

Dense retrieval using word embeddings

CS6101

Learning word and text

representations

- Unsupervised (AD)
 - Word2vec, GloVe (next)
- Contextualized
 - ConvNet, RNN, Transformer (later)
- Pretrained then fine-tuned

Whole doc → shorter windows

- Whole doc co-occurrence signal is quite noisy
- Instead slide a window through text
 - Center word is w
 - Context words in window are C
- If we hide w , can we predict it from C ?
- Model for cooccurrence of w and $c \in C$
 - Word2vec (Google)
 - GloVe (Stanford)

word2vec

- Doc shortened to sliding window ‘passage’ centered on focus word f
- Cooccurring with other context words c
- For each word w , fit two vectors $\mathbf{u}_w, \mathbf{v}_w$
 - \mathbf{u}_w for focus ‘role’, \mathbf{v}_w for context ‘role’
- Model probability that
 - f and c cooccur as $\Pr(f, c) = \sigma(\mathbf{u}_f \cdot \mathbf{v}_c)$
 - f and a negative sample word \bar{c} do not cooccur as $\Pr(f, \neg \bar{c}) = 1 - \sigma(\mathbf{u}_f \cdot \mathbf{v}_{\bar{c}}) = \sigma(-\mathbf{u}_f \cdot \mathbf{v}_{\bar{c}})$
- Fit U, V to maximize $\prod_f \prod_c \sigma(\mathbf{u}_f \cdot \mathbf{v}_c) \prod_{\bar{c}} \sigma(-\mathbf{u}_f \cdot \mathbf{v}_{\bar{c}})$
- Nonconvex, sensitive to how (many) \bar{c} are sampled per f

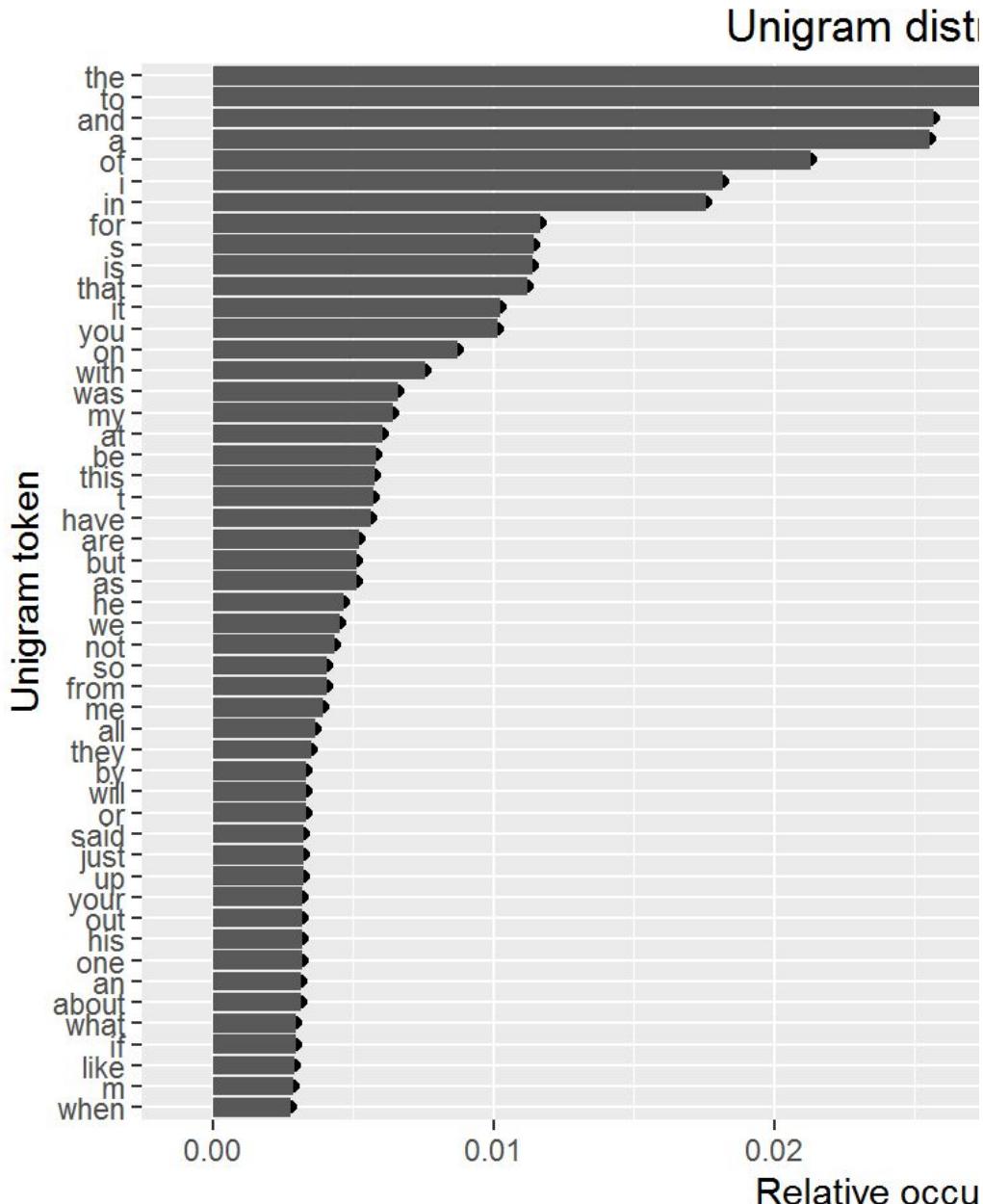
1

is the

There are a few variations on the form of the model, but they show comparable behavior

Negative sampling for word2vec

- For each $\langle f, c \rangle$, collect K samples $\langle f, \bar{c} \rangle$ where \bar{c} does not occur in current context
- Note, negative \bar{c} in one context can be positive c in another
- Natural to sample unigram distrib q_w
- Any f, f' similar to the, to, and, a, etc. therefore $f \approx f'$
- Helps to flatten unigram distribution to $p(w) \propto q_w^{0.75}$ (plus some magic)
- No theory to justify the new gravitational constant!



