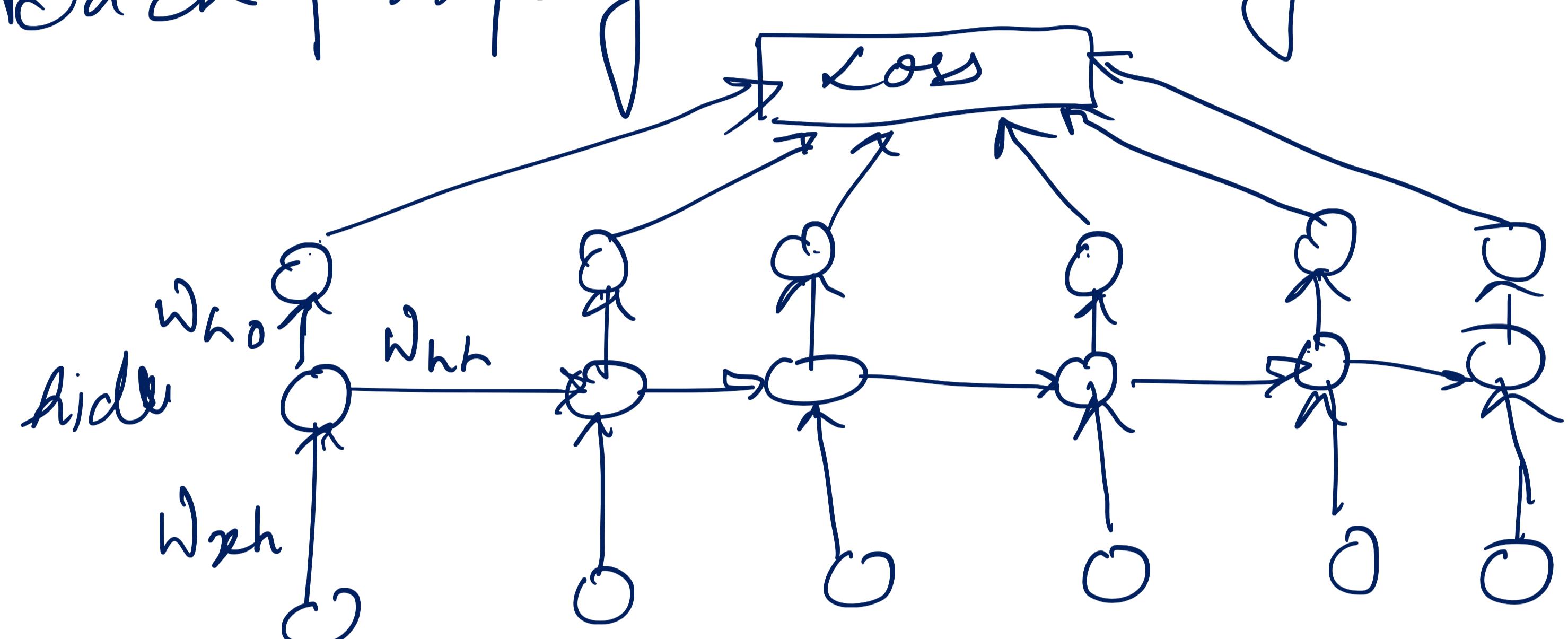


RNN's

Back propagation through time.



- 1) Vanishing gradient
- 2) Exploding gradient
 - ↳ One possible way of handling is Gradient Clipping

Difficult to handle

E.g. You are training a model on ~~a~~ set of wikipedia documents.

