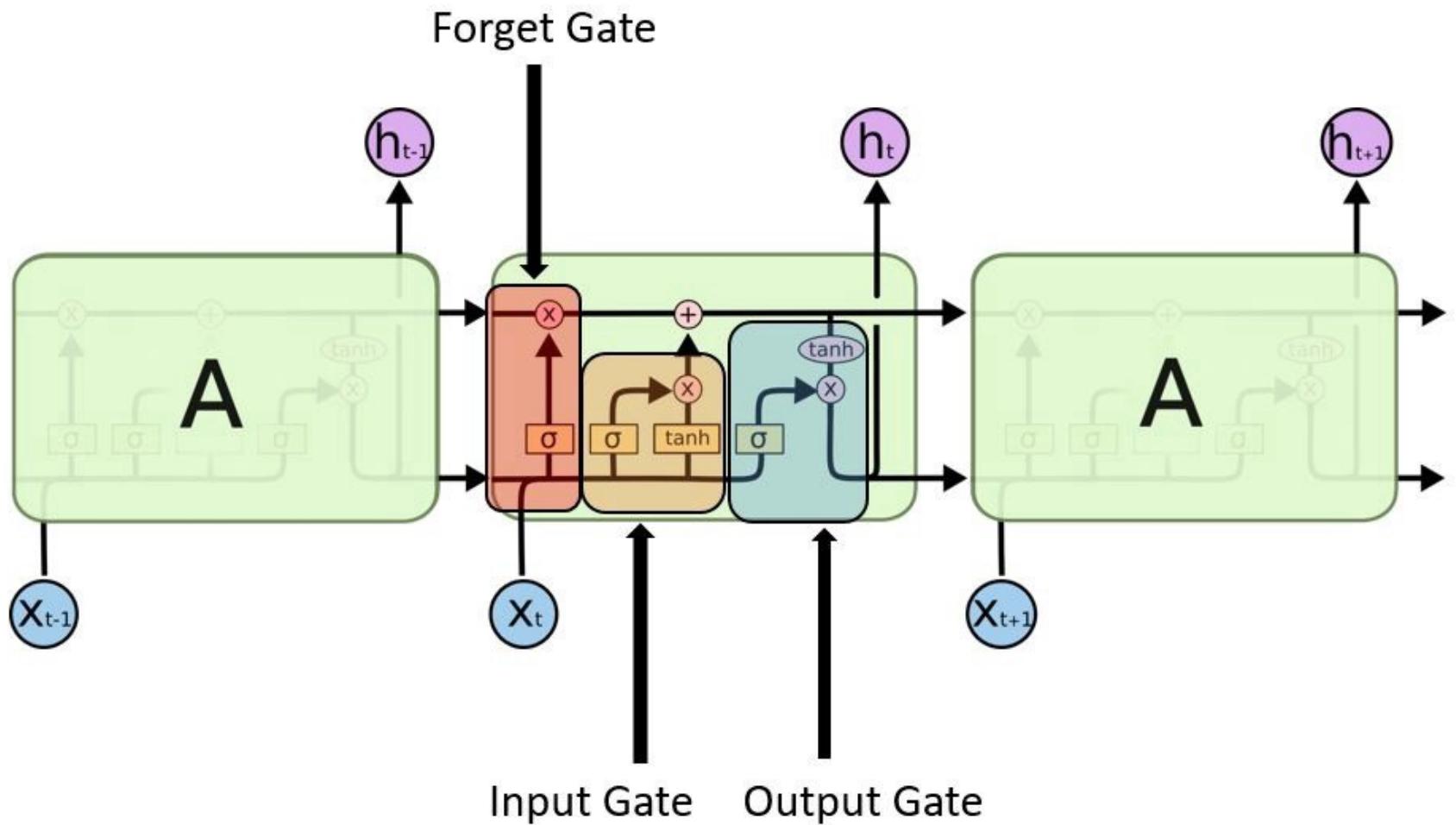
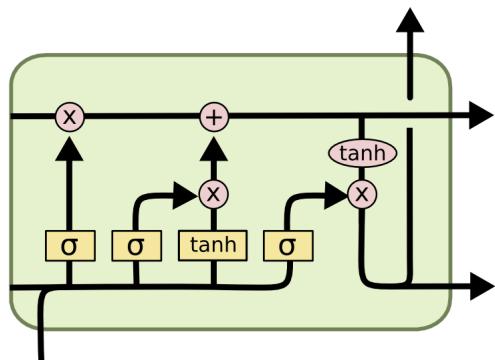


Deep Learning

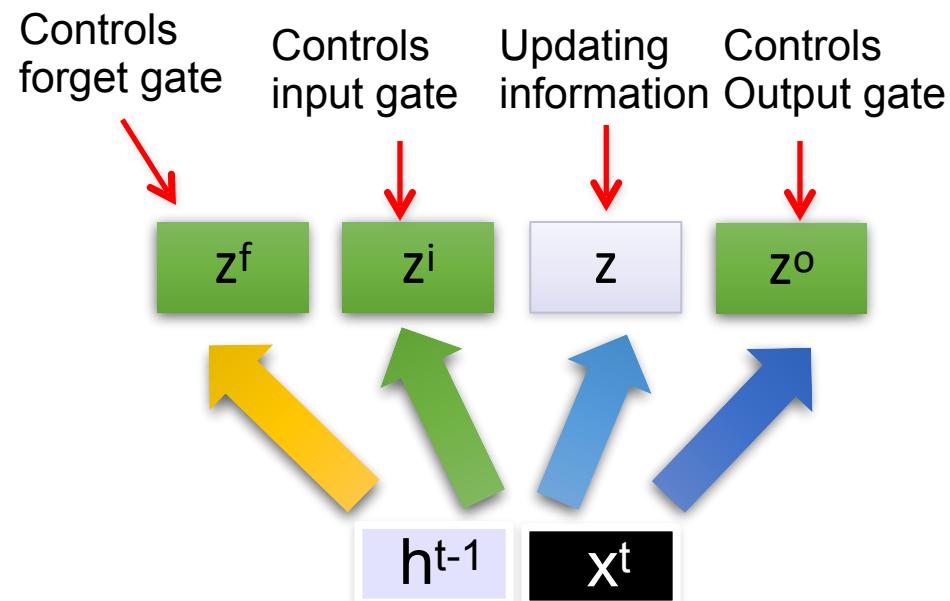
Mayank Vatsa



These 4 matrix computation should be done concurrently.



