



When you don't use LSTM  
for a long sequence RNN

RNN

LSTM



## Agenda

→ 11:15

1. Data Shape in RNNs, multi-layer RNN
2. Encoder Decoder Architecture
3. LSTM
4. GRU



Hardly used in industry

Transformers

# Data Shallow into RNN

tent

→ ~ 100 won

~ 1K won

~ 10 words

target → multi-class

{Batch-size, no-of-words} → Final Value

Product now → 6000 character level

[n rows, 6000] → [n, 6000, 128] dimension of emball.

200) 64gb Ram

128 gb Ram

median length  $\rightarrow$

lose some data

lose  $\rightarrow$  minimum

Delay + hind + Err

V.U. good Args

LSTM

LPADS

99.7%

g

love

Lasagne

LPADS

z fin

Dotsize  $\rightarrow$  [1000, 10, 128]

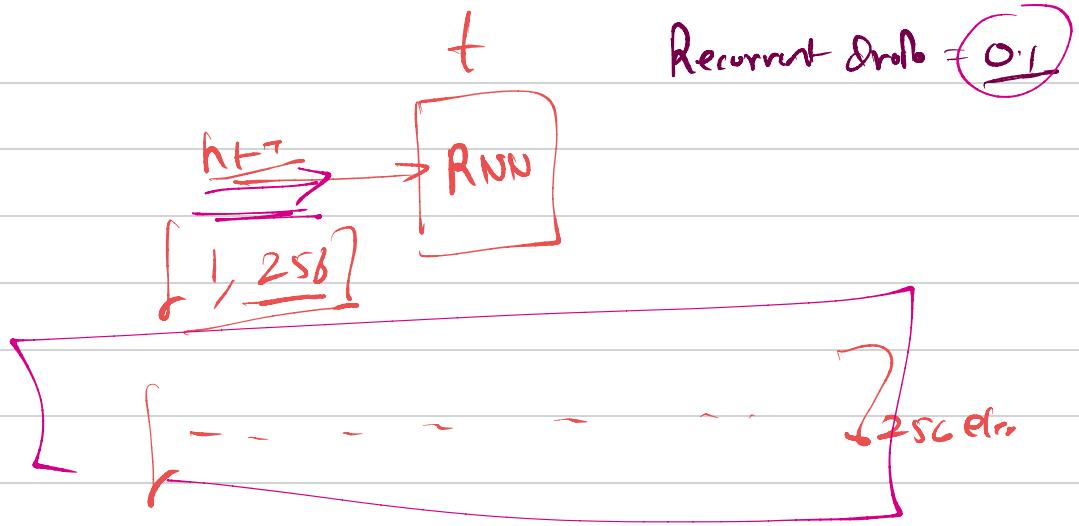
{ word } - 128 emb  
b.size  
1, 10, 128  
max length  
embedding



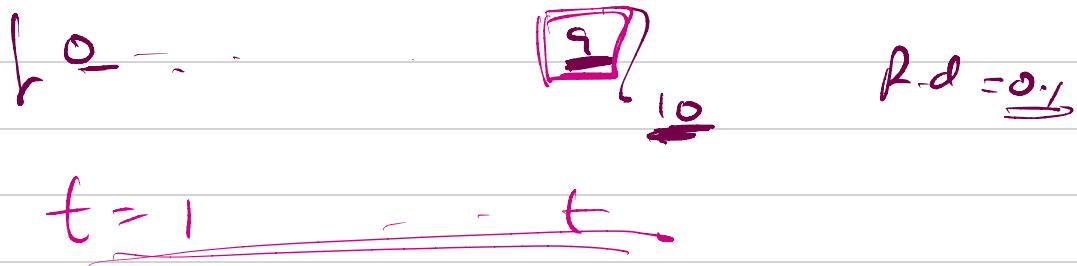
(128)

