
1097 Figure 4 details the ablation study results of the
1098 MC-MAS framework on Bartle player type clas-
1099 sification (ACC_{Bartle}) and cultural attribute predic-
1100 tion (P_{Culture}) tasks. The experiment aims to verify
1101 the effectiveness of multimodal inputs and core
1102 architectural components. The bar charts repre-
1103 sent performance scores, while the purple dashed
1104 lines illustrate the performance gap between the
1105 full model and the ablated variants.



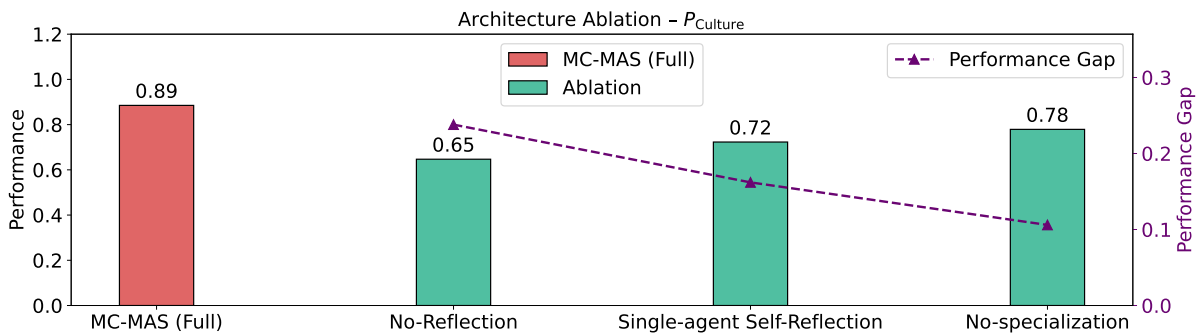
(a) Effect of removing expression, action, or dialogue modalities on Bartle type accuracy



(b) Effect of removing modalities on cultural profile prediction



(c) Effect of architectural modifications on Bartle type accuracy



(d) Effect of architectural modifications on cultural profile prediction

Figure 4: MC-MAS ablation study: Bartle type accuracy and cultural profile prediction with performance gap indicated.