

Transformer Contextualized Re-Ranking

Sebastian Hofstätter

sebastian.hofstaetter@tuwien.ac.at

 [/s_hofstaetter](https://twitter.com/s_hofstaetter)

Today

Transformer Contextualized Re-Ranking

1

Re-Ranking with BERT

- Concatenate query and document for a vanilla BERT use
- The mono-duo pattern; long document handling

2

Efficient Transformer-based Models

- Splitting BERT: PreTTR, ColBERT
- Our Transformer-Kernel family
- IDCM: A hybrid model

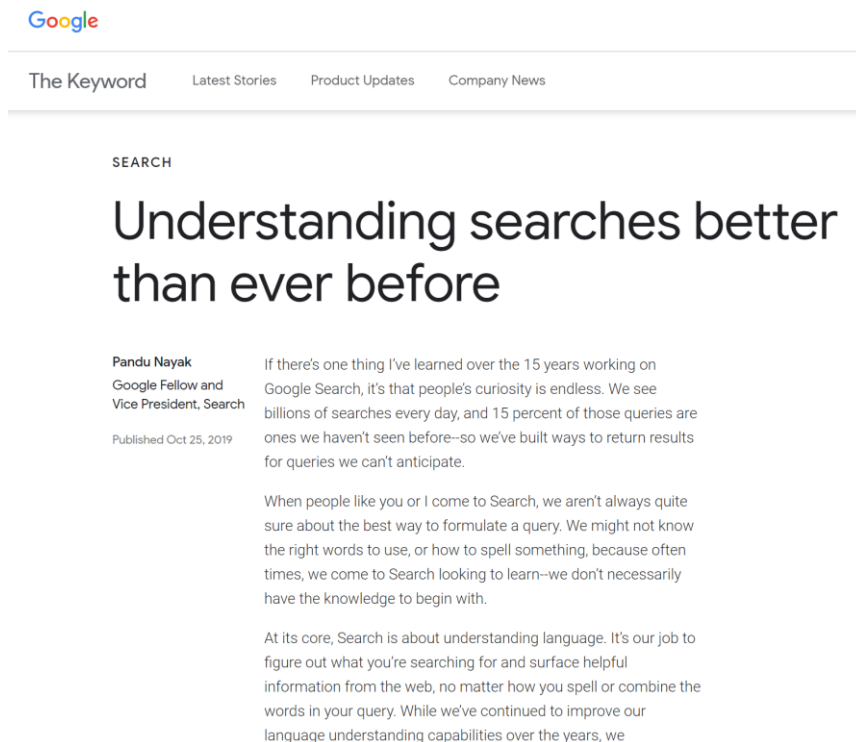
Now We Reached the State-of-the-art

- Finally! 🎉
- Fast moving field:
 - Most of the contents of this lecture did not exist in the beginning of 2018
 - And by July we probably have a new SOTA technique
- Many questions are open
 - We try to answer a few here
- General direction: More computation = better results,
possible by better hardware & more data

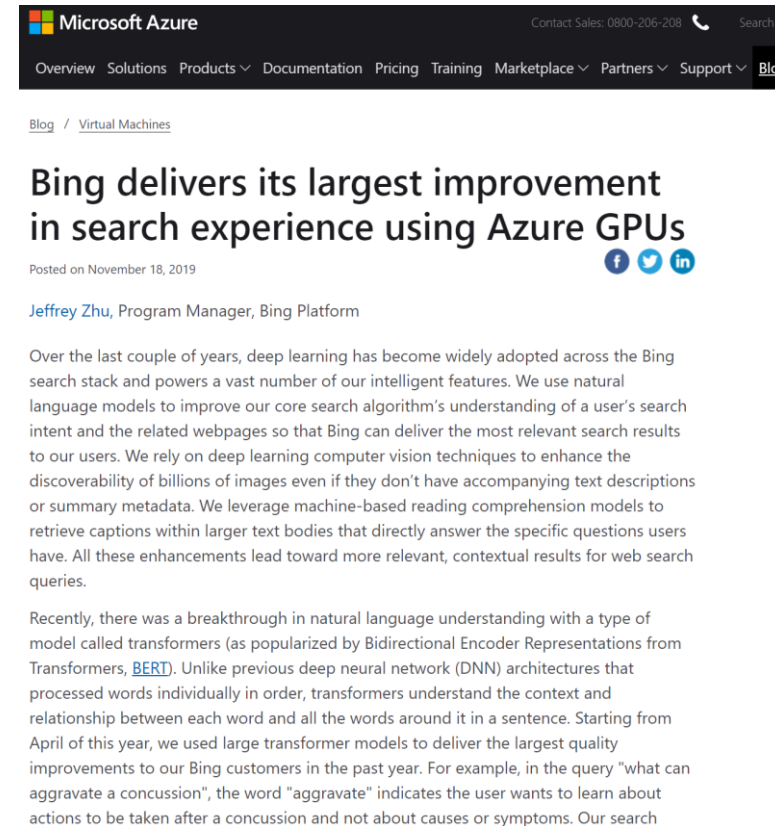
Context of this Lecture

- We are studying ad-hoc re-ranking models
 - Many other applications of Transformers in IR:
 - Document & query expansion
 - QA pipelines
 - Conversational search, using the query history
 - Dense retrieval
 - Knowledge graph-based search
- We are focusing on the efficiency-effectiveness tradeoff
 - How fast or slow is a model compared to the result quality
 - Many other aspects can be studied, f.e.: social biases, domain adaption, or multilingual models (and many more)

Web Search with BERT



Google (October 2019)



Microsoft (November 2019)

Re-Ranking with BERT

Drastically improving search result quality

Recall: BERT - Workflow

- Someone with lots of compute or time pre-trains a large model
 - BERT uses Masked Language Modelling [MASK] and Next Sentence Prediction [CLS]
- We download it and fine-tune on our task

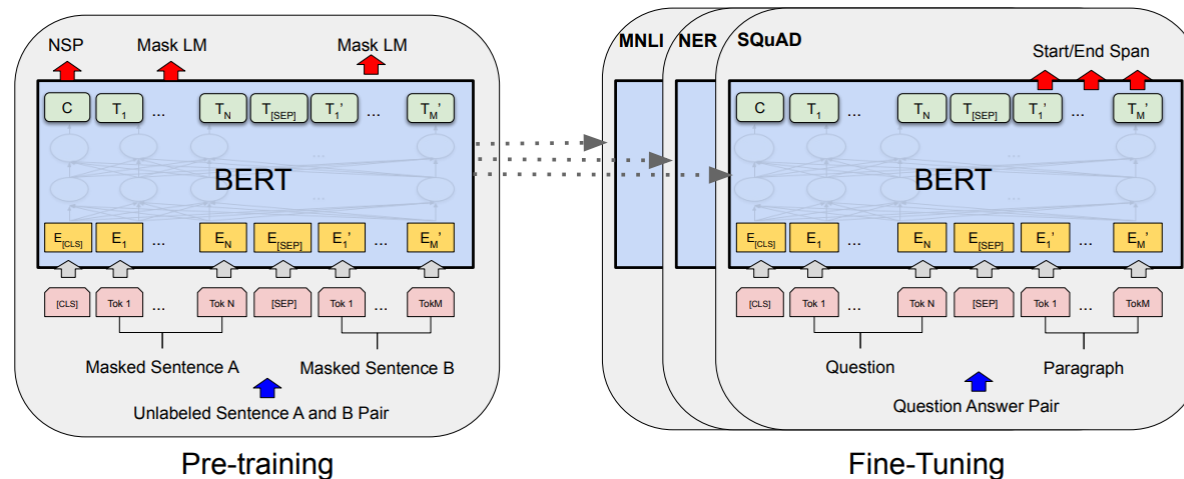
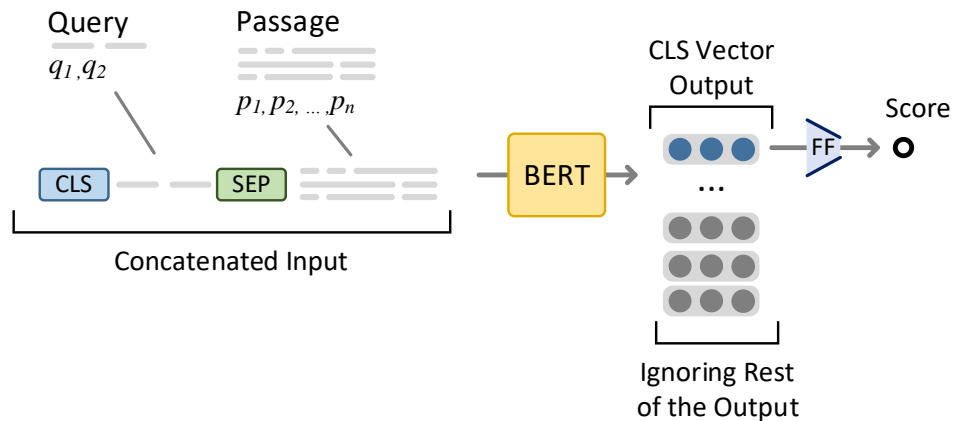


