

# Sequence Learning

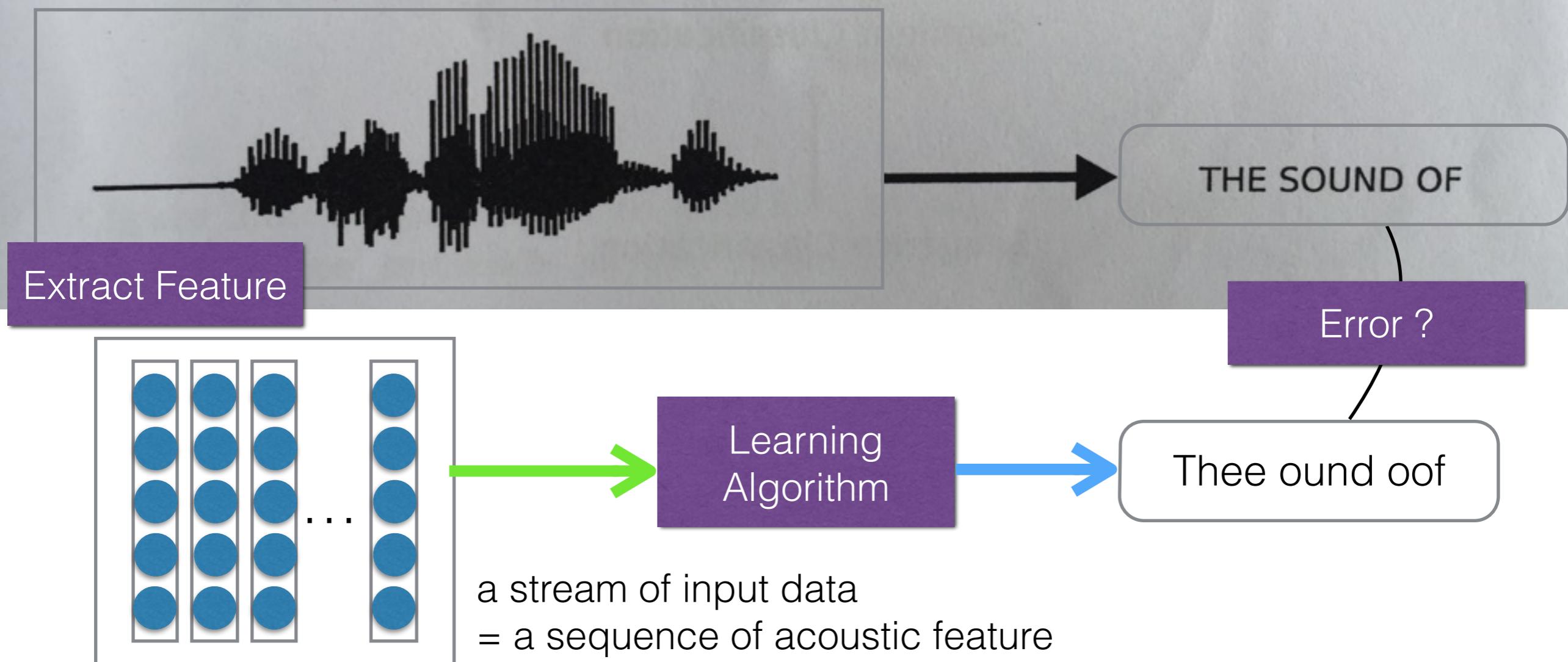
from character based acoustic model to  
End-to-End ASR system

替代役男 王俊豪

[chuchuhao831@gmail.com](mailto:chuchuhao831@gmail.com)

# Sequence Learning

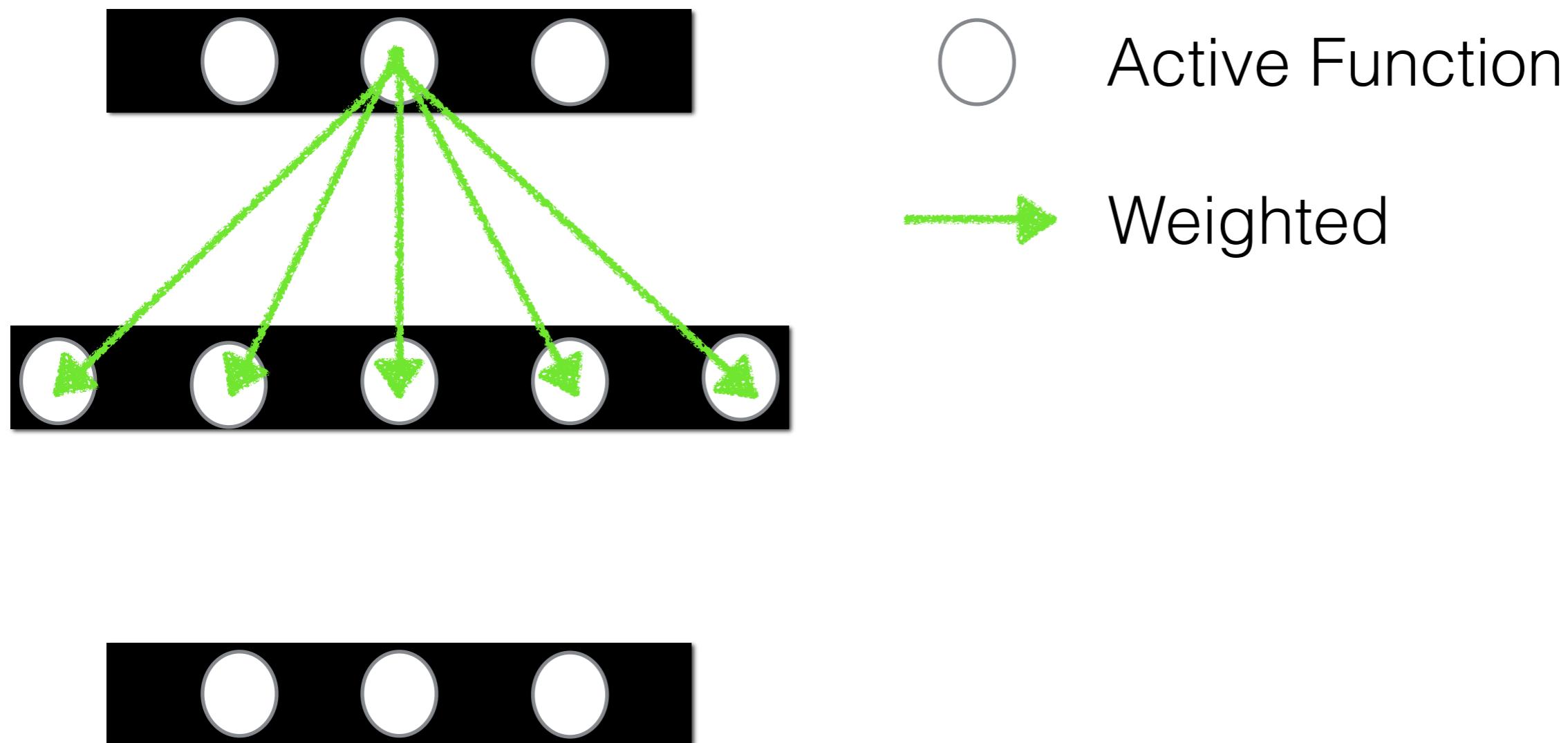
Foreign minister. → FOREIGN MINISTER.



# Outline

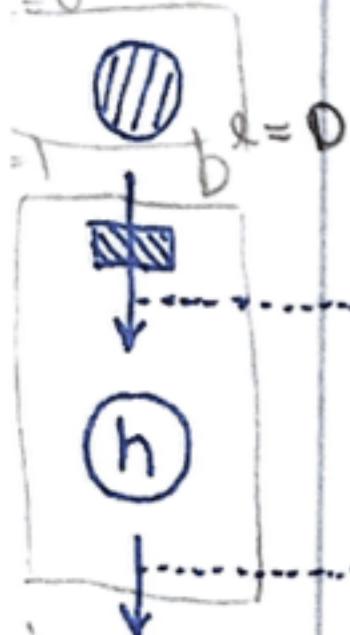
1. **F**eedforward **N**eural **N**etwork
2. **R**ecurrent **N**eural **N**etwork
3. **B**idirectional **R**ecurrent **N**eural **N**etwork
4. **L**ong **S**hort-Term **M**emory
5. Connectionist Temporal Classification
6. Build a End-to-End System
7. Experiment on Low-resource Language

