

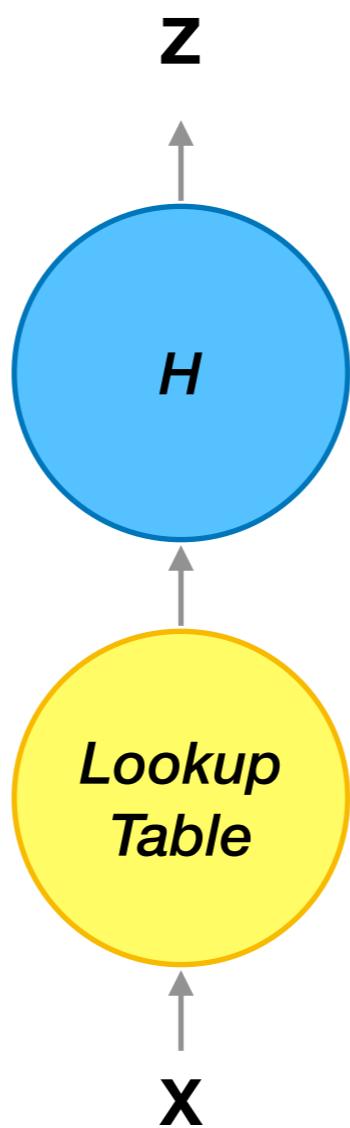
Recurrent Neural Network

Boris Zubarev

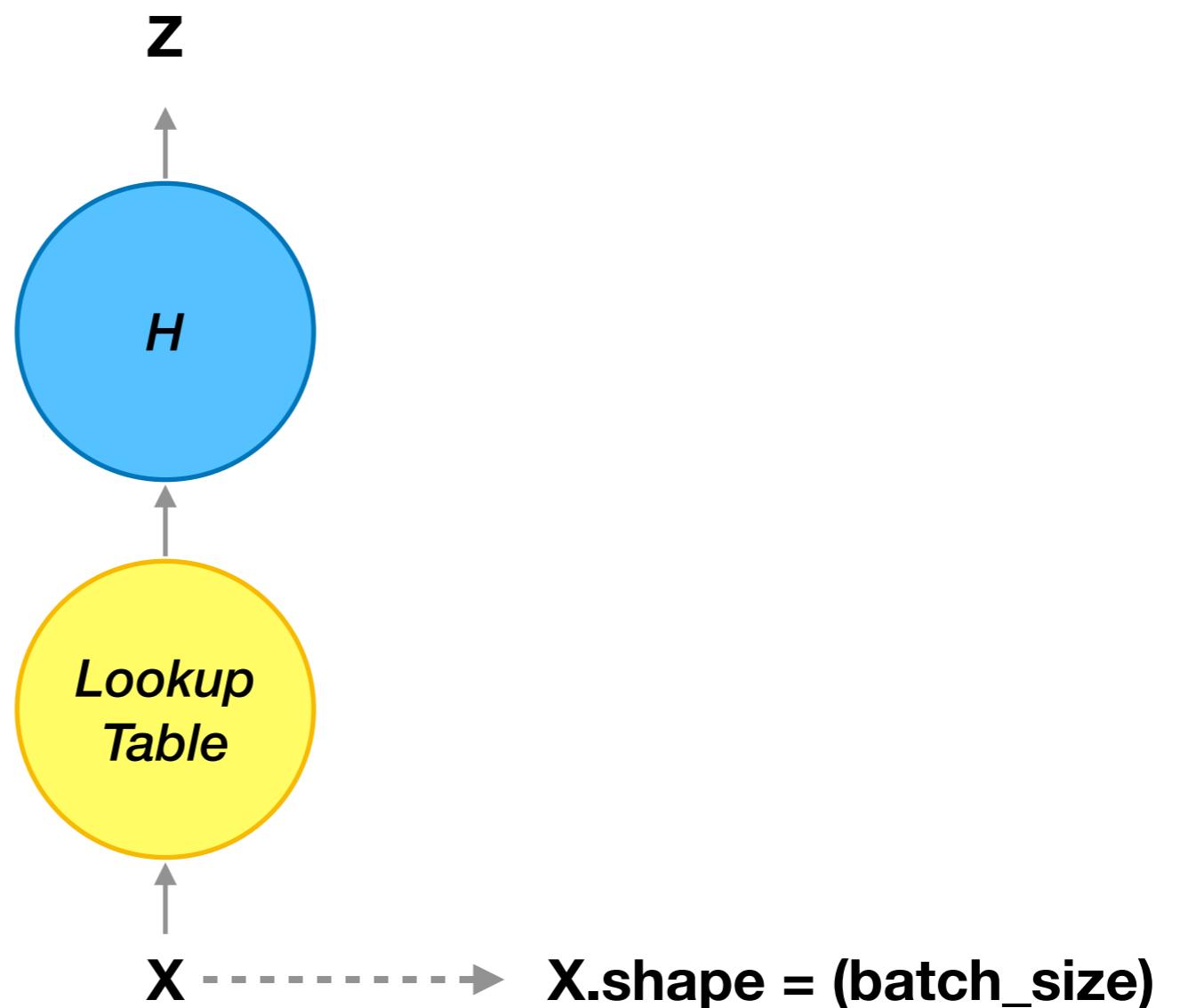


@bobazooba

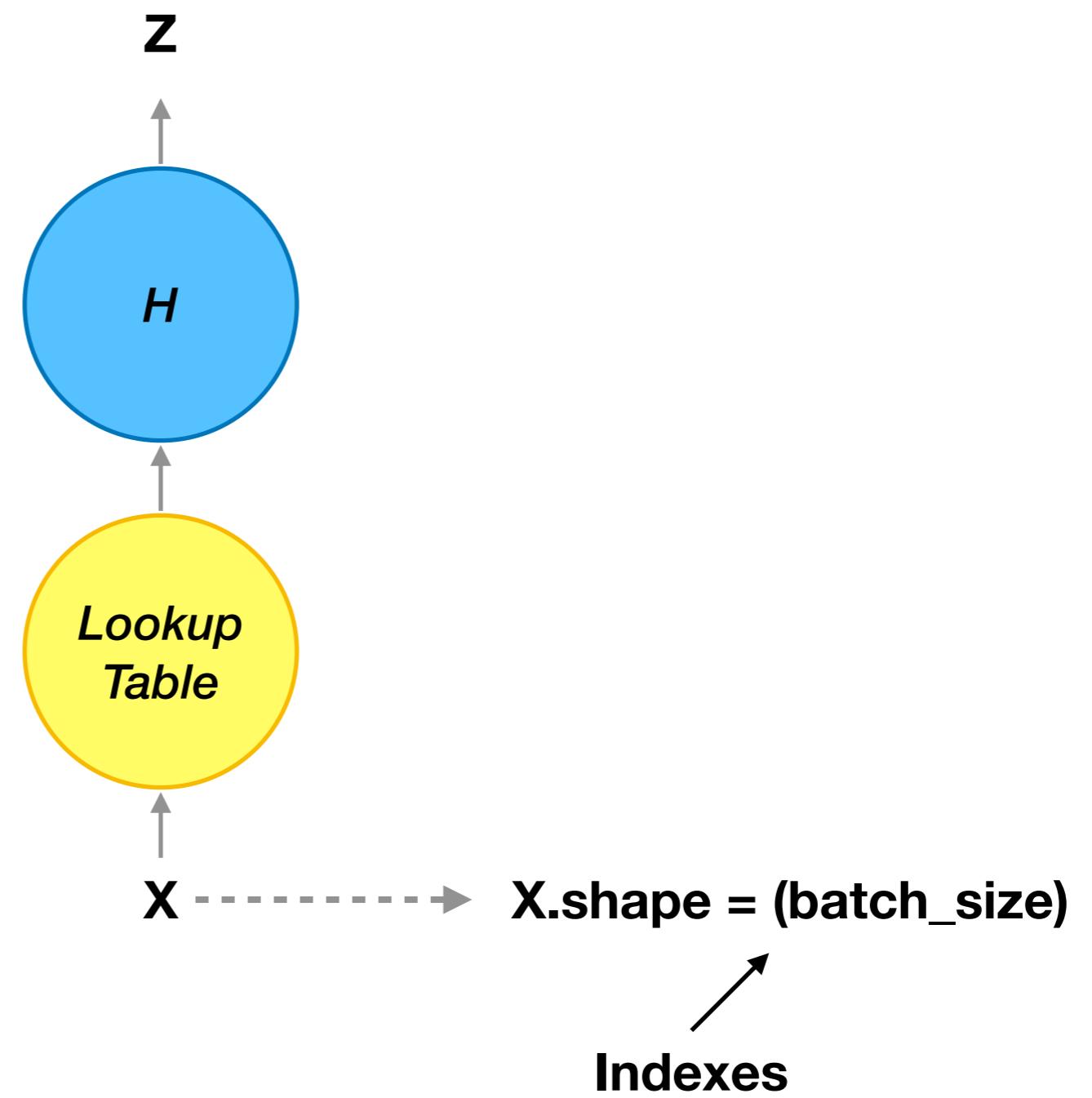
Neural Network



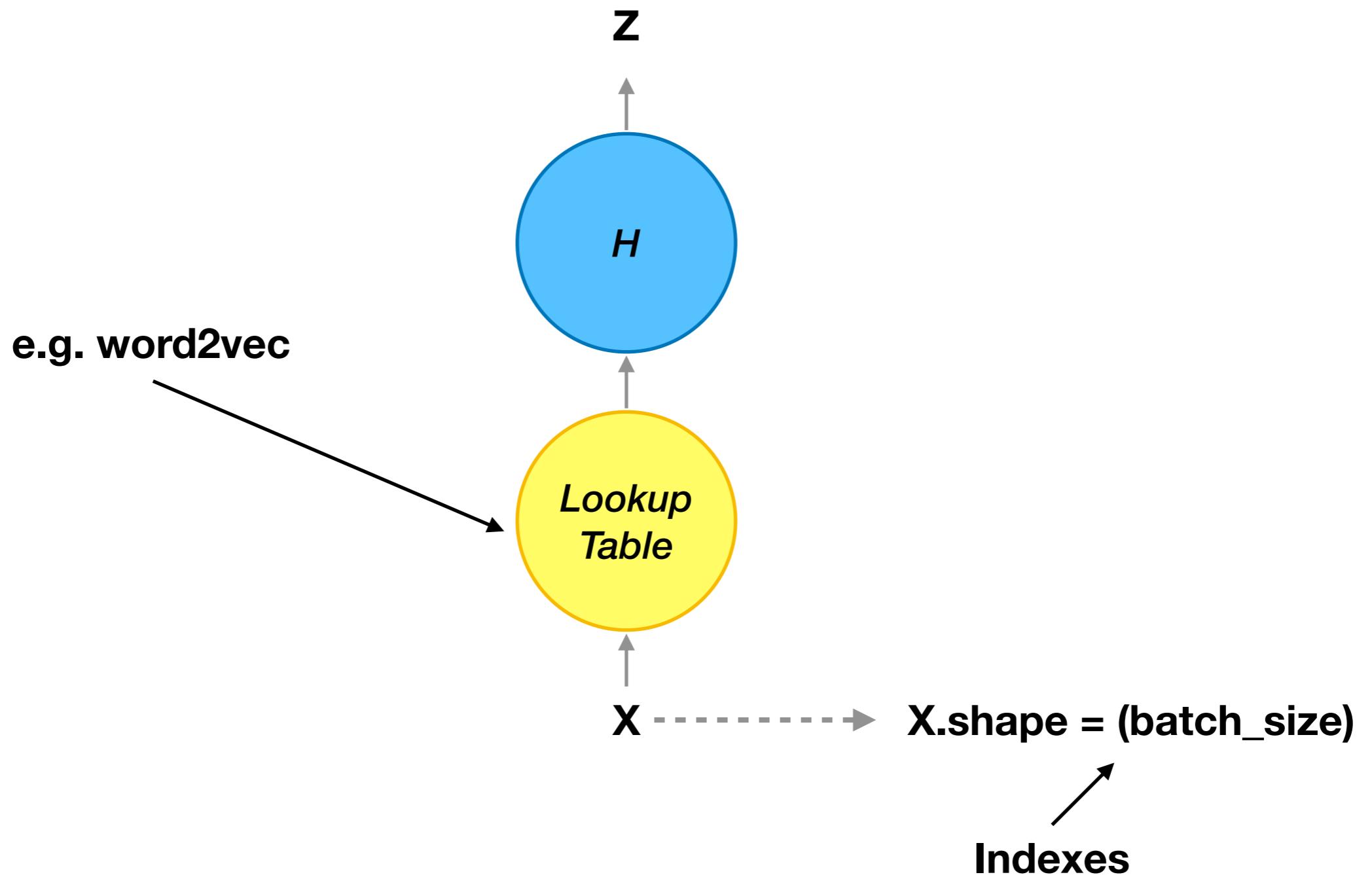
Neural Network



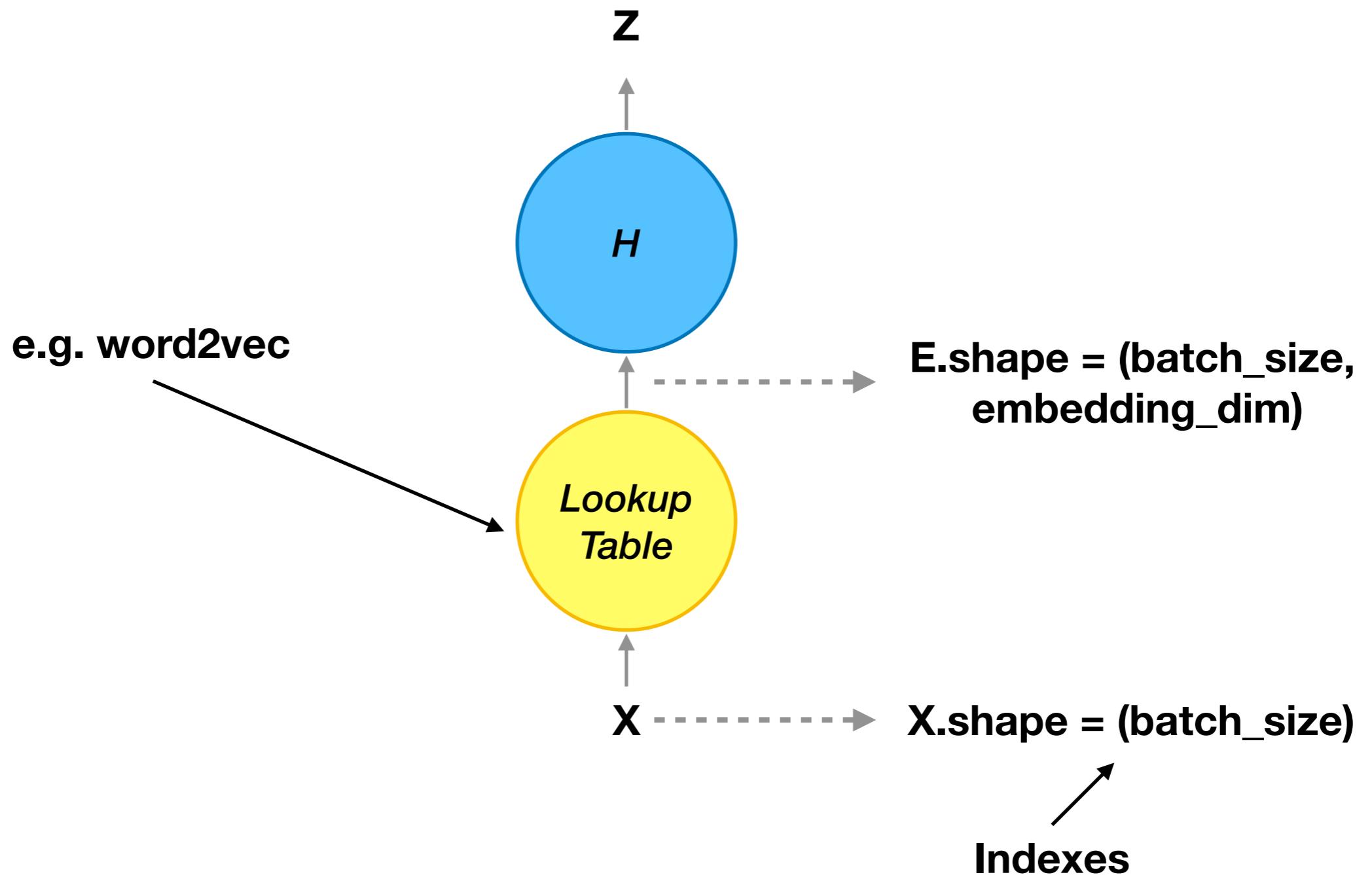
Neural Network



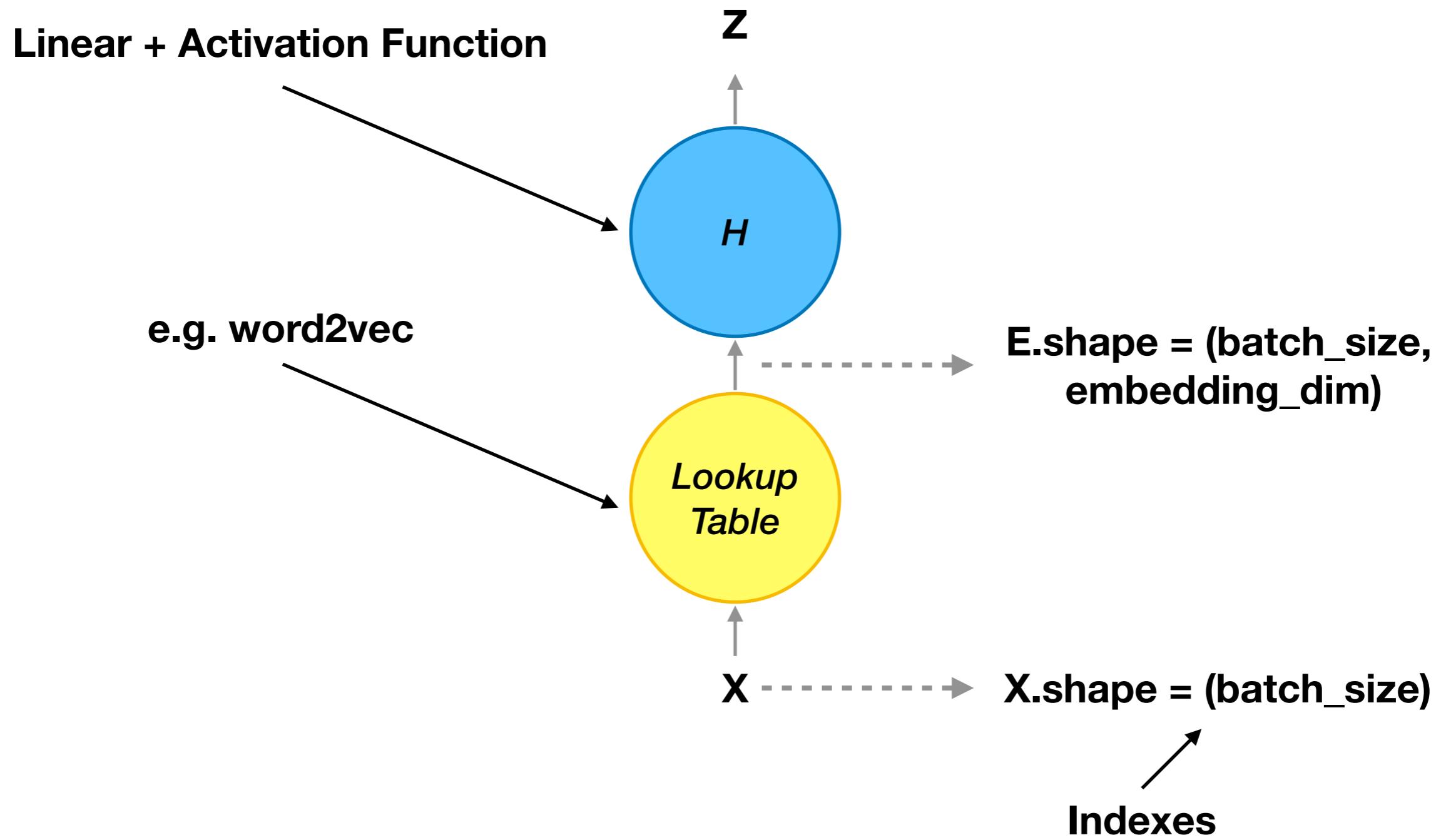
Neural Network



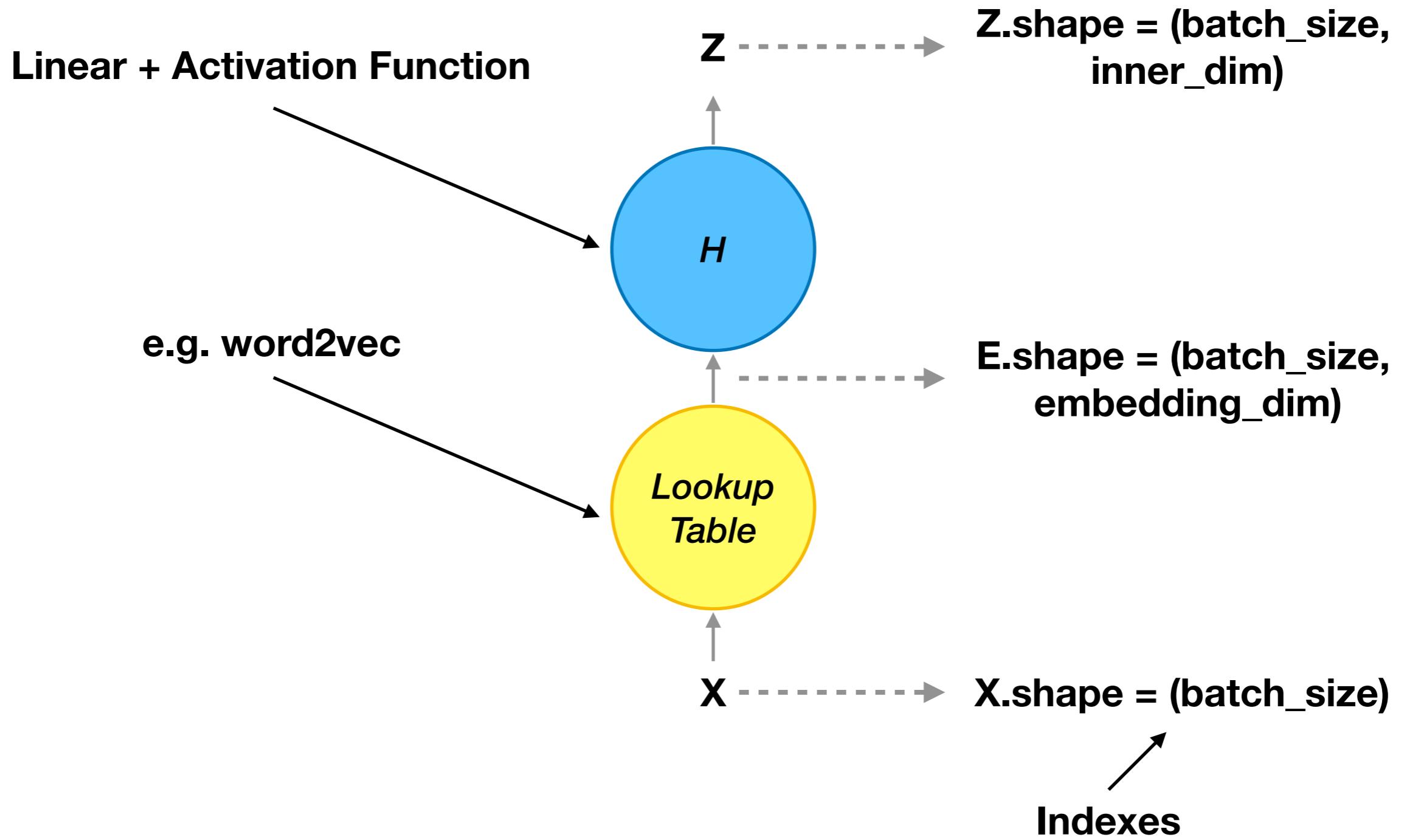
Neural Network



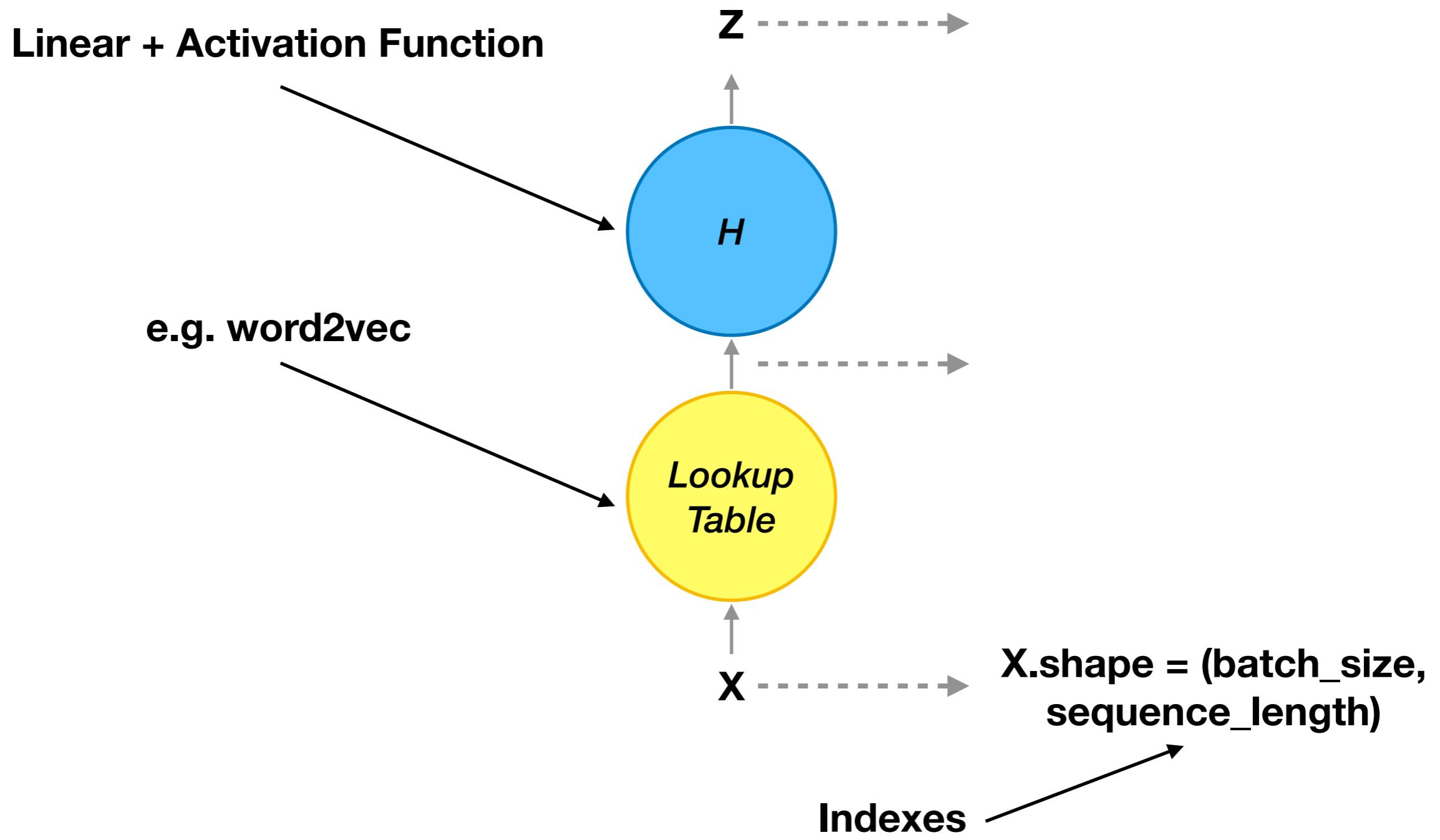
Neural Network



Neural Network

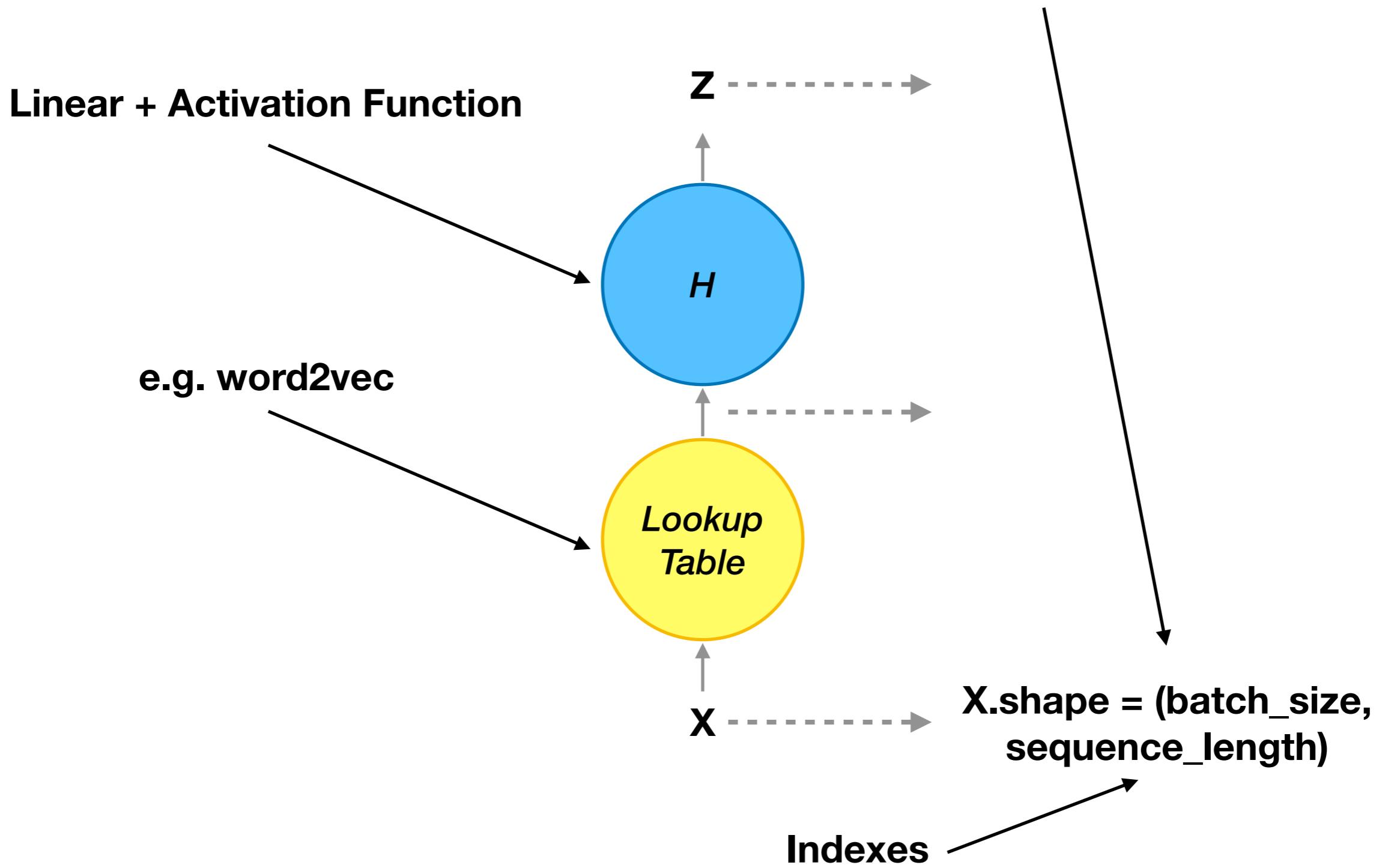


Neural Network



Neural Network

What if we have different lengths?



Padding

[Мама, мыла, раму]

[Покупать, новый, телефон, затратно]

[Bay]

[Мне, очень, не, нравится ваш, банк]

Padding

[Мама, мыла, раму]

[Покупать, новый, телефон, затратно]

[Bay]

[Мне, очень, не, нравится ваш, банк]



`max_sequence_length = 8`

Depend of your dataset statistics



Padding

[Мама, мыла, раму]

[Покупать, новый, телефон, затратно]

[Bay]

[Мне, очень, не, нравится ваш, банк]



max_sequence_length = 8

Depend of your dataset statistics



[Мама, мыла, раму, PAD, PAD, PAD, PAD, PAD]

[Покупать, новый, телефон, затратно, PAD, PAD, PAD, PAD]

[Bay, PAD, PAD, PAD, PAD, PAD, PAD]

[Мне, очень, не, нравится, ваш, банк, PAD, PAD]

Padding

[Мама, мыла, раму, PAD, PAD, PAD, PAD, PAD]

[Покупать, новый, телефон, затратно, PAD, PAD, PAD, PAD, PAD]

[Bay, PAD, PAD, PAD, PAD, PAD, PAD, PAD]

[Мне, очень, не, нравится, ваш, банк, PAD, PAD]

Padding

[Мама, мыла, раму, PAD, PAD, PAD, PAD, PAD, PAD]

[Покупать, новый, телефон, затратно, PAD, PAD, PAD, PAD, PAD]

[Bay, PAD, PAD, PAD, PAD, PAD, PAD, PAD]

[Мне, очень, не, нравится, ваш, банк, PAD, PAD]

Indexing

Collect word2index vocabulary



Padding

[Мама, мыла, раму, PAD, PAD, PAD, PAD, PAD]

[Покупать, новый, телефон, затратно, PAD, PAD, PAD, PAD]

[Bay, PAD, PAD, PAD, PAD, PAD, PAD]

[Мне, очень, не, нравится, ваш, банк, PAD, PAD]

