

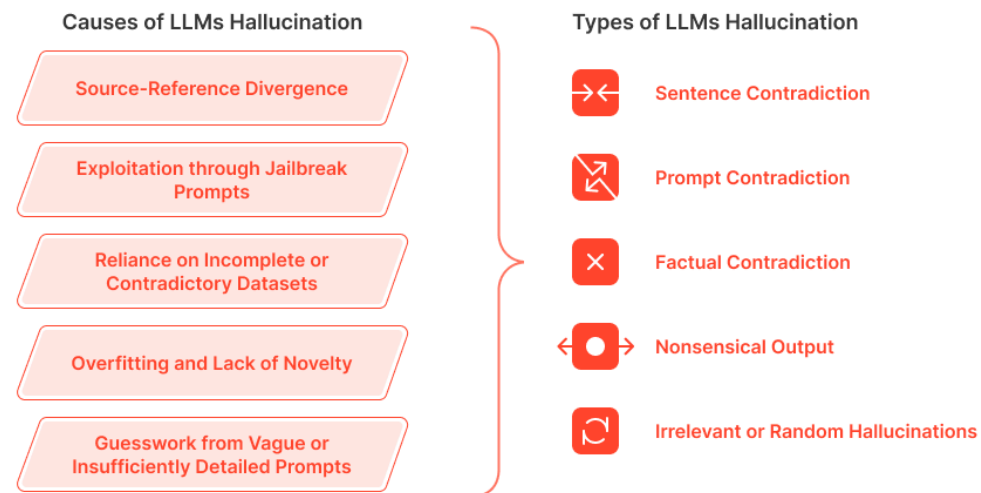
# Retrieval Augmented LM - Attempts to improve retrieval

HCC Lab  
Heo Dong Seok

# Retrieval Augmented Generation

- LM 's big problem : Hallucination / Imaginary content

## Causes and Types of LLMs Hallucination



# Retrieval Augmented Generation

- RAG : Uses external docs/knowledge to generate output

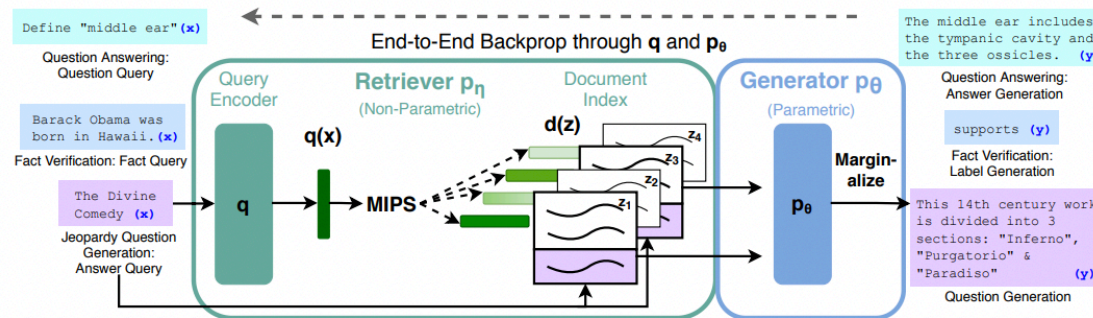


Figure 1: Overview of our approach. We combine a pre-trained retriever (*Query Encoder + Document Index*) with a pre-trained seq2seq model (*Generator*) and fine-tune end-to-end. For query  $x$ , we use Maximum Inner Product Search (MIPS) to find the top-K documents  $z_i$ . For final prediction  $y$ , we treat  $z$  as a latent variable and marginalize over seq2seq predictions given different documents.

## Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks (2020)

<https://arxiv.org/abs/2005.11401>

# Retrieval Augmented Generation

---

- **Step 1 : Retrieval**

- Retrieve top k passages for given query from external knowledge base (Wiki, Bing, etc)
- (ex) MIPS / BM25 / DPR

- **Step 2 : Generation**

- Generate answer output with (query + passages) input
- (ex) LLM (BART / GPT / LoRa ... )

