

# DPR 을 활용한 검색 솔루션

- 다차원 벡터 활용을 통한 문서 유사도 검색 (Wisenut)
- 

팀명 현명한호두

팀장 홍세일

팀원 정재훈 이정호 유예지 조

은상

# 목차

table of contents

- 1 프로젝트 배경
- 2 **Dense Passage Retrieval**
- 3 **Pretraining Model**
- 4 유사도 엔진
- 5 검색 엔진 연동
- 6 **Conclusion**



# 1

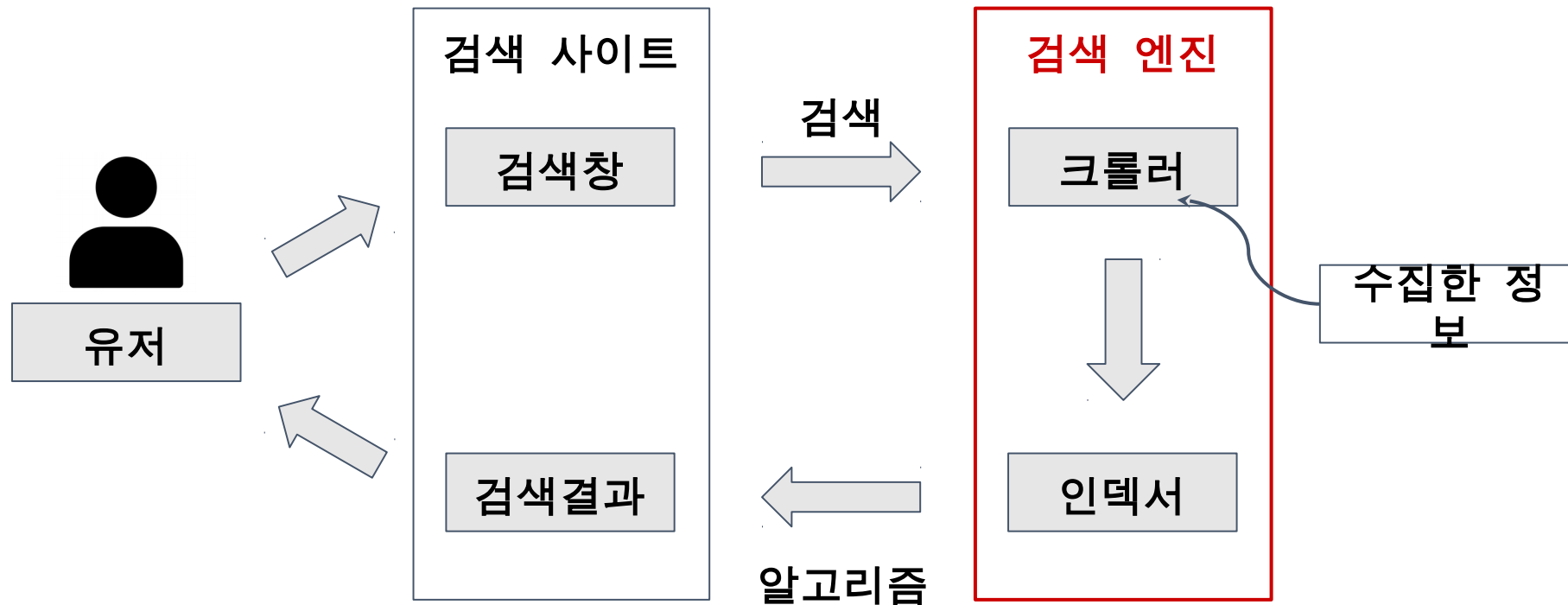
## 프로젝트 배경

---

# # ‘검색 엔진’ 이란 ?

- 검색 엔진 (Search Engine)

→ 사용자가 웹에서 특정 정보를 찾을 때 , 검색어를 입력하면 그 검색어와 관련된 웹페이지를 찾아주는 프로그램



# # 배경 및 문제점

## 1. 프로젝트 배경



# 2

## Dense Passage Retrieval

---

참고 논문

Dense Passage Retrieval for Open-Domain Question  
Answering

<https://arxiv.org/pdf/2004.04906.pdf>

# # 'DPR (Dense Passage Retrieval)' 이란 ?

## 2. Dense Passage Retrieval

### < Retrieval 의 Backbone Model >

‘사전훈련된 BERT model’ 로 이루어진  
‘2 개의 Encoder’ (Question Encoder, Passage Encoder)

활용

### < 유사도 측정 방법 >

question 벡터와 passage 벡터의 내적 (dot product)

$$\text{sim}(q, p) = E_Q(q)^T E_P(p).$$

### < Model Training >

loss 를 구하는 식 :

NLL (Negative Log Likelihood)

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

### < Negative sample >

loss 가 감소하는 방향으로 학습하기 위해

- Random : 코퍼스 내의 random 한 passage 를 선택
- BM25 : 코퍼스 내에서 BM25 기준으로 top-k 의 문서 사용
- Gold : 학습셋 내의 다른 질의의 positive passage 선택

### < In-batch negatives >

모델 학습 시 같은 batch 내에서 (In-batch) gold passage 를 negative passage 로 활용하는 것 (재사용)

계산적인 효율 뿐만 아니라, 좋은 성능을 낼 수

$$S = QP^T S = QP^T \text{ (BxB)}$$

< test set 의 retrieval 정

Trainin g0	Retriev er	top20					top-100				
		NQ	TriviaQ A	WQ	TREC	SQuAD	NQ	TriviaQ A	WQ	TREC	SQuAD
None	BM25	59.1	66.9	55.0	70.9	68.8	73.7	76.7	71.1	84.1	80.0
Singl e	DPR	78.4	79.4	73.2	79.8	63.2	85.4	85.0	81.4	89.1	77.2
	BM25 + DPR	78.4	79.8	71.0	85.2	71.5	83.8	84.5	80.5	92.7	81.3
Multi	DPR	79.4	78.8	75.0	89.1	51.6	86.0	84.7	82.9	93.9	67.6
	BM25 + DPR	78.0	79.9	74.7	88.5	66.2	83.9	84.4	82.3	94.1	78.6

- Single : 각각의 Dataset 에 대하여 학습시킴
- Multiple : SQuAD 를 제외한 4 개의 Dataset 을 합쳐서 학습시킴
- 40 epoch (NQ, TriviaQA, SQuAD)
- 100 epoch (TREC, WQ)
- dropout ratio : 0.1
- optimizer : Adam
- top k : 상위 k 개의 retrieval 정확도를 가지고 측정



# # DPR 학습 결과

## 2. Dense Passage Retrieval

< NQ 데이터셋에 대해 서로 다른 학습 방법을 테스트한 결과 >

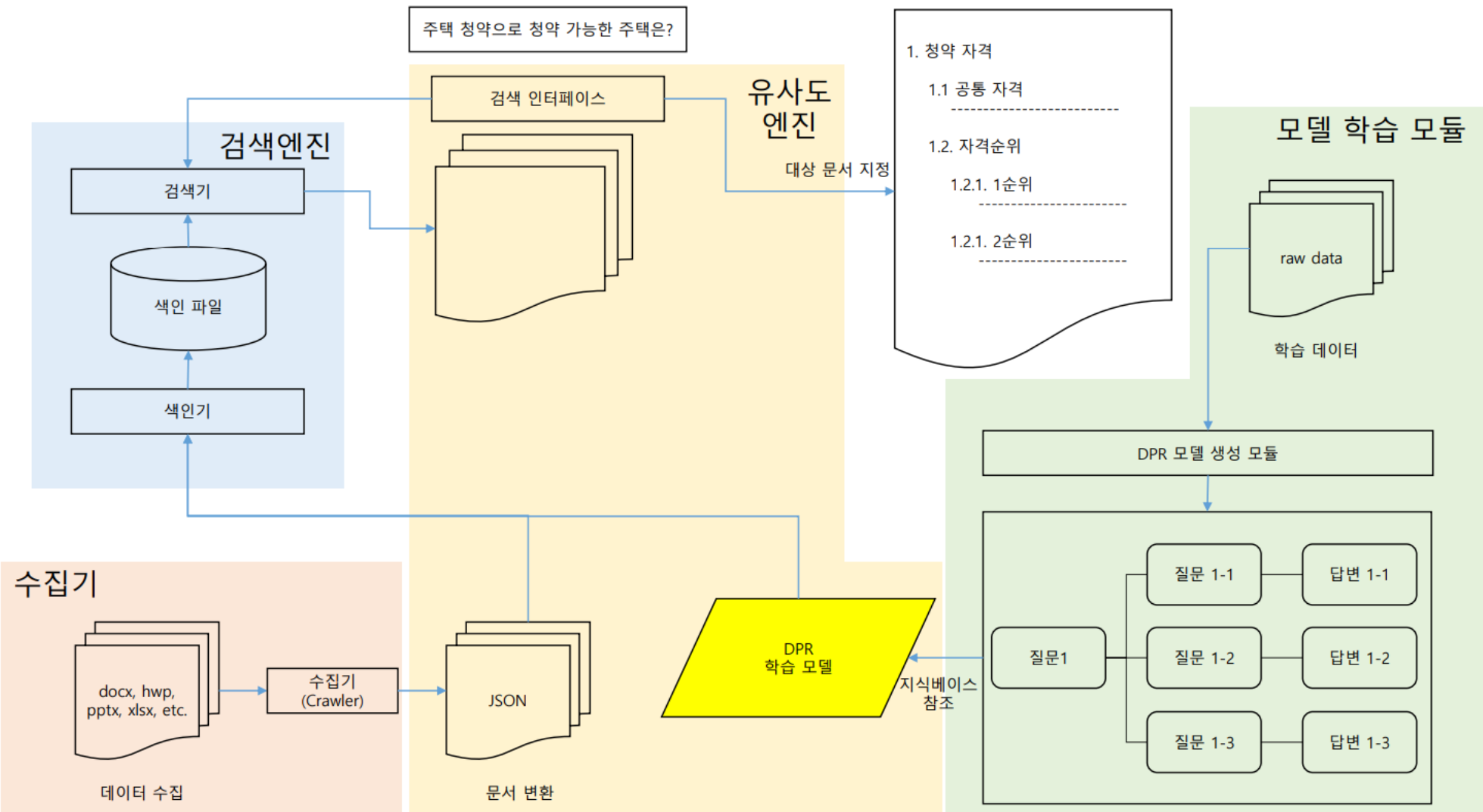
Type	#N	IB	Top-5	Top-20	Top-100
Random	7	✗	47.0	64.3	77.8
BM25	7	✗	50.0	63.3	74.8
Gold	7	✗	42.6	63.1	78.3
Gold	7	✓	51.1	69.1	80.8
Gold	31	✓	52.1	70.8	82.1
Gold	127	✓	55.8	73.0	83.1
G.+BM25 <sup>(1)</sup>	31+32	✓	65.0	77.3	84.4
G.+BM25 <sup>(2)</sup>	31+64	✓	64.5	76.4	84.0
G.+BM25 <sup>(1)</sup>	127+128	✓	<b>65.8</b>	<b>78.0</b>	<b>84.9</b>

