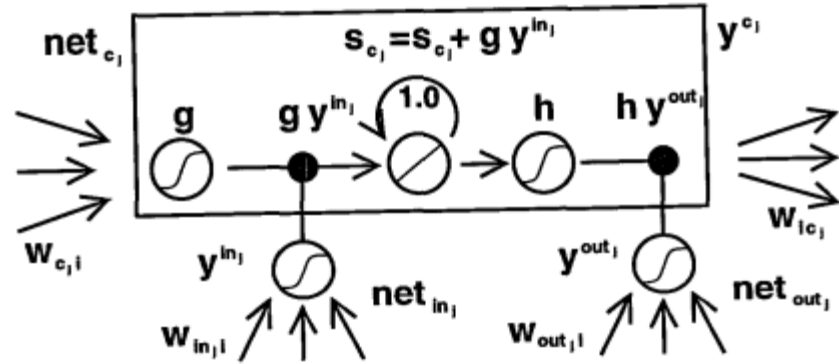


Long Short Term Memory Networks

Brian Cheung

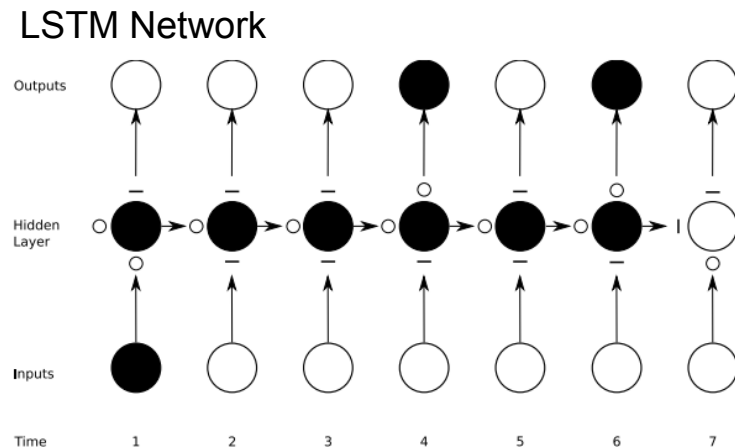
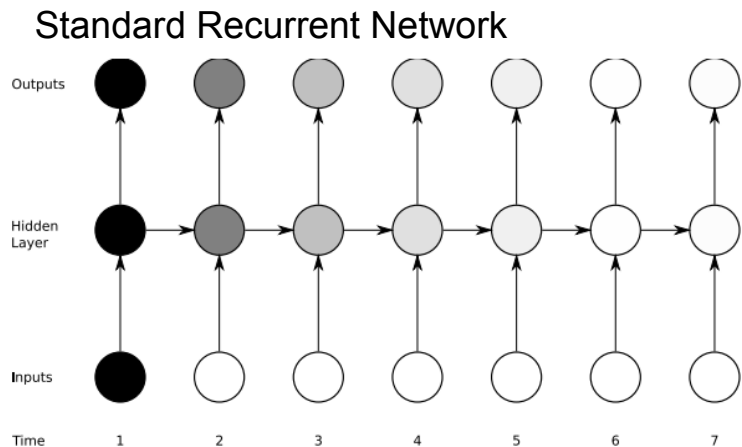
LSTM



- Proposed by Sepp Hochreiter in 1997
- Originally used approximate error gradient with Real Time Recurrent Learning and truncated backpropagation through time
- Used for processing long range contextual information

Hochreiter & Schmidhuber 1997

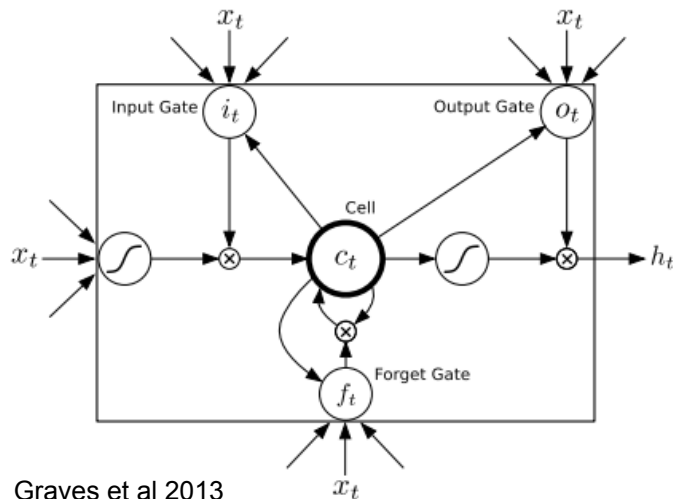
LSTMs reduce vanishing gradient problem



Graves et al 2013

- The darker the shade, the greater the sensitivity
- The sensitivity decays exponentially over time as new inputs overwrite the activation of hidden unit and the network 'forgets' the first input

LSTMs reduce vanishing gradient problem



Graves et al 2013

- Memory cells and gating units allow information to be stored for long periods of time.
- Memory cells are additive in time
 - Gradients also additive in time which alleviates vanishing gradient

Backpropagation through time

- Derivative of objective function, O , with respect to linear activations, a

$$a_j = \sum_{i=1} w_{ij} b_i$$








$$\delta_j^t \stackrel{\text{def}}{=} \frac{\partial O}{\partial a_j^t}$$

- Gradient descent weight update

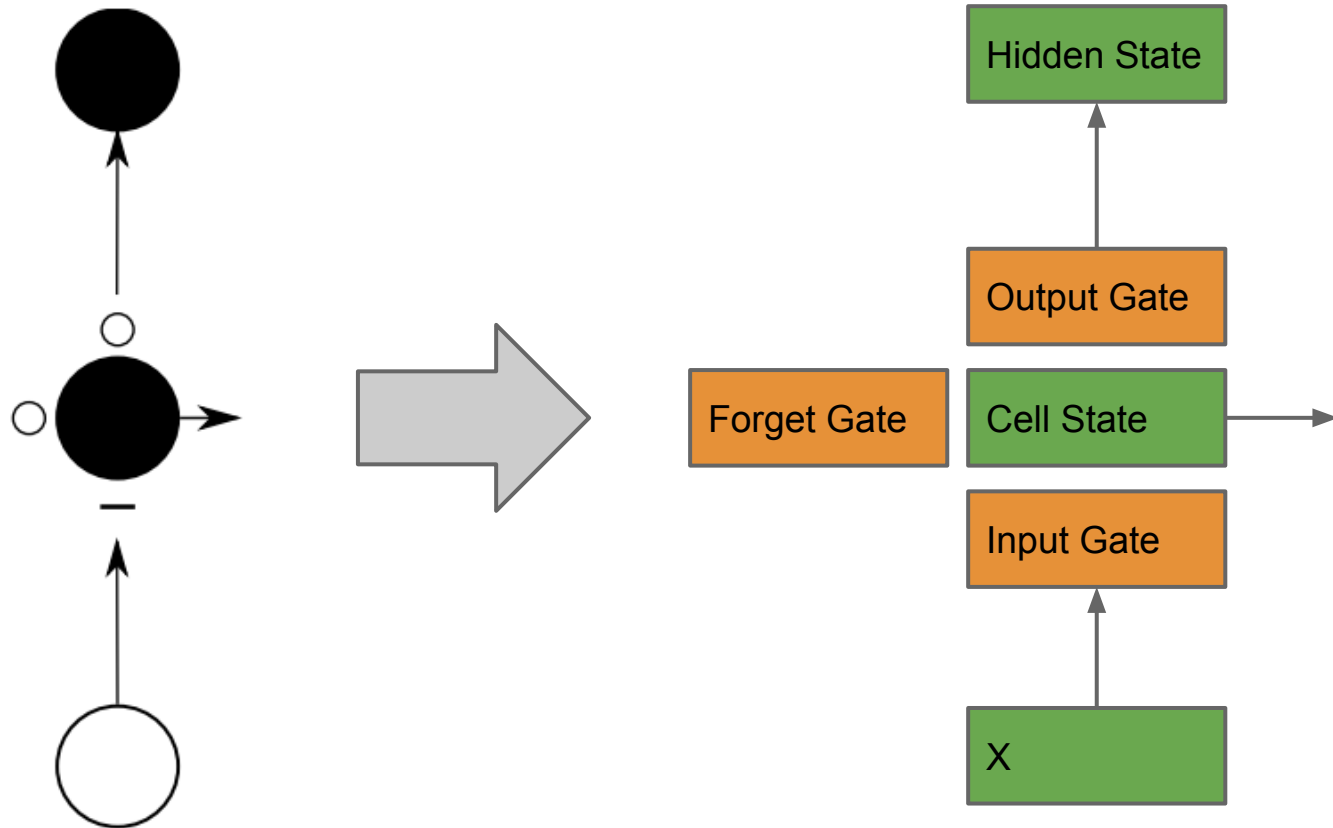
$$\frac{\partial O}{\partial w_{ij}} = \sum_{t=1}^T \frac{\partial O}{\partial a_j^t} \frac{\partial a_j^t}{\partial w_{ij}} = \sum_{t=1}^T \delta_j^t b_i^t$$

