# DUET AT TREC 2019 DEEP LEARNING TRACK

**Bhaskar Mitra**
Microsoft, University College London
Montreal, Canada
bmitra@microsoft.com

**Nick Craswell**
Microsoft
Redmond, USA
nickcr@microsoft.com

## ABSTRACT

This report discusses three submissions based on the Duet architecture to the Deep Learning track at TREC 2019. For the document retrieval task, we adapt the Duet model to ingest a "multiple field" view of documents—we refer to the new architecture as Duet with Multiple Fields (DuetMF). A second submission combines the DuetMF model with other neural and traditional relevance estimators in a learning-to-rank framework and achieves improved performance over the DuetMF baseline. For the passage retrieval task, we submit a single run based on an ensemble of eight Duet models.

***Keywords*** Deep learning · Neural information retrieval · Ad-hoc retrieval

## 1 Introduction

The Duet architecture was proposed by Mitra et al. [2017] for document ranking. Fig. 7 from The original paper show that the retrieval effectiveness of the model is still improving as the size of the training data approaches $2^{17}$ samples. The training data employed in that paper is a proprietary dataset from Bing. A similar plot was later reproduced on a public benchmark by Nanni et al. [2017], but in the context of a passage ranking dataset with synthetic queries. Variations of the Duet model [Mitra and Craswell, 2019, Mitra et al., 2019, Cohen et al., 2018] have since then been evaluated on other public passage ranking datasets. However, the lack of large scale training data prevented the public evaluation of Duet for document ranking.

The deep learning track at TREC 2019 makes large training datasets—suitable for traininig deep models with large number of learnable parameters—publicly available in the context of a document ranking and a passage ranking tasks. We benchmark the Duet model on both tasks.

In the context of the document ranking task, we adapt the Duet model to ingest a "multiple field" view of the documents, based on findings from Zamani et al. [2018]. We refer to this new architecture as Duet with Multiple Fields (DuetMF) in the paper. Furthermore, we combine the relevance estimates from DuetMF with several other traditional and neural retrieval methods in a learning-to-rank (LTR) [Liu, 2009] framework.

For the passage ranking task, we submit a single run based on an ensemble of eight Duet models. The architecture and the training scheme resembles that of the "Duet V2 (Ensembled)" baseline listed on the MS MARCO leaderboard[1].

## 2 TREC 2019 deep learning track

The TREC 2019 deep learning track introduces: (i) a document retrieval task and (ii) a passage retrieval task. For both tasks, participants are provided a set of candidates—100 documents and 1000 passages, respectively—per query that should be ranked. Participants can choose to either rerank provided candidates or retrieve from the full collection.

For the passage retrieval task, the track reuses the set of 500K+ manually-assessed binary training labels released as part of the Microsoft Machine Reading COmprehension (MS MARCO) challenge [Bajaj et al., 2016]. For the document retrieval task, the passage-level labels are transferred to their corresponding source documents—producing a training dataset of size close to 400K labels.

---

[1] http://www.msmarco.org/leaders.aspx

Table 1: Official TREC results. The recall metric is computed at position 100 for the document retrieval task and at position 1000 for the passage retrieval task.

| Run description | Run ID | Subtask | MRR | NDCG@10 | MAP | Recall |
|---|---|---|---|---|---|---|
| **Document retrieval task** | | | | | | |
| LTR w/ DuetMF as feature | ms_ensemble | fullrank | 0.876 | 0.578 | 0.237 | 0.368 |
| DuetMF model | ms_duet | rerank | 0.810 | 0.533 | 0.229 | 0.387 |
| **Passage retrieval task** | | | | | | |
| Ensemble of 8 Duet models | ms_duet_passage | rerank | 0.806 | 0.614 | 0.348 | 0.694 |

For evaluation, a shared test set of 200 queries is provided for both tasks, of which two different overlapping set of 43 queries were later selected for manual NIST assessments corresponding to the two tasks.

Full details of all datasets is available on the track website[2] and in the track overview paper [Craswell et al., 2019].

## 3 Methods and results

The Duet model proposed by Mitra et al. [2017] employs two deep neural networks trained jointly towards a retrieval task: (i) The "distributed" sub-model learns useful representations of text for matching and (ii) the "local" sub-model estimates relevance based on patterns of exact term matches between query and document. Mitra and Craswell [2019] propose several modifications to the original Duet model that show improved performance on the MS MARCO passage ranking challenge. We adopt the updated Duet model from Mitra and Craswell [2019] and incorporate additional modifications, in particular to consider multiple fields for the document retrieval task. Table 1 summarizes the official evaluation results for all three runs.

