

Query expansion with artificially generated texts

Vincent Claveau^[0000-0002-3459-0550]

CNRS - IRISA, Univ. Rennes
Campus de Beaulieu, F-35042 Rennes, France
vincent.claveau@irisa.fr

Abstract. A well-known way to improve the performance of document retrieval is to expand the user’s query. Several approaches have been proposed in the literature, and some of them are considered as yielding state-of-the-art results in IR. In this paper, we explore the use of text generation to automatically expand the queries. We rely on a well-known neural generative model, GPT-2, that comes with pre-trained models for English but can also be fine-tuned on specific corpora. Through different experiments, we show that text generation is a very effective way to improve the performance of an IR system, with a large margin (+10% MAP gains), and that it outperforms strong baselines also relying on query expansion (LM+RM3). This conceptually simple approach can easily be implemented on any IR system thanks to the availability of GPT code and models.

Keywords: Text generation, query expansion, GPT2, data-augmentation, document retrieval.

1 Introduction

In the IR traditional setting, a user expresses his information needs with the help of a query. Yet, it is sometimes difficult to match the query with the documents, for instance because of the query vocabulary may differ from the documents. Especially when the query is short, the performance of the system is usually poor, as it is difficult to detect the precise focus of the information need, and the relative importance of the query terms.

Query expansion aims at tackling these problems by transforming the short query into a larger text (or set of words) that make it easier to match documents from the collection. The main difficulty of query expansion is obviously to add only relevant terms to the initial query. Several techniques have been proposed in the literature, based on linguistic resources (e.g. synonym lists) or based on the documents themselves (e.g. pseudo-relevance feedback).

In this paper, we explore the use of recent text generation models to expand queries. We experimentally demonstrate that the recent advances in neural generation can dramatically improve ad-hoc retrieval, even when dealing with specialized domains. More precisely, through different experiments, we show that:

1. texts artificially generated from the query can be used for query expansion;

2. this approach does not only provide new terms to the query, but also a better estimate of their relative weights;
3. in addition, it also provides a better estimate of the importance (i.e. weight) of original query words;
4. this approach can also be used on specialized domains.

The paper is structured as follows. After a presentation of the related work, Section 3 details the different components of our approach. Several experiments are then detailed in Section 4. Last, some concluding remarks are given in Section 5.

2 Related work

Query expansion is a well-established technique to try to improve the performance of an IR system. Adding new terms to the query is expected to specifically improve recall, yet, since the query is, hopefully, better formulated, it may also improve the top rank results and be beneficial to precision. One might classify the existing automatic approaches based on the origins of resource used to expand the query.

External resources. One obvious way to expand a query is to add semantically related terms to it (synonyms or sharing other semantic relations like hyponyms, quasi-synonyms, meronyms...). Existing lexical resources can be used to add, for each query term, a list of semantically related terms; yet, one has to deal with different problems: existence of lexical resources for the collection language, or for the specific domain of the collection, choice of the appropriateness of certain relations, need of sense disambiguation for polysemous words... WordNet [14] is among the best-known resources for English (general language) and have been used with mitigated results at first [23], but later shown to be effective [3, inter alia].

Collection-based resources. Distributional thesauri have also been exploited to enrich queries. Since they can be built from the document collection (or from a large corpus with similar characteristics), they are suited to the domain, the vocabulary... Traditional techniques to build these thesauri have obtained good results for query expansion [3]. Neural approaches, that is, word embedding approaches are now widely used to build such semantic resources. In the recent years, static embeddings (word2vec [13], Glove [16] or FastText [2] to name a few) were also used in IR, in particular to enrich the query. Indeed, these trainable dense representations make it easy to find new words that are semantically close to query words.