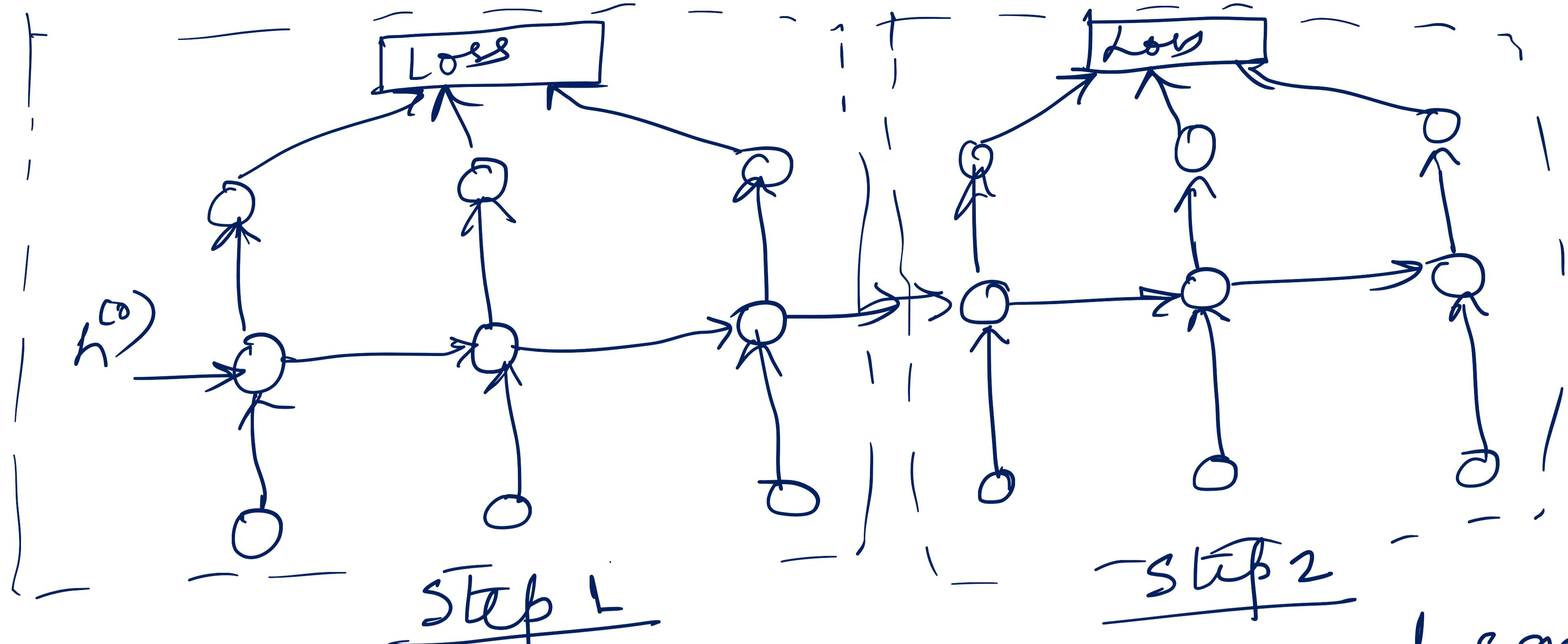
Each input unit receives a word from a wikipedia document (in vector format)

For one update of the weights, run the forward propagation to compute the loss and then backward through the entire sequence to compute the gradient

- Even for a single update the number of gradients to be calculated equals the no. of words in the document.

→ Result is slow computation

- Solution is to use a truncated back propagation of the sequence instead of whole sequence.
- Run forward and backward through chunks



