

CROSS-JEM: Accurate and Efficient Cross-encoders for Short-text Ranking Tasks

Bhawna Paliwal*[†]
bhawna@microsoft.com
Microsoft Research
Bengaluru, India

Deepak Saini*[†]
desaini@microsoft.com
Microsoft
Redmond, USA

Mudit Dhawan[‡]
muditd2000@gmail.com
Carnegie Mellon University
Pittsburgh, USA

Siddarth Asokan[†]
Siddarth.Asokan@microsoft.com
Microsoft Research
Bengaluru, India

Nagarajan Natarajan
Nagarajan.Natarajan@microsoft.com
Microsoft Research
Bengaluru, India

Surbhi Aggarwal
Surbhi.Aggarwal@microsoft.com
Microsoft
Bengaluru, India

Pankaj Malhotra
pamalhotra@microsoft.com
Microsoft
Bengaluru, India

Jian Jiao
Jian.Jiao@microsoft.com
Microsoft
Redmond, USA

Manik Varma
manik@microsoft.com
Microsoft Research
Bengaluru, India

Abstract

Ranking a set of items based on their relevance to a given query is a core problem in search and recommendation. Transformer-based ranking models are the state-of-the-art approaches for such tasks, but they score each query-item independently, ignoring the joint context of other relevant items. This leads to sub-optimal ranking accuracy and high computational costs. In response, we propose Cross-encoders with Joint Efficient Modeling (CROSS-JEM), a novel ranking approach that enables transformer-based models to jointly score multiple items for a query, maximizing parameter utilization. CROSS-JEM leverages (a) redundancies and token overlaps to jointly score multiple items, that are typically short-text phrases arising in search and recommendations, and (b) a novel training objective that models ranking probabilities. CROSS-JEM achieves state-of-the-art accuracy and over 4x lower ranking latency over standard cross-encoders. Our contributions are threefold: (i) we highlight the gap between the ranking application’s need for scoring thousands of items per query and the limited capabilities of current cross-encoders; (ii) we introduce CROSS-JEM for joint efficient scoring of multiple items per query; and (iii) we demonstrate state-of-the-art accuracy on standard public datasets and a proprietary dataset. CROSS-JEM opens up new directions for designing tailored early-attention-based ranking models that incorporate strict production constraints such as item multiplicity and latency.

Keywords

sponsored search, information retrieval, keyword scoring, keyword ranking, cross encoders, efficiency, production systems, online systems, industrial applications, large-scale learning

*Equal contribution

[†]Corresponding authors

[‡]Work done while at Microsoft Research

Reference Format:

Bhawna Paliwal, Deepak Saini, Mudit Dhawan, Siddarth Asokan, Nagarajan Natarajan, Surbhi Aggarwal, Pankaj Malhotra, Jian Jiao, and Manik Varma. . CROSS-JEM: Accurate and Efficient Cross-encoders for Short-text Ranking Tasks. In *arXiv*. 2024, 15 pages.

1 Introduction

We consider the problem of ranking that arises in search and recommendation pipelines, wherein the goal is to rank a set of items based on their relevance to a given query. Our work is in the context of two-stage *retrieve-then-rank* pipelines in modern search engines [8, 15, 16, 25, 43, 45] as depicted in Figure 1. Given a *query*, i.e., a search phrase such as “*patagonia japan*”, the retrieval stage pares the *items*, i.e., keywords such as “*clothing store japan*”, “*patagonia homes shimoda*” that are bid against for displaying ads, from billions to a few hundreds [11, 23] of potential interest. These items are subsequently provided as the input to the ranking stage. In this work, we focus on the ranking model, given a black-box retriever. We consider short-text items (as in the above example), which appear in a myriad of recommendation systems applications such as product recommendation, query to advertiser bid phrase recommendation, and Wikipedia category tagging [5, 35, 44]. In designing the ranking model, two key axes are the model architecture, and the choice of the loss function, while the key performance metrics for such systems are accuracy and inference latency.

Ranking architectures and limitations: Along the architecture vertical, *encoders* that employ stacked attention layers to encode a query-item pair, followed by a *classifier* to predict ranking scores are widely adopted for ranking [25, 27, 45]. Recently, sequence-to-sequence models with encoder-decoder and decoder-only architectures have also been proposed for ranking. These models provide a ranking score based on particular vocabulary-token logits [26, 42, 46]. However, these approaches model ranking as a *pointwise* task, providing ranking scores for a given query and item pair. But ranking is inherently a list-based task that requires scoring query-item pairs relative to one another, and not in isolation. In

particular, the other items in the list to be ranked for a given query provide crucial context for scoring. Pointwise models neglect the list context, produce independent scores that may not reflect the optimal ranking order, and are difficult to calibrate across items for sorting to provide final rankings [32].

Table 1 illustrates this by juxtaposing the top-5 ranked items obtained using our proposed approach (listwise modeling) and a baseline pointwise ranking model [25]. We observe that relatively more *generic* items such as “*mexican cuisine*”, although relevant to the query “*different foods of oaxaca mexico*”, are ranked higher in baseline predictions, owing to (a) their frequency in training data, (b) token-level matching and other biases which are difficult to mitigate in pointwise modeling. On the other hand, our proposed listwise approach evaluates all the items to be ranked holistically, and subsequently ranks more specific (not just relevant) items higher.

Furthermore, pointwise transformer based models [25, 46] are computationally expensive and impractical for real-time ranking systems that need to handle large-scale traffic requiring low latency and high throughput. Therefore, many industrial systems resort to using simpler sparse neural networks [3] or late-interaction models [14, 18] for online ranking, sacrificing accuracy for latency.

Another line of research has been along the loss function, wherein training with listwise loss functions [2, 9, 46] have shown accuracy gains. These listwise losses optimize the models for a list of items to be ranked, rather than for individual query-item pairs, and can enhance the ranking performance without increasing the model size or complexity. However, these model architectures still operate at the query-item (pointwise) level, and produce independent scores for each item in the list, without explicitly modeling the inter-item dependencies or the query context. Recent works also incorporate listwise modeling via pre-trained LLMs for ranking [29, 31, 38, 42]. These approaches typically work with pre-trained large-scale models, focus on specifying all ranking items in the prompt and employing prompt engineering methods to improve accuracy and minimize LLM inference calls. Due to the huge parameter count (running into a few billions), these models cannot be deployed in large-scale online ranking systems. **In this paper, we aim to bridge this gap by proposing a ranking model that works at the list level, explicitly models inter-item interactions, and achieves a superior latency-accuracy tradeoff, making it deployable in real-time scenarios.**

The Proposed Approach: We propose an end-to-end joint ranking approach that models the listwise ranking of the query and given items, capturing both the query-item and the item-item interactions and satisfying the strict latency requirements of industry-scale recommendation systems. Our approach, entitled **CROSS-JEM** (**CROSS**-Encoder with **Joint Efficient Modeling**), leverages the list structure of the input in both the encoder and the classifier components, as well as in the training objective. Specifically, our encoder is a transformer-based architecture that processes a query and a list of items to rank in a single pass, generating ‘*list-aware*’ context vectors for each input token. To achieve this, each item token attends to all other items in the list, as well as to the query, allowing the encoder to capture the joint relevance of all items. This is followed by a classifier which is jointly trained with the encoder. While any standard loss function (e.g., binary cross entropy) could be used

to train a CROSS-JEM model, in this work, we propose a novel variant of the listwise loss functions [2, 30], Ranking Probability Loss (RPL); which can be interpreted as a divergence between the predicted probabilities and the estimated ranking probabilities (as a function of model logits and ground truth rankings) of items to be ranked. The proposed listwise loss works much better in conjunction with our joint modeling than standard pointwise or listwise losses [2, 4, 30]. To the best of our knowledge, we are the first to propose a joint ranking approach that can effectively model listwise ranking in both the model architecture and the training objective with real-time latency constraints.

CROSS-JEM is able to support low latency applications (few ms) by significantly reducing the computational costs of cross-attention across query and the list of ranking items. This is achieved by exploiting the presence of token duplicates in the retrieved set of items for a given query. CROSS-JEM exploits this redundancy to sidestep processing long sequences, thereby keeping the inference latency small. Further, jointly obtaining the ranking scores for a list of items avoids multiple calls to the expensive encoder (unlike pointwise approaches) making CROSS-JEM significantly faster.

In summary, our key contribution is the introduction of a novel **joint ranking approach CROSS-JEM**, that scores multiple items per query in a single pass, exploiting token interaction across items for better accuracy and token redundancies for better efficiency. We demonstrate the effectiveness of CROSS-JEM on two public benchmark datasets for short-text re-ranking, wherein it outperforms the best-performing baselines by at-least 3% in terms of MRR. When applied to large-scale search-based recommendation, CROSS-JEM demonstrated a 13% higher accuracy than state-of-the-art models, while being over 6× faster than standard cross-encoders [25]. We also deploy CROSS-JEM for real-time ranking on live traffic, where it reduces the quick-back-rate by 1.8%, indicating improved relevance of ads to users. Our work presents a general, scalable framework for joint ranking of multiple items across various domains and short-text tasks, accounting for ranking under task-specific constraints such as item multiplicity and latency.

2 Background and Related Work

Ranking Architecture: State-of-the-art transformer based models with stacked attention layers are typically encoder based, such as monoBERT and Birch [1, 25]. These models encode query-passage pairs with a bi-directional attention encoder and use a classifier to obtain a ranking score. Alternatively, sequence-to-sequence models, such as monoT5 [26] leverage the pre-training knowledge of generative models for ranking. These models generate a ranking score from a specific vocabulary token in the decoder output. Another line of work explores decoder-only models, such as llama and GPT-based models and rely on their extensive pre-training knowledge and parameter count for ranking [19, 42]. Despite their advantages, decoder-only models fine-tuned on large ranking corpora [42] do not outperform the fine-tuned encoder-based and sequence-to-sequence models on short-text ranking tasks (cf. Section 5). CROSS-JEM is the first joint-ranking approach that effectively incorporates listwise ranking into the model architecture and training objective while maintaining latency constraints.