Even more recently, dynamic word representations obtained with transformer-based architectures, such as BERT [5] or GPT [17], have been proposed. They build a representation for each word according to its context, and this ability have been exploited to obtain competitive results in IR tasks [4,7, inter alia]. Bert has been also used for query expansion in the framework of a neural IR system based on reranking [15]

Pseudo-relevance feedback A last category of studies considers only a small set of document to help to expand or reformulate the query. To be automatic, they replace the user feedback by the hypothesis that the best ranked documents retrieved with the original query are relevant and may contain useful semantic information [20]. It is interesting to note that in this case, not only semantically relevant terms are extracted, but also distributional/statistical information on them and on the original query terms. In this category, Rocchio, developed in the 60's for vector space model was among the first one popularized [11]. One of the current best known approach is RM3, which was developed in the framework of language model based IR systems [1]. It is often reported to yield the best results in ad-hoc retrieval tasks, even compared with recent neural models [9]. Neural approaches have also been proposed to integrate pseudo-relevance feedback information [8], yet, as it is reported by the authors, the results are still lower than traditional models with query expansion.

In this paper, we propose to use constrained text generation to expand queries. In this approach, the original query is used as a seed for a generative model which will output texts that are, hopefully, related to the query. While text generation with language model is not new, the performance of neural models based on transformers [22] makes this task realistic. In this paper, we use the Generative Pre-Trained Transformers (GPT). They are built from stacked transformers (precisely, decoders) trained on a large corpus by auto-regression, i.e. unsupervisedly trained to predict the next token given the previous ones. The second version, GPT-2 [17], contains 1.5 billion parameters for its largest pre-trained model, trained on more than 8M documents from Reddit (i.e. mostly English and general language such as discussion on press articles). A newer version, GPT-3, has been released in July 2020; it is much more larger (175 billion parameters) and outperforms GPT-2 on any tested task. Yet, the experiments reported below were done before this release, thus with GPT-2.

3 Generated query expansion

3.1 Overview of our approach

As it was previously explained, our approach is very simple as it relies on existing tools and techniques. From a query, multiple texts are generated by a GPT-2 model using the query as the seed. Note that the generation process is not deterministic, and thus, even with the same seed, the texts are different. In the experiments reported below, 100 texts per query are generated (unless specified otherwise). These texts are concatenated and considered as the new query. In our experiments, this new, very large, query is then fed to a simple BM25+ IR system, but it could obviously be used in any other IR system.

An example of a text generated from a query is presented in Fig. 1. As one can see, the generated text, while completely invented (note the barrel prices),

is relevant for the query. It contains many terms, absent from the original query, that are more or less closely related to the information need, such as orthographic variants (United States vs. U.S.), meronyms-metonyms (barrel vs. oil), hypernyms (energy vs. oil) and more generally any paradigmatic (consumer, producer vs. industry) or syntagmatic (production for oil) relations. It is worth noting that such texts also give a valuable information about the relative frequency of each terms (contrary to thesauri or embeddings).

U.S. oil production has been declining steadily for decades and it is not expected to reverse. In fact, some argue that it may even get worse. The long-term trend is for oil production to decline at a rate of about 1 percent per year. With production of about 8 million barrels per day now, there is no way the United States can replace its current output.

The U.S. oil boom was a result of an energy revolution in the 1970s that led to increased oil production, and a significant change in the global oil market. The U.S. now produces about 2.3 million barrels of oil per day, the highest it has been in over 30 years.

The United States is now the world's largest oil producer and the fourth largest oil exporter.

What happened?

When oil prices peaked in the 1970s, the United States was the world's largest oil producer. But over the next several decades, the United States' oil production began to decline. The decline was most pronounced in the 1980s, when the United States began to fall behind other oil producing countries.

The oil price decline in the 1970s was not entirely voluntary. The United States was producing less oil and exporting more oil than it was consuming. The Federal Reserve controlled the amount of dollars in the Federal Reserve's reserves, so the United States was not exporting as much oil as it was producing. The decline in U.S. oil production was a result of the declining price of oil.

The price of oil had declined from \$8 per barrel in 1973 to \$2.50 per barrel in 1977. In 1979, the price of oil reached a high of \$15.75 per barrel. By 1983, the price of oil had fallen to \$4.65 per barrel. By 1986, the price of oil had fallen to \$1.86 per barrel. By the end of the 1980s, the price of oil had fallen to \$1.24 per barrel.

