

Question Answering Systems

Reading Comprehension and Open-domain QA

Rishiraj Saha Roy

Max Planck Institute for Informatics, Germany

8

Question of the day

How do we (generally) design question answering systems over text (today)?

You'll find this covered in

- 1 ■ Reading Wikipedia to Answer Open-Domain Questions
 - Chen et al. *DrQA*
 - ACL 2017
 - <https://www.aclweb.org/anthology/P17-1171.pdf> ✓
- 2 ■ Simple and Effective Multi-Paragraph Reading Comprehension
 - Clark and Gardner *Document QA*
 - ACL 2018
 - ■ <https://www.aclweb.org/anthology/P18-1078.pdf>

Research paper 1

Reading Wikipedia to Answer Open-Domain Questions

[Reading wikipedia to answer open-domain questions](#) [PDF] [arxiv.org](#)

[D Chen](#), [A Fisch](#), [J Weston](#), [A Bordes](#) - arXiv preprint

[arXiv:1704.00051](#), [2017](#) - [arxiv.org](#)

This paper proposes to tackle open-domain question answering using Wikipedia as the unique knowledge source: the answer to any factoid question is a text span in a Wikipedia article. This task of machine reading at scale combines the challenges of document retrieval (finding the relevant articles) with that of machine comprehension of text (identifying the answer spans from those articles). Our approach combines a search component based on bigram hashing and TF-IDF matching with a multi-layer recurrent neural network model ...

☆ [🔖](#) [Cited by 615](#) [Related articles](#) [All 20 versions](#) [↔](#)

Open-domain question answering

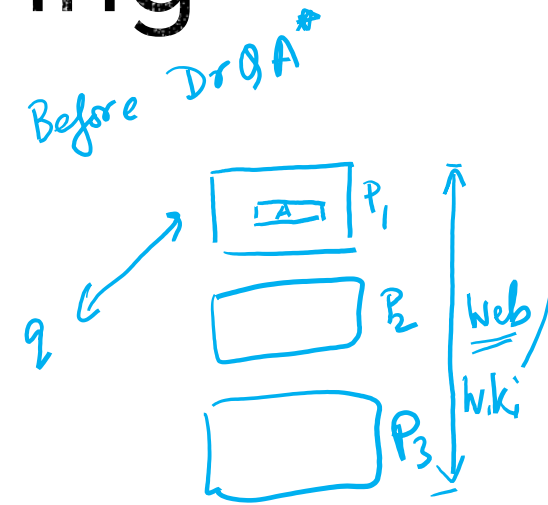
Can we **scale up** neural reading comprehension systems and **generalize** them to open-domain?



DrQA

(Chen et al, 2017)

An open-domain QA system which uses the full English Wikipedia as the knowledge source



from where
do we get
these P_i 's?
AS/NLP \rightarrow IR

Thanks to Danqi Chen
for the slide

DrQA: An open-domain QA system

Q: How many of Warsaw's inhabitants spoke Polish in 1933?

mr c
TBSK?



①
Document
Retriever

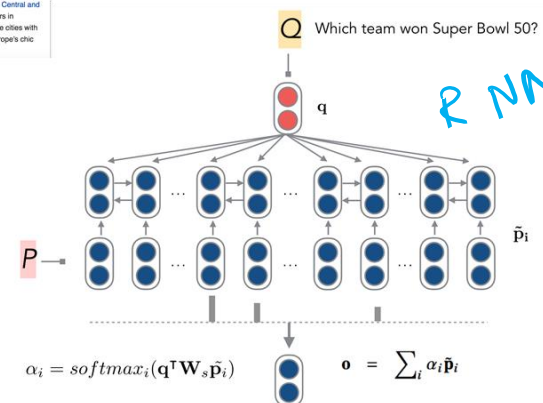


②
Document
Reader

833,500



DrQA = Information Retrieval +



Thanks to Danqi Chen
for the slide

DrQA: What's new?

- Before DrQA: MRC \leftarrow Text QA
- MRC is cool 😊 AI/NLU
- But how effective is it?
- Can we make the goal more realistic / general?
- DrQA: IR (DR) + MRC = MRS \rightarrow Open domain QA
- Previous open systems: Watson, YodaQA, QuASE, AskMSR
- Learning from Wikipedia: Low-redundancy (analogous to KGs)

Machine Reading Comprehension
MRC/MC/RC

(in practice)

Web
 \downarrow
redundancy of evidence
 \downarrow
aggregate

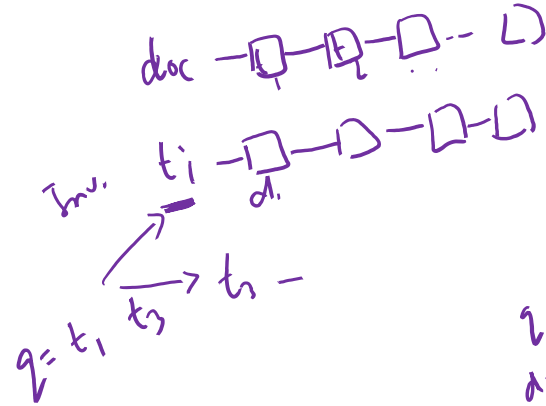
strict

DrQA: In a nutshell

- Document retriever *IR*
 - Bigram hashing and TF-IDF scoring
- Main* | ■ Document reader
 - Multi-layer RNN
- Retriever outperforms Wikipedia search (ElasticSearch)
- | ■ Reader comparable to state-of-the-art on SQuAD *RC*
- Why was it so successful? Visit: *popular* <https://github.com/facebookresearch/DrQA>

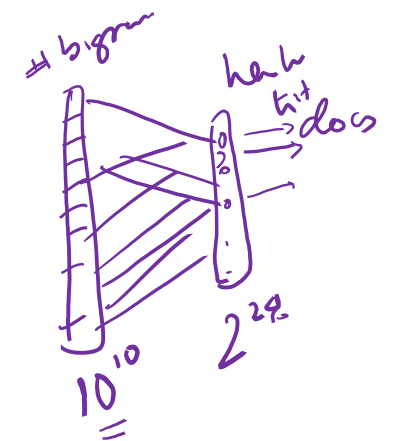
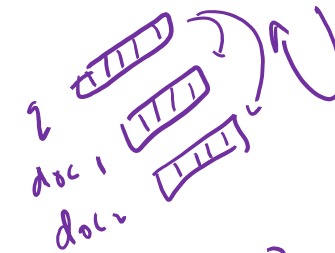
Document retriever

- Unsupervised approach
- Inverted index lookup
- Local word order with n-gram features
- Use bigram counts
- Hashing to map bigrams to 2^{24} bins with murmur3 hash [Recall Bloom filters? ☺]



↓ 2-gram

who is the wife of donald trump?