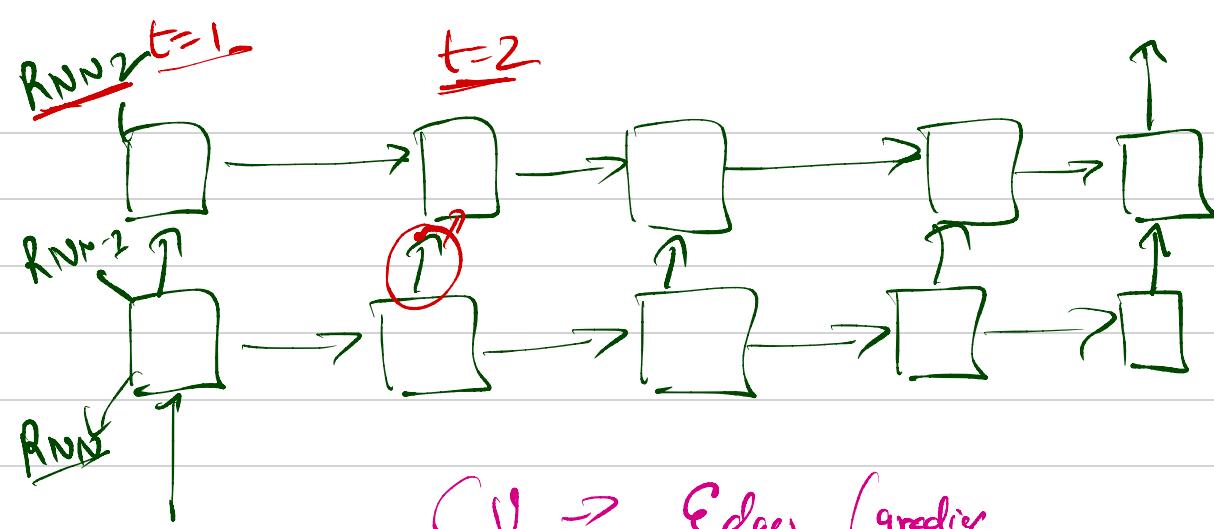
2D

(128)



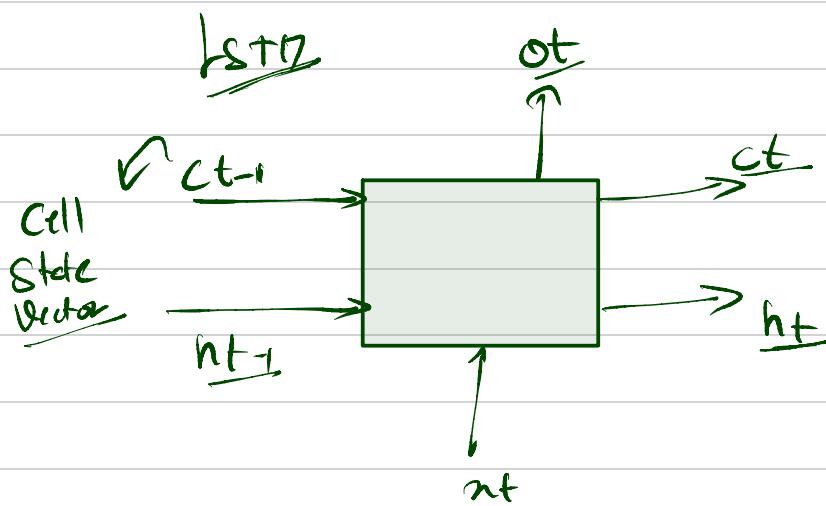
For each time step  $\rightarrow$  the same probability  
of  $\frac{1}{10}$  of elements  
in hidden state  
 $\rightarrow$  drop

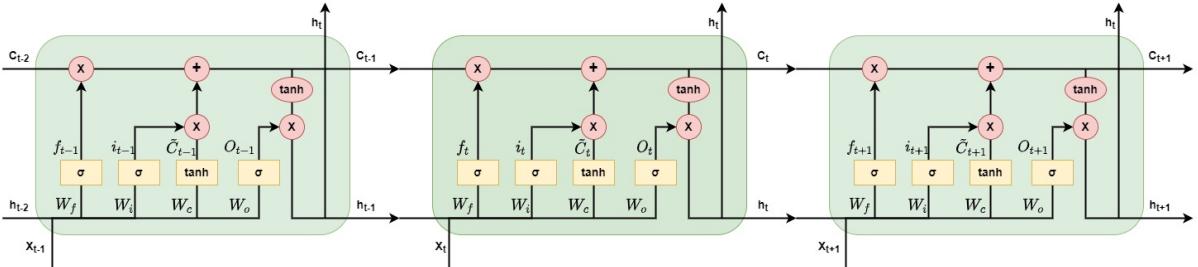
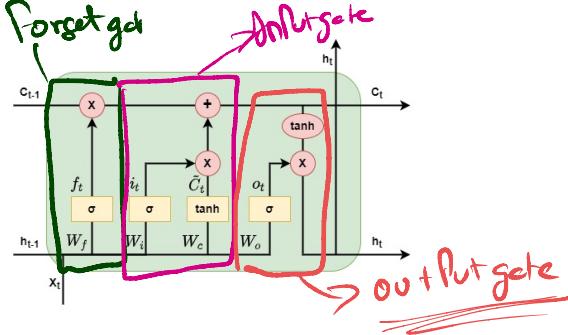