$$\Delta \mathbf{w}(n) = -\alpha \frac{\partial O}{\partial \mathbf{w}(n)}$$

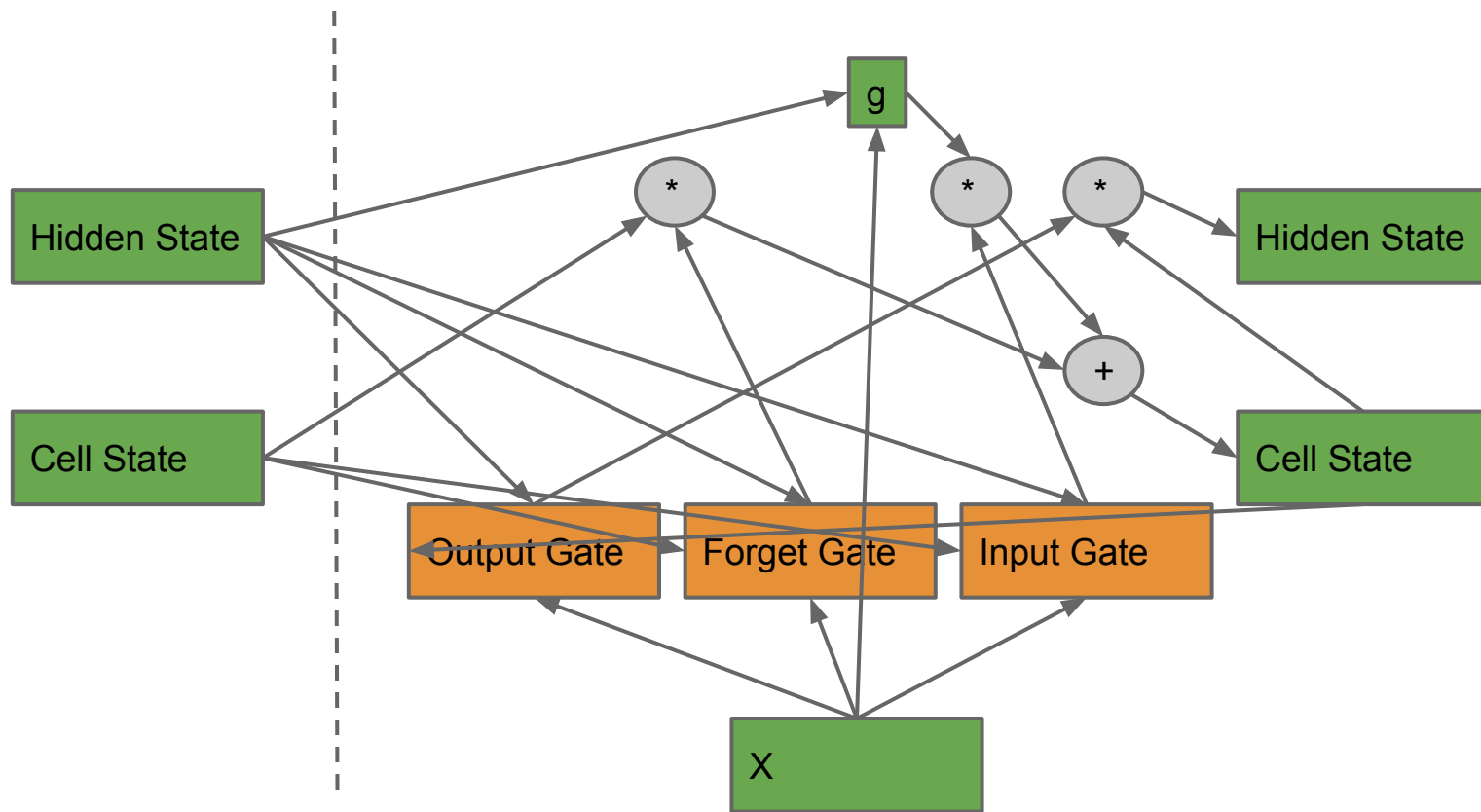
Anatomy of an LSTM node

x		Input to the LSTM node at a particular time point (character, phoneme, word)
i		Determines whether the input will be stored in the memory cell
ϕ		Determines whether current contents of memory will be forgotten (erased)
ω		Determines if current memory contents will be output
c		Memory Cell
h		Output of LSTM node
g		Input nonlinearity of Memory Cell (i.e. tanh, sigmoid)

