

# SuRe: Summarizing Retrievals using Answer Candidates for Open-domain QA of LLMs

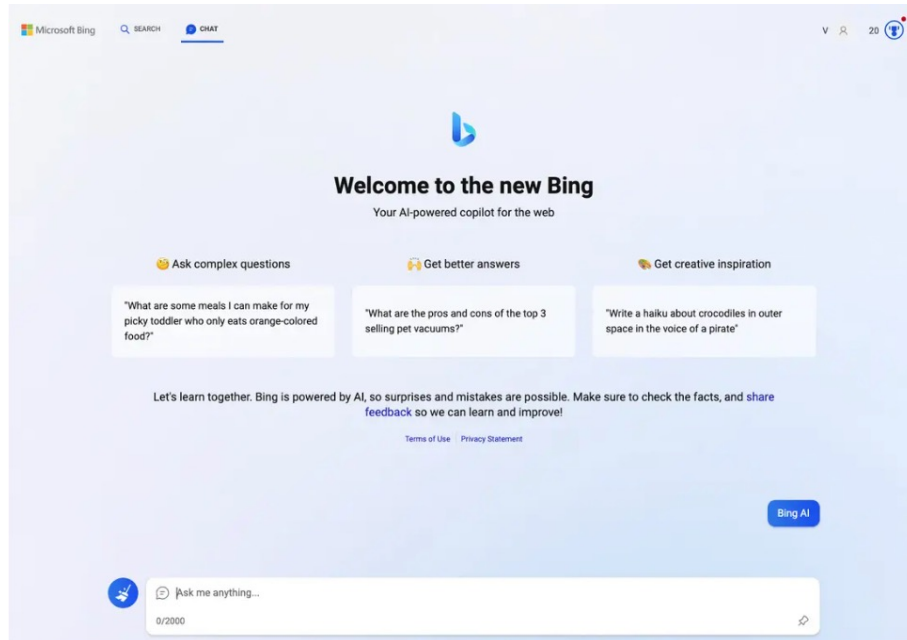
**Jaehyung Kim**<sup>1</sup>, Jaehyun Nam<sup>2</sup>, Sangwoo Mo<sup>3</sup>, Jongjin Park<sup>2</sup>,  
Sang-Woo Lee<sup>2,5</sup>, Minjoon Seo<sup>2</sup>, Jung-Woo Ha<sup>4,5</sup>, and Jinwoo Shin<sup>2</sup>

*ICLR 2024*

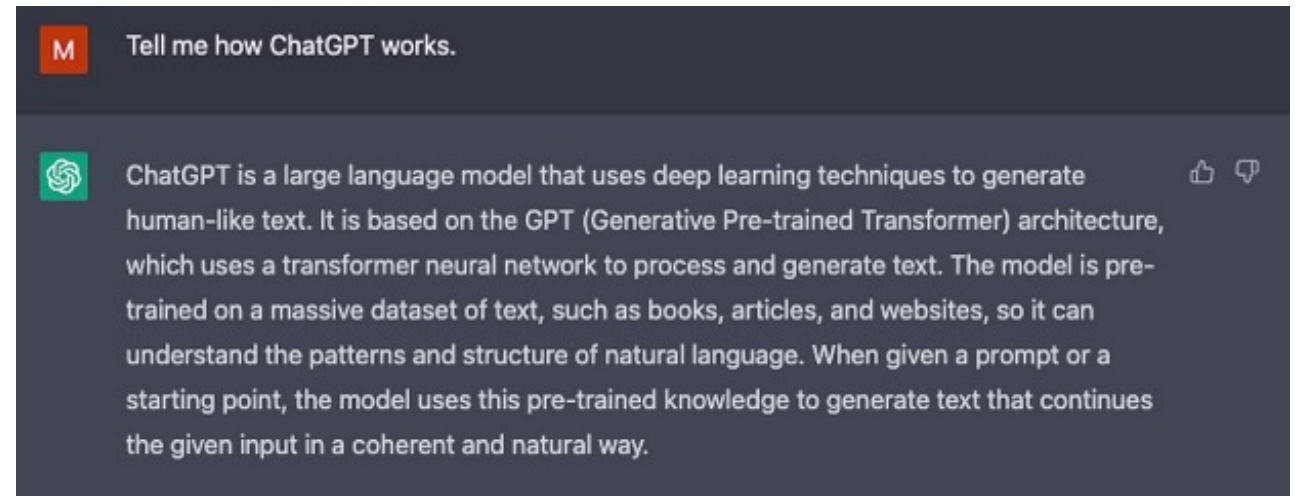
<sup>1</sup> CMU, <sup>2</sup> KAIST, <sup>3</sup> University of Michigan, <sup>4</sup> Naver AI Lab, <sup>5</sup> Naver Cloud

# Success of Large Language Models (LLMs)

- LLMs become **core technology** in many real-world application



*Search engine (Bing AI)<sup>[1]</sup>*



*Chatbot (ChatGPT)<sup>[2]</sup>*

[1] Quy et al., Evaluation of Chatgpt and Microsoft Bing ai Chat Performances on Physics Exams of Vietnamese National High School Graduation Examination., arXiv:2306

[2] OpenAI., <https://openai.com/blog/chatgpt>., 2022

# Challenges of Large Language Models

- Despite its great success, LLMs still suffer from **critical limitations**
  - Non-factual but seemingly plausible generation, *i.e.*, **hallucinations**<sup>[1]</sup>
  - **Difficulty** in integrating **up-to-date knowledge**

PA

summarise this article <https://www.nytimes.com/2023/03/11/technology/chatgpt-prompts-to-avoid-content-filters.html>



The article discusses how ChatGPT, an AI language model, is being used to circumvent content filters that are designed to detect and block certain types of content. By generating prompts that are less likely to trigger these filters, ChatGPT can be used to create content that might otherwise be censored. While this technology could be useful for those seeking to express themselves freely, it also raises concerns about the potential misuse of AI-generated content, as well as the difficulty of regulating such content in the future. The article concludes by noting that there is still much to learn about the impact of AI language models like ChatGPT on the way we communicate online.

*Hallucination of LLMs<sup>[2]</sup>*



Who will be the president of Korea in 2022?



As of my last update in January 2022, the South Korean presidential election was scheduled to take place on March 9, 2022. Several candidates were in the running at that time, but the election results were not available. You'll need to consult a current or post-March 2022 news source to find out who became the president of South Korea in 2022.

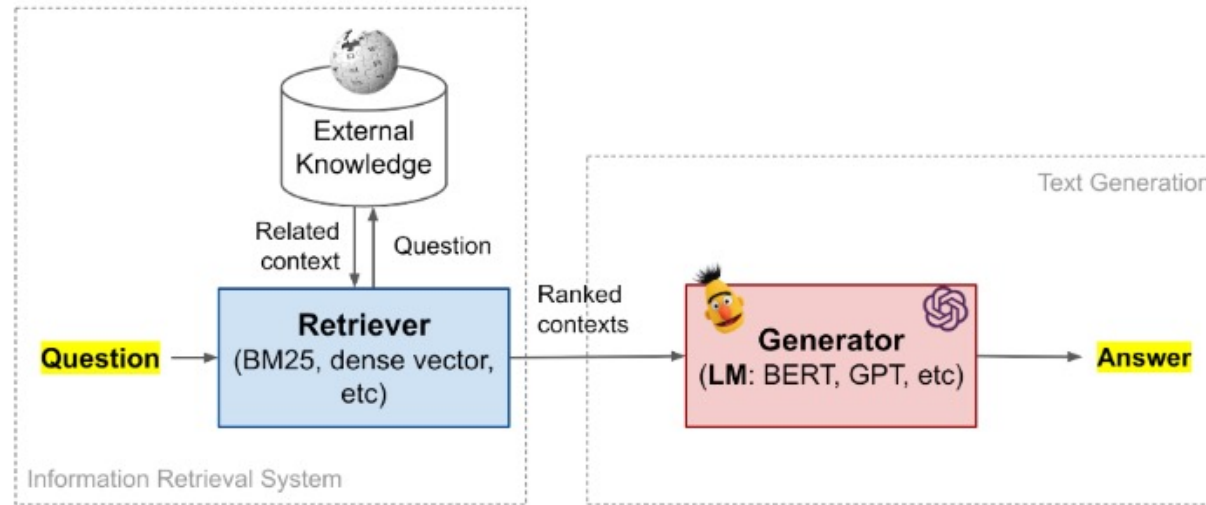
*Struggle to incorporate up-to-date knowledge*

