

Shallow Cross-Encoders for Low-Latency Retrieval

Aleksandr V. Petrov, Sean MacAvaney, and Craig Macdonald

University of Glasgow, Glasgow, UK
a.petrov.1@research.gla.ac.uk
{sean.macavaney;craig.macdonald}@glasgow.ac.uk

Abstract. Transformer-based Cross-Encoders achieve state-of-the-art effectiveness in text retrieval. However, Cross-Encoders based on large transformer models (such as BERT or T5) are computationally expensive and allow for scoring only a small number of documents within a reasonably small latency window. However, keeping search latencies low is important for user satisfaction and energy usage. In this paper, we show that weaker shallow transformer models (i.e. transformers with a limited number of layers) actually perform *better* than full-scale models when constrained to these practical low-latency settings, since they can estimate the relevance of more documents in the same time budget. We further show that shallow transformers may benefit from the generalised Binary Cross-Entropy (gBCE) training scheme, which has recently demonstrated success for recommendation tasks. Our experiments with TREC Deep Learning passage ranking queriesets demonstrate significant improvements in shallow and full-scale models in low-latency scenarios. For example, when the latency limit is 25ms per query, MonoBERT-Large (a cross-encoder based on a full-scale BERT model) is only able to achieve NDCG@10 of 0.431 on TREC DL 2019, while TinyBERT-gBCE (a cross-encoder based on TinyBERT trained with gBCE) reaches NDCG@10 of 0.652, a +51% gain over MonoBERT-Large. We also show that shallow Cross-Encoders are effective even when used without a GPU (e.g., with CPU inference, NDCG@10 decreases only by 3% compared to GPU inference with 50ms latency), which makes Cross-Encoders practical to run even without specialised hardware acceleration.

1 Introduction

The introduction of the Transformer [35] neural network architecture, and especially pre-trained language models that use Transformers (such as BERT [7]), has been transformative for the IR field; for example, Nogueira et al. [27] improved MRR@10 on the MS-MARCO dev set by 31% with the help of a BERT-based model. Although there are a variety of ranking architectures used within IR (e.g., dense Bi-Encoders [14,18,40], sparse Bi-Encoders [9,22], and late interaction models [13,15]), the best results for document re-ranking are typically achieved with the help of *Cross-Encoders* [14] – a family of models which encode both the query and the document simultaneously as a single textual input [41]. Aside from their high in-domain precision, Cross-Encoders tend to be more robust when generalising across retrieval tasks/domains [33]. Although Cross-Encoders can only practically be used as re-ranking models, limitations in their first-stage recall can be efficiently mitigated using pseudo-relevance feedback [23]. Further, Cross-Encoders can typically be fine-tuned from scratch (i.e., starting from the checkpoint of a foundational model, such as BERT).

Despite these benefits, the application of Cross-Encoders in production retrieval systems is still limited. Cross-Encoders require a model inference for each query-document pair and, therefore, struggle with high computational complexity and high latency [24]. In real-world search systems, high latency negatively affects key performance metrics, such as the number of clicks, revenue, and user satisfaction [16, Ch. 5]. Further, high latencies tend to be correlated with higher energy usage, resulting in negative impacts on the climate [32].

The high computational complexity and resulting latency of Cross-Encoder models motivated researchers to investigate Bi-Encoder [14] models. These models separately encode the query and the document, and then estimate relevance score using an inexpensive operation over the encoded representations (e.g. cosine similarity [30] or the MaxSim operation [15]). By pre-computing the document representations offline and using a variety of approaches to accelerate retrieval [17], Bi-Encoders can achieve low retrieval latency. However, this comes at other costs. For instance, Bi-Encoders are markedly more complicated to train than Cross-Encoders, typically relying on knowledge distillation from other models (e.g., [19]), training data balancing (e.g., [12]), and/or hard negative mining (e.g., [40]). Further, Bi-Encoders must pre-encode all documents in the collection and keep the encoded versions of all documents in memory.

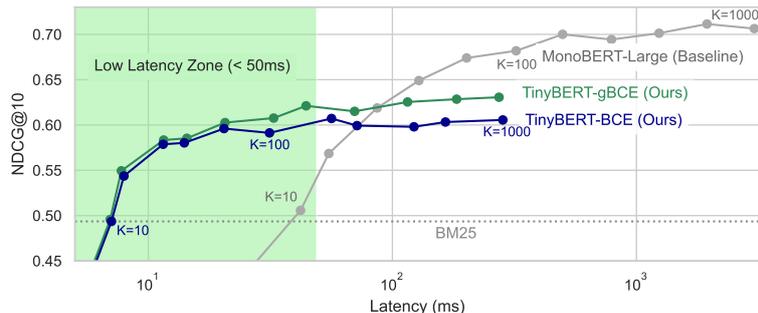


Fig. 1: Latency/NDCG tradeoffs on the TREC-DL2020 queryset when varying the number of retrieved candidates K .

This may be an inefficient strategy given the long-tail distribution of document popularity in a corpus. Indeed, if most documents are never retrieved, building their dense representations and keeping them in memory wastes resources. Moreover, document encoding costs for Bi-Encoders must also be incurred every time the retrieval model changes (e.g., when a model is re-trained to reflect new search trends.) Finally, Bi-Encoders often struggle to transfer across retrieval tasks and domains [33].

Therefore, in this paper, we investigate *shallow Cross-Encoders* (Cross-Encoders with a limited number of transformer layers) as a solution for low-latency search. Shallow Cross-Encoders are much smaller than full-scale models and require much fewer computations than full-scale models. Therefore, these models can score many times more documents within the low-latency window¹ than the full-scale models; this means that in low-latency scenarios, they can rerank more candidates and, ultimately, have better effectiveness than the full-scale models.