Joint (Listwise) Ranking: Ranking inherently involves comparing

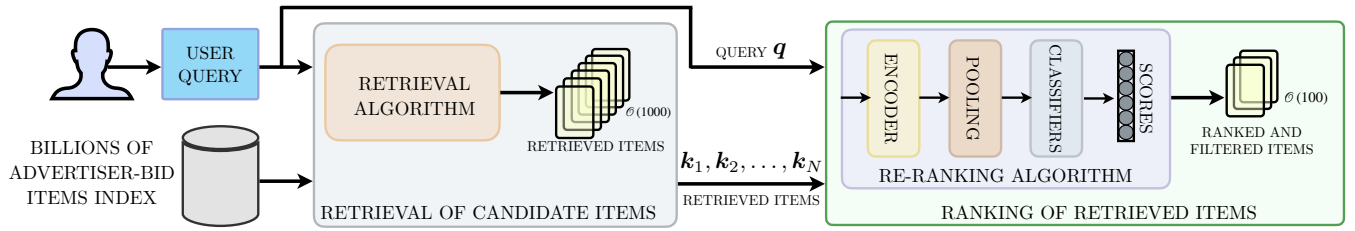


Figure 1: The two-stage large-scale search and recommendation pipeline comprises: (i) candidate selection from billions of items; (ii) re-ranking of retrieved items. CROSS-JEM works at stage (ii), to output the score for all N item in a single pass.

Table 1: A comparison of CROSS-JEM’s listwise modeling and a pointwise ranking model [25]: Relatively more *generic* (but relevant) items are ranked higher in baseline predictions, owing to biases such as high frequency in the training data. Our listwise model evaluates all the shortlisted items in a single forward pass, and ranks more specific (and relevant) items higher.

Query	Top-5 ranked items in the Proposed Approach (CROSS-JEM)	Top-5 ranked items in the baseline Cross Encoder
different foods of oaxaca mexico	'the foods of oaxaca ', 'oaxacan cuisine', 'exploring oaxacan food ', 'authentic recipes from oaxaca, mexico ', '6 things you'll love about oaxaca '	'culinary tales: the kinds of food mexicans eat', 'mexican christmas foods', 'mexican cuisine', 'culture: food and eating customes in mexico', 'popular food in mexico',
what is the bovine growth hormone	' recombinant bovine growth hormone ', 'what is rbgh? ', 'rbgh', 'bovine growth hormone and milk : what you need to know', 'what is rbst? '	'growth hormone', ' human growth hormone ', ' alternative names for growth hormone ', ' human growth hormone and insulin are friends ', ' growth hormone (somatotropin) '

a list of items; thus, recent works employ listwise loss functions in encoder or encoder-decoder models to enhance ranking accuracy. Gao et al. [10] devised a listwise multi-class cross-entropy loss to optimize ranking probabilities in encoder models. Zhuang et al. [46] introduced a ranking-specific listwise cross-entropy loss to improve performance of Seq2Seq models for ranking. They also demonstrated that an expanded form of cross-entropy loss (poly-1) achieved superior performance across various ranking metrics. Although listwise losses improved ranking accuracy, the model architectures of these rankers remain pointwise and do not fully capture the ranking task. CROSS-JEM ranks multiple items per query in one pass, capturing item-item interactions in the ranking list to improve accuracy. It is optimized with Ranking Probability Loss (RPL), a novel variant of the listwise ranking loss (as considered in ListNet [2]) aligned with the CROSS-JEM architecture, and estimates the ranking probabilities using target relevance labels and predicted model logits.

LLMs for Ranking: LLMs have emerged as a powerful tool for ranking tasks; they can leverage pre-trained knowledge, large parameter count, and effective prompting to achieve superior performance. Existing works have adopted two main approaches to exploit LLMs for ranking: (a) Using LLMs directly as re-rankers by designing novel prompting schemes [21, 29, 38] and sorting strategies [31, 47]; the high-level idea is to encode the ranking items and the query into a single input and generate a ranked list as output, using a sliding window technique to handle long inputs; (b) Using LLMs for learning more accurate smaller ranking models [20, 28, 38, 42] via standard distillation or via augmenting the training data with synthetic samples generated by the LLM. These works have demonstrated that LLMs can outperform small-scale supervised methods [25, 26, 31] on various ranking benchmarks. Deploying LLMs at scale for real-time serving scenarios is still challenging due to their high computational costs and memory requirements. But using LLMs to enhance smaller ranking models,

as in (b), is more practical and scalable – this can be applied to our proposed CROSS-JEM as well.

Efficiency and Scaling: Another closely related area is the focus on improving the efficiency and scaling of transformer based rankers (and retrievers) using light-weight architectures. Approaches employing early-attention (such as monoBERT [25]) have shown high ranking accuracies, but cannot support the low-latency requirements of online raking applications in industry-scale recommendation systems. This can be attributed to the requirement of making multiple calls to the expensive transformer-based encoder to rank a list of items per query, which is infeasible in a few milliseconds. Therefore, online production systems use a variation of sparse neural networks, namely MEB [3], for ranking thousands of items in real-time. Late-interaction models such as ColBERT [14], Baleen [13] and TwinBERT [18] are also used for online ranking due to their computational efficiency. These models reduce the computational costs by applying a late-interaction layer over query-item embeddings. This comes at the price of accuracy, owing to the lack of interactions between query and item tokens. They also incur high storage overheads in online settings, as they need to store and retrieve token-level embeddings. Another line of work that focuses on efficient retrieval and ranking is based on dual-encoder architectures, such as ANCE [41], DPR [12], and INSTRUCTOR [36]. These methods use contrastive-style training to learn a query and item encoder and metrics such as cosine similarity to rank query-item pairs. They can be made highly efficient via a nearest neighbor search [22, 37] over pre-computed embeddings, but lose out on accuracy compared to cross-encoder based approaches [25, 34].

3 CROSS-encoder with Joint Efficient Modeling

We now present CROSS-JEM, our proposed approach to accurate and low-latency ranking via efficient scoring of multiple items. Before we describe CROSS-JEM in detail, we develop the notation used in the subsequent sections of this manuscript.

Notation: We denote queries by q and the corresponding set of N candidate items retrieved for q by $\mathbb{K}_q = \{k_1, k_2, \dots, k_N\}$. We denote the dataset of queries and items used for training by \mathbb{Q}_{tr} and \mathbb{I}_{tr} respectively, and that of the test datasets by \mathbb{Q}_{te} and \mathbb{I}_{te} . We drop the subscripts when the meaning is clear from the context. The ground-truth scores are given by $y_i \in \mathbb{R}^N$, $[y_i]_j = y_{ij}$ denoting the score for the item k_j from set \mathbb{K}_q , associated with the query q_i . The queries q and items k_j are tokenized via $\mathcal{T}(\cdot)$ to obtain d -dimensional representations (tokens), given by $\mathbb{T}_q = \{q^1, q^2, \dots, q^{L_q}\}$, and $\mathbb{T}_{k_j} = \{k_j^1, k_j^2, \dots, k_j^{L_{k_j}}\}$, respectively, where $q^l, k_j^l \in \mathbb{R}^d$.

Research Problem: We seek to learn a ranking model that, given a query q , assigns scores $\mathbb{S}_q = \{s_1, s_2, \dots, s_N\}$ for all items in \mathbb{K}_q , such that the ranking induced by the scores is accurate.

3.1 The CROSS-JEM Architecture

CROSS-JEM comprises an encoder to obtain representations of a given query q and all its candidate items \mathbb{K}_q . The representations are pooled and passed to a classification head, that outputs a score corresponding to each item $k_j \in \mathbb{K}_q$ associated with the query.

CROSS-JEM’s architecture is primarily inspired by the observation that the candidate items in the set \mathbb{K}_q , for a given query q , has significant token overlap amongst themselves. A more in-depth exploration of this phenomenon, in the context of efficient-scoring, is provided in Section 5.3. Given a query q and its candidate items \mathbb{K}_q , the core idea is we form the union of tokens \mathbb{T}_{U_q} of all items in \mathbb{K}_q . Subsequently, the representations of query and set \mathbb{T}_{U_q} can be obtained in a *single pass* of the encoder, and scored via a *single pass* over the classifiers.

While at first glance, it might appear that using \mathbb{T}_{U_q} (with a potential loss of ordering of the item tokens) could adversely affect performance, we observed in preliminary experimentation that this is not the case when scoring short-text items. In particular, we compared the performance of two cross-encoder models on search engine logs, one with items as-is, and another comprising items with alphabetically sorted tokens. Both the mean average precision (MAP) and accuracy of the latter model was found to be within 1% of the score obtained when the sequence information is retained. Additional discussions are provided in Appendix F.

We now describe the CROSS-JEM encoder and classifier in detail.

Encoder: CROSS-JEM employs a trainable encoder \mathcal{E}_θ , which takes as input a sequence of tokens $\mathbb{T} = [t^1, t^2, \dots, t^n]$ (of length n), and generates as output another sequence of d -dimensional contextual embeddings $\mathcal{E}_\theta(\mathbb{T}) = \mathbf{E} = [e_1, e_2, \dots, e_n]$. These embeddings provide context-dependent representations of the input tokens, and can be used for downstream tasks such as classification, generation, and retrieval. Given the tokenization of the query (\mathbb{T}_q), and that of an item (\mathbb{T}_{k_j}), the contextual embeddings of the tokens

$$[t^{[\text{CLS}]}, q^1, q^2, \dots, q^{L_q}, t^{[\text{SEP}]}, k_j^1, k_j^2, \dots, k_j^{L_{k_j}}]$$

in baseline variants is given by

$$[e^{[\text{CLS}]}, e^{q^1}, e^{q^2}, \dots, e^{q^{L_q}}, e^{[\text{SEP}]}, e^{k_j^1}, e^{k_j^2}, \dots, e^{k_j^{L_{k_j}}}]$$

where $e^{[\text{CLS}]}$ and $e^{[\text{SEP}]}$ denote the embeddings of the $t^{[\text{CLS}]}$ and $t^{[\text{SEP}]}$ tokens, respectively. These contextual embeddings are pooled

to obtain a single d -dimensional embedding for each pair (q, k_j) by means of a sum or mean pooling layer, or by taking $e^{[\text{CLS}]}$. This process is computationally expensive due to the need for N forward passes of the encoder to compute the scores for each item in \mathbb{K}_q .

