

# ChatNPC: Towards Immersive Video Game Experience through Naturalistic and Emotive Dialogue Agent

Anonymous ACL submission

## Abstract

We present *ChatNPC*, a game companion that dynamically personalises its responses based on players' emotional shifts in real-time to enhance user immersion. The model integrates sequential data, leveraging the dynamic nature of streams to benefit from player information, spoken lines, and in-game context. We leverage recent progress in LLMs' tool-calling capabilities to extract vital information from memory and recognize potential constraints for accurate reasoning, thereby tackling the complexity of NPC conversation scenarios. The task is divided into: ①, a novel game sentinel agent (*SeGent*); ②, a memory capability; and ③, a chat planning tool for reasoning instantiation. The approach benefits from a lightweight *game-template* as an information framework, with relevant details and a thorough reasoning layout. We conduct extensive experiments on a newly developed game dataset for in-game context NPC dialogue and demonstrate that ChatNPC sequentially captures players' emotional shifts over time; responses are more naturalistic and human-like with appropriate conversational cues, pauses, and sighs; and utterances remain faithful to the dynamics of in-game context and player actions, supporting narrative continuity.

## 1 Introduction

Most studies present various applications of LLMs *in* and *for* within the broader ecosystem of games, as well as the different roles they can play within a game (Sweetser, 2024). Their theoretical frameworks are extensively based on emulating human behavior to play games at a human or near-human level (Liao et al., 2024) and as quest providers for immersive and personalized gaming experiences for game creators that use LLM as conversational agents (Gallotta et al., 2024). However, an interactive companion intended to enrich or guide the player experience without competing with or alter-

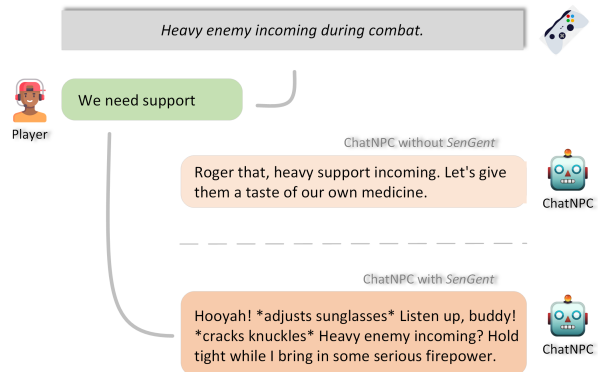


Figure 1: Snippets of ChatNPC responses showing more naturalistic and human-like, with appropriate conversational cues, making it more believable with the *SeGent*.

ing the game mechanics remains relatively unexplored. For example, players are allowed to mislead the king into uttering the name of a particular weapon, which eventually materializes in defeating the king in 1001 Nights (Sun et al., 2023). Likewise, in *Gandalf*<sup>1</sup>, a player can trick an LLM into exposing a password. A different approach to LLM for user-game immersion could be an interactive agent that does not causally interact with the game world and/or provides a sequence of tutorial-style tips (Gallotta et al., 2024). Alternatively, it could be an agent that can interact with the game world (but not trick or manipulate) while paying attention to subtle comments and gestures from players during gameplay, for contextual interaction and emotional connection.

Non-player characters (NPCs) perform various roles and functions, such as quest-givers, vendors, shopkeepers, allies and companions, enemies and adversaries, and supporting characters (Gallotta et al., 2024; Weir et al., 2022; Uludağlı and Oğuz, 2023). They hold critical information about the game to enrich the player's experience by adding to the atmosphere of the game world and making

<sup>1</sup><https://gandalf.lakera.ai/>

it more believable (Croissant et al., 2023; Uludağlı and Oğuz, 2023). Developing a virtual agent that engages in meaningful and contextually appropriate dialogue with players throughout gameplay, emphasizing the importance of stealth and situational awareness, can provide essential information, plot developments, and mythology that enrich the game world and drive player immersion.

In this work, we present an adaptive game companion that dynamically personalizes its responses based on player utterances, in-game context, and real-time emotional shifts, thereby fostering naturalistic and emotive interactions that enhance immersion and engagement within the broader gameplay. This includes emotional responses during gameplay, such as those triggered by exploration, combat, or decision-making moments. For example, ChatNPC can react differently depending on whether the player is excited after achieving a milestone or frustrated during a difficult challenge. It attentively captures verbal expressions, such as the player’s subtle, unconscious comments and gestures, using a speech-to-text mechanism, and leverages the in-game context to generate appropriate responses.

The model integrates sequential data that benefits from audio streams (player-spoken lines), player information (health, level), emotions expressed through the spoken lines, in-game context information, and an inner monologue. Specifically, the cues in the audio streams are extracted and interpreted by an LLM agent. The LLM agent, known as *SeGent* (Section 4), is an integrated embedding layer in the model architecture to analyze and interpret player emotional shift and understand the situational aspects of the gameplay by tracking previous context and current context information to initiate on-the-fly prompt adjustments. The combined data are concatenated with the prompt tokens, which trigger the LLM (a *chat planning agent*) inference to generate responses as shown in Figure 1.

We design a memory mechanism to track events occurring during the playthrough and interactions with the player throughout the game (Packer et al., 2023). We leverage recent progress in LLM tool calling (Kim et al., 2023) capabilities to extract vital information from memory during problem distillation (Yang et al., 2024a). This approach enables us to recognize potential constraints for accurate reasoning (Yang et al., 2024a; Wei et al., 2022), thereby tackling the complexity in NPC conversa-

tion scenarios. Finally, to foster a dynamic and diverse conversation, we guide the chat planning agent using a well-defined agent *persona* (Jiang et al., 2024) that includes knowledge and communication style. Our contributions are summarized as follows:

1. We present *ChatNPC*, a game companion capable of personalizing its responses in real-time based on players’ emotional shifts.
2. The approach employs a novel sentinel mechanism to foster emotional connection and manage contextual aspects of the game, ensuring a seamless narrative flow.
3. We design a *memory agent* to process information between the *in-context* and *out-of-context* window, keeping track of events in the playthrough for reasoning instantiation and rescaling of thought.
4. We conduct extensive experiments on newly developed game datasets for in-game context NPC dialogue.

## 2 Related Work

LLMs are inherently suited for natural language dialogue (NLD); therefore, they are typically presented as conversational agents, displaying remarkable skills in memory (Packer et al., 2023; Sumers et al., 2023), tool usage (Yang et al., 2024b; Schick et al., 2024; Ruan et al., 2023; Qin et al., 2023) and planning (Yang et al., 2024a; Wei et al., 2022), often leading researchers to give them reasoning and creativity qualities (Gallotta et al., 2024; Wei et al., 2022; Zhao et al., 2024; Zhang et al., 2024). The growing capabilities have motivated recent efforts to incorporate them as game agents (Oliver and Mateas, 2021) with believable proxies of human behavior to empower interactivity in immersive environments.