#N : Negative Sample 개수 , IB : In - Batch 적용 여부

< end-to-end QA accuracy >

Training	Model	NQ	TriviaQA	WQ	TREC	SQuAD
Single	BM25+BERT (Lee et al., 2019)	26.5	47.1	17.7	21.3	33.2
Single	ORQA (Lee et al., 2019)	33.3	45.0	36.4	30.1	20.2
Single	HardEM (Min et al., 2019a)	28.1	50.9	-	-	-
Single	GraphRetriever (Min et al., 2019b)	34.5	56.0	36.4	-	-
Single	PathRetriever (Asai et al., 2020)	32.6	-	-	-	<b>56.5</b>
Single	REALM <sub>Wiki</sub> (Guu et al., 2020)	39.2	-	40.2	46.8	-
Single	REALM <sub>News</sub> (Guu et al., 2020)	40.4	-	40.7	42.9	-
Single	BM25	32.6	52.4	29.9	24.9	38.1
	DPR	<b>41.5</b>	56.8	34.6	25.9	29.8
	BM25+DPR	39.0	57.0	35.2	28.0	36.7
Multi	DPR	<b>41.5</b>	56.8	<b>42.4</b>	49.4	24.1
	BM25+DPR	38.8	<b>57.9</b>	41.1	<b>50.6</b>	35.8

# # 프로젝트 구조도



## # 사용한 tools



# 3

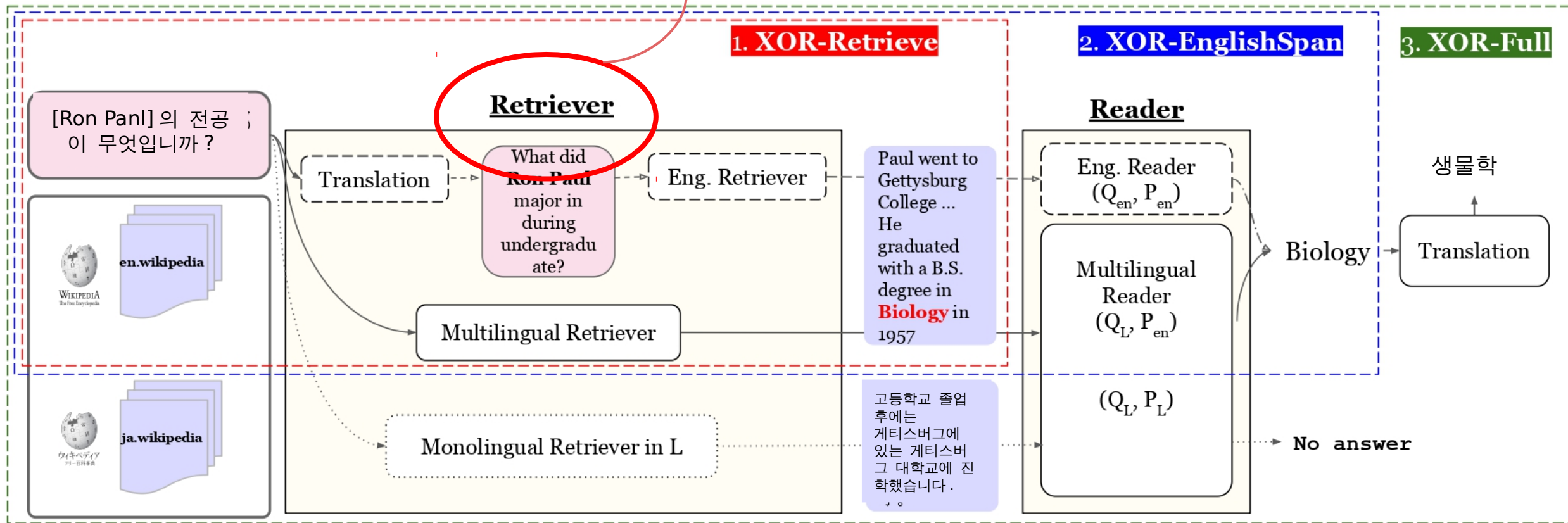
**Pretraining**  
~~**Model**~~

---

# # DPR 구조

## 3. Pretraining Model

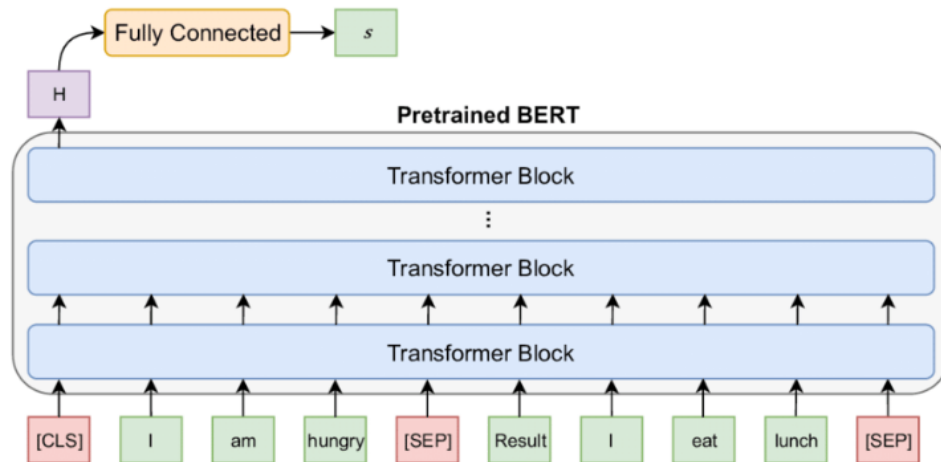
Bert 모델 사용 → 한국어 처리를 위해 **Kobert**로 변경



# # BERT , KoBERT Model

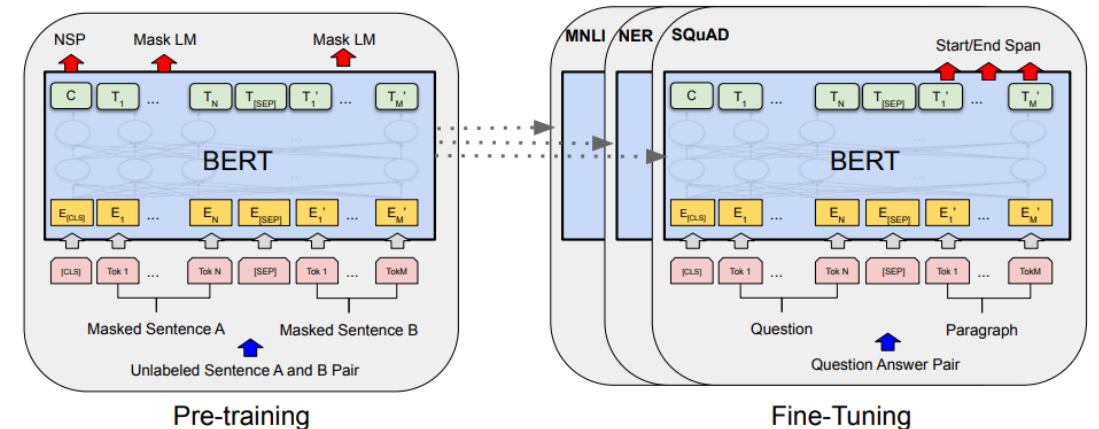
## 3. Pretraining Model

### < BERT Model >



- 약 33억개의 단어로 pretrain 되어있는 기계 번역 모델
- 사용목적에 따라 fine-tuning 이 가능
- 텍스트를 양방향 ( 앞뒤 ) 로 확인하여 자연어 처리
- 영어의 정확도는 높지만 , 한국어의 정확도는 낮음