Batching for word2vec

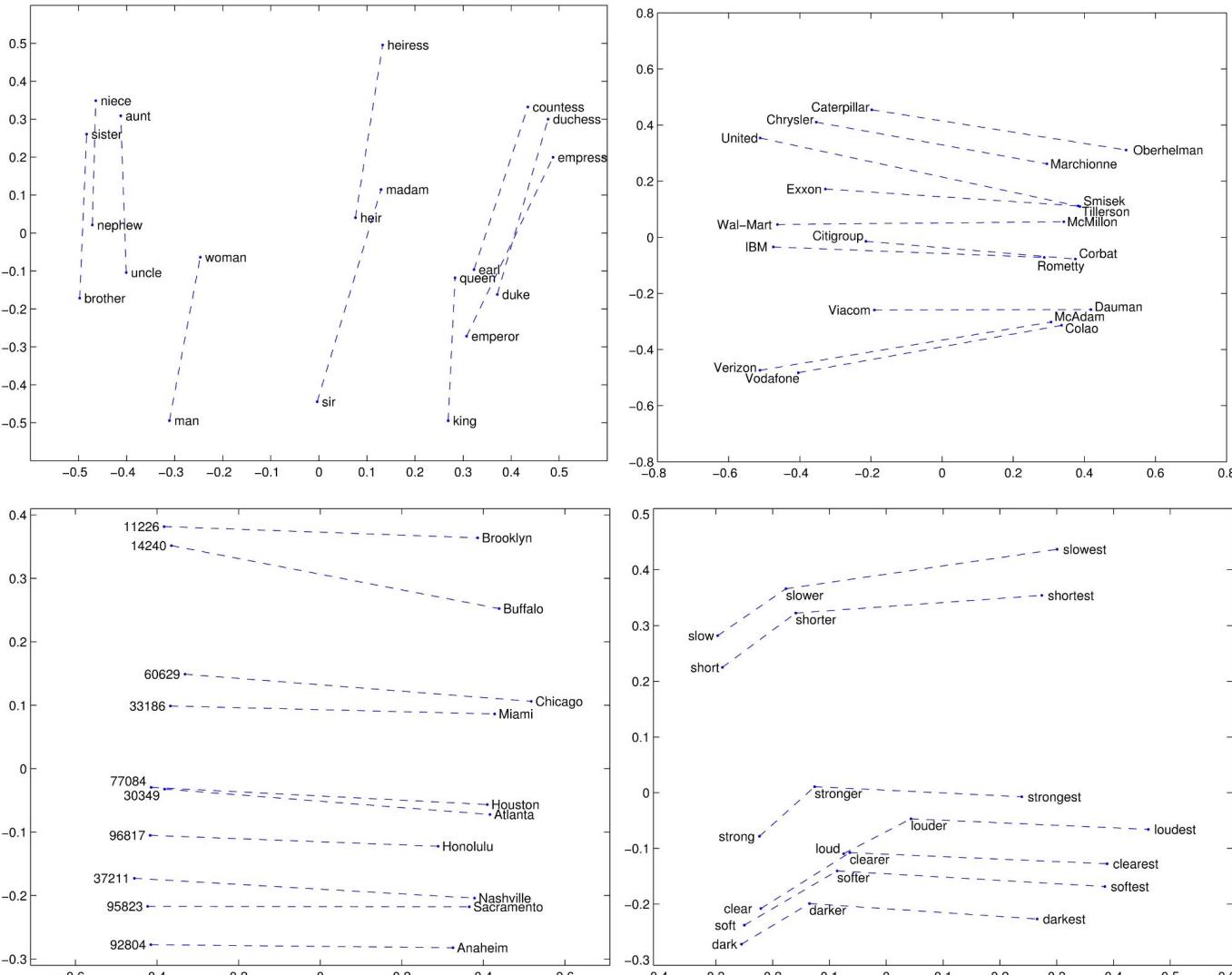
- Repeated window scans through corpus
 - (Not really necessary, as shown in GloVe)
- A batch is collected from a focus word f and several context group, each having
 - One positive context word c
 - Several (K) negative context words \bar{c}
- Notation $\langle f, \{\langle c, \{\bar{c}\} \rangle\} \rangle$; $\{\bar{c}\} = \{\bar{c}_1, \dots, \bar{c}_K\}$
- Find gradient of objective wrt these words, and update with small step size
- Flattening to $\langle f, c, \bar{c} \rangle$ and shuffling destroys optimization!

GloVe

- Collect global word cooccurrence (in window) counts $n(w, w')$
 - Collect top counts (“heavy hitters”) in $o(W^2)$ space, one/few scans over corpus
- Parameter space with four components
 - Compared to word2vec’s two
 - U, V as in word2vec (direction of word)
 - Plus biases $a_w, b_w \in \mathbb{R}$ for each word (marginal popularity)
- Model $\log n(w, w') = \mathbf{u}_w \cdot \mathbf{v}_{w'} + a_w + b_{w'}$, with loss
$$\sum_{w, w'} \text{squash}(n(w, w')) [\mathbf{u}_w \cdot \mathbf{v}_{w'} + a_w + b_{w'} - \log n(w, w')]^2$$
- Parameters fitted by gradient descent
- But no corpus scan or negative sampling during optimization

Surprising property of word vectors

- Word embeddings were known since LSI
- Novelty of word2vec and GloVe partly in pointing out that abstract relations are represented by translations in space
- (Surprising, given objective involved only dot products of unnormalized vectors)
- Analogy games
 - India : Delhi :: Australia : ?
- Also played a role in early knowledge graph embeddings