# Feedforward Neural Network



# FNN

## [Forward Pass]



$I = \text{input units}$

$$a_h^l = \sum_{h'=1}^{H^{l-1}} w_{hh'}^l \cdot b_{h'}^{l-1}$$

$$b_h^l = \sigma_h(a_h)$$

# of hidden unit's

hidden unit

$$a_k^l = \sum_{h=1}^{H^{l-1}} w_{hk}^l \cdot b_h^{l-1}$$

Softmax

$$y_k = p(C_k | x) = \frac{e^{a_k}}{\sum_{k'=1}^K e^{a_{k'}}}$$

Multiclass

$l = \text{Layer}$

## [Backward Pass]

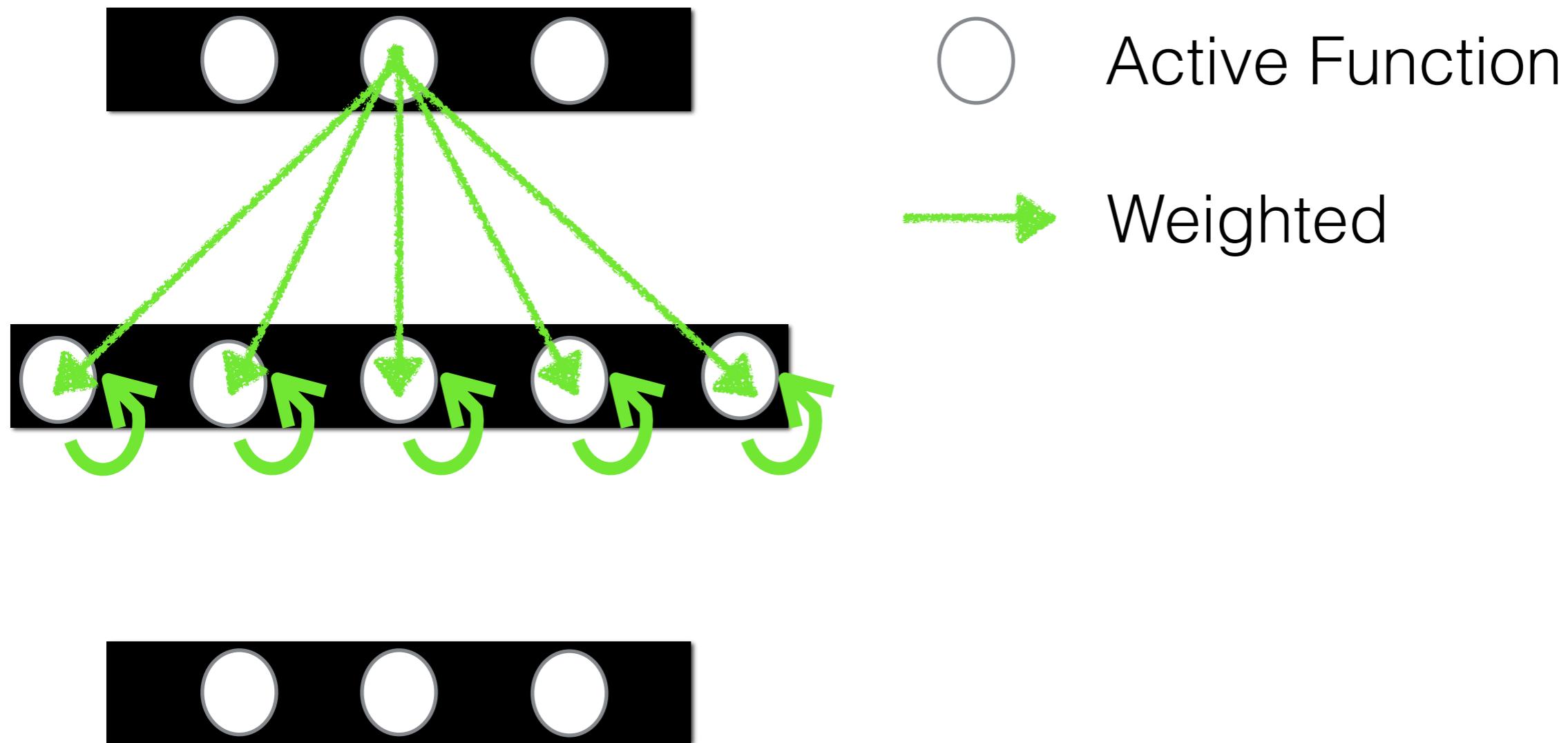
$$\frac{\partial O}{\partial b_h} = \sum_{k=1}^K \frac{\partial O}{\partial a_k} \frac{\partial a_k}{\partial b_h}$$

$$\delta_h = \frac{\partial O}{\partial a_h} = \left( \frac{\partial O}{\partial b_h} \right) \left( \frac{\partial b_h}{\partial a_h} \right) = \sigma'(a_h) \sum_{k=1}^K \delta_k w_{hk}$$

$$\frac{\partial O}{\partial w_{ij}} = \frac{\partial O}{\partial a_i} \frac{\partial a_i}{\partial w_{ij}} = \delta_j b_i = \delta_j b_i$$

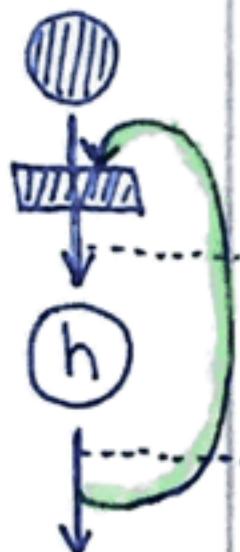
$$O = - \sum_{(x,z) \in S} \sum_{k=1}^K z_k \ln y_k$$

# Recurrent Neural Network



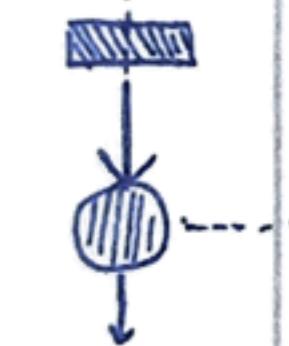
# RNN

## [Forward Pass]



$$a_h^{t,s} = \sum_{h=1}^{H-1} W_{hh} \cdot b_h^t + \sum_{h=1}^H W_{hh} \cdot b_h^{t-1}$$

$$b_h^t = \sigma_h(a_h^t)$$



$$a_k^t = \sum_{h=1}^{H-1} W_{hk} \cdot b_h^t$$

Softmax

## [Backward Pass]

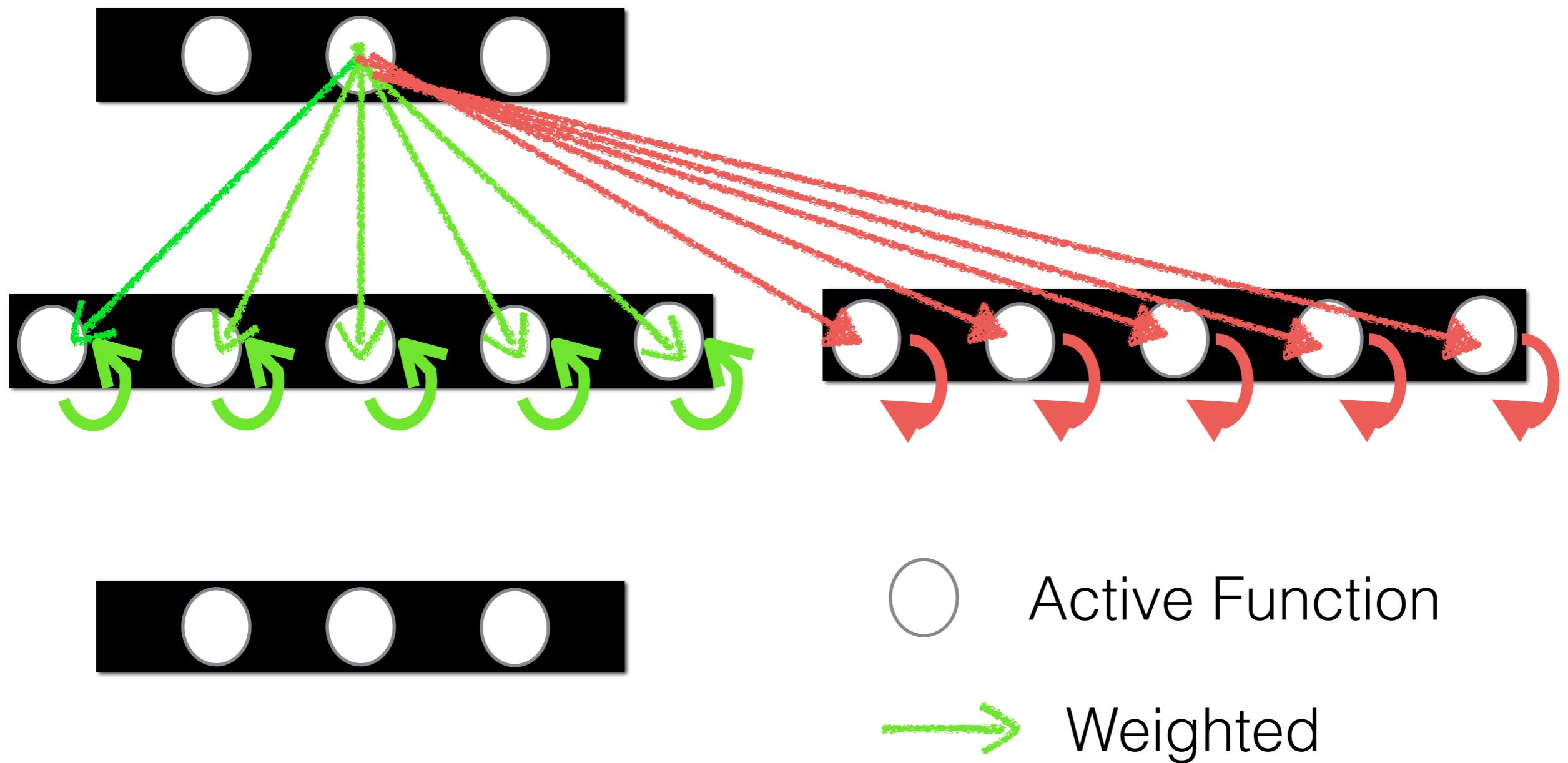
$$\frac{\partial O}{\partial b_h^t} = \sum_{k=1}^K \frac{\partial O}{\partial a_k^t} \frac{\partial a_k^t}{\partial b_h^t} + \sum_{h'=1}^H \frac{\partial O}{\partial a_h^{t+1}} \frac{\partial a_h^{t+1}}{\partial b_h^t} \frac{\partial a_h^t}{\partial b_h^t}$$

$$\delta_h^t = \sum_{k=1}^K \delta_k^t W_{hk} + \sum_{h'=1}^H \delta_{h'}^{t+1} W_{hh'}$$

$$\delta_h^t = \frac{\partial O}{\partial a_h^t} = \frac{\partial O}{\partial b_h^t} \frac{\partial b_h^t}{\partial a_h^t} = \sigma'(a_h^t)$$

$$\frac{\partial O}{\partial W_{ij}} = \sum_{t=1}^T \left( \frac{\partial O}{\partial a_i^t} \right) \frac{\partial a_i^t}{\partial a_j^t} \frac{\partial a_j^t}{\partial W_{ij}} = \sum_{t=1}^T \delta_i^t b_i^t$$

# Bidirectional RNN



```
for  $t = 1$  to  $T$  do
    Do forward pass for the forward hidden layer, storing activations at each timestep
for  $t = T$  to  $1$  do
    Do forward pass for the backward hidden layer, storing activations at each timestep
for  $t = 1$  to  $T$  do
    Do forward pass for the output layer, using the stored activations from both hidden layers
```

### Algorithm 3.1: BRNN Forward Pass

Similarly, the backward pass proceeds as for a standard RNN trained with BPTT, except that all the output layer  $\delta$  terms are calculated first, then fed back to the two hidden layers in opposite directions:

```
for  $t = T$  to  $1$  do
    Do BPTT backward pass for the output layer only, storing  $\delta$  terms at each timestep
for  $t = T$  to  $1$  do
    Do BPTT backward pass for the forward hidden layer, using the stored  $\delta$  terms from the output layer
for  $t = 1$  to  $T$  do
    Do BPTT backward pass for the backward hidden layer, using the stored  $\delta$  terms from the output layer
```