c^{t-1}



$$z = \tanh(W h^{t-1} + x^t)$$

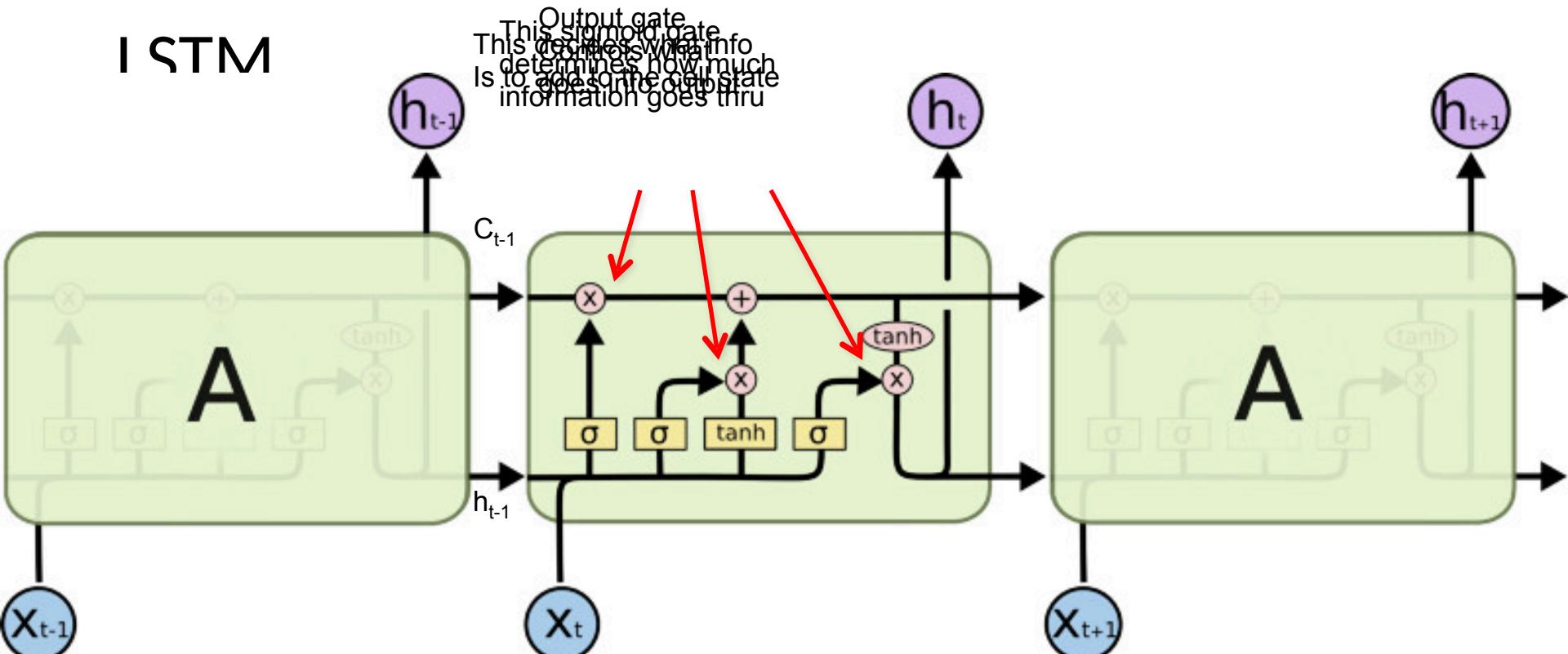
$$z^i = \sigma(W_i h^{t-1} + x^t)$$

$$z^f = \sigma(W_f h^{t-1} + x^t)$$

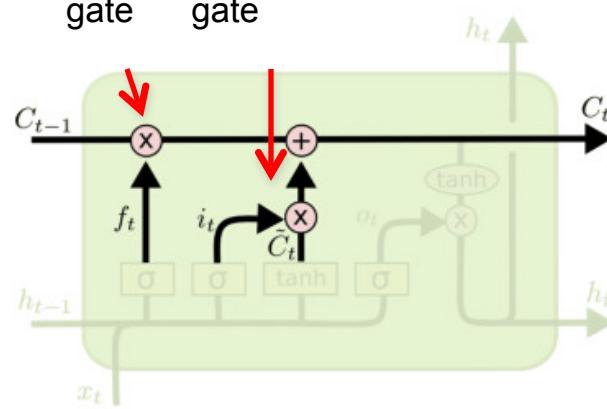
$$z^o = \sigma(W_o h^{t-1} + x^t)$$

Information flow of LSTM

I LSTM

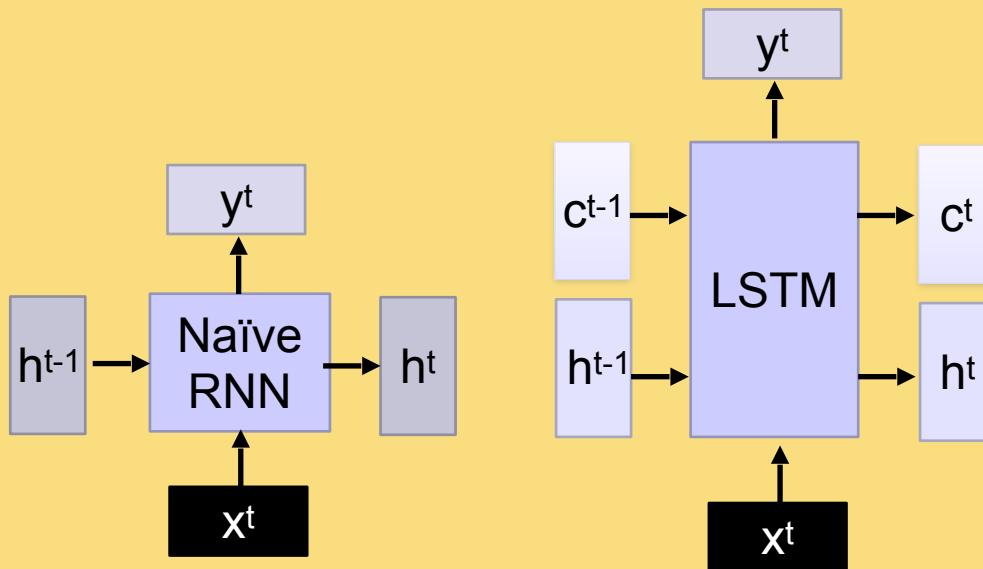


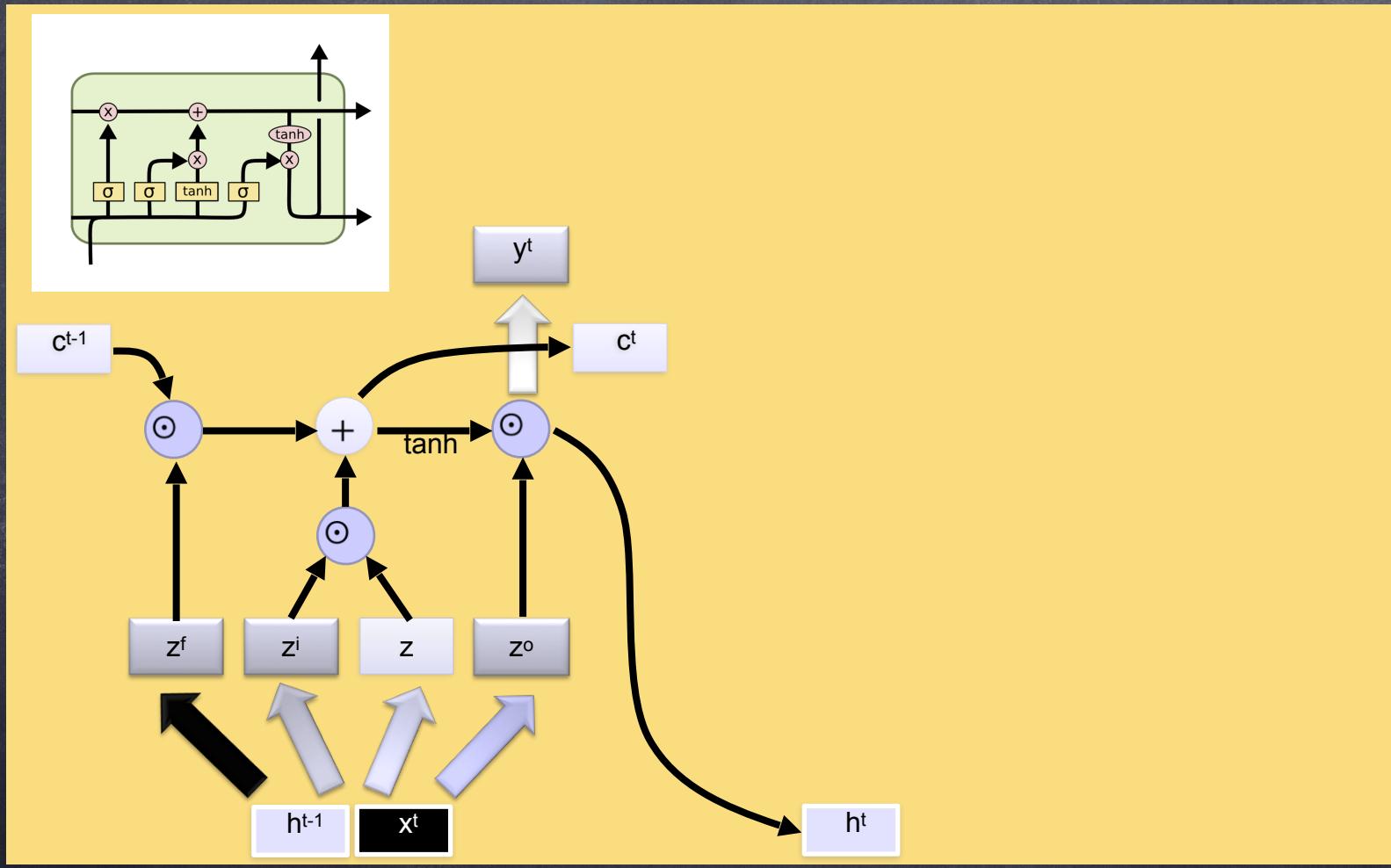
The core idea is this cell state C_t , it is changed slowly, with only minor linear interactions. It is very easy for vanishing gradient problem in information to flow along it unchanged. LSTM is handled already. ReLU replaces tanh

