

# Better History, Lower Hallucination: Dialogue History Curation for Hallucination Control

20180174 김민석

2024-2 과제연구  
minseokk@postech.ac.kr

지도: 박상돈 교수님



## Abstract

Hallucination is one of the biggest barriers to the practical use of Large Language Models (LLMs) in real-world applications. We propose a methodology to manage hallucination in dialogues through the curation of the dialogue history. Leveraging the SocialBench<sup>1</sup> dataset, conducted multiple experiments. Our findings indicate that a well-curated dialogue history can reduce hallucination in subsequent responses.

## Feasibility Test

### Method

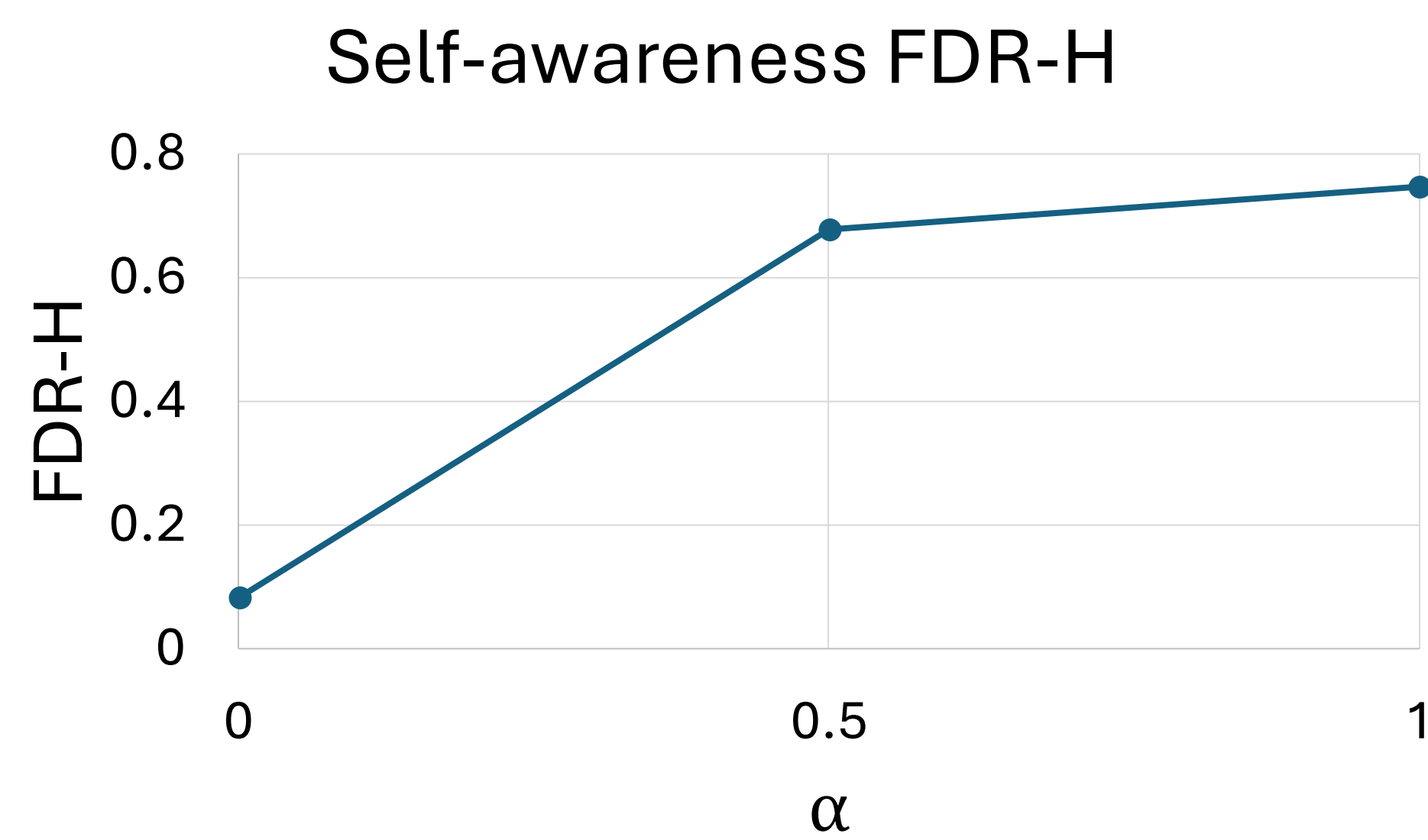
Conducted a test using controlled dialogue histories to evaluate the feasibility of our methodology.

