

# TUTORIAL – 2

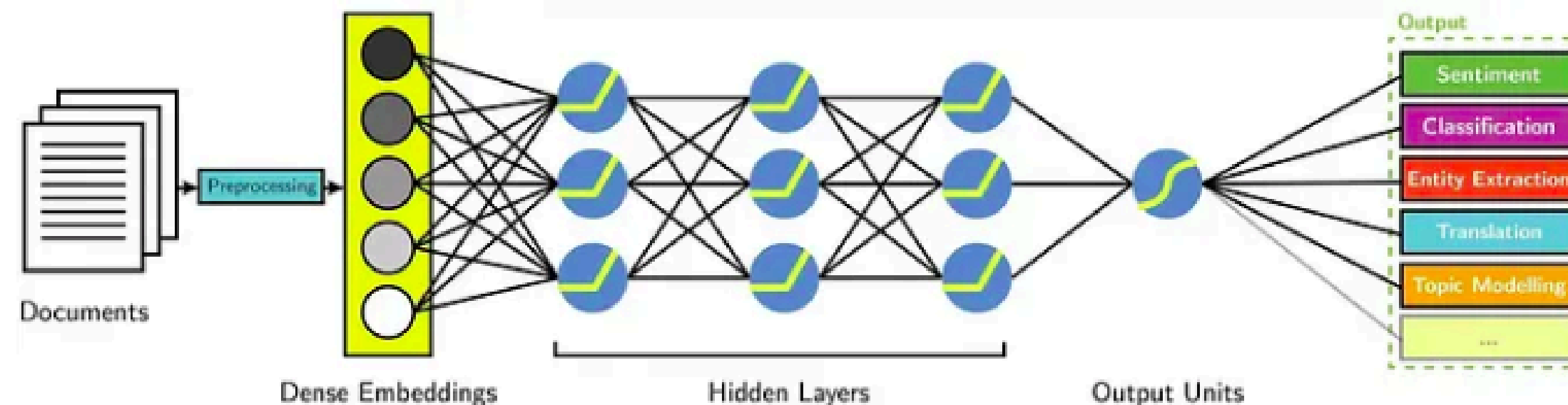
# NEXT WORD PREDICTION

- For example, if we have some text  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$ , then the probability of this text (according to the Language Model) is:

$$\begin{aligned} P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}) &= P(\mathbf{x}^{(1)}) \times P(\mathbf{x}^{(2)} | \mathbf{x}^{(1)}) \times \dots \times P(\mathbf{x}^{(T)} | \mathbf{x}^{(T-1)}, \dots, \mathbf{x}^{(1)}) \\ &= \prod_{t=1}^T P(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(1)}) \end{aligned}$$

  
This is what our LM provides

# FFNN



Source: Aylien

Input data: We need a Mathematical representation of the textual data for computational purposes.

Simple Idea????

Simple Idea: One hot encoding.

Can we do better? Can we learn the embeddings of the words as we train our model. We can use  
pytorch for embeddings