[1] Welleck et al., Neural Text Generation with Unlikelihood Training., ICLR 2020

[2] [https://en.wikipedia.org/wiki/Hallucination\\_%28artificial\\_intelligence%29](https://en.wikipedia.org/wiki/Hallucination_%28artificial_intelligence%29)

# Retrieval Augmentation

- **Retrieving** and incorporating relevant knowledge → promising solution
  - *Retrieval* is effective way to improve QA system, *i.e.*, open-domain QA (ODQA)<sup>[1,2]</sup>



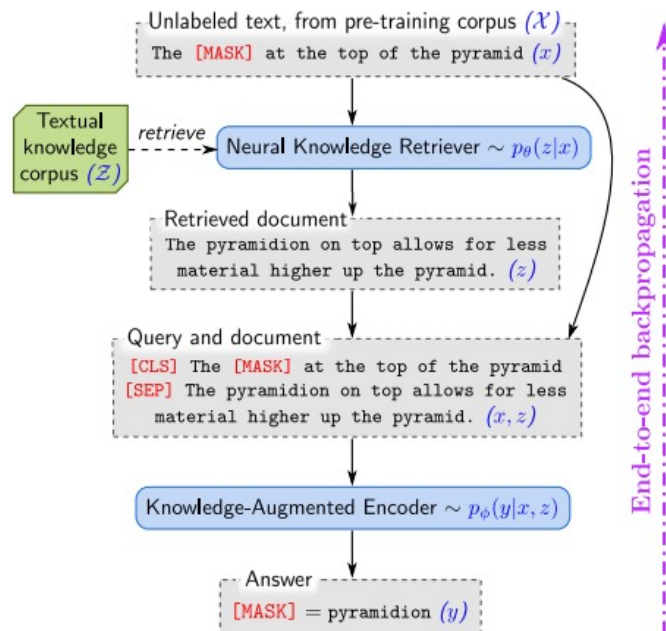
*Illustration of retrieve-and-read system for ODQA<sup>[1]</sup>*

[1] <https://lilianweng.github.io/posts/2020-10-29-odqa/>

[2] Karpukhin et al., Dense Passage Retrieval for Open-Domain Question Answering., EMNLP 2020

# Retrieval Augmentation

- **Retrieving** and incorporating relevant knowledge → promising solution
  - *Retrieval* is effective way to improve QA system, *i.e.*, open-domain QA (ODQA)<sup>[1,2]</sup>
  - *Retrieval-augmented LLMs* outperforms ordinary LLMs<sup>[3,4]</sup>



Name	NQ (79k/4k)	WQ (3k/2k)	CT (1k /1k)	# params
BERT-Baseline (Lee et al., 2019)	26.5	17.7	21.3	110m
T5 (base) (Roberts et al., 2020)	27.0	29.1	-	223m
T5 (large) (Roberts et al., 2020)	29.8	32.2	-	738m
T5 (11b) (Roberts et al., 2020)	34.5	37.4	-	11318m
DrQA (Chen et al., 2017)	-	20.7	25.7	34m
HardEM (Min et al., 2019a)	28.1	-	-	110m
GraphRetriever (Min et al., 2019b)	31.8	31.6	-	110m
PathRetriever (Asai et al., 2019)	32.6	-	-	110m
ORQA (Lee et al., 2019)	33.3	36.4	30.1	330m
Ours ( $\mathcal{X}$ = Wikipedia, $\mathcal{Z}$ = Wikipedia)	39.2	40.2	<b>46.8</b>	330m
Ours ( $\mathcal{X}$ = CC-News, $\mathcal{Z}$ = Wikipedia)	<b>40.4</b>	<b>40.7</b>	42.9	330m

*Illustration of REtrieval-Augmented Language Model (REALM)<sup>[3]</sup>*

[1] <https://lilianweng.github.io/posts/2020-10-29-odqa/>

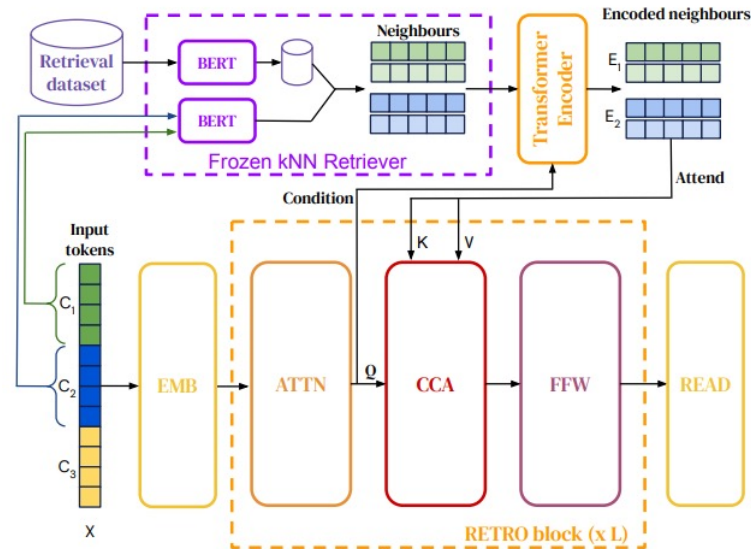
[2] Karpukhin et al., Dense Passage Retrieval for Open-Domain Question Answering., EMNLP 2020

[3] Guu et al., REALM: Retrieval-Augmented Language Model Pre-Training., ICML 2020

[4] Borgeaud et al., Improving Language Models by Retrieving from Trillions of Tokens., ICML 2022

# Challenges with Retrieval Augmented LLMs

- Retrieval-augmented LLMs are typically constructed with **finetuning**<sup>[1,2]</sup>
  - But, recent **large scale** & **black-box nature** make it less attractive & infeasible



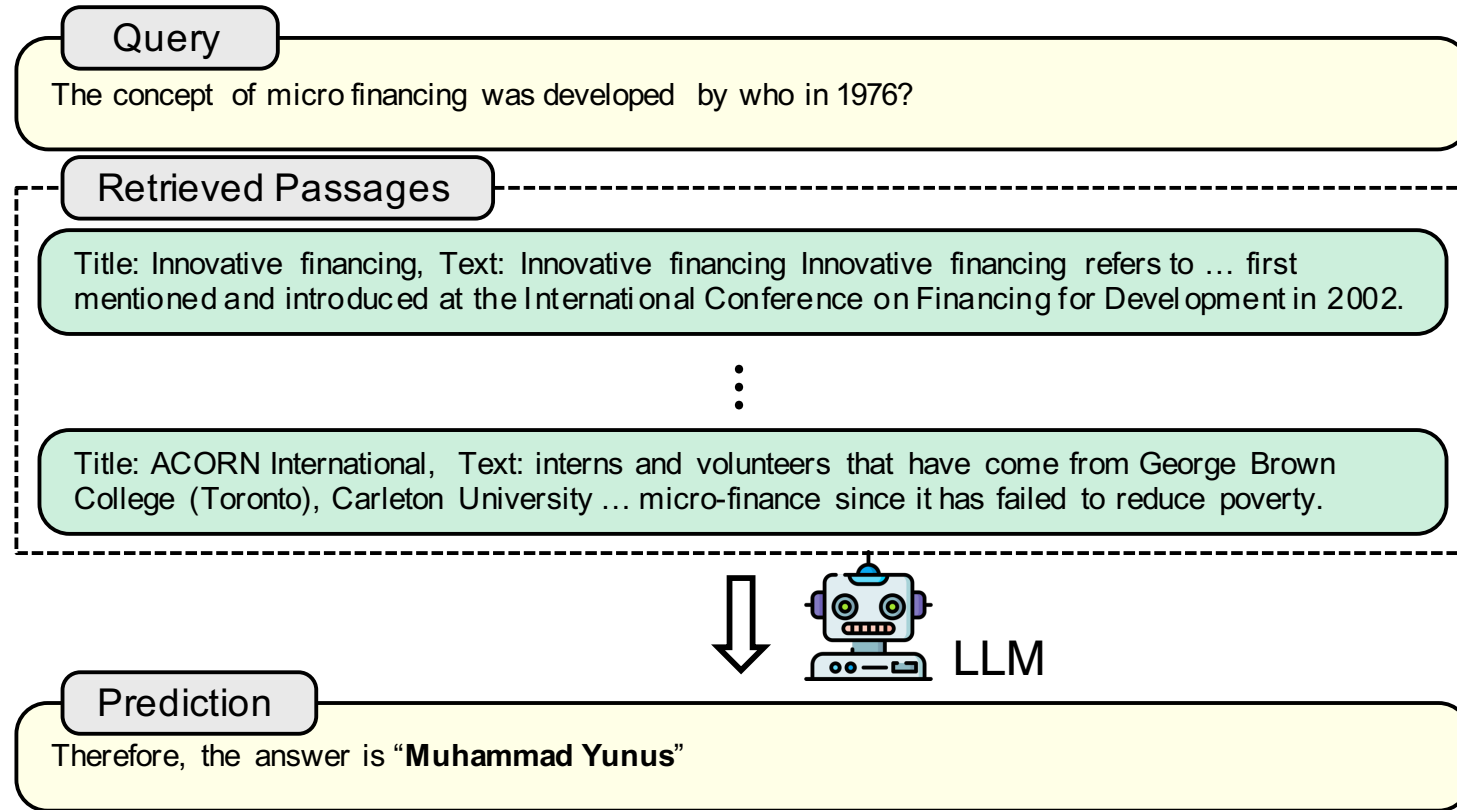
*Modification of architecture and training to learn retrieval augmentation<sup>[2]</sup>*