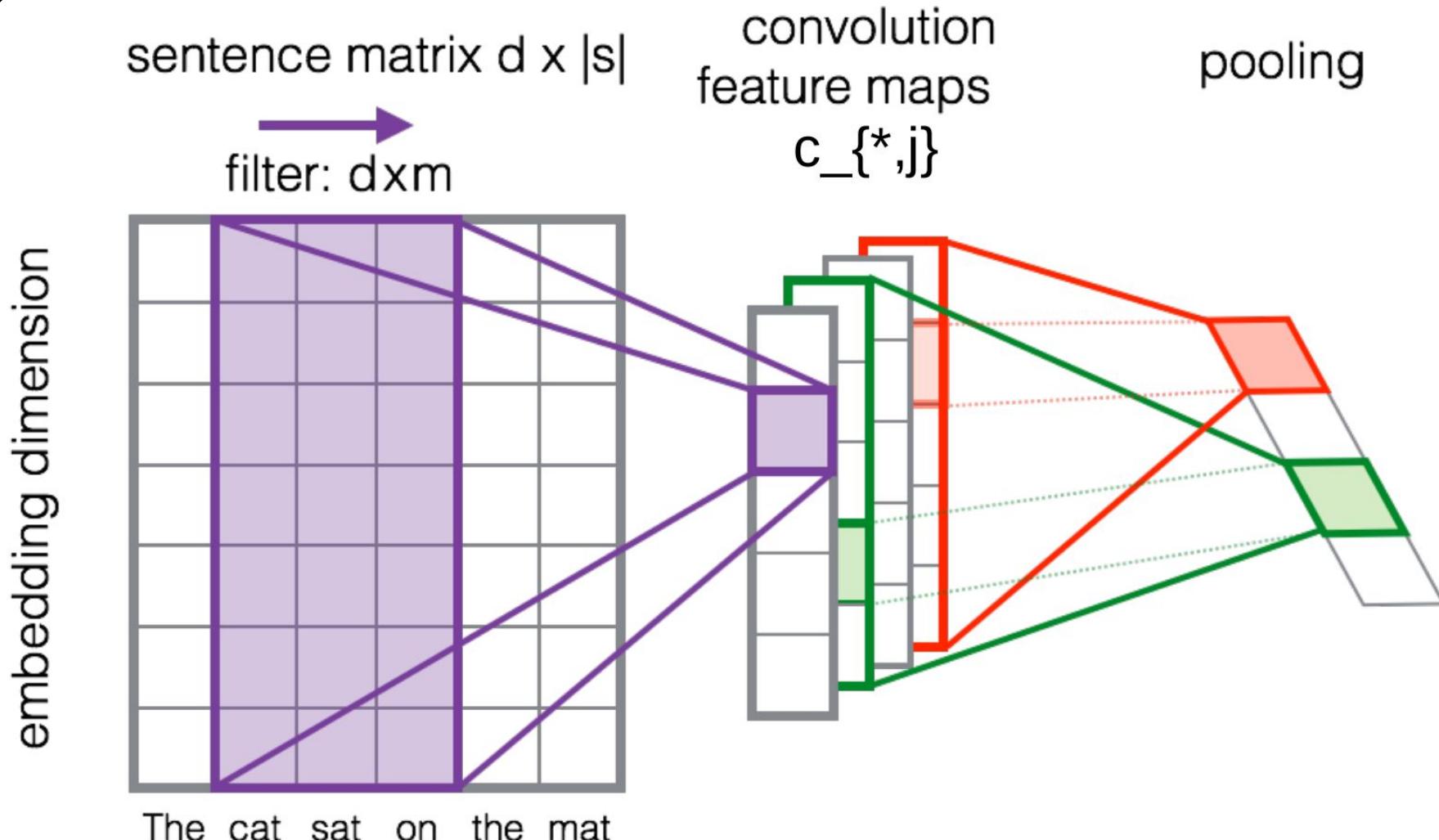


Sentence/passage/doc representation

- LSI and PSI simultaneously computed embeddings for each word and doc
 - New docs have to be “folded in” via EM
 - Cart before horse vs. horse before cart
 - Word embeddings are context independent
- Easiest: weighted average of word vectors
- Sliding window and convolution
- Memory (RNN) and attention (transformers)

1d convnet for text

- Hierarchical position-aware aggregation without parsing ~



1d convnet, multiple channels

- $d \times |s|$ word embedding matrix S for sentence with $|s|$ words, here 7×6
- 7×3 ‘stencil’ slides through four positions
- Weight matrix is $W \in \mathbb{R}^{7 \times 3}$
- When left edge of W is positioned at column k of S , we compute
$$c_k = \sum_{i=0}^6 \sum_{j=0}^2 W[i, j] S[i, j + k]$$
- This gives c , a vector with 4 numbers
- In place of a single W , use ‘channels’ or ‘heads’ $\textcolor{red}{W}, \textcolor{blue}{W}, \textcolor{green}{W}, \textcolor{yellow}{W}, \textcolor{purple}{W}$
- Leading to $\textcolor{red}{c}, \textcolor{blue}{c}, \textcolor{green}{c}, \textcolor{yellow}{c}, \textcolor{purple}{c}$ i.e., matrix $C \in \mathbb{R}^{4 \times 5}$
- Then another ‘layer’ (vector) $u \in \mathbb{R}^5$, leading to output $f = Cu \in \mathbb{R}^4$

Expectation: different channels can be trained to pick up different kinds of signals in text