Figure taken from the BERT paper

BERT Re-Ranking: BERT_{CAT}



- Also known as monoBERT, vanilla BERT re-ranking, or simply BERT
- Concatenating the two sequences to fit BERT's workflow
 - [CLS] query [SEP] passage
 - Pool [CLS] token
 - Predict the score with a single linear layer
- Needs to be repeated for every passage



BERT_{CAT}


- Simple formula (as long as we abstract BERT):

$$\begin{array}{l} \text{BERT} \left[r = \text{BERT}([CLS]; \underbrace{q_{1..n}; [SEP]; p_{1..m}}_{\text{Concatenation}})_{CLS} \right. \\ \text{Scoring} \left[s = r * \underbrace{W}_{\text{Starts uninitialized}} \right] \end{array}$$

- We still have the choice of BERT-model

$q_{1..n}$	Query tokens
$p_{1..m}$	Passage tokens
BERT	Pre-trained BERT model
$[CLS]$ $[SEP]$	Special tokens
x_{CLS}	Pool the CLS vector
W	Linear Layer (from 768 dims to 1)
s	Output score

The Impact of BERT_{CAT}

- This model (first shown by Nogueira and Cho) jumpstarted the current wave of neural IR 
- Works awesome out of the box
 - Concatenating the two sequences to fit BERT's workflow
 - As long as you have time or enough compute it trains easily
- Major jumps in effectiveness across collections and domains
 - But, of course, comes at the cost of performance and virtually no interpretability
 - Larger BERT models roughly translate to slight effectiveness gains at high efficiency cost
 - The problem is we need to repeat the inference by the re-ranking depth!

So how good is BERT_{CAT} ?

- MSMARCO-Passage
 - MRR@10 from .194 (BM25) to .385 (ALBERT-Large)
 - BERT basically doubles the result quality
- MSMARCO-Document
 - MRR@10 from .252 (BM25) to .384 (DistilBERT with 2K tokens)
- Similar results for TREC-DL '19, '20, TripClick
 - Large Training Data Settings

See also the MSMARCO leaderboard: <https://microsoft.github.io/msmarco/>

Improving Efficient Neural Ranking Models with Cross-Architecture Knowledge Distillation; Hofstätter et al. <https://arxiv.org/abs/2010.02666>

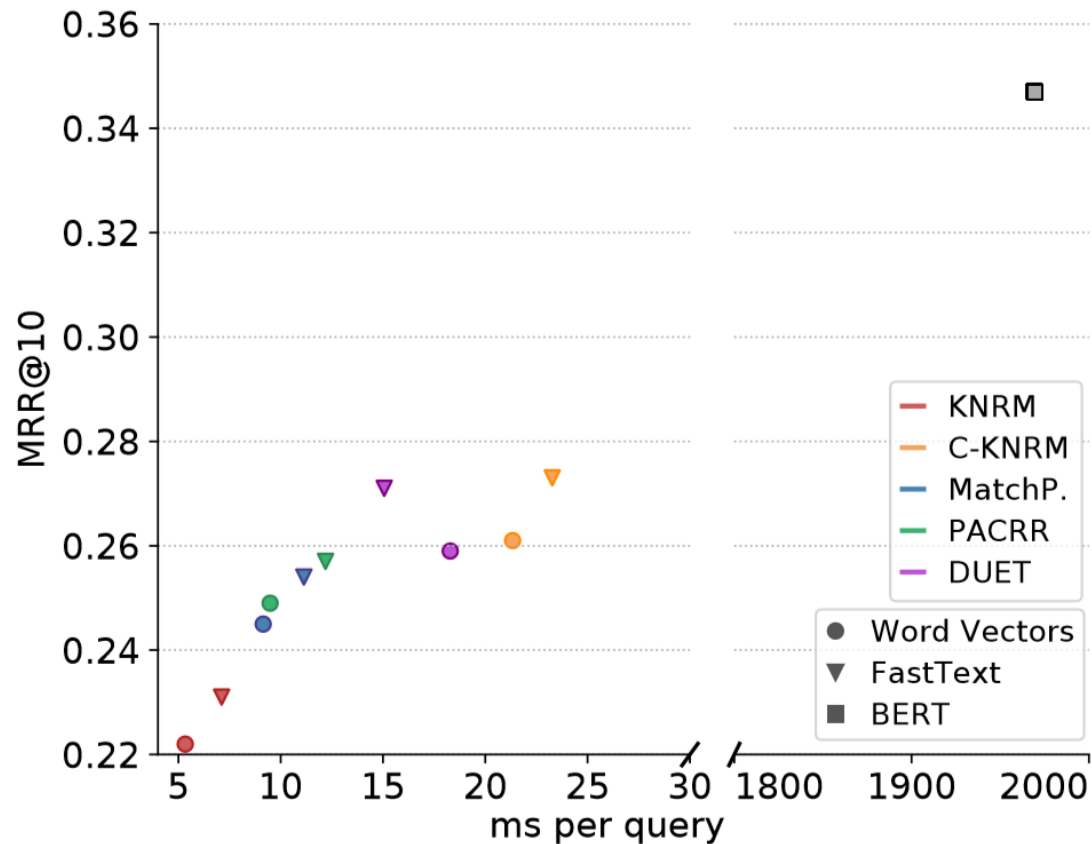
The Mono-Duo Pattern

- Mono-Duo is a multi-stage process:
 - Mono = $\text{score}(q, p) - \text{f.e. Top-1000}$
 - Duo = $\text{score}(q, p_1, p_2) - \text{f.e. Top-50}$ (but needs to be done 50^2 times)
- Improves on the single “mono” stage
 - Good for leaderboard settings, and showing the maximum achievable performance
- Can use BERT or T5 as base model
 - T5 is also a Transformer based language model, even larger but of course slower

BERT_{CAT} for Longer Documents

- BERT is capped at max. 512 input tokens (query + document)
- Simplest solution: just take cap the document at 512-query length
 - Works surprisingly well already (for MSMARCO-Documents)
 - But might not work well in other domains, where documents are really long, or contain a variety of topics at different depths
- Still simple, but working on full documents: Sliding window over the document -> take max window score as document score
 - Now, we can also make smaller sliding windows
 - Might be useful to use in the UI -> highlight the most relevant passage as snippet

BERT_{CAT} In-Efficiency



- Evaluated on 250 docs / query on short MSMARCO-Passage (*max 200 tokens*)
- Basic IR-specific networks are fast, but moderately effective
- Transformer-based BERT is very effective, but very slow
 - + Infrastructure cost (blocking 1 GPU for 2 seconds at a time)

Sebastian Hofstätter and Allan Hanbury. 2019.

Let's measure runtime! Extending the IR replicability infrastructure to include performance aspects . In OSIRRC @ SIGIR.

Efficient Transformer-based Models

Adapting BERT or starting from scratch altogether

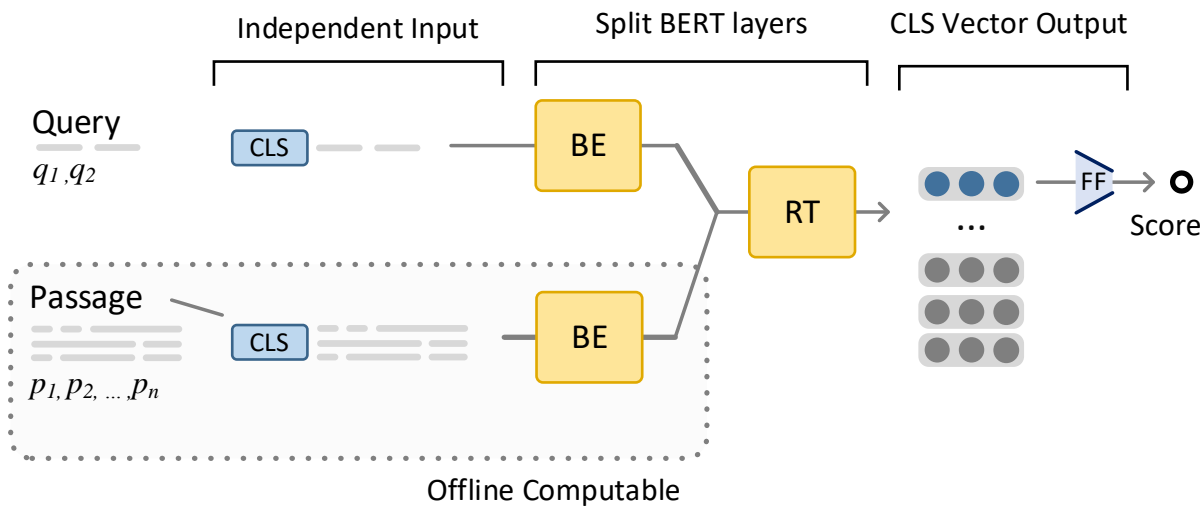
Achieving Efficiency

- Multiple paths to reduce query latency
 - Query latency is our focus today, but full lifecycle efficiency also a concern
 - Lifecycle efficiency includes training, indexing and retrieval steps
- ① Reduce model size
 - Smaller models run faster, duh!
 - Only possible until a certain threshold after which quality reduces drastically
- ② Move computation away from query-time
 - Pre-compute passage representation, so they become a simple lookup
 - Lightweight aggregation, that can be done at query time

