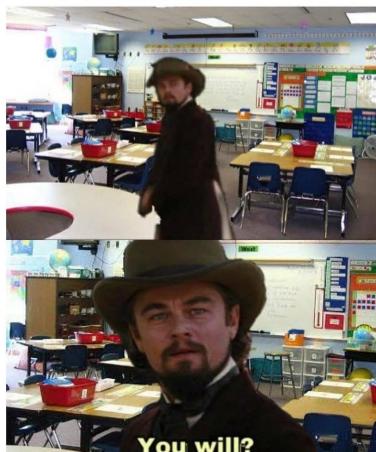


Teacher: If you don't want to pay attention I'll honestly end this class right now.

The Students:



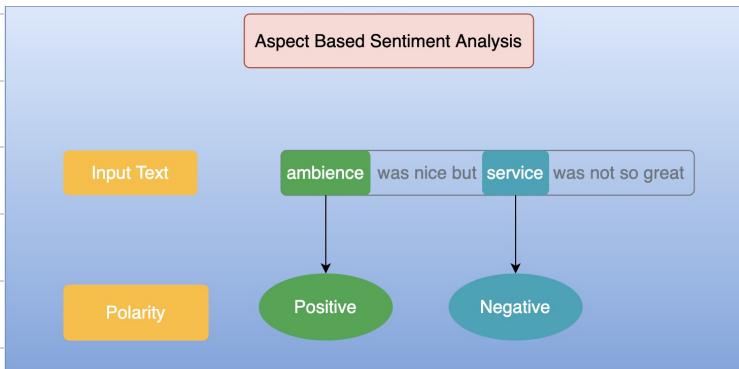
Agenda

1. Case-study, how can we solve it?
2. What's wrong with RNNs?
3. How attention helps?
4. How attention works
5. Different kinds of attention
6. Self Attention and Multi-head attention
7. Code Walkthrough.

Aspect Based Sentiment Analysis (ABSA)

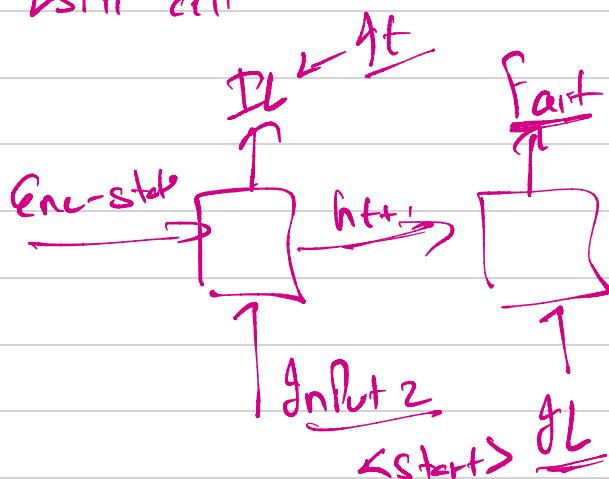
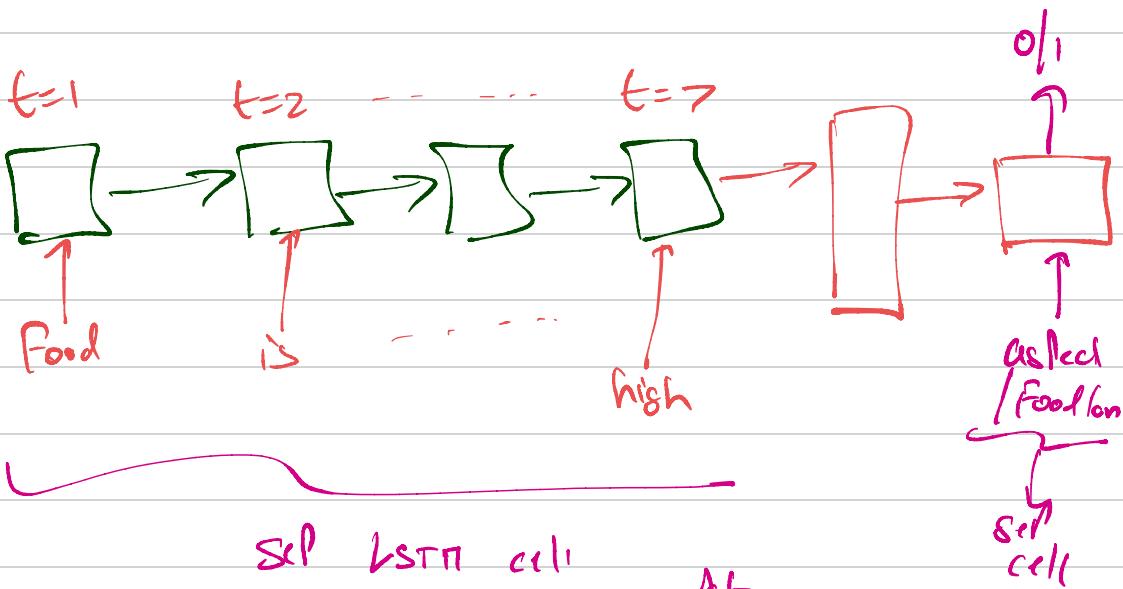
Food is delicious but price is high) Overall neutral