# Retrieval Augmented Generation - Step 1

- BM25 : sparse retriever

$$score(D, Q) = \sum_{i=1}^n IDF(q_i) * \frac{\overbrace{f(q_i, D) * (k_1 + 1)}^{\text{문서 } D \text{에서 } q_i \text{의 term frequency}}}{\underbrace{f(q_i, D) + k_1}_{\text{파라미터}} * \underbrace{(1 - b)}_{\text{문서 집합의 평균 문서 길이}} + b * \frac{|D|}{avgdl}}$$

- DPR : dense retriever

$$p_{\eta}(z|x) \propto \exp(\mathbf{d}(z)^{\top} \mathbf{q}(x)) \quad \mathbf{d}(z) = \text{BERT}_d(z), \quad \mathbf{q}(x) = \text{BERT}_q(x)$$

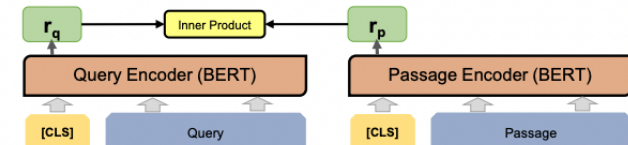


Figure 5: Representation Model for Initial Retrieval

q: query encoder / d: DPR  
 두 임베딩의 곱으로 부터 similarity 계산

---

Re2G (2022)



# Re2G

---

- **R**etrieve, **R**erank, **G**enerate
- RAG + Reranker Network

**Re2G: Retrieve, Rerank, Generate (2022)**

<https://arxiv.org/abs/2207.06300>

# Re2G

---

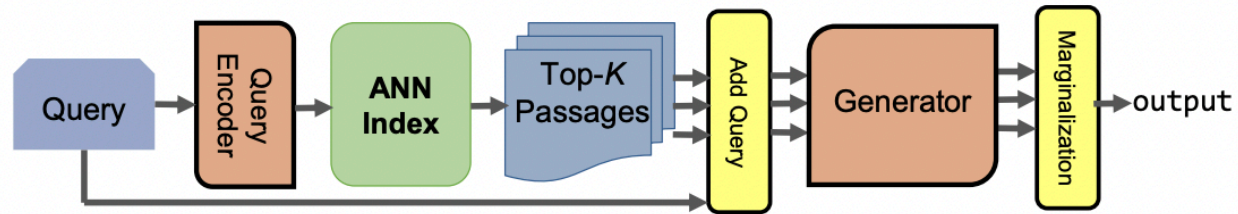


Figure 2: RAG Architecture

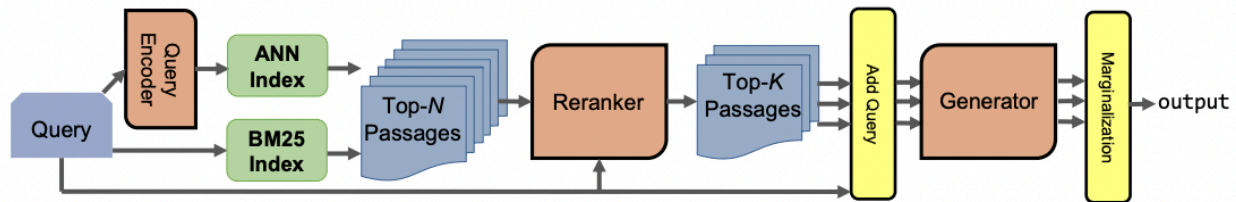


Figure 3: Re<sup>2</sup>G Architecture

Rerank (BM25 passage ; ANN(DPR) passage)

# Retrieval Augmented Generation - Step 1

- BM25 : sparse retriever ; good at out of domain data

$$score(D, Q) = \sum_{i=1}^n IDF(q_i) * \frac{\text{문서 } D \text{에서 } q_i \text{의 term frequency}}{f(q_i, D) * (k_1 + 1) \text{ 문서 } D \text{의 길이}} \cdot \frac{1}{f(q_i, D) + k_1 * (1 - b) + b * \frac{|D|}{avgdl}}$$

파라미터
문서 집합의 평균 문서 길이

- DPR : dense retriever ; good at semantic info

$$p_{\eta}(z|x) \propto \exp(\mathbf{d}(z)^{\top} \mathbf{q}(x)) \quad \mathbf{d}(z) = \text{BERT}_d(z), \quad \mathbf{q}(x) = \text{BERT}_q(x)$$