Note that there are a number of approaches for reducing the latency of ranking models, such as *dynamic pruning* or the use of *approximate nearest neighbour* indices (an overview of these methods can be found in [3]); however, most of these methods are not applicable to cross-encoder models. Indeed, in this paper, we focus on two main ways of reducing latency in the cross-encoder: (i) reducing the model’s size and (ii) reducing the number of candidate documents, and we evaluate which of these ways achieves better effectiveness in the low-latency scenario.

Training effective Shallow Cross-Encoders, despite their promise of good efficiency, presents a significant challenge (this is in contrast to full-scale cross-encoders, which, as we previously mentioned, are relatively easy to train). MacAvaney et al. [24] first explored limiting the depths of a transformer network for ranking, but was only able to reduce the depth of the network to 5 layers without substantial effectiveness degradation. Other successful attempts to train shallower Cross-Encoder models have required applying complicated knowledge-distillation techniques. For example, the popular Sentence-Transformer package [30] provides several shallow models² and reports model performance competitive to larger models. Unfortunately, no academic paper is associated with these checkpoints, and the exact details of training these models are unclear. Our analysis of the training code shows that these models were trained using a knowledge distillation setup from an ensemble of full-scale models, loosely following the process described by Hofstätter et al. [11], which assumes the existence of such an ensemble in the first place. While resulting in effective checkpoints for the MSMARCO dataset, we argue that a training strategy that requires training an ensemble of full-scale language models before training the shallow model is hard to replicate for other settings (e.g. for different languages or other datasets). Indeed, in [37], the challenges of reproducing knowledge distillation models with “dependency chains” of models was found to be challenging - and hence, we argue that training shallow Cross-Encoders this way is rather art than science.

In contrast, this paper proposes a direct training method for shallow Cross-Encoders, which does not resort to complex techniques such as Knowledge Distillation. Our approach based on the gBCE training scheme [29] has been recently shown to improve the effectiveness of transformer-based models for recommender systems. The training scheme consists of two key components: (i) an increased number of sampled negative instances for each labelled positive instance and (ii) the gBCE loss function, which counters the effects of negative sampling (which is typically used to train Cross-Encoder models). Our

¹ In [6], Google researchers argued that for a smooth user experience, total search latency should be kept under 100ms. This includes time for network round-trips, page rendering, and other overheads. Therefore, this paper uses a 50ms cutoff for defining *low-latency retrieval*, leaving the remainder of the time to these other overheads.

² <https://www.sbert.net/docs/pretrained-models/ce-msmarco.html>

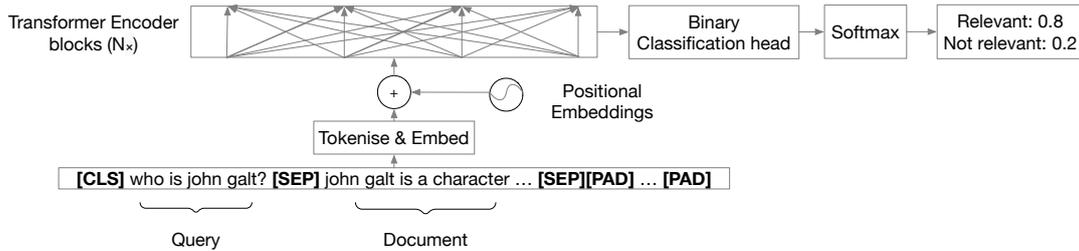


Fig. 2: A typical BERT-based [7] Cross-Encoder. Note that the structure can be adapted to other transformers with slight modification.

experiments show that both of these components positively affect model training. Figure 1 summarises some of the main findings of this paper. The figure illustrates efficiency/effectiveness tradeoffs of two very small cross-encoders (2 layers, TinyBERT [34]) and a full-scale MonoBERT-Large model [27]. One of the small Cross-Encoders is trained using the gBCE training scheme, and the other is trained using traditional BCE. As the figure shows, while MonoBERT-Large is more effective when allowed latency is high (i.e. it is allowed to re-rank many documents), within the low latency zone ($< 50\text{ms}$ latency), shallow Cross-Encoders are more effective. Moreover, TinyBERT trained with gBCE is more effective compared to TinyBERT trained with traditional BCE.

Overall, the contributions of this paper can be summarised as follows: (i) We propose a simple and replicable method for training shallow Cross-Encoders based on the gBCE training scheme, which does not rely on knowledge distillation; (ii) We analyse the efficiency/effectiveness tradeoffs of Cross-Encoders of different sizes and demonstrate that shallow Cross-Encoders are preferable to full-size models under low-latency constraints; (iii) We demonstrate that shallow Cross-Encoders are efficient and effective even when used without a GPU.

We note that outside academia, there is interest in shallow Cross-Encoders, which is expressed in several industrial blog posts.³⁴ This interest makes us believe that our research has high potential to be adopted by industry and serve as a basis for further study.

The rest of the paper is organised as follows: Section 2 provides an overview of the Cross-Encoder architecture and Efficiency/Effectiveness tradeoffs arising from this architecture, Section 3 describes training scheme for shallow Cross-Encoders, Section 4 contains experimental evaluation of the efficiency/effectiveness tradeoffs for shallow Cross-Encoders, and Section 5 contains final remarks.