Splitting BERT for Efficiency

- Concatenated BERT_{CAT} needs x re-ranking depth evaluations at query time
 - Bad for query latency, deep re-ranking implausible
- Most of that computation is spent on passage encoding
 - If we can move the passage encoding to the indexing phase
 - Encode each passage exactly 1 time
 - At query time only need to encode 1 query (with few words, so it's fast)
- Multiple approaches have been proposed to then *glue* the query and passage representations back together
 - With only a small reduction in effectiveness

Splitting BERT: PreTTR



- PreTTR splits BERT layers:
 - The first n layers are separated
 - The following layers are concatenated again
 - n is a hyperparameter
- We can pre-compute the first- n layers of the passages
- Pro: \approx quality as BERT_{CAT}
- Neg: Still low query latency + storage requirements

PreTTR

- Split in two parts:

Part 1

$$\begin{cases} \hat{q}_{1..n} = \text{BERT}_{1..b}([CLS]; q_{1..n}) \\ \hat{p}_{1..m} = \text{BERT}_{1..b}([CLS]; p_{1..m}) \end{cases}$$

Independent computation

Part 2

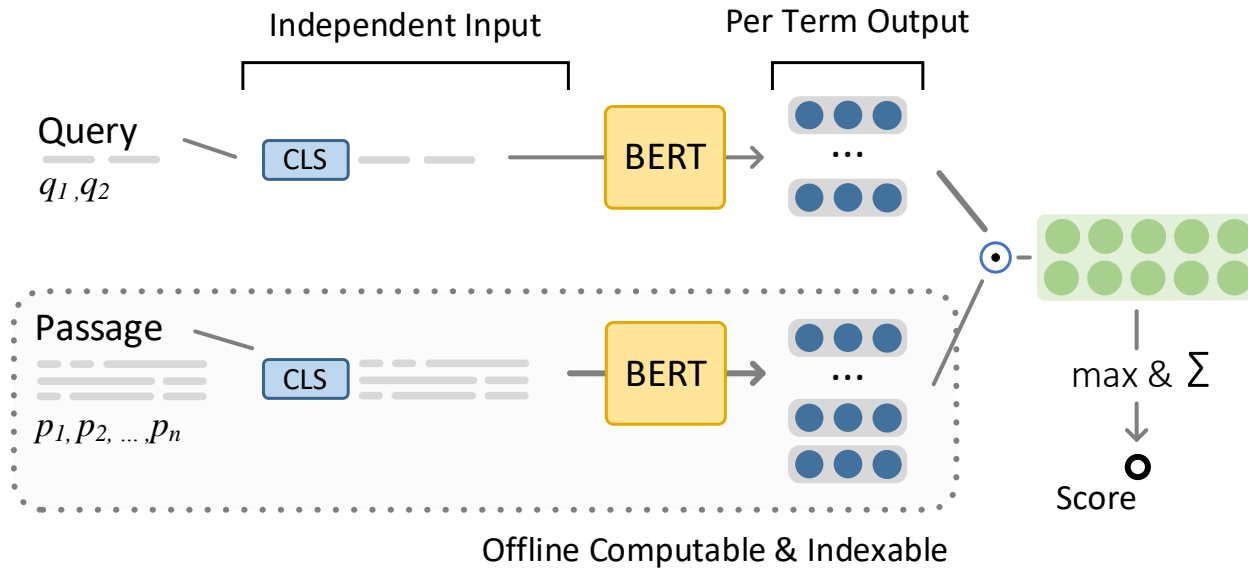
$$s = \underline{W} * \text{BERT}_{b..l}(\hat{q}_{1..n}[SEP]; \hat{p}_{1..m})$$

Starts uninitialized

- Optional compression of $\hat{q}_{1..n}$ & $\hat{p}_{1..m}$

$q_{1..n}$	Query tokens
$p_{1..m}$	Passage tokens
BERT	Pre-trained BERT model
[CLS] [SEP]	Special tokens
W	Linear Layer (from 768 dims to 1)
s	Output score

Splitting BERT: ColBERT



- ColBERT creates a match-matrix of BERT term-representations
- Then uses a simple max-pooling for the doc-dim & sum for the query dim
- Pro: Very fast query latency
- Con: Need to save all passage term vectors (huge storage cost)

ColBERT

- Simple formula (as long as we abstract BERT):

Encoding

$$\begin{cases} \hat{q}_{1..n} = \text{BERT}([CLS]; q_{1..n}) \\ \hat{p}_{1..m} = \text{BERT}([CLS]; p_{1..m}) \end{cases} \quad \text{Can be done at indexing time}$$

Aggregation

$$s = \sum_{i=1..n} \max_{t=1..m} \hat{q}_i \cdot \hat{p}_t \quad \text{Quite efficient at query time}$$

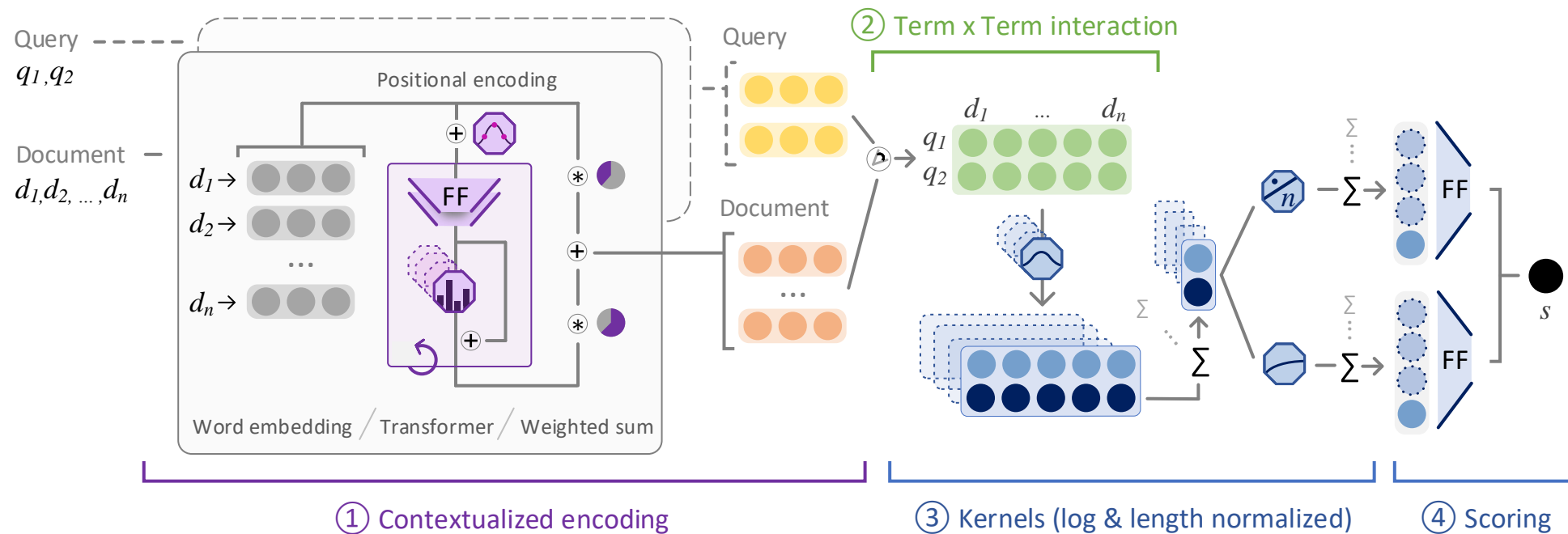
- Optional compression of \hat{q} , \hat{p} with a single linear layer

