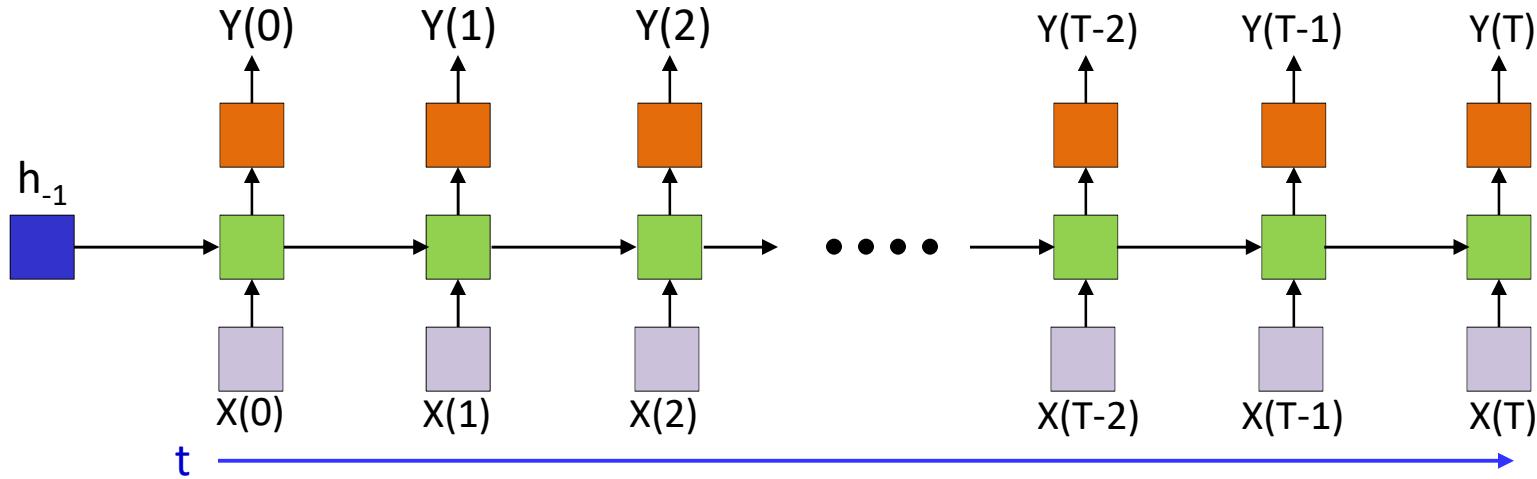


# How do we *train* the network



- Back propagation through time (BPTT)
- Given a collection of *sequence* inputs
  - $(\mathbf{X}_i, \mathbf{D}_i)$ , where
  - $\mathbf{X}_i = X_{i,0}, \dots, X_{i,T}$
  - $\mathbf{D}_i = D_{i,0}, \dots, D_{i,T}$
- Train network parameters to minimize the error between the output of the network  $\mathbf{Y}_i = Y_{i,0}, \dots, Y_{i,T}$  and the desired outputs
  - This is the most generic setting. In other settings we just “remove” some of the input or output entries

- ① Note that  $D(0), D(1), \dots, D(T)$  are the desired output sequence.
- ②  $Y(0), Y(1), \dots, Y(T)$  is the output sequence generated by the RNN.
- ③ Thus, the loss will capture distance between two sequences and is denoted as  $DIV(\{Y(0), Y(1), \dots, Y(T)\}, \{D(0), D(1), \dots, D(T)\})$ .
- ④ An example of the loss is where  $D(t)$  is compared to only  $Y(t)$  and the overall loss is a sum of local loss values as follows.

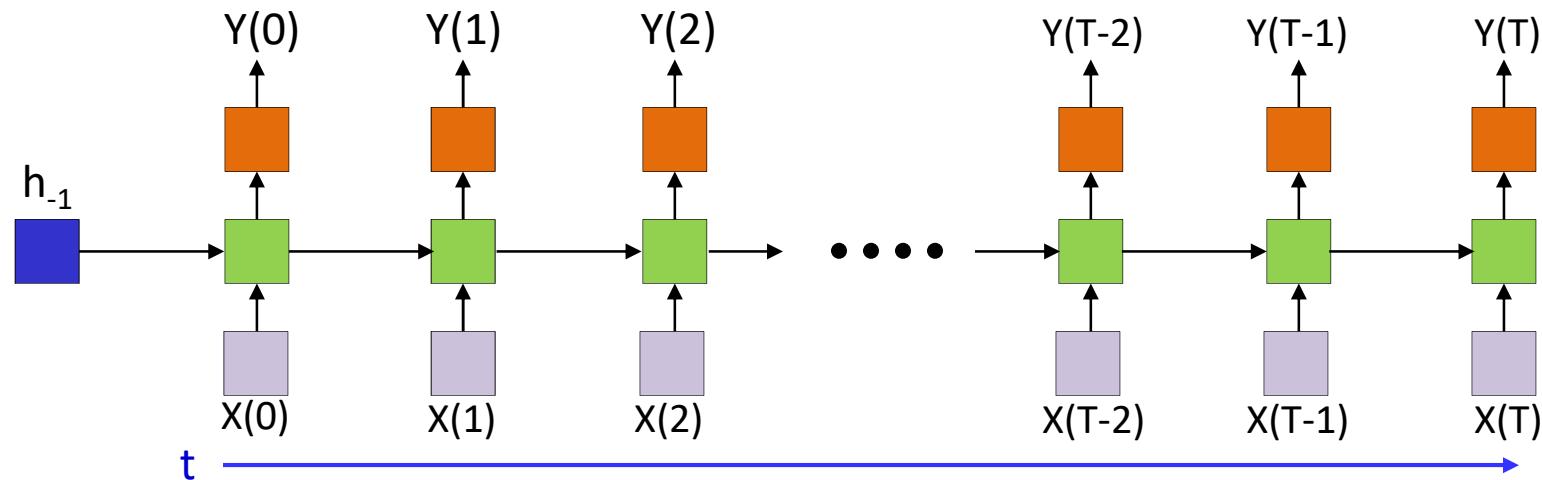
$$DIV(\{Y(0), Y(1), \dots, Y(T)\}, \{D(0), D(1), \dots, D(T)\})$$

$$= \sum_{t=0}^T Div(Y(t), D(t))$$

$$\Rightarrow \frac{\partial DIV(\{Y(0), Y(1), \dots, Y(T)\}, \{D(0), D(1), \dots, D(T)\})}{\partial Y(t)}$$

$$= \frac{\partial Div(Y(t), D(t))}{\partial Y(t)}$$

# Training: Forward pass



- For each training input:
- Forward pass: pass the entire data sequence through the network, generate outputs

- ① Let the input  $X(t)$  has dimension  $N$  for all  $t = 0, \dots, T$ .
- ② Let the hidden layer has  $N_1$  nodes.
- ③ Then the pre-activation value  $Z_i^{(1)}(t)$  of the  $i^{th}$  node in the hidden layer at time  $t$  is computed as follow. Note that it takes  $X(t)$  and  $\mathbf{h}(t-1)$  as input.

$$Z_i^{(1)}(t) = \sum_{j=1}^N w_{ji}^{(1)} X_j(t) + \sum_{j=1}^{N_1} w_{ji}^{(11)} h_j^{(1)}(t-1) + b_i^{(1)}, \quad i = 1 \dots N_1$$

- ④ Output of  $i^{th}$  node in the hidden layer is found by applying activation function  $f^{(1)}$  on  $Z_i^{(1)}(t)$ .

$$h_i^{(1)}(t) = f^{(1)} \left( Z_i^{(1)}(t) \right), \quad i = 1 \dots N_1$$

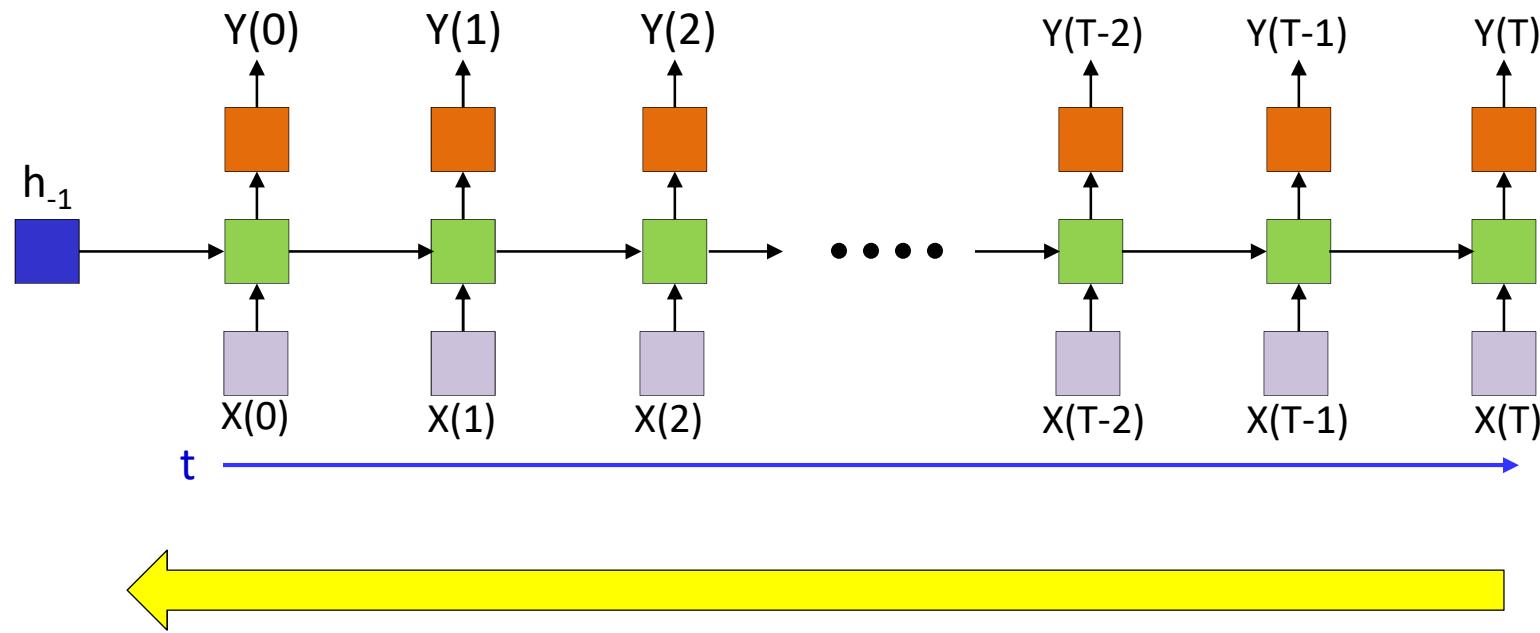
- ① Let the output  $Y(t)$  is an  $M$ -dimensional vector.
- ② Then the pre-activation value  $Z_i^{(2)}(t)$  of the  $i^{th}$  node in the output layer at time  $t$  is computed as follow. Note that it takes  $\mathbf{h}(t)$  as input.

$$Z_i^{(2)}(t) = \sum_{j=1}^{N_1} w_{ji}^{(2)} h_j(t) + b_i^{(2)}, \quad i = 1 \dots M$$

- ③ Output of  $i^{th}$  node in the output layer is found by applying activation function  $f_i^{(2)}$  on  $Z_i^{(2)}(t)$ .

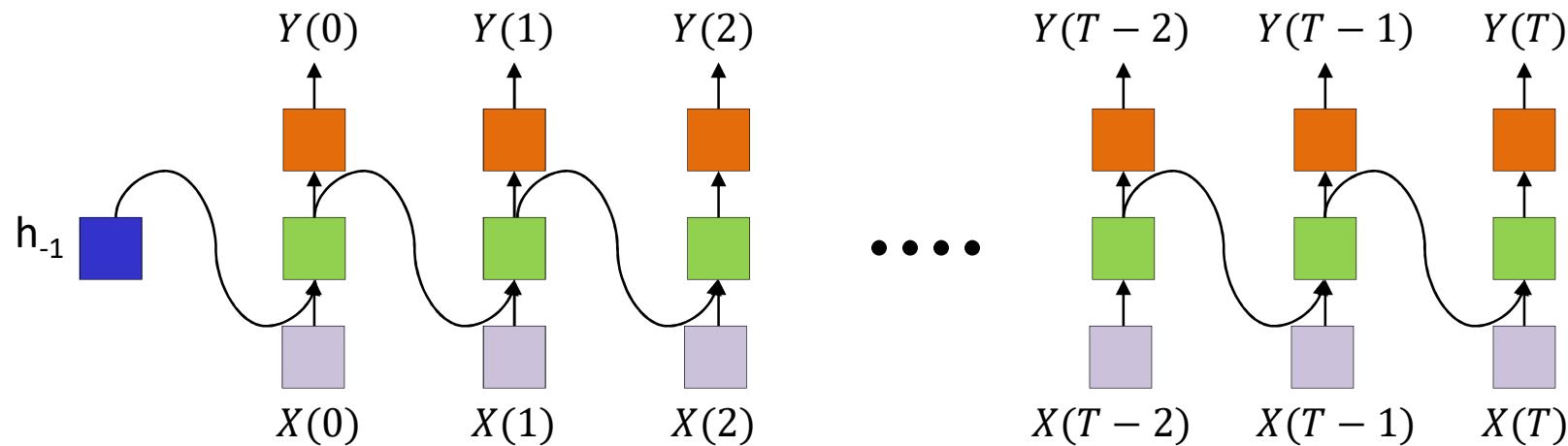
$$Y_i(t) = f_i^{(2)} \left( Z_1^{(2)}(t), \dots, Z_M^{(2)}(t) \right), \quad i = 1 \dots M$$

# Training: Computing gradients



- For each training input:
- Backward pass: Compute gradients via backpropagation
  - *Back Propagation Through Time*

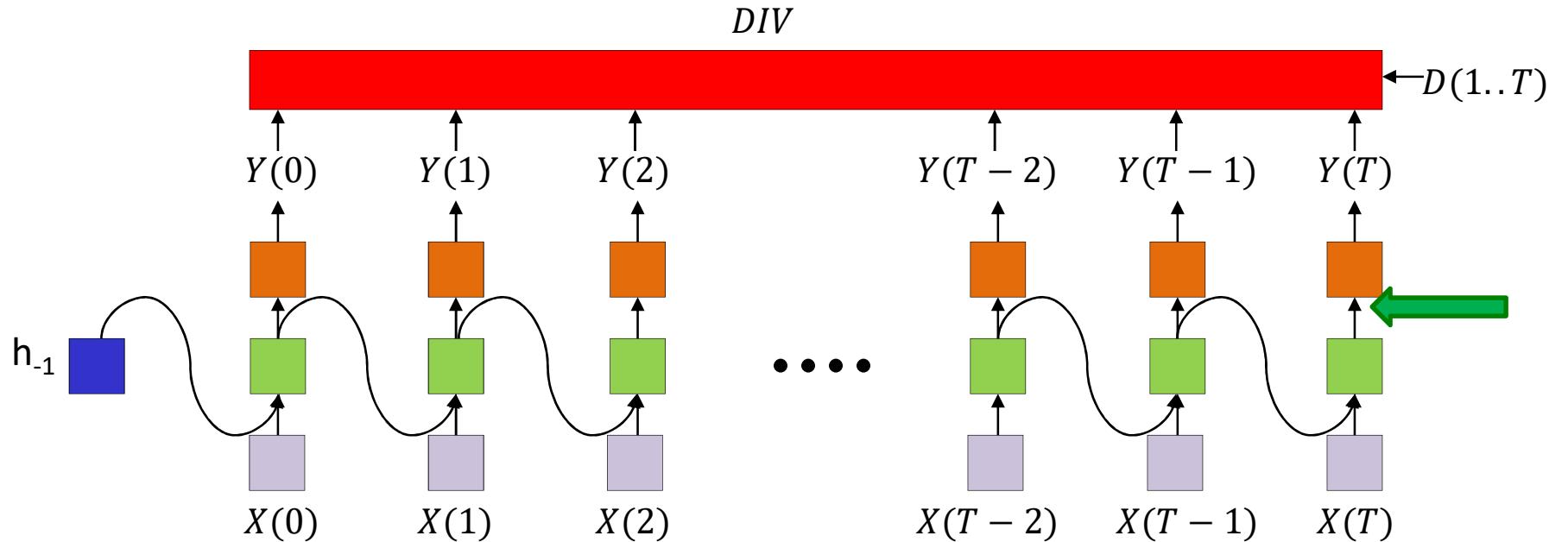
# Back Propagation Through Time



Will only focus on one training instance

All subscripts represent components and not training instance index

# Back Propagation Through Time



First step of backprop: Compute  $\frac{dDIV}{dY_i(T)}$  for all i

$$\nabla_{Z^{(2)}(T)} DIV = \nabla_{Y(T)} DIV \nabla_{Z^{(2)}(T)} Y(T)$$

Vector output activation

$$\frac{dDIV}{dZ_i^{(2)}(T)} = \frac{dDIV}{dY_i(T)} \frac{dY_i(T)}{dZ_i^{(2)}(T)}$$