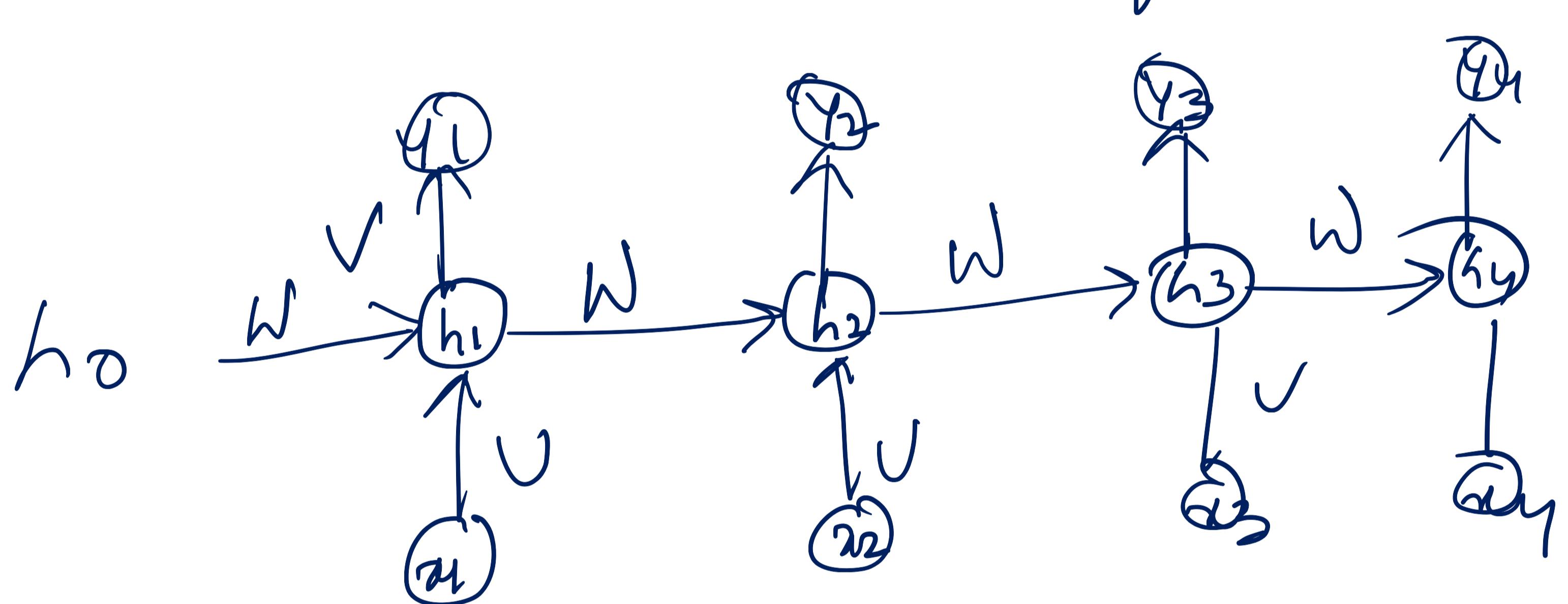
- Step 1

 - Generate subsequences and compute loss only over this subsequence of data.

Step 2

 - Use a batch of data as subsequence for calculating the loss
 - Carry the hidden states forward in time forever, but only back propagate for some smaller no. of steps -
 - i.e back propagate through the next ~~back~~ batch.

Derivation of the Gradient for BPTT



$$h_t = \psi(x_t) + w\phi(h_{t-1}) - \theta$$

$$y_t = \cancel{f(t)} + \phi(h_t) - \boxed{2}$$

Let the max stages be C
and let t be the current no. of stages.

