

NEURAL NETWORKS AS COMPUTATIONAL GRAPHS

L8-1

$$\begin{array}{ccccc}
 X & z^{(0)} & a^{(0)} & z^{(1)} & p_{y|x} \\
 \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\
 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & - 0 \\
 0 & 0 & 0 & 0 & - 0 \\
 0 & 0 & 0 & 0 & - 0 \\
 0 & 0 & 0 & 0 & - 0 \\
 \end{array}
 \rightarrow \hat{y} = \text{ARGMAX } p_{y|x}$$

$\sigma_{\text{SOFTMAX}}(z^{(1)})$

\downarrow $XW^{(0)T} + b^{(0)}$ \downarrow $\max(z^{(0)}, \alpha z^{(0)})$ \downarrow $a^{(0)}W^{(1)T} + b^{(1)}$

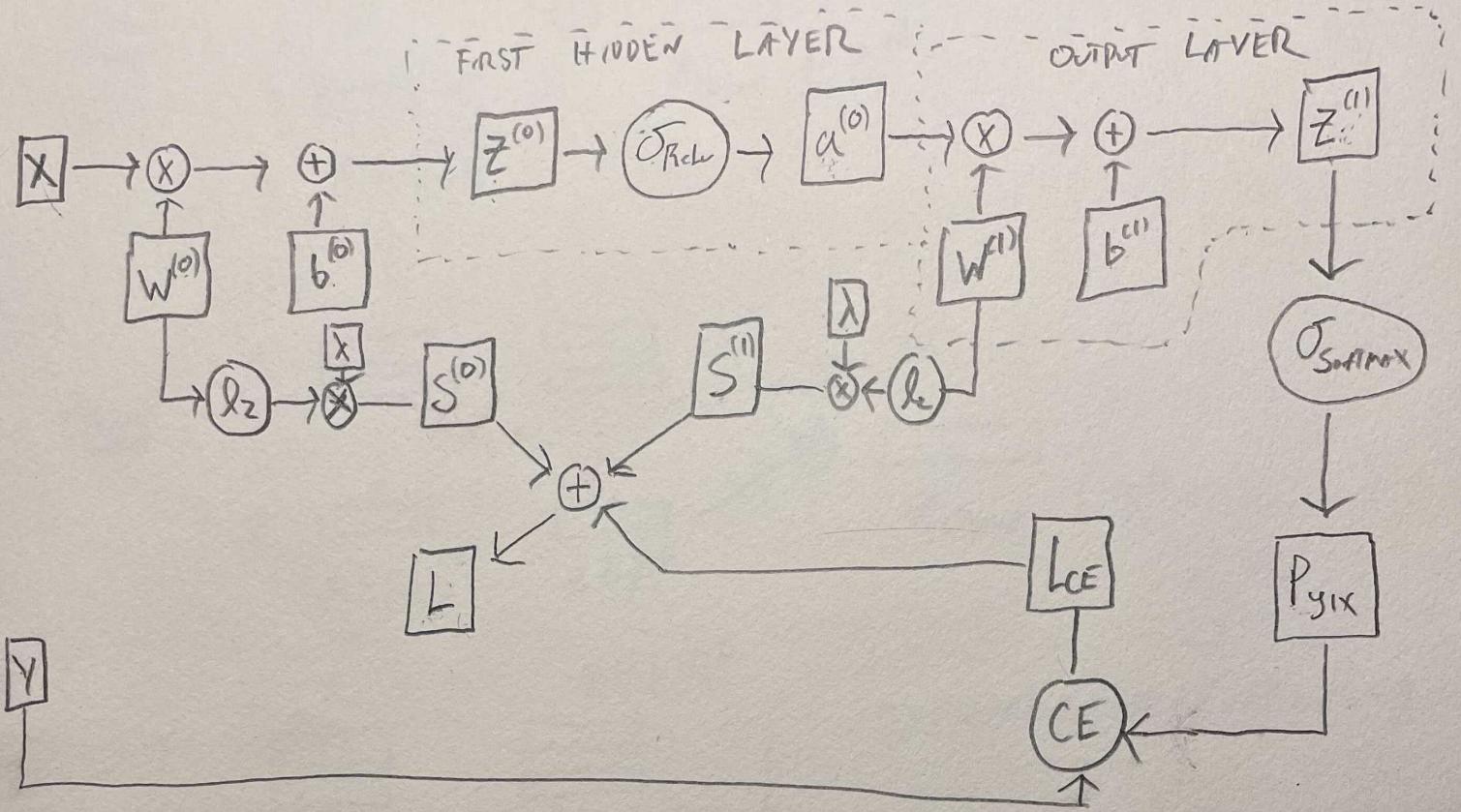
$$\begin{aligned}
 X &\in \mathbb{Z}_+^N & W^{(0)} &\in \mathbb{R}^{D \times N} & W^{(1)} &\in \mathbb{R}^{K \times D} \\
 z^{(0)} &\in \mathbb{R}^D & b^{(0)} &\in \mathbb{R}^D & b^{(1)} &\in \mathbb{R}^K \\
 z^{(1)} &\in \mathbb{R}^K & & & & \\
 p_{y|x} &\in [0,1]^K & & & &
 \end{aligned}$$