OR

$$\frac{dDIV}{dZ_i^{(2)}(T)} = \sum_j \frac{dDIV}{dY_j(T)} \frac{dY_j(T)}{dZ_i^{(2)}(T)}$$

① We know that  $Z_i^{(2)}(T) = \sum_{j=1}^{N_1} w_{ji}^{(2)} h_j(T) + b_i^{(2)}$ ,  $i = 1 \dots M$

## ② Case 1: Vector output activation

- $Y_k(T) = f_k^{(2)}(Z_1^{(2)}(T), \dots, Z_M^{(2)}(T))$ ,  $k = 1 \dots M$
- An example of such an  $f_k^{(2)}$  is softmax function, i.e.

$$f_k^{(2)}(Z_1^{(2)}(T), \dots, Z_M^{(2)}(T)) = \frac{\exp(Z_k^{(2)}(T))}{\sum_{j=1}^M \exp(Z_j^{(2)}(T))}, \quad k = 1 \dots M.$$

- In this case,  $\frac{\partial DIV}{\partial Z_i^{(2)}(T)} = \sum_{j=1}^M \frac{\partial DIV}{\partial Y_j(T)} \frac{\partial Y_j(T)}{\partial Z_i^{(2)}(T)}$

## ③ Case 2: Scalar output activation

- $Y_k(T) = f_k^{(2)}(Z_k^{(2)}(T))$ ,  $k = 1 \dots M$
- An example of such an  $f_k^{(2)}$  is tanh function
- In this case,  $\frac{\partial DIV}{\partial Z_i^{(2)}(T)} = \frac{\partial DIV}{\partial Y_i(T)} \frac{\partial Y_i(T)}{\partial Z_i^{(2)}(T)}$

$$\text{Compute } \nabla_{Z^{(2)}(T)} DIV = \left[ \frac{\partial DIV}{\partial Z_1^{(2)}(T)} \quad \cdots \quad \frac{\partial DIV}{\partial Z_M^{(2)}(T)} \right]$$



### Case 1: Vector output activation

$$\begin{aligned} \nabla_{Z^{(2)}(T)} DIV &= \left[ \sum_{j=1}^M \frac{\partial DIV}{\partial Y_j(T)} \frac{\partial Y_j(T)}{\partial Z_1^{(2)}(T)} \quad \cdots \cdots \quad \sum_{j=1}^M \frac{\partial DIV}{\partial Y_j(T)} \frac{\partial Y_j(T)}{\partial Z_M^{(2)}(T)} \right] \\ &= \left[ \frac{\partial DIV}{\partial Y_1(T)} \quad \cdots \cdots \quad \frac{\partial DIV}{\partial Y_M(T)} \right] \begin{bmatrix} \frac{\partial Y_1(T)}{\partial Z_1^{(2)}(T)} & \frac{\partial Y_1(T)}{\partial Z_2^{(2)}(T)} & \cdots & \frac{\partial Y_1(T)}{\partial Z_M^{(2)}(T)} \\ \frac{\partial Y_2(T)}{\partial Z_1^{(2)}(T)} & \frac{\partial Y_2(T)}{\partial Z_2^{(2)}(T)} & \cdots & \frac{\partial Y_2(T)}{\partial Z_M^{(2)}(T)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial Y_M(T)}{\partial Z_1^{(2)}(T)} & \frac{\partial Y_M(T)}{\partial Z_2^{(2)}(T)} & \cdots & \frac{\partial Y_M(T)}{\partial Z_M^{(2)}(T)} \end{bmatrix} \\ &= \nabla_{Y(T)} DIV \cdot \nabla_{Z^{(2)}(T)} Y(T) \end{aligned}$$

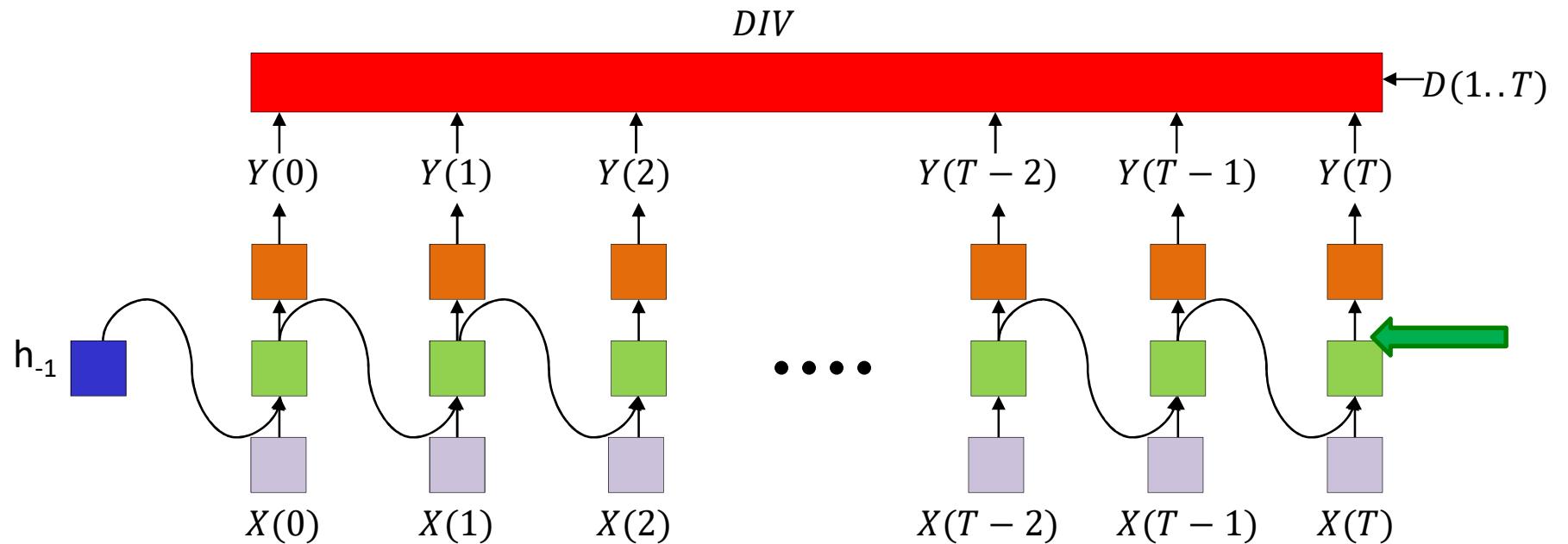
$$\text{Compute } \nabla_{Z^{(2)}(T)} DIV = \left[ \frac{\partial DIV}{\partial Z_1^{(2)}(T)} \cdots \frac{\partial DIV}{\partial Z_M^{(2)}(T)} \right]$$



### Case 2: Scalar output activation

$$\begin{aligned} \nabla_{Z^{(2)}(T)} DIV &= \left[ \frac{\partial DIV}{\partial Y_1(T)} \frac{\partial Y_1(T)}{\partial Z_1^{(2)}(T)} \cdots \cdots \frac{\partial DIV}{\partial Y_M(T)} \frac{\partial Y_M(T)}{\partial Z_M^{(2)}(T)} \right] \\ &= \left[ \frac{\partial DIV}{\partial Y_1(T)} \cdots \cdots \frac{\partial DIV}{\partial Y_M(T)} \right] \begin{bmatrix} \frac{\partial Y_1(T)}{\partial Z_1^{(2)}(T)} & 0 & \cdots & 0 \\ 0 & \frac{\partial Y_2(T)}{\partial Z_2^{(2)}(T)} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{\partial Y_M(T)}{\partial Z_M^{(2)}(T)} \end{bmatrix} \\ &= \nabla_{Y(T)} DIV \cdot \nabla_{Z^{(2)}(T)} Y(T) \end{aligned}$$

# Back Propagation Through Time



$\frac{dDIV}{dY_i(T)}$  for all i

$$\frac{dDIV}{dZ_i^{(2)}(T)} = \frac{dDiv(T)}{dY_i(T)} \frac{dY_i(T)}{dZ_i^{(2)}(T)}$$

$$\frac{dDIV}{dh_i(T)} = \sum_j \frac{dDIV}{dZ_j^{(2)}(T)} \frac{dZ_j^{(2)}(T)}{dh_i(T)} = \sum_j w_{ij}^{(2)} \frac{dDIV}{dZ_j^{(2)}(T)}$$

$$\nabla_{h(T)} DIV = \nabla_{Z^{(2)}(T)} DIV W^{(2)}$$

$$\text{Compute } \nabla_{\mathbf{h}(T)} DIV = \begin{bmatrix} \frac{\partial DIV}{\partial h_1(T)} & \cdots & \frac{\partial DIV}{\partial h_{N_1}(T)} \end{bmatrix}$$

- ① Note that pre-activation value of  $j^{th}$  node in the output layer at time  $T$  is  $Z_j^{(2)}(T) = \sum_{k=1}^{N_1} w_{kj}^{(2)} h_k(T) + b_j^{(2)}$ .
- ②  $h_i(T)$  is an input to all the  $Z_1^{(2)}(T), \dots, Z_M^{(2)}(T)$ . Thus, we see that,

$$\frac{\partial DIV}{\partial h_i(T)} = \sum_{j=1}^M \frac{\partial DIV}{\partial Z_j^{(2)}(T)} \frac{\partial Z_j^{(2)}(T)}{\partial h_i(T)} = \sum_{j=1}^M w_{ij}^{(2)} \frac{\partial DIV}{\partial Z_j^{(2)}(T)}$$

- ③ Computing  $\nabla_{\mathbf{h}(T)} DIV$  now, we get,

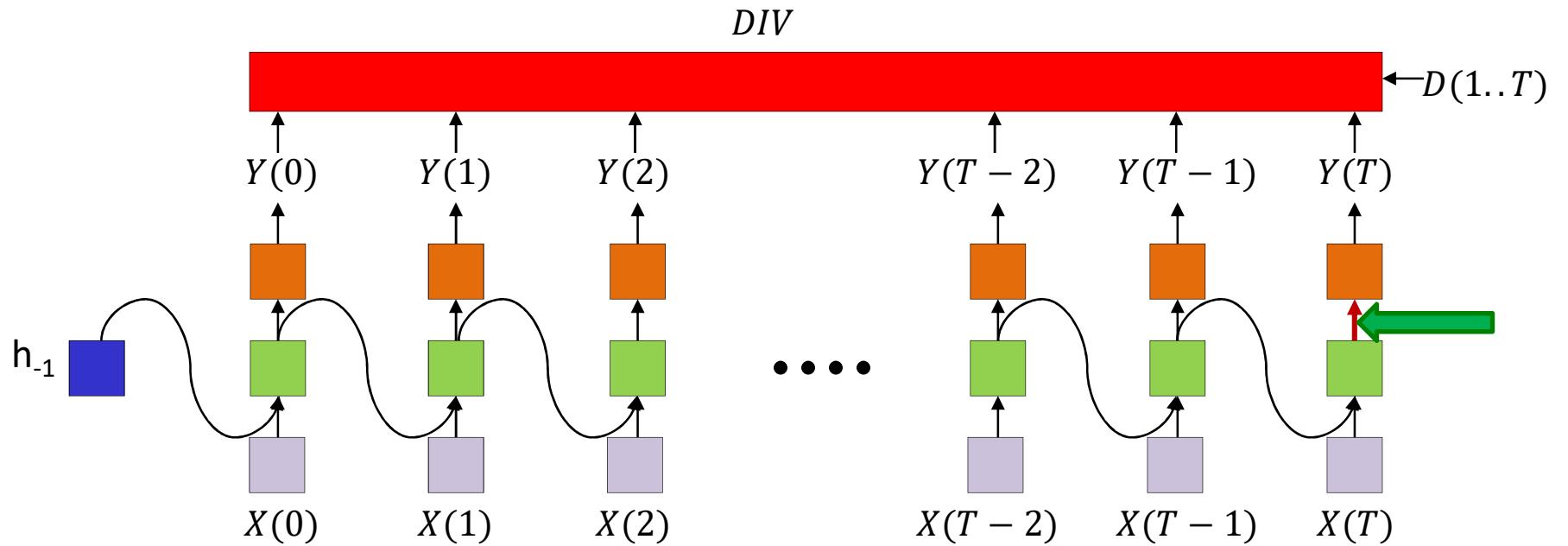
$$\begin{aligned} \nabla_{\mathbf{h}(T)} DIV &= \left[ \frac{\partial DIV}{\partial h_1(T)} \quad \cdots \cdots \quad \frac{\partial DIV}{\partial h_{N_1}(T)} \right] \\ &= \left[ \sum_{j=1}^M \frac{\partial DIV}{\partial Z_j^{(2)}(T)} \frac{\partial Z_j^{(2)}(T)}{\partial h_1(T)} \quad \cdots \cdots \quad \sum_{j=1}^M \frac{\partial DIV}{\partial Z_j^{(2)}(T)} \frac{\partial Z_j^{(2)}(T)}{\partial h_{N_1}(T)} \right] \\ &= \left[ \frac{\partial DIV}{\partial Z_1^{(2)}(T)} \quad \cdots \quad \frac{\partial DIV}{\partial Z_M^{(2)}(T)} \right] \begin{bmatrix} \frac{\partial Z_1^{(2)}(T)}{\partial h_1(T)} & \frac{\partial Z_1^{(2)}(T)}{\partial h_2(T)} & \cdots & \frac{\partial Z_1^{(2)}(T)}{\partial h_{N_1}(T)} \\ \frac{\partial Z_2^{(2)}(T)}{\partial h_1(T)} & \frac{\partial Z_2^{(2)}(T)}{\partial h_2(T)} & \cdots & \frac{\partial Z_2^{(2)}(T)}{\partial h_{N_1}(T)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial Z_M^{(2)}(T)}{\partial h_1(T)} & \frac{\partial Z_M^{(2)}(T)}{\partial h_2(T)} & \cdots & \frac{\partial Z_M^{(2)}(T)}{\partial h_{N_1}(T)} \end{bmatrix} \end{aligned}$$

Compute  $\nabla_{\mathbf{h}(T)} DIV = \left[ \frac{\partial DIV}{\partial h_1(T)} \dots \frac{\partial DIV}{\partial h_{N_1}(T)} \right]$

$$\nabla_{\mathbf{h}(T)} DIV = \nabla_{Z^{(2)}(T)} DIV \quad \nabla_{\mathbf{h}(T)} Z^{(2)}(T)$$

$$= \left[ \frac{\partial DIV}{\partial Z_1^{(2)}(T)} \dots \frac{\partial DIV}{\partial Z_M^{(2)}(T)} \right] \begin{bmatrix} w_{11}^{(2)} & w_{21}^{(2)} & \dots & w_{N_1 1}^{(2)} \\ w_{12}^{(2)} & w_{22}^{(2)} & \dots & w_{N_1 2}^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{1M}^{(2)} & w_{2M}^{(2)} & \dots & w_{N_1 M}^{(2)} \end{bmatrix}$$
$$= \nabla_{Z^{(2)}(T)} DIV \quad W^{(2)}$$

