

Université Grenoble Alpes / LIG at TREC Deep Learning Track 2022

Petra Galuščáková, Lucas Alberede, Gabriela Nicole Gonzalez Saez, Aidan Mannion, Philippe Mulhem, Georges Quénot

Univ. Grenoble Alpes, CNRS, Inria, Grenoble INP**, LIG, 38000 Grenoble, France

Abstract. This paper describes the submissions of the MRIM research Group/Laboratoire d’Informatique de Grenoble from the Université Grenoble Alpes at the TREC Deep Learning Track 2022. We participated in both fullranking and reranking sub-tasks in the passages ranking task. We proposed several ways to combine neural and non-neural runs using reciprocal-rank based fusion.

1 Introduction

This report describes the submissions of the Modeling and Retrieval of Multimedia Information research Group at Université Grenoble Alpes team at the TREC Deep Learning Track 2022. We have participated in the passages ranking task, in both fullranking and reranking sub-tasks.

For the **fullranking sub-task**, we used several ranking models for the first stage retrieval and an index of the full collection. We subsequently applied reranking on top of the outputs of this ranking. Last, we combined the outputs of multiple rankers and re-rankers into a single output. We comprehensively experimented with a range of first stage retrieval models, rerankers, and architectures and selected 3 runs which we submitted and which were used in the pooling and 4 additional runs, which we submitted but which were not used in the pooling.

In the **reranking sub-task**, we applied three reranking approaches to the provided list of passages and we either combined them or we chained them on the top of each other.

In the submissions selection process, we used the test collection used at TREC Deep Learning 2021 Passage ranking task, which in terms of collection statistics, should well correspond to this year track’s collection. This collection contains 477 queries, for 53 out of which there are provided the relevance judgements. As this number is relatively small, in the selection we aimed for both good performance and robustness of the submitted system.

In the following section, we cover the basic blocks which we used in the submissions. We review the fullranking models used for getting the top relevant passages from the full passage collection and the models for reranking which we applied to the passages retrieved by these fullranking models. We then review

** Institute of Engineering Univ. Grenoble Alpes

our fullranking submissions in the Section 3 and the reranking submissions in Section 4. All the scripts used in our submissions are freely available¹.

2 Basic Building Blocks

We use full ranking models for ranking all the passages available in the collection as the first stage of the retrieval. We experimented with a range of full ranking models from which we selected 5 well performing and diverse models, as such systems might be expected to perform well in a combination. We eventually use following full ranking models in all our submissions:

- **DPH:** We indexed and retrieved the segments using Terrier v5.0 framework [8]. DPH was used with stemming and with blind relevance feedback², as this is expected to increase a diversity of the performance of the models. DPH [4] is a well performing and a parameter free scoring technique derived from the Divergence from Randomness model.
- **BM25:** The standard BM25 retrieval model with default parameters, using the index and implementation provided by Anserini [10]. In addition to the index of passages, we also use the index of documents, which are further used in the reranking using the Graph Model (see below). If not stated otherwise, we further refer to the retrieval of passages in this report.
- **Expanded BM25:** BM25 with document expansion using docT5query [6] and RM3 blind relevance feedback. We used the index and implementation provided by Anserini.
- **Colbert:** A BERT-based passage retrieval model [2]. Specifically, we used the TCT-ColBERTv2 model fine-tuned on MS MARCO v2 provided by Anserini.
- **uniCOIL:** A model that integrates BERT embeddings into a sparse representation [3]. Specifically, we have used Castorini uncoil-msmarco-passages model pre-trained on MS MARCO v2 Passage Corpus, along with on-the-fly docT5query query inference. We used an index provided by Anserini.

Passages returned by these models which were ranked at the top positions are in some cases further reranked. Reranking is done by calculating the similarity between the query and retrieved passage and reordering the passages based on this similarity. Specifically, we use following approaches:

- **BERT:** We used the monoBERT model [7] from Castorini, available in Pygaggle³, which was finetuned on MS MARCO v1.
- **T5:** We used the monoT5 model [5] from Castorini available in Pygaggle which was finetuned on MS MARCO v1.