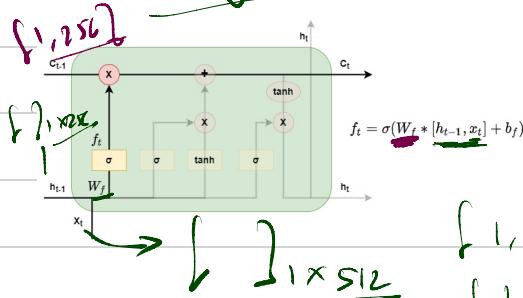
$CV \rightarrow$  Edges / gradients  
Parts of objects

$RNN \rightarrow L_1$  - Small range defn  
 $L_2 \rightarrow$  Bigger Phrase interact





Forget gate



what information in  
"cell state" should retain  
or discard

$$\begin{cases} 1, 256 \\ 1, 512 \end{cases} \downarrow \text{dimis=0}$$

$$= [1, 768] \times [768, 256] = [1, 256]$$

w.r.t

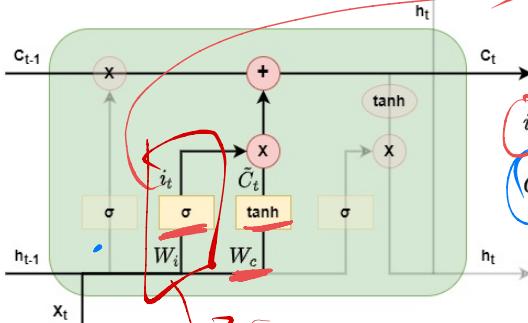
$$\begin{cases} 5, 6, 8, 10 \end{cases} \otimes$$

$$\{0, 0, 1, 2\}$$

$$= \{0, 0, 8, 20\}$$

$$\begin{cases} 1, 256 \end{cases} \text{ o }$$

Input Gate



controlling how much of  $c_t$  to be added to  $c_{t-1}$

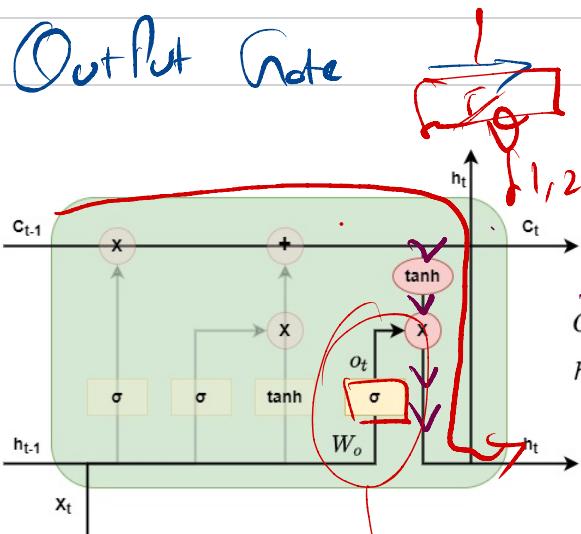
$$i_t = \sigma(W_i * [h_{t-1}, x_t] + b_i)$$

$$\tilde{c}_t = \tanh(W_C * [h_{t-1}, x_t] + b_c)$$

Potential cell state vector

$i_t (0-1)$   
on/off switch

Output Gate



what information should be passed  
 $[1, 256] \times 9^5 \rightarrow [256]$  to next hidden state

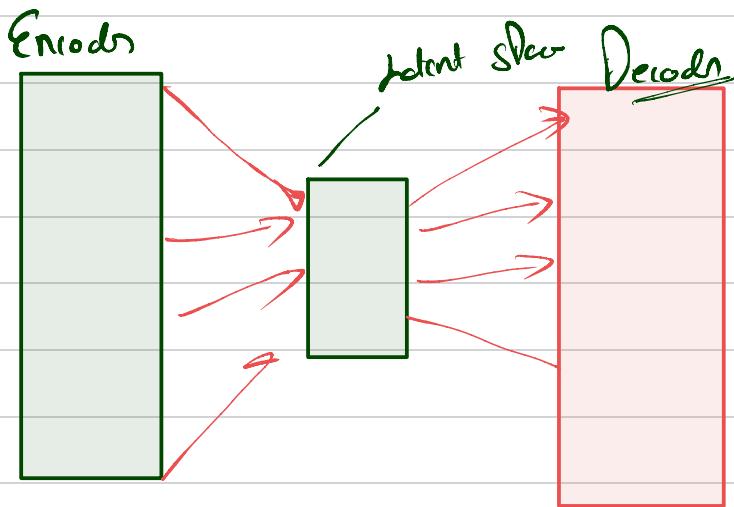
$$o_t = \sigma(W_o * [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(c_t)$$