In CROSS-JEM, we leverage the short-text nature of the items, item-item interactions, and redundancy of tokens amongst items in \mathbb{K}_q . This is done by computing the contextual embeddings for all *distinct tokens in the retrieved item set* \mathbb{K}_q in one pass over the sequence of the query tokens \mathbb{T}_q combined with *item token union set* $\mathbb{T}_{U_q} = \{t^{u^1}, t^{u^2}, \dots, t^{u^M}\}$, where M is the total number of tokens in the union set. The contextual embeddings of all input tokens in CROSS-JEM are computed as

$$\begin{aligned} \mathbf{E} &= \mathcal{E}_\theta \left([q^1, \dots, q^{L_q}, t^{[\text{SEP}]}, t^{u^1}, \dots, t^{u^M}] \right) \\ &= [e^{[\text{CLS}]}, e^{q^1}, \dots, e^{q^{L_q}}, e^{[\text{SEP}]}, e^{t^{u^1}}, \dots, e^{t^{u^M}}]. \end{aligned}$$

Since the number of tokens in the item union set is significantly smaller than the sum of tokens of all items in \mathbb{K}_q (cf. Section 5.3), the proposed token-union-based inference enables highly efficient computation of contextual embeddings. Figure 2 (a) illustrates the difference between the CROSS-JEM encoder, and standard encoders such as monoBERT.

Selective Pooling Layer: Given the contextual embeddings \mathbf{E} for all query tokens (\mathbb{T}_q) and union over keyword tokens (\mathbb{T}_{U_q}), CROSS-JEM employs a selective pooling layer to *jointly model* the per-item relevance score. Given $k_j \in \mathbb{K}_q$, a pooled representation for pair (q, k_j) is computed as the mean of the contextual embeddings for all tokens in \mathbb{T}_q and those tokens in \mathbb{T}_{U_q} which are present in \mathbb{T}_{k_j} . The set of selected tokens for the pooling layer is given by $\mathbb{P}_{qk_j} = \mathbb{T}_q \cup \{t^{[\text{SEP}]}\} \cup \{\mathbb{T}_{U_q} \cap \mathbb{T}_{k_j}\}$. The selectively pooled representation $e^{qk_j} \in \mathbb{R}^d$ is obtained via selective mean pooling:

$$e^{qk_j} = \frac{1}{|\mathbb{P}_{qk_j}|} \sum_{t^j \in \mathbb{P}_{qk_j}} e^{t^j}. \quad (1)$$

Ablations on designing the pooling layer are provided in Section 5.

Classifier: The final stage in CROSS-JEM is a d -dimensional shared linear classifier $\mathbf{w} \in \mathbb{R}^d$ which computes the relevance score associated with each pair (q, k_j) . The selectively pooled representations e^{qk_j} obtained for all $k_j \in \mathbb{K}_q$ are batched together ($e^{qk} \in \mathbb{R}^{N \times d}$) allowing for the computation of all logits $[f_q]_j = \langle \mathbf{w}, e^{qk_j} \rangle$ in a single shot. The scores are defined over these logits (cf. Section 4.1).

4 The CROSS-JEM Algorithm

For **training** CROSS-JEM, we jointly learn the encoder model and classifier parameters $\{\theta, \mathbf{w}\}$ with target scores obtained from a teacher model. The teacher is a large cross-encoder-based model which is highly accurate but computationally expensive (cf. Section 5). The encoder followed by the selective pooling layer output representations associated with the query and each candidate item, while the linear classifier generates logits for each query-item pair. During **inference**, CROSS-JEM predicts the logits for a set of retrieved shortlist items jointly (for a given query), as in the training phase. The SoftMax over the logits give the scores, and in turn, the ranking. The inference algorithm is provided in Appendix B.

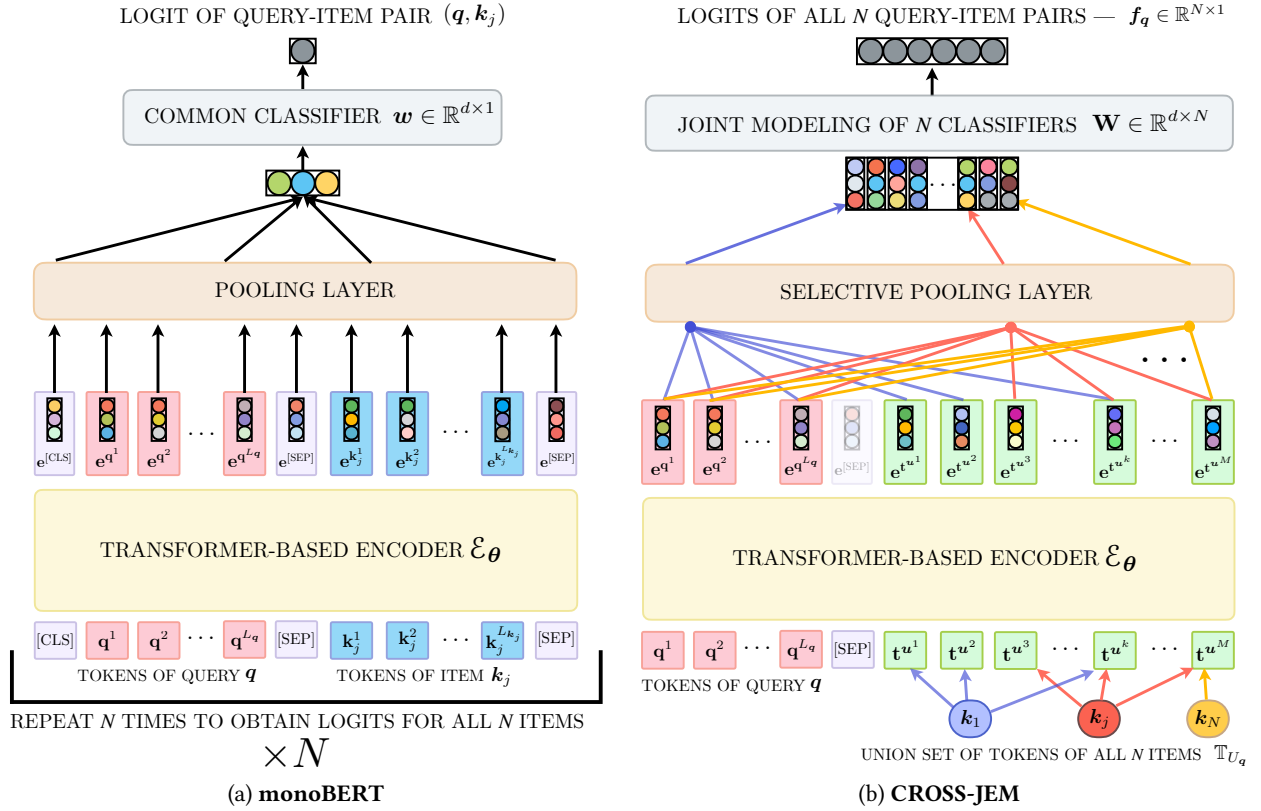


Figure 2: (Best viewed in color) Computing the relevance scores for query q and retrieved set of N items $\{k_1, k_2, \dots, k_N\}$ using (a) monoBERT and (b) CROSS-JEM: (a) monoBERT (standard cross-encoders) infers the scores for each (q, k_j) pair with query and item tokens individually (pointwise). The inference is repeated N times for N items. (b) CROSS-JEM jointly inputs query tokens and the set of union of tokens from all items to be ranked (listwise). The logits for an item k_j can be computed by selecting the contextual embeddings corresponding to its tokens. The selected embeddings are pooled, and input to the linear classifier jointly (i.e., $W = [w, w, \dots, w] \in \mathbb{R}^{d \times N}$), thus obtaining the N scores in a single encoder and classifier pass.

The Ranking Probability Loss (RPL) used to train CROSS-JEM is a novel variant of existing list-based loss functions, and is designed to take advantage of CROSS-JEM architecture, wherein all item scores are available jointly from a single forward pass.

4.1 The CROSS-JEM Ranking Probability Loss

The standard loss to train cross-encoders is the binary cross-entropy loss, defined over (q_i, k_j) , given by $\mathcal{L}^{BCE}(\theta, w)$ is

$$-\sum_{i=1}^{|\mathcal{Q}_{tr}|} \sum_{j=1}^N \left[y_{ij} \log(f_{w, q_i, j}) + (1 - y_{ij}) \log(1 - f_{w, q_i, j}) \right], \quad (2)$$

where $f_{w, q_i, j} = \langle w, \mathcal{E}_\theta(q_i, k_j) \rangle$ is the score of item j , associated with query i , computed by means of an inner product with the classifier w . However, such cross-entropy-based pointwise losses fail to account for the list of items available for ranking. *List-based loss functions* [2], in contrast, leverage the task-specific ranking information, help learn a scoring function for a list of items to be ranked, rather than for individual query-item pairs. As an example,

consider the ListNet [2] loss $\mathcal{L}^{LN}(\theta, w)$, given by:

$$\mathcal{L}^{LN}(\theta, w) = - \sum_{i=1}^{|\mathcal{Q}_{tr}|} \sum_{j=1}^N P_{y, j} \log(P_{f, j}), \quad (3)$$

where $P_{x, j} = \frac{\Phi([x]_j)}{\sum_{\ell=1}^N \Phi([x]_\ell)}$, and x is set either to the targets $[y_i]$, or the output logits $[f_{q_i}]$, and Φ is a normalizing function, typically the exponential operation, leading to P being a SoftMax function. However, this formulation is still centered around obtaining the pointwise logits, and subsequently computing the top-one probability $P_{x, j}$ using a normalization term that accounts for all pairs.

In CROSS-JEM, we design a novel version of the ListNet loss, one that factors in the availability of all logits $[f_{q_i}]$ computed by taking into account the item-item interactions. Given f_{q_i} , the ground-truth scores y_i , and a candidate item k_j , we define the set $\mathbb{L}_j = \{k \in \{1, N\} : [y_i]_k < [y_i]_j\}$, i.e., \mathbb{L}_j comprises the indices k for which the ground truth score at location k is **lower** than $[y_i]_j$, the ground truth score at location j . We now define the RPL as:

$$\mathcal{L}^{RPL} = - \sum_{i=1}^{|\mathcal{Q}_{tr}|} \sum_{j=1}^N \left(\sum_{k \in \mathbb{L}_j} [y_i]_k \right) \log \left(\text{SoftMax} \left(\sum_{k \in \mathbb{L}_j} [f_{q_i}]_k \right) \right). \quad (4)$$

The loss \mathcal{L}^{RPL} represents a cross-entropy loss over functions of the target \mathbf{y}_i and scores computed as a function of the logits \mathbf{f}_{q_j} . The following Lemma sheds light on the relationship between RPL and the ranking probability distribution.