### Algorithm 3.2: BRNN Backward Pass

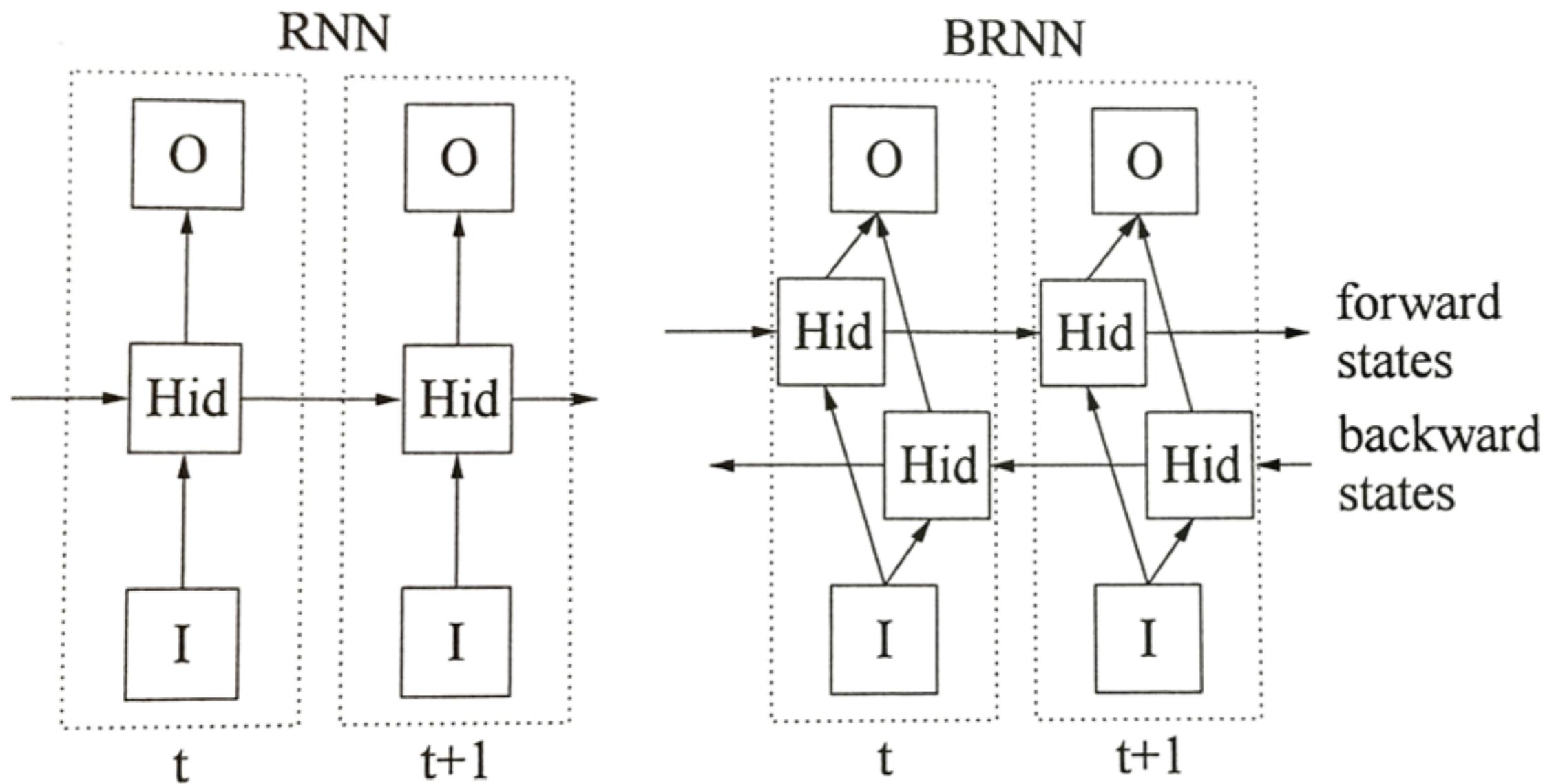


Figure 3.4: Standard and bidirectional RNNs

G

Gate

W

Weighted

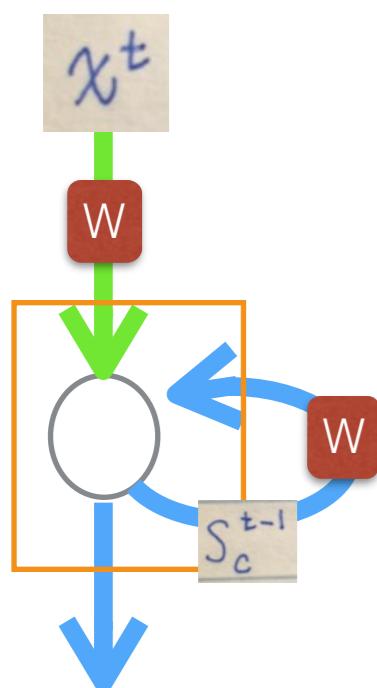
M

Multiplier

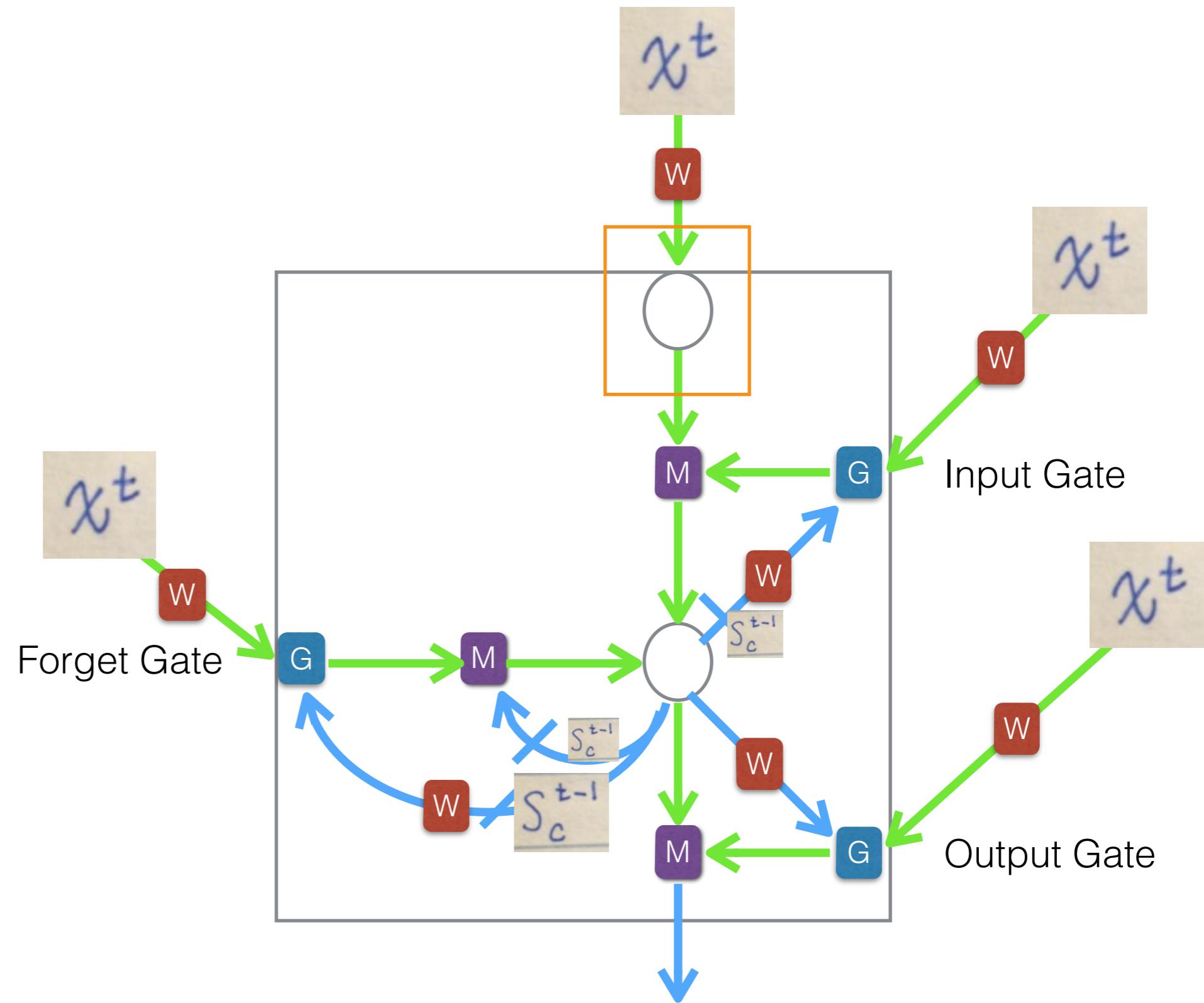


Active func

RNNs



# LSTM

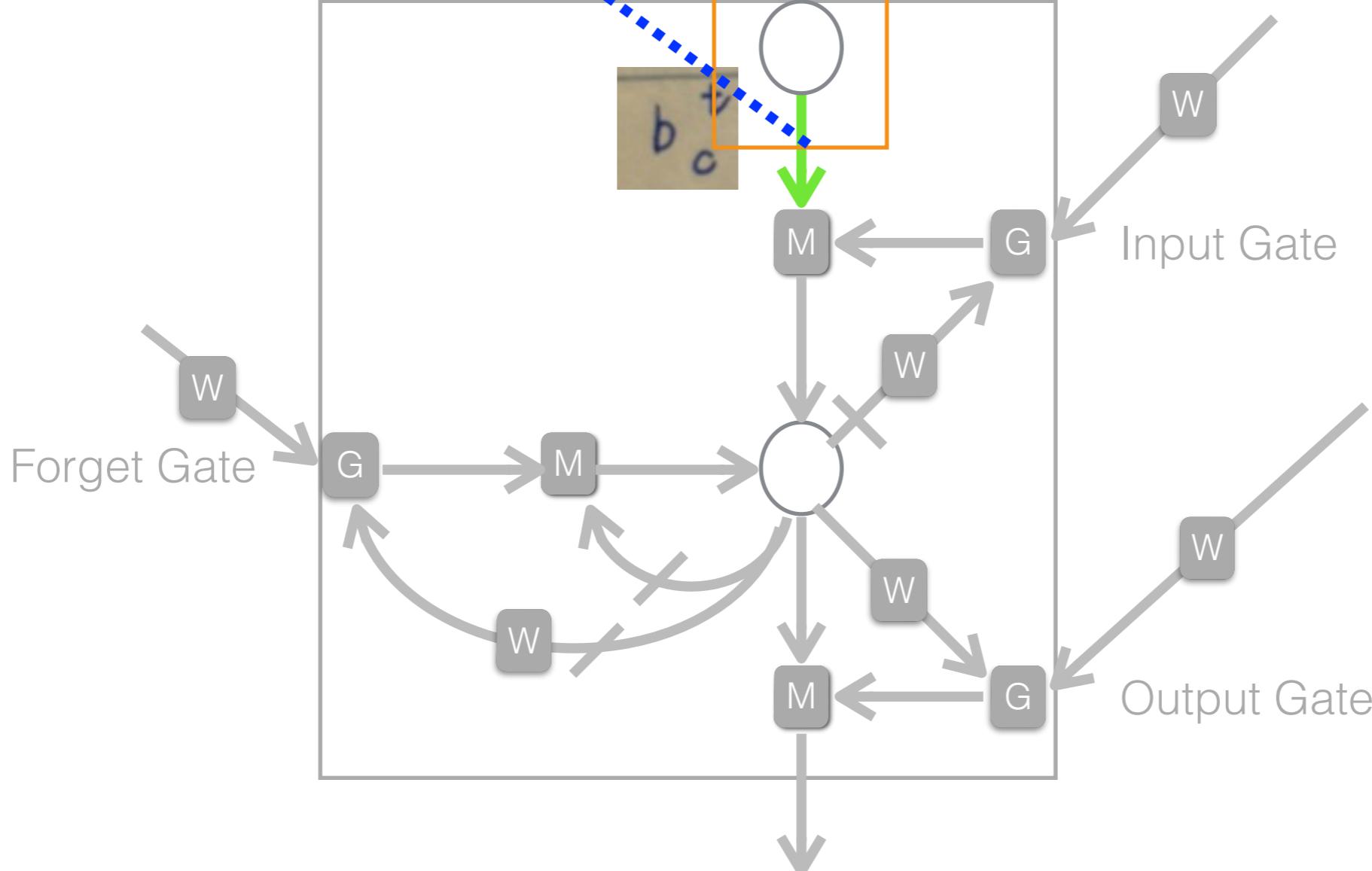


# LSTM Forward Pass, Previous Layer Input

(1)  $a_c^t = \sum_{i=1}^I w_{ic} x_i^t + \left( \sum_{h=1}^H w_{hc} b_h^{t-1} \right)$

Previous layer

$b_c^t = g(a_c^t)$



# LSTM Forward Pass, Input Gate

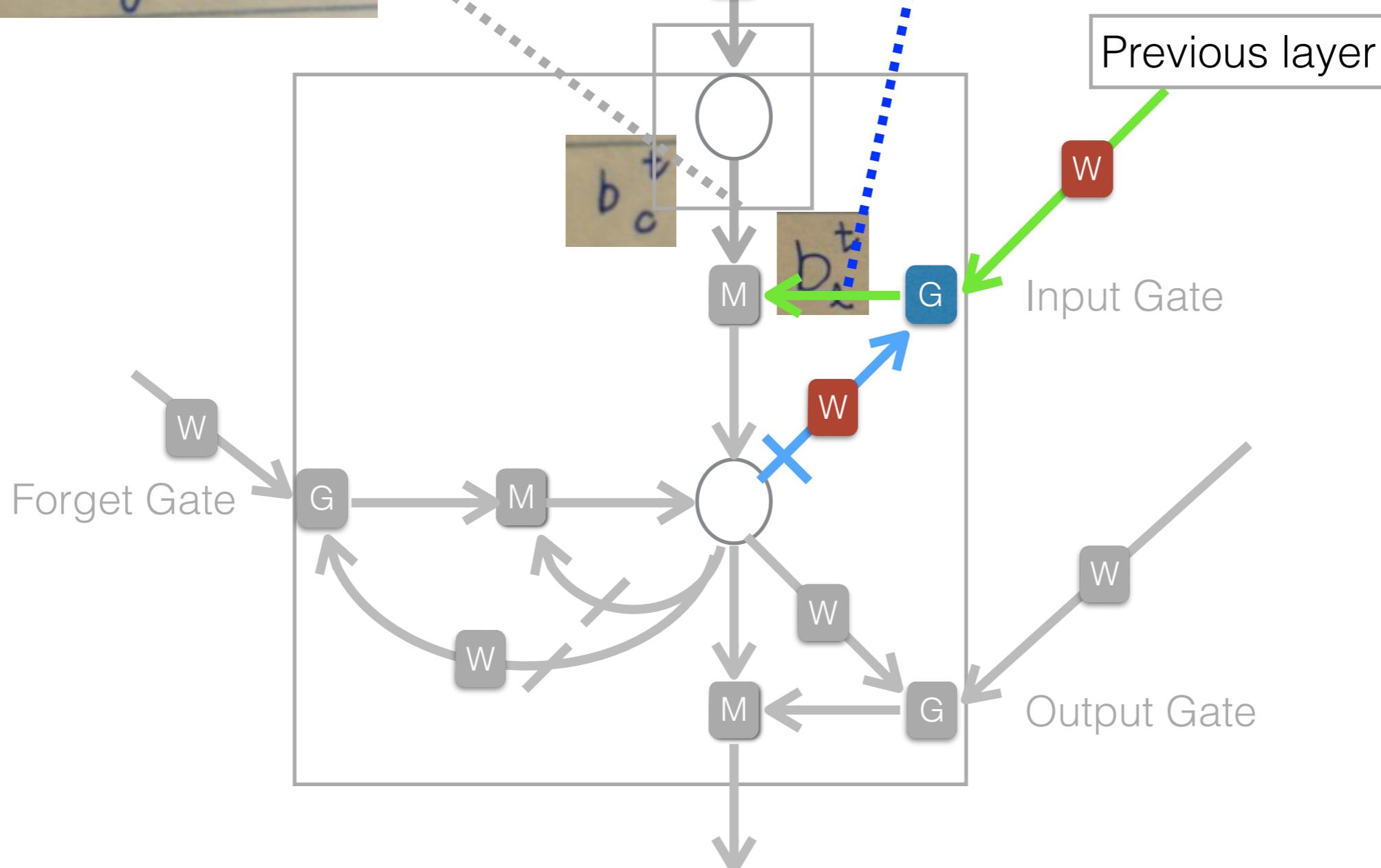
$$1) \quad a_c^t = \sum_{i=1}^I w_{ic} x_i^t + \left( \sum_{h=1}^H w_{hc} b_h^{t-1} \right)$$

$$b_c^t = g(a_c^t)$$

2)  $a_e^t = \sum_{i=1}^I w_{ie} x_i^t + \sum_{c=1}^C w_{ce} s_c^{t-1}$

