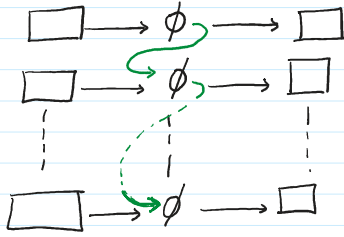


## Problem with RNN

Formula  $h_t = \phi(W_h h_{t-1} + W_x x_t + b)$

some  
Activation  
function

At time  $t$ , take input  $x_t$  & previous hidden state  $h_{t-1}$  to produce  $h_t$

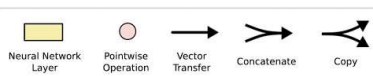
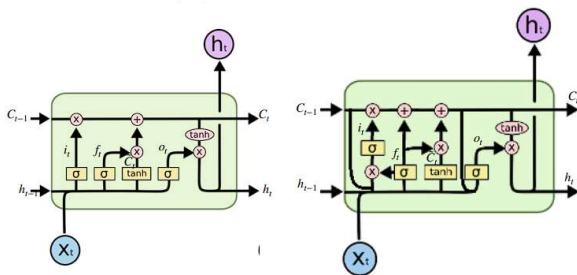


During Backpropagation  
you keep multiplying  
numbers that  $< 1$ ,  
gradient  $\rightarrow 0$

(Vanishing Gradient  
Problem)

Recurrent Neural Network

A standard LSTM as discussed in the lectures is shown below in left figure. We proposed a new LSTM which is shown block in the right figure



Write the equations gate ( $\hat{C}_t$ ), input gate ( $i_t$ ), output gate ( $o_t$ ), forget gate ( $f_t$ ), cell update ( $C_t$ ) and state update ( $h_t$ )

## LSTMS

At  $t$ , given

①  $x_t \in \mathbb{R}^D$

②  $h_{t-1} \in \mathbb{R}^h$

Compute

① Forget Gate :-  $f_t \in (0, 1)^h$

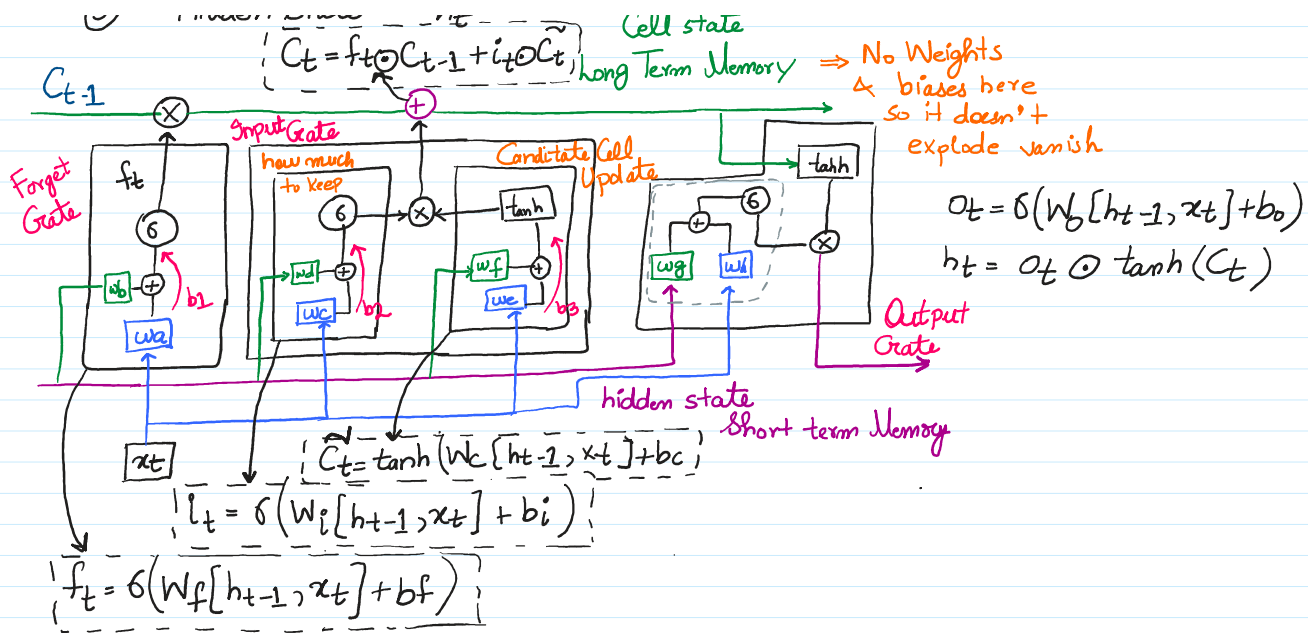
② Input Gate :-  $i_t \in (0, 1)^h$

③ Candidate Gate :-  $\hat{C}_t \in (-1, +1)^h$

④ Output Gate :-  $o_t \in (0, 1)^h$

⑤ Hidden State  $h_t$

$C_t = f_t \odot C_{t-1} + i_t \odot \hat{C}_t$  Cell state  
long Term Memory  $\Rightarrow$  No Weights & biases here



$$\sigma = \frac{e^x}{e^x + 1}$$

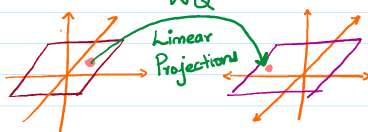
$$\tanh = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