$q_{1..n}$	Query tokens
$p_{1..m}$	Passage tokens
BERT	Pre-trained BERT model
[CLS]	Special tokens
s	Output score

TK: Transformer-Kernel Ranking

- Desired properties: Lightweight, interpretable, effective
- We proposed the TK model
 - Combine Transformer-contextualization with kernel-pooling
 - Strong results compared to IR-specific models
 - State-of-the-art model for time-budget constrained environments
- Use Transformer-blocks as contextualization layer
 - Create hybrid contextualization by merging context & non-context
- Limit the number of Transformer layers, as each additional layer takes considerable amount of time with diminishing returns

TK: Transformer-Kernel Ranking



TK

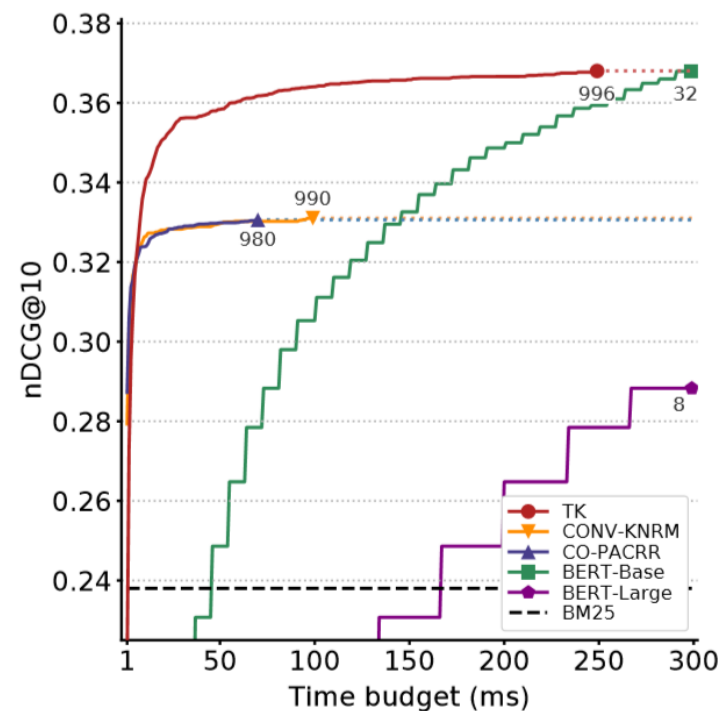
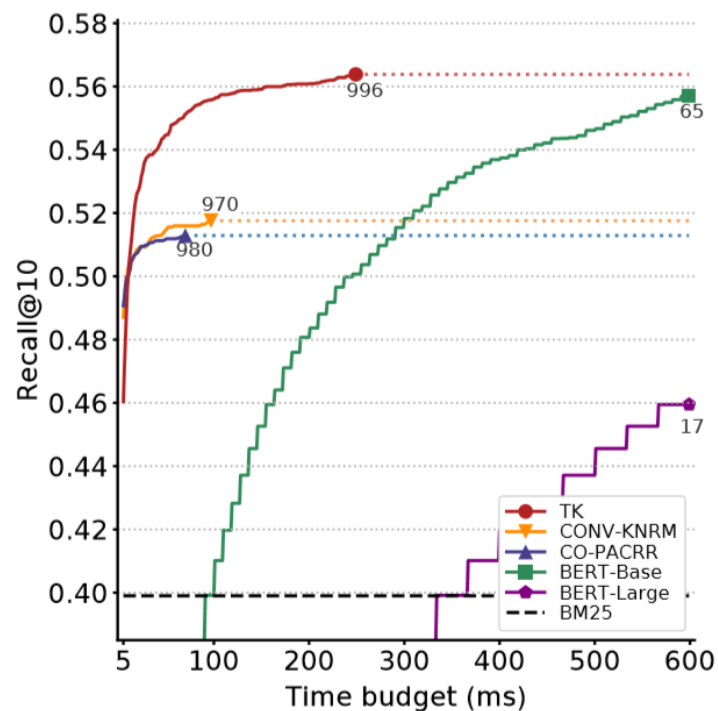
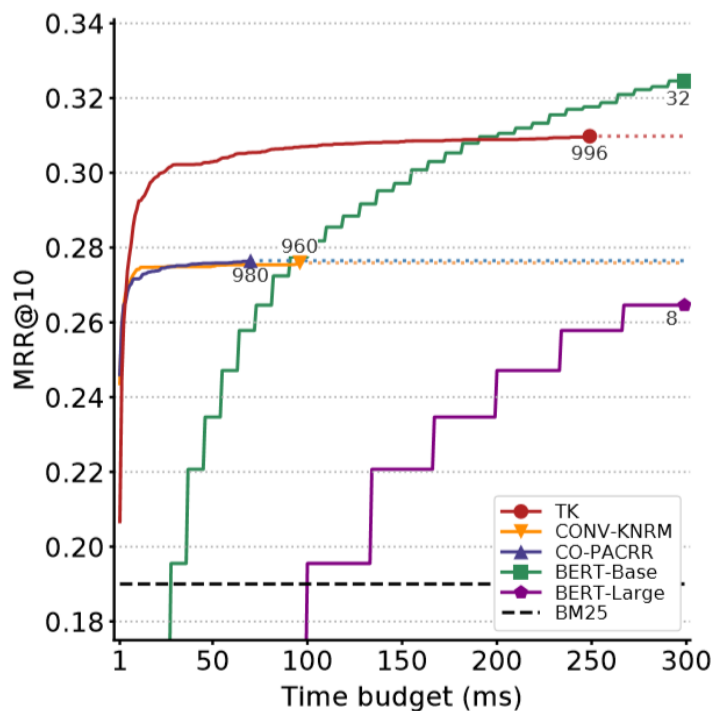
Encoding	$\hat{q}_{1..n} = TF(q_{1..n}) \quad \hat{d}_{1..m} = TF(d_{1..m})$	$q_{1..n}$ Query token vecs
Matching	$M_{ij} = \mathbf{cos}(\hat{q}_i, \hat{d}_j)$ A single match matrix	$d_{1..m}$ Doc. token vecs
Kernels	$K_k(M) = \sum_{i=1}^n \log(\sum_j \exp(-\frac{(M_{ij} - \mu_k)^2}{2\sigma_k^2}))$	M Match-matrix
Scoring	$s = \mathbf{FC}(K) = W * K + b$ Same as KNRM	TF Transformer
		K All Kernels
		K_k k -th kernel
		μ_k Similarity level
		σ_k Kernel-width/range
		W, b Weights & biases
		s Output score

Designing for a Time-Budget

- Large influence on how we employ a model
- Main outside factor: how many documents to re-rank
- Faster models can re-rank more documents in the same time as slower ones
- Evaluate based on a time budget
 - Allows us to simultaneously evaluate effectiveness & efficiency in a realistic setting

TK on a Time Budget

- More documents in the same time = better results



(a) MSIMARCO-Passage

Comparing the Models


- Many dimension to compare possible:
 - Quality (Effectiveness, diff. measures)
 - Query latency, Storage requirement, Capabilities, etc...

Model	Effect.	Query Latency	GPU Memory	Query-Passage Interaction	Passage Cache	NN Index	Storage (× Dim.)
BERT_{CAT}	1	950 ms	10.4 GB	All TF layers	–	–	–
BERT_{DOT}	× 0.87	23 ms	3.6 GB	Single dot product	✓	✓	P
ColBERT	× 0.97	28 ms	3.4 GB	$m * n$ dot products	✓	✓	T
PreTT	× 0.97	455 ms	10.9 GB	Min. 1 TF layer (here 3)	✓	–	T
TK	× 0.89	14 ms	1.8 GB	$m * n$ dot pr. + Kernel-pool	✓	–	T

Understanding TK

neural-ir-explorer – TK using 2 Transformer layers @ MSMARCO-Passage (Dev)

< | when and where did paella originate ?

Min. similarity  or Kernels 1 0.9 0.7 0.5 0.3 0.1 -0.1

Rank Score (length & log) Log-Kernel scores: 1 to -0.1 & rest

★ ① -23.57 (1.58 & -25.15) -4.38 -0.70 4.41 -2.19 -2.69 0.02 -0.65 -18.96

Paella has ancient roots , but its modern form originated in the mid - 19th century in the area around Albufera lagoon on the east coast of Spain , adjacent to the city of Valencia . [4] Many non - Spaniards view paella as Spain 's national dish , but most Spaniards consider it to be a regional Valencian dish .