optional  $\rightarrow \left( + \sum_{h=1}^H w_{he} b_h^{t-1} \right)$

$$b_e^t = f(a_e^t)$$



# LSTM Forward Pass, Foget Gate

$$(1) \quad a_c^t = \sum_{i=1}^I w_{ic} x_i^t + \left( \sum_{h=1}^H w_{hc} b_h^{t-1} \right)$$

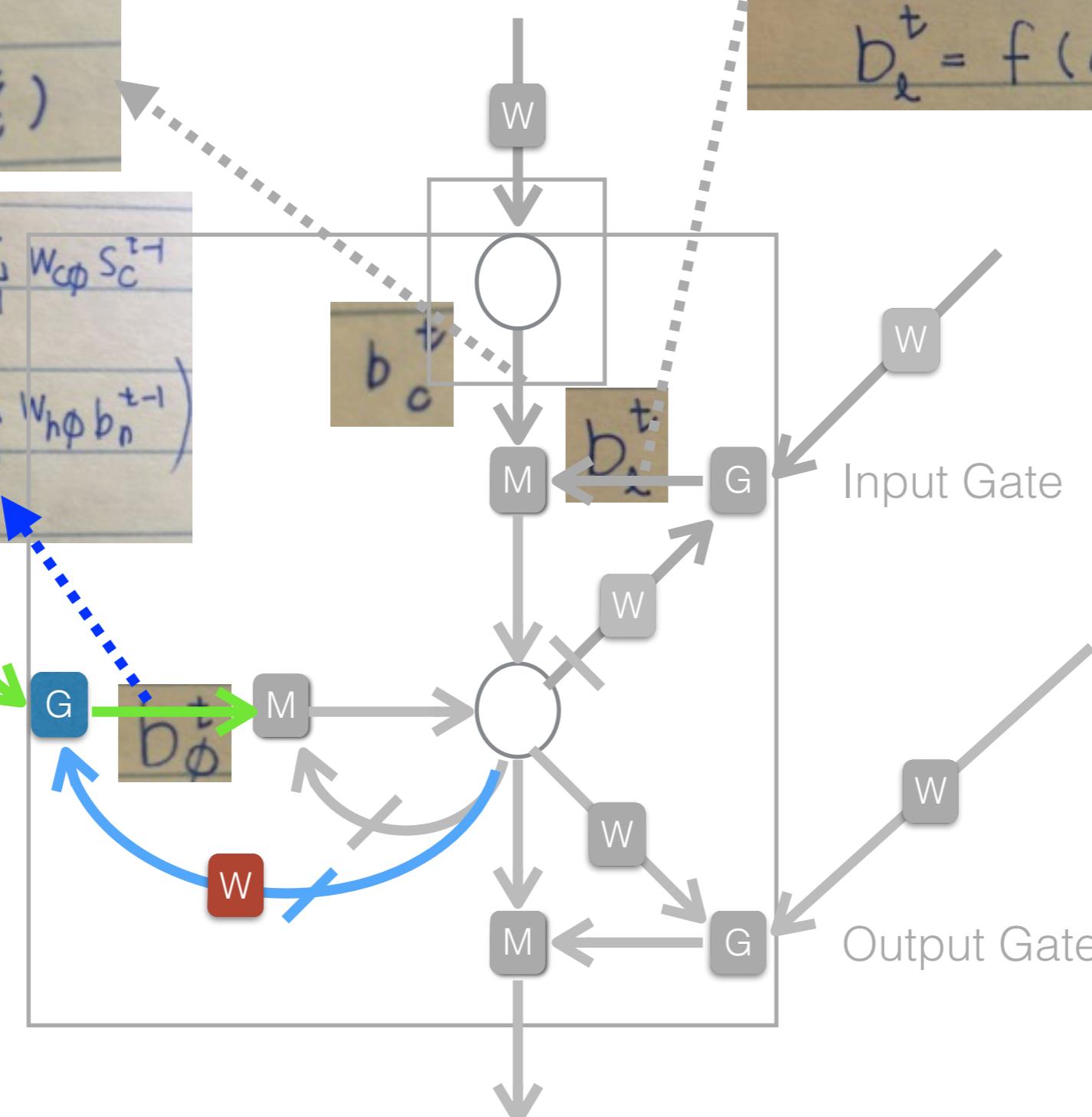
$$b_c^t = g(a_c^t)$$

$$(3) \quad a_\phi^t = \sum_{i=1}^I w_{i\phi} x_i^t + \sum_{c=1}^C w_{c\phi} s_c^{t-1} \\ \left( + \sum_{h=1}^H w_{h\phi} b_h^{t-1} \right)$$

$$b_\phi^t = f(a_\phi^t)$$

Previous layer

Forget Gate



$$2) \quad a_\ell^t = \sum_{i=1}^I w_{i\ell} x_i^t + \sum_{c=1}^C w_{c\ell} s_c^{t-1}$$

optional  $\rightarrow$

$$\left( + \sum_{h=1}^H w_{h\ell} b_h^{t-1} \right)$$

$$b_\ell^t = f(a_\ell^t)$$

# LSTM Forward Pass, Cell

$$(1) \quad a_c^t = \sum_{i=1}^I w_{ic} x_i^t + \left( \sum_{h=1}^H w_{hc} b_h^{t-1} \right)$$

