

# Hybrid Graphs for Table-and-Text based Question Answering using LLMs

Ankush Agarwal, Ganesh S, Chaitanya Devaguptapu

Fujitsu Research of India

Contact: ankush.agarwal@fujitsu.com

NAACL 2025

ALBUQUERQUE, NEW MEXICO  
APRIL 29-MAY 4, 2025

## Problem Definition

**Table-Text QA:** The goal is to answer complex questions that require reasoning over both unstructured text (documents) and structured data (tables).

**Question:** What **system** was the most recent game to be released by the company known for the **Street Fighter** franchise released on ?

### Context

**Table : Video Games in 1996**

Release	Title	System	Developer	Notes
Mar 22	Resident Evil	PS1	Capcom	Survival Horror
Mar 29	DBZ: Hyper Dimension	SNES	Bandai	Third X installment

**Document Title: Capcom**  
Capcom Co., Ltd. is a Japanese video game developer and publisher known for creating numerous multi-million selling game franchises, including Mega Man, **Street Fighter**, Resident Evil, Devil May Cry, Dino Crisis, Dead Rising, Sengoku Basara, Ghosts 'n Goblins, Monster Hunter, Breath of Fire, and Ace Attorney as well as games based on Disney animated properties.

**Document Title: Bandai**  
Bandai Co., Ltd. (Kabushiki-gaisha Bandai) is a Japanese toy maker and a producer of many plastic model kits as well as a former video game company. It was the world's third-largest producer of toys in 2008 after Mattel and Hasbro. Its headquarters is located in Tokyo. It produced video games for systems such as **PS1**, **SNES**, **Xbox 360**, etc.

Figure 1: An illustrative example of the Table-Text QA

## Results

### Findings:

- Our method achieves the best performance in a zero-shot setting across various LLMs.
- For Hybrid-QA, our method performs comparably to fine-tuning-based approaches and outperforms them on the OTT-QA dataset.

Datasets	Hybrid-QA					OTT-QA				
Methods	EM (↑)	F1 (↑)	P (↑)	R (↑)	B (↑)	EM (↑)	F1 (↑)	P (↑)	R (↑)	B (↑)
<i>Reader: gpt-4-1106-preview (zero-shot)</i>										
Base	4.60	12.44	12.08	12.25	62.10	4.85	12.44	12.25	14.24	64.30
Base w/ Table & Text (Zhang et al., 2023)	55.40	68.84	68.92	71.79	85.54	58.86	72.28	72.16	74.28	87.51
Base w/ Table & Summarized Text <sup>2</sup>	45.29	58.72	58.78	61.39	81.14	48.31	60.90	61.47	63.12	82.02
Our Method w/o hopwise	58.20	71.54	71.75	74.35	86.30	61.00	72.64	73.60	74.27	88.06
Our Method w/ hopwise	<b>58.40</b>	<b>71.80</b>	<b>71.62</b>	<b>74.22</b>	<b>86.53</b>	<b>62.02</b>	<b>73.02</b>	<b>73.40</b>	<b>75.13</b>	<b>88.18</b>
<i>Reader: gpt-3.5-turbo-1106 (zero-shot)</i>										
Base	4.20	11.54	11.62	12.45	65.05	5.27	12.20	12.35	13.44	66.17
Base w/ Table & Text (Zhang et al., 2023)	40.22	53.47	54.18	55.57	81.18	42.41	54.06	54.04	55.8	81.63
Base w/ Table & Summarized Text <sup>2</sup>	41.19	51.64	52.03	53.12	81.05	37.34	49.58	49.80	51.73	79.87
Our Method w/o hopwise	41.8	52.37	52.82	53.72	81.31	42.19	53.61	54.18	55.30	81.58
Our Method w/ hopwise	<b>44.2</b>	<b>55.82</b>	<b>55.28</b>	<b>56.90</b>	<b>83.98</b>	<b>44.30</b>	<b>54.08</b>	<b>55.04</b>	<b>54.67</b>	<b>81.73</b>
<i>Reader: Llama3-8B (zero-shot)</i>										
Base	2.0	7.07	6.91	7.07	59.05	0.64	7.00	6.77	8.72	59.71
Base w/ Table & Text	28.6	37.05	37.22	48.07	74.01	33.12	43.43	43.53	45.12	76.75
Base w/ Table & Summarized Text <sup>2</sup>	30.33	39.42	39.60	41.30	75.06	31.22	41.72	42.42	42.88	75.57
Our Method w/o hopwise	33.2	41.37	41.77	42.95	75.15	36.08	45.75	46.60	45.75	77.04
Our Method w/ hopwise	<b>37.0</b>	<b>46.43</b>	<b>46.56</b>	<b>48.78</b>	<b>77.55</b>	<b>37.13</b>	<b>47.38</b>	<b>48.24</b>	<b>48.31</b>	<b>77.62</b>