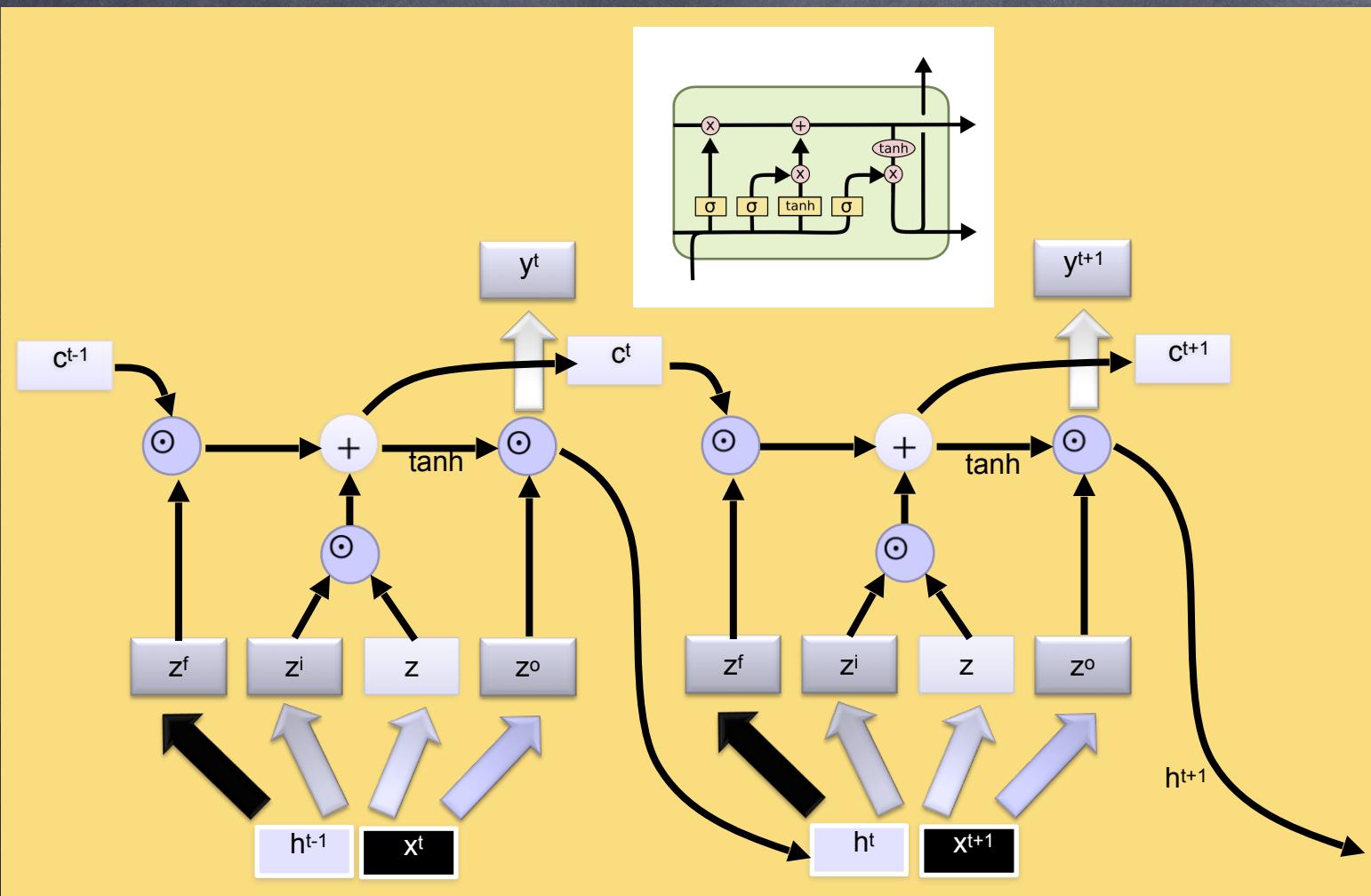


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

RNN vs LSTM







Applications of LSTM / RNN

a train traveling down a track next to a forest.



a group of young boys playing soccer on a field.



Image caption generation using attention (From CY Lee lecture)

z^0 is initial parameter, it is also learned

A vector for each region

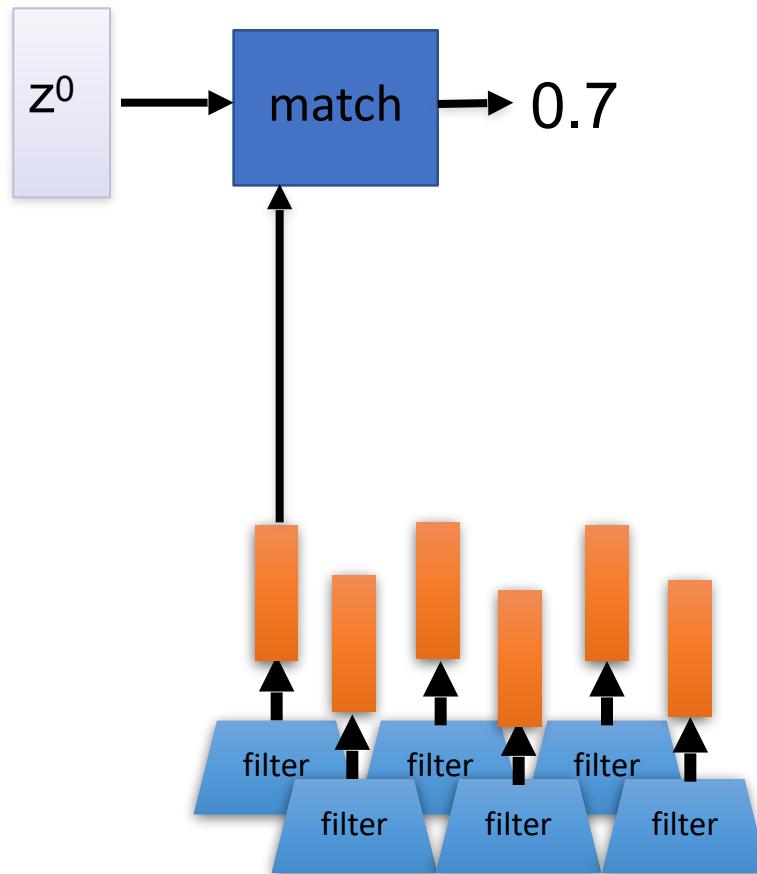
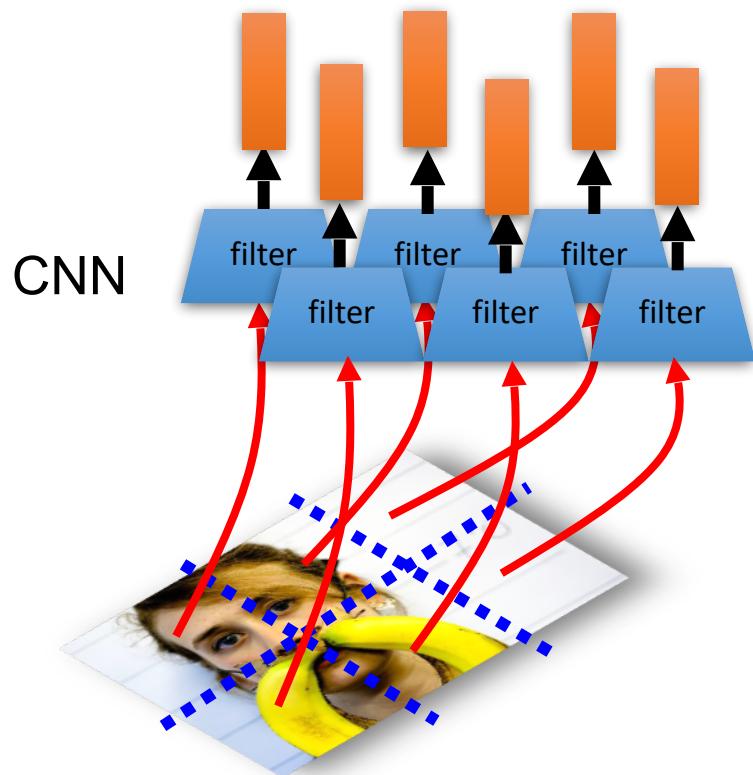


Image Caption Generation

A vector for each region

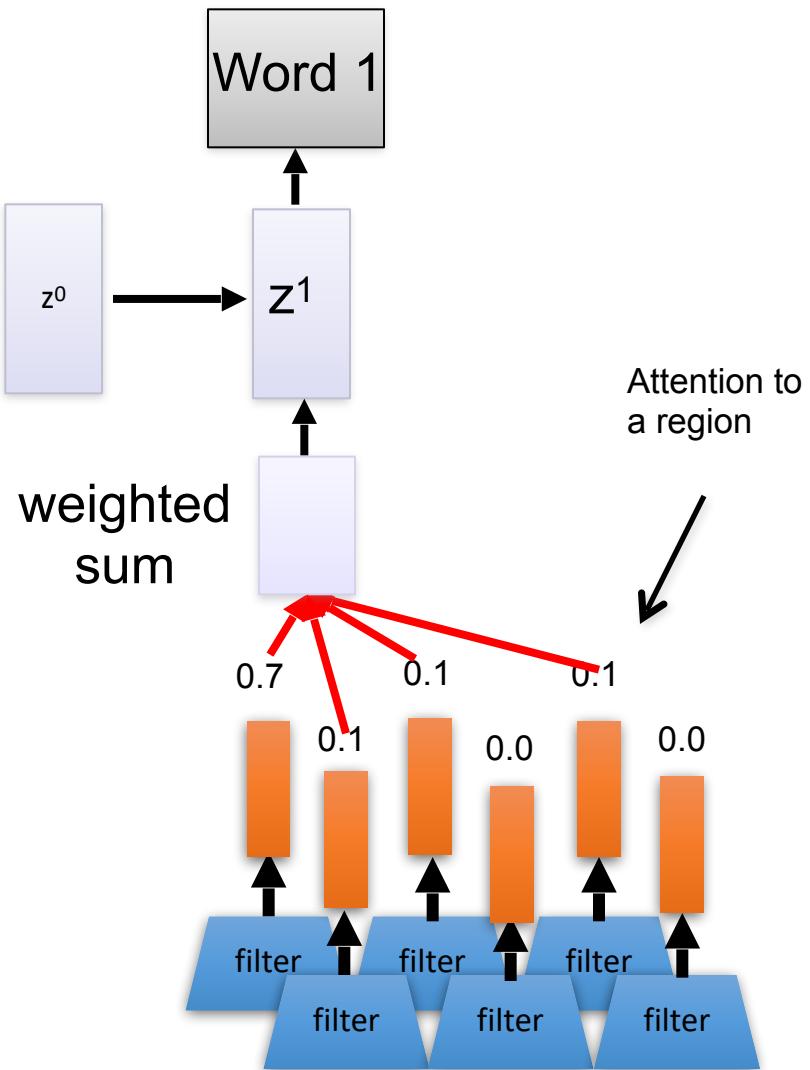
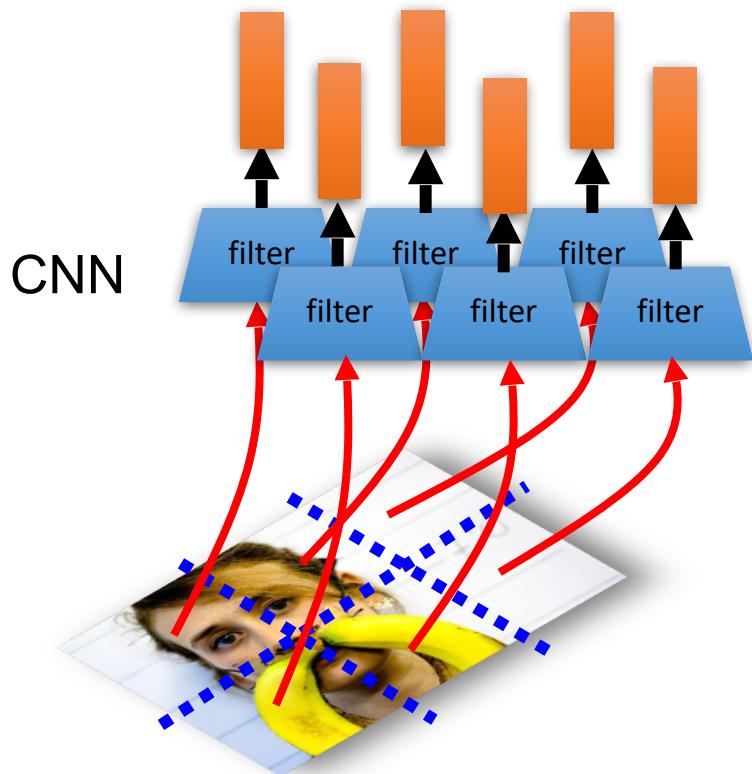