LEMMA 1. (Ranking Probability Loss) *Assume without loss of generality that $\mathbf{f}_{q_i} \in [0, 1]^N$. Let $\mathbf{P} \in \mathbb{R}^{N \times N}$ denote a matrix with entries p_{jk} given by $p_{jk} = \text{Prob}(\text{ranking item } \mathbf{k}_j \text{ at location } k) \triangleq C \sum_{\ell \in \mathbb{L}_k} [\mathbf{f}_{q_i}]_{\ell}$, where C is a normalizing constant. Then, the Ranking Probability Loss maximizes the probability of ranking items \mathbb{K}_{q_i} in the ordering of the ground-truth \mathbf{y}_i .*

The detailed proof is given in Appendix A. The sketch of the proof follows by evaluating the matrix \mathbf{P} for the predicted logits \mathbf{f}_{q_i} defined above, and \mathbf{P}^* defined over the ground-truth scores \mathbf{y}_{ij} . Minimizing the distance between \mathbf{P} and \mathbf{P}^* is equivalent to minimizing the KL divergence between the predicted, and ground truth ranking distributions, which yields the Ranking Probability Loss defined in Equation (4).

The following Corollary presents an equivalence between the ListNet loss [2] and the proposed Ranking Probability Loss.

COROLLARY 2. (RPL and the ListNet loss) *Minimizing the Ranking Probability Loss is equivalent to optimizing for the ListNet top-1 probability loss (Equation (3) [2]) defined over modified scores $\sum_{k \in \mathbb{L}_j} [\mathbf{f}_{q_i}]_k$ and modified ground-truth scores $\tilde{\mathbf{y}}_{ij} = \sum_{k \in \mathbb{L}_j} [\mathbf{y}_i]_k$.*

Intuitively, defining the modified scores in terms of the sum of all logits in \mathbb{L}_j , the set of indices of ground-truth scores lower than the logit at j , ensures that the loss takes into account the interactions between the contextual embeddings contribution to the different logits. This is unique to the CROSS-JEM setting, and Corollary 2 shows that all guarantees derived for ListNet loss also hold for RPL.

We show in Section 5 that CROSS-JEM trained with RPL yields significantly better ranking accuracy than existing pointwise and listwise loss functions, and leads to state-of-the-art performance on public and proprietary ranking benchmarks.

5 Experimental Validation

Datasets: We evaluate CROSS-JEM on the publicly available on **Stack Overflow Duplicate Questions** [17] (**SODQ**) and a short-text version of **MS MARCO** [6] datasets. While the SODQ dataset is used as is, for MS MARCO, a ‘short-text’ variant of the query-webpage click dataset is constructed by exclusively considering webpage titles (subsequently referred to as **MS MARCO-Titles**). This narrows the dataset’s length distribution, aligning it more closely with the typical item lengths seen in sponsored search. We note that the standard metrics for MS MARCO passage ranking do not apply to MS MARCO-Titles, as they rely on the passage content as well as the title for ranking. Therefore, we report updated numbers for MS MARCO-Titles in Table 2, as appropriate. See Appendix C for more details on the datasets.

Evaluation Metrics: To validate the efficacy of CROSS-JEM for ranking tasks, we consider Mean Average Precision (MAP) and Mean Reciprocal Rank (MRR) as the evaluation metrics on both MS MARCO and SODQ datasets (see Appendix D for metric definitions). Note that MAP is a more comprehensive version of MRR, which is specifically designed for scenarios where the test set contains only one positive item per test query. Hence, on MS MARCO, where

there is only one relevant item per query, MAP@K and MRR@K are equivalent for any K.

Baselines: We compare CROSS-JEM with transformer based ranking approaches as well as sparse BoW methods such as **BM25**. Within transformer based ranking baselines, we consider state-of-the-art encoder based methods such as **monoBERT** [25]. These approaches employ *early interaction* across query-item tokens (pointwise) in the form of self-attention and are widely used for ranking. We also compare CROSS-JEM against more efficient encoders such as **ColBERT** [14]. *Late-interaction* Models like ColBERT apply self-attention within query or item separately and combine the contextual embeddings across query and item later via a computationally cheap cross-interaction layer. We compare CROSS-JEM with a ColBERT model consisting of 12 layers (109M parameters). Dual Encoders such as **ANCE** [41], **DPR** [12], and **INSTRUCTOR** [36] use contrastive learning to train the encoder and obtain dense representations of queries and items. Typically, a dot product between the query and item representations is used as a measure of relevance. INSTRUCTOR is a dual encoder approach, which obtains task-specific embeddings using instruction prompts to the model. Following Su et al. [36], we use pre-trained and instruction-tuned INSTRUCTOR model for zero-shot evaluation on MS MARCO and SODQ. We also compare against state-of-the-art sequence-to-sequence class of models containing encoder-decoder architecture, **RankT5** [46]. RankT5-base is a 24 layer model trained on MS MARCO passage ranking task and is taken as it is for evaluating as a baseline here. For a fair comparison with 6 layer CROSS-JEM model, we also fine-tune our own RankT5 model with 6 layers (encoder+decoder) on both MS MARCO-Titles and SODQ datasets. For completeness, we also compare against a decoder only ranking model with Llama2-7B architecture, referred to as **RankingGPT-Llama2-7B**, available from Zhang et al. [42]. These models have been pre-trained as well as fine-tuned specifically for ranking tasks with huge amount of training data generated from larger and more accurate LLMs.

Hyper-parameters: CROSS-JEM’s tunable hyper-parameters include maximum item union length, L_u and number of items per query, N . These hyper-parameter values are dictated by the application requirements and efficiency constraints. We use $N = 10$ and $L_u = 360$ on MS MARCO and $N = 30$ and $L_u = 242$ on SODQ. Hyperparameters used for the baselines are given in Appendix E.

5.1 Accuracy Results on Public Benchmarks

Table 2 shows the results of CROSS-JEM and several competitive baselines on the SODQ and MS MARCO datasets. We observe that CROSS-JEM achieves superior performance over the encoder based methods on both datasets, demonstrating the effectiveness of its listwise ranking approach. We also compare CROSS-JEM with encoder-decoder based methods such as RankT5-6L, which encode queries and generate item rankings using a decoder. Here, we note that RankT5-6L, which is fine-tuned with a listwise ranking loss, outperforms similar sized encoder based methods (monoBERT and ColBERT), highlighting the benefits of listwise modeling for ranking. However, we find that RankT5-base, a larger encoder-decoder model (24 layers, 223M parameters) that is fine-tuned on a long-text passage ranking task [46], performs worse than smaller models that are specifically fine-tuned on short-text ranking. A similar

Table 2: Performance of CROSS-JEM and the baseline methods on the SODQ and MS MARCO-Titles ranking datasets: All baselines and our method, CROSS-JEM, use ~60-100M parameter base models and are fine-tuned on the corresponding datasets, except for the large pre-trained models (indicated with asterisk (*)), which are used as is without any further fine-tuning on the two datasets. CROSS-JEM surpasses similar-sized state-of-the-art methods fine-tuned for short-text ranking as well as large pre-trained models by at least 3%.

	Method	Parameters	SODQ		MS MARCO-Titles		
			MAP@5	MAP@10	MRR@5	MRR@10	
Sparse Models	BM25	–	32.80	39.26	23.71	24.57	
Encoder	Early-interaction	monoBERT	66M	46.79	48.04	30.89	32.47
	Late-interaction	ColBERT	109M	36.10	37.68	30.25	32.00
	Dual Encoders	DPR	66M	47.32	48.48	28.78	30.87
		ANCE	66M	48.31	49.41	28.48	30.53
INSTRUCTOR*		335M	49.47	50.81	28.84	30.55	
Encoder-Decoder	Seq2Seq	RankT5-6L	74M	49.50	50.75	30.73	32.52
		RankT5-base*	223M	45.66	49.47	27.87	29.75
Decoder	Ranking LLMs	RankingGPT-Llama2-7B*	7B	47.64	50.62	28.66	30.47
Ours	Joint Ranking	CROSS-JEM-6L	66M	52.40	53.05	33.82	35.45

observation is made for RankingGPT-llama2-7b, a large decoder only model (7B parameters) that is pre-trained and fine-tuned on large-scale ranking datasets [42], but performs similarly to much smaller models on short-text ranking. These results suggest that the Seq2Seq models are sensitive to the domain and length of the ranking tasks, and require careful fine-tuning and adaptation. To verify this, we fine-tune a RankT5-base model on the short-text ranking datasets and observe a significant improvement in performance (cf. Appendix E). We also report that CROSS-JEM, which uses a 6-layer BERT as the base encoder, has the same number of parameters as monoBERT, but is much faster and more accurate. Specifically, CROSS-JEM can perform joint (listwise) inference for ranking over 4× faster than monoBERT, which requires multiple pointwise computations (cf. Section 5.3). Moreover, CROSS-JEM performs up to 5% more accurately than dual encoder methods and up to 20% more accurately than sparse models like BM25. We perform more experiments on the comparison of pointwise and listwise loss functions in CROSS-JEM in Section 5.4.

5.2 Case Study on Sponsored Search Ads

In large-scale search and recommendation systems like sponsored search, the ranking model serves to weeding out bad retrievals and rank the prediction pool of different retrievers to select the top-k. We evaluate the effectiveness of CROSS-JEM on this real-world task of matching user queries to relevant advertiser-bid keywords. A large scale dataset consisting of 1.8B query-keyword pairs was created by mining search engine logs (detailed in Appendix C).

Accuracy comparison: As shown in Table 3 (details on baselines in Appendix F), CROSS-JEM improves over the in-production sparse neural model MEB [3] in MAP by over 13%. Further, CROSS-JEM also outperforms ANCE and TwinBERT by large margins in MAP, Precision, and Recall. We also assess CROSS-JEM’s ability to eliminate irrelevant items while preserving relevant ones. Table 3 presents the negative and overall accuracy when retaining top 80%

of positive items per query. CROSS-JEM achieves 99.45% negative accuracy, removing nearly all irrelevant items.

Efficiency Gains: We observe that CROSS-JEM takes only 9.8 ms to score 700 keywords for a query on a A100 GPU. In contrast, monoBERT takes 41.3 ms for the same task, rendering it unsuitable for online deployment. This represents a *more than 4-fold reduction in latency* compared to standard cross-encoder models. The latency gains are because CROSS-JEM scores multiple items for a query in one shot by passing their concatenated tokens through the model. On the other hand, monoBERT scores each query-item pair independently necessitating 700 passes compared to CROSS-JEM’s 7 passes. Additionally, CROSS-JEM provides 3× lower latency on GPUs than MEB on CPUs. This highlights CROSS-JEM’s ability to leverage GPU acceleration for efficiency, unlike sparse models.