¹ <https://github.com/galuscakova/trec-dl-2022>

² As recommended in <http://terrier.org/docs/current/javadoc/org/terrier/matching/models/DPH.html>

³ <https://github.com/castorini/pygaggle>

- **MiniLM:** Cross-encoders are transformer models [9] fine-tuned for sentence similarity calculation, which are widely used in semantic search as they compute attention weights across the query and the document together. We used the MiniLM architecture provided by the `sentence-transformers` library⁴, which have been pre-trained on the MS MARCO v2 passage reranking dataset.
- **Graph Model:** We use Heterogeneous Graph Attention Networks (HGATs). Document retrieval (BM25) is applied as a first stage retrieval returning the top 1000 documents. Every passage from these documents is then reranked using HGATs to leverage its relations with other passages in a graph-based document representation [1]. Passage and query embeddings are computed using a multiple-representation embedding encoder similar to the one in the Colbert architecture [2].

Finally, for the combination of the ranked and reranked runs, we used reciprocal rank-based (RR) fusion implemented in `TrecTools`⁵.

3 Fullranking Submissions

All our submissions consisted of a combination of fullranking and reranking models. We submitted three runs which were further used in the pooling: hierarchical combination, flat combination and double reranking. A graphical view of the submissions is displayed in Figure 1.

Hierarchical Combination: First, Colbert, UniCOIL and Expanded BM25 rankers were combined and then reranked by BERT [top 10 documents]. Then, this combination was further combined with UniCOIL reranked by dense MiniLM model [top 100 documents] and with BM25 reranked by the Graph Model [top 100 document]. This system had the highest performance on our training data in terms of nDCG@10, nDCG@1000 and MAP. The official run name used in the submissions is *hierarchical combination*.

Flat Combination: This combination consists of 6 rankers, rerankers and combinations which all together form the final combination. Due to the diversity of the combined models, this system is also supposed to be robust. Following models are combined:

- UniCOIL
- UniCOIL reranked by a MiniLM model [top 100 documents]
- Colbert reranked by T5 [top 100 documents]
- Expanded BM25 reranked by T5 [top 100 documents]
- BM25 reranked by the Graph Model [top 100 documents]

⁴ <https://www.sbert.net/docs/pretrained-models/ce-msmarco.html>

⁵ <https://github.com/joaopalotti/trectools>

- BM25 and DPH first combined, then reranked by T5 [top 100 documents]

The official run name used in the submissions is *6systems*.