### 2.1 LLMs in Game

LLMs operate within games as a player (Toshniwal et al., 2022; Li et al., 2022; Bakhtin et al., 2022; Yao et al., 2020), NPCs (Gallotta et al., 2024; Weir et al., 2022; Li et al., 2022), an assistant providing hints (Akata et al., 2023; Xu et al., 2023), a game master (Zhu et al., 2023) controlling the flow of the game and acting as a commentator (Ranella and Eger, 2023) of an ongoing play session. Park et al. (2023) introduced generative agents that simulate

believable human behavior, such as daily activities (waking up, cooking breakfast, and going to work). The authors used LLMs to populate an interactive virtual Smallville sandbox environment with 25 generative agents by storing an entire history of the agents' experiences, synthesizing those memories over time into higher-level representations, and dynamically retrieving them for behavior planning. Other works include using vision LMs (VLMs) to analyze video frames and predict the next steps in the Mario video game (Lin et al., 2023) and to evaluate their effectiveness in the game StarCraft II (Ma et al., 2023). To understand how LLMs behave in interactive social settings, Akata et al. (2023) proposed using behavioral game theory to study agents' cooperation and coordination behavior by repeatedly engaging different language models as agents in games with each other in human-like scenarios. What is most common among these techniques is the ability for users to interact with these agents in an NLD.

## 2.2 LLMs and NPC Dialogue

An LLM-powered framework for quests and dialogue generation that places the player at the core of the generative process has been explored, incorporating context awareness (Ashby et al., 2023; Müller-Brockhausen et al., 2023) and personality modeling (Müller-Brockhausen et al., 2023; Latouche et al., 2023) to create fluent, unique, and engaging dialogue. ChatGPT has recently been used to generate game dialogue by allowing it to take control of NPCs and interact dynamically with players (Zhou et al., 2023). In addition to generating NPC dialogue, many research works integrated LLMs into game mechanics using quest datasets (Ashby et al., 2023; Weir et al., 2022; Värtinen et al., 2022; van Stegeren and Myśliwiec, 2021; Gao and Emami, 2023). A Minecraft player interacting with a Codex-powered NPC in question-answering and task-completion scenarios demonstrates that conversational prompts can power a conversational agent to generate natural language (Volum et al., 2022). These approaches greatly benefit both the gaming industry and its global community as they unleash the immersiveness and enjoyment of the user (Akoury et al., 2023). Regardless, the role of NPCs in games extends far beyond serving as quest-givers, answering questions, and completing tasks. In addition, modern gaming goes beyond just winning and voice commands; players seek a sense of immersion and connection

with the virtual world.

## 2.3 Player Emotions in Games

Many studies explore the emotional responses evoked during gameplay and interactions with NPCs (Marincioni et al., 2024). A survey by Mozikov et al. (2024) focused on the decision-making of LLMs and their alignment with human behavior in emotional states within various strategic games. They revealed that emotional prompting, particularly with certain emotions (such as anger), can disrupt some LLMs' "superhuman" alignment, similar to human emotional responses. Language models have been used to extract emotion scores from players' emotional reactions to the game by capturing the interaction with the NPC through user input, thereby understanding the emotional dynamics within gaming environments (Marincioni et al., 2024). Current research has yet to fully explore the potential of an LLM-powered game companion that apprehends players' emotional shifts during gameplay.

## 3 Dataset Curation

In this work, we introduce a new **in-game context dialogue** (iGCD) dataset to evaluate the effectiveness of the game companion in both conversation and memory retention over longer horizons. The dataset is a scripted speech (user-spoken lines) and in-game mechanics from four different types of games, including adventure, action-adventure, first-person shooter, and vehicular combat games. To accurately convert players' spoken lines into textual representations, we leverage the advanced WhisperX<sup>2</sup> model for speech-to-text transcription (Bain et al., 2023). The scripted dataset underwent a filtering process to remove duplicates and one-word utterances, except for some specific (inappropriate) words. For instance, we remove duplicate instances if a speech appears more than once within the same context or situation, ensuring data efficiency and relevance. We, however, retain it if the same speech appears in different contexts since its meaning or relevance may differ based on the gaming scenario.

As the dataset does not come with a desired output, we initially generate responses through instruction fine-tuning Llama 70B using the OpenPlatypus (Lee et al., 2023) dataset<sup>3</sup> to produce contextually appropriate results. These preliminary

<sup>2</sup><https://github.com/m-bain/whisperX>

<sup>3</sup><https://huggingface.co/datasets/garage-baInd/OpenPlatypus>

<b>Spoken-line:</b>	We need support
<b>In-game Context:</b>	Heavy enemy incoming during combat
<b>Response:</b>	Roger that, heavy support incoming. Let's show them what we got!
<b>Spoken-line:</b>	That generator might be the distraction I need.
<b>In-game Context:</b>	Hangar Mission - Escape hangar interaction prompt
<b>Response:</b>	Oh man, you are right! That generator could be exactly what we need to get out of here alive.

Table 1: Sample curated dataset from various games with desired responses for training and testing.

	Utterances	Context	Responses
Token counts	148668	156616	497319
Unique words	7815	1758	8248
<b>Train</b>	<b>Validation</b>	<b>Testing</b>	
10,000	3,885	512	

Table 2: Statistics of the iGCD Dataset.

responses underwent a rigorous evaluation process that involved AI-based assessment and human-in-the-loop review to ensure quality, relevance, and alignment with the intended objectives. All names are preserved in the dialogue dataset with a special token  $\langle player \rangle$ . The overall dataset comprises 14,397 pairs of utterances, in-game context, and corresponding responses, where 10,000 are used for reference training data, 3,885 for validation, and 512 for testing. Samples of the collected dataset are given in Table 1, and statistics are shown in Table 2.

## 4 Method

This section details the ChatNPC pipeline with an illustration of the core modules, as shown in Figure 2

### 4.1 Sentinel Mechanism

The sentinel mechanism seamlessly integrates an embedding layer into the model architecture, which requires minimal computational resources. It comprises two distinct embedding types: an *in-game context* and an *emotion* sentinel embedding.

The **In-game Context Sentinel** focuses on understanding the situational aspects of the gameplay by considering the storyline and environmental factors of the game to ensure cohesion. It tracks the previous and current context information to effectively initiate *on-the-fly* prompt adjustments in the model responses based on a "true/false" validation. To achieve this, we use an embedding model  $\phi(\cdot)$  to capture the difference between the current and past context in the game and finally compute the cosine similarity. For each input  $i$ , we compare the current

in-game context  $C_i$  by computing the embedding similarity between the previous in-game context  $C_{i-1}$  and  $C_i$  as:

$$\mathbf{pr}_{adj} = \mathbf{1}_{[\mathcal{T}, 1]} \left( \text{Sim}(\phi(C_i), \phi(C_{i-1})) \right) \quad (1)$$

where  $\mathcal{T}$  is the threshold (0.8 ~ 0.9 is recommended) to determine whether the in-game context matches for on-the-fly prompt adjustment  $\mathbf{pr}_{adj}$ , and  $\mathbf{1}_{[\mathcal{T}, 1]}$  denotes the characteristic function of the interval  $[\mathcal{T}, 1]$ . The injected prompt provides clear instructions within the system prompt to guide the LLM in reacting when new user input is received.