Anatomy of an LSTM node



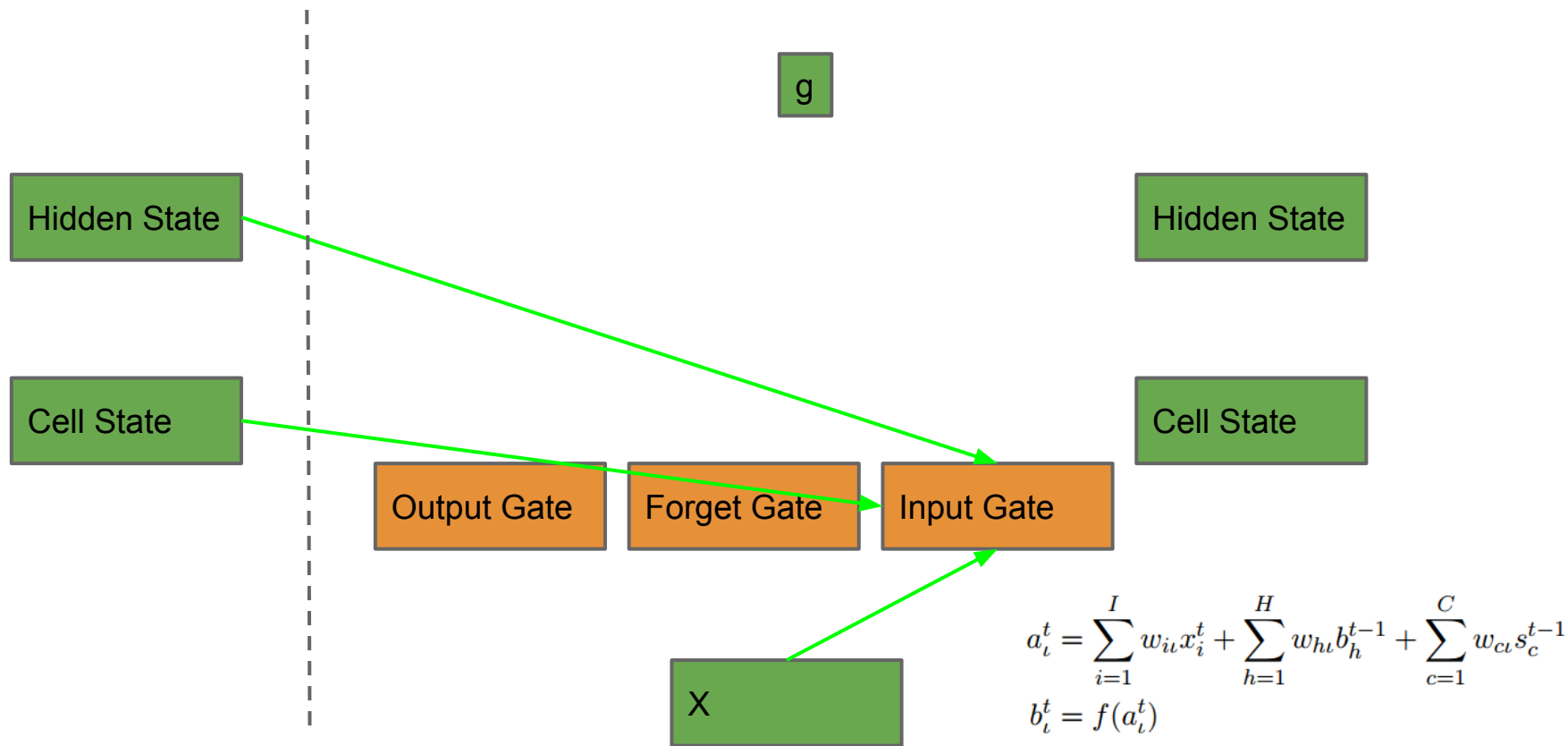
LSTM Forward Propagation



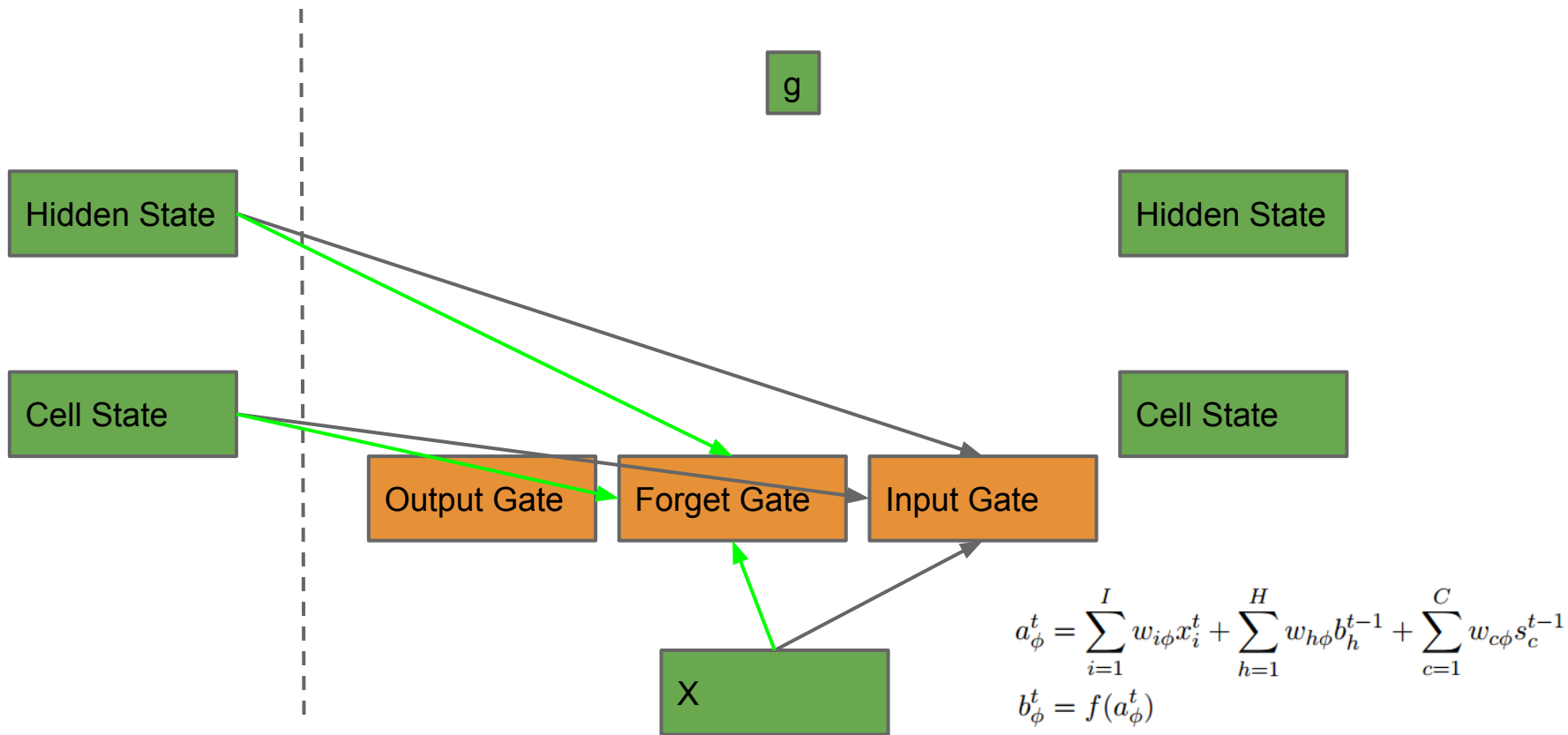
Time: $t-1$

t

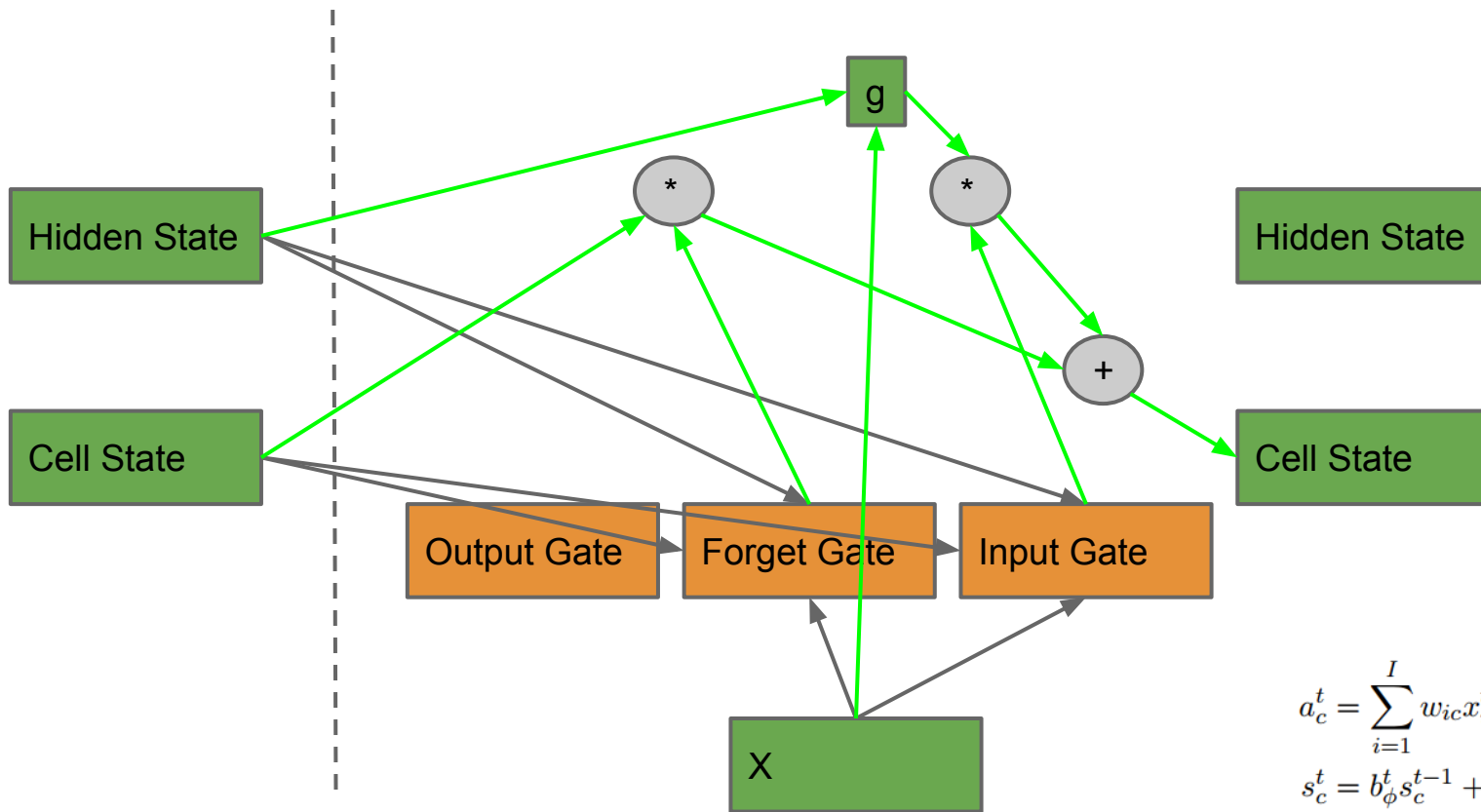
1. Input Gate



2. Forget Gate