CROSS-JEM achieves a high throughput of 17,200 query-keyword pairs per second. This is over 5× more than the 3,350 pairs per second for monoBERT. The massive throughput and latency gains show CROSS-JEM’s ability to meet the computational demands of large-scale industrial systems without sacrificing accuracy.

Online A/B testing: CROSS-JEM was deployed in the ranking stage of a premier search engine to conduct A/B tests on live traffic. The ranking stage receives an average of 700 keywords and up to 1400 keywords in the 99th percentile, from a suite of retrieval algorithms. The control group consisted of a proprietary combination of late interaction, dense retrieval, and sparse-neural-network algorithms. CROSS-JEM demonstrated a decrease in the quick-back-rate (users who close the ad quickly, indicating non-relevance) by over 1.8%. Furthermore, as judged by expert judges, CROSS-JEM improved the proportion of accurate predictions by 10.2%.

5.3 Interpreting CROSS-JEM’s Performance

To better understand the **efficiency gains** in CROSS-JEM, we analyze the effect of the significant token overlap amongst candidate items in the set \mathbb{K}_q for a given query q . We trained an ANCE [41]

Table 3: Comparison of CROSS-JEM with production baselines on Sponsored Search Ads Dataset for ranking advertiser-bid keywords for a user query. CROSS-JEM outperforms baseline methods (with a latency small enough to be deployed for real-time ranking) by 13% in MAP. CROSS-JEM filters >90% irrelevant predictions at a threshold which retains 80% of good predictions.

Method	MAP@100	P@50	R@50	AUC	Negative Accuracy	Overall Accuracy
ANCE	78.39	41.94	94.02	89.84	86.19	85.55
TwinBERT	83.56	43.50	95.36	92.10	90.60	88.58
MEB	84.38	42.94	94.65	91.77	84.59	84.40
CROSS-JEM	97.48	45.76	99.07	99.41	99.45	95.27

Table 4: (a) Ablation on loss function in CROSS-JEM; (b) Adding item token sequence information in CROSS-JEM via positional encodings in the pooling layer.

Method	BCE	CE	ListNet	RPL
MRR @ 10	31.46	32.03	30.27	35.45

(a)

Method	Without Positional Encodings	With Positional Encodings
MRR@10	35.45	35.68
MRR@5	33.82	33.98

(b)

Table 5: A comparison of latency between CROSS-JEM and various baselines. The mean latency for scoring 700 items per query was computed on A100 GPUs for all models. We observe that CROSS-JEM takes about 4× lower inference time than standard cross-encoders (monoBERT).

Method	ANCE	ColBERT	monoBERT	RankT5-6L	CROSS-JEM
Latency (ms)↓	4.0	4.5	41.3	41.3	9.8

dense retriever on the same train set as above. For a query $q \in \mathbb{Q}_{te}$, recall that $\mathbb{K}_q = \{k_1, k_2, \dots, k_N\}$ are the top N ($=100$ in our experiments) items retrieved using ANCE. Let $\mathbb{T}_{k_j}^L$ denote word-piece tokens in k_j with a max-length of L . Let \mathbb{T}_U denote the union of all tokens of items $k_j \in \mathbb{K}_q$. We compute the following statistics:

$$m \text{ (mean total tokens)} = \frac{1}{|\mathbb{Q}_{te}|} \sum_{q \in \mathbb{Q}_{te}} \left(\sum_{j=1}^N |\mathbb{T}_{k_j}^L| \right), \text{ and } \quad (5)$$

$$N_u \text{ (mean size of the set } \mathbb{T}_U) = \frac{1}{|\mathbb{Q}_{te}|} \sum_{q \in \mathbb{Q}_{te}} \left(\left| \bigcup_{j=1}^N \mathbb{T}_{k_j}^L \right| \right). \quad (6)$$

Intuitively, m is the the sum of item token lengths on average, while N_u is the cardinality of the union set, averaged over the queries q for which the candidate items \mathbb{K}_q were obtained. We hypothesize that, if the items k_j have significant overlap, $m \gg N_u$. Statistically, the union size N_u is found to be at least 5× smaller than m , indicating high redundancy, and correlates with observations on the *sponsored search case study* (cf. Table 8 (b) in Appendix).

To further analyze this effect, we characterize the **time complexity** of CROSS-JEM in terms of number of query tokens $L_q = |\mathbb{T}_q|$,

number of item tokens $L_k = |\mathbb{T}_k|$, number of transformer layers L , number of candidate items N , and item union compression factor C (approximated as $\frac{L_k \times N}{|\mathbb{T}_U|}$). The time complexity for scoring all N items jointly is $(L_q + L_k N / C)^2 L$. On the other hand, for the standard cross-encoder, the corresponding time complexity is $(L_q + L_k)^2 L N$. In practice, the inference time depends on factors such as the implementation of the model, the hardware employed, and optimizations used (such as quantization). For the sponsored search setting considered in Section 5, assuming $L_q \approx L_k$, letting $N = 100$, the maximum tokens in an item is 24 and maximum item union length used is 220, we have $C \approx (24 * 100) / 220 \approx 11$. Then, the standard cross-encoder time complexity is $400 L_q^2 L$, while that of CROSS-JEM is $101.8 L_q^2 L$, which is 3.9× lower. These inference time gains are also reflected in the latency comparison of CROSS-JEM against pointwise approaches such as monoBERT (cf. Table 5).

5.4 Ablations

Loss Function: We demonstrate performance comparison of pointwise and listwise loss functions with CROSS-JEM architecture in Table 4. While listwise loss functions (Cross-Entropy and ListNet [2]) either perform similar to or slightly better than pointwise losses (such as Binary Cross-Entropy), CROSS-JEM trained with RPL performs much better than any pointwise or listwise loss.

Incorporating Token Sequence Order Information: CROSS-JEM encodes the union of tokens from all ranking items, which enables fast inference by reducing the input sequence length. However, this also discards the original token order information within each item, treating them as bags of tokens. Though in our preliminary experiments, we observe that the loss of sequence information in the encoder has a negligible impact on the accuracy particularly in short-text ranking tasks (within 1%, refer Appendix F). However, token sequence information could be crucial in many ranking scenarios, and could get completely ignored in CROSS-JEM. To address this limitation, we propose an extension of CROSS-JEM that leverages sinusoidal positional encodings [40] to inject the item sequence information back into the model via the selective pooling layer. The key idea is to preserve the original position indices of the tokens for each item in the ranking list, remove them during the encoding process, and then add them back to the corresponding token context vectors after the CROSS-JEM encoder. The positional encodings are computed based on the original position indices and are summed with the token context vectors. The resulting vectors are then pooled together for classification as described in Section 3.

This simple yet effective technique improves the MRR@10 by 0.2% on the MS MARCO dataset (cf. Table 4 (b)), without any additional latency. This technique could also be useful for improving CROSS-JEM performance on long texts, which we leave for future work.

6 Conclusion: Limitations and Future Work

We introduced CROSS-JEM, an accurate and efficient approach for joint ranking of a set of short-text items for a given query based on relevance. CROSS-JEM effectively addresses the two major challenges in existing ranking architectures – sub-optimal accuracy due to pointwise inference, and significantly higher computational cost. Our extensive evaluations on publicly available ranking benchmarks as well as large-scale sponsored search datasets reveal that CROSS-JEM significantly outperforms the baselines, establishing a new state of the art. The scope of this work is primarily focused on the ranking of short texts, a common requirement in both industrial Sponsored Search applications and academic benchmarks, including tasks like matching queries with webpage titles and ranking duplicate questions. While the current work demonstrates significant gains on such short-text ranking tasks, the proposed approach could be adapted for long-text ranking by incorporating post-hoc positional encodings (cf. Section 5.4) and more sophisticated attention mechanisms that account for longer document lengths. Exploring these adaptations is an area for future research. CROSS-JEM opens up new directions for designing accurate ranking architectures and algorithms, accounting for application-specific constraints.

7 Ethical Considerations

Our data usage and service provision practices have received approval from our legal and ethical boards. Socially, our research is significantly enhancing the efficiency and user experience for millions of people searching for goods and services online. This improvement is crucial in today’s context, as it enables contactless and time-efficient purchasing and delivery. Additionally, our work is boosting the revenue of numerous small and medium-sized businesses by expanding their market reach and lowering customer acquisition costs.