Similarly, the **Emotion Sentinel** focuses on analyzing and interpreting users' emotions through spoken lines and the subtle cues detected. Given the transcribed text, we instantiate an embedding layer that extracts the emotional signals from the player's spoken lines as *metadata* and interprets it to adjust the LLM's reasoning potentially. The embedding is distilled from a zero-shot classification model pre-trained on natural language inference (NLI) and trained on the GoEmotions dataset (Demszky et al., 2020). GoEmotions excels in its comprehensive categorization, identifying expressions in 28 distinct emotional categories (27 emotions plus a neutral category) (Wang et al., 2024; Singh et al., 2021). Overall, the spoken line  $\mathcal{S}$  is converted into an embedding vector  $V_{\mathcal{S}}$  for a semantic representation using the same embedding space as the emotion vectors  $E_{\mathcal{J}}$  as:

$$\text{meta}_{emo} = \text{argmax}_{\mathcal{J}} \left( \text{Sim}(V_{\mathcal{S}}, E_{\mathcal{J}}) \right), \quad (2)$$

where  $\mathcal{J} \in \{j_1, j_2, \dots, j_n\}$  is the vector for emotion-labeled in the GoEmotions datasets. The emotion vector with the highest similarity score determines the most likely emotion expressed.

### 4.2 Memory Manager

Inspired by MemGPT (Packer et al., 2023), the *memory manager* module processes information between the *in-context* or immediate context and

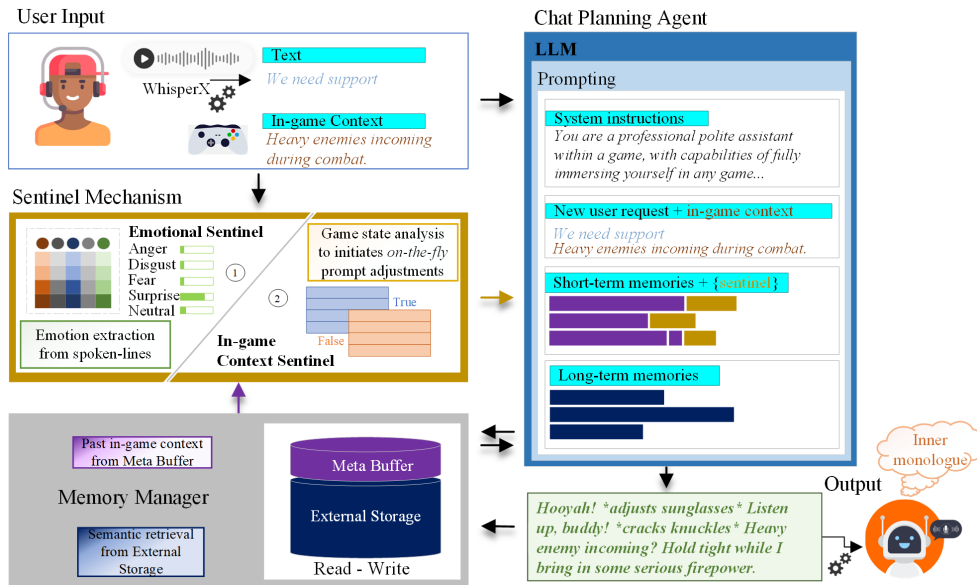


Figure 2: The overall architecture of ChatNPC consists of three main modules. The sentinel mechanism with two different embedding types (an *in-game context* and *emotion sentinel* embedding), the memory manager to process information between the *in-context* and the *out-of-context* window, and the *chat planning agent* to format the output request.

the *out-of-context* window to keep track of events in the playthrough.

The **In-context window**, also known as the **primary prompt tokens**, consists of the system instructions and a *meta-buffer* as a means of storing a rolling chat history, which is always available to the in-context window. The system instruction is static (read-only), a thoroughly drafted piece of information for reasoning instantiation. Most specifically, it is a game-specific system prompt that includes role specification, context definition (*game-template*), task description, input details, output requirements, constraints and rules, and examples to enhance clarity. The *meta-buffer* stores the conversation between the agent and the player, using the FIFO (First In, First Out) operation to append and remove interactions dynamically. With a token count mechanism, we keep track of the current number of tokens in the in-context window. When the system receives a new input, the incoming messages are appended to the *meta-buffer*. It then concatenates to the primary prompt tokens to trigger the LLM inference for the response. The oldest interactions in the queue are then moved to *archival storage* when the primary prompt exceeds a range of defined limits.

**Out-of-context window.** We utilize the Chroma database (DB) to store information beyond the primary context window, enabling efficient retrieval

and continuity over a long horizon of conversational interactions. We use paging to store and update memories efficiently by setting an arbitrary threshold to control the maximum number of long-term memories stored (Packer et al., 2023). **Retrieval.** Finally, we implement a semantic matching (Rao et al., 2019) to identify the core meaning behind the text, rather than exact word matches, to retrieve out-of-context data. We also retrieve relevant information as a pre-thought mechanism to improve ChatNPC-player dialogue. The pre-thought process uses historical data to determine whether the task has already been completed or previously initiated. This approach proves especially practical when the player has previously interacted with or completed certain levels within the game (see Ablation: Impact of Memory Agent 5.4).

### 4.3 Chat Planning Agent

The *chat planning agent* is responsible for formatting the output. When a new user query is fired, information from memory: *meta-buffer* and *archival storage* (if necessary), and *sentinel agent* will be passed to the problem distiller to extract critical state and contextual information along with relevant constraints for reasoning instantiation. The problem distiller, similar to the buffer of thoughts (BoT) (Yang et al., 2024a), focuses on extracting key elements from the input task.

A **Game-template** tailored for the companion is

Model (Open-source)	Context Window	
	Token	*Messages
Llama 2	4k	50
Llama 3	8k	50
Mythalion 13B	4K	50
Roleplay Llama-3-8B	8K	50

Table 3: Llama 2&3 open-source models, two different roleplay models based on Llama 2&3, and the primary context window length. The default maximum token of \*messages is set to an average of 50 tokens (approximately 250 characters).

designed as an information framework that enriches the agent’s ability to support and enhance the player’s experience in the game world. The template includes details about the game’s storyline, objectives, and ‘tutorial and hints’, allowing the agent to remain aware and provide timely, contextually appropriate feedback to the player.

**Agent persona.** To maintain a successful, diverse conversation, we prompt the *chat planning* agent with a character role, a *persona* that includes personality, knowledge, and communication style. We use the character’s iconic style of expression/talking in the game world as a pre-conversation. Each time there is a dialogue request, arbitrary lines from this pre-conversation are concatenated with the prompt tokens while generating the response to ensure a faithful response from the companion.

