

Re²G: Retrieve, Rerank, Generate

Glass, Michael, et al. "Re2G: Retrieve, rerank, generate."
arXiv preprint arXiv:2207.06300 (2022).

Introduction

- GPT-3, T5 등은 텍스트 생성에 있어서 좋은 성능을 보임.
파라미터 크기가 클수록 모델의 성능도 높아짐
- 최근, non-parametric knowledge를 사용한 transformer 연구가 이뤄짐(REALM,RAG)
-> 모델이 사용 가능한 정보량 확장 가능
- 본 논문에서는
neural initial retrieval과 re-ranking을 결합하여
BART기반 seq2seq generation 방법인 **Re²G(Retrieval, Rerank, Generate)** 접근법을 제안한다
 1. **Re-Ranking 접근방식 개선**
-> 점수를 비교하기 어려운 것들의 병합을 가능하게 함
(ex) BM25, neural initial retrieval)
 2. **End-to-End 학습을 위한 접근방식**
Knowledge distillation을 통해 ground truth data를 중심으로
initial retriever, re-ranker, generation을 통합적으로 학습

--> Reranking 효과 재입증, 초기 리트리버 방법을 ensembling으로 확장(neural+tradition)

Introduction

1. Re-Ranking 접근방식 개선

- > 비교하기 어려운 score의 병합을 가능하게 함
(ex) BM25, neural initial retrieval)

2. End-to-End 학습을 위한 접근방식

Knowledge distillation을 통해 ground truth data를 중심으로
initial retriever, re-ranker, generation을 통합적으로 학습

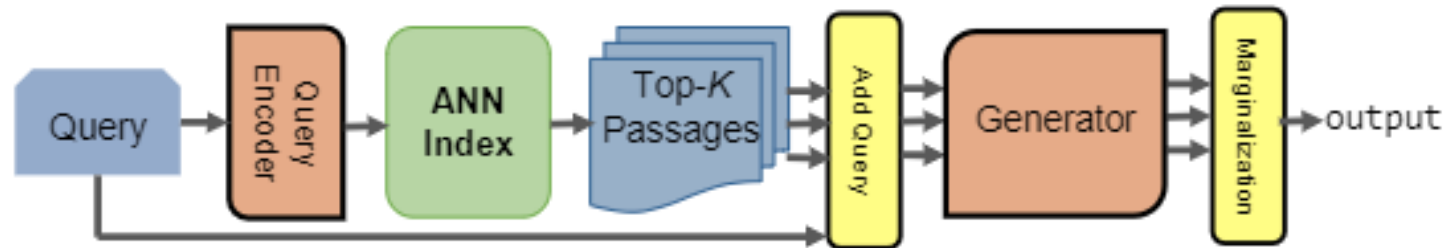


Figure 2: RAG Architecture

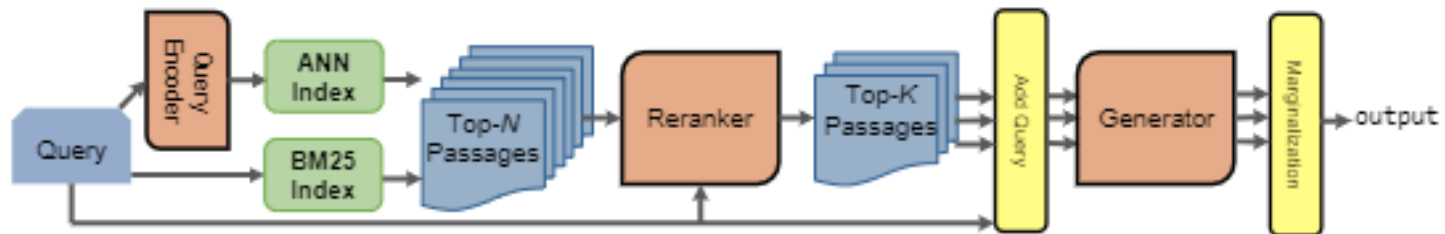


Figure 3: Re²G Architecture

Method

ReRanker

: 초기 검색 단계에서 반환된 문서나 결과 목록을 재정렬(re-ranking)

1) **initial retriever**은 Ranker의 사용을 통해 크게 개선될 수 있다.