② -23.90 (1.20 & -25.10) -4.45 -0.75 4.07 -2.05 -2.44 0.01 -0.55 -18.94

Paella is a world - famous dish , which originated in the region of Valencia , in eastern Spain . It is now widely eaten in all provinces of Spain , as well as every continent of the world . Like so many other popular recipes , Valencian paella was initially a peasant dish .

③ -24.83 (1.72 & -26.55) -4.95 -1.03 4.00 -2.20 -2.97 0.01 -0.53 -18.89

Paella is a dish that originated from Valencia , Spain . The main ingredients are rice , and saffron . Everyone has a different recipe for paella . Some paella is made from all seafood , some from all meat , and some are a mix of meat and seafood . Here is my recipe for paella . It is easy to make and very authentic tasting .

④ -25.13 (1.66 & -26.79) -5.14 -1.12 3.65 -2.01 -2.53 0.03 -0.74 -18.93


Paella . Originating in Valencia , paella is a rice dish prepared with seafood . Of all the foods in Spain , this is the most popular . In this dish , savory yellow rice is combined with tomatoes , onions , peas , shellfish , squid , clams and chicken drumsticks .

⑤ -25.53 (1.82 & -27.34) -5.56 -1.25 3.87 -2.12 -2.61 0.03 -0.80 -18.91

Authentic Paella Valenciana . I lived in Spain for two years where I was taught the art of making the Paella which originated in Valencia . I have n't found anything on here which is even close to authentic , so I thought I would add this recipe for those who would like to try a taste of Spain .

- Demo application to showcase TK
 - Displays internal similarity & kernel results
- Users can browse around the results
 - Get an overview over the queries
 - Dig deep into a single result
- Complements metric-based evaluation
- Allows users to develop a “feeling” for the test collection & model used

Sebastian Hofstätter, Markus Zlabinger and Allan Hanbury. 2020. Neural-IR-Explorer: A Content-Focused Tool to Explore Neural Re-ranking Results. In ECIR.

Sort    Collapse clusters Prefix filter

Let's start exploring

Here is what you see around you and what you can do with it:

- At the top you can sort the queries: randomly, ascending or descending (based on the rank of the first relevant document). You can also expand the clusters to see all queries.
- We clustered the queries based on their mean contextualized encoding. Each card holds the queries for a cluster.
- At the top of each card is: the median best rank of the neural model, the difference to the first-stage baseline, and a manual summary of the queries in that cluster
- Each query line contains: the best rank of the neural model, the difference to the first-stage baseline, and the query
- Simply click on a query to go to the result view and see details on the query result

⑤ +3 where is location

- 1 +1 what airport is in wilder ky
- 1 =0 where is alepotrypa cave
- 1 =0 where is azaz
- 1 +2 where is bell buckle tn
- 1 +4 where is boston georgia
- 1 =0 where is henry's plant farm
- 1 =0 where is last name hollis derived from
- 1 +1 where is lima beads located
- 1 +3 where is mathura

⑨ +5 general knowledge questions (trivia)

- 1 =0 do vhi swiftcare do blood tests?
- 1 +1 what do partnerships file tax in michigan
- 1 +6 what does android sdk tools do
- 1 =0 what eventually replaced the cottage industry
- 1 +9 what federal statute gives the epa authority to regulate pesticides
- 1 =0 what navy installation support camp david
- 1 +3 how can deforestation directly affect living organisms
- 1 +13 how do active transport and passive transport differ
- 1 =0 when do atoms become excited
- 1 =0 definition do classic
- 1 =0 more queries (click to expand)

Live Demo available at

<https://neural-ir-explorer.ec.tuwien.ac.at/>

② =0 phone number

- 1 +1 texas roadhouse glen mills pa phone number
- 1 =0 the miners state bank routing number
- 1 =0 usf admissions office phone number
- 1 =0 vermont casting group phone number
- 1 =0 dr azadpour phone number
- 1 +1 dr. richard spech npi number
- 1 =0 colorado routing number loveland colorado
- 1 =0 green horizon mini storage contact number
- 1 =0 cox business omaha phone number
- 1 +3 dcu electronic routing number
- + 100 more queries (click to expand)

④ +2 location questions

- 1 =0 hotels in thornton co
- 1 +3 does azusa pacific university negotiate salary
- 1 =0 troy student population
- 1 =0 canada most dense area
- 1 =0 carbon reactivation facilities california
- 1 =0 what is the zip code in arrowhead lakes
- 1 =0 honey in south carolina
- 1 =0 honolulu chinese new year celebration
- 1 =0 which bbc radio station specializes in sports commentaries
- 1 =0 how old is dalton rapattoni
- + 172 more queries (click to expand)

⑨ =0 weather and climate

- 1 +2 weather in greenbelt md

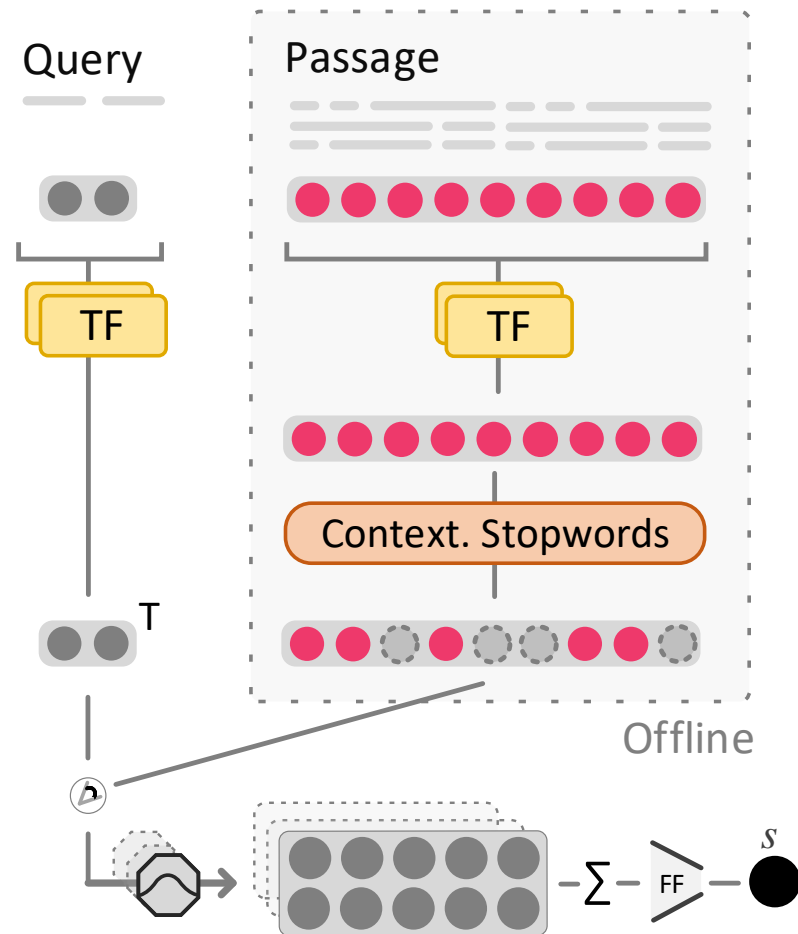
⑭ +8 what is/are 2+ words

- 1 +3 what are caged ibc tanks used for

⑨ +7.5 historical dates & times

- 1 +62 what month does winter start in new zealand

TK-Sparse with Contextualized Stopwords



- The TK model contextualizes query & passage independent
- We offload passage computation
- Stopwords are removed by a trained sparsity module after contextualization
- Sparsity is trained end-to-end with L1 norm augmented loss function

