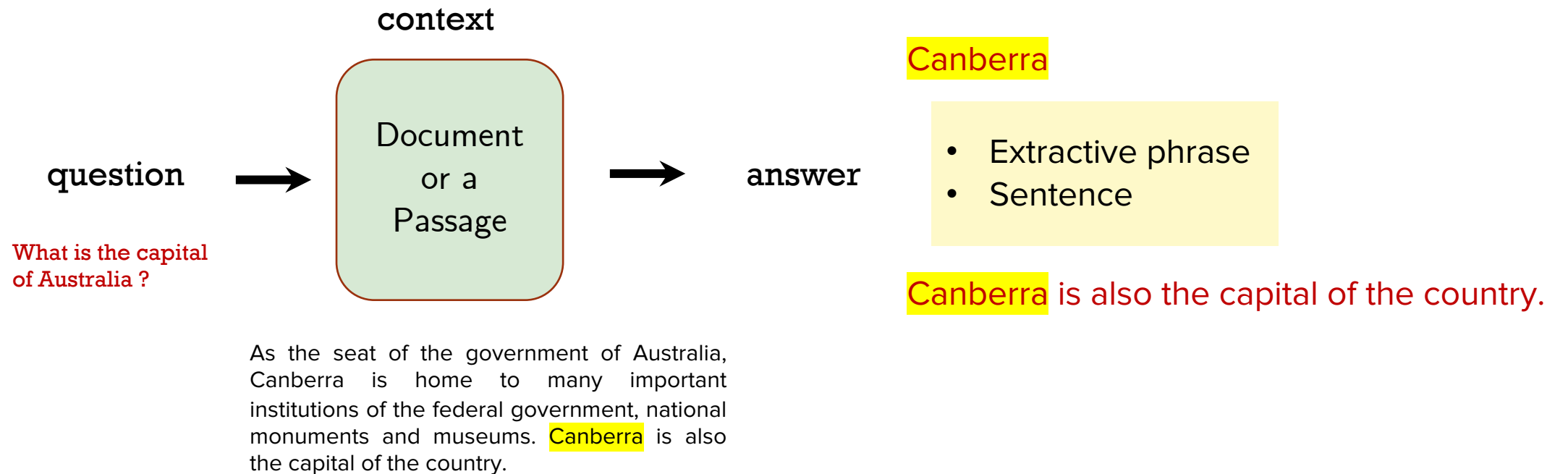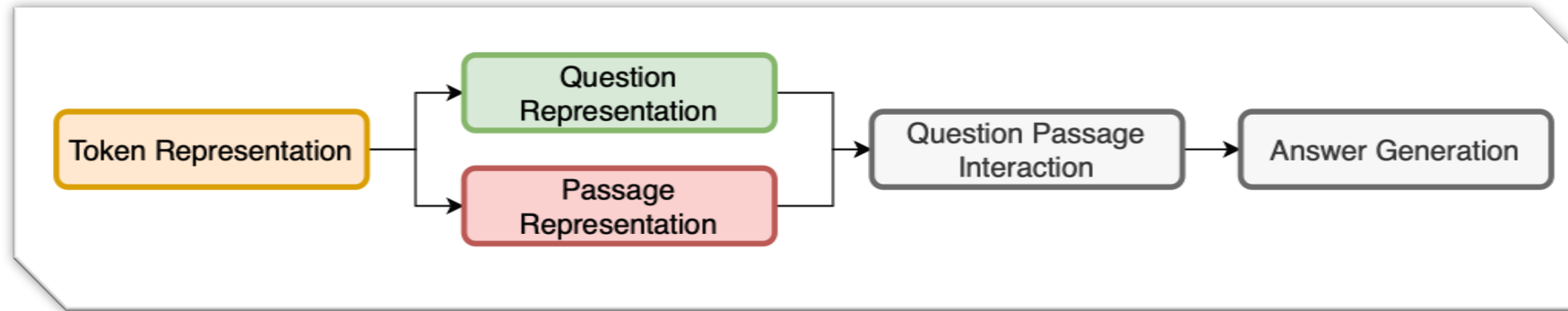# MACHINE COMPREHENSION

# Machine Comprehension

Chris Burges 2013

"A machine **comprehends** a passage of **text** if, for any **question** regarding that text that can be **answered** correctly by a majority of native speakers, that machine can provide a string which those speakers would agree both answers that question, and does not contain information irrelevant to that question."

# Problem Setting

context

**question** $\longrightarrow$

Document or a Passage

$\longrightarrow$ **answer**

What is the capital of Australia ?

As the seat of the government of Australia, Canberra is home to many important institutions of the federal government, national monuments and museums. Canberra is also the capital of the country.