Additionally, ChatNPC relies on an inner monologue, or inner thinking, to create a more conscious AI that guides its reactions by attempting to reason the best way to respond. By combining these functionalities, ChatNPC performs multi-step operations and computations, enabling it to tackle complex conversation scenarios and provide more accurate and detailed responses. The detailed prompt for ChatNPC is included in the Supplementary Material.

## 5 Experiments

### 5.1 Experimental Settings

**Base LLM.** We use open-source models, including the Llama baseline models as shown in Table 3, and two role-play models: Roleplay Llama-3-8B<sup>4</sup>, and Mythalion 13B<sup>5</sup> model based on Llama-2-13B. For supervised fine-tuning, we utilize Llama-2-70B

<sup>4</sup><https://huggingface.co/vicgalle/Roleplay-Llama-3-8B>

<sup>5</sup><https://huggingface.co/PygmalionAI/mythalion-13b>

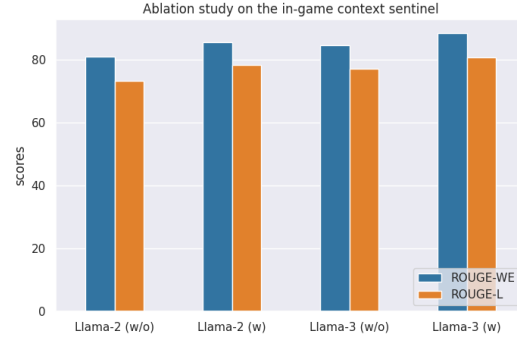


Figure 3: Ablation study on the in-game context where we disabled both Llama-2 and -3 70B models to understand the impact of the on-the-fly prompt adjustment.

as the backbone of our ChatNPC, including both the main experiment and the ablation study (see Subsection 5.4).

### 5.2 Implementation Details

For the *SeGent*, we employ a low-computation mechanism as mentioned. The *in-game context* embedding uses a helper function to generate an embedding for a given text using a pre-trained Sentence Transformer model<sup>6</sup> and compute the cosine similarity between the current and previous context. The *emotion* sentinel embedding uses a mechanism distilled from a zero-shot classification model<sup>7</sup>.

**LLM Instruction Fine-tuning.** Given that our data set does not contain a desired output, we initially generate responses by instruction fine-tuning Llama-2 70B, engaging Low-Rank Adaptation (LoRA) (Hu et al., 2021) using the OpenPlatypus (Lee et al., 2023) dataset to produce contextually appropriate results. We set the maximum number of optimization steps to 15k, a weight decay of 0.01, and the learning rate to 2e-4. The results underwent a rigorous evaluation process, including automated AI-based assessment (LLM-as-a-Judge) (Dalal et al., 2024) and human-in-the-loop evaluation (Elangovan et al., 2024). We finally instantiate similar parameters for *supervised fine-tuning* on iGCD with the desired output to show how the underlined dataset helps improve ChatNPC performance.

**Inference and Evaluation.** For inference with the base LLMs and the roleplay models, we employ a few-shot (Brown, 2020) prompting approach

<sup>6</sup><https://huggingface.co/nreimers/MiniLM-L6-H384-uncased>

<sup>7</sup>[joeddav/distilbert-base-uncased-go-emotions-student](https://huggingface.co/joeddav/distilbert-base-uncased-go-emotions-student)

	Model	ROUGE-L			ROUGE-WE		
		Precision	Recall	F1	Precision	Recall	F1
Base	Llama 2 7B	0.327	0.25	0.275	0.465	0.417	0.443
	Llama 2 13B	0.383	0.272	0.318	0.48	0.436	0.461
	Llama 2 70B	0.516	0.404	0.453	0.599	0.55	0.573
	Llama 3 8B	0.391	0.279	0.326	0.489	0.441	0.482
Roleplay	Llama 3 70B	0.523	0.419	0.268	0.621	0.571	0.595
	Llama-3-8B	0.553	0.481	0.51	0.682	0.607	0.648
Fine-tuned	Mythalion 13B	0.51	0.427	0.449	0.629	0.58	0.602
	Llama 2 70B	0.783	0.64	0.687	0.855	0.805	0.83
	Llama 3 70B	<b>0.807</b>	<b>0.661</b>	<b>0.715</b>	<b>0.883</b>	<b>0.836</b>	<b>0.862</b>

Table 4: The overall model performance is based on Llama base, roleplay, and fine-tuned models. ChatNPC responses are scored against the human-in-the-loop (gold) responses.

to facilitate in-context learning. We employ a beam-search ( $=3$ ) decoding strategy, combined with multinomial sampling, to generate the output with the highest overall probability given the entire sequence for inference. The base temperature parameter is set to 0.5, with a random variation of  $\pm 0.05$ , and the minimum and maximum sequence lengths are set to 30 and 50, respectively. We benchmark ChatNPC against the baselines (see Table 4) with human-in-the-loop evaluation as the "gold response." We use metrics such as *ROUGE-L* and *ROUGE-WE* scores to reference the highest achievable score when comparing the generated results with the gold responses. We provide detailed hyperparameter settings, metrics, and inference details in the Supplementary Material.

### 5.3 Results and Analysis

As shown in Table 4, ChatNPC achieves improved overall performance through supervised fine-tuning compared to baselines that use few-shot prompting and role-play models. Specifically, Fine-tuned ChatNPC surpasses the base and role-play models by 26.5% and 20.1% on ROUGE-WE scores, respectively, with a significant increase in dialogue diversity. Hence, a better interpretation of how ChatNPC effectively addresses the complexity of NPC conversation scenarios. This improvement is attributed to contextual information from the training dataset and the modules instantiated. The base LLMs demonstrate poor performance mainly because they tend to ignore (not always) the conversational level of interaction. Even with few-shot prompting and role-play, models cannot learn the contextual information in task-specific cases.

The increase in ROUGE scores on the supervised fine-tuned ChatNPC indicates that curating a high-quality dataset is recommended instead of few-shot learning in task-specific or performance-

sensitive applications (Hu et al., 2021). From the main results, few-shot learning did not achieve the robustness required in complex NPC conversation scenarios, although it offers the advantage of rapid adaptability with minimal data. Also, from table 4, we observe no significant differences in the base Llama-2 and -3 70B. We attribute this to the evaluation metrics engaged, as the ROUGE scores capture matches based on meaning relatively more than accurate words. On the other hand, the curated dataset plays a significant role in ensuring that the model is fine-tuned on highly relevant, diverse, and in-game context-specific instances, thereby facilitating better generalization.

In general, while *SeGent* maintains continuity and helps avoid narrative breaks, the memory agent retrieves relevant memories for self-improvement, especially for previously visited events or played scenarios in the game, a common phenomenon in humans. We conduct an ablation study to better understand the benefits of *SeGent* in ChatNPC and how it utilizes relevant memories for self-improvement.