Indexing

Collect word2index vocabulary

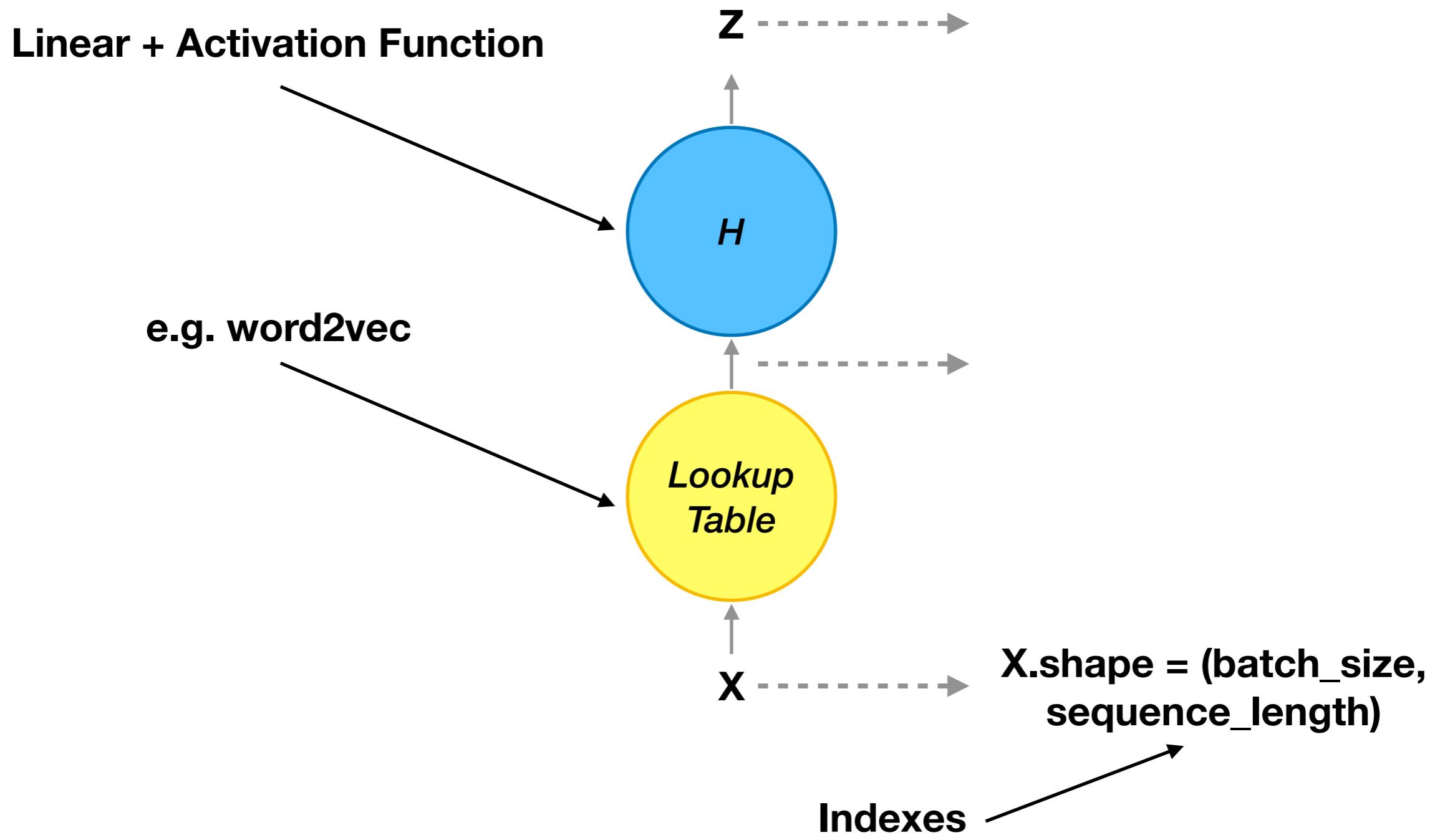
34 2 3 0 0 0 0 0

64 23 1515 45 0 0 0 0

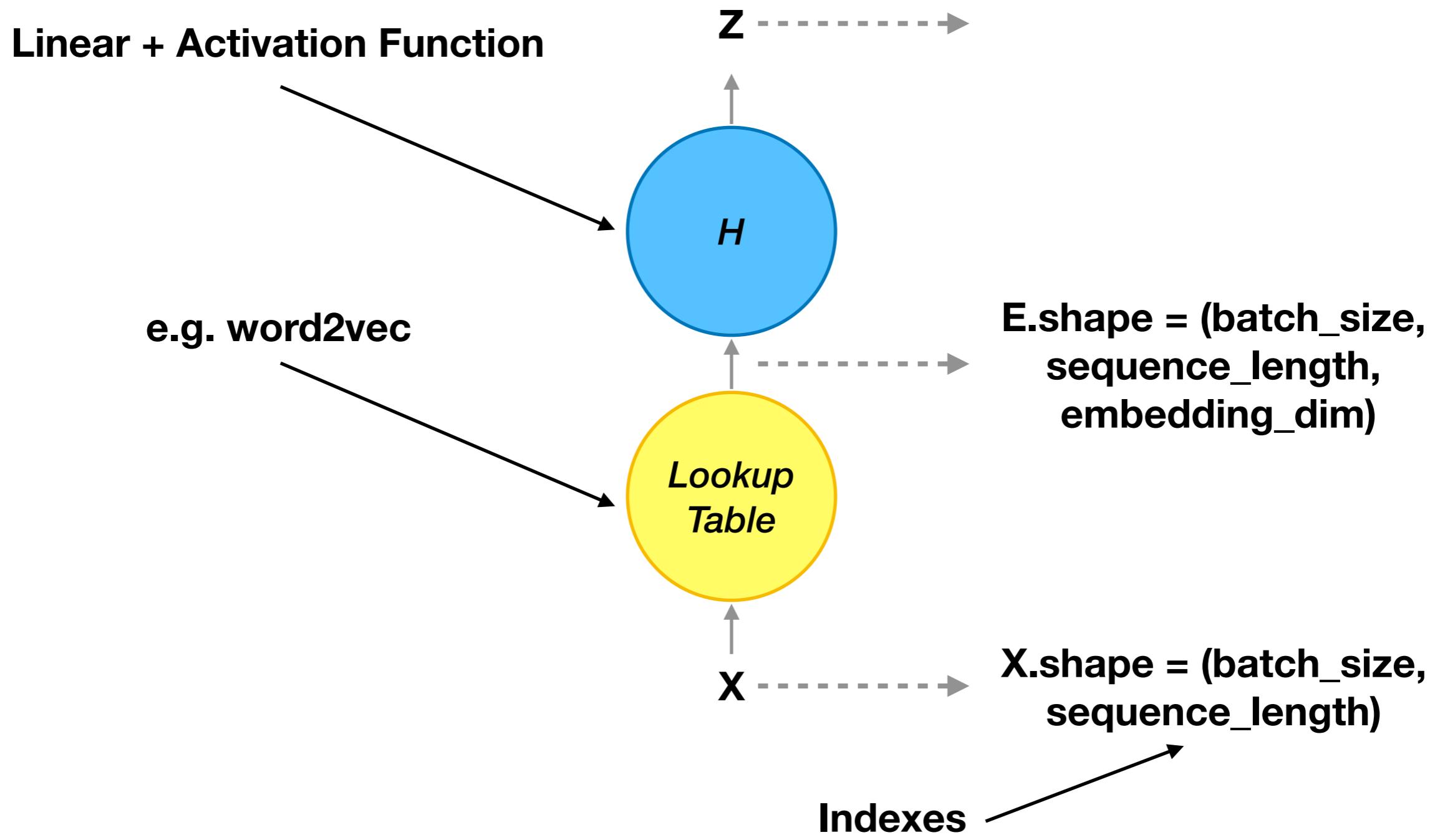
55 0 0 0 0 0 0 0

36 43 22 7 124 356 0 0

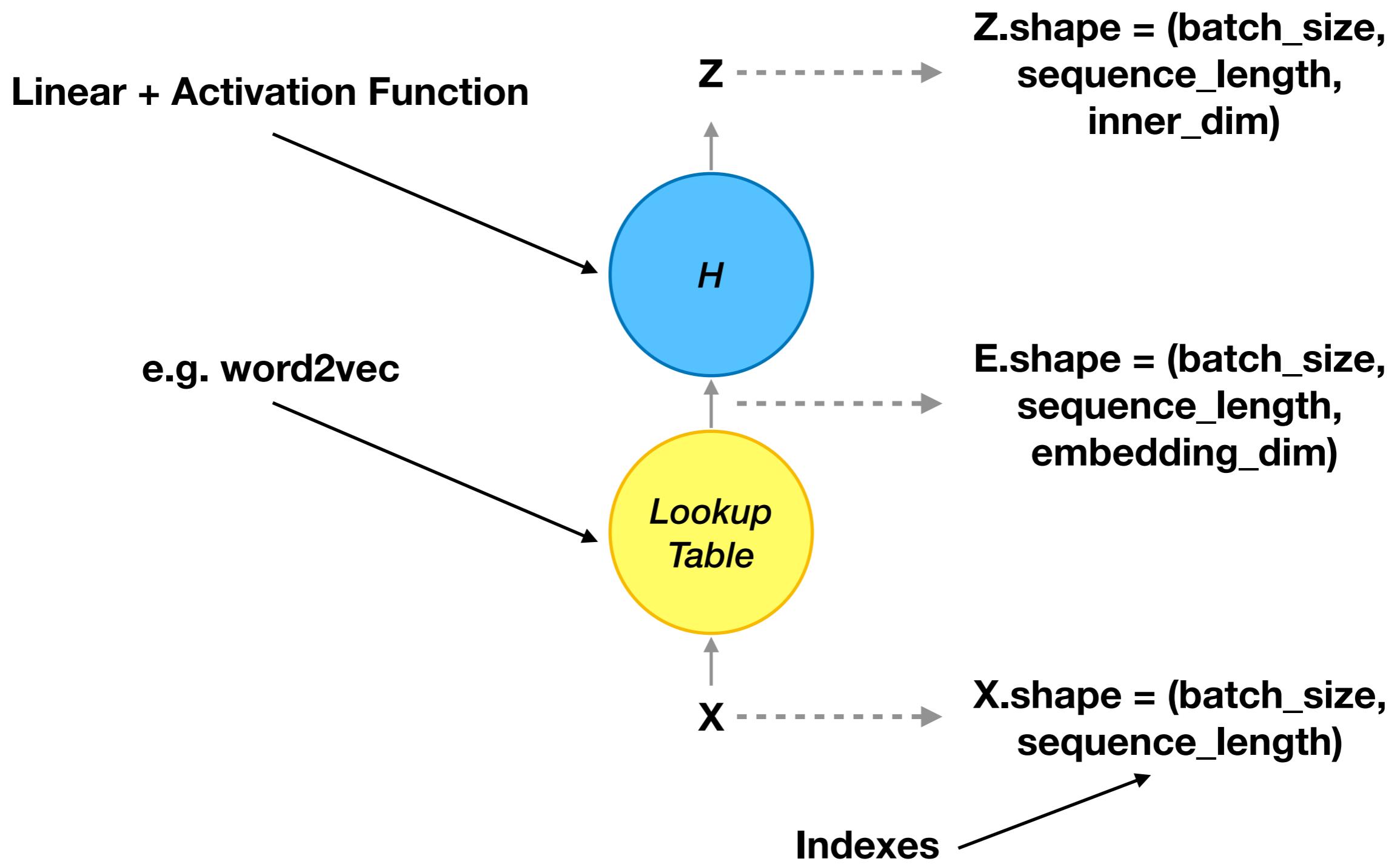
Neural Network



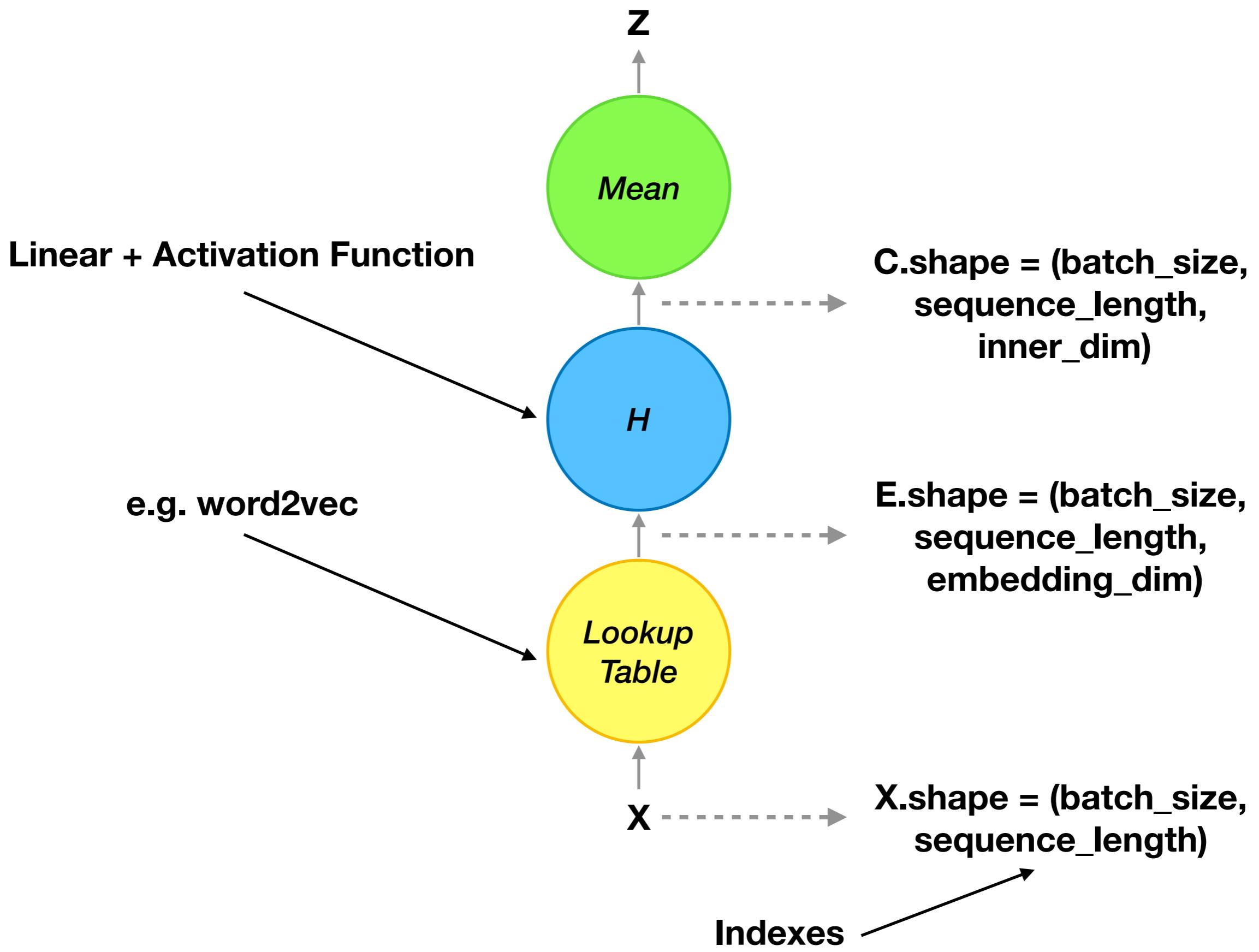
Neural Network



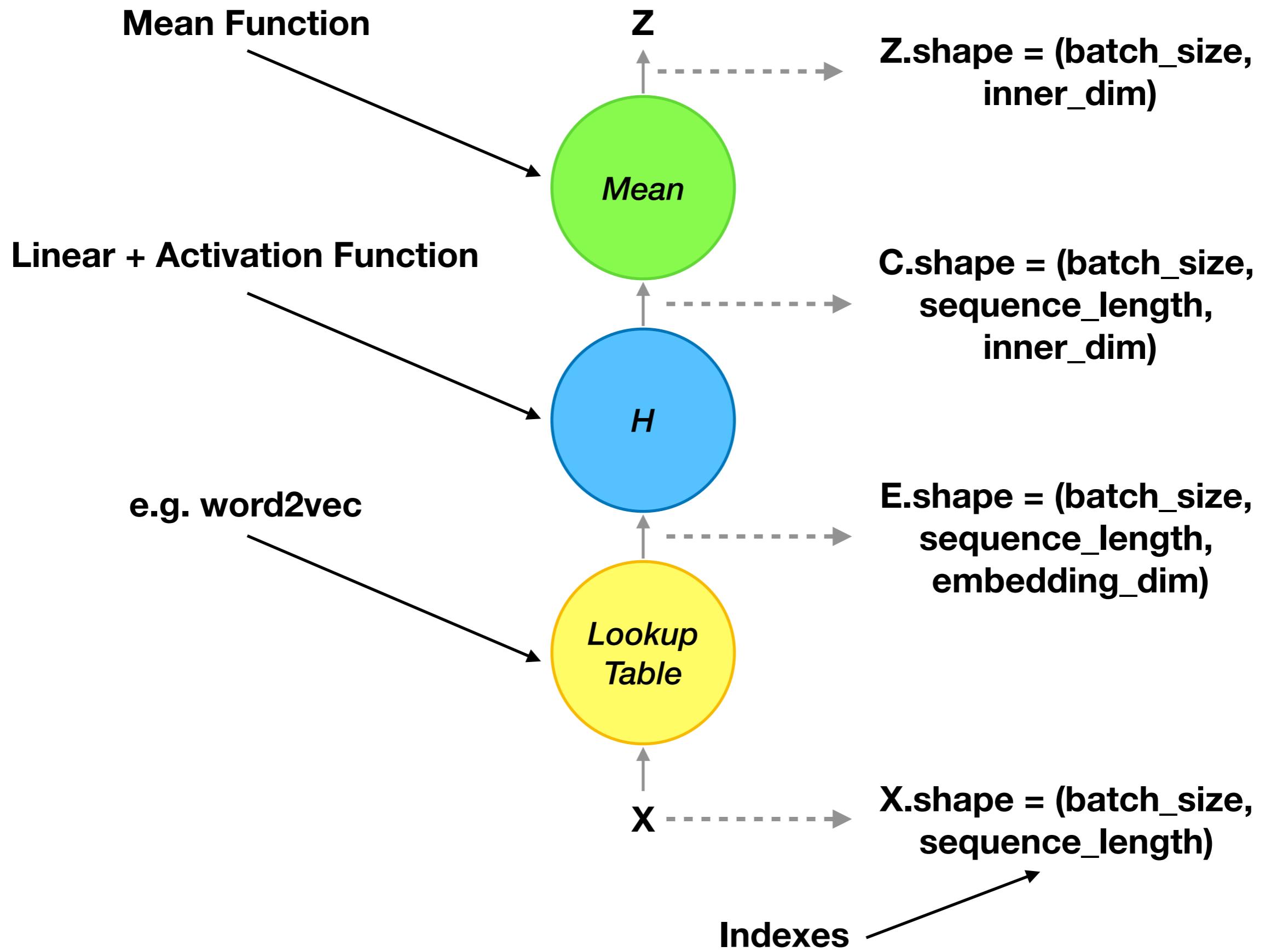
Neural Network



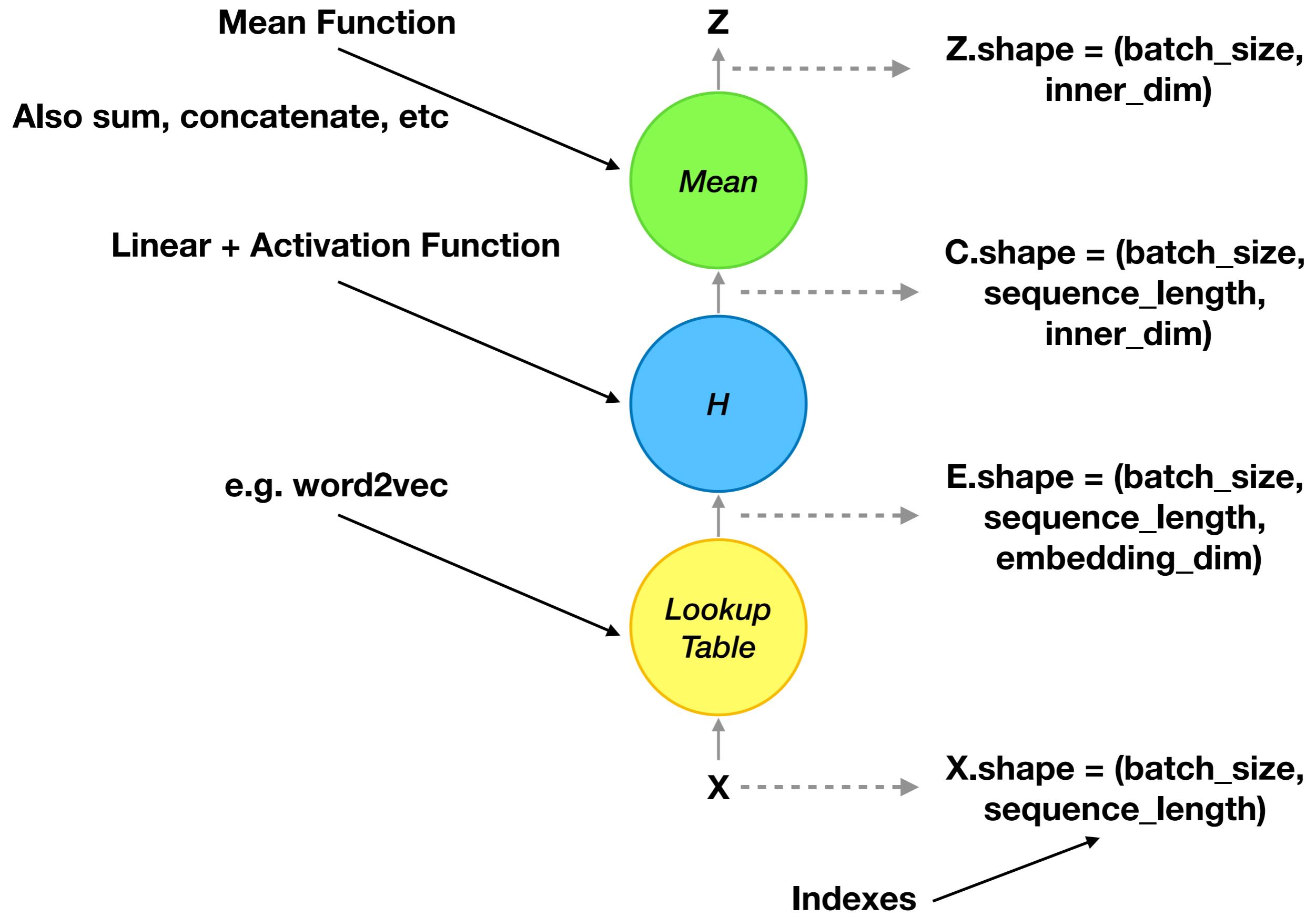
Neural Network



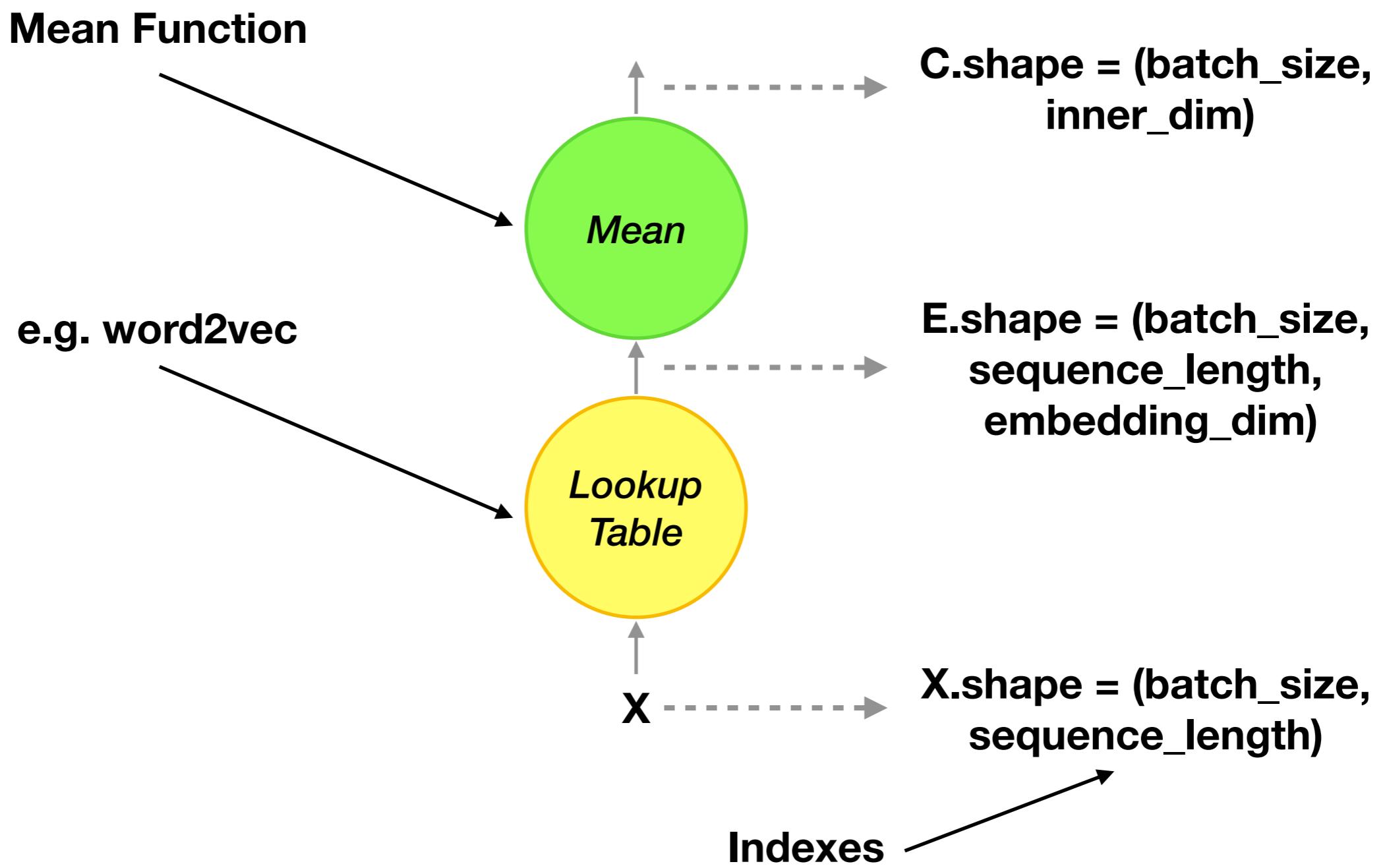
Neural Network



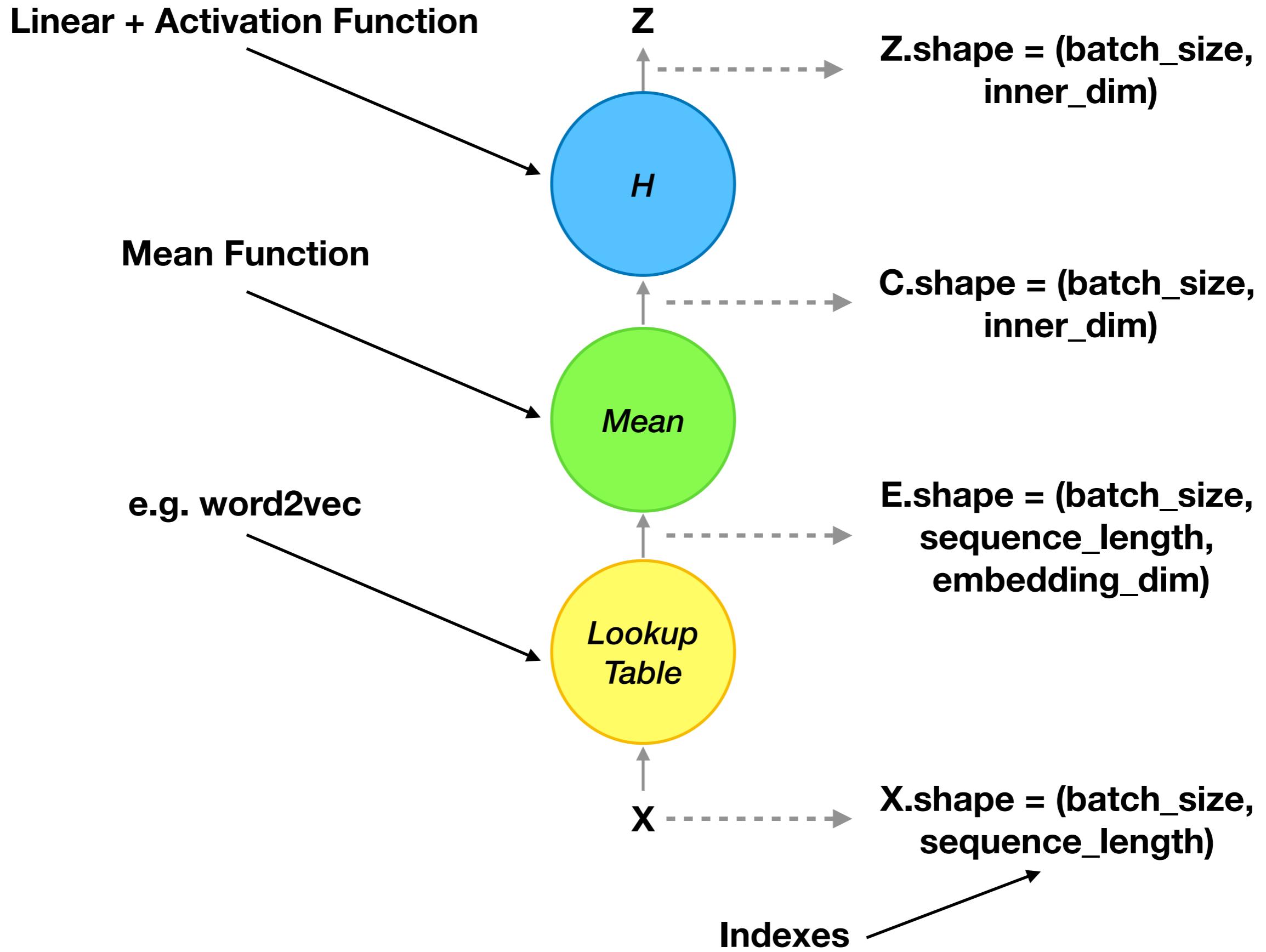
Neural Network



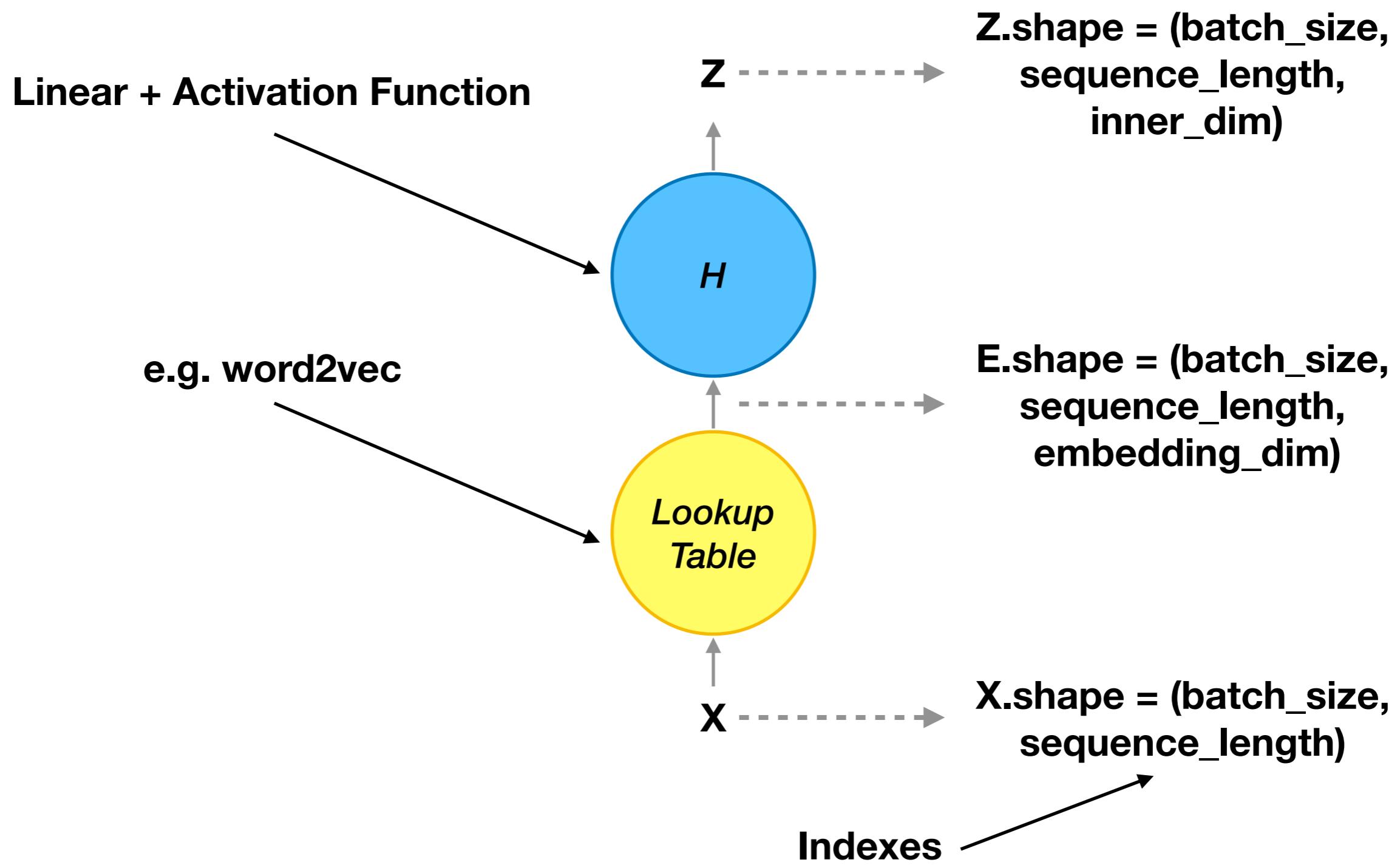
Neural Network



Neural Network

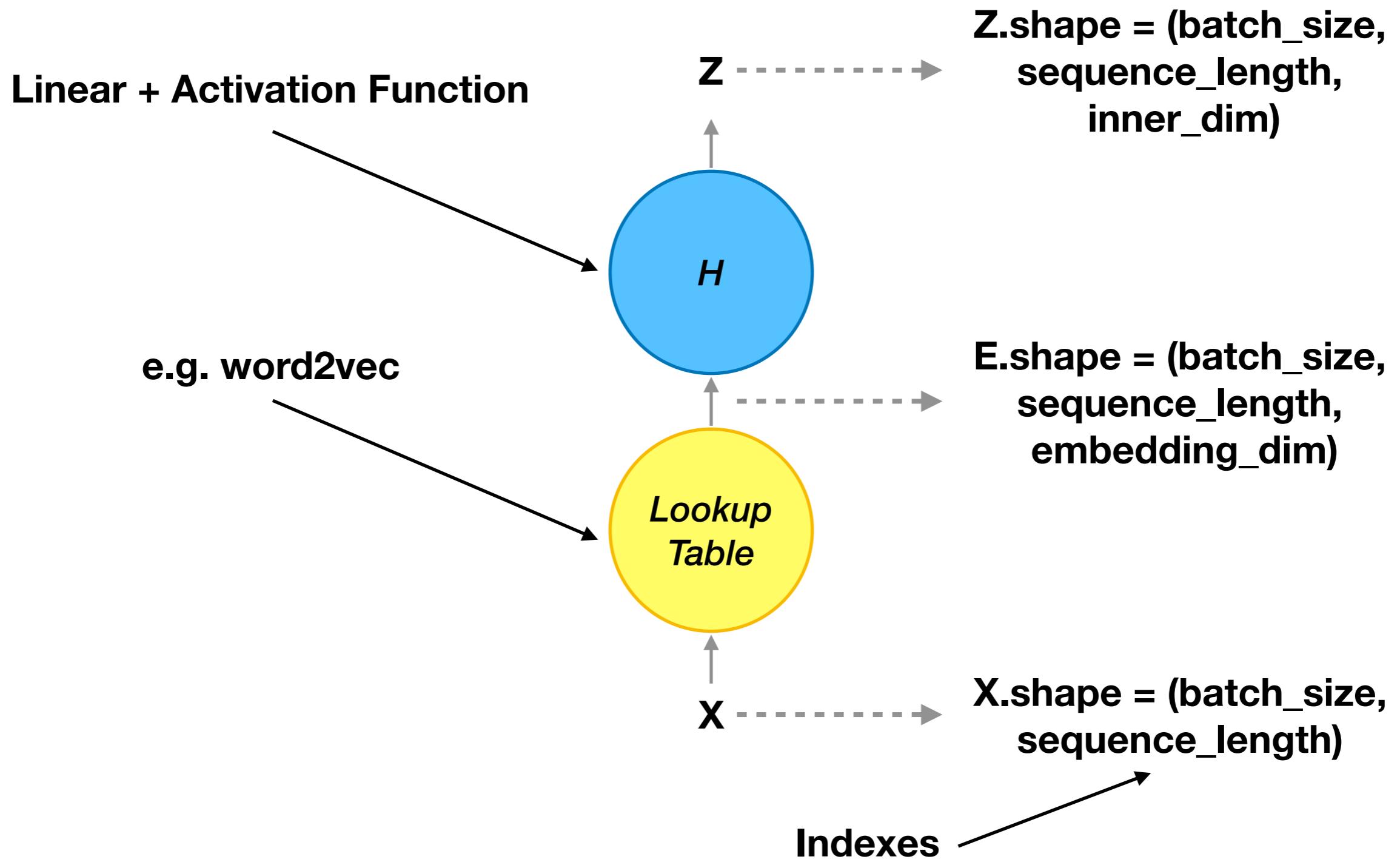


Recurrent Neural Network



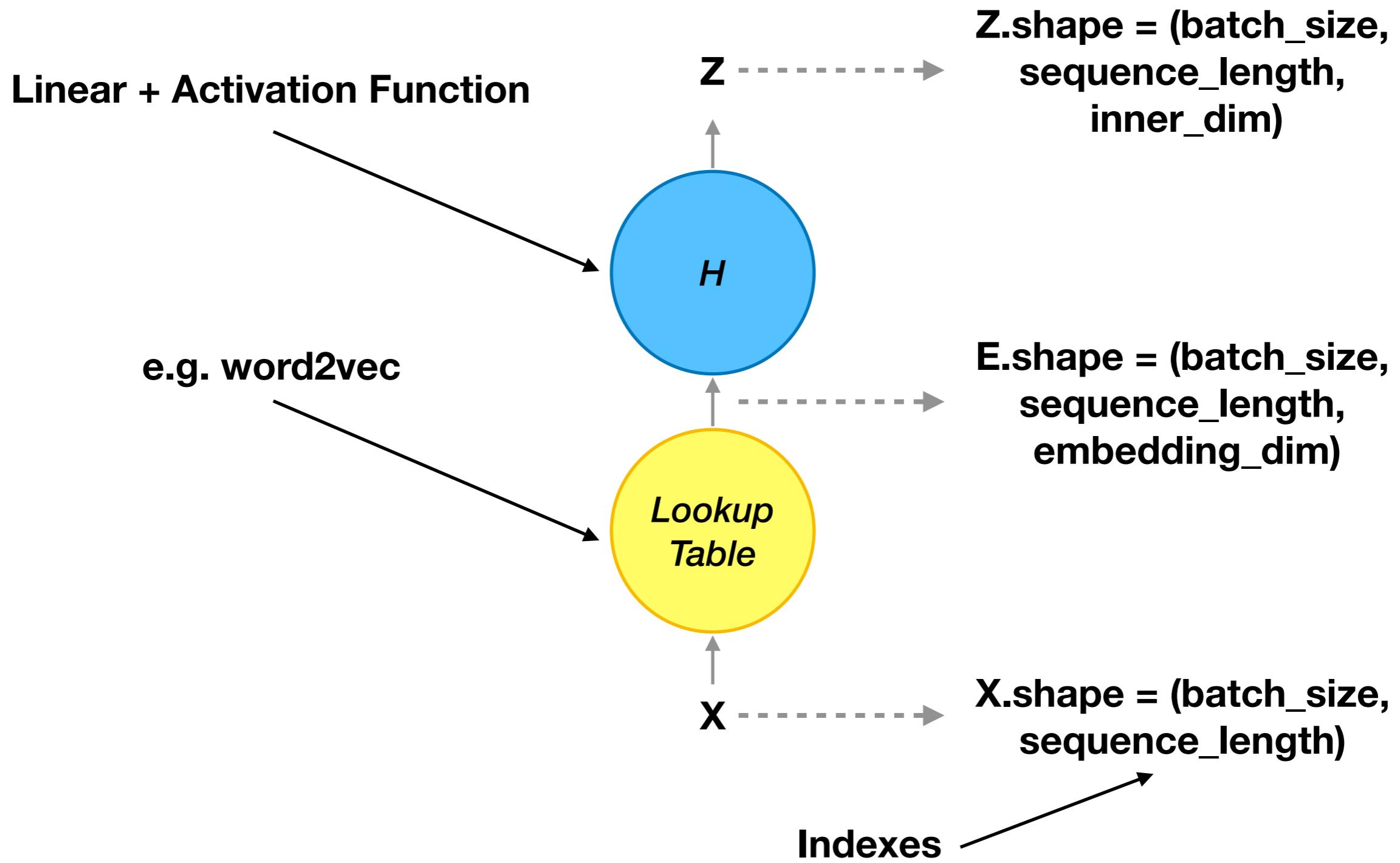
Recurrent Neural Network

Maybe we can take last token?