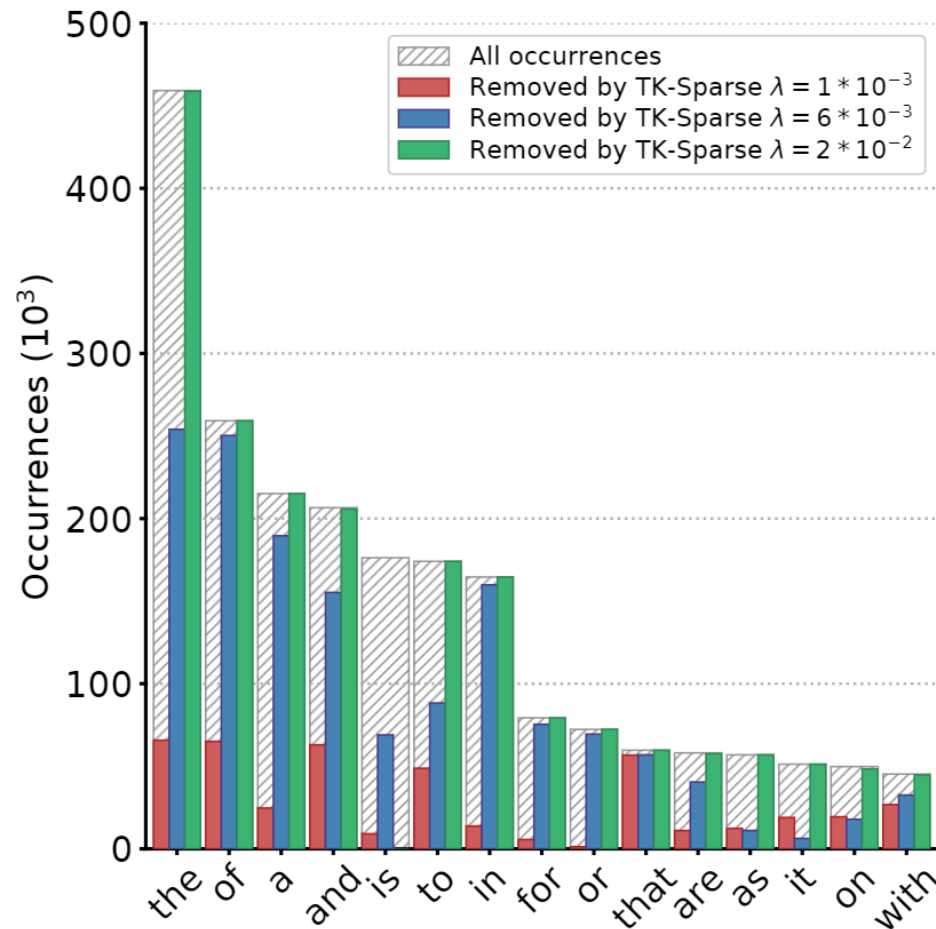
Effectiveness Results (MSMARCO-Passage)

Sig. Model	Stop	DEV (limited judgments)			TREC-2019 (dense judgments)				
		nDCG@10	MRR@10	R@10	nDCG@3	nDCG@10	MRR@10	R@10	MAP@1K
<i>a</i> TK	-	<i>bcdj</i> 0.369	<i>bcdj</i> 0.311	<i>bcdj</i> 0.564	0.655	0.649	0.821	0.242	0.396
<i>b</i> TK /w LuceneStop	24 %	<i>cdj</i> 0.359	<i>cdj</i> 0.301	<i>cdj</i> 0.555	0.661	0.630	<i>dg</i> 0.795	0.220	0.392
<i>c</i> TK /w CollectionTop25	35 %	<i>cdj</i> 0.353	<i>dj</i> 0.296	<i>j</i> 0.547	0.649	0.624	0.756	0.214	<i>d</i> 0.388
<i>d</i> TK /w CollectionTop50	41 %	0.350	0.293	0.543	0.635	0.627	0.745	0.230	0.376
<i>e</i> TK-Sparse $\lambda = 8 * 10^{-4}$	3 %	<i>abcdghij</i> 0.373	<i>abcdhij</i> 0.314	<i>abcdj</i> 0.569	<i>d</i> 0.669	0.638	0.789	<i>c</i> 0.230	0.384
<i>f</i> TK-Sparse $\lambda = 1 * 10^{-3}$	19 %	<i>abcdhj</i> 0.373	<i>abcdhj</i> 0.314	<i>abcdhj</i> 0.570	<i>cd</i> 0.691	<i>bcdeg</i> 0.658	<i>dceg</i> 0.840	<i>c</i> 0.232	<i>bcdegh</i> 0.400
<i>g</i> TK-Sparse $\lambda = 3 * 10^{-3}$	12 %	<i>abcdj</i> 0.372	<i>abcdj</i> 0.314	<i>abcdj</i> 0.568	<i>d</i> 0.665	0.632	0.758	0.220	<i>d</i> 0.382
<i>h</i> TK-Sparse $\lambda = 6 * 10^{-3}$	26 %	<i>abcdj</i> 0.371	<i>bcdj</i> 0.312	<i>bcdj</i> 0.567	0.681	0.653	0.821	0.231	0.391
<i>i</i> TK-Sparse $\lambda = 9 * 10^{-3}$	18 %	<i>abcdhj</i> 0.371	<i>abcdj</i> 0.312	<i>bcdj</i> 0.567	<i>bcdeg</i> 0.714	<i>bcdeg</i> 0.657	<i>dcg</i> 0.827	0.227	<i>bdegh</i> 0.400
<i>j</i> TK-Sparse $\lambda = 2 * 10^{-2}$	43 %	0.350	0.293	0.542	<i>bcdgh</i> 0.705	<i>bcdeg</i> 0.657	<i>dcg</i> 0.832	<i>c</i> 0.239	<i>d</i> 0.395

Effectiveness Results (MSMARCO-Passage)

Sig. Model	Stop	DEV (limited judgments)			TREC-2019 (dense judgments)				
		nDCG@10	MRR@10	R@10	nDCG@3	nDCG@10	MRR@10	R@10	MAP@1K
<i>a</i> TK	-	<i>bcdj</i> 0.369	<i>bcdj</i> 0.311	<i>bcdj</i> 0.564	0.655	0.649	0.821	0.242	0.396
<i>b</i> TK /w LuceneStop	24 %	<i>cdj</i> 0.359	<i>cdj</i> 0.301	<i>cdj</i> 0.555	0.661	0.630	<i>dg</i> 0.795	0.220	0.392
<i>c</i> TK /w CollectionTop25	35 %	<i>cdj</i> 0.353	<i>dj</i> 0.296	<i>j</i> 0.547	0.649	0.624	0.756	0.214	<i>d</i> 0.388
<i>d</i> TK /w CollectionTop50	41 %	0.350	0.293	0.543	0.635	0.627	0.745	0.230	0.376
<i>e</i> TK-Sparse $\lambda = 8 * 10^{-4}$	3 %	<i>abcdghij</i> 0.373	<i>abcdhj</i> 0.314	<i>abcdj</i> 0.569	<i>d</i> 0.669	0.638	0.789	<i>c</i> 0.230	0.384
<i>f</i> TK-Sparse $\lambda = 1 * 10^{-3}$	19 %	<i>abcdhj</i> 0.373	<i>abcdhj</i> 0.314	<i>abcdhj</i> 0.570	<i>cd</i> 0.691	<i>bcdeg</i> 0.658	<i>dceg</i> 0.840	<i>c</i> 0.232	<i>bcdegh</i> 0.400
<i>g</i> TK-Sparse $\lambda = 3 * 10^{-3}$	12 %	<i>abcdj</i> 0.372	<i>abcdj</i> 0.314	<i>abcdj</i> 0.568	<i>d</i> 0.665	0.632	0.758	0.220	<i>d</i> 0.382
<i>h</i> TK-Sparse $\lambda = 6 * 10^{-3}$	26 %	<i>abcdj</i> 0.371	<i>bcdj</i> 0.312	<i>bcdj</i> 0.567	0.681	0.653	0.821	0.231	0.391
<i>i</i> TK-Sparse $\lambda = 9 * 10^{-3}$	18 %	<i>abcdhj</i> 0.371	<i>abcdj</i> 0.312	<i>bcdj</i> 0.567	<i>bcdeg</i> 0.714	<i>bcdeg</i> 0.657	<i>dcg</i> 0.827	0.227	<i>bdegh</i> 0.400
<i>j</i> TK-Sparse $\lambda = 2 * 10^{-2}$	43 %	0.350	0.293	0.542	<i>bcdgh</i> 0.705	<i>bcdeg</i> 0.657	<i>dcg</i> 0.832	<i>c</i> 0.239	<i>d</i> 0.395