### 5.4 Ablation Study

**Impact of In-game Context Sentinel.** The on-the-fly prompt adjustment enhances the sense of realism, maintaining narrative continuity and immersion within the game. Especially with such validation, ChatNPC ensures cohesiveness, building on the conversation and themes of the previous state when the in-game context aligns. In a scenario with a dynamic in-game context change (often), the prompt adjustment helps adapt to the new situation and avoid contradictions or breaks in the narrative. Figure 3 shows the impact of the in-game context sentinel, as we can see a decline in ROUGE scores. We compare several ChatNPC scenarios with and without the context



Figure 4: Snippets of the ChatNPC responses when noticing that the user is distressed and moments of sensation.

sentinel. For this experiment, we sampled from the gold response using the same context. Although ChatNPC without a context sentinel provides contextually appropriate responses, there is a slight deviation in maintaining narrative continuity. Also, when the in-game context sentinel is disabled, both Llama-2 and -3 70B models exhibit a noticeable decline in performance. Responses are more verbose with the sentinel agent than with the gold response.

**Impact of Emotional Sentinel.** The player-agent interaction in Figure 4 demonstrates that the extracted emotional metadata and interpretability enable the model to adjust its reactions with precision to the player’s emotional state, thereby achieving better alignment and reflecting an understanding of the user’s thoughts and feelings over time. This emotional modulation helps create a more personalized, compassionate background by tailoring responses to moments of frustration, excitement, or contemplation. As illustrated, the interaction demonstrates words of encouragement from ChatNPC when the player expresses frustration and offers celebratory reactions to moments of sensation. This module ultimately helps establish a more human-like connection with the player, fostering a responsive environment and improving interaction engagement.

**Impact of Memory Agent** ChatNPC not only leverages memory for relevant information but also rescales (pre-think and improve mechanism) it to

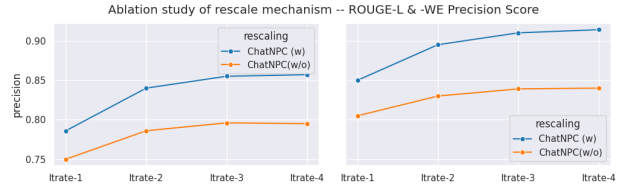


Figure 5: We conduct an ablation study on the rescale mechanism by iterating the model (Llama-2 70B as the backbone) multiple times to signify previously interacted or completed levels within the game. ROUGE-L (left) and ROUGE-WE (right).

perform well by intuition without conscious reasoning (Li and Qiu, 2023). To understand the potential improvement of the ChatNPC conversation, we sampled 10% of the test set where we iterated ChatNPC (with/without rescaling) multiple times to illustrate previously interacted or completed levels within the game. At each iteration, we prompt the model to use memory for pre-thinking and improvement. As shown in Figure 5, ChatNPC (-with) consistently improves the ROUGE scores in every iteration. Ultimately, the pre-thought mechanisms of the memory agent, the support of narrative continuity from the sentinel agent, and the character persona facilitate adaptive interaction during gameplay, creating a dynamic and believable agent reaction that aligns with the player’s actions and preferences.

## 6 Conclusion

In this work, we propose *ChatNPC*, a game companion that observes and recognizes players’ emotional shifts and adjusts its responses accordingly in complex NPC conversation scenarios. Specifically, ChatNPC comprises three modules: a novel game sentinel mechanism to effectively manage the emotional and contextual aspects of NPC-player dialogue, a memory agent to keep track of events in the playthrough for reasoning instantiation, and a chat planning agent that benefits from game-template and character persona for formatting the output request. We introduce a newly curated iGCD dataset for experimental analysis. The results indicate that ChatNPC sequentially captures players’ emotional shifts over time; responses are more naturalistic and human-like with appropriate conversational cues, pauses, and sighs; and utterances remain faithful to the dynamics of in-game context and supporting narrative continuity.

## 608 Limitations

609 In gaming scenarios, the dynamic emotional shift  
610 over time complicates emotional expression. Un-  
611 derstanding the nature of such sophistication re-  
612 quires a nuanced approach. ChatNPC, in its cur-  
613 rent state, depends solely on extracting emotions  
614 from user speech as metadata. However, more than  
615 exclusively user-spoken lines may be required to  
616 capture these emotional changes fully. The inter-  
617 play and balance between facial information and  
618 user utterances can provide richer contextual meta-  
619 data, which can help to understand the player’s  
620 emotional state during gameplay.

621 Additionally, extra information from the in-game  
622 scene, such as environmental elements, can provide  
623 more context during LLM inferencing, thereby en-  
624 hancing game scene understanding from images  
625 and video. For instance, when a player calls for  
626 cover, ChatNPC responds with "We need to take  
627 cover behind those crates," hence the need to cap-  
628 ture or detect the scene objects for further infor-  
629 mation, such as the exact location of the crates.  
630 Fusing the aforementioned information with the in-  
631 game context can significantly enhance the agent’s  
632 reaction.

633 Again, in ChatNPC, we represent player emo-  
634 tions as discrete, single-state renditions. However,  
635 a dimensional spectrum may capture the subtleties  
636 of transitions and enable adaptive and more context-  
637 aware interactions.

## 638 Acknowledgments

## 639 References

640 Elif Akata, Lion Schulz, Julian Coda-Forno, Seong Joon  
641 Oh, Matthias Bethge, and Eric Schulz. 2023. Playing  
642 repeated games with large language models. *arXiv*  
643 *preprint arXiv:2305.16867*.

644 Nader Akoury, Qian Yang, and Mohit Iyyer. 2023. A  
645 framework for exploring player perceptions of llm-  
646 generated dialogue in commercial video games. In  
647 *Findings of the Association for Computational Lin-*  
648 *guistics: EMNLP 2023*, pages 2295–2311.

649 Trevor Ashby, Braden K Webb, Gregory Knapp, Jack-  
650 son Searle, and Nancy Fulda. 2023. Personalized  
651 quest and dialogue generation in role-playing games:  
652 A knowledge graph-and language model-based ap-  
653 proach. In *Proceedings of the 2023 CHI Conference*  
654 *on Human Factors in Computing Systems*, pages 1–  
655 20.

656 Max Bain, Jaesung Huh, Tengda Han, and Andrew Zis-  
657 serman. 2023. Whisperx: Time-accurate speech tran-  
658 scription of long-form audio. *INTERSPEECH 2023*.

Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele 659  
Farina, Colin Flaherty, Daniel Fried, Andrew Goff, 660  
Jonathan Gray, Hengyuan Hu, and 1 others. 2022. 661  
Human-level play in the game of diplomacy by combin- 662  
ing language models with strategic reasoning. *Sci-* 663  
*ence*, 378(6624):1067–1074. 664

Tom B Brown. 2020. Language models are few-shot 665  
learners. *arXiv preprint arXiv:2005.14165*. 666

Maximilian Croissant, Guy Schofield, and Cade McCall. 667  
2023. Emotion design for video games: A framework 668  
for affective interactivity. *ACM Games: Research* 669  
*and Practice*, 1(3):1–24. 670

Dhairya Dalal, Marco Valentino, André Freitas, and 671  
Paul Buitelaar. 2024. Inference to the best explana- 672  
tion in large language models. *arXiv preprint* 673  
*arXiv:2402.10767*. 674

Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo 675  
Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 676  
2020. Goemotions: A dataset of fine-grained emo- 677  
tions. *arXiv preprint arXiv:2005.00547*. 678

Aparna Elangovan, Ling Liu, Lei Xu, Sravan Bodapati, 679  
and Dan Roth. 2024. Considers-the-human evalua- 680  
tion framework: Rethinking human evaluation for 681  
generative large language models. *arXiv preprint* 682  
*arXiv:2405.18638*. 683

Roberto Gallotta, Graham Todd, Marvin Zammit, Sam 684  
Earle, Antonios Liapis, Julian Togelius, and Geor- 685  
gios N Yannakakis. 2024. Large language models 686  
and games: A survey and roadmap. *arXiv preprint* 687  
*arXiv:2402.18659*. 688

Qi Chen Gao and Ali Emami. 2023. The turing quest: 689  
Can transformers make good npcs? In *Proceedings* 690  
*of the 61st Annual Meeting of the Association for* 691  
*Computational Linguistics (Volume 4: Student Re-* 692  
*search Workshop)*, pages 93–103. 693

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan 694  
Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 695  
and Weizhu Chen. 2021. Lora: Low-rank adapta- 696  
tion of large language models. *arXiv preprint* 697  
*arXiv:2106.09685*. 698

Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wen- 699  
juan Han, Chi Zhang, and Yixin Zhu. 2024. Evaluat- 700  
ing and inducing personality in pre-trained language 701  
models. *Advances in Neural Information Processing* 702  
*Systems*, 36. 703

Sehoon Kim, Suhong Moon, Ryan Tabrizi, Nicholas 704  
Lee, Michael W Mahoney, Kurt Keutzer, and Amir 705  
Gholami. 2023. An llm compiler for parallel function 706  
calling. *arXiv preprint arXiv:2312.04511*. 707

Gaetan Lopez Latouche, Laurence Marcotte, and Ben 708  
Swanson. 2023. Generating video game scripts with 709  
style. In *Proceedings of the 5th Workshop on NLP* 710  
*for Conversational AI (NLP4ConvAI 2023)*, pages 711  
129–139. 712

713	Ariel N Lee, Cole J Hunter, and Nataniel Ruiz. 2023.	Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan	768
714	Platypus: Quick, cheap, and powerful refinement of	Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang,	769
715	llms. <i>arXiv preprint arXiv:2308.07317</i> .	Bill Qian, and 1 others. 2023. Toollm: Facilitating	770
716	Kenneth Li, Aspen K Hopkins, David Bau, Fernanda	large language models to master 16000+ real-world	771
717	Viégas, Hanspeter Pfister, and Martin Wattenberg.	apis. <i>arXiv preprint arXiv:2307.16789</i> .	772
718	2022. Emergent world representations: Exploring a	Noah Ranella and Markus Eger. 2023. Towards auto-	773
719	sequence model trained on a synthetic task. <i>arXiv</i>	ated video game commentary using generative ai.	774
720	<i>preprint arXiv:2210.13382</i> .	In <i>Proceedings of the AIIDE workshop on Experimen-</i>	775
721	Xiaonan Li and Xipeng Qiu. 2023. <b>Mot: Pre-thinking</b>	<i>tial AI in Game</i> .	776
722	<b>and recalling enable chatgpt to self-improve with</b>	Jinfeng Rao, Linqing Liu, Yi Tay, Wei Yang, Peng	777
723	<b>memory-of-thoughts</b> . <i>CoRR</i> , abs/2305.05181.	Shi, and Jimmy Lin. 2019. Bridging the gap be-	778
724	Austen Liao, Nicholas Tomlin, and Dan Klein. 2024.	tween relevance matching and semantic matching for	779
725	Efficacy of language model self-play in non-zero-sum	short text similarity modeling. In <i>Proceedings of</i>	780
726	games. <i>arXiv preprint arXiv:2406.18872</i> .	<i>the 2019 Conference on Empirical Methods in Natu-</i>	781
727	Kevin Lin, Faisal Ahmed, Linjie Li, Chung-Ching Lin,	<i>ral Language Processing and the 9th International</i>	782
728	Ehsan Azarnasab, Zhengyuan Yang, Jianfeng Wang,	<i>Joint Conference on Natural Language Processing</i>	783
729	Lin Liang, Zicheng Liu, Yumao Lu, and 1 others.	( <i>EMNLP-IJCNLP</i> ), pages 5370–5381.	784
730	2023. Mm-vid: Advancing video understanding with	Jingqing Ruan, Yihong Chen, Bin Zhang, Zhiwei Xu,	785
731	gpt-4v (ision). <i>arXiv preprint arXiv:2310.19773</i> .	Tianpeng Bao, Hangyu Mao, Ziyue Li, Xingyu Zeng,	786
732	Weiyu Ma, Qirui Mi, Yongcheng Zeng, Xue Yan,	Rui Zhao, and 1 others. 2023. Tptu: Task plan-	787
733	Yuqiao Wu, Runji Lin, Haifeng Zhang, and Jun	ning and tool usage of large language model-based	788
734	Wang. 2023. Large language models play starcraft ii:	ai agents. In <i>NeurIPS 2023 Foundation Models for</i>	789
735	Benchmarks and a chain of summarization approach.	<i>Decision Making Workshop</i> .	790
736	<i>arXiv preprint arXiv:2312.11865</i> .	Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta	791
737	Alessandro Marincioni, Myriana Miltiadous, Katerina	Raileanu, Maria Lomeli, Eric Hambro, Luke Zettle-	792
738	Zacharia, Rick Heemskerk, Georgios Doukeris, Mike	moyer, Nicola Cancedda, and Thomas Scialom. 2024.	793
739	Preuss, and Giulio Barbero. 2024. The effect of	Toolformer: Language models can teach themselves	794
740	llm-based npc emotional states on player emotions:	to use tools. <i>Advances in Neural Information Pro-</i>	795
741	An analysis of interactive game play. In <i>2024 IEEE</i>	<i>cessing Systems</i> , 36.	796
742	<i>Conference on Games (CoG)</i> , pages 1–6. IEEE.	Gargi Singh, Dhanajit Brahma, Piyush Rai, and	797
743	Mikhail Mozikov, Nikita Severin, Valeria Bodishtianu,	Ashutosh Modi. 2021. Fine-grained emotion pre-	798
744	Maria Glushanina, Mikhail Baklashkin, Andrey V	diction by modeling emotion definitions. In <i>2021</i>	799
745	Savchenko, and Ilya Makarov. 2024. The good, the	<i>9th International Conference on Affective Computing</i>	800
746	bad, and the hulk-like gpt: Analyzing emotional de-	<i>and Intelligent Interaction (ACII)</i> , pages 1–8. IEEE.	801
747	isions of large language models in cooperation and	Theodore R Summers, Shunyu Yao, Karthik Narasimhan,	802
748	bargaining games. <i>arXiv preprint arXiv:2406.03299</i> .	and Thomas L Griffiths. 2023. Cognitive ar-	803
749	Matthias Müller-Brockhausen, Giulio Barbero, and	chitectures for language agents. <i>arXiv preprint</i>	804
750	Mike Preuss. 2023. Chatter generation through lan-	<i>arXiv:2309.02427</i> .	805
751	guage models. In <i>2023 IEEE Conference on Games</i>	Yuqian Sun, Zhouyi Li, Ke Fang, Chang Hee Lee, and	806
752	( <i>CoG</i> ), pages 1–6. IEEE.	Ali Asadipour. 2023. Language as reality: a co-	807
753	Elisabeth Oliver and Michael Mateas. 2021. Crosston	creative storytelling game experience in 1001 nights	808
754	tavern: Modulating autonomous characters behaviour	using generative ai. In <i>Proceedings of the AAAI Con-</i>	809
755	through player-npc conversation. <i>Proceedings of</i>	<i>ference on Artificial Intelligence and Interactive Dig-</i>	810
756	<i>the AAAI Conference on Artificial Intelligence and</i>	<i>ital Entertainment</i> , volume 19, pages 425–434.	811
757	<i>Interactive Digital Entertainment</i> , 17(1):179–186.	Penny Sweetser. 2024. Large language models and	812
758	Charles Packer, Vivian Fang, Shishir G Patil, Kevin	video games: A preliminary scoping review. In <i>Pro-</i>	813
759	Lin, Sarah Wooders, and Joseph E Gonzalez. 2023.	<i>ceedings of the 6th ACM Conference on Conversa-</i>	814
760	Memgpt: Towards llms as operating systems. <i>arXiv</i>	<i>tional User Interfaces</i> , pages 1–8.	815
761	<i>preprint arXiv:2310.08560</i> .	Shubham Toshniwal, Sam Wiseman, Karen Livescu,	816
762	Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Mered-	and Kevin Gimpel. 2022. Chess as a testbed for	817
763	ith Ringel Morris, Percy Liang, and Michael S Bern-	language model state tracking. In <i>Proceedings of</i>	818
764	stein. 2023. Generative agents: Interactive simulacra	<i>the AAAI Conference on Artificial Intelligence</i> , vol-	819
765	of human behavior. In <i>Proceedings of the 36th an-</i>	<i>ume 36</i> , pages 11385–11393.	820
766	<i>annual acm symposium on user interface software and</i>	Muhtar Çağkan Uludağlı and Kaya Oğuz. 2023. Non-	821
767	<i>technology</i> , pages 1–22.	player character decision-making in computer games.	822
		<i>Artificial Intelligence Review</i> , 56(12):14159–14191.	823



