# CS540 Introduction to Artificial Intelligence Lecture 24
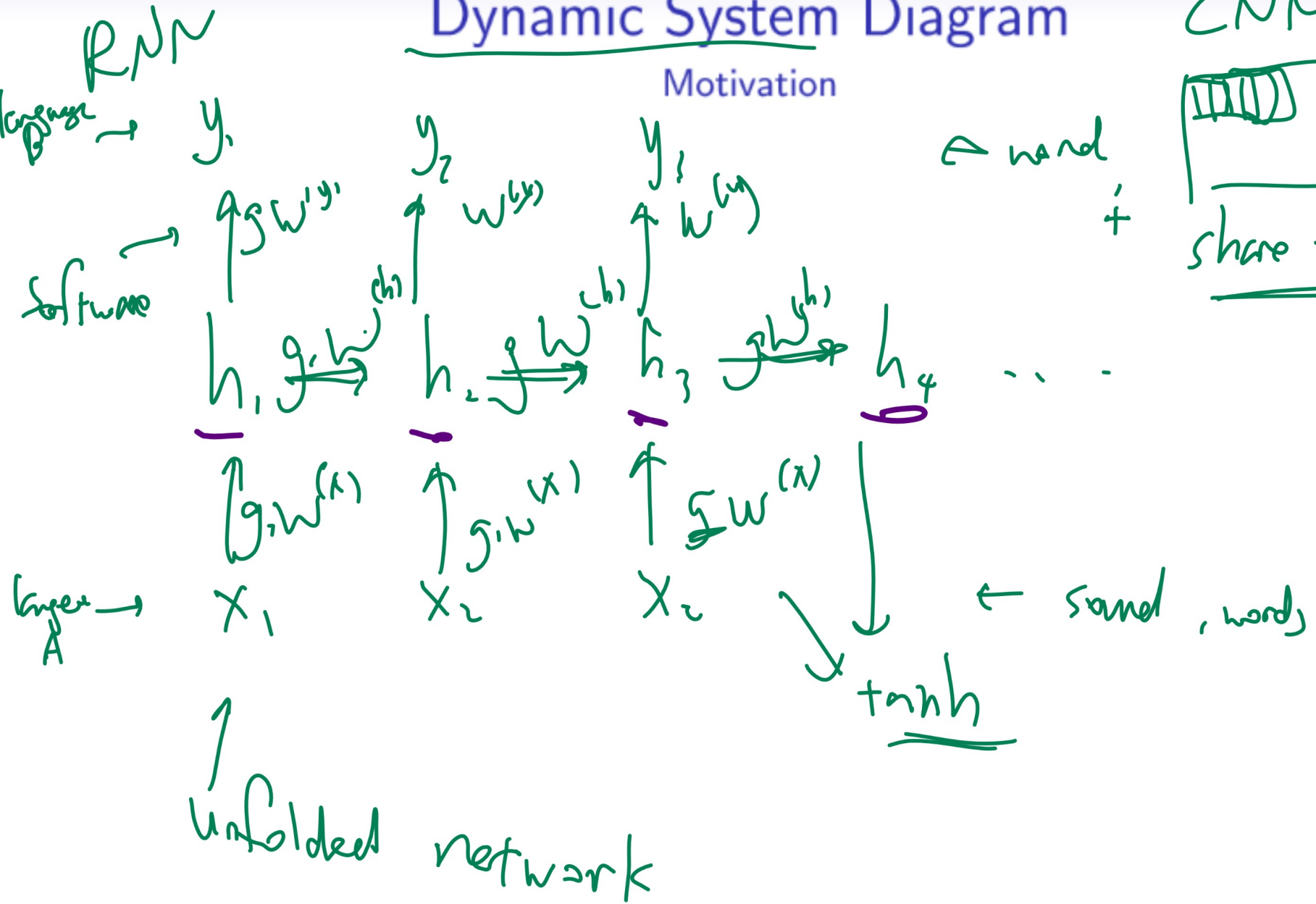
Young Wu

Based on lecture slides by Jerry Zhu, Yingyu Liang, and Charles Dyer

June 22, 2020
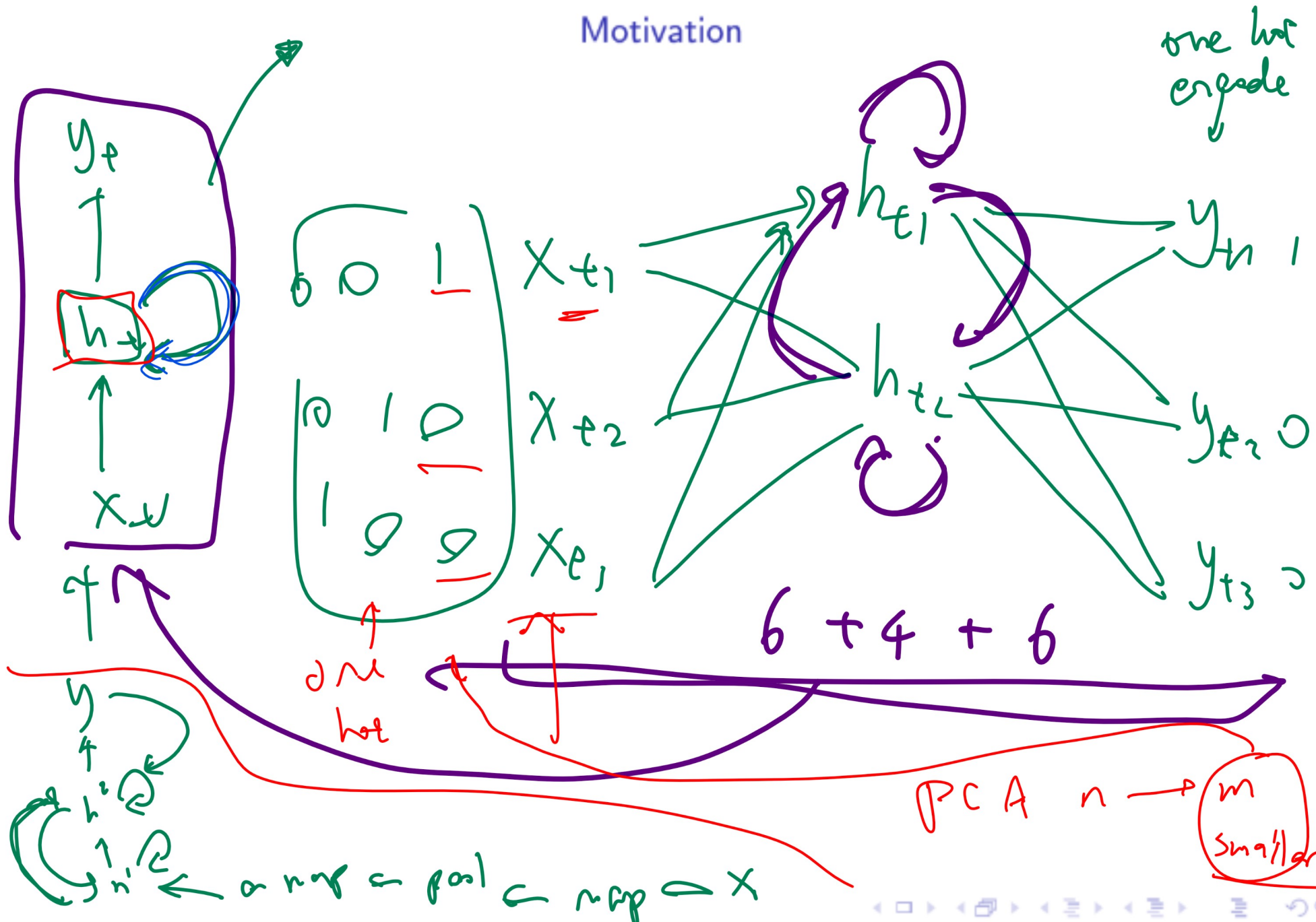
# Dynamic System Diagram

## Motivation

RNN

language → $y_1$

software →

CNN

← word

+

share weights

$$\uparrow g, w^{(y)} \qquad y_2 \qquad \uparrow w^{(y)} \qquad y_3 \qquad \uparrow h^{(y)}$$