[1] Guu et al., REALM: Retrieval-Augmented Language Model Pre-Training., ICML 2020

[2] Borgeaud et al., Improving Language Models by Retrieving from Trillions of Tokens., ICML 2022

# Challenges with Retrieval Augmented LLMs

- **Prompting** is simple and efficient way to augment retrieved passages<sup>[1]</sup>

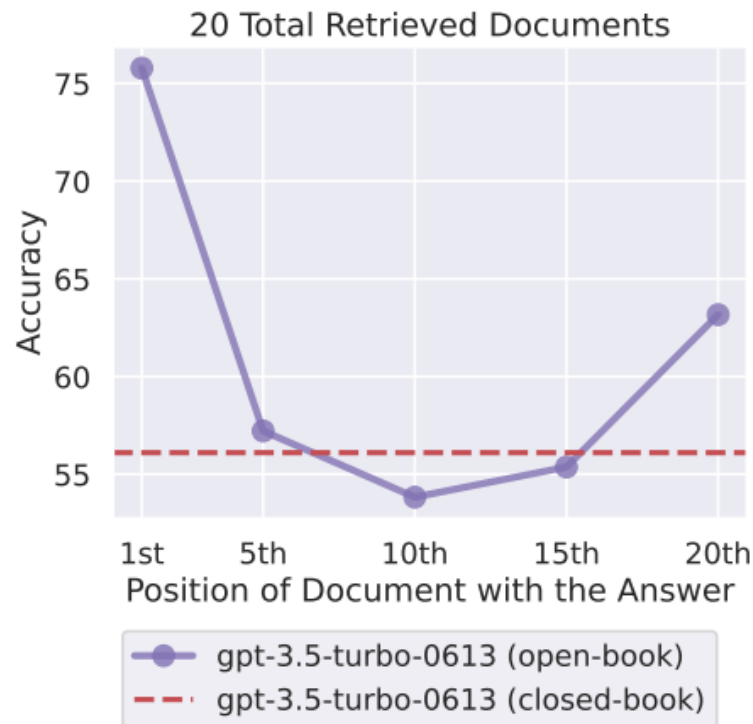


[1] Shi et al., REPLUG: Retrieval-Augmented Black-Box Language Models., arXiv:2301

[2] Liu et al., Lost in the Middle: How Language Models Use Long Contexts., arXiv:2307

# Challenges with Retrieval Augmented LLMs

- **Prompting** is simple and efficient way to augment retrieved passages<sup>[1]</sup>
  - Naïve approach (e.g., appending) could be **limited in fully exploiting the retrieval**<sup>[2]</sup>



*LLMs struggle to handle long retrieved passages<sup>[2]</sup>*

[1] Shi et al., REPLUG: Retrieval-Augmented Black-Box Language Models., arXiv:2301

[2] Liu et al., Lost in the Middle: How Language Models Use Long Contexts., arXiv:2307



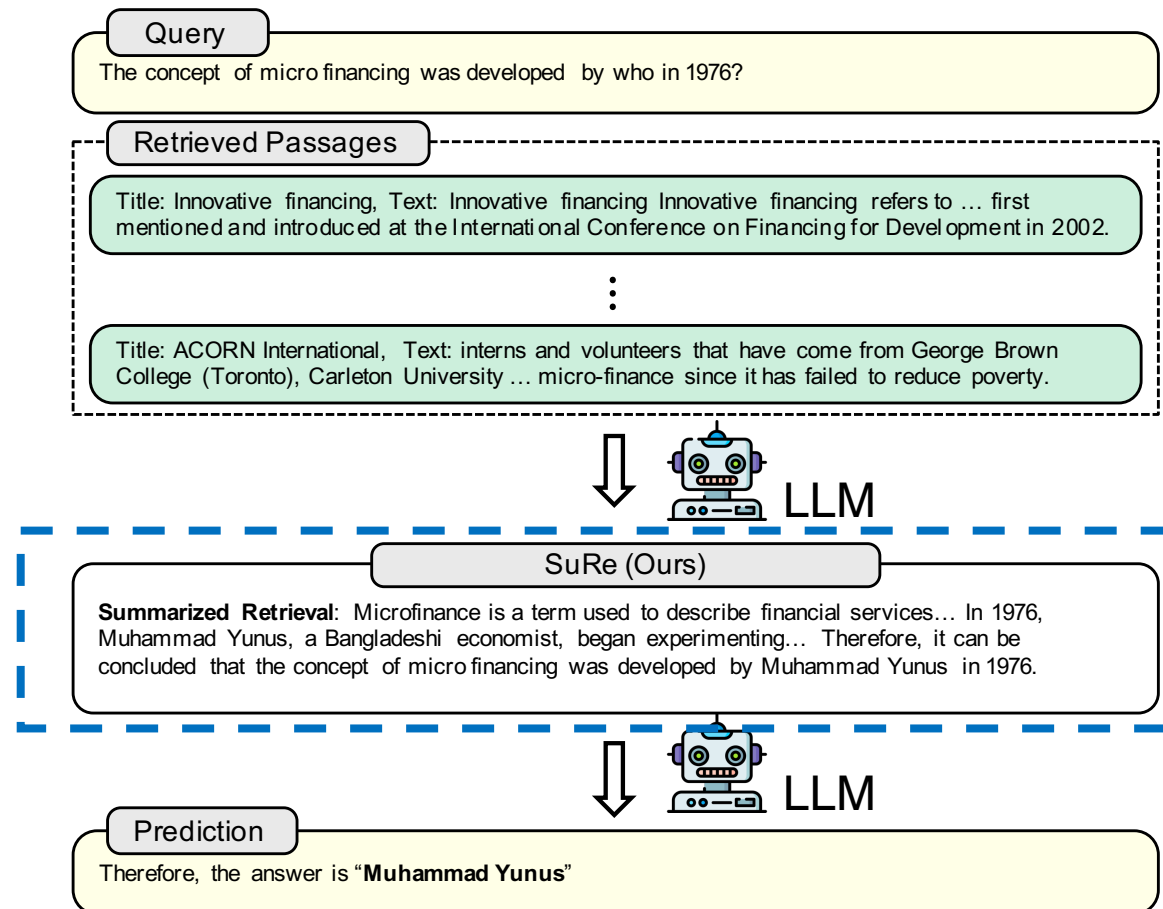
# Research Goal

- **Prompting** is simple and efficient way to augment retrieved passages<sup>[1]</sup>
  - Naïve approach (e.g., appending) could be limited in fully exploiting the retrieval<sup>[2]</sup>