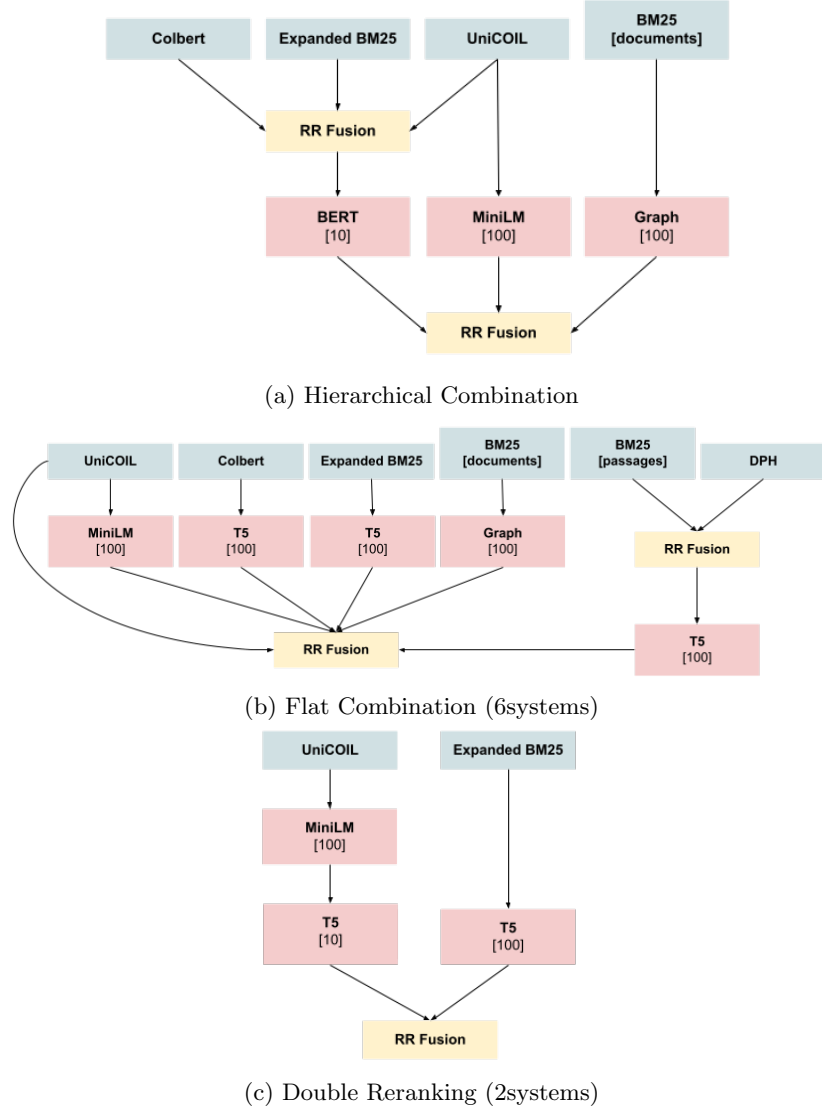


Fig. 1: Full ranking architectures. Ranking models are in blue, reranking models in rose, and fusions in yellow. Number of reranked passages/documents is in the brackets.

Double Reranking: This submission follows the architecture of duoT5 model with two different rerankers chained. The first stage ranker is UniCOIL model. Top 100 returned passages are first reranked using the MiniLM model. Subsequently, the top 10 reranked documents are further reranked by T5 model. Last, this double reranked model is further combined with the expanded BM25 model reranked by T5 model [top 100 documents]. The official run name used in the submissions is *2systems*.

In addition to these runs, we submitted other four runs which were officially evaluated by the task organizers, but haven't been included into an evaluation pool.

Graph: This system is simply the Graph Model described in Section 2 which uses a BM25 ranking of the documents followed by reranking done using the Graph Attention Networks. The official run name used in the submissions is *graph colbert*.

Reranked Unicoil: This system also does not use any reciprocal rank-based combination. Instead, it only uses double reranking of the UniCOIL system, first by MiniLM model [top 100 documents] and then by T5 model [top 10 documents]. The system is thus the same as the Double Reranking system which was submitted for the pooling, but without the expanded BM25 model applied. The official run name used in the submissions is *unicoil reranked*.

4 Systems with Graph: This combination contains a subset of 4 systems which are included in the Flat Combination:

- UniCOIL
- UniCOIL reranked by a MiniLM model [top 100 passages]
- Expanded BM25 reranked by T5 [top 100 passages]
- BM25 reranked by the Graph Model [top 100 passages]

The official run name used in the submissions is *4systems*.

4 Systems with DPH: This combination again contains a subset of 4 systems which are included in the Flat Combination. Comparing with 4 Systems with Graph, the BM25 model reranked by the graph model is here replaced by a combination of BM25 and DPH models:

- UniCOIL
- UniCOIL reranked by a MiniLM model [top 100 documents]
- Expanded BM25 reranked by T5 [top 100 documents]
- BM25 and DPH first combined, then reranked by T5 [top 100 documents]

The official run name used in the submissions is *c47*.