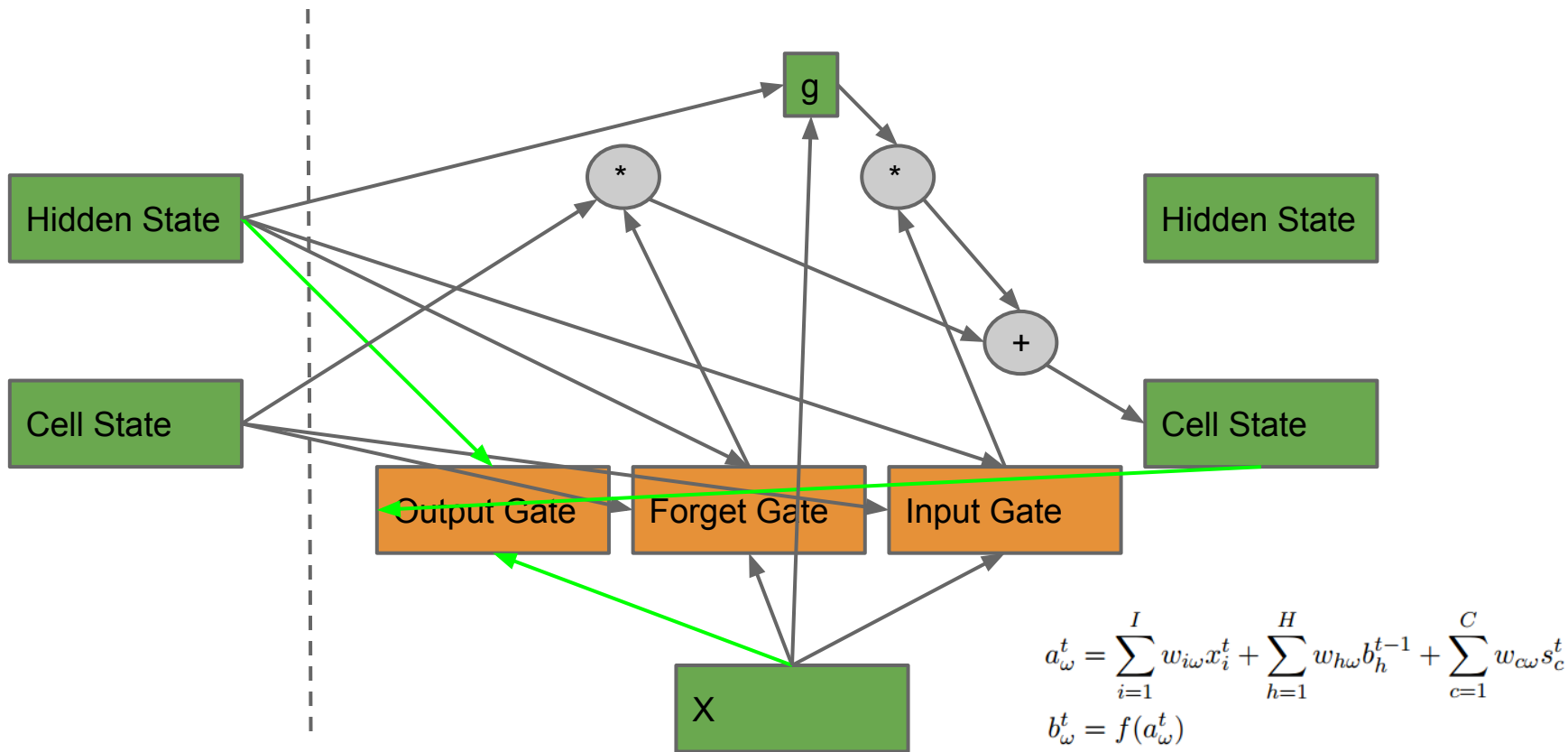


3. Cell State

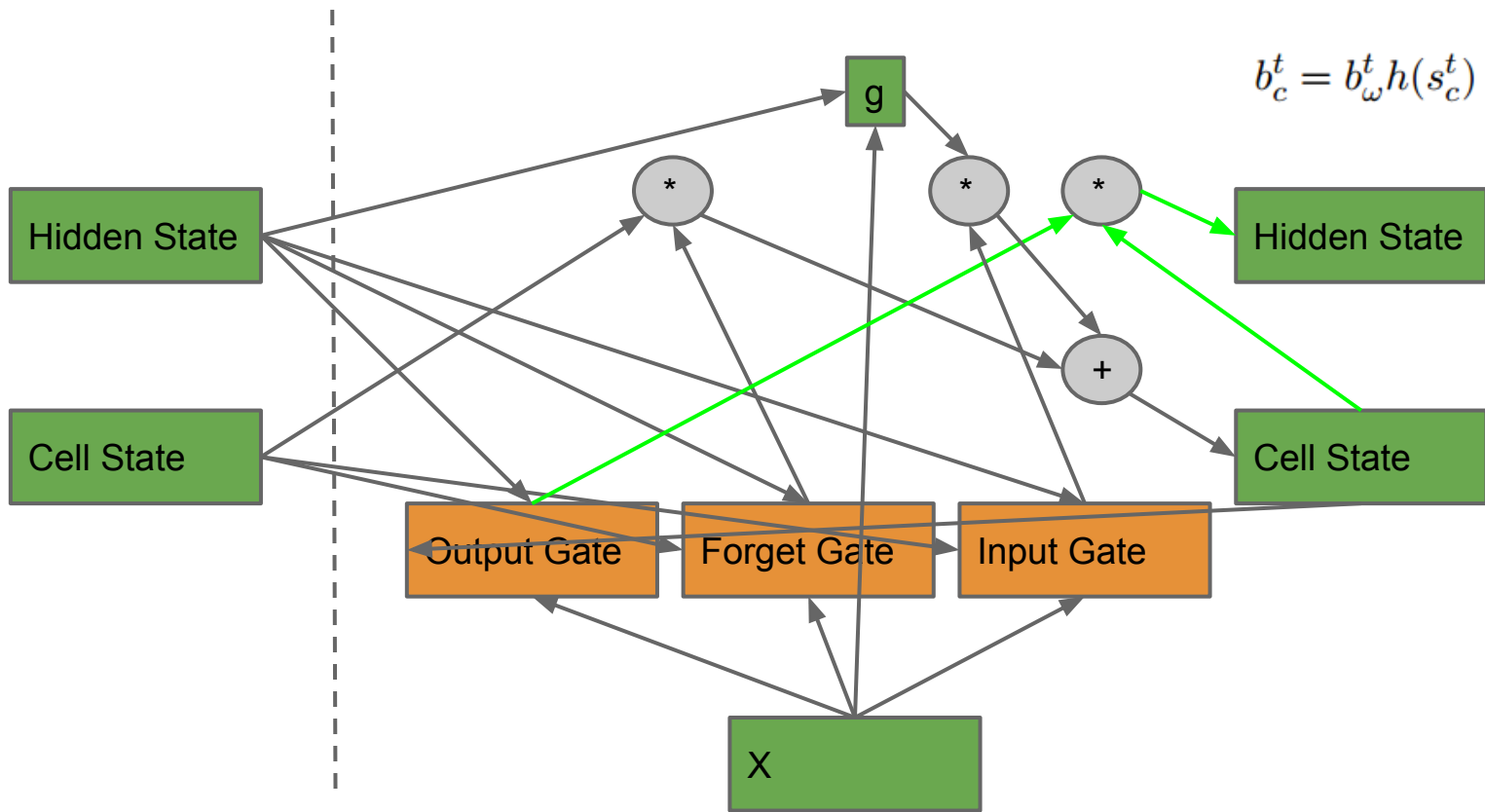


$$\begin{aligned} a_c^t &= \sum_{i=1}^I w_{ic} x_i^t + \sum_{h=1}^H w_{hc} b_h^{t-1} \\ s_c^t &= b_\phi^t s_c^{t-1} + b_l^t g(a_c^t) \end{aligned}$$