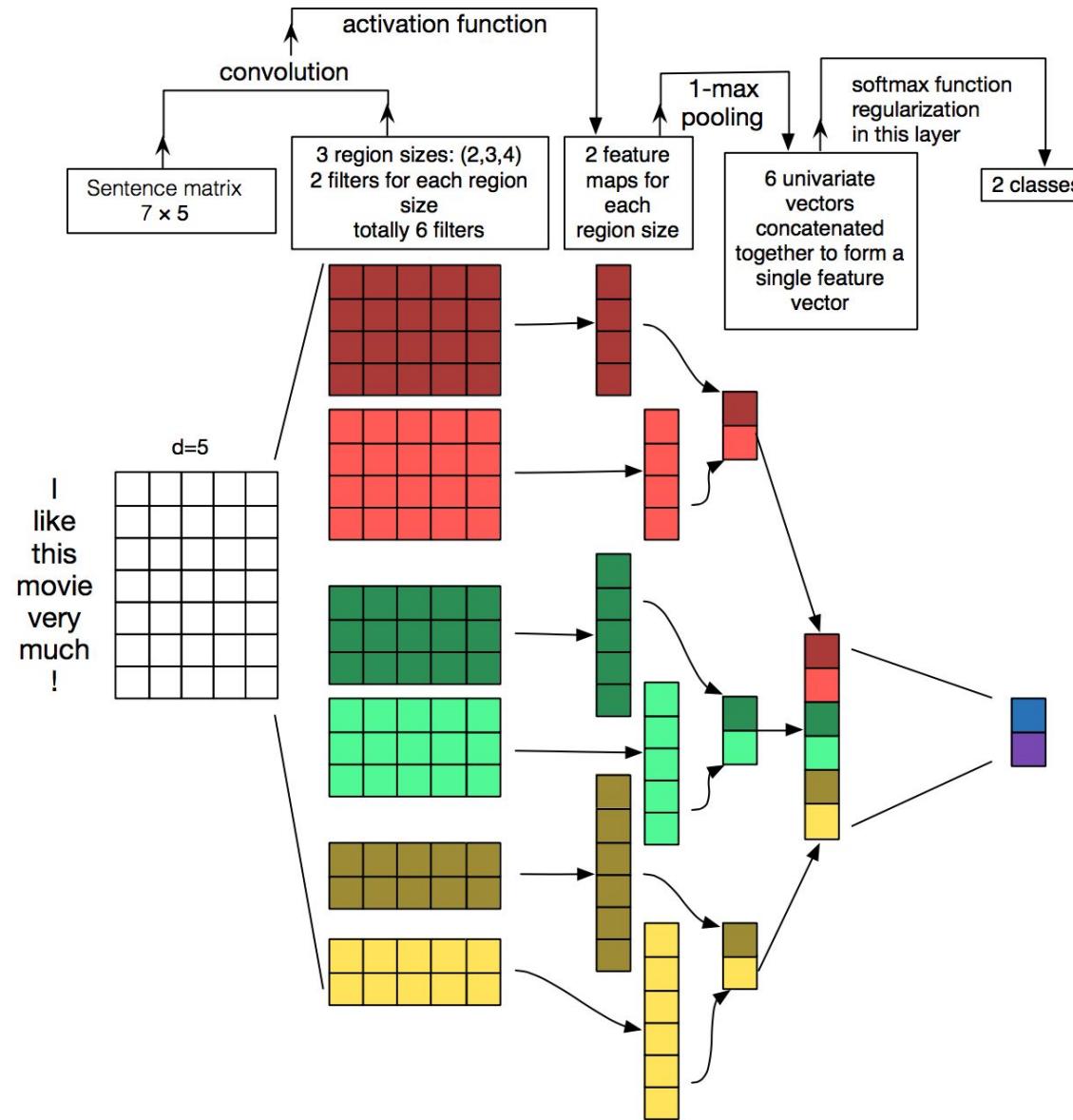
Applications of passage

representation

• (Not just for convnet)

- Learn a classifier whose input is f , to predict labels like topic, sentiment, product rating, spam/not-spam, etc.
- If there are M labels, use matrix $V \in \mathbb{R}^{4 \times M}$; output $\text{softmax}(fV)$
- Compare two passages s_1, s_2 for ‘semantic’ similarity $s_1 \cong s_2$ or entailment $s_1 \Rightarrow s_2$
 - Common use in search
- If unsure of best shape of sliding stencil, use several stencils of various sizes

1d convnet, multiple stencils



Contextual word embeddings

He swam near the bank because the current was swift.

‘Encoder’

His bank offered a low interest on current accounts.

‘Encoder’

Collective ‘meaning’
disambiguation

Two paradigms:

- Left to right influence (RNN, LSTM)
- All to all influence (attention)

HMM, MEMM

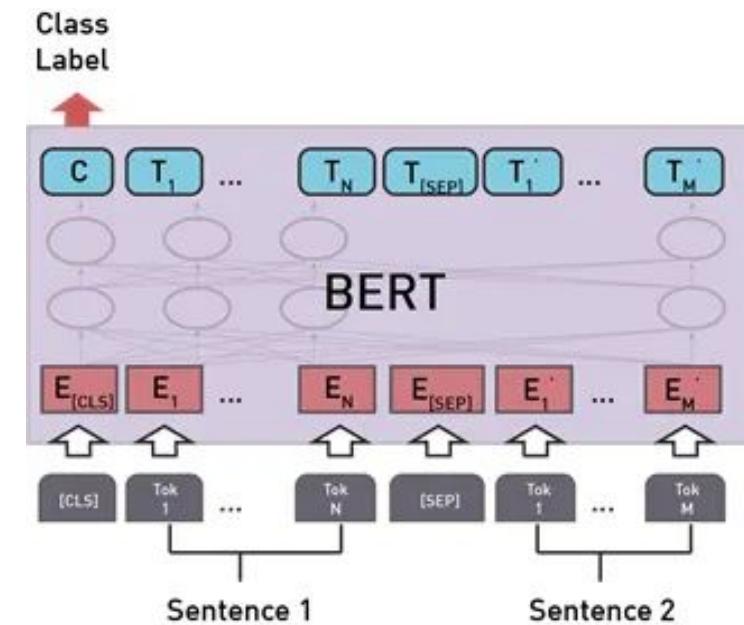
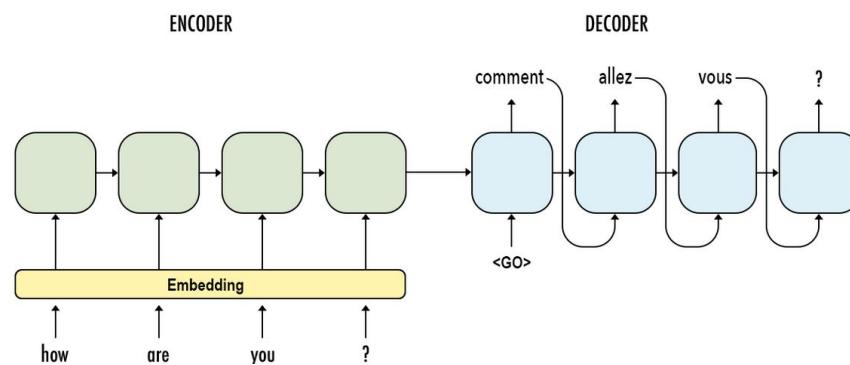
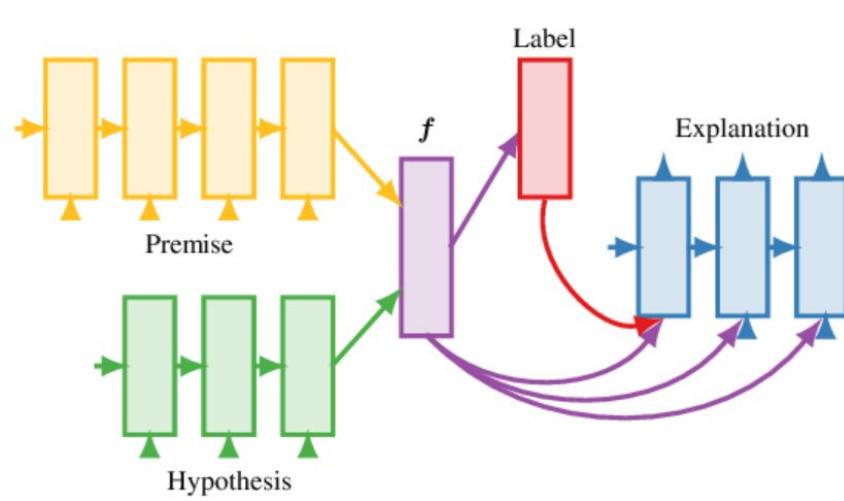
- Observed input symbols $x_{1:T}, x_t \in [F]$
- Inferred output ‘states’ $y_{0:T}$ (sentinel y_0); $y_t \in [M]$
- Classic applications: POS, NE tagging
- Hidden Markov model
 - $\Pr(x_{1:T}, y_{1:T}) = \prod_{t=1}^T \Pr(y_{t-1} \rightarrow y_t) \Pr(y_t \uparrow x_t)$
 - State transition ‘ \rightarrow ’ and emission ‘ \uparrow ’
- Maximum entropy Markov model
 - $\Pr(y_{1:T}|x_{1:T}) = \prod_{t=1}^T \Pr(y_{t-1} \rightarrow y_t | x_t)$
 - Each term modeled as logistic regression
- RNNs are continuous state counterparts to MEMMs