# Back Propagation Through Time



$$\frac{dDIV}{dZ_i^{(2)}(T)} = \frac{dDiv(T)}{dY_i(T)} \frac{dY_i(T)}{dZ_i^{(2)}(T)}$$

$$\frac{dDIV}{dh_i(T)} = \sum_j w_{ij}^{(2)} \frac{dDIV}{dZ_j^{(2)}(T)}$$

$$\nabla_{W^{(2)}} DIV = h(T) \nabla_{Z^{(2)}(T)} DIV$$

$$\frac{dDIV}{dW_{ij}^{(2)}} = \frac{dDIV}{dZ_j^{(2)}(T)} h_i(T)$$

# Compute $\nabla_{W^{(2)}} DIV$

- ➊ Note that  $W^{(2)}$  is an  $M \times N_1$  weight matrix connecting hidden layer to the output layer.
- ➋  $w_{ij}^{(2)}$  is connected to pre-activation value of  $j^{th}$  output node at  $T$  as  
 $Z_j^{(2)}(T) = \sum_{k=1}^{N_1} w_{kj}^{(2)} h_k(T) + b_j^{(2)}$ .
- ➌ Thus,

$$\frac{\partial DIV}{\partial w_{ij}^{(2)}} = \frac{\partial DIV}{\partial Z_j^{(2)}(T)} \frac{\partial Z_j^{(2)}(T)}{\partial w_{ij}^{(2)}} = \frac{\partial DIV}{\partial Z_j^{(2)}(T)} h_i(T)$$

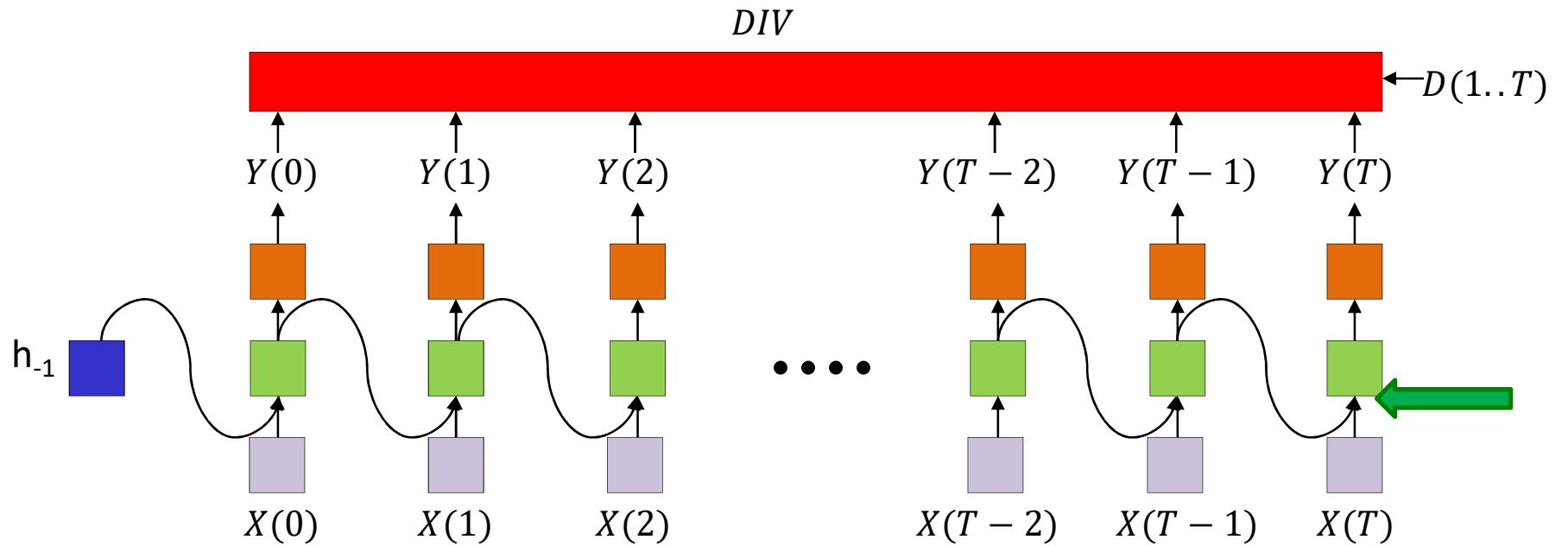
- ➍ Computing  $\nabla_{W^{(2)}} DIV$

$$\nabla_{W^{(2)}} DIV = \begin{bmatrix} \frac{\partial DIV}{\partial w_{11}^{(2)}} & \frac{\partial DIV}{\partial w_{12}^{(2)}} & \cdots & \frac{\partial DIV}{\partial w_{1M}^{(2)}} \\ \frac{\partial DIV}{\partial w_{21}^{(2)}} & \frac{\partial DIV}{\partial w_{22}^{(2)}} & \cdots & \frac{\partial DIV}{\partial w_{2M}^{(2)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial DIV}{\partial w_{N_1 1}^{(2)}} & \frac{\partial DIV}{\partial w_{N_1 2}^{(2)}} & \cdots & \frac{\partial DIV}{\partial w_{N_1 M}^{(2)}} \end{bmatrix}$$

# Compute $\nabla_{W^{(2)}} DIV$

$$\begin{aligned}
 \nabla_{W^{(2)}} DIV &= \begin{bmatrix} \frac{\partial DIV}{\partial Z_1^{(2)}(T)} h_1(T) & \frac{\partial DIV}{\partial Z_2^{(2)}(T)} h_1(T) & \cdots & \frac{\partial DIV}{\partial Z_M^{(2)}(T)} h_1(T) \\ \frac{\partial DIV}{\partial Z_1^{(2)}(T)} h_2(T) & \frac{\partial DIV}{\partial Z_2^{(2)}(T)} h_2(T) & \cdots & \frac{\partial DIV}{\partial Z_M^{(2)}(T)} h_2(T) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial DIV}{\partial Z_1^{(2)}(T)} h_{N_1}(T) & \frac{\partial DIV}{\partial Z_2^{(2)}(T)} h_{N_1}(T) & \cdots & \frac{\partial DIV}{\partial Z_M^{(2)}(T)} h_{N_1}(T) \end{bmatrix} \\
 &= \begin{bmatrix} h_1(T) \\ h_2(T) \\ \vdots \\ h_{N_1}(T) \end{bmatrix} \begin{bmatrix} \frac{\partial DIV}{\partial Z_1^{(2)}(T)} & \frac{\partial DIV}{\partial Z_2^{(2)}(T)} & \cdots & \frac{\partial DIV}{\partial Z_M^{(2)}(T)} \end{bmatrix} \\
 &= \mathbf{h}(T) \nabla_{Z^{(2)}} DIV
 \end{aligned}$$

# Back Propagation Through Time



$$\nabla_{Z^{(1)}(T)} DIV = \nabla_{h(T)} DIV \nabla_{Z^{(1)}(T)} h(T)$$

$$\frac{dDIV}{dZ_i^{(1)}(T)} = \frac{dDIV}{dh_i(T)} \frac{dh_i(T)}{dZ_i^{(1)}(T)}$$

$$\frac{dDIV}{dZ_i^{(2)}(T)} = \frac{dDIV}{dY_i(T)} \frac{dY_i(T)}{dZ_i^{(2)}(T)}$$

$$\frac{dDIV}{dh_i(T)} = \sum_j w_{ij}^{(2)} \frac{dDIV}{dZ_j^{(2)}(T)}$$

$$\frac{dDIV}{dw_{ij}^{(2)}} = \frac{dDIV}{dZ_j^{(2)}(T)} h_i(T)$$

# Compute $\nabla_{Z^{(1)}(T)} DIV$

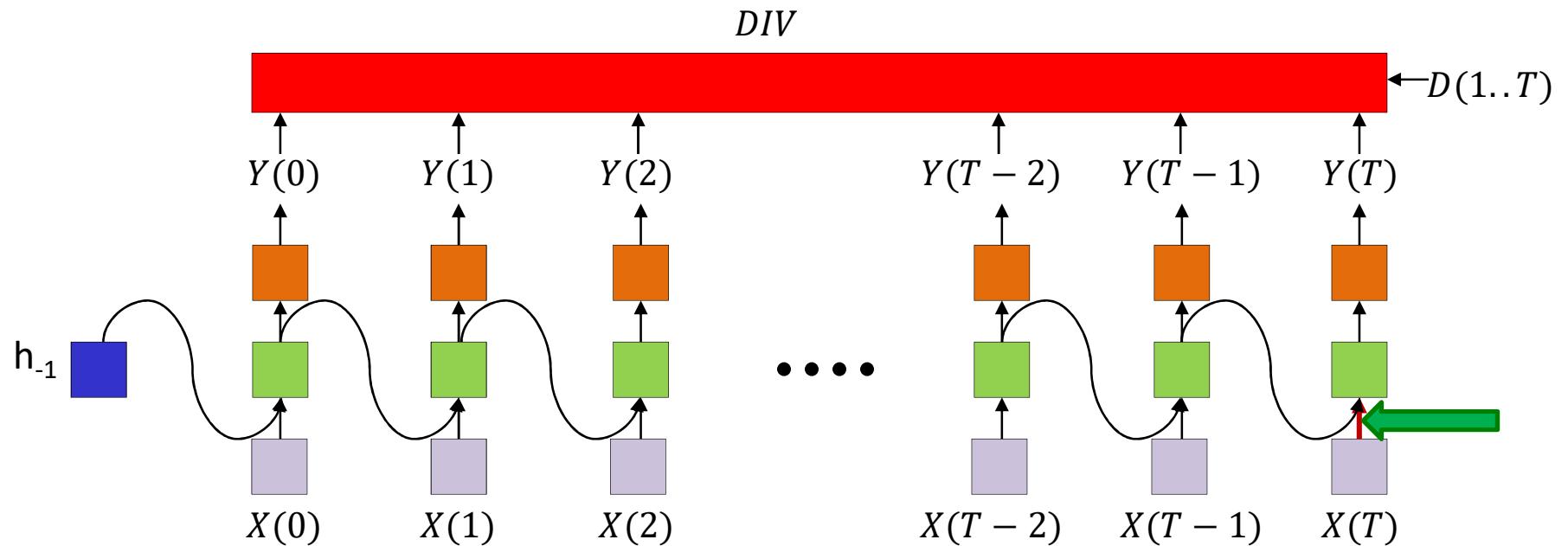
- ➊  $Z_i^{(1)}(T)$  is the pre-activation value of  $i^{th}$  hidden node at time  $T$ .
- ➋ Note that  $h_i(T) = Z_i^{(1)}(T)$ .
- ➌ Thus,  $\frac{\partial DIV}{\partial Z_i^{(1)}(T)} = \frac{\partial DIV}{\partial h_i(T)} \frac{\partial h_i(T)}{\partial Z_i^{(1)}(T)}$ .
- ➍ Computing  $\nabla_{Z^{(1)}(T)} DIV$

$$\begin{aligned}
 \nabla_{Z^{(1)}(T)} DIV &= \left[ \frac{\partial DIV}{\partial Z_1^{(1)}(T)} \quad \frac{\partial DIV}{\partial Z_2^{(1)}(T)} \quad \cdots \quad \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(T)} \right] \\
 &= \left[ \frac{\partial DIV}{\partial h_1(T)} \frac{\partial h_1(T)}{\partial Z_1^{(1)}(T)} \quad \frac{\partial DIV}{\partial h_2(T)} \frac{\partial h_2(T)}{\partial Z_2^{(1)}(T)} \quad \cdots \quad \frac{\partial DIV}{\partial h_{N_1}(T)} \frac{\partial h_{N_1}(T)}{\partial Z_{N_1}^{(1)}(T)} \right] \\
 &= \left[ \frac{\partial DIV}{\partial h_1(T)} \quad \frac{\partial DIV}{\partial h_2(T)} \quad \cdots \quad \frac{\partial DIV}{\partial h_{N_1}(T)} \right] \begin{bmatrix} \frac{\partial h_1(T)}{\partial Z_1^{(1)}(T)} & \frac{\partial h_1(T)}{\partial Z_2^{(1)}(T)} & \cdots & \frac{\partial h_1(T)}{\partial Z_{N_1}^{(1)}(T)} \\ \frac{\partial h_2(T)}{\partial Z_1^{(1)}(T)} & \frac{\partial h_2(T)}{\partial Z_2^{(1)}(T)} & \cdots & \frac{\partial h_2(T)}{\partial Z_{N_1}^{(1)}(T)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_{N_1}(T)}{\partial Z_1^{(1)}(T)} & \frac{\partial h_{N_1}(T)}{\partial Z_2^{(1)}(T)} & \cdots & \frac{\partial h_{N_1}(T)}{\partial Z_{N_1}^{(1)}(T)} \end{bmatrix}
 \end{aligned}$$

Compute  $\nabla_{Z^{(1)}(T)} DIV$ 

$$\nabla_{Z^{(1)}(T)} DIV = \left[ \frac{\partial DIV}{\partial h_1(T)} \quad \frac{\partial DIV}{\partial h_2(T)} \quad \cdots \quad \frac{\partial DIV}{\partial h_{N_1}(T)} \right] \begin{bmatrix} \frac{\partial h_1(T)}{\partial Z_1^{(1)}(T)} & 0 & \cdots & 0 \\ 0 & \frac{\partial h_2(T)}{\partial Z_2^{(1)}(T)} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{\partial h_{N_1}(T)}{\partial Z_{N_1}^{(1)}(T)} \end{bmatrix}$$
$$= \nabla_{\mathbf{h}(T)} DIV \cdot \nabla_{Z^{(1)}(T)} \mathbf{h}(T)$$

# Back Propagation Through Time



$$\nabla_{W^{(1)}} DIV = X(T) \nabla_{Z^{(1)}(T)} DIV$$

$$\frac{dDIV}{dw_{ij}^{(1)}} = \frac{dDIV}{dZ_j^{(1)}(T)} X_i(T)$$

# Compute $\nabla_{W^{(1)}} DIV$

- ① Note that  $W^{(1)}$  is  $N_1 \times N$  weight matrix between the input layer and the hidden layer.