### < KoBERT Model >



- BERT Model 에서 “한국어 데이터” 를 추가로 학습시킨 모델 ( 한국어 위키에서 5 백만개의 문장과 5,400 만개의 단어를 학습시킨 모델 )
- 한국어 버전 **BERT Model**  
( 한국어 데이터에서 높은 정확도 )

# # Input data

## 3. Pretraining Model

### < input\_train.json >

```
[{'question': '바그너는 괴테의 파우스트를 읽고 무엇을 쓰고자 했는가?',  
  'answers': ['교향곡'],  
  'positive_ctxs': [{'title': '파우스트 서곡',  
    'text': '1839년 바그너는 괴테의 파우스트를 처음 읽고 그 내용에  
    마음이 끌려 이를 소재로 해서 하나의 교향곡을 쓰려는 뜻을 갖는다. 이  
    시기 바그너는 1838년에 빛 독촉으로 산전수전을 다 겪은 상황이라 좌절과  
    실망에 가득했으며 메피스토펠레스를 만나는 파우스트의 심경에 공감했다고  
    한다. 또한 파리에서 아브네크의 지휘로 파리 음악원 관현악단이 연주하는  
    베토벤의 교향곡 9번을 듣고 깊은 감명을 받았는데, 이것이 이듬해 1월  
    에 파우스트의 서곡으로 쓰여진 이 작품에 조금이라도 영향을 끼쳤으리라는  
    것은 의심할 여지가 없다. 여기의 라단조 조성의 경우에도 그의 전기에 적  
    혀 있는 것처럼 단순한 정신적 피로나 실의가 반영된 것이 아니라 베토벤의  
    합창교향곡 조성의 영향을 받은 것을 볼 수 있다. 그렇게 교향곡 작곡을  
    1839년부터 40년에 걸쳐 파리에서 착수했으나 1악장을 쓴 뒤에 중단했  
    다. 또한 작품의 완성과 동시에 그는 이 서곡(1악장)을 파리 음악원의  
    연주회에서 연주할 파트보까지 준비하였으나, 실제로는 이루어지지  
    않았다. 결국 초연은 4년 반이 지난 후에 드레스덴에서 연주되었고 재연도 이  
    루어졌지만, 이후에 그대로 방치되고 말았다. 그 사이에 그는 리엔치와  
    방황하는 네덜란드인을 완성하고 탄호이저에도 착수하는 등 분주한 시간을  
    보냈는데, 그런 바쁜 생활이 이 곡을 잊게 한 것이 아닌가 하는 의견도 있  
    다.'}]},  
  ...
```

### < input\_dev.json >

```
[{'question': '임종석이 여의도 농민 폭력 시위를 주도한 혐의로 지명  
수배 된 날은?',  
  'answers': ['1989년 2월 15일'],  
  'positive_ctxs': [{'title': '임종석',  
    'text': '1989년 2월 15일 여의도 농민 폭력 시위를 주도한 혐  
    의 ( 폭력행위등처벌에관한법률위반 ) 으로 지명수배되었다. 1989년 3월  
    12일 서울지방검찰청公安부는 임종석의 사전구속영장을 발부받았다. 같  
    은 해 6월 30일 평양축전에 임수경을 대표로 파견하여 국가보안법위반  
    혐의가 추가되었다. 경찰은 12월 18일~20일 사이 서울 경희대학교에  
    서 임종석이 성명 발표를 추진하고 있다는 첩보를 입수했고, 12월 18일  
    오전 7시 40분 경 가스총과 전자봉으로 무장한 특공조 및 대공과 직원  
    12명 등 22명의 사복 경찰을 승용차 8대에 나누어 경희대학교에 투입했  
    다. 1989년 12월 18일 오전 8시 15분 경 서울청량리경찰서는 호위  
    학생 5명과 함께 경희대학교 학생회관 건물 계단을 내려오는 임종석을 발  
    견, 검거해 구속을 집행했다. 임종석은 청량리경찰서에서 약 1시간 동  
    안 조사를 받은 뒤 오전 9시 50분 경 서울 장안동의 서울지방검찰청 공  
    안분실로 인계되었다.'}]},  
  {'question': '1989년 6월 30일 평양축전에 대표로 파견 된 인물  
    은?',  
    'answers': ['임수경'],  
    'positive_ctxs': [{'title': '임종석',  
      'text': '1989년 2월 15일 여의도 농민 폭력 시위를 주도한 혐  
      의 ( 폭력행위등처벌에관한법률위반 ) 으로 지명수배되었다. 1989년 3월  
      12일 서울지방검찰청公安부는 임종석의 사전구속영장을 발부받았다.  
      ...
```

# # KoBERT code

## 1. RESOURCE\_MAP 에 한국어 데이터 추가

```
"data.retriever.input-dev": {
  "s3_url":
  "https://dl.fbaipublicfiles.com/dpr/data/retriever/biencoder-squad1-
train.json.gz",
  "original_ext": ".json",
  "compressed": False,
  "desc": "SQUAD 1.1 train subset with passages pools for the
Retriever training",
},
"data.retriever.input-train": {
  "s3_url":
  "https://dl.fbaipublicfiles.com/dpr/data/retriever/biencoder-squad1-
train.json.gz",
  "original_ext": ".json",
  "compressed": False,
```

## 2. encoder train default.yaml 과 retriever default.yaml 수정

```
},
input_train:
  _target_: dpr.data.biencoder_data.JsonQADataset
  file: data.retriever.input-train

input_dev:
  _target_: dpr.data.biencoder_data.JsonQADataset
  file: data.retriever.input-dev
```

## 3. Pretraining Model

### 3. hf\_bert.yaml 에서 하이퍼 파라미 터 수정

```
# @package _group_

# model type. One of [hf_bert, pytext_bert,
fairseq_roberta]
encoder_model_type: hf_bert

# HuggingFace's config name for model
initialization
pretrained_model_cfg: skt/kobert-base-v1

# Some encoders need to be initialized from a file
pretrained_file:

# Extra linear layer on top of standard
bert/roberta encoder
projection_dim: 0

# Max length of the encoder input sequence
sequence_length: 512

dropout: 0.1

# whether to fix (don't update) context encoder
during training or not
fix_ctx_encoder: False

# if False, the model won't load pre-trained BERT
weights
```



# # KoBERT code

## 3. Pretraining Model

### 4. hf\_model.py 의 kobert 변환

```
def get_bert_tokenizer(pretrained_cfg_name: str,
do_lower_case: bool = True):
    return
    KoBERTTokenizer.from_pretrained(pretrained_cfg_name,
do_lower_case=do_lower_case)

class BertTensorizer(Tensorizer):
    def __init__(self, tokenizer: KoBERTTokenizer, max_length:
int, pad_to_max: bool = True):
        self.tokenizer = tokenizer
        self.max_length = max_length
        self.pad_to_max = pad_to_max

if transformers.__version__.startswith("4"):
    from transformers import BertConfig, BertModel
    from transformers import AdamW
    from kobert_tokenizer import KoBERTTokenizer
    from transformers import RobertaTokenizer
```

### 5. hf\_model.py 의 truncation 전략 수정

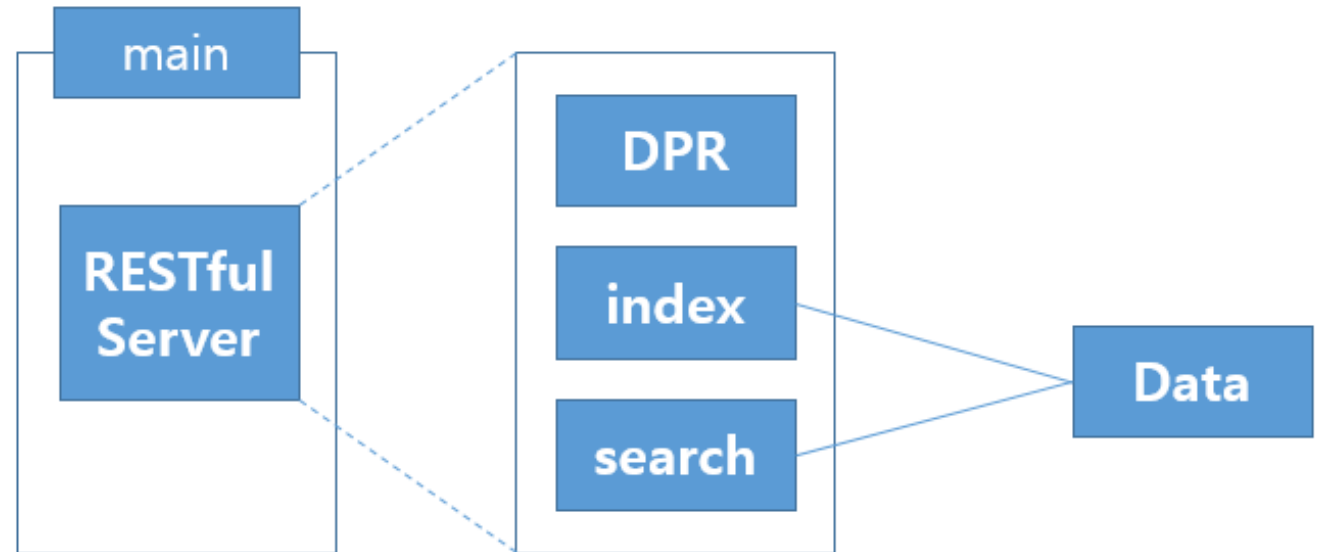
```
class BertTensorizer(Tensorizer):
    def __init__(self, tokenizer: KoBERTTokenizer,
max_length: int, pad_to_max: bool = True):
        self.tokenizer = tokenizer
        self.max_length = max_length
        self.pad_to_max = pad_to_max

        if title:
            token_ids = self.tokenizer.encode(
                title,
                text_pair=text,
                add_special_tokens=add_special_tokens,
                max_length=self.max_length if apply_max_len
            else 10000,
                pad_to_max_length=False,
                truncation="only_second",
            )
        else:
            token_ids = self.tokenizer.encode(
                text,
                add_special_tokens=add_special_tokens,
                max_length=self.max_length if apply_max_len
            else 10000,
                pad_to_max_length=False,
                truncation="only_second",
            )
```