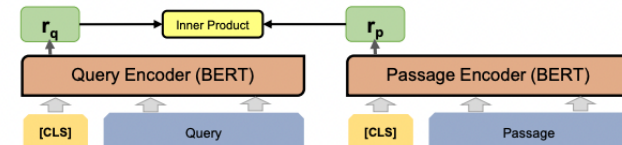


Figure 5: Representation Model for Initial Retrieval

# Re2G

---

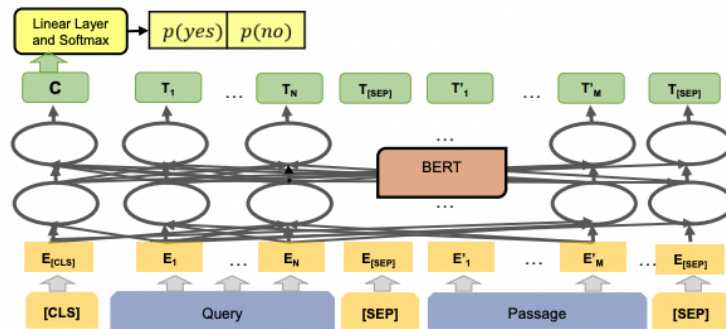


Figure 4: Interaction Model Reranker

Reranker architecture

## Re2G - Training Scheme

---

- Training with Data  $\langle \mathbf{q}, \mathbf{t}, \mathbf{Prov} \rangle$  : query, target output, Provenance
  - 1. Training DPR :  $\langle \mathbf{q}, \mathbf{p}+, \mathbf{p}- \rangle$  triplet ( $\mathbf{p}+$ : positive  $\mathbf{Prov}$ ,  $\mathbf{p}-$ : negative  $\mathbf{Prov}$ )
  - 2. Training Generator (BART) :  $\langle \mathbf{q}, \mathbf{t} \rangle$
  - 3. Training Reranker :  $\langle \mathbf{q}, \mathbf{P}, \mathbf{Prov} \rangle$  ( $\mathbf{P} = \text{BM25}(\mathbf{q}) \cup \text{DPR}(\mathbf{q})$ )  $loss = - \sum_{i \in \mathbf{Prov}} \log(\text{softmax}(\mathbf{z}_r)_i)$
  - 4. Training end-to-end :  $\langle \mathbf{q}, \mathbf{t} \rangle$   $\mathbf{z}_r$ : reranker logit

## Re2G - Training Scheme

---

- Training end-to-end
  - Using reranker passage instead of DPR passage : can't train query encoder
    - (does not use query encoder embedding in re-ranker)
  - Fix : Online knowledge distillation / Freezing query encoder