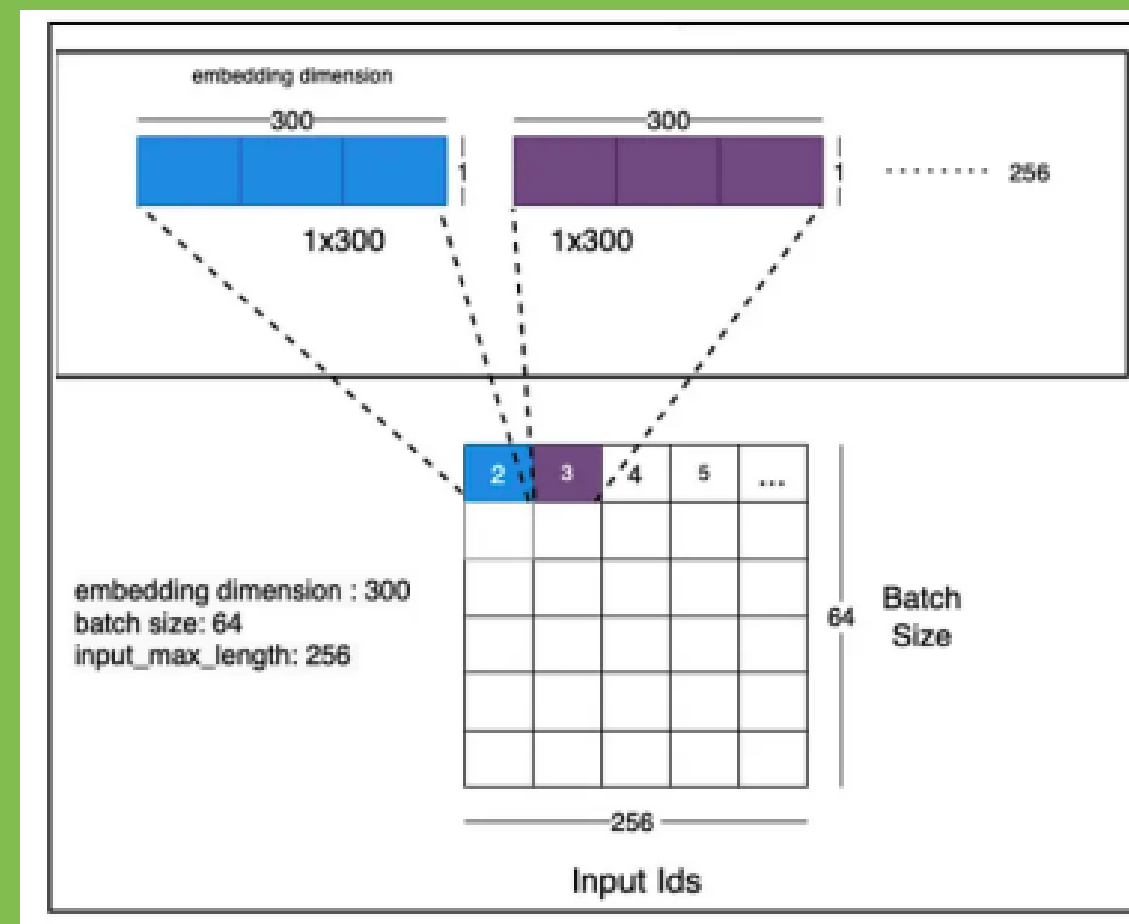
Input data: We need a Mathematical representation of the textual data for computational purposes.

Simple Idea????

Simple Idea: One hot encoding.

Can we do better? Can we learn the embeddings of the words as we train our model.

Embedding layer is also a Matrix, and what is the Matrix dimension? PyTorch provides a Module for Embedding layer, which we can initialize with appropriate dimensions