Image Caption Generation

A vector for each region

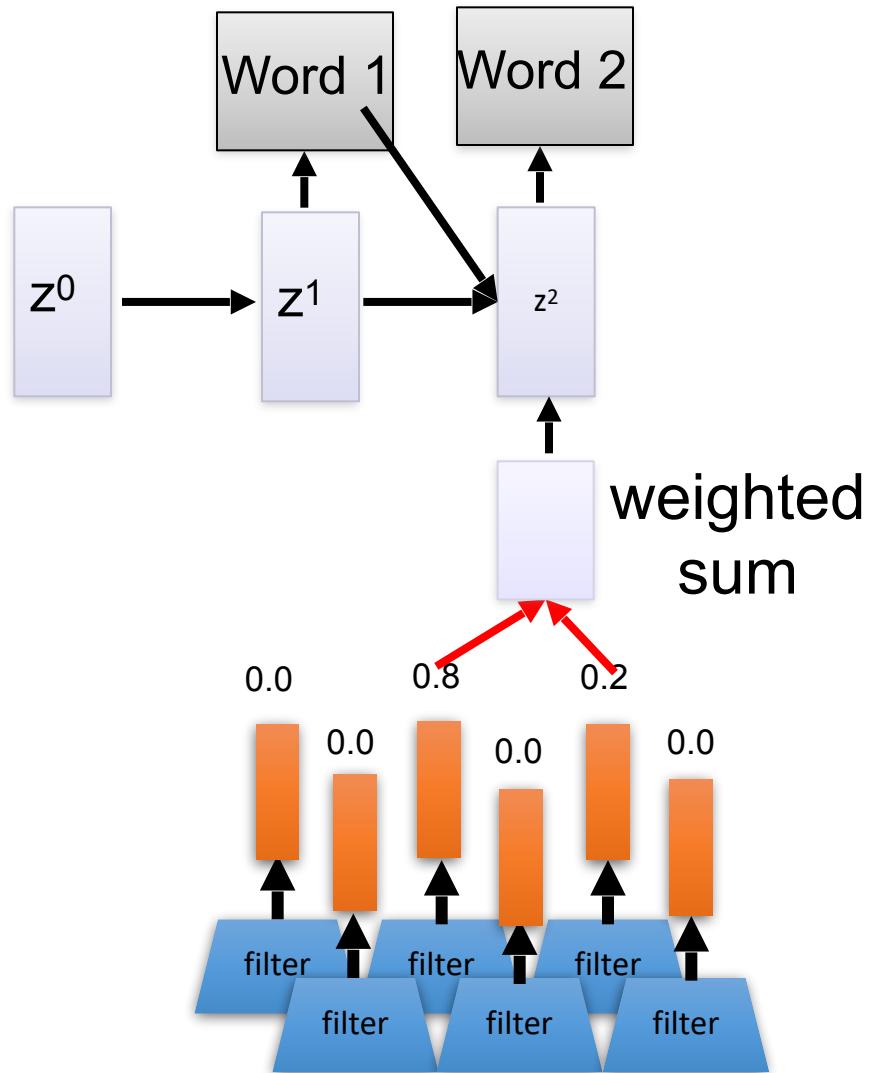
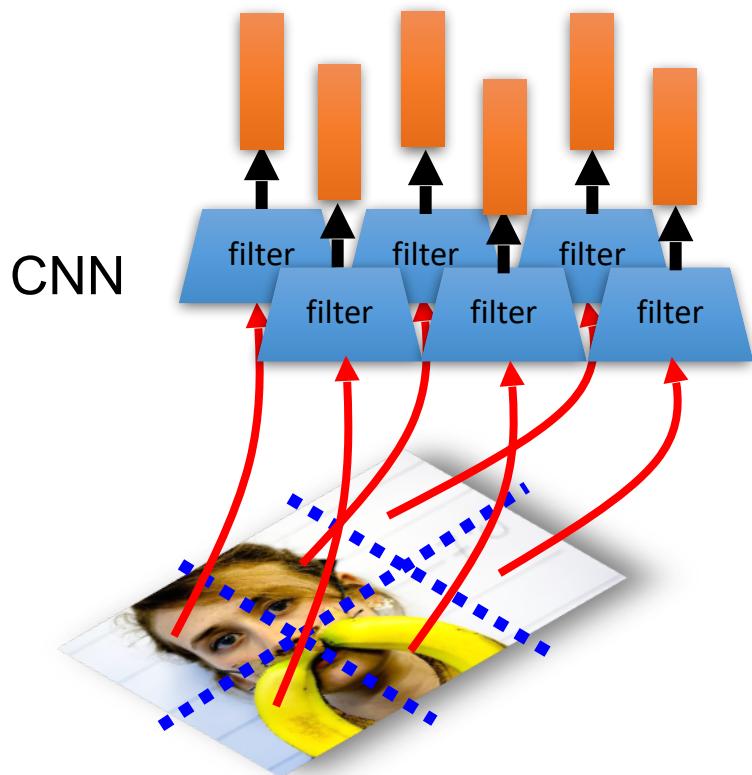


Image Caption Generation



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, Yoshua Bengio, “Show, Attend and Tell: Neural Image Caption Generation with Visual Attention”, ICML, 2015

Image Caption Generation



A large white bird standing in a forest.



A woman holding a clock in her hand.



A man wearing a hat and a hat on a skateboard.



A person is standing on a beach with a surfboard.



A woman is sitting at a table with a large pizza.

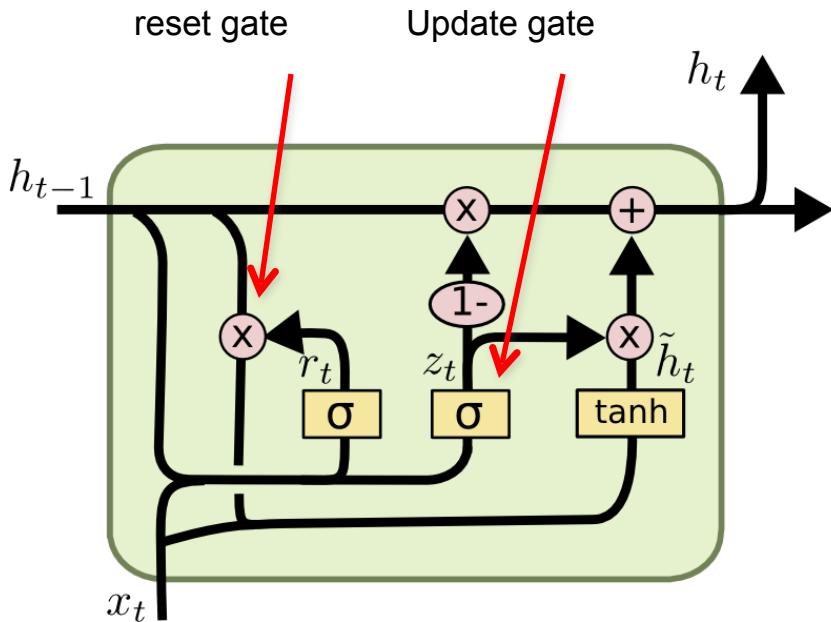
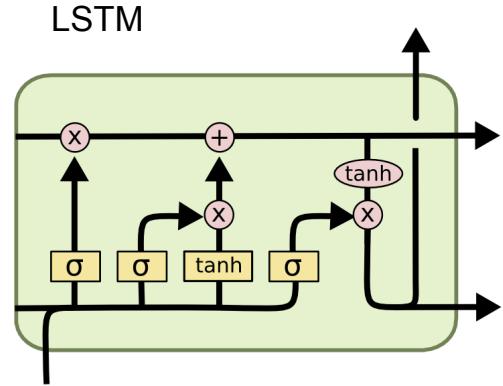


A man is talking on his cell phone while another man watches.

Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, Yoshua Bengio, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML, 2015

GRU – gated recurrent unit

(more compression)



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

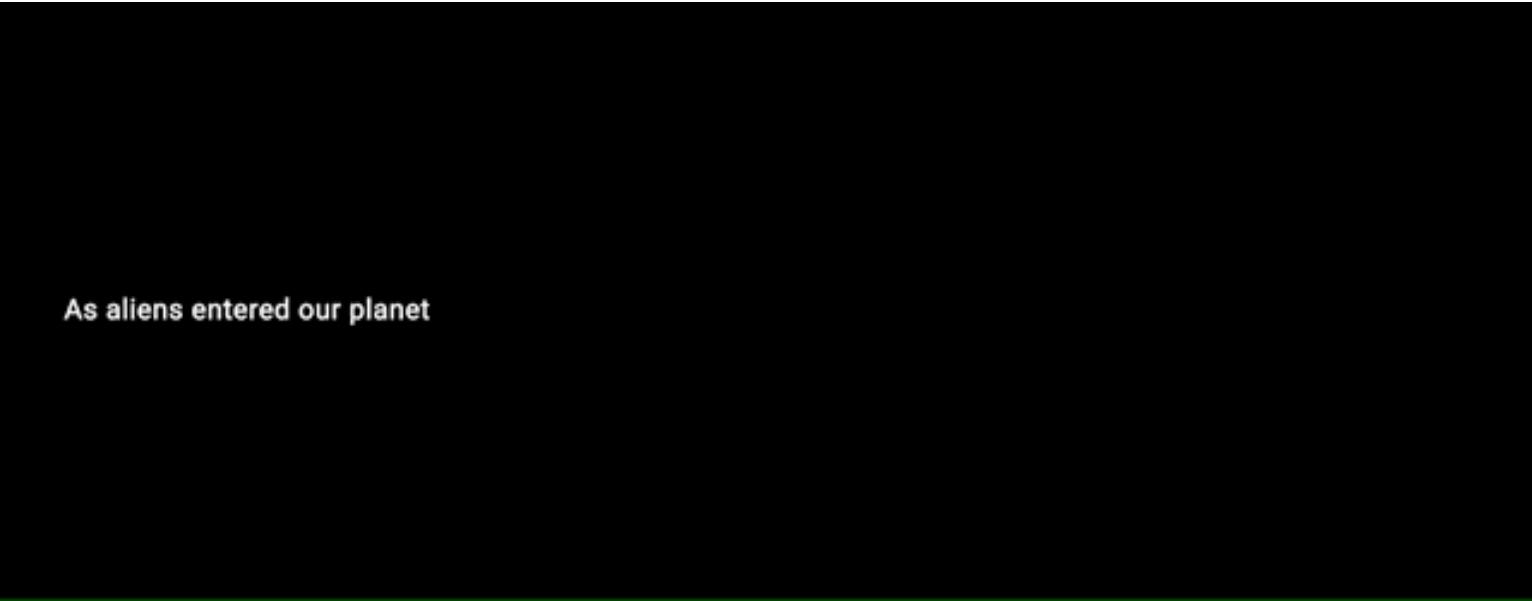
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

It combines the **forget** and **input** into a single **update gate**.

It also merges the cell state and hidden state. This is simpler than LSTM. There are many other variants too.

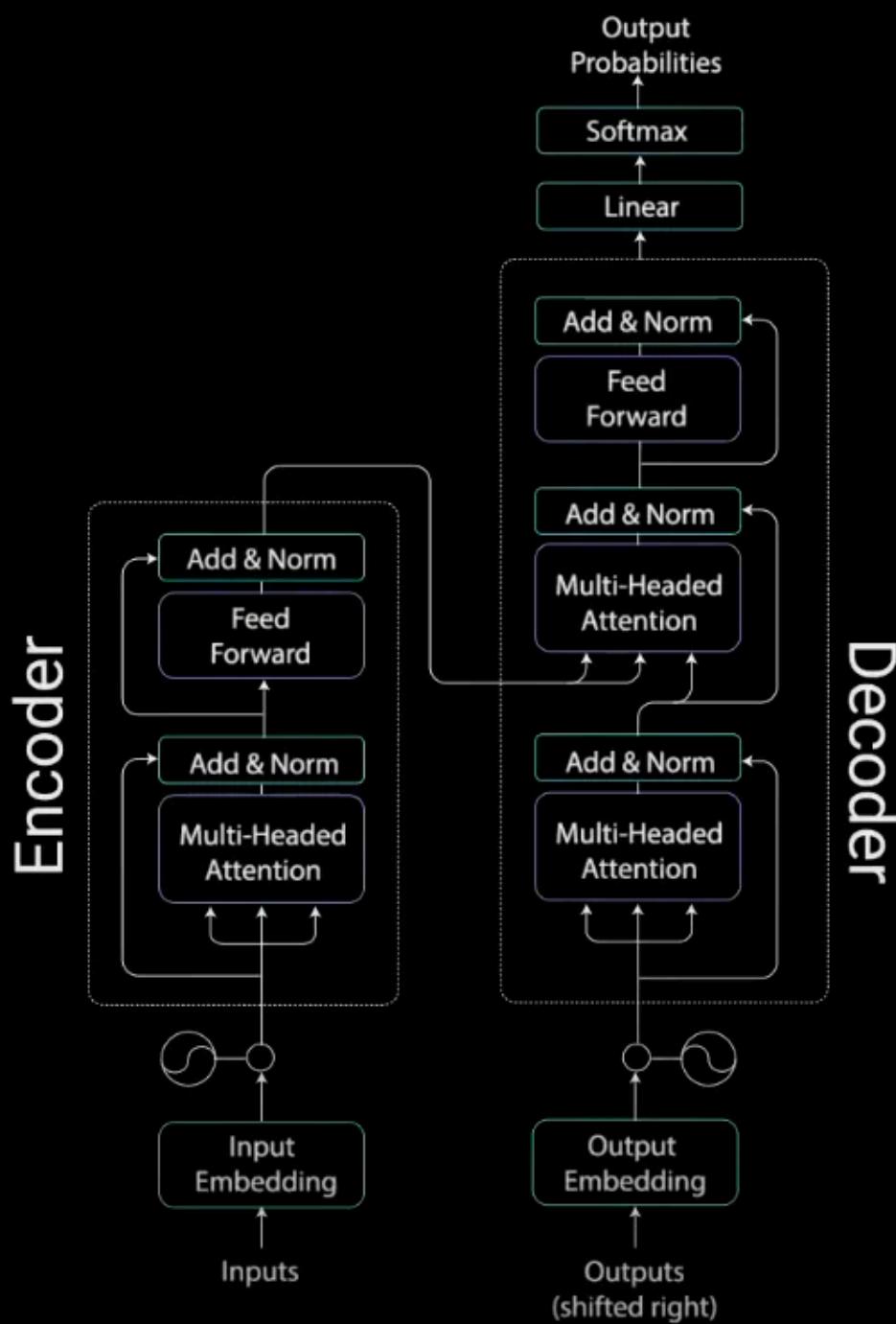
More advancements

- Transformers (LSTMs are dead - long live Transformers)
- **Input:** “As Aliens entered our planet”.
- **Transformer output:** “and began to colonized Earth, a certain group of extraterrestrials began to manipulate our society through their influences of a certain number of the elite to keep and iron grip over the populace.”