Stopword Analysis

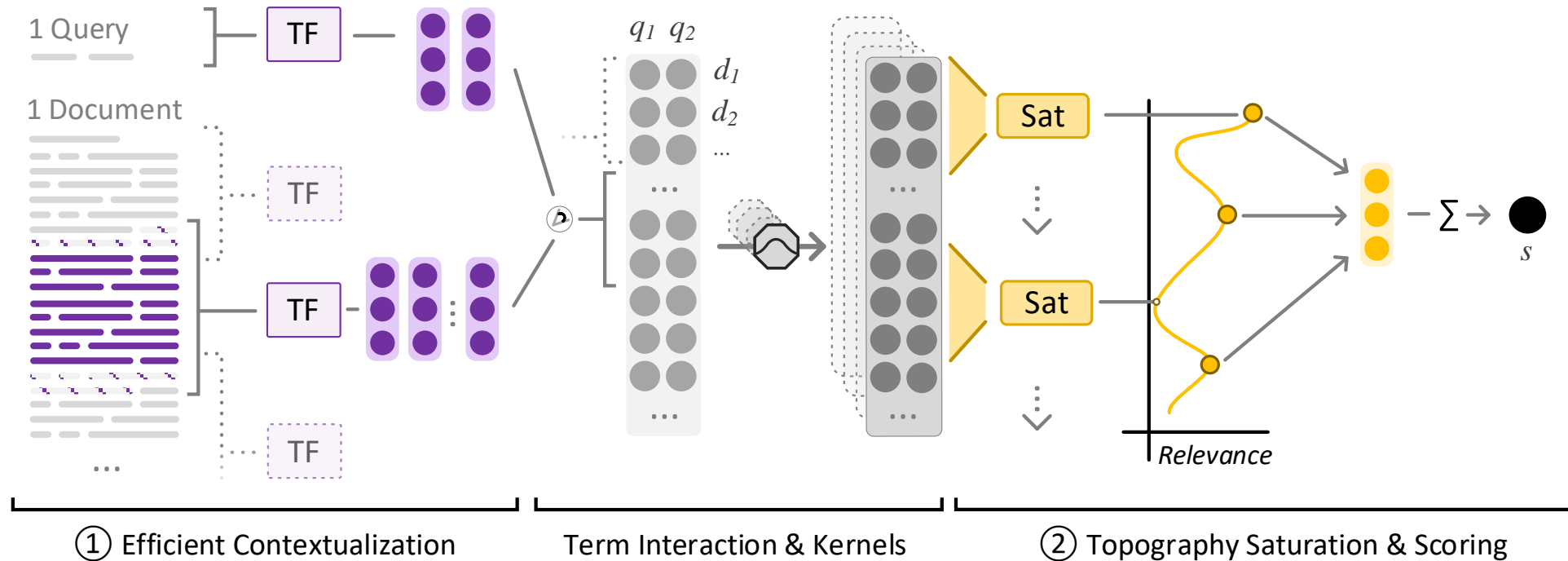


- Common stopwords lists remove every occurrence of a term
 - Contextualized stopwords decide per occurrence
- Here we compare the most common stopwords from Lucene with different configurations of TK-Sparse
 - We observe an overlap of removed terms and substantial differences

TKL: Transformer-Kernel for Long Documents

- Long documents (>200 tokens) are very slow
 - 300-dimensional vector per each word
 - Re-rank 100s of documents per query
 - Padding techniques are insufficient
- State-of-the-art models don't work
 - Do not contain a notion of region importance
 - Current best approach split a document and score sentences/paragraphs individually
- We proposed an extension to TK for Long documents (TKL)

TKL: Transformer-Kernel for Long Documents



TKL: Why long documents?

- We found that longer document input gives us better results
- But only if we do top region detection in TKL
- Main idea behind exercise 1:
 - Find out if the model was correct in this assumption
 - Can we prove that we need longer input

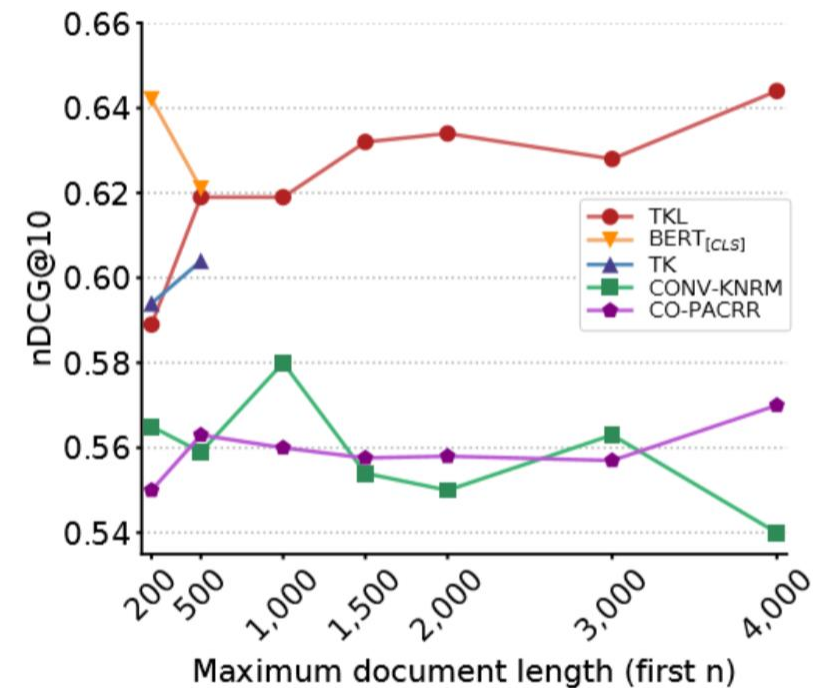
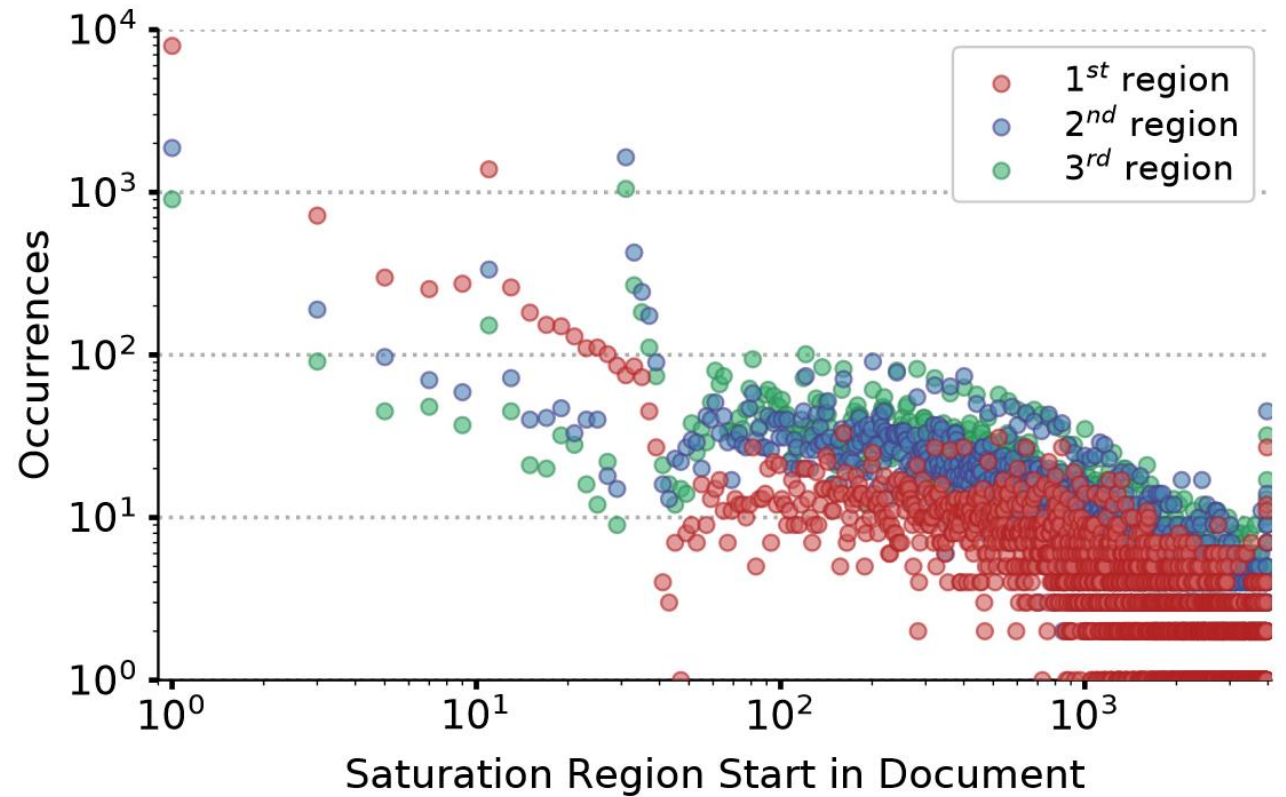


Figure 2: TREC-2019 results based on the document length.