# 4

유사도 엔진

---



### 1.DPR : 두 문장 간의 유사도 파악

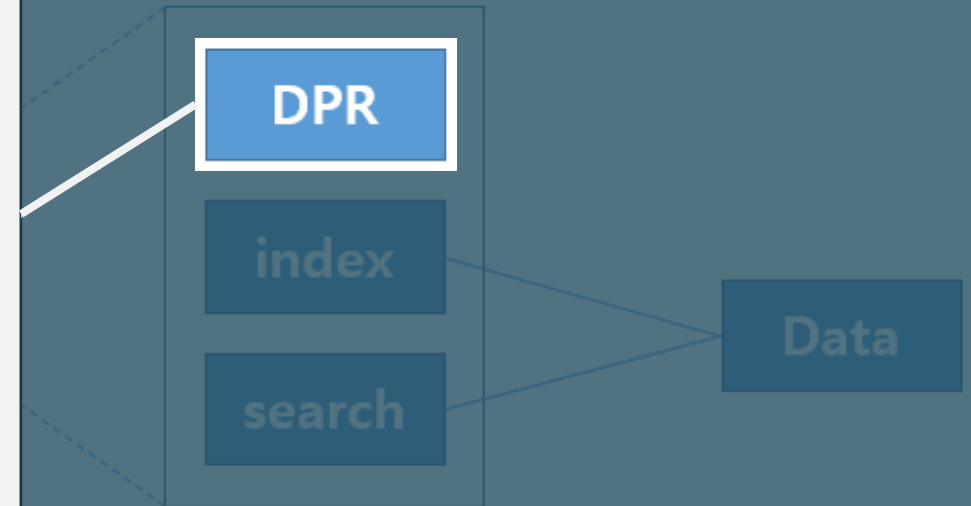
문장 유사도 찾기

← → ↻ ⚠ 주의 요함 | 172.30.1.53:5000/similarity/172.30.1.83

CAKD Home **similarity** index search

문장 1:

문장 2:



# # DPR code : similarity.py

```
def get_idx(*args):
    idxs = []
    for _ in args:
        idxs.append(tokenizer(_, return_tensors="pt")["input_ids"])
    return idxs

def get_pooleroutput(List):
    embeddings = []
    for _ in List:
        embeddings.append(model(_).pooler_output)
    return embeddings

def dot_product_scores(q_vectors: T, ctx_vectors: T) -> T:
    """
    calculates q->ctx scores for every row in ctx_vector
    :param q_vector:
    :param ctx_vector:
    :return:
    """
    # q_vector: n1 x D, ctx_vectors: n2 x D, result n1 x n2
    r = torch.matmul(q_vectors, torch.transpose(ctx_vectors, 0, 1))
    return r
```

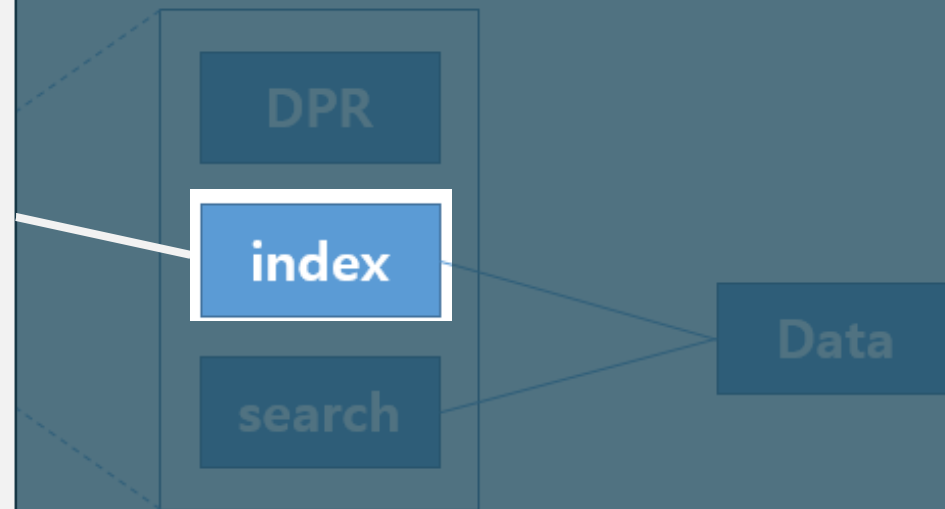
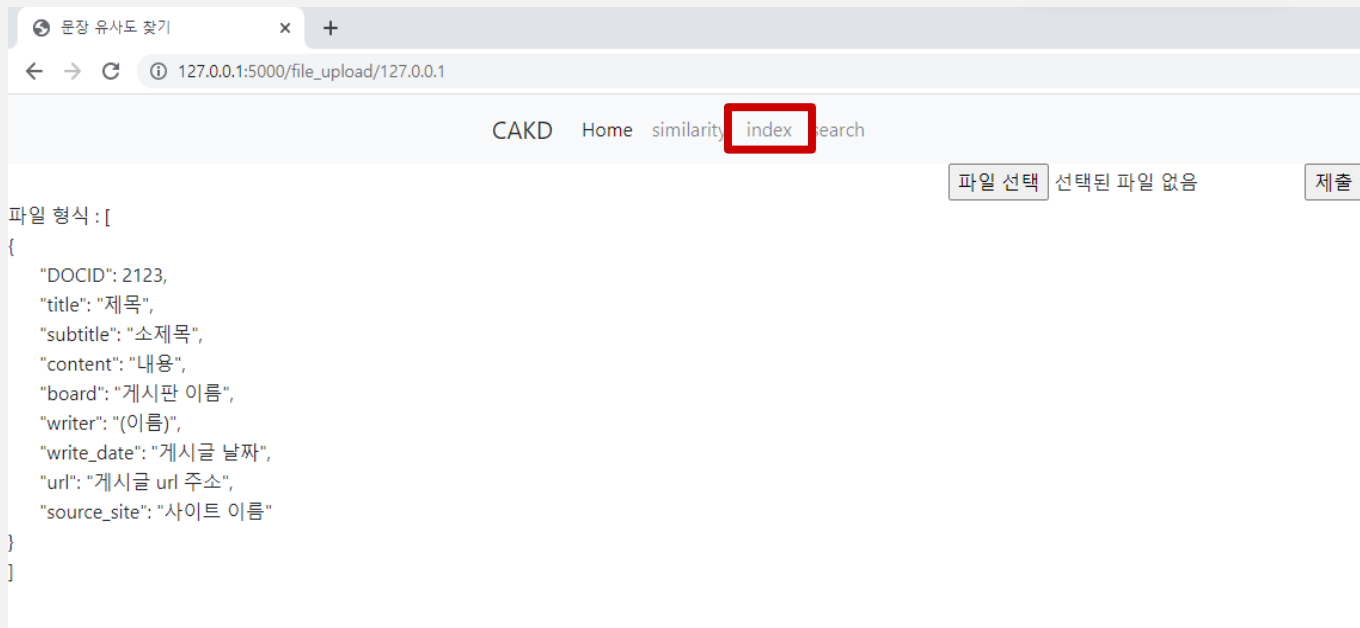
```
def cosine_scores(q_vector: T, ctx_vectors: T):
    # q_vector: n1 x D, ctx_vectors: n2 x D, result n1 x n2
    return F.cosine_similarity(q_vector, ctx_vectors, dim=1)

def get_total_scores(q_vector: T, ctx_vectors: T):
    cos_score=cosine_scores(q_vector, ctx_vectors)
    dot_pdt_score=dot_product_scores(q_vector, ctx_vectors)
    total_score = cos_score
    return round(total_score.item(), 4)

def get_title_dpr(item):
    _ = tokenizer(item, return_tensors="pt")["input_ids"]
    input_data = _.to(device)
    res = model(input_data).pooler_output.detach().numpy().tolist()
    return res

def get_content_dpr(item):
    _ = tokenizer(item, return_tensors="pt", truncation=True, max_length=512)
    ["input_ids"]
    input_data = _.to(device)
    res = model(input_data).pooler_output.detach().numpy().tolist()
    return res
```