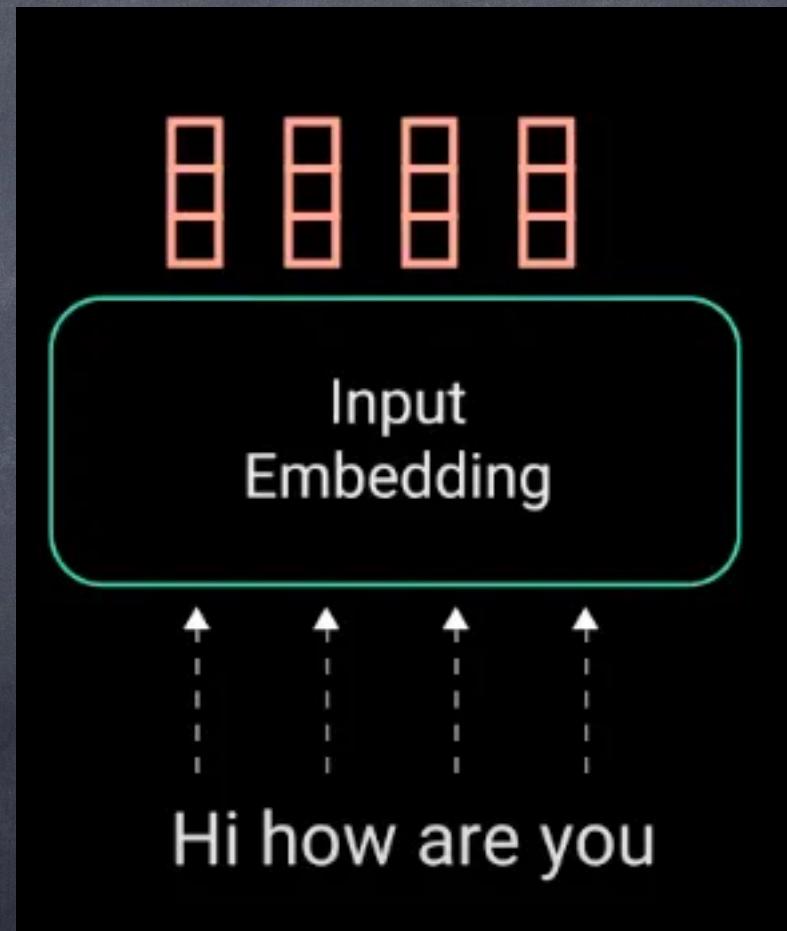
As aliens entered our planet

Transformer



Input

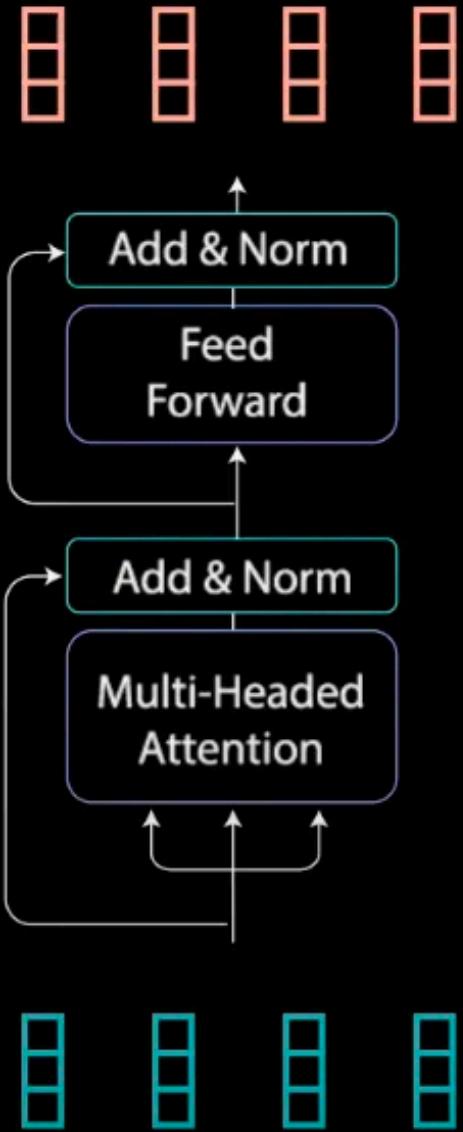
- Input embedding
- Positional encoding



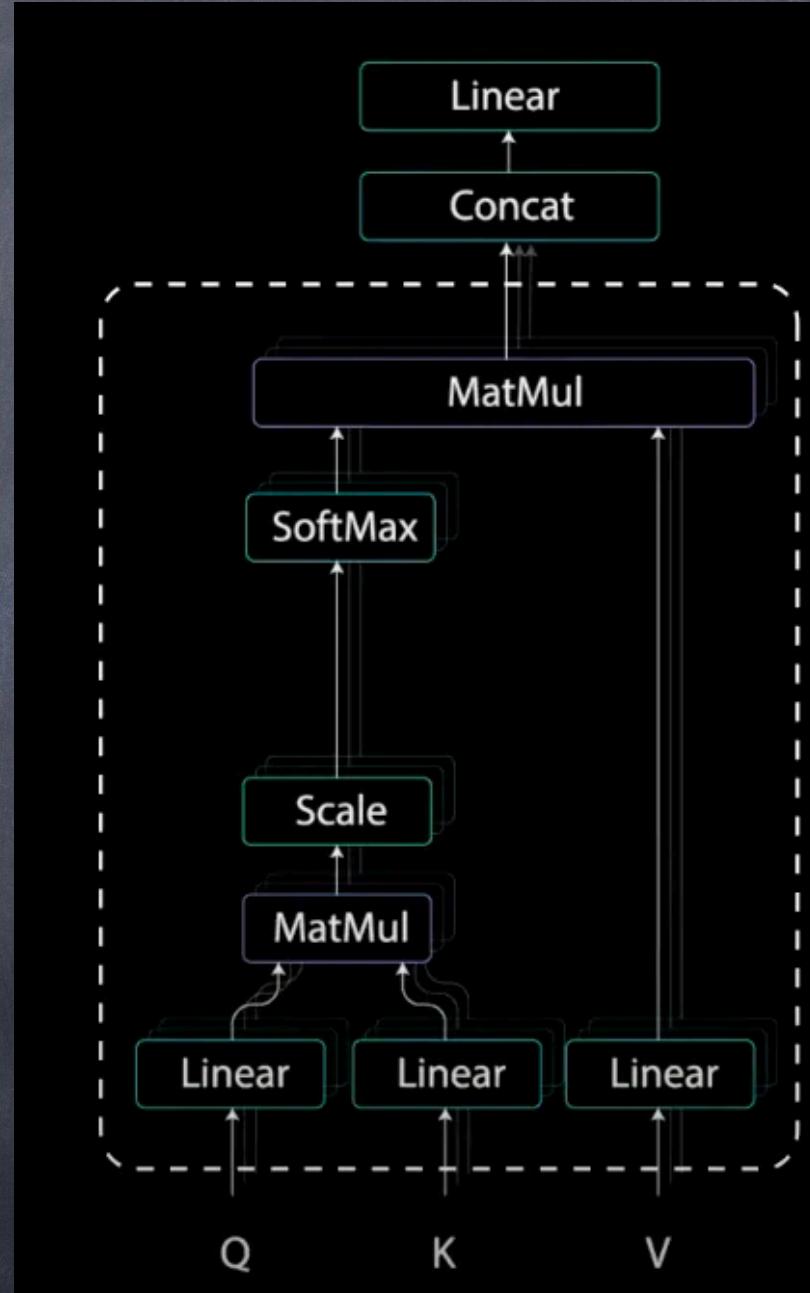
- Encoder
- multi-headed attention, followed by a fully connected network
- There are residual connections around each of the two sublayers followed by a layer normalization

Encoder Input Representation

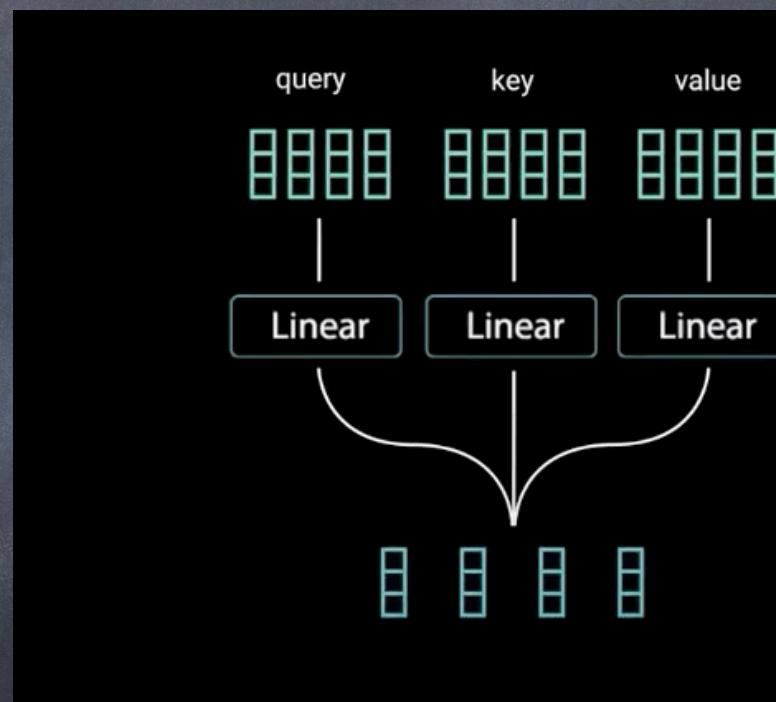
Positional Input Embedding



- Multi-Headed Attention
- Query, Key, and Value Vectors
- "The query key and value concept come from retrieval systems. For example, when you type a query to search for some video on Youtube, the search engine will map your query against a set of keys (video title, description etc.) associated with candidate videos in the database, then present you the best matched videos (values)."



- Dot Product of Query and Key and score is obtained
- The score matrix determines how much focus should a word be put on other words. So each word will have a score that corresponds to other words in the time-step. The higher the score the more focus.



Hi how are you

Hi

98

27

10

12

how

27

89

31

67

are

10

31

91

54

you

12

67

54

92

- Scaling Down the Attention Scores
- square root of the dimension of query and key

$$\frac{\text{Red Grid}}{\sqrt{d_k}} = \text{Scaled Scores}$$

The diagram illustrates the scaling of attention scores. It shows a red 4x4 grid representing raw attention scores. This grid is divided by a horizontal line, and the resulting value is multiplied by the square root of d_k (indicated by a cyan square root symbol) to produce a cyan scaled scores grid. The text "Scaled Scores" is written next to the cyan grid.

- Softmax of the Scaled Scores

- By doing a softmax the higher scores get heightened, and lower scores are depressed. This allows the model to be more confident about which words to attend too.