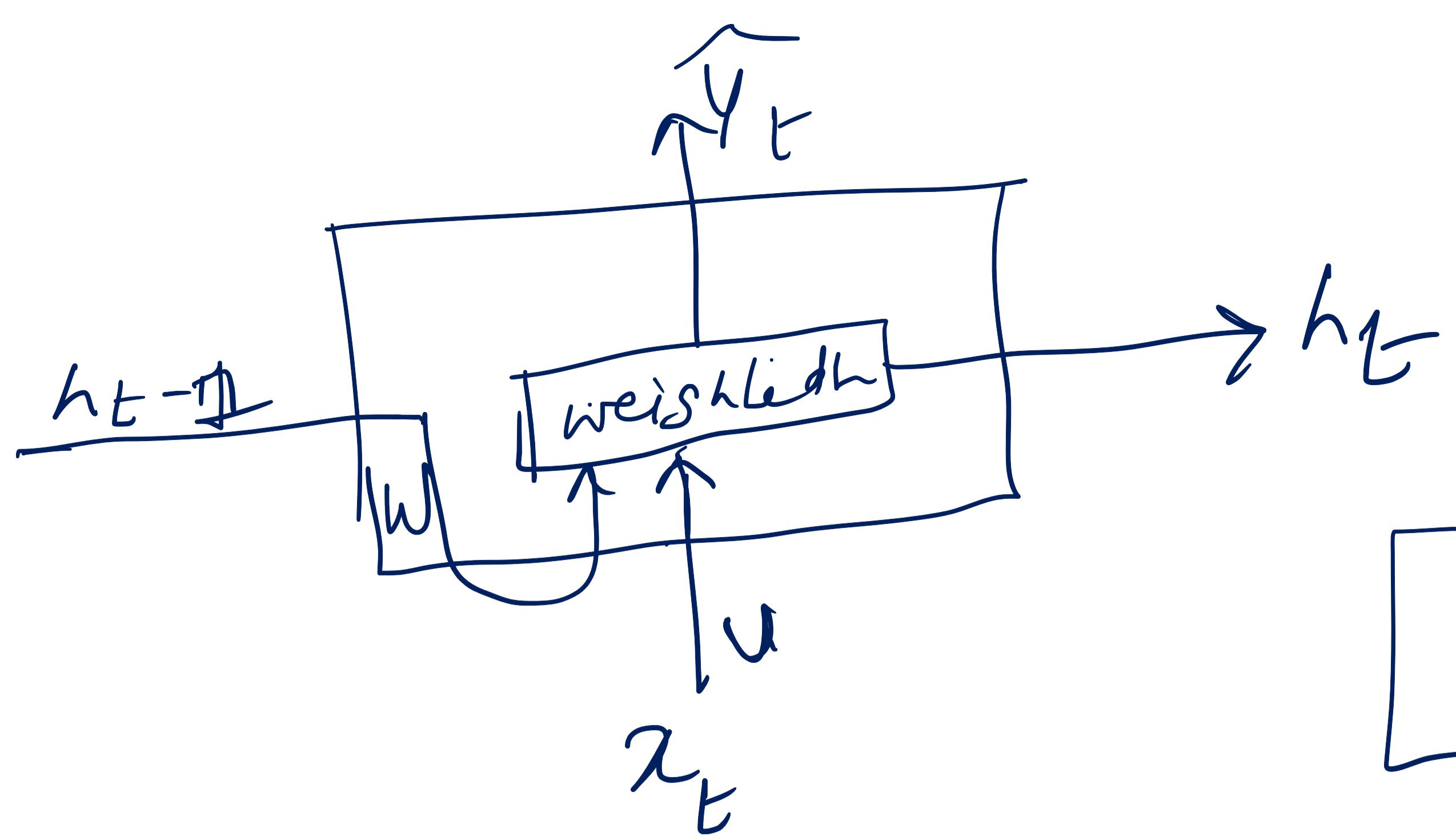
$$\frac{\partial E}{\partial w} = \sum_{t=1}^T \frac{\partial E}{\partial y_t} \cdot \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial w}$$

$$\frac{\partial h_t}{\partial h_k} = \frac{\partial h_t}{\partial h_{t-1}} \cdot \frac{\partial h_{t-1}}{\partial h_{t-2}} \cdot \dots \cdot \frac{\partial h_{k+1}}{\partial h_k}.$$

$$= \cancel{\pi} h_k = \frac{\partial h_i +'}{j_{h_i}} \stackrel{i=k}{=} \pi^{t-1} \circ 3$$

$$\text{from Eq ①} \quad \frac{\partial h_i^+}{\partial \omega_i^-} = \cancel{\omega^T \phi(h_i^+)}$$

$$\text{so } \frac{\partial h^t}{\partial h_k} = \prod_{l=k}^{t-1} w^T \text{diag}(\phi(h_i)).$$



