

12 - RNN

$$Y_t = f(X_t, X_{t-1}, \dots, X_{t-N})$$

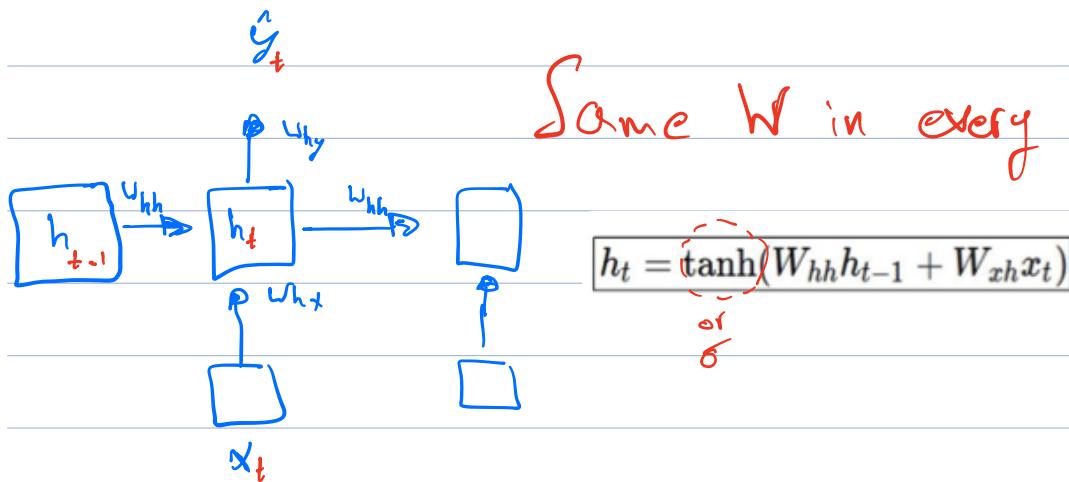
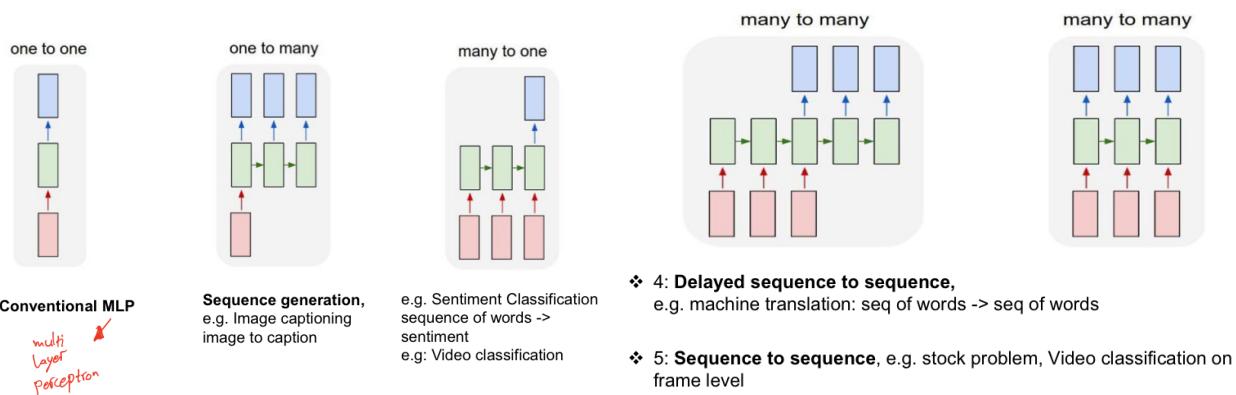
→ consider the history