Chain RULE Refresh

$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial X}$$

Computational Graph

L8-2



From Lecture 05 - Softmax Regression

$$\frac{\partial L}{\partial z^{(1)}} = P_{y|x} - Y \in \mathbb{R}^K$$

$$\begin{array}{c} \text{LReLU} \\ \text{---} \\ \text{---} \end{array} \quad \begin{array}{c} \text{PRELU} \\ \text{---} \\ \text{---} \end{array}$$

$$\frac{\partial L}{\partial b^{(1)}} = \frac{\partial L}{\partial z^{(1)}} \cdot \frac{\partial z^{(1)}}{\partial b^{(1)}} = P_{y|x} - Y \in \mathbb{R}^K$$

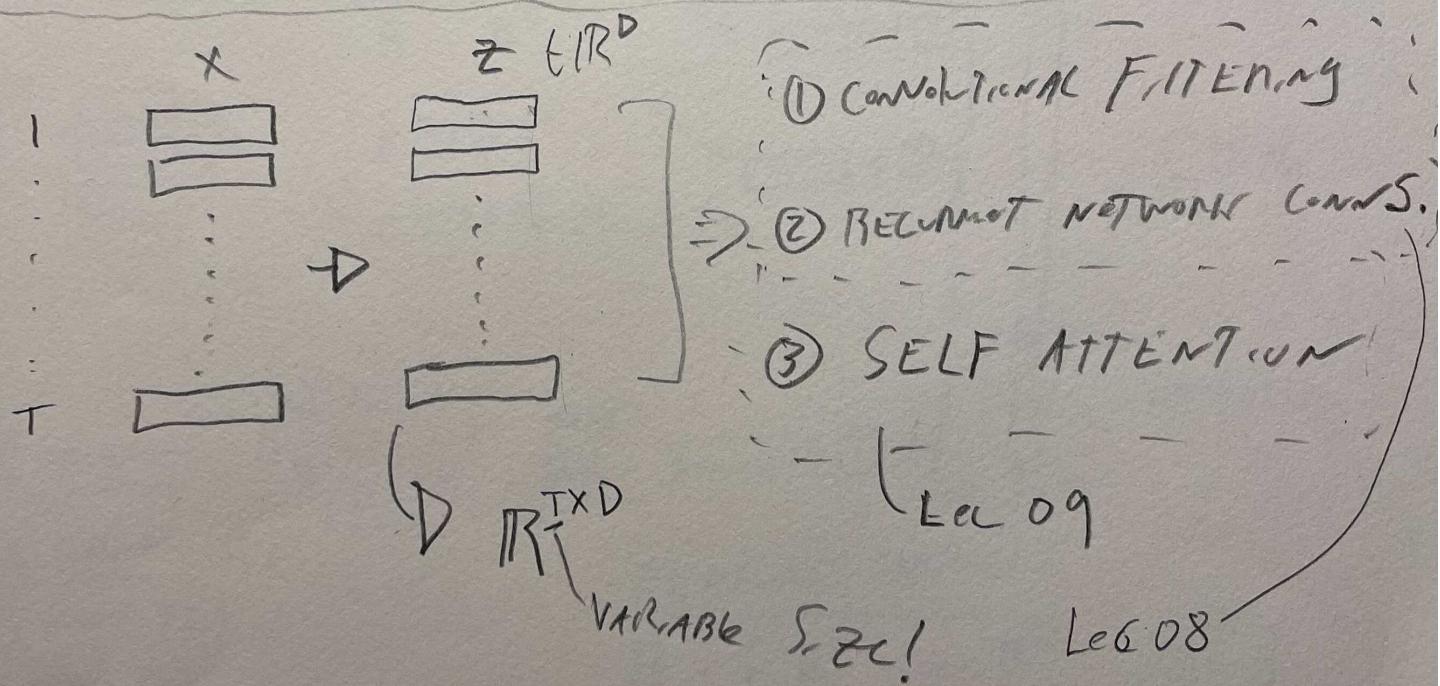
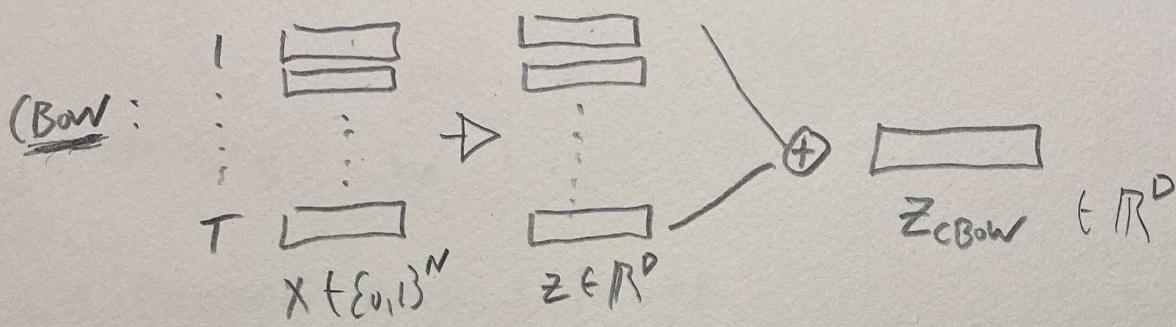
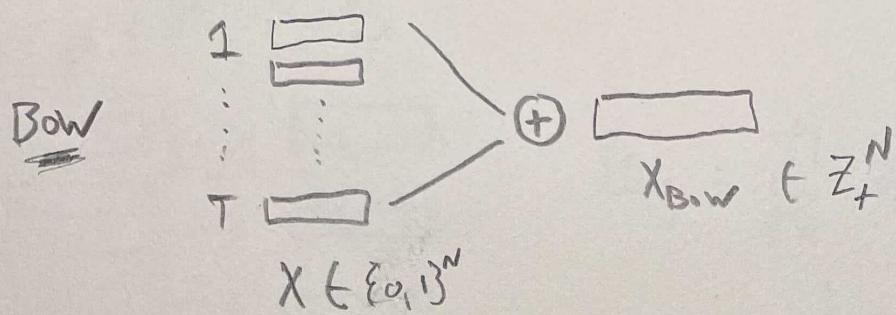
$$\frac{\partial L}{\partial w^{(1)}} = \frac{\partial L}{\partial z^{(1)}} \cdot \frac{\partial z^{(1)}}{\partial w^{(1)}} + \frac{\partial L}{\partial s^{(1)}} \cdot \frac{\partial s^{(1)}}{\partial w^{(1)}} = (P_{y|x} - Y) \cdot a^{(0)} + \lambda w^{(1)} \in \mathbb{R}^{K \times D}$$