Canberra

- Extractive phrase
- Sentence

Canberra is also the capital of the country.
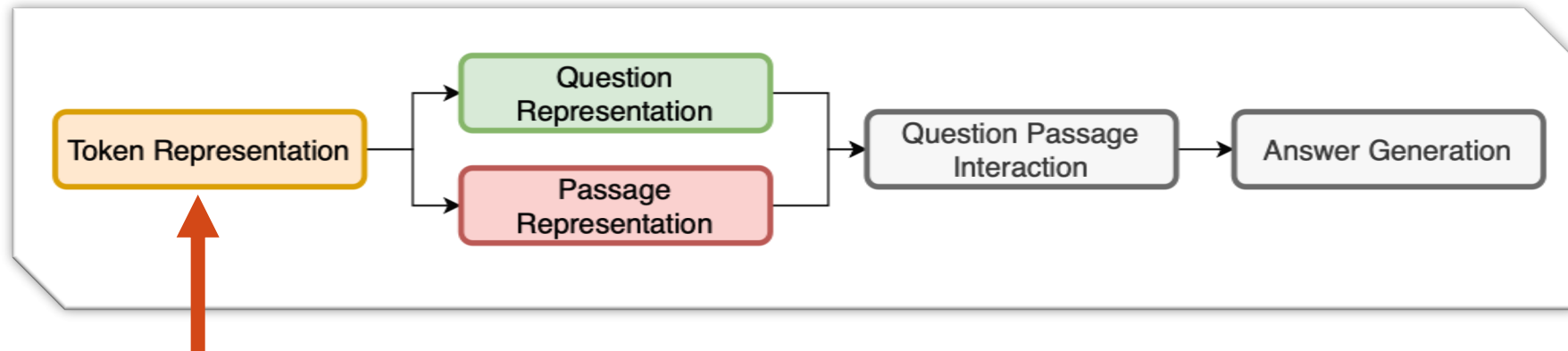
# The MRC Pipeline



- Words, characters, sub-words embeddings
- Contextual Embeddings
- Other features – Matching, Alignment, Language structure

- Sequential representation
- Contextual representation
- Attentive reading

- Attentive reading
- Attention flows
- Multiple input passes inputs
- Re-representation of question and passages

- Token prediction
- Span prediction
- Free-form generation

# Token Representation



| Conventional | Linguistic | Contextual |
|---|---|---|
| • One Hot | • POS | • Bi-LSTM |
| • Word Embeddings | • Named Entity | • BERT |
| • Sub-word, Character Embeddings | • Query Category | • ELMO |

# Question/Passage Representation



| CNN | RNN | Transformers |
|---|---|---|
| • CNN. for doc rep.<br>• Cross-attention<br>• … | • LSTM<br>• Bi-LSTM<br>• Bi-GRU | • GPT2<br>• BERT<br>• XLNET |

# Question And Passage Interaction



| Attention | Interaction Type |
|---|---|
| • Uni directional<br><br>• Bi-directional<br><br>• Cross-attention | <br>• Single Hop Interaction<br><br>• Multi hop Interaction |

# Answer Generation



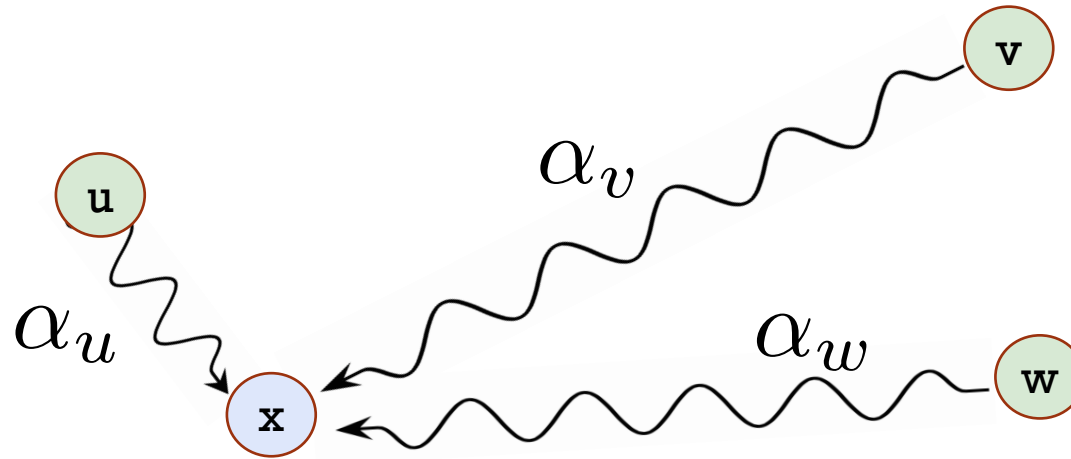| Cloze | MCQ | Free Text | Span Pred. |
|---|---|---|---|
| • Word-level prediction (BC) | • Choosing one answer among many (MC) | • Generative models | • Predict begin and end of sequence |

# Attention Mechanism

- Attention is used to represent tokens, question and passages
    - How do we re-represent otherwise independent token representations ?
    - How do we leverage contextualization ?