2 Efficiency/Effectiveness Tradeoffs in Cross-Encoders

A *Cross-Encoder* [14] is a model that jointly encodes a query-document pair using a single language model. Figure 2 illustrates a typical Transformer [35] encoder-based Cross-Encoder. The input to the Cross-Encoder model consists of a concatenation of the query with the document, joined with the help of some special tokens. In our example, there are three different special tokens: (i) $[CLS]$ token is added to the beginning of the input; a contextualised representation of this token is then used for classification; (ii) $[SEP]$ token is added at the end of both text and the document; these tokens separate help the model to separate different groups of text; (iii) a series of $[PAD]$ tokens is added at the end of the sequence to equalise the input length of each document in the batch. The input is then encoded using a standard Transformer encoder network, which consists of an embedding layer, positional embeddings, and a stack of Transformer blocks. For brevity, we omit details of the Transformer encoder network and refer to the original papers [7,35]. The output of the Transformer encoder consists of a sequence of embedding, where each embedding is a contextualised representation of each input token. In particular, for classification tasks, the representation of the $[CLS]$ token is usually used to represent the whole input sequence. The task is usually cast as a binary classification for information retrieval with two possible outcomes (relevant/not relevant). Typically, these probabilities are obtained by passing the $[CLS]$ representation through a simple Feed-Forward network with two outputs and then apply a Softmax operation over these two outputs.

³ <https://blog.vespa.ai/pretrained-transformer-language-models-for-search-part-4/>

⁴ <https://towardsdatascience.com/tinybert-for-search-10x-faster-and-20x-smaller-than-bert-74cd1b6b5aec>

As Cross-Encoders jointly encode the query and the document, the scoring of K candidate documents requires applying the model for inference K times⁵. Assuming each inference requires λ milliseconds, we can infer an upper bound on the number of documents that can be scored within the latency window ω :

$$K \leq \frac{\omega}{\lambda} \quad (1)$$

Equation (1) defines the tradeoff between the latency window ω and the number of scored documents K . In practice, this tradeoff limits the use of Cross-Encoders to a *re-ranking* scenario, where the Cross-Encoder is applied to a relatively small number of candidates retrieved with a lightweight first-stage model, such as BM25 [31]. Moreover, reducing the latency window ω decreases the number of retrieved documents that can be scored, decreasing the result’s recall and reducing overall model effectiveness. Overall, we can say that there exists a tradeoff between model performance $Q(K)$ and the latency window ω :

$$Q(K) \bowtie \omega \quad (2)$$

where the \bowtie symbol denotes a dependency between Q and ω . The exact form of this dependency is unknown, and investigating the properties of this dependency is one of the main goals of this paper. Equation (1) also shows that the number of scored documents K can be increased if we decrease model inference time λ . In the case of Transformer-based Cross-Encoders, we can decrease the inference time by limiting the number of transformer blocks in the model and/or by limiting the computational complexity of each block by reducing embedding sizes and the number of attention heads (i.e. making the Cross-Encoder *shallow*). On the other hand, many recent publications [27,28,38] show that larger models are more effective for document re-ranking (without considering the latency).

In summary, decreasing the Cross-Encoder size has two effects on overall model effectiveness:

- **Positive effect:** Decreasing model size allows for scoring more documents, which leads to **increased model effectiveness** according to Equation (1).
- **Negative effect:** Decreasing model size hinders accuracy on each individual (query/document) pair; as a result, overall **model effectiveness decreases**.

To the best of our knowledge, no published research has analysed which of these two effects dominates in a low-latency scenario. This work aims to close this gap. In particular, we aim to verify the following research hypothesis:

Hypothesis H1

When reducing the latency window ω , the ability of shallow Cross-Encoder to score and rank more candidate documents within the given latency window results in higher effectiveness than when ranking less candidate documents with more accurate full-scale models.

Note that hypothesis **H1** is somewhat counter-intuitive. Indeed, it has been consistently shown that larger language models are more effective for documents re-ranking [28,27,38] compared to the smaller ones. However, contrary to these findings, hypothesis **H1** states that smaller language models can be **more effective** compared to the larger ones when the latency window is limited.

To analyse whether or not this hypothesis holds, we need to establish an effective training scheme for shallow Cross-Encoders, which we do in the next section. We then analyse Hypothesis **H1** experimentally in Section 4.

3 Training Shallow Cross-Encoders with gBCE

Usually, Cross-Encoder models output estimated relevance *scores* (logits), which can be converted to *probabilities*. For example, the architecture of a BERT-based Cross-Encoder on Figure 2 outputs positive and negative scores s^+ and s^- , and the probability of the document being relevant to the query is then computed using the Softmax(\cdot) transformation:

$$p^+ = \frac{e^{s^+}}{e^{s^+} + e^{s^-}} \quad (3)$$

⁵ For simplicity, we omit batching in this reasoning. With a slight tweaking, it remains valid for batching as well (e.g. inference time should be divided by the batch size).

As the model outputs can be converted into probabilities, Binary Cross-Entropy Loss can be used to train the model:

$$\mathcal{L}_{BCE} = - [y \cdot \log(p^+) + (1 - y) \cdot \log(1 - p^+)] \quad (4)$$

where y is the ground truth relevance judgment for a given query-document pair, BCE loss drives the model to minimise the KL divergence between ground truth relevancy and predicted probabilities $D_{KL}(y||p)$. Therefore, p^+ will converge to the “real” probability (the probability that a user finds the document relevant before the judgement has been made; see [29] for proofs). BCE is a very popular choice and has been used to train many effective Cross-Encoder models [28,36]. However, an underlying assumption of BCE is that the distribution of positive/negative samples during training and inference match each other. In practice, there are many more negative documents than positives. Due to computational and memory constraints, it is not feasible to score all possible negatives with large-scale models and therefore, researchers employ *negative sampling* techniques. For example, the popular MonoT5 model [28] samples just one negative (non-relevant) document for each query during training.

A recent publication [29] has shown that negative sampling coupled with the BCE loss leads to the *overconfidence* problem: the estimated probability for positives becomes too high, and the model training becomes unstable, leading to poorer performance compared to the models trained without negative sampling.