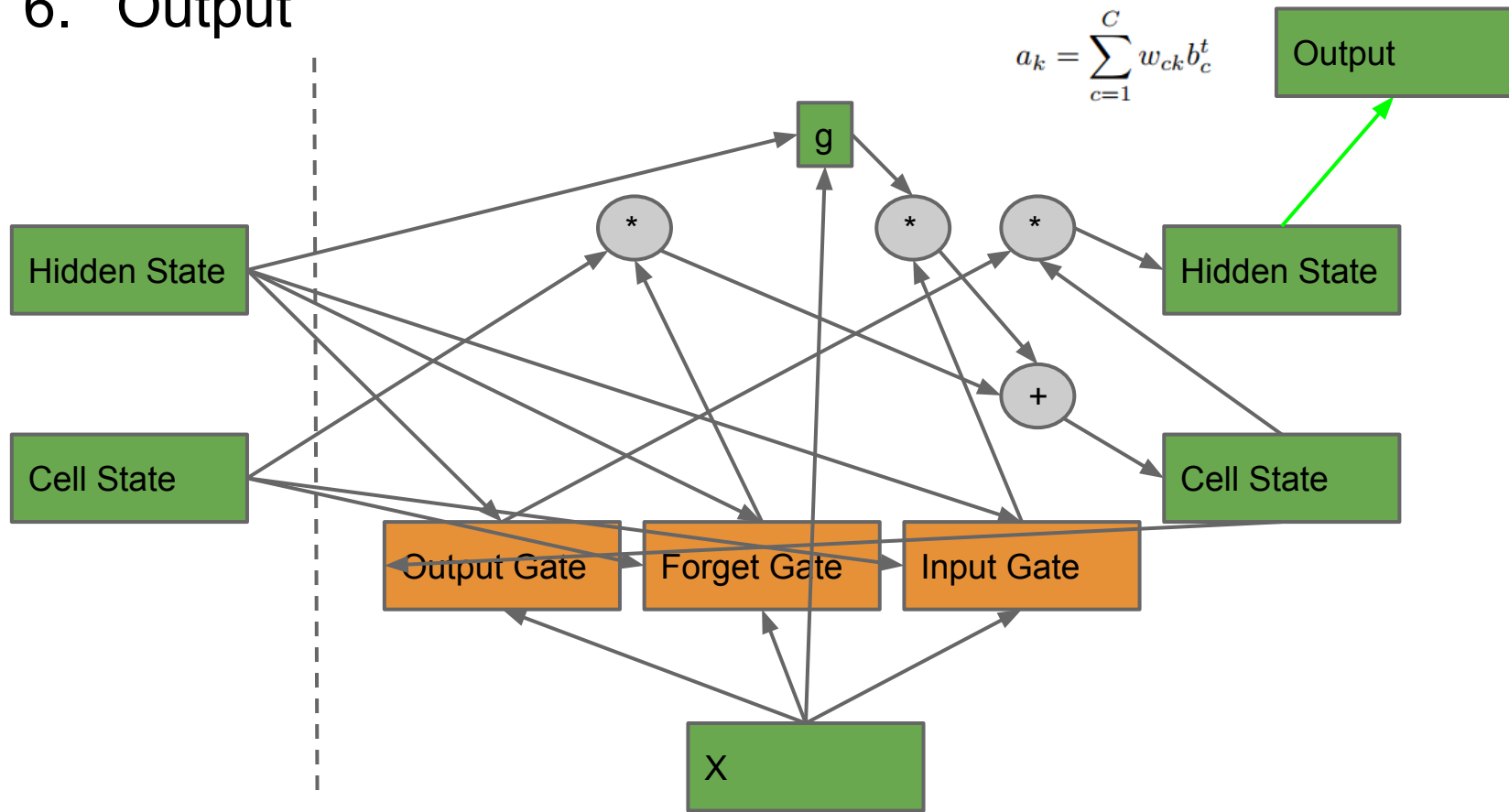
4. Output Gate



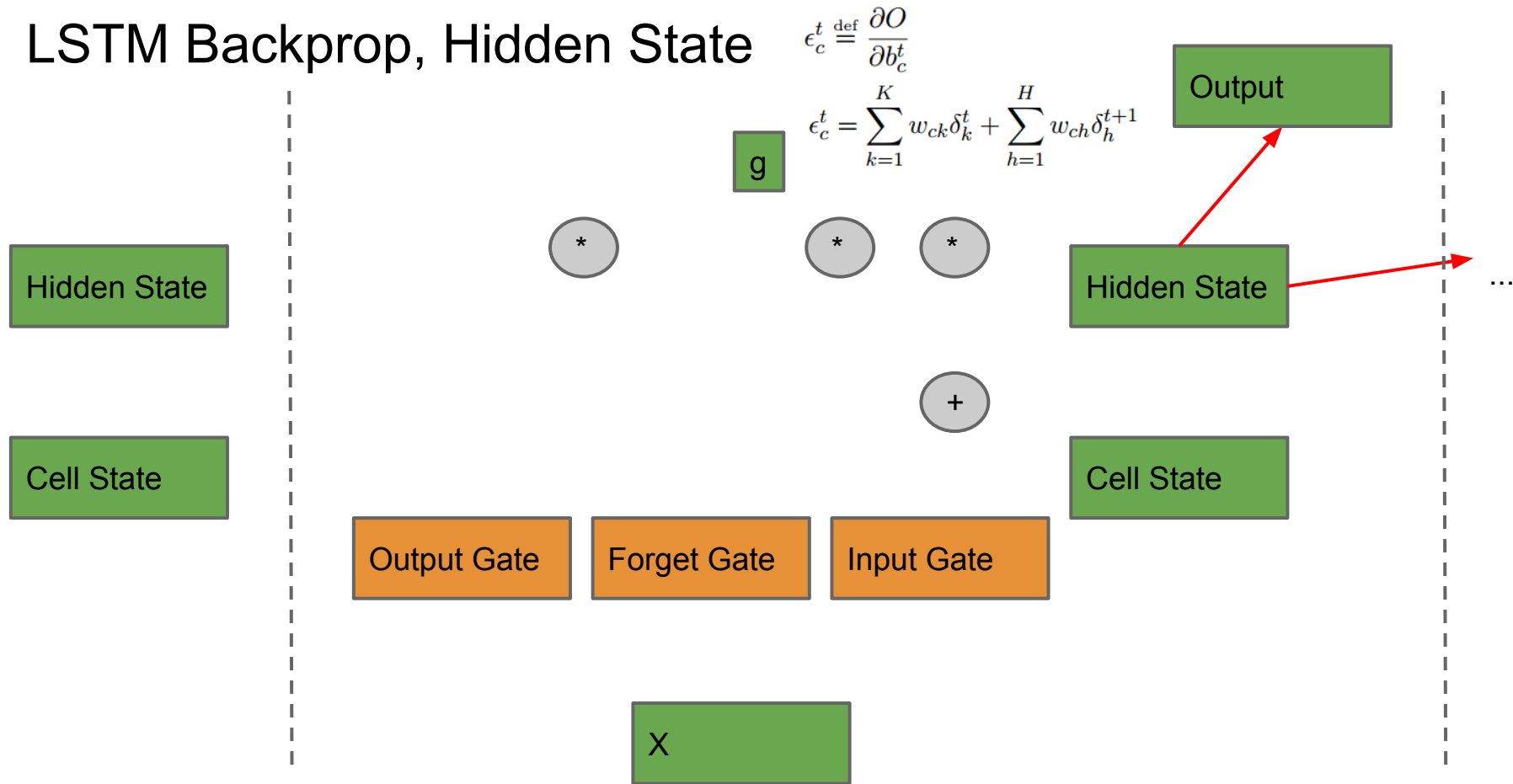
5. Hidden State



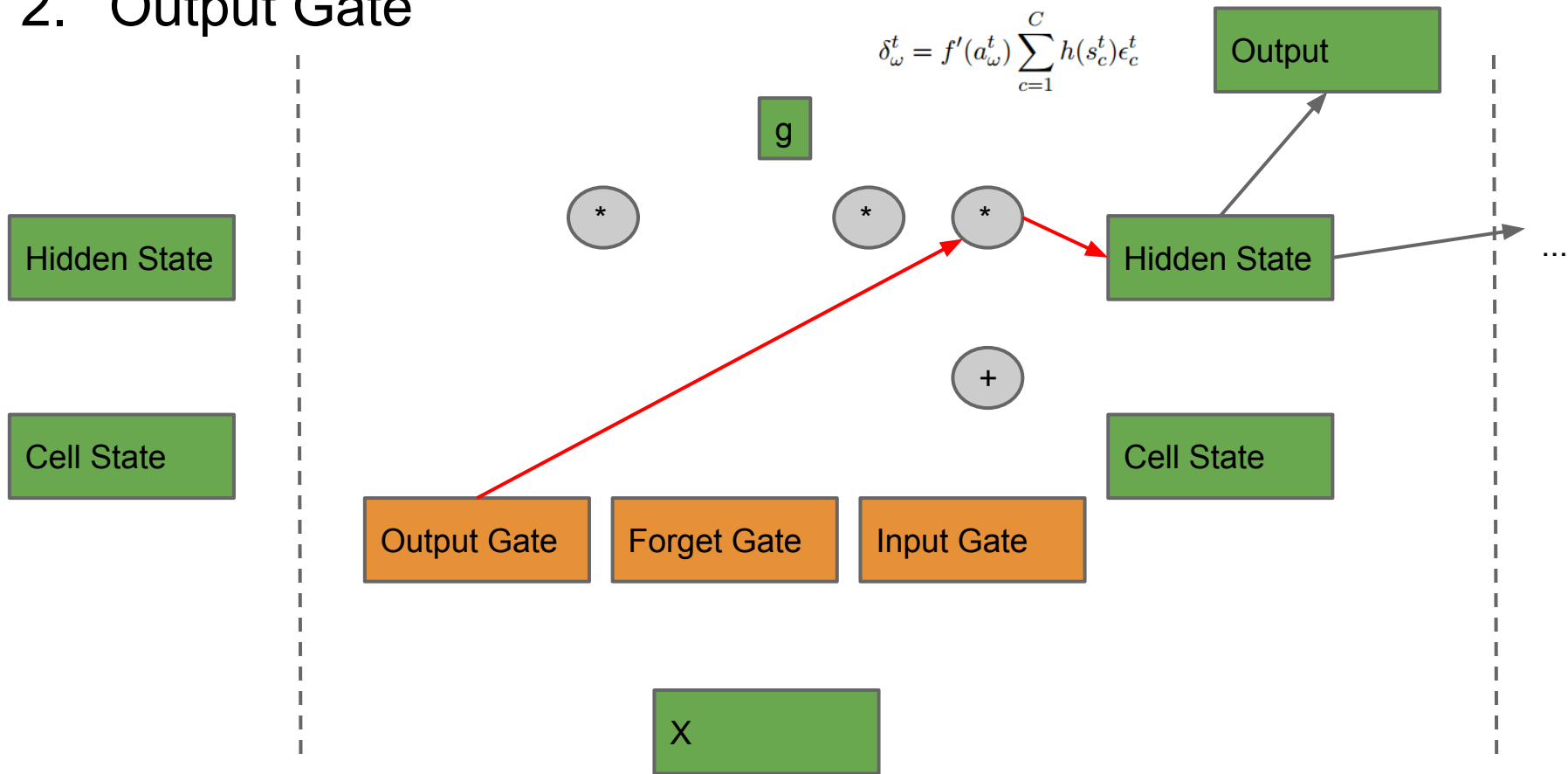
6. Output



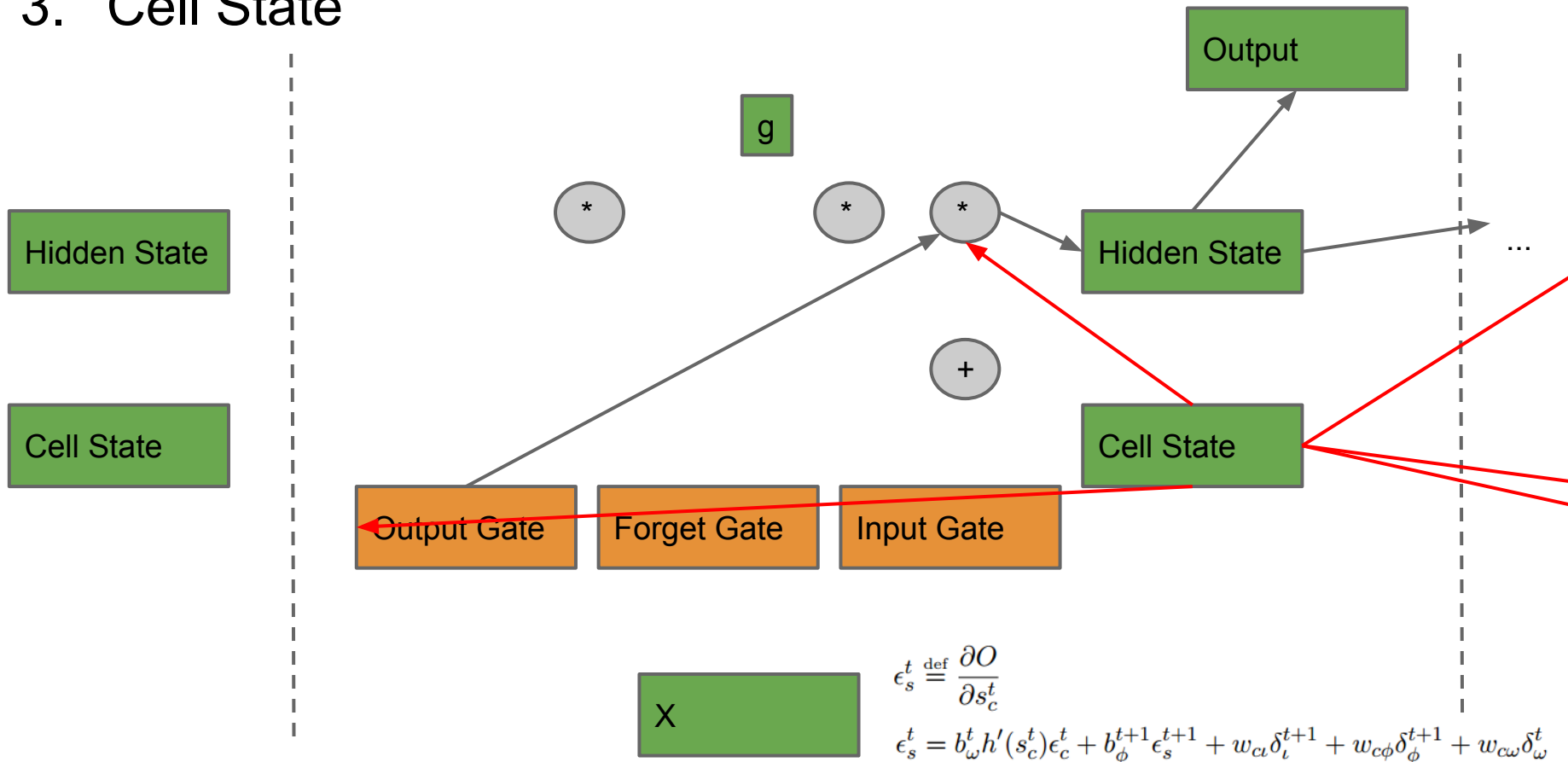