- Hard attention
- **Soft Attention**
- Co-attention
- Self-attention

# Attention − Influence Point Of View

$$\alpha_u = \left( \frac{e^{\mathbf{x} \cdot \mathbf{u}}}{e^{\mathbf{x} \cdot \mathbf{u}} + e^{\mathbf{x} \cdot \mathbf{v}} + e^{\mathbf{x} \cdot \mathbf{w}}} \right)$$
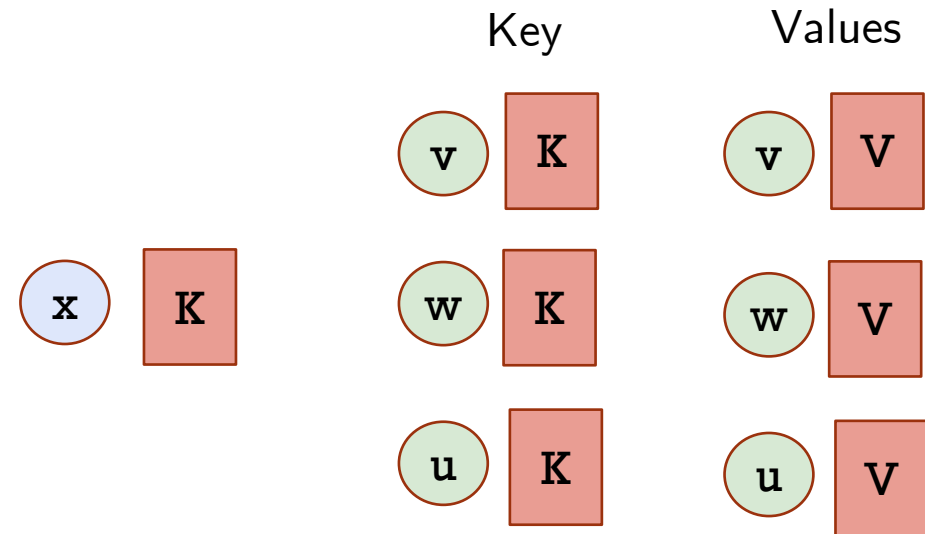


$$\mathbf{x}' = \alpha_u \mathbf{u} + \alpha_v \mathbf{v} + \alpha_w \mathbf{w}$$

- Typically x and context vectors are first projected through a learnable matrix W

# Attention Mechanism – Memory Point Of View

Attention retrieves values from a continuous memory using fuzzy matching

Key Values

- Assume vectors are stored in memory referenced by Key matrix K

- Thought expt: for 1-hot vectors = hashmaps

- Instead Kx retrieves from this continuous memory as a weighted sum over all values

Attention weight
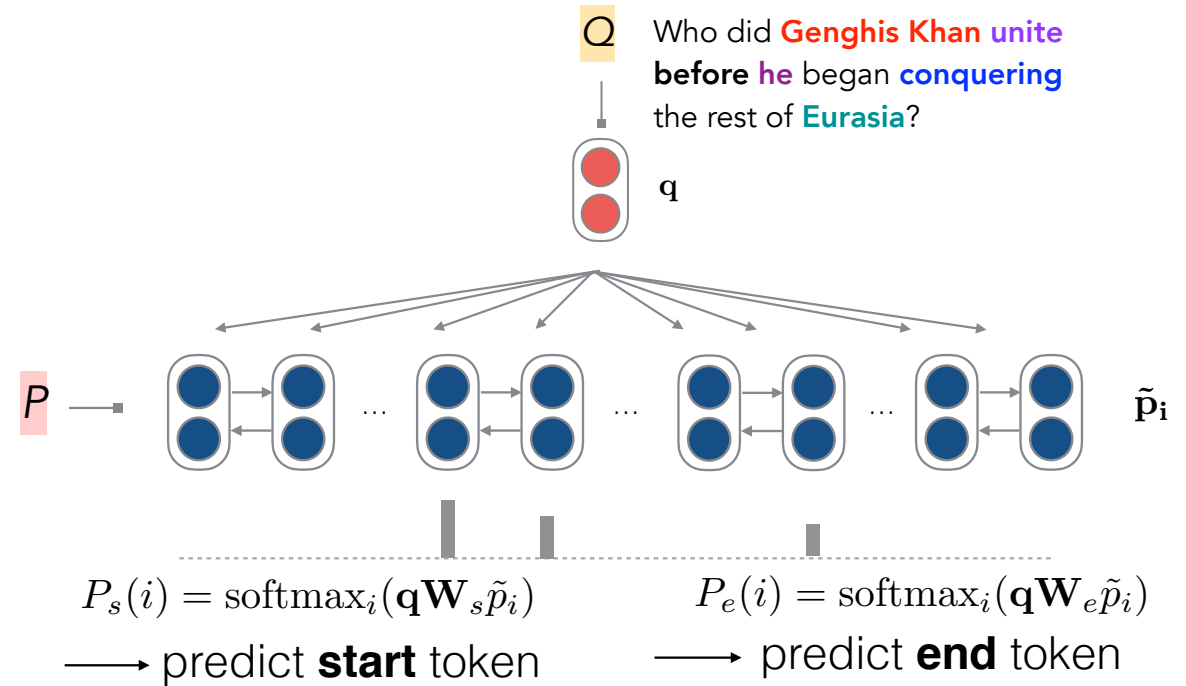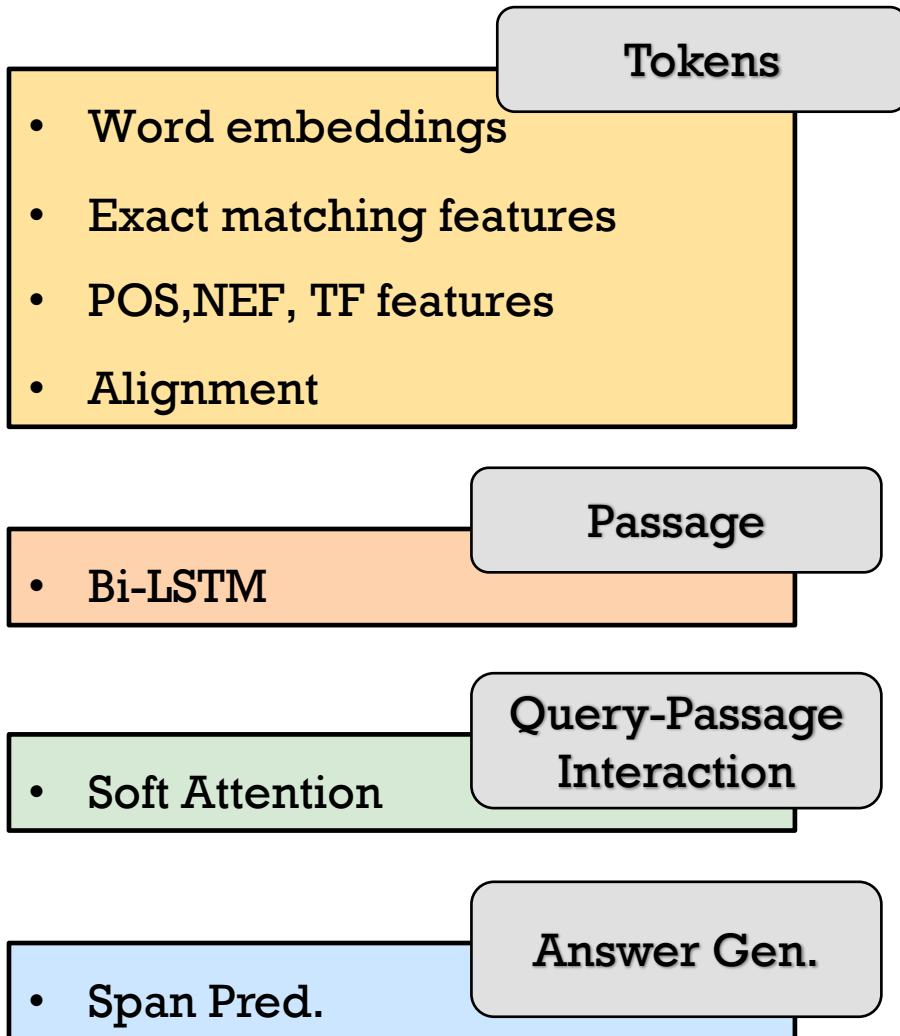
$$\alpha_u = \frac{e^{K\mathbf{x} \cdot K\mathbf{u}}}{e^{K\mathbf{x} \cdot K\mathbf{u}} + e^{K\mathbf{x} \cdot K\mathbf{v}} + e^{K\mathbf{x} \cdot K\mathbf{w}}}$$