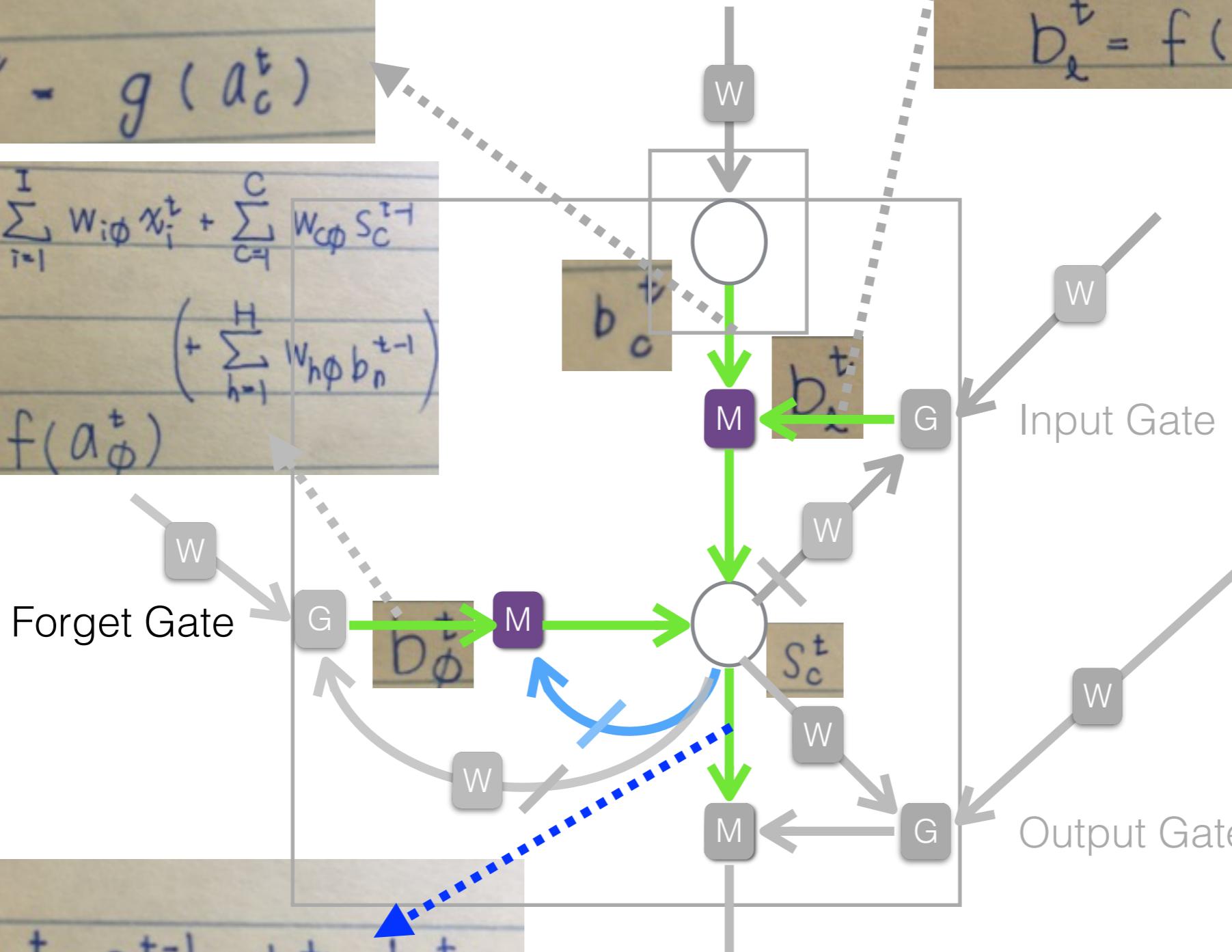
$$b_c^t = g(a_c^t)$$

$$(3) \quad a_\phi^t = \sum_{i=1}^I w_{i\phi} x_i^t + \sum_{c=1}^C w_{c\phi} s_c^{t-1} \\ \left( + \sum_{h=1}^H w_{h\phi} b_h^{t-1} \right)$$

$$b_\phi^t = f(a_\phi^t)$$

2)  $a_e^t = \sum_{i=1}^I w_{ie} x_i^t + \sum_{c=1}^C w_{ce} s_c^{t-1}$   
 optional  $\rightarrow \left( + \sum_{h=1}^H w_{he} b_h^{t-1} \right)$

$$b_e^t = f(a_e^t)$$



$$4) \quad s_c^t = b_\phi^t \cdot s_c^{t-1} + b_e^t \cdot b_c^t$$

(3) forget Gate      (2) Input Gate

# LSTM Forward Pass, Output Gate

$$(1) \quad a_c^t = \sum_{i=1}^I w_{ic} x_i^t + \left( \sum_{h=1}^H w_{hc} b_h^{t-1} \right)$$

$$b_c^t = g(a_c^t)$$

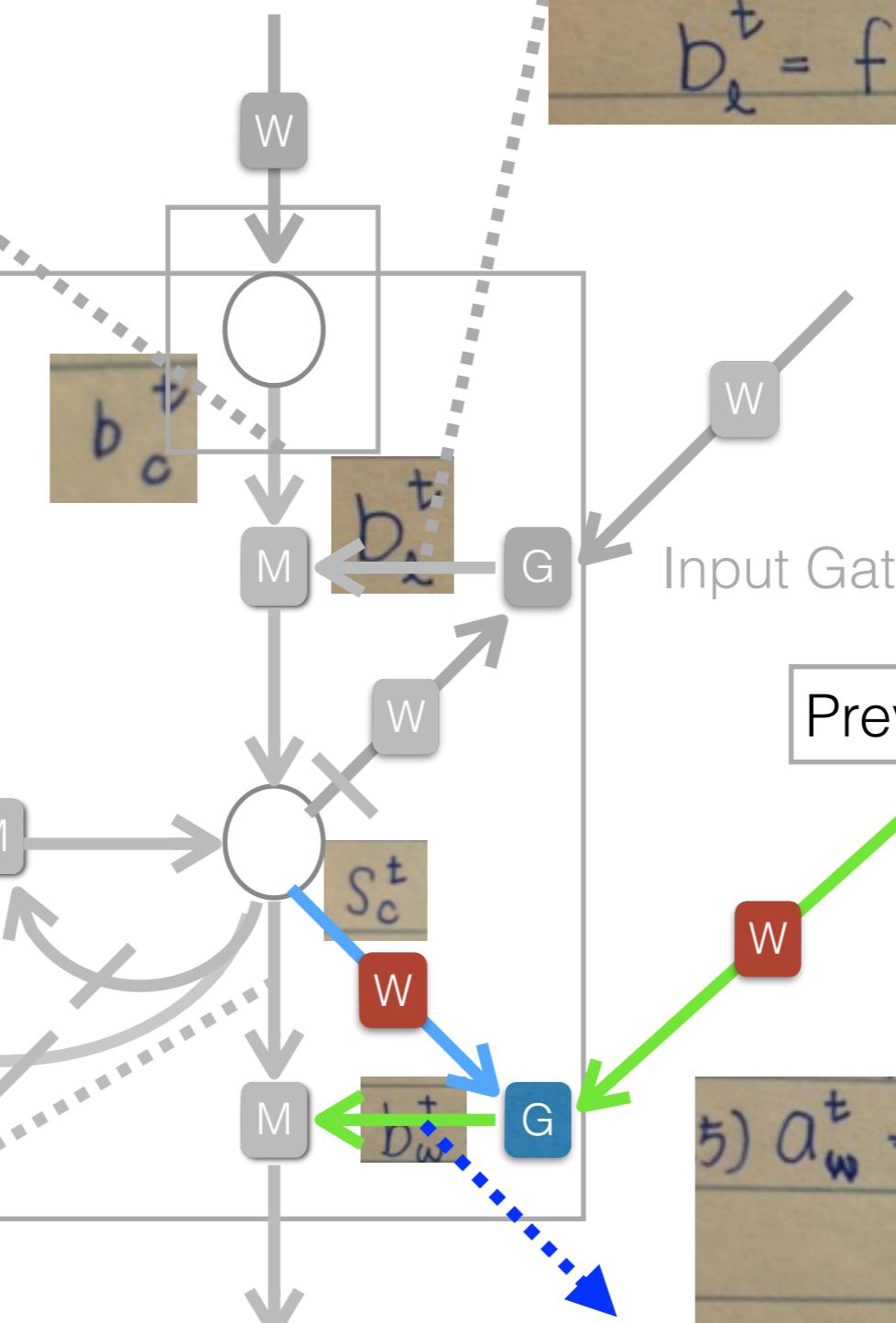
$$(3) \quad a_\phi^t = \sum_{i=1}^I w_{i\phi} x_i^t + \sum_{c=1}^C w_{c\phi} s_c^{t-1} \\ \left( + \sum_{h=1}^H w_{h\phi} b_h^{t-1} \right)$$

$$b_\phi^t = f(a_\phi^t)$$

Forget Gate

$$S_c^t = \underline{b_\phi^t \cdot S_c^{t-1}} + \underline{b_e^t \cdot b_c^t}$$

(3) forget Gate      (2) Input Gate



$$2) \quad a_e^t = \sum_{i=1}^I w_{ie} x_i^t + \sum_{c=1}^C w_{ce} s_c^{t-1}$$

optional  $\rightarrow$

$$\left( + \sum_{h=1}^H w_{he} b_h^{t-1} \right)$$

$$b_e^t = f(a_e^t)$$

$$5) \quad a_w^t = \sum_{i=1}^I w_{iw} x_i^t + \sum_{c=1}^C w_{cw} s_c^t \\ \left( + \sum_{h=1}^H w_{hw} b_h^{t-1} \right)$$