### 2. Index : 파일 인덱싱



# # Index code : index.py

```
mecab = Komoran()

def file_save(f,fname,upload_folder):
    try :
        file_path = os.path.join(upload_folder, fname)
        print(file_path)
        f.save(file_path)
        if fname in os.listdir(upload_folder):
            print('success')
            return (True, file_path)
        else :
            print('failed')
            return (False, file_path)

    except Exception as e:
        print('-----')
        print(e)
        return (False, file_path)
        new_dict['title'] = item['title']
        new_dict['title_dpr'] = title_dpr

return_tensors="pt",truncation=True,max_length=512)
["input_ids"]).pooler_output.detach().numpy().tolist()
    content_dpr = get_content_dpr(item['content'])
    new_dict['content'] = item['content']
    new_dict['content_dpr'] = content_dpr
    new_dict['content_morphs'] = get_morphs(item['content'])
    new_dict['title_morphs'] = get_morphs(item['title'])
    new_dict['DOCID'] = item['DOCID']
    es.index(index='my_index', body=new_dict,id = new_dict['DOCID'])
    end = time.time()
    count += 1
    print(count,':',end-start)

    #bulk(es,json_file,index = 'my_index')
return True
```

```
def file_upload_in_db(db, document_obj):
    try:
        db.session.add(document_obj)
        db.session.commit()
        return True
    except:
        return False

def get_morphs(item):
    res = mecab.morphs(item.replace('[^가-힣a-zA-Z0-9]', ''))
    return res

def file_indexing(path):
    with open(path,'r',encoding='utf-8') as f:
        json_file=json.load(f)
        count = 0
        for item in json_file:
            start = time.time()
            new_dict = {}
            #title_dpr = model(tokenizer(item['title'], return_tensors="pt")
["input_ids"]).pooler_output.detach().numpy().tolist()
            title_dpr = get_title_dpr(item['title'])
            content_dpr = get_content_dpr(item['content'])
            new_dict['content'] = item['content']
            new_dict['content_dpr'] = content_dpr
            new_dict['content_morphs'] = get_morphs(item['content'])
            new_dict['title_morphs'] = get_morphs(item['title'])
            new_dict['DOCID'] = item['DOCID']
            es.index(index='my_index', body=new_dict,id = new_dict['DOCID'])
            end = time.time()
            count += 1
            print(count,':',end-start)

            #bulk(es,json_file,index = 'my_index')
return True
```

### 3. Search : 유사도 top 문장 출력

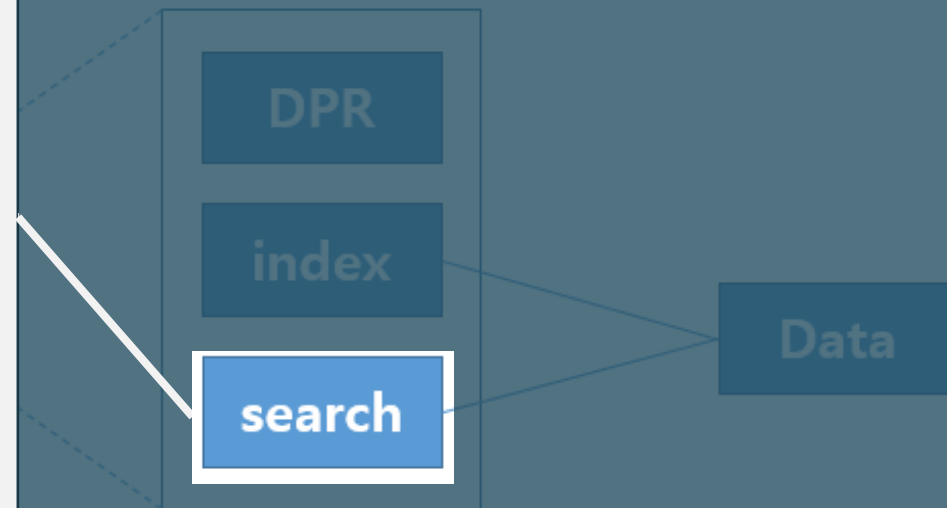
문장 유사도 찾기 x +

← → ↻ 주의 요함 | 172.30.1.53:5000/search/172.30.1.83

CAKD Home similarity index **search**

query :

input n :





## # Search code : data\_util.py

```
def get_similarity_in_document(input_text, input_nbest, json_objs):
    for json_obj in json_objs:
        indicies = get_idx(input_text, json_obj['text'])
        embeddings = get_pooleroutput(indicies)
        json_obj['similarity'] = get_total_scores(embeddings[0], embeddings[1])
    json_objs.sort(key=lambda x : -x['similarity'])
    for i, json_obj in enumerate(json_objs):
        json_obj['nbest'] = i+1
        if (i+1) == input_nbest: break;
    return json_objs

def search_in_es(q, n):
    que = q.replace('[^A-Za-z0-9가-힣]', '')
    query_with_tag = mecab.pos(que)
    print(query_with_tag)
    print(mecab.tagset)
    query = []
    for item in query_with_tag:
        if (item[1] == 'XR') | (item[1] == 'NNG') | (item[1] == 'VV') | (item[1] ==
'NNB') | (item[1] == 'NNP') | (item[1] == 'VA') | (item[1] == 'NP') | (item[1] ==
'SN') | (item[1] == 'NA') | (item[1] == 'NR') | (item[1] == 'SL'):
            query.append(item[0])
    print(mecab.pos(que)[0][0])
    print(query)
    data = [] # 반환할 데이터를 담을 리스트를 초기화합니다
    for item in query:
        query = {
            "query": {
                "multi_match": {
                    "query": item,
                    "fields": ["title", "content", "title_morphs", "content_morphs"]
                }
            }
        }
    res = es.search(index = 'my_index', body = query)
    hits = res.get('hits', {}).get('hits', []) # 실제 데이터가 있는 hits 리스트를 추출
합니다
```

```
# 각각의 문서에 대해서 필요한 정보만 추출하여 data 리스트에 추가합니다
for hit in hits:
    source = hit.get('_source', {})
    item = {
        'title': source.get('title', ''),
        'content': source.get('content', ''),
        'title_dpr': source.get('title_dpr', []),
        'content_dpr': source.get('content_dpr', []),
    }
    if item not in data:
        data.append(item)
for item in data:
    question = model(tokenizer(q, return_tensors="pt", truncation=True, max_length=512)
["input_ids"]).pooler_output
    ctx = torch.tensor(item['content_dpr'])
    title = torch.tensor(item['title_dpr'])
    score1 = get_total_scores(question, ctx)
    score2 = get_total_scores(question, title)
    item.setdefault('score', score1 + score2)
sorted_list = sorted(data, key=lambda x: x['score'], reverse=True)
print(len(sorted_list))
return sorted_list[:n]
```

# 5

## 검색 엔진 연

동

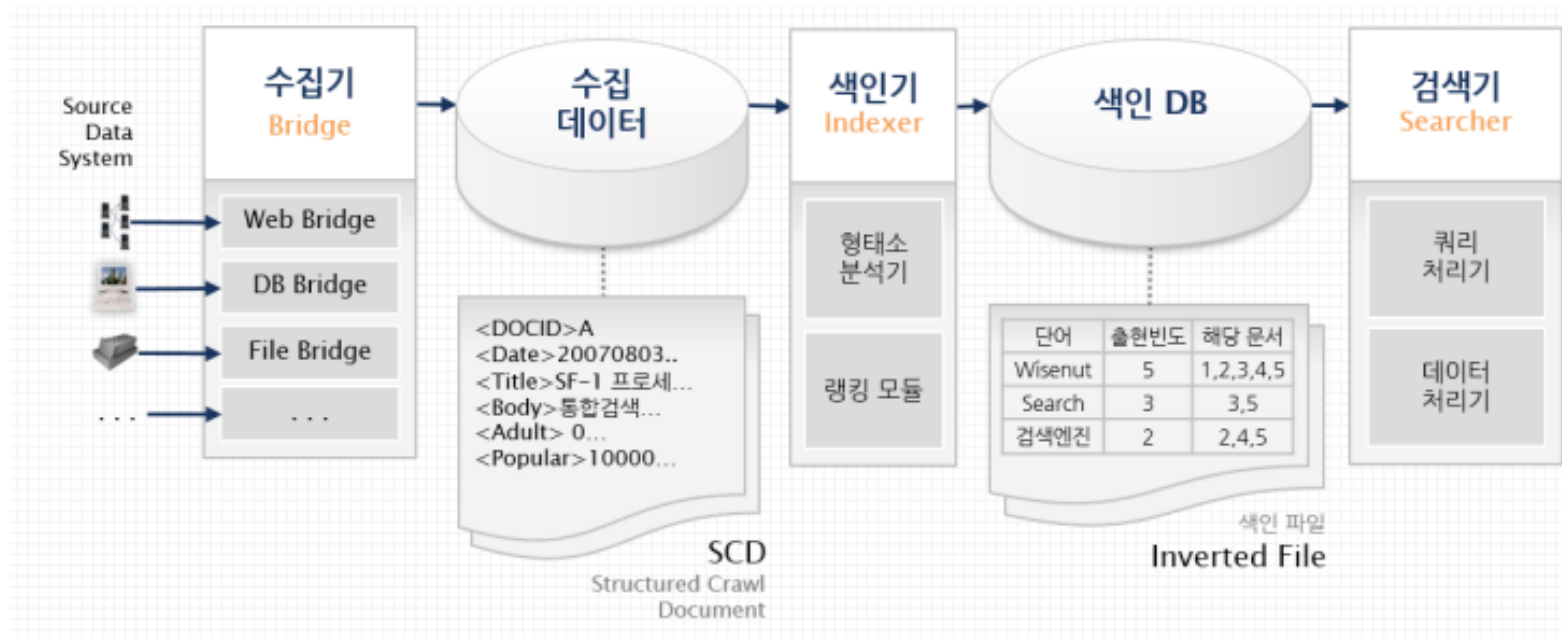
색인 / 색인어 및 검색어 분석 / 역색인 /  
Elastic search