- * Learning to Rank (학습간 순위 생성 / 2009, Liu)
initial retriever 후 training된 모델로 각 passag를 평가(Ranking)

- * Cascade Ranking Model (계단식 순위 모델 / 2011, Wang)

| | | | | |
|-------------|----|--------------------|----|----------------|
| 초기검색 | -> | 중간검색 | -> | 최종검색 |
| BM25등으로 빠르게 | | n-gram 등 확률적 언어모델로 | | neural network |

2) 비교할 수 없는 **score 결과**를 병합할 수 있음.

ex) BM25와 DPR의 retriever 결과

1. ReRanker

Initial Retrieval

- 초기 후보군 검색
- q,p 독립적으로 임베딩

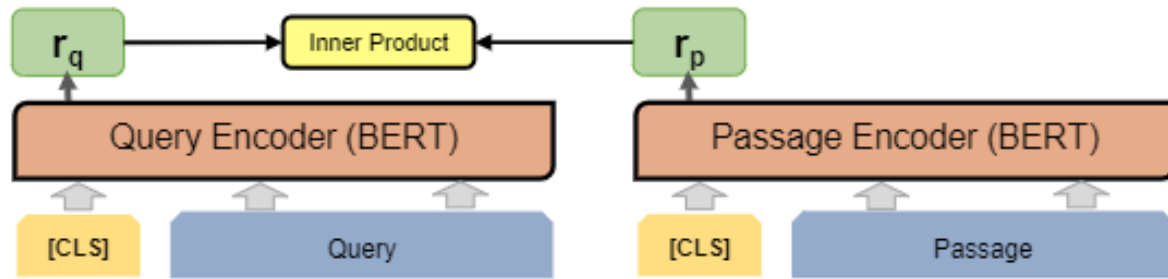


Figure 5: Representation Model for Initial Retrieval

ReRanker

- 후보군 재정렬
- q,p 동시에 입력 -> cross-attention

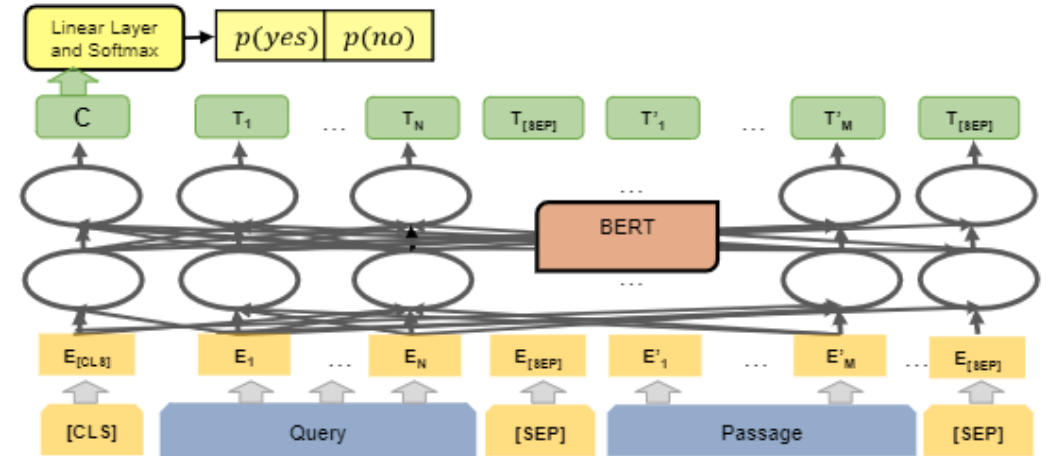


Figure 4: Interaction Model Reranker

-> Accuracy, scalability

2. Training

KILT Downstream: $\langle q, t, \text{prov} \rangle$

Prov ground truth(정답 근거 출처) : 관련 문서 검색

Target : 최종 사용자 질문에 대한 정확한 답변 생성

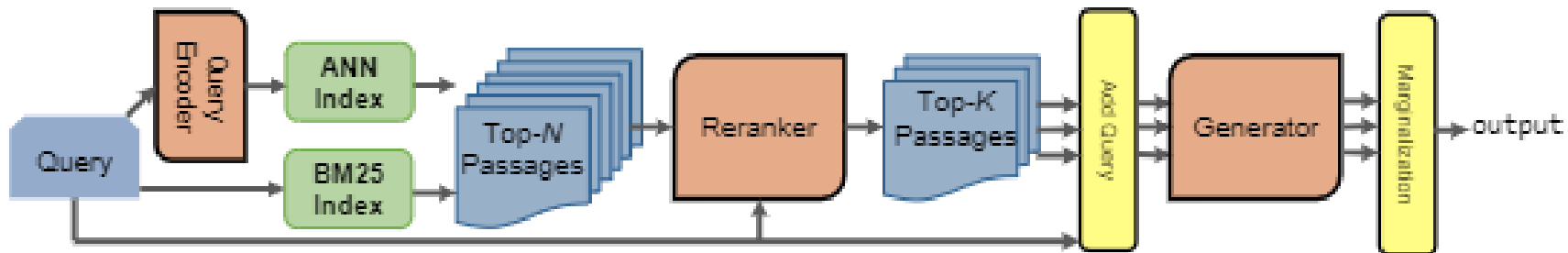
1. DPR training

-> KGI의 기본구조를 따름

2. Generator training

3. Re-Ranking training

4. end-to-end training



Related Work

KILT Benchmark(Knowledge Intensive Language Tasks)

: 지식 집약적인 작업을 다루며, 모델이 외부 지식을 검색하고 이를 활용하는 능력. 즉, 지식 검색-활용-생성을 통합적으로 평가

예) 질의응답, 사실 검증, 슬롯 채우기, 대화 생성 등.

