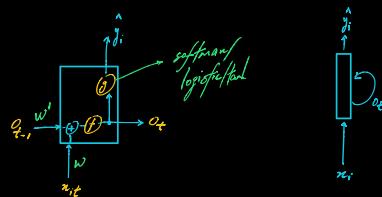
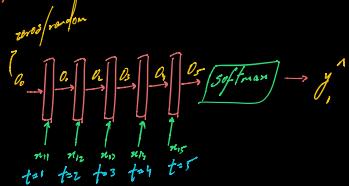
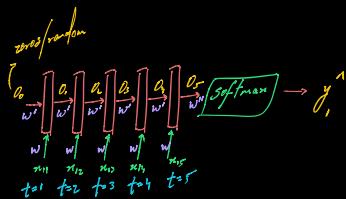


3 set of weights
we want \$o_i\$ to depend on \$x_{i,t}\$ & \$o_{i-1}

\$w \rightarrow\$ current input
\$w' \rightarrow\$ previous output each output depends on \$\rightarrow\$ current word
\$w'' \rightarrow\$ in the end + output of previous word



Backprop over Time



$$o_i = f(wx_{i,t} + w'o_{i-1})$$

$$o_i = f(wx_{i,t} + w'o_{i-1})$$

⋮

$$o_r = f(wx_{i,r} + w'o_{i-1})$$

$$\hat{y}_i = g(w''o_r)$$

$$\frac{\partial L}{\partial w''} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w''}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial o_r} \cdot \frac{\partial o_r}{\partial w}$$

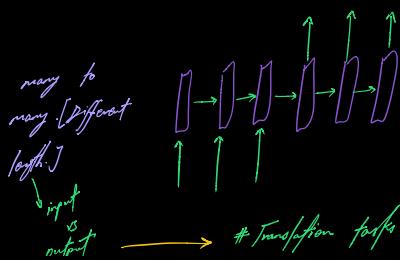
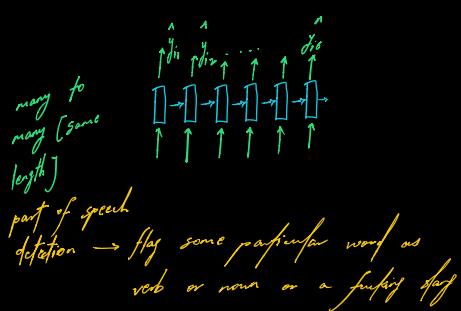
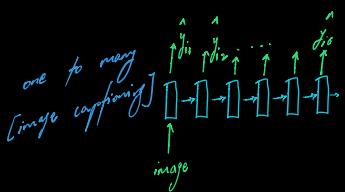
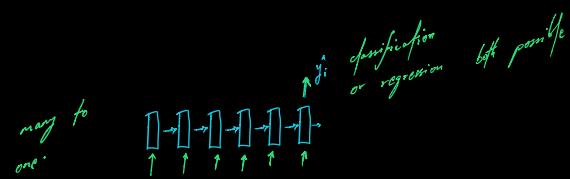
↑
\$w\$ on
the last
time

$$\frac{\partial L}{\partial w'} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial o_r} \cdot \frac{\partial o_r}{\partial w'}$$

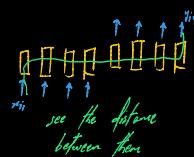
This backprop happens over time, on
the same weights.

The longer the sequence, the larger the
backprop. prompt. \therefore Vanishing/Exploding
gradient problem will persist.


LSTM / GRU
avoid these problems
through modifications to
RNN

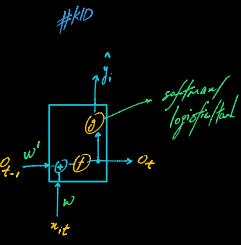


- # Single words such
 - ① Vanishing gradients occur
both forward & backward
passes
 - (if first & last word are
referred to each other, you will
not effectively affect that.)
- # RNNs break at long term dependencies
in input sequences



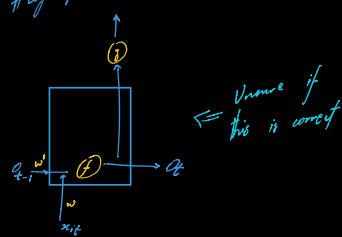
↳ we want the network
to have complete input
before outputting something

LSTM



RNN

RNN /



$$o_t = f([\omega', \omega], [x_{it}, h_{t-1}])$$

converge (not off the path) & off path

check by a
print for LSTM

GRU

→ optimized

→ less params to train

→ faster training

→ almost same performances → apply
as LSTM.