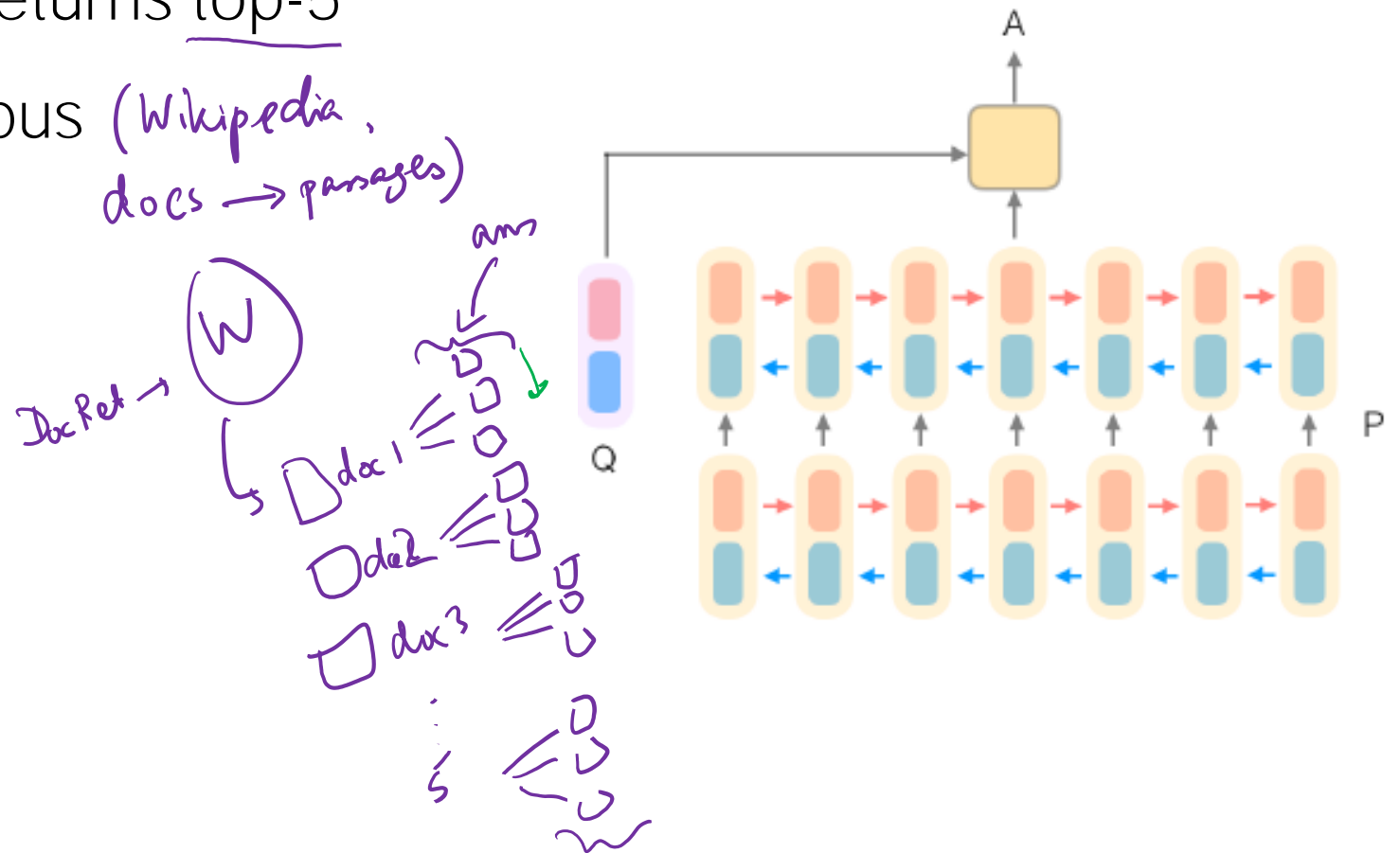
index at bigram level

$$2^{24} \ll 10^{10}$$

unique bigram : too high!!
 $\sim 10^5$
 $\sim 10^{10}$
 $\gg 10^5$

Document reader

- Document retriever returns top-5 ~~paragraphs~~ ^{articles} from corpus (Wikipedia, docs → passages)
- 1 Paragraph encoding
- 2 Question encoding
- 3 Prediction of answer



Reader: Paragraph encoding

- Uses sequence models
- Recurrent neural networks *RNN*
- Specifically, bi-LSTMs *family*

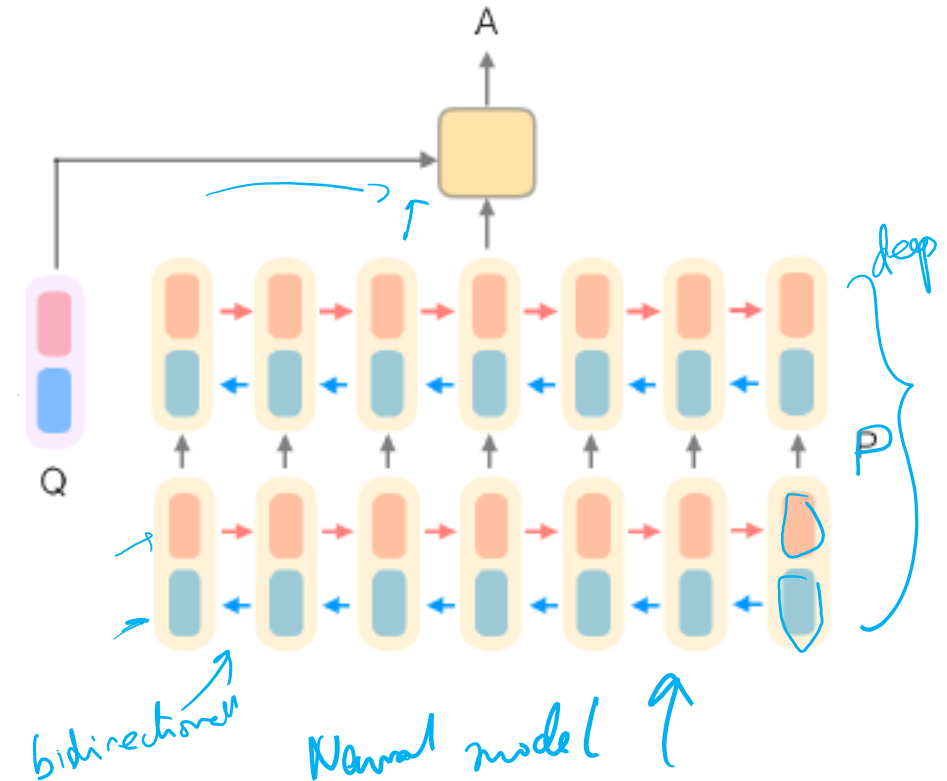
$$\{\vec{p_1}, \dots, \vec{p_m}\} = \text{RNN}(\{\tilde{p_1}, \dots, \tilde{p_m}\})$$

encodings

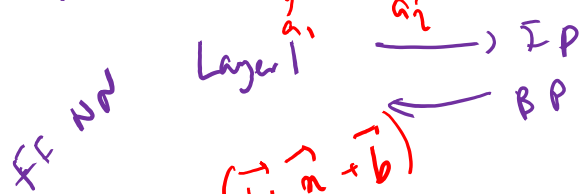
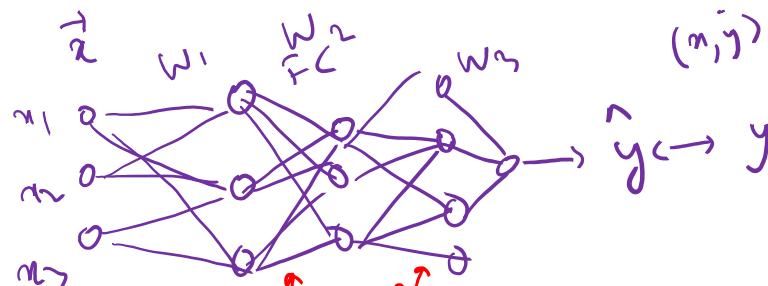
$$Q = \{q_1, q_2, q_3, \dots, q_l\}$$

feature vectors of words

$$P = \{p_1, p_2, p_3, \dots, p_m\}$$



Digression: RNN



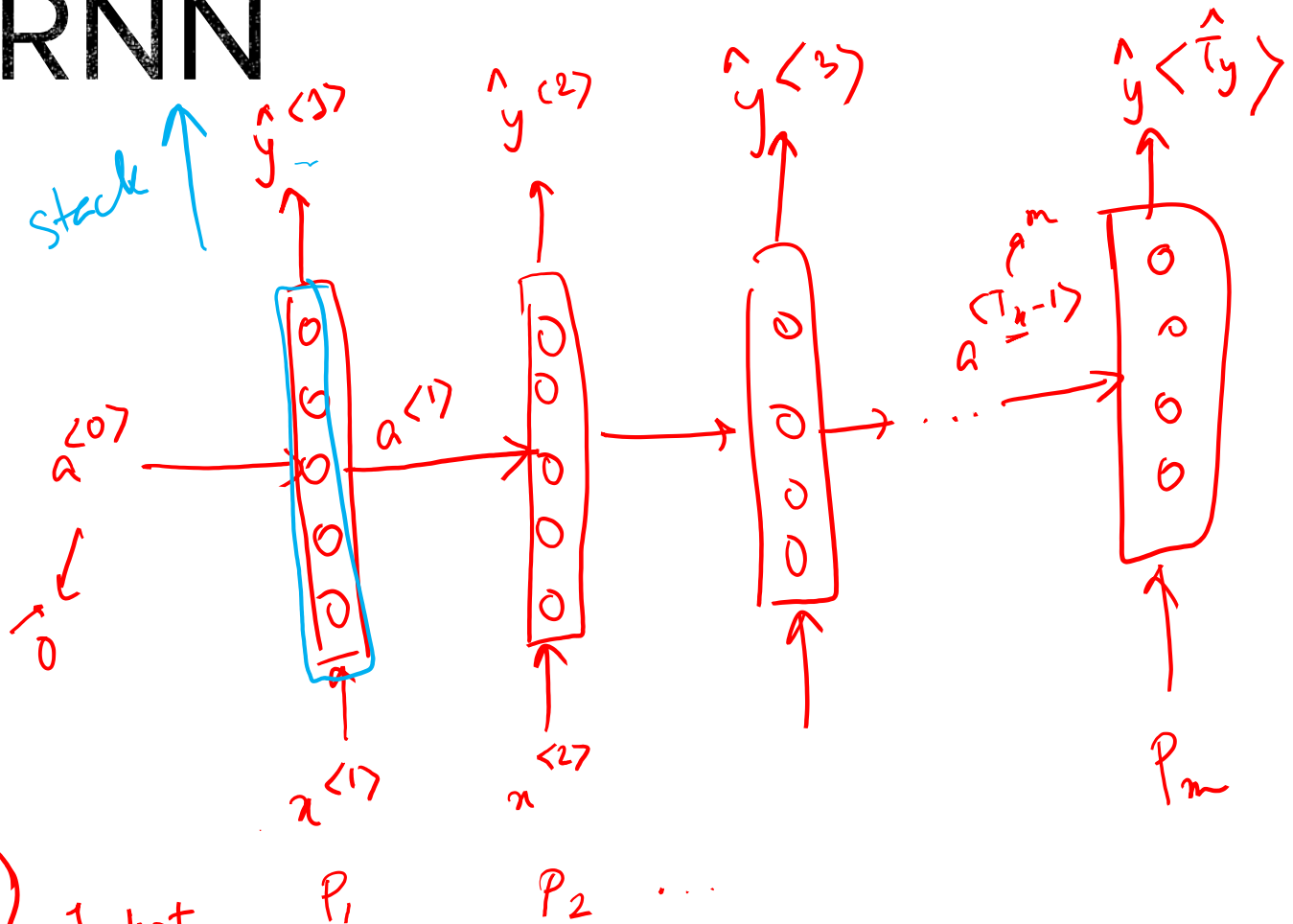
$$\vec{a} = \sigma(\vec{w} \cdot \vec{x} + \vec{b})$$