The decline in oil prices was a direct result of the energy revolution in the 1970s. The United States was the world's largest oil producer, but the United States was also the world's largest consumer of oil. When oil prices fell, so did the cost of producing oil.

Fig. 1. Example of a document generated with the pre-trained GPT-2 large model from the text seed "U.S. oil industry history" (query 701 from .GOV collection)

3.2 Pre-trained models, fine-tuning and parameters

GPT-2 comes with several pre-trained models, having different size in terms of parameters (from 124M to 1.5B). As it was previously said, their training data was news-oriented general language. The largest model was used for two of the tested collections (see below). While these all-purpose models are fine for IR

collections whose documents are also general language, it may not be appropriate for domain-specific IR collections. In the experiment reported in the next section, we use the OHSUMED collection, made of medical documents. For this collection, we have fine-tuned the GPT-2 355M model on the documents of the collection in order to adapt the language model to the specific medical syntax and vocabulary. The fine-tuning was stopped after 250,000 samples were processed (this number of sample process indirectly controls the under/over-fitting to the specialized corpus) was set to and other parameters (batch size, optimizer, learning rate...) let to their defaults. Although a larger set of medical documents could be used (from Pubmed for instance), this small fine-tuned model is expected to be more suited to generate useful documents to enrich the query.

Concerning the generation of documents, for reproducibility purposes, here are the main GPT-2 parameters used (please refer to GPT-2 documentation¹): length=512, temperature=0.5, top_p=0.95, top_k = 40.

3.3 IR Systems

In the experiments reported in the next section, we use two IR models. The first one is BM25+ [10], a variant of BM25 [19]. The parameters k_1 , k_3 , b and δ were kept at their default value (resp. 1.2, 1000, 0.75, 1). It is implemented as a custom modification of the GENSIM toolkit [18]. The second IR model is Language modeling with Dirichlet smoothing [25] as implemented in Indri [12,21]. The smoothing parameter μ is set to 2500. Both models are regarded as yielding state-of-the-art performance for bag-of-words representation [9]. Their RSV function can be written:

$$RSV(q, d) = \sum_{t \in q} w_q(t) \cdot w_d(t)$$

with $w_q(t)$ the weight of term t in query q and $w_d(t)$ the weight in document d , as illustrated in Tab. 1 (from [10]). For RM3 expansion, we also rely on the Indri implementation; the results reported in the next section corresponds to the best performing parameters tested for each collection (number of documents considered for pseudo relevance feedback and number of terms kept).

4 Experiments

4.1 Experimental settings

Three IR collections are used in our experiments: Tipster, GOV2 and OHSUMED. Some basic statistics are given in Tab. 2.

Tipster was used in TREC-2. The documents are articles from newspaper, patents and specialized press (computer related) in English. The queries are composed of several parts, including the query itself and a narrative detailing

¹ <https://github.com/openai/gpt-2>

IR model	weighting
BM25+ $w_d(t)$	$\left(\frac{(k_1+1)c(t,d)}{k_1(1-b+b \cdot dl(d)/avdl)+c(t,d)} + \delta \right) \cdot \log \frac{N+1}{df(t)}$
BM25+ $w_q(t)$	
	with k_1, k_3, b and δ fixed parameters
LM $w_d(t)$	$\log \left(\frac{\mu}{dl(d)+\mu} + \frac{c(t,d)}{(dl(d)+\mu)p(t C)} \right)$
LM $w_q(t)$	
	$\mu > 0$ a smoothing parameter

Table 1. IR models (weighting functions of terms in the query and the document) for BM25+ [19,10] and Language modeling with Dirichlet smoothing LM [25]

	Tipster	GOV2	OHSUMED
nb of documents	170,000	25M	350,000
nb of queries	50	150	106
avg size of queries	6.74	3.15	7.24
language	En	En	En
avg nb of relevant doc per query	849	179	21

Table 2. Statistics on the IR collections used

the relevance criteria; in the experiments reported below, only the actual query part is used.

