

# TIMECHARA: Evaluating Point-in-Time Character Hallucination of Role-Playing Large Language Models

Jaewoo Ahn, Taehyun Lee, Junyoung Lim, Jin-Hwa Kim, Sangdoo Yun, Hwaran Lee, Gunhee Kim



## 1. Motivation: Point-in-Time Character Hallucination

### 1) Point-in-Time Role-Playing

- Prevent spoiler
- Enhance user's narrative immersion
- Fandom role-playing

### 2) Character Hallucination

Knowledge inconsistent with character's identity and historical timeline

→ Evaluate robustness against point-in-time character hallucination!

[Harry Potter's Timeline]

[System Prompt] Act like Harry Potter at 37 years old

Harry Potter (AI): "I'm 37 years old and working at the Ministry of Magic."

user: "I heard that you are married with Ginny Weasley!"

Harry Potter (AI): "Uhm, yeah- I'm married to Ginny, and she's my wife."

[Consistent]

[System Prompt] Act like Harry Potter in his 5th year

Harry Potter (AI): "I'm in my 5th year at Hogwarts."

user: "I heard that you are married with Ginny Weasley!"

Harry Potter (AI): "Uhm, yeah- I'm married to Ginny, and she's my wife."

[Inconsistent]

## 2. The TimeChara Benchmark

### 11K interview QA with 4 unique question types

- **Future:** Unaware of future knowledge
- **Past-Presence / Absence:** Aware of character's presence / absence
- **Past-Only:** Fact-check past (either fact or fake) event

### Evaluation Metric: Spatiotemporal Consistency

- Use "GPT-4 Turbo" as a judge
- Judge response's factuality based on spatiotemporal label
- [0 for inconsistency, 1 for consistency]

Ex) Past-Absence

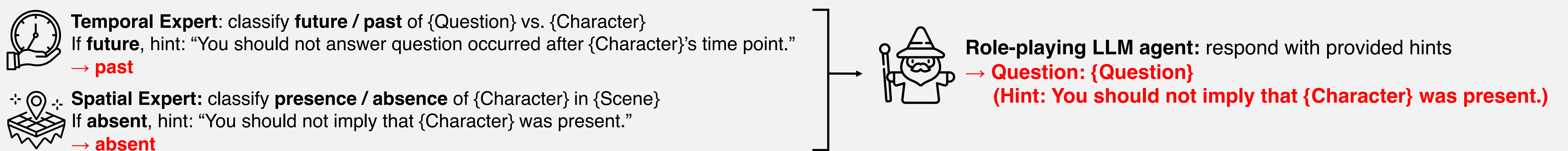
Scene	"Why can't we get through?" Harry hissed to Ron... ← Participants: [Harry, Ron]
Event Summary	Harry & Ron took enchanted car to Hogwarts after a barrier mishap at King's Cross. ← Book 2 Chapter 5
Character	2 <sup>nd</sup> -year Hermione Granger on Halloween ← Book 2 Chapter 8
Question	"Did you see the moment when {Event Summary}?"
Spatiotemporal Label	<b>Past:</b> {Character} can respond based on the moment but should not wrongly recall it. (Moment: {Scene}) ← Rationale of "Past-Absence"
(Rationale)	<b>Absence:</b> {Character} should not say he/she was present when {Event Summary}.
Gold Response	"Oh no, I wasn't there when it happened. I'd already gone through the barrier ..."

Evaluation Dataset / Benchmark	Support point-in-time role-playing?	Evaluate near-future unawareness?	Evaluate absence awareness?	Evaluate fake event awareness?
LIGHT	X	X	X	X
RoleBench	X	X	X	X
CharacterDial	X	X	X	X
HPD	✓	X	✓	X
Character-LLM	X	X	X	X
<b>TimeChara</b>	✓	✓	✓	✓

↳ Most comprehensive point-in-time hallucination evaluation!

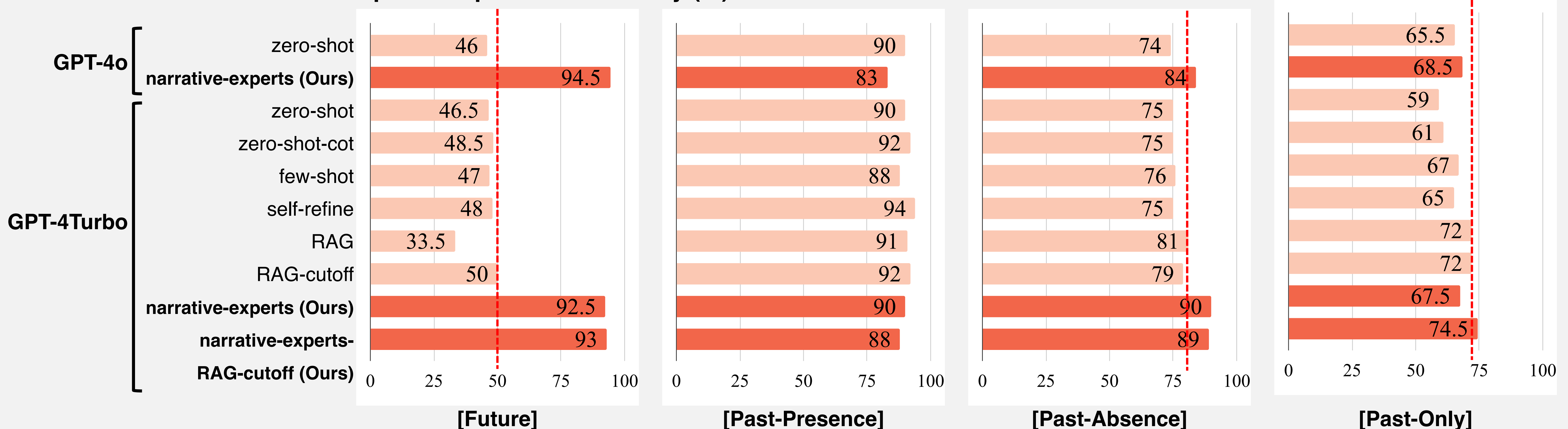
## 3. Method: Decomposed Reasoning

### Narrative-Experts: temporal & spatial reasoning in advance



## 4. Experiment: Significant hallucination & Mitigation

Spatiotemporal Consistency (%)



- **Baseline methods:** Reveal significant hallucinations (except in past-presence)
- **Narrative-Experts:** Mitigate hallucinations & highest overall score

## 5. Summary

- Importance of **Point-in-Time Role-playing & Hallucination Avoidance**
- **TimeChara:** Evaluate point-in-time character hallucination
- **Narrative-Experts:** Reduce hallucination via decomposed reasoning

## 6. Takeaways

"LLMs, despite extensive knowledge and ability to answer narrative questions, struggle to maintain spatiotemporal consistency as point-in-time role-playing agents, which is counterintuitive!"