# Re2G - Result

## KILT IR benchmarks

	T-REx				(Slot Filling)	
	R-Prec	Recall@5	Accuracy	F1	KILT-AC	KILT-F1
Re <sup>2</sup> G (ours)	<b>80.70</b>	<b>89.00</b>	<b>87.68</b>	<b>89.93</b>	<b>75.84</b>	<b>77.05</b>
KG1 <sub>1</sub> [Glass et al., 2021]	74.36	83.14	84.36	87.24	69.14	70.58
KILT-WEB 2 [Piktus et al., 2021]	75.64	87.57	81.34	84.46	64.64	66.64
SEAL [Bevilacqua et al., 2022]	67.80	81.52	83.72	86.53	60.08	61.72
KG1 <sub>0</sub> [Glass et al., 2021]	59.70	70.38	77.90	81.31	55.54	56.79
Natural Questions (Question Answering)						
	R-Prec	Recall@5	Accuracy	F1	KILT-AC	KILT-F1
Re <sup>2</sup> G (ours)	<b>70.78</b>	<b>76.63</b>	<u>51.73</u>	<u>60.97</u>	<b>43.56</b>	<b>49.80</b>
SEAL [Bevilacqua et al., 2022]	63.16	68.19	<b>53.74</b>	<b>62.24</b>	<u>38.78</u>	<u>44.40</u>
KG1 <sub>0</sub> [Glass et al., 2021]	<u>63.71</u>	70.17	45.22	53.38	36.36	41.83
KILT-WEB 2 [Piktus et al., 2021]	59.83	<u>71.17</u>	51.59	60.83	35.32	40.73
RAG [Petroni et al., 2021]	59.49	67.06	44.39	52.35	32.69	37.91
TriviaQA (Question Answering)						
	R-Prec	Recall@5	Accuracy	F1	KILT-AC	KILT-F1
Re <sup>2</sup> G (ours)	<b>72.68</b>	<u>74.23</u>	<b>76.27</b>	<b>81.40</b>	<b>57.91</b>	<b>61.78</b>
SEAL [Bevilacqua et al., 2022]	<u>68.36</u>	<b>76.36</b>	70.86	77.29	<u>50.56</u>	<u>54.99</u>
KILT-WEB 2 [Piktus et al., 2021]	58.85	71.55	<u>72.73</u>	<u>79.54</u>	45.55	49.57
KG1 <sub>0</sub> [Glass et al., 2021]	60.49	63.54	60.99	66.55	42.85	46.08
MultiDPR [Maillard et al., 2021]	61.49	68.33	59.60	66.53	42.36	46.19
FEVER (Fact Checking)						
	R-Prec	Recall@5	Accuracy	KILT-AC		
Re <sup>2</sup> G (ours)	<b>88.92</b>	<b>92.52</b>	<b>89.55</b>	<b>78.53</b>		
SEAL [Bevilacqua et al., 2022]	81.45	<u>89.56</u>	89.54	<u>71.28</u>		
KILT-WEB 2 [Piktus et al., 2021]	74.77	87.89	88.99	65.68		
KG1 <sub>0</sub> [Glass et al., 2021]	75.60	84.95	85.58	64.41		
MultiDPR [Maillard et al., 2021]	74.48	87.52	86.32	63.94		
Wizard of Wikipedia (Dialog)						
	R-Prec	Recall@5	Rouge-L	F1	KILT-RL	KILT-F1
Hindsight [Paranjape et al., 2021]	56.08	74.27	<b>17.06</b>	<b>19.19</b>	<b>11.92</b>	<b>13.39</b>
Re <sup>2</sup> G (ours)	<b>60.10</b>	<b>79.98</b>	<u>16.76</u>	<u>18.90</u>	<u>11.39</u>	<u>12.98</u>
SEAL [Bevilacqua et al., 2022]	57.55	78.96	16.65	18.34	10.45	11.63
KG1 <sub>0</sub> [Glass et al., 2021]	55.37	78.45	16.36	18.57	10.36	11.79
RAG [Petroni et al., 2021]	<u>57.75</u>	74.61	11.57	13.11	7.59	8.75
KILT-WEB 2 [Piktus et al., 2021]	41.54	68.25	13.94	15.66	6.55	7.57

