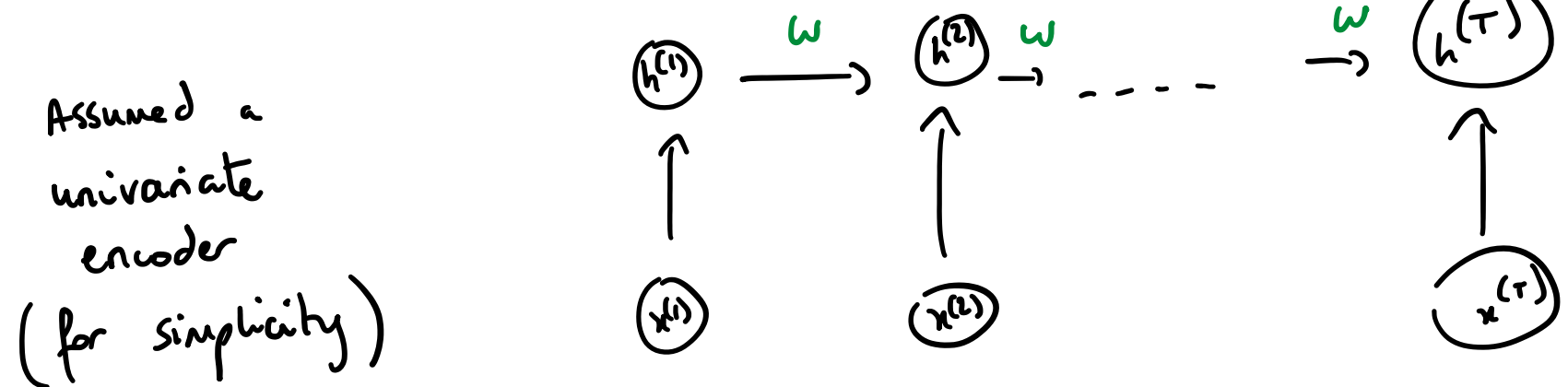


long temporal gap  
 => Difficult to learn long-distance dependencies

Recall backpropagation through time:



Assume  $\bar{h}^{(T)}$  has already been computed (by backpropagating through the decoder)

Encoder backpropagation rules:

from video on backpropagation through time

$$\begin{cases} \bar{h}^{(t)} = \bar{z}^{(t+1)} w \\ \bar{z}^{(t)} = \bar{h}^{(t)} \phi'(z^{(t)}) \end{cases}$$

By iteration:

$$\begin{aligned} \bar{h}^{(1)} &= \bar{z}^{(2)} w = \bar{h}^{(2)} \phi'(z^{(2)}) w \\ &= \bar{z}^{(3)} w \phi'(z^{(2)}) w \\ &= \bar{h}^{(3)} \phi'(z^{(3)}) \phi'(z^{(2)}) w^2 \end{aligned}$$

$$\begin{aligned} &\vdots \\ &= w^{T-1} \phi'(z^{(T)}) \dots \phi'(z^{(2)}) \bar{h}^{(T)} \end{aligned}$$

computed by backpropagating through the decoder.

for simplicity, if  $\phi$  is linear:  $\frac{\partial \bar{h}^{(T)}}{\partial \bar{h}^{(1)}} = w^{T-1}$

if  $w$  is not close to 1,  $w^{T-1}$  will either explode or vanish

$|w| > 1$   $\Rightarrow w^{T-1} = 117.4$  (if  $w = -1.1, T = 50$ )  
 $0 < |w| < 1$   $\Rightarrow w^{T-1} = 0.0015$  (if  $w = 0.9, T = 50$ )