$$\mathbf{x}' = \alpha_u \mathbf{u} + \alpha_v \mathbf{v} + \alpha_w \mathbf{w}$$

answering?  **A** a compu

**Tokens**
- Word embeddings
- Exact matching features
- POS,NEF, TF features
- Alignment

**A** 42

**Passage**
- Bi-LSTM

**iever + Document P**

**Query-Passage Interaction**
- Soft Attention

(reading comprehension

**Answer Gen.**
- Span Pred.

Bidirectional LSTMs

Who did **Genghis Khan** **unite** **before** he began **conquering** the rest of **Eurasia**?

$q$

$P$

$\tilde{p}_i$
$\mathbf{p}_i$

$$P_s(i) = \text{softmax}_i(\mathbf{q}\mathbf{W}_s\tilde{p}_i)$$

$\longrightarrow$ predict **start** token

$$P_e(i) = \text{softmax}_i(\mathbf{q}\mathbf{W}_e\tilde{p}_i)$$

$\longrightarrow$ predict **end** token

**Data:** SQuAD + **Distantly Supervised** Data

**Data:** SQuAD + **Distantly Supervised** Data

**Data:** SQuAD + **Distantly Supervised** Data

(Q', A) → (P, Q, A) if P is retrieved and A can be found in P

**Q:** What part of the atom did Chadwick discover?

WebQuestions

# Reader

puter science
etrieval and natural al

**Reader**

les from 5 million

on system):

articles

Who did **Genghis Khan** **unite**
**before** **he** began **conquering**
the rest of **Eurasia**?

**Bidirectional LSTMs**

$\mathbf{q}$

$P$

$\tilde{p}_i$

$$P_s(i) = \mathrm{softmax}_i(\mathbf{q}\mathbf{W}_s\tilde{p}_i)$$

⟶ predict **start** token

$$P_e(i) = \mathrm{softmax}_i(\mathbf{q}\mathbf{W}_e\tilde{p}_i))$$

⟶ predict **end** token

$$\Pr(a|q, p_i) = P_s(a_s)P_e(a_e)$$

**Data:** SQuAD + **Distantly Supervised** Data

Logistic regressio
Fine-Grained Gating (C
Match-LSTM (Singa
Management U
DCN (Salesforce
BiDAF (UW & Allen Ins
Ours
r-net (MSR Asia
State-of-the-art (July
Human performan