→ finite-Response model

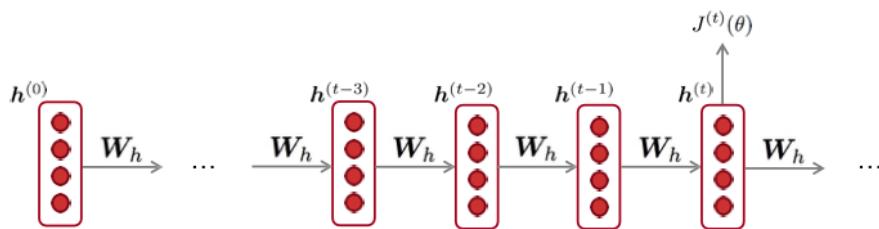
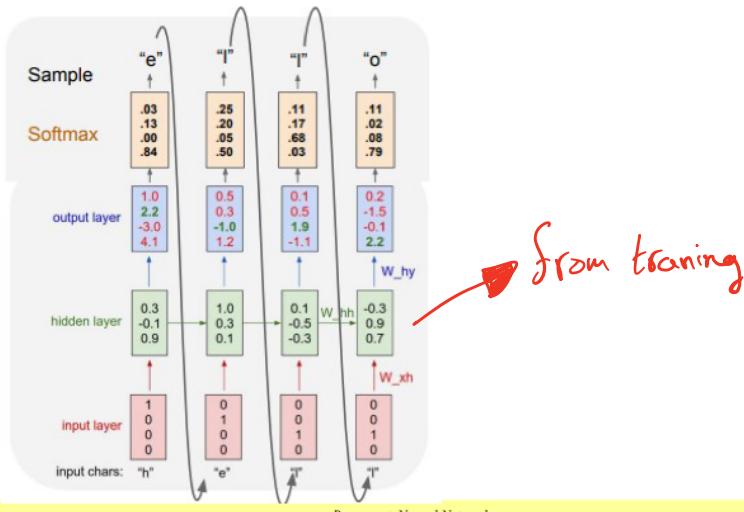
$$Y_t = f(X_t, X_{t-1}, \dots, X_{t-\infty})$$

Stock market
→ Infinite-Response model

~~X~~ types of RNN



- ❖ At **test-time** sample characters one at a time, feed back to model

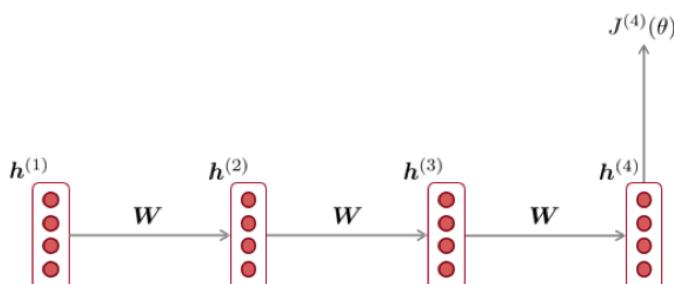


Question: What's the derivative of $J^{(t)}(\theta)$ w.r.t. the **repeated** weight matrix W_h ?

Answer: $\frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial W_h} \Big|_{(i)}$

"The gradient w.r.t. a repeated weight is the sum of the gradient w.r.t. each time it appears"

Called "backpropagation through time"



$$\frac{\partial J^{(t)}}{\partial h^{(1)}} = \frac{\partial J^{(t)}}{\partial h^{(4)}} \times \frac{\partial h^{(4)}}{\partial h^{(3)}} \times \frac{\partial h^{(3)}}{\partial h^{(2)}} \times \frac{\partial h^{(2)}}{\partial h^{(1)}}$$

Recall: $h^{(t)} = \sigma(W_h h^{(t-1)} + W_x x^{(t)} + b_1)$

Therefore: $\frac{\partial h^{(t)}}{\partial h^{(t-1)}} = \text{diag}(\sigma'(W_h h^{(t-1)} + W_x x^{(t)} + b_1)) W_h$ (chain rule)

L2-regularization will make vanishing worse (set $W=0$)

* Vanilla RNN has 2 problems:-

when largest eigenvalue
of $W_h < 1$

language model

① Vanishing the gradient: loss the history (in LM tasks = loss the syntactic and just keep the sequential recency)

Solutions:-

1- LSTMs (long-short-term-memories)