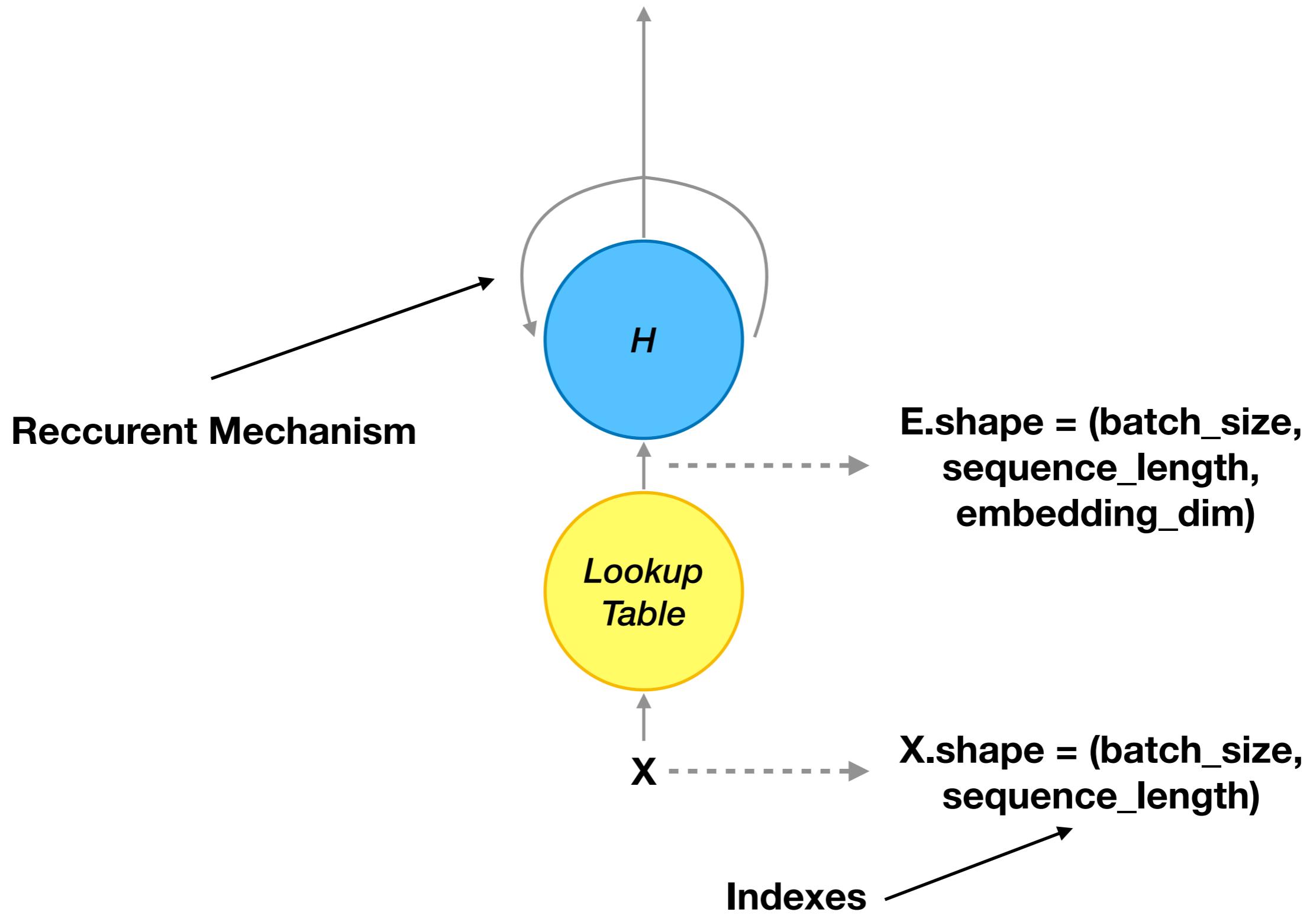


Recurrent Neural Network

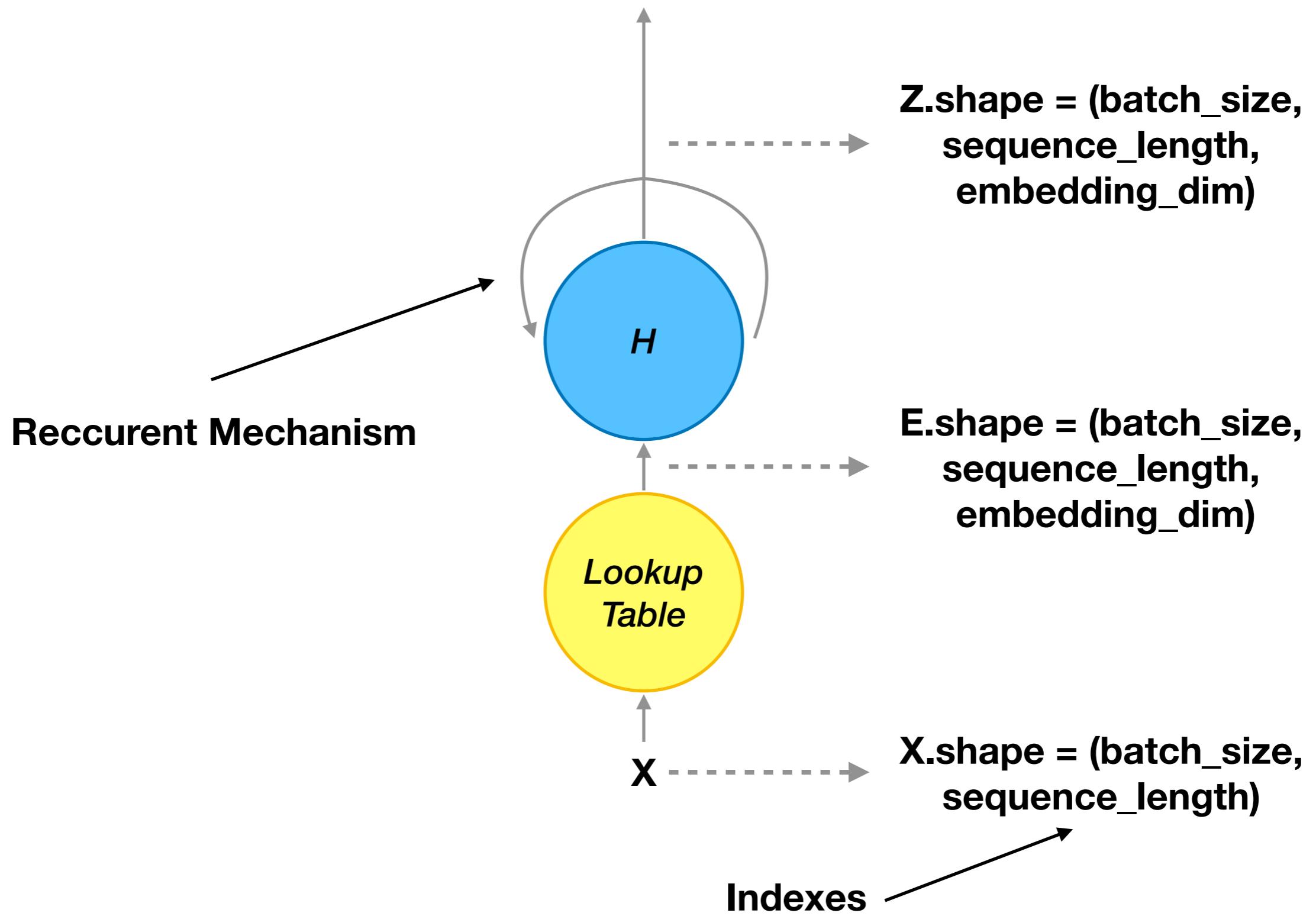
No, because we don't have information about whole sentence