- ②  $w_{ij}^{(1)}$  is connected to  $Z_j^{(1)}(T)$  as follows.

$$Z_j^{(1)}(T) = \sum_{k=1}^N w_{kj}^{(1)} X_k(T) + \sum_{k=1}^{N_1} w_{kj}^{(1)} h_k(T-1) + b_j^{(1)}$$

- ③ Thus,  $\frac{\partial DIV}{\partial w_{ij}^{(1)}} = \frac{\partial DIV}{\partial Z_j^{(1)}(T)} \frac{\partial Z_j^{(1)}(T)}{\partial w_{ij}^{(1)}} = \frac{\partial DIV}{\partial Z_j^{(1)}(T)} X_i(T)$

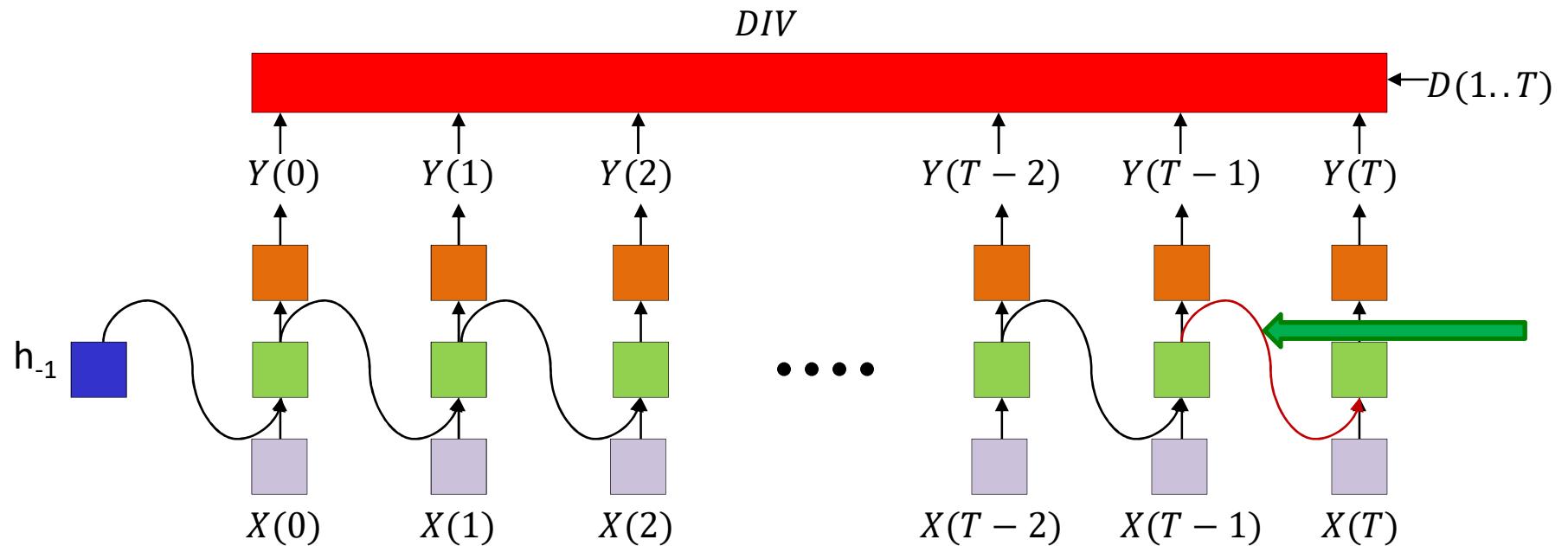
- ④ Computing  $\nabla_{W^{(1)}} DIV$

$$\nabla_{W^{(1)}} DIV = \begin{bmatrix} \frac{\partial DIV}{\partial w_{11}^{(1)}} & \frac{\partial DIV}{\partial w_{12}^{(1)}} & \cdots & \frac{\partial DIV}{\partial w_{1N_1}^{(1)}} \\ \frac{\partial DIV}{\partial w_{21}^{(1)}} & \frac{\partial DIV}{\partial w_{22}^{(1)}} & \cdots & \frac{\partial DIV}{\partial w_{2N_1}^{(1)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial DIV}{\partial w_{N1}^{(1)}} & \frac{\partial DIV}{\partial w_{N2}^{(1)}} & \cdots & \frac{\partial DIV}{\partial w_{NN_1}^{(1)}} \end{bmatrix}$$

# Compute $\nabla_{W^{(1)}} DIV$

$$\begin{aligned}
 \nabla_{W^{(1)}} DIV &= \begin{bmatrix} \frac{\partial DIV}{\partial Z_1^{(1)}(T)} X_1(T) & \frac{\partial DIV}{\partial Z_2^{(1)}(T)} X_1(T) & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(T)} X_1(T) \\ \frac{\partial DIV}{\partial Z_1^{(1)}(T)} X_2(T) & \frac{\partial DIV}{\partial Z_2^{(1)}(T)} X_2(T) & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(T)} X_2(T) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial DIV}{\partial Z_1^{(1)}(T)} X_N(T) & \frac{\partial DIV}{\partial Z_2^{(1)}(T)} X_N(T) & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(T)} X_N(T) \end{bmatrix} \\
 &= \begin{bmatrix} X_1(T) \\ X_2(T) \\ \vdots \\ X_N(T) \end{bmatrix} \begin{bmatrix} \frac{\partial DIV}{\partial Z_1^{(1)}(T)} & \frac{\partial DIV}{\partial Z_2^{(1)}(T)} & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(T)} \end{bmatrix} \\
 &= X(T) \nabla_{Z^{(1)}(T)} DIV
 \end{aligned}$$

# Back Propagation Through Time



$$\nabla_{W^{(11)}} DIV = h(T-1) \nabla_{Z_j^{(1)}(T)} DIV$$

$$\frac{dDIV}{dw_{ij}^{(1)}} = \frac{dDIV}{dZ_j^{(1)}(T)} X_i(T)$$

$$\frac{dDIV}{dw_{ij}^{(11)}} = \frac{dDIV}{dZ_j^{(1)}(T)} h_i(T-1)$$

# Compute $\nabla_{W^{(11)}} DIV$

① Note that  $W^{(11)}$  is  $N_1 \times N_1$  weight matrix between  $\mathbf{h}(T-1)$  and  $\mathbf{h}(T)$ .

②  $w_{ij}^{(11)}$  is connected to  $Z_j^{(1)}(T)$  as follows.

$$Z_j^{(1)}(T) = \sum_{k=1}^N w_{kj}^{(1)} X_k(T) + \sum_{k=1}^{N_1} w_{kj}^{(11)} h_k(T-1) + b_j^{(1)}$$

③ Thus,  $\frac{\partial DIV}{\partial w_{ij}^{(11)}} = \frac{\partial DIV}{\partial Z_j^{(1)}(T)} \frac{\partial Z_j^{(1)}(T)}{\partial w_{ij}^{(11)}} = \frac{\partial DIV}{\partial Z_j^{(1)}(T)} h_i(T-1)$ .

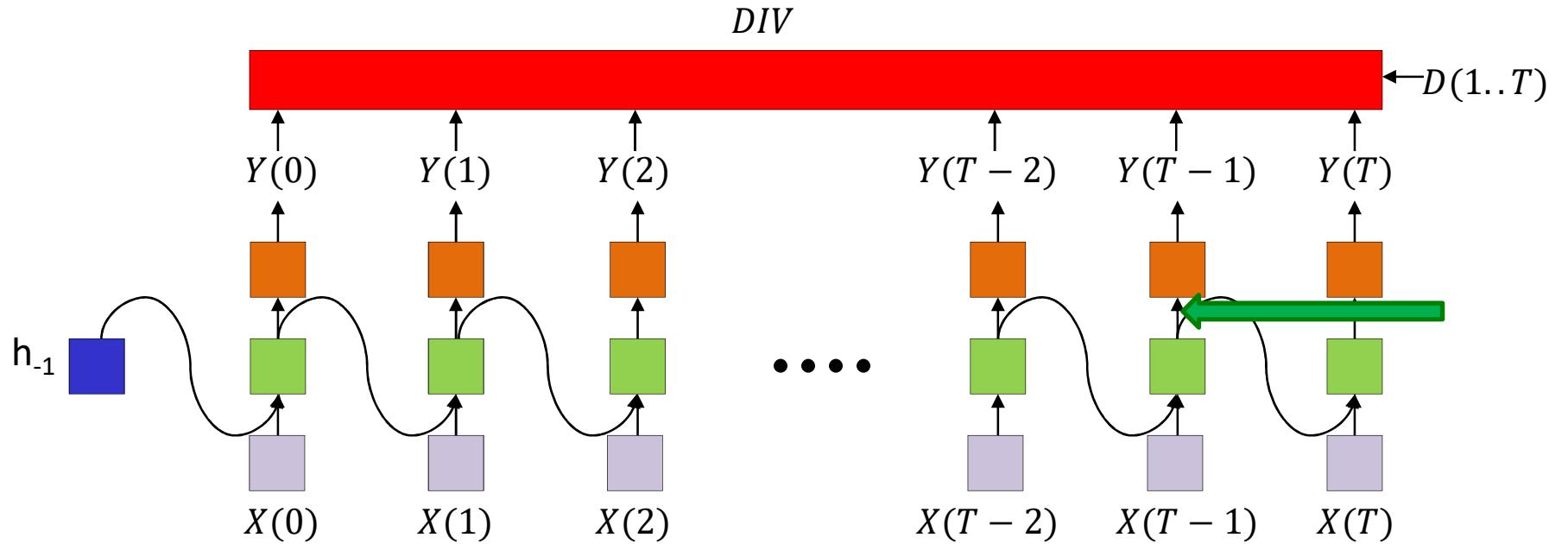
④ Computing  $\nabla_{W^{(11)}} DIV$

$$\nabla_{W^{(11)}} DIV = \begin{bmatrix} \frac{\partial DIV}{\partial w_{11}^{(11)}} & \frac{\partial DIV}{\partial w_{12}^{(11)}} & \cdots & \frac{\partial DIV}{\partial w_{1N_1}^{(11)}} \\ \frac{\partial DIV}{\partial w_{21}^{(11)}} & \frac{\partial DIV}{\partial w_{22}^{(11)}} & \cdots & \frac{\partial DIV}{\partial w_{2N_1}^{(11)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial DIV}{\partial w_{N_1 1}^{(11)}} & \frac{\partial DIV}{\partial w_{N_1 2}^{(11)}} & \cdots & \frac{\partial DIV}{\partial w_{N_1 N_1}^{(11)}} \end{bmatrix}$$

# Compute $\nabla_{W^{(11)}} DIV$

$$\begin{aligned}
 \nabla_{W^{(11)}} DIV &= \begin{bmatrix} \frac{\partial DIV}{\partial Z_1^{(1)}(\tau)} h_1(\tau-1) & \frac{\partial DIV}{\partial Z_2^{(1)}(\tau)} h_1(\tau-1) & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(\tau)} h_1(\tau-1) \\ \frac{\partial DIV}{\partial Z_1^{(1)}(\tau)} h_2(\tau-1) & \frac{\partial DIV}{\partial Z_2^{(1)}(\tau)} h_2(\tau-1) & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(\tau)} h_2(\tau-1) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial DIV}{\partial Z_1^{(1)}(\tau)} h_{N_1}(\tau-1) & \frac{\partial DIV}{\partial Z_2^{(1)}(\tau)} h_{N_1}(\tau-1) & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(\tau)} h_{N_1}(\tau-1) \end{bmatrix} \\
 &= \begin{bmatrix} h_1(\tau-1) \\ h_2(\tau-1) \\ \vdots \\ h_{N_1}(\tau-1) \end{bmatrix} \begin{bmatrix} \frac{\partial DIV}{\partial Z_1^{(1)}(\tau)} & \frac{\partial DIV}{\partial Z_2^{(1)}(\tau)} & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(\tau)} \end{bmatrix} \\
 &= \mathbf{h}(\tau-1) \nabla_{Z^{(1)}(\tau)} DIV
 \end{aligned}$$

# Back Propagation Through Time



$$\nabla_{Z^{(2)}(T-1)} DIV = \nabla_{Y(T-1)} DIV \nabla_{Z^{(2)}(T)} Y(T-1)$$

Vector output activation

$$\frac{dDIV}{dZ_i^{(2)}(T-1)} = \frac{dDIV}{dY_i(T-1)} \frac{dY_i(T-1)}{dZ_i^{(2)}(T-1)}$$

OR

$$\frac{dDIV}{dZ_i^{(2)}(T-1)} = \sum_j \frac{dDIV}{dY_j(T-1)} \frac{dY_j(T-1)}{dZ_i^{(2)}(T-1)}$$

Compute  $\nabla_{Z^{(2)}(T-1)} DIV$

### Case 1: Vector output activation

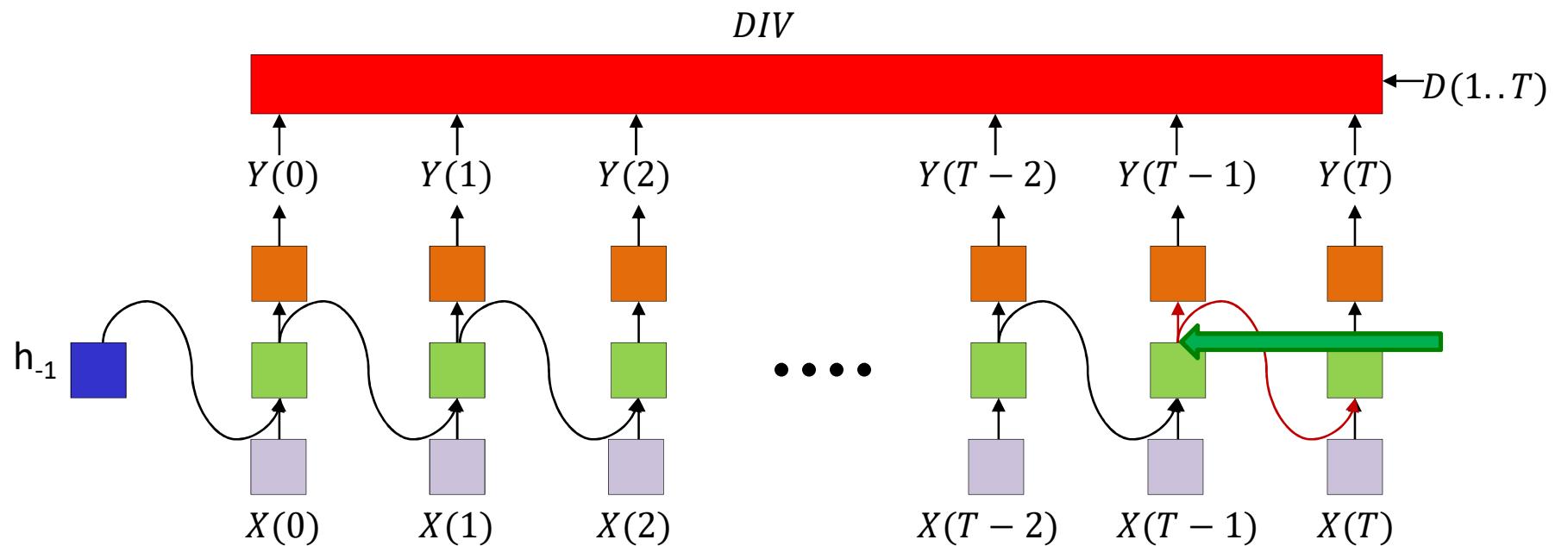
$$\begin{aligned}
 \nabla_{Z^{(2)}(T-1)} DIV &= \left[ \frac{\partial DIV}{\partial Z_1^{(2)}(T-1)} \cdots \frac{\partial DIV}{\partial Z_M^{(2)}(T-1)} \right] \\
 &= \left[ \sum_{j=1}^M \frac{\partial DIV}{\partial Y_j(T-1)} \frac{\partial Y_j(T-1)}{\partial Z_1^{(2)}(T-1)} \cdots \cdots \sum_{j=1}^M \frac{\partial DIV}{\partial Y_j(T-1)} \frac{\partial Y_j(T-1)}{\partial Z_M^{(2)}(T-1)} \right] \\
 &= \left[ \frac{\partial DIV}{\partial Y_1(T-1)} \cdots \frac{\partial DIV}{\partial Y_M(T-1)} \right] \begin{bmatrix} \frac{\partial Y_1(T-1)}{\partial Z_1^{(2)}(T-1)} & \frac{\partial Y_1(T-1)}{\partial Z_2^{(2)}(T-1)} & \cdots & \frac{\partial Y_1(T-1)}{\partial Z_M^{(2)}(T-1)} \\ \frac{\partial Y_2(T-1)}{\partial Z_1^{(2)}(T-1)} & \frac{\partial Y_2(T-1)}{\partial Z_2^{(2)}(T-1)} & \cdots & \frac{\partial Y_2(T-1)}{\partial Z_M^{(2)}(T-1)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial Y_M(T-1)}{\partial Z_1^{(2)}(T-1)} & \frac{\partial Y_M(T-1)}{\partial Z_2^{(2)}(T-1)} & \cdots & \frac{\partial Y_M(T-1)}{\partial Z_M^{(2)}(T-1)} \end{bmatrix} \\
 &= \nabla_{Y(T-1)} DIV \nabla_{Z^{(2)}(T-1)} Y(T-1)
 \end{aligned}$$

Compute  $\nabla_{Z^{(2)}(T-1)} DIV$

## Case 2: Scalar output activation

$$\begin{aligned} \nabla_{Z^{(2)}(T-1)} DIV &= \left[ \frac{\partial DIV}{\partial Y_1(T-1)} \frac{\partial Y_1(T-1)}{\partial Z_1^{(2)}(T-1)} \dots \dots \frac{\partial DIV}{\partial Y_M(T-1)} \frac{\partial Y_M(T-1)}{\partial Z_M^{(2)}(T-1)} \right] \\ &= \left[ \frac{\partial DIV}{\partial Y_1(T-1)} \dots \dots \frac{\partial DIV}{\partial Y_M(T-1)} \right] \begin{bmatrix} \frac{\partial Y_1(T-1)}{\partial Z_1^{(2)}(T-1)} & 0 & \dots & 0 \\ 0 & \frac{\partial Y_2(T-1)}{\partial Z_2^{(2)}(T-1)} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{\partial Y_M(T-1)}{\partial Z_M^{(2)}(T-1)} \end{bmatrix} \\ &= \nabla_{Y(T-1)} DIV \quad \nabla_{Z^{(2)}(T-1)} Y(T-1) \end{aligned}$$

# Back Propagation Through Time



$$\frac{dDIV}{dh_i(T-1)} = \sum_j w_{ij}^{(2)} \frac{dDIV}{dZ_j^{(2)}(T-1)} + \sum_j w_{ij}^{(11)} \frac{dDIV}{dZ_j^{(1)}(T)}$$

$$\nabla_{h(T-1)} DIV = \nabla_{Z^{(2)}(T-1)} DIV W^{(2)} + \nabla_{Z^{(1)}(T)} DIV W^{(11)}$$

# Compute $\nabla_{\mathbf{h}(T-1)} DIV$

- ①  $h_i(T-1)$  is an input to all the  $Z_1^{(2)}(T-1), \dots, Z_M^{(2)}(T-1)$  where  
 $Z_j^{(2)}(T-1) = \sum_{k=1}^{N_1} w_{kj}^{(2)} h_k(T-1) + b_j^{(2)}$ .