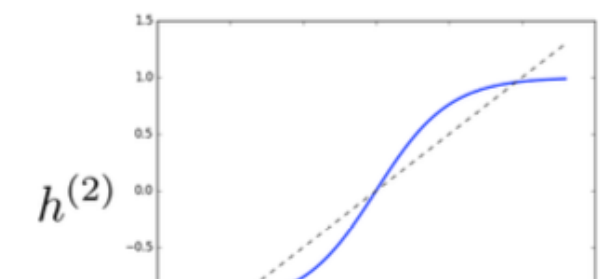
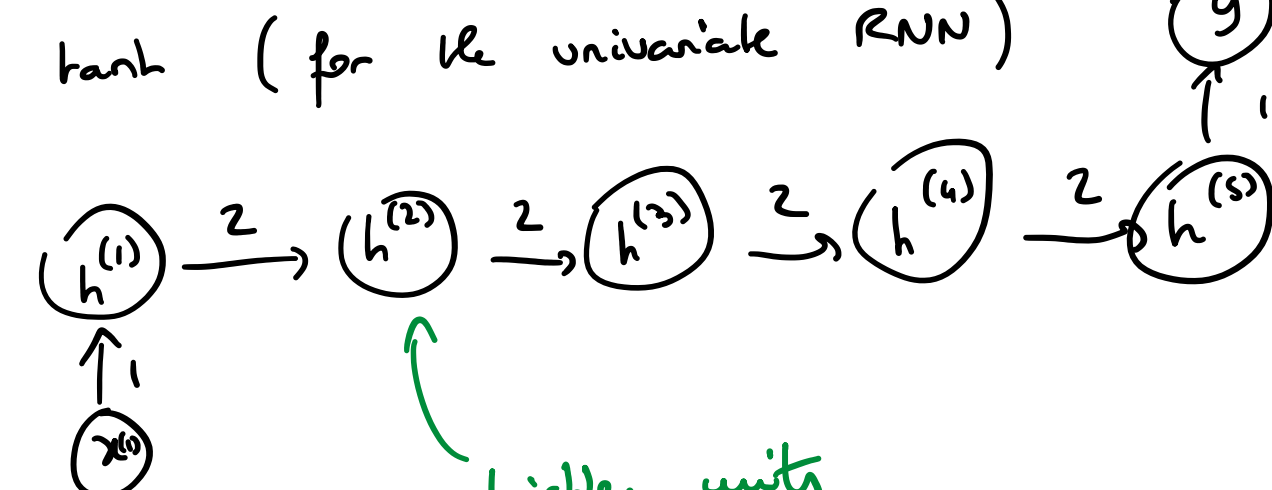
In the multivariate case:

$$\frac{\partial \vec{h}^{(T)}}{\partial \vec{h}^{(1)}} = \frac{\partial \vec{h}^{(T)}}{\partial \vec{h}^{(T-1)}} \dots \frac{\partial \vec{h}^{(2)}}{\partial \vec{h}^{(1)}} = W^{T-1}$$

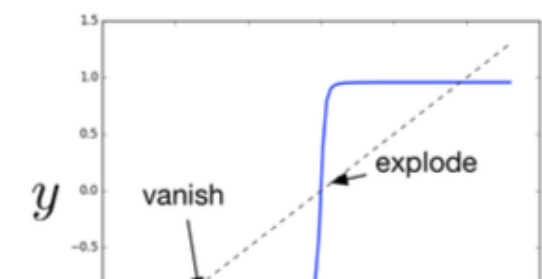
assume that activation is linear

Exploding/vanishing gradients

Example with a tanh (for the univariate RNN)