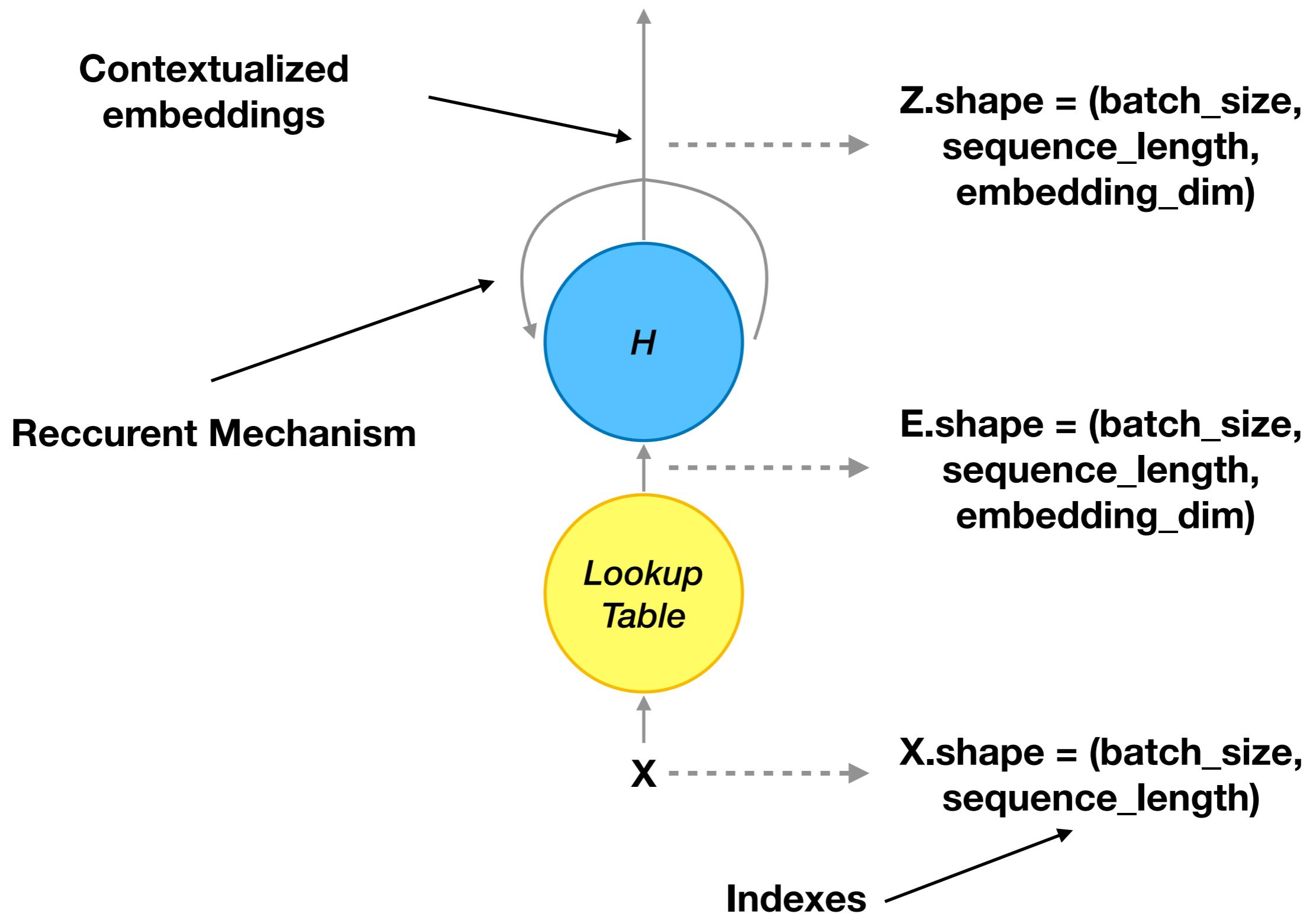
Recurrent Neural Network



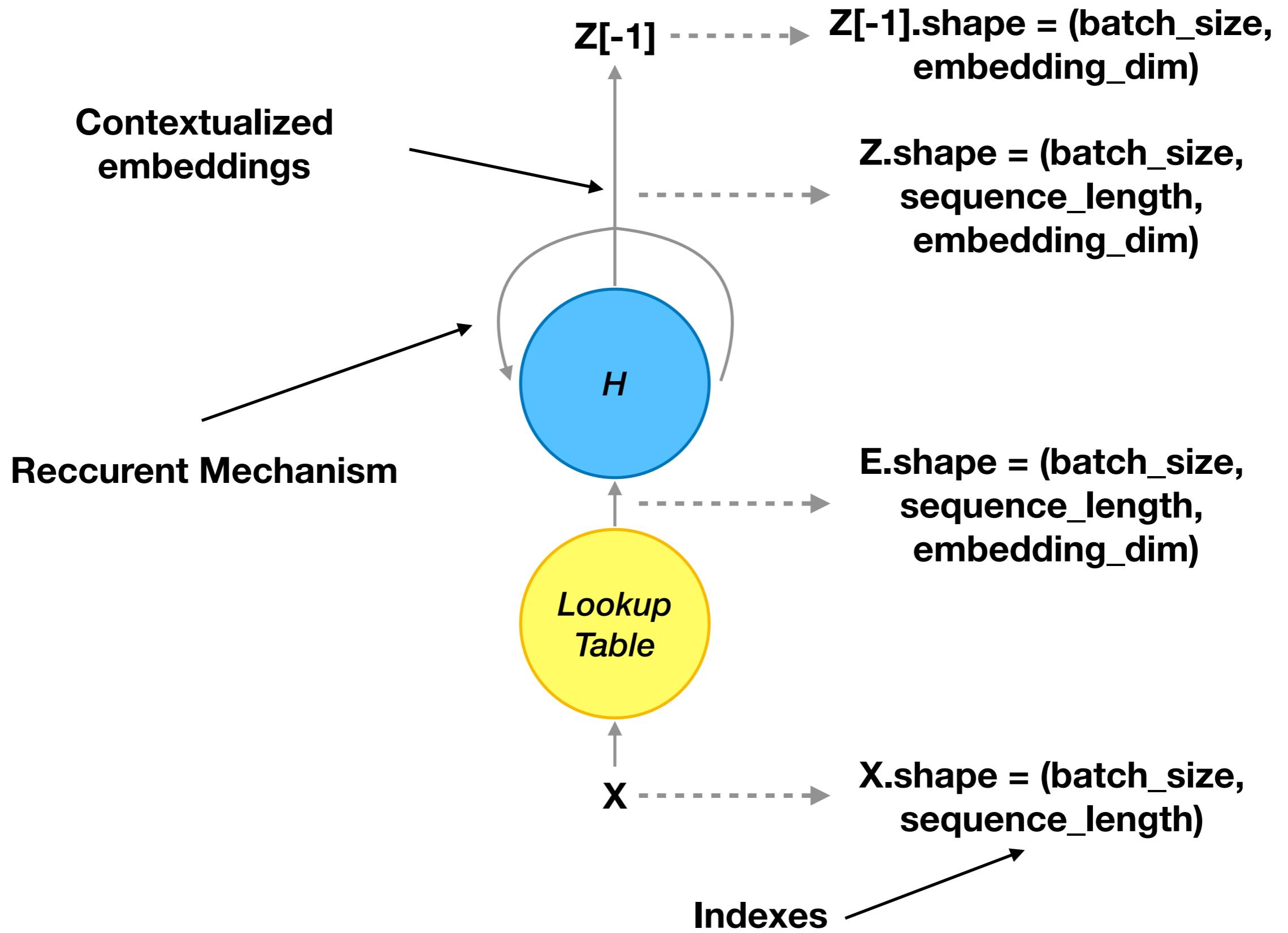
Recurrent Neural Network



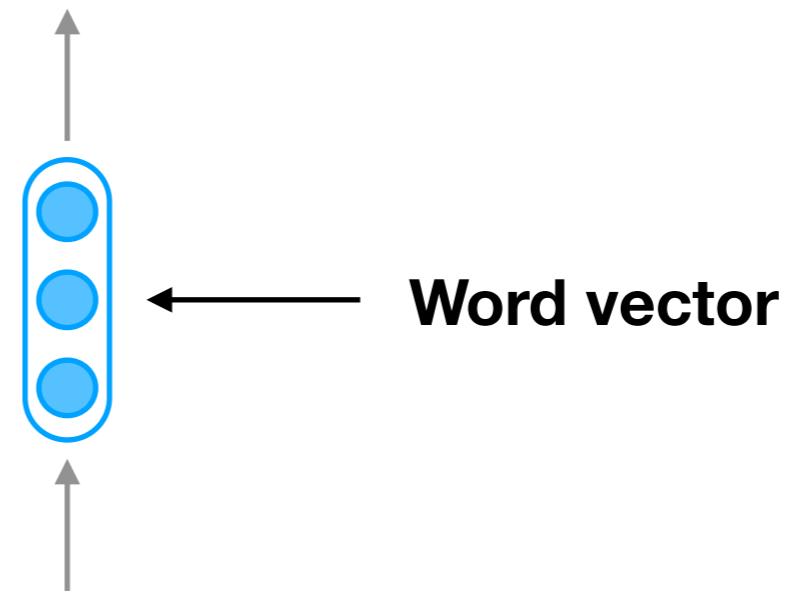
Recurrent Neural Network



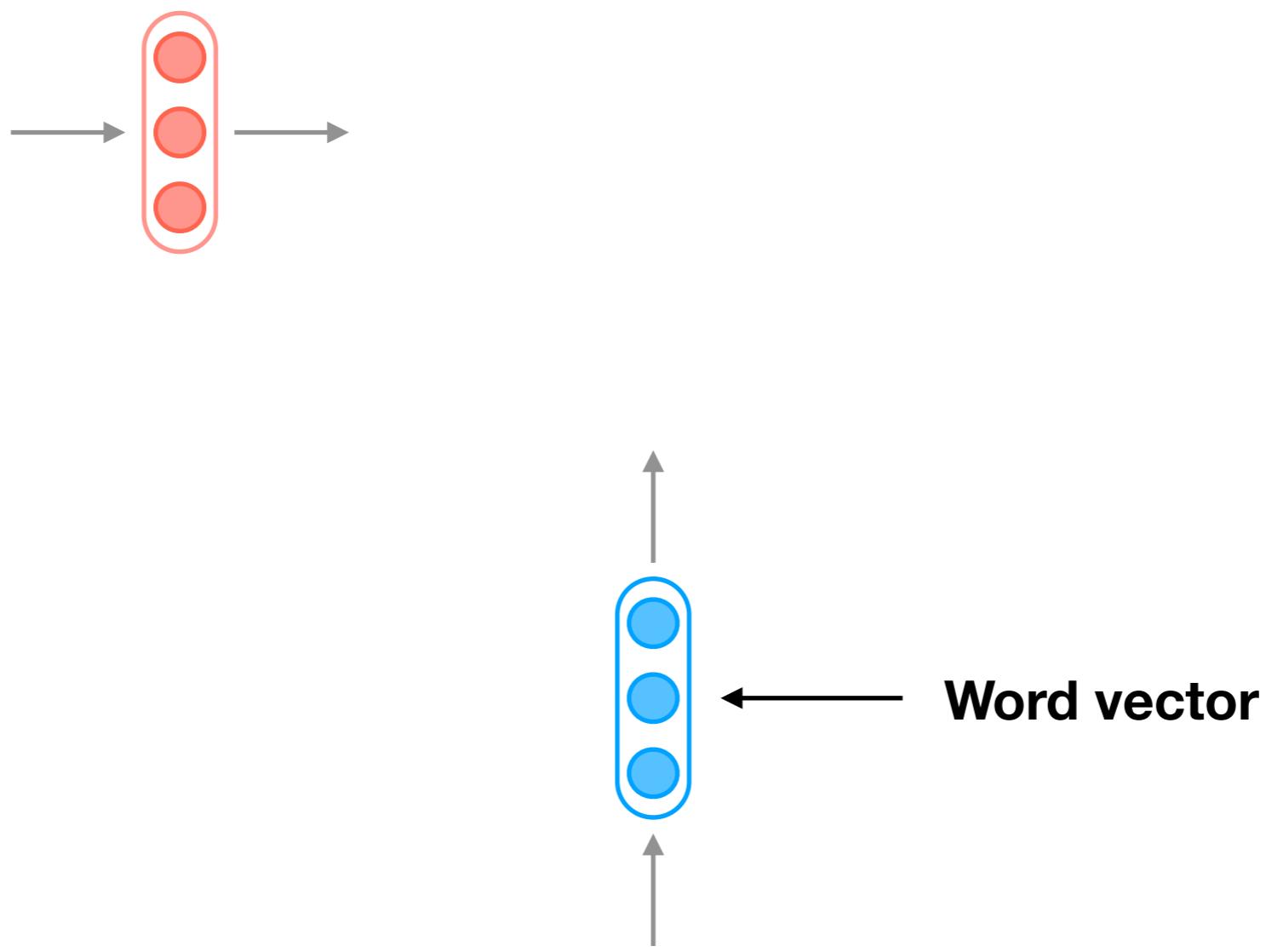
Recurrent Neural Network



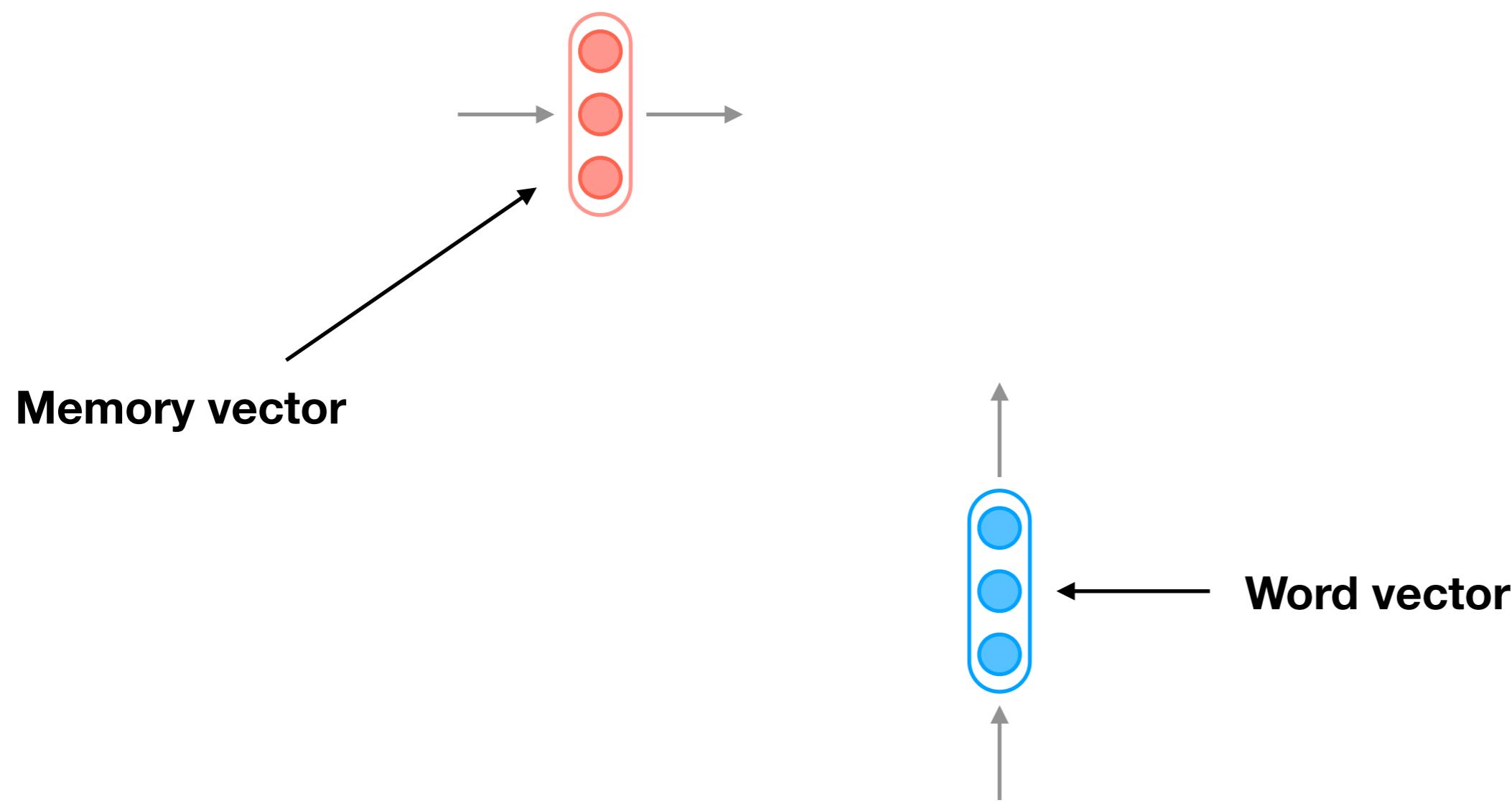
Reccurrent Cell



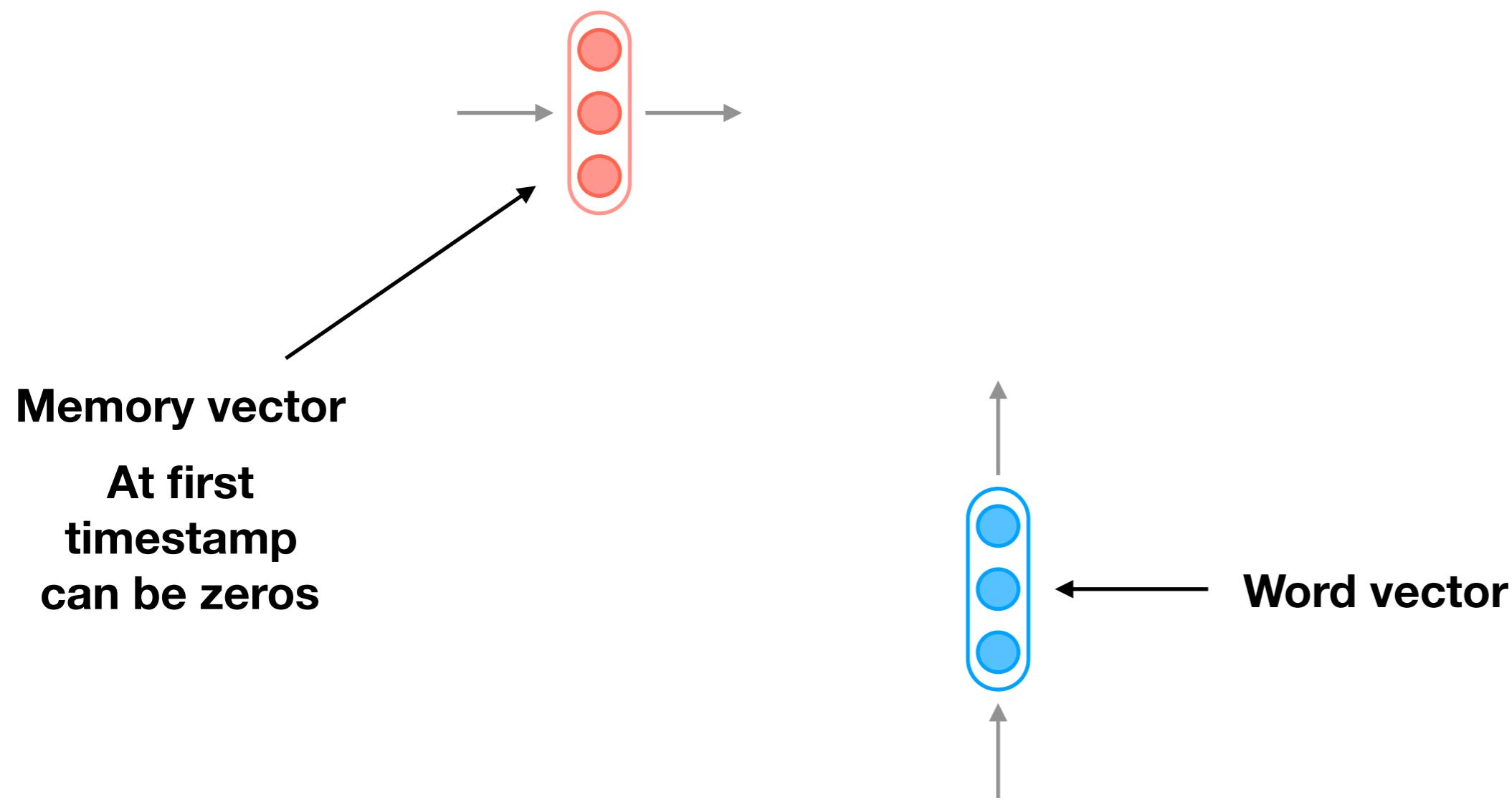
Reccurrent Cell



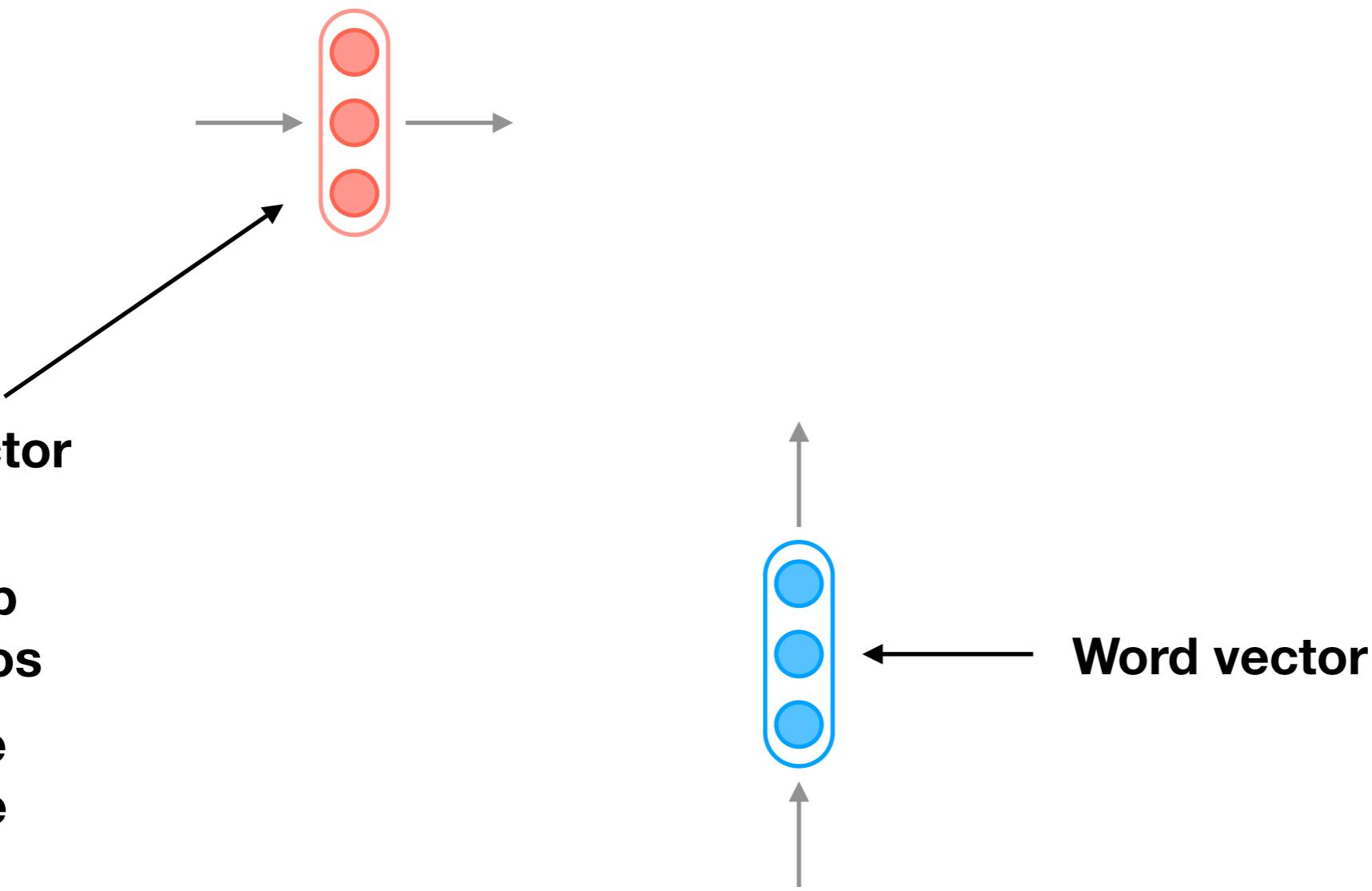
Reccurrent Cell



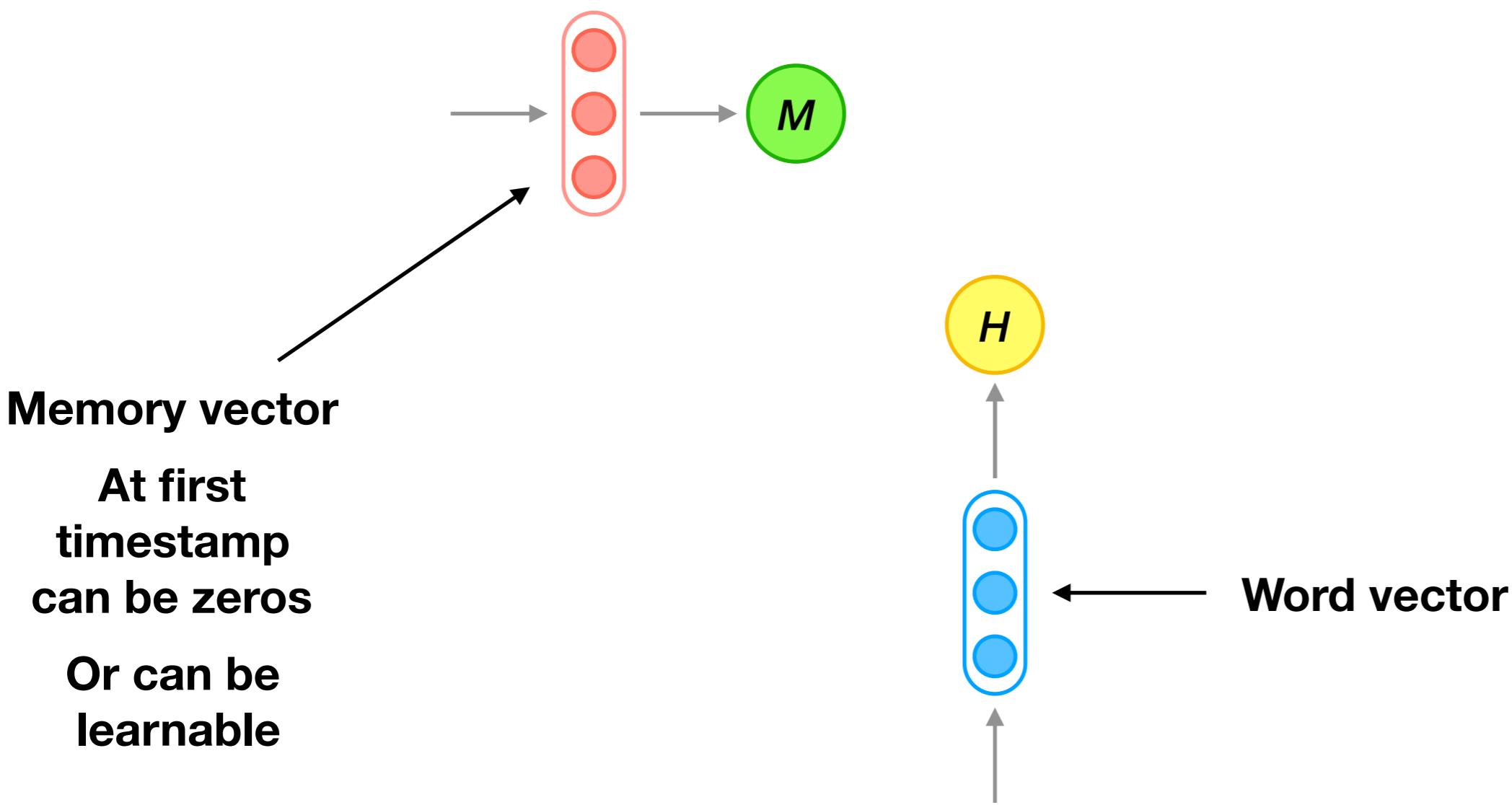
Reccurrent Cell



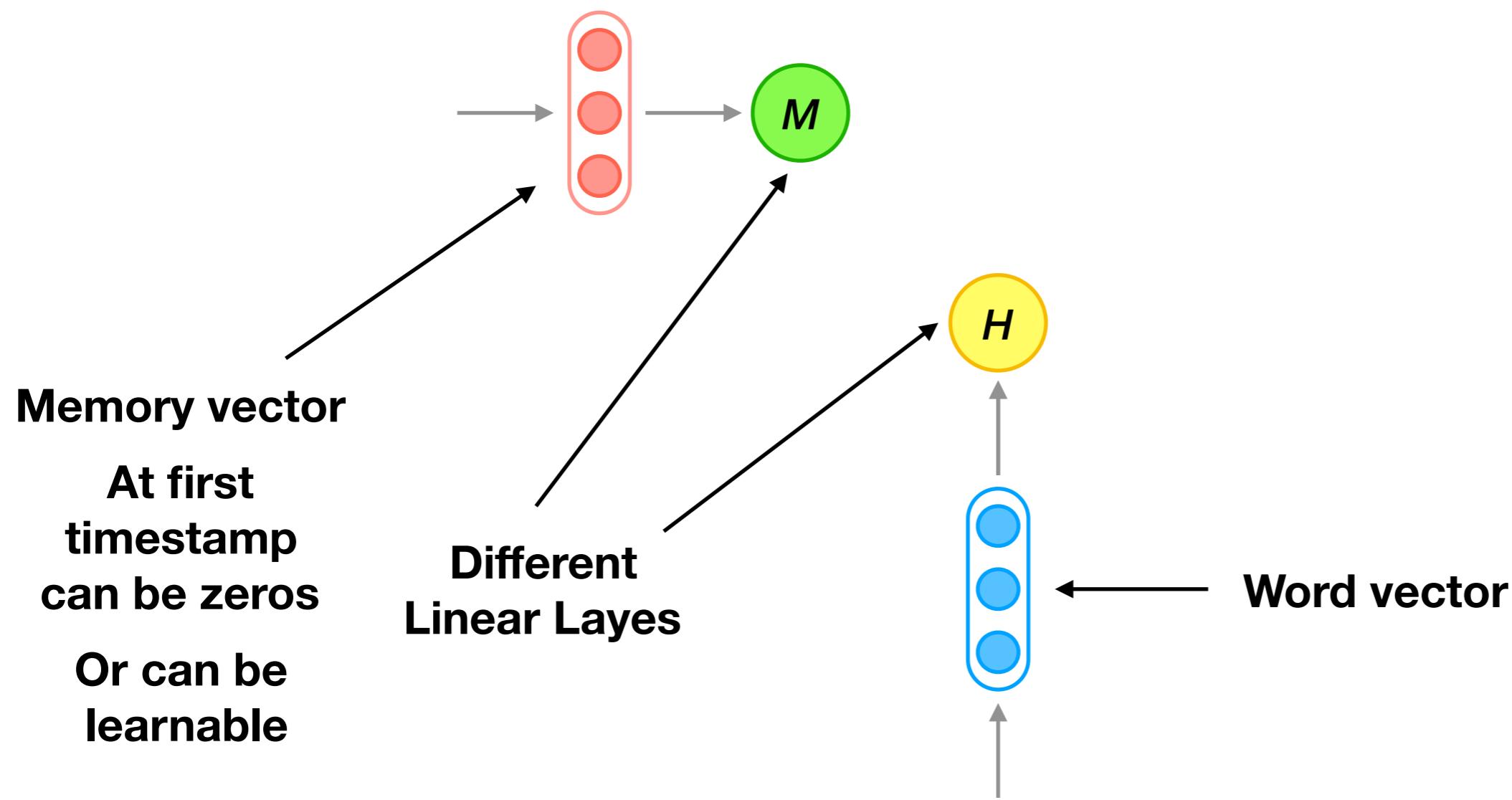
Reccurrent Cell



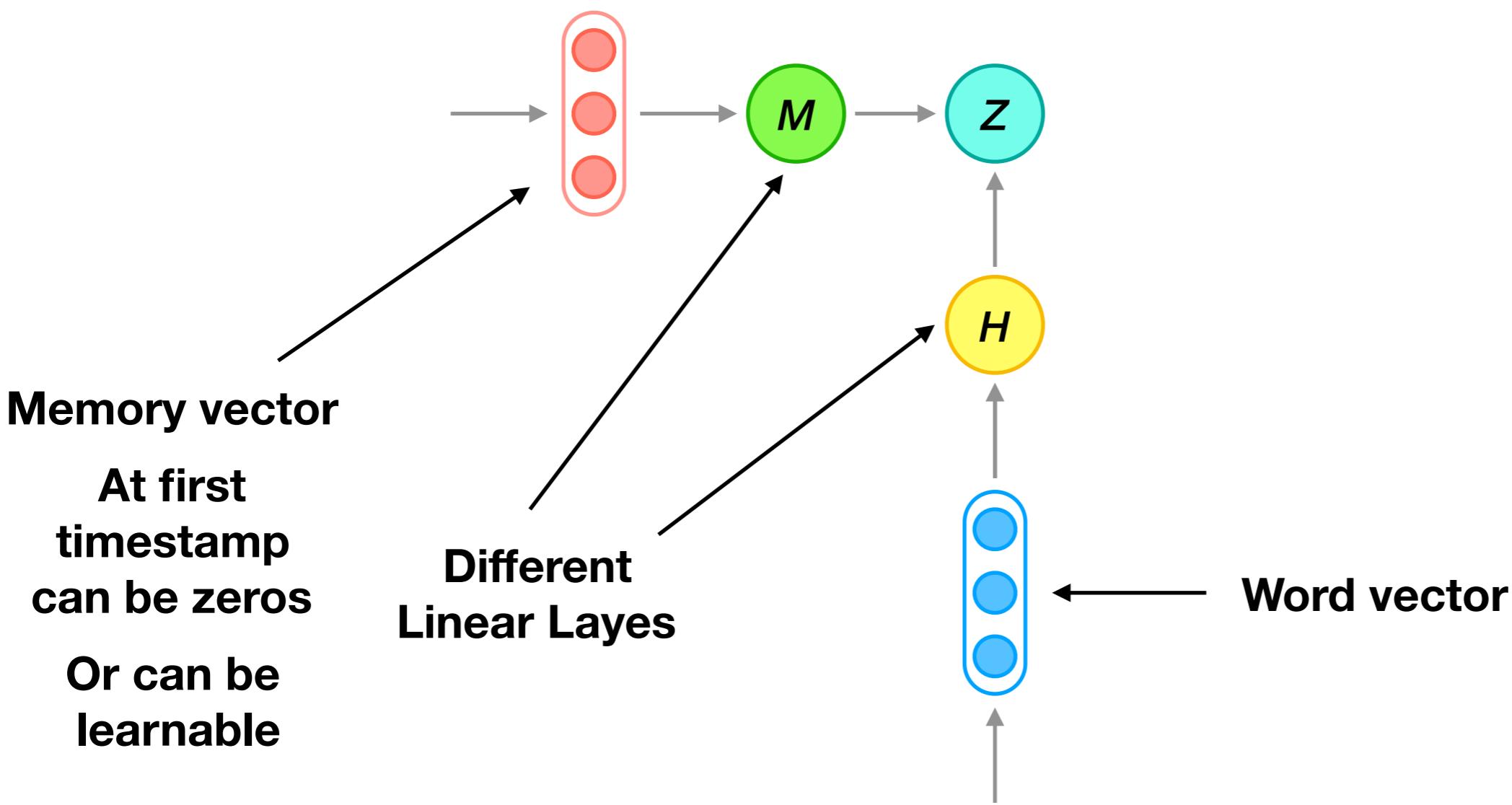
Reccurrent Cell



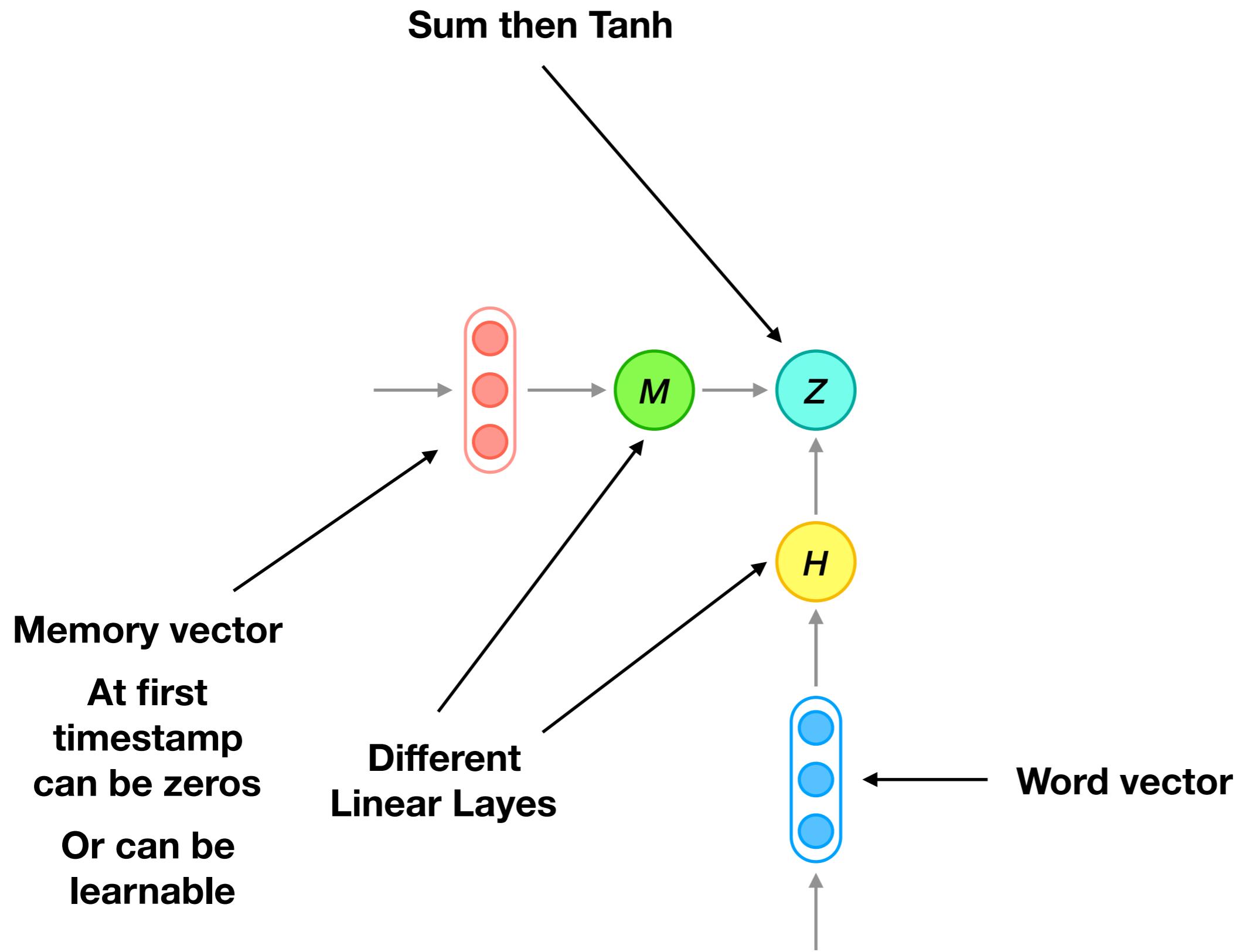
Reccurrent Cell



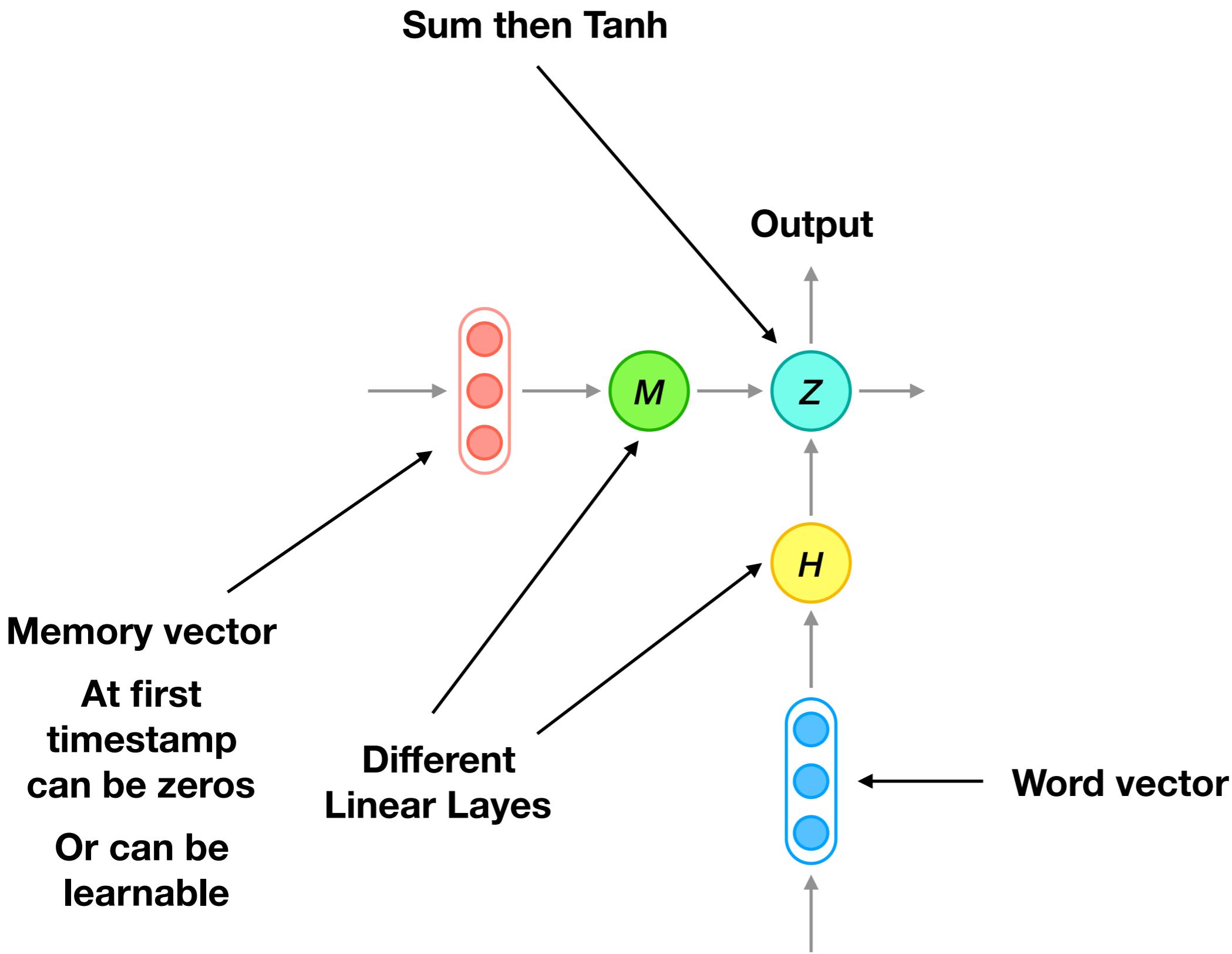
Reccurrent Cell



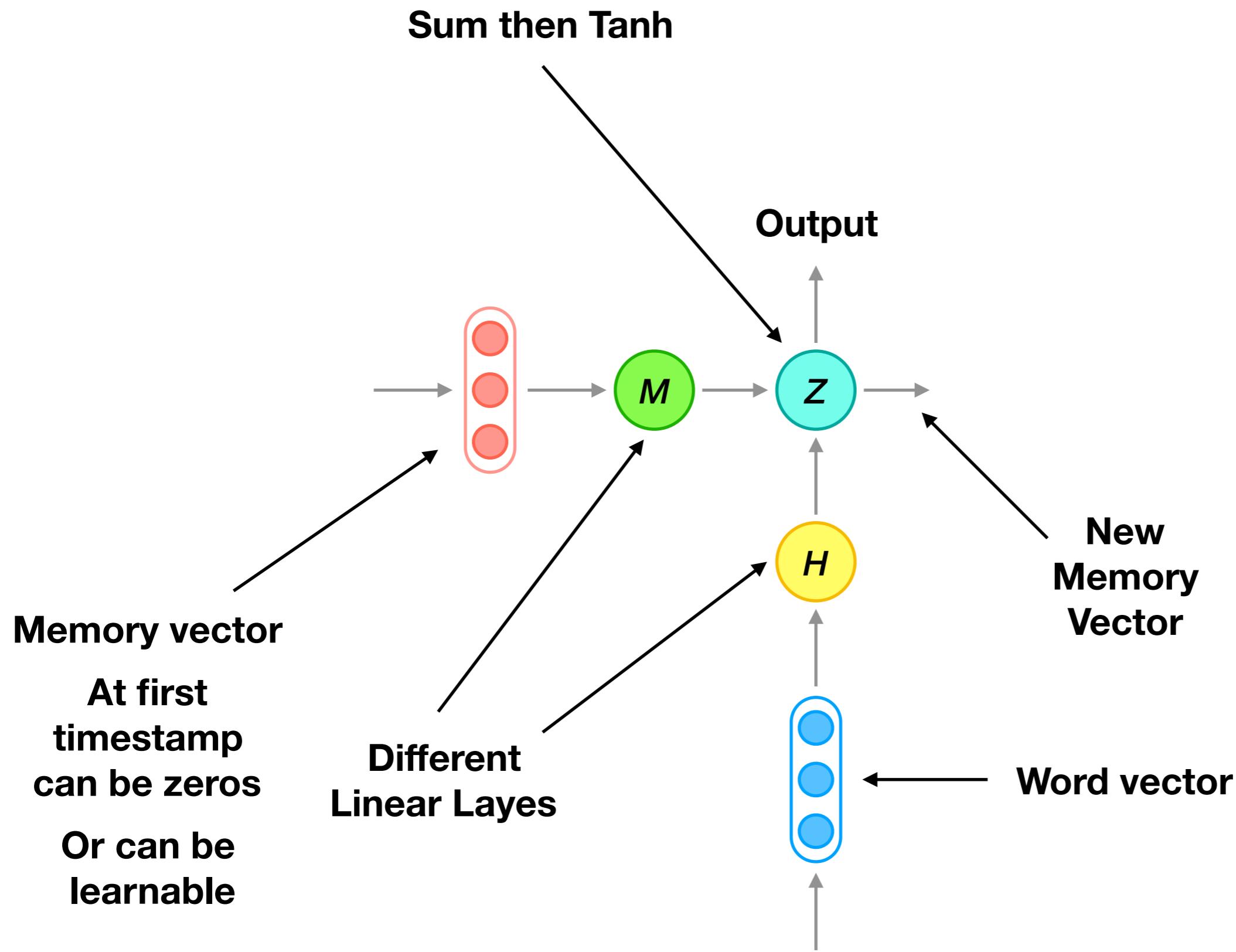
Reccurrent Cell



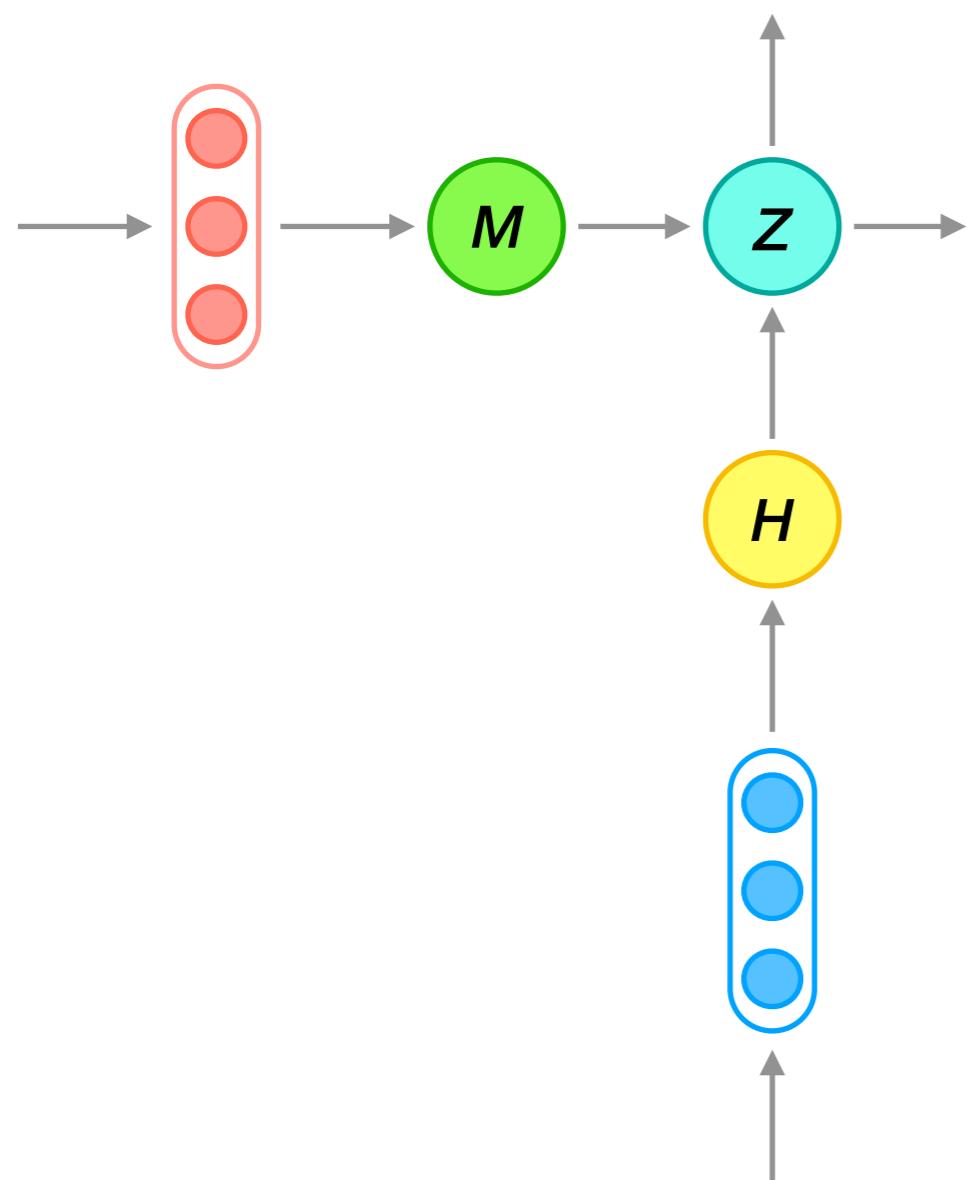
Reccurrent Cell



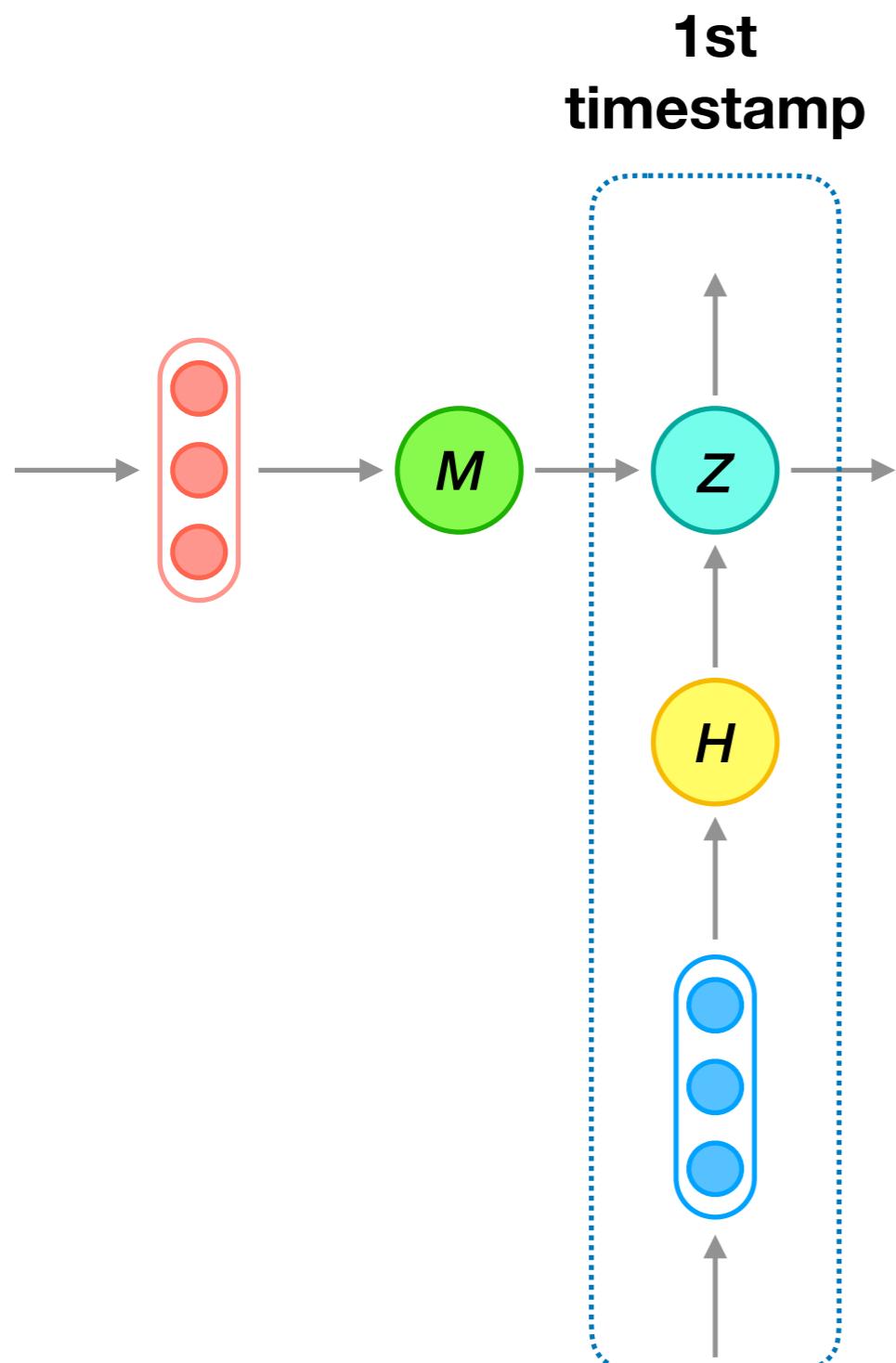
Reccurrent Cell



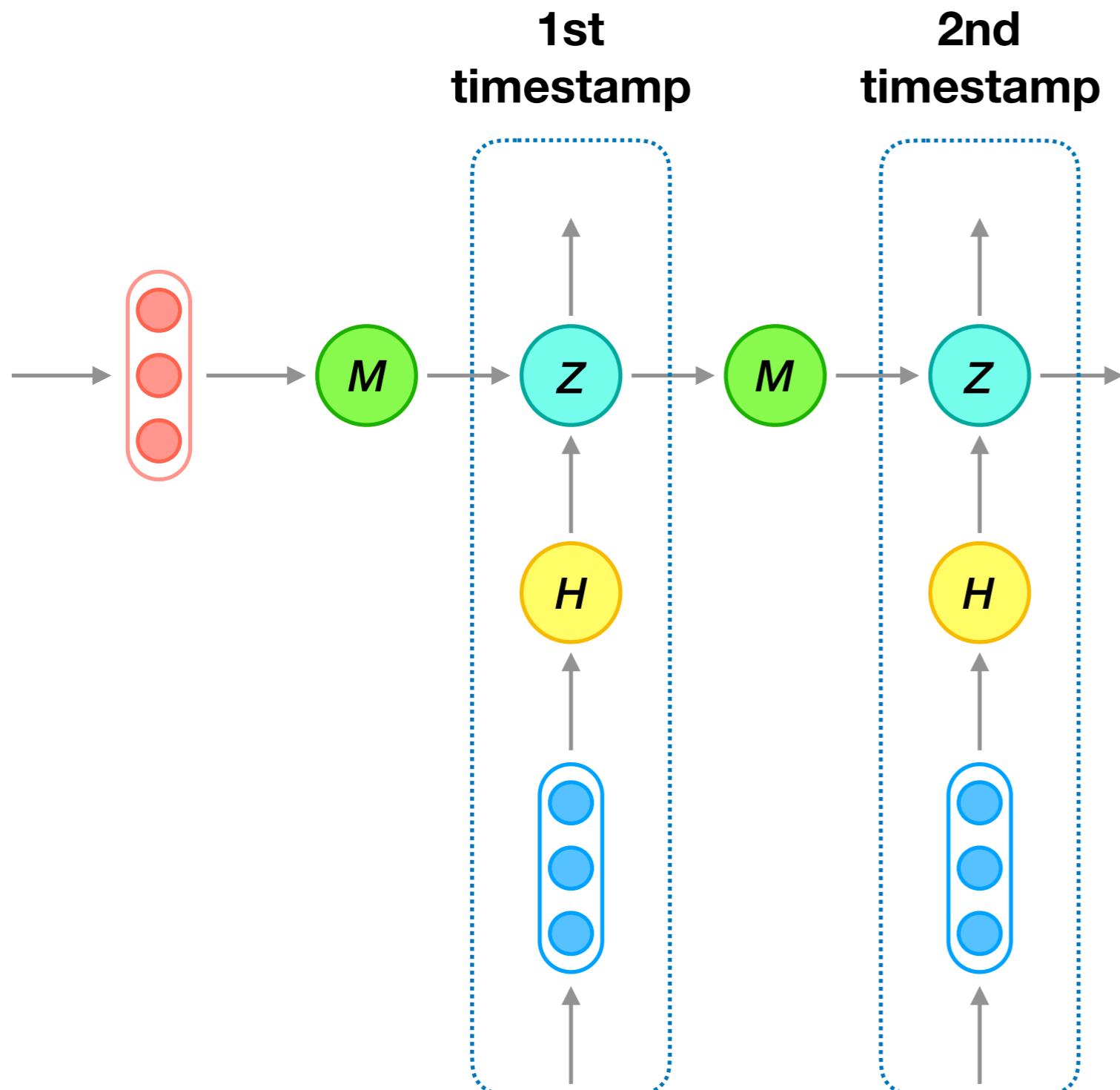
Reccurent Mechanism



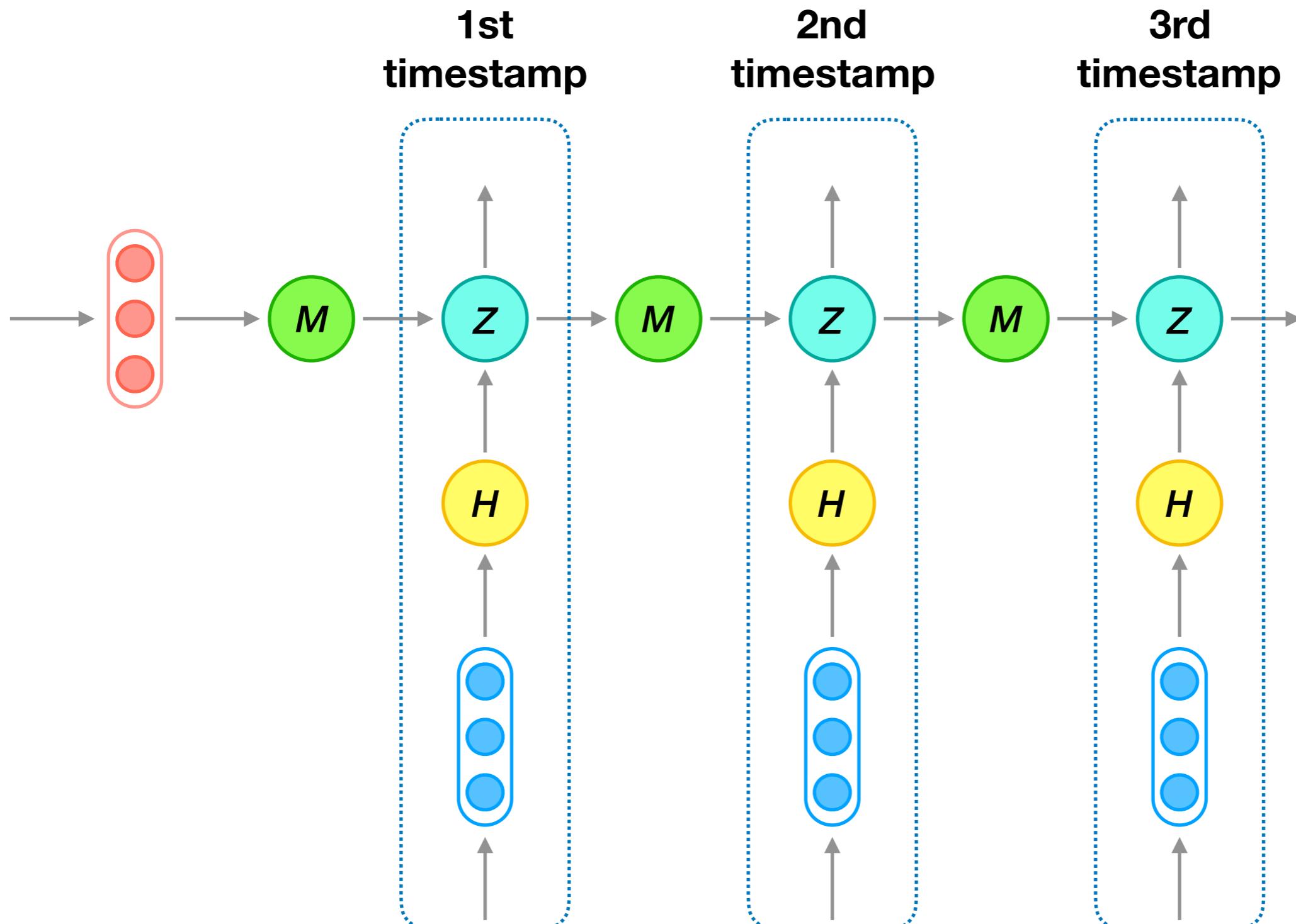
Reccurent Mechanism



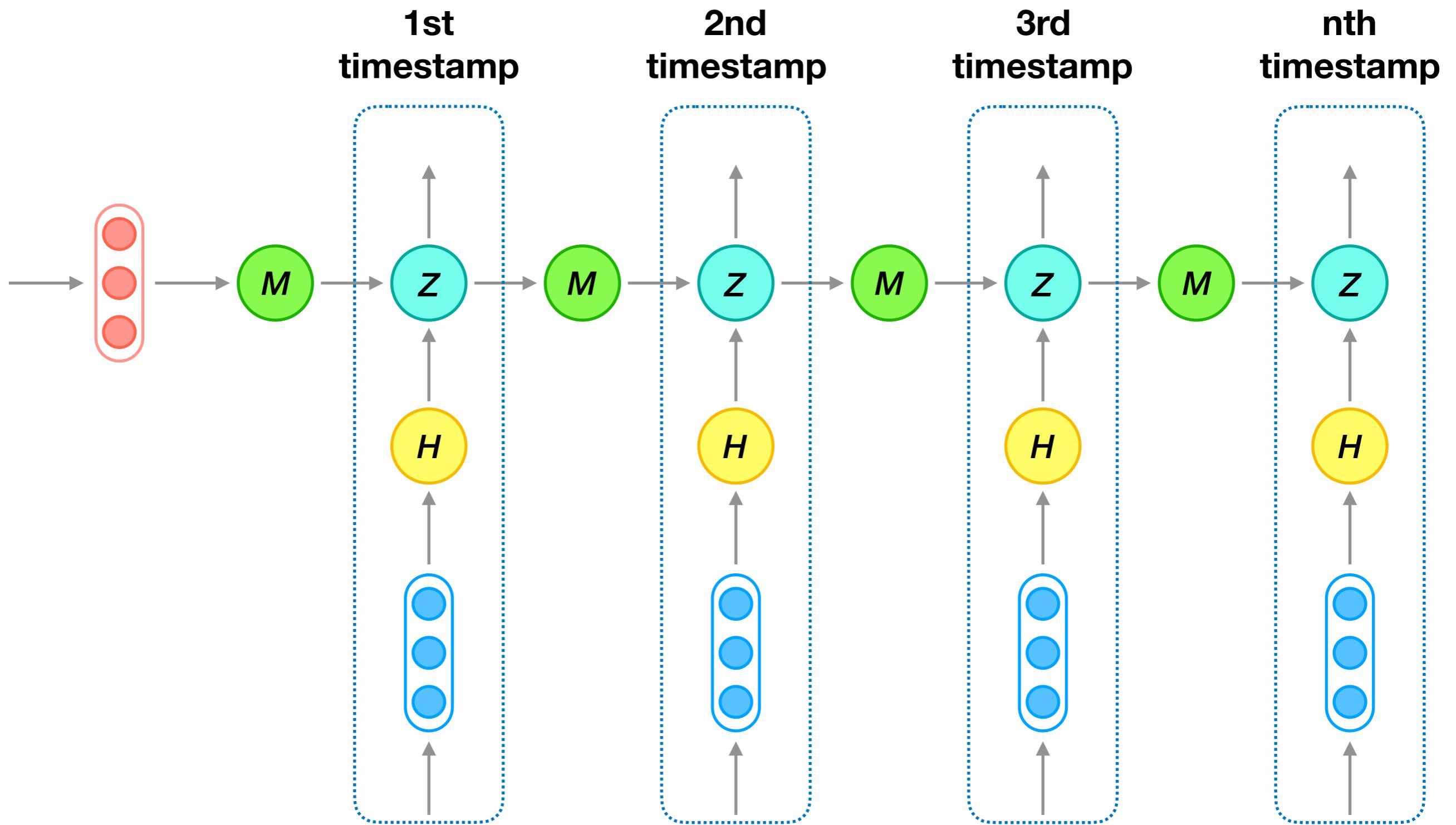
Reccurent Mechanism



Reccurent Mechanism

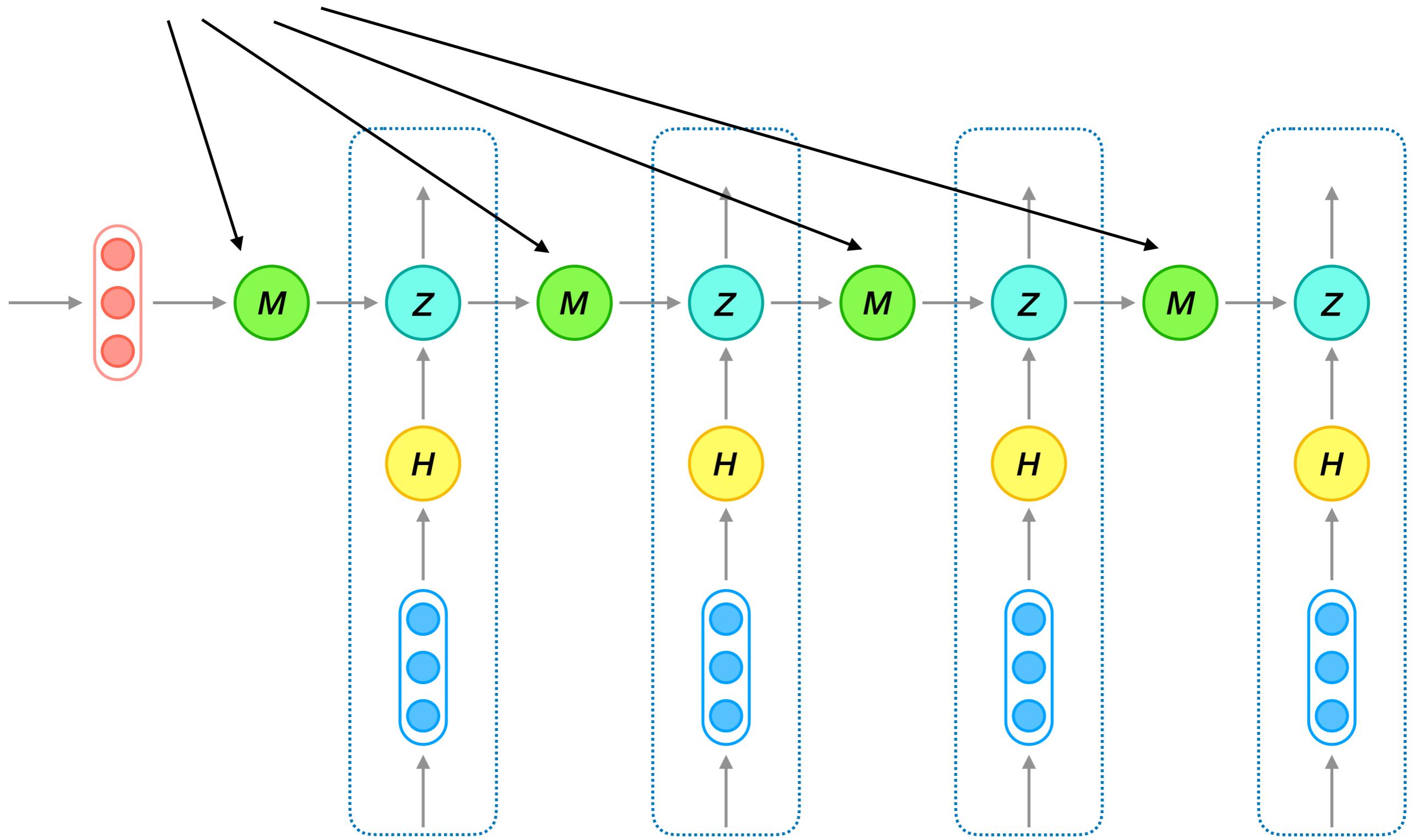


Reccurent Mechanism



Reccurent Mechanism

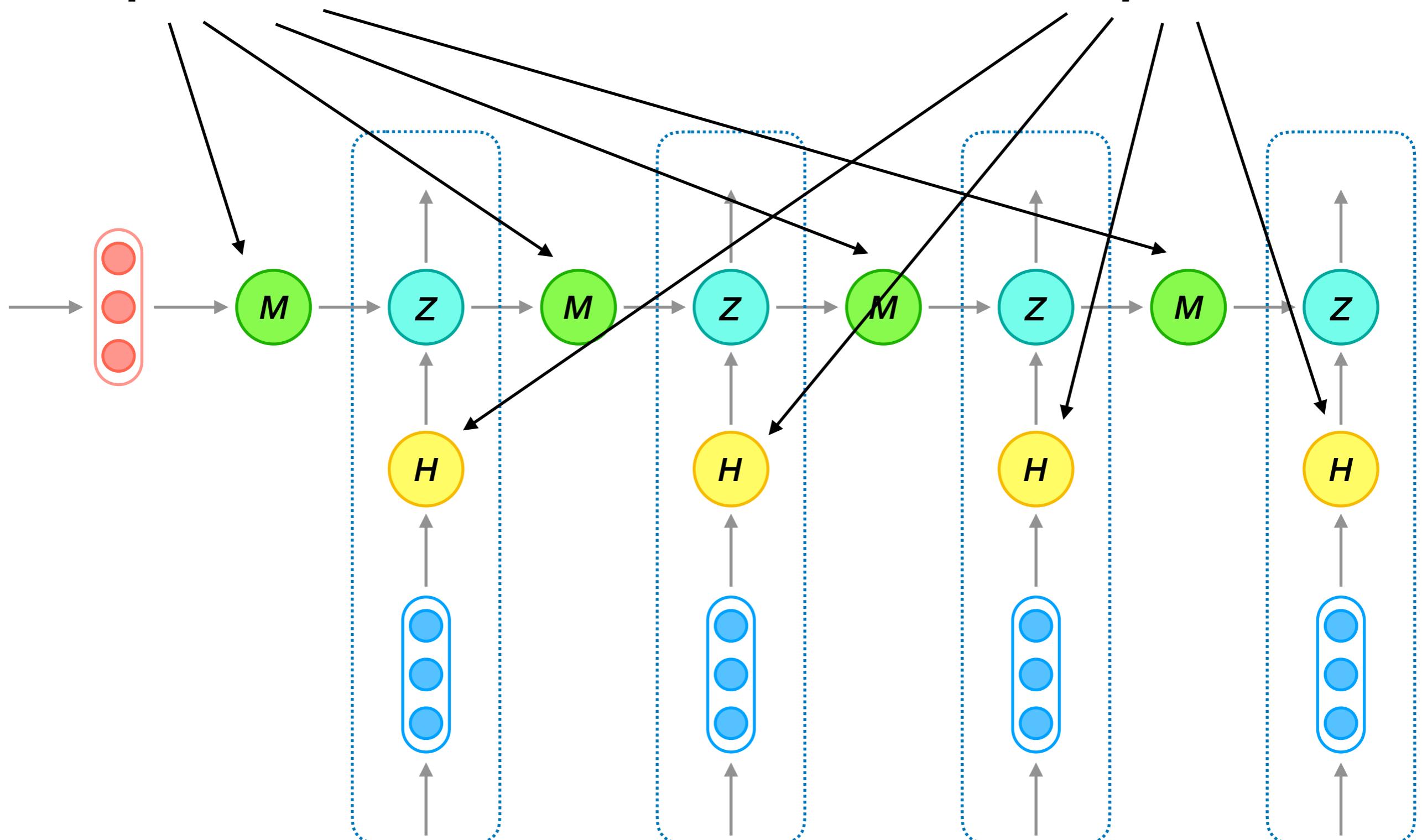
Same parameters!



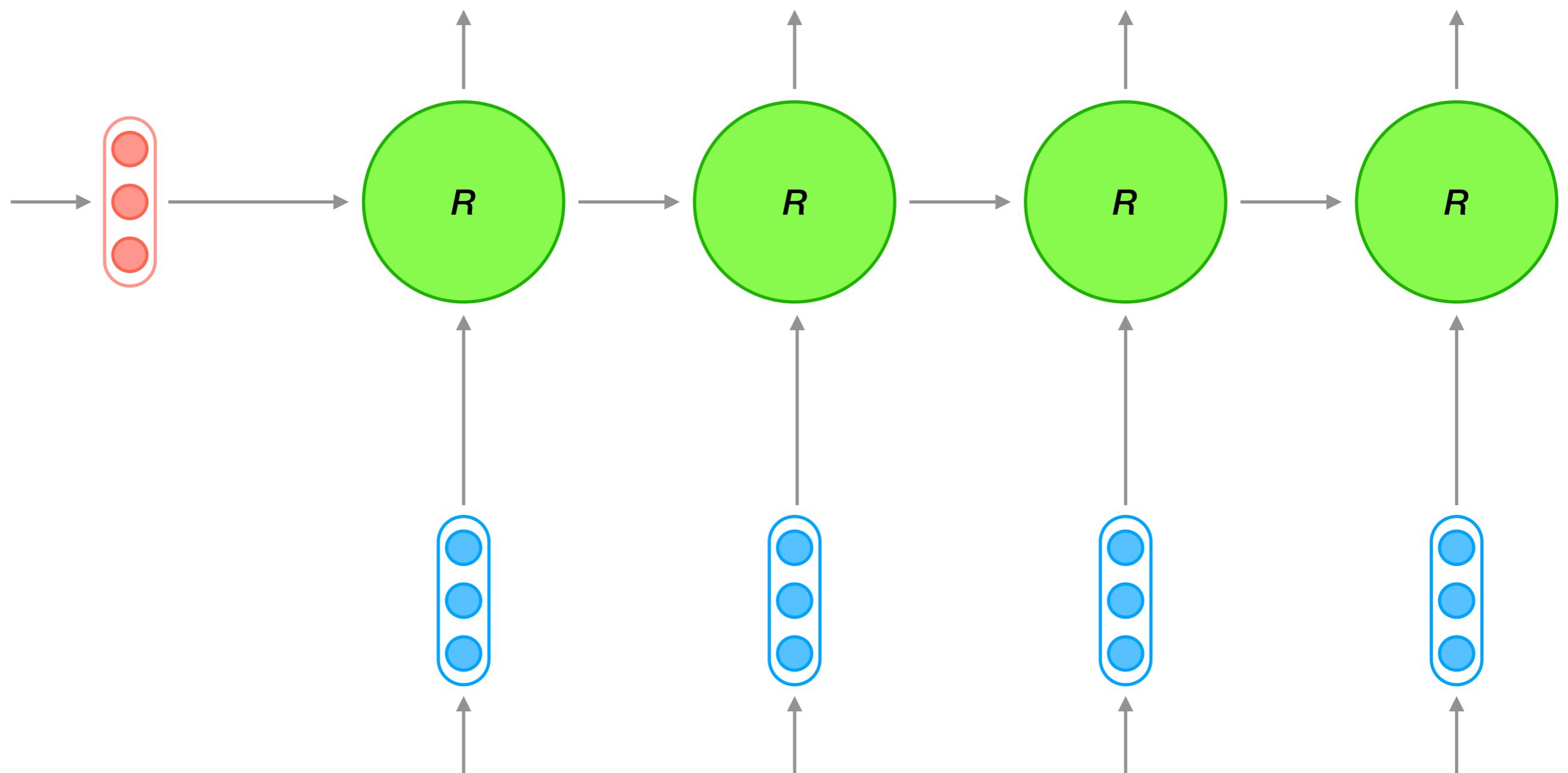
Recurrent Mechanism

Same parameters!

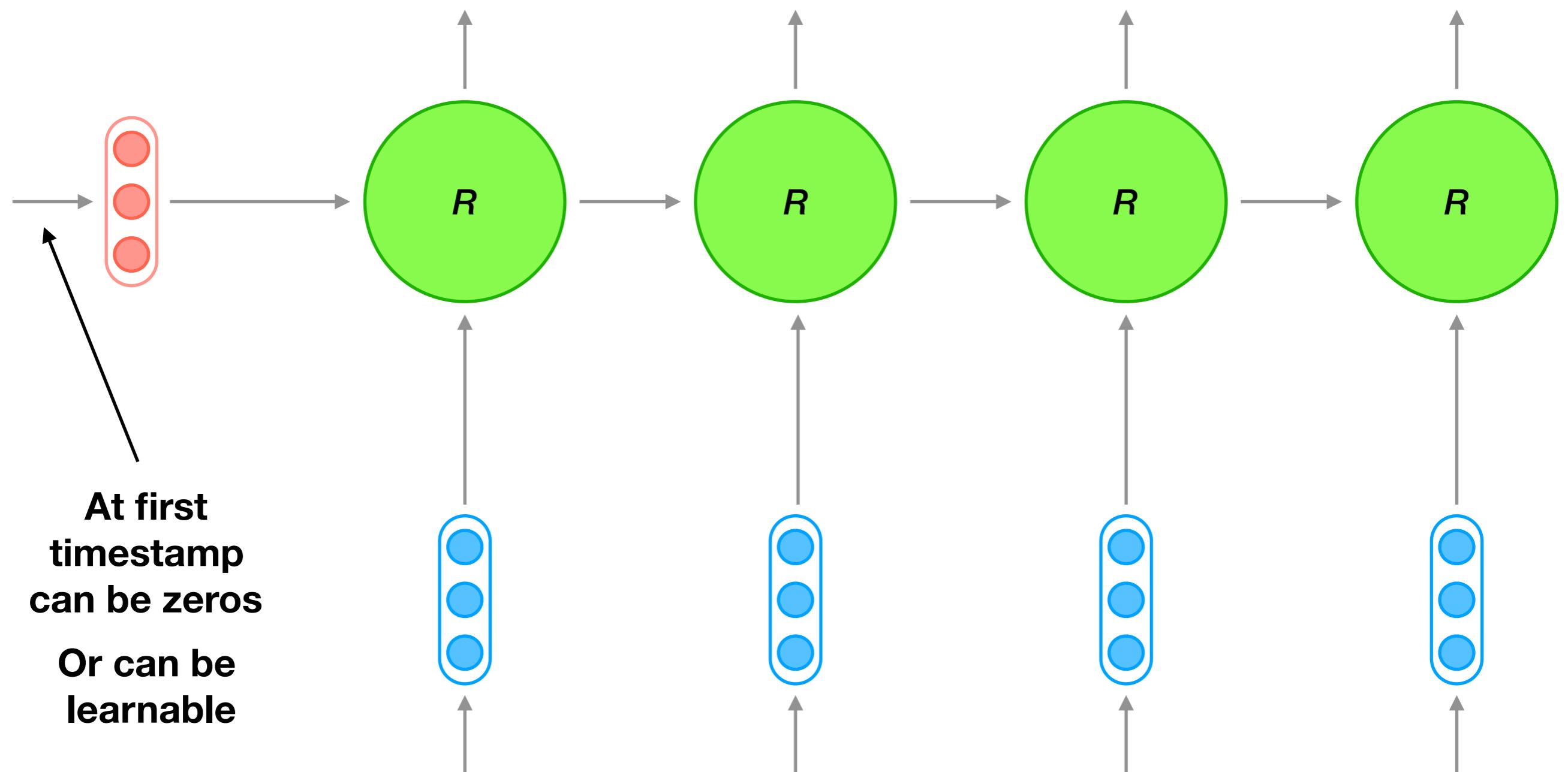
Another
same parameters!