nonlinear linear

tanh/relu

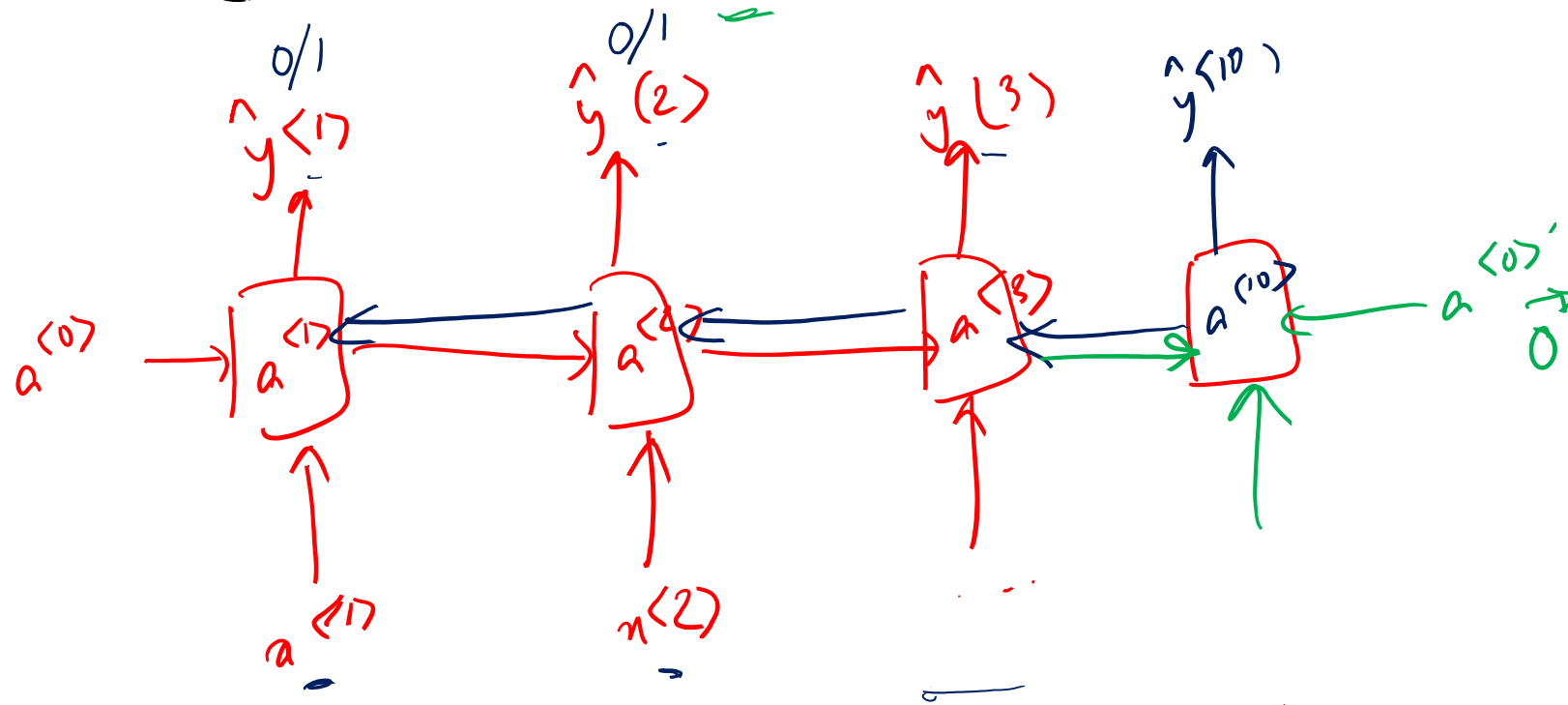
$$a^{(1)} = g(w_{ac}c^{(0)} + w_{an}x^{(1)} + b_a) \quad \text{1-hot}$$

$$\xrightarrow{\text{sigmoid}} y^{(1)} = g_2(w_{ya}a^{(1)} + b_y) \quad \text{or word2vec slope}$$