$$h_t = \sigma(\mathbf{W}_1 \mathbf{x}_t + \mathbf{b}_1)$$

$$\hat{y}_t = \text{softmax}(\mathbf{W}_2 \mathbf{h}_t + \mathbf{b}_2)$$

$W_1$  and  $b_1$  are the weight matrix and bias vector for the hidden layer, respectively.

They are learned during the training process during back propagation

Cross-Entropy Loss:

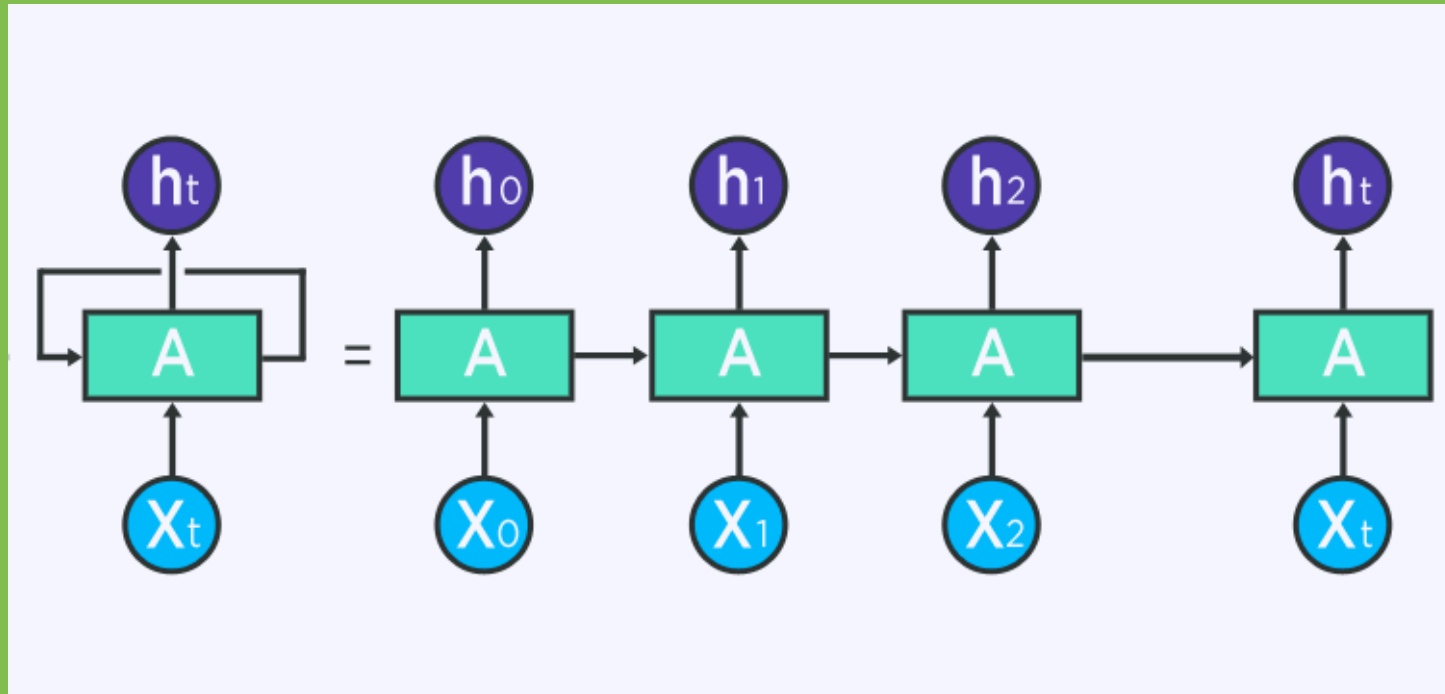
$$\text{Loss} = - \sum_{i=1}^N y_i \cdot \log(\hat{y}_i)$$