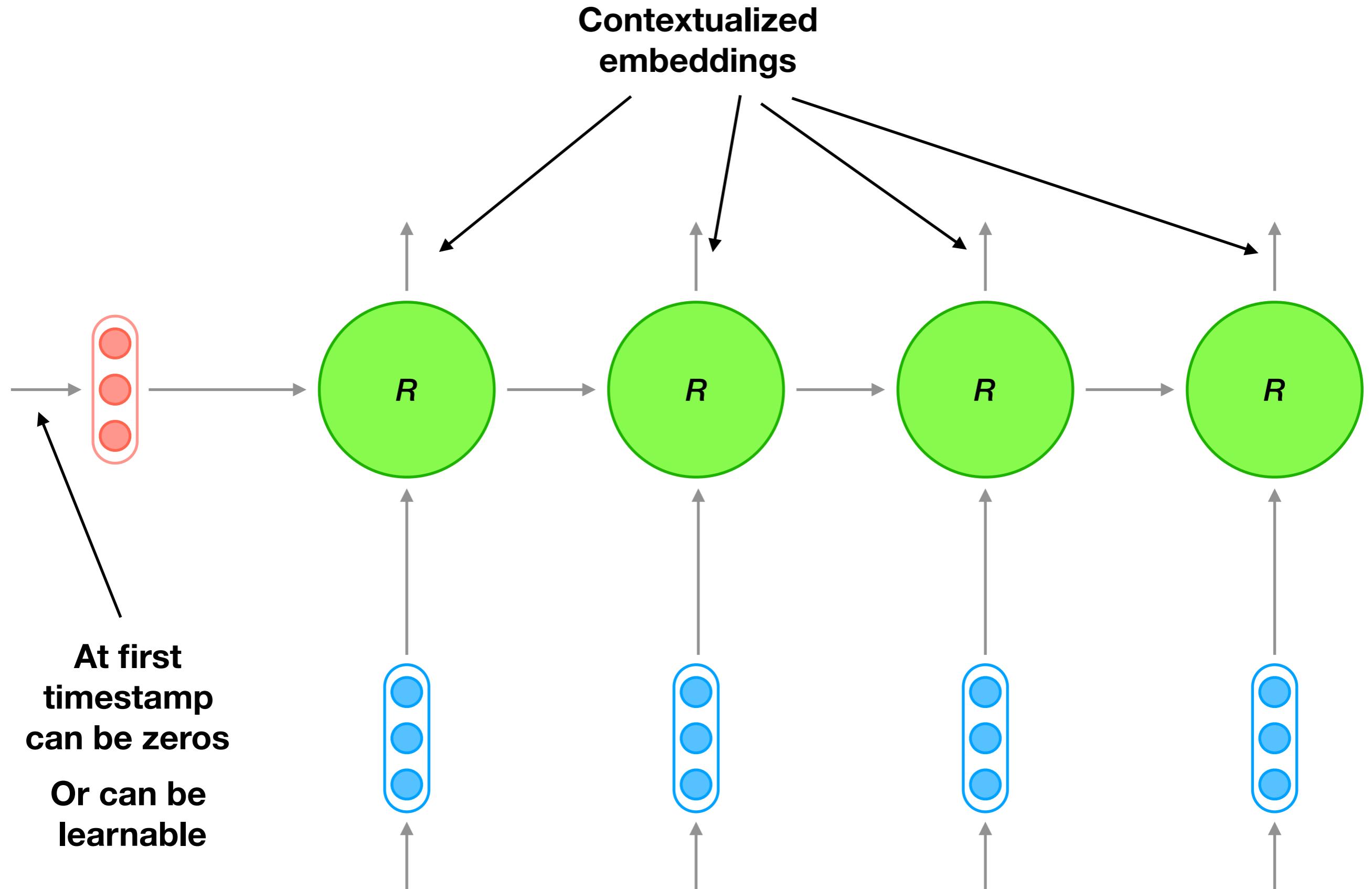
Reccurent Mechanism



Reccurent Mechanism

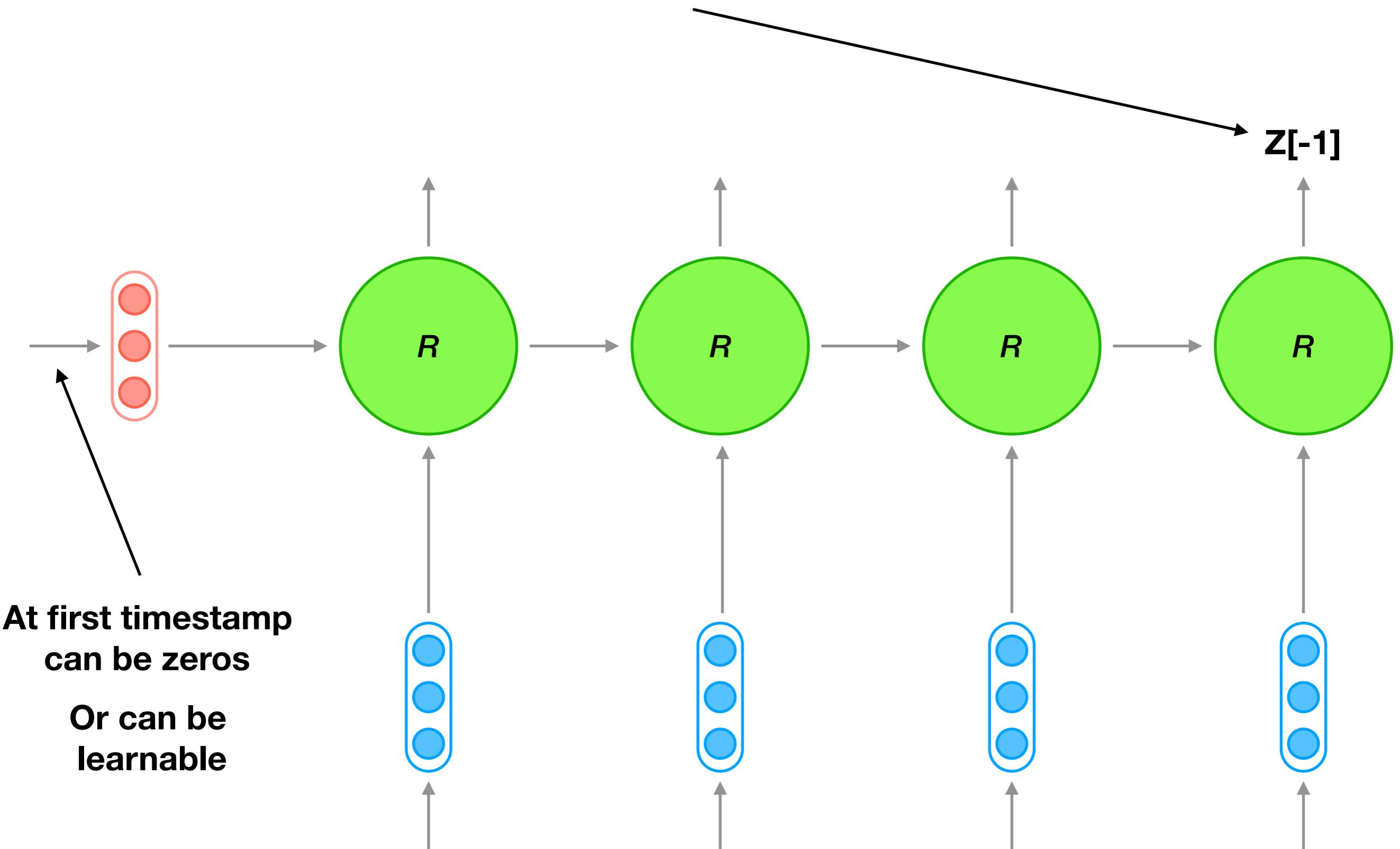


Reccurent Mechanism

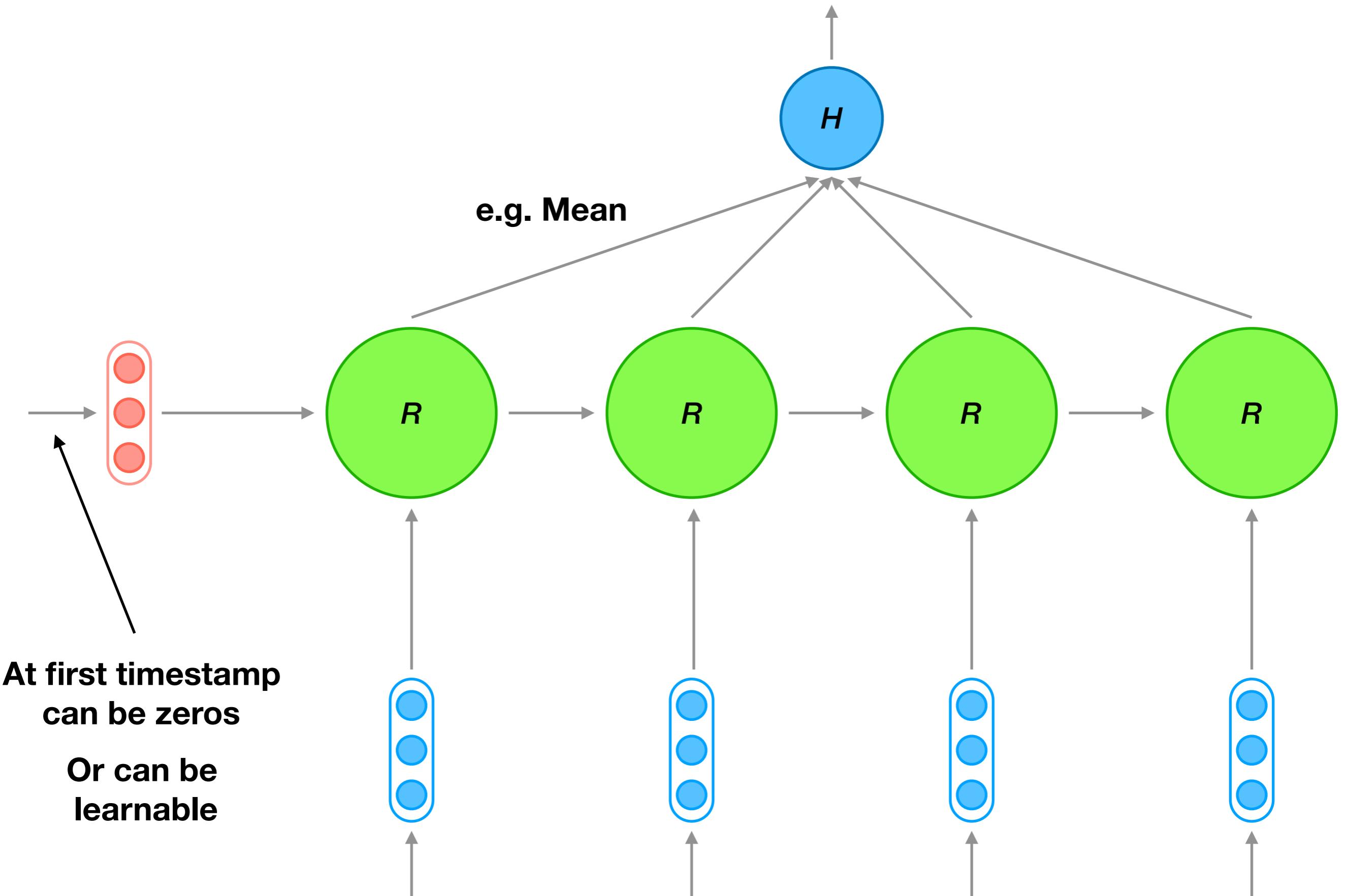


Reccurent Mechanism

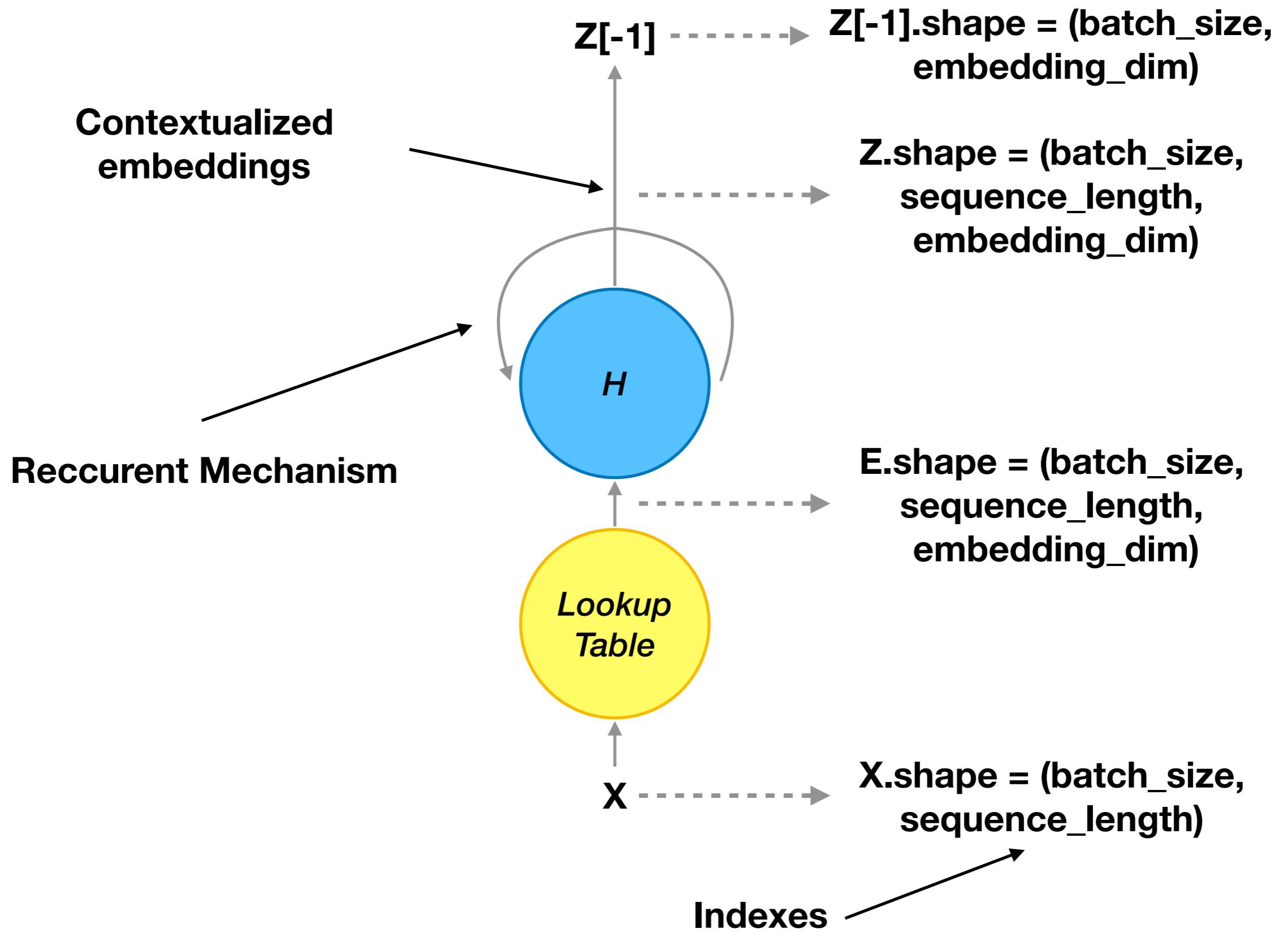
Now we can take last vector from sequence



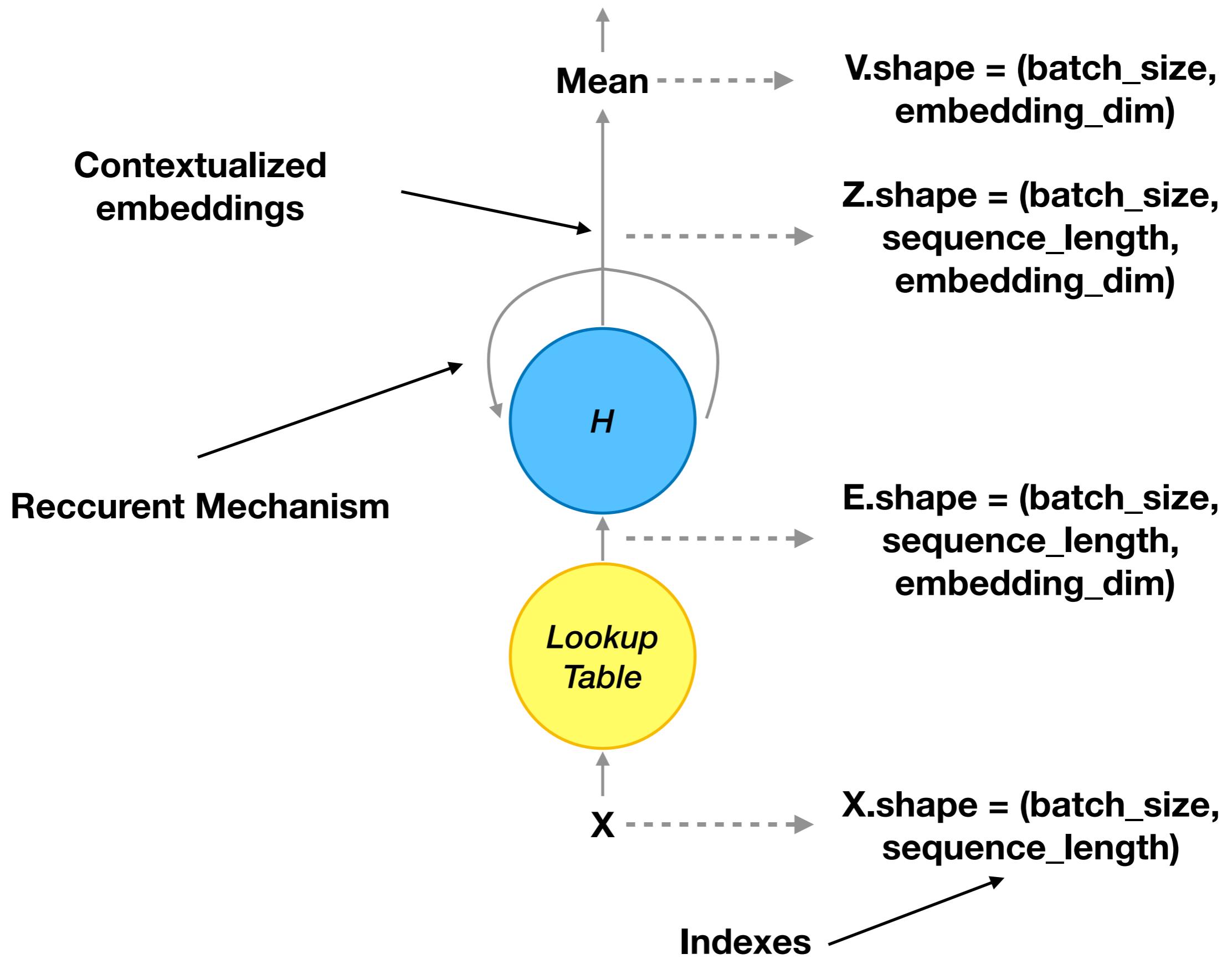
Recurrent Mechanism



Recurrent Neural Network



Recurrent Neural Network



```
class RNNCell(nn.Module):

    def __init__(self, in_features, hidden_features):
        super().__init__()

        self.in_features = in_features
        self.hidden_features = hidden_features

        self.input_linear = nn.Linear(in_features=in_features, out_features=hidden_features)
        self.memory_linear = nn.Linear(in_features=hidden_features, out_features=hidden_features)

    def init_memory(self, batch_size):
        return torch.zeros((batch_size, self.hidden_features))

    def forward(self, x, memory=None):
        if memory is None:
            memory = self.init_memory(batch_size=x.shape[0])

        x = torch.tanh(self.input_linear(x) + self.memory_linear(memory))

        memory = x

    return x, memory
```

```
class RNN(nn.Module):

    def __init__(self, rnn_cell, output_last=False):

        super().__init__()

        self.rnn_cell = rnn_cell

        self.output_last = output_last

    def forward(self, x, memory=None):

        if memory is None:
            memory = self.rnn_cell.init_memory(batch_size=x.shape[0])

        hiddens = []

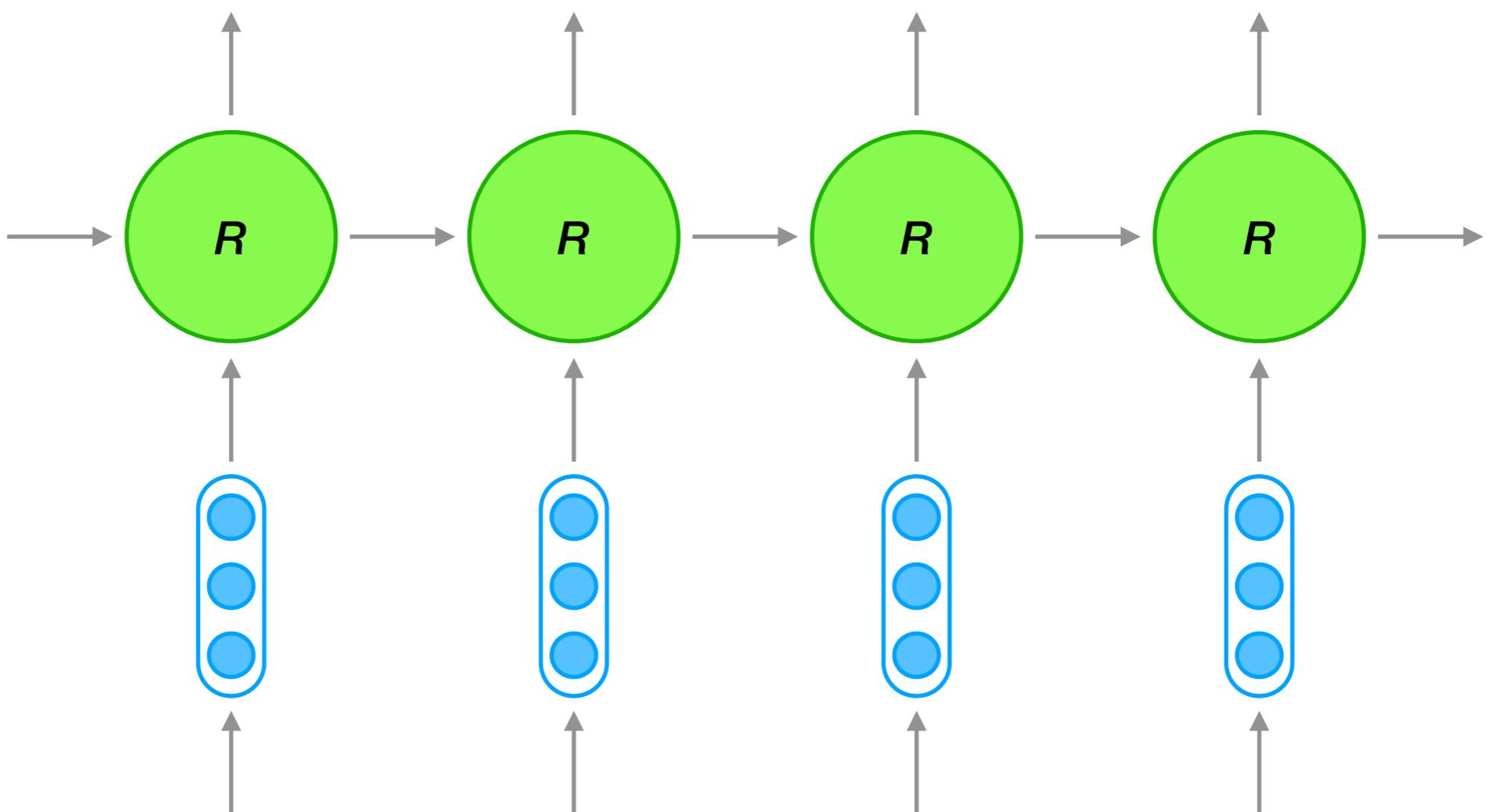
        for timestamp in range(x.size(1)):
            current_timestamp = x[:, timestamp, :]
            current_hidden, memory = self.rnn_cell(current_timestamp, memory)
            hiddens.append(current_hidden.unsqueeze(1))

        if self.output_last:
            return hiddens[-1].squeeze()

        hiddens = torch.cat(hiddens, dim=1)

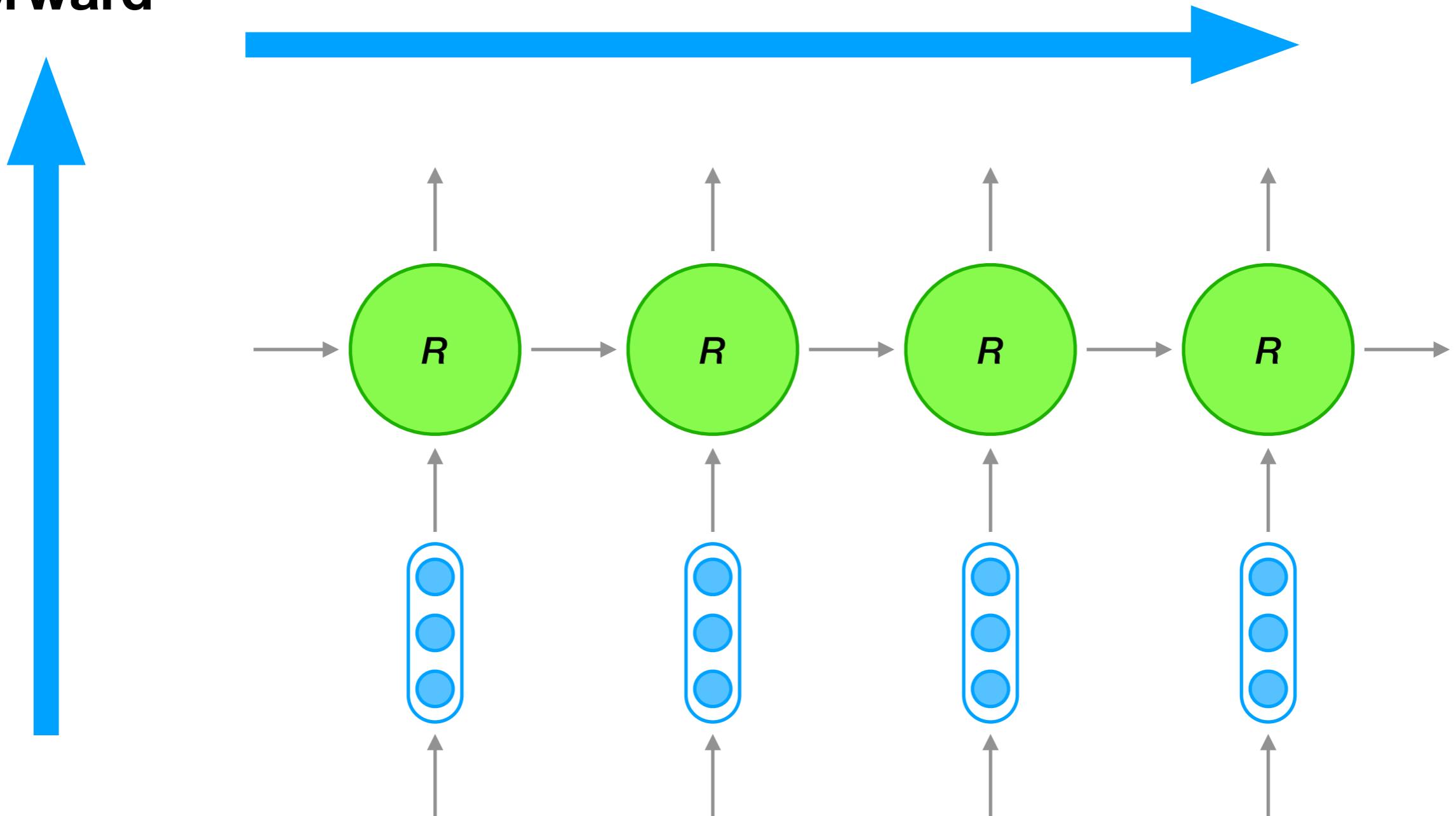
    return hiddens
```