Softmax()=

	Hi	how	are	you
Hi	0.7	0.1	0.1	0.1
how	0.1	0.6	0.2	0.1
are	0.1	0.3	0.6	0.1
you	0.1	0.3	0.3	0.3

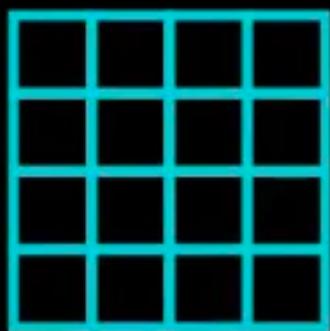
$$\text{softmax}(x)_i = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$

- Multiply Softmax Output with Value vector
- The higher softmax scores will keep the value of words the model learns is more important
- The lower scores will drown out the irrelevant words
- Then you feed the output of that into a Linear layer to process.

attention weights

value

output



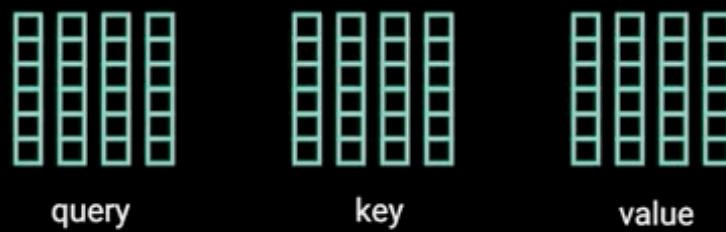
x



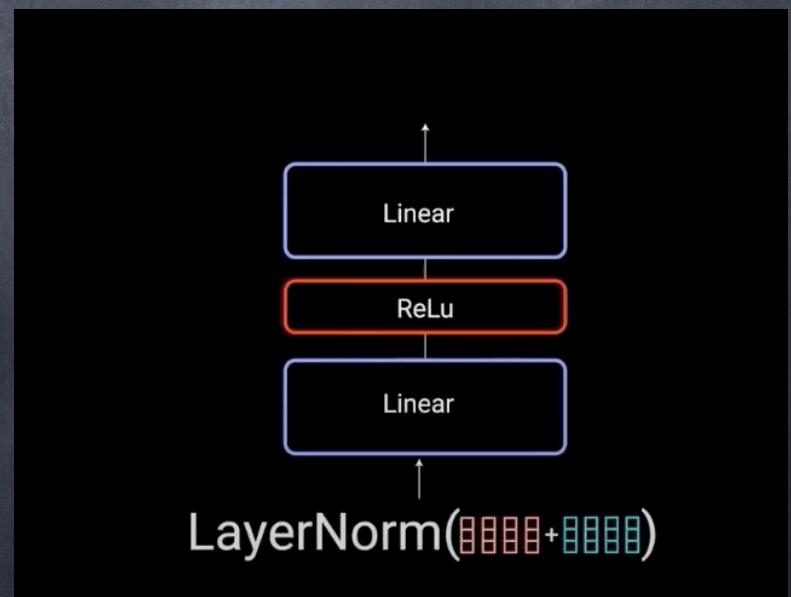
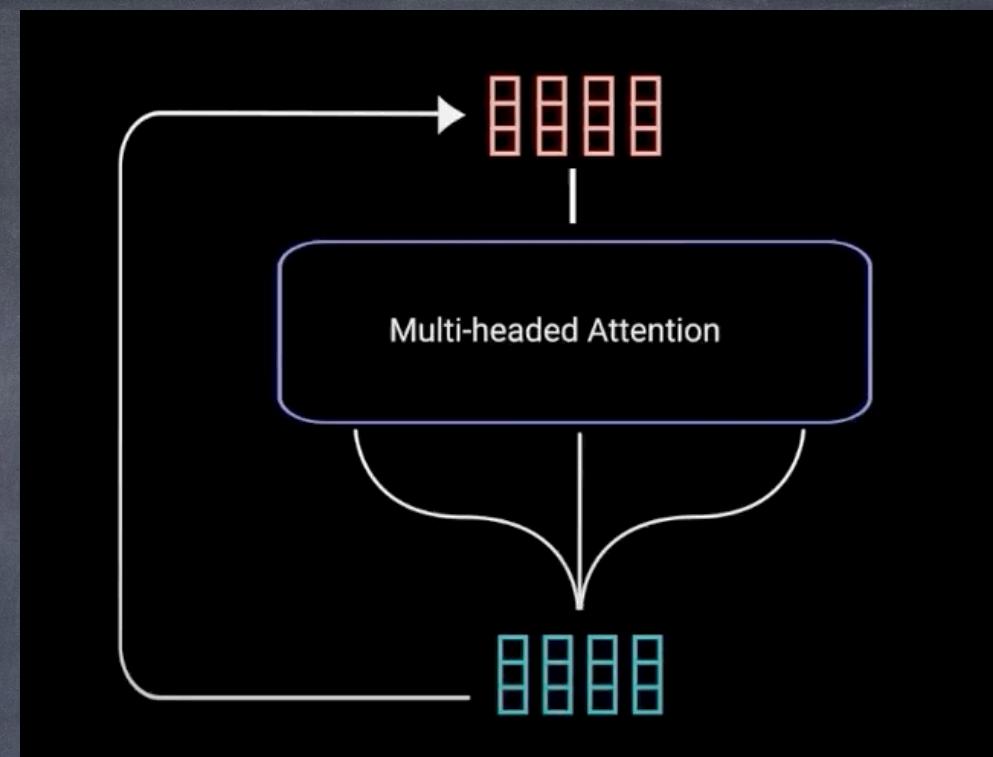
=