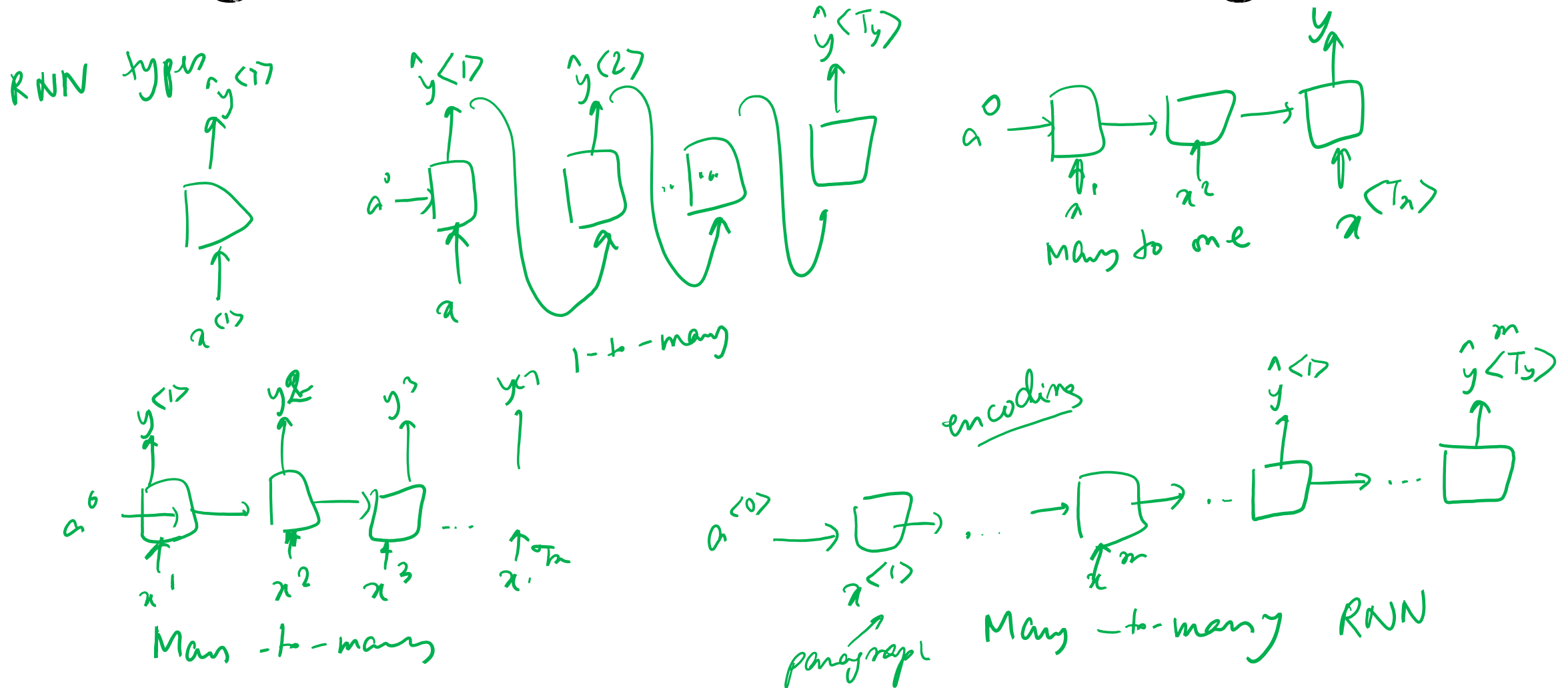
Andrew Ng

Digression: Bidirectional RNN

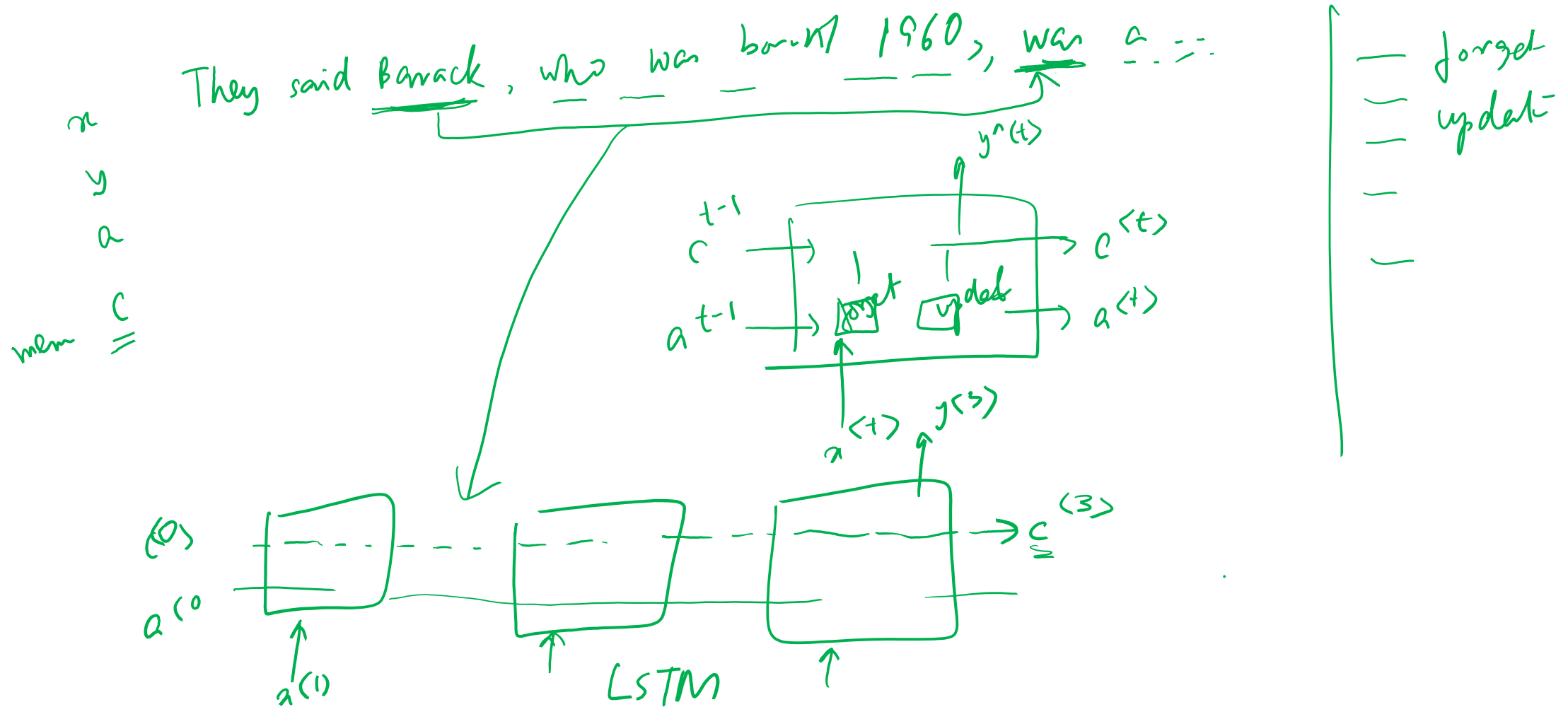


They said that Bonack Obama was a good president.
 NE NE ← Tom Cruise went
 Tom cruised thru L

Digression: RNNs encodings

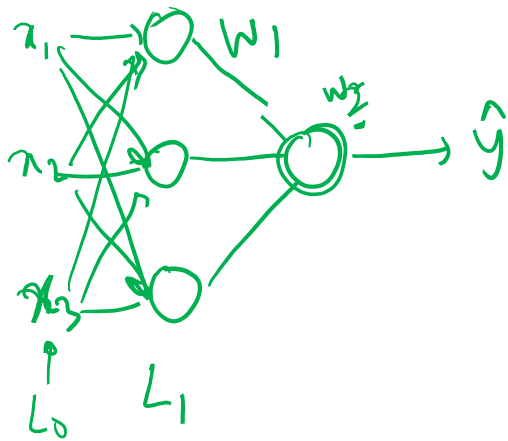


Digression: LSTM



Digression: ReLU

Activation functions



differentiable

$$z^{[1]} = W^{[1]} x + b^{[1]} \text{ \& linear}$$

$$a^{[1]} = g(z^{[1]})$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} / a^{[2]} = g(z^{[2]}) \dots$$

sigmoid $a = \frac{1}{1 + e^{-z}}$

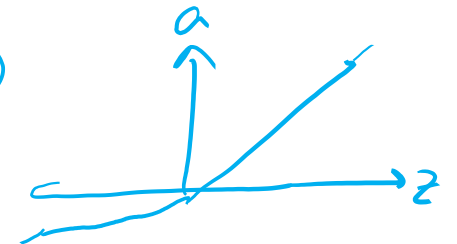
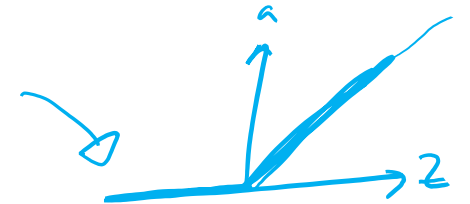
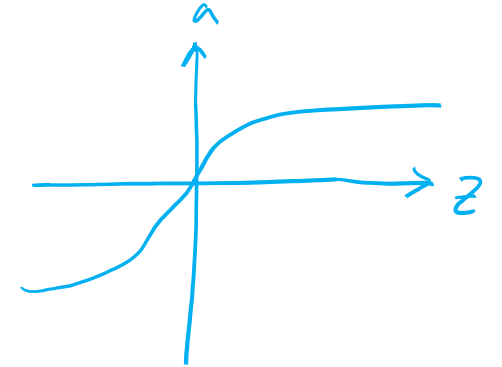
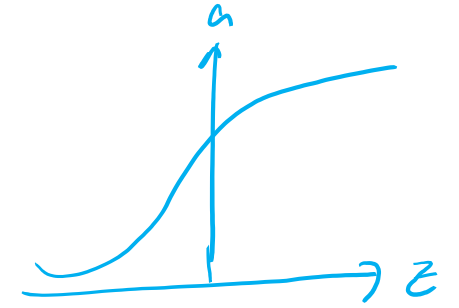
$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

tanh : $a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

$$1 - \tanh^2(z)$$

ReLU : $a = \max(0, z)$

Leaky ReLU : $a = \max(0.01z, z)$



Paragraph encoding: Features

1 ■ Word embeddings *GloVe*

2 ■ Exact match

3 ■ Token features

4 ■ Aligned question features

capture → *Attention*
→ *Synonymy*

$$f_{emb}(\underline{p_i}) = \underline{\mathbf{E}}(\underline{p_i}) \quad \begin{matrix} \text{word} & \text{embeddings} \\ & 300D \end{matrix}$$

$$f_{exact_match}(\underline{p_i}) = \mathbb{I}(\underline{p_i} \in \underline{q}) \quad \begin{matrix} 0/1 \end{matrix}$$

$$f_{token}(\underline{p_i}) = (\overset{\text{tag}}{\text{POS}}(\underline{p_i}), \overset{\text{tag}}{\text{NER}}(\underline{p_i}), \text{TF}(\underline{p_i}))$$

$$f_{align}(\underline{p_i}) = \sum_j \overset{\text{attention}}{a_{i,j}} \underline{\mathbf{E}}(\underline{q_j})$$

$$\underline{a_{i,j}} = \frac{e^{\underset{\text{ReLU}}{\alpha(\underline{\mathbf{E}}(\underline{p_i})) \cdot \alpha(\underline{\mathbf{E}}(\underline{q_j}))}}}{\sum_{j'} \exp(\alpha(\underline{\mathbf{E}}(\underline{p_i})) \cdot \alpha(\underline{\mathbf{E}}(\underline{q_{j'}})))}$$

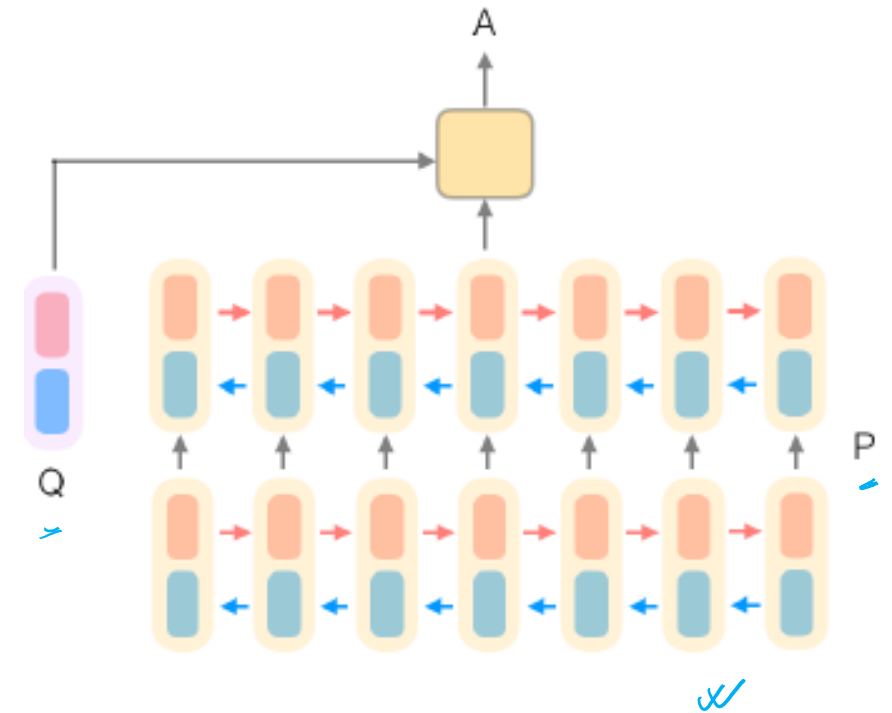
Question encoding

- Another RNN
- Over word embeddings of question tokens

$$\begin{array}{c} \xrightarrow{\quad} \xrightarrow{\quad} \\ \{q_1, q_2, \dots\} \rightarrow \vec{q} \end{array} \quad \left| \quad \vec{q} = \sum_j b_j \vec{q}_j \right.$$