BPTT



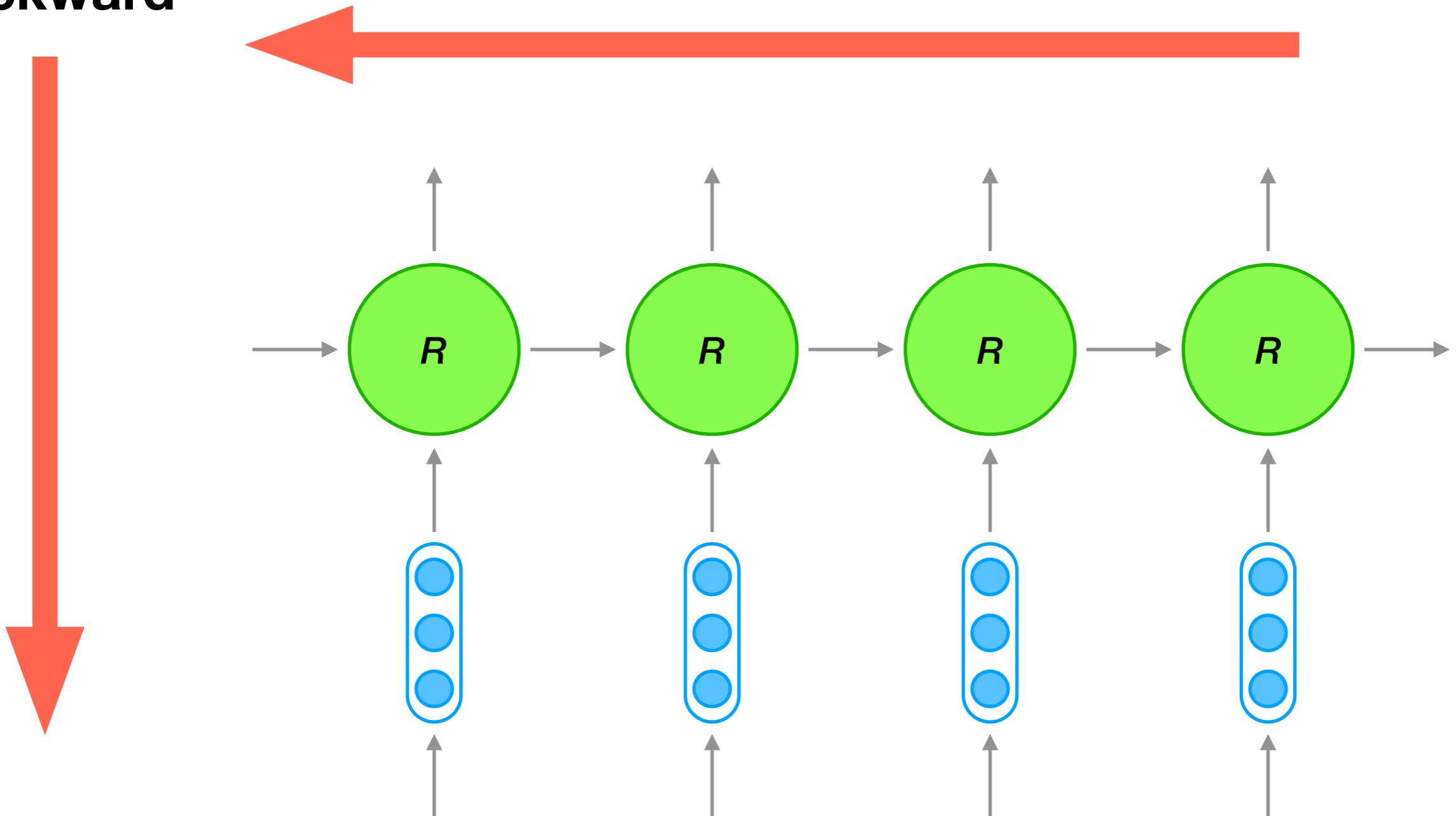
BPTT

Forward

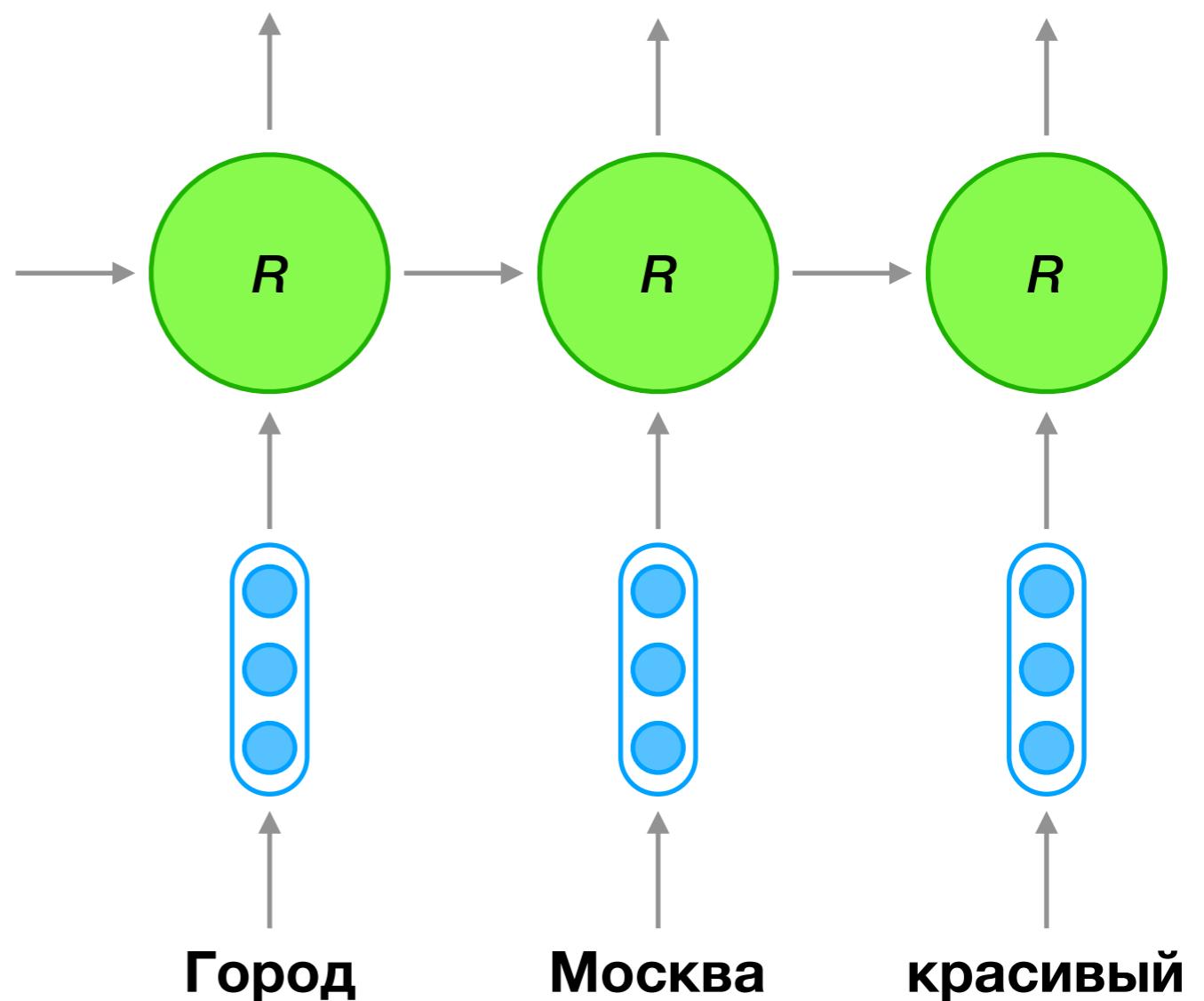


BPTT

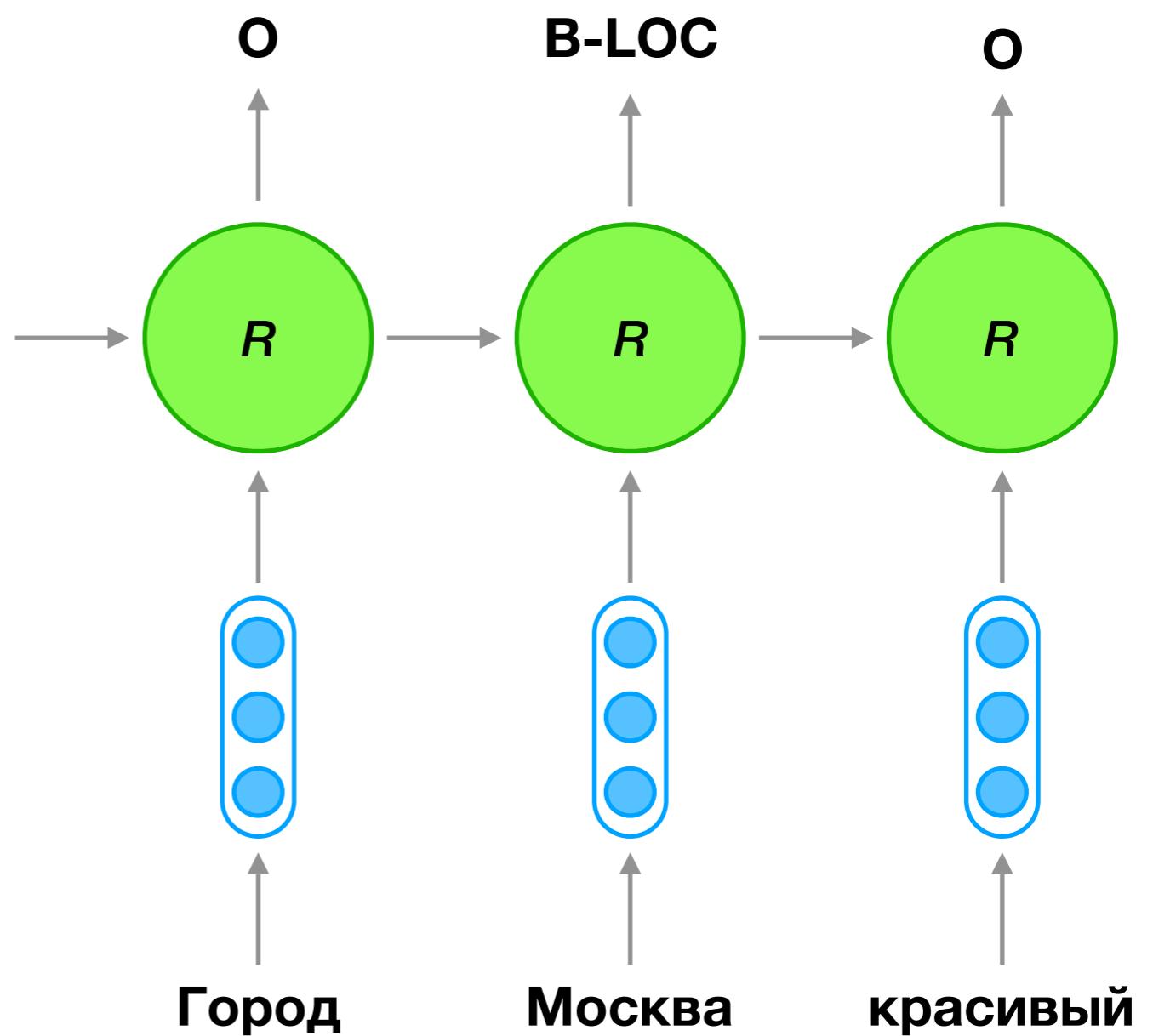
Backward



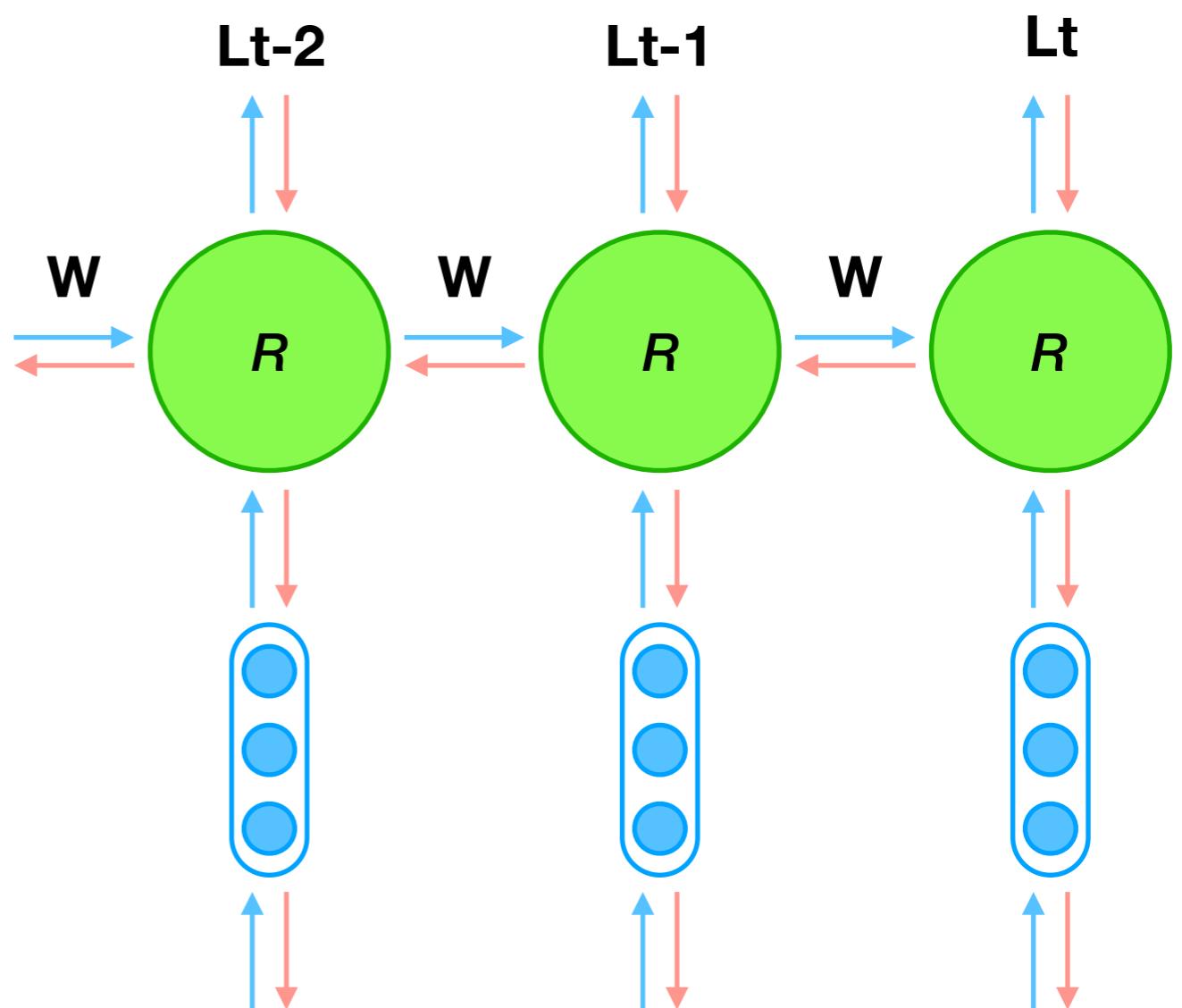
ВРТГ



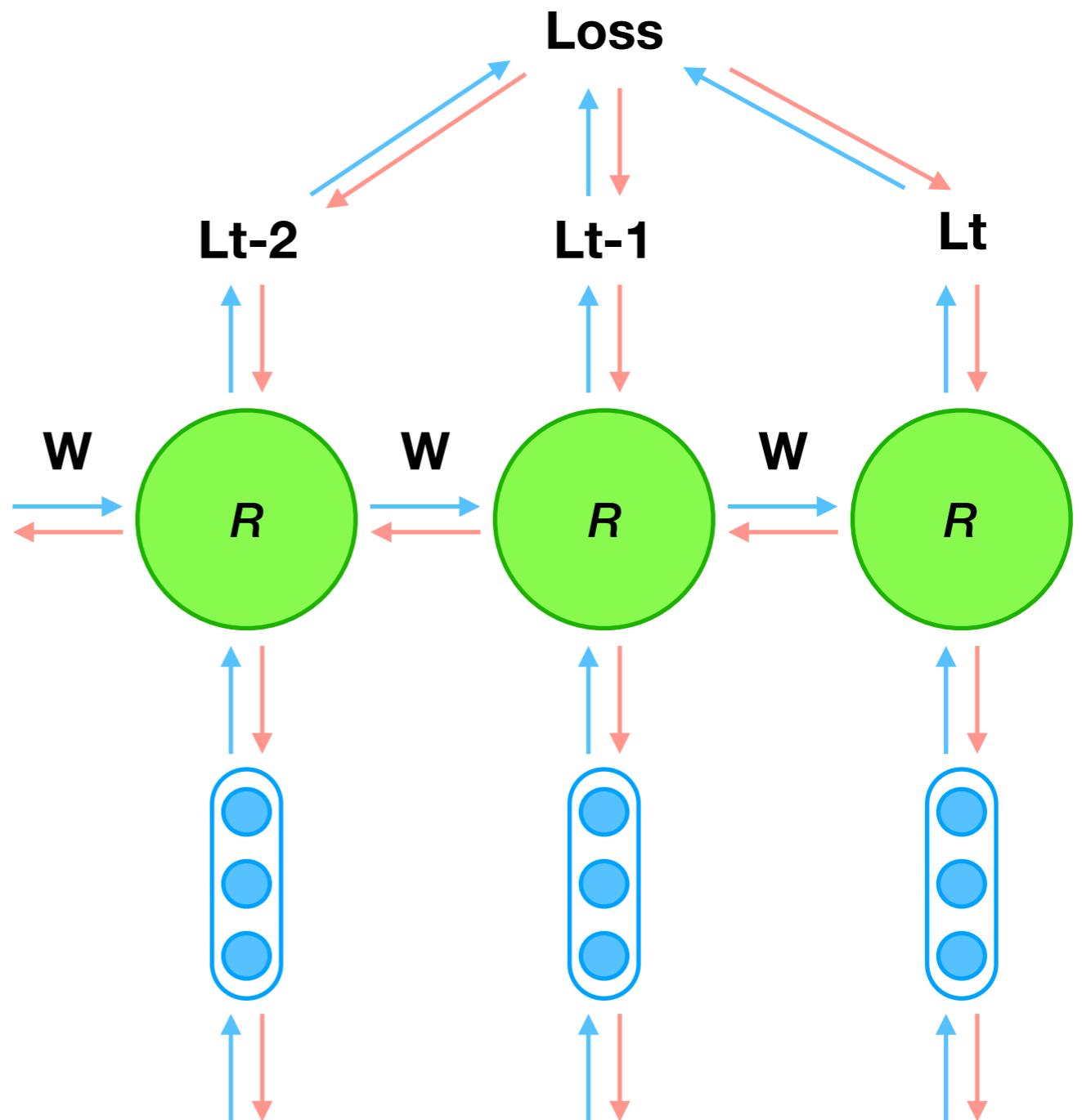
ВРТГ



BPTT

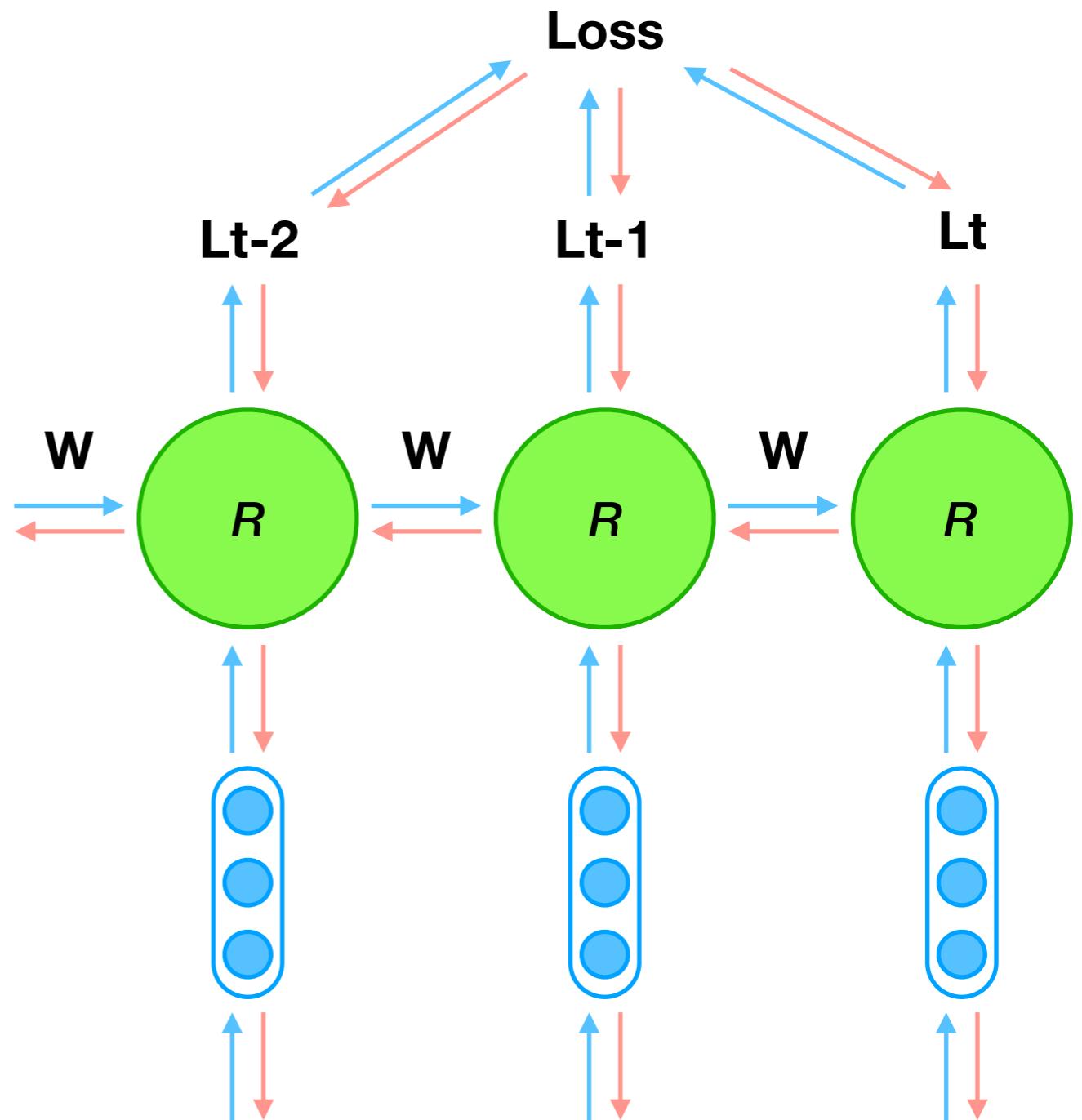


BPTT



BPTT

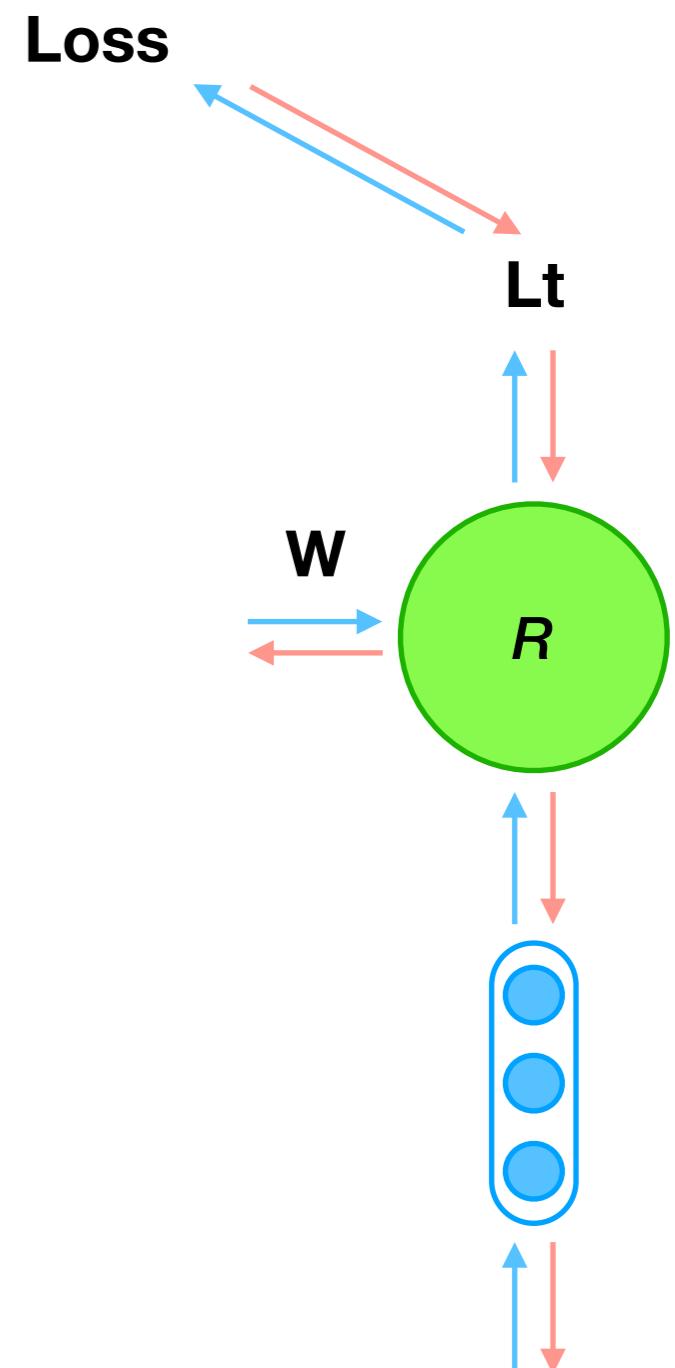
$$\frac{\partial L}{\partial W} = \sum_{i=0}^T \frac{\partial L_i}{\partial W}$$



BPTT

$$\frac{\partial L}{\partial W} = \sum_{i=0}^T \frac{\partial L_i}{\partial W}$$

$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial W}$$



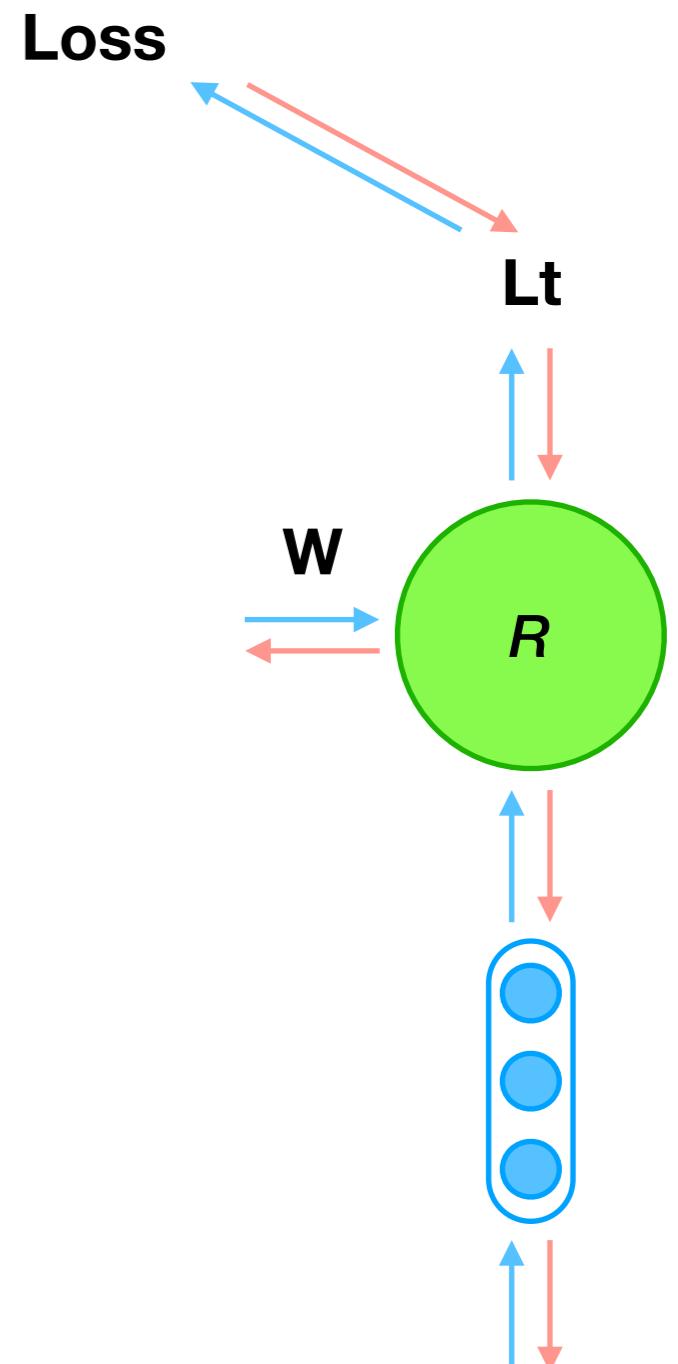
BPTT

$$\frac{\partial L}{\partial W} = \sum_{i=0}^T \frac{\partial L_i}{\partial W}$$

$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial W}$$

$$h_t = f_h(Vx_t + \boxed{W}h_{t-1} + b_h)$$

This is NOT the only dependence!



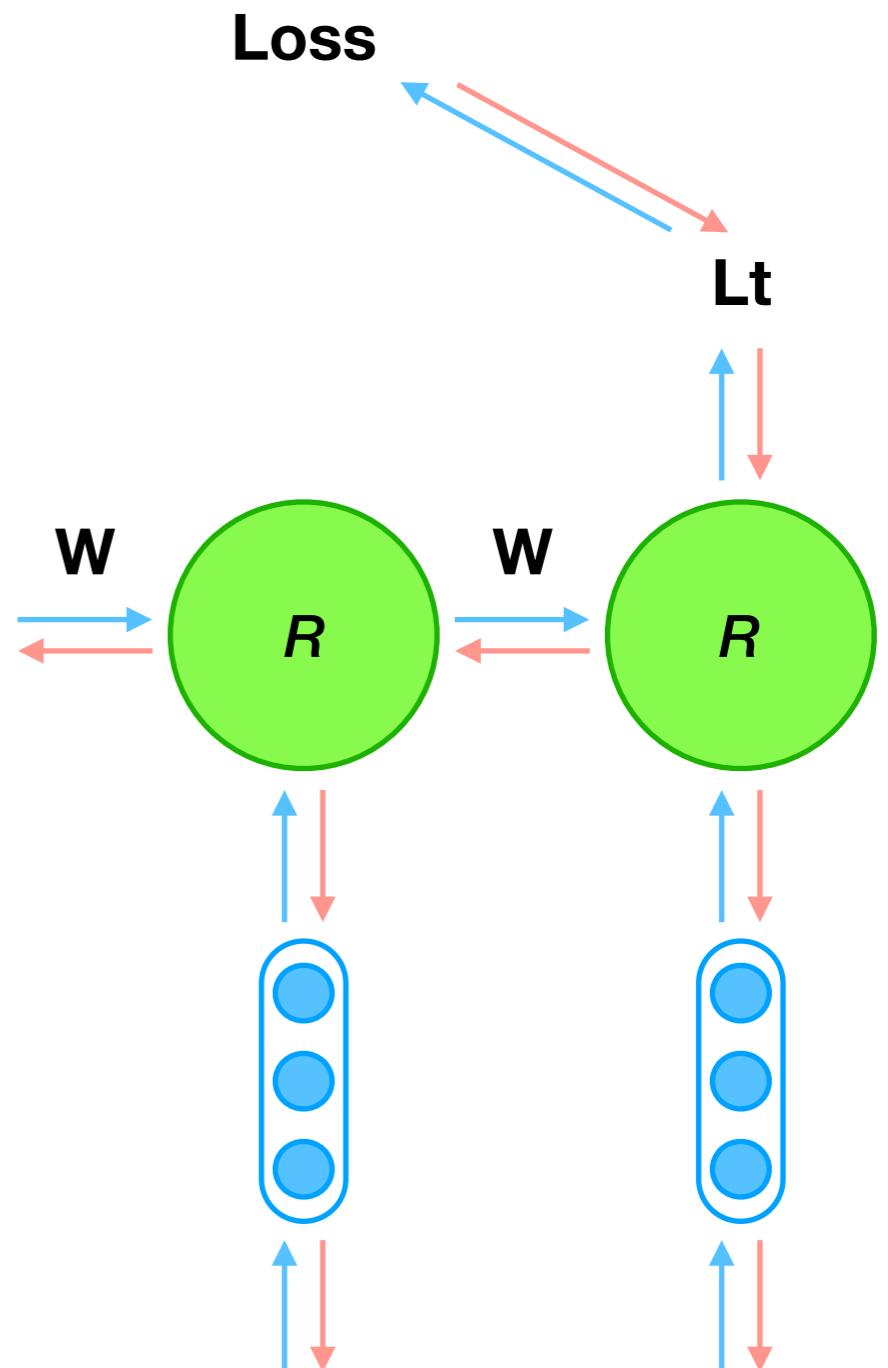
BPTT

$$\frac{\partial L}{\partial W} = \sum_{i=0}^T \frac{\partial L_i}{\partial W}$$

$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial W}$$

$$h_t = f_h(Vx_t + Wh_{t-1} + b_h)$$

This is **NOT** the only dependence!



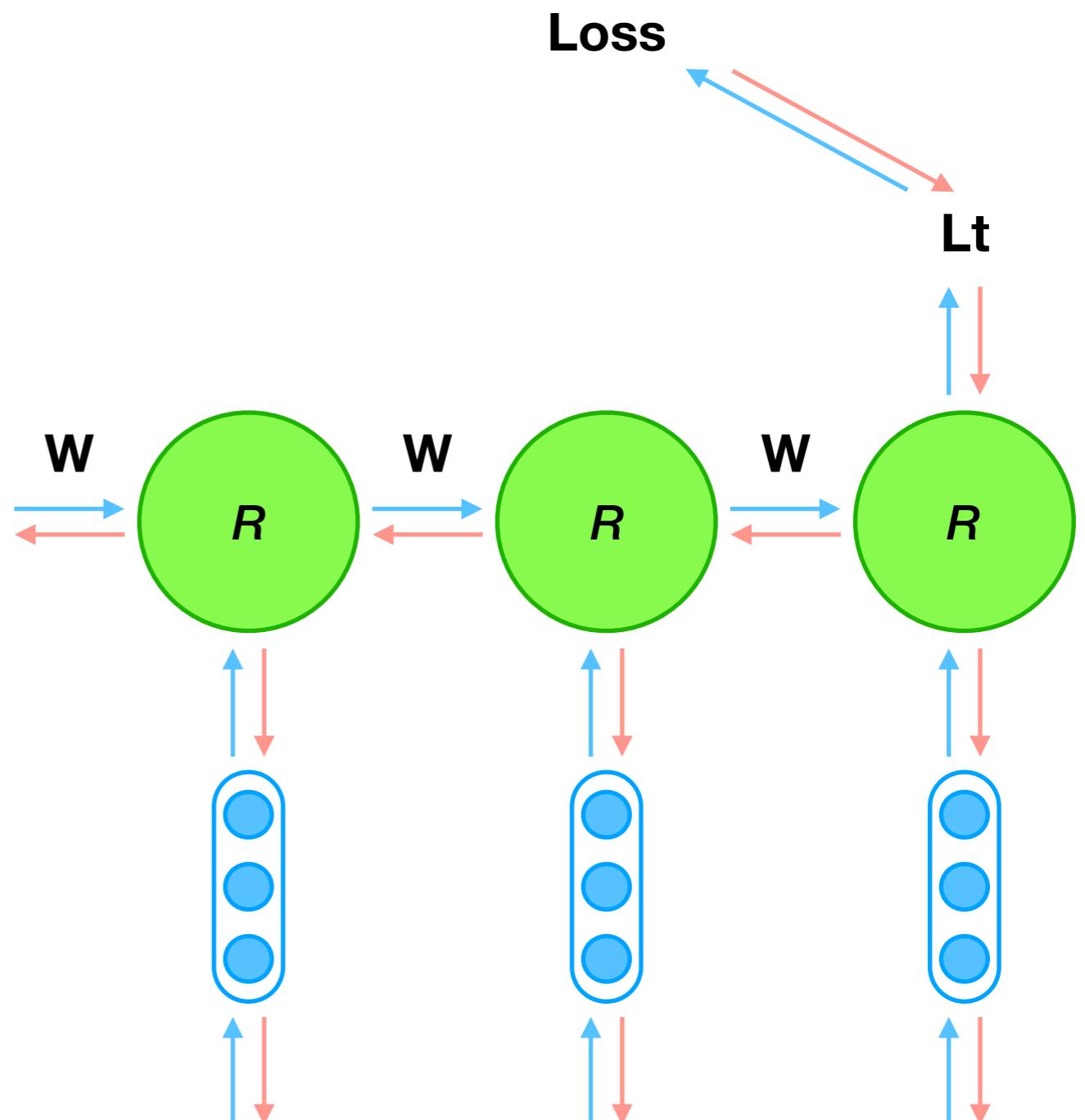
BPTT

$$\frac{\partial L}{\partial W} = \sum_{i=0}^T \frac{\partial L_i}{\partial W}$$

$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial W}$$

$$h_t = f_h(Vx_t + Wh_{t-1} + b_h)$$

This is **NOT** the only dependence!



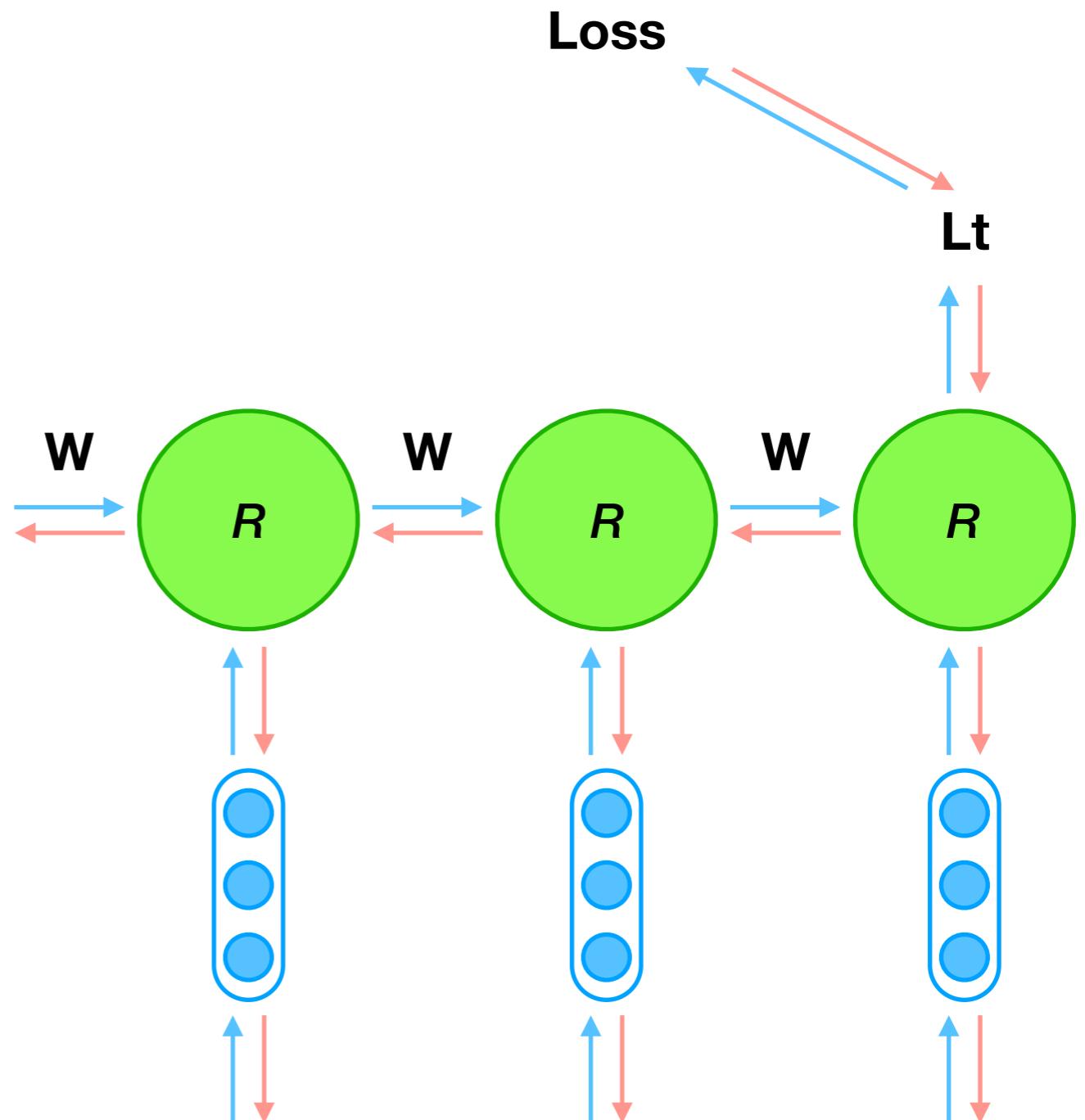
BPTT

$$\frac{\partial L}{\partial W} = \sum_{i=0}^T \frac{\partial L_i}{\partial W}$$

$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial W}$$

$$h_t = f_h(Vx_t + Wh_{t-1} + b_h)$$

This is **NOT** the only dependence!



$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \left(\frac{\partial h_t}{\partial W} + \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial W} + \dots \right)$$

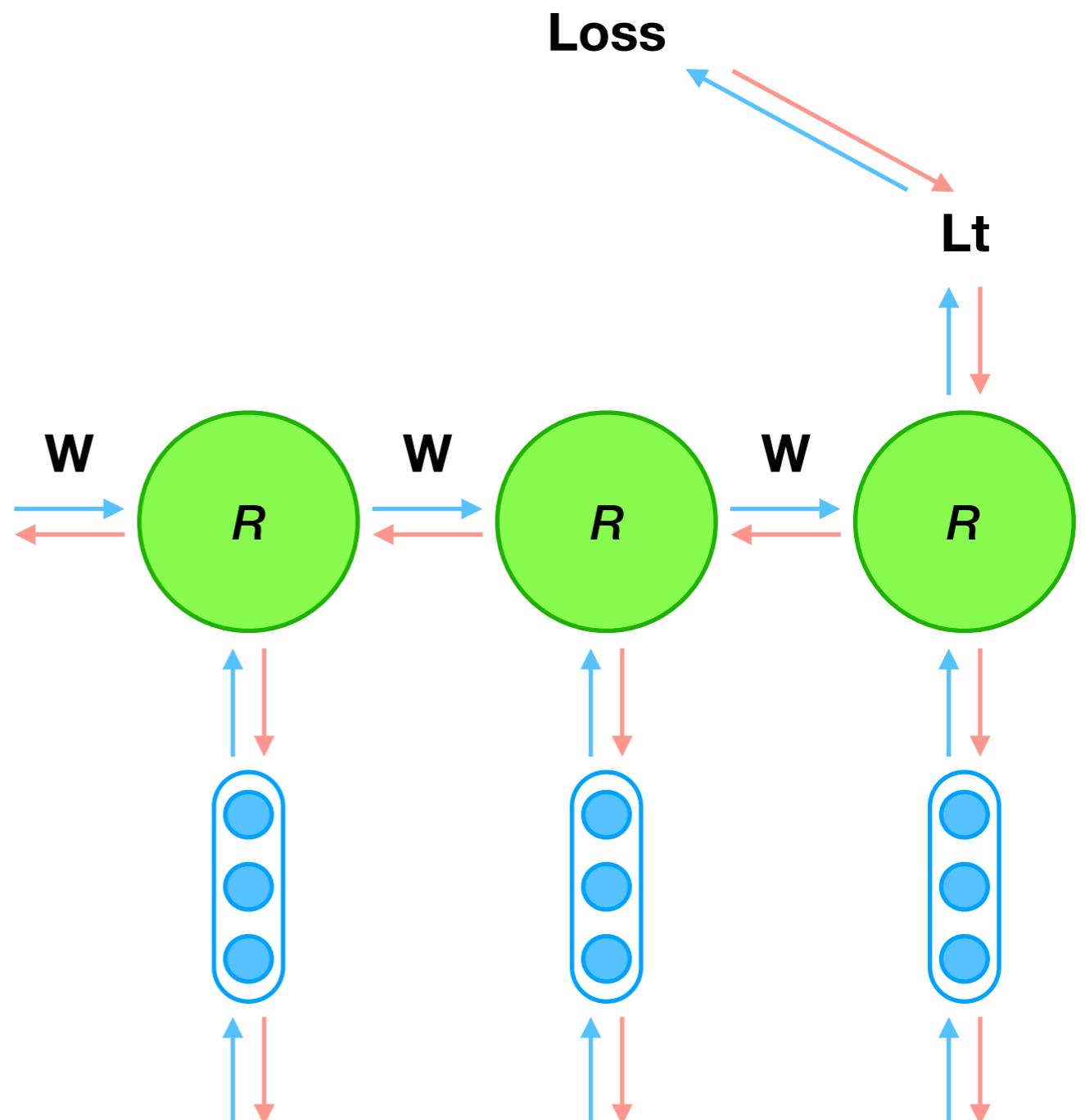
BPTT

$$\frac{\partial L}{\partial W} = \sum_{i=0}^T \frac{\partial L_i}{\partial W}$$

$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial W}$$

$$h_t = f_h(Vx_t + Wh_{t-1} + b_h)$$

This is **NOT** the only dependence!



$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \sum_{k=0}^t \left(\prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}} \right) \frac{\partial h_k}{\partial W}$$

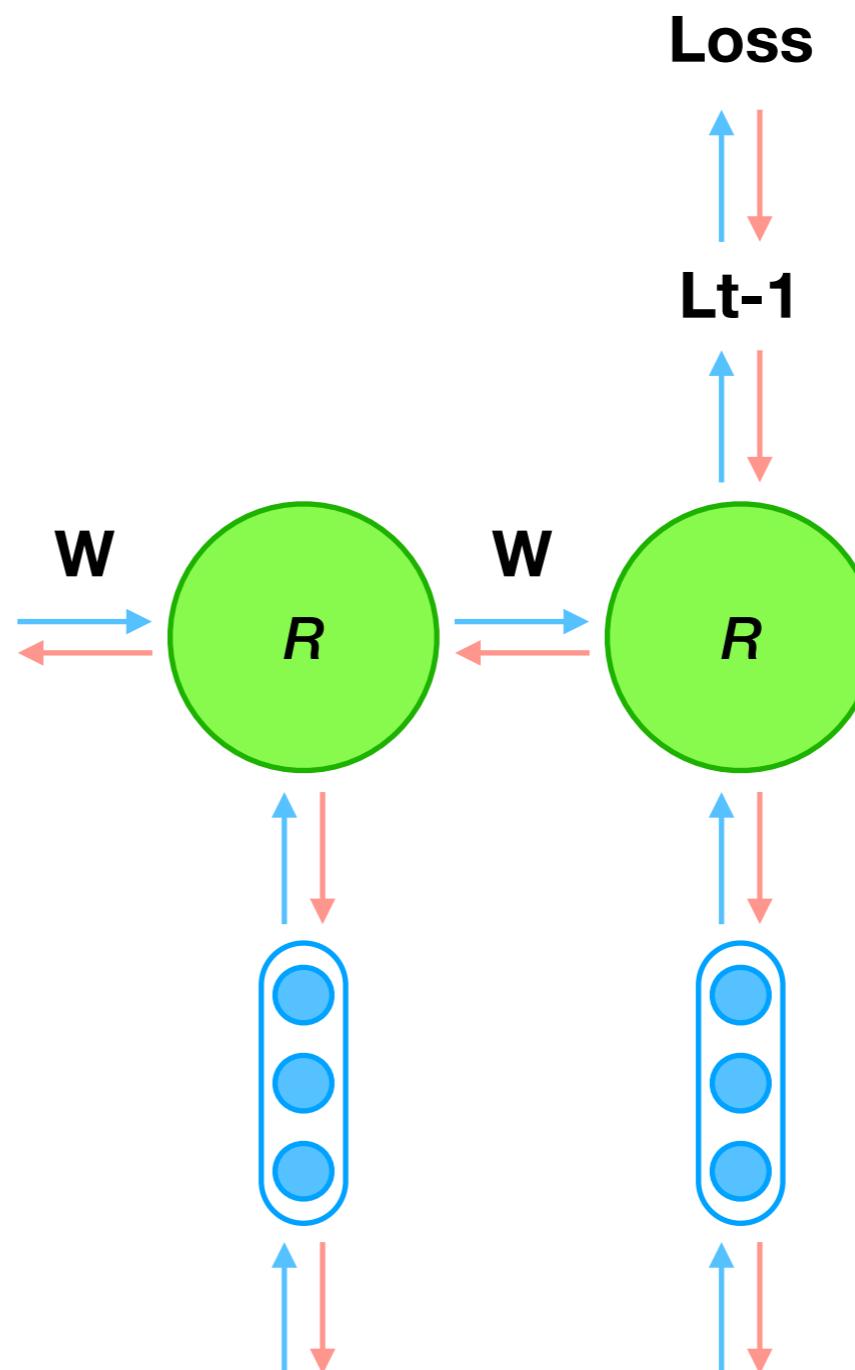
BPTT

$$\frac{\partial L}{\partial W} = \sum_{i=0}^T \frac{\partial L_i}{\partial W}$$

$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial W}$$

$$h_t = f_h(Vx_t + Wh_{t-1} + b_h)$$

This is **NOT** the only dependence!



$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \sum_{k=0}^t \left(\prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}} \right) \frac{\partial h_k}{\partial W}$$

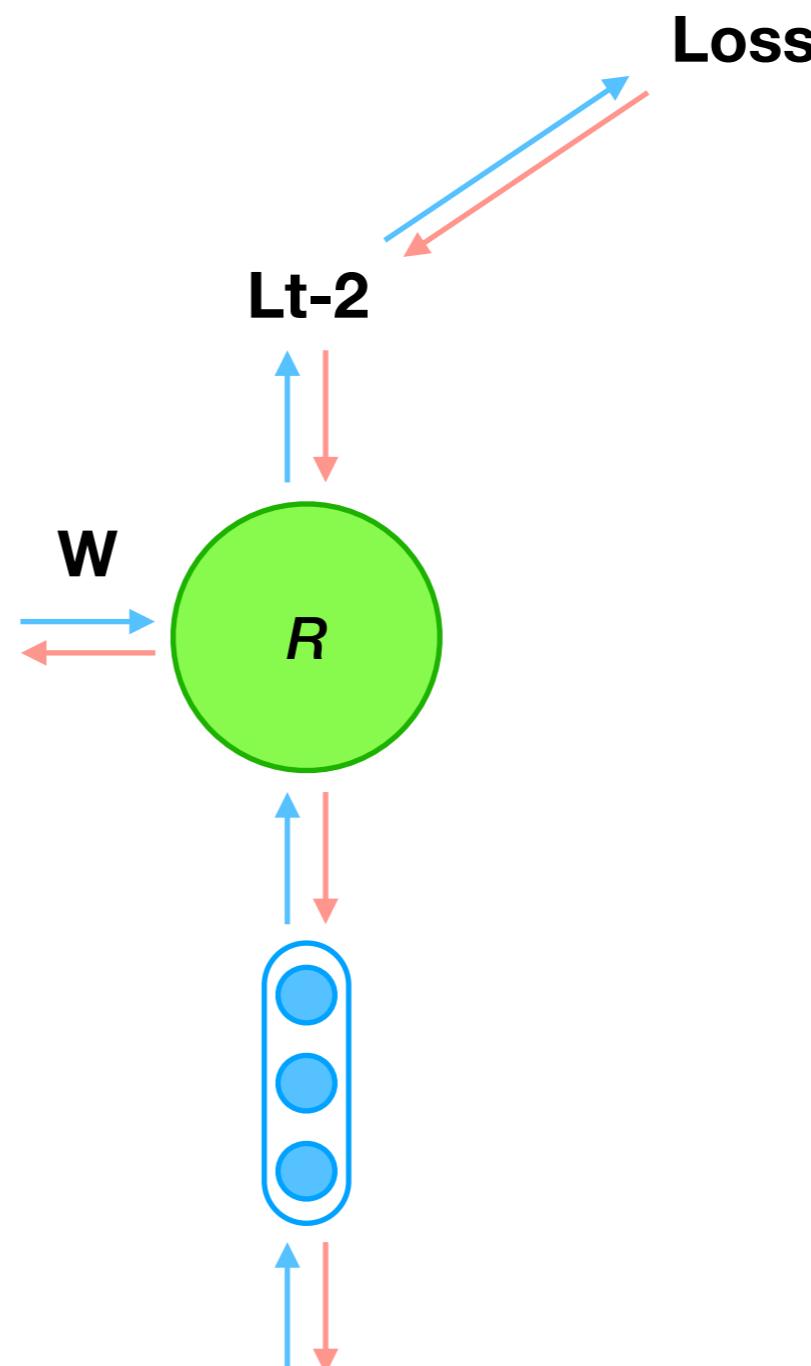
BPTT

$$\frac{\partial L}{\partial W} = \sum_{i=0}^T \frac{\partial L_i}{\partial W}$$

$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial W}$$

$$h_t = f_h(Vx_t + \boxed{W}h_{t-1} + b_h)$$

This is NOT the only dependence!