References

- [1] Zeynep Akkalyoncu Yilmaz, Wei Yang, Haotian Zhang, and Jimmy Lin. 2019. Cross-Domain Modeling of Sentence-Level Evidence for Document Retrieval. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 3490–3496. <https://doi.org/10.18653/v1/D19-1352>
- [2] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. 2007. Learning to rank: from pairwise approach to listwise approach. In *Proceedings of the 24th international conference on Machine learning (Corvallis, Oregon) (ICML '07)*. ACM, New York, NY, USA, 129–136. <https://doi.org/10.1145/1273496.1273513>
- [3] Junyan Chen, Frédéric Dubut, Jason (Zengzhong) Li, and Rangan Majumder. 2023. Make Every feature Binary: A 135B parameter sparse neural network for massively improved search relevance. <https://tinyurl.com/y9amxjau>
- [4] Wei Chen, Tie-yan Liu, Yanyan Lan, Zhi-ming Ma, and Hang Li. 2009. Ranking Measures and Loss Functions in Learning to Rank. In *Advances in Neural Information Processing Systems*, Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta (Eds.), Vol. 22. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2009/file/2f55707d4193dc27118a0f19a1985716-Paper.pdf
- [5] K. Dahiya, D. Saini, A. Mittal, A. Shaw, K. Dave, A. Soni, H. Jain, S. Agarwal, and M. Varma. 2021. DeepXML: A Deep Extreme Multi-Label Learning Framework Applied to Short Text Documents. In *WSDM*.
- [6] Zhuyun Dai and Jamie Callan. 2020. Context-aware document term weighting for ad-hoc search. In *Proceedings of The Web Conference 2020*. 1897–1907.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *CoRR* abs/1810.04805 (2018). arXiv:1810.04805 <http://arxiv.org/abs/1810.04805>
- [8] Yixing Fan, Xiaohui Xie, Yingqiong Cai, Jia Chen, Xinyu Ma, Xiangsheng Li, Ruqing Zhang, Jiafeng Guo, et al. 2022. Pre-training methods in information retrieval. *Foundations and Trends® in Information Retrieval* 16, 3 (2022), 178–317.
- [9] Luyu Gao, Zhuyun Dai, and Jamie Callan. 2021. Rethink Training of BERT Rerankers in Multi-stage Retrieval Pipeline. In *Advances in Information Retrieval: 43rd European Conference on IR Research, ECIR 2021, Virtual Event, March 28–April 1, 2021, Proceedings, Part II*. Springer-Verlag, Berlin, Heidelberg, 280–286. https://doi.org/10.1007/978-3-030-72240-1_26
- [10] Luyu Gao, Zhuyun Dai, and Jamie Callan. 2021. Rethink training of BERT rerankers in multi-stage retrieval pipeline. In *Advances in Information Retrieval: 43rd European Conference on IR Research, ECIR 2021, Virtual Event, March 28–April 1, 2021, Proceedings, Part II 43*. Springer, 280–286.
- [11] Jiafeng Guo, Yingqiong Cai, Yixing Fan, Fei Sun, Ruqing Zhang, and Xueqi Cheng. 2022. Semantic models for the first-stage retrieval: A comprehensive review. *ACM Transactions on Information Systems (TOIS)* 40, 4 (2022), 1–42.
- [12] V. Karpukhin, B. Oguz, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, and W.-T. Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *EMNLP*.
- [13] Omar Khattab, Christopher Potts, and Matei Zaharia. 2021. Baleen: Robust multi-hop reasoning at scale via condensed retrieval. *Advances in Neural Information Processing Systems* 34 (2021), 27670–27682.
- [14] Omar Khattab and Matei Zaharia. 2020. Colbert: Efficient and effective passage search via contextualized late interaction over bert. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*. 39–48.
- [15] J. Lin, R. Nogueira, and Yates A. 2021. Pretrained Transformers for Text Ranking: BERT and Beyond. In *NAACL*.
- [16] Shichen Liu, Fei Xiao, Wenwu Ou, and Luo Si. 2017. Cascade ranking for operational e-commerce search. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 1557–1565.
- [17] Xueqing Liu, Chi Wang, Yue Leng, and ChengXiang Zhai. 2018. LinkSO: a dataset for learning to retrieve similar question answer pairs on software development forums. In *4th ACM SIGSOFT International Workshop on NLP for Software Engineering*. 2–5. <https://doi.org/10.1145/3283812.3283815>
- [18] Wenhao Lu, Jian Jiao, and Ruofei Zhang. 2020. TwinBERT: Distilling Knowledge to Twin-Structured Compressed BERT Models for Large-Scale Retrieval. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management (Virtual Event, Ireland) (CIKM '20)*. Association for Computing Machinery, New York, NY, USA, 2645–2652. <https://doi.org/10.1145/3340531.3412747>
- [19] Xueguang Ma, Liang Wang, Nan Yang, Furu Wei, and Jimmy Lin. 2024. Fine-Tuning LLaMA for Multi-Stage Text Retrieval. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval (Washington DC, USA) (SIGIR '24)*. Association for Computing Machinery, New York, NY, USA, 2421–2425. <https://doi.org/10.1145/3626772.3657951>
- [20] Xueguang Ma, Liang Wang, Nan Yang, Furu Wei, and Jimmy Lin. 2024. Fine-tuning llama for multi-stage text retrieval. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2421–2425.
- [21] Xueguang Ma, Xinyu Zhang, Ronak Pradeep, and Jimmy Lin. 2023. Zero-shot listwise document reranking with a large language model. *arXiv preprint arXiv:2305.02156* (2023).
- [22] Yu A Malkov and Dmitry A Yashunin. 2018. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and machine intelligence* 42, 4 (2018), 824–836.
- [23] Irina Matveeva, Chris Burges, Timo Burkard, Andy Laucius, and Leon Wong. 2006. High accuracy retrieval with multiple nested ranker. In *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*. 437–444.
- [24] Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2022. MTEB: Massive Text Embedding Benchmark. *arXiv preprint arXiv:2210.07316* (2022). <https://doi.org/10.48550/ARXIV.2210.07316>
- [25] Rodrigo Nogueira and Kyunghyun Cho. 2020. Passage Re-ranking with BERT. arXiv:1901.04085 [cs.IR]
- [26] Rodrigo Nogueira, Zhiying Jiang, Ronak Pradeep, and Jimmy Lin. 2020. Document Ranking with a Pretrained Sequence-to-Sequence Model. In *Findings of the Association for Computational Linguistics: EMNLP 2020, Trevor Cohn, Yulan He, and Yang Liu (Eds.)*. Association for Computational Linguistics, Online, 708–718. <https://doi.org/10.18653/v1/2020.findings-emnlp.63>
- [27] Rodrigo Nogueira, Wei Yang, Kyunghyun Cho, and Jimmy Lin. 2019. Multi-stage document ranking with BERT. *arXiv preprint arXiv:1910.14424* (2019).
- [28] Ronak Pradeep, Sahel Sharifmoghadam, and Jimmy Lin. 2023. Rankvicuna: Zero-shot listwise document reranking with open-source large language models. *arXiv preprint arXiv:2309.15088* (2023).
- [29] Ronak Pradeep, Sahel Sharifmoghadam, and Jimmy Lin. 2023. RankZephyr: Effective and Robust Zero-Shot Listwise Reranking is a Breeze! *arXiv preprint arXiv:2312.02724* (2023).

- [30] Tao Qin, Tie-Yan Liu, and Hang Li. 2008. Query-level loss functions for information retrieval. *Information Processing Management* 44, 2 (January 2008), 838–855. <https://www.microsoft.com/en-us/research/publication/query-level-loss-functions-for-information-retrieval/>
- [31] Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, et al. 2024. Large Language Models are Effective Text Rankers with Pairwise Ranking Prompting. In *Findings of the Association for Computational Linguistics: NAACL 2024*. 1504–1518.
- [32] Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, Xuanhui Wang, and Michael Bendersky. 2024. Large Language Models are Effective Text Rankers with Pairwise Ranking Prompting. In *Findings of the Association for Computational Linguistics: NAACL 2024*, Kevin Duh, Helena Gomez, and Steven Bethard (Eds.). Association for Computational Linguistics, Mexico City, Mexico, 1504–1518. <https://doi.org/10.18653/v1/2024.findings-naacl.97>
- [33] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108* (2019).
- [34] Keshav Santhanam, Omar Khattab, Jon Saad-Falcon, Christopher Potts, and Matei Zaharia. 2022. ColBERTv2: Effective and Efficient Retrieval via Lightweight Late Interaction. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (Eds.). Association for Computational Linguistics, Seattle, United States, 3715–3734. <https://doi.org/10.18653/v1/2022.naacl-main.272>
- [35] Peter Schonhofen. 2006. Identifying Document Topics Using the Wikipedia Category Network. In *2006 IEEE/WIC/ACM International Conference on Web Intelligence (WI 2006 Main Conference Proceedings)(WI'06)*, 456–462. <https://doi.org/10.1109/WI.2006.92>
- [36] Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorf, Wen-tau Yih, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. 2023. One Embedder, Any Task: Instruction-Finetuned Text Embeddings. In *Findings of the Association for Computational Linguistics: ACL 2023*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 1102–1121. <https://doi.org/10.18653/v1/2023.findings-acl.71>
- [37] S. J. Subramanya, Devvrit, R. Kadekodi, R. Krishnaswamy, and H. Simhadri. 2019. DiskANN: Fast Accurate Billion-point Nearest Neighbor Search on a Single Node. In *NeurIPS*.
- [38] Weiwei Sun, Lingyong Yan, Xinyu Ma, Shuaiqiang Wang, Pengjie Ren, Zhumin Chen, Dawei Yin, and Zhaochun Ren. 2023. Is ChatGPT Good at Search? Investigating Large Language Models as Re-Ranking Agents. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. 14918–14937.
- [39] Andrew Trotman, Antti Puurula, and Blake Burgess. 2014. Improvements to BM25 and Language Models Examined. In *Proceedings of the 19th Australasian Document Computing Symposium* (Melbourne, VIC, Australia) (ADCS '14). Association for Computing Machinery, New York, NY, USA, 58–65. <https://doi.org/10.1145/2682862.2682863>
- [40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [41] L. Xiong, C. Xiong, Y. Li, K.-F. Tang, J. Liu, P. Bennett, J. Ahmed, and A. Overwijk. 2021. Approximate nearest neighbor negative contrastive learning for dense text retrieval. In *ICLR*.
- [42] Longhui Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, Meishan Zhang, and Min Zhang. 2024. A Two-Stage Adaptation of Large Language Models for Text Ranking. In *Findings of the Association for Computational Linguistics ACL 2024*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand and virtual meeting, 11880–11891. <https://aclanthology.org/2024.findings-acl.706>
- [43] Wayne Xin Zhao, Jing Liu, Ruiyang Ren, and Ji-Rong Wen. 2024. Dense Text Retrieval Based on Pretrained Language Models: A Survey. *ACM Trans. Inf. Syst.* 42, 4, Article 89 (feb 2024), 60 pages. <https://doi.org/10.1145/3637870>
- [44] H. Zhou, M. Huang, Y. Mao, C. Zhu, P. Shu, and X. Zhu. 2019. Domain-Constrained Advertising Keyword Generation. In *WWW*.
- [45] Yucheng Zhou, Tao Shen, Xiubo Geng, Chongyang Tao, Can Xu, Guodong Long, Binxing Jiao, and Daxin Jiang. 2023. Towards Robust Ranker for Text Retrieval. In *Findings of the Association for Computational Linguistics: ACL 2023*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 5387–5401. <https://doi.org/10.18653/v1/2023.findings-acl.332>
- [46] Honglei Zhuang, Zhen Qin, Rolf Jagerman, Kai Hui, Ji Ma, Jing Lu, Jianmo Ni, Xuanhui Wang, and Michael Bendersky. 2023. RankT5: Fine-Tuning T5 for Text Ranking with Ranking Losses. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Taipei, Taiwan) (SIGIR '23). Association for Computing Machinery, New York, NY, USA, 2308–2313. <https://doi.org/10.1145/3539618.3592047>
- [47] Shengyao Zhuang, Honglei Zhuang, Bevan Koopman, and Guido Zuccon. 2024. A setwise approach for effective and highly efficient zero-shot ranking with large language models. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 38–47.

A Proofs of Lemma 1 and Corollary 2

To derive the proof, we recall both Lemma 1 and the Ranking Probability Loss:

$$\mathcal{L}^{RPL} = - \sum_{i=1}^{|\mathbb{Q}_{tr}|} \sum_{j=1}^N \left(\sum_{k \in \mathbb{L}_j} [\mathbf{y}_i]_k \right) \log \left(\text{SoftMax} \left(\sum_{k \in \mathbb{L}_j} [\mathbf{f}_{q_i}]_k \right) \right),$$

where $\tilde{s}_{ij} = \sum_{k \in \mathbb{L}_j} [\mathbf{f}_{q_i}]_k$ is the modified score associated with the query-item pair (q_i, k_j) . The following Lemma sheds light on the relationship between RPL and ranking probability distribution.