Our initial experiments have shown that overconfidence does not cause effectiveness degradation of full-scale Cross-Encoder models, such as MonoT5; we speculate that the use of pre-trained checkpoints works as a strong model regulariser (see [10, Ch. 15] for an intuition). However, our experiments show (see Section 4) that overconfidence is indeed a problem in shallow cross-encoders, and it has to be mitigated to achieve high effectiveness.

To mitigate the overconfidence problem for recommender systems, the authors [29] propose the *gBCE* training scheme. The gBCE training scheme consists of two components:

1. Increased number of negative samples per positive.
2. Generalised Binary-Crossentropy (gBCE) loss function instead of classic BCE.

gBCE, parametrised by a parameter β , is defined as:

$$\begin{aligned} \mathcal{L}_{gBCE}^\beta &= - [y \cdot \log((p^+)^\beta) + (1 - y) \cdot \log(1 - p^+)] \\ &\quad (\beta \text{ can be taken out of the log}) \\ &= - [y \cdot \beta \cdot \log(p^+) + (1 - y) \cdot \log(1 - p^+)] \end{aligned} \quad (5)$$

Parameter β controls model confidence: when $\beta = 1$, gBCE becomes regular BCE; however, decreasing β decreases the model’s tendency to predict high probabilities. The authors of [29] proposed to control β indirectly with the help of a calibration parameter t :

$$\beta(t) = \alpha \left(t \left(1 - \frac{1}{\alpha} \right) + \frac{1}{\alpha} \right) \quad (6)$$

where α is the negative sampling rate: $\alpha = \frac{|D_k^-|}{|D^-|}$ (number of sampled negatives as a fraction of the overall number of retrieved negative candidates). For example, if, at each training step, we use 10 sampled negatives for each positively labelled relevant document, and the overall number of retrieved negatives from the first stage retriever is 1000, then the sampling rate is $\alpha = 0.01$.

In summary, at each training step, we sample a batch of B positive (query/ document) pairs, then use the first stage retriever model (we use BM25) to retrieve 1000 negatives. After that, we sample K negatives out of 1000 candidates and add these K negative (query/document) pairs for each query to the training batch. Overall, each training batch contains $B \cdot (K + 1)$ (query/document) pairs with B positives and $B \cdot K$ negatives. We then perform a standard gradient descent update step using the gBCE loss function.

This concludes the description of the gBCE training scheme for Cross-Encoder models. We now turn to the experimental evaluation of shallow Cross-Encoders for low-latency retrieval.

4 Experiments

We design our experiments to answer the following research questions:

Table 1: Salient characteristics of experimental models. * For the MonoT5 model, “Transformer Layers” is the combined number of encoder and decoder layers.

Model Type	Model	Transformer Layers	Embedding Size	Attention Heads	Vocab Size	Number of Parameters	Chekpont File Size
Shallow Cross-Encoders	TinyBERT	2	128	2	30522	4,386,178	17 Mb
	MiniBERT	4	256	4	30522	11,171,074	43 Mb
	SmallBERT	4	512	8	30522	28,764,674	110 Mb
Full-Size baselines	MonoT5-Base	24*	768	12	32128	222,903,552	850 Mb
	MonoBERT-Large	24	1024	16	30522	335,143,938	1.2 Gb

RQ1 How does the model size affect efficiency and effectiveness of Cross-Encoders?

RQ2 What are the effects of an increased number of negatives and calibration parameter t in the gBCE training strategy applied to shallow Cross-Encoders?

RQ3 What is the efficiency/effectiveness tradeoff of shallow Cross-Encoders when using CPU-only inference?

4.1 Experimental Setup

Frameworks We use PyTerrier [25] as our main experimental framework, the PyTerrier-Pisa [20,26] plugin for a low-latency BM25 first-stage retriever (making use of a memory-mapped BlockMax-WAND index [8]), and the ir-measures [21] library for evaluation metrics. We use model implementations from the HuggingFace Transformers [39] library v4.30.2.⁶

Datasets For training shallow Cross-Encoders, we use the large-scale MSMARCO dataset [1]. For training, we select queries from the *train* section of the dataset, for which the BM25 retriever retrieves at least one relevant document within the top 1000 results⁷. After pre-filtering, the dataset contains 436299 queries. Out of this pre-filtered dataset, we randomly select 200 queries into a *validation set*, which are held out from training and used for monitoring effectiveness metrics during training and for early stopping. For evaluation, we use TREC Deep Learning track (denoted TREC-DL) queries from the 2019 [5] and 2020 [4] tracks. Both querysets are based on queries from MSMARCO but with comprehensive assessments for reliable evaluation.

Hardware All our experiments use a computer with a Ryzen 5950x CPU, 128Gb DDR-4 memory, NVIDIA RTX 4090 GPU, and a Samsung 980 SSD.

Models As the backbone architecture for shallow Cross-Encoders, we use pre-trained versions of BERT [7]. Namely, we use TinyBERT, MiniBERT, and SmallBERT checkpoints⁸ provided by Turc et al. [2,34]. The smallest model, TinyBERT, has just two Transformer Encoder layers and an embedding size of 128, and the largest model, SmallBERT, has four Transformer Layers and an embedding size of 512. We also use two full-size pre-trained Cross-Encoders as the baselines: MonoBERT-Large [27] and MonoT5-Base [28]. For both models, we use official pre-trained checkpoints provided by the authors^{9,10}. Table 1 describes salient characteristics of all experimental models. For the inference of all BERT-based models, we use a batch size of 8 – the largest batch size with which we did not experience memory-related issues with the MonoBERT-Large model. For MonoT5, we use a batch size of 64.

During training, following best practices [10, Ch.7], we use an early stopping mechanism to ensure model convergence. In particular, after every 600 training batches, we measure NDCG@10 on the validation set and stop training when validation does not improve for 200 validation steps.

