

# Unsupervised Document Expansion for Information Retrieval with Stochastic Text Generation

Soyeong Jeong<sup>1</sup>, Jinheon Baek<sup>2</sup>, ChaeHun Park<sup>1</sup>, Jong C. Park<sup>1\*</sup>

School of Computing<sup>1</sup>, Graduate School of AI<sup>2</sup>, KAIST<sup>1,2</sup>, South Korea

{syjeong, ddehun, park}@nlp.kaist.ac.kr, jinheon.baek@kaist.ac.kr

## Motivation

One of the challenges in information retrieval (IR) is the *vocabulary mismatch problem*, which happens when the terms between queries and documents are lexically different but semantically similar. While recent work has proposed to expand the queries or documents by enriching their representations with additional relevant terms or densely represent documents to address this challenge, they usually require a large volume of query-document pairs to train an expansion model.

## UDEG Framework

To tackle the vocabulary mismatch problem, we propose an Unsupervised Document Expansion with Generation (UDEG) framework with a pre-trained language model, which generates diverse supplementary sentences for the original document without using labels on query-document pairs for training. We first generate document-related sentences with a pre-trained language model, already fine-tuned on a summarization task. However, such a framework generates only one static sentence at a time, so we further propose to stochastically generate multiple sentences which reflect diverse points of view for the given document and minimize the vocabulary mismatch cases (See Figure 1).