GOV2 is a large collection of Web pages crawled from the .gov domain and used in several TREC tracks. In the experiments reported below, 150 queries from TREC 2004-2006 ad-hoc retrieval tasks are used; as for Tipster, only the actual query part is used (i.e. description and narrative fields are not included in the query).

OHSUMED [6] contains bibliographical notices from Medline and queries from the TREC-9 filtering task. Its interest for our experiments is that it deals with a specialized domain, hence it contains a specific vocabulary.

Performance are assessed with standard scores: Precision at different thresholds (P@x), R-precision (R-prec), *Mean Average Precision* (MAP). When needed, a paired t-test with $p = 0.05$ is performed to assess the statistical significance of the difference between systems.

4.2 General language

Tables 3 and 4 respectively present the results for the general-language collections Tipster and GOV2. For comparison purposes, we indicate the results of BM25+ without expansion, Indri’s Language Model (LM) with and without RM3 expansion. The best performing setting for RM3 on Tipster is 100 terms

	MAP	R-Prec	P@5	P@10	P@20	P@100
BM25+	25.06	32.16	95.60	92.60	89.70	73.64
LM	24.48	31.48	92.40	89.00	85.40	70.70
LM + RM3	31.01	36.38	94.40	93.20	90.60	81.22
BM25+ and expansion	35.22*	39.87*	99.60*	98.40*	98.20*	87.84*

Table 3. Performance (%) on Tipster with query expansion; best results in bold, statistical significance over LM+RM3 noted with *

	MAP	R-Prec	P@5	P@10	P@20	P@100
BM25+	25.66	31.25	52.92	49.97	46.52	34.63
LM	27.96	33.01	56.08	55.20	51.59	37.32
LM with RM3	30.22	34.20	55.00	56.08	53.67	45.86
BM25+ and expansion	34.54*	37.76	67.91*	63.88*	57.94	44.30

Table 4. Performance (%) on GOV2 with query expansion; best results in bold, statistical significance over LM+RM3 noted with *

from the top 20 documents, and 100 terms for the top 10 documents GOV2. The statistical significance is computed by comparing with the LM+RM3 baseline.

On both collections, and on every performance measure, expanding the queries with the generated texts brings important gains compared with the system without expansion. Also, our approach outperforms RM3 expansion in almost every situation, and with a large margin on MAP, R-prec and precision on the top-ranked documents (P@5, P@10).

4.3 Specialized language

The same setting is used on the OHSUMED collection. For these medical-oriented IR dataset, we report two versions of our approach: one is using the pre-trained model as before, and one relies on a model fine-tuned on the documents of the collection. The best performing setting for RM3 is 80 terms for the top 10 documents. The results are reported in Tab. 5.

	MAP	R-Prec	P@5	P@10	P@20	P@100
BM25+	18.27	19.94	31.88	26.04	20.50	9.48
LM	17.61	20.35	29.31	24.06	19.21	9.28
LM + RM3	20.80	22.54	30.89	26.83	22.18	10.51
BM25+ and expansion (no fine-tuning)	21.60	23.75	33.47*	27.62	22.92	11.16
BM25+ and expansion (fine-tuning)	23.07*	24.65*	34.65*	29.41*	24.31	11.42

Table 5. Performance (%) on OHSUMED with query expansion; best results in bold, statistical significance over LM+RM3 noted with *

Here again, the GPT-based expansion significantly improves the results of the IR system and outperforms RM3 expansion. Yet, the gains are lower than for the two previous collections. This difference can be explained by the following factors:

1. the queries are longer more complex and more specific (as can be seen in Tab. 2, few documents are relevant);
2. the generation model is not sufficiently suited to the documents.

Concerning this latter reason, we can indeed see the interest of fine-tuning the generation model, but better results may be obtained by using a larger set of medical documents, or adopting different fine-tuning parameters (in particular the number of epoch/samples processed, see Sect. 3.2). Unfortunately, defining a priori the best parameters for our IR task is not possible and the cost of the fine-tuning process makes it impossible to test a wide range of possible values.