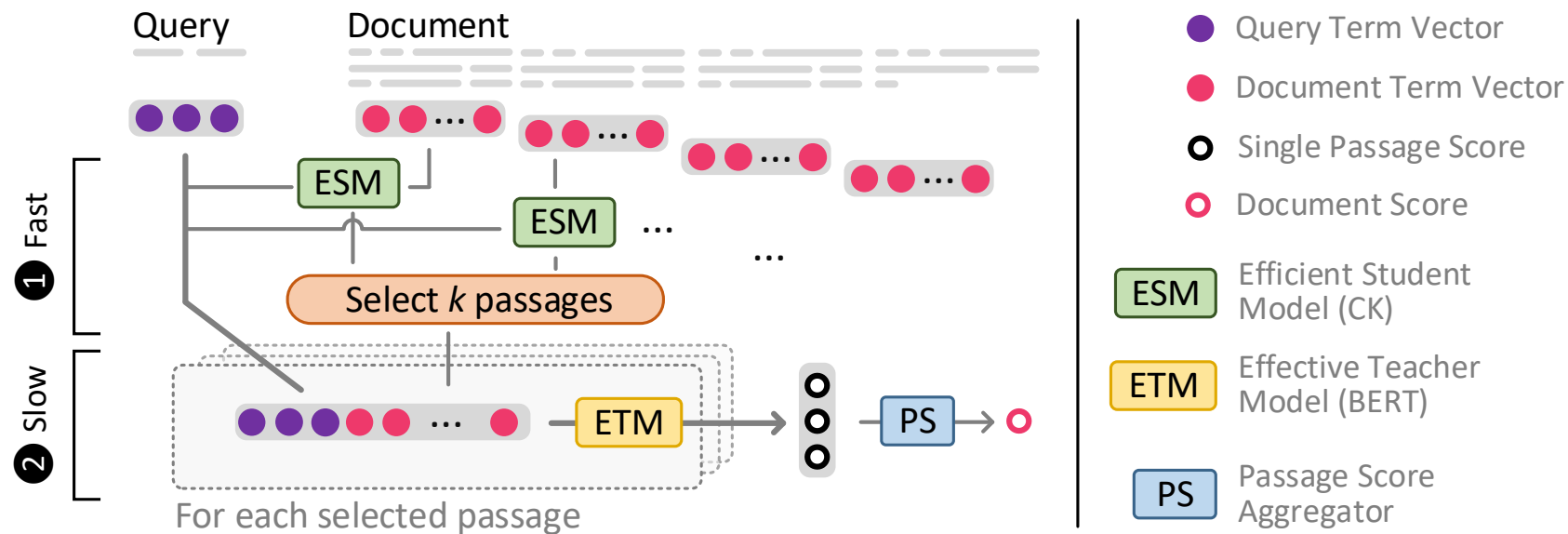
TKL: Where is the relevance?

- TKL provides relevant regions location in the document
- Occurrences follow a Zipfian-Distribution
 - In the TREC-DL 2019 Document collection
- Could be used as snippet generation
 - Better user interfaces

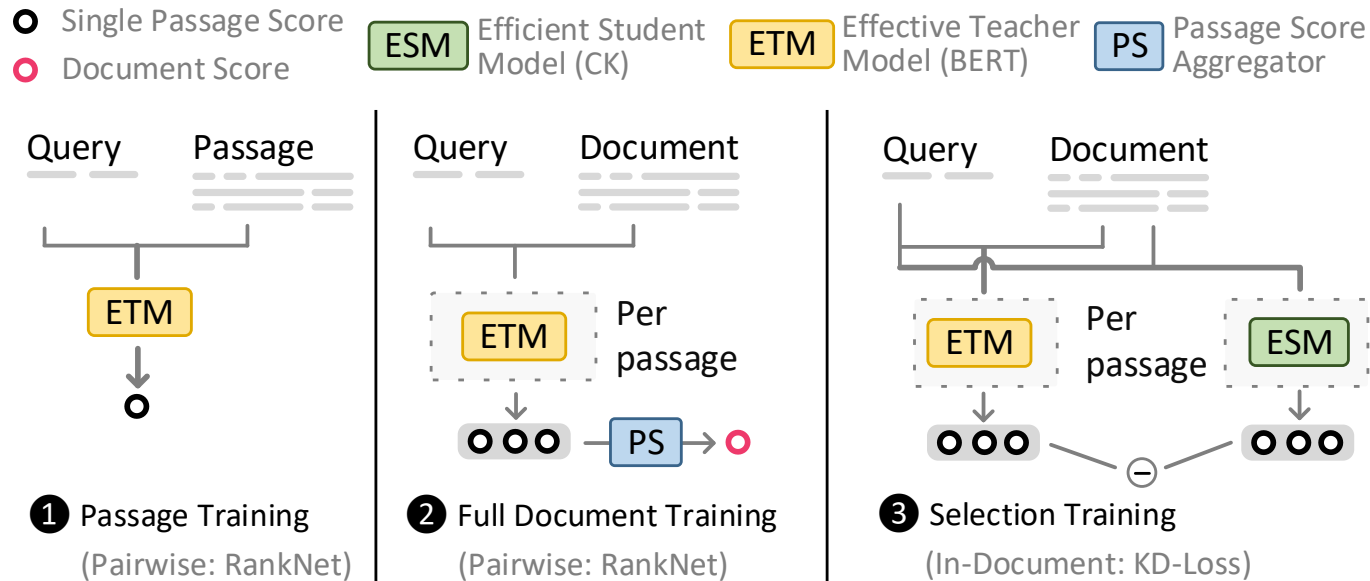


IDCM: A Hybrid Approach

- Combining a slow and fast module in one model to provide an intra-document cascade (IDCM)

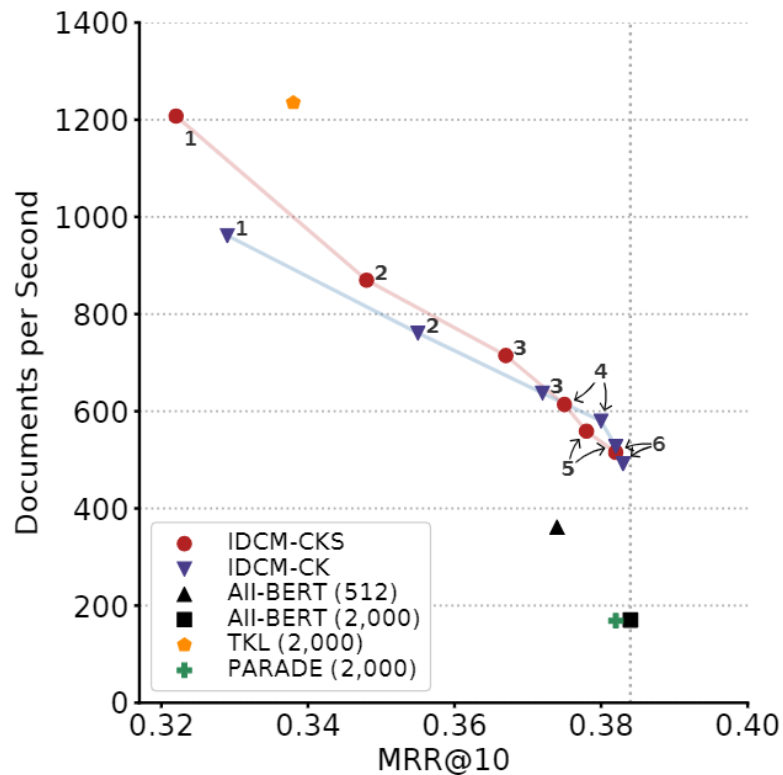


IDCM: Distilled Training



- IDCM is trained in 3 steps:
 - 1. & 2. Prepare slow, but effective BERT model
 - 3. Distill BERT passage-ranks into fast student (CK)
- We want a slow model, that gets the same passages in the top spot, for BERT to fine-grained score
 - Easier the more relevant a document is

IDCM: Improving Throughput



- We can choose different cascade settings:
 - How many passages to score with the costly model
 - The more we score, the better the results, but at the cost of efficiency
- IDCM offers you to choose exactly where you want to be on that curve
 - Ends with up to ~4x lower query latency, but same results

More Resources

- Pretrained Transformers for Text Ranking: BERT and Beyond
Jimmy Lin et al <https://arxiv.org/abs/2010.06467>
 - Great Survey on different techniques using Transformers in Ranking
- ECIR 2021 Tutorial - IR From Bag-of-words to BERT and Beyond through Practical Experiments <https://github.com/terrier-org/ecir2021tutorial>
 - Hands on Tutorial with many Google Colab Notebooks
- Running list of new papers in the field: <https://arxiv.org/list/cs.IR/recent>
 - ~10 new papers a day; the fastest source to get info about pre-prints

Summary: Transformer Contextualized Re-Ranking

- 1 The concatenated BERT_{CAT} opened a new era of information retrieval
- 2 BERT provides enormous effectiveness jumps at the cost of speed
- 3 Combining Transformers and Kernels leads to a good compromise

- ① The concatenated BERTCAT opened a new era of information retrieval
- ② BERT provides enormous effectiveness jumps at the cost of speed
- ③ Combining Transformers and Kernels leads to a good compromise

Thank You