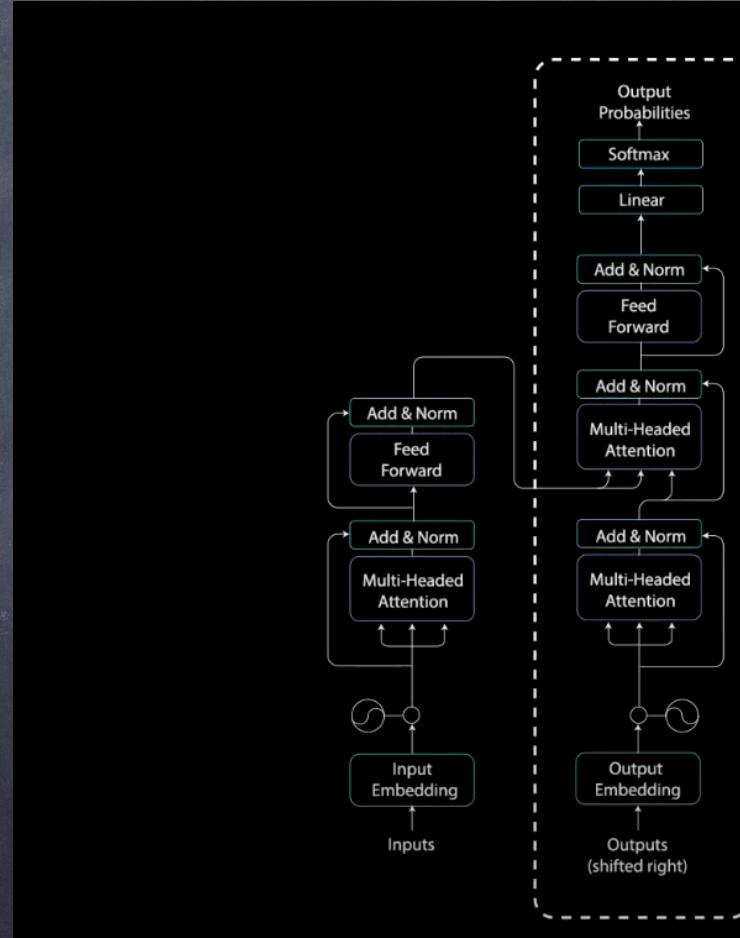
• Computing Multi-headed Attention



• The Residual Connections, Layer Normalization, and Feed Forward Network

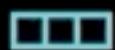


Decoder



HI, |||
how |||
are |||
you? |||

Transformers Decoder



<start>

Decoder Input Embeddings & Positional Encoding

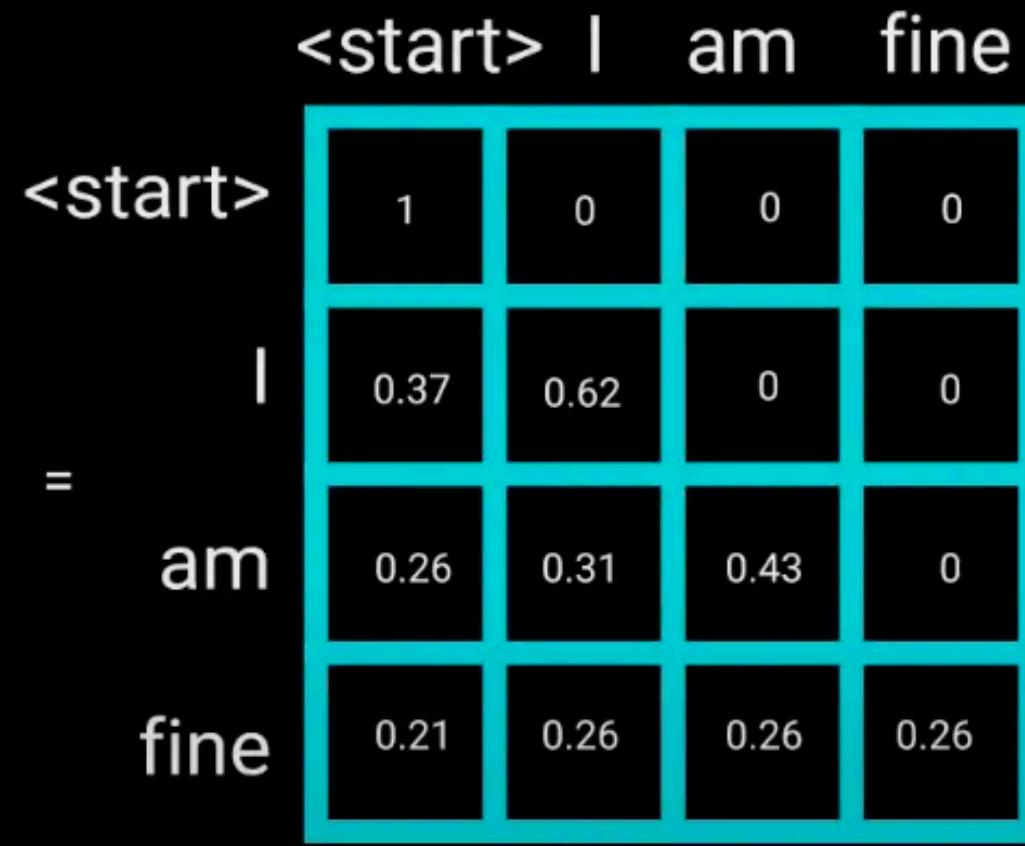
Decoders First Multi-Headed Attention

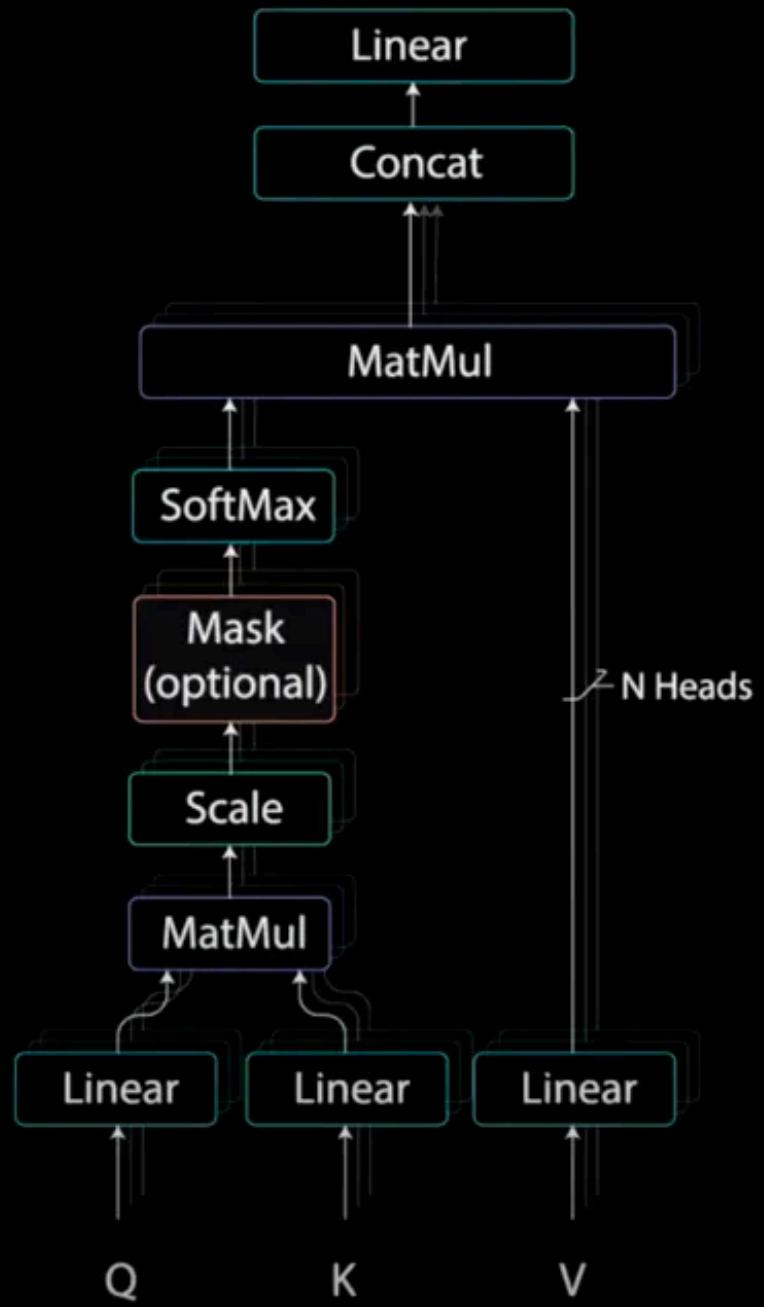
	<start>	I	am	[fine]
<start>	0.7	0.1	0.1	0.1
I	0.1	0.6	0.2	0.1
am	0.1	0.3	0.6	0.1
fine	0.1	0.3	0.3	0.3

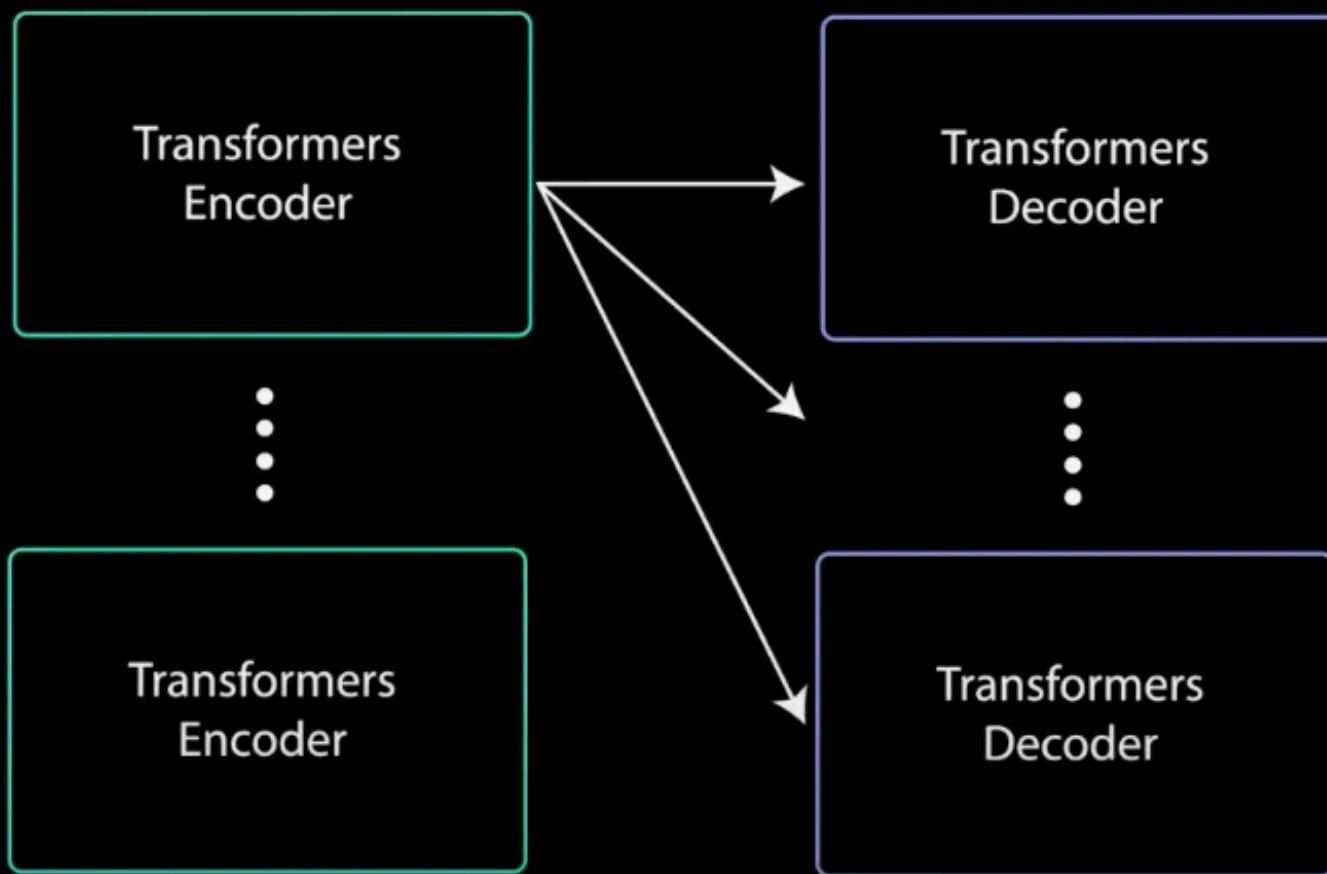
$$\begin{array}{c} \text{Scaled Scores} \\ \begin{matrix} 0.7 & 0.1 & 0.1 & 0.1 \\ 0.1 & 0.6 & 0.2 & 0.1 \\ 0.1 & 0.3 & 0.6 & 0.1 \\ 0.1 & 0.3 & 0.3 & 0.3 \end{matrix} \\ + \\ \text{Look-Ahead Mask} \\ \begin{matrix} 0 & -\inf & -\inf & -\inf \\ 0 & 0 & -\inf & -\inf \\ 0 & 0 & 0 & -\inf \\ 0 & 0 & 0 & 0 \end{matrix} \\ = \\ \text{Masked Scores} \\ \begin{matrix} 0.7 & -\inf & -\inf & -\inf \\ 0.1 & 0.6 & -\inf & -\inf \\ 0.1 & 0.3 & 0.6 & -\inf \\ 0.1 & 0.3 & 0.3 & 0.3 \end{matrix} \end{array}$$

$$\text{Softmax} \left(\begin{array}{cccc} 0.7 & -\inf & -\inf & -\inf \\ 0.1 & 0.6 & -\inf & -\inf \\ 0.1 & 0.3 & 0.6 & -\inf \\ 0.1 & 0.3 & 0.3 & 0.3 \end{array} \right)$$

) =







And That's it ...

- Next class: Multi-task Learning and Transfer Learning