$$\frac{\partial L}{\partial a^{(0)}} = \frac{\partial L}{\partial z^{(1)}} \cdot \frac{\partial z^{(1)}}{\partial a^{(0)}} \cdot \frac{\partial a^{(0)}}{\partial z^{(0)}} \cdot \frac{\partial z^{(0)}}{\partial b^{(0)}} = (P_{y|x} - Y) \cdot W^{(1)} \cdot \left\{ \begin{array}{l} \frac{\partial z^{(0)}}{\partial b^{(0)}} > 0 \\ \frac{\partial z^{(0)}}{\partial b^{(0)}} < 0 \end{array} \right\} \cdot 1 \in \mathbb{R}^D$$

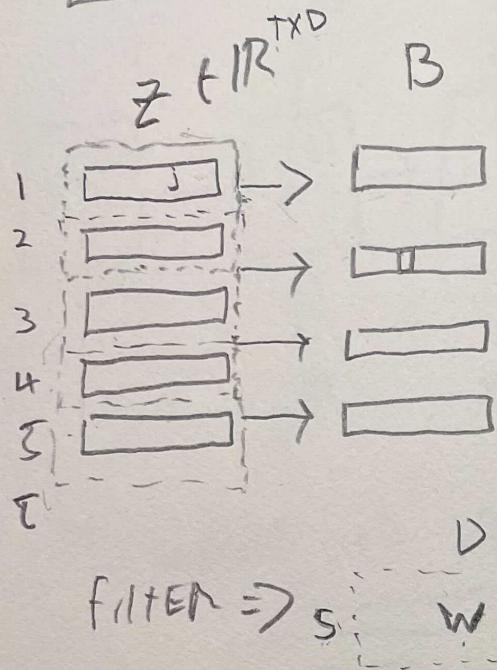
$$\frac{\partial L}{\partial w^{(0)}} = \frac{\partial L}{\partial z^{(1)}} \cdot \frac{\partial z^{(1)}}{\partial a^{(0)}} \cdot \frac{\partial a^{(0)}}{\partial z^{(0)}} \cdot \frac{\partial z^{(0)}}{\partial w^{(0)}} + \frac{\partial L}{\partial s^{(0)}} \cdot \frac{\partial s^{(0)}}{\partial w^{(0)}} = (P_{y|x} - Y) \cdot W^{(1)} \cdot \left\{ \begin{array}{l} \frac{\partial z^{(0)}}{\partial b^{(0)}} > 0 \\ \frac{\partial z^{(0)}}{\partial b^{(0)}} < 0 \end{array} \right\} \cdot X + \lambda w^{(0)} \in \mathbb{R}^{D \times N}$$