RNNs (GRU, LSTM) and bidirectional RNNs

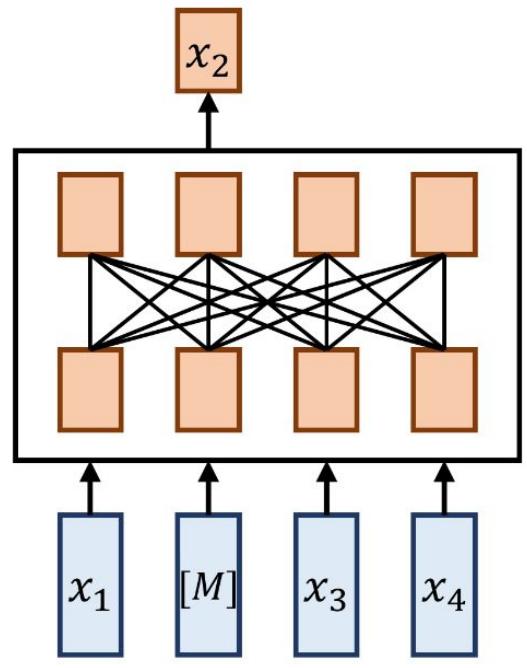
- (Brief review of cell logic)
- In left-to-right RNN, \vec{h}_t is a compressed representation of $x_{1:t}$
- Similarly \overleftarrow{h}_t for right-to-left RNN
- Together, $\vec{h}_t = [\vec{h}_t, \overleftarrow{h}_t]$ give a compressed representation of the whole text tailored to position t
 - Can choose to include or exclude info about x_t itself
- Various tasks can be solved using $\vec{h}_{1:T}$
 - Pooled $\bigoplus_t \vec{h}_t$ can be used for classification; extend to sentence pair labeling
 - Can mask x_t and try to predict it using \vec{h}_t
 - Can give up on right-to-left signals and generate new text

Sentence pair tasks

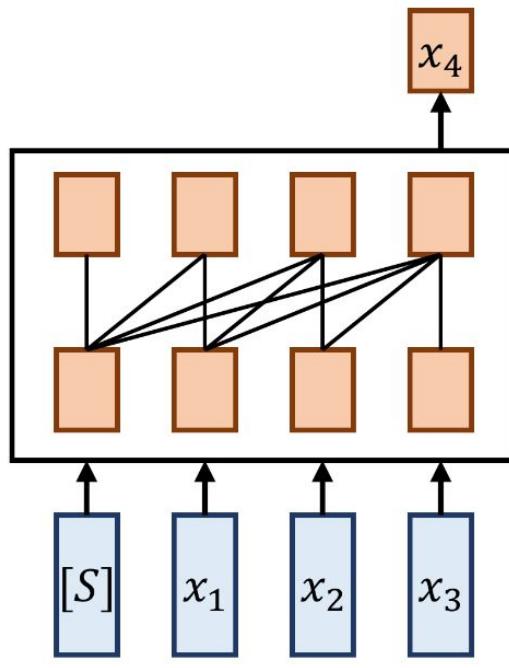
- Pair classification (e.g. implication, next sentence)
- Translation (primed generation)
- Question answering (can be regarded as either classification or generation)