LEMMA. (Ranking Probability Loss) Let $\mathbf{P} \in \mathbb{R}^{N \times N}$ denote a matrix with entries p_{jk} given by $p_{jk} = \text{Prob}(\text{ranking item } k_j \text{ at location } k) \triangleq C \sum_{\ell \in \mathbb{L}_k} [\mathbf{f}_{q_i}]_\ell$, where C is a normalizing constant. Then, the Ranking Probability Loss maximizes the probability of ranking queries \mathbb{K}_{q_i} in the ordering of the ground-truth ranking \mathbf{y}_i .

PROOF. Given \mathbf{f}_{q_i} , the ground-truth scores \mathbf{y}_i , and a candidate item k_j , we define the set $\mathbb{L}_j = \{k \in \{1, N\} : [\mathbf{y}_i]_k < [\mathbf{y}_i]_j\}$, i.e., \mathbb{L}_j comprises the indices k for which the ground truth score of $[\mathbf{y}_i]_j$ is larger than the ground truth score at location k . Additionally, without loss of generality, we assume that \mathbf{f}_{q_i} and \mathbf{y}_i are normalized to have entries that lie in $[0, 1]$. We have:

$$p_{jk} = \text{Prob}(\text{ranking item } k_j \text{ at location } k) \triangleq C \sum_{\ell \in \mathbb{L}_k} [\mathbf{f}_{q_i}]_\ell,$$

Where C is an appropriately chosen constant, such that \mathbf{P} is a matrix with entries summing to one. Further, in the context of ranking with the logits/scores, we have:

$$\begin{aligned} p_{jk} &= \text{Prob}(\text{ranking item } k_j \text{ at location } k) \\ &= \text{Prob}(\text{logit of item } k_j > \{\text{logits of all items } k_k \text{ such that } k \in \mathbb{L}_k\}). \end{aligned}$$

Since the entire analysis is presented in the context of a each query q_i , for convenience, we abuse notation, and denote $[\mathbf{f}_{q_i}]_j = [\mathbf{f}]_j = f_j$. Then, we have

$$\begin{aligned} p_{jk} &= \text{Prob}(\text{logit of item } k_j > \{\text{logits of all items } k_k \text{ such that } k \in \mathbb{L}_k\}) \\ &= \text{Prob}(f_j > \{f_k \text{ for all } k \in \mathbb{L}_k\}) \\ &= \text{Prob}\left(f_j > \max_{k \in \mathbb{S}_k} \{f_k\}\right) \\ &\approx \text{Prob}\left(f_j > \sum_{\ell \in \mathbb{L}_k} f_\ell\right), \end{aligned} \tag{7}$$

where the last step is a consequence of the norm-bound $\|\mathbf{f}\|_1 \leq \|\mathbf{f}\|_\infty$. First, we note that, given the condition in Equation (7), for all k ranked higher than j , $p_{jk} = 0$. Further, for all k ranked lower than j , it suffices to show that the probability p_{jk} is non-zero, for $j = k$. To verify this, consider two probabilities p_{ik_1} and p_{ik_2} , with k_1 ranked higher than k_2 . Let $\mathbf{1}_{ik_1}$ denote the indicator of the event associated with p_{ik_1} . We have

$$p_{ik_2} \approx \text{Prob}\left(f_j > \sum_{\ell \in \mathbb{L}_{k_2}} f_\ell \mid \mathbf{1}_{ik_1}\right) \text{Prob}(\mathbf{1}_{ik_1}) = p_{ik_1}$$

To build intuition for the this, assume without loss of generality that \mathbf{f} have been sorted in a non-increasing manner. Let $j = 1$, $k_1 = 1$ and $k_2 = 3$. Then, it is clear to see that

$$\begin{aligned} &\text{Prob}\left(f_1 > f_4 + f_5 + \dots \mid f_1 > f_2 + f_3 + f_4 + f_5 + \dots\right) \text{Prob}(f_1 > f_2 + f_3 + \dots) \\ &= \text{Prob}(f_1 > f_2 + f_3 + \dots). \end{aligned}$$

Given \mathbf{P} built as described above, and \mathbf{P}^* defined similarly over scores \mathbf{y} , from the discussion above, we see that minimizing the distance between \mathbf{P} and \mathbf{P}^* is equivalent to minimizing the distance between the vectors $\mathbf{p} = \text{Diag}(\mathbf{P})$ and $\mathbf{p}^* = \text{Diag}(\mathbf{P}^*)$. We observe that, when normalized, the entries of \mathbf{p} correspond to the probability of ranking item k_j and rank j , given sorted vectors $\tilde{\mathbf{f}}$ and $\tilde{\mathbf{y}}$. Then, the RPL is the binary cross entropy loss defined between \mathbf{p} and \mathbf{p}^* , and represents minimizing the KL-divergence between the predicted ranking probability distribution $\tilde{\mathbf{f}}$ and the ground-truth ranking distribution $\tilde{\mathbf{y}}$, up to a normalizing constant factor. This completes the proof of Lemma 1. \square

Proof of Corollary 2: To link the Ranking Probability Loss to the ListNet loss [2], we recall that:

$$\mathcal{L}^{LN} = - \sum_{i=1}^{|\mathbb{Q}_{tr}|} \sum_{j=1}^N P_{\mathbf{y},j} \log(P_{s,j}),$$

where $P_{x,j} = \frac{\Phi([\mathbf{x}]_j)}{\sum_{\ell=1}^N \Phi([\mathbf{x}]_\ell)}$ Setting $[\tilde{\mathbf{s}}_{q_i}]_j = \tilde{s}_{ij} = \sum_{\ell \in \mathbb{L}_j} [f_{q_i}]_\ell$, and $[\tilde{\mathbf{y}}_i]_j = \tilde{y}_{ij} = \sum_{\ell \in \mathbb{L}_j} [y_i]_\ell$ and a Φ that result in mapping $\tilde{\mathbf{s}}$ and $\tilde{\mathbf{y}}$ to be valid probabilities distributions, we see the equivalence between the ListNet loss, and the proposed Ranking Probability Loss.

B The CROSS-JEM Inference

The procedure for forward pass of a CROSS-JEM model during training as well as inference is outlined in Algorithm 2, while the procedure to obtain the attention map over the item union set \mathbb{T} is described in Algorithm 1.

Algorithm 1 Method to create attention masks for item k_i . **Input:** T_q : Tokenized Query, T_U : Tokens in item Union Set, k : tokenized k_i keyword. **Output:** AttMask: Attention Mask for the keyword

```

1: procedure GETKUATTENTIONMASK( $T_q, T_U, k$ )
2:   AttMask  $\leftarrow$  []
3:   for  $i$  from 0 to LEN( $T_q$ ) - 1 do
4:     AttMask.ADDITEM(1)
5:   end for
6:   for  $i$  from 0 to LEN( $T_U$ ) - 1 do
7:     if  $T_U[i]$  in  $k$  then
8:       AttMask.ADDITEM(1)
9:     else
10:      AttMask.ADDITEM(0)
11:    end if
12:   end for
13:   return AttMask
14: end procedure

```

Algorithm 2 Getting relevance scores for a query and retrieved set of items using CROSS-JEM. **Input:** Query q , retrieved set of N items for q where $K_q = \{k_0, k_1, \dots, k_{N-1}\}$. **Output:** Scores $S = \{s_0, s_1, \dots, s_{N-1}\}$

```

1:  $T_q \leftarrow$  TOKENIZE( $Q$ ) ▷ Tokenize the query
2:  $T_U \leftarrow$  {}
3:  $K_{tokens} \leftarrow$  [] ▷ Store tokenized items
4: for  $k$  in  $K_q$  do
5:    $k_{tokens} \leftarrow$  TOKENIZE( $k$ )
6:    $T_U \leftarrow$  UNION( $T_U, k_{tokens}$ )
7:    $K_{tokens}.ADDITEM(k_{tokens})$ 
8: end for
9:  $T_U \leftarrow$  SORTED( $T_U$ )
10:  $KUAttMask \leftarrow$  [] ▷  $KUAttMask$ : item Union Attention Mask
11: for  $i$  from 0 to  $N - 1$  do
12:   AttMask  $\leftarrow$  GETKUATTENTIONMASK( $T_q, T_U, K_{tokens}[i]$ ) (cf. Algorithm 1)
13:    $KUAttMask.ADDITEM(AttMask)$ 
14: end for
15: sepToken  $\leftarrow$  TOKENIZE([SEP]) ▷ Token id for [SEP] token
16: encInpToks  $\leftarrow$   $T_Q$  ▷ Tokens to be passed through the Encoder
17: encInpToks.ADDITEM(sepToken)
18: for  $t_U$  in  $T_U$  do
19:   encInpToks.ADDITEM( $t_U$ )
20: end for
21:  $E \leftarrow$  ENCODER(encInpToks)
22:  $S \leftarrow$  SELECTIVEPOOLING( $E, KUAttMask$ ) ▷  $S \in \mathbb{R}^{N \times d}$ 
23:  $S \leftarrow$  CLASSIFIER( $S$ ) ▷  $S \in \mathbb{R}^N$ 
24: return  $S$ 

```

C Datasets

MS MARCO: MS MARCO document re-ranking dataset contains queries and their clicked web pages consisting of webpage title, URL, and passage. We create a short-text version of the dataset by considering only the titles of the webpages, making the length statistics of

Table 6: Improvement in performance of a large Seq2Seq model, RankT5-base after fine-tuning on short-text ranking benchmarks

Method	SODQ		MS MARCO	
	MAP@5	MAP@10	MRR@5	MRR@10
RankT5-base (pre-trained)	45.66	49.47	27.87	29.75
RankT5-base (fine-tuned)	55.9	56.8	33.72	35.14

the dataset better aligned with the real-world applications of ranking in sponsored search. We experiment with HDCT [6] retriever based training dataset consisting of 0.37M training queries and top 10 predictions from HDCT along with their ground truth click labels. We use the dev set to report our metrics as this set has ground truth labels available for evaluation. We use ~3.7M training pairs available in MS MARCO HDCT dataset sourced from [6] as described above for training CROSS-JEM and baselines using DistilBERT [33]. The target scores for all training pairs are obtained from a monoBERT model trained on binary ground truth click data with BERT-Base [7] as the base encoder for MS MARCO dataset.