Pre-trained SQuAD model
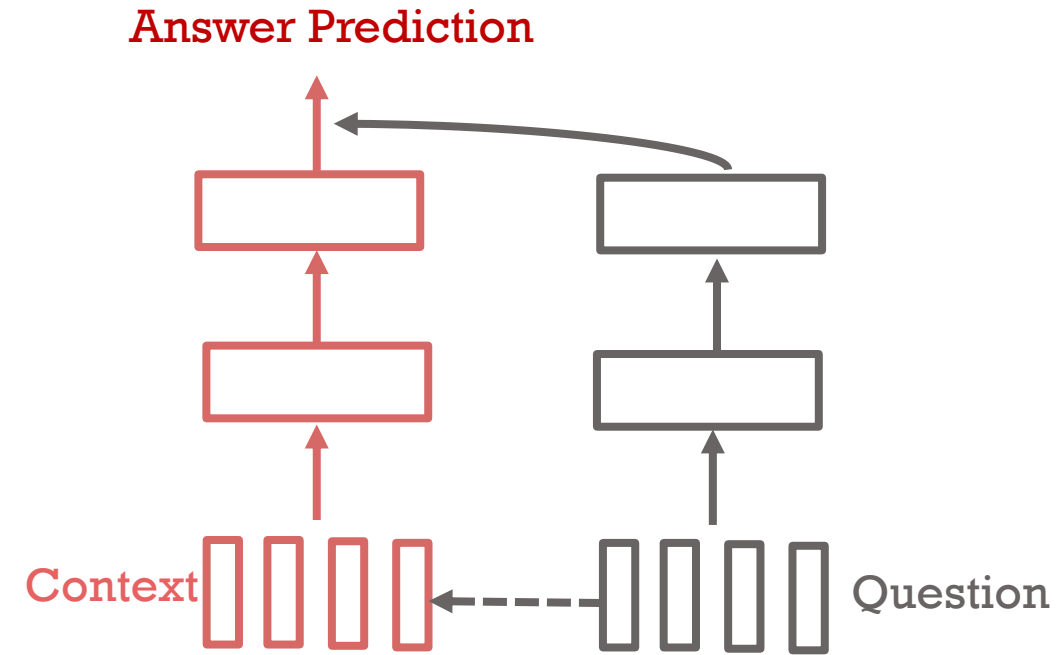
# Input Representation

- Context words are represented based on similarity with the query

- Semantic similarity
  - Word embeddings

- Matching similarity
  - Direct word-level matching
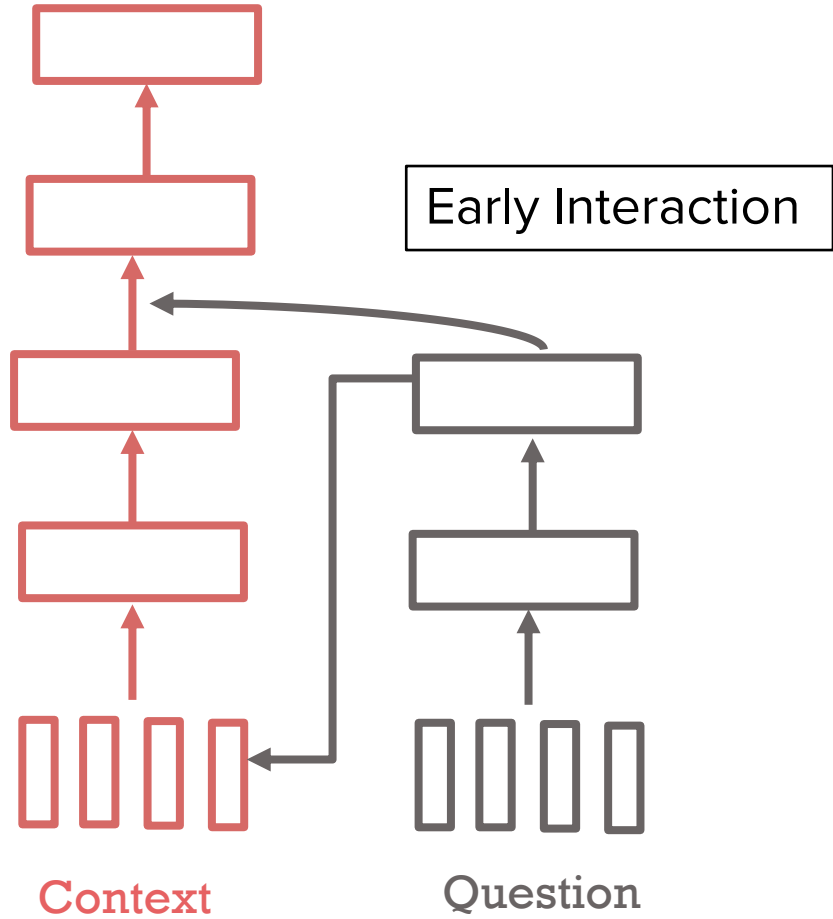  - Weighted matching
  - Attention mechanism