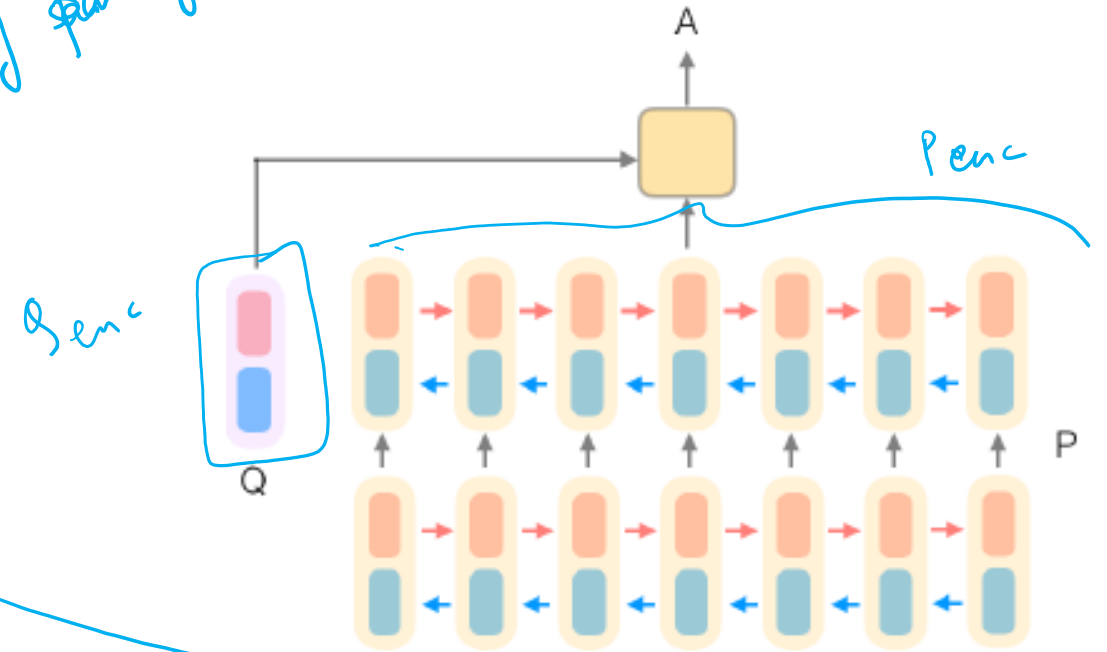
learned

$$b_j = \frac{\exp(\underline{\mathbf{w}} \cdot \underline{\mathbf{q}}_j)}{\sum_{j'} \exp(\underline{\mathbf{w}} \cdot \underline{\mathbf{q}}_{j'})}$$



Answer prediction

- Predict answer span
- Two independent classifiers
- Use bilinear term
- Choose span from indices i to i'
- $i \leq i' \leq i + 15$
- Maximize $P_{start}(i) \times P_{end}(i')$



$$P_{start}(i) \propto \frac{\exp(\mathbf{p}_i \mathbf{W}_s \mathbf{q})}{\exp(\mathbf{p}_i \mathbf{W}_e \mathbf{q})}$$

$$P_{end}(i) \propto \frac{\exp(\mathbf{p}_{i'} \mathbf{W}_s \mathbf{q})}{\exp(\mathbf{p}_{i'} \mathbf{W}_e \mathbf{q})}$$

Annotations:
 - \mathbf{p}_i : start
 - $\mathbf{p}_{i'}$: end
 - \mathbf{q} : max 15 words
 - \mathbf{W}_s : unnormalized exponentiation
 - \mathbf{W}_e : avg max over all
 - \mathbf{q} : score
 - \mathbf{q} : Q_enc

Representative examples

Dataset	Example	Article / Paragraph
SQuAD	Q: How many provinces did the Ottoman empire contain in the 17th century? A: <u>32</u>	Article: Ottoman Empire Paragraph: ... At the beginning of the 17th century the empire contained <u>32</u> provinces and numerous vassal states. Some of these were later absorbed into the Ottoman Empire, while others were granted various types of autonomy during the course of centuries.
CuratedTREC	Q: What U.S. state's motto is "Live free or Die"? A: <u>New Hampshire</u>	Article: Live Free or Die Paragraph: "Live Free or Die" is the official motto of the U.S. state of <u>New Hampshire</u> , adopted by the state in 1945. It is possibly the best-known of all state mottos, partly because it conveys an assertive independence historically found in American political philosophy and partly because of its contrast to the milder sentiments found in other state mottos.
WebQuestions	Q: What part of the atom did Chadwick discover? [†] A: neutron	Article: Atom Paragraph: ... The atomic mass of these isotopes varied by integer amounts, called the whole number rule. The explanation for these different isotopes awaited the discovery of the <u>neutron</u> , an uncharged particle with a mass similar to the proton, by the physicist James Chadwick in 1932. ...
WikiMovies	Q: Who wrote the film Gigli? A: <u>Martin Brest</u>	Article: Gigli Paragraph: Gigli is a 2003 American romantic comedy film written and directed by <u>Martin Brest</u> and starring Ben Affleck, Jennifer Lopez, Justin Bartha, Al Pacino, Christopher Walken, and Lainie Kazan.

Training

Dr. QA
vs
Yoda QA

SQuAD
Web Questions

Research paper 2

Simple and Effective Multi-Paragraph Reading Comprehension

Simple and effective multi-paragraph reading
comprehension

[\[PDF\] arxiv.org](#)

[C Clark, M Gardner](#) - arXiv preprint arXiv:1710.10723, 2017 - arxiv.org

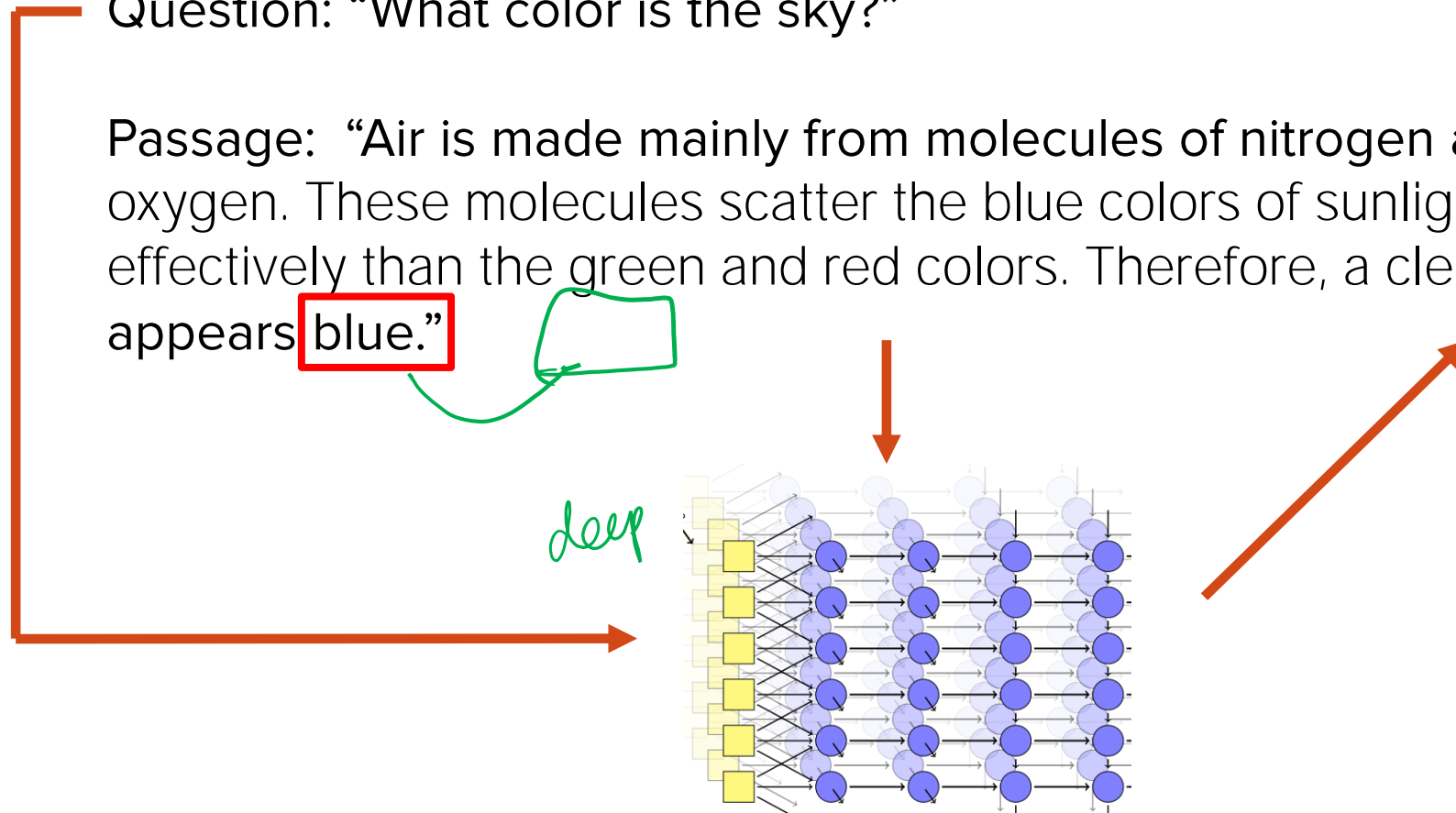
We consider the problem of adapting neural paragraph-level question answering models to the case where entire documents are given as input. Our proposed solution trains models to produce well calibrated confidence scores for their results on individual paragraphs. We sample multiple paragraphs from the documents during training, and use a shared-normalization training objective that encourages the model to produce globally correct output. We combine this method with a state-of-the-art pipeline for training models on ...