LSTM Backprop, Hidden State



2. Output Gate



3. Cell State



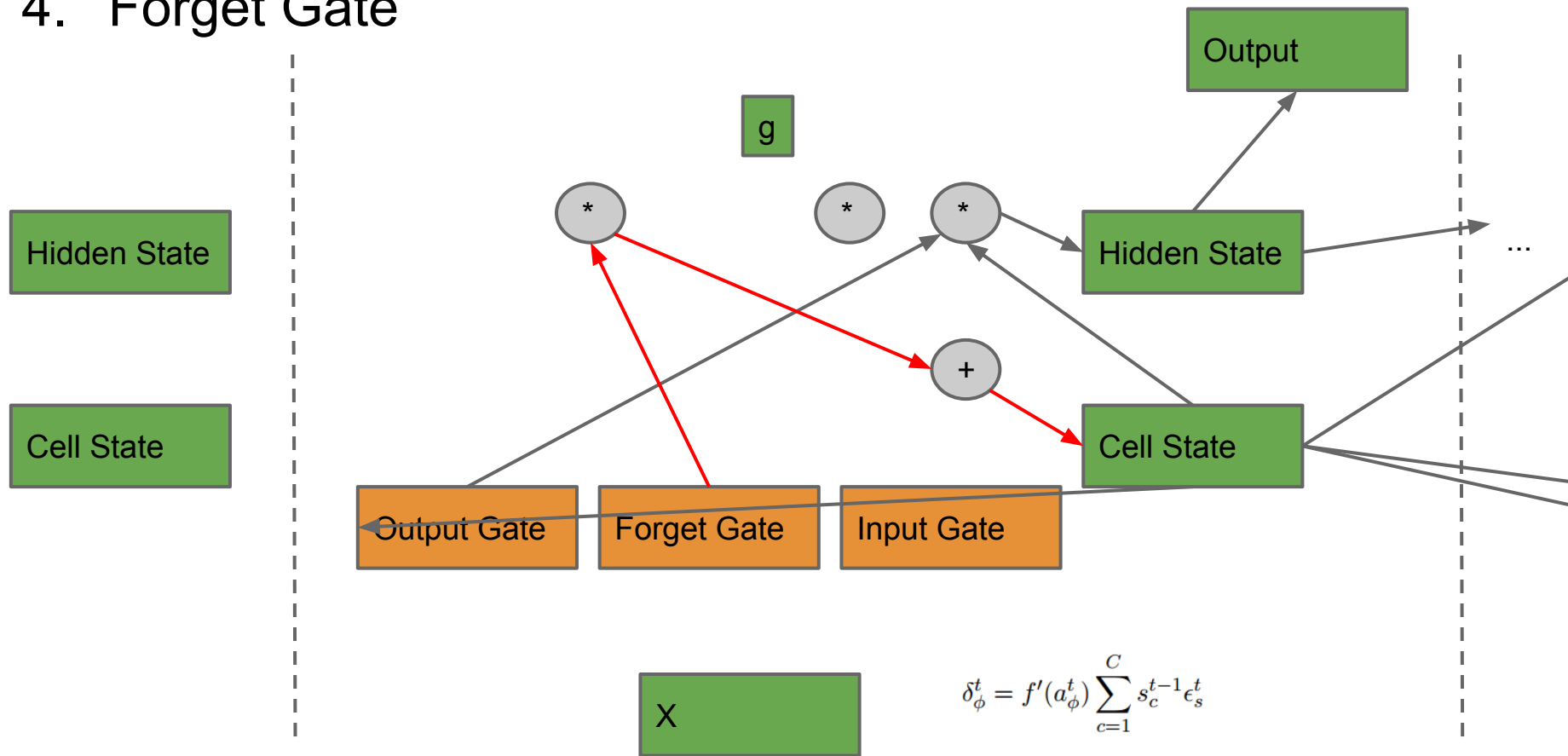
$$X$$

$$\epsilon_s^t \stackrel{\text{def}}{=} \frac{\partial O}{\partial s_c^t}$$