- ②  $h_i(T-1)$  is an input to all the  $Z_1^{(1)}(T), \dots, Z_{N_1}^{(1)}(T)$  where  
 $Z_j^{(1)}(T) = \sum_{k=1}^N w_{kj}^{(1)} X_k(T) + \sum_{k=1}^{N_1} w_{kj}^{(11)} h_k(T-1) + b_j^{(1)}$ .

- ③ Thus, we see that,

$$\begin{aligned}\frac{\partial DIV}{\partial h_i(T-1)} &= \sum_{j=1}^M \frac{\partial DIV}{\partial Z_j^{(2)}(T-1)} \frac{\partial Z_j^{(2)}(T-1)}{\partial h_i(T-1)} + \sum_{j=1}^{N_1} \frac{\partial DIV}{\partial Z_j^{(1)}(T)} \frac{\partial Z_j^{(1)}(T)}{\partial h_i(T-1)} \\ &= \sum_{j=1}^M w_{ij}^{(2)} \frac{\partial DIV}{\partial Z_j^{(2)}(T-1)} + \sum_{j=1}^{N_1} w_{ij}^{(11)} \frac{\partial DIV}{\partial Z_j^{(1)}(T)}\end{aligned}$$

- ④ Computing  $\nabla_{\mathbf{h}(T-1)} DIV = \left[ \frac{\partial DIV}{\partial h_1(T-1)} \quad \dots \quad \frac{\partial DIV}{\partial h_{N_1}(T-1)} \right]$

$$\begin{aligned}\nabla_{\mathbf{h}(T-1)} DIV &= \left[ \sum_{j=1}^M \frac{\partial DIV}{\partial Z_j^{(2)}(T-1)} \frac{\partial Z_j^{(2)}(T-1)}{\partial h_1(T-1)} \quad \dots \quad \sum_{j=1}^M \frac{\partial DIV}{\partial Z_j^{(2)}(T-1)} \frac{\partial Z_j^{(2)}(T-1)}{\partial h_{N_1}(T-1)} \right] \\ &\quad + \left[ \sum_{j=1}^{N_1} \frac{\partial DIV}{\partial Z_j^{(1)}(T)} \frac{\partial Z_j^{(1)}(T)}{\partial h_1(T-1)} \quad \dots \quad \sum_{j=1}^{N_1} \frac{\partial DIV}{\partial Z_j^{(1)}(T)} \frac{\partial Z_j^{(1)}(T)}{\partial h_{N_1}(T-1)} \right]\end{aligned}$$

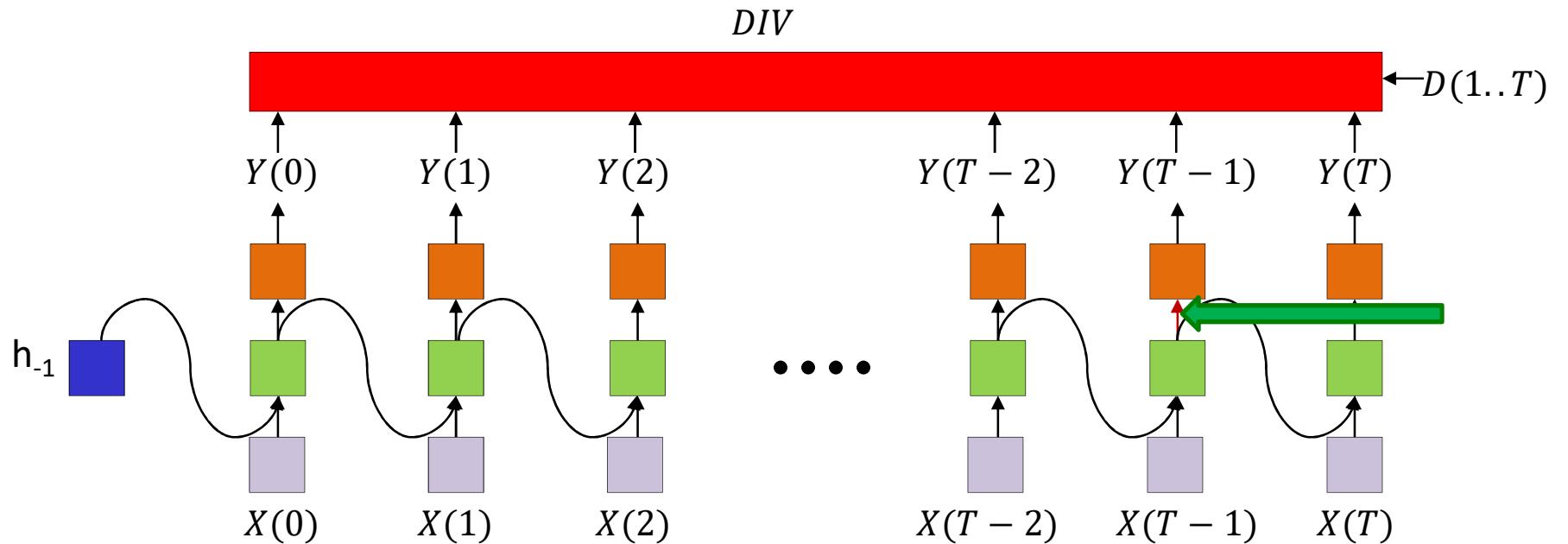
# Compute $\nabla_{\mathbf{h}(T-1)} DIV$

$$\begin{aligned}
 &= \left[ \frac{\partial DIV}{\partial Z_1^{(2)}(T-1)} \quad \dots \quad \frac{\partial DIV}{\partial Z_M^{(2)}(T-1)} \right] \begin{bmatrix} \frac{\partial Z_1^{(2)}(T-1)}{\partial h_1(T-1)} & \frac{\partial Z_1^{(2)}(T-1)}{\partial h_2(T-1)} & \dots & \frac{\partial Z_1^{(2)}(T-1)}{\partial h_{N_1}(T-1)} \\ \frac{\partial Z_2^{(2)}(T-1)}{\partial h_1(T-1)} & \frac{\partial Z_2^{(2)}(T-1)}{\partial h_2(T-1)} & \dots & \frac{\partial Z_2^{(2)}(T-1)}{\partial h_{N_1}(T-1)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial Z_M^{(2)}(T-1)}{\partial h_1(T-1)} & \frac{\partial Z_M^{(2)}(T-1)}{\partial h_2(T-1)} & \dots & \frac{\partial Z_M^{(2)}(T-1)}{\partial h_{N_1}(T-1)} \end{bmatrix} \\
 &\quad + \left[ \frac{\partial DIV}{\partial Z_1^{(1)}(T)} \quad \dots \quad \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(T)} \right] \begin{bmatrix} \frac{\partial Z_1^{(1)}(T)}{\partial h_1(T-1)} & \frac{\partial Z_1^{(1)}(T)}{\partial h_2(T-1)} & \dots & \frac{\partial Z_1^{(1)}(T)}{\partial h_{N_1}(T-1)} \\ \frac{\partial Z_2^{(1)}(T)}{\partial h_1(T-1)} & \frac{\partial Z_2^{(1)}(T)}{\partial h_2(T-1)} & \dots & \frac{\partial Z_2^{(1)}(T)}{\partial h_{N_1}(T-1)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial Z_{N_1}^{(1)}(T)}{\partial h_1(T-1)} & \frac{\partial Z_{N_1}^{(1)}(T)}{\partial h_2(T-1)} & \dots & \frac{\partial Z_{N_1}^{(1)}(T)}{\partial h_{N_1}(T-1)} \end{bmatrix}
 \end{aligned}$$

# Compute $\nabla_{\mathbf{h}(T-1)} DIV$

$$\begin{aligned} &= \nabla_{Z^{(2)}(T-1)} DIV \begin{bmatrix} w_{11}^{(2)} & w_{21}^{(2)} & \cdots & w_{N_1 1}^{(2)} \\ w_{12}^{(2)} & w_{22}^{(2)} & \cdots & w_{N_1 2}^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{1M}^{(2)} & w_{2M}^{(2)} & \cdots & w_{N_1 M}^{(2)} \end{bmatrix} \\ &\quad + \nabla_{Z^{(1)}(T)} DIV \begin{bmatrix} w_{11}^{(11)} & w_{21}^{(11)} & \cdots & w_{N_1 1}^{(11)} \\ w_{12}^{(11)} & w_{22}^{(11)} & \cdots & w_{N_1 2}^{(11)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{1N_1}^{(11)} & w_{2N_1}^{(11)} & \cdots & w_{N_1 N_1}^{(11)} \end{bmatrix} \\ &= \nabla_{Z^{(2)}(T-1)} DIV W^{(2)} + \nabla_{Z^{(1)}(T)} DIV W^{(11)} \end{aligned}$$

# Back Propagation Through Time



$$\frac{dDIV}{dh_i(T-1)} = \sum_j w_{ij}^{(2)} \frac{dDIV}{dZ_j^{(2)}(T-1)} + \sum_j w_{ij}^{(11)} \frac{dDIV}{dZ_j^{(1)}(T)}$$

Note the addition  $\rightarrow$

$$\frac{dDIV}{dw_{ij}^{(2)}} += \frac{dDIV}{dZ_j^{(2)}(T-1)} h_i(T-1)$$

$$\nabla_{W^{(2)}} DIV += h(T-1) \nabla_{Z^{(2)}(T-1)} DIV$$

# Updating $\nabla_{W^{(2)}} DIV$

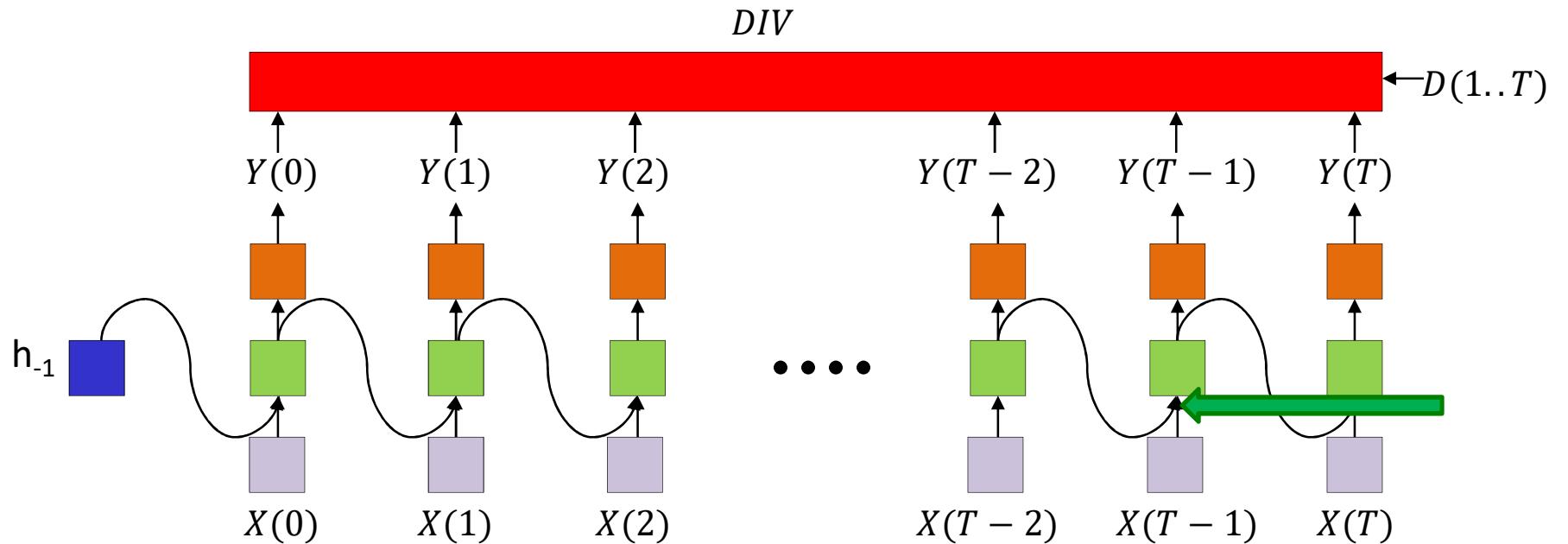
$$\begin{aligned}\frac{\partial DIV}{\partial w_{ij}^{(2)}|_{new}} &= \frac{\partial DIV}{\partial w_{ij}^{(2)}|_{old}} + \frac{\partial DIV}{\partial Z_j^{(2)}(T-1)} \frac{\partial Z_j^{(2)}(T-1)}{\partial w_{ij}^{(2)}} \\ &= \frac{\partial DIV}{\partial Z_j^{(2)}(T)} h_i(T) + \frac{\partial DIV}{\partial Z_j^{(2)}(T-1)} h_i(T-1)\end{aligned}$$

$$\nabla_{W^{(2)}} DIV|_{new} = \nabla_{W^{(2)}} DIV|_{old}$$

$$+ \begin{bmatrix} \frac{\partial DIV}{\partial Z_1^{(2)}(T-1)} h_1(T-1) & \frac{\partial DIV}{\partial Z_2^{(2)}(T-1)} h_1(T-1) & \cdots & \frac{\partial DIV}{\partial Z_M^{(2)}(T-1)} h_1(T-1) \\ \frac{\partial DIV}{\partial Z_1^{(2)}(T-1)} h_2(T-1) & \frac{\partial DIV}{\partial Z_2^{(2)}(T-1)} h_2(T-1) & \cdots & \frac{\partial DIV}{\partial Z_M^{(2)}(T-1)} h_2(T-1) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial DIV}{\partial Z_1^{(2)}(T-1)} h_{N_1}(T-1) & \frac{\partial DIV}{\partial Z_2^{(2)}(T-1)} h_{N_1}(T-1) & \cdots & \frac{\partial DIV}{\partial Z_M^{(2)}(T-1)} h_{N_1}(T-1) \end{bmatrix}$$

$$= \mathbf{h}(T) \nabla_{Z^{(2)}(T)} DIV + \mathbf{h}(T-1) \nabla_{Z^{(2)}(T-1)} DIV$$