**Duet model with Multiple Fields (DuetMF) for document ranking.** Zamani et al. [2018] study neural ranking models in the context of documents with multiple fields. In particular, they make the following observations:

Obs. 1: It is more effective to summarize the match between query and individual document fields by a vector—as opposed to a single score—before aggregating to estimate full document relevance to the query.

Obs. 2: It is better to learn different query representations corresponding to each document field under consideration.

Obs. 3: Structured dropout (*e.g.*, field-level dropout) is effective for regularization during training.

We incorporate all of these ideas to modify the Duet model from Mitra and Craswell [2019]. The updated model is shown in Fig. 1.

Documents in the deep learning track dataset contains three text fields: (i) URL, (ii) title, and (iii) body. We employ the Duet architecture to match the query against each individual document fields. In line with Obs. 1 from [Zamani et al., 2018], the field-specific Duet architecture outputs a vector instead of a single score. We do not share the parameters of the Duet architectures between the field-specific instances based on Obs. 2. Following Obs. 3, we introduce structured dropouts at different stages of the model. We randomly dropout each of the local sub-models for $50\%$ of the training samples. Similarly, we also dropout different combinations of field-level models uniformly at random—taking care that at least one field-level model is always retained.

We consider the first 20 terms for queries and for document URLs and titles. For document body text, we consider the first 2000 terms. Similar to Mitra and Craswell [2019], we employ pretrained word embeddings as the input text representation for the distributed sub-models. We train the word embeddings using a standard word2vec [Mikolov et al., 2013] implementation in FastText [Joulin et al., 2016] on a combination of the MS MARCO document corpus and training queries.

Similar to previous work [Mitra et al., 2017, Mitra and Craswell, 2019], the query and document field embeddings are learned by deep convolutional-pooling layers. We set the hidden layer size at all stages of the model to 300 and dropout rate for different layers to 0.5. For training, we employ the RankNet loss [Burges et al., 2005] over $< q, d_{pos}, d_{neg} >$ triples and the Adam optimizer [Kingma and Ba, 2014]—with a minibatch size of 128 and a learning rate of 0.0001 for training. We sample $d_{neg}$ uniformly at random from the top 100 candidates provided that are not positively labeled. When employing structured dropout, the same sub-models are masked for both $d_{pos}$ and $d_{neg}$.

In light of the recent success of large pretrained language models—*e.g.*, [Nogueira and Cho, 2019]—we also experiment with an unsupervised pretraining scheme using the MS MARCO document collection. The pretraining is performed
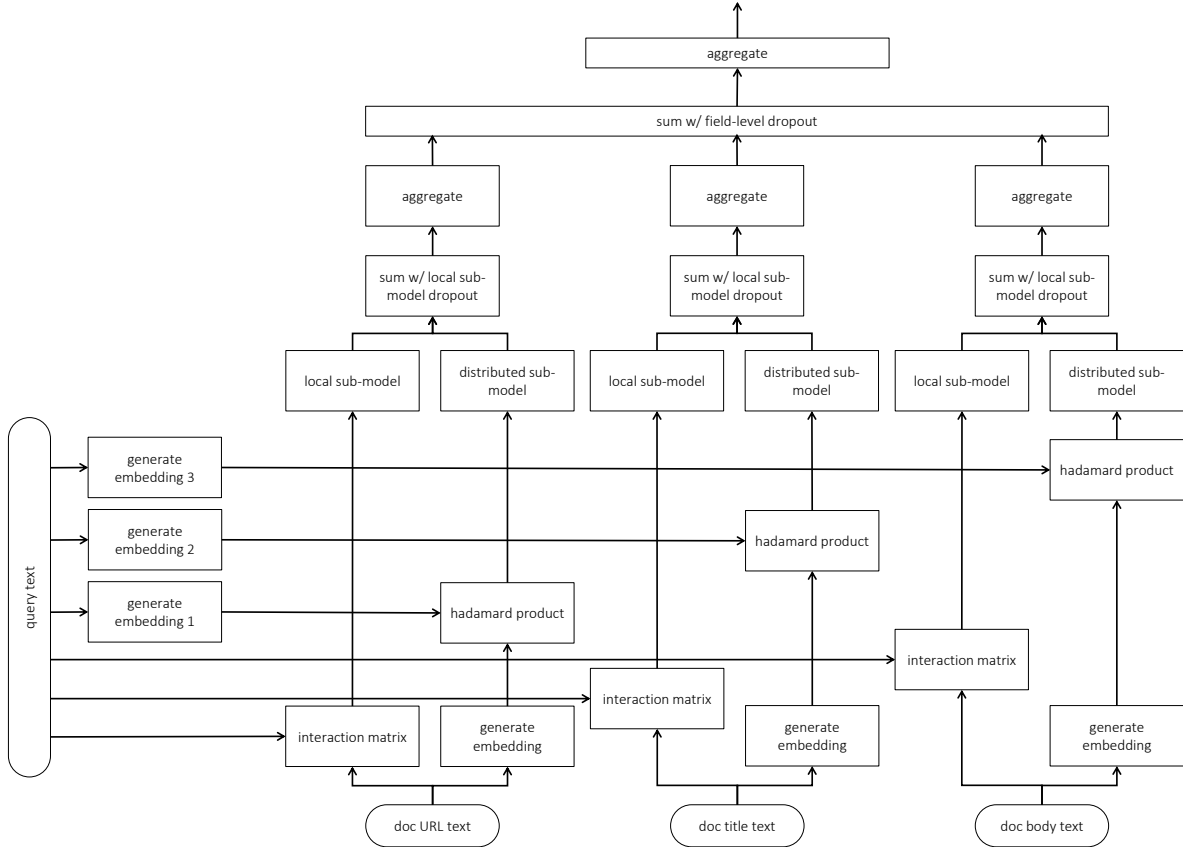
---