Table 1: **Table-Text QA Evaluation:** We analyze Exact Match (EM), F1-Score, Precision (P), Recall (R), and BERTScore-F1 (B) in (%) to compare our method against baselines in a zero-shot setting using Llama3-8B, GPT-3.5, and GPT-4. The results consistently demonstrate significant improvements across datasets, metrics, and various language models. Base (only reader LLM); w/ Table & Text (table and passages relevant to the question); w/ Table & Summarized Text (table with summarized supporting passages); w/o hopwise (pruned information without considering hop-wise extraction).

## Contributions

- A novel approach, **ODYSSEY**, jointly distills information from structured and unstructured data sources to construct a Hybrid Graph.
- An increase in performance over the current SoTA fine-tuning-free approach, improving EM and F1 scores by 7.3% and 20.9% for Hybrid-QA using GPT-4.
- A significant reduction in the input token size. Our Hybrid Graph based approach uses up to 45% and 53% fewer tokens than the original table and text for the Hybrid-QA and OTT-QA datasets respectively.

## ODYSSEY: A Hybrid Graph Approach

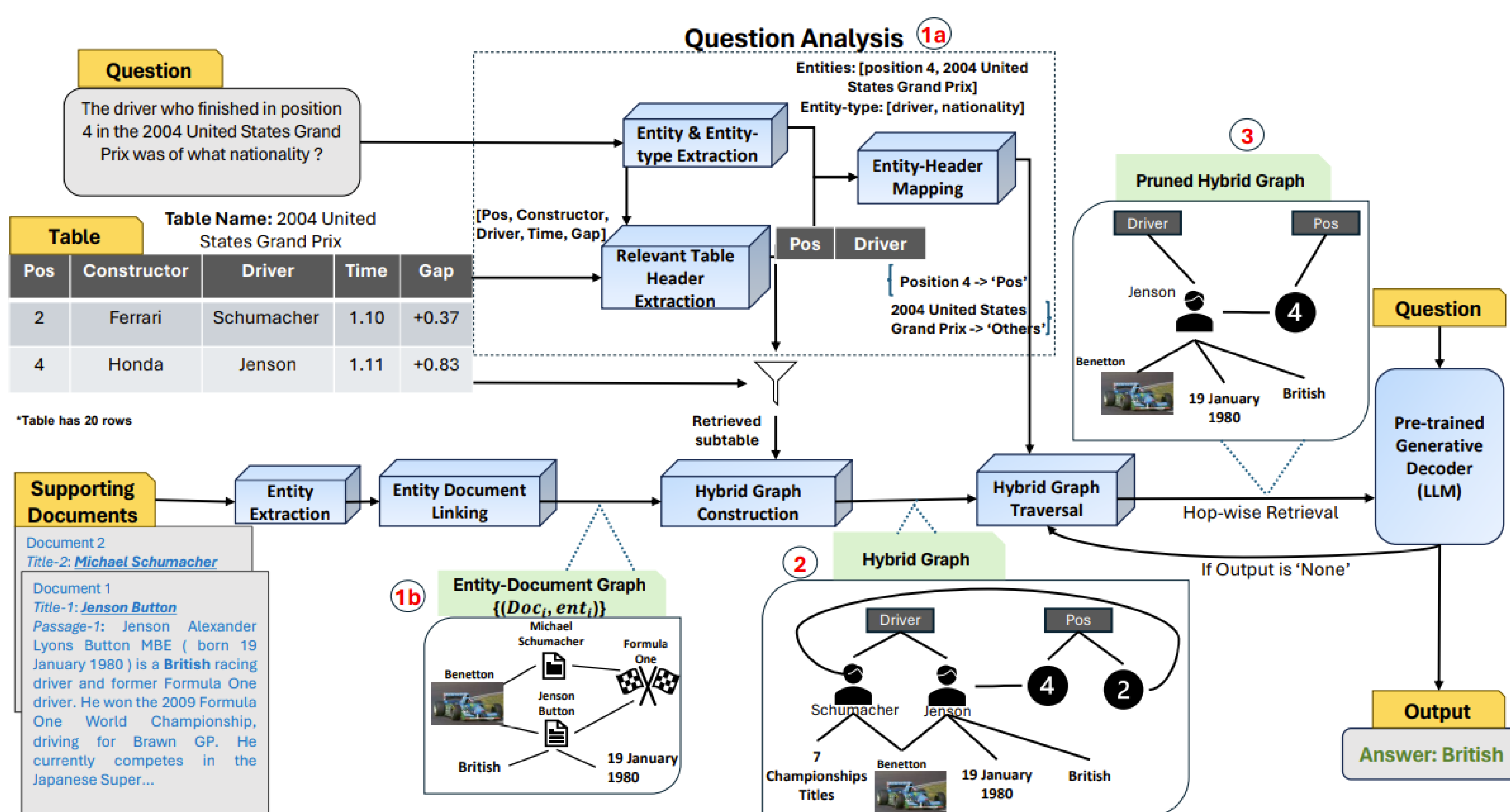


Figure 2: **Overview of the ODYSSEY framework.** Our method comprises of 3 steps: i) Question Analysis, ii) Hybrid Graph Construction, and iii) Hybrid Graph Traversal. First, we begin with Question Analysis (1a in the figure) from where we get question entities, retrieved sub-table, and entity-header mapping. Next, we construct the Entity-Doc Graph (1b in the figure). Using entity-doc graph and retrieved sub-table, we construct the Hybrid Graph (2 in the figure). At last, we perform Hybrid Graph Traversal (3 in the figure) to get the pruned graph which serves as input for the LLM.

Method	EM (↑)	F1 (↑)