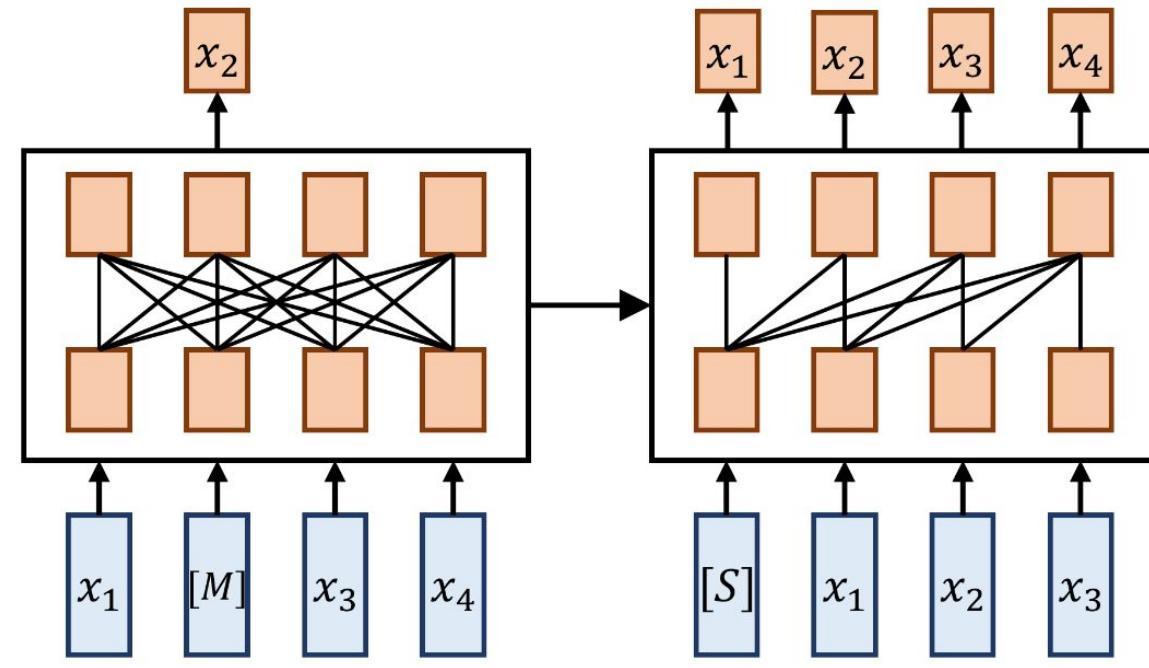
Sequence network nomenclature



Encoder-only



Decoder-only



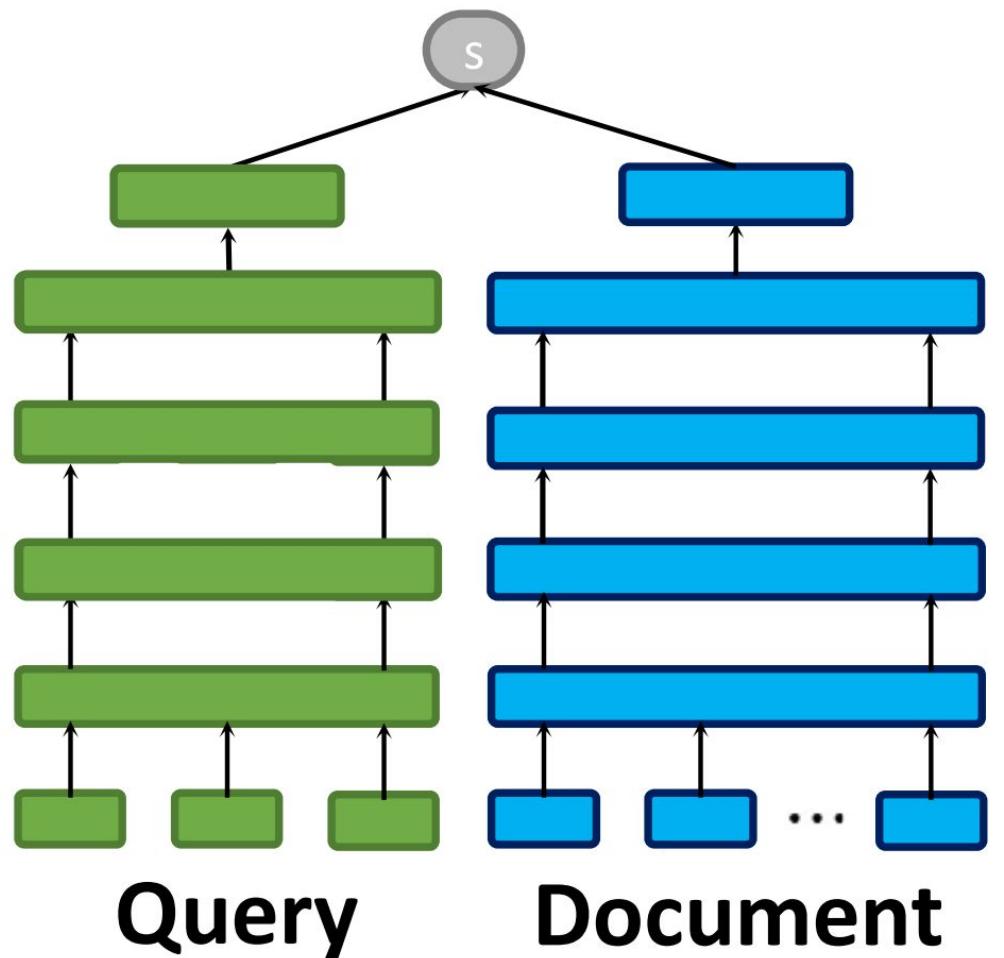
Encoder-decoder

- Classification/tagging, generation, translation
- Encoder can look ahead, decoder cannot

Next....

- Late interaction single vector retrieval
 - LSH for various distance/similarity scores
- Other ANN techniques ([HNSW](#))
- Trainable (dense) retriever
 - aka “dense passage retrieval” or [DPR](#)
- [Late interaction multi-vector retrieval](#)
 - Middle ground between single-vector and early interaction
- [Generative rerankers](#)
- Generative retrievers ([DSI](#))

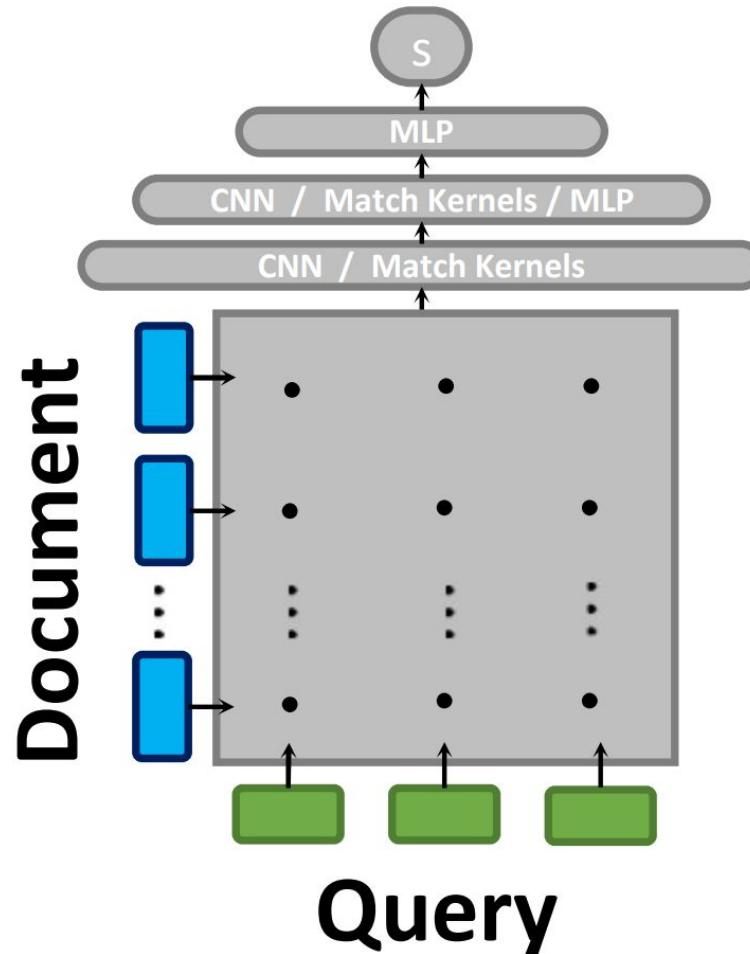
Deep retrieval: late interaction



(a) Representation-based Similarity
(e.g., DSSM, SNRM)