[2]`https://microsoft.github.io/TREC-2019-Deep-Learning/`

Figure 1: The modified Duet model (DuetMF) that considers multiple document fields.

over $< q_{\text{pseudo}}, d_{\text{pos}}, d_{\text{neg}} >$—where $d_{\text{pos}}$ and $d_{\text{neg}}$ are randomly sampled from the collection and a pseudo-query $q_{\text{pseudo}}$ is generated by picking the URL or the title of $d_{\text{pos}}$ randomly (with equal probability) and masking the corresponding field on the document side for both $d_{\text{pos}}$ and $d_{\text{neg}}$. We see faster convergence during supervised training when the DuetMF model is pretrained in this fashion on the MS MARCO document collection. We posit that a more formal study should be performed in the future on pretraining Duet models on large collections, such as Wikipedia and the BookCorpus [Zhu et al., 2015].

**Learning-to-rank model for document ranking.** We train a neural LTR model with two hidden layers—each with 1024 hidden nodes. The LTR run reranks a set of 100 document candidates retrieved by query likelihood (QL) [Ponte and Croft, 1998] with Dirichlet smoothing ($\mu = 1250$) [MacKay and Peto, 1995]. Several ranking algorithms based on neural and inference networks act as features: (i) DuetMF, (ii) Sequential Dependence Model (SDM) [Metzler and Croft, 2005], and (iii) Pseudo-Relevance Feedback (PRF) [Lavrenko and Croft, 2001, Lavrenko, 2008], (iv) BM25, [Robertson et al., 2009], and (v) Dual Embedding Space Model (DESM) [Nalisnick et al., 2016, Mitra et al., 2016].

We employ SDM with an order of 3, combine weight of 0.90, ordered window weight of 0.034, and an unordered window weight of 0.066 as our base candidate scoring function. We use these parameters to retrieve from the target corpus as well as auxiliary corpora of English language Wikipedia (`enwiki-20180901-pages-articles-multistream.xml.bz2`), LDC Gigaword (`LDC2011T07`). For PRF, initial retrievals—from either of the target, wikipedia, or gigaword corpora—adopted the SDM parameters above, however are used to rank 75-word passages with a 25-word overlap. These passages are then interpolated using the top $m$ passages and standard relevance modeling techniques, from which we select the top 50 words to use as an expanded query for the final ranking of the target candidates. We do not explicitly adopt RM3 [Abdul-Jaleel et al., 2004] because our LTR model implicitly combines our initial retrieval score and score from the expanded query. All code for the SDM and PRF feature computation is available at `https://github.com/diazf/indri`.

We evaluate two different BM25 models with hyperparameters $< k_1 = 0.9, b = 0.4 >$ and $< k_1 = 3.44, b = 0.87 >$.

Corresponding to each of the DuetMF, SDM, PRF, and BM25 runs we generate two features based on the score and the rank that the model predicts for a document *w.r.t.* the target query.

We generate eight features by comparing the query against two different document fields (title and body) and using different DESM similarity estimates (INxIN, INxOUT, OUTxIN, OUTxOUT).

Lastly, we add couple of features based on query length and domain quality—where the latter is defined simply as a ratio between how often documents from a given domain appear in the positively labeled training data and in the overall document collection.