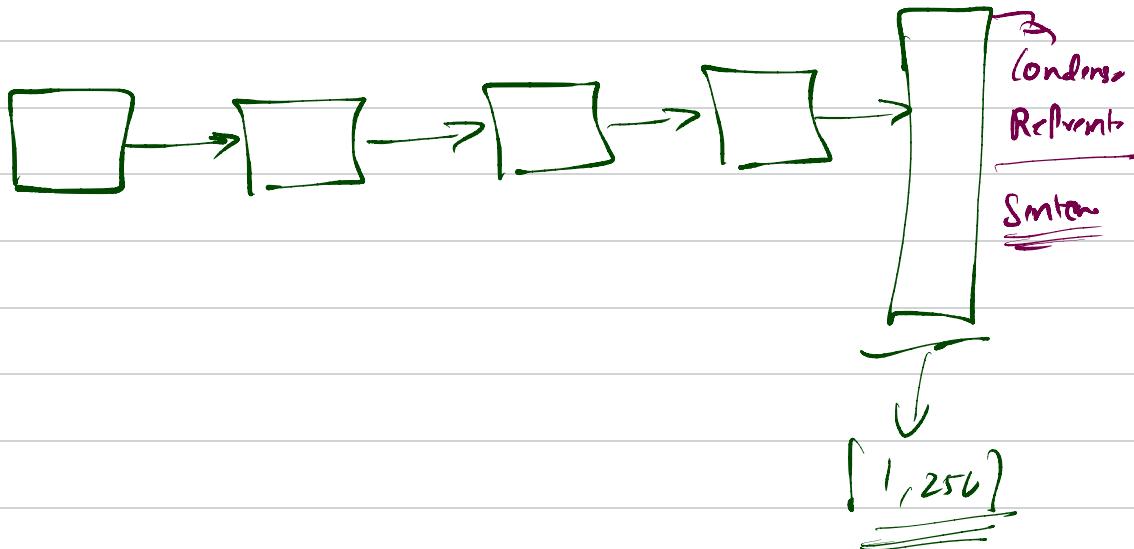
on out put  
of  
cell  
state  
vector

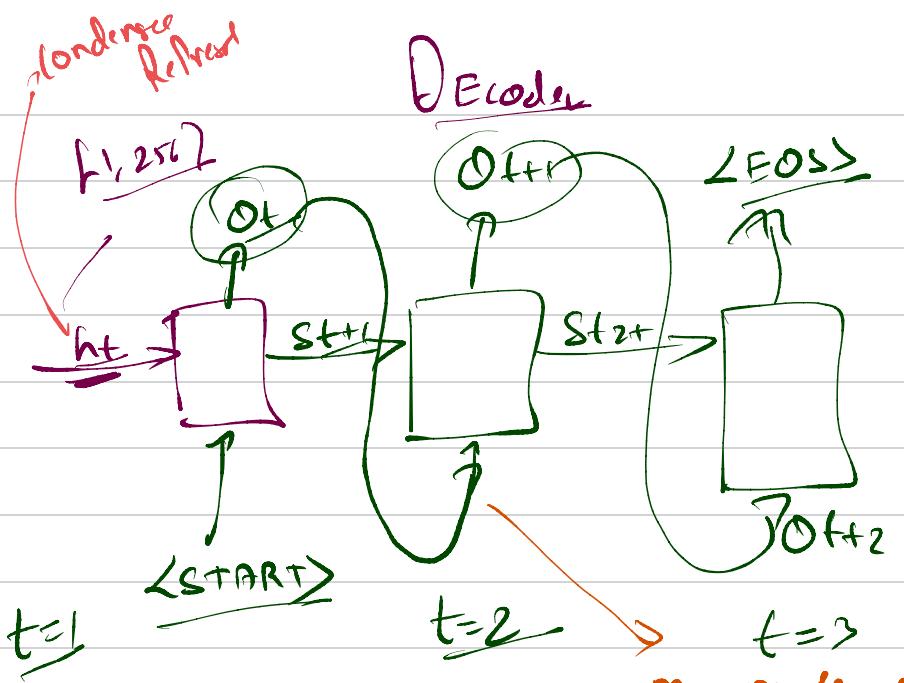
$$c_t = f_t \times c_{t-1} + i_t \times \tilde{c}_t$$

# ENCODER- DECODER



"g loved Roller-coaster at wonderl"   
 t<sub>1</sub> t<sub>2</sub> t<sub>3</sub> t<sub>4</sub> t<sub>5</sub> t<sub>6</sub>





my model needs to  
 know, what it has  
 spoken