One application of the tanh



Five applications of the tanh

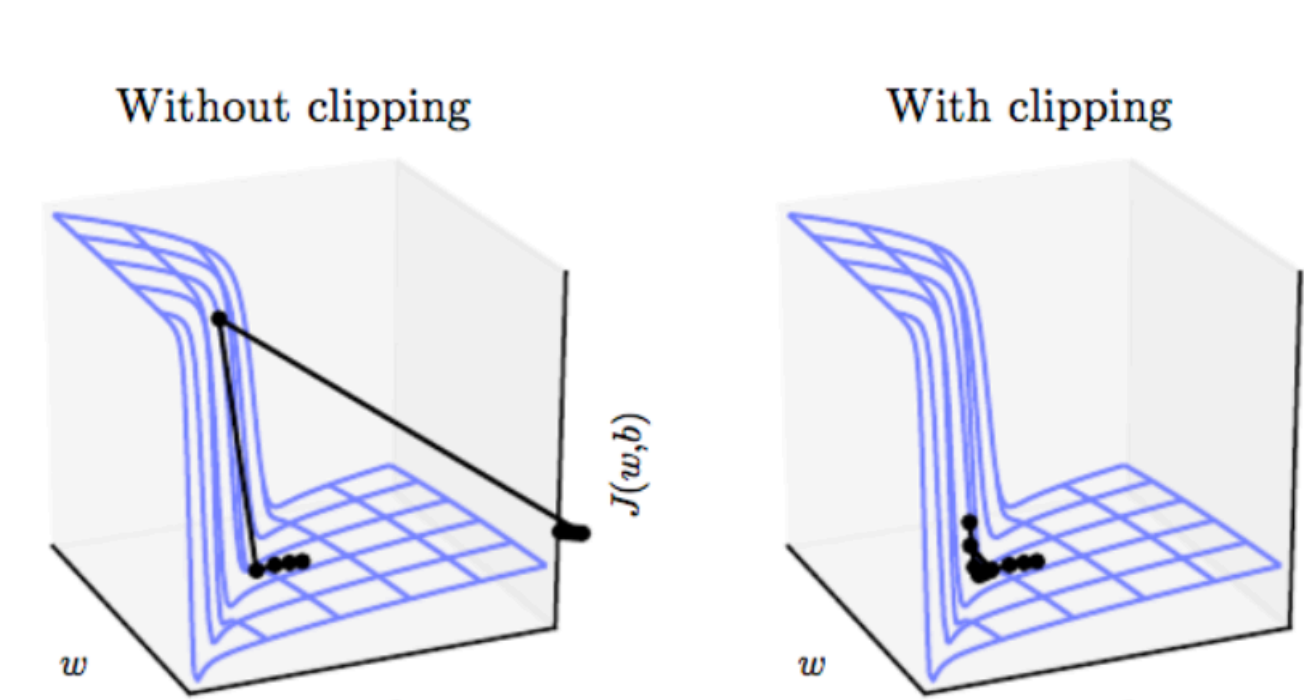
Solution 1: Gradient clipping

Prevent gradients from exploding by rescaling them:

$$\vec{g} \leftarrow \frac{\vec{g}}{\|\vec{g}\|}$$

constant used to rescale gradients

Biases training procedure but can make it more stable.