☆ [🔖](#) Cited by 198 [Related articles](#) [All 5 versions](#) [↗](#)

Neural question answering

Question: “What color is the sky?”

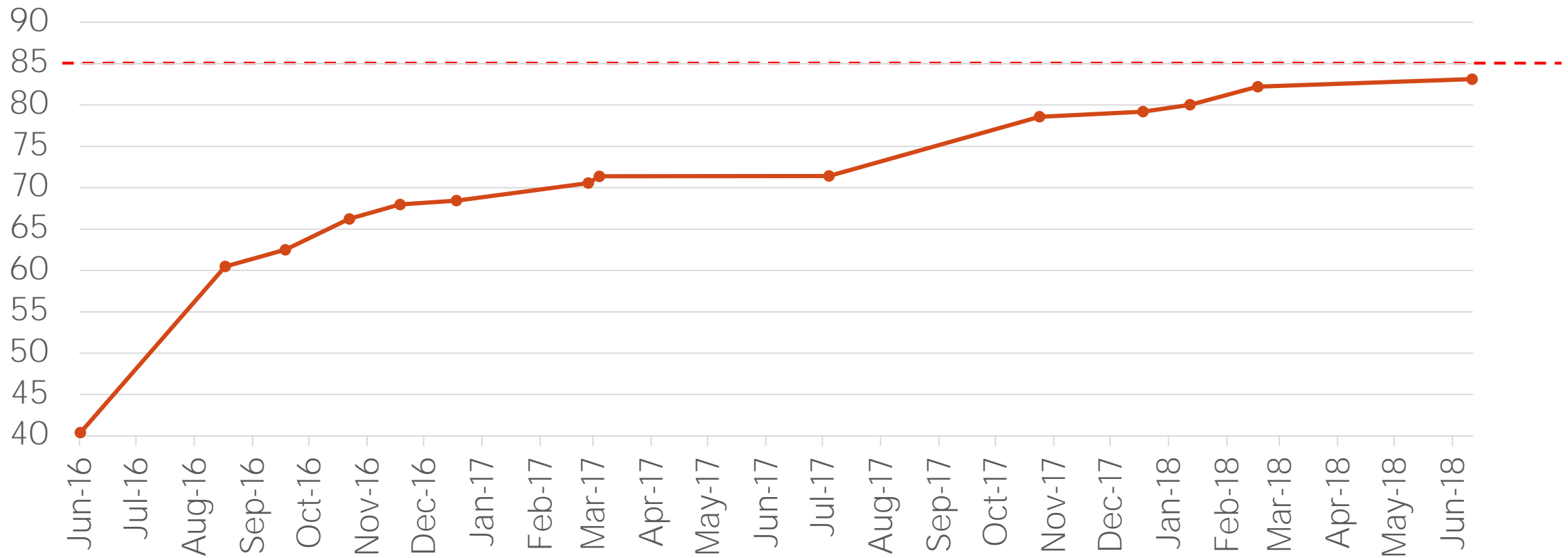
Passage: “Air is made mainly from molecules of nitrogen and oxygen. These molecules scatter the blue colors of sunlight more effectively than the green and red colors. Therefore, a clean sky appears **blue.**”



Thanks to Christopher Clark
for the slides

Fast progress on paragraph datasets

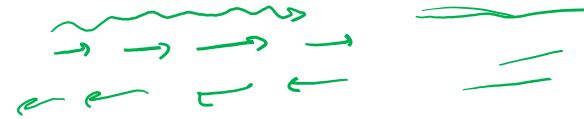
Accuracy on SQuAD 1.1



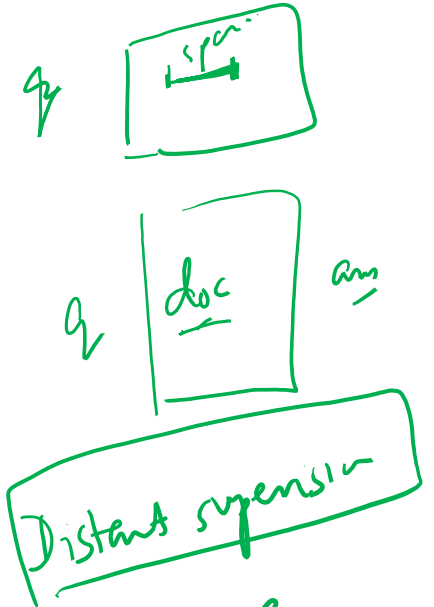
Challenge: Scaling Models to Documents

2018 July

- Modern reading comprehension models have many layers and parameters
- The trend is continuing in this direction, for example with the use of large language models
- Reduced efficiency as the paragraph length increases due to long RNN chains or transformers/self-attention modules
- Limits the model to processing short paragraphs