<q,t,Prov>

t : target output

Prov: 목표출력을 지원하는 출처 문서의 집합

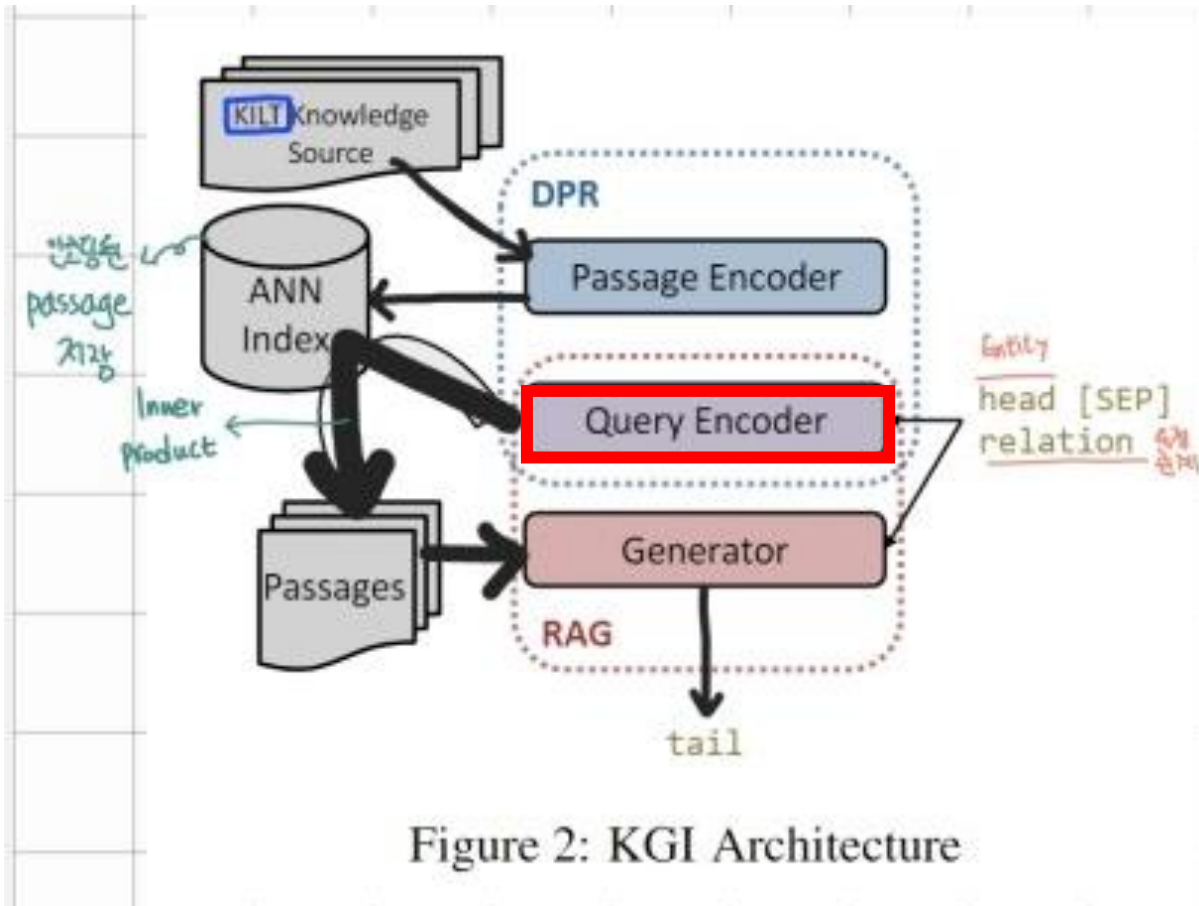
| | | |
|---|--|---|
| T-REx | Dracula (7923) Dracula is an 1897 Gothic horror novel by Irish author Bram Stoker. It introduced the character of Count Dracula, and established many conventions of subsequent vampire fantasy. The novel tells the story of Dracula's attempt to move from Transylvania to England so that he may find new blood and spread the undead curse, and of the battle between Dracula and a small group of men and a woman led by Professor Abraham Van Helsing. | Wizard of Wikipedia |
| Input: Dracula [SEP] narrative location Output: Transylvania Provenance: 7923-2 | | Input: <ul style="list-style-type: none">• I really like vampires!!• Vampires are intense and based on European folklore. Do you have any favorite vampires?• I think dracula is the best one!!! Output: He's one of the best! He's based on the character from the 1897 horror book of the same name. Provenance: 7923-1 |
| Natural Questions | | |
| Input: when did bram stoker's dracula come out Output: 1897 Provenance: 7923-1 | | |
| FEVER | | |
| Input: Dracula is a novel by a Scottish author. Output: REFUTES Provenance: 7923-1 | | |