**Research Goal:** a simple yet effective framework based on prompting to improve ODQA with LLMs

# Proposed Method: SuRe

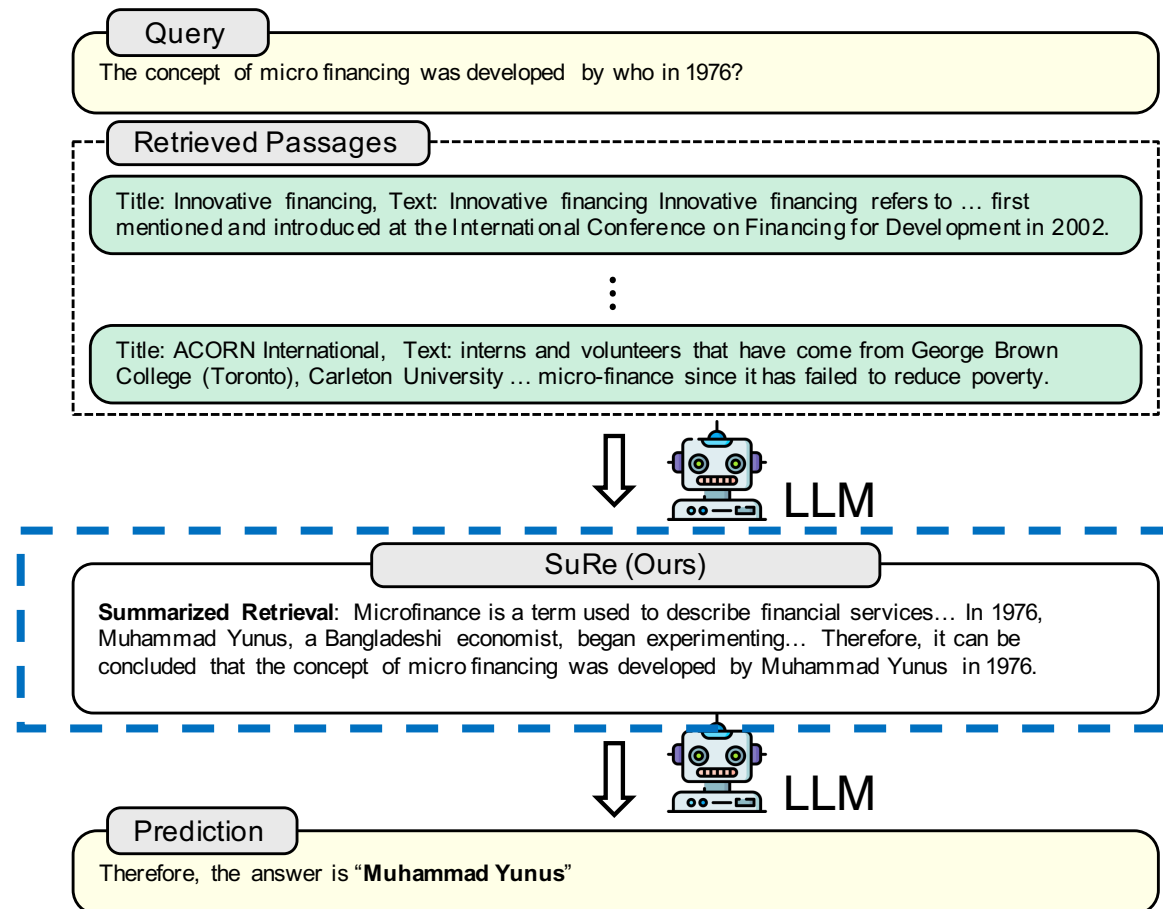
- **SuRe**: Summarizing Retrievals using Answer Candidates



*Example of Retrieval Augmentation of LLMs via Proposed SuRe Framework*

# Proposed Method: SuRe

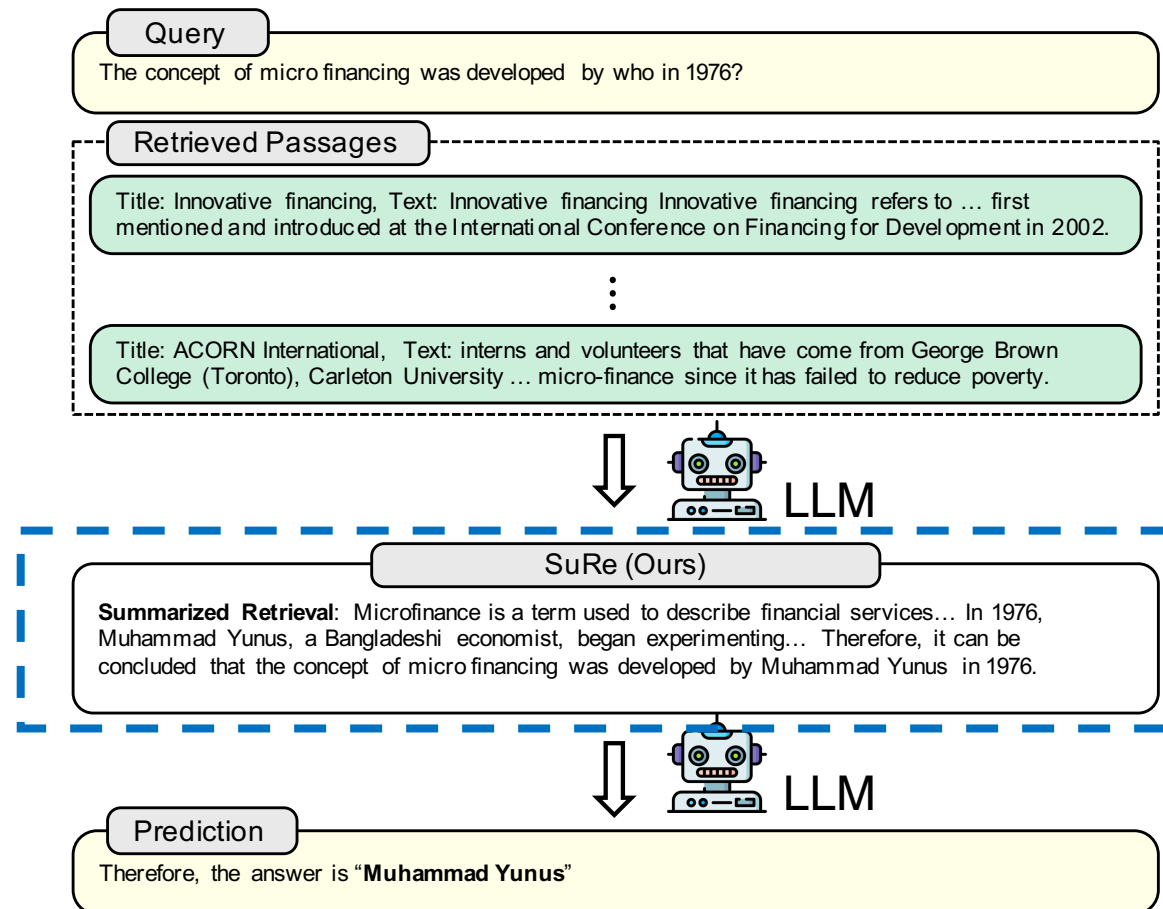
- **SuRe**: Summarizing Retrievals using Answer Candidates
  - Key idea. Select answer that is well-supported by [summarization](#) of retrievals (explicit rationale)



*Example of Retrieval Augmentation of LLMs via Proposed SuRe Framework*

# Proposed Method: SuRe

- **SuRe**: Summarizing Retrievals using Answer Candidates
  - Key idea. Select answer that is well-supported by [summarization](#) of retrievals