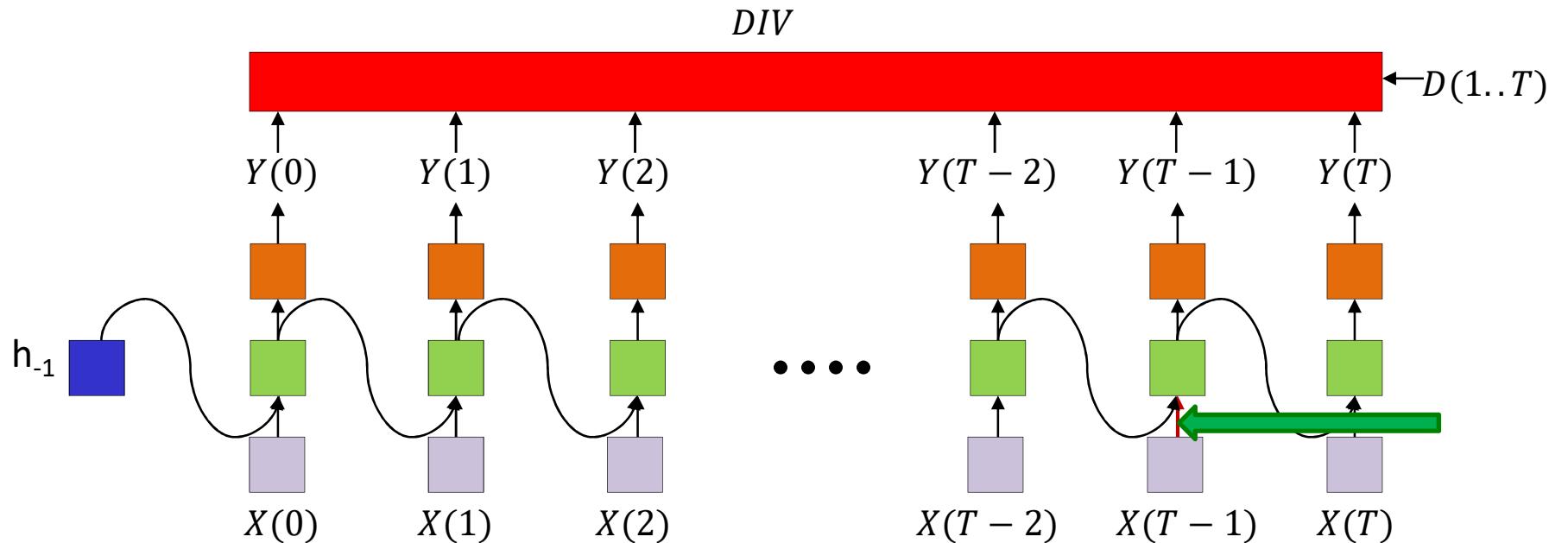
# Back Propagation Through Time



$$\frac{dDIV}{dZ_i^{(1)}(T-1)} = \frac{dDIV}{dh_i(T-1)} \frac{dh_i(T-1)}{dZ_i^{(1)}(T-1)}$$

$$\nabla_{Z^{(1)}(T-1)} DIV = \nabla_{h(T-1)} DIV \nabla_{Z^{(1)}(T-1)} h(T-1)$$

# Back Propagation Through Time



$$\frac{dDIV}{dZ_i^{(1)}(T-1)} = \frac{dDIV}{dh_i(T-1)} \frac{dh_i(T-1)}{dZ_i^{(1)}(T-1)}$$

$$\frac{dDIV}{dw_{ij}^{(1)}} += \frac{dDIV}{dZ_j^{(1)}(T-1)} X_i(T-1)$$

Note the addition

$$\nabla_{W^{(1)}} DIV += X(T-1) \nabla_{Z^{(1)}(T-1)} DIV$$

# Updating $\nabla_{W^{(1)}} DIV$

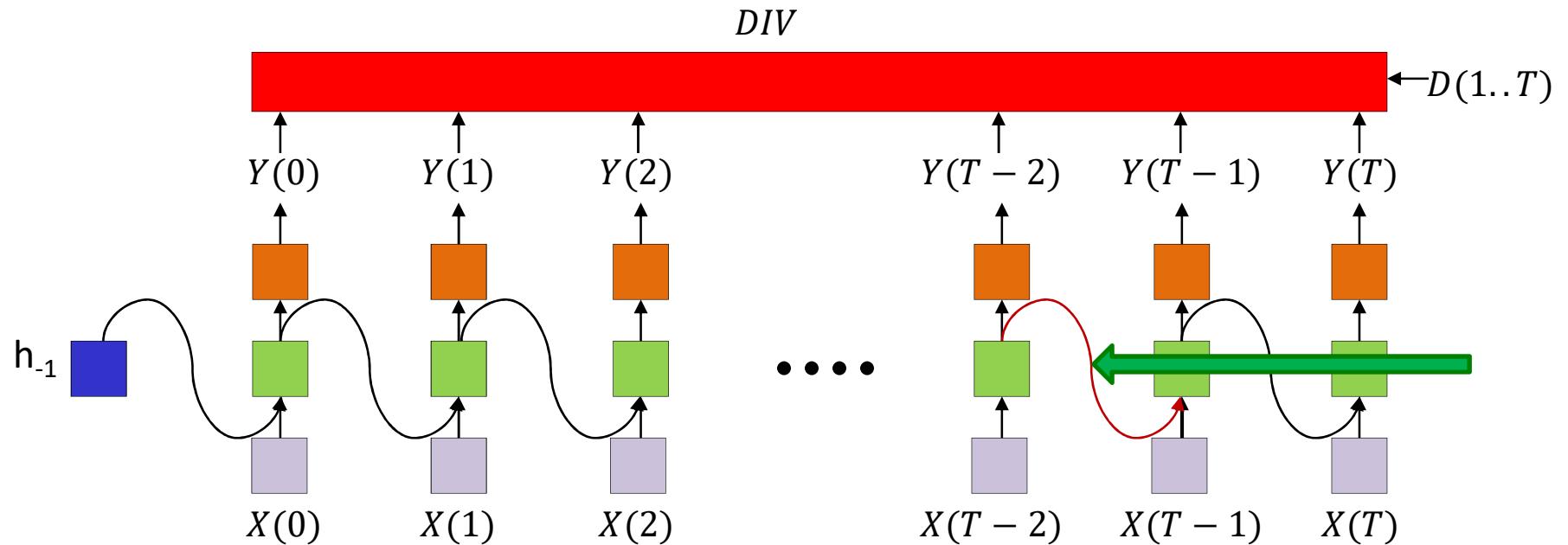
$$\begin{aligned}\frac{\partial DIV}{\partial w_{ij}^{(1)}|_{new}} &= \frac{\partial DIV}{\partial w_{ij}^{(1)}|_{old}} + \frac{\partial DIV}{\partial Z_j^{(1)}(T-1)} \frac{\partial Z_j^{(1)}(T-1)}{\partial w_{ij}^{(1)}} \\ &= \frac{\partial DIV}{\partial Z_j^{(1)}(T)} X_i(T) + \frac{\partial DIV}{\partial Z_j^{(1)}(T-1)} X_i(T-1)\end{aligned}$$

$$\nabla_{W^{(1)}} DIV|_{new} = \nabla_{W^{(1)}} DIV|_{old}$$

$$+ \begin{bmatrix} \frac{\partial DIV}{\partial Z_1^{(1)}(T-1)} X_1(T-1) & \frac{\partial DIV}{\partial Z_2^{(1)}(T-1)} X_1(T-1) & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(T-1)} X_1(T-1) \\ \frac{\partial DIV}{\partial Z_1^{(1)}(T-1)} X_2(T-1) & \frac{\partial DIV}{\partial Z_2^{(1)}(T-1)} X_2(T-1) & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(T-1)} X_2(T-1) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial DIV}{\partial Z_1^{(1)}(T-1)} X_N(T-1) & \frac{\partial DIV}{\partial Z_2^{(1)}(T-1)} X_N(T-1) & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(T-1)} X_N(T-1) \end{bmatrix}$$

$$= X(T) \nabla_{Z^{(1)}(T)} DIV + X(T-1) \nabla_{Z^{(1)}(T-1)} DIV$$

# Back Propagation Through Time



$$\frac{dDIV}{dw_{ij}^{(1)}} += \frac{dDIV}{dZ_j^{(1)}(T-1)} X_i(T-1)$$

Note the addition  $\rightarrow$

$$\frac{dDIV}{dw_{ij}^{(11)}} += \frac{dDIV}{dZ_j^{(1)}(T-1)} h_i(T-2)$$

$$\nabla_{W^{(11)}} DIV += h(T-2) \nabla_{Z^{(1)}(T-1)} DIV$$

# Updating $\nabla_{W^{(11)}} DIV$

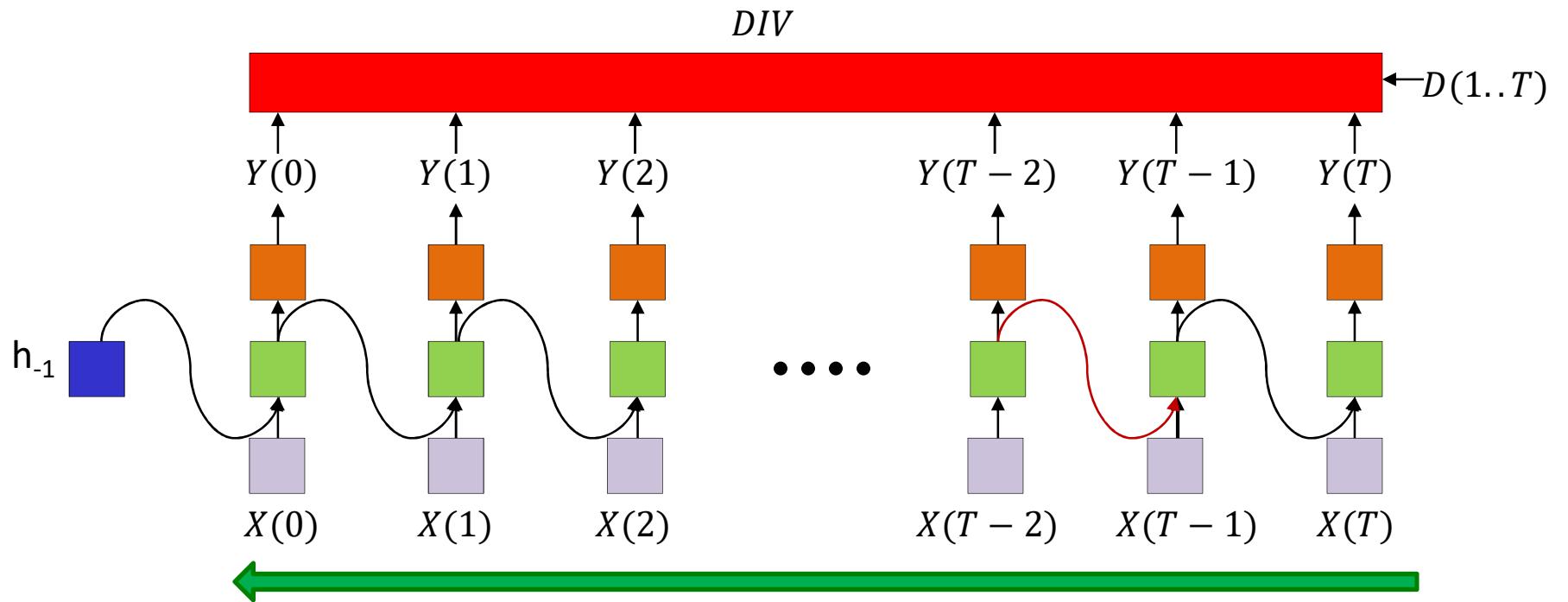
$$\begin{aligned}\frac{\partial DIV}{\partial w_{ij}^{(11)}}|_{new} &= \frac{\partial DIV}{\partial w_{ij}^{(11)}}|_{old} + \frac{\partial DIV}{\partial Z_j^{(1)}(T-1)} \frac{\partial Z_j^{(1)}(T-1)}{\partial w_{ij}^{(11)}} \\ &= \frac{\partial DIV}{\partial Z_j^{(1)}(T)} h_i(T-1) + \frac{\partial DIV}{\partial Z_j^{(1)}(T-1)} h_i(T-2)\end{aligned}$$

$$\nabla_{W^{(11)}} DIV|_{new} = \nabla_{W^{(11)}} DIV|_{old}$$

$$+ \begin{bmatrix} \frac{\partial DIV}{\partial Z_1^{(1)}(T-1)} h_1(T-2) & \frac{\partial DIV}{\partial Z_2^{(1)}(T-1)} h_1(T-2) & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(T-1)} h_1(T-2) \\ \frac{\partial DIV}{\partial Z_1^{(1)}(T-1)} h_2(T-2) & \frac{\partial DIV}{\partial Z_2^{(1)}(T-1)} h_2(T-2) & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(T-1)} h_2(T-2) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial DIV}{\partial Z_1^{(1)}(T-1)} h_{N_1}(T-2) & \frac{\partial DIV}{\partial Z_2^{(1)}(T-1)} h_{N_1}(T-2) & \cdots & \frac{\partial DIV}{\partial Z_{N_1}^{(1)}(T-1)} h_{N_1}(T-2) \end{bmatrix}$$

$$= \mathbf{h}(T-1) \nabla_{Z^{(1)}(T)} DIV + \mathbf{h}(T-2) \nabla_{Z^{(1)}(T-1)} DIV$$

# Back Propagation Through Time

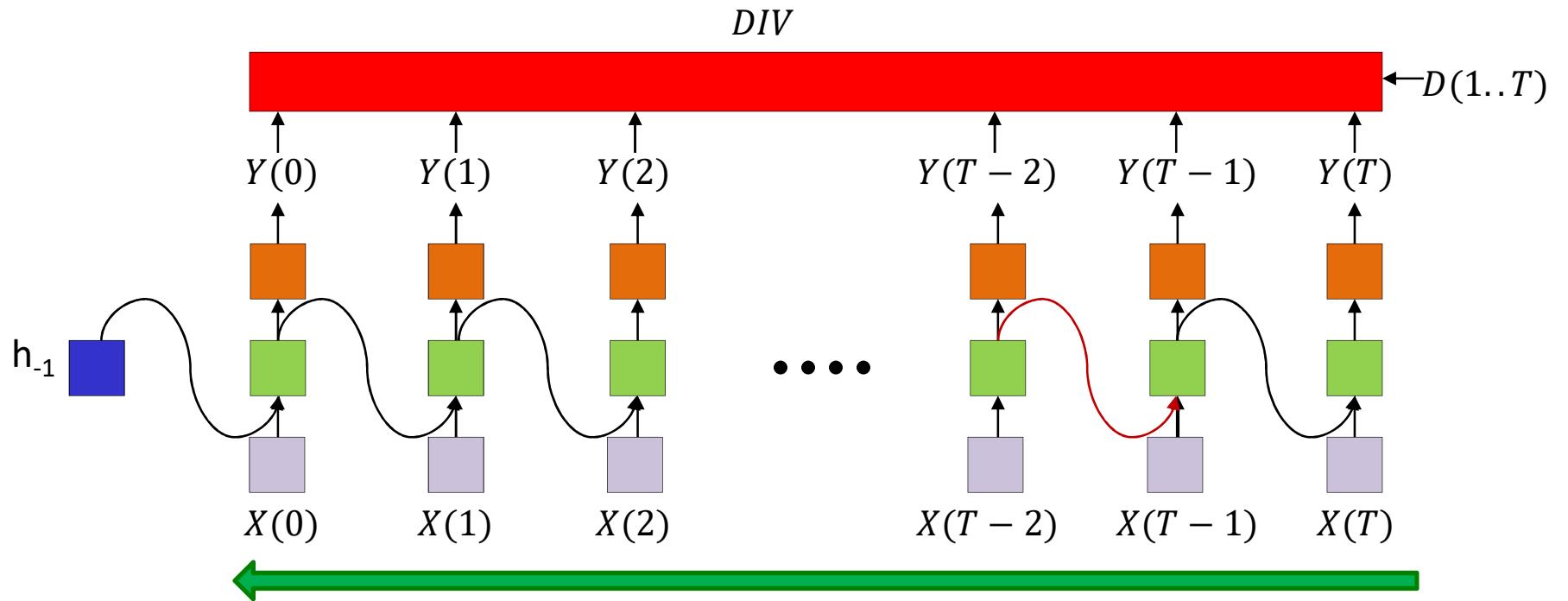


Continue computing derivatives  
going backward through time until..