**Ensemble of Duet models for passage ranking.** For the passage ranking task, we adopt the exact same model and training procedure from [Mitra and Craswell, 2019]. Our final submission is an ensemble of eight Duet models.

## 4 Discussion and conclusion

One of the main goals of the deep learning track is to create a public reusable dataset for benchmarking the growing body of neural information retrieval literature [Mitra and Craswell, 2018]. We submit three runs based on the Duet architecture for the two—document and passage—retrieval tasks. Our main goal is to enrich the set of pooled documents for NIST assessments with documents that a Duet based architecture is likely to rank highly. As a secondary goal, we are also interested in benchmarking Duet against other state-of-the-art neural and traditional methods. A more detailed comparison of the performance of these Duet runs with other TREC submissions is provided in the track overview paper [Craswell et al., 2019].

## References

Nasreen Abdul-Jaleel, James Allan, W Bruce Croft, Fernando Diaz, Leah Larkey, Xiaoyan Li, Mark D Smucker, and Courtney Wade. Umass at trec 2004: Novelty and hard. 2004.

Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, et al. Ms marco: A human generated machine reading comprehension dataset. *arXiv preprint arXiv:1611.09268*, 2016.

Chris Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Greg Hullender. Learning to rank using gradient descent. In *Proceedings of the 22nd international conference on Machine learning*, pages 89–96. ACM, 2005.

Daniel Cohen, Bhaskar Mitra, Katja Hofmann, and W Bruce Croft. Cross domain regularization for neural ranking models using adversarial learning. In *Proc. SIGIR*, pages 1025–1028. ACM, 2018.

Nick Craswell, Bhaskar Mitra, Emine Yilmaz, and Daniel Campos. Overview of the trec 2019 deep learning track. In *TREC (to appear)*, 2019.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*, 2016.

Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Victor Lavrenko. *A generative theory of relevance*, volume 26. Springer Science & Business Media, 2008.

Victor Lavrenko and W Bruce Croft. Relevance based language models. In *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 120–127. ACM, 2001.

Tie-Yan Liu. Learning to rank for information retrieval. *Foundation and Trends in Information Retrieval*, 3(3):225–331, March 2009.

David JC MacKay and Linda C Bauman Peto. A hierarchical dirichlet language model. *Natural language engineering*, 1(3):289–308, 1995.

Donald Metzler and W Bruce Croft. A markov random field model for term dependencies. In *Proc. SIGIR*, pages 472–479. ACM, 2005.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Proc. NIPS*, pages 3111–3119, 2013.

Bhaskar Mitra and Nick Craswell. An introduction to neural information retrieval. *Foundations and Trends® in Information Retrieval (to appear)*, 2018.

Bhaskar Mitra and Nick Craswell. An updated duet model for passage re-ranking. *arXiv preprint arXiv:1903.07666*, 2019.

Bhaskar Mitra, Eric Nalisnick, Nick Craswell, and Rich Caruana. A dual embedding space model for document ranking. *arXiv preprint arXiv:1602.01137*, 2016.

Bhaskar Mitra, Fernando Diaz, and Nick Craswell. Learning to match using local and distributed representations of text for web search. In *Proc. WWW*, pages 1291–1299, 2017.

Bhaskar Mitra, Corby Rosset, David Hawking, Nick Craswell, Fernando Diaz, and Emine Yilmaz. Incorporating query term independence assumption for efficient retrieval and ranking using deep neural networks. *arXiv preprint arXiv:1907.03693*, 2019.

Eric Nalisnick, Bhaskar Mitra, Nick Craswell, and Rich Caruana. Improving document ranking with dual word embeddings. In *Proc. WWW*, 2016.

Federico Nanni, Bhaskar Mitra, Matt Magnusson, and Laura Dietz. Benchmark for complex answer retrieval. In *Proc. ICTIR*. ACM, 2017.

Rodrigo Nogueira and Kyunghyun Cho. Passage re-ranking with bert. *arXiv preprint arXiv:1901.04085*, 2019.

Jay M Ponte and W Bruce Croft. A language modeling approach to information retrieval. In *Proc. SIGIR*, pages 275–281. ACM, 1998.

Stephen Robertson, Hugo Zaragoza, et al. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389, 2009.

Hamed Zamani, Bhaskar Mitra, Xia Song, Nick Craswell, and Saurabh Tiwary. Neural ranking models with multiple document fields. In *Proc. WSDM*, pages 700–708. ACM, 2018.

Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pages 19–27, 2015.