We cannot parallelize an LSTM  
Very long Range dependencies fade

## Attention

①  $X \in \mathbb{R}^{n \times d}$   $\rightarrow$  Model dimension

②  $W_Q, W_K, W_V \in \mathbb{R}^{d \times d_k}$

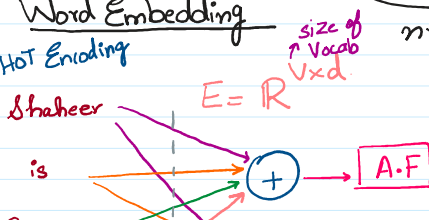


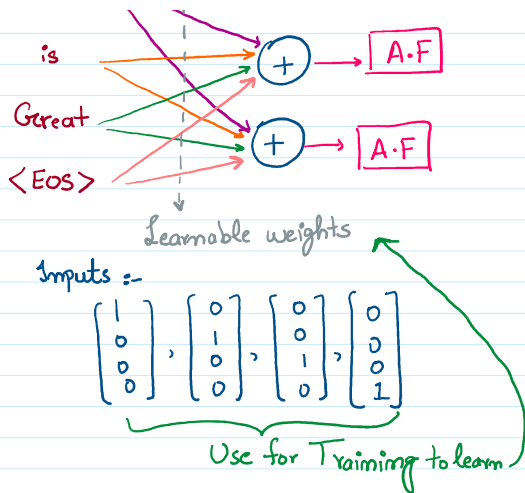
③  $Q = W_Q X \quad K = X W_K \quad V = X W_V$

$Q, K, V \in \mathbb{R}^{n \times d_k}$   
because  $n \times d \quad d \times d_k$   $\xrightarrow{\text{Matrix Multiplication Rules}}$   $n \times d_k$

## Word Embedding

ONE HOT Encoding





## Positional Embedding

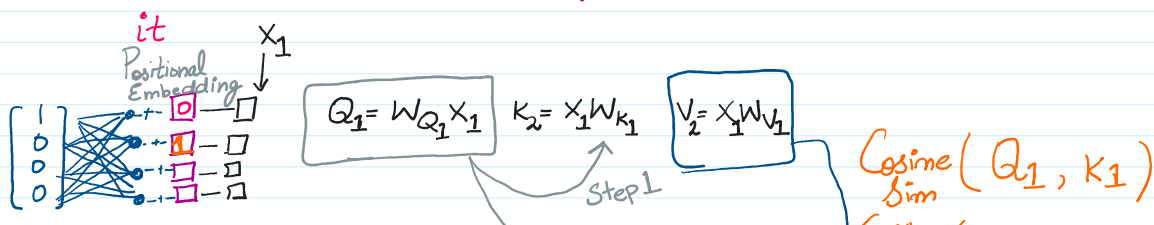
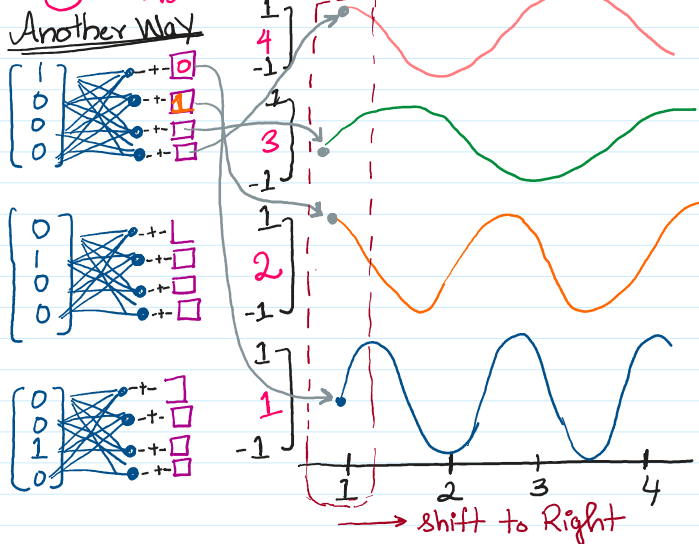
One Way is Sinusoidal