SEQUENCE Processing with NN's

L 8-3



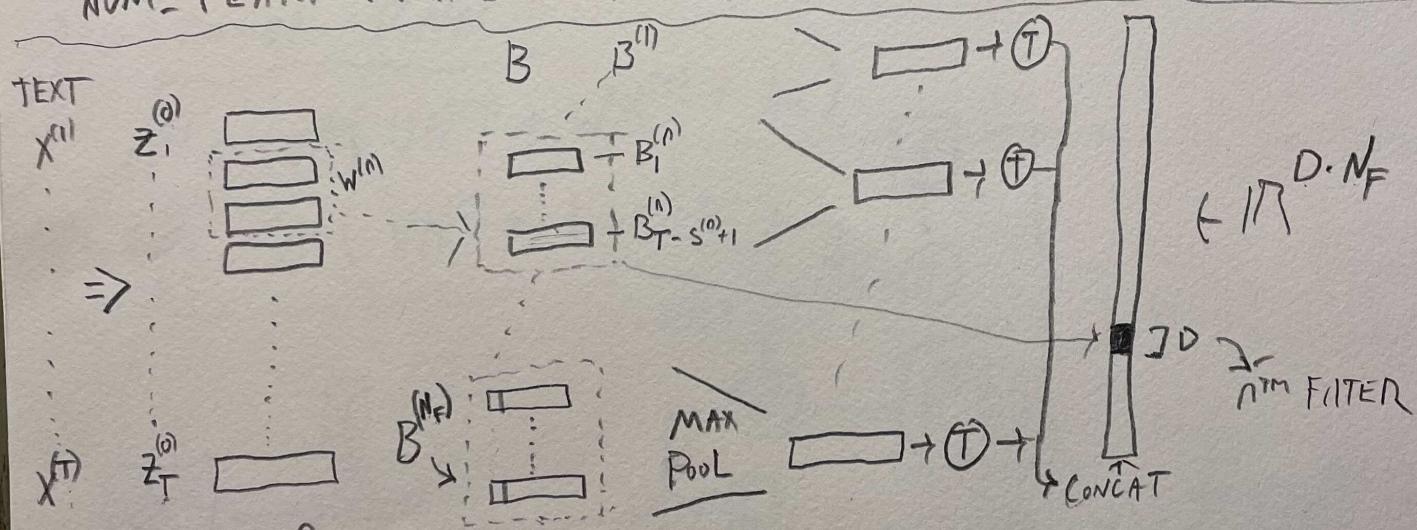
Convolutional Filtering



$$B_{tj} = \sum_{s=1}^S W_{tj}^s \cdot Z_{t+s-1,j}$$

$\forall t \in \{1, \dots, T-S+1\}$

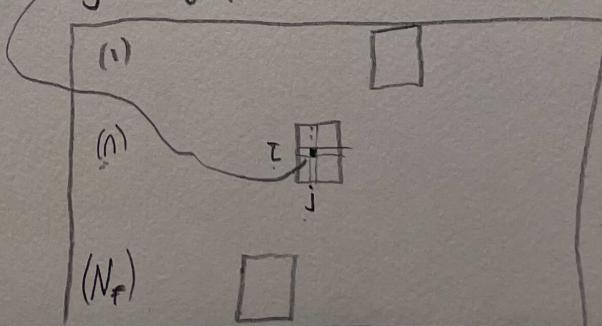
$$\text{NUM_FEATURE_MAPS} = T-S+1$$



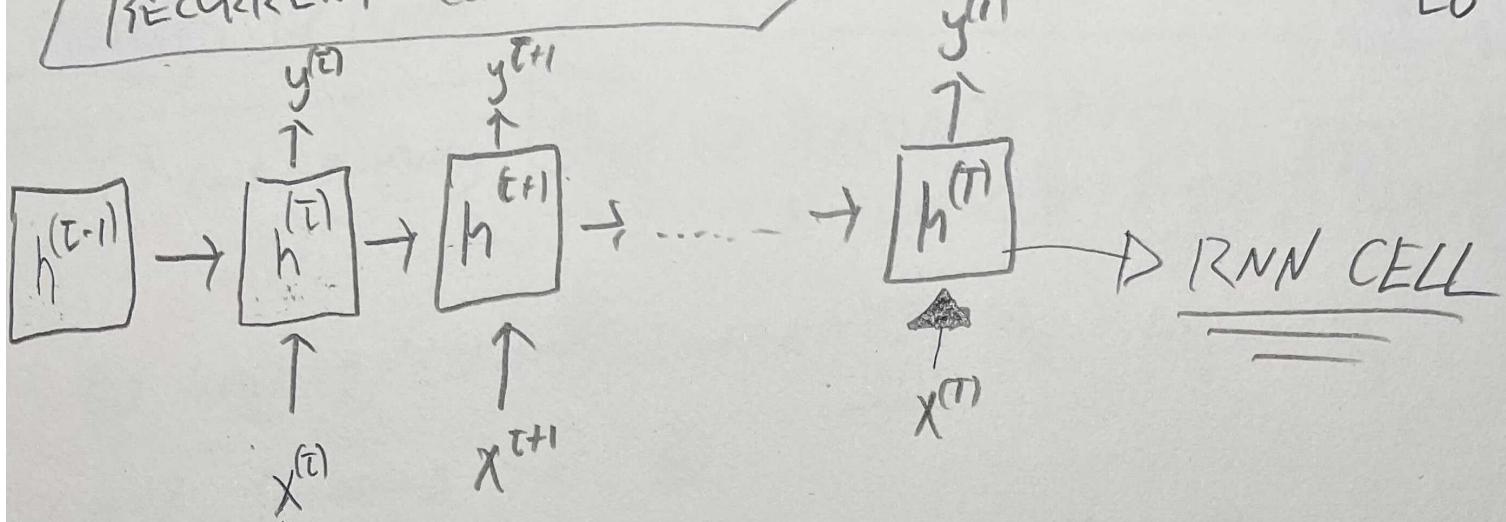
FEATURE MAP

$$B = \{B^{(1)}, \dots, B^{(n)}, \dots, B^{(N_f)}\}$$

$$B_{tj}^{(n)} = \sum_{s=1}^{S^{(n)}} W_{tj}^{(n)} Z_{t+s-1,j}^{(0)}$$



RECURRENT CONNECTIONS