**Answer Prediction**

Context ← Question

Embedding-based representation

Query-based representation

# Late Interaction

- First encode question and context sufficiently

- Choice of encoders
  - Bi-LSTMs
  - Conv Nets

- Most popular Model
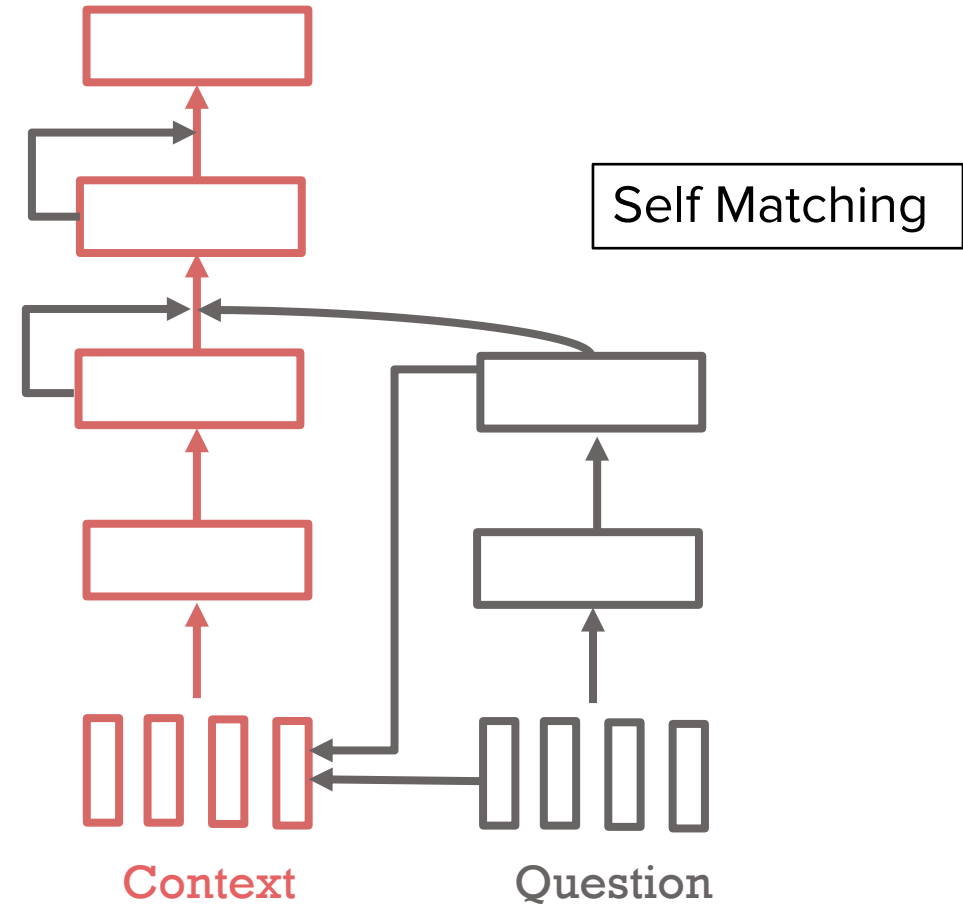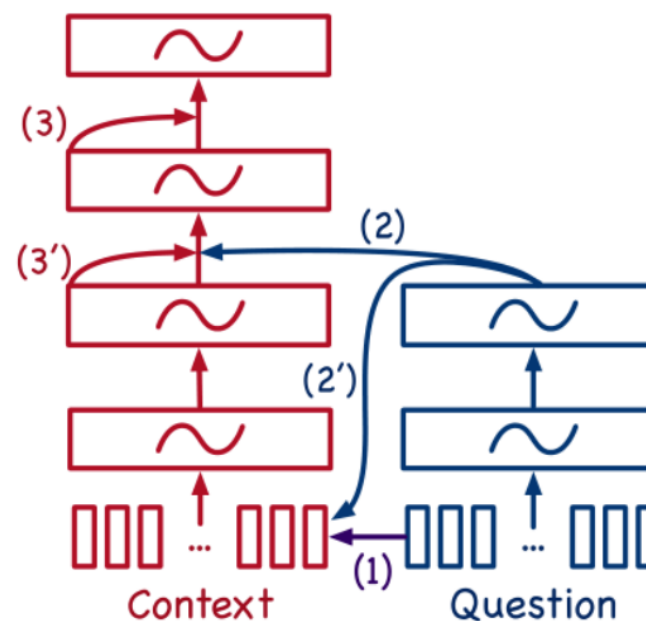  - Bi-directional attention flow [Seo '17]

**Answer Prediction**

Context — — — Question

# Other Variants



Answer Prediction

Early Interaction

Context    Question

Answer Prediction

Self Matching

Context    Question

# Attention Based Architectures

- 2016 – 2017 – Multitude of attention based architecture

| Architectures | (1) | (2) | (2') | (3) | (3') |
|---|---|---|---|---|---|
| Match-LSTM (Wang and Jiang, 2016) | | ✓ | | | |
| DCN (Xiong et al., 2017) | | ✓ | | | ✓ |
| FastQA (Weissenborn et al., 2017) | ✓ | | | | |
| FastQAExt (Weissenborn et al., 2017) | ✓ | ✓ | | ✓ | |
| BiDAF (Seo et al., 2017) | | ✓ | | | ✓ |
| RaSoR (Lee et al., 2016) | ✓ | | ✓ | | |
| DrQA (Chen et al., 2017) | ✓ | | | | |
| MPCM (Wang et al., 2016) | ✓ | ✓ | | | |
| Mnemonic Reader (Hu et al., 2017) | ✓ | ✓ | | ✓ | |
| R-net (Wang et al., 2017b) | | ✓ | | ✓ | |