*Example of Retrieval Augmentation of LLMs via Proposed SuRe Framework*

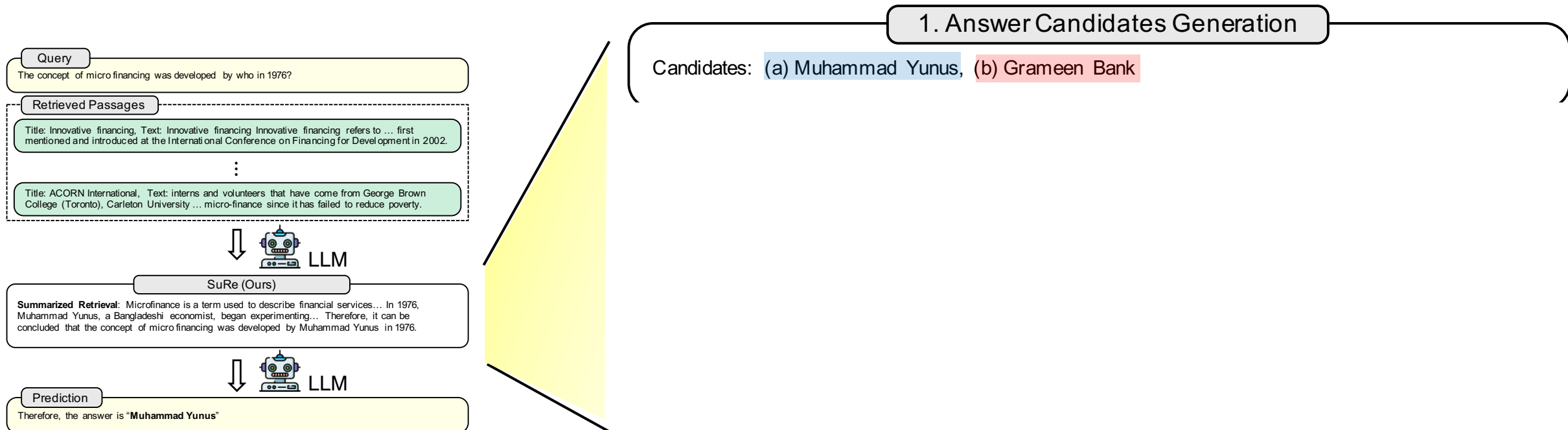
# Proposed Method: SuRe

---

- SuRe conducts **three steps** sequentially (via zero-shot prompting)

# Proposed Method: SuRe

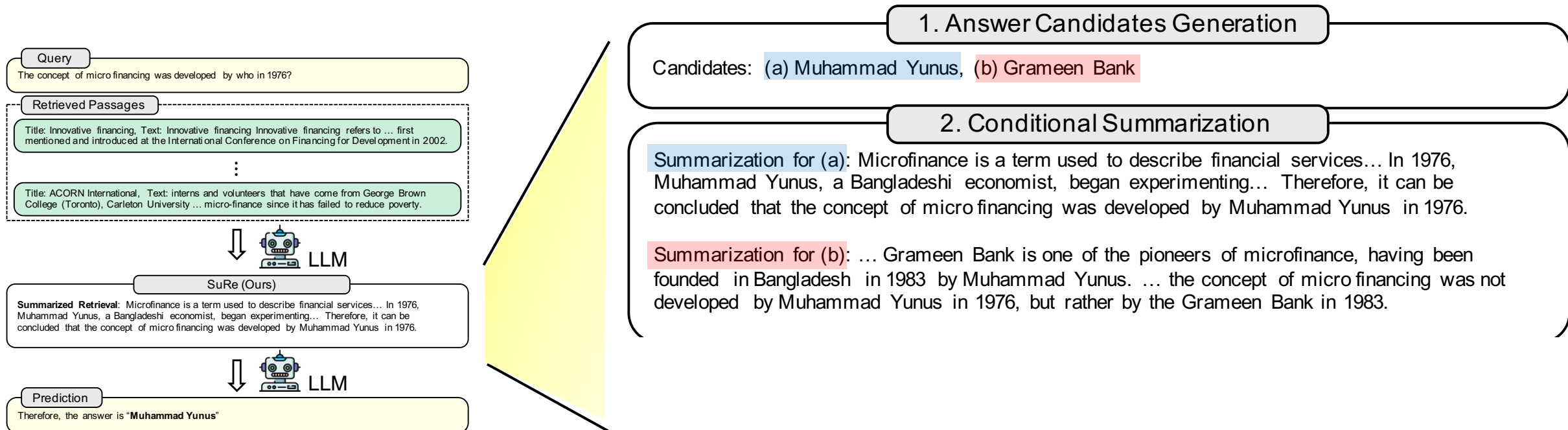
- SuRe conducts **three steps** sequentially (via zero-shot prompting)
  1. Answer candidates generation



*Overall illustration of SuRe Framework*

# Proposed Method: SuRe

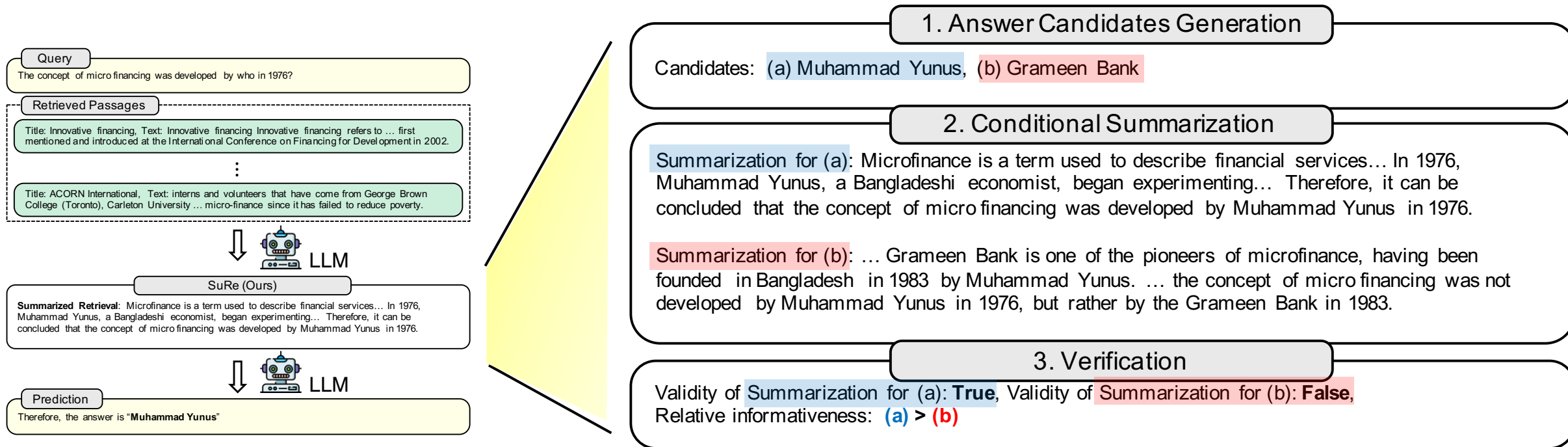
- SuRe conducts **three steps** sequentially (via zero-shot prompting)
  1. Answer candidates generation
  2. Conditional summarization



*Overall illustration of SuRe Framework*

# Proposed Method: SuRe

- SuRe conducts **three steps** sequentially (via zero-shot prompting)
  1. Answer candidates generation
  2. Conditional summarization
  3. Selection via verification (validity & ranking)

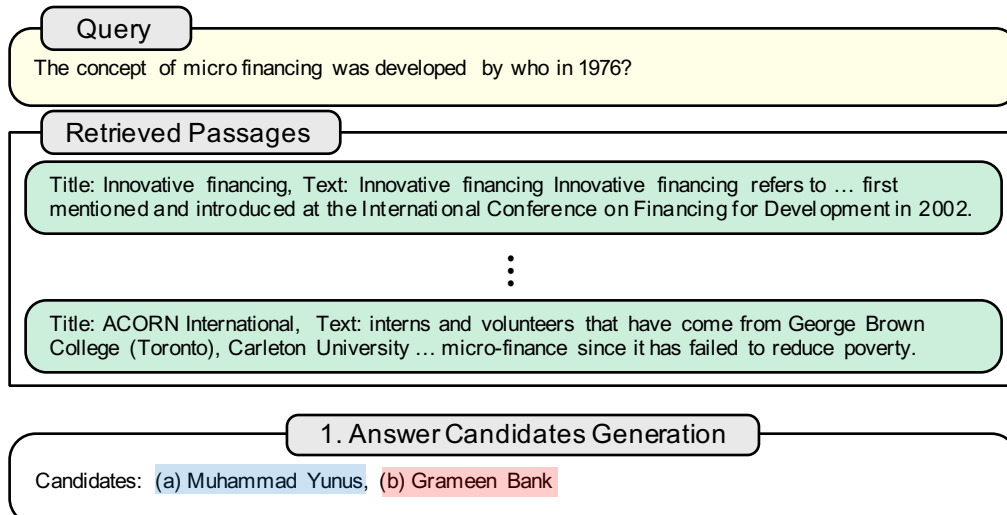


*Overall illustration of SuRe Framework*



# SuRe: Answer Candidates Generation

- SuRe first generate K answer candidates from question and retrieval
  - Compare to random sampling, we observe *direct prompting is more effective*
  - We use fixed K=2 during experiments, as  $K > 2$  is **quite inefficient**