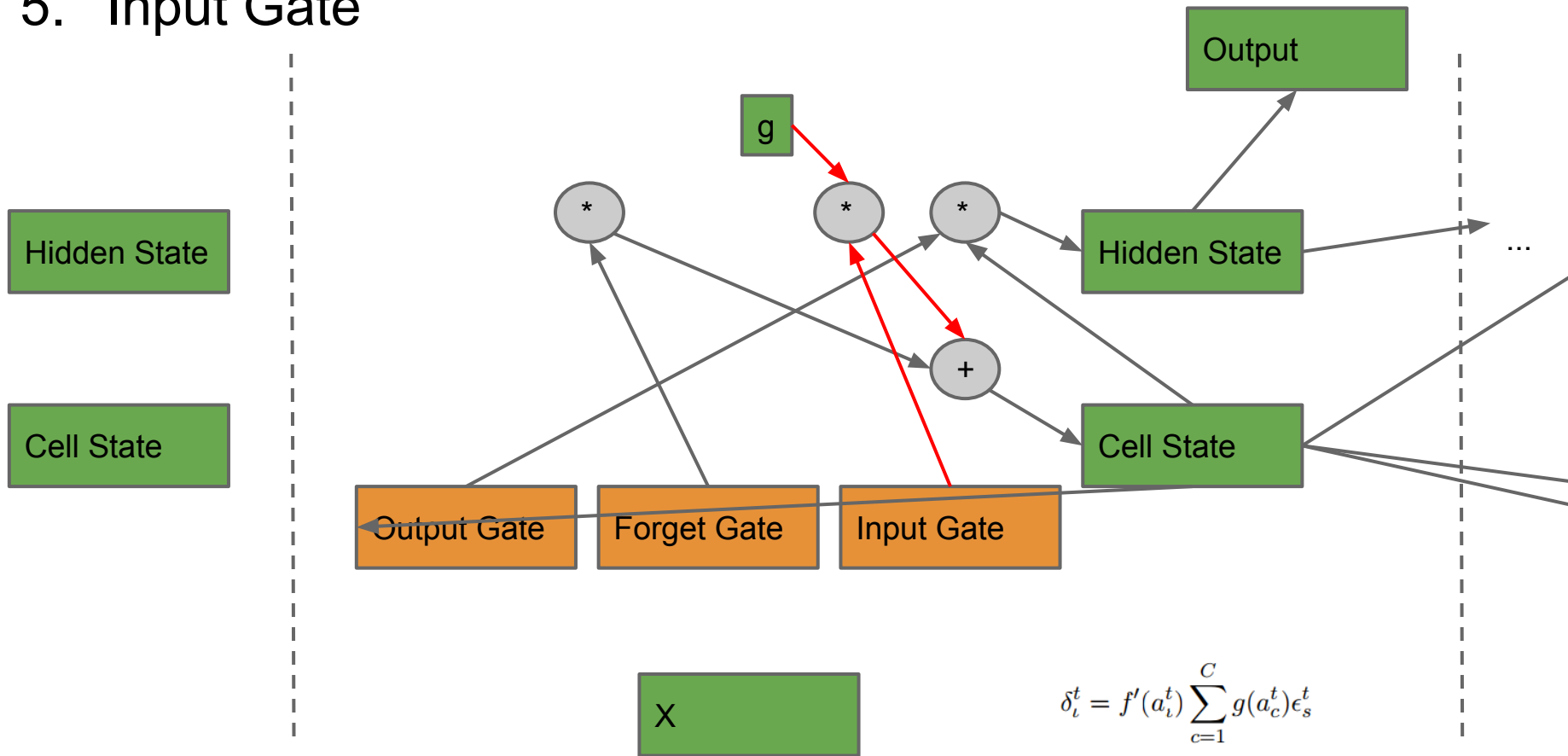
$$\epsilon_s^t = b_\omega^t h'(s_c^t) \epsilon_c^t + b_\phi^{t+1} \epsilon_s^{t+1} + w_{ci} \delta_i^{t+1} + w_{c\phi} \delta_\phi^{t+1} + w_{c\omega} \delta_\omega^t$$