Figure 1: KILT tasks of slot filling, question answering, fact checking and dialog

Related Work

KGI(Knowledge Graph Induction)

: RAG의 검색 및 생성 성능을 확장한 모델 = RAG+DPR



> > 지식그래프와 유사.

- 1) 입출력 형식 유사(head-relation-tail)
- 2) target외에도 prov 요구

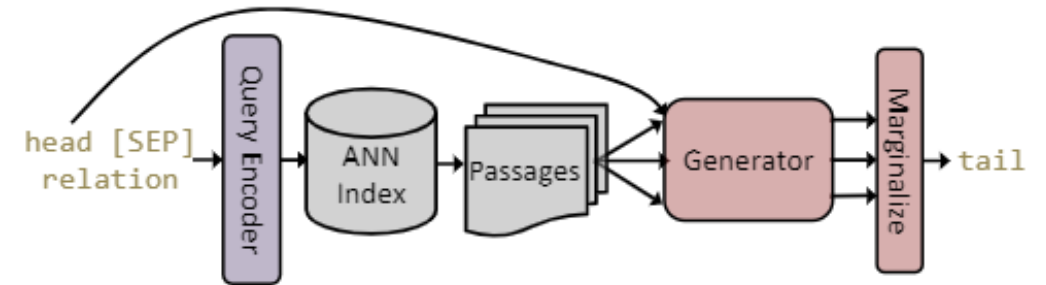


Figure 4: RAG Architecture

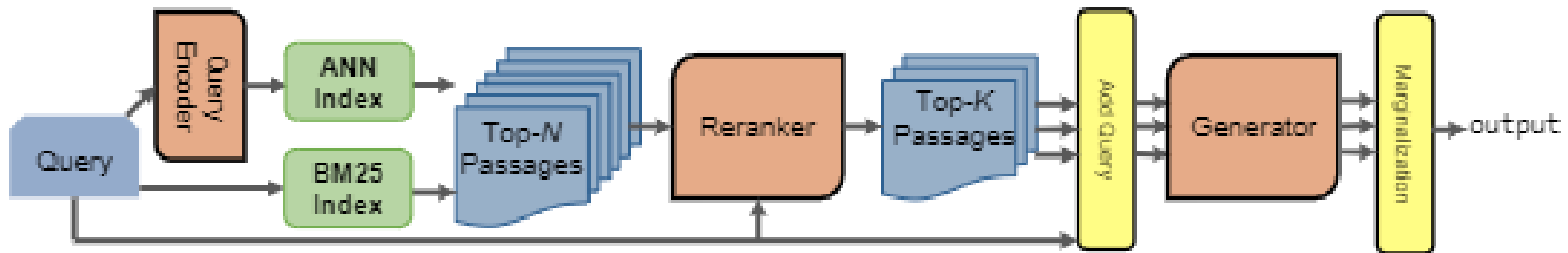
2. Training

KILT Downstream: $\langle q, t, \text{prov} \rangle$

Prov ground truth(정답 근거 출처) : 관련 문서 검색

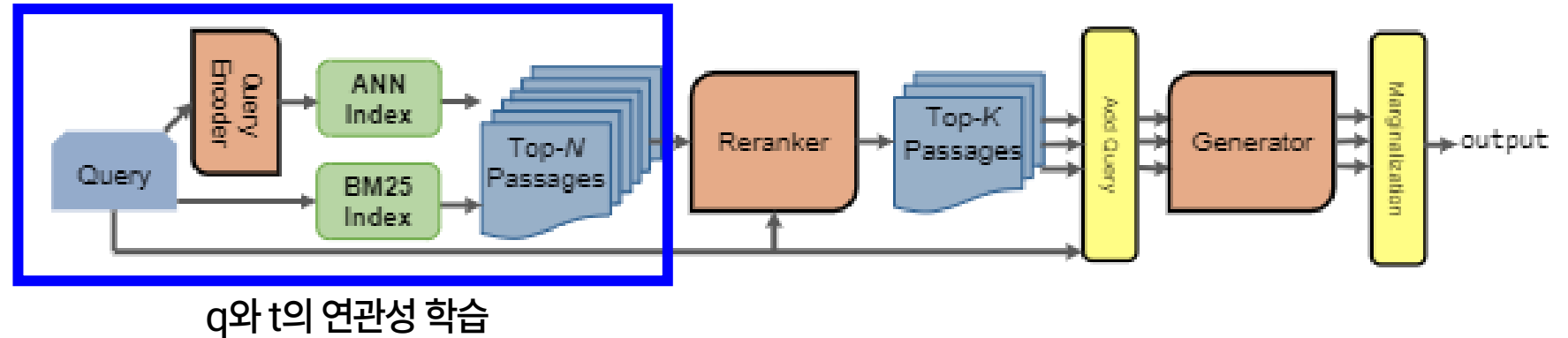
Target : 최종 사용자 질문에 대한 정확한 답변 생성

1. DPR training $\langle q, p+, p- \rangle$
2. Generator training $\langle q, t \rangle$
3. Re-Ranking training $\langle q, p, \text{Prov} \rangle$
4. end-to-end training (DPR, Reranker, Generator)



2-1. DPR Training

1. DPR training
2. Generator training
3. Re-Ranking training
4. end-to-end training



$\langle q, p^+, p^- \rangle$