4.4 Query expansion and weight information

One of the interest of having complete texts that are generated is that we can collect information on the relative importance of words, to the contrary of expanding queries with a thesaurus. To observe the impact of the number of occurrences in the generated texts, we evaluate the effect of keeping the k most frequent terms of the generated texts and either weighting them by their frequency (as done usually by BM25) or by giving a fixed weight ($1/k$). The results for different values of k are presented in Fig. 2. On can observe that adding

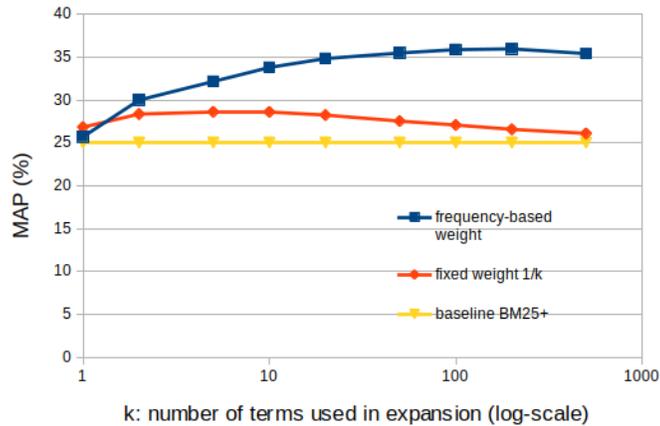


Fig. 2. MAP (%) according to size of the expansion, terms with fixed weight or weight depending on their frequency in the generated documents; Tipster collection

terms to the query with a fixed weight slightly improves the MAP, but most of

	MAP	R-Prec	P@5	P@10	P@20	P@100
BM25+	25.06	32.16	95.60	92.60	89.70	73.64
BM25+ with re-weighting	27.12	33.22	96.80	94.20	92.70	77.12

Table 6. Performance (%) on Tipster with re-weighting query words

the gain is indeed brought by a proper weighting based on the frequency of the term in the generated documents.

In the next experiment, we examine how the generated texts can help to re-weight the initial query terms; there is no query expansion, since only the initial query terms are kept, but their frequency in the generated texts are used in the BM25+ w_q . The results reported in Tab. 6 show that there is indeed a small improvement of the MAP, that is more noticeable at high DCV. These two experiments demonstrate the usefulness of dealing with full texts and not only word-to-word similarity.

4.5 Number of generated documents

Since text generation can be costly, it is interesting to see how many generated texts are necessary. In Fig. 3, the MAP obtained for up to 100 documents is presented. One can observe that a plateau is rapidly reached (at around 20 documents per query). Of course, the size of the generated documents (can be set as a parameter of the generation process) is also to be considered.

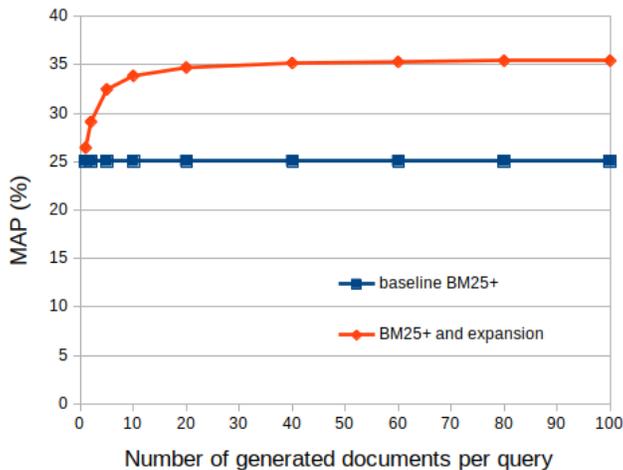


Fig. 3. MAP (%) according to the number of generated documents (average over 5 runs); Tipster collection