# # ‘색인’이란?

## ● 색인 (Indexing)

- 데이터베이스나 검색 엔진에서 검색 속도를 높이기 위해 사용되는 기술
- 특정 칼럼에 대한 정렬된 데이터 구조를 생성하여 검색 속도를 높임

## ● 색인 과정



# # ‘색인어 및 검색어 분석’ 이란 ?

## ● 검색의 기본

- “문서” 혹은 “검색어” 를 어떠한 방식으로 주요 단어를 추출할지 결정
- ‘어떻게 추출하는가’ 에 따라 검색의 품질 및 정확성이 보장

## ● 색인어 및 검색어 분석 과정

문서 → 레굴레이션 → 색인어 추출 → 색인 DB ← 검색어 추출 ← 레굴레이션  
← 검색어

## ● 레굴레이션 (Regulation) - 구분자 및 병합자 지저

예시 ) “안녕하세요 ? # 중앙정보 ! 홍길동 + 입니다 .” → [ 안녕하세요 ], [ 중앙정보 ],  
[ 홍길동 ], [ 입니다 ]

# # ‘역색인’ 이란?

## ● 역색인 (Inverted index)

- 매우 빠른 풀텍스트 검색을 할 수 있도록 설계된 것
- 문서에 나타나는 모든 고유한 단어의 목록을 만들고, 각 단어가 발생하는 모든 문서를 식별
- java는 컴파일 언어이며, python보다 좀 더 빠르기 때문에 java를 사용

## ● 역색인 코드

```
1 .io.*;
2 .util.*;
3 .util.Map.Entry;
4 .nio.file.*;
5
6
7 s aufgabel {
8 ic static void main(String[] args) throws IOException {
9
10 system.out.println("loading "+args[0]+" file");
11 ong start = System.currentTimeMillis();
12 /바이트사인 포맷크기해 알람을 받을 포맷크기는 약 35000byte이고 전체를 받으려면 파일을것으로 생각됨
13 BufferedReader reader = new BufferedReader(
14     new FileReader("C:\\Users\\admin\\Downloads\\"+args[0]),40960
15 );
16 Path path = Paths.get("C:\\Users\\admin\\Downloads\\"+args[0]);
17 //문장의 개수를 미리 알면 속도향상해 도움이 될것으로 생각해서 문장의 개수 확인
18 int lineCount = (int) Files.lines(path).count();
19 String[] saetze = new String[lineCount];
20 final int [] count = {0};
21 //읽은 문장들을 배열해 저장
22 Files.lines(path).forEach(line -> saetze[count[0]++] = line);
23 reader.close();
24
25 long end = System.currentTimeMillis();
26 double time = (end - start)/1000.0;
27 System.out.println(String.format("Complete! (%.3fs)", time));
28
29 //유틸리티를 사용하면 기존보다 더 편리할것으로 생각됨
30 //containskey 보다 get함수가 빠르지만 큰 차이 없음
31 long start1 = System.currentTimeMillis();
32
33
34 //해쉬맵을 2중으로 사용
35 HashMap<String, HashMap<Integer,Integer>> doc = new HashMap<String, HashMap<Integer,Integer>>();
36
37 //문장들을 배열처리
38 Iterator<String> iterator = Arrays.stream(saetze).parallel().iterator();
39 while (iterator.hasNext()) {
40     String satz = iterator.next();
41     String[] temp = satz.toLowerCase().replaceAll("[^a-z0-9]", " ").split(" ");
42     temp = Arrays.stream(temp).filter(s -> !s.isEmpty()).toArray(String[]::new);
```

```
40 String satz = iterator.next();
41 String[] temp = satz.toLowerCase().replaceAll("[^a-z0-9]", " ").split(" ");
42 temp = Arrays.stream(temp).filter(s -> !s.isEmpty()).toArray(String[]::new);
43 int doc_id = Integer.parseInt(temp[0]);
44 for (int i = 1; i < temp.length; i++) {
45     String word = temp[i];
46
47     if(!doc.containsKey(word)) {
48
49         doc.put(word, new HashMap<Integer,Integer>());
50
51     }
52
53     HashMap<Integer,Integer> innerMap = doc.get(word);
54     innerMap.put(doc_id, innerMap.getOrDefault(doc_id, 0) + 1);
55
56 }
57
58 }
59
60
61 }
62 Comparator<HashMap.Entry<Integer, Integer>> valueComparator = new Comparator<HashMap.Entry<Integer, Integer>>() {
63     @Override
64     public int compare(Entry<Integer, Integer> o1, Entry<Integer, Integer> o2) {
65         int compare = o2.getValue().compareTo(o1.getValue()); // comparing values in descending order
66         if (compare == 0) {
67             return o1.getKey().compareTo(o2.getKey()); // comparing keys in ascending order
68         }
69         return compare;
70     }
71 }
72
73 HashMap<String, LinkedHashMap<Integer,Integer>> sortedDoc = new HashMap<>();
74 Set<String> keySet = doc.keySet();
75 List<String> keyList = new ArrayList<>(keySet);
76 Collections.sort(keyList);
77 for (String key : keyList) {
78     HashMap<Integer,Integer> innerMap = doc.get(key);
79     List<Map.Entry<Integer, Integer>> list = new ArrayList<>(innerMap.entrySet());
80     // comparator로 정렬
81     Collections.sort(list,valueComparator);
82     // comparator로 정렬된 결과를 새로운 해쉬맵에 저장
83     LinkedHashMap<Integer,Integer> innerMapSorted = new LinkedHashMap<>();
84     for(HashMap.Entry<Integer, Integer> innerEntry : list) {
```

```
81 Collections.sort(list,valueComparator);
82 // comparator로 정렬된 결과를 새로운 해쉬맵에 저장
83 LinkedHashMap<Integer,Integer> innerMapSorted = new LinkedHashMap<>();
84 for(HashMap.Entry<Integer, Integer> innerEntry : list) {
85     innerMapSorted.put(innerEntry.getKey(), innerEntry.getValue());
86 }
87 sortedDoc.put(key, innerMapSorted);
88 }
89 /*
90 Set<String> li = sortedDoc.keySet();
91 for (String w : li) {
92     LinkedHashMap<Integer, Integer> inn = sortedDoc.get(w);
93     Set<Integer> inner = inn.keySet();
94     for (Integer docId : inner) {
95         System.out.println(w+" "+ docId + " " + inn.get(docId) + " ");
96     }
97 }
98 */
99
100
101 BufferedWriter bw = null;
102 try {
103     bw = new BufferedWriter(new FileWriter("sortedDoc.txt",1024);
104     List<String> sortedDocKeySet = keyList;
105
106     for (String word : sortedDocKeySet) {
107         LinkedHashMap<Integer, Integer> innerMap = sortedDoc.get(word);
108         Set<Integer> innerMapKeySet = innerMap.keySet();
109         for (Integer docId : innerMapKeySet) {
110             bw.write(word+" "+ docId + " " + innerMap.get(docId) + " ");
111         }
112         bw.newLine();
113         bw.flush();
114     }
115 } catch (IOException e) {
116     e.printStackTrace();
117 } finally {
118     try {
119         if (bw != null)
120             bw.close();
121     } catch (IOException e) {
122         e.printStackTrace();
123     }
124 }
```