$p^+ \in \text{Prov}$

$p^- \in \text{BM25}(q) \cap p^- \notin \text{Prov}$

Hard negative Batch negative

- NLL 사용

$$L(q_i, p_i^+, p_{i,1}^-, \dots, p_{i,n}^-) = -\log \frac{e^{\text{sim}(q_i, p_i^+)}}{e^{\text{sim}(q_i, p_i^+)} + \sum_{j=1}^n e^{\text{sim}(q_i, p_{i,j}^-)}}$$

- train 완료 시, passage를 FAISS를 사용해 HNSW로 Indexing

Hard Negative

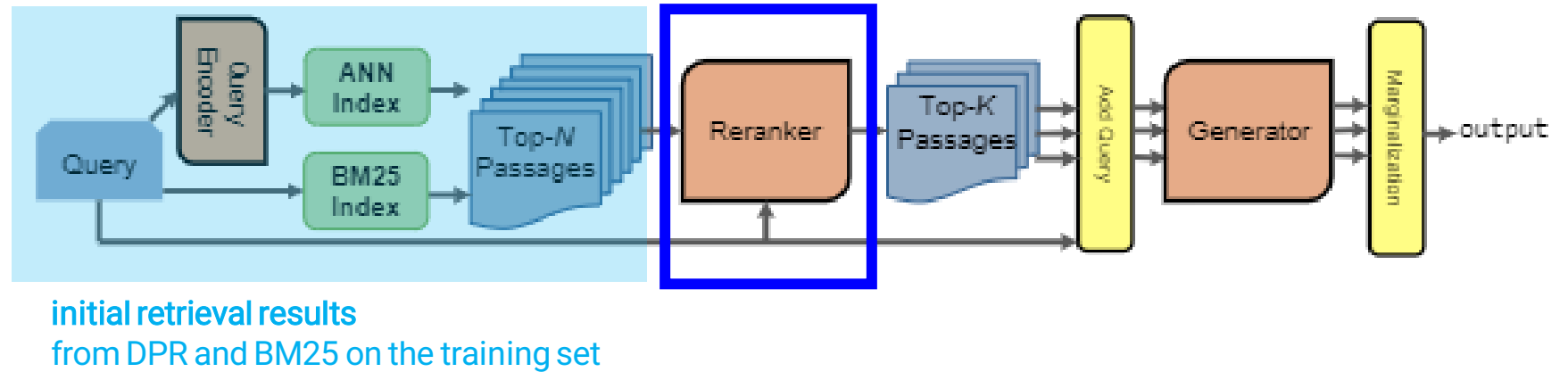
: 쿼리와 유사하지만 정답이 포함되지 않은 문서. (BM25 등으로 유사한 passage 택 -> 추가 retriever 필요)

Batch Negative

: 동일한 batch 내 다른 query의 정답 passage (추가 retriever 필요 X)

2-3. ReRanking Training

1. DPR training
2. Generator training
3. Re-Ranking training
4. end-to-end training



$\langle q, p, \text{Prov} \rangle$

- Reranking training begins with the application of DPR and BM25, producing tuples: $\langle q, \mathbf{P}, \text{Prov} \rangle$ where $\mathbf{P} = \text{BM25}(q) \cup \text{DPR}(q)$