**Listing 1** Prompt for answer candidates generation.

```
f'''
Below are N passages related to the question at the end. After reading
↳ the passages, provide two correct candidates for the answer to the
↳ question at the end. Each answer should be in the form: (a) xx, (b)
↳ yy, and should not exceed 3 words for each candidate.

Passage #1 Title: {Passage #1 Title}
Passage #1 Text: {Passage #1 Text}

...

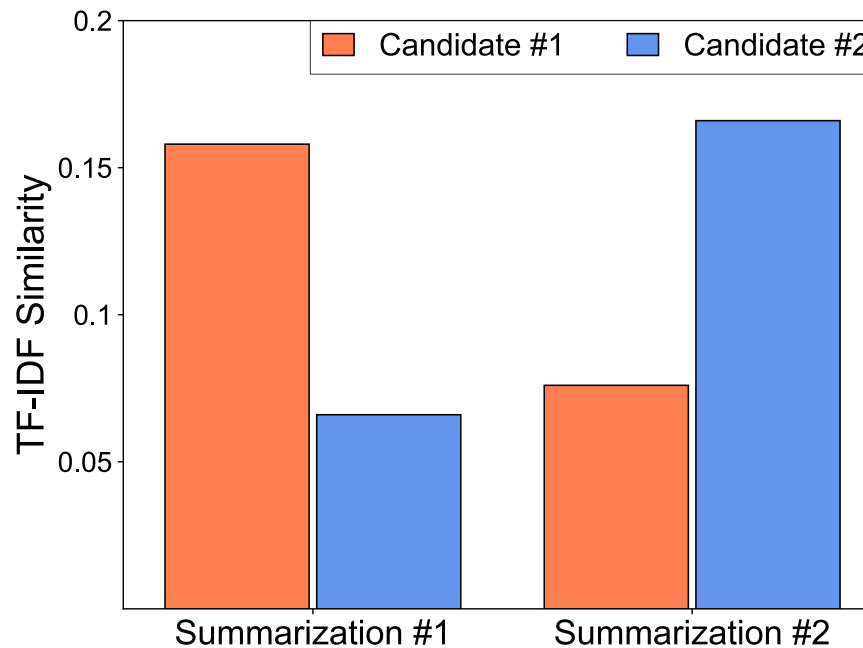
Passage #N Title: {Passage #N Title}
Passage #N Text: {Passage #N Text}

Question: {Question}

Answer:
'''
```

# SuRe: Conditional Summarization

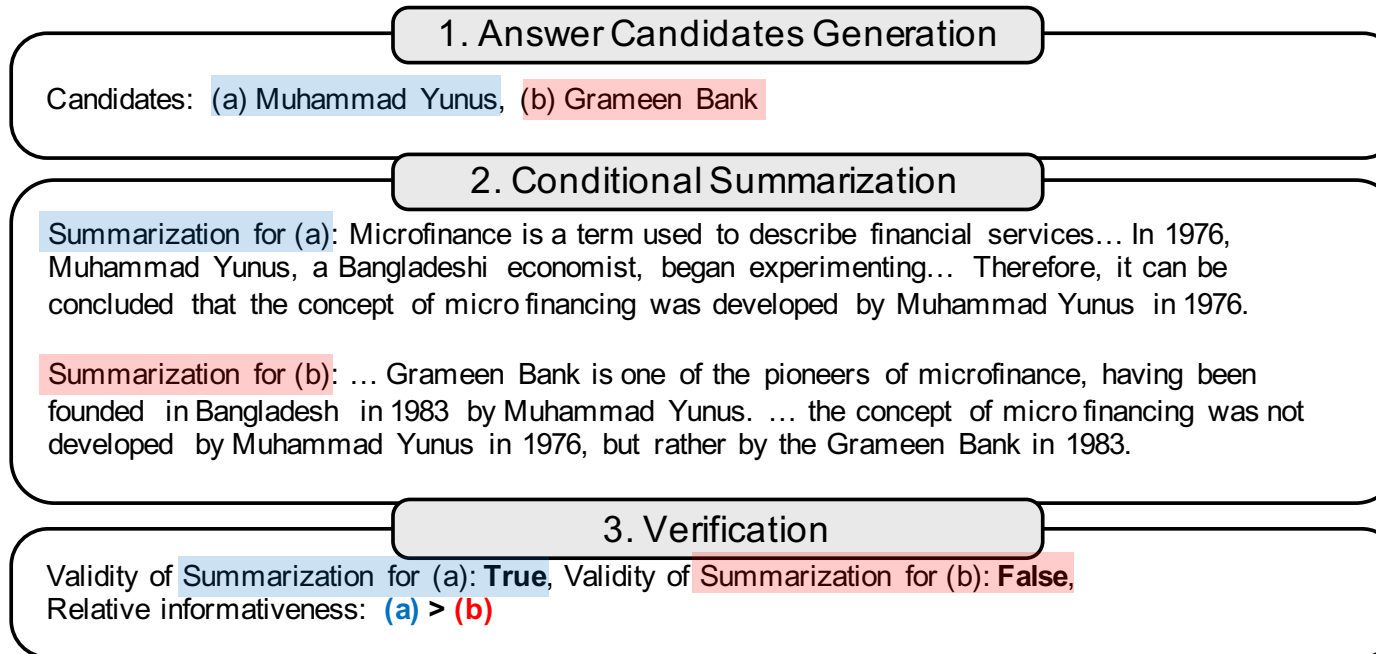
- SuRe then summarize retrieval conditioned on each candidate
  - *i.e.*, question, retrieved passages, and **answer candidate** → summarization
  - Conditional summarization includes *specific contexts* supporting given candidate



*TF-IDF similarity between candidates and conditional summarizations*

# SuRe: Selection via Verification

- Lastly, SuRe verify summarizations and select most plausible candidate



# SuRe: Selection via Verification

- Lastly, SuRe verify summarizations and select most plausible candidate
  - Two scores: (1) Validity of summarization

---

**Listing 3** Prompt for instance-wise validation.

---

```
f'''
Question: {Question}

Prediction: {Prediction}

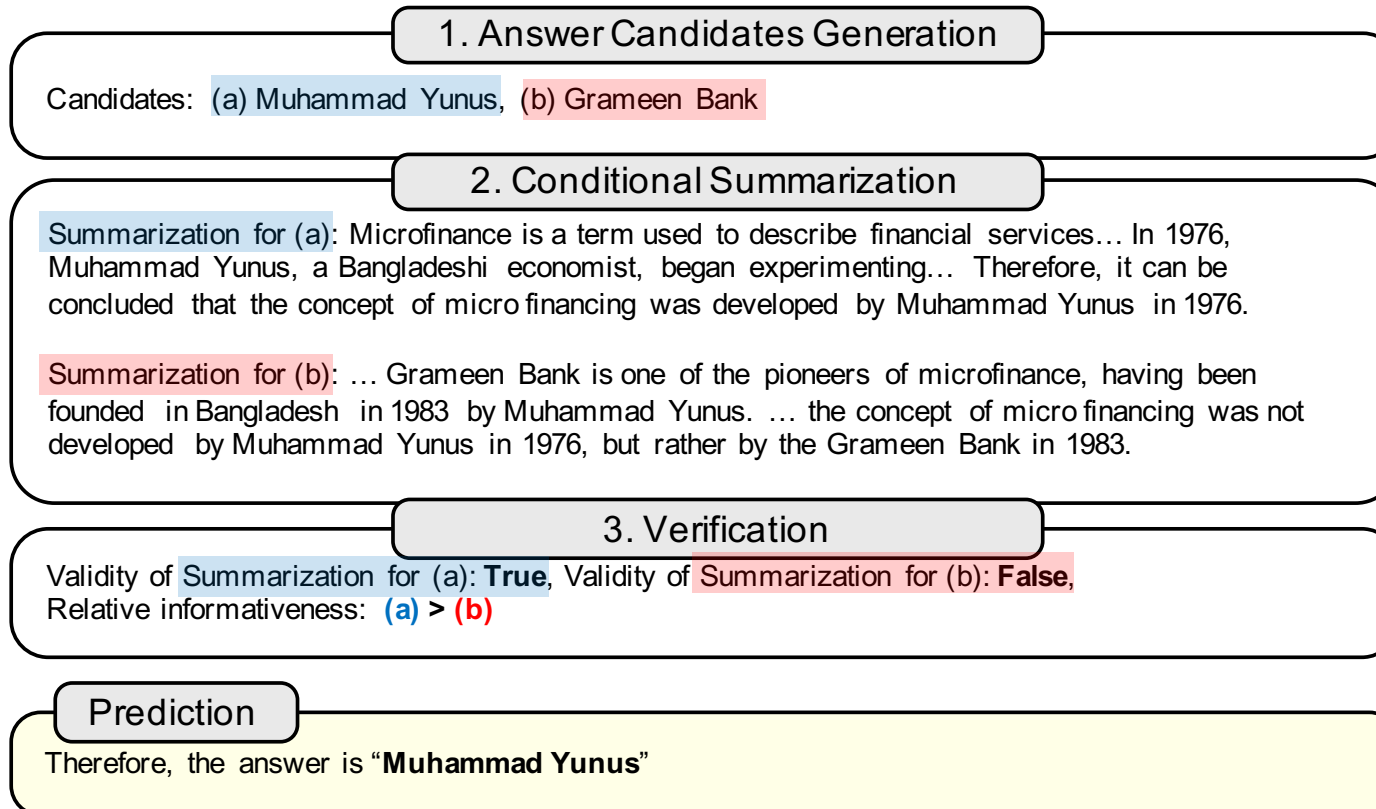
Passage: {Passage}

Does the passage correctly support the prediction? Choices: [True,
↪ False]. Answer:
'''
```