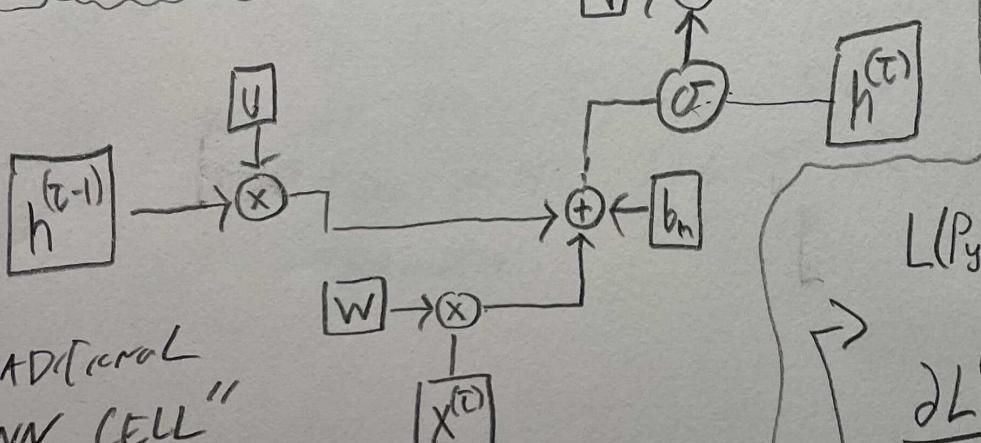


What is $h^{(t)}$?

TRADITIONAL RNN

$$h^{(t)} = \sigma(Vh^{(t-1)} + WX^{(t)} + b_h)$$

$$y^{(t)} = \sigma(Vh^{(t)} + b_y)$$



"TRADITIONAL
RNN CELL"

$$L(P_{y|x}, y) = \sum_{t=1}^T L(P_{y|x_t}, y^{(t)})$$

$$\frac{\partial L^{(T)}}{\partial W} = \sum_{t=1}^T \frac{\partial L^{(T)}}{\partial W^{(t)}}$$

How are RNNs OPTIMIZED?

PROBLEMS w/ RNNs

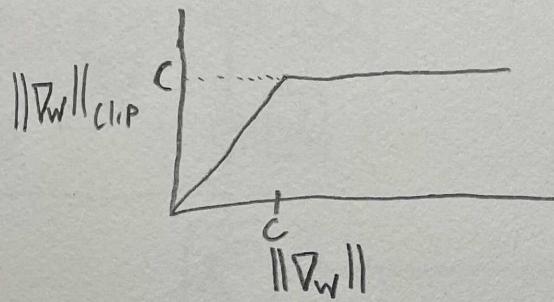
L8-6

① computationally slow \Rightarrow Backprop through
NETWORK + TIME!