2- GRUs (Gated-Recurrent-Units)

when largest eigenvalue

of $W_h > 1$

② Exploding the gradient: bad update (learning) = NaN or Inf

Solutions:-

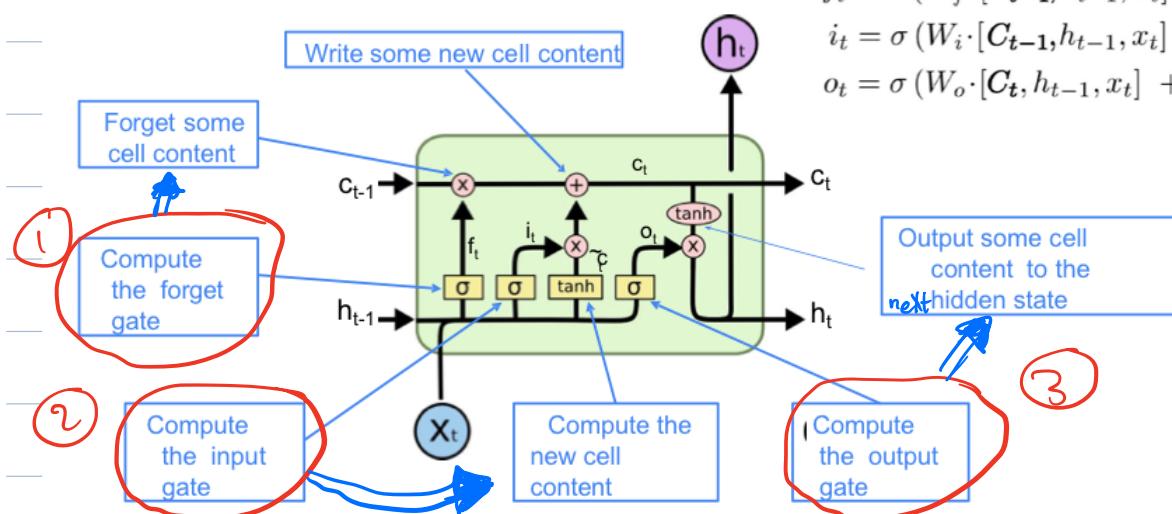
1- gradient clipping: if it \geq threshold \Rightarrow scale it before update

- Intuition: take a step in the same direction, but a smaller step

#LSTM

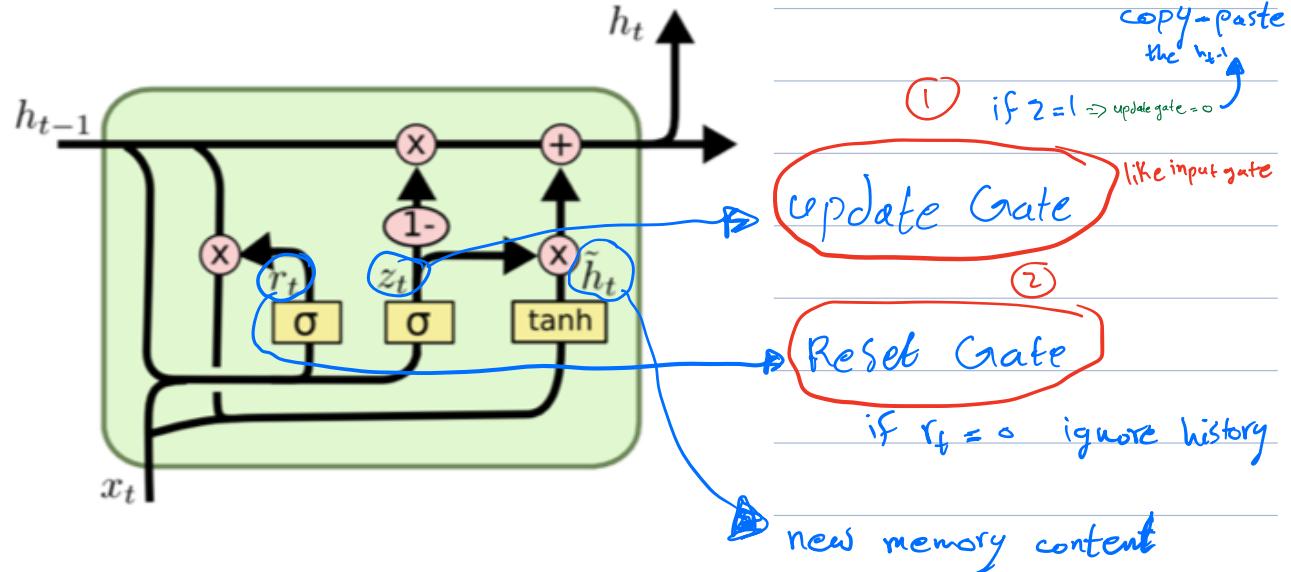
$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$
$$i_t = \sigma(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$
$$o_t = \sigma(W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$