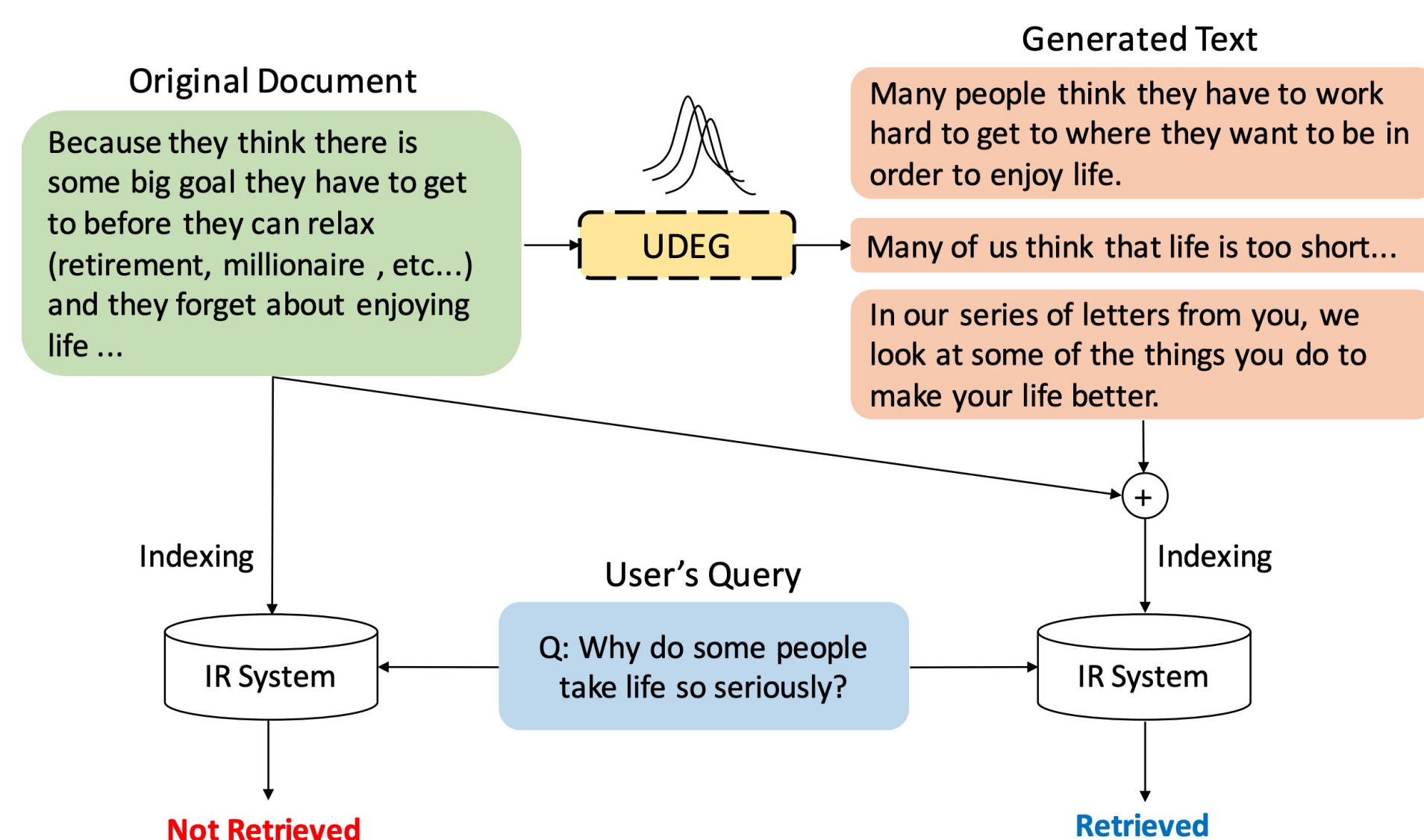


Figure 1. UDEG framework, where the example is generated from UDEG. Given an original document (green box), our UDEG framework stochastically generates several sentences (orange box) relevant to the given document, and augments the generated sentences to the input document.

## Baseline Expansion Models

- No Expansion
- Query Expansion Model
  - RM3
- Extractive Document Expansion Models
  - MP-rank
  - LexRank
  - Pegasus<sub>ext</sub>
- Abstractive Document Expansion Models
  - LexRank + paraphrase
  - UDEG

## Main Results

UDEG outperforms other baselines on the ANTIQUE dataset and sampled MS MARCO dataset (See Table 1 and Table 2).

		No Expan.	MP-rank	LexRank	Lex.+Para.	Pegasus <sub>ext</sub>	UDEG
MRR	BM25	0.595	0.584	0.571	0.561	0.585	<b>0.645</b>
	BM25+RM3	0.558	0.579	0.542	0.567	0.555	<b>0.616</b>
	QL	0.499	0.534	0.567	0.518	0.562	<b>0.650</b>
	QL+RM3	0.396	0.447	0.456	0.432	0.504	<b>0.583</b>
R@10	BM25	0.218	0.220	0.208	0.209	0.207	<b>0.237</b>
	BM25+RM3	0.217	0.221	0.208	0.204	0.213	<b>0.226</b>
	QL	0.189	0.199	0.203	0.196	0.205	<b>0.232</b>
	QL+RM3	0.159	0.179	0.182	0.162	0.191	<b>0.211</b>
P@3	BM25	0.378	0.381	0.346	0.351	0.356	<b>0.431</b>
	BM25+RM3	0.361	0.355	0.360	0.373	0.366	<b>0.433</b>
	QL	0.301	0.333	0.340	0.315	0.358	<b>0.418</b>
	QL+RM3	0.240	0.281	0.275	0.271	0.301	<b>0.386</b>
MAP	BM25	0.211	0.212	0.199	0.202	0.201	<b>0.238</b>
	BM25+RM3	0.212	0.213	0.203	0.203	0.207	<b>0.234</b>
	QL	0.172	0.191	0.192	0.181	0.199	<b>0.230</b>
	QL+RM3	0.150	0.168	0.170	0.158	0.180	<b>0.212</b>
NDCG@3	BM25	0.437	0.442	0.417	0.425	0.419	<b>0.478</b>
	BM25+RM3	0.424	0.434	0.423	0.433	0.426	<b>0.470</b>
	QL	0.356	0.389	0.400	0.375	0.413	<b>0.471</b>
	QL+RM3	0.277	0.324	0.319	0.306	0.350	<b>0.424</b>

Table 1. Retrieval results on the ANTIQUE dataset.

		No Expan.	LexRank	UDEG
MRR	BM25	0.427	0.441	<b>0.463</b>
	BM25+RM3	0.366	0.385	<b>0.415</b>
	QL	0.402	0.420	<b>0.454</b>
	QL+RM3	0.319	0.337	<b>0.382</b>
R@10	BM25	0.636	0.646	<b>0.679</b>
	BM25+RM3	0.600	0.617	<b>0.651</b>
	QL	0.611	0.633	<b>0.671</b>
	QL+RM3	0.552	0.579	<b>0.629</b>
P@1	BM25	0.311	0.324	<b>0.344</b>
	BM25+RM3	0.248	0.265	<b>0.291</b>
	QL	0.289	0.302	<b>0.334</b>
	QL+RM3	0.202	0.215	<b>0.255</b>
MAP	BM25	0.422	0.435	<b>0.457</b>
	BM25+RM3	0.361	0.380	<b>0.409</b>
	QL	0.398	0.414	<b>0.448</b>
	QL+RM3	0.315	0.333	<b>0.377</b>

Table 2. Retrieval results on the MS MARCO dataset.