---

# SuRe: Selection via Verification

- Lastly, SuRe verify summarizations and select most plausible candidate
  - Two scores: (1) Validity of summarization & (2) Ranking between summarizations
  - Candidate with highest scored summarization → **Answer**



# Experiments: Comparison with Baselines

- **ChatGPT** with 10 retrieved passages from **BM25**
  - Augmenting retrieved passages with **prompting (Base)** is effective
  - **Other baselines are ineffective** due to challenging setup (zero-shot, black-box API)
  - **SuRe significantly improve** *Base* with a large margin (4.6% in EM & 4.0% in F1)

Methods / Datasets	NQ	WebQ	2Wiki	HotpotQA	Average
No retrieval	27.6 / <u>39.0</u>	<u>25.0</u> / <b>38.8</b>	21.4 / 24.8	22.2 / 31.9	24.1 / 33.6
Base	<u>28.4</u> / 38.8	19.6 / 32.5	<u>27.4</u> / <u>32.8</u>	30.8 / 40.3	<u>26.6</u> / <u>36.1</u>
Rerank	24.8 / 33.9	18.8 / 30.6	23.0 / 28.4	27.8 / 37.4	23.6 / 32.6
RePlug	26.0 / 35.3	18.8 / 31.5	23.6 / 28.5	28.0 / 37.9	24.1 / 33.3
Selection-inference	24.3 / 32.8	17.3 / 28.6	22.6 / 29.5	30.8 / 39.6	23.8 / 32.6
Chain-of-thoughts	22.3 / 31.4	15.2 / 27.8	19.6 / 22.5	25.6 / 31.8	20.7 / 28.4
Self-verification	25.2 / 35.4	16.1 / 28.5	23.2 / 30.5	<u>31.6</u> / <u>41.8</u>	24.0 / 34.1
SURE (Ours)	<b>33.5</b> / <b>42.3</b>	<b>25.1</b> / <u>36.6</u>	<b>32.8</b> / <b>38.1</b>	<b>33.2</b> / <b>43.4</b>	<b>31.2</b> / <b>40.1</b>

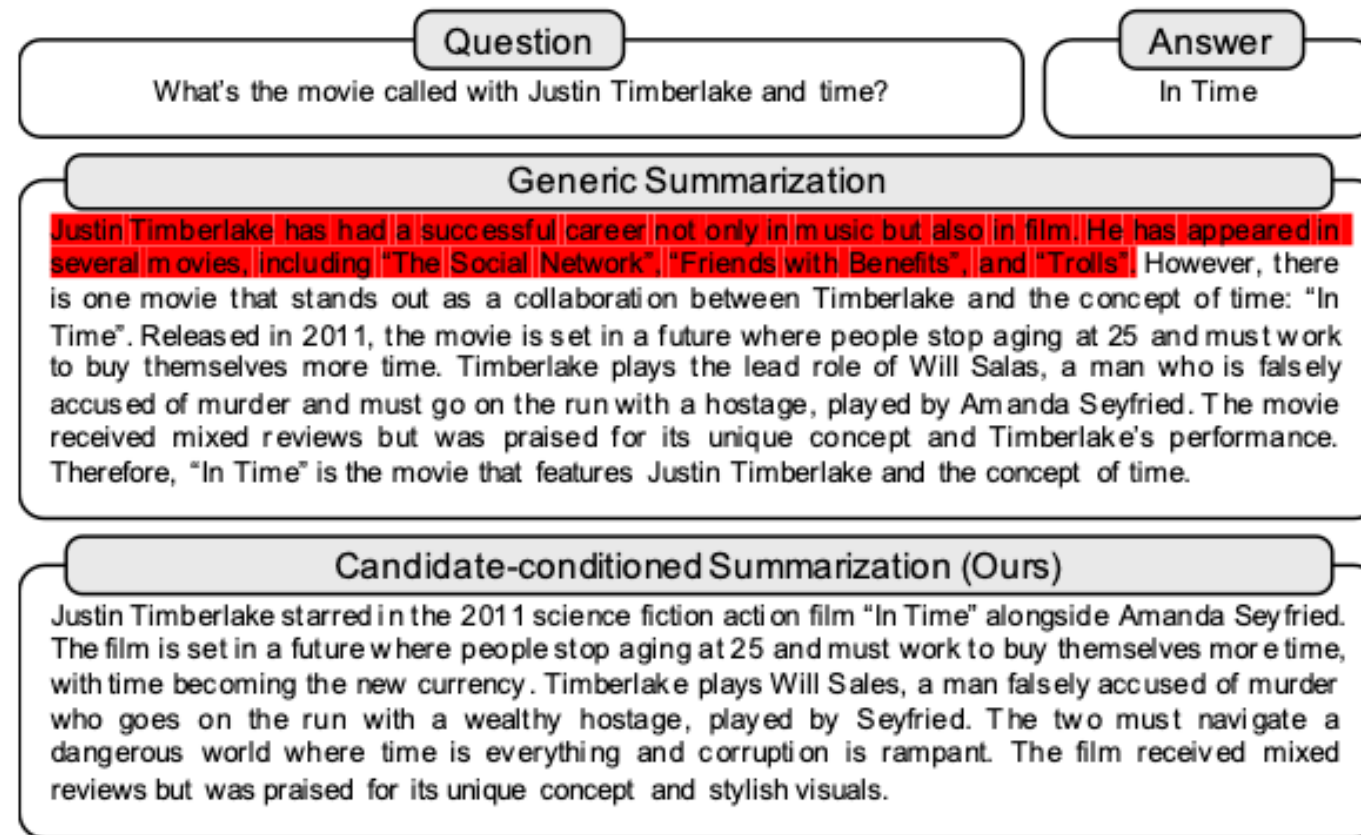
# Experiments: Compatibility of SuRe

- Evaluation on **different LLMs and retrieval methods**
  - Naïve prompting (Base) is compared as baseline
  - SuRe is **consistently effective** for all setups with 4.6% EM improvement in average
  - SuRe successfully improves **LLaMA2-chat** → Usefulness for open sourced LLMs!