$$b_w^t = f(a_w^t)$$

# LSTM Forward Pass, Output

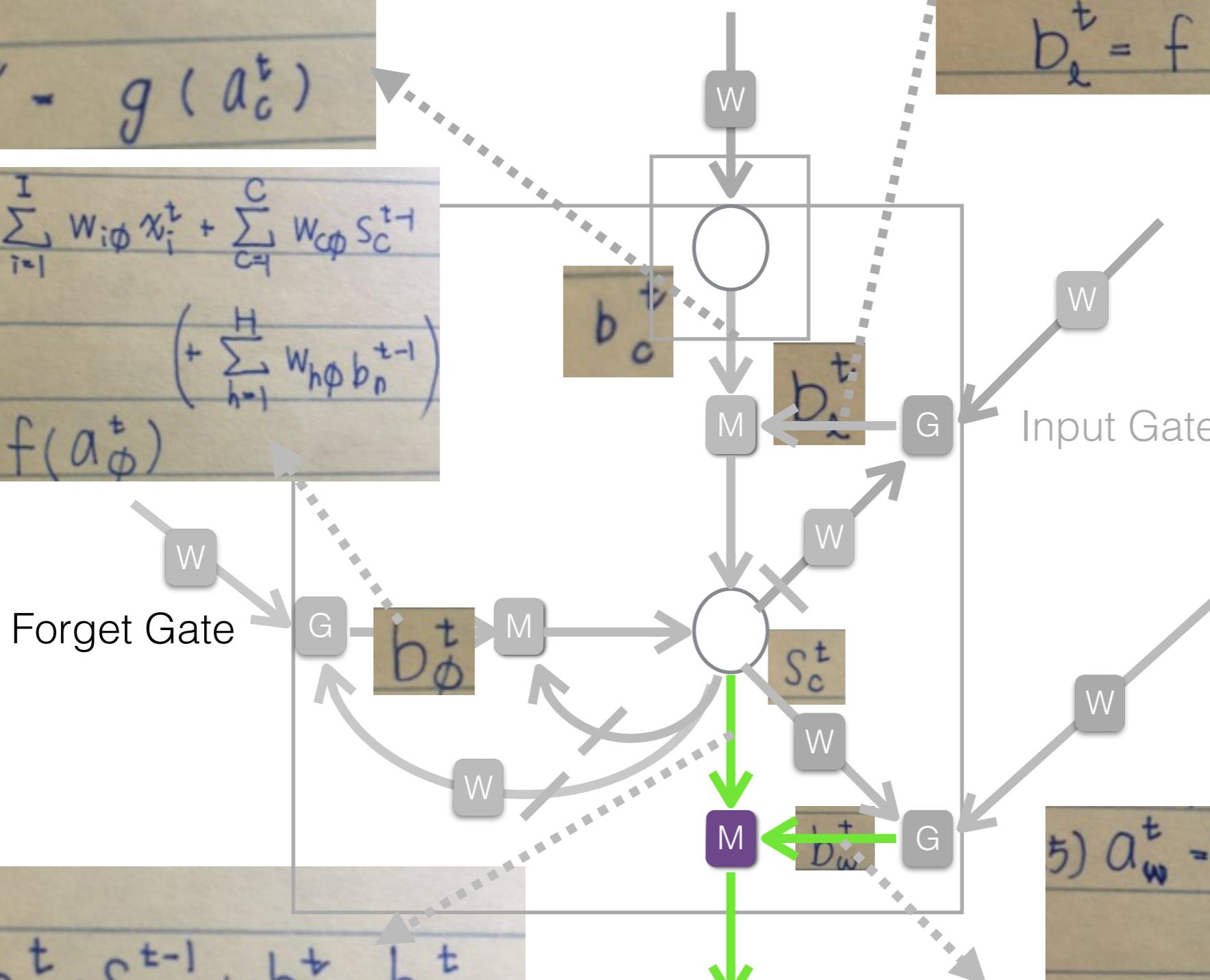
$$(1) \quad a_c^t = \sum_{i=1}^I w_{ic} x_i^t + \left( \sum_{h=1}^H w_{hc} b_h^{t-1} \right)$$

$$b_c^t = g(a_c^t)$$

$$(3) \quad a_{\emptyset}^t = \sum_{i=1}^I w_{i\emptyset} x_i^t + \sum_{C=1}^C w_{C\emptyset} s_C^{t-1}$$

$$\left( + \sum_{h=1}^H w_{h\emptyset} b_h^{t-1} \right)$$

$$b_{\emptyset}^t = f(a_{\emptyset}^t)$$



$$S_c^t = \underbrace{b_\phi^t \cdot S_c^{t-1}}_{(3) \text{ Forget Gate}} + \underbrace{b_i^t \cdot b_c^t}_{(2) \text{ Input Gate}}$$

### (3) forget Gate

$$(6). b_c^t = b_w^t \cdot \underline{h(s_c^t)}$$

$$b_{(t)}^t = f(a_{(t)}^t)$$

$$2) \quad a_e^t = \sum_{i=1}^I w_{iel} x_i^t + \sum_{c=1}^C w_{ct} s_c^{t-1}$$

Optional  $\rightarrow$

$$\left( + \sum_{h=1}^H w_{he} b_h^{t-1} \right)$$

# LSTM Backward Pass, Output Gate

Cells

$$\delta_c^t = b_c^t g'(a_c^t) \epsilon_c^t$$

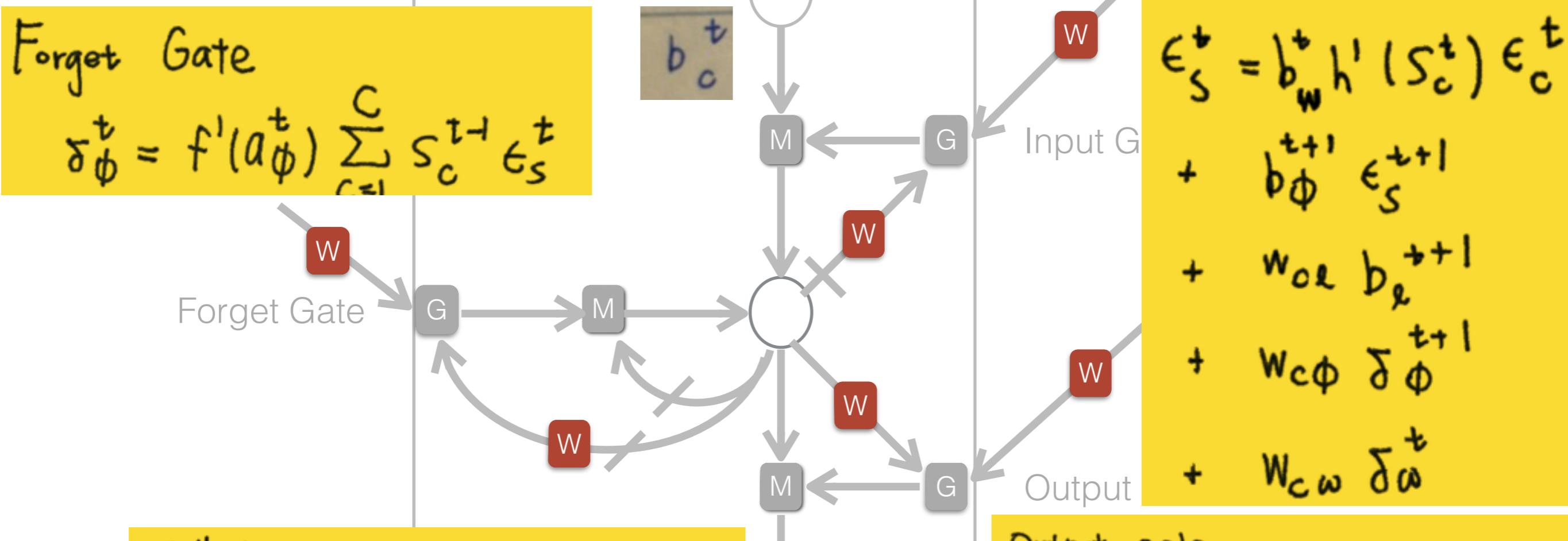
Input Gate

$$\delta_i^t = f'(a_i^t) \sum_{c=1}^C g(a_c^t) \epsilon_s^t$$

Forget Gate

$$\delta_\phi^t = f'(a_\phi^t) \sum_{c=1}^C s_c^{t-1} \epsilon_s^t$$

Forget Gate



Cell Output

$$\epsilon_c^t = \frac{\partial O}{\partial b_c^t} = \sum_{k=1}^K w_{ck} \delta_k^t + \sum_{h=1}^H w_{ch} \delta_h^{t+1}$$