# # ‘역색인’ 이란?

## 5. 검색 엔진 연동

### ● 역색인 (Inverted index) < 역색인 결과 >

- 매우 빠른 풀텍스트
- 문서에 나타나는 모든
- 서를 식별
- java 는 컴파일 언어

### ● 역색인 코드

```
1 .io.*;
2 .util.*;
3 .util.Map.Entry;
4 .nio.file.*;
5
6
7 s aufgabel {
8 ic static void main(String[] args) throws IOException {
9
10 ystem.out.println("loading "+args[0]+" file");
11 ong start = System.currentTimeMillis();
12 /비파괴적인 파일의 내용을 불러 올때까지는 약 35000byte이고 내용을 불러오면
13 BufferedReader reader = new BufferedReader(
14     new FileReader("C:\\Users\\admin\\Downloads\\"+args[0],40960
15 ));
16 Path path = Paths.get("C:\\Users\\admin\\Downloads\\"+args[0]);
17 //문서의 개수를 읽어 보기 속도향상책 도움이 됨안으로 생각해서 문장의 경우
18 int lineCount = (int) Files.lines(path).count();
19 String[] saetze = new String[lineCount];
20 final int [] count = {0};
21 //읽은 문장들을 배열에 저장
22 Files.lines(path).forEach(line -> saetze[count[0]++] = line);
23 reader.close();
24
25 long end = System.currentTimeMillis();
26 double time = (end - start)/1000.0;
27 System.out.println(String.format("Complete! (%.3fs)", time));
28
29 //문장들을 사분할까? 아니면 다 할까? 일단은 사분할
30 //containskey 보다가 get해서가 빠르거든 큰 차이 없음
31 long start1 = System.currentTimeMillis();
32
33
34 //해시맵을 이용해서 저장
35 HashMap<String, HashMap<Integer, Integer>> doc = new HashMap<String, HashMap<Integer, Integer>>();
36
37 //문장들을 배열에서
38 Iterator<String> iterator = Arrays.stream(saetze).parallel().iterator();
39 while (iterator.hasNext()) {
40     String satz = iterator.next();
41     String[] temp = satz.toLowerCase().replaceAll("[^a-z0-9]", " ").split(" ");
42     temp = Arrays.stream(temp).filter(s -> !s.isEmpty()).toArray(String[]::new);
```

```
loading input.big file
Complete! (0.176s)
Complete! (5.762s)
total Complete! (5.938s)
input : frowning
Doc: 11396, Freuquency: 1
Doc: 82860, Freuquency: 1
Doc: 104083, Freuquency: 1
Doc: 249603, Freuquency: 1
```

```
input : happy
Doc: 64017, Freuquency: 3
Doc: 8929, Freuquency: 2
Doc: 68177, Freuquency: 2
Doc: 71409, Freuquency: 2
Doc: 122598, Freuquency: 2
Doc: 125899, Freuquency: 2
Doc: 143714, Freuquency: 2
Doc: 236782, Freuquency: 2
Doc: 262076, Freuquency: 2
Doc: 911, Freuquency: 1
Doc: 1214, Freuquency: 1
Doc: 1428, Freuquency: 1
Doc: 4452, Freuquency: 1
Doc: 6080, Freuquency: 1
Doc: 6859, Freuquency: 1
Doc: 7005, Freuquency: 1
Doc: 7902, Freuquency: 1
Doc: 8201, Freuquency: 1
Doc: 8333, Freuquency: 1
Doc: 9214, Freuquency: 1
Doc: 9686, Freuquency: 1
Doc: 11735, Freuquency: 1
Doc: 12324, Freuquency: 1
Doc: 14368, Freuquency: 1
Doc: 14374, Freuquency: 1
Doc: 14519, Freuquency: 1
Doc: 15320, Freuquency: 1
Doc: 16165, Freuquency: 1
Doc: 16361, Freuquency: 1
```

```
119 Collections.sort(list,valueComparator);
120 // comparator로 정렬된 결과를 새로운 해시맵에 저장
121 LinkedHashMap<Integer, Integer> innerMapSorted = new LinkedHashMap<>();
122 for(HashMap.Entry<Integer, Integer> innerEntry : list) {
```

```
actions.sort(list,valueComparator);
120 comparator로 정렬된 결과를 새로운 해시맵에 저장
121 LinkedHashMap<Integer, Integer> innerMapSorted = new LinkedHashMap<>();
122 for(HashMap.Entry<Integer, Integer> innerEntry : list) {
123     innerMapSorted.put(innerEntry.getKey(), innerEntry.getValue());
124 }
125 }
126 }
127 }
128 }
129 }
130 }
131 }
132 }
133 }
134 }
135 }
136 }
137 }
138 }
139 }
140 }
141 }
142 }
143 }
144 }
145 }
146 }
147 }
148 }
149 }
150 }
151 }
152 }
153 }
154 }
155 }
156 }
157 }
158 }
159 }
160 }
161 }
162 }
163 }
164 }
165 }
166 }
167 }
168 }
169 }
170 }
171 }
172 }
173 }
174 }
175 }
176 }
177 }
178 }
179 }
180 }
181 }
182 }
183 }
184 }
185 }
186 }
187 }
188 }
189 }
190 }
191 }
192 }
193 }
194 }
195 }
196 }
197 }
198 }
199 }
200 }
```

# # Elastic search 란 ?

## 5. 검색 엔진 연동



Elastic Search

- 텍스트, 숫자, 위치 기반 정보, 정형 및 비정형 데이터 등 모든 유형의 데이터를 위한 무료 검색 및 분석 엔진

### 실시간 검색 플랫폼

문서가 색인될 때부터 검색 가능해질 때까지의 대기 시간이 매우 짧음

### 분산적

Elastic search에 저장된 문서는 '샤드' (여러 다른 컨테이너)에 걸쳐 분산되며,

샤드는 복제되어 하드웨어 장애 시에 중복되는 데이터 사본을 제공함

### 강력한 기본기능

속도, 확장성, 복원력 뿐만 아니라 '데이터 롤업, 인덱스 수명 주기 관리' 등과 같이 데이터를 더욱 효율적으로 저장 및 검색할 수 있는 강력한 기본기능 탑재

### 데이터 처리 간소화

Beats 와 Logstash 는 Elastic search 로 색인하기 전 데이터를 쉽게 처리할 수 있게 함



1

REST API  
&  
elastic search  
&  
DPR  
연동



2

raw data  
embedding  
&  
elastic server  
es.index



3

elastic search  
에서 형태소분석  
기를 사용한  
question  
sentence 의  
주요 명사 추출



4

es.search 사용  
하여  
문서 가져오기  
&  
유사도 비교



5

rank 로  
n-best 추출





# Web Crawling - www.musinsa.com

5. 검색 엔진 연동

품목

브랜드

인기 Best

니트/스웨터 (19,093)

후드 집업 (4,490)

데님 팬츠 (12,544)

백팩 (5,321)

후드 티셔츠 (21,500)

메신저/크로스백 (8,395)

패션스니커즈화 (2,305)

맨투맨/스웨트셔츠 (35,510)

슈트 팬츠/슬랙스 (7,986)

카디건 (9,639)

스타디움 재킷 (1,082)

숏패딩/숏에비 아우터 (7,445)

상의 Top

아우터 Outer

바지 Pants

원피스 Onepiece

스커트 Skirt

스니커즈 Sneakers

신발 Shoes

가방 Bag

여성 가방 Women's bag

스포츠/용품 Sports/Goods

모자 Headwear

양말/레그웨어 Socks/Legwear

속옷 Underwear

선글라스/안경테 Eyewear

엑세서리 Accessory

시계 Watch

주얼리 Jewelry

뷰티 Beauty

무신사 스토어 > 상품 랭킹

all

Ranking Shop

무신사 랭킹은 상품 매출, 판매 수량, 상품 조회 수, 작성 후기 수를 반영한 공식에 의해 선정됩니다. 무신사 스토어는 광고 목적으로 랭킹을 절대 임의 조작하지 않으므로 믿고 구매하셔도 됩니다.

상품

브랜드

검색어

무신사 랭킹을 알려드립니다.

자세히 보기

기간 구분	실시간	일간	주간	월간	3개월					
대분류	전체	상의	아우터	바지	원피스	스커트	스니커즈	신발	가방	여성 가방
		스포츠/용품	모자	양말/레그웨어	속옷	선글라스/안경테	엑세서리	시계	주얼리	뷰티
		디지털/테크	리빙	책/음악/티켓	번려동물					

가격

전체

5만원 이하

5~10만원

10~20만원

20~30만원

30만원 이상

원 ~ 원

검색

NEW X

5분 전 갱신

11

NEW

단독

세일

품질 포함

골프

키즈

상품 수 : 182,035개

100 페이지 중 1 페이지

1위

한정 판매

분크  
Occam Doux Shoulder XL (오리 두 숏더 엑스라)  
475,000원  
쿠폰 -47,500원  
MEMBERSHIP PRICE  
★★★★★ 8  
685

W

OPTION ▼

10위

부티크

2위

한정 판매

테어드  
2 / 28 배송 PURPLE TWEED BLOUSON  
268,000원  
MEMBERSHIP PRICE  
29

W

OPTION ▼

11위

3위

무신사 단독

아식스  
젤-1090 V2 SMU - 폴라 웨이드크림 /  
99,000원  
MEMBERSHIP PRICE  
★★★★★ 328  
5,997

MW

OPTION ▼

12위

부티크

4위

부티크

디스커버리 익스페디션  
라이크 에어 시프트 백팩 (BLACK)  
169,000원  
MEMBERSHIP PRICE  
★★★★★ 342  
1,976

MW

OPTION ▼

13위

5위

부티크

질 샌디  
새벽배송 여성 스몰 카슬로 숏더백 - 올리브 /  
2,583,888원  
1,036,000원  
MEMBERSHIP PRICE  
13