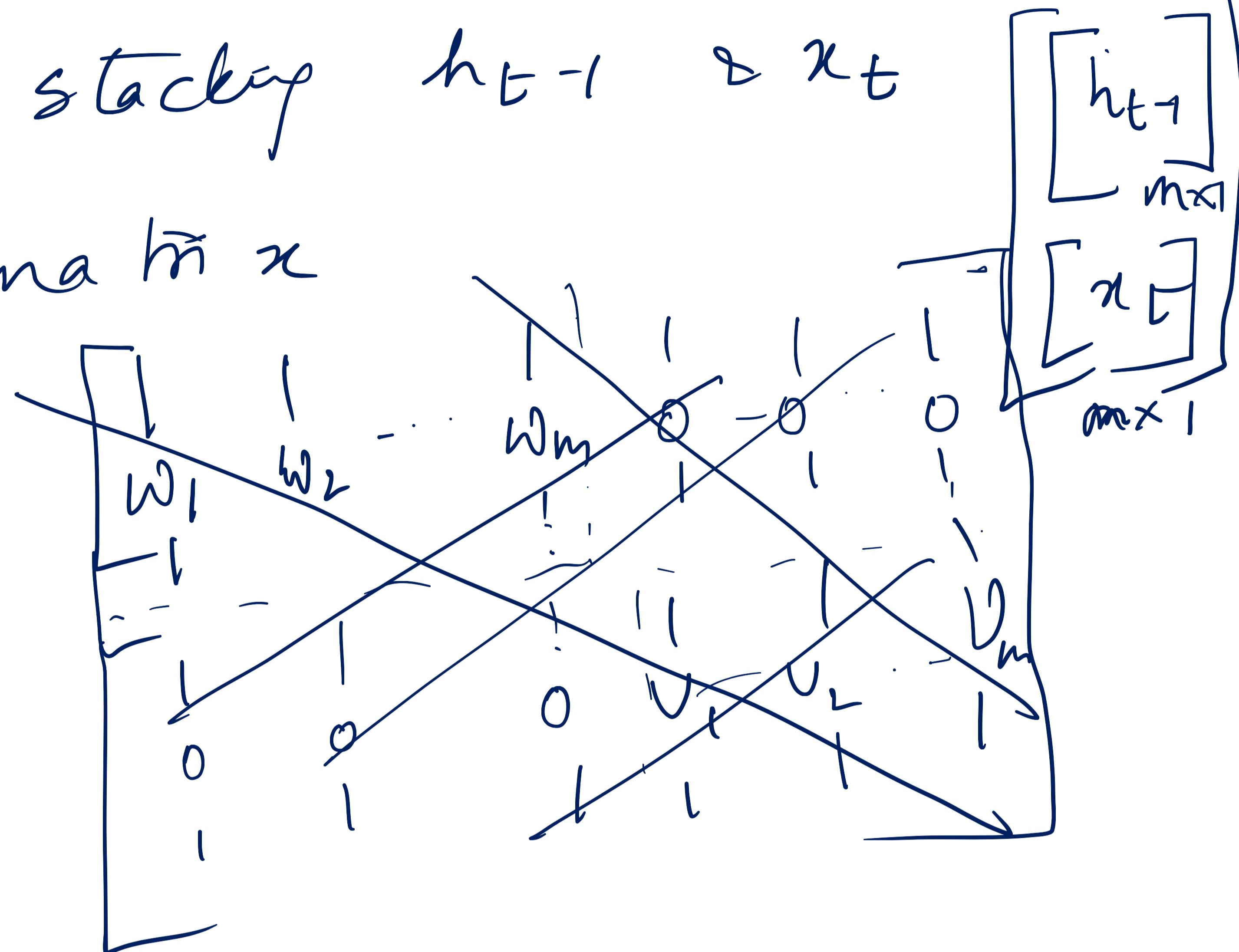
$$h_t = W h_{t-1} + U x_t$$

Create a new input & say $[h_{t-1}, x_t]$

i.e. by stacking h_{t-1} & x_t

Create a new weight matrix \tilde{w}

$$\tilde{w}_a = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ w_1 & w_2 & -w_m & -U_1 & U_2 & -U_m \\ 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix}$$



$$\tilde{w}_a [h_{t-1}, x_t] = a_t$$

$$h_t = \phi(\tilde{w}_a [h_{t-1}, x_t] + b_a)$$

For the detailed derivation refer to slides 138'4.

Applications of RNN

Language model and Sequence Generation.

Generate text by training a model.

What is language modeling?

Answer: Consider a language with a vocabulary of words v_1, v_2, \dots, v_T

Given a given sequence of words y_1, y_2, \dots, y_n "I am a good boy"

For a language model, the task is