3.1 Official Evaluation Results

The evaluation of these submitted systems is displayed in Table 1. We see from this table, that the single Graph-based reranking and the Hierarchical Combination outperform the other combinations for the measures taking into account top 1000 results (MAP, nDCG@1000 and Recall@1000). The 4-way combinations perform the worst according to these measures. As the Hierarchical Combination also includes the Graph-based reranking as one of the combined systems and as they perform similarly in terms of retrieving the top 1000 documents, we can claim, that the additional systems in the combination add no extra improvement to the Graph system. However, the Hierarchical Combination performs better than the Graph in terms of the measures working only with top 10 documents (nDCG@10 and P@10). The additional systems added to the Graph system thus do not change the performance of the retrieval on the top 1000 documents, but lead to more than 5% relative improvement for the top 10 documents retrieval. However, clearly the best performance in terms of retrieving top 10 documents is achieved by the Double Reranking system (in terms of nDCG@10 and P@10 measures). This might confirm the usefulness of the duo architecture and it should lead to further experimentation with different cascading setups and rerankers.

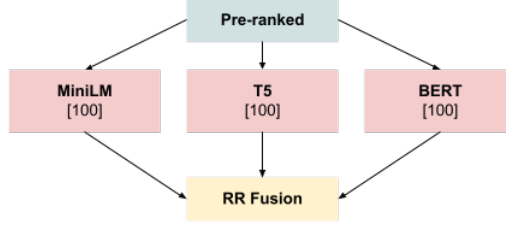
| Run name | Pooling | MAP | nDCG@10 | nDCG@1000 | P@10 | Recall@1000 |
|--------------------------|---------|--------------|--------------|--------------|--------------|--------------|
| Hierarchical Combination | Yes | 0.170 | 0.562 | 0.358 | 0.480 | 0.332 |
| Flat Combination | Yes | 0.166 | 0.567 | 0.339 | 0.487 | 0.309 |
| Double Reranking | Yes | 0.168 | 0.590 | 0.328 | 0.497 | 0.280 |
| Reranked Unicoil | No | 0.167 | 0.583 | 0.327 | 0.487 | 0.280 |
| Graph | No | 0.171 | 0.537 | 0.355 | 0.450 | 0.334 |
| 4 Systems w. Graph | No | 0.159 | 0.565 | 0.322 | 0.484 | 0.280 |
| 4 Systems w. DPH | No | 0.155 | 0.561 | 0.321 | 0.467 | 0.280 |

Table 1: Full ranking sub-task UGA official results. The highest scores for each measure are highlighted.

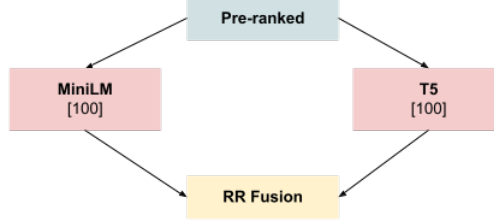
4 Reranking submissions

In the reranking submissions, we only used the reranking methods described above and we either combined them using a fusion or using a chain of double reranking. The graphical overview of all submissions is displayed in Figure 2.

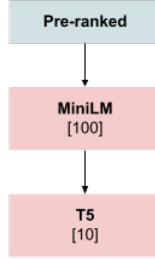
3-way Combination: We reranked the original run by the MiniLM model, T5 model and BERT model. Then we combined all three reranked runs into a single run. The official run name for this submission is *fused_3runs*.



(a) 3-way Combination (fused_3runs)



(b) 2-way Combination (fused_2runs)



(c) Double Reranking (hierarchical_2runs)

Fig. 2: Reranking architectures.

2-way Combination: We reranked the original run by the MiniLM model and T5 model. Then we combined both reranked runs into a single run. This run might be directly compared with the double reranking run, which uses the same rerankers, but combined together in a different architecture and to the 3-way combination that uses one additional reranker. The official run name for this submission is *fused_2runs*.

Double Reranking: Similarly, as in the full reranking run, we used the duo architectures and stacked two rerankers on a top of each other. We first reranked all 100 passages by the MiniLM model and then reranked again the top 10 passages by the T5 model. The official run name for this submission is *hierarchical_2runs*.