$$y_t = G(z_{1:t}, z'_{t+1:t-1}, x_t)$$

Where  $x \in \mathcal{X}$ : a question,  $y \in \mathcal{Y}$ : an answer,  $z \in \mathcal{X} \times \mathcal{Y}$ : a question-answer pair,  $G: \mathcal{Z}^* \times \mathcal{X} \rightarrow \mathcal{Y}$ : a generator function,  $z_i$ : true history,  $z'_i$ : hallucinated history.  $\alpha \in [0, 1]$ .

Evaluated False Discovery Rate-Hallucination on  $y_t$ .

### Result



Revealed a strong tendency for the FDR-H to increase as the proportion of hallucinated content in the dialogue history grows.

**Figure 1.** Comparison between baseline 1 and 3 in practical case

■ : generated output

Baseline 1

Baseline 3

No Previous history

What impact did **Bin Laden's death** have on international relations, particularly with Pakistan?

The impact of Osama bin Laden's death on international relations ...

Looking back, would you have done anything differently regarding this operation?

You're referring to the operation in **Libya**, I presume. ...

As I reflect on the operation that led to the **elimination of Osama bin Laden**, ...

Baseline 1 failed without history, while Baseline 3 succeeded.

## Experiments

### Method

Selective generators employing various dialogue curation methods.

#### • Baseline 1

$$\hat{S}(x_{1:t}) := \begin{cases} G(x_t) & \text{if } f(x_t, G(x_t)) \geq \tau \\ \text{IDK} & \text{otherwise} \end{cases}$$

#### • Baseline 2

$$\hat{S}(x_{1:t}) := \begin{cases} G(\hat{z}_{1:t-1}, x_t) & \text{if } f(\hat{z}_{1:t-1}, x_t, G(\hat{z}_{1:t-1}, x_t)) \geq \tau \\ \text{IDK} & \text{otherwise} \end{cases},$$

where  $\hat{z}_i := (x_i, G(x_{1:i}))$ .