4.2 RQ1. Efficiency/Effectiveness tradeoffs

Our first research question aims to verify Hypothesis H1 and test whether or not shallow Cross-Encoders provide better efficiency/effectiveness tradeoffs compared to the full-scale models.

⁶ Source code for this paper can be found at

<https://github.com/asash/shallow-cross-encoders>

⁷ We do not use the standard MSMARCO triplets file because it only contains one negative per query, and for gBCE training scheme we need up to 128 negatives.

⁸ <https://huggingface.co/prajjwall1/bert-tiny>

⁹ <https://huggingface.co/castorini/monot5-base-msmarco-10k>

¹⁰ <https://huggingface.co/castorini/monobert-large-msmarco-finetune-only>

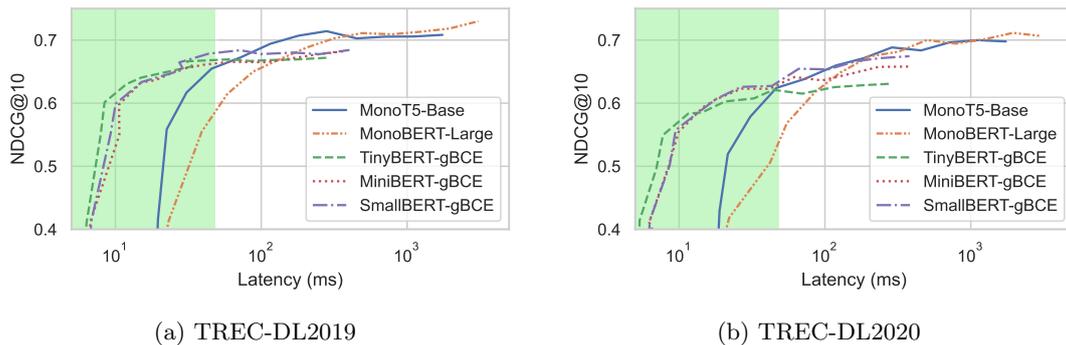


Fig. 3: Latency/NDCG tradeoffs of experimental models when varying the number of candidates from BM25 between 1 and 1000. The shaded area represents the low-latency zone (latency less than 50ms).

To evaluate the tradeoffs, we train TinyBERT, MiniBERT and SmallBERT using the gBCE training scheme. Following the recommendations in [29], we use 128 negatives per positive, and we set calibration parameter $t = 0.75$ for training. We then evaluate shallow models and pre-trained full-scale models using a variable number of candidates from the first-stage BM25 retriever. In our experiments, we vary the number of candidate documents between 1 and 1000. For each number of candidates, we measure model effectiveness using the NDCG@10 metric and the efficiency by measuring latency in milliseconds.

Figure 3 illustrates the efficiency/effectiveness tradeoffs for all experimental models when varying the number of retrieved BM25 candidates, K . Note that the latency includes overheads, such as the BM25 retrieval time and tokenization. From the figure, we see that on both the TREC-DL2019 and TREC-DL2020 querysets, shallow Cross-Encoders (TinyBERT-gBCE, MiniBERT-gBCE, SmallBERT-gBCE) outperform the larger full-scale models in the low-latency zone (latency less than 50ms). For example, on the TREC-DL2019 queryset, for the maximum latency of 25ms, TinyBERT-gBCE achieves NDCG@10 of 0.652, a +13.7% improvement over MonoT5 (NDCG@10 0.573) and +51% improvement over MonoBERT-Large (NDCG@10 0.4316). However, if we allow large latencies, full-scale models outperform shallow Cross-Encoders: for example, the best NDCG@10 on TREC-DL2020 achieved by MonoBERT-Large is 0.711, whereas the best NDCG@10 achieved by TinyBERT-gBCE is 0.630 (-12%). Note that the difference between shallow models in the low-latency zone is relatively small; for example, at TREC-DL2019 at 10ms latency, all shallow models achieve an NDCG@10 of 0.60. However, the smallest model (TinyBERT-gBCE, which has just two transformer layers) consistently has better effectiveness than the larger models with latencies less than 10ms; therefore, we argue that at very small latency requirements, the usage of the smallest model is preferable, as in addition to the best possible performance it also has lowest memory requirements.

Note that these effectiveness/efficiency tradeoffs are specific to our hardware. Also, the tradeoffs can be improved using better-optimised versions of the transformers or using engineering techniques, such as document pre-tokenisation. Nevertheless, these improvements are likely to take effect on all models, and overall, we still expect shallow models to outperform full-size models for low-latency retrieval (however, the point where they start to perform better may change).

Overall, for RQ1, we conclude that Hypothesis H1 holds, and shallow cross-encoders are more effective compared to full-size models for low-latency retrieval.

4.3 RQ2. Effect of the Training Scheme

To answer RQ2, we first analyse whether or not overconfidence is present in shallow Cross-Encoders and whether or not gBCE helps to mitigate it.

To do that, we analyse predicted probabilities at different ranks on queries from TREC-DL2019 and TREC-DL2020. Figure 4 shows the results of this analysis on a sample query from TREC-DL2019 (results for other queries looked similar). As we can see from the figure, TinyBERT-BCE (trained with BCE loss and 1 negative) predicts probabilities very close to 1 for ranks from 1 to 5. For this query, there was only one relevant (label = 3) document identified by the TREC assessors (among 138); therefore, most of these high probabilities are false positives – a clear sign of overconfidence. If such a false positive (query/ document) pair appears in the training data, according to Equation (4), the BCE loss function on the false positive sample will be computed as $\log(1 - p^+)$, which tends to $-\infty$ when p^+ tends to 1;

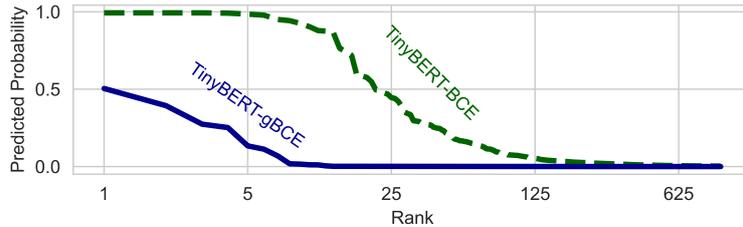


Fig. 4: Predicted probabilities at different ranks for TREC-DL2019 query 146187 “difference between a mcdouble and a double cheeseburger”.