- Too much compression
- + Fast similarity search

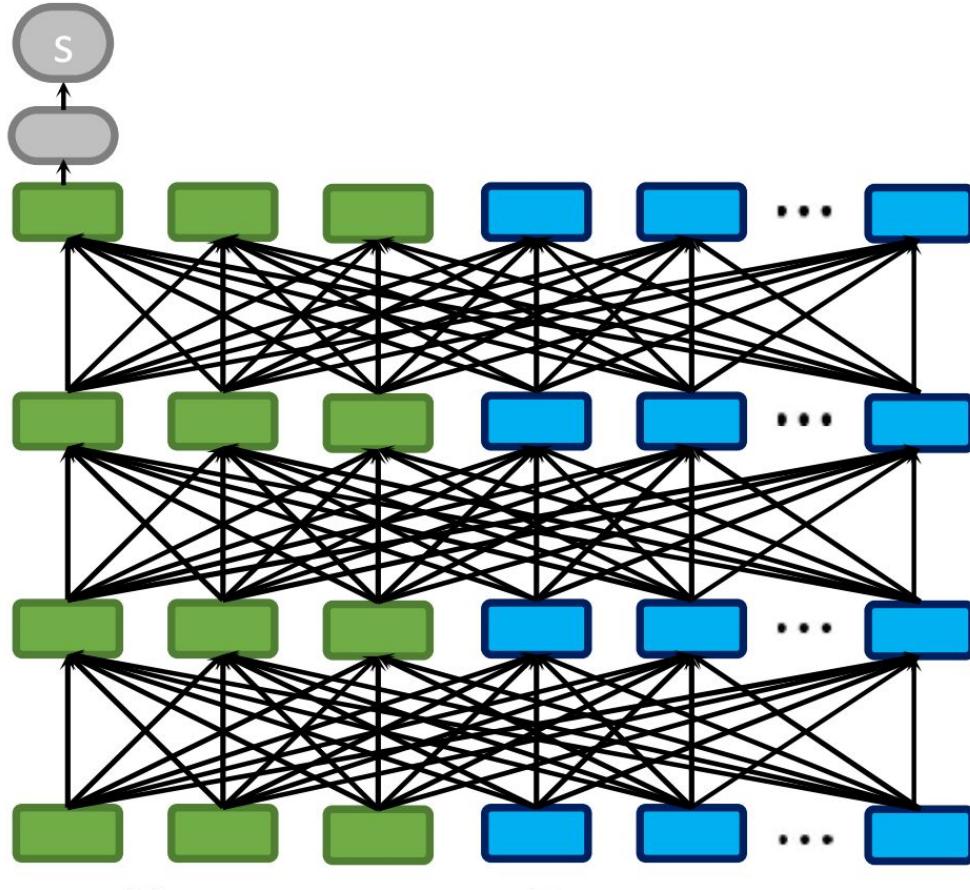
Deep retrieval: crossbar



- All (promising) docs to be checked; slow
- No self-attention

(b) Query-Document Interaction
(e.g., DRMM, KNRM, Conv-KNRM)