SODQ: Stack Overflow Duplicate Questions dataset involves ranking questions on Stack Overflow as duplicates or not with the tags Java, JavaScript and Python [17]. It is also one of the re-ranking datasets on the popular MTEB benchmark [24]. StackOverflowDupQuestions on MTEB benchmark is the only short text re-ranking dataset with both training and evaluation data available, and is hence used in our experiments. Similar to our experiments on MS MARCO dataset, the target scores for training CROSS-JEM as well as baselines are obtained from a BERT [7] based monoBERT model trained on binary relevance of duplicate questions.

Sponsored Search Dataset: The training dataset for sponsored search query to advertiser matching task is created using a BERT-Large based monoBERT model trained on manually labeled and good-click data as the teacher model. A query-item (advertiser bid keyword) pair in the good click data is obtained when the user clicked on the ad corresponding to an advertiser keyword in response to their query, and did not close the ad quickly indicating they found it relevant. This BERT-Large teacher model was used to score 100 predicted items each for 18.6M queries on the search engine during a time period. This resulted in around 1.8B query-item pairs with scores in 0 to 1 range as training data for CROSS-JEM and all baselines in Table 3.

D Metrics

- **Mean Average Precision (MAP):** This is a ranking metric defined as the mean of Average Precision (AP) over the positive and negative detected classes:

$$\frac{1}{|Q|} \sum_{u=1}^{|U|} \left(\frac{1}{m} \sum_{k=1}^N P_u(k) \cdot rel_u(k) \right) \quad (8)$$

where $|Q|$ is the total number of queries, $P_u(k)$ is the precision at cut-off k in the list, $rel_u(k)$ is an indicator function equaling 1 if the item at rank k is a relevant document, otherwise zero.

- **Mean Reciprocal Rank (MRR):** Rank is defined as the position of the first relevant item in the ranked list. MRR is hence defined as below:

$$\frac{1}{|Q|} \sum_{i=1}^{|Q|} \left(\frac{1}{rank_i} \right) \quad (9)$$

- **Accuracy:** Positive, Negative, and Overall Accuracy denote the proportion of positive, negative, and overall instances, respectively, in the test set that are accurately identified.
- **Area Under the ROC Curve (AUC-ROC):** The ROC curve is a plot of True Positive Rate (TPR) or sensitivity against False Positive Rate (FPR) at different thresholds.

E Baselines and Hyperparameters

For baselines monoBERT, DPR, and ANCE, we tune the following hyperparameters based on the validation set accuracy: learning rate, weight decay, number of epochs. ColBERT is trained and evaluated with default hyperparameters provided by the authors¹. We use the pre-trained checkpoint² and code-base³ provided by the authors for INSTRUCTOR model, and test its zero-shot performance using the instructions mentioned in the paper. CROSS-JEM is trained with exactly same setting as monoBERT: learning rate of 1e-4, linear learning rate scheduler, and AdamW optimizer. Hyperparameters specific to CROSS-JEM (N and L_u) are provided in Table 7. The metrics for BM25 on SODQ are taken from [17], while they are computed following Trotman et al. [39] on MS MARCO.

Fine-tuning RankT5-base (24L) for Short-text Ranking: We fine-tune the model checkpoint available from Zhang et al. [42] on short-text ranking benchmarks and note the performance improvements in Table 6.

¹<https://github.com/stanford-futuredata/ColBERT>

²<https://huggingface.co/hkunlp/instructor-base>

³<https://github.com/xlang-ai/instructor-embedding>

Table 7: Hyperparameters used in CROSS-JEM.

Hyperparam	SODQ	MS MARCO	Sponsored Search
N	30	10	100
L_u	265	360	262

Table 8: Sponsored search dataset statistics motivating CROSS-JEM architecture design. (a) Cross-encoder trained with ordered tokens (\mathcal{E}_{CE}) and alphabetically sorted tokens for all items (\mathcal{E}'_{CE}). The small performance delta (between rows) indicates that sequence ordering is not critical for short-text. (b) The mean total tokens m in a retrieved item set Q_{te} and the mean size N_u of tokens in the union of Q_{te} . The ratio of m to N_u being ~ 5 indicates strong token overlap.

Algorithm	MAP@100	Max-length in word-piece (L)	m	N_u
\mathcal{E}_{CE}	93.03	12	498.38	93.70
		16	499.34	94.02
\mathcal{E}'_{CE}	92.76	32	501.52	94.60

(a) (b)

Table 9: Variation in accuracy on varying the number of items scored per query by CROSS-JEM. We observe only minor variation in changing the number of items to be ranked at inference time. This observation is useful in real-world ranking which receive item candidates from a set of retrieval algorithms and hence the number of items to be scored can vary with the query.

N	Negative Accuracy	Overall Accuracy	AUC
10	99.18	94.94	99.24
20	99.37	94.96	99.34
50	99.52	94.82	99.40
80	99.56	94.71	99.51
100	99.45	95.27	99.42

E.1 Compute

All baselines on MS MARCO and SODQ datasets including CROSS-JEM were trained on 8 V100 GPUs. Experiments on proprietary Sponsored Search dataset were conducted on larger GPU cluster with 16 V100s.

F Experiments: Sponsored Search Dataset

We compare CROSS-JEM against methods that can be deployed for real-time ranking including ANCE, MEB, and TwinBERT. TwinBERT is a lighter version of ColBERT. It applies an MLP layer to individual query and keyword embeddings, unlike ColBERT which considers interactions along all token embeddings. This makes TwinBERT more efficient and practical for real-world systems due to lower storage requirements. Table 3 shows the comparison of CROSS-JEM with baseline methods in production where CROSS-JEM outperforms existing methods by large margins.

Ablation on Number of items per Query (N): From Table 9, we vary the number of items scored per query from 10 to 100. With more items, the token overlap increases, providing CROSS-JEM more opportunity for joint modeling. Correspondingly, we observe gains in negative accuracy and AUC as items per query increase.

Ablation on Encoder \mathcal{E}_θ in CROSS-JEM: Table 10 shows that even with a smaller encoder, CROSS-JEM provides significant accuracy gains over sparse models like MEB while having low latency.

Ablation on Sequence Information: We analyze the effect of the ordering/sequencing of the item text on classification accuracy using a cross-encoder model; and a train dataset consists of about 100M query-item pairs (q, k) , drawn from $Q_{tr} \times I_{tr}$, mined from proprietary search engine logs. Given (q, k) , we compare two cross-encoder \mathcal{E} models: 1) \mathcal{E}_{CE} : Standard cross-encoder scoring the pairs (q, k) , and 2) \mathcal{E}'_{CE} : Cross-encoder trained to score pairs (q, k') , where the item k' is obtained by sorting the tokens in k alphabetically.

The hypothesis is that, if the cross-encoders \mathcal{E}_{CE} and \mathcal{E}'_{CE} have similar scoring accuracy, then the sequence ordering is relatively less informative for this task. While testing \mathcal{E}'_{CE} is evaluated on (q, k') pairs k' is drawn from I_{te} , with its tokens sorted alphabetically. Table 8(a) shows how the variants perform on a test set of 10M pairs. We observe that, when the sequence information in k is discarded, both the mean average precision (MAP) and accuracy are within 1% of the case when the sequence information is retained.

Table 10: Variation in accuracy with base encoder \mathcal{E}_θ in CROSS-JEM.

Encoder	Negative Accuracy	Overall Accuracy	AUC	Latency CPU
TinyBERT - 2 layer	96.52	92.76	96.67	114.3
DistilBERT - 6 layer	99.45	95.27	99.61	744.7

G Qualitative Analysis

Table 11: A comparison of the top-5 ranked items obtained using CROSS-JEM’s listwise modeling and a pointwise ranking model [25]. Relatively more *generic* (but still relevant) items are ranked higher in baseline predictions, owing to their frequency in training data, token-level matching and other biases, which are circumvented when using listwise architectures capable of evaluating all the shortlisted items in a single forward pass, and rank the more relevant (and specific) items higher.

Query	Top-5 Ranked Item in the Proposed Approach (CROSS-JEM)	Top-5 Ranked Item in a Baseline Cross Encoder
different foods of oxaca mexico	<ul style="list-style-type: none"> 'the foods of oxaca' 'oaxacan cuisine' 'exploring oxacan food' 'authentic recipes from oxaca, mexico' '6 things you'll love about oxaca' 	<ul style="list-style-type: none"> 'culinary tales: the kinds of food mexicans eat' 'mexican christmas foods' 'mexican cuisine' 'culture: food and eating customes in mexico' 'popular food in mexico'
what is the bovine growth hormone	<ul style="list-style-type: none"> 'recombinant bovine growth hormone' 'what is rbgh?' 'rbgh' 'bovine growth hormone and milk : what you need to know' 'what is rbst?' 	<ul style="list-style-type: none"> 'growth hormone' 'human growth hormone' 'alternative names for growth hormone' 'human growth hormone and insulin are friends' 'growth hormone (somatotropin)'
what is the state nickname of new mexico	<ul style="list-style-type: none"> 'what are the nicknames of the state of new mexico?' 'state nicknames new mexico - south carolina and their explanation' 'what is new mexico's nickname?' 'new mexico state names (etymology of names)' 'the state of new mexico' 	<ul style="list-style-type: none"> meh'the state of new mexico' meh'state of new mexico' 'new mexico' 'new mexico state university' 'state of mexico'
is the elliptical bad for your knees	<ul style="list-style-type: none"> 'does an elliptical make bad knees worse?' 'are the elliptical machines bad for your knees?' 'is an elliptical the best machine for knees that are chronically painful?' 'why does my knee hurt on an elliptical machine?' 'elliptical machine is good for osteoarthritis of the knee!' 	<ul style="list-style-type: none"> 'what exercises can help relieve knee pain?' '4 bad exercises for bad knees' 'how to treat a knee sprain' 'how to strengthen legs with bad knees' 'yoga bad for your knees, indian doctor warns'
what do you use for oxygen facial machines	<ul style="list-style-type: none"> 'oxygen facials and other skincare services' 'oxygenating facial treatment' 'oxygen facial : home kits, beauty benefits, side effects, process' 'benefits of oxygen facial' '4 beauty - boosting benefits of oxygen facials' 	<ul style="list-style-type: none"> 'using oxygen safely' 'the oxygen machine' 'shop cpap and oxygen' 'oxygen concentrators and generators' 'oxygen concentrators & generators'