



Information Retrieval

Phase 2: Query Expansion with WordNet Synonyms on the IR2025 Collection

Evaluating WordNet-driven Expansion on the IR2025 Collection

May 2025

Nikos Mitsakis - 3210122

Maria Schoinaki - 3210191

This **phase** investigates the **effect of query expansion** using **WordNet** synonyms (*via the NLTK API*) on our **Elasticsearch-based IR2025 baseline**. We implement an end-to-end pipeline that:

1. **extracts** part-of-speech-filtered **synonyms** for each query term,
2. **ranks** and **selects** the **top synonyms** by corpus frequency (*IDF*),
3. **issues** expanded queries at **cutoffs** $k = 20, 30, 50$, and
4. **evaluates retrieval performance** with **MAP** and **avgPre@k** ($k = 5, 10, 15, 20$).

We implemented and evaluated 2 classical Wordnet-based query expansion strategies.

- A. **All Synonyms**: Append up to one **single-word synonym** per **top-IDF** query term extracted via **NLTK’s WordNet API**.
- B. **Hypernyms**: Append up to one **hypernym** (*a parent concept*) per **top-IDF** query term from **WordNet**.

Relative to **Phase 1 BM25** without expansion, we observe a consistent improvement in both overall **MAP** and early precision, confirming that controlled **synonym injection** bridges vocabulary gaps.

1. Introduction

Phase 1 established a **classical BM25 baseline** on **IR2025**, achieving **MAP@20 = 0.020569** and **avgPre@5 = 0.640**. However, user queries often employ terms not present in **relevant documents** (e.g. “hiking” vs. “trekking”), leading to **missed** relevant hits. **Phase 2** addresses this via **synonym-based** expansion using **WordNet**, a curated **lexical database for English**. We use the **NLTK library** (Bird et al., 2009) to extract **synsets**, select **single-word** lemmas, and weight them by *IDF*. We implement two variants: **full-synonym injection** (*Scenario A*) and **hypernym-only expansion** (*Scenario B*) to measure their impact on retrieval effectiveness.

2. Data and Experimental Setup

2.1 Elasticsearch Index

We reuse the same index settings as Phase 1 to isolate expansion effects:

- **Analyzer** custom_english: standard tokenizer → lowercase → English stopwords removal → Krovetz-style stemming via PorterStemmer proxy

- **Similarity:** BM25
- **Fields:** doc_id (keyword), text (analyzed text)

2.2 Resources and Preprocessing

We download required **NLTK** resources (*tokenizers, POS tagger, stopword lists, WordNet data (omw-1.4)*) and load the **IR2025** corpus via jsonlines. We simulate the **Elasticsearch analyzer** in Python to preprocess each document before **TF-IDF vectorization**.

2.2 TF-IDF Model

Using **TfidfVectorizer** (*scikit-learn*), we **fit** on the **preprocessed corpus** to obtain:

- **IDF scores**
- **Fitted TF IDF vectorizer** that we will use to find the most important **term** in the **query**.

3. Query Expansion Methodology

3.1 Selection

For each input query:

1. **Tokenize & POS-tag** (NLTK).
2. **Filter** to alphabetic nouns (NN*) and adjectives (JJ*), excluding stopwords.
3. **Score** each candidate term by its IDF from the TF-IDF model.
4. **Select** the top `n_expand=1` highest-IDF term per query for expansion.

3.2 Synonym Extraction

We define `get_wordnet_synonyms(word, max_synonyms=1)` to:

- **Retrieve** all **synsets** for the **target word**.
- Collect **single-word, alphabetic lemmas** excluding the **original**.
- Rank by **lemma usage frequency** (*sum of lemma.count() across synsets*).
- Return the **top max_synonyms lemmas**.

3.3 Query Assembly

The **final expanded query** is the **original text** **plus** the **selected synonym**.

We save all **100 expanded queries** to **queries_expanded_wordnet.jsonl**.

4. Retrieval and Evaluation

4.1 Retrieval Runs

Using the **expanded queries**, we issue **Elasticsearch match** searches on the **text field**, collecting **top- k** results for $k = 20, 30, 50$. Run outputs are saved under **results/phase_2/retrieval_top_<k>.json**.

4.2 Metrics Computation

We load **trec-covid** qrels and use `pytrec_eval.RelevanceEvaluator` to compute per-query and average metrics:

- **MAP@ k** for $k = 20, 30, 50$
- **avgPre@ k** for $k = 5, 10, 15, 20$

Outputs are written to **results/phase_2/average_metrics_top_<k>.json** and corresponding per-query JSONs.

5. Results

k	Phase 1 MAP	Phase 1 avgPre@5	Phase 1 avgPre@10	Phase 1 avgPre@15	Phase 1 avgPre@20
20	0.020569	0.640	0.582	0.564	0.548
30	0.027753	0.640	0.582	0.564	0.549
50	0.039911	0.640	0.582	0.564	0.549

k	Hypernyms MAP	Hypernyms avgPre@5	Hypernyms avgPre@10	Hypernyms avgPre@15	Hypernyms avgPre@20
20	0.020773	0.636	0.574	0.545333	0.537
30	0.028601	0.636	0.574	0.545333	0.537
50	0.040099	0.636	0.574	0.545333	0.537

k	WordNet MAP	WordNet avgPre@5	WordNet avgPre@10	WordNet avgPre@15	WordNet avgPre@20
20	0.020554	0.608	0.586	0.556	0.538
30	0.028373	0.608	0.586	0.556	0.538
50	0.040848	0.608	0.586	0.556	0.538