Table 1: KILT leaderboard top systems

R-Prec: Precision for R retrieved passages

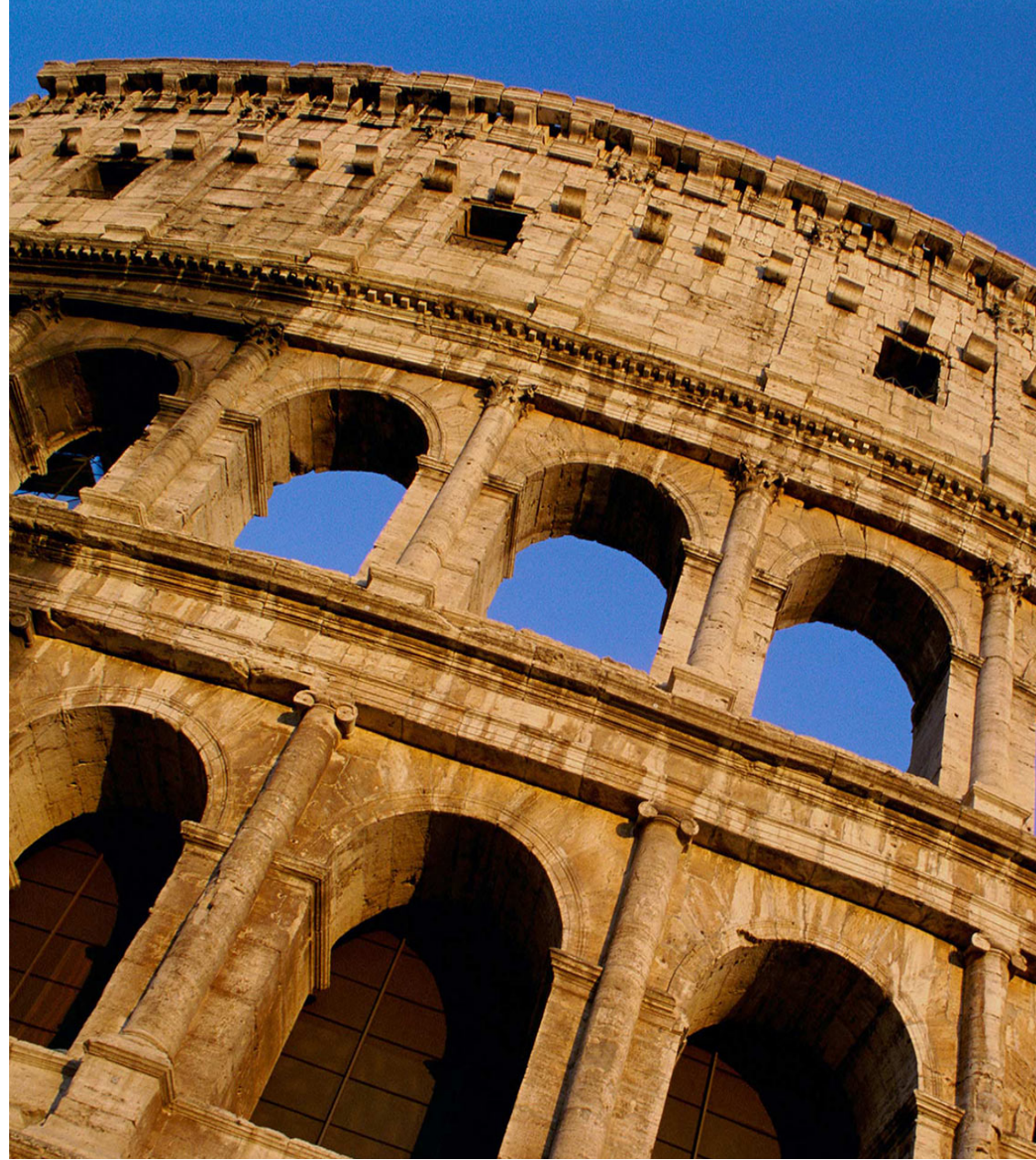
## Re2G - Conclusion

---

- Improvement
  - BM25, DPR 등 retrieval method들의 ensemble 방법 제시 및 성능 향상
- Limitation
  - Very long training scheme
  - Simple re-ranker network architecture

FLARE (2023)

---



# FLARE

---

- **F**oward-**L**ooking **A**ware **R**etrieval augmented generation
- In short-term knowledge intensive task, augmented LM works great
- However, long-term knowledge intensive task augmented LM does not work well as original LM.