$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \sum_{k=0}^t \left(\prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}} \right) \frac{\partial h_k}{\partial W}$$

Gradients

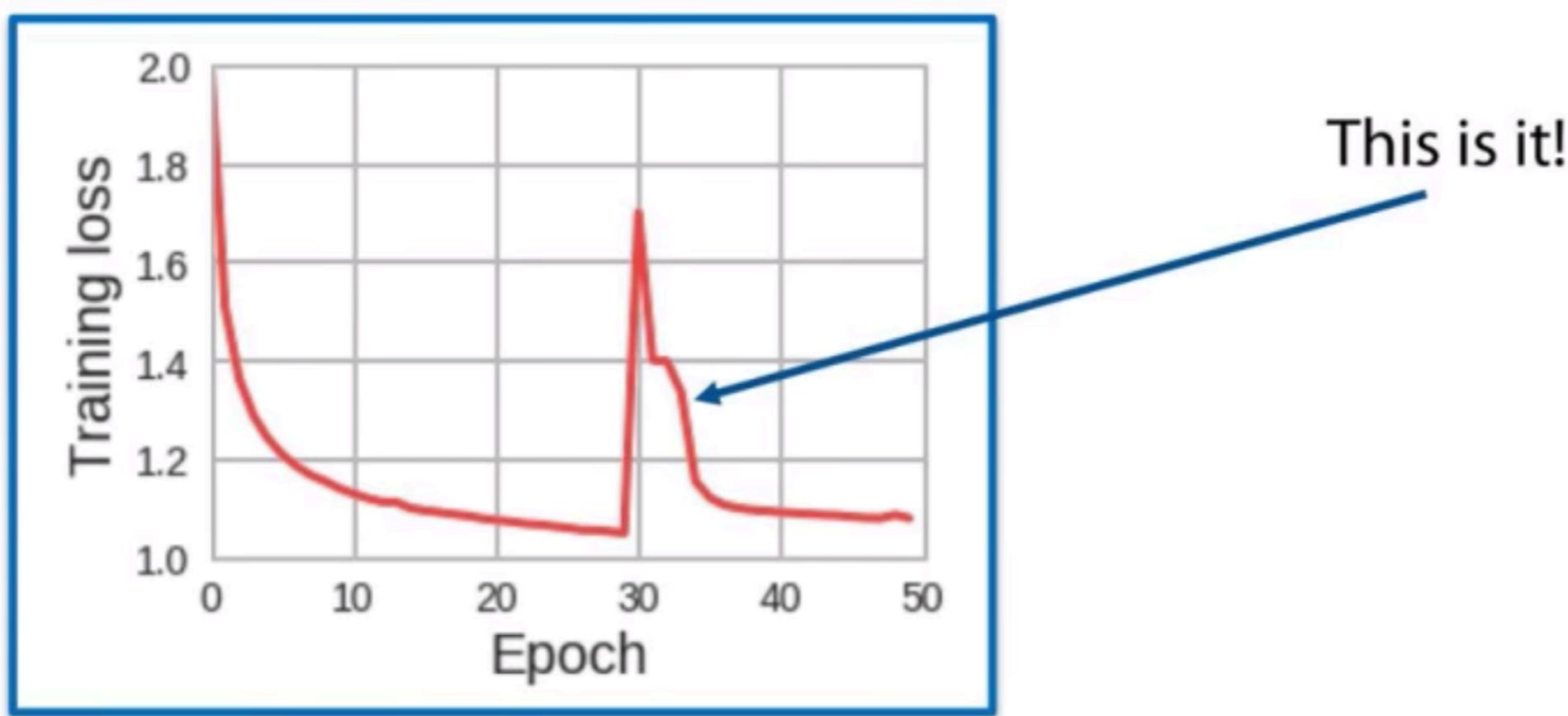
$$\frac{\partial L_t}{\partial W} \propto \sum_{k=0}^t \left(\prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}} \right) \frac{\partial h_k}{\partial W}$$

$\left\| \frac{\partial h_i}{\partial h_{i-1}} \right\|_2 < 1 \longrightarrow \text{Vanishing gradients}$

$\left\| \frac{\partial h_i}{\partial h_{i-1}} \right\|_2 > 1 \longrightarrow \text{Exploding gradients}$

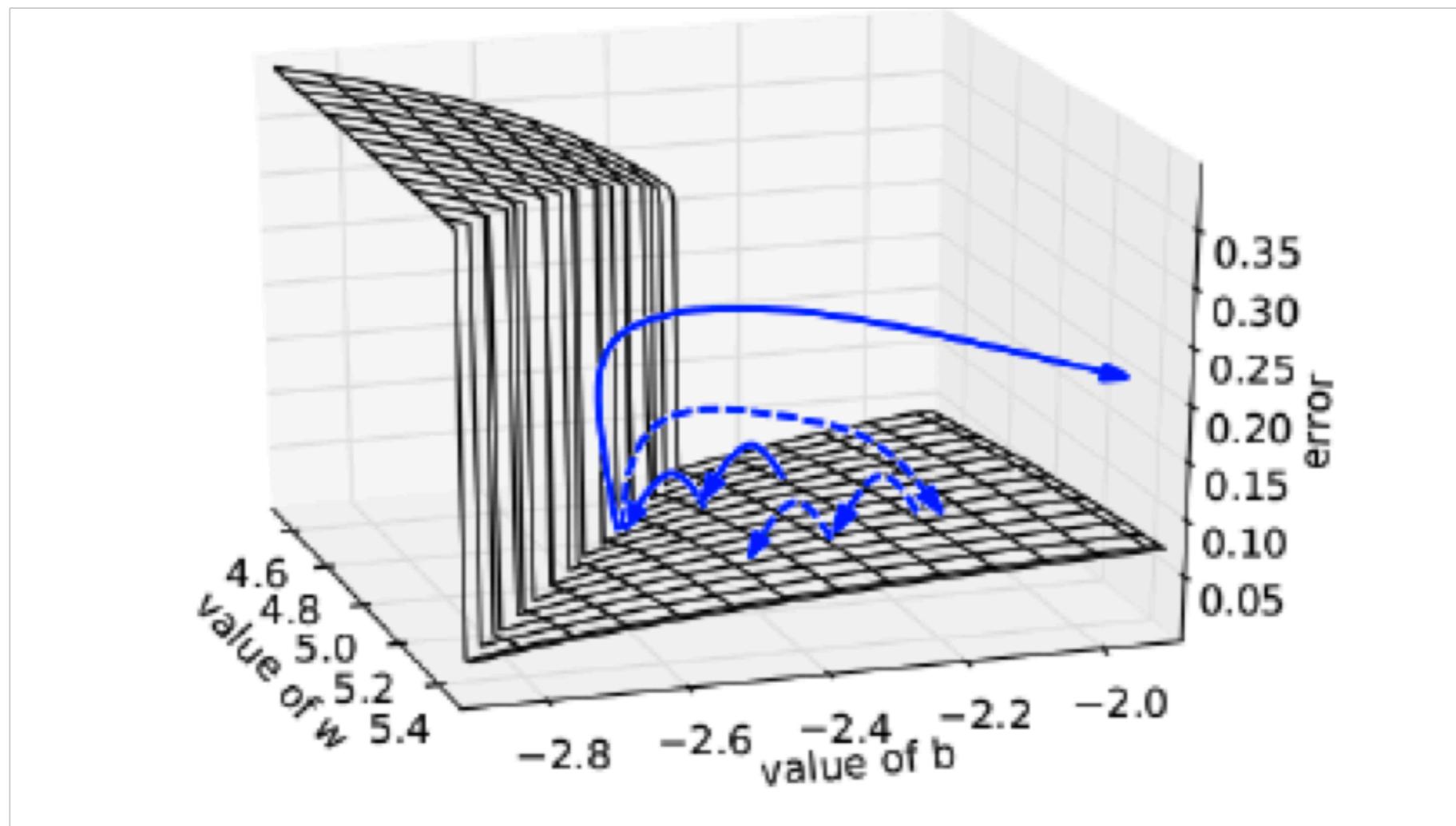
Exploding gradients

Unstable learning curve



If the gradients contain NaNs you end up
with NaNs in the weights

Exploding gradients



Gradient Clipping

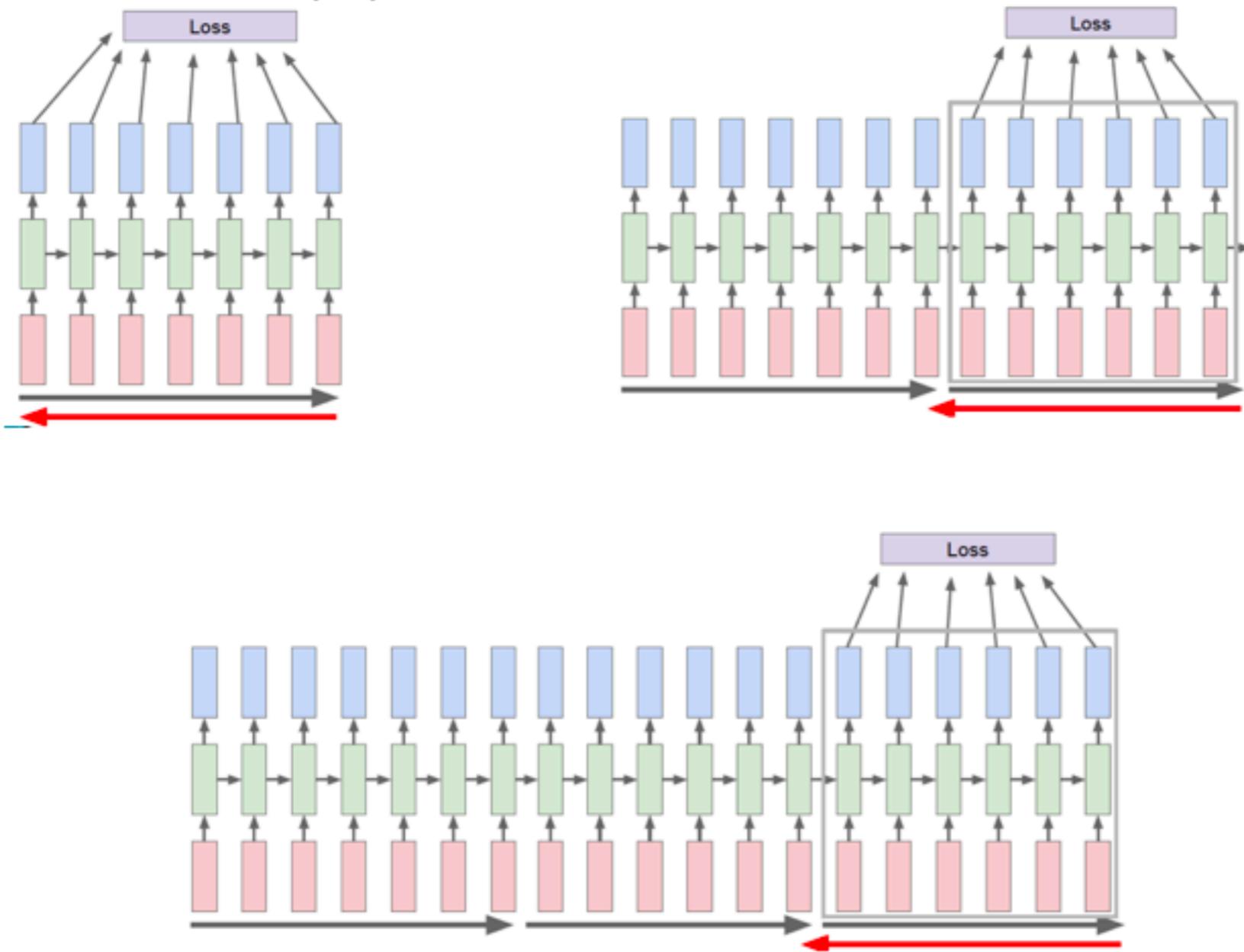
Gradient $g = \frac{\partial L}{\partial \theta}$, θ - all the network parameters

If $\|g\| > \text{threshold}$:

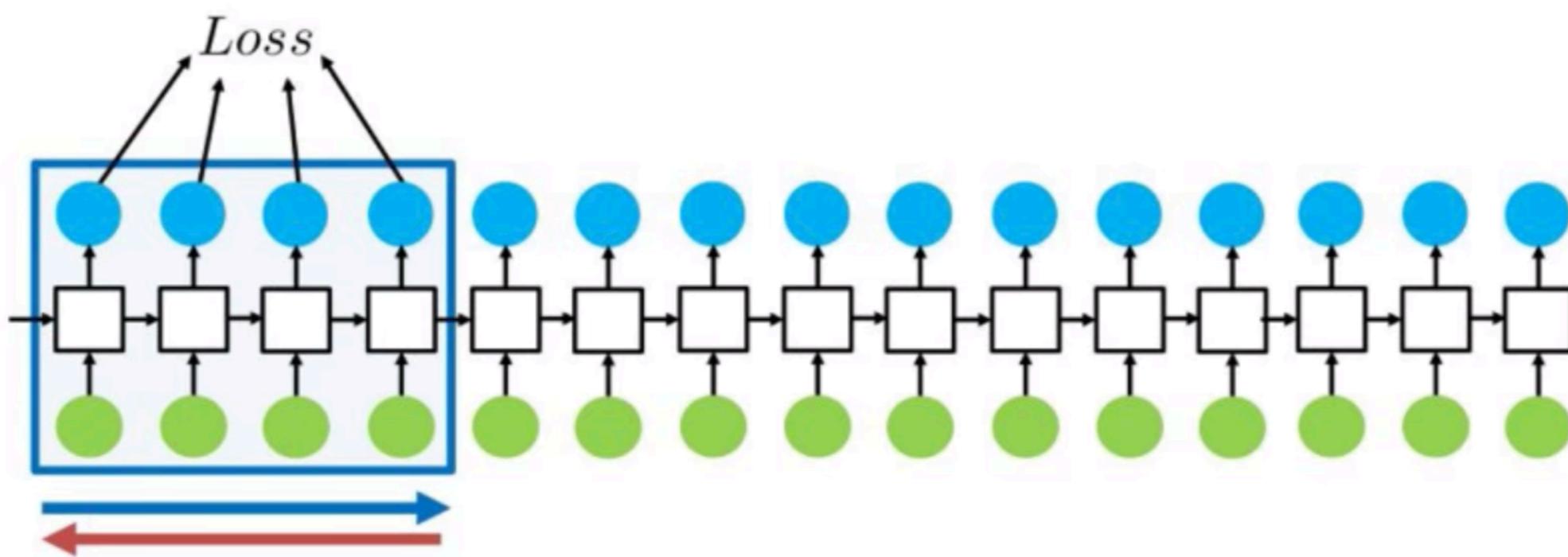
$$g \leftarrow \frac{\text{threshold}}{\|g\|} g$$

Simple but still very effective!

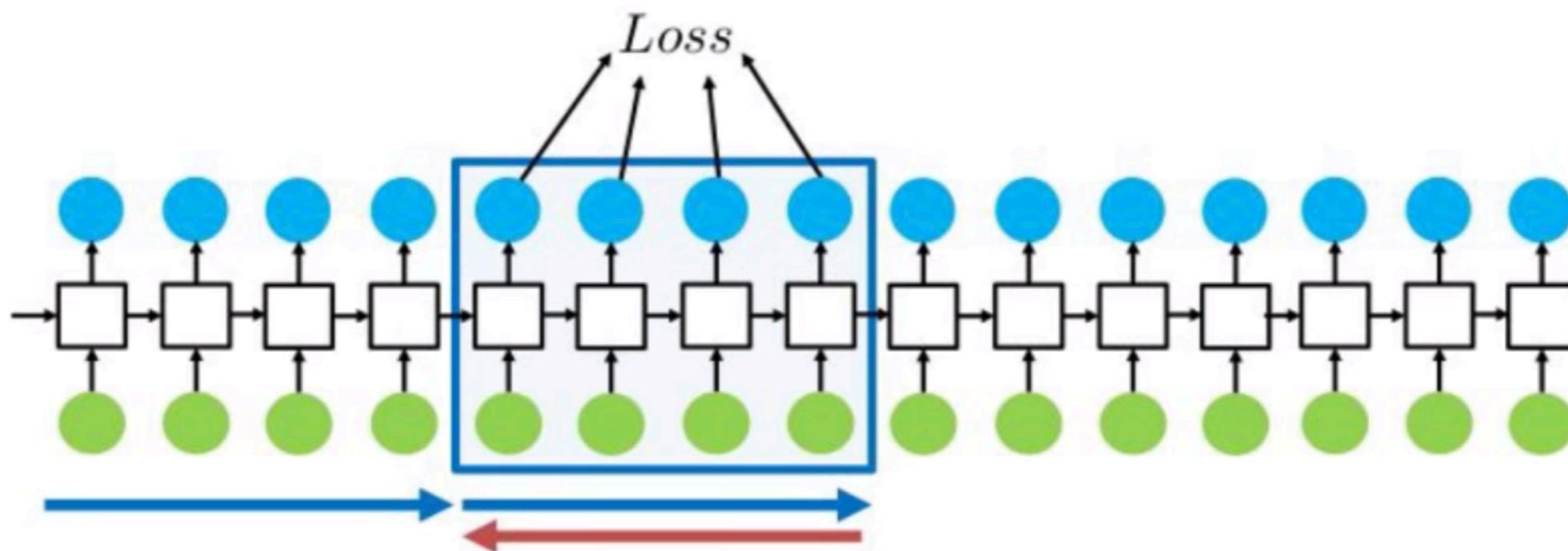
Truncated BPTT



Truncated BPTT



Truncated BPTT



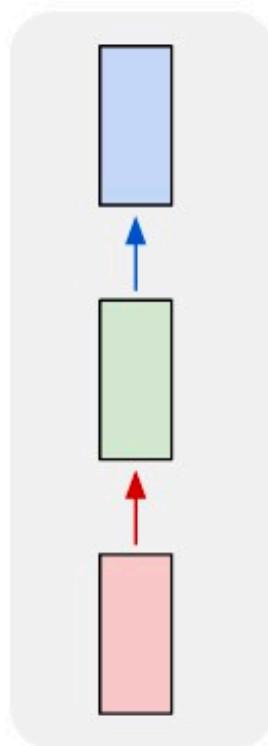
Recurrent Neural Network

+

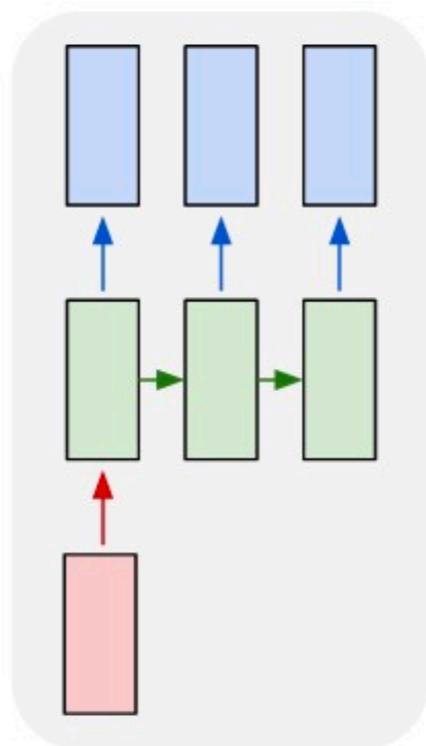
-

Recurrent Neural Network

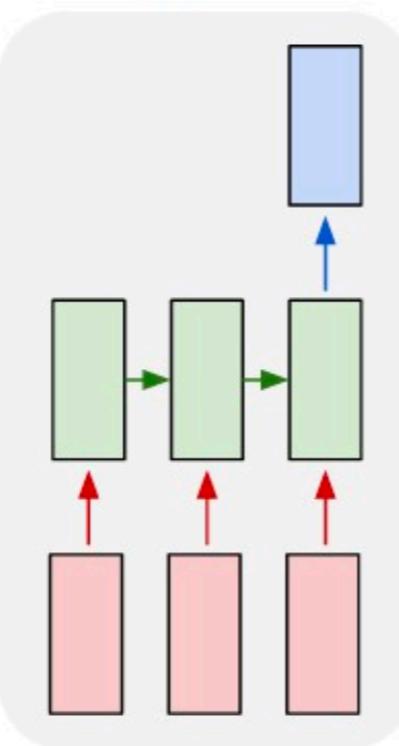
one to one



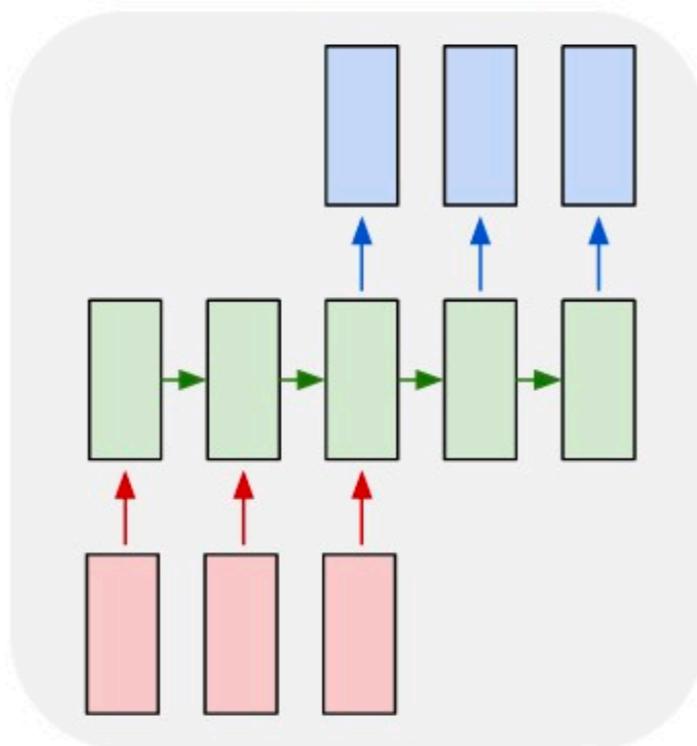
one to many



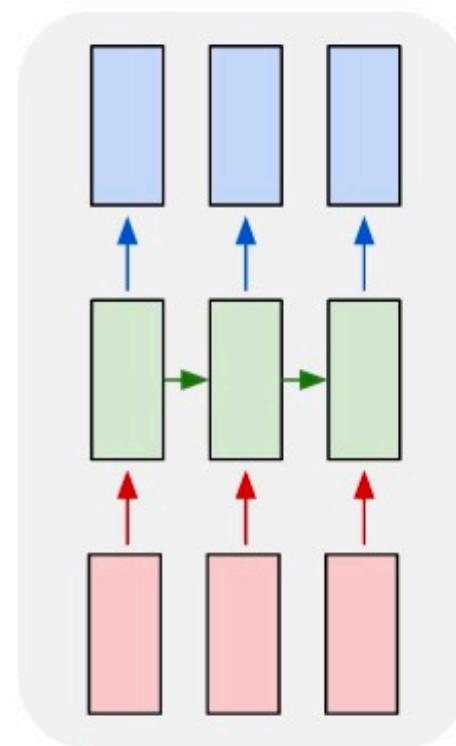
many to one



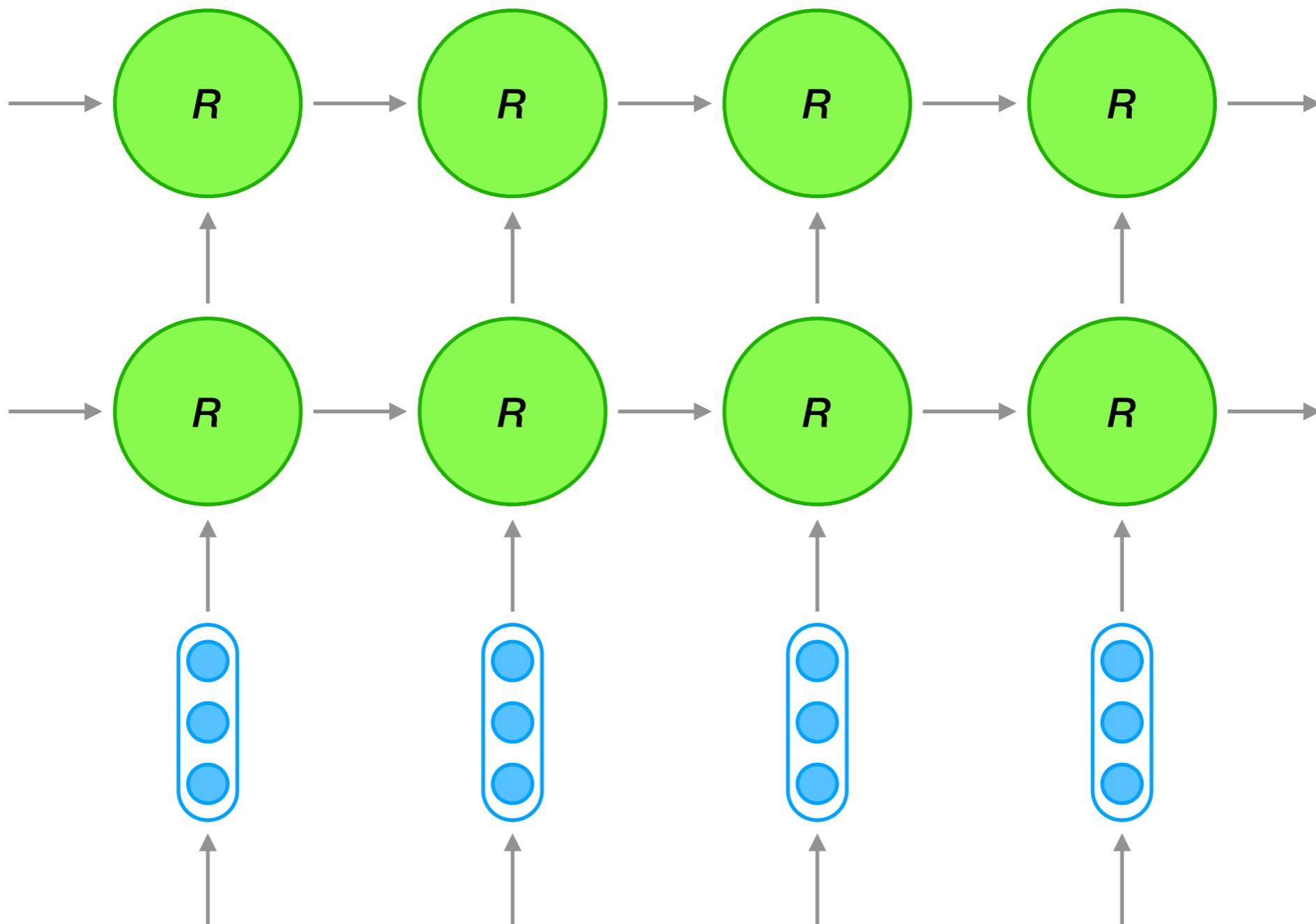
many to many



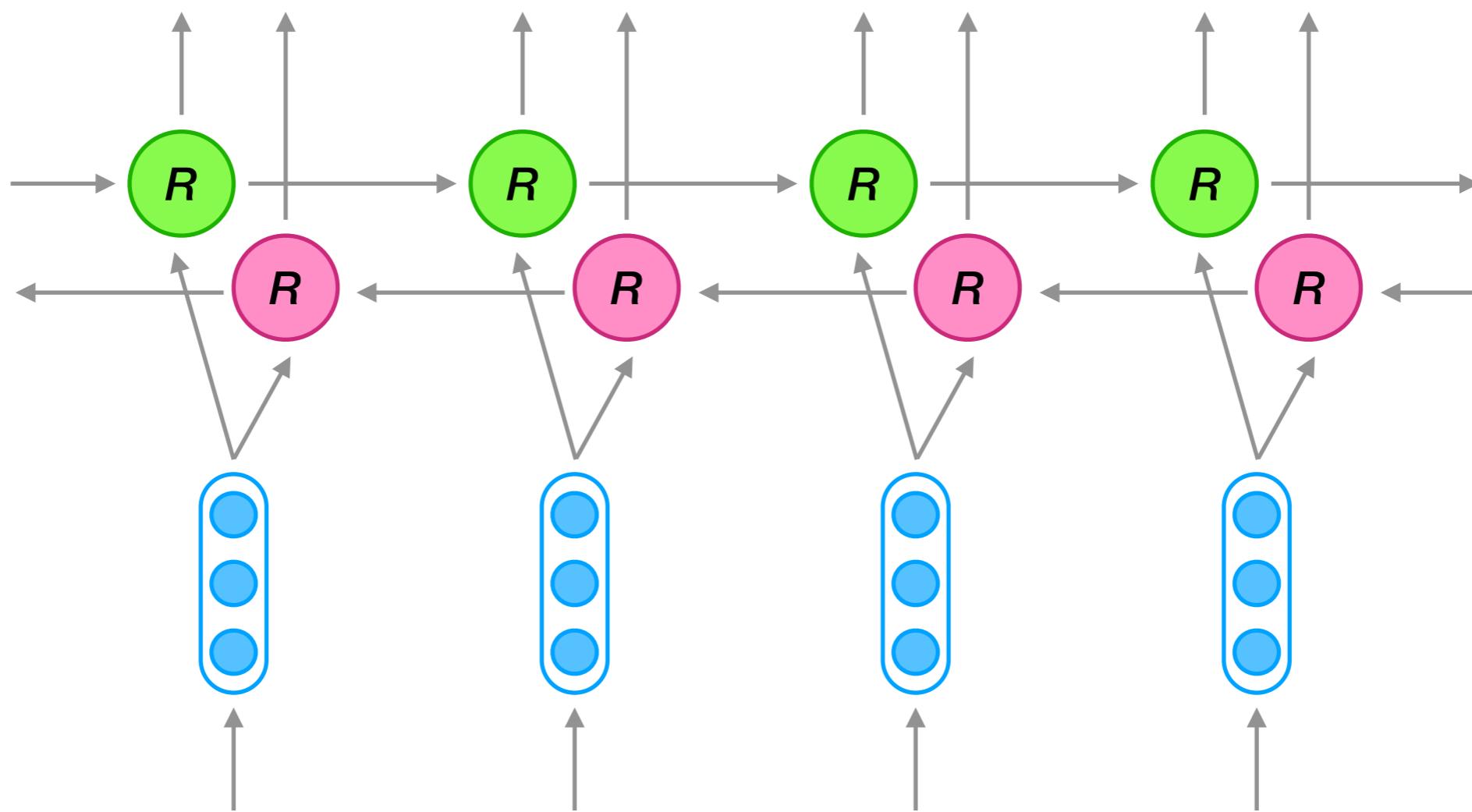
many to many



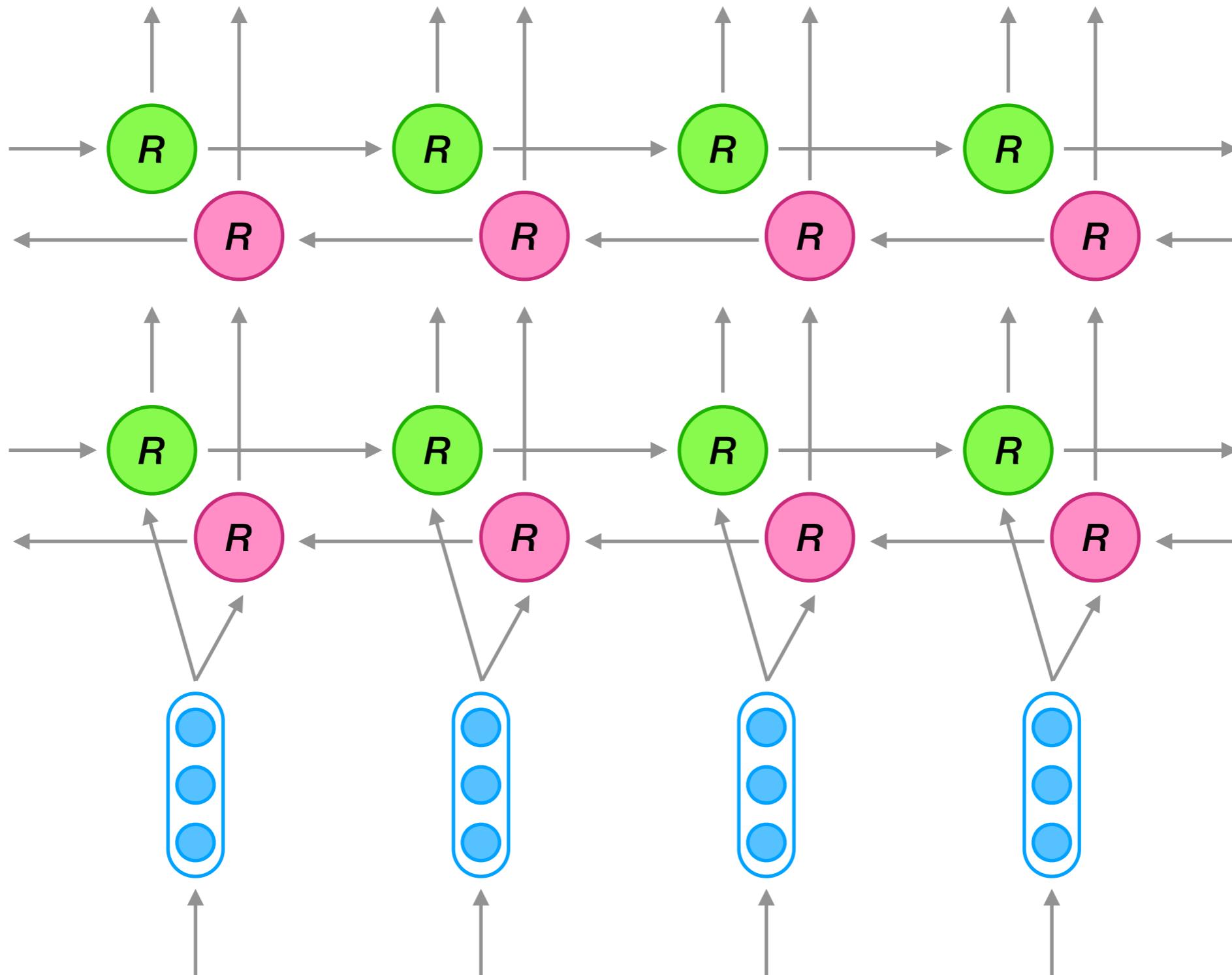
Deeper



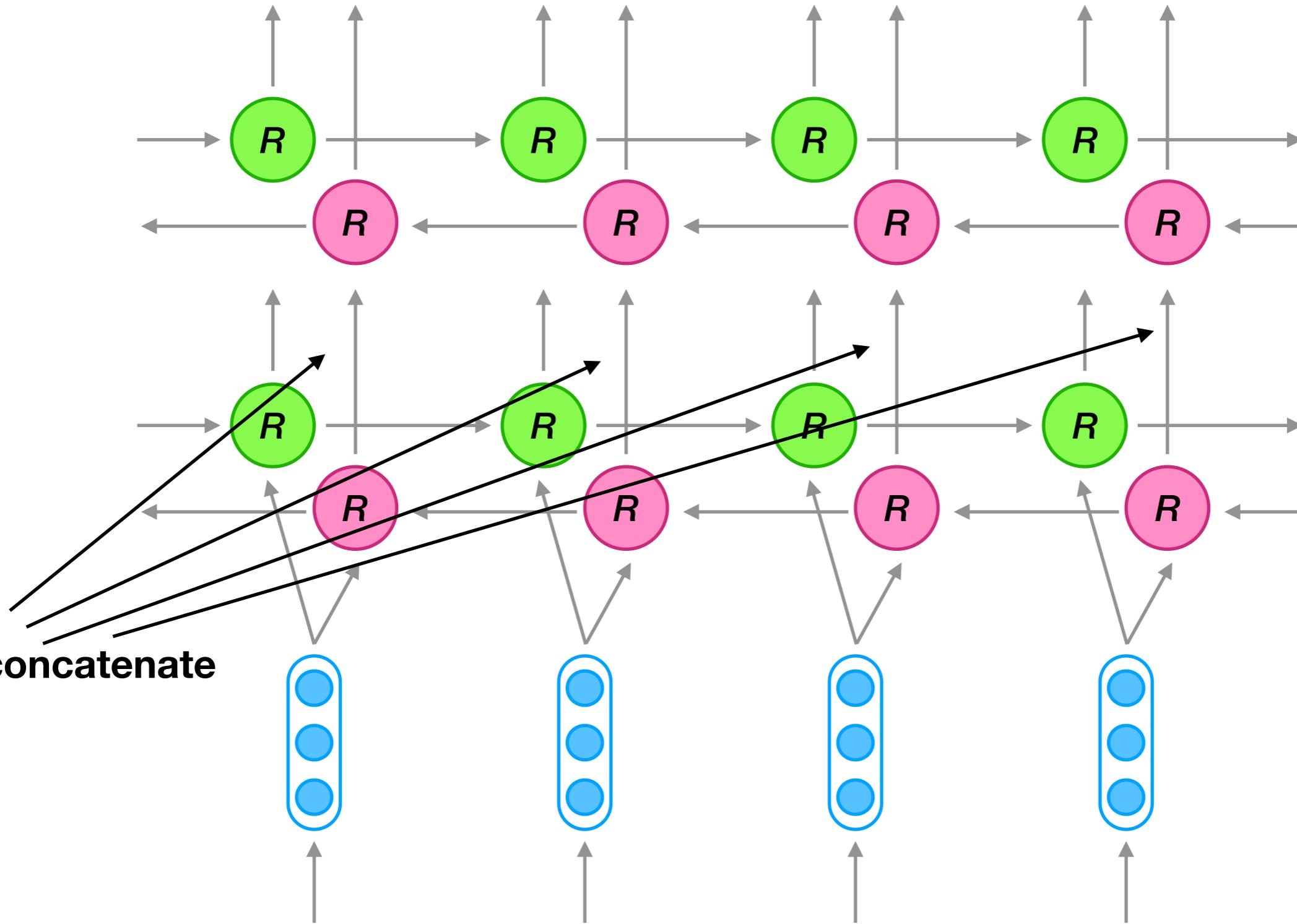
Bidirectional



Deeper Bidirectional



Deeper Bidirectional



Bidirectional

Source <START> The poor don't have any money <END>

Bidirectional

Source <START> The poor don't have any money <END>

Forward

Backward

Bidirectional

Source <START> The poor don't have any money <END>

Forward <START>

Backward <END>

Bidirectional

Source <START> The poor don't have any money <END>

Forward <START> The

Backward <END> money

Bidirectional

Source <START> The poor don't have any money <END>

Forward <START> The poor

Backward <END> money any

Bidirectional

Source <START> The poor don't have any money <END>

Forward <START> The poor don't

Backward <END> money any have

Bidirectional

Source <START> The poor don't have any money <END>

Forward <START> The poor don't have

Backward <END> money any have don't

Bidirectional

Source <START> The poor don't have any money <END>

Forward <START> The poor don't have any

Backward <END> money any have don't poor

Bidirectional

Source <START> The poor don't have any money <END>

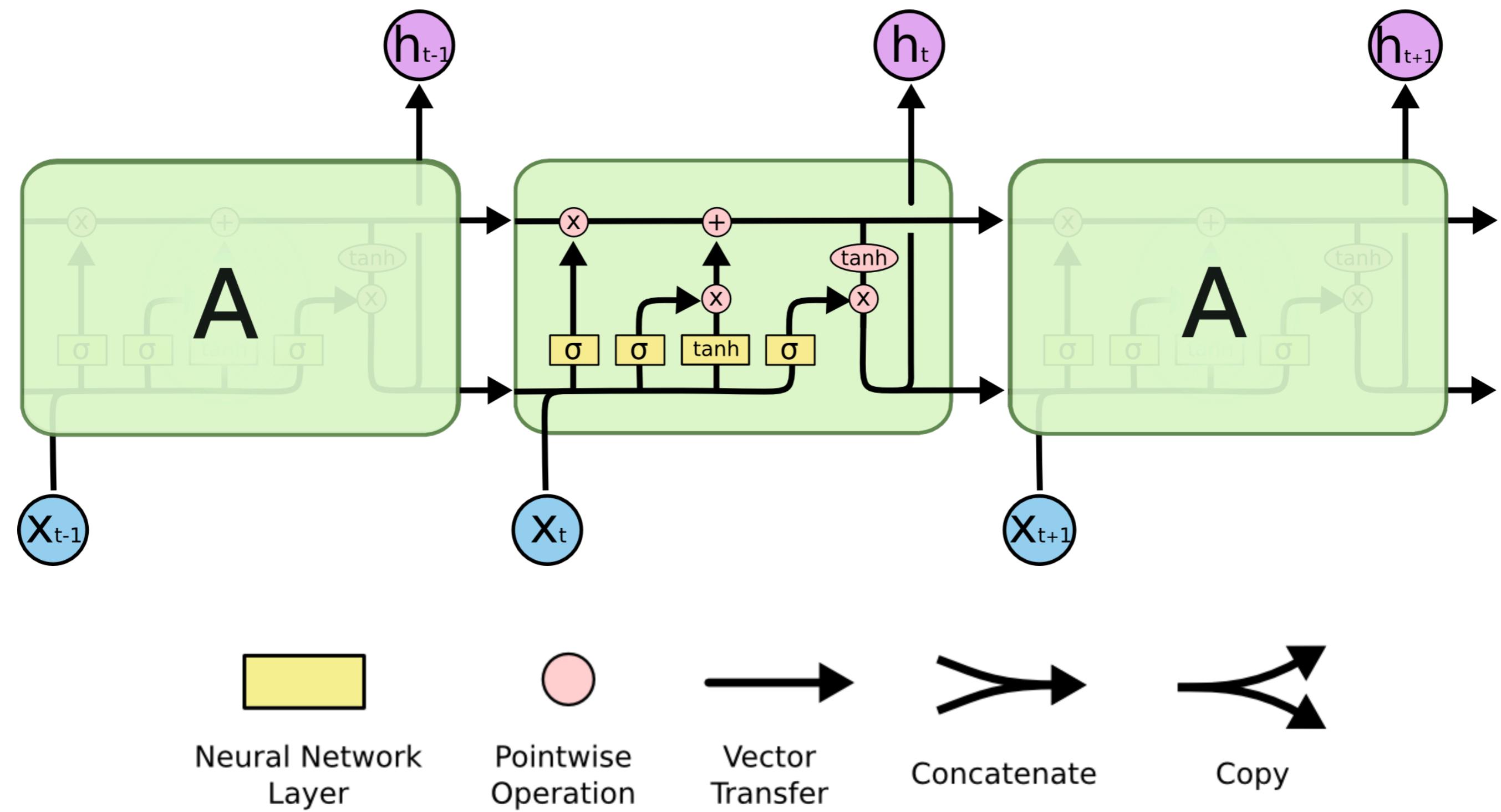
Forward <START> The poor don't have any money