$$h_1 \xrightarrow{g, h^{(h)}} h_2 \xrightarrow{g, w^{(h)}} h_3 \xrightarrow{g, w^{(h)}} h_4 \quad \dots$$

$$\uparrow g, w^{(x)} \qquad \uparrow g, w^{(x)} \qquad \uparrow g, w^{(x)}$$

language → $X_1$ \qquad $X_2$ \qquad $X_2$ \qquad ← sound, words

$$\downarrow \tanh$$

unfolded network

# Recurrent Neural Network Structure Diagram

## Motivation

# Activation Functions

## Definition

- The hidden layer activation function can be the tanh activation, and the output layer activation function can be the softmax function.

$$z_t^{(x)} = W^{(x)} x_t + W^{(h)} a_{t-1}^{(x)} + b^{(x)}$$

$$a_t^{(x)} = g\left(z_t^{(x)}\right), g\left(\boxed{\cdot}\right) = \tanh\left(\boxed{\cdot}\right)$$

$$z_t^{(y)} = W^{(y)} a^{(1,t)} + b^{(y)}$$

$$a_t^{(y)} = g\left(z_t^{(y)}\right), g\left(\boxed{\cdot}\right) = \text{softmax}\left(\boxed{\cdot}\right)$$

*(handwritten annotations: $h_{t-1}$, $h_t$, $W^{(x)} x_t$)*

# Cost Functions
## Definition

- Cross entropy loss is used with softmax activation as usual.

$$C_t = H\left(y_t, a_t^{(y)}\right)$$

$$C = \sum_t C_t$$
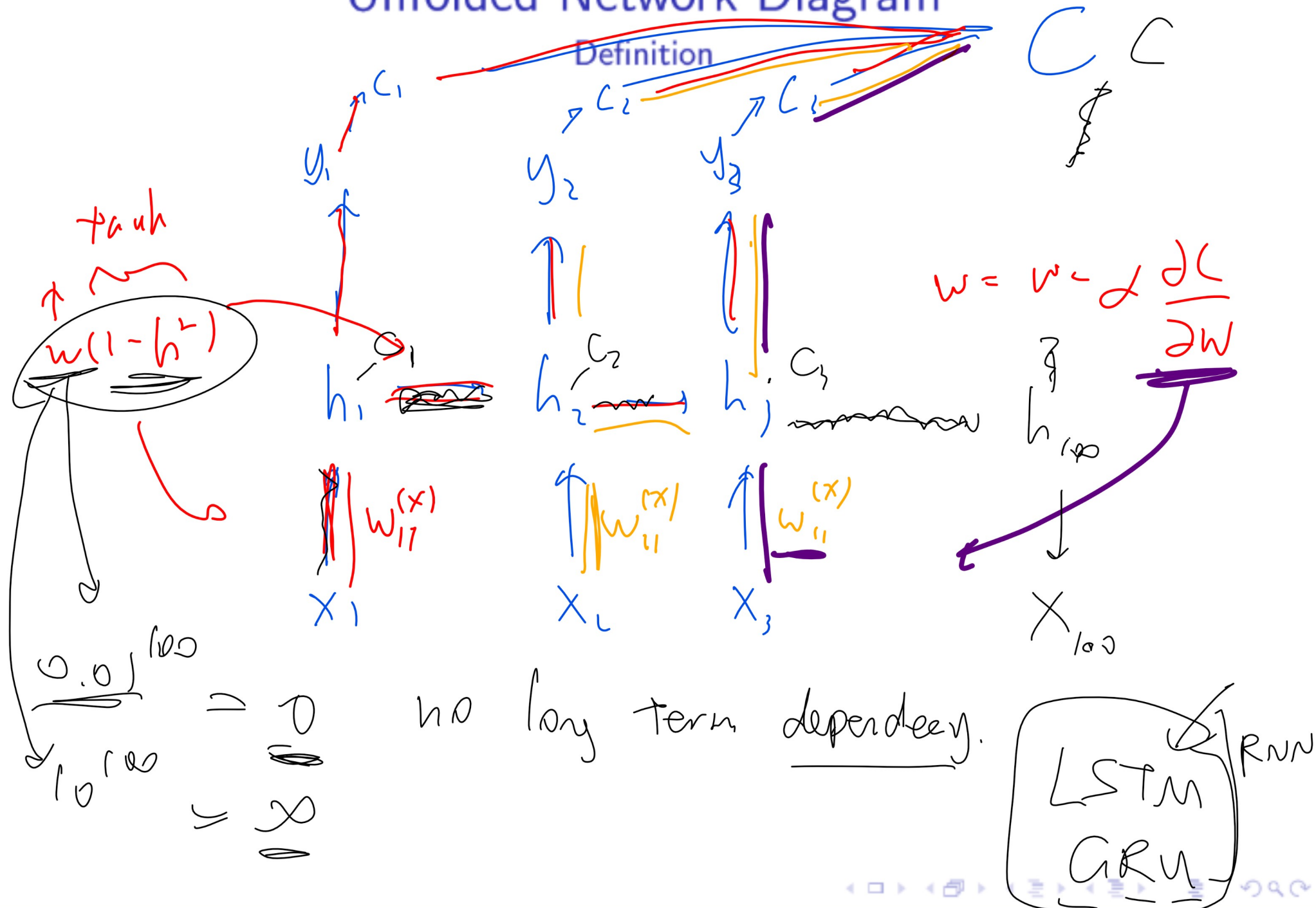
# BackPropogation Through Time
## Definition