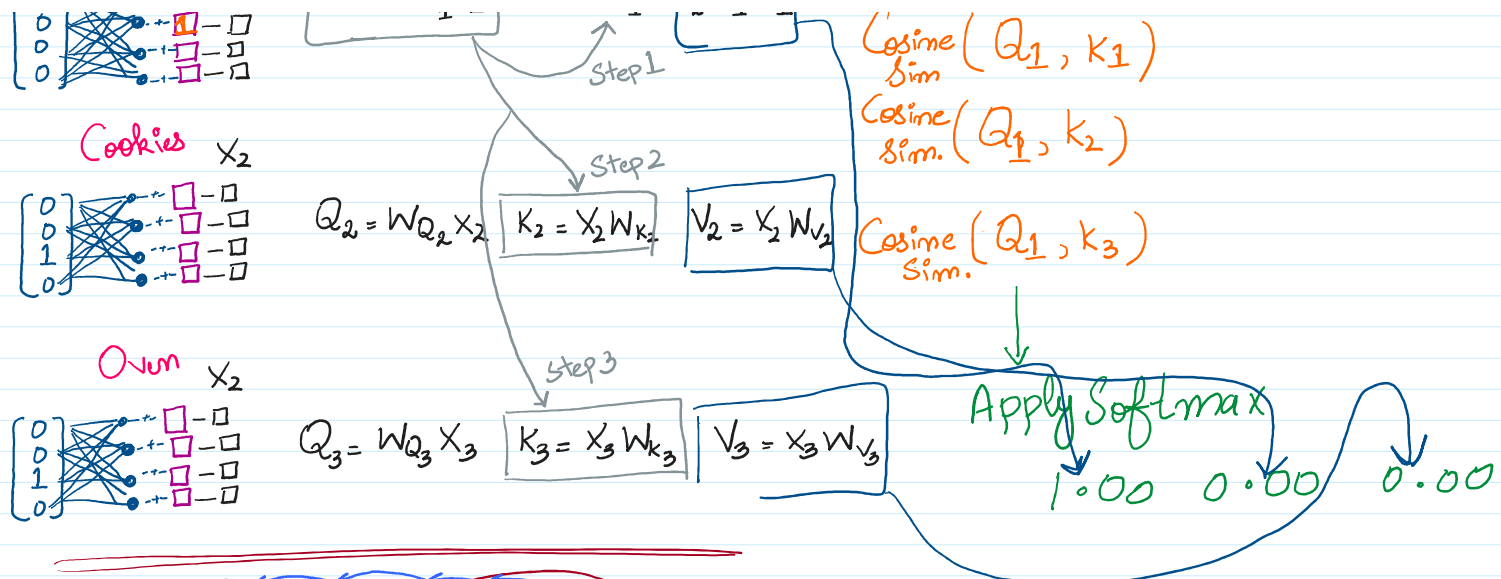
For Positions  $t = 0, 1, 2 \dots n-1$   
& dimension  $d = 0, 1, 2 \dots d_{model}-1$

$$PE_{t,i} = \begin{cases} \sin\left(\frac{t}{10000^{2i/d_{model}}}\right); & i \text{ even} \\ \cos\left(\frac{t}{10000^{2(i-1)/d_{model}}}\right); & i \text{ odd} \end{cases}$$

$$PE \in \mathbb{R}^{n \times d}$$

- ① Unique Patterns
- ② Smooth functions
- ③ No Extra Parameters





Someone took some cookies out of the oven & it  
smelt so Yum

Self attention

'it' & 'cookies' should have a higher cosine similarity

$$S = QK^T$$

$$\hat{S} = \frac{S}{\sqrt{d_k}}$$

$$A_{ij} = \frac{e^{\hat{S}_{ij}}}{\sum_j \exp(\hat{S}_{ij})}$$

$$Z = AV$$

$$Z_i = \sum_j A_{ij} V_j$$

## Diffusion Models

### ① Forward Process

$$q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t} x_{t-1}, \beta_t I)$$

Base Case =

$$q(x_1 | x_0) = \mathcal{N}(x_1; \sqrt{\alpha_1} x_0, \beta_1 I)$$

Where does  $q(x_t | x_0)$

$$q(x_{t-1} | x_0) = \mathcal{N}(x_{t-1}; \sqrt{\alpha_{t-1}} x_0, (1 - \alpha_{t-1}) I)$$

$$q(x_{t-1} | x_0) = \mathcal{N}(x_{t-1}; \sqrt{\bar{\alpha}_{t-1}} x_0, (1 - \bar{\alpha}_{t-1}) I)$$

where

$$\bar{\alpha}_{t-1} = \prod_{s=1}^{t-1} \alpha_s$$

$$= \alpha_1 \alpha_2 \dots \alpha_{t-1}$$