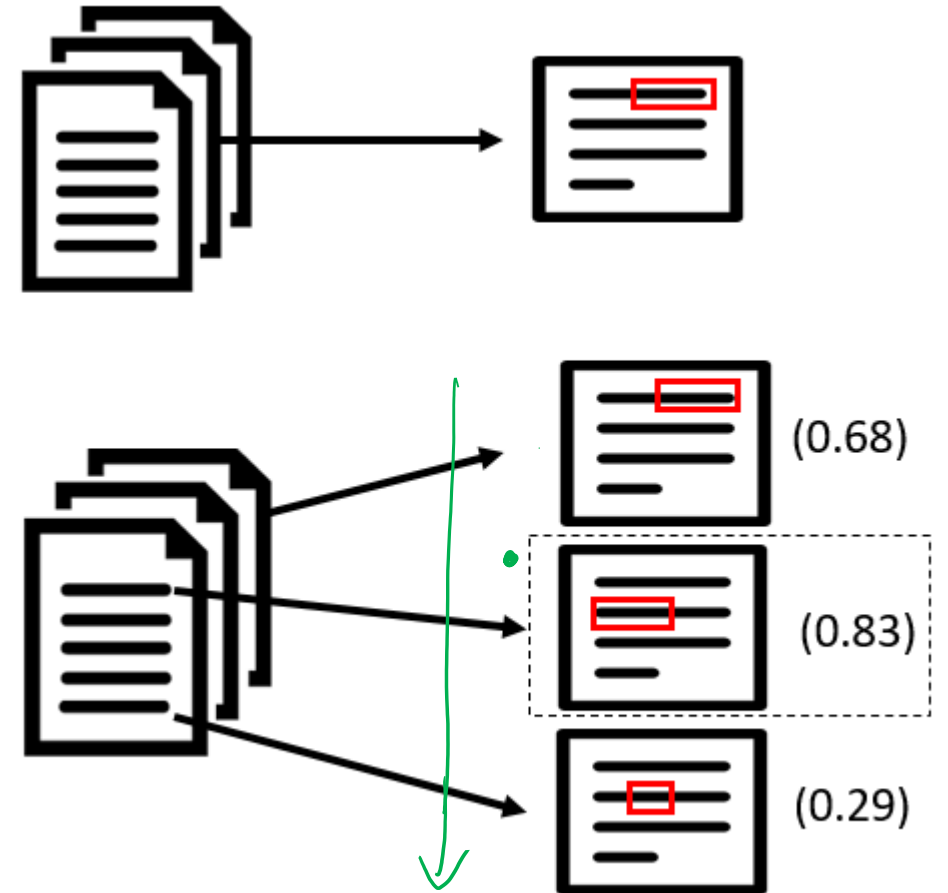
① Training's (q, doc) pairs



② Shared normalization
multipassage

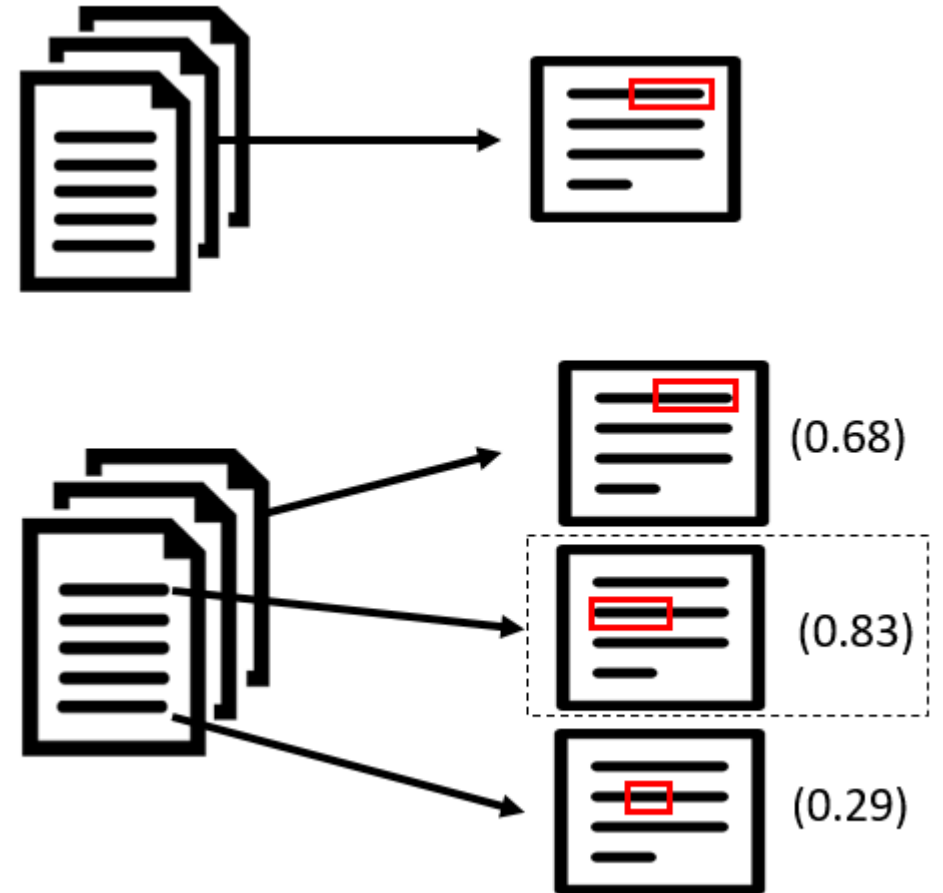
Challenge: Scaling Models to Documents

- ① ■ Pipelined Systems
 - Select a single paragraph from the input, and run the model on that paragraph
- ② ■ Confidence Systems *DrQA*
 - Run the model on many paragraphs from the input, and have it assign a confidence score to its results on each paragraph



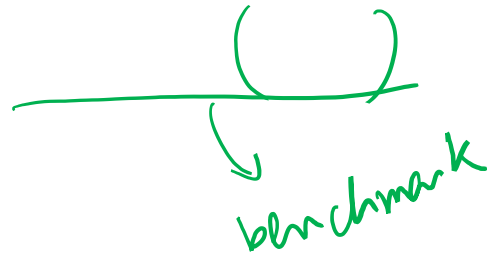
DocumentQA

- Improved Pipeline Method
 - Improve several of key design decisions that arise when training on document-level data
- Improved Confidence Method
 - Study ways to train models to produce correct confidence scores



Pipeline Method: Paragraph Selection

- Train a shallow linear model to select the best paragraphs *from corpus*
- Features include TF-IDF, word occurrences, and its position within document
- If there is just one document TF-IDF alone is effective
- Improves chance of selecting an answer-containing paragraph from 83.0 to 85.1 on TriviaQA Web



Pipeline Method: Noisy Supervision

Document level data can be expected to be distantly supervised:

Q Question: Which British general was killed at Khartoum in 1885?

A Gordon

Passage:

In February 1884 Gordon returned to the Sudan to evacuate Egyptian forces. Rebels broke into the city, killing Gordon and the other defenders. The British public reacted to his death by acclaiming ' Gordon of Khartoum', a saint. However, historians have since suggested that Gordon defied orders and....

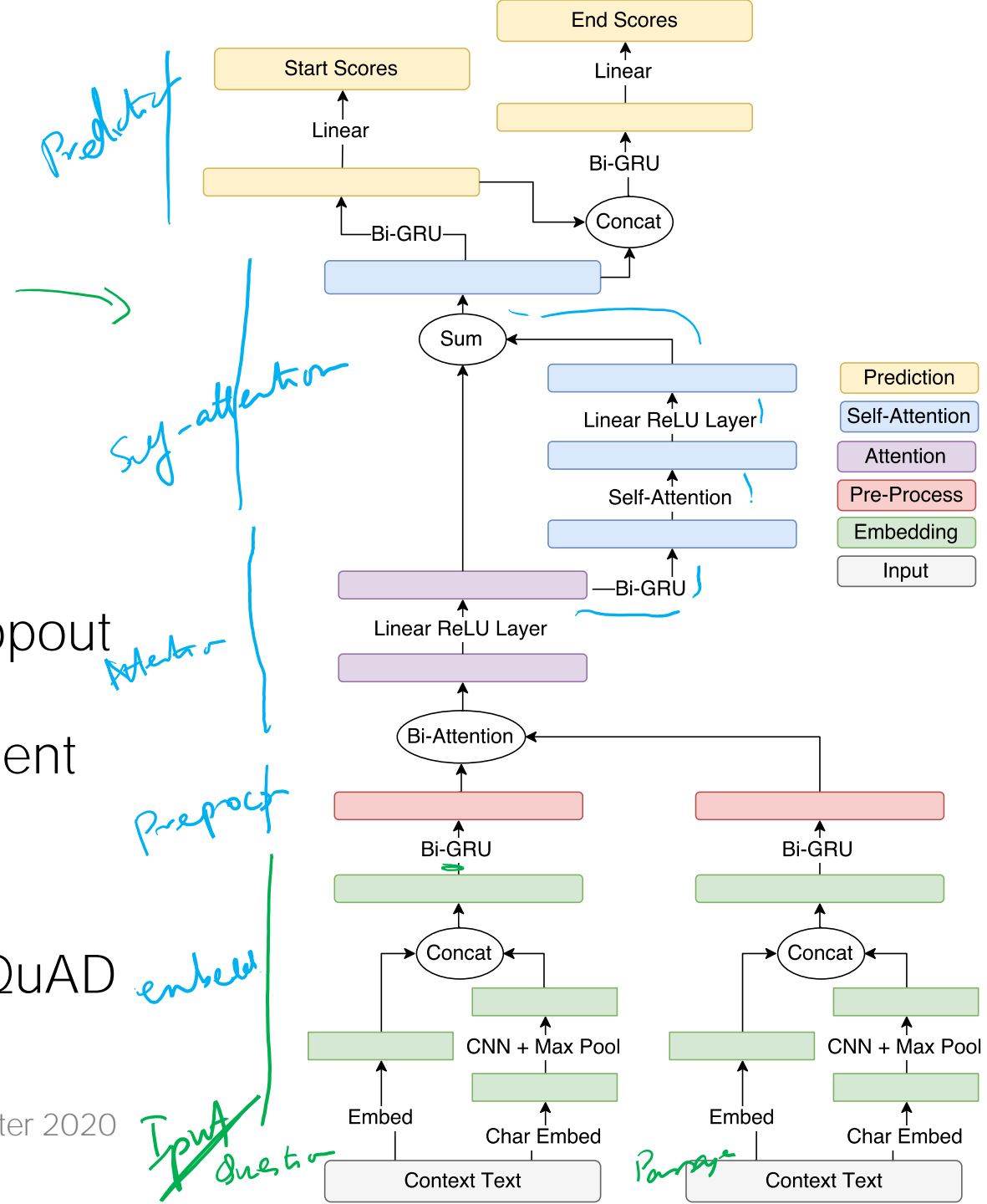
Pipeline Method: Noisy Supervision

- Need a training objective that can handle multiple (noisy) answer spans
- Use the summed objective from Kadlec et al (2016), that optimizes the log sum of the probability of all answer spans
- Remains agnostic to how probability mass is distributed among the answer spans