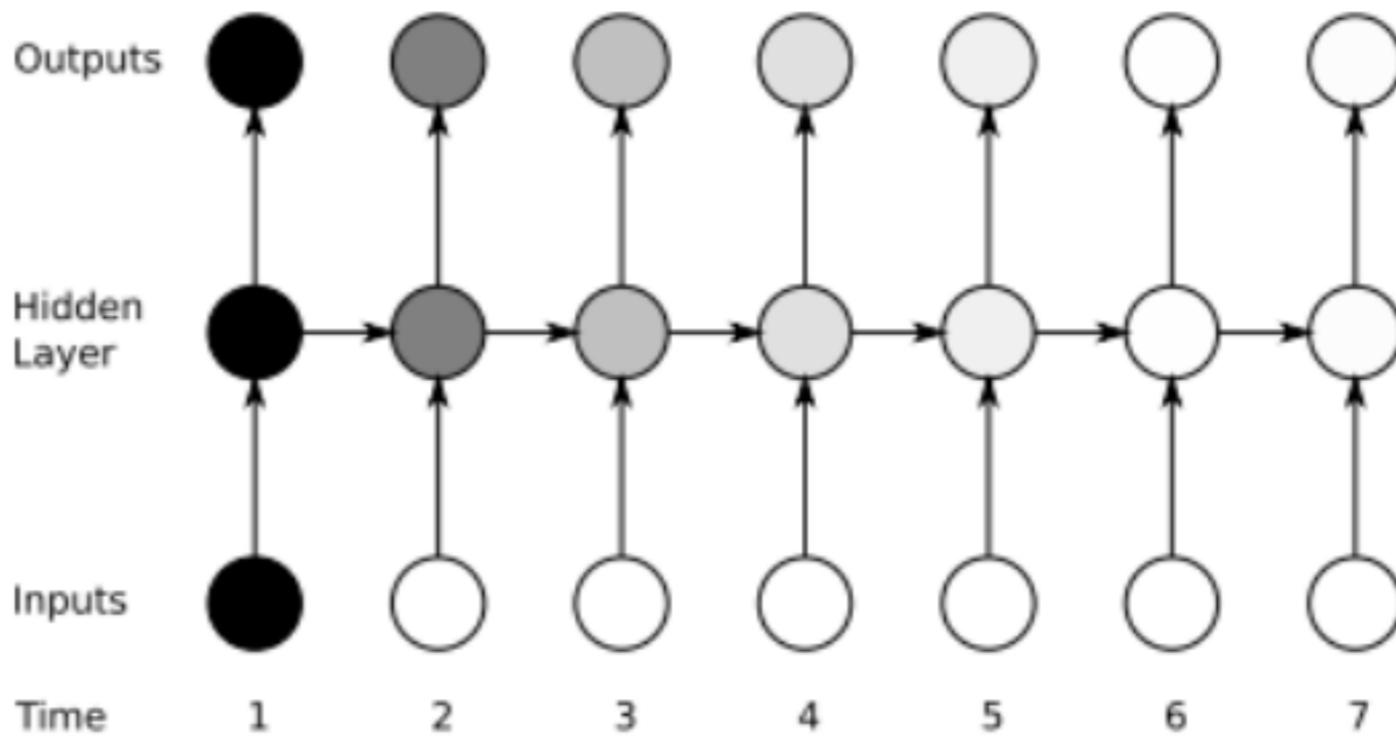
Output gate

$$\delta_o^t = f'(a_o^t) \sum_{c=1}^C h(s_c^t) \epsilon_c^t$$

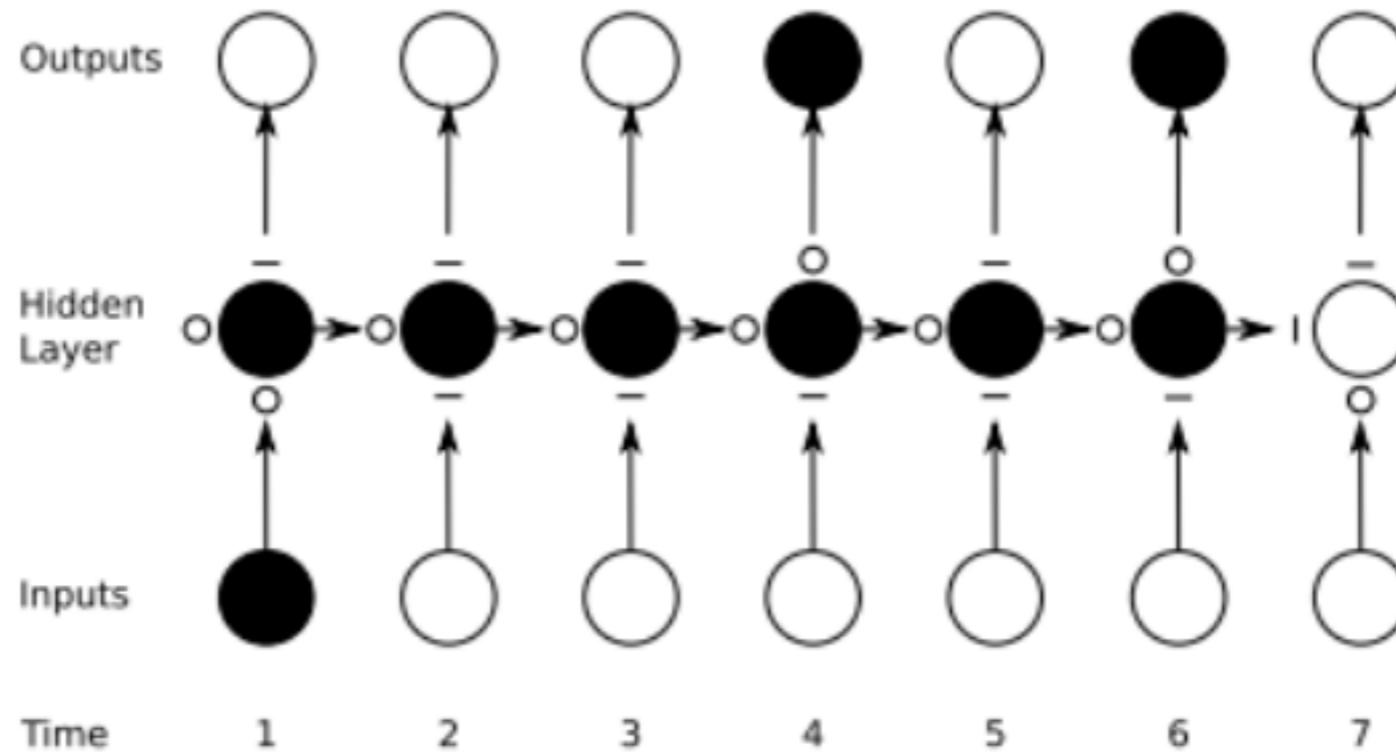
S states

$$\begin{aligned} \epsilon_s^t &= b_s^t h'(s_c^t) \epsilon_c^t \\ &+ b_\phi^{t+1} \epsilon_s^{t+1} \\ &+ w_{o\phi} b_\phi^{t+1} \\ &+ w_{c\phi} \delta_\phi^{t+1} \\ &+ w_{co} \delta_o^t \end{aligned}$$

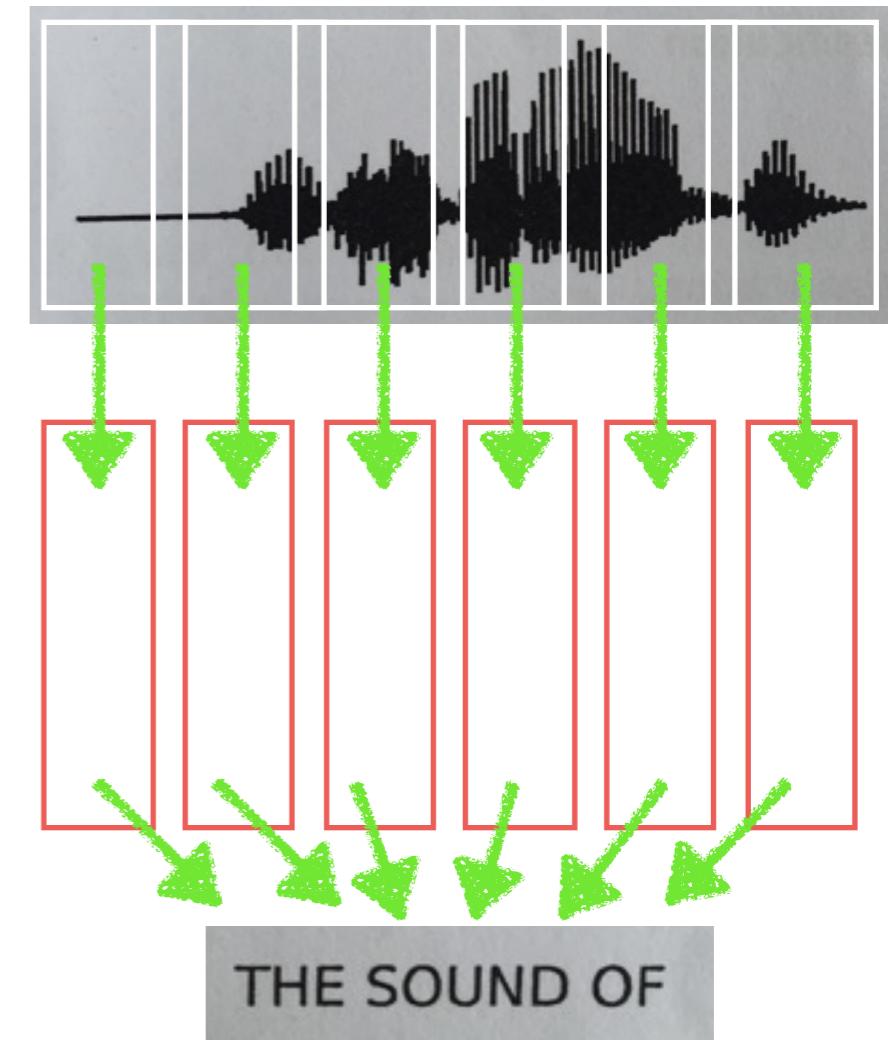
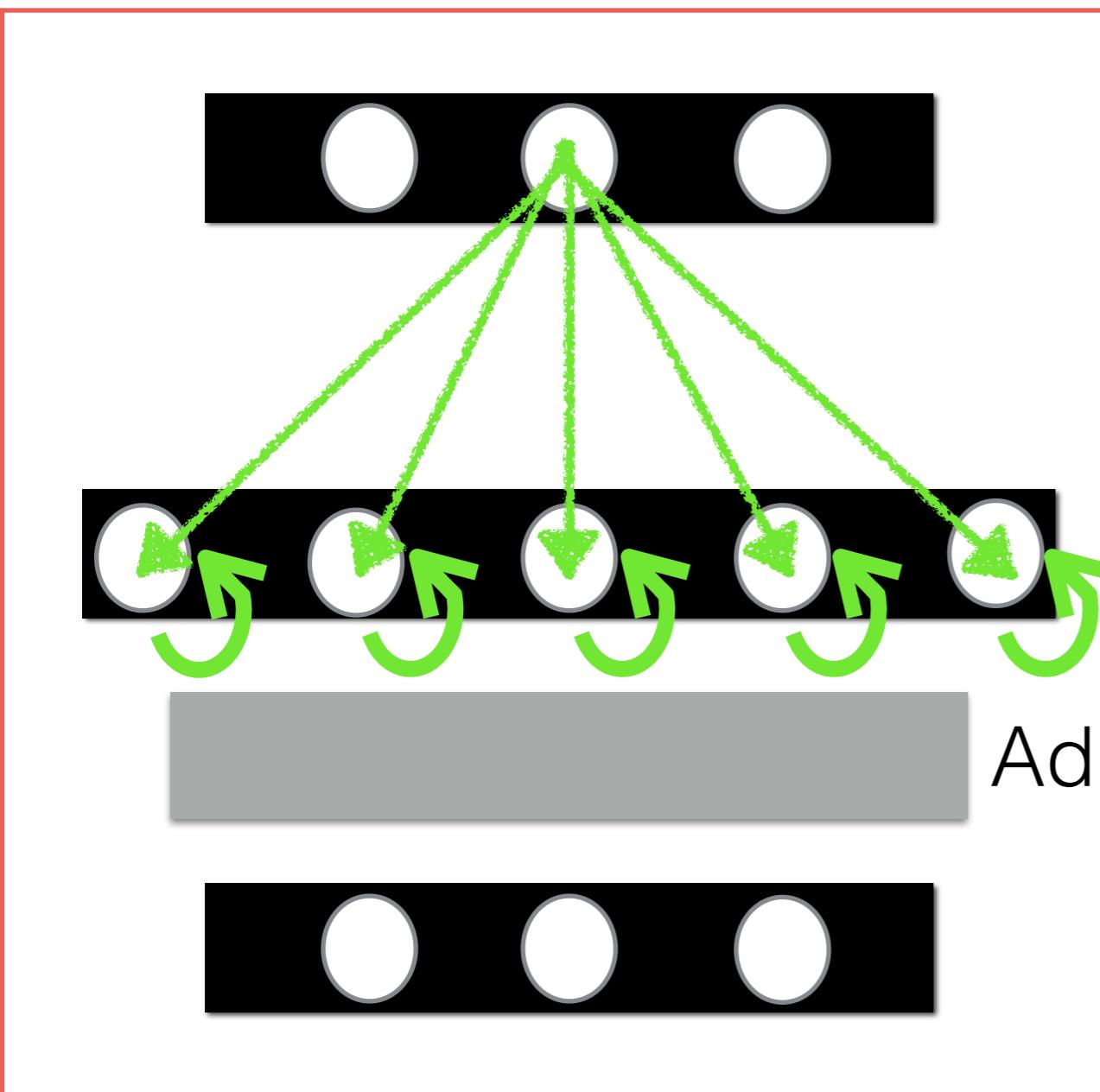
# RNN