<i>Hybrid-QA Fine-Tuning</i>		
HYBRIDER (Chen et al., 2020b)	43.5	50.6
HYBRIDER-LARGE (Chen et al., 2020b)	44.0	50.7
DocHopper (Sun et al., 2021)	47.7	55.0
MuGER <sup>2</sup> (Wang et al., 2022)	57.1	67.3
<b>S<sup>3</sup>HQA (Lei et al., 2023) [SoTA]</b>	<b>68.4</b>	<b>75.3</b>
<i>w/o Fine-Tuning</i>		
Unsupervised-QG (Pan et al., 2021)	25.7	30.5
GPT-4 <sup>†</sup> w/ Retriever (Shi et al., 2024)	24.5	30.0
GPT-4 <sup>†</sup> + CoT (Wei et al., 2022)	48.5	63.0
HProPro <sup>†</sup> (Shi et al. (2024), ACL 2024)	48.0	54.6
<b>ODYSSEY<sup>†</sup> (Our Method)</b>	<b>51.5</b>	<b>66.0</b>

Table 2: **Performance comparison of ODYSSEY with fine-tuning-based and fine-tuning-free approaches.** We evaluate our method against state-of-the-art fine-tuning-based methods as well as approaches without fine-tuning using GPT-4.

Method	EM (↑)	F1 (↑)
<i>OTT-QA Fine-Tuning</i>		
BM25-HYBRIDER (Chen et al., 2020a)	10.3	13.0
Fusion+Cross-Reader (Chen et al., 2020a)	28.1	32.5
CARP (Zhong et al., 2022)	33.2	38.6
CORE (Ma et al., 2022)	49.0	55.7
<b>COS (Ma et al., 2023) [SoTA]</b>	<b>56.9</b>	<b>63.2</b>
<i>w/o Fine-Tuning</i>		
GPT-4 <sup>†</sup> + CoT (Wei et al., 2022)	61.0	72.3
<b>ODYSSEY<sup>†</sup> (Our Method)</b>	<b>62.02</b>	<b>73.02</b>

## Analysis

### Efficient Query Context handling

Method	Input Token Size (↓)	Input Token Cost (↓)
<i>Dataset: Hybrid-QA</i>		
Original Context	7195	\$71.95
Summarized	3923	\$39.23
<b>Our Method</b>	<b>3857</b>	<b>\$38.57</b>
<i>Dataset: OTT-QA</i>		
Original Context	5866	\$58.66
Summarized	3778	\$37.78
<b>Our Method</b>	<b>2745</b>	<b>\$27.45</b>

Table 3: **Reader Input Token Count and Cost:** We compare our method with baselines on average reader input token size and its pricing in dollars w.r.t. GPT-4 Turbo OpenAI pricing for 1000 samples.