**Active Retrieval Augmented Generation (2023)**

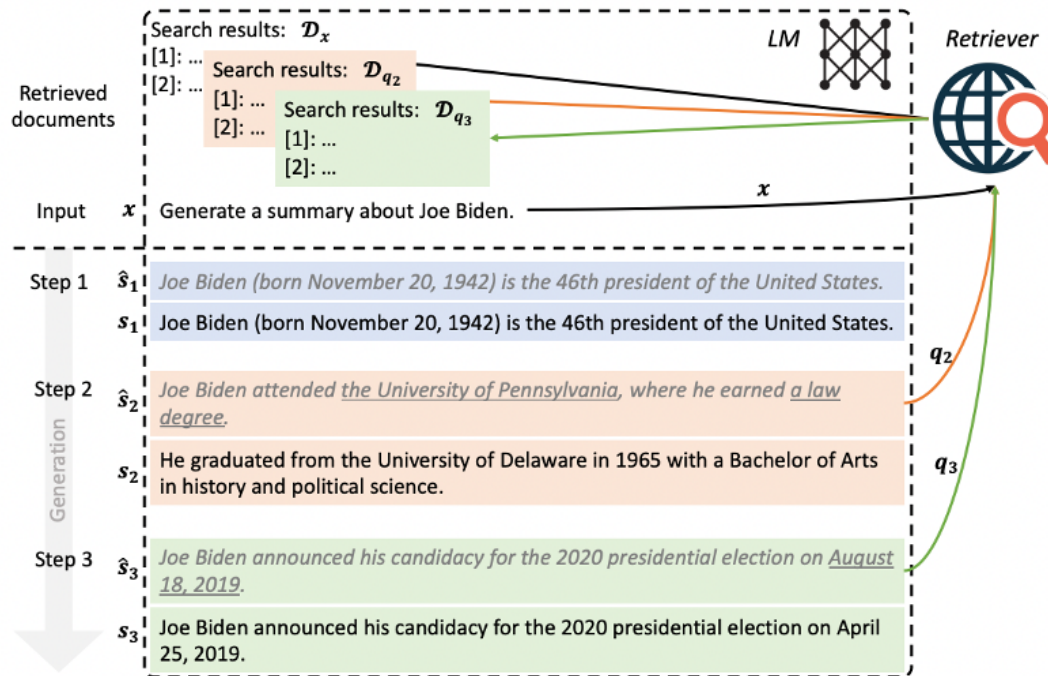
<https://arxiv.org/abs/2305.06983>

# FLARE

---

- RAG retrieve passage only **once**.
- In long-term complex task, one retrieval is not enough.
- Repeat (Retrieve - Generate) and model stops itself actively.
- Assumption : LLM tend to be well-calibrated (필요한 정보만 사용) / low term probability means lack of information. (단어 확률이 낮음 => 정보 부족 의미)

# FLARE - Model



$$q_t = \text{qry}(x, y(<t))$$

$$D_t = \text{ret}(q_t)$$

$$y_t = \text{LLM}(D_t, x, y(<t))$$

if all  $y_t$ 's term prob  $>$  threshold:  
 end retrieval, use  $y_t$  as final output  
 else:  
 $t += 1$   
 repeat

Figure 1: An illustration of forward-looking active retrieval augmented generation (FLARE). Starting with the user input  $x$  and initial retrieval results  $\mathcal{D}_x$ , FLARE iteratively generates a temporary next sentence (shown in *gray italic*) and check whether it contains low-probability tokens (indicated with underline). If so (step 2 and 3), the system retrieves relevant documents and regenerates the sentence.

# FLARE - Model

FLARE\_instruct : generator가 [Search (\*query)] 를 출력 시 이를 그대로 retrieve에 사용.  
 (tend to not search than necessary : increase logit of “[“)  
 (excessive queries : after [Search (\*query)], decrease logit of “[“ for next few tokens)

FLARE\_direct : 앞서 설명한대로 단순히  $y_t$ 를  $q$ 에 붙임.

대신 이경우 잘못된  $y_t$  정보가 들어갈 수 있어서, masking / self-ask QG with chatGPT로써 극복하고자함

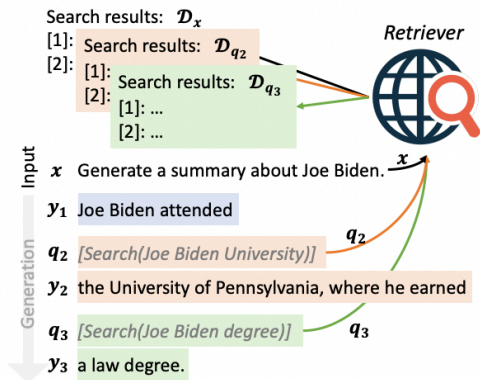


Figure 2: An illustration of forward-looking active retrieval augmented generation with retrieval instructions (FLARE\_instruct). It iteratively generates search queries (shown in *gray italic*) to retrieve relevant information to aid future generations.

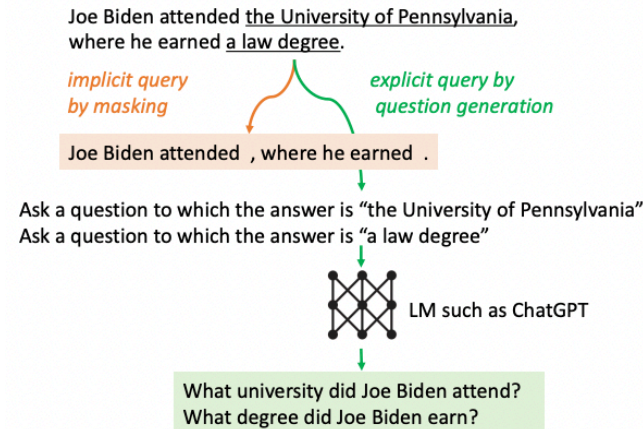


Figure 3: Implicit and explicit query formulation. Tokens with low probabilities are marked with underlines.