② VANISHING / EXPLODING GRADIENTS

Common Solutions to ②

④ EXPLODING GRADIENTS: GRADIENT CLIPPING



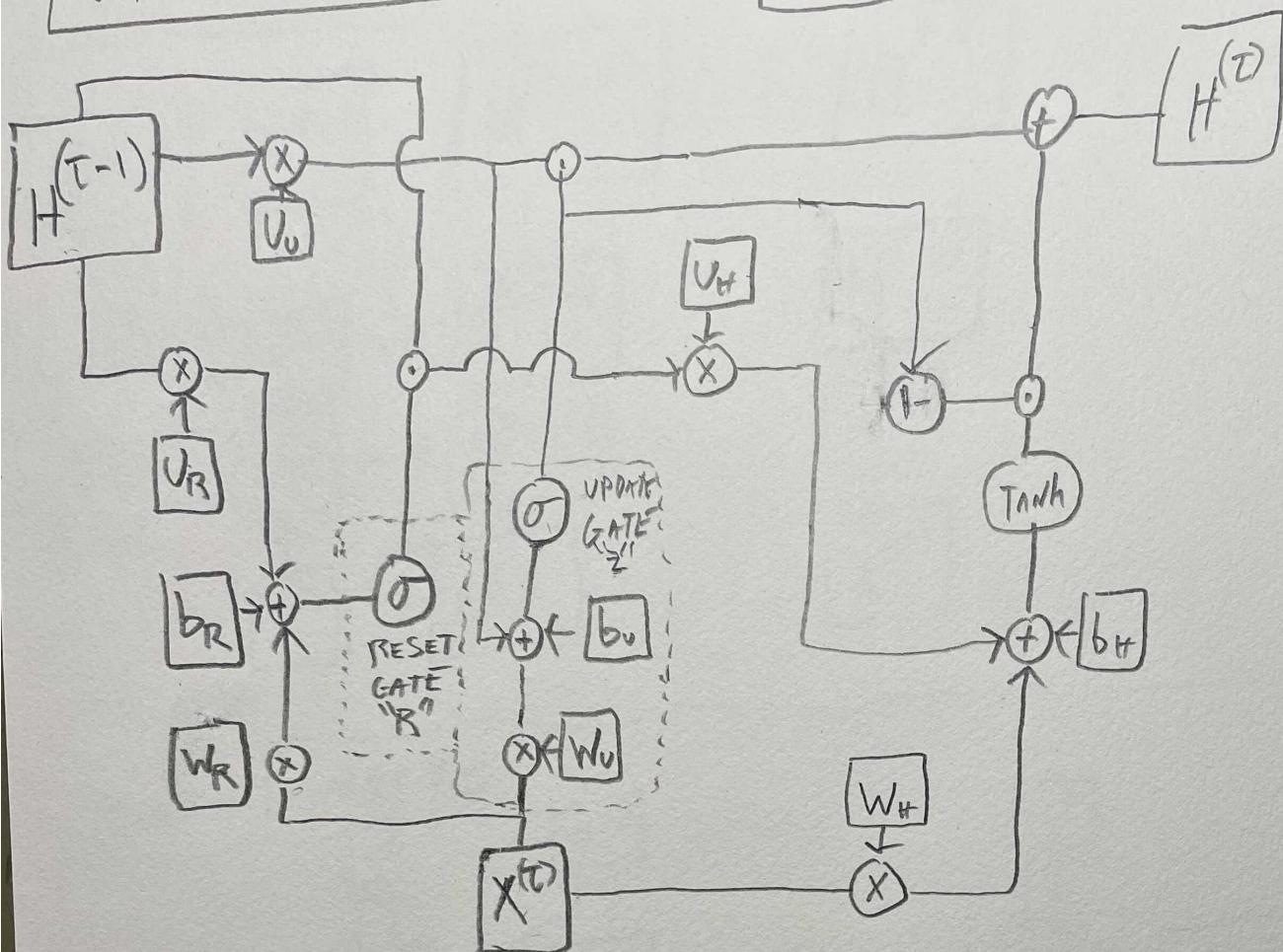
④ VANISHING GRADIENT: GATING

④ GATE: ONLY GATES TO THE HIDDEN STATE
TO LIMIT THE FORWARD/BACKWARD PROPAGATION
OF THE SIGNAL SUCH THAT ONLY
THE MOST SALIENT FEATURES PASS
THROUGH.