$$\frac{dDIV}{dh_i(-1)} = \sum_j w_{ij}^{(11)} \frac{dDIV}{dZ_j^{(1)}(0)}$$

$$\nabla_{h_{-1}} DIV = \nabla_{Z^{(1)}(0)} DIV W^{(11)}$$

# Back Propagation Through Time

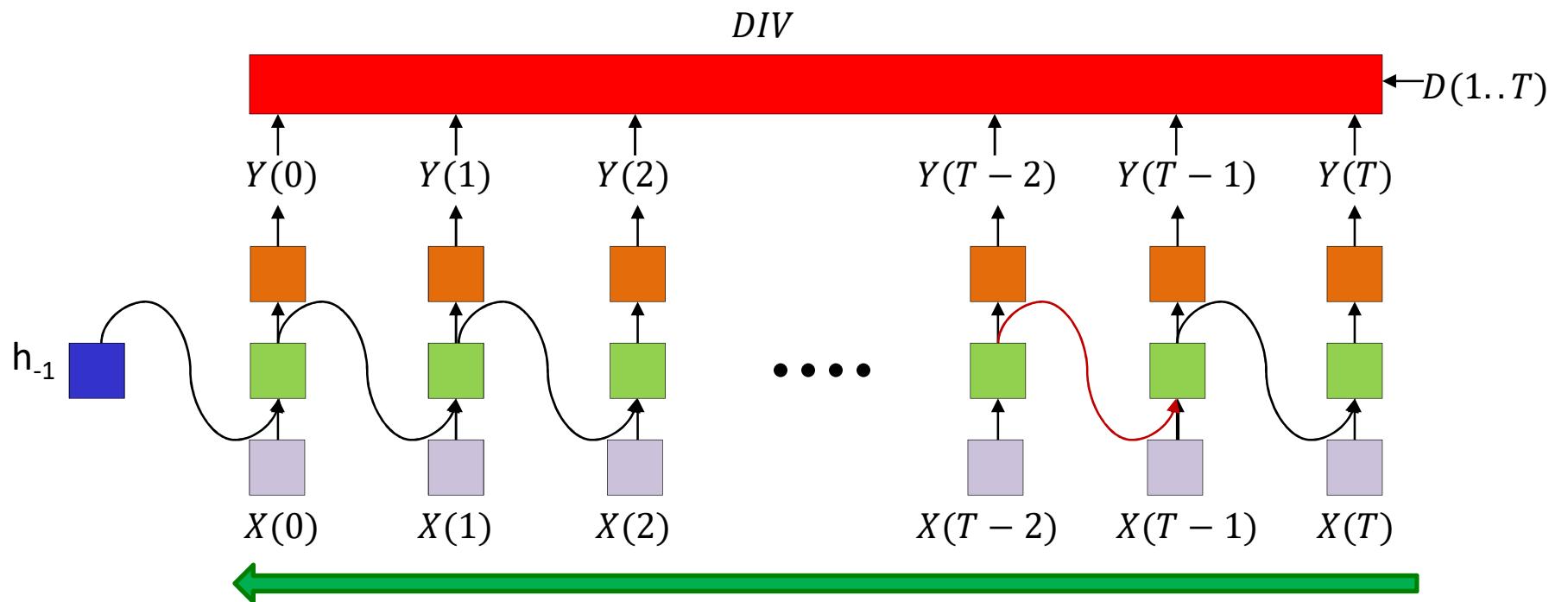


$$\frac{dDIV}{dh_i^{(k)}(t)} = \sum_j w_{i,j}^{(k+1)} \frac{dDIV}{dZ_j^{(k+1)}(t)} + \sum_j w_{i,j}^{(k,k)} \frac{dDIV}{dZ_j^{(k)}(t+1)}$$

Not showing derivatives  
at output neurons

$$\frac{dDIV}{dZ_i^{(k)}(t)} = \frac{dDIV}{dh_i^{(k)}(t)} f'_k(Z_i^{(k)}(t))$$

# Back Propagation Through Time



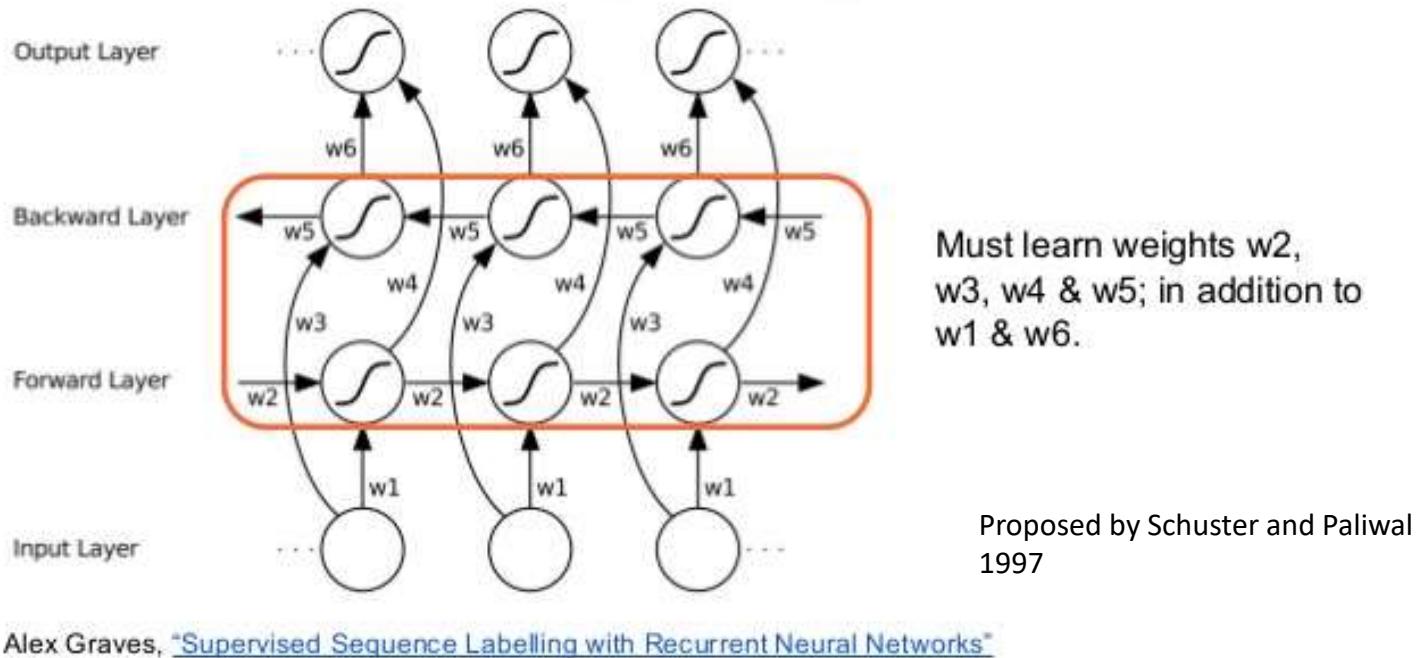
$$\frac{dDIV}{dh_i(-1)} = \sum_j w_{ij}^{(11)} \frac{dDIV}{dZ_j^{(1)}(0)}$$

$$\frac{dDIV}{dw_{ij}^{(1)}} = \sum_t \frac{dDIV}{dZ_j^{(1)}(t)} X_i(t)$$

$$\frac{dDIV}{dw_{ij}^{(11)}} = \sum_t \frac{dDIV}{dZ_j^{(1)}(t)} h_i(t-1)$$

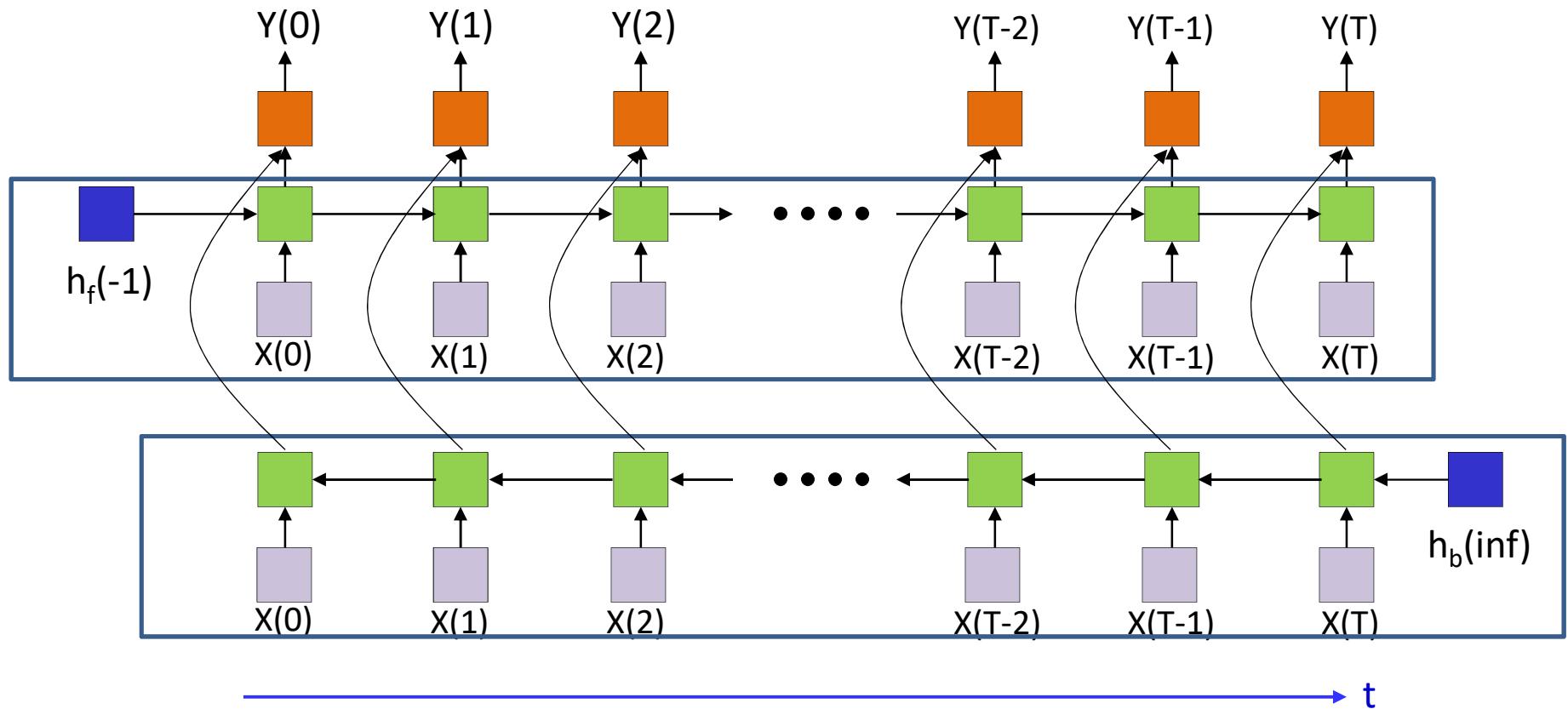
# Extensions to the RNN: *Bidirectional RNN*

## Bidirectional RNN (BRNN)



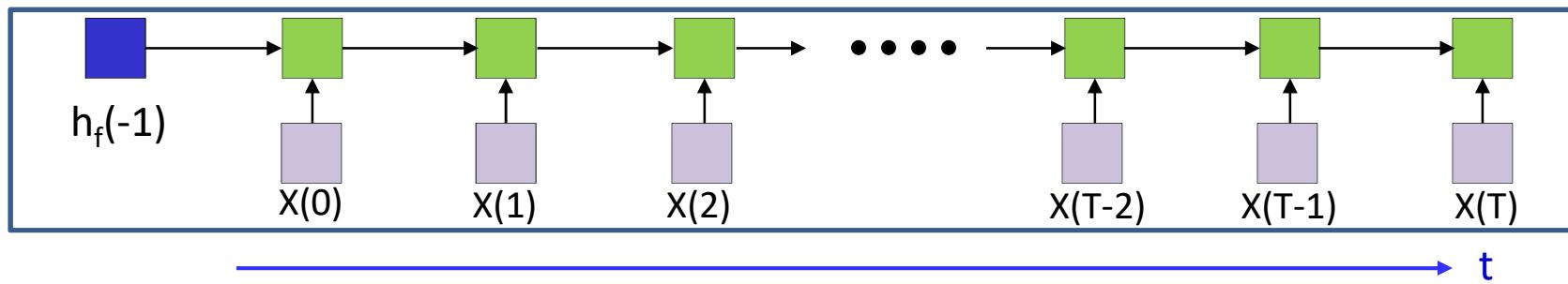
- RNN with both forward and backward recursion
  - Explicitly models the fact that just as the future can be predicted from the past, the past can be deduced from the future

# Bidirectional RNN



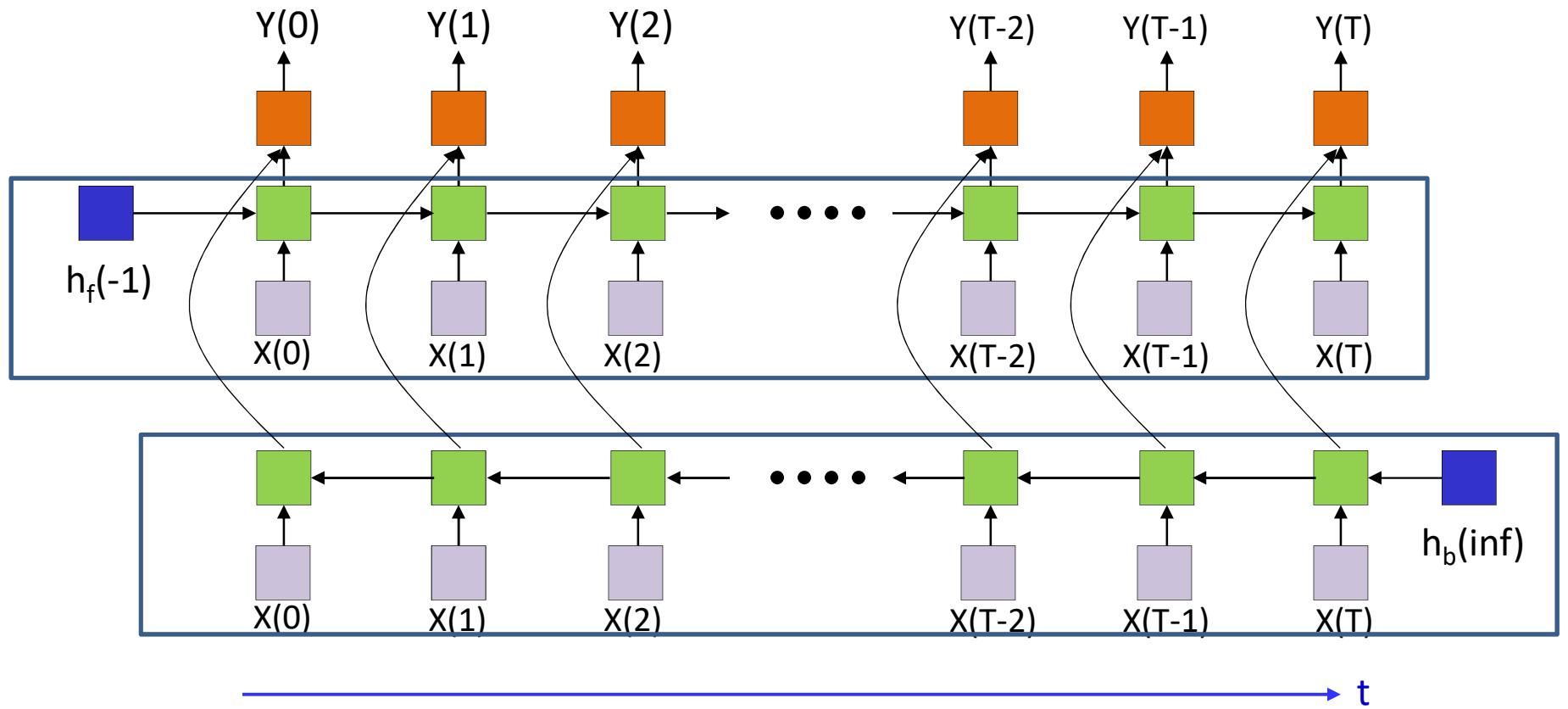
- A forward net process the data from  $t=0$  to  $t=T$
- A backward net processes it backward from  $t=T$  down to  $t=0$