taste → U. Positive "taste"
Cost → U. Negative.



- ① Use NER to remove / select other entries
↳ Personalized NER model

Encoder - Decoder.



Bidirectional



What are some problems with LSTM??

1. If I have a very long paragraph, let's say like 100 words, then the specific word like "taste" maybe present in the beginning of the word.

Now cell state vector takes care of that little bit, but still information is lost.

2. The weight matrix, for forget, input, and output gate are the same wt matrix, being multiplied to all the words.
3. Lack of parallel training for LSTM networks.

ATTENTION

① Key, Query, Value \rightarrow matrices

② Attention function / Energy / Alignment

③ Content Vector

Value matrix $f(10, 512) \leftarrow w_k \cdot f(512, 512)$

1

Key matrix
↓

$\rightarrow [w_1, \dots, w_{10}]$

$\rightarrow f(2, 10 \times 512) \leftarrow b_0$

Representation of tokenized
Sentences

Food is delicious but price is high

Key
matrix

	2	-2.3	3.6	-1.5	0.76	-2.3	-0.31	-3
o	0.64	-0.69	-1.3	-0.69	2.8	-0.69	1.2	-2
r	1.1	-0.81	-0.91	0.95	-0.89	0.81	1.3	-1
m	-1	0.051	0.66	-0.47	-0.6	0.051	1.3	0
a	-0.78	-0.36	0.38	1.2	0.21	-0.36	-0.25	-2
	food	is	delicious	but	price	is	high	

$\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$
 $K_1 \quad K_2 \quad K_3 \quad \dots \quad K_d$

128d or 512d → PCA

sd ??

Never do this in Pract.

2

Query Matrix

$\rightarrow [1, 512] \odot w_q$

$[512]$

↳ Representation of tokenized aspects $[512]$

↓

1 word
2 word
!

1 word → "cost"



③

Value Vector



In our R-station

Store separately in RAR

Value Vector = Key Vector

Initially same

Key matrix: 5×7 / 10, 512

Query matrix: 5×1 / 1, 512

Value matrix: 5×7 / 10, 512

Attention \rightarrow Google

↓
Search Engin

Key: Value



1 trillion Keys

Query

most relevant Key

Key matrix: 5×7 / $(10, 512)$

Query matrix: 5×1 / $(1, 512)$

Value matrix: 5×7 / $(10, 512) = (10, 1)$

$$\frac{2 \times 5}{\downarrow} \quad \frac{5 \times 1}{\longrightarrow} = \frac{2 \times 1}{\downarrow}$$