— Goodfellow et al., Deep Learning

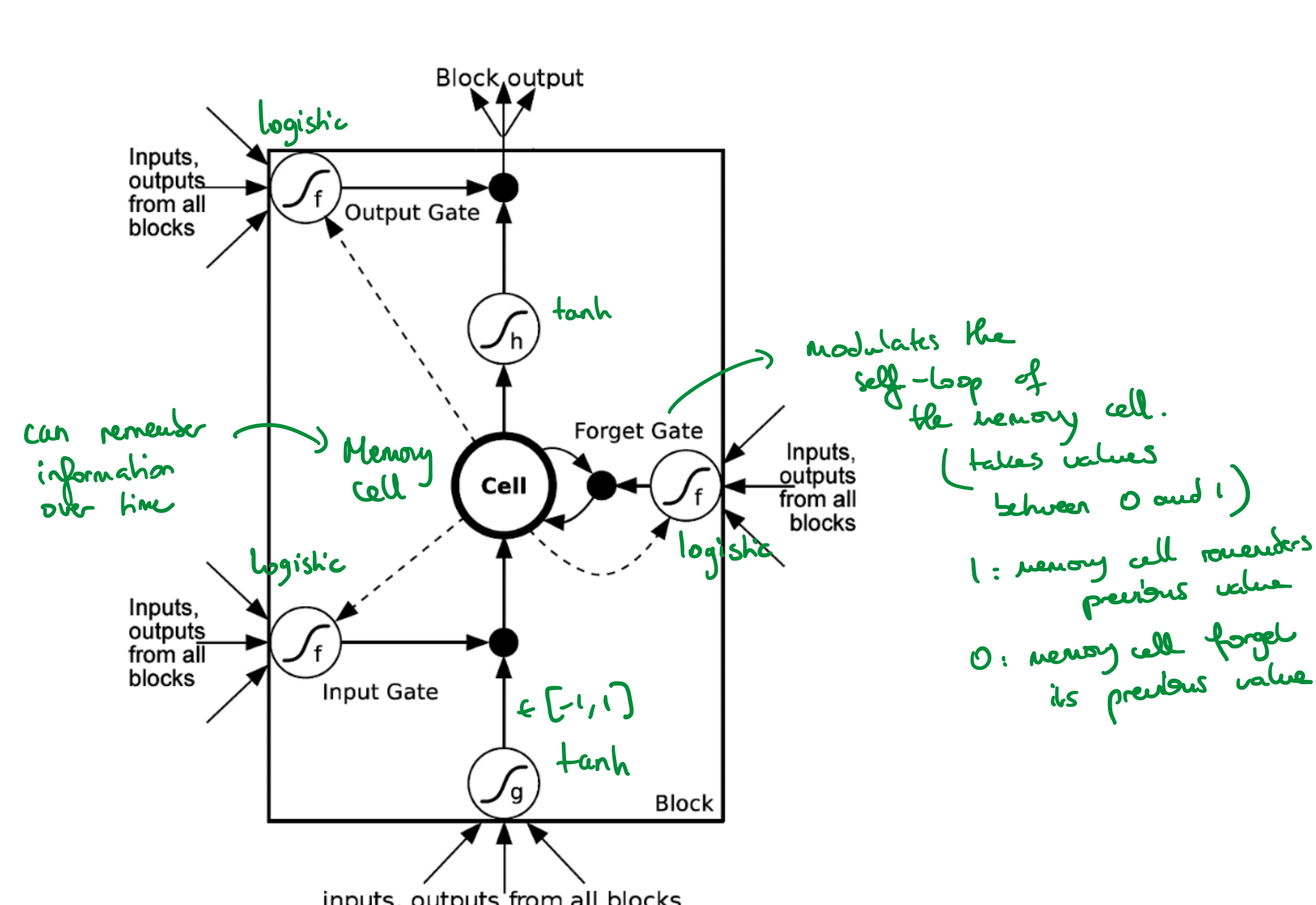
Solution 2: Long-Term Short-Term Memory (LSTM)

Goal: remember information over long time periods.

- long-term memory: weights
- short-term memory: activation

Architecture uses an LSTM unit which is called a block

Each hidden layer of an LSTM architecture is made up of many of these blocks



Input gate	Forget gate	Behavior of the memory cell
0	1	remember the previous value
1	1	add to previous value
0	0	erase the value
1	0	overwrite the value.

Implementation  $\rightarrow$  values of the memory cell the different gates

input gates of all blocks  $\rightarrow \begin{pmatrix} \vec{i}^{(t)} \\ \vec{f}^{(t)} \\ \vec{o}^{(t)} \\ \vec{g}^{(t)} \end{pmatrix}$

forget gate  $\rightarrow$

output gate  $\rightarrow$

weighted input of all blocks  $\rightarrow$

$$\vec{c}^{(t)} = \vec{f}^{(t)} \cdot \vec{c}^{(t-1)} + \vec{i}^{(t)} \cdot \vec{g}^{(t)}$$

memory cells for all blocks  $\rightarrow \vec{h}^{(t)} = \vec{o}^{(t)} \cdot \tanh(\vec{c}^{(t)})$

each of the gates and the weighted input computations involve a set of independent weights