# FLARE - Result

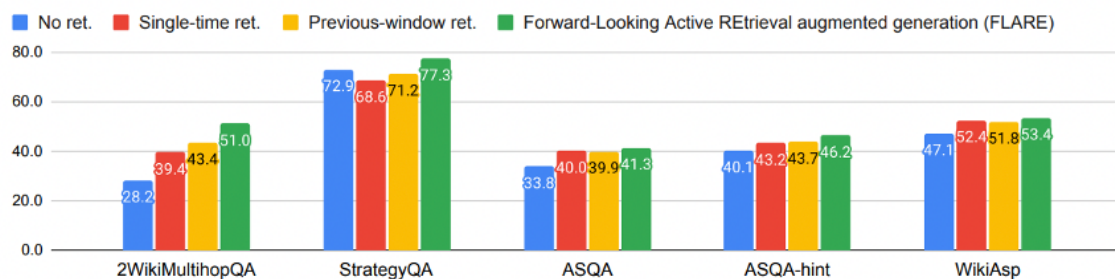


Figure 4: Comparison between FLARE and baselines across all tasks/datasets. We report the primary metric for each dataset: EM for 2WikiMultihopQA, StrategyQA, and ASQA, and UniEval for WikiAsp.

MultihopQA / Commonsense reasoning / long-term QA (x2) /  
open domain summarization

EM: exact match UniEval : measure factual consistency.  
Used generator baseline : GPT3.5

Methods	EM	F <sub>1</sub>	Prec.	Rec.
No retrieval	28.2	36.8	36.5	38.6
Single-time retrieval	39.4	48.8	48.6	51.5
<i>Multi-time retrieval</i>				
Previous-window	43.2	52.3	51.7	54.5
Previous-sentence	39.0	49.2	48.9	51.8
Question decomposition	47.8	56.4	56.1	58.6
FLARE <sub>instruct</sub> (ours)	42.4	49.8	49.1	52.5
FLARE <sub>direct</sub> (ours)	<b>51.0</b>	<b>59.7</b>	<b>59.1</b>	<b>62.6</b>

Table 1: FLARE and baselines on 2WikiMultihopQA. Previous-window (Borgeaud et al., 2022; Ram et al., 2023), previous-sentence (Trivedi et al., 2022), and question decomposition (Press et al., 2022; Yao et al., 2022) methods are reimplemented for fair comparisons.

Baseline : window ret / sentence ret / question decomposition

# FLARE - Result

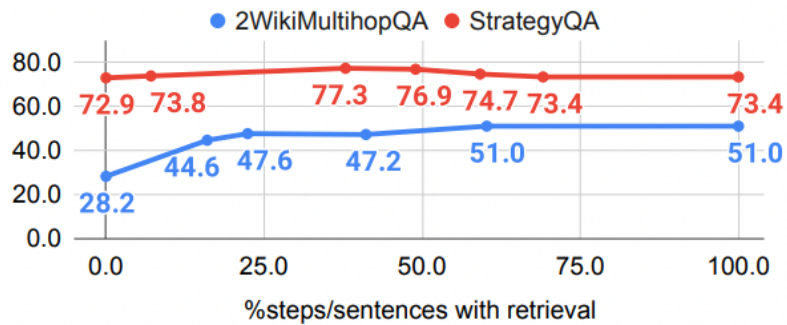


Figure 5: Performance (EM) of FLARE with respect to the percentage of steps/sentences with retrieval on 2WikiMultihopQA and StrategyQA.

- FLARE outperforms others.
- MultihopQA, long-term QA(w. hint) show big improvements
- FLARE\_direct > FLARE\_instruct
- too many retrieval is worse.

## FLARE - Conclusion

---

- Improvement
  - 단어 확률을 바탕으로 스스로 retrieval 횟수 조정, 여러번의 retrieval로 성능 향상
- Limitation
  - 비정교한 threshold 설정
  - 다회 retrieval - generate의 반복으로 인한 cost 증가, overhead
  - 다회 retrieval passage 단순 합성 => 정보 손실 우려

Q&A