Table 2: Effect of the loss function and the number of negatives training on Tiny BERT-based Cross-Encoder NDCG@10. Bold indicates the best result, and * indicates a statistically significant difference ($pvalue < 0.05$) compared to the baseline (BCE loss, one negative).

(a) TREC-DL2019					(b) TREC-DL2020						
Num Negatives→	1	2	8	32	128	Num Negatives→	1	2	8	32	128
Loss↓						Loss↓					
BCE	0.6386	0.6640	0.6735*	0.6622	0.6701	BCE	0.6056	0.6203	0.6340	0.6424*	0.6323
gBCE	0.6593	0.6607	0.6700*	0.6619	0.6721	gBCE	0.6193	0.6150	0.6270	0.6381	0.6306

this causes numerical instability (i.e. large gradients) during training. This is less problematic for full-size models: recall from Section 3 that the use of pre-trained checkpoints in large models works as a strong regulariser, which stabilises model training. Full-size models are also more robust because they have more other regularisations in their architecture, such as dropouts and layer normalisations in each transformer block.

In contrast, TinyBERT-gBCE (model with gBCE loss and 128 negatives) never predicts a probability higher than 0.5, showing more variation in scores in the top predicted results. This confirms that the gBCE training scheme is able to mitigate overconfidence in shallow Cross-Encoders.

We now analyse the effect of the gBCE training scheme on shallow Cross-Encoders. As we discussed in Section 3, the gBCE scheme has two key components: (i) a large number of negatives and (ii) gBCE loss. To better understand the effect of gBCE, we analyse these components independently: we train the TinyBERT model, selecting the number of negatives from (1, 128) and the loss function from (BCE, gBCE).

Table 2 summarises the results of our evaluation. As we can see from the table, compared to the standard TinyBERT-BCE model (BCE loss, one negative), both components have a positive effect. For example, on TREC-DL2019, an increased number of negatives improves the NDCG@10 metric from 0.638 to 0.670 (+4.9%), and gBCE loss improves the result to 0.6593. The combination of these improvements leads to an improvement over the baseline (0.6721, +5.2%). As we can see, with 128 negatives, changing the loss function from BCE to gBCE does not have a big effect +0.3% on TREC-DL2019 and -0.3% on TREC-DL2020. This is in line with the original gBCE paper [29], which suggests that gBCE loss is less important with many negatives. We also observe from Table 2 that the effectiveness of the model only increases up to a certain number of negatives, after which the effectiveness fluctuates. Indeed, the best effectiveness is achieved with 8 negatives on TREC-DL2019 and with 32 negatives on TREC-DL2020. At this point, switching from BCE to gBCE is already unnecessary, and all the improvements come from the increased number of negatives. Note that these numbers depend on how many candidate documents are re-ranked. However, our analysis shows that TinyBERT trained with gBCE and 128 negatives consistently outperforms TinyBERT with BCE and one negative with a different number of candidates, and therefore with different latencies; the same observation can also be seen in Figure 1 for varying sizes of K .

We also performed the analysis of the gBCE training scheme on the MS-MARCO dev queryset and found that the evaluation results are in line with the results on the TREC querysets, but due to the larger size of the queryset, we achieved statistically significant improvements compared to the baseline in most of the experiments. For example, evaluation of TinyBERT on the dev queryset showed statistically significant improvements when switching to gBCE loss with one negative sample; see Appendix A for more details.

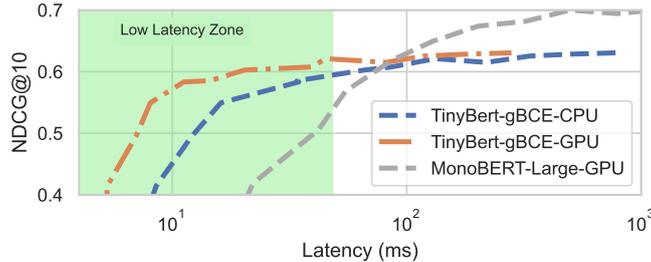


Fig. 5: Comparison of tradeoffs on CPU and GPU, TREC-DL2020 queryset.

Overall, in answer to RQ2, we summarise that the gBCE training scheme leads to improvements on both experimental TREC-DL querysets; a large number of negatives is the most important component of the scheme for shallow Cross-Encoders.

4.4 RQ3. CPU inference

Finally, we analyse the effectiveness/efficiency tradeoffs of a shallow Cross-Encoder without use of a GPU. For this experiment, we use the smallest model, TinyBERT-gBCE, because it has shown the best effectiveness in the low-latency zone in RQ1.

Figure 5 compares the tradeoffs of the model when it is used with CPU inference and with GPU inference. The plot also shows the tradeoff of a full-scale MonoBERT-Large model with GPU inference. The figure shows that GPU inference is better than CPU inference, especially within the low latency zone. For example, with a 10ms latency window, the model with CPU inference only achieves NDCG@10 of 0.447, whereas the model with GPU inference achieves NDCG@10 of 0.573 (+28%). However, with a larger allowed latency, the difference decreases. At 50ms (the upper bound of our “low latency zone”), the difference in effectiveness between GPU and CPU is just 3%. As we can see from the plot, in both cases, TinyBERT provides a better tradeoff than MonoBERT with GPU inference, and all 3 models intersect at approximately 100ms latency.