$$(Q^T \cdot Q)$$

Dot Product of

$$Vec(A) = P \cdot B \rightarrow \text{Similarity}$$

$A \rightarrow B$ of how similar

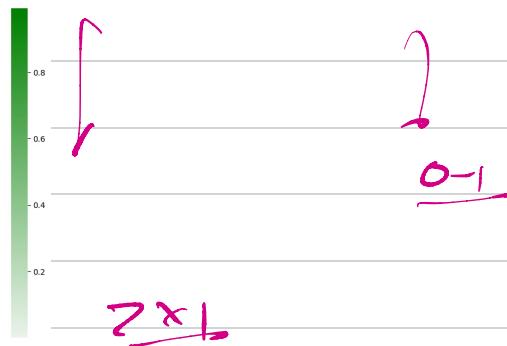
what do they represent

Apply softmax on

similarity of Query w.r.t. each token in Key

Softmax $(Q^T \cdot Q)$

ATTN_WEIGHTS	
food	0.0018
is	4.1e-05
delicious	0.00082
but	0.00012
price	0.99
is	4.1e-05
high	0.0049
cost	



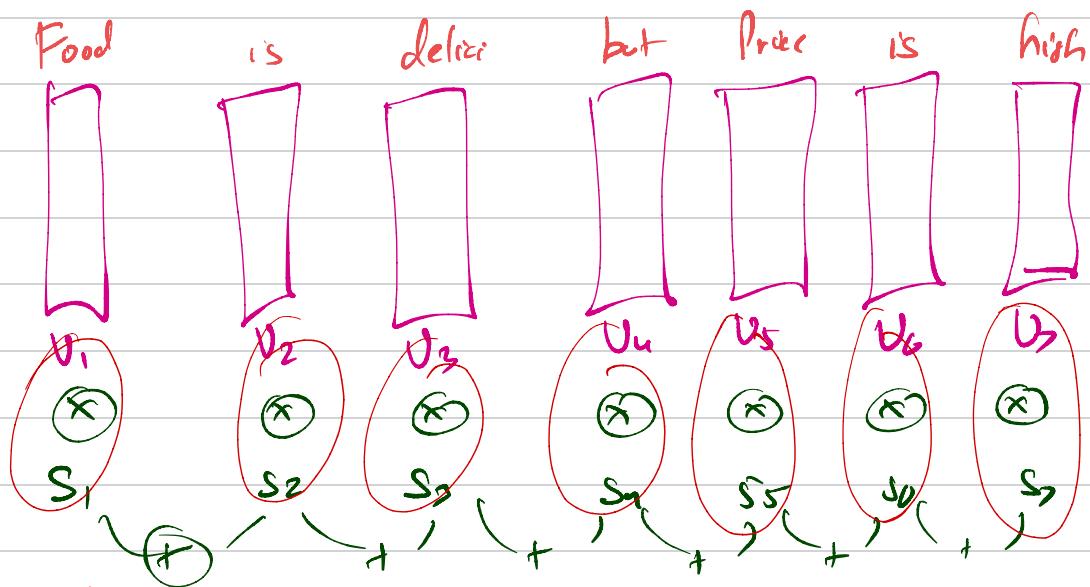
3

Contract Vector

→ Here enters the Value matrix

Element wise multiplication of software outlet
with Value matrix

→ 5×7

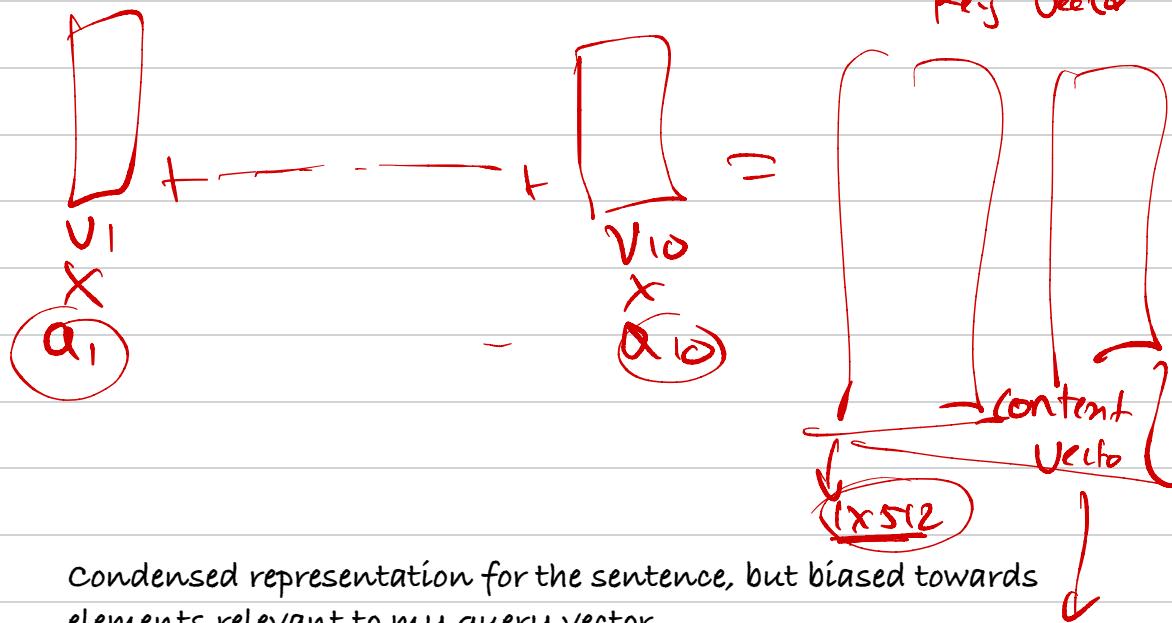