# Bidirectional RNN: Processing an input string



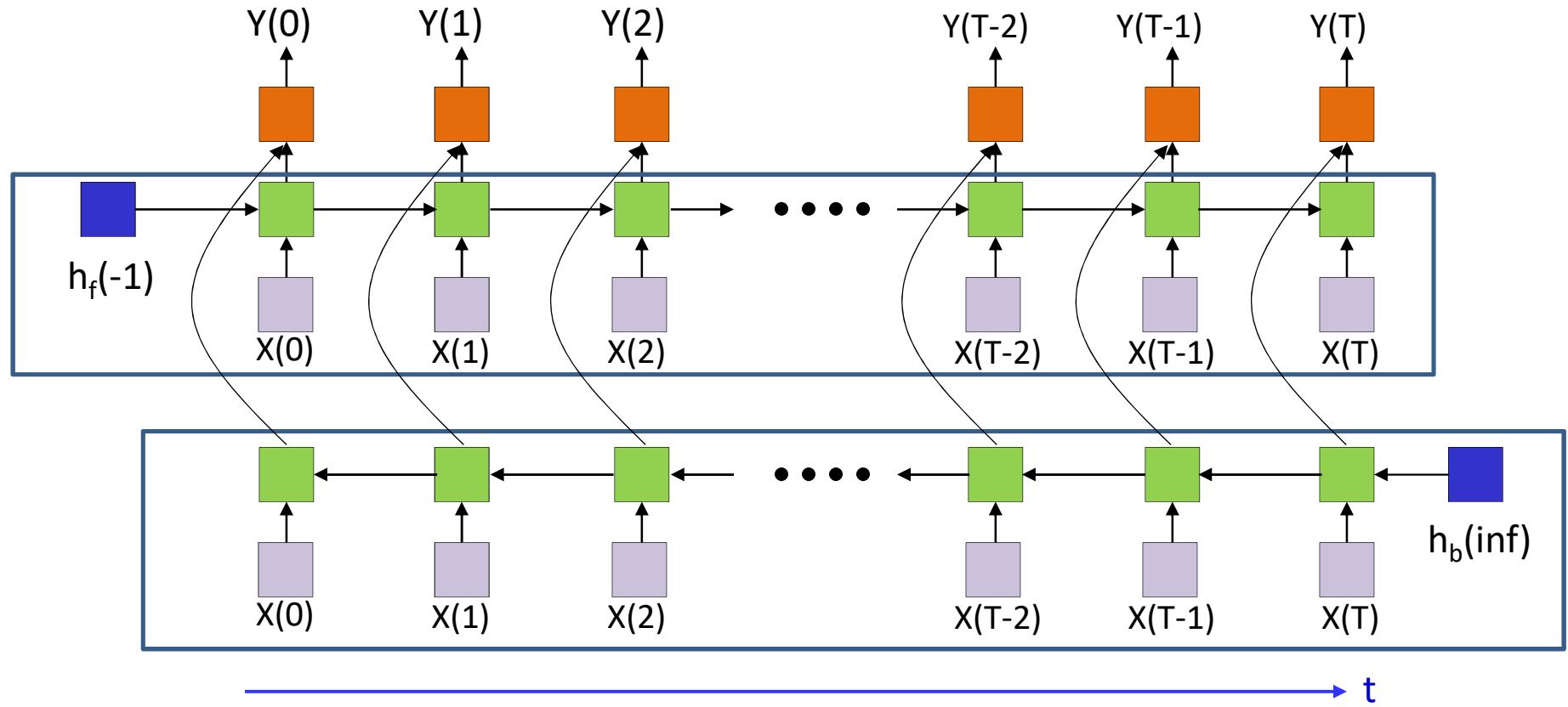
- The forward net process the data from  $t=0$  to  $t=T$ 
  - Only computing the hidden states, initially

# Bidirectional RNN: Processing an input string



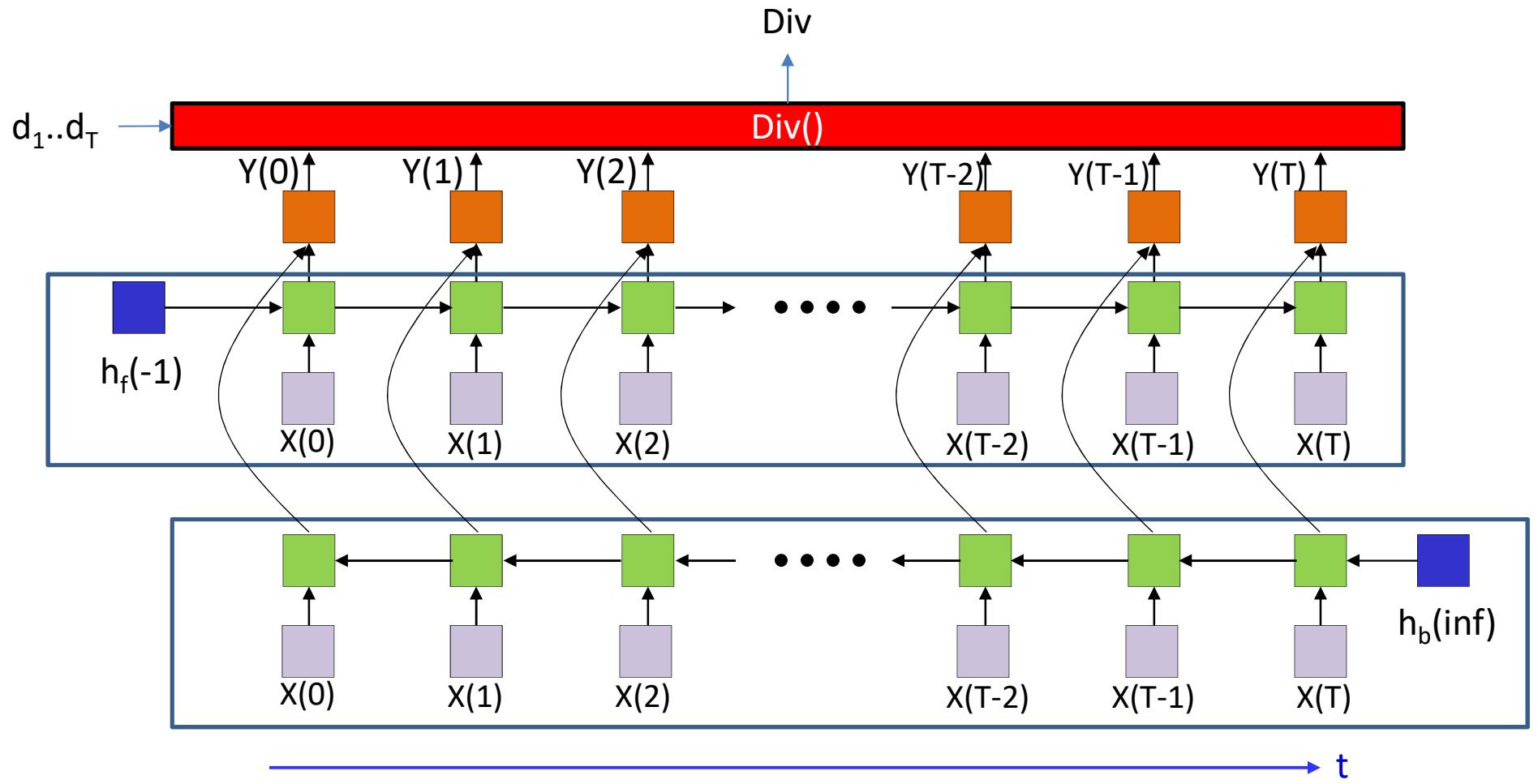
- The computed states of both networks are used to compute the final output at each time

# Backpropagation in BRNNs



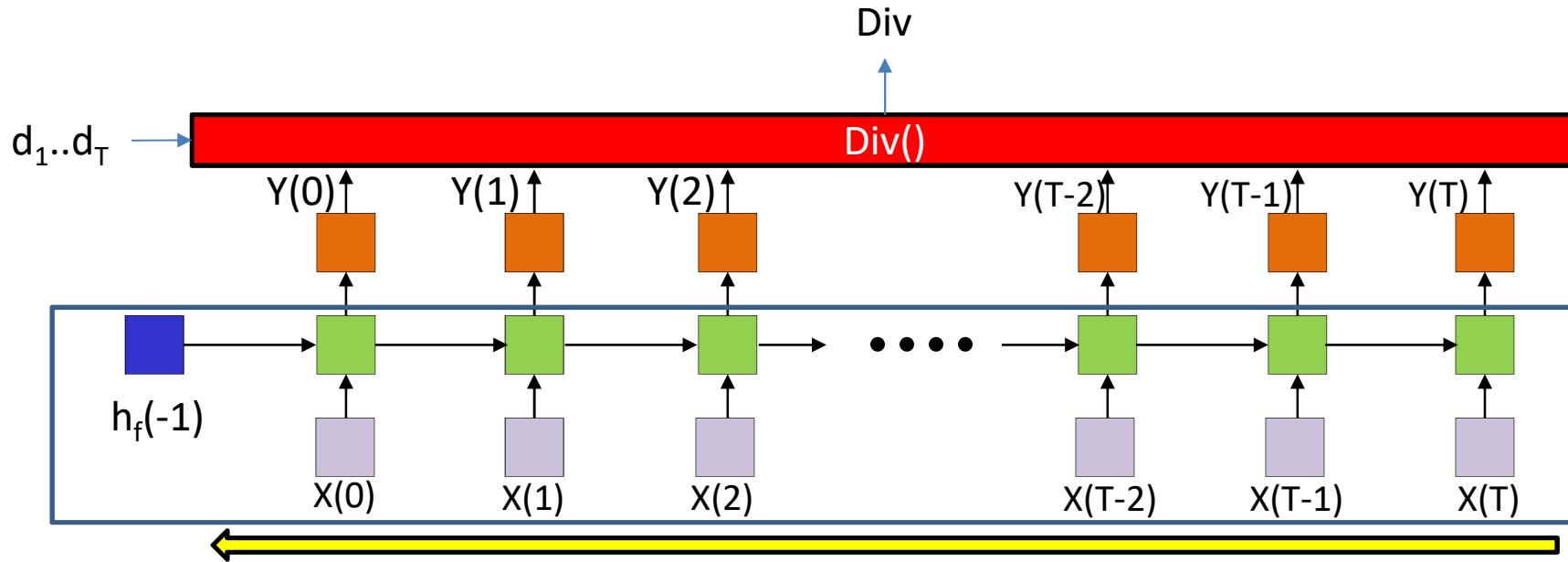
- Forward pass: Compute both forward and backward networks and final output

# Backpropagation in BRNNs



- Backward pass: Define a divergence from the desired output

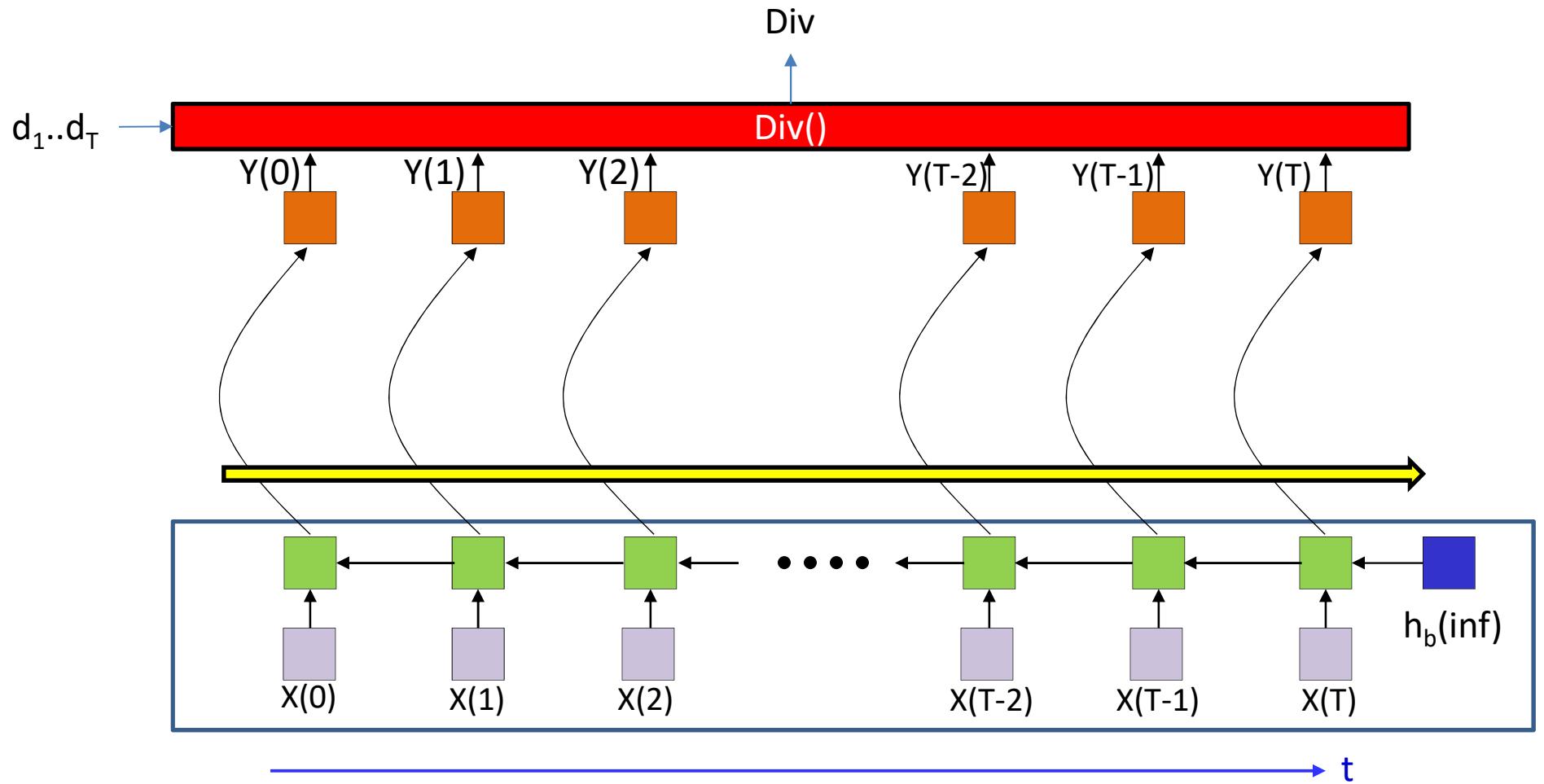
# Backpropagation in BRNNs



- Backward pass: Define a divergence from the desired output
- Separately perform back propagation on both nets
  - **From  $t=T$  down to  $t=0$  for the forward net**

$t$

# Backpropagation in BRNNs



- Backward pass: Define a divergence from the desired output
- Separately perform back propagation on both nets
  - From  $t=T$  down to  $t=0$  for the forward net
  - **From  $t=0$  up to  $t=T$  for the backward net**

- ① Consider the following three tasks of filling in the blanks in a text:
  - I am \_\_\_\_\_.
  - I am \_\_\_\_\_ very hungry.
  - I am \_\_\_\_\_ very hungry, I could eat half a pig.
- ② We might fill the blanks with very different words such as “happy”, “not”, and “very”.
- ③ End of the phrase (if available) conveys significant information about which word to pick.
- ④ A sequence model should be capable of taking advantage of this information.
- ⑤ For instance, to recognize whether “Green” refers to “Mr. Green” or to the color, longer-range context is equally vital.