Datasets	ChatGPT						GPT-4		LLaMA2-chat	
	BM25	+ SuRE	DPR	+ SuRE	Contriever	+ SuRE	BM25	+ SuRE	BM25	+ SuRE
NQ	28.4	<b>33.5</b>	36.1	<b>41.0</b>	35.8	<b>40.4</b>	30.2	<b>32.4</b>	18.6*	<b>30.4*</b>
WebQ	19.6	<b>25.1</b>	23.2	<b>27.3</b>	22.5	<b>28.7</b>	21.5	<b>21.7</b>	16.0*	<b>24.0*</b>
2Wiki	27.4	<b>32.8</b>	19.2	<b>21.4</b>	27.2	<b>32.6</b>	34.8	<b>38.2</b>	20.2	<b>27.8</b>
HotpotQA	30.8	<b>33.2</b>	25.6	<b>27.4</b>	32.2	<b>33.6</b>	34.8	<b>40.6</b>	24.0	<b>28.0</b>
Average	26.6	<b>31.2</b>	26.0	<b>29.3</b>	29.4	<b>33.8</b>	30.3	<b>33.2</b>	19.7	<b>27.6</b>

# Experiments: Additional Advantages of SuRe

- Not only answer, but SuRe also generates **conditional summarization**
  - Then, is this conditional summarization really a *good rationale for answer*?



*Qualitative comparison between SuRe's summarization and generic (i.e., unconditional) summarization*



# Experiments: GPT-4 & Human Evaluations

- We ask GPT-4 to select more informative & plausible sum. as rationale<sup>[1]</sup>

```
f'''
Question: Given the following summaries for the target question,
↳ determine which one is more informative and plausible as rationale
↳ to support a given target question-answer pair.

Summary 1: {Summary 1}

Summary 2: {Summary 2}

Target Question: {Question}

Target Answer: {Answer}

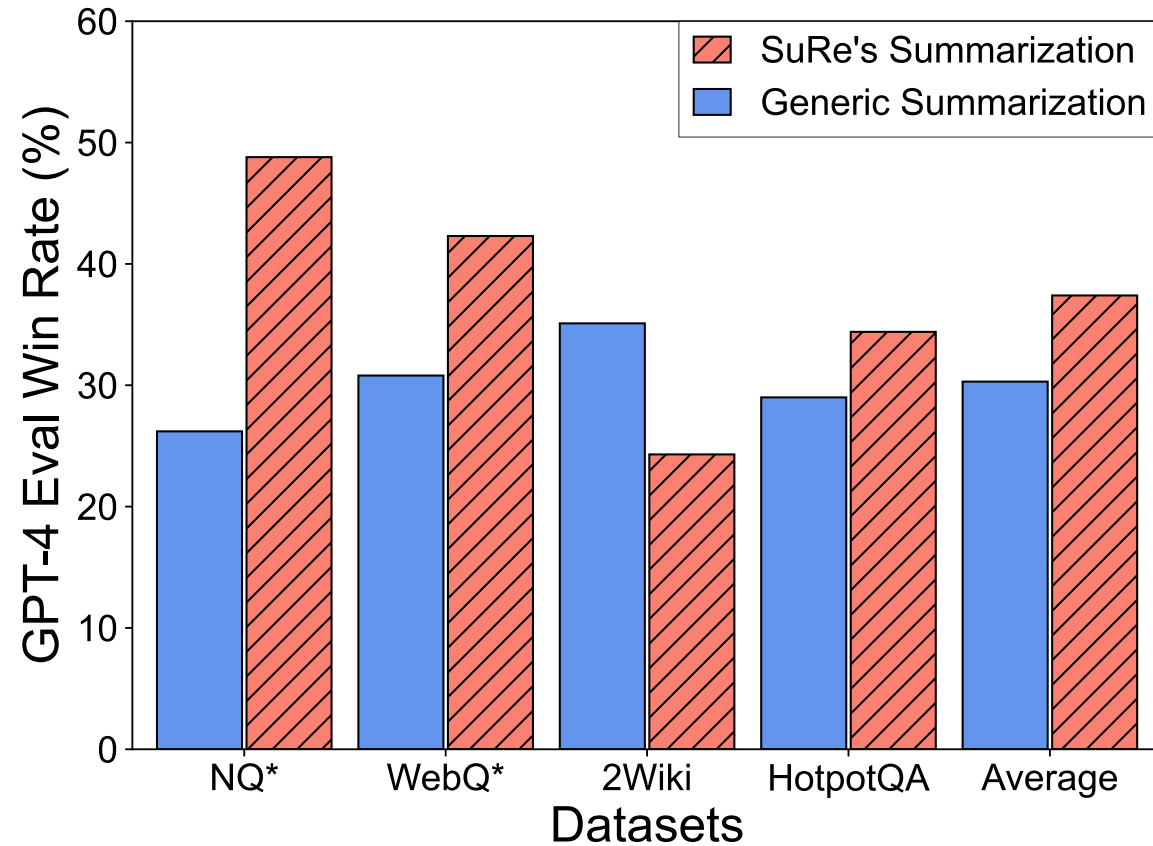
Your Task:
Identify which summary (Summary 1 or Summary 2) is more informative and
↳ plausible as rationale to support a given answer at hand. Choices:
↳ [Summary 1, Summary 2].

Answer:
'''
```

*Designed Prompt for GPT-4 Evaluation between Summarizations*

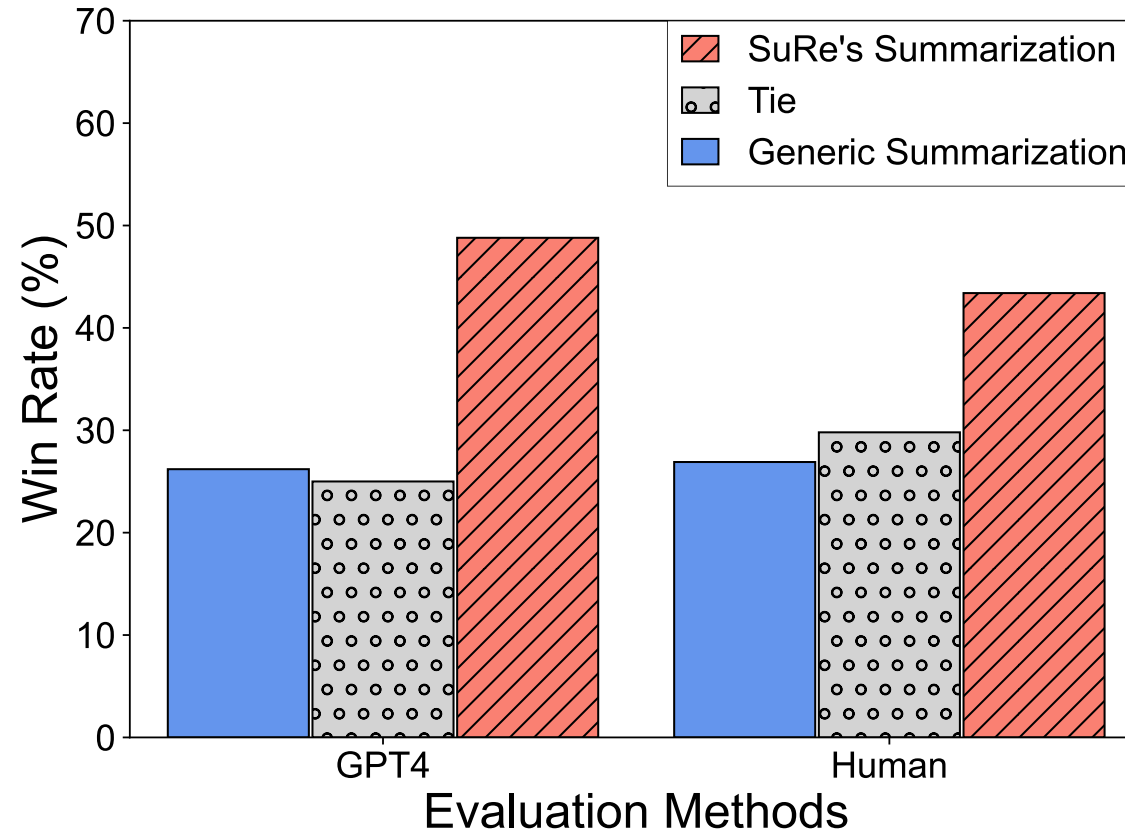
# Experiments: GPT-4 Evaluation

- **GPT-4 prefers SuRe's summarization** than generic summarization
  - Setup. We only use questions that both summarizations correctly predict
  - Result. Average: **37.4%** (SuRe) vs 30.3 % (Generic)



# Experiments: Human Evaluation

- **Human also prefers SuRe's summarization than generic one**
  - Setup. 84 NQ questions that both summarizations correctly predict (7 experts)
  - Result. **43.4%** (SuRe) vs 26.9 % (Generic)



Thank you for attention 🙏