Some Kind of Condense representation

Aggregated information from all parts of data but biased towards elements which have more weightage/important.

$\{2 \text{ } 5 \times 1\}$

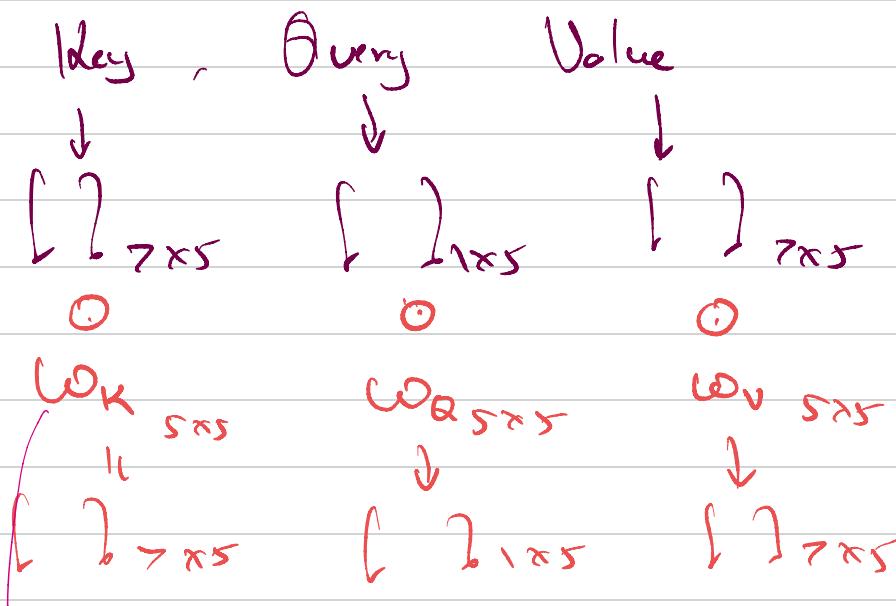
Softmax $\left(\frac{12_{10}, 1}{\sqrt{512}} \right)$ ← Scaled dot product
 \downarrow
 $\left[\begin{array}{c} \\ \\ \end{array} \right] \frac{10 \times 1}{\downarrow}$
 a_1 ← weights
 for each Key vector



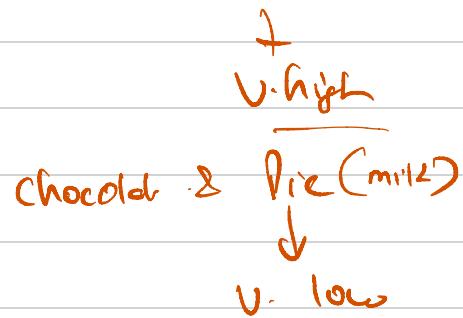
Condensed representation for the sentence, but biased towards elements relevant to my query vector.