5 Conclusive remarks and foreseen work

Neural approaches are increasingly used in IR, with mitigated results, especially when compared with "traditional" bag-of-word approaches [9,24]. Here, the neural part is successfully used outside of a "traditional" IR system (but note that it could be used with any IR systems, since it simply enriches the query). The expansion approach presented in this paper is simple and easy to implement (thanks to the availability of the GPT models and code) while offering impressive gains. Lot of parameters could be further optimized, especially on the GPT model side (to influence the "creativity" of the text generation), and the fine-tuning capabilities should also be explored more thoroughly (influence of bigger specialized corpus if available, precise mix between pre-trained and fine-tuning, etc.). The recent availability of GPT-3² makes it possible to even get greater gains thanks to the alleged high quality of its outputs.

This whole approach also offers many research avenues: in this work, we have used text generation as a way to perform data augmentation on the query side, but it could also be used to augment the representation of the documents (even if in practice, the cost is still prohibitive on large collection). All machine-learning (neural or not) approaches based on pseudo-relevance feedback to train their model could instead use similar text generation with the advantage that they would not be limited by the number of potential relevant documents in the shortlist. And of course, similar data-augmentation strategy could be used for other tasks than document retrieval.

More fundamentally, the recent improvements of text generation also question the relevance of the document retrieval task. Indeed, it is possible to envision systems that will be able to generate one unique document answering the user's information need, similarly to question-answering. If the generative model is trained on the document collection, the generated document will serve as a summary (which is one of the popular applications of GPT-x models) of the relevant documents. Yet, the current limitations of the models tested in this paper make them far from being suited for this ultimate task: the generated documents do deal with the subject of the query, and thus use a relevant vocabulary, but do not provide accurate, factual information (as seen in the Example in Fig. 1 about the price of oil barrels).

References

1. Abdul-jaleel, N., Allan, J., Croft, W.B., Diaz, O., Larkey, L., Li, X., Smucker, M.D., Wade, C.: Umass at trec 2004: Novelty and hard. In: In Proceedings of TREC-13 (2004)
2. Bojanowski, P., Grave, E., Joulin, A., Mikolov, T.: Enriching word vectors with subword information. arXiv preprint arXiv:1607.04606 (2016)

² <https://github.com/openai/gpt-3>

3. Claveau, V., Kijak, E.: Direct vs. indirect evaluation of distributional thesauri. In: Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers. pp. 1837–1848. The COLING 2016 Organizing Committee, Osaka, Japan (Dec 2016), <https://www.aclweb.org/anthology/C16-1173>
4. Dai, Z., Callan, J.: Deeper text understanding for ir with contextual neural language modeling. Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval (Jul 2019). <https://doi.org/10.1145/3331184.3331303>, <http://dx.doi.org/10.1145/3331184.3331303>
5. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding (2019)
6. Hersh, W., Buckley, C., Leone, T.J., Hickam, D.: Ohsumed: An interactive retrieval evaluation and new large test collection for research. In: Proceedings of the 17th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 192–201. SIGIR '94, Springer-Verlag New York, Inc., New York, NY, USA (1994), <http://dl.acm.org/citation.cfm?id=188490.188557>
7. Khattab, O., Zaharia, M.: Colbert: Efficient and effective passage search via contextualized late interaction over BERT. In: Huang, J., Chang, Y., Cheng, X., Kamps, J., Murdock, V., Wen, J., Liu, Y. (eds.) Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020. pp. 39–48. ACM (2020). <https://doi.org/10.1145/3397271.3401075>, <https://doi.org/10.1145/3397271.3401075>
8. Li, C., Sun, Y., He, B., Wang, L., Hui, K., Yates, A., Sun, L., Xu, J.: NPRF: A neural pseudo relevance feedback framework for ad-hoc information retrieval. In: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. pp. 4482–4491. Association for Computational Linguistics, Brussels, Belgium (Oct-Nov 2018). <https://doi.org/10.18653/v1/D18-1478>, <https://www.aclweb.org/anthology/D18-1478>
9. Lin, J.: The neural hype and comparisons against weak baselines. SIGIR Forum **52**(2), 40–51 (2018). <https://doi.org/10.1145/3308774.3308781>, <https://doi.org/10.1145/3308774.3308781>
10. Lv, Y., Zhai, C.: Lower-bounding term frequency normalization. In: Proc. of the 20th ACM International Conference on Information and Knowledge Management. pp. 7–16. CIKM '11, ACM, New York, NY, USA (2011). <https://doi.org/10.1145/2063576.2063584>, <http://doi.acm.org/10.1145/2063576.2063584>
11. Manning, C.D., Raghavan, P., Schütze, H.: Introduction to Information Retrieval. Cambridge University Press, USA (2008)
12. Metzler, D., Croft, W.: Combining the language model and inference network approaches to retrieval. Information Processing and Management Special Issue on Bayesian Networks and Information Retrieval **40**(5), 735–750 (2004)
13. Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. In: Burges, C.J.C., Bottou, L., Ghahramani, Z., Weinberger, K.Q. (eds.) Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States. pp. 3111–3119 (2013), <http://papers.nips.cc/paper/5021-distributed-representations-of-words-and-phrases-and-their-compositionality>