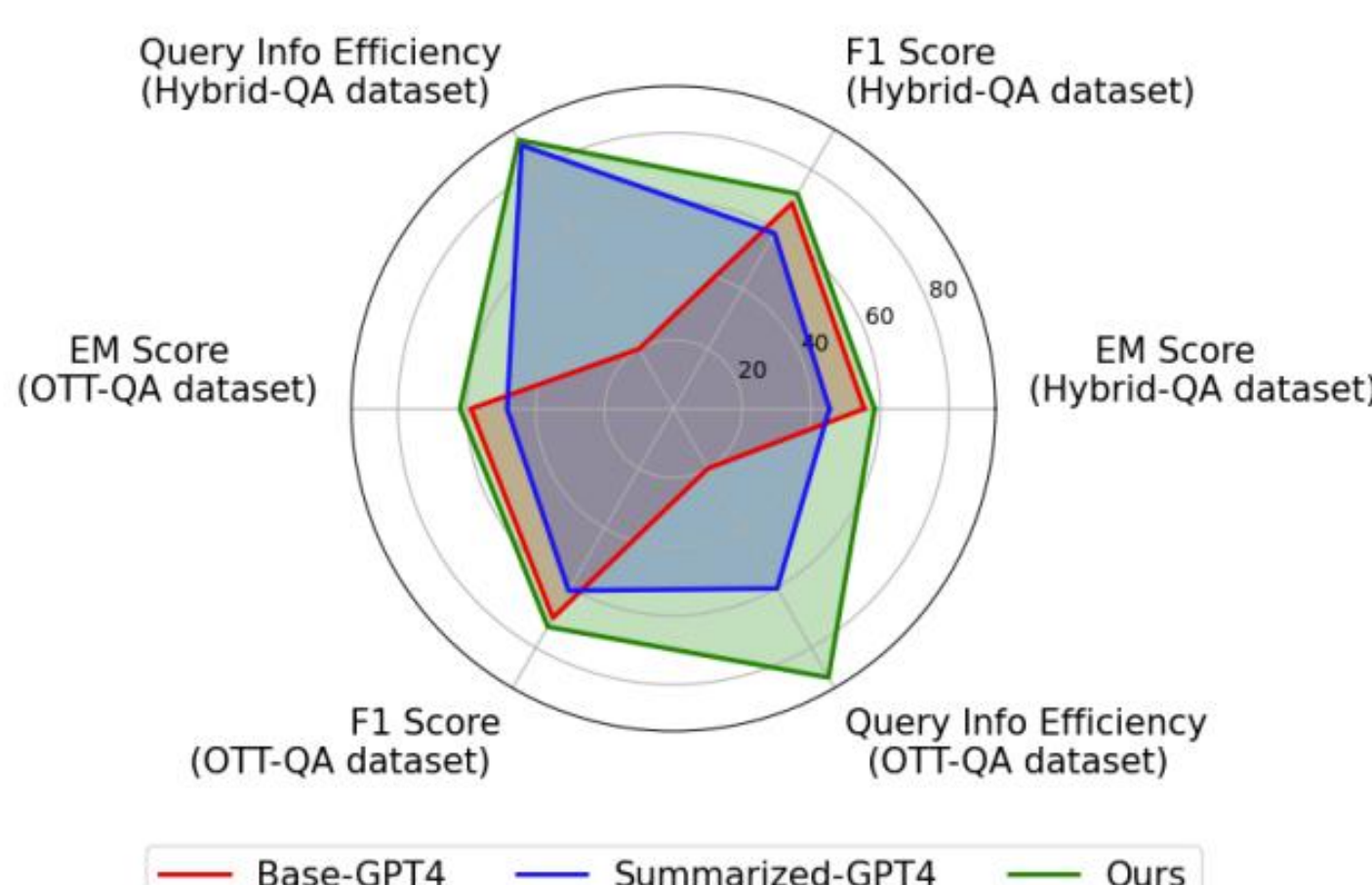


Figure 3: **Multi-Dimensional Improvements:** Our method (with GPT-4 as reader LLM) demonstrates superior results on Hybrid-QA and OTT-QA.