$\frac{2 \times 512}{\downarrow}$
 512

$\{ \}_{5 \times 1} \rightarrow \text{MLP} \rightarrow \text{Sigmoid} \rightarrow 0/1$



→ can be learned so that it
can modify my embedding (chocolate & deer)
which benefits my
use-case.

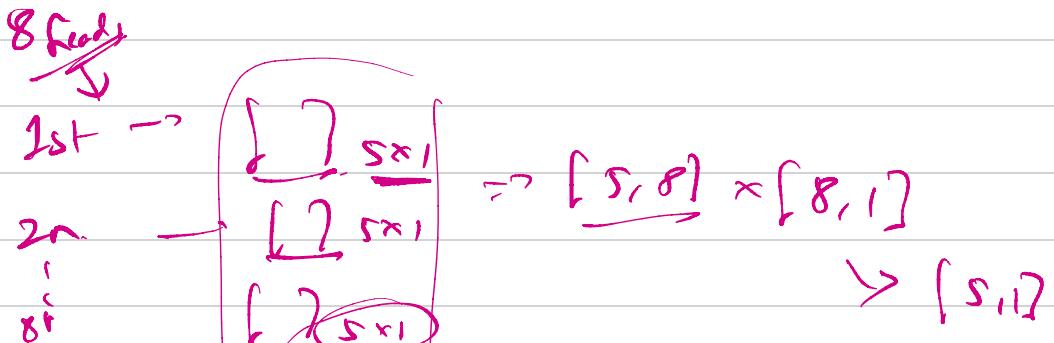
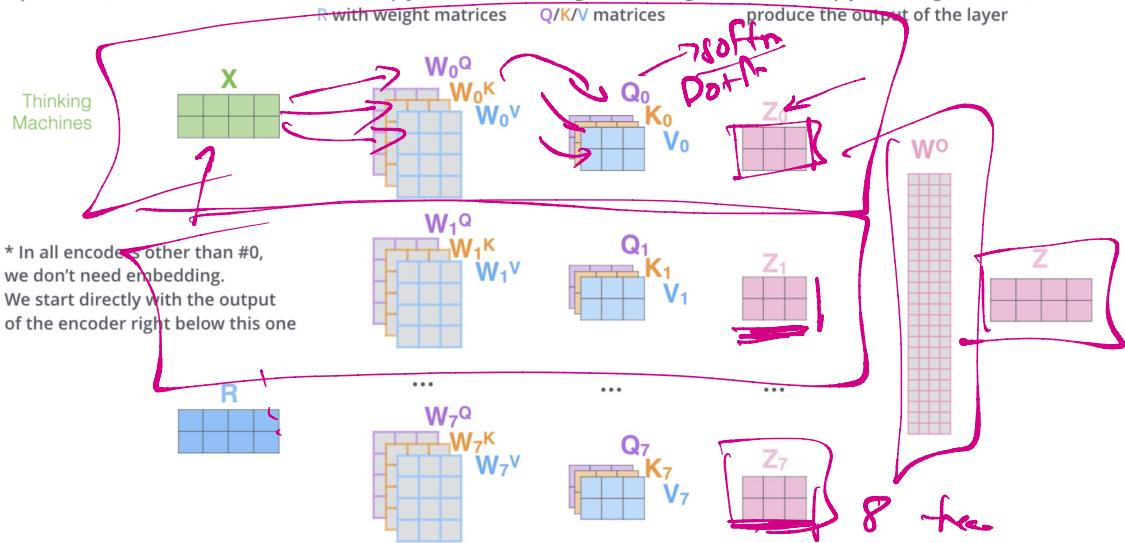