(1) Word-level fusion, (2) high-level fusion, (2') high-level fusion (alternative), (3) self-boosted fusion, and (3') self-boosted fusion (alternative).
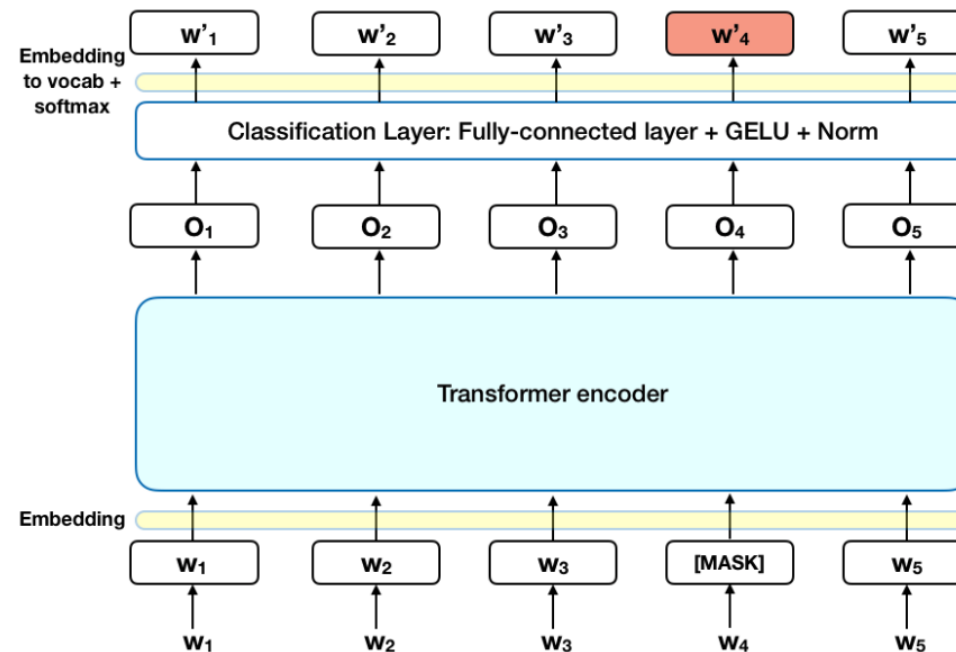
# Contextual Language Models