$$\delta_c^t = b_i^t g'(a_c^t) \epsilon_s^t$$

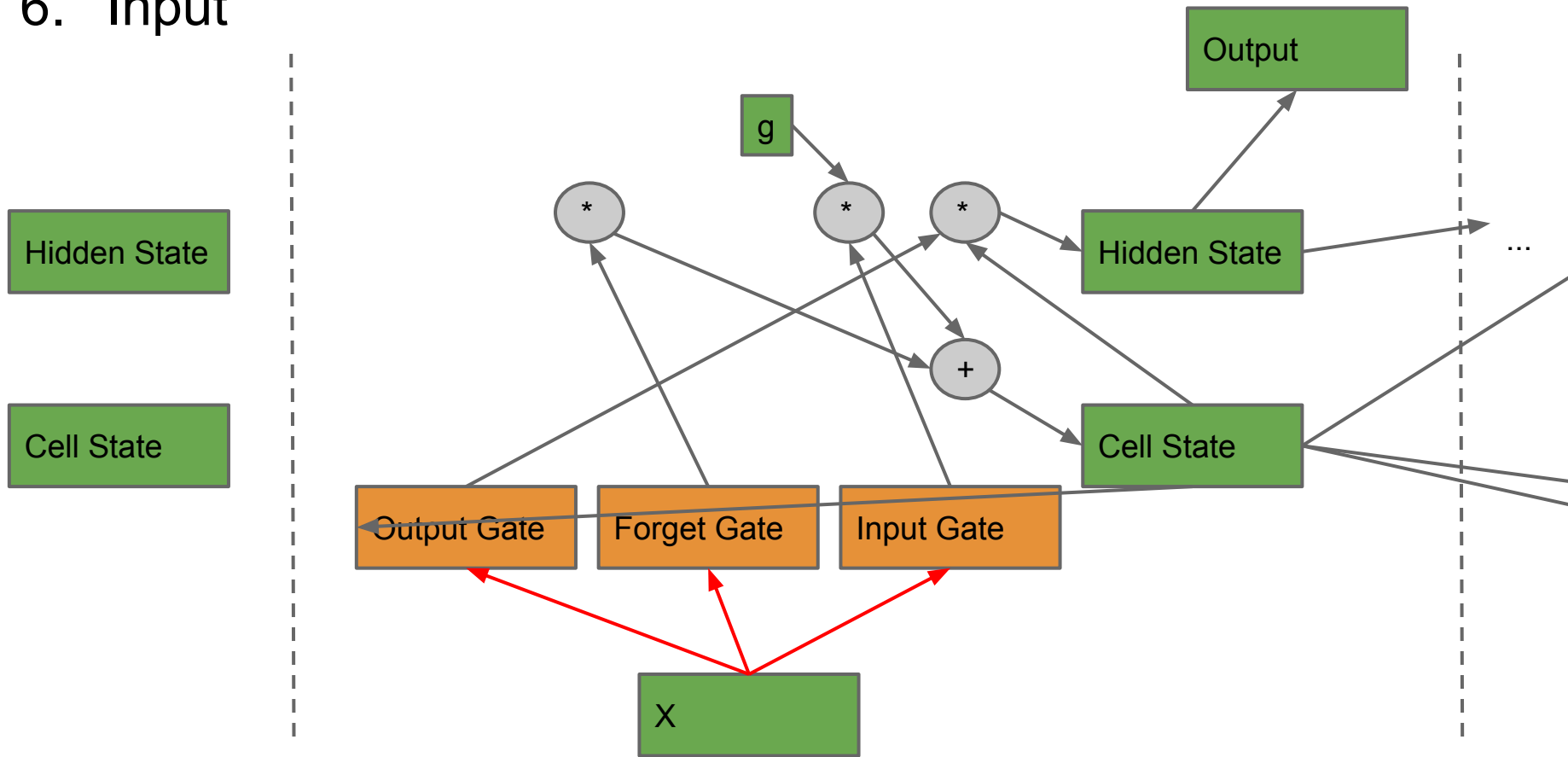
4. Forget Gate



5. Input Gate

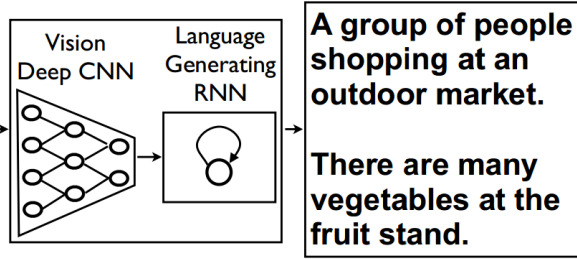


6. Input

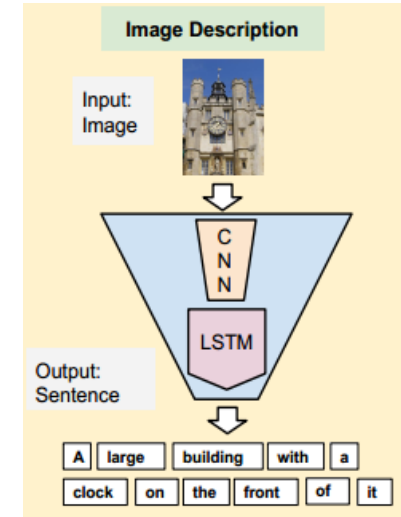


LSTM Applications

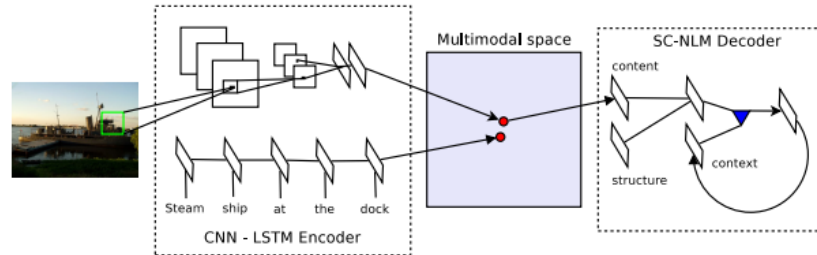
Image Captioning



Vinyals et al 2014



Donahue et al 2014



Kiros et al 2014

LSTM Applications

Handwriting Synthesis

more of national temperament

more of national temperament

more of national temperament

more of national temperament

more of national temperament

more of national temperament

LSTM Applications

- Speech recognition (Graves et al 2013)
- Neural Machine Translation (Sutskever et al 2014)
- Neural Turing Machine (Graves et al 2014)