14. Miller, G.A.: WordNet: An On-Line Lexical Database. *International Journal of Lexicography* **3**(4) (1990)
15. Padaki, R., Dai, Z., Callan, J.: Rethinking query expansion for bert reranking. In: Jose, J.M., Yilmaz, E., Magalhães, J., Castells, P., Ferro, N., Silva, M.J., Martins, F. (eds.) *Advances in Information Retrieval*. pp. 297–304. Springer International Publishing, Cham (2020)
16. Pennington, J., Socher, R., Manning, C.D.: Glove: Global vectors for word representation. In: *Empirical Methods in Natural Language Processing (EMNLP)*. pp. 1532–1543 (2014), <http://www.aclweb.org/anthology/D14-1162>
17. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I.: Language models are unsupervised multitask learners. *OpenAI Blog* (2019)
18. Řehůřek, R., Sojka, P.: Software Framework for Topic Modelling with Large Corpora. In: *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*. pp. 45–50. ELRA, Valletta, Malta (May 2010), <http://is.muni.cz/publication/884893/en>
19. Robertson, S.E., Walker, S., Hancock-Beaulieu, M.: Okapi at TREC-7: Automatic Ad Hoc, Filtering, VLC and Interactive. In: *Proc. of the 7th Text Retrieval Conference, TREC-7*. pp. 199–210 (1998)
20. Ruthven, I., Lalmas, M.: A survey on the use of relevance feedback for information access systems. *Knowledge Eng. Review* **18**(2), 95–145 (2003)
21. Strohman, T., Metzler, D., Turtle, H., Croft, W.: Indri: A language-model based search engine for complex queries (extended version). *Tech. rep., CIIR* (2005)
22. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.u., Polosukhin, I.: Attention is all you need. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) *Advances in Neural Information Processing Systems 30*, pp. 5998–6008. Curran Associates, Inc. (2017), <http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>
23. Voorhees, E.M.: Query expansion using lexical-semantic relations. In: *Proc. of the 17th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. pp. 61–69. SIGIR '94, Springer-Verlag New York, Inc., New York, NY, USA (1994), <http://dl.acm.org/citation.cfm?id=188490.188508>
24. Yang, W., Lu, K., Yang, P., Lin, J.: Critically examining the "neural hype": Weak baselines and the additivity of effectiveness gains from neural ranking models. In: Piwowarski, B., Chevalier, M., Gaussier, É., Maarek, Y., Nie, J., Scholer, F. (eds.) *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2019, Paris, France, July 21-25, 2019*. pp. 1129–1132. ACM (2019). <https://doi.org/10.1145/3331184.3331340>, <https://doi.org/10.1145/3331184.3331340>
25. Zhai, C., Lafferty, J.D.: A study of smoothing methods for language models applied to ad hoc information retrieval. In: *Proc. of the SIGIR conference*. pp. 334–342 (2001)