- BERT – No Recurrence, only attention

- Re-representing each token based on the context

- Shows the most promising performance

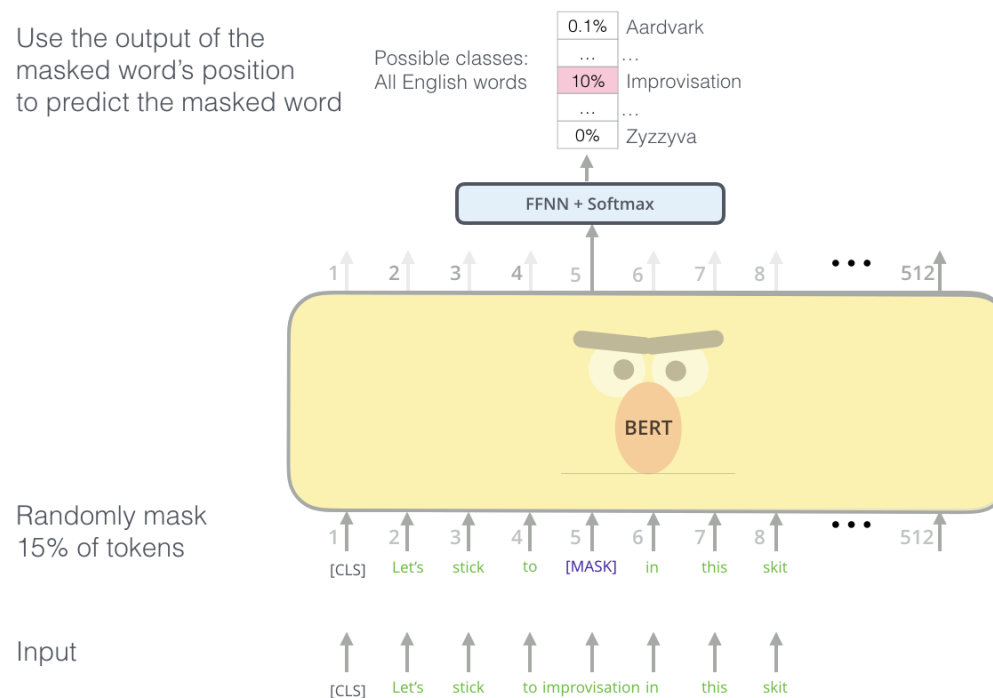**Answer Prediction**



Context    Question

# BERT

- Bi-directional : Transformer encoder reads the entire sequence of words at once.
  - Learns the context of a word based on all of its surroundings (left and right of the word).
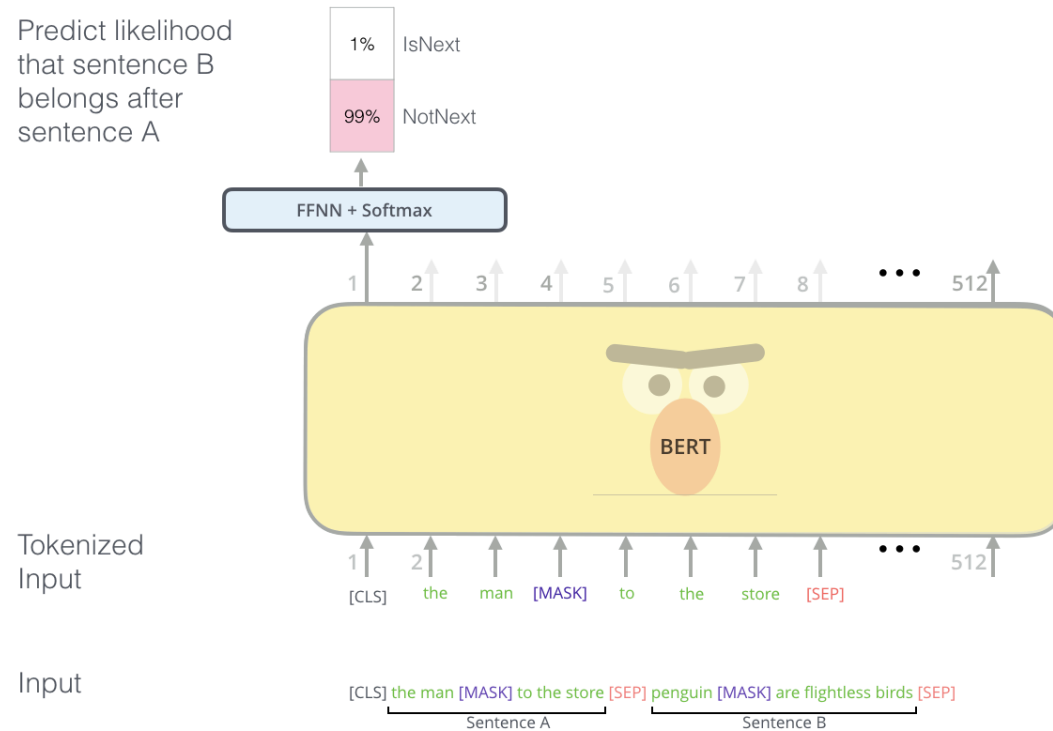
# BERT– Masked Language Model

## Masked word prediction

- Given a sentence with some words masked at random, can we predict them?
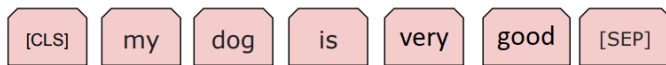- Randomly select 15% of tokens to be replaced with "<MASK>"

# Next Sentence Prediction

- Given two sentences, does the first follow the second? Teaches BERT about relationship between two sentences
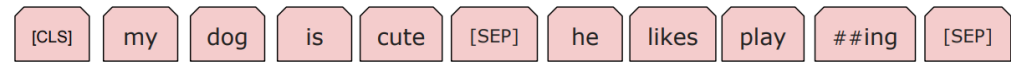- 50% of the time the actual next sentence, 50% random

Predict likelihood that sentence B belongs after sentence A

| 1% | IsNext |
| 99% | NotNext |

FFNN + Softmax

1  2  3  4  5  6  7  8  ···  512

BERT

Tokenized Input

1  2  ···  512

[CLS]  the  man  [MASK]  to  the  store  [SEP]

Input

[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flightless birds [SEP]

Sentence A          Sentence B

# BERT Fine Tuning

- Inputs to BERT – [CLS] <token embeddings> [SEP] ...

| [CLS] | my | dog | is | very | good | [SEP] |
|---|---|---|---|---|---|---|

### Single sentence input

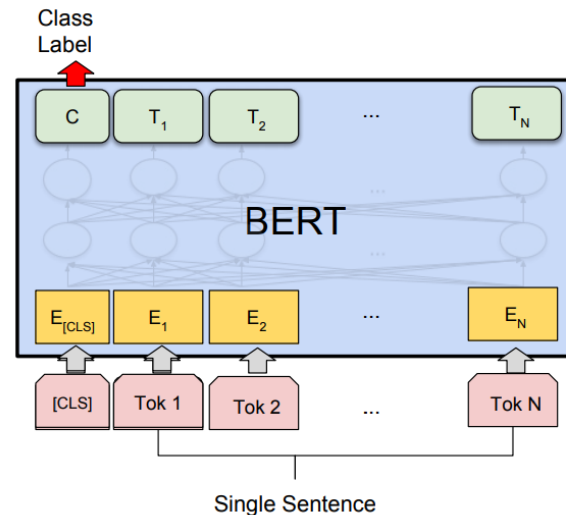| [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|

### Single sentence input

- Classification tasks such as sentiment analysis are done similarly to Next Sentence classification, by adding a classification layer on top of the Transformer output for the [CLS] token.

# BERT Fine Tuning

Q&A model can be trained by learning two extra vectors that mark the beginning and the end of the answer.