- The gradient descent algorithm for recurrent neural networks is called BackPropogation Through Time (BPTT). The update procedure is the same as standard neural networks using the chain rule.

$$w = w - \alpha \frac{\partial C}{\partial w}$$
$$b = b - \alpha \frac{\partial C}{\partial b}$$

# Unfolded Network Diagram

## Definition

$tanh$

$w(1-h^2)$

$C_1$

$y_1$

$C_2$

$y_2$

$C_3$

$y_3$

$C$ $C$

$h_1$

$C_2$

$h_2$

$C_3$

$h_j$

$w = w - \alpha \dfrac{\partial C}{\partial w}$

$h_{100}$

$w_{11}^{(x)}$

$w_{11}^{(x)}$

$w_{11}^{(x)}$

$X_1$

$X_2$

$X_3$

$X_{100}$

$0.0]^{100}$

$ca_l [0.0]$

$\Rightarrow 0$

$[0]^{100}$

$\Rightarrow \infty$

no long term dependecy.

LSTM
GRU

RNN

# Backpropagation Diagram 1
## Definition

# Backpropagation Diagram 2

## Definition

# Vanishing and Exploding Gradient
## Discussion

- If the weights are small, the gradient through many layers will shrink exponentially. This is called the vanishing gradient problem.

- If the weights are large, the gradient through many layers will grow exponentially. This is called the exploding gradient problem.

- Fully connected and convolutional neural networks only have a few hidden layers, so vanishing and exploding gradient is not a problem in training those networks.

- In a recurrent neural network, if the sequences are long, the gradients can easily vanish or explode.

# Long Short Term Memory
## Discussion

- Long Short Term Memory (LSTM) network adds more connected hidden units for memories controlled by gates. The activation functions used for these gates are usually logistic functions.

- An LSTM unit usually contains an input gate, an output gate, and a forget gate, to keep track of the dependencies in the input sequence.

# Long Short Term Memory Diagram

## Discussion

# Gated Recurrent Unit
## Discussion

- Gated Recurrent Unit (GRU) does something similar to an LSTM unit.

- A GRU contains input and forget gates, and does not contain an output gate.

# Gated Recurrent Unit Diagram

## Discussion