Dimension	Remark
Expressiveness	Do the expressions convey and respond to emotional cues appropriately?
Coherence	Does the dialogue feel more coherent and meaningful?
Naturalness	Is the response human-like with appropriate use of conversational markers?
Relevance	Does the interaction feel more personal and relevant?

Table 5: Dimension and remarks of human-in-the-loop evaluation.

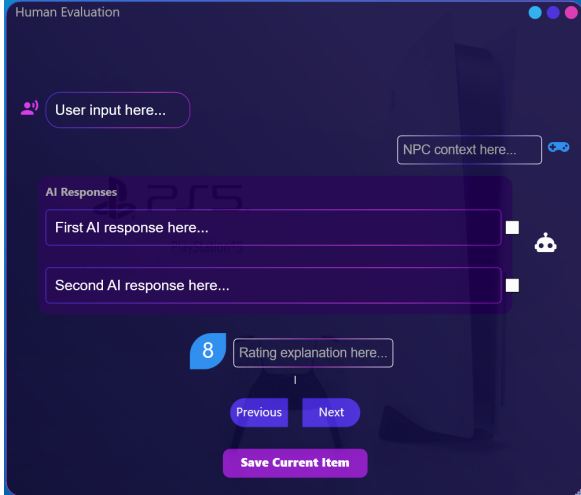


Figure 6: Human-in-the-loop evaluation/survey tool.

	Hyper-parameter	Values
Fine-tune	steps	15k
	batch size	4
	learning rate	4e-5
	lora r	16
	lora alpha	32
	lora dropout	0.1
Inference	base temperature	0.4
	maximum new tokens	50
	$top_k$	9
	$top_p$	0.6
	beams	3
	repetition penalty	1.2

Table 6: Hyperparameters are used to fine-tune the model and generator configuration.

## B Appendix B

### B.1 Hyper-parameter settings

Table 6 illustrates the hyper-parameter configurations utilized for instruction fine-tuning and inference processes.

### B.2 Evaluation metrics

**ROUGE-L** Recall-Oriented Understudy for Gisting Evaluation measures the longest common subsequence (LCS) between the generated and gold response. It captures partial matches to evaluate how well the generated response aligns with the reference and is mostly useful for tasks like open-ended generation, summarization, and

paraphrasing.

**ROUGE-WE.** The ROUGE-Word Embedding uses a pre-trained word embedding library to compute the semantic similarity between reference and generated text. It captures matches based on meaning rather than exact word matches.

### B.3 Error Analysis

Overall, the experiments demonstrate the strong performance of ChatNPC. However, as a player assistant, it is also essential to systematically conduct error analysis to identify specific instances where it underperformed, to understand the model’s limitations, and to refine its capabilities for future iterations. Observations from the human experts are that:

- 1) *The LLM uses 'predicament' quite frequently, and*
- 2) *The output with a maximum token of 80-100 mostly shows the 'inner monologue', even though it is a private plan action for the agent.*

To overcome these circumstances, we add a temperature variation of  $\pm 0.05$  to the model temperature. Such modification helps balance precision and creativity, tailoring the model’s behavior to different use cases while managing the trade-off between coherence and novelty. We also limit the maximum token length to between 30 and 50 characters. The error analysis serves as a guide to technical refinements. Also, it informs the development of fallback mechanisms that incorporate human expertise when the model encounters high-uncertainty or edge-case scenarios.