Official evaluation results: The evaluation of our reranking submitted systems is displayed in Table 2. Compared to Table 1, we do not present the Recall@1000 values as reranking has no effect on it. Clearly, these results are lower than the results for the fullranking, in terms of all considered measures. This shows the importance of using adequate and varied inputs for the reranking step. What we also notice is that the fused 3 runs and fused 2 runs achieve similar results. The hierarchical combination of 2 runs (i.e. the Double Reranking) underperforms the two other systems, especially when evaluating top 10 retrieved documents. This is contrary to a good performance of the Double Reranking run in the fullranking task. Using multiple rerankers in parallel followed by a fusion thus in this case seems to be a slightly better strategy than performing several rerankings in sequence. This, especially the conditions for an effective stacking of the rankers (i.e. their performance, the performance of the initial run), thus should be further explored.

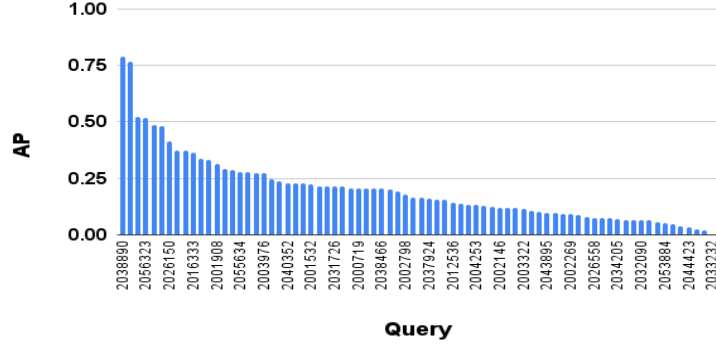
| Run name | Pooling | MAP | nDCG10 | nDCG@1000 | P@10 |
|-------------------|---------|--------------|--------------|--------------|--------------|
| 3-way Combination | Yes | 0.096 | 0.499 | 0.231 | 0.370 |
| 2-way Combination | Yes | 0.096 | 0.499 | 0.232 | 0.366 |
| Double Reranking | Yes | 0.095 | 0.489 | 0.231 | 0.347 |

Table 2: Reranking sub-task UGA official results. The highest scores for each measure are highlighted.

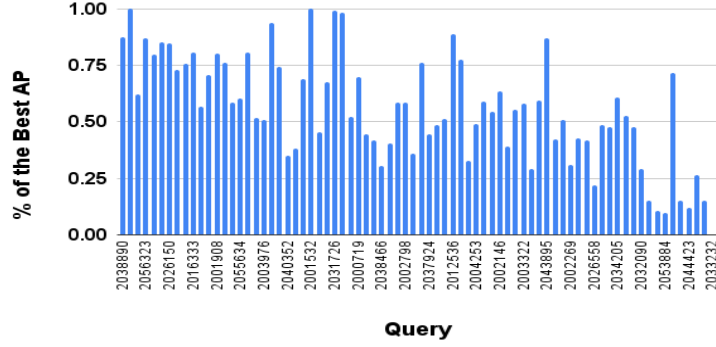
5 Results per Query

Last, we explored the results per query for the pools of all our submitted systems for fullranking and reranking. Overall, our systems performed the worst for the query 2033232 (“people who have fled the us because trump”), for which we achieved zero score for all AP, NDCG@10 and Precision measures. However, this query was hard for all the systems submitted to the pool. Both NDCG@10 and Precision scores are equal to 0 and the highest achieved AP is 0.0009 in the pool of all available systems. Comparing the pool of all submitted systems and the pool of our systems, we achieved the highest score for two queries (query 2001532: “example of what a family advocate does” and query 2013306: “who is gehan homes”) in terms of AP, for other two queries (query 2001908: “how do I replace the burners on a ducane grill” and query 2039908: “what is the name of a baby nurse”) in terms of NDCG@10 and for 14 queries in terms of Precision. We show the overall best achieved AP scores by any of our submitted systems in Figure 3 (a). Figure (b) shows the relative comparison of our best performing system for the query and the best performing system submitted into the pool by any team. Though there is some correlation between the absolute performance in (a) and the relative comparison in (b), for example for the query 2028378 (‘when is trial by jury used’), the relative performance is 0.0461 AP, the overall

best performing system achieves 0.0643 AP and our score is thus 72% of this value.