Prov : 정답 문서 집합

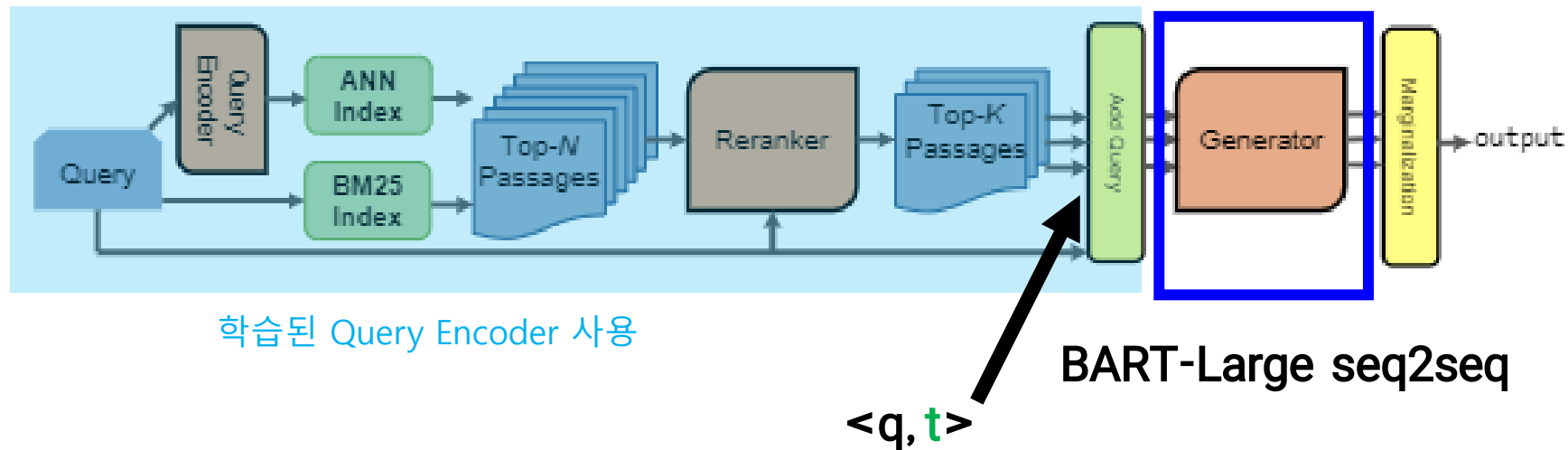
P : 후보 문서 집합 ($P^+ \cup P^-$)

- Loss Function
(여러 개의 Positive Passage가 존재하는 상황)

$$\text{loss} = - \sum_{i \in \text{Prov}} \log(\text{softmax}(\mathbf{z}_r)_i)$$

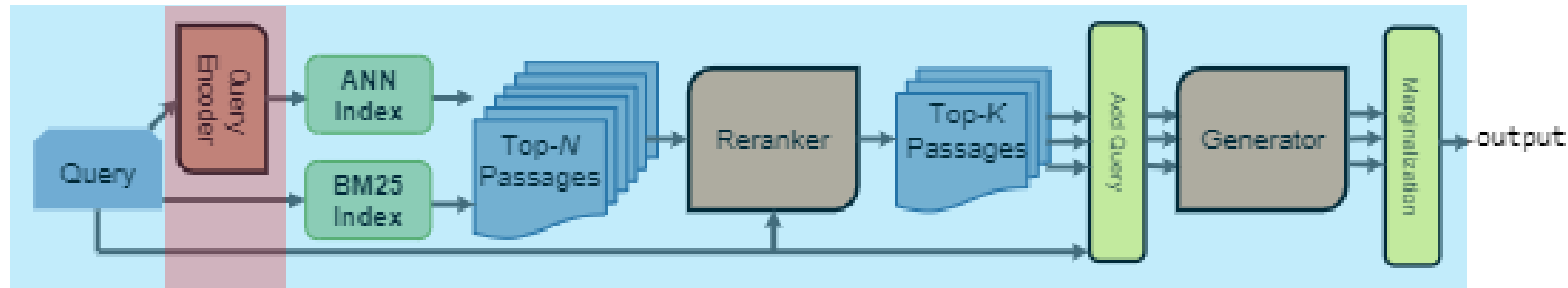
2-2. Generator Training

1. DPR training
2. Generator training
3. Re-Ranking training
4. end-to-end training



2-4. End-to-End Training

1. DPR training
2. Generator training
3. Re-Ranking training
4. end-to-end training



학습 단절

End-to-End(DPR, Reranker, Generator)