W

OPTION ▼

14위

한정 판매

6위

무신사 단독

아디다스  
센테니얼 85 로우 - 화이트:블루 / IF5419  
139,000원  
쿠폰 -13,900원  
MEMBERSHIP PRICE  
★★★★★ 169  
18,398

MW

OPTION ▼

15위

7위

한정 판매

인사일런스  
커브드 라인 울 블레이지 MOCHA  
222,988원  
188,000원  
쿠폰 -9,400원  
MEMBERSHIP PRICE  
★★★★★ 3  
1,012

M

OPTION ▼

16위

8위

한정 판매

에프씨엔엠  
2 / 24 배송 [오링X에프씨엔엠] 헤비웨이트 후디 -  
75,000원  
쿠폰 -7,500원  
MEMBERSHIP PRICE  
269

MW

OPTION ▼

17위

9위

한정 판매

가민  
인스팅트 2 솔라 택티컬 코요테탄  
599,000원  
MEMBERSHIP PRICE  
19

L

OPTION ▼

18위

전체

남성

여성

장바구니

로그인

회원가입

하위메뉴

더보기

# Web Crawling - www.musinsa.com

5. 검색 엔진 연동

품목

브랜드

무신사 스토어 > 상품 랭킹

Ranking Shop

무신사 랭킹은 상품 매출, 판매 수량, 상품 조회 수, 적성 후기 수를 반영한 공식에 의해 선정됩니다.  
무신사 스토어는 광고 목적으로 랭킹을 절대 임의 조작하지 않으므로 믿고 구매하셔도 됩니다.

상품

신상품

검색어

인기 Best

니트/스웨터 (5,093)

무드 집업 (4,490)

데님 팬츠 (12,544)

백팩 (5,321)

후드 티셔츠 (21,500)

메신저/크로스백 (8,395)

상의 Top

아우터 Outer

바지 Pants

원피스 Onepiece

스커트 Skirt

스니커즈 Sneakers

신발 Shoes

가방 Bag

여성 가방 Women's bag

스포츠/용품 Sports/Goods

모자 Headwear

양말/레그웨어 Socks/Legwear

속옷 Underwear

선글라스/안경테 Eyewear

액세서리 Accessory

시계 Watch

주얼리 Jewelry

뷰티 Beauty

(23SS) 2 TONE ARCH TEE WHITE

XL 구매

★★★★★

적당한 길이라 좋아요  
예일 글자 포인트도 좋습니다  
목이 좀 높아났지만 요정도만 괜찮은 정도?

사이즈 보통이에요

밝기 보통이에요

색감 보통이에요

두께

도움돼요 0

LV 4 mind3

남성 · 170cm · 65kg · 신고

(23SS) 2 TONE ARCH TEE WHITE

L 구매

★★★★★

예일 화이트 반팔티 기본으로 입기 편하고 좋아요~

사이즈 보통이에요

밝기 밝아요

색감 선명해요

두께감 보통이에요

도움돼요 0

LV 3 빙빙빙글

남성 · 173cm · 68kg · 신고

(23SS) 2 TONE ARCH TEE WHITE

L 구매

★★★★★

깨끗하게 하얀 느낌이에요. 시원해보이고 사이즈도 좋아요

분크

Occam Doux Shoulder XL (오리 두 손디 엑스라 475,000원

MEMBERSHIP PRICE 47,500원

★★★★★ 685

테어드

2/28 배송 PURPLE TWEED BLOUSON 268,000원

MEMBERSHIP PRICE 29

아식스

젤-1090 V2 SMU - 폴라 웨이드크림 / 99,000원

MEMBERSHIP PRICE 5,997

디스커버리 엑스페디션 라이크 에어 시프트 백팩 (BLACK) 169,000원

MEMBERSHIP PRICE 1,976

질 샌디

새벽배송 여성 스몰 카놀로 손디백 - 올리브 / 2,583,000원

MEMBERSHIP PRICE 13

10위

부티크

11위

부티크

12위

부티크

13위

부티크

14위

안정 판매

15위

부티크

16위

부티크

17위

부티크

18위

부티크

원 검색

100 페이지 중 1 페이지

7위

안정 판매

8위

안정 판매

9위

전체

남성

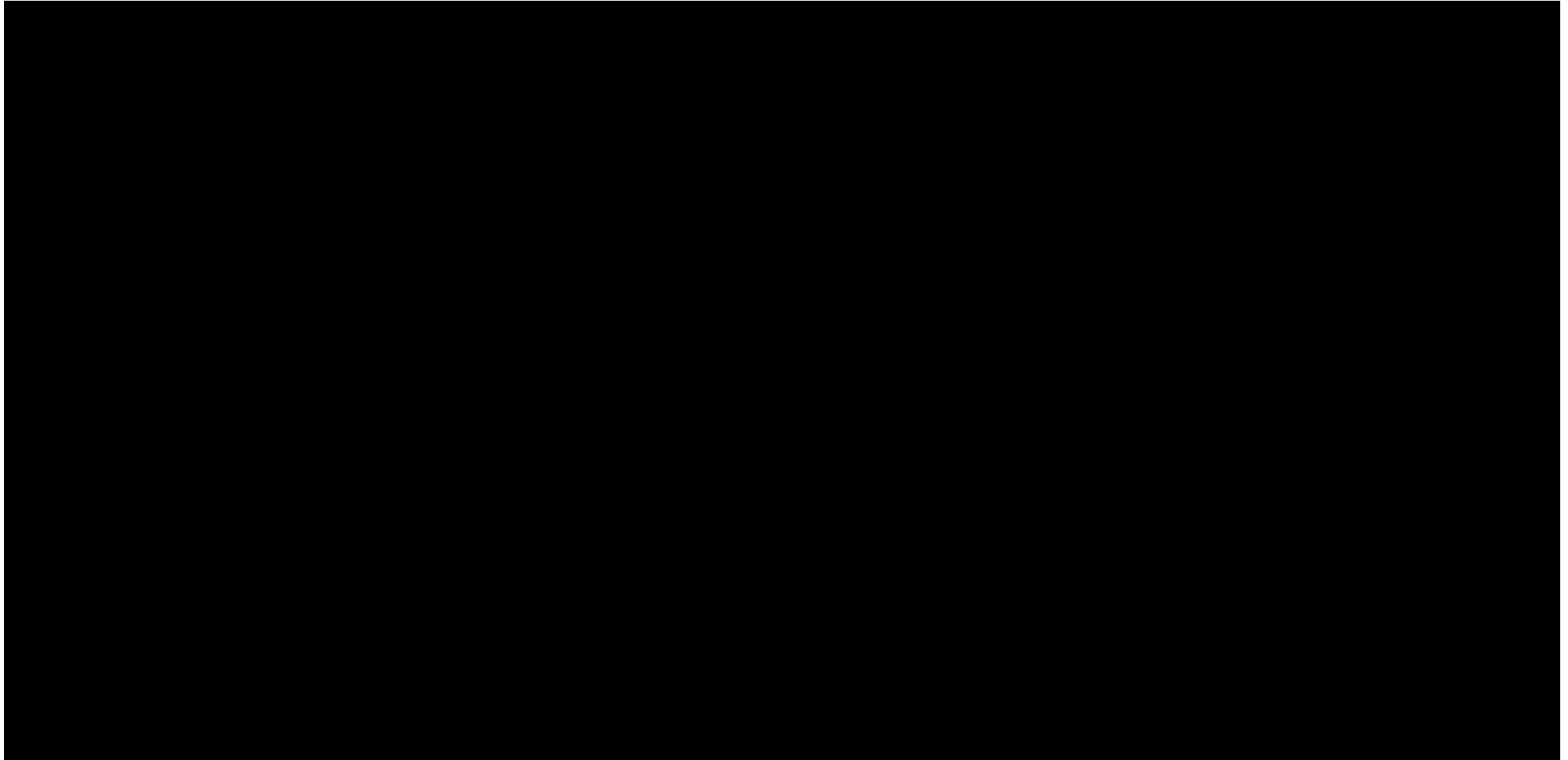
여성

옵션

옵션

## # Demonstration

## 5. 검색 엔진 연동



# 6

## Conclusion

---

### 기대효과

- 검색 품질의 개선 : dense retrieval 모델을 이용한 품질개선
- 검색 효율성 향상 : 여러개의 문서를 병렬로 처리 가능
- 기존 시스템 개선 : 기존 시스템의 검색 품질 및 효율성 개선
- 새로운 비즈니스 모델 개발 : 검색 엔진 서비스 제공

### 한계점

- gpu 서버가 없어 모델 훈련에 많은시간이 소모됨
- 한국어 데이터셋 부족
- 인터프리터 언어인 파이썬은 컴파일 언어에 비해 느린 속도
- gpu 사용시에도 문장 임베딩 시 오랜시간 소요

# Conclusion

### 개선점

- 로그인 기능을 구현 및 클라이언트 세션 구현
- 모델 하이퍼 파라미터 튜닝
- 엘라스틱 서치 / 플라스크 서버 보안설정
- 유사도 함수를 종합해서 사용하여 엔진 성능 개선

**Q & A**

# Thank You