(a) AP achieved by our best submitted system per query, sorted by AP performance.



(b) Ratio of performance, for the same query, of our best performing system over the overall best performing system submitted to the task. The queries are sorted according to the AP value in Figure (a).

Fig. 3: Performance of our best performing system per query in terms of Average Precision (AP) measure.

6 Conclusion

We described our submissions for the fullranking and reranking sub-tasks of TREC Deep Learning 2022 Passage Retrieval Task. For the fullranking sub-task, a graph-based reranking of the passages identified using document ranking performs especially well for retrieving top 1000 documents. Its combination with additional well performing systems then further leads to a strong performance

on the top 10 documents. The best performing system for the retrieval of the top 10 documents combines double reranked UniCOIL system with reranked BM25 system with query and document expansions applied. However, double reranking was in the reranking task outperformed by a reciprocal rank-based combination of two or three rerankers. Our experiments thus lead to further questions about behaviour of multiple rankers and rerankers, which we plan to explore in future.

Acknowledgments

This work is supported by the ANR Kodicare bi-lateral project, grant ANR-19-CE23-0029 of the French Agence Nationale de la Recherche, and by the Austrian Science Fund (FWF).

References

1. L. Albarede, P. Mulhem, L. Goeuriot, C. L. Pape-Gardeux, S. Marie, and T. Chardin-Segui. Passage retrieval on structured documents using graph attention networks. In M. Hagen, S. Verberne, C. Macdonald, C. Seifert, K. Balog, K. Nørsvåg, and V. Setty, editors, *Advances in Information Retrieval - 44th European Conference on IR Research, ECIR 2022, Stavanger, Norway, April 10-14, 2022, Proceedings, Part II*, volume 13186 of *Lecture Notes in Computer Science*, pages 13–21. Springer, 2022.
2. O. Khattab and M. Zaharia. Colbert: Efficient and effective passage search via contextualized late interaction over bert. In *Proceedings of SIGIR 2020*, SIGIR '20, page 39–48, New York, NY, USA, 2020. Association for Computing Machinery.
3. J. Lin and X. Ma. A Few Brief Notes on DeepImpact, COIL, and a Conceptual Framework for Information Retrieval Techniques, 2021. <https://arxiv.org/abs/2106.14807>.
4. R. McCreadie, C. Macdonald, I. Ounis, J. Peng, and R. L. T. Santos. University of glasgow at TREC 2009: Experiments with terrier. In E. M. Voorhees and L. P. Buckland, editors, *Proceedings of The Eighteenth Text REtrieval Conference, TREC 2009, Gaithersburg, Maryland, USA, November 17-20, 2009*, volume 500-278 of *NIST Special Publication*. National Institute of Standards and Technology (NIST), 2009.
5. R. Nogueira, Z. Jiang, R. Pradeep, and J. Lin. Document ranking with a pre-trained sequence-to-sequence model. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 708–718, Online, Nov. 2020. Association for Computational Linguistics.
6. R. Nogueira and J. Lin. From doc2query to docTTTTTquery. Technical report, University of Waterloo, 2019.
7. R. Nogueira, W. Yang, K. Cho, and J. Lin. Multi-Stage Document Ranking with BERT, 2019. <https://arxiv.org/abs/1910.14424>.
8. I. Ounis, G. Amati, P. V., B. He, C. Macdonald, and Johnson. Terrier Information Retrieval Platform. In *Proceedings of the 27th European Conference on IR Research (ECIR 2005)*, volume 3408 of *Lecture Notes in Computer Science*, pages 517–519. Springer, 2005.

9. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
10. P. Yang, H. Fang, and J. Lin. Anserini: Enabling the use of lucene for information retrieval research. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '17, page 1253–1256, New York, NY, USA, 2017. Association for Computing Machinery.