[기존 loss]

$$P(s_j) = \text{softmax}(\mathbf{z}_r)_j$$

$$P(t_i|s_j) = \text{softmax}(\text{BART}(s_j)_i)_{t_i}$$

$$\text{loss} = - \sum_{i,j} \log(P(t_i|s_j) \cdot P(s_j))$$

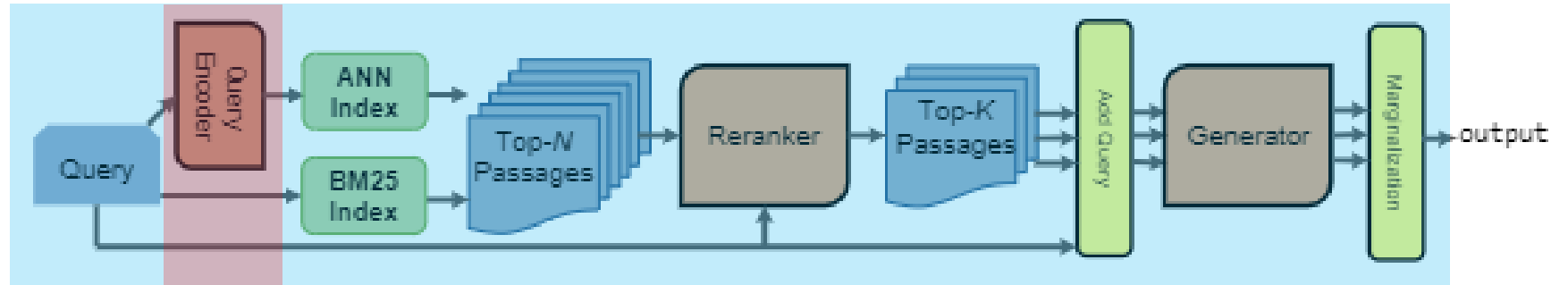


Query Encoder(DPR)의 학습 단절 문제 발생

1. Gradient 전파의 단절
2. DPR과 ReRanking 모델의 분리

2-4. End-to-End Training

1. DPR training
2. Generator training
3. Re-Ranking training
4. end-to-end training



학습 단절

End-to-End(DPR, Reranker, Generator)

[기존 loss]

$$\begin{aligned} P(s_j) &= \text{softmax}(\mathbf{z}_r)_j \\ P(t_i|s_j) &= \text{softmax}(\text{BART}(s_j)_i)_{t_i} \\ \text{loss} &= - \sum_{i,j} \log(P(t_i|s_j) \cdot P(s_j)) \end{aligned}$$



Query Encoder(DPR)의 학습 단절 문제 발생



제안

• Combine the DPR and reranker scores

• Freeze the query encoder

• Online Knowledge Distillation



[본 논문의 loss]

$$\text{loss} = D_{KL} \left(\text{softmax} \left(\frac{\mathbf{z}_s}{T} \right) \parallel \text{softmax} \left(\frac{\mathbf{z}_t}{T} \right) \right) \cdot T^2$$

Student 모델
점수분포

Teacher 모델
점수분포



Experiments

Table 1: KILT leaderboard top systems

| | T-REx (Slot Filling) | | | | | |
|---------------------------------------|--|-----------------|-----------------|--------------|----------------|----------------|
| | R-Prec | Recall@5 | Accuracy | F1 | KILT-AC | KILT-F1 |
| Re ² G (ours) | 80.70 | 89.00 | 87.68 | 89.93 | 75.84 | 77.05 |
| KGI ₁ [Glass et al., 2021] | 74.36 | 83.14 | <u>84.36</u> | <u>87.24</u> | <u>69.14</u> | <u>70.58</u> |
| KILT-WEB 2 [Piktus et al., 2021] | <u>75.64</u> | <u>87.57</u> | 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 |
| KGI ₀ [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 ² G (ours) | 70.78 | 76.63 | <u>51.73</u> | <u>60.97</u> | 43.56 | 49.80 |
| SEAL [Bevilacqua et al., 2022] | 63.16 | 68.19 | 53.74 | 62.24 | <u>38.78</u> | <u>44.40</u> |
| KGI ₀ [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 |