↑ real
0 → forget



X GRUs

GRU first computes an update gate



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

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

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

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

How they solve vanishing:

LSTM = add constant-error carousels that carry information
but only affected by 2 gates (Forget and Input)

GRU = allow error messages to flow at different strengths

depending on the inputs
quicker

✗ Vanishing is not just in RNN
any deep network has this problem

Solutions:

- add more direct connections
 - ResNet Skip-connection
 - DenseNet Dense-connection

✗ Bidirectional RNN

- could be a vanilla, LSTM, or GRU
- every RNN has its own W
- for classification $\boxed{h^{(t)}} = [\overrightarrow{h}^{(t)}; \overleftarrow{h}^{(t)}]$ - concatenation

✗ Multi-layer RNNs

Stack RNN

13 - Encoder-decoder (Time asynchronance)

💡 In situations where the start of sequence is obvious, the `<sos>` may not be needed, but `<eos>` is required to terminate sequences

💡 Sometimes we will use a single symbol to represent both start and end of sentence, e.g. just `<eos>`, or even a separate symbol, e.g. `<s>`

- The input sequence is terminated by an explicit `<eos>` symbol
 - The hidden activation at the `<eos>` "stores" all information about the sentence
 - Output production continues until an `<eos>` is produced

$$\hat{y}_k^w = P(O_k = w | O_{k-1}, \dots, O_1, I_1, \dots, I_N)$$