GATED RECURRENT UNIT

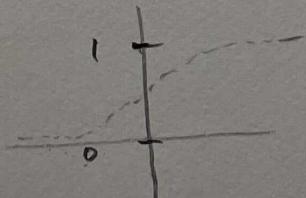
GRU

L8-7



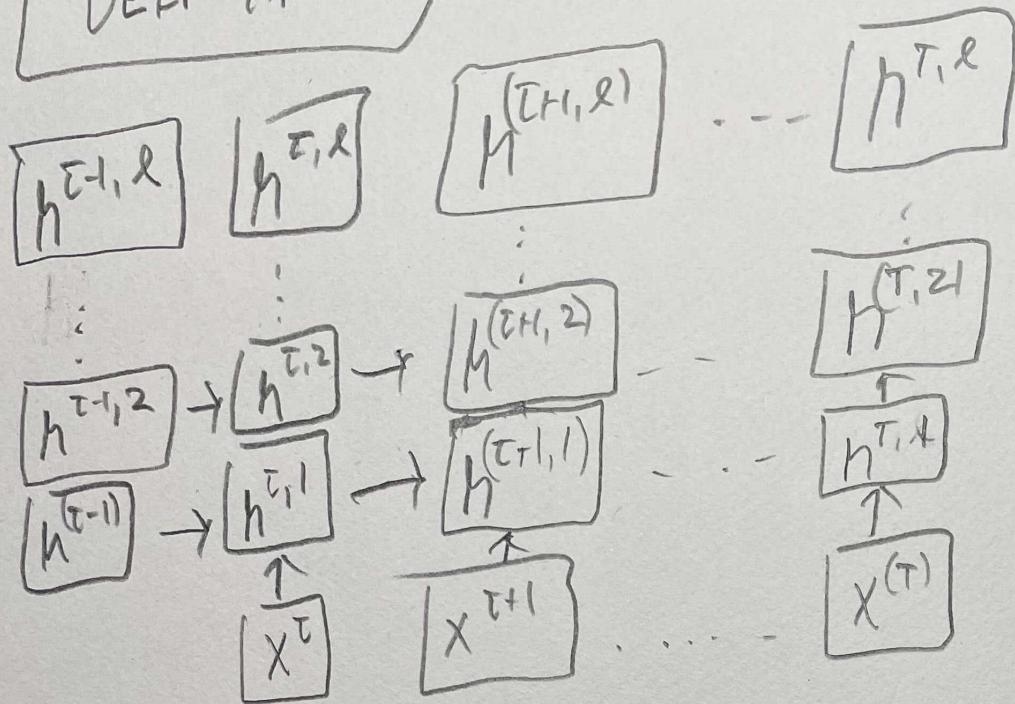
$$\begin{aligned}
 H &\in \mathbb{R}^H \\
 X &\in \mathbb{R}^D \\
 R &\in \mathbb{R}^H \\
 Z &\in \mathbb{R}^H \\
 V_H &\in \mathbb{R}^{H \times H} \\
 b_H &\in \mathbb{R}^H \\
 W_H &\in \mathbb{R}^{H \times D} \\
 V_R &\in \mathbb{R}^{H \times H} \\
 U_Z &\in \mathbb{R}^{H \times H} \\
 b_R &\in \mathbb{R}^H \\
 W_R &\in \mathbb{R}^{H \times D} \\
 b_Z &\in \mathbb{R}^H \\
 W_Z &\in \mathbb{R}^{H \times D}
 \end{aligned}$$

Sigmoid function looks like this:



DEEP RNNs

L8-8



BIDIRECTIONALITY