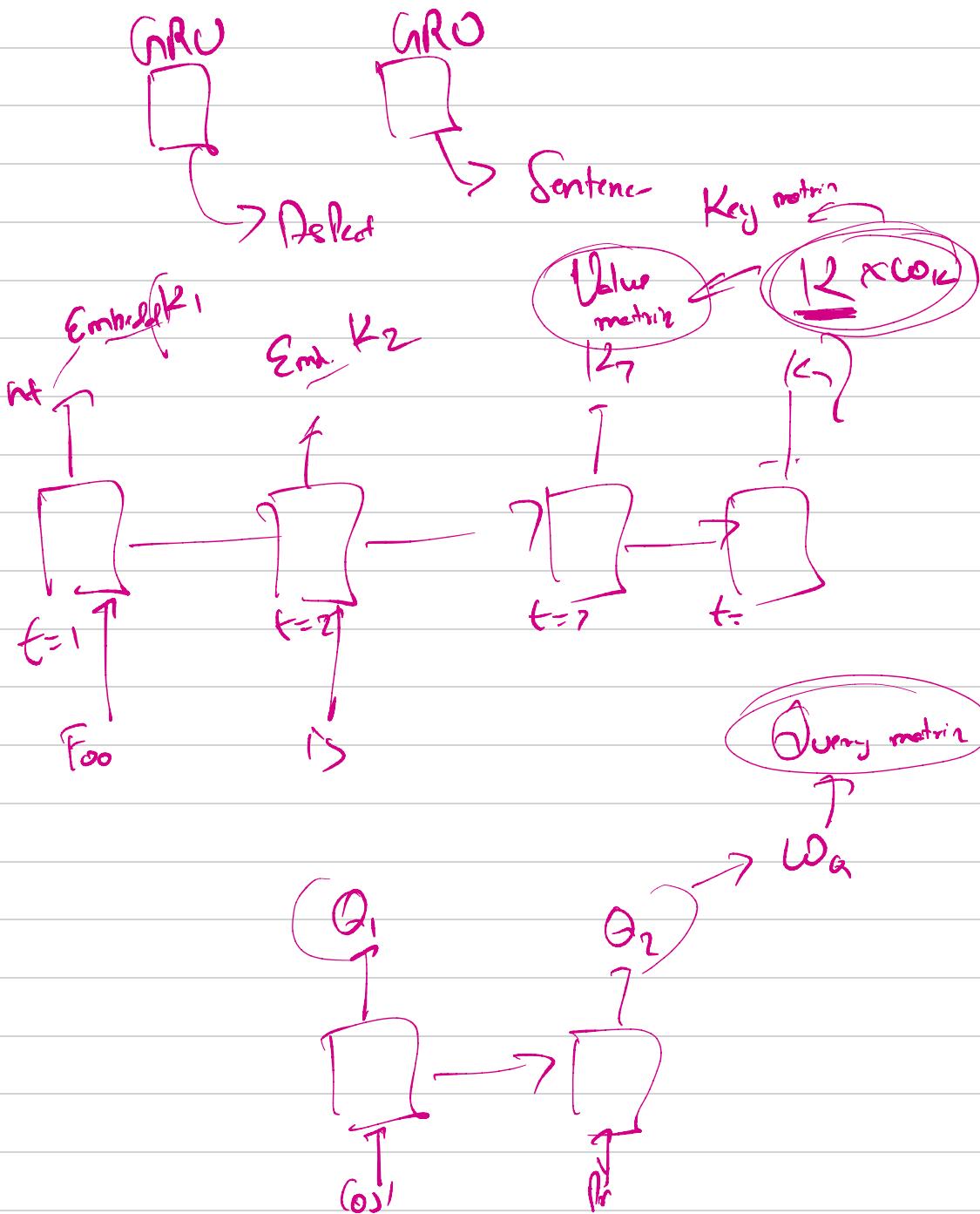
beau

beautiful



- 1) This is our input sentence*
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer





I want food

I want paratha

aspect hunger