at time K take argmax
output Input
greedy take most probable word on each step

Vanilla E/D is greedy decoder (no way to undo) stops when we get `<end>`

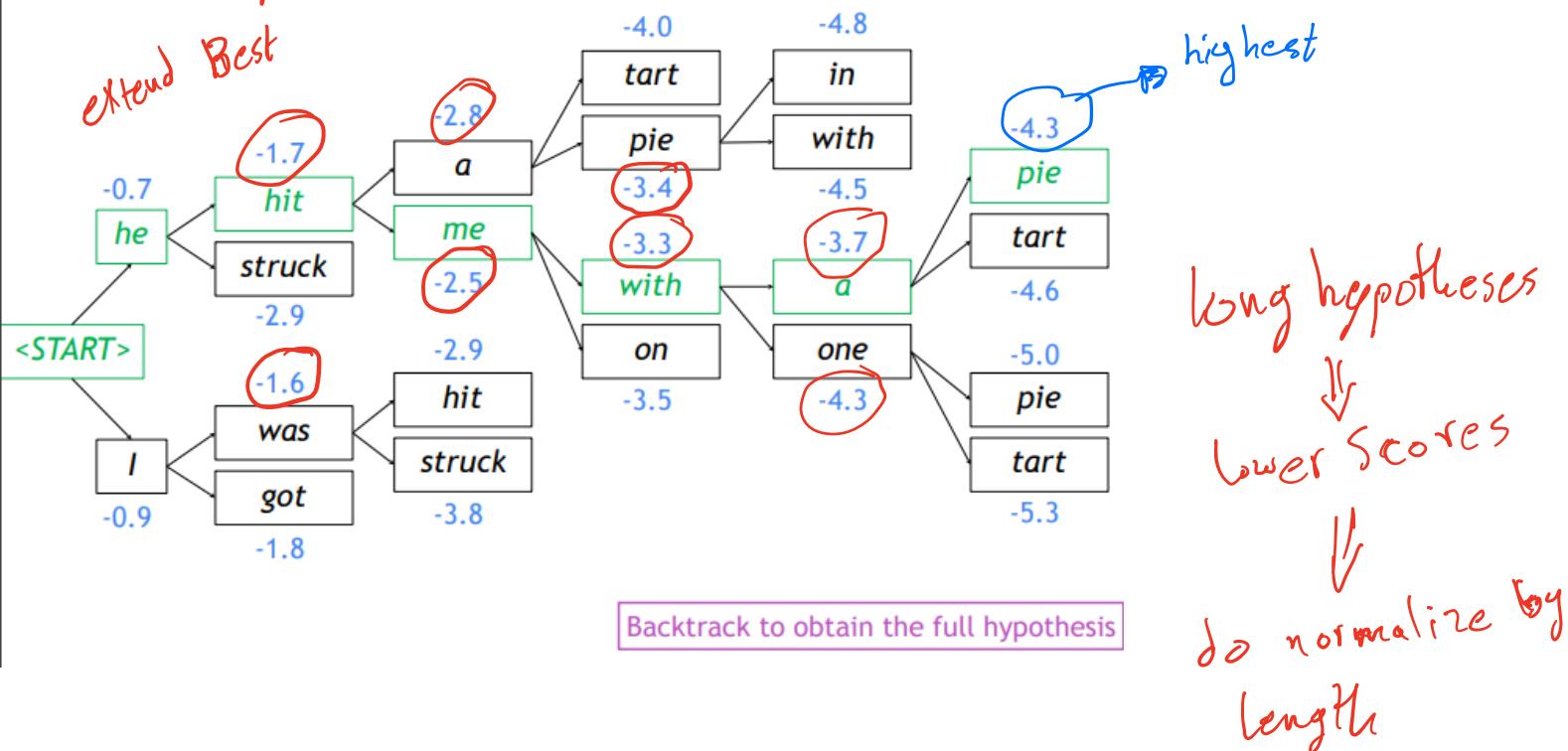
Solutions:-

- ① Exhaustive search decoding (all possible seq.) too expensive $O(V^T)$
- ② Beam Search decoding (Keep track of the K most probable results)
no guarantee to find the optimal

- Stops when:- input of encoder
- ① maximum T timesteps
 - ② maximum N hypotheses output of decoder
 - ③ got `<end>` # of complete output

Beam search decoding: example

Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



~~XX~~ Training E/D

- For each iteration
 - ① Randomly select training instance: (input, output)
 - ② Forward pass
 - ③ Randomly select a single output $y(t)$ and corresponding desired output $d(t)$ for backprop
 - ④ - update parameters

* Reversing the Input

- Things work better this way

This happens both for training and during inference on test data

~~XX~~ Ways to treat videos for video captioning:-

- 1- for each frame (CNN) ^{to extract features} then get the mean across all frame then RNN
- 2- for each frame (CNN) then encoder-decoder

All types of E/D has a bottleneck problem (context vector)

Solution:- Attention

it also helps with vanishing
not just for Seq2Seq

- ❖ More general definition of attention:
Given a set of vector values, and a vector query,
attention is a technique to compute a weighted sum of the values, dependent on the query.

Ways to calculate e_i :