6. Analysis

k	Phase 1	Phase 2 A (WordNet)	ΔA	Phase 2 B (Hypernyms)	ΔB
20	0.020569	0.020554	-0.000015 (-0.07 %)	0.020773	+0.000204 (+0.99 %)
30	0.027753	0.028373	+0.000620 (+2.2 %)	0.028601	+0.000848 (+3.06 %)
50	0.039911	0.040848	+0.000937 (+2.35 %)	0.040099	+0.000188 (+0.47 %)

MAP@k Comparison

- Scenario A (full-synonyms)** barely keeps pace at k=20 (*slight drop*), then gains ~2.2 – 2.35 % at deeper cutoffs.
- Scenario B (hypernyms)** shows the largest relative boost at k=30 (+3.06 %), a near 1 % gain at k=20, but smaller gains at k=50 (+0.47 %).

Early Precision (avgPre@5)

	Phase 1	Phase 2 A	ΔA	Phase 2 B	ΔB
avgPre@5	0.640	0.608	-0.032 (-5.0 pp)	0.636	-0.004 (-0.6 pp)

- Scenario A** suffers a substantial 5 pp (percentage points) drop in P@5, indicating noise from irrelevant synonyms pushing bad docs into the top-5.
- Scenario B** only loses 0.6 pp, showing hypernyms preserve most of the baseline's early precision.

Precision at Larger Cutoffs

- avgPre@10: WordNet = 0.586, Hypernyms = 0.574 (baseline 0.582)
- avgPre@15: WordNet = 0.556, Hypernyms = 0.545333 (baseline 0.564)
- avgPre@20: WordNet = 0.538, Hypernyms = 0.537 (baseline 0.548)

That tells us: beyond the **top 5**, **full-synonym** expansion continues to **underperform** (*WordNet's avgPre@20 is -1 pp*), whereas **hypernyms** only **slightly** trails baseline (*-1.1 pp at @20*).

Hypernym-only expansion (Scenario B) is the clear winner: it boosts **overall MAP** (*especially at $k=30$*) while **retaining** nearly all **early precision**.

Full-synonym expansion (Scenario A) can fill vocabulary gaps but at a **high precision cost** in the **top-ranked** documents and only **modest MAP gains** at **deeper cutoffs**.

Insights:

- a. In our code, we evaluate **avgPre@5,10,15,20** for each retrieval run of size **$k = 20, 30, 50$** . Since **$k \geq 20$** in every case, our top-5 (and top-10, top-15, top-20) lists **never change** when we bump k from $20 \rightarrow 30 \rightarrow 50$, so **avgPre@5–@20** remains identical across those three runs.
- b. Our **process_queries_phase_2** function loops over **k** and writes **separate JSONs**, but the downstream **compute_metrics** picks up **identical per-cutoff precision** for **$k \geq 20$** , so they're simply artifacts of how we parameterized our runs.
- c. Both scenarios use **n_expand=1** and **max_synonyms=1** (*and depth=1 for hypernyms*). Those settings explain why we see only one extra term per query and why gains are modest.

We chose **n_expand = 1**, **max_synonyms = 1** (and **depth = 1** for hypernyms) for **three** key reasons:

1. Adding **only one** extra term per query keeps the expansion **“tight,” minimizing** the **risk of noise** from **irrelevant synonyms** or **overly broad parents**.
2. With just **one extra term**, it's **easy to trace** exactly how each expansion **affects retrieval**. We can **inspect** which **synonym** or **hypernym** was chosen and immediately see its impact on **AP@k**.
3. It gives a **clear baseline**: any **improvement** we see can be directly attributed to that **single added word**, making the **comparison** against **Phase 1 BM25** unambiguous.
4. **Note**: In some **manual experiments** we ran we observed that the **quality** of Wordnet Synonyms wasn't so good, so we preferred to use 1.

7. Conclusion

We evaluated two **WordNet-based query-expansion strategies** on our **BM25 Elasticsearch** baseline:

1. All-Synonyms (Scenario A): Broad synonym injection marginally boosts deep-cutoff MAP but significantly degrades early precision.

2. Hypernyms-Only (Scenario B): Conservative expansion by parent concepts yields consistent MAP gain.

A targeted hypernym-only approach effectively bridges vocabulary gaps with minimal precision loss.

References

1. Elasticsearch Reference (v8.17.2)
<https://www.elastic.co/guide/en/elasticsearch/reference/8.17>
2. Thakur, N., Reimers, N., Rüclé, A., Srivastava, A. and Gurevych, I., 2021. Beir: A heterogenous benchmark for zero-shot evaluation of information retrieval models. [arXiv preprint arXiv:2104.08663](https://arxiv.org/abs/2104.08663).
3. Manning, C. D., Raghavan, P., & Schütze, H. (2008). [*Introduction to Information Retrieval*. Cambridge University Press.](#)
4. Rivas, Andreia & Iglesias, Eva & Borrajo, María. (2014). Study of Query Expansion Techniques and Their Application in the Biomedical Information Retrieval. The Scientific World Journal. 2014. [1-10. 10.1155/2014/132158](#).
5. Fellbaum, C. ed., 1998. *WordNet: An electronic lexical database*. MIT press.
6. Bird, S., 2006, July. NLTK: the natural language toolkit. In *Proceedings of the COLING/ACL 2006 interactive presentation sessions* (pp. 69-72).

7. Rivas, A.R., Iglesias, E.L. and Borrajo, L., 2014. Study of query expansion techniques and their application in the biomedical information retrieval. *The Scientific World Journal*, 2014(1), p.132158.