### ODYSSEY Hopwise Analysis

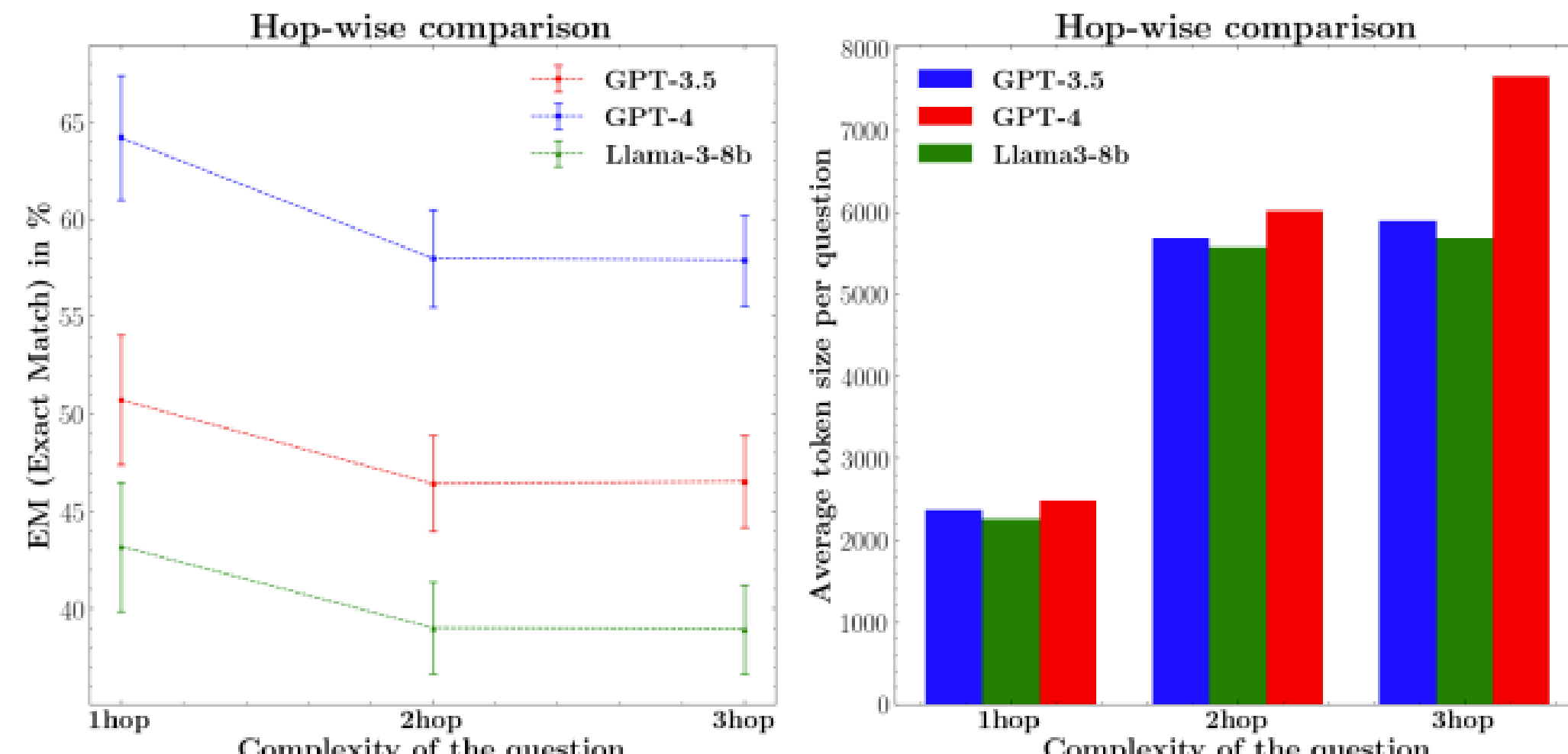


Figure 4: **Hopwise analysis:** For ODYSSEY (our method w/ hopwise), we calculate the cumulative EM score (left-side in figure) and average token size (right-side in figure) utilized after each hop for Llama3-8B, GPT3.5, and GPT-4 on Hybrid-QA.

### Findings:

- Our method effectively reduces and prunes the input token size for LLMs, enhancing efficiency.
- The reduction in token size directly correlates with a decrease in computational cost.