CE Loss and Adam optimizer can be used for back propagating the loss

$\hat{y}_t$  is the predicted probability distribution over all the possible output words at time step  $t$ .

# RNN

concept of memory - introduces a hidden state that carries information from previous inputs.



## A RNN Language Model

output distribution

$$\hat{y}^{(t)} = \text{softmax}(U h^{(t)} + b_2) \in \mathbb{R}^{|V|}$$

hidden states

$$h^{(t)} = \sigma(W_h h^{(t-1)} + W_e e^{(t)} + b_1)$$

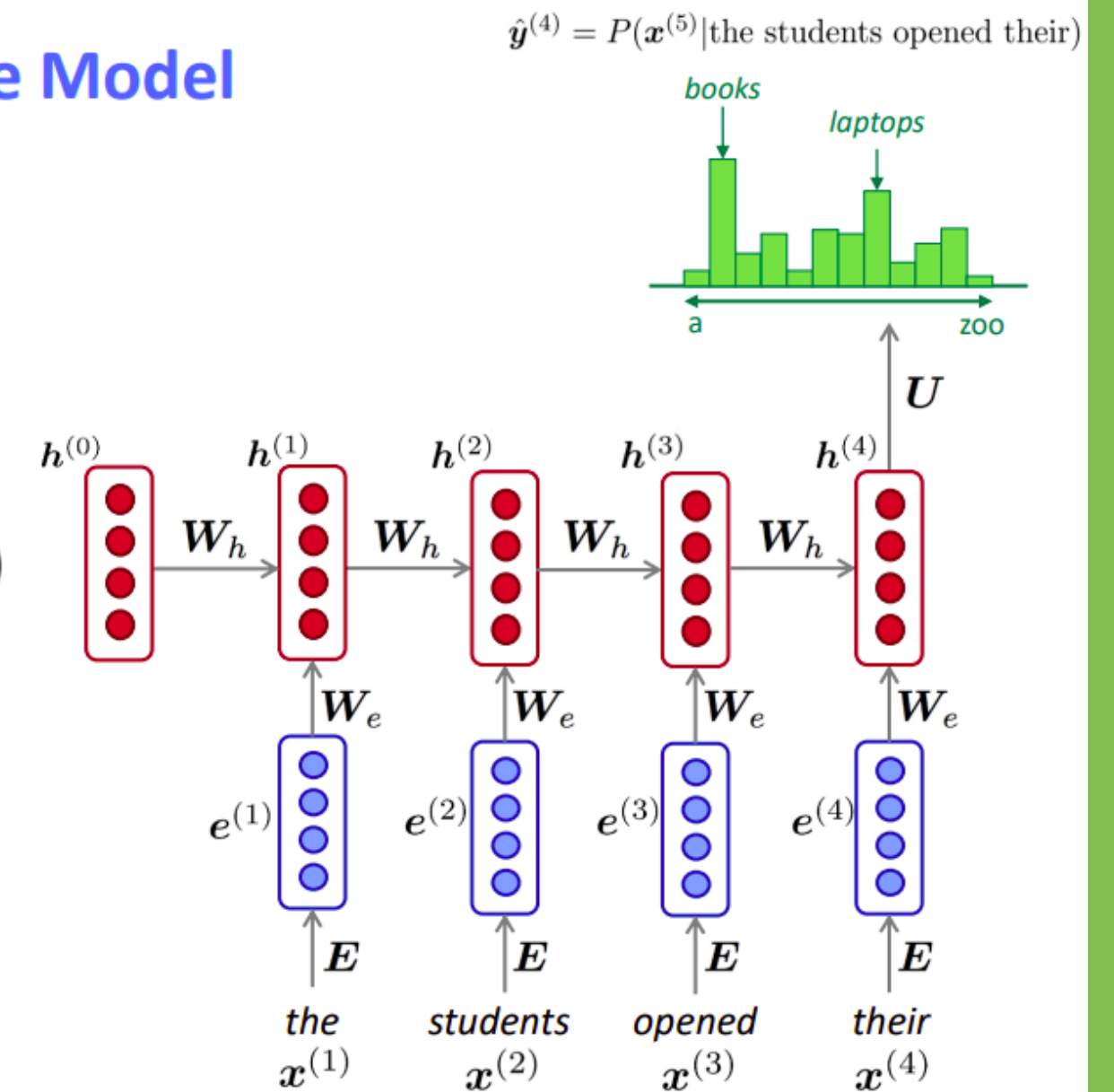
$h^{(0)}$  is the initial hidden state

word embeddings

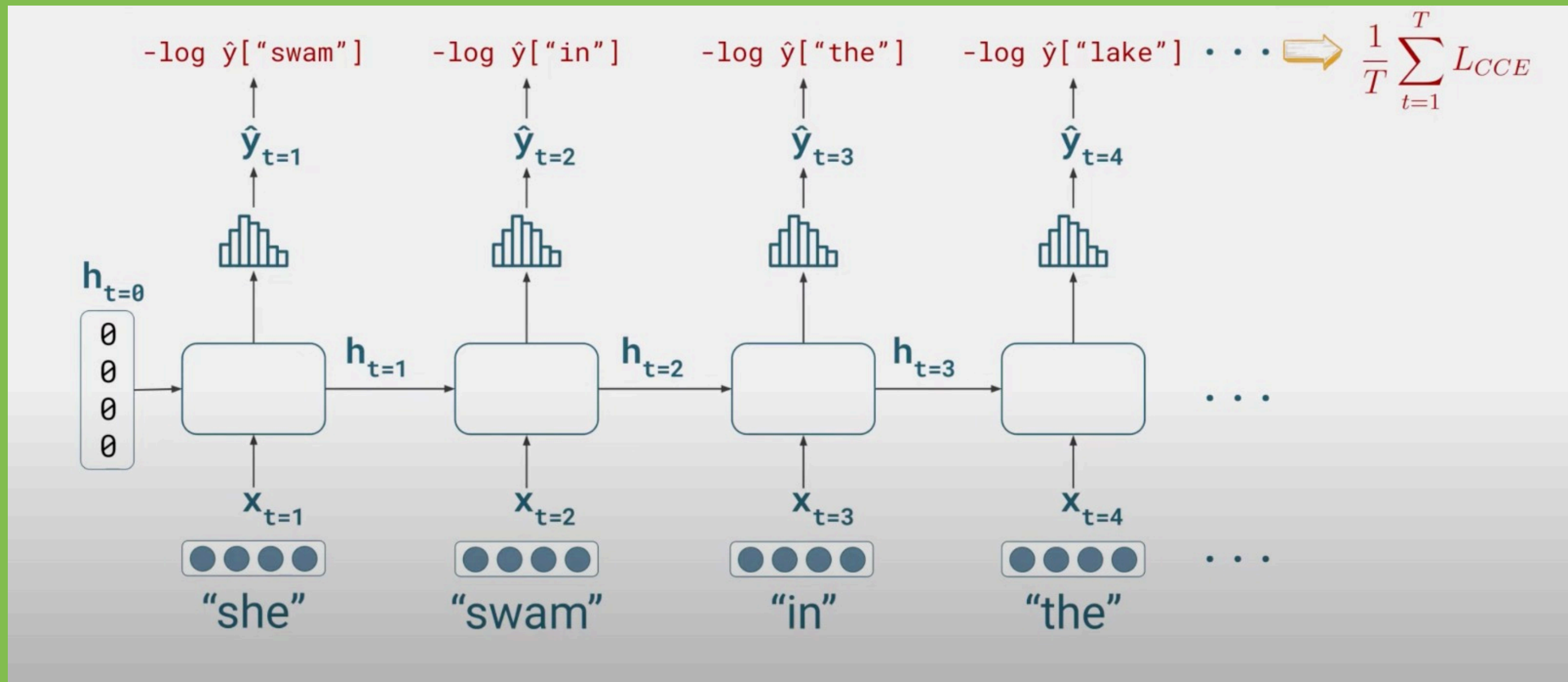
$$e^{(t)} = E x^{(t)}$$