## Ablation Study

- Robustness on different LMs (See Figure 2).
- Comparison of Stochastic Generation Strategy (See Figure 3).
- Varying the Number of Expanded Sentences (See Figure 3).

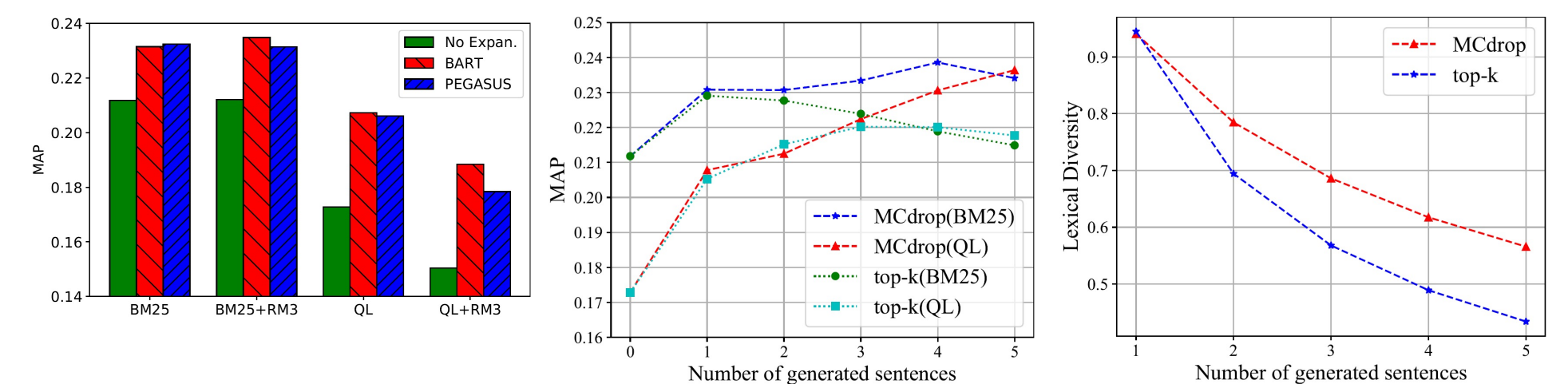


Figure 2. Comparison of BART and PEGASUS.

Figure 3. MC dropout vs top-k sampling (Left) with a varying number of generated sentences and (Right) with a lexical diversity.

## Case Study

UDEG successfully retrieves documents (See Table 3).

Query	How is the chemistry is a basic of science?
Relevant Document	Chemistry is a basic because all matter can be broken down into elements (i.e., hydrogen, oxygen, nitrogen, etc.); without matter, nothing could be studied.
Generated Sentences	1) Chemistry is the study of <b>atoms</b> and <b>molecules</b> . 2) Chemistry is the study of matter and how it is made. 3) Chemistry is the <b>study</b> of matter. 4) Chemistry is a basic <b>science</b> .
	Original Document Rank: 104      Expanded Document Rank: 5
Query	How is the library consider as a heart of university?
Relevant Document	Whatever you are studying has to be found somewhere for you to learn it. That's where the library comes into focus.
Generated Sentences	1) If you're studying at <b>university</b> , you'll need a library. 2) A library is a <b>place</b> where you can <b>find</b> out more about the <b>subject</b> you are studying. 3) If you're studying, you'll be studying. 4) There are <b>many different ways</b> you can study.
	Original Document Rank: 636      Expanded Document Rank: 32
Query	What do doctors do when a patient has a Do Not Resuscitate Order?
Relevant Document	All healthcare professionals involved in the care of that patient will not do anything to prolong the patient's life if in case patient deteriorates/dies. DNR orders may be modified, some may choose mechanical ventilation, or drugs. Usually when a pt is DNR, comfort measures is provided only.
Generated Sentences	1) DNR is not <b>life-support</b> . 2) When a patient is in a "do not resuscitated" (DNR) state, that patient's life will not be <b>saved</b> . 3) A DNR is a <b>decision</b> made by the patient's <b>family</b> or health care provider to prolong the life of the patient. 4) A "do not resuscitate"(DNR) order does not <b>mean</b> that a patient should be put on <b>life support</b> .
	Original Document Rank: 40      Expanded Document Rank: 1

Table 3. Examples of generated sentences on the ANTIQUE dataset.

## Contributions

- To mitigate the vocabulary mismatch problem, we present UDEG framework that augments a document with abstractly generated sentences without paired query-document data.
- Under an unsupervised document expansion framework, we generate document-related sentences with a pre-trained LM, and further stochastically generate diverse sentences.
- We show that our framework achieves outstanding performances on benchmark datasets for IR tasks.