The fact that TinyBERT allows us to achieve relatively high performance even using CPU-only inference allows us to use it in cases where GPU inference is not feasible, such as for on-device search in applications. The ability to efficiently perform inference on a CPU can also amount to cost savings. For instance, the AWS’s GPU-equipped `g5.4xlarge` instance costs \$1.624/hr at the time of writing, while a roughly equivalent instance without a GPU (`m6gd.4xlarge`) costs less than half as much, at \$0.7232/hr.

Overall, in answer to RQ3, we conclude that while CPU inference is less efficient for shallow Cross-Encoders compared to GPU inference, it allows us to achieve relatively high effectiveness. Considering its low memory footprint (checkpoint is only 17 Mb), we argue that a TinyBERT-based encoder may be an effective solution for systems without a GPU (such as on-device search).

5 Conclusion

In this paper, we proposed shallow Cross-Encoders as a solution for low-latency information retrieval. We showed that shallow Cross-Encoders are more effective than full-size when latency is limited (e.g. TinyBERT model achieved +51% NDCG@10 on TREC-DL2019 compared to MonoBERT-Large with latency limited by 25ms; see Figure 1). We adapted the gBCE training scheme to shallow Cross-Encoders and showed that it improves the effectiveness of shallow Cross-Encoder models (e.g. +5.2% NDCG@10 on TREC-DL2019; see Table 2). We also showed that shallow Cross-Encoders can be effective even for CPU-only inference (e.g., on TREC-DL2020, the difference in NDCG@10 is only 3% with 50ms latency; see Figure 5). We believe that shallow Cross-Encoders can be further optimised by applying engineering techniques, such as pre-tokenisation.

A Effect of gBCE training scheme on Tiny BERT-based Cross-Encoder on the MS MARCO dev set

Table 3 reports the effectiveness of a Tiny BERT model on a 6,980 queries sub-set of the MS MARCO dev set (dataset `irds:msmarco-passage/dev/small` in PyTerrier). The evaluation follows the scheme

Table 3: Effect of the loss function and the number of negatives training on Tiny BERT-based Cross-Encoder MRR@10 on the MS MARCO dev set. Bold indicates the best result, and * indicates a statistically significant difference ($pvalue < 0.05$) compared to the baseline (BCE loss, one negative).

Num Negatives→	1	2	8	32	128
Loss↓					
BCE	0.2942	0.3032*	0.3057*	0.3128*	0.3172*
gBCE	0.3035*	0.2974	0.3109*	0.3183*	0.3200*

described in Section 4.3, with the exception of using MRR@10, which is the official metric for this queryset, instead of NDCG@10. As we can see from the table, the overall trends follow the observations in Section 4.3. In particular, an increased number of negatives is more important than the loss function; gBCE loss improves results with a small number of negatives but has a moderate effect when the number of negatives increases. However, we observe that overall gBCE in this experiment is better than BCE loss in 5 out of 6 cases. With 1 negative, the improvement over BCE loss is statistically significant. Overall, the combination of gBCE loss and 128 number of negatives provides a significant improvement of MRR@10, from 0.2942 to 0.3200 (+8.76%), compared to the “standard” training scheme with 1 negative and BCE loss. Note that this result is lower compared to the larger models – e.g. Nogueira et al. [27] achieved MRR@10 of 0.36 on this queryset with a BERT-Large model. Lower effectiveness compared to the full-scale models is an expected result, as we do not control for latency in this experiment. When latency is limited, shallow Cross-Encoders are more effective (see Figure 1).

References

- Bajaj, P., Campos, D., Craswell, N., Deng, L., Gao, J., Liu, X., Majumder, R., McNamara, A., Mitra, B., Nguyen, T., Rosenberg, M., Song, X., Stoica, A., Tiwary, S., Wang, T.: MS MARCO: A Human Generated MACHine Reading COMprehension Dataset. In: Proc. NeurIPS (2018)
- Bhargava, P., Drozd, A., Rogers, A.: Generalization in NLI: Ways (Not) To Go Beyond Simple Heuristics (2021), <http://arxiv.org/abs/2110.01518>
- Bruch, S., Lucchese, C., Nardini, F.M.: Efficient and Effective Tree-based and Neural Learning to Rank. Foundations and Trends® in Information Retrieval **17**(1), 1–123 (2023)
- Craswell, N., Mitra, B., Yilmaz, E., Campos, D.: Overview of the TREC 2020 deep learning track. In: Proc. TREC (2020)
- Craswell, N., Mitra, B., Yilmaz, E., Campos, D., Voorhees, E.M.: Overview of the TREC 2019 deep learning track. In: Proc. TREC (2019)
- Dean, J., Barroso, L.A.: The tail at scale. Communications of the ACM **56**(2), 74–80 (2013)
- Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In: Proc. of NAACL-HLT. pp. 4171–4186 (2019)
- Ding, S., Suel, T.: Faster top-k document retrieval using block-max indexes. In: Proc. SIGIR. pp. 993–1002 (2011)
- Formal, T., Piwowarski, B., Clinchant, S.: SPLADE: Sparse Lexical and Expansion Model for First Stage Ranking. In: Proc. SIGIR. pp. 2288–2292 (2021)
- Goodfellow, I., Bengio, Y., Courville, A., Bach, F.: Deep Learning. MIT Press (2017)
- Hofstätter, S., Althammer, S., Schröder, M., Sertkan, M., Hanbury, A.: Improving Efficient Neural Ranking Models with Cross-Architecture Knowledge Distillation (2021), <http://arxiv.org/abs/2010.02666>
- Hofstätter, S., Lin, S.C., Yang, J.H., Lin, J., Hanbury, A.: Efficiently Teaching an Effective Dense Retriever with Balanced Topic Aware Sampling. In: Proc. SIGIR. pp. 113–122 (2021)
- Hofstätter, S., Zlabinger, M., Hanbury, A.: Interpretable & Time-Budget-Constrained Contextualization for Re-Ranking. In: Proc. ECAI (2020)
- Humeau, S., Shuster, K., Lachaux, M.A., Weston, J.: Poly-encoders: Transformer Architectures and Pre-training Strategies for Fast and Accurate Multi-sentence Scoring (2020), <http://arxiv.org/abs/1905.01969>
- Khattab, O., Zaharia, M.: ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. In: Proc. SIGIR. pp. 39–48 (2020)
- Kohavi, R., Tang, D., Xu, Y.: Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing. Cambridge University Press (2020)
- Kulkarni, H., MacAvaney, S., Goharian, N., Frieder, O.: Lexically-Accelerated Dense Retrieval. In: Proc. SIGIR. pp. 152–162 (2023)
- Lin, S.C., Yang, J.H., Lin, J.: In-Batch Negatives for Knowledge Distillation with Tightly-Coupled Teachers for Dense Retrieval. In: Proc. Repl4NLP. pp. 163–173 (2021)
- Lu, W., Jiao, J., Zhang, R.: TwinBERT: Distilling Knowledge to Twin-Structured Compressed BERT Models for Large-Scale Retrieval. In: Proc. CIKM. pp. 2645–2652 (2020)
- MacAvaney, S., Macdonald, C.: A Python Interface to PISA! In: Proc. SIGIR. pp. 3339–3344 (2022)