Backward <END> money any have don't poor The

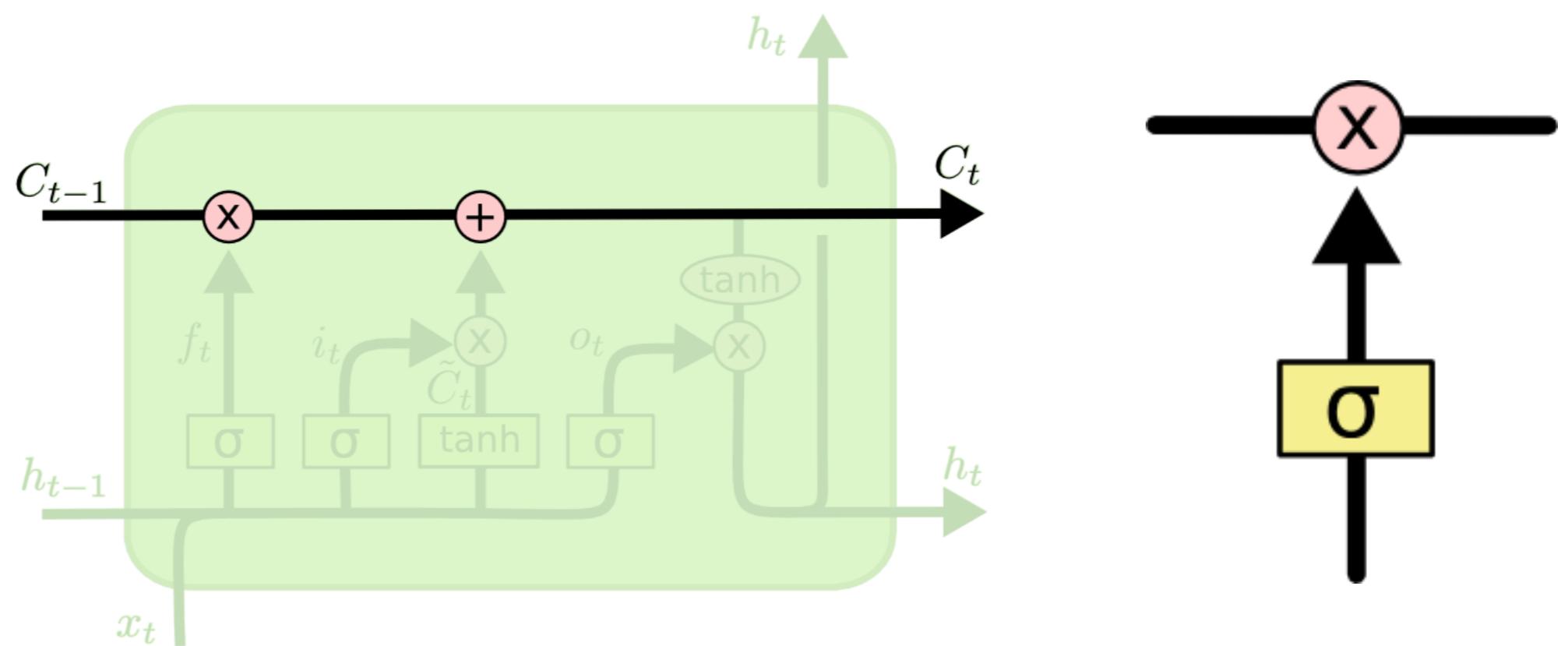
Bidirectional

Source	<START>	The	poor	don't	have	any	money	<END>
Forward	<START>	The	poor	don't	have	any	money	<END>
Backward	<END>	money	any	have	don't	poor	The	<START>

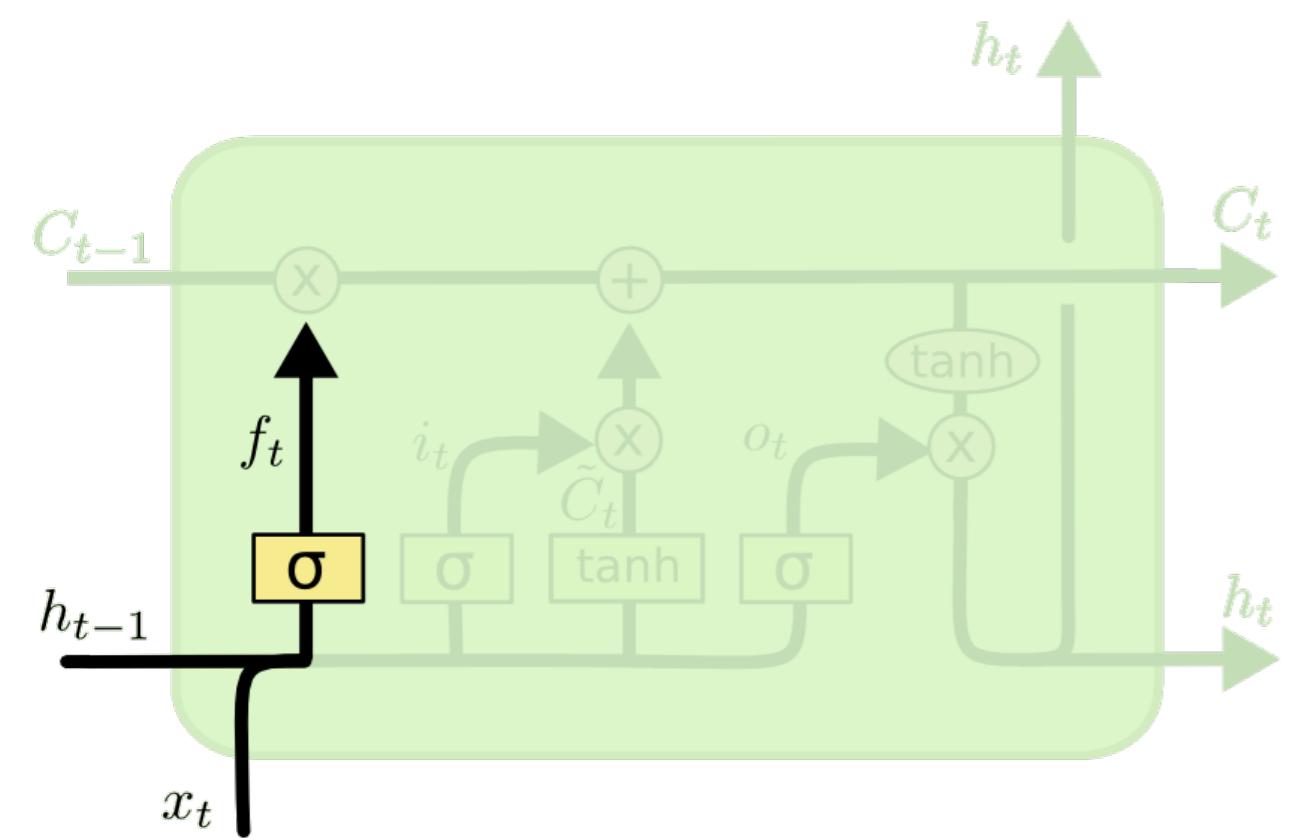
Long Short Term Memories



Long Short Term Memories

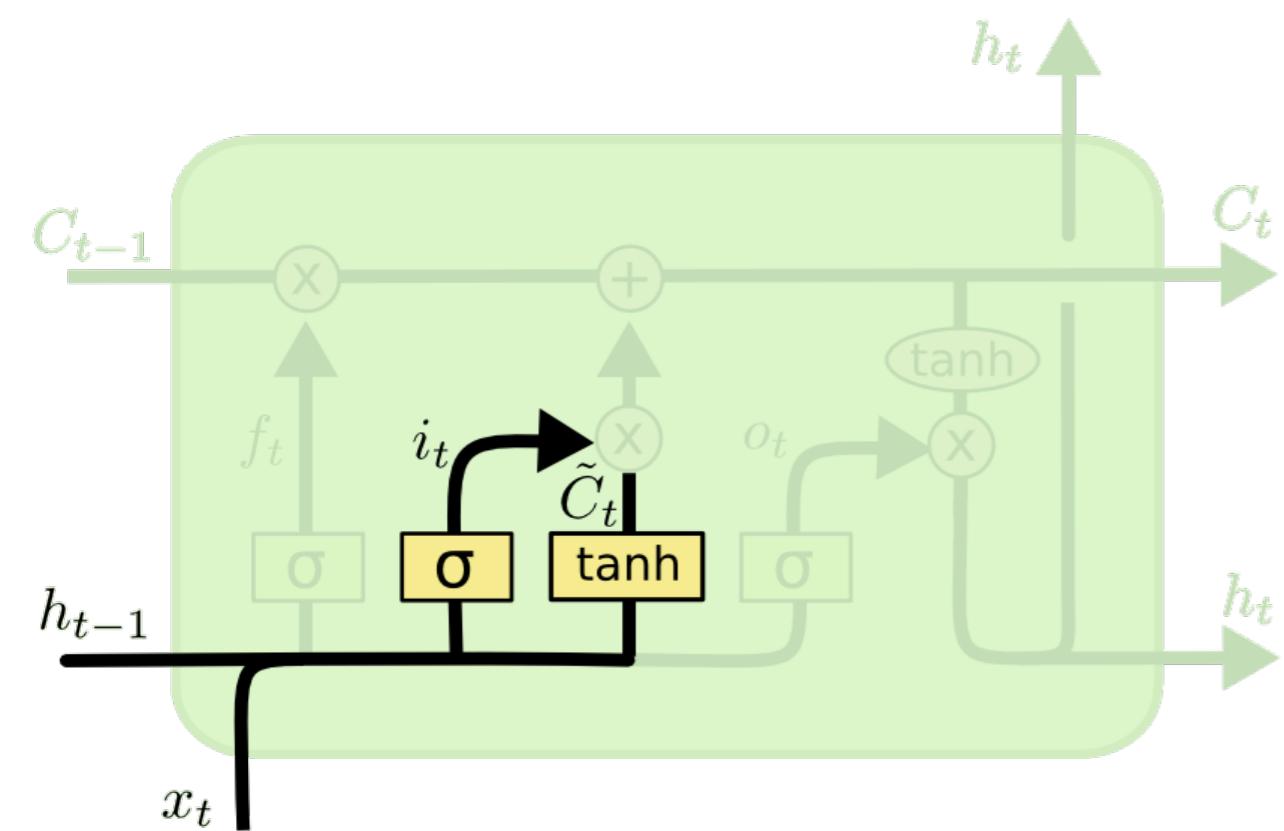


Long Short Term Memories



$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

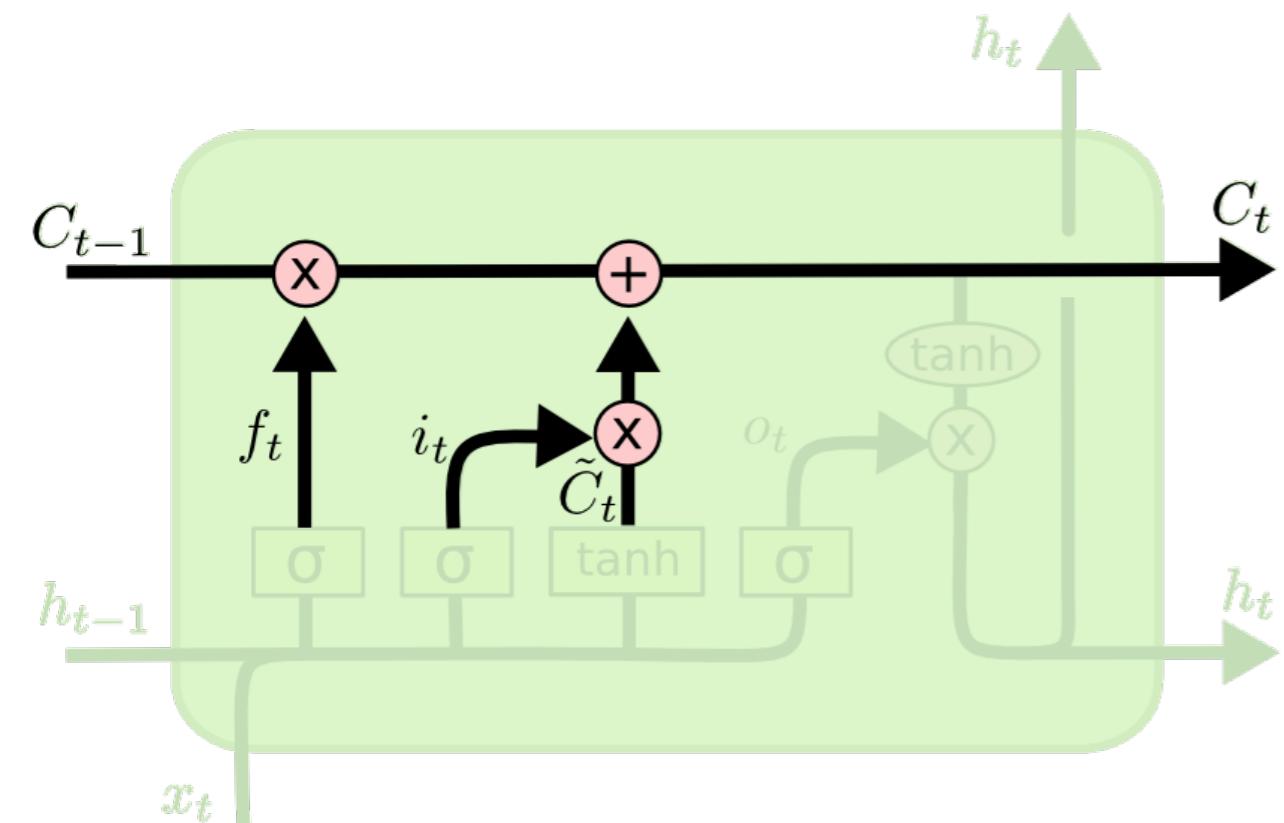
Long Short Term Memories



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

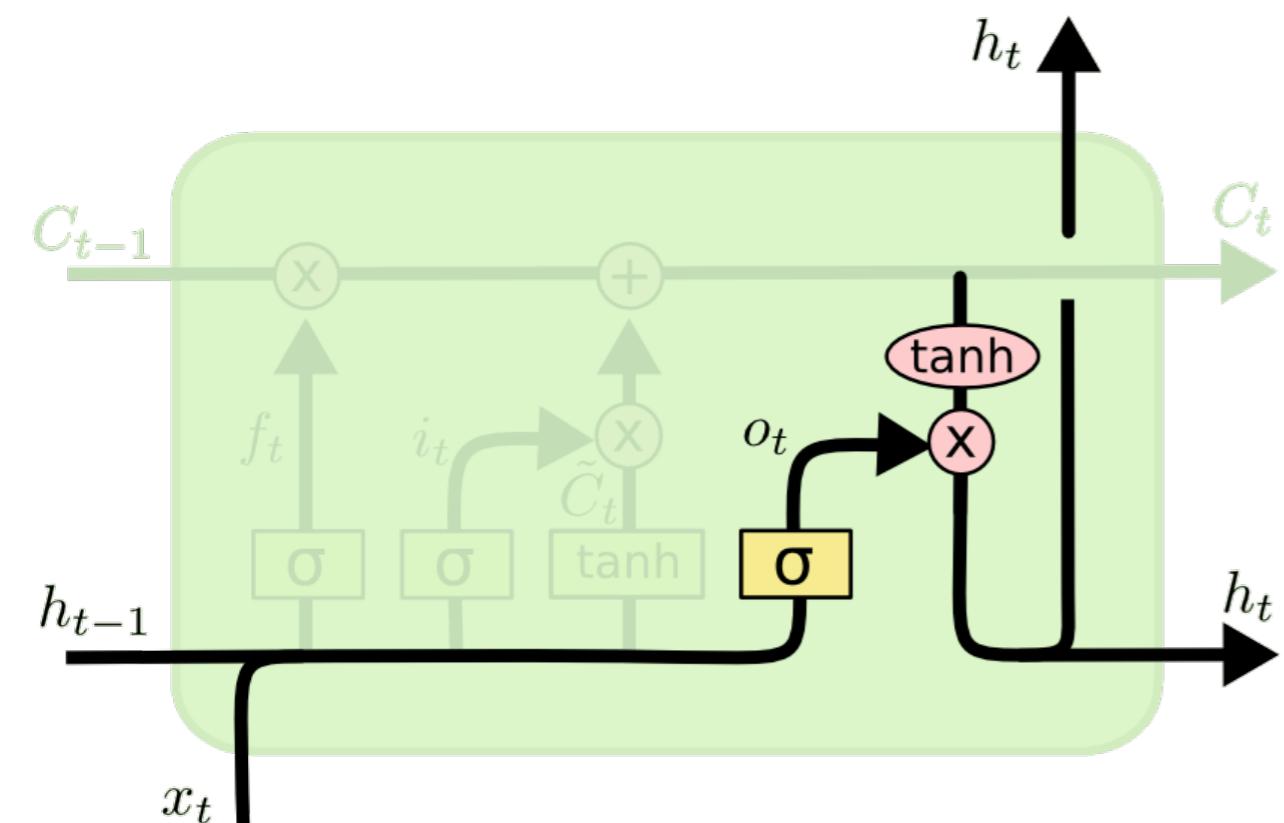
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Long Short Term Memories



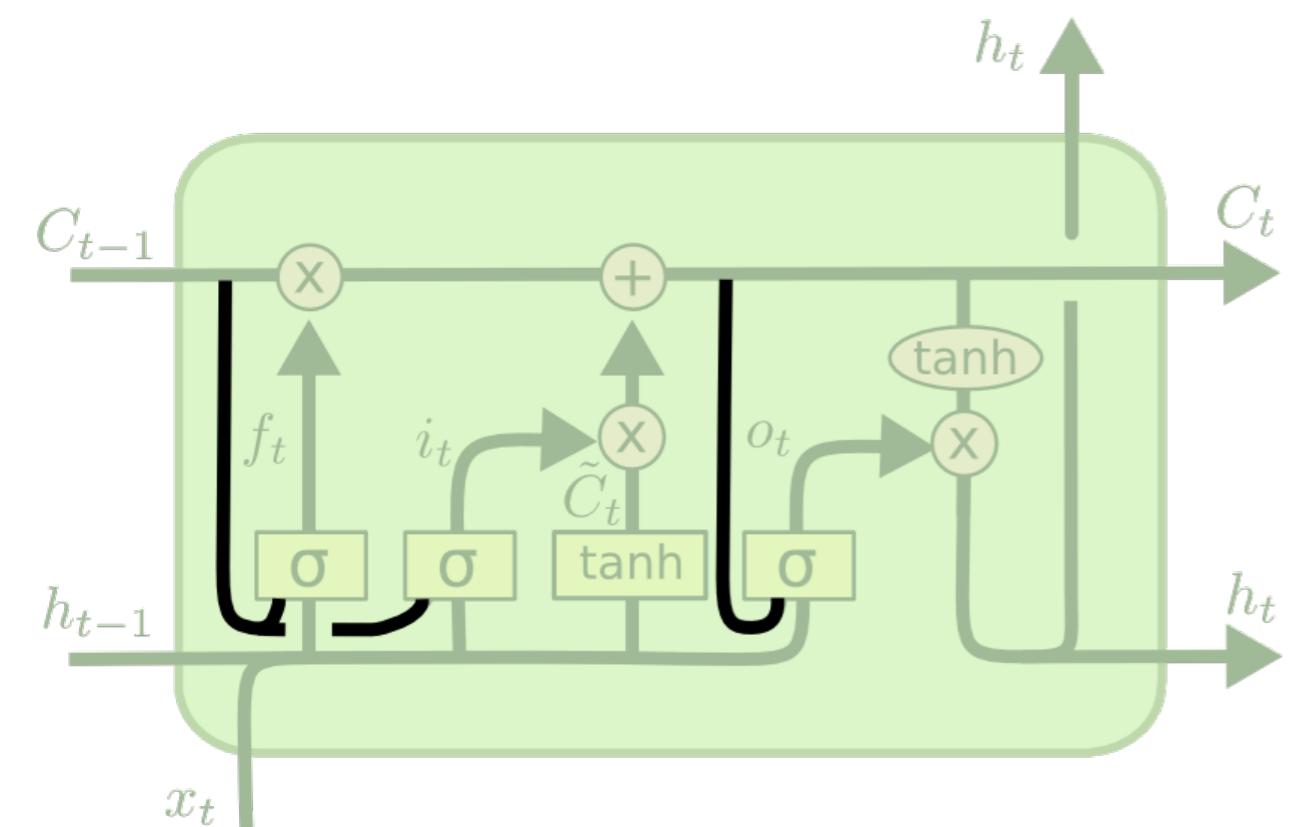
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Long Short Term Memories



$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$
$$h_t = o_t * \tanh (C_t)$$

Long Short Term Memories

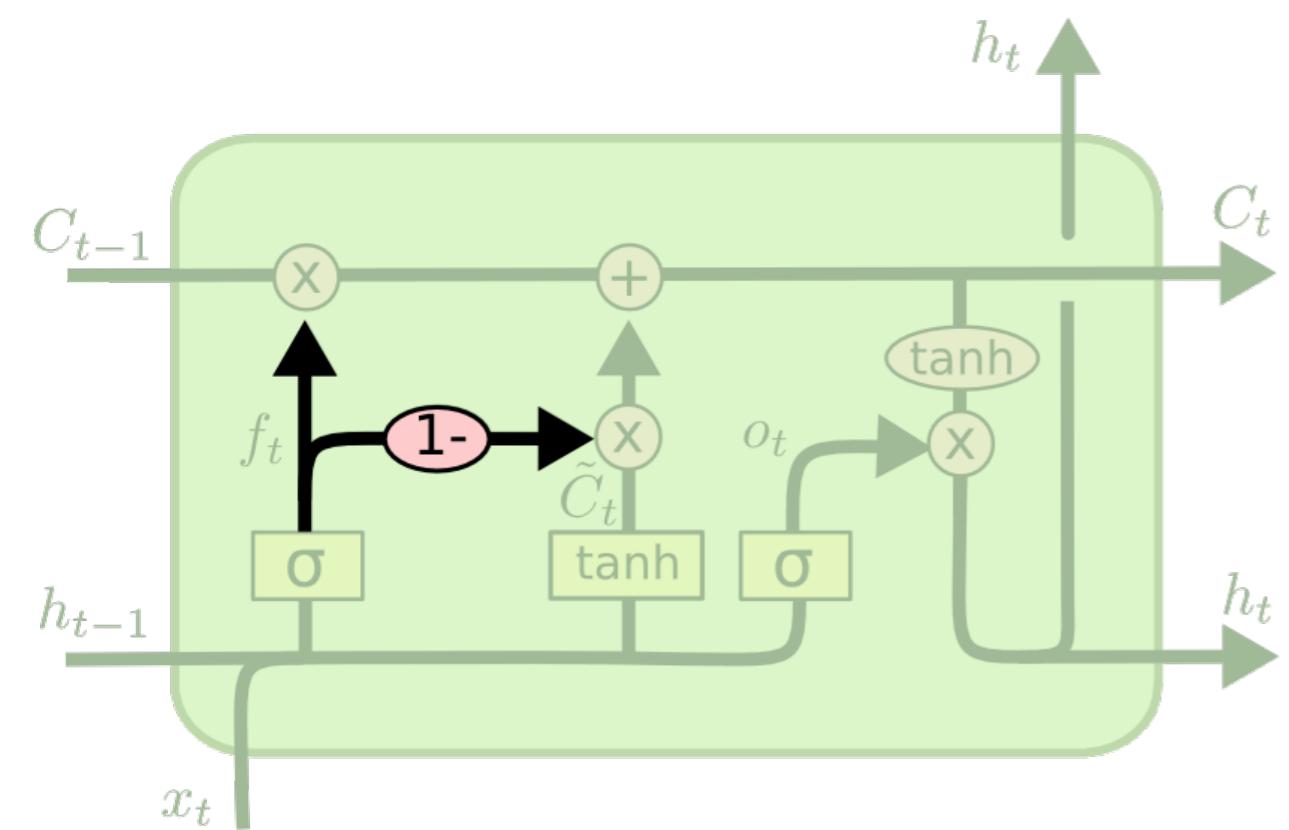


$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$

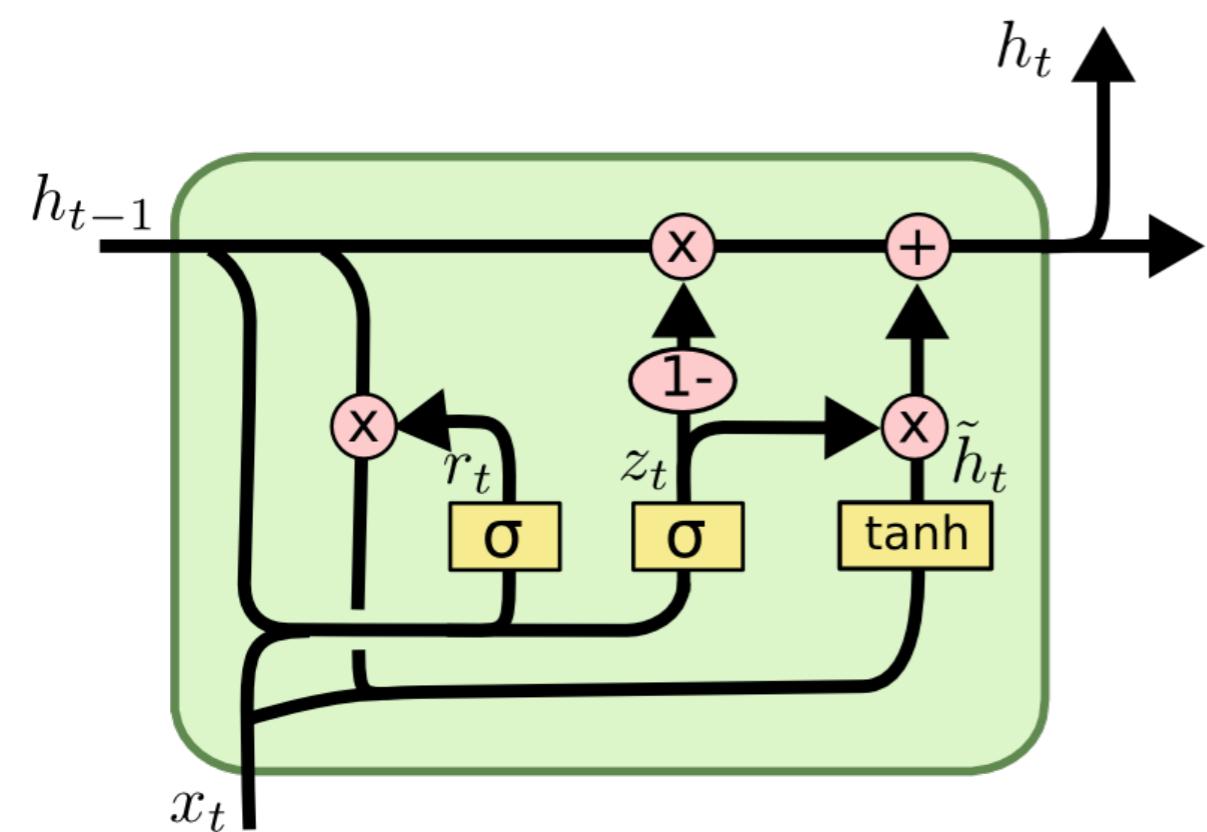
$$o_t = \sigma(W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$

Long Short Term Memories



$$C_t = f_t * C_{t-1} + (1 - f_t) * \tilde{C}_t$$

Gated Recurrent Unit



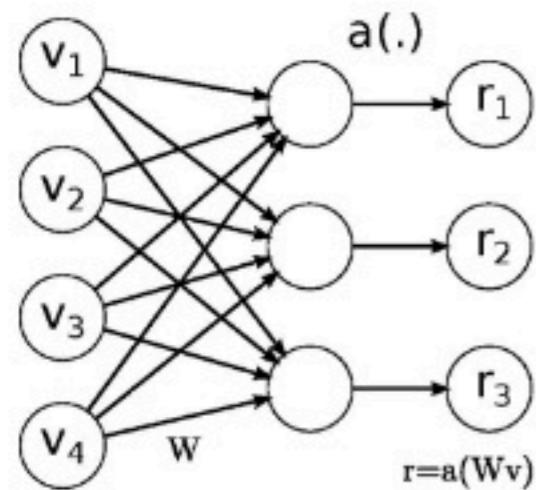
$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

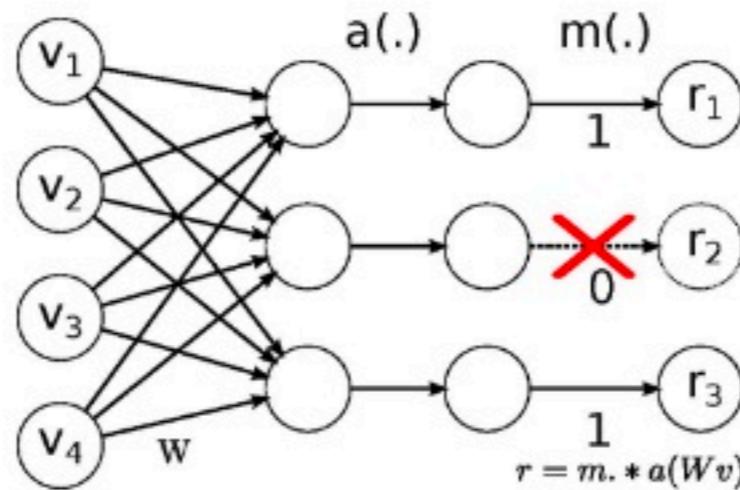
$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

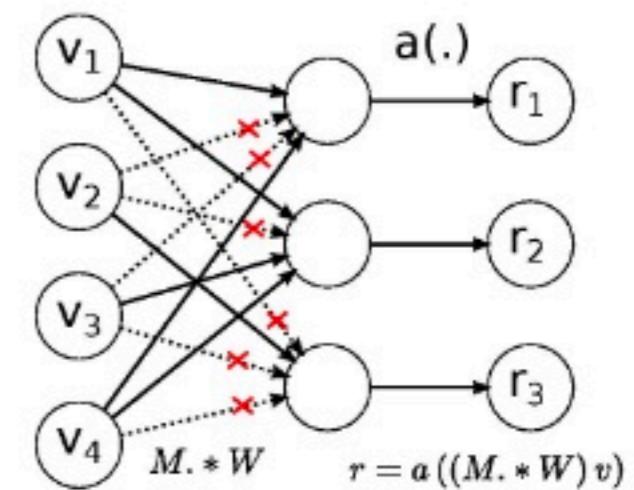
Dropout in RNN



No-Drop Network



DropOut Network



DropConnect Network

Layer Normalization

Batch Normalization

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_{ij}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_{ij} - \mu_j)^2$$

$$\hat{x}_{ij} = \frac{x_{ij} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Layer Normalization

Batch Normalization



Layer Normalization

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_{ij}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_{ij} - \mu_j)^2$$

$$\hat{x}_{ij} = \frac{x_{ij} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

$$\mu_i = \frac{1}{m} \sum_{j=1}^m x_{ij}$$

$$\sigma_i^2 = \frac{1}{m} \sum_{j=1}^m (x_{ij} - \mu_i)^2$$

$$\hat{x}_{ij} = \frac{x_{ij} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

Layer Normalization

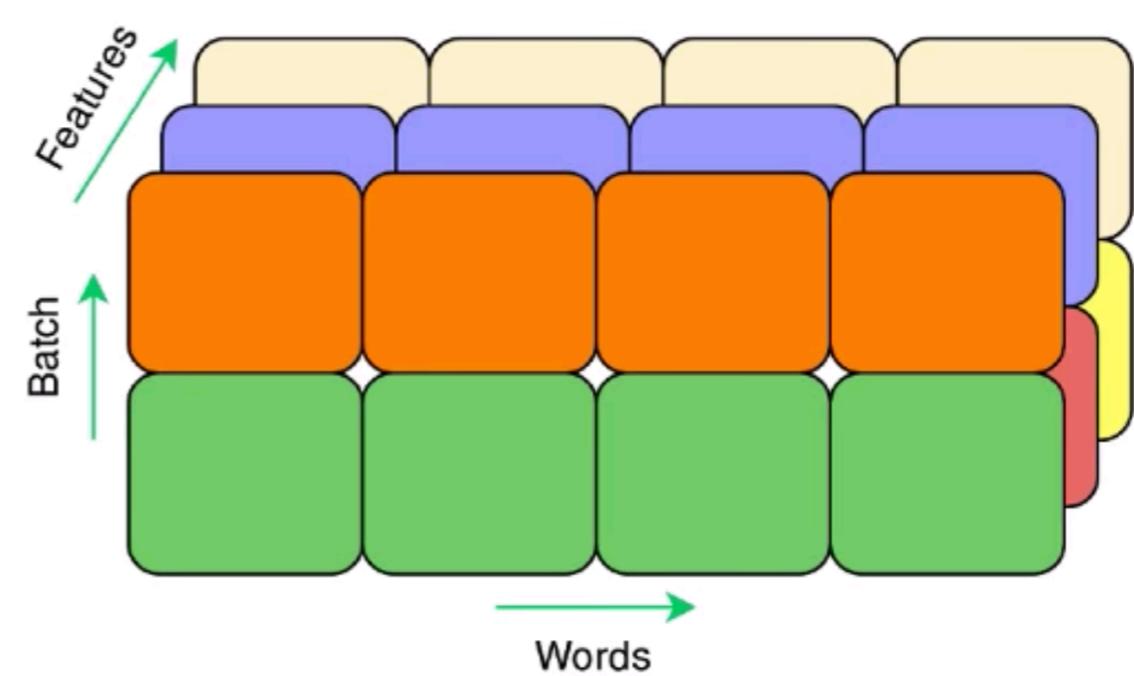
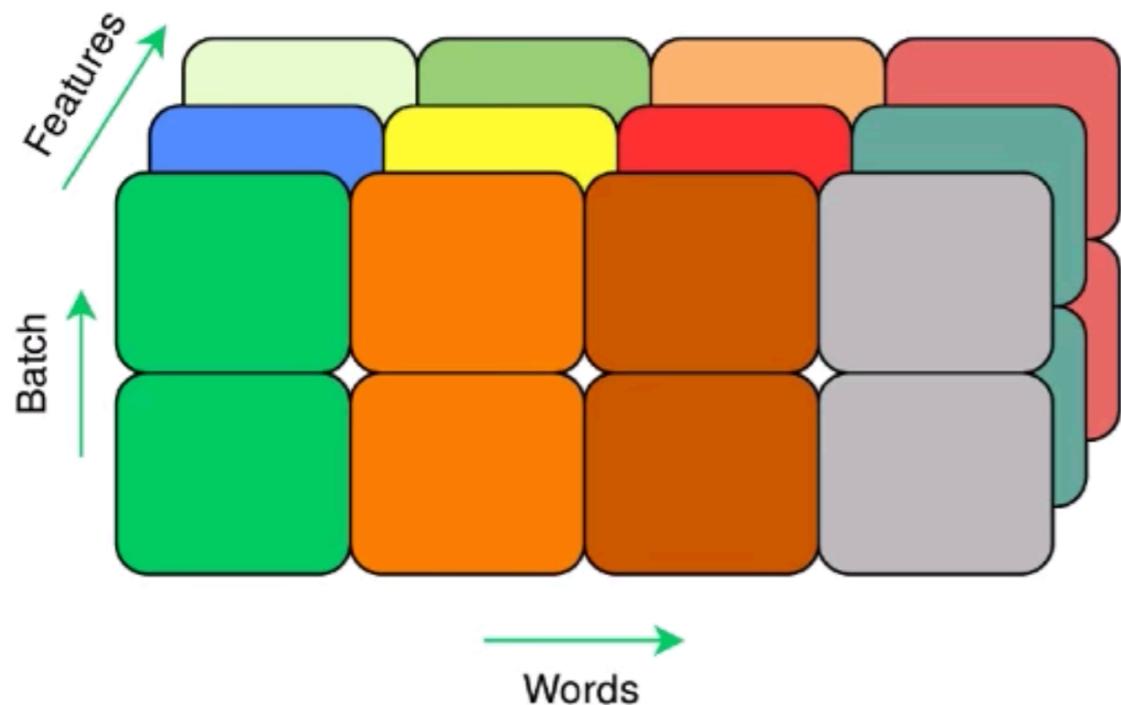
Batch Normalization



Layer Normalization

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_{ij}$$
$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_{ij} - \mu_j)^2$$
$$\hat{x}_{ij} = \frac{x_{ij} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

$$\mu_i = \frac{1}{m} \sum_{j=1}^m x_{ij}$$
$$\sigma_i^2 = \frac{1}{m} \sum_{j=1}^m (x_{ij} - \mu_i)^2$$
$$\hat{x}_{ij} = \frac{x_{ij} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$



Thanks for your Attention!

Boris Zubarev



@bobazooba