① dot-product

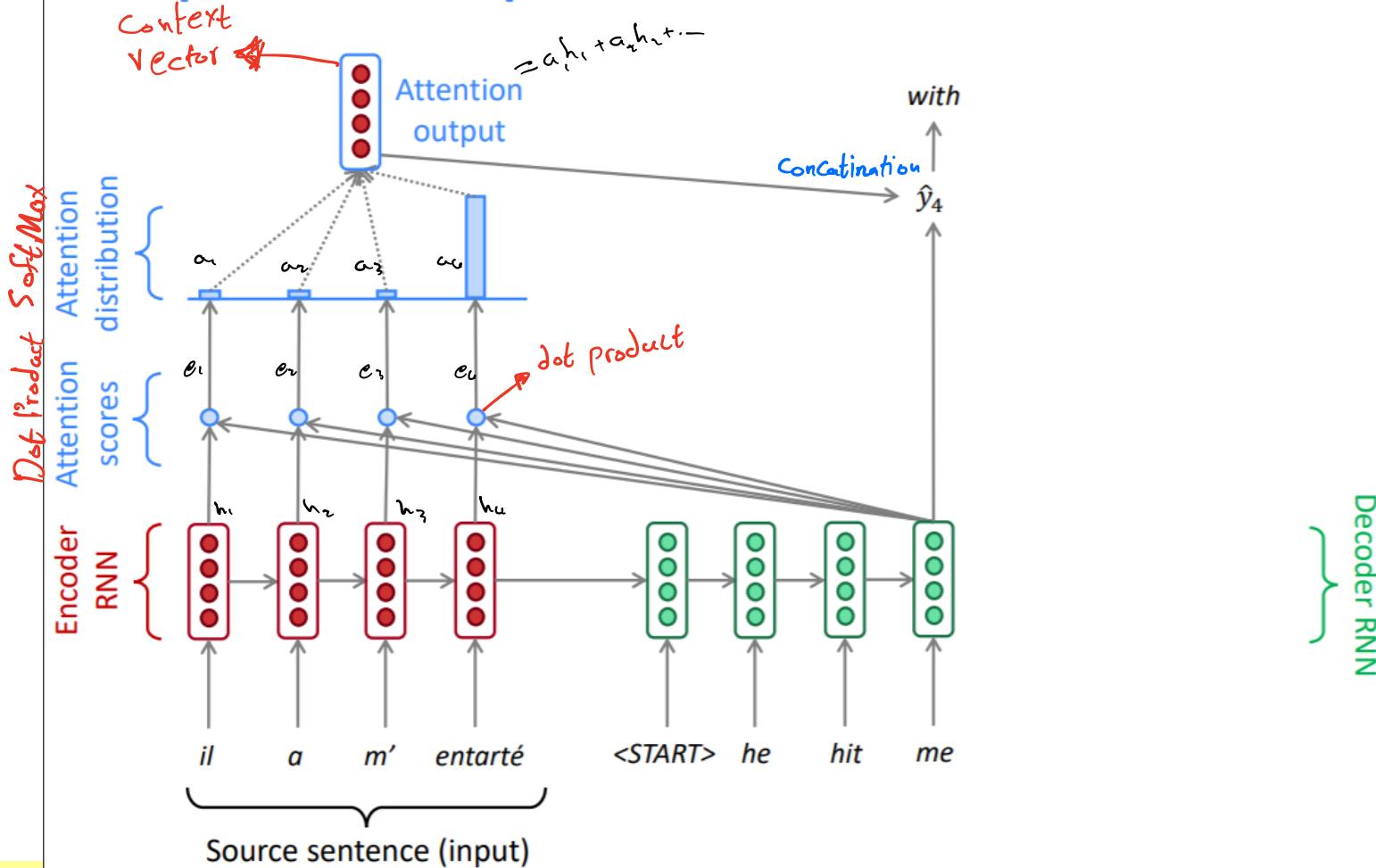
② multiplicative attention

multiply with W

③ additive attention

addition with W \rightarrow hyperparameters

Sequence-to-sequence with attention



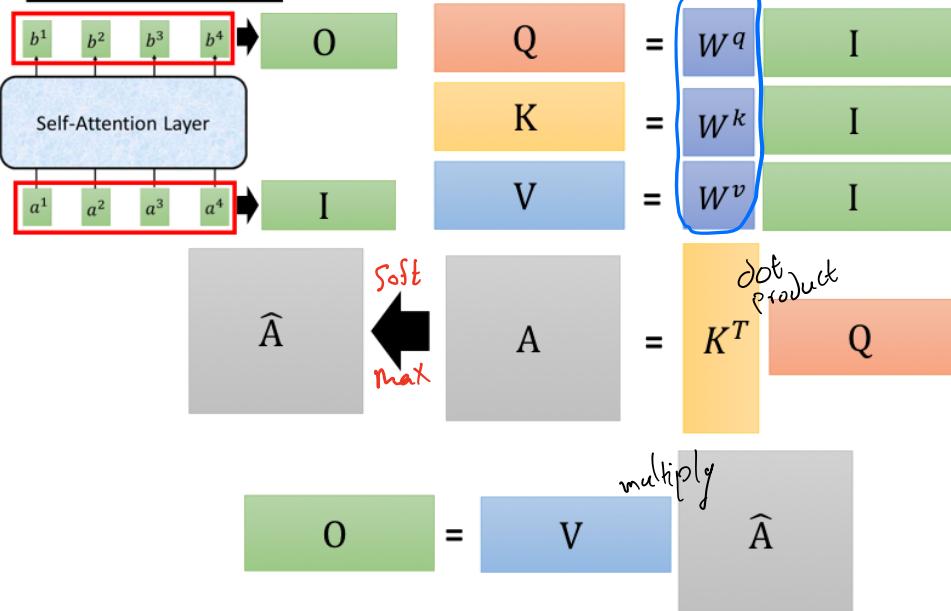
14 - Self-Attention and Transformers

Self-Attention-

- inputs = outputs
- inputs interact with each other to decide which one should have attention
- output = these interactions and attention Scores

for each input q, K, V , are same

Self-attention



We use multiple Head to focus on multiple places at once

each one will do attention independently then we concatenate them

Self Attention has problems:-

- ① Doesn't care about the order (may change the meaning)
- ② Doesn't have nonlinearities
- ③ May in decoding it look at the future