words / one-hot vectors

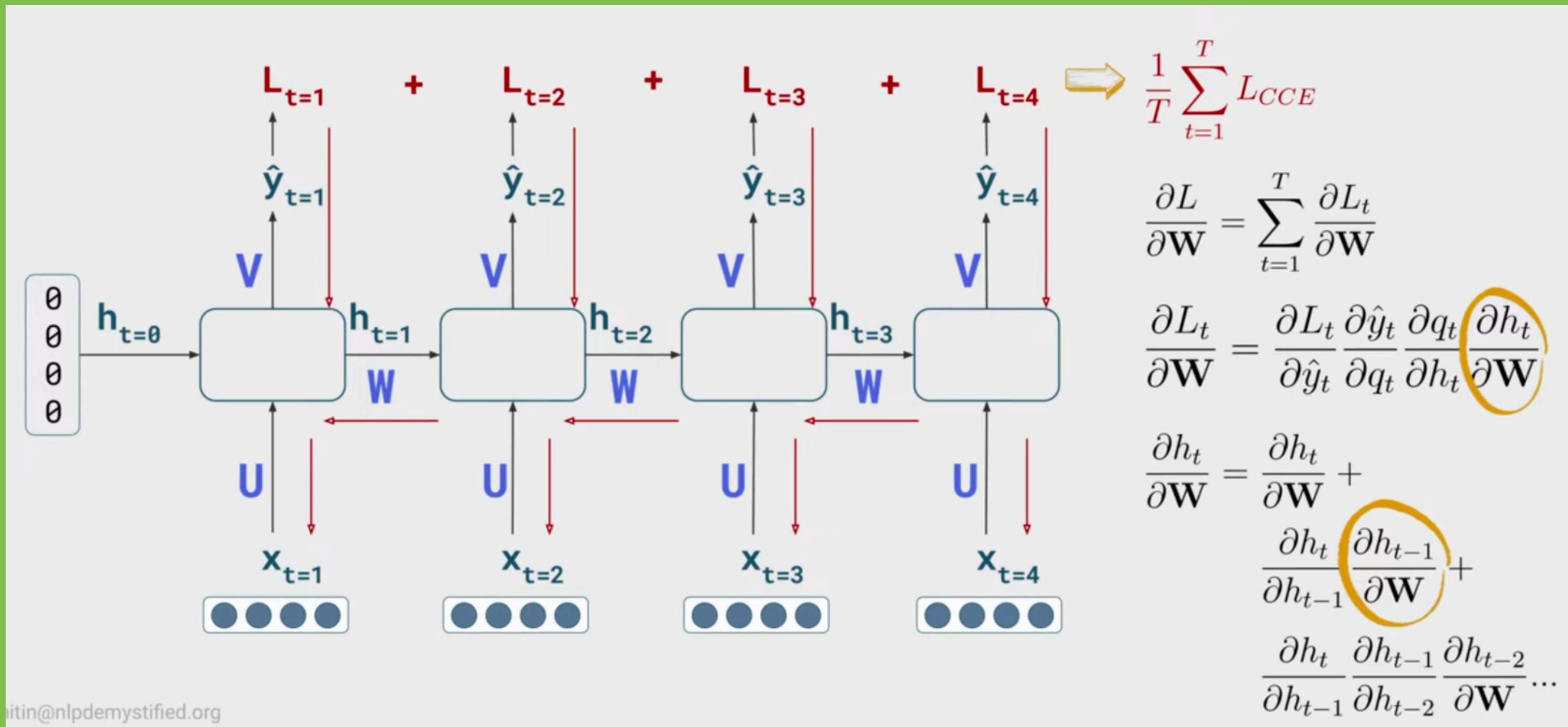
$$x^{(t)} \in \mathbb{R}^{|V|}$$



# TRAINING RNN



# BACKPROPOGATION



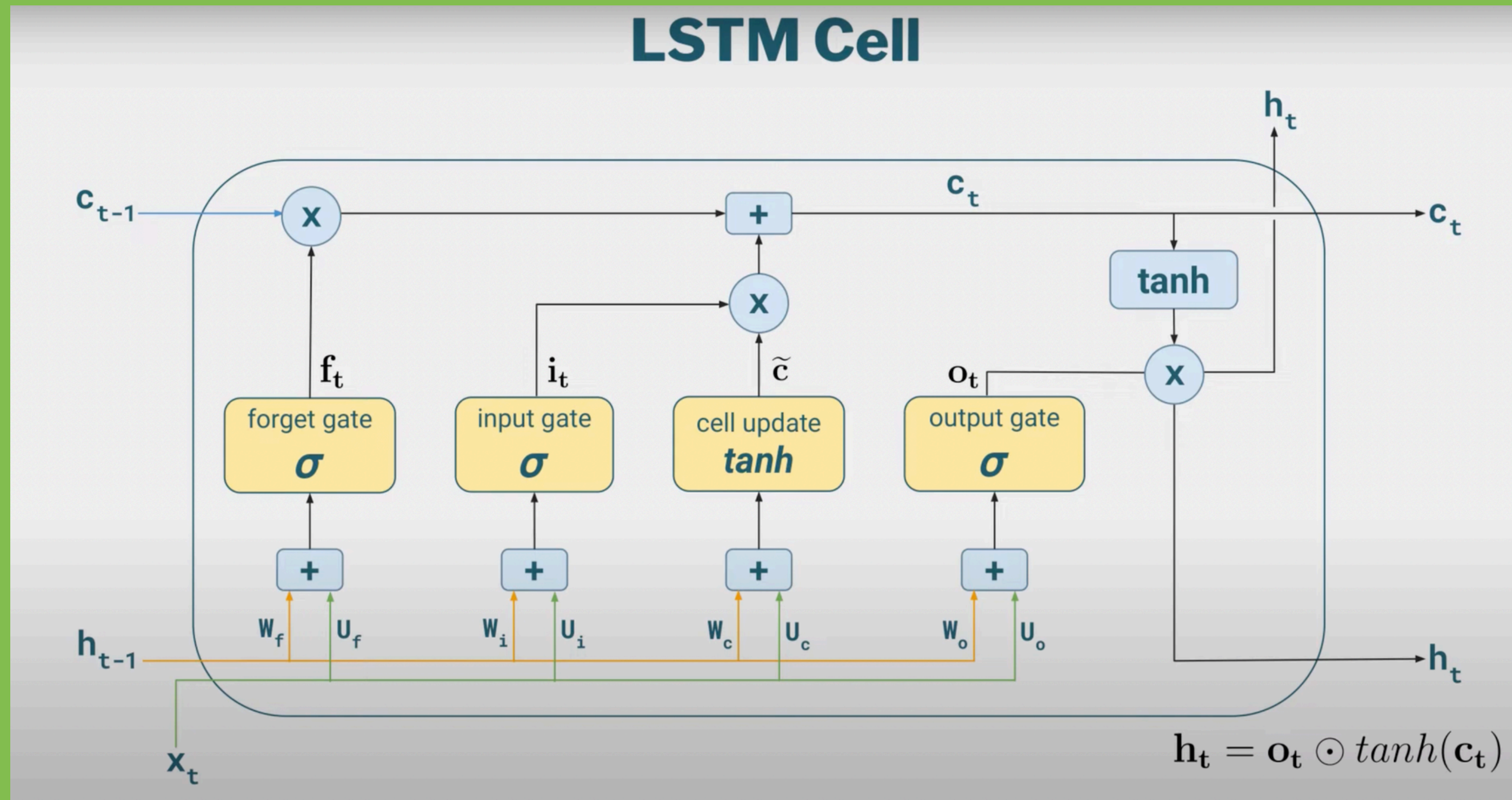


# DISADVANTAGE

He was supposed to meet his friends at the show but something came up at work, so he called them



# LSTM



**DOUBTS?**