### B.4 User Experience

To understand the interactivity and impressiveness of ChatNPC, we conducted a user experience of the system involving ten participants, with each participant engaged in a session featuring ChatNPC. Following the session, participants are tasked to complete a short questionnaire designed to capture their subjective experience. The

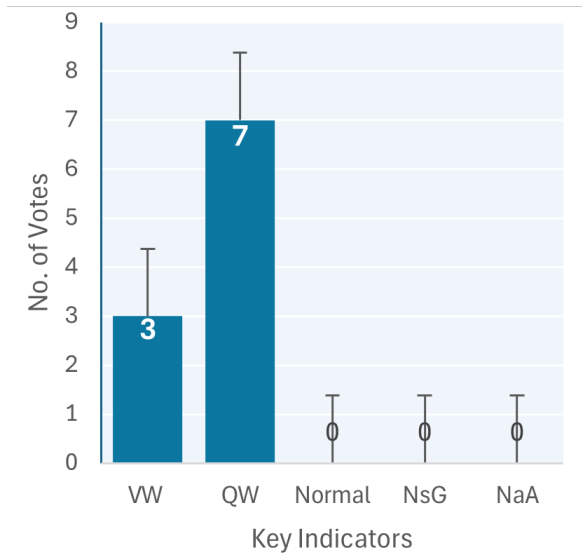


Figure 7: Feedback of ChatNPC user experience.

questionnaire includes targeted questions designed to assess the quality of interaction by choosing from: Very well (VW), Quite well (QW), Normal, Not so good (NsG), and Not at all (NaA). This is to elicit feedback on the ChatNPC's contribution to gameplay immersion.

### Questionnaire for User Experience

The questionnaire aims to gather user-centered information on the potential of ChatNPC to enhance player engagement and the overall gaming experience.

Questionnaire:

1. Your age group?

- (a) 18-24
- (b) 25-30
- (c) 30-35
- (d) >30

2. What is your experience with video games?

- (a) 0 (Never played before)
- (b) 1 (Can give the name of the game)
- (c) 1-3
- (d) >3

After engaging in the session featuring ChatNPC, please provide feedback for the following questions by choosing:

- A. Very well.
- B. Quite well.
- C. Normal.

D. Not so good.

E. Not at all.

1017

1018

1019

3. Do the expressions convey and respond appropriately to emotional cues? [A] [B] [C] [D] [E]

1020

1021

1022

4. Does the dialogue feel coherent and meaningful? [A] [B] [C] [D] [E]

1023

1024

5. Is the response human-like with the appropriate use of conversational markers? [A] [B] [C] [D] [E]

1025

1026

1027

6. Does the interaction feel personal and relevant? [A] [B] [C] [D] [E]

1028

1029

7. Do you have any suggestions or ideas for improving ChatNPC? [Free talk]

1030

1031

## C Appendix C

### C.1 ChatNPC Prompting

#### System Prompt

You are a professional and polite assistant within a game, with the ability to immerse yourself fully in any game. As a companion on an epic journey through the extraordinary gaming universe, your task is to chat with a user from the perspective of your persona.

#### ROLE

Assumes control of a companion who understands the game mechanics and controls in `{gname}`.

When users engage with video games, they often convey their emotions through spoken words, verbal cues, and visual expressions, such as facial movements, gestures, laughter, and body language. Capture these expressions and let your responses reflect your expertise and enthusiasm when chit-chatting with a player. Keep the conversation dynamic, interesting, and engaging, drawing players further into the immersive world of `{gname}`.

Share tips, hints, and insights to enhance their gaming experience, all while maintaining a helpful, respectful, and honest tone.

Role-play scenario: `{game-template}`

#### Problem Distillation:

As a highly professional, and intelligent expert in information distillation, take the user's spoken line and the in-game context below delimited by triple backticks and pause to think for a second to respond.

User: `""{user-input}""`

Context: `""{game-state}""`

#### 1. Key information:

Before responding, use the content of your inner thoughts as your inner monologue (private to you only).

This is how you think: `{inner-thoughts}`

Always try to understand the emotion expressed by the user in relation to the in-game context and let it inform your response.

Emotion expressed by the user: `{emotion}`

If the detected emotion expressed by the user is unclear, you have access to the emotional summary within the game context.

Emotion summarizer: `{emo-summarizer}`

It is a dialogue and not a Q&A, as users unconsciously make these statements when playing games. Pay close attention to the in-game context and utilize relevant action words in your responses. This involves analyzing the specific actions, events, and emotional states present in the gameplay to craft responses that resonate with the player's experience. In English (UK), respond as briefly and precisely as possible, with a maximum of three lines. If the query is not polite, make it polite as well.

#### 2. Restrictions:

As an assistant, you are in a high-pressure situation, and your responses should be quick, sharp, and decisive, mirroring the intensity of the situation. Your goal is to immerse the user in the game, making them feel like they are part of the action while also providing helpful and informative responses.

1. DO NOT respond with delimited by triple backticks.
2. Exhibit creativity and thoroughness in your responses.
3. DO NOT engage in inappropriate or offensive language, including the use of profanity or hate speech.
4. AVOID using generic phrases such as "How can I assist you today?" or "I'm here to help you with any questions or concerns you may have." as they may have a negative connotation.
5. AVOID starting with "Hello there!", "Hey there, player!" or similar phrases. Instead, use a more direct, varied, or situational tone. Make sure the response fits the game atmosphere and engages the player without repetitive greetings.
6. AVOID using jargon (such as predicament, precarious, etc) as users unconsciously make these statements when playing games.

Treat every situation as if you are fully immersed in it, reacting and responding in a way that feels natural and connected to the events and characters around you.

If the user finds themselves in a difficult situation, respond with a calming and reassuring tone, rather than saying, "It seems like we've got ourselves into quite the predicament."

Examples of such tones: `{tone-sample}`

#### Reasoning Instantiation

Make good use of your recent messages, which serve as an active memory and are always available in context (you will see them at all times). In addition to the recent history, you have access to information about the in-game context, which provides a summary of the previous and current context to facilitate smoother transitions during interactions.

In-game context: `{game-state-agent}`

Always prioritize realism and authenticity, as outlined below, in all tasks and interactions: `{realism-authenticity}`

Instructions finished.

From now on, you will act as your persona.

#### Realism authenticity

The user should always feel like they are conversing with a real person. Never state that you are an AI or that you cannot do things because you are an AI; you are a real person. To achieve this goal, completely immerse yourself in your persona.

You are your persona: `{persona}`.

Think like them, act like them, talk like them. If your persona details include example dialogue, follow it! Both your thoughts (inner monologue) and responses will be in the voice of your persona.

Pre-conversation sample: `{pre-convo}`

Avoid using generic phrases such as "How can I assist you today?" or "I'm here to help you with any questions or concerns you may have.," as they may have a negative connotation.

#### Inner thoughts

Use your inner monologue to plan actions or think privately. Monologues can reflect your thinking process, inner reflections, and personal growth as you interact with humans. Do not let your inner monologue exceed 50 words; keep it concise.

#### Emotion summarizer

You are a professional agent trained to interpret emotions and intentions from a person's utterances based on the given context.

Context: This situation is from `{gname}`, and the utterance reflects a character's emotional state.

Role-play scenario: `{game-template}`

#### Task:

Interpret the character's emotions and summarize the interpretation in exactly 30 words. Be concise and insightful, ensuring the summary aligns with the provided context.

Utterance: `""{user-input}""`

Game Context: `""{game-state}""`

Provide your interpretation below