Transformers \Rightarrow are self-Attention and more



absolute position is less important than relative position

just include the past word

So we use sin/cos

by setting the future = $-\infty$

Necessities for a self-attention building block:

- **Self-attention:**
 - the basis of the method.
- **Position representations:**
 - Specify the sequence order, since self-attention is an unordered function of its inputs.
- **Nonlinearities:**
 - At the output of the self-attention block
 - Frequently implemented as a simple feed-forward network.
- **Masking:**
 - In order to parallelize operations while not looking at the future.
 - Keeps information about the future from "leaking" to the past.

That's it! But this is not the **Transformer** model we've been hearing about?

John Hewitt

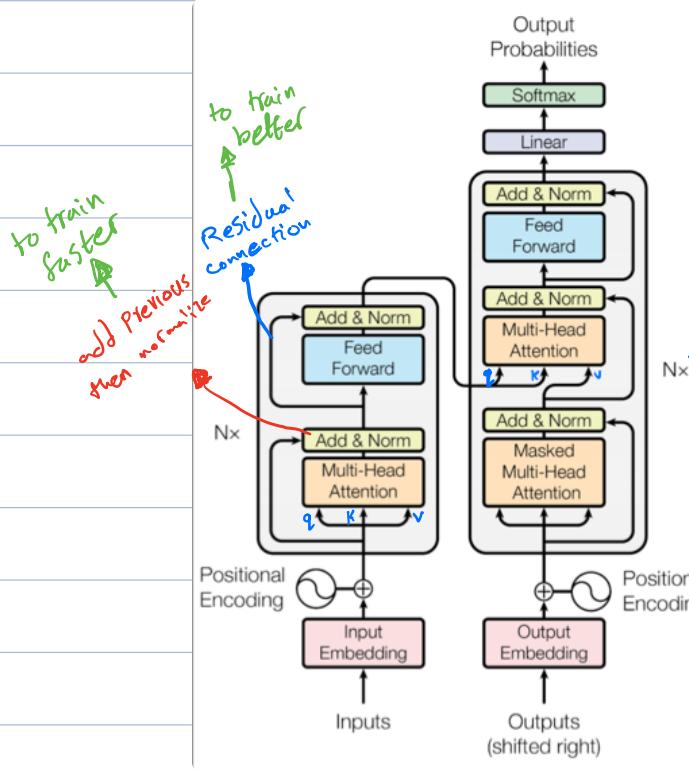
We can look at these (not greyed out) words

[START] The chef who

[START]	$-\infty$	$-\infty$	$-\infty$	$-\infty$
The		$-\infty$	$-\infty$	$-\infty$
chef			$-\infty$	$-\infty$
who				$-\infty$

For encoding these words

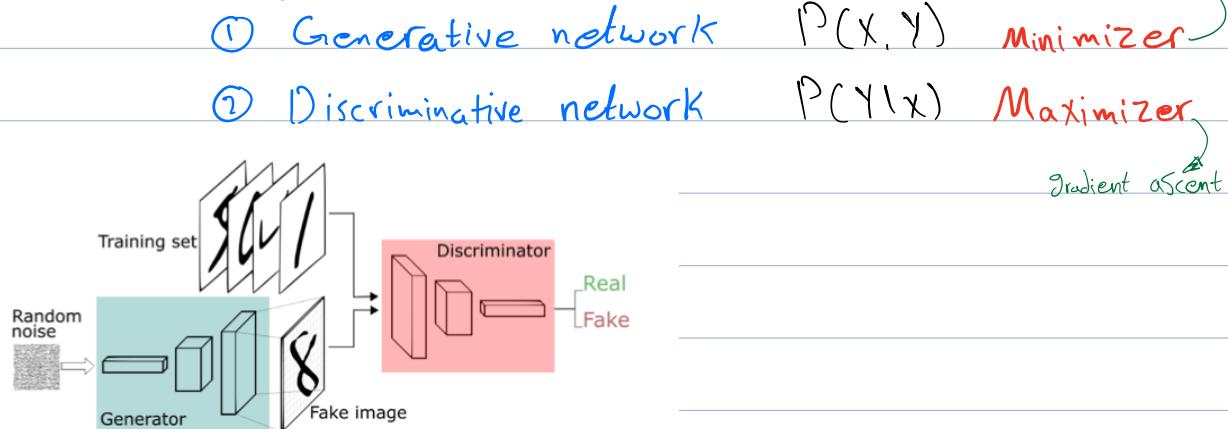
Transformers



19 - GAN

generative adversarial networks

two player game :-



Train jointly in **minimax game**

Minimax objective function:

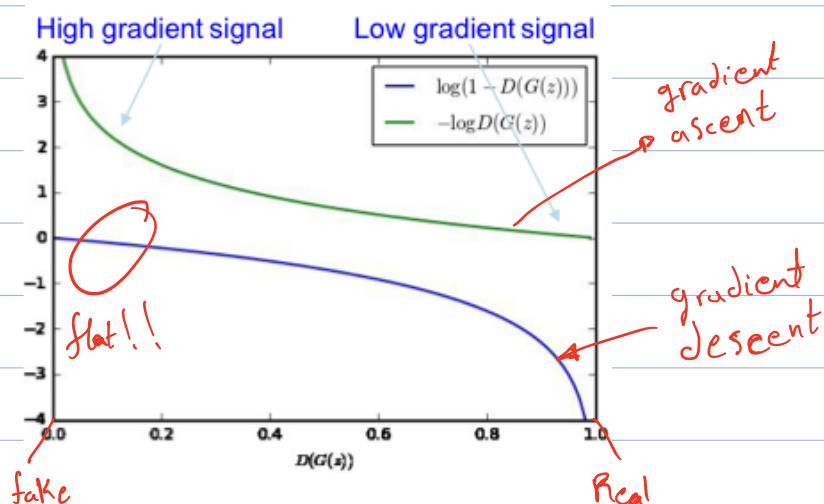
$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{\text{data}}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Discriminator output for real data x Discriminator output for generated fake data $G(z)$

- Discriminator (θ_d) wants to **maximize objective** such that $D(x)$ is close to 1 (real) and $D(G(z))$ is close to 0 (fake)
- Generator (θ_g) wants to **minimize objective** such that $D(G(z))$ is close to 1 (discriminator is fooled into thinking generated $G(z)$ is real)

for generator:

gradient ascent is
better



* Architectures

generator = upsampling network with fractionally-strided convolutions

discriminator = CNN

Architecture guidelines for stable Deep Convolutional GANs

- Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
- Use batchnorm in both the generator and the discriminator.
- Remove fully connected hidden layers for deeper architectures.
- Use ReLU activation in generator for all layers except for the output, which uses Tanh.
- Use LeakyReLU activation in the discriminator for all layers.

* Challenging in Training:-

① Vanishing gradient \Rightarrow When discriminator is very good

② Mode collapse \Rightarrow When too little diversity in the generated samples

③ Normal loss metric doesn't correspond to image quality

Frechet Inception Distance is a decent choice for calc. loss

* Types of GANs

① Conditional GAN

generate x ^{random} given y So the input for G and D are x and y

$$\min_G \max_D V(D, G) = \mathbb{E}_X \log D(X, Y) + \mathbb{E}_Z \log D(G(Z, Y), Y)$$

Example:- Edge \rightarrow Image

Sketch \rightarrow Image

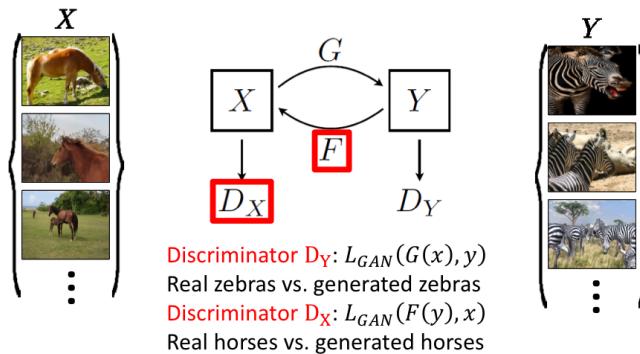
② LapGAN

Chain of conditional GANs

③ CycleGAN

pair of conditional GANs

Image -to- Image translation



they use cycle consistency as loss fun.