21. MacAvaney, S., Macdonald, C., Ounis, I.: Streamlining Evaluation with ir-measures. In: Proc. ECIR. pp. 305–310 (2022)
22. MacAvaney, S., Nardini, F.M., Perego, R., Tonello, N., Goharian, N., Frieder, O.: Expansion via Prediction of Importance with Contextualization. In: Proc. SIGIR. pp. 1573–1576 (2020)
23. MacAvaney, S., Tonello, N., Macdonald, C.: Adaptive Re-Ranking with a Corpus Graph. In: Proc. SIGIR. pp. 1491–1500 (2022)
24. MacAvaney, S., Yates, A., Cohan, A., Goharian, N.: CEDR: Contextualized Embeddings for Document Ranking. In: Proc. SIGIR. pp. 1101–1104 (2019)
25. Macdonald, C., Tonello, N., MacAvaney, S., Ounis, I.: PyTerrier: Declarative Experimentation in Python from BM25 to Dense Retrieval. In: Proc. CIKM. pp. 4526–4533 (2021)
26. Mallia, A., Siedlaczek, M., Mackenzie, J.M., Suel, T.: PISA: Performant indexes and search for academia. In: Proc. OSIRRC@SIGIR 2019. vol. 2409, pp. 50–56 (2019)
27. Nogueira, R., Cho, K.: Passage Re-ranking with BERT (2020), <http://arxiv.org/abs/1901.04085>
28. Nogueira, R., Jiang, Z., Lin, J.: Document Ranking with a Pretrained Sequence-to-Sequence Model (2020), <http://arxiv.org/abs/2003.06713>
29. Petrov, A.V., Macdonald, C.: gSASRec: Reducing Overconfidence in Sequential Recommendation Trained with Negative Sampling. In: Proc. RecSys. pp. 116–128 (2023)
30. Reimers, N., Gurevych, I.: Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In: Proc. EMNLP (2019)
31. Robertson, S., Walker, S., Jones, S., Hancock-Beaulieu, M., Gatford, M.: Okapi at TREC 3. In: Proc. TREC (1994)
32. Scells, H., Zhuang, S., Zuccon, G.: Reduce, Reuse, Recycle: Green Information Retrieval Research. In: Proc. SIGIR. pp. 2825–2837 (2022)
33. Thakur, N., Reimers, N., Rücklé, A., Srivastava, A., Gurevych, I.: BEIR: A Heterogenous Benchmark for Zero-shot Evaluation of Information Retrieval Models (2021), <http://arxiv.org/abs/2104.08663>
34. Turc, I., Chang, M.W., Lee, K., Toutanova, K.: Well-Read Students Learn Better: On the Importance of Pre-training Compact Models (2019), <http://arxiv.org/abs/1908.08962>
35. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is All you Need. In: Proc. NeurIPS (2017)
36. Wallat, J., Berlinger, F., Anand, A., Anand, A.: Probing BERT for Ranking Abilities. In: Proc. ECIR. pp. 255–273 (2023)
37. Wang, X., MacAvaney, S., Macdonald, C., Ounis, I.: An Inspection of the Reproducibility and Replicability of TCT-ColBERT. In: Proc. SIGIR. pp. 2790–2800 (2022)
38. Wang, X., Macdonald, C., Tonello, N., Ounis, I.: Reproducibility, Replicability, and Insights into Dense Multi-Representation Retrieval Models: From ColBERT to Col*. In: Proc. SIGIR. pp. 2552–2561 (2023)
39. Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Scao, T.L., Gugger, S., Drame, M., Lhoest, Q., Rush, A.M.: HuggingFace’s Transformers: State-of-the-art Natural Language Processing (2020), <http://arxiv.org/abs/1910.03771>
40. Xiong, L., Xiong, C., Li, Y., Tang, K.F., Liu, J., Bennett, P., Ahmed, J., Overwijk, A.: Approximate Nearest Neighbor Negative Contrastive Learning for Dense Text Retrieval (2020)
41. Zhuang, H., Qin, Z., Jagerman, R., Hui, K., Ma, J., Lu, J., Ni, J., Wang, X., Bendersky, M.: RankT5: Fine-Tuning T5 for Text Ranking with Ranking Losses. In: Proc. SIGIR. pp. 2308–2313 (2023)