All-to-all early interaction

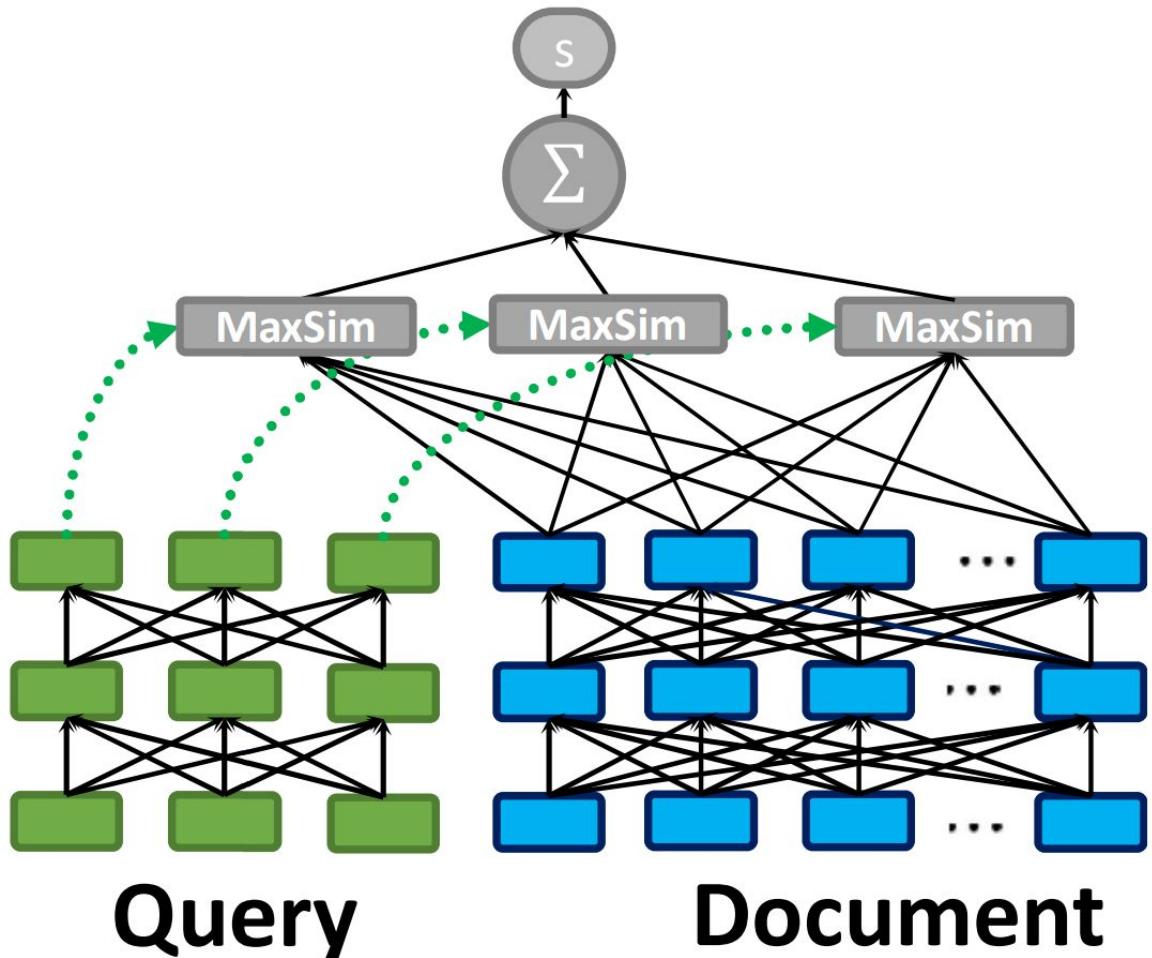


Query Document

(c) All-to-all Interaction
(e.g., BERT)

- Even slower
- Soft vs hard matches?
- Position embeddings, layout tricks

Late, sparse interaction (ColBERT)



(d) Late Interaction
(i.e., the proposed ColBERT)

Separate stacks for q, d

$$E_q := \text{Normalize}(\text{CNN}(\text{BERT}("[Q]q_0q_1\dots q_l\#\#\#")))$$

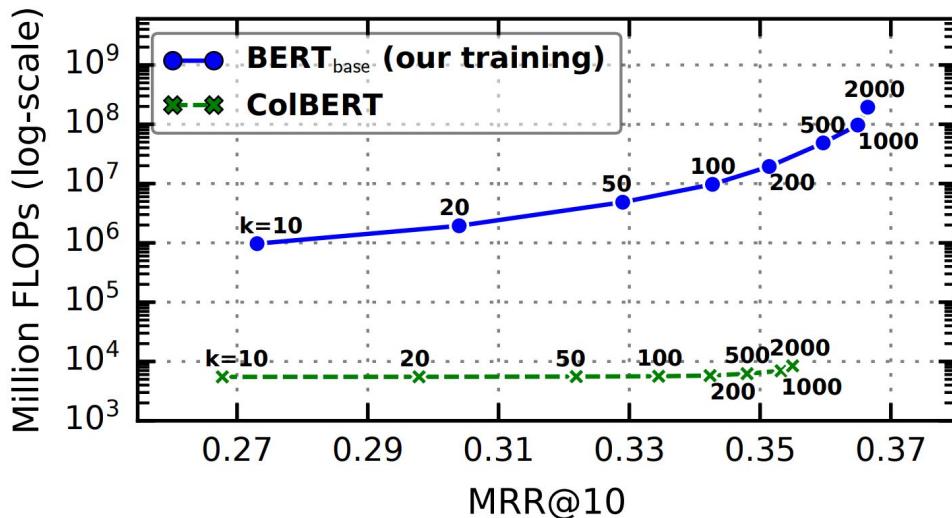
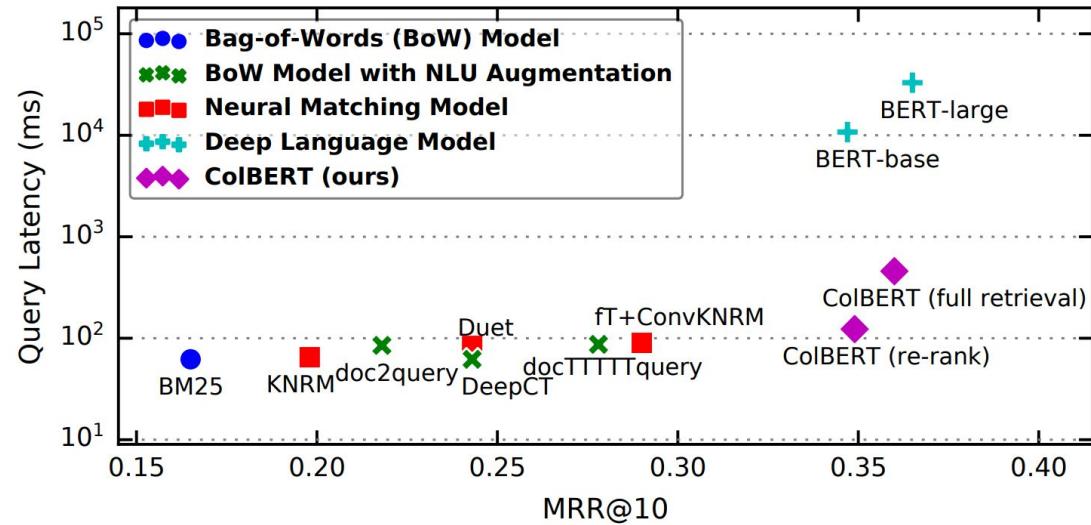
$$E_d := \text{Filter}(\text{Normalize}(\text{CNN}(\text{BERT}("[D]d_0d_1\dots d_n"))))$$

Scoring is not one-shot
but sum over all query
tokens

$$S_{q,d} := \sum_{i \in [|E_q|]} \max_{j \in [|E_d|]} E_{q_i} \cdot E_{d_j}^T$$

Each q token picks best
 d cover

Performance summary



Deep retrieval data set

- <https://microsoft.github.io/msmarco/>
- 100,000 real Bing questions
- Human generated answer
- Another 1,000,000 question dataset
- Tasks
 - Document ranking [leaderboard](#)
 - Passage retrieval [leaderboard](#)
 - Also QA and NL generation

References

- [Huang+2013](#) bag-of-word-embeddings, symmetric aggregation, dot product similarity
- [Shen+2014](#) use convnet instead of bag
- [Palangi+2015](#) use LSTM instead of convnet
- [Gillick+2018](#) avg word vectors; sampling and loss
- [Nogueira+2019](#) use BERT instead of ...
- Bi-encoder vs cross-encoder
- [Karpukhin+2020](#) difficult negatives from BM25
- [Xiong+2020](#) difficult negatives from model itself
- [Humeau+2019](#) [Khattab+2020](#) address cross-encoder inefficiency