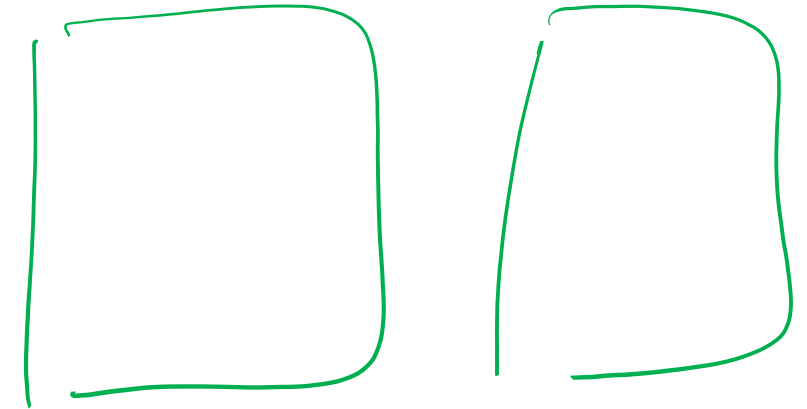
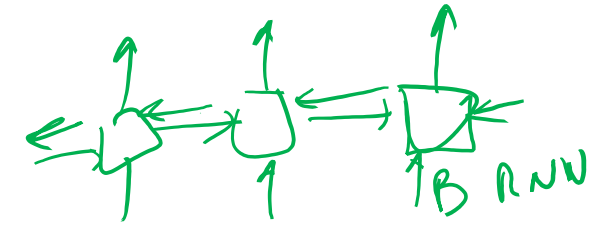
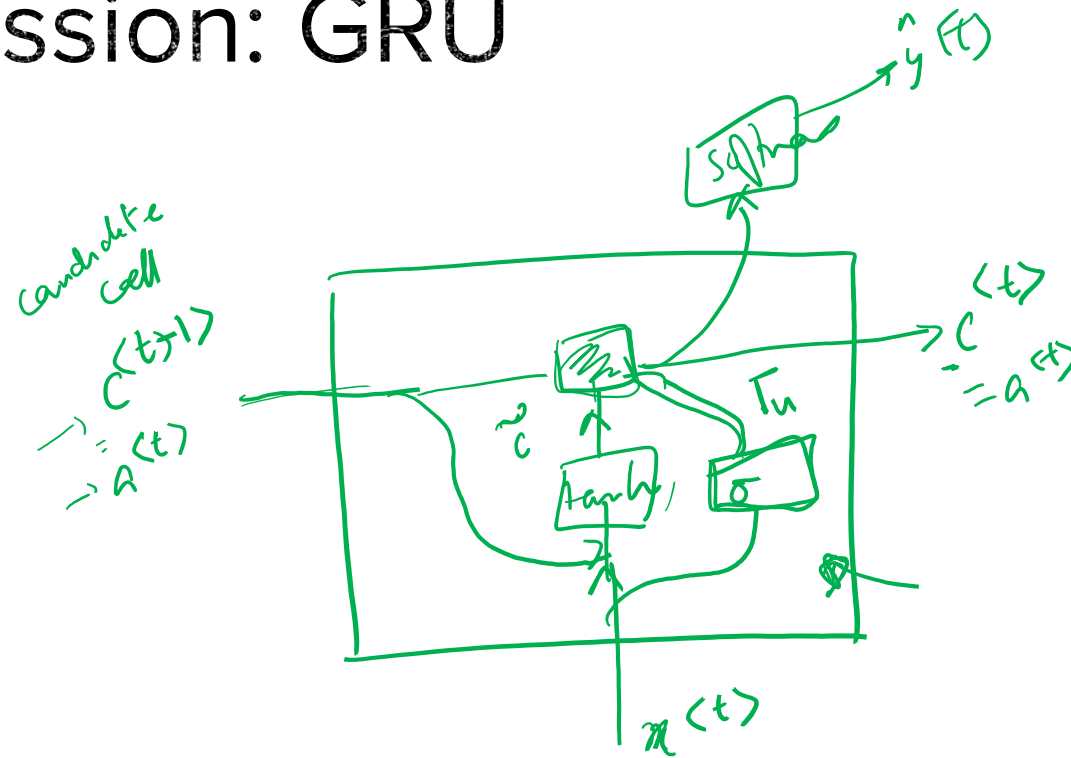
Pipeline Method: Model

- Construct a fast, competitive model
- Use some key ideas from prior work, bidirectional-attention, self-attention, character-embeddings, variational dropout
- Also added learned tokens for document and paragraphs starts
- < 5 hours to train for 26 epochs on SQuAD

start $\langle w_i \rangle \langle v_c \rangle$



Digression: GRU



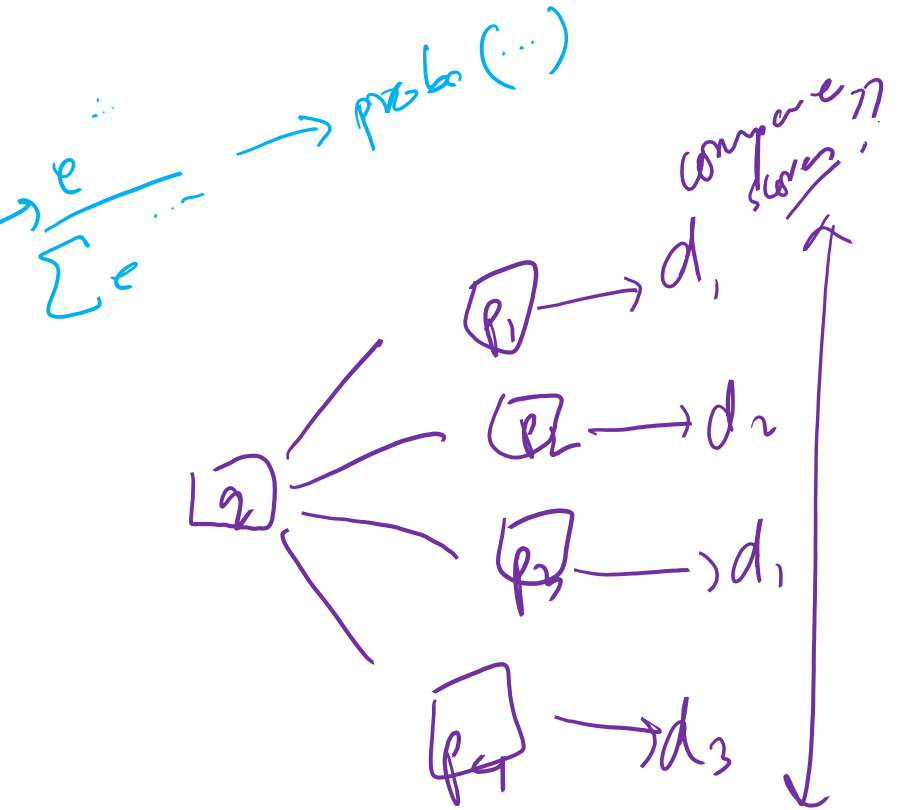
GRU ~ 2014
LSTM ~ 1997

They said that Obama, man, was a great president.



Confidence methods

- ^{then} We can derive confidence scores from the ^{logistic $\sigma()$} logit scores given to each span by the model, i.e., the scores given before the softmax operator is applied
- Without re-training this can work poorly



Example from SQuAD

Question: “When is the Members Debate held?”

Model Extraction: “..majority of the Scottish electorate voted for it in a referendum to be held on 1 March 1979 that represented at least... ”

Correct Answer: “Immediately after Decision Time a “Members Debate” is held, which lasts for 45 minutes...”

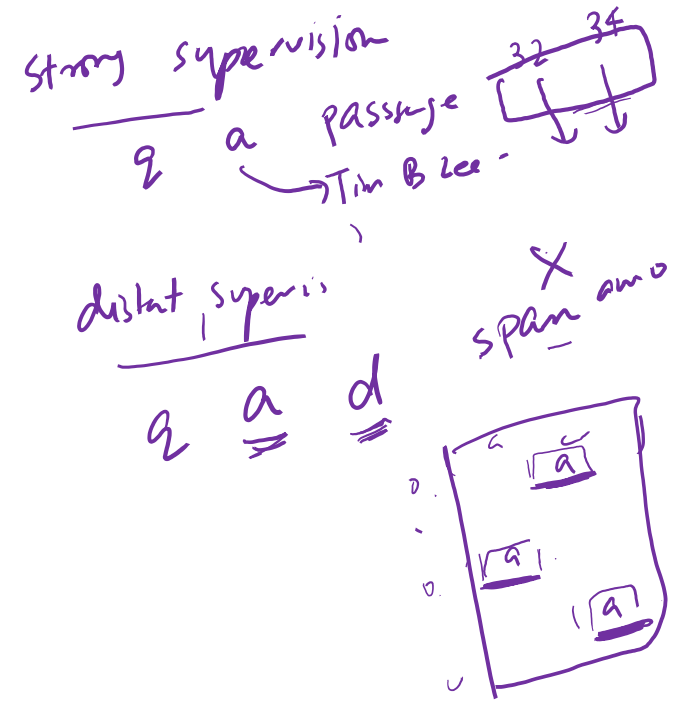
Learning well-calibrated confidence scores *for span*

- Train the model on both answer-containing and non-answer-containing paragraph and use a modified objective function
 - 1 ■ Merge: Concatenate sampled paragraphs together
 - 2 ■ No-Answer: Process paragraphs independently, and allow the model to place probability mass on a “no-answer” output
 - 3 ■ Sigmoid: Assign an independent probability on each span using the sigmoid operator
 - 4 ■ Shared-Norm: Process paragraphs independently, but compute the span probability across spans in all paragraphs
- $$\frac{e^{s_{ap}}}{\sum_{j \in \mathcal{P}} \sum_{i=1}^{n_j} e^{s_{ij}}}$$

over all passages with same context

Conclusions

- Neural machine reading comprehension systems now form the bulk of text-QA *today*
- MRC systems coupled with retrieval pipeline result in so-called “open domain QA systems” *IR+MRC*
- Retrieval is still primarily TF-IDF based
- *Reader* Answers from text are predicted as free-form spans *arbitrary*
- Often designed as classifiers for start and end positions in passage text *not restricted to entities*



Thank
you