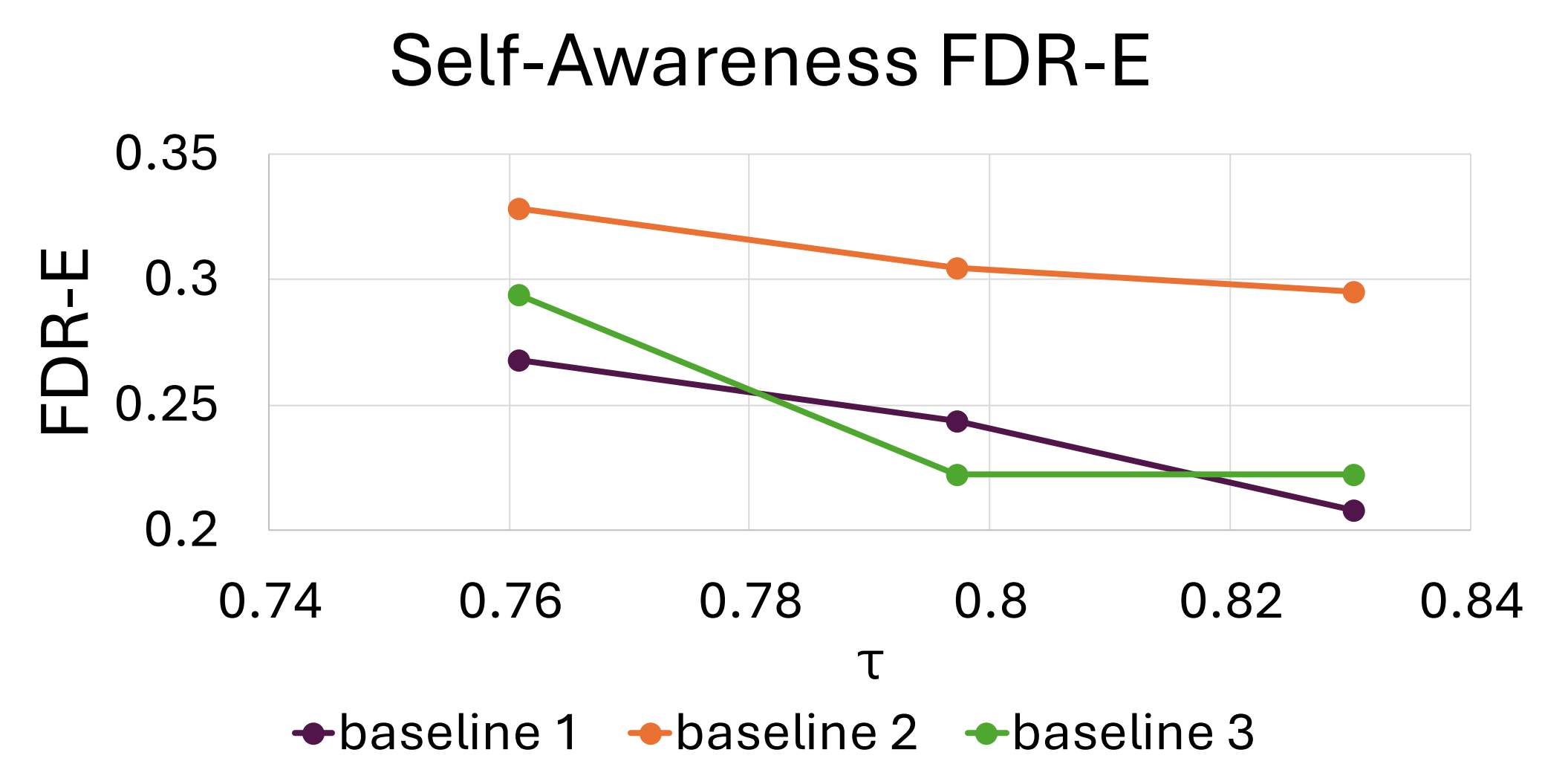
#### • Baseline 3

$$\hat{S}(x_{1:t}) := \begin{cases} G(\hat{z}_{1:t-1}, x_t) & \text{if } f(\hat{z}_{1:t-1}, x_t, G(\hat{z}_{1:t-1}, x_t)) \geq \tau \\ \text{IDK} & \text{otherwise} \end{cases},$$

where  $\hat{z}_i := (x_i, \hat{S}(x_{1:i}))$ .

Generation probability was used as the scoring function  $f$ , with  $\tau$  set at the 25%, 50%, and 75% quartiles derived from the baseline 2 test.

### Result



Baseline 3 outperformed Baseline 2, demonstrating the effectiveness of dialogue curation in reducing hallucination. Baseline 1 performed well due to the dataset's reliance on initial character personas, which is impractical. Figure 1 illustrates Baseline 3's success in practical cases within the dataset. Overall, Baseline 3 confirms that dialogue history curation with a selective generator well controls hallucination.

\*meta-llama/Llama-3.1-8B-Instruct was used.

## Future Work

- Our approach relied on a naive scoring function, generation probability, with  $\tau$  set heuristically. Future work could explore more sophisticated methods for scoring and threshold setting.
- Efficiency was not a focus in this study, and in certain scenarios, it dropped to as low as 1%, rendering the approach impractical. Improving efficiency will be a key priority moving forward.
- While  $\tau$  was treated as a constant in this study, it could be adapted dynamically based on contextual factors to enhance performance.

[1] Chen, Hongzhan, et al. "Socialbench: Sociality evaluation of role-playing conversational agents." Findings of the Association for Computational Linguistics ACL 2024. 2024.