Experiments

| | | | | TriviaQA (Question Answering) | | | | | | | |
|--|--|--|--|---------------------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--|
| | | | | R-Prec | Recall@5 | Accuracy | F1 | KILT-AC | KILT-F1 | | |
| | | | | Re ² G (ours) | 72.68 | <u>74.23</u> | 76.27 | 81.40 | 57.91 | 61.78 | |
| | | | | SEAL [Bevilacqua et al., 2022] | 68.36 | 76.36 | 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 | |
| | | | | KGI ₀ [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 ² G (ours) | 88.92 | 92.52 | 89.55 | 78.53 | | | |
| | | | | SEAL [Bevilacqua et al., 2022] | <u>81.45</u> | <u>89.56</u> | <u>89.54</u> | <u>71.28</u> | | | |
| | | | | KILT-WEB 2 [Piktus et al., 2021] | 74.77 | 87.89 | 88.99 | 65.68 | | | |
| | | | | KGI ₀ [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 | 17.06 | 19.19 | 11.92 | 13.39 | |
| | | | | Re ² G (ours) | 60.10 | 79.98 | <u>16.76</u> | <u>18.90</u> | <u>11.39</u> | <u>12.98</u> | |
| | | | | SEAL [Bevilacqua et al., 2022] | 57.55 | <u>78.96</u> | 16.65 | 18.34 | 10.45 | 11.63 | |
| | | | | KGI ₀ [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 | |

| T-REx | | | |
|---------------------------------------|--------------|--------------|--------------|
| | R-Prec | Recall@5 | Accuracy |
| Re ² G (ours) | 80.70 | 89.00 | 87.68 |
| KGI ₁ [Glass et al., 2021] | 74.36 | 83.14 | <u>84.36</u> |
| KILT-WEB 2 [Piktus et al., 2021] | <u>75.64</u> | <u>87.57</u> | 81.34 |
| SEAL [Bevilacqua et al., 2022] | 67.80 | 81.52 | 83.72 |
| KGI ₀ [Glass et al., 2021] | 59.70 | 70.38 | 77.90 |
| Natural Questions | | | |
| | R-Prec | Recall@5 | Accuracy |
| Re ² G (ours) | 70.78 | 76.63 | <u>51.73</u> |
| SEAL [Bevilacqua et al., 2022] | 63.16 | 68.19 | 53.74 |
| KGI ₀ [Glass et al., 2021] | <u>63.71</u> | 70.17 | 45.22 |
| KILT-WEB 2 [Piktus et al., 2021] | 59.83 | <u>71.17</u> | 51.59 |
| RAG [Petroni et al., 2021] | 59.49 | 67.06 | 44.39 |

Table 1: KILT leaderboard top systems

Experiments

| | T-REx | | NQ | | TriviaQA | | FEVER | | WoW | |
|----------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | R-Prec | R@5 | R-Prec | R@5 | R-Prec | R@5 | R-Prec | R@5 | R-Prec | R@5 |
| BM25 | 46.88 | 69.59 | 24.99 | 42.57 | 26.48 | 45.57 | 42.73 | 70.48 | 27.44 | 45.74 |
| DPR Stage 1 | 49.02 | 63.34 | 56.64 | 64.38 | 60.12 | 64.04 | 75.49 | 84.66 | 34.74 | 60.22 |
| KGI ₀ DPR | 65.02 | 75.52 | 64.65 | 69.60 | 60.55 | 63.65 | 80.34 | 86.53 | 48.04 | 71.02 |
| Re ² G DPR | 67.16 | 76.42 | 65.88 | 70.90 | 62.33 | 65.72 | 84.13 | 87.90 | 47.09 | 69.88 |
| KGI ₀ DPR+BM25 | 60.48 | 80.06 | 36.91 | 66.94 | 40.81 | 64.79 | 65.95 | 90.34 | 35.63 | 68.47 |
| Reranker Stage 1 | 81.22 | 87.00 | 70.78 | 73.05 | 71.80 | 71.98 | 87.71 | 92.43 | 55.50 | 74.98 |
| Re ² G Reranker | 81.24 | 88.58 | 70.92 | 74.79 | 60.37 | 70.61 | 90.06 | 92.91 | 57.89 | 74.62 |

Table 2: Development Set Results for Retrieval

Conclusion

Reranker

BM25와 같은 여러 검색 소스를 통합하여 정확도 개선

Knowledge Distillation

DPR 점수와 직접적으로 의존하지 않는 학습 방식임에도 불구하고 긍정적인 효과를 보임

Ensembling effective

Reranker를 통해 DPR과 BM25의 ensembling이 가능해져, 5개 데이터셋 중 4개에서 성능 향상을 달성.

성능 개선

- Re2G는 KILT 데이터셋의 5개 태스크에서 SOTA(State-of-the-Art)를 크게 향상시켰으며, 그중 4개 태스크에서 여전히 최고 성능을 유지.
- 이전 모델(RAG, KGI) 대비 검색 성능과 엔드-투-엔드 성능(슬롯 채우기, 질문 답변, 사실 검증, 대화) 모두에서 크게 개선됨.

향후 연구 방향

질문 답변과 대화 태스크에서 Re2G의 도메인 적응에 대한 추가 실험