# LSTM



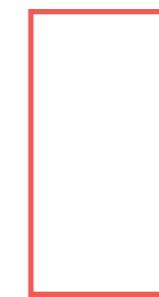
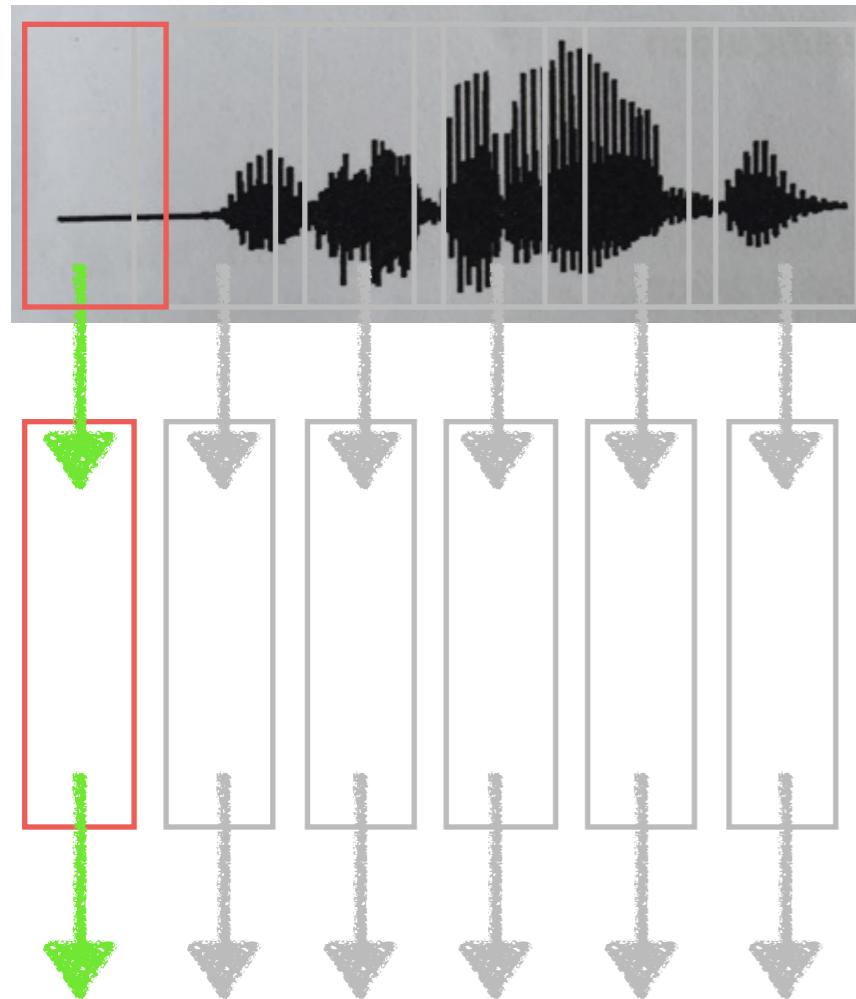
# Connectionist Temporal Classification



Additional CTC Output Layer

- (1) Predict label at any timestep
- (2) Probability of Sentence

Predict label  
at any timestep



Framewise Acoustic Feature



Label

A B C D E F G I J  
K L M N O P Q R  
S T U V W X Y Z  
{space}  
{blank}  
others like “ ” .

$$\mathbf{y} = \mathcal{N}_w(\mathbf{x}) \quad \mathcal{N}_w : (\mathbb{R}^m)^T \mapsto (\mathbb{R}^n)^{\bar{T}}$$

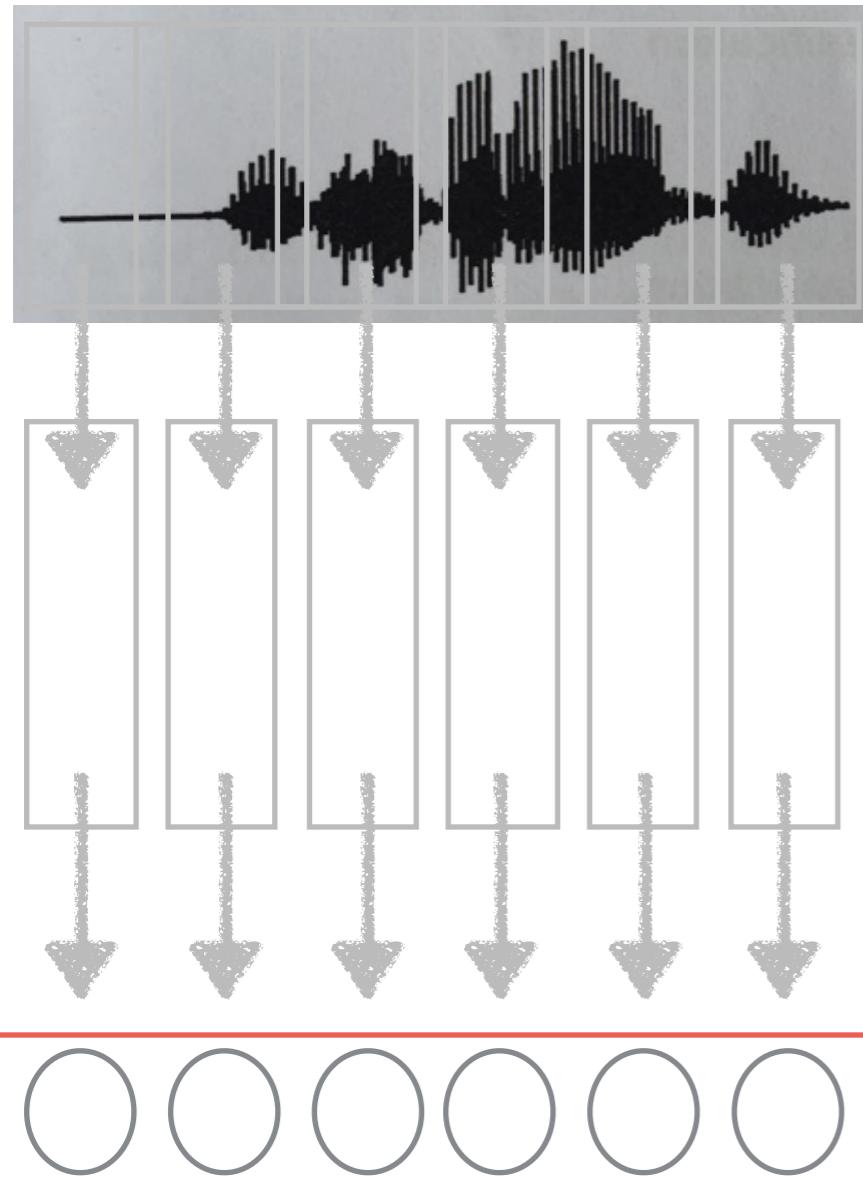


$$p(\pi|\mathbf{x}) = \prod_{t=1}^{\bar{T}} y_{\pi_t}^t, \quad \forall \pi \in L'^{\bar{T}}$$

THE SOUND OF

# Probability of Sentence

$$\mathcal{B} : L'^T \mapsto L^{\leq T}$$



paths:  $\pi \in L'^T$  **Thh..e S..outtd ooof**

Step 1, remove repeated labels  
Th.e S.outd of

Step 2, remove all blanks  
Predict  
**The Soudt of**

labelling:  $l \in L^{\leq T}$

$$p(l|x) = \sum_{\pi \in \mathcal{B}^{-1}(l)} p(\pi|x)$$

**Target** THE SOUND OF  $\leftarrow - -$

Now, we can do  $\dashrightarrow$   
the maximum likelihood training

## Output Decoding

Choose the labelling :  $h(\mathbf{x}) = \arg \max_{\mathbf{l} \in L^{\leq T}} p(\mathbf{l}|\mathbf{x})$

Best Path, the most probable path will correspond to most probable labelling

$$h(\mathbf{x}) \approx \mathcal{B}(\pi^*)$$

$$\text{where } \pi^* = \arg \max_{\pi \in N^t} p(\pi|\mathbf{x})$$

Something wrong here. Prefix Search would be better

$$\underset{\text{ML}}{O}(S, N_w) = - \sum_{(x_i, z_i) \in S} \ln(p(z_i | x_i))$$

↓  
training set  
 ↑  
weights  
 ↑  
target  
 ↑  
acoustic

objective function = 0

$$\frac{\partial O}{\partial y_k^t} = - \frac{\partial \ln(p(\underline{z} | \underline{x}))}{\partial y_k^t}$$

# How to connect target label and RNNs Output

network output unit  $k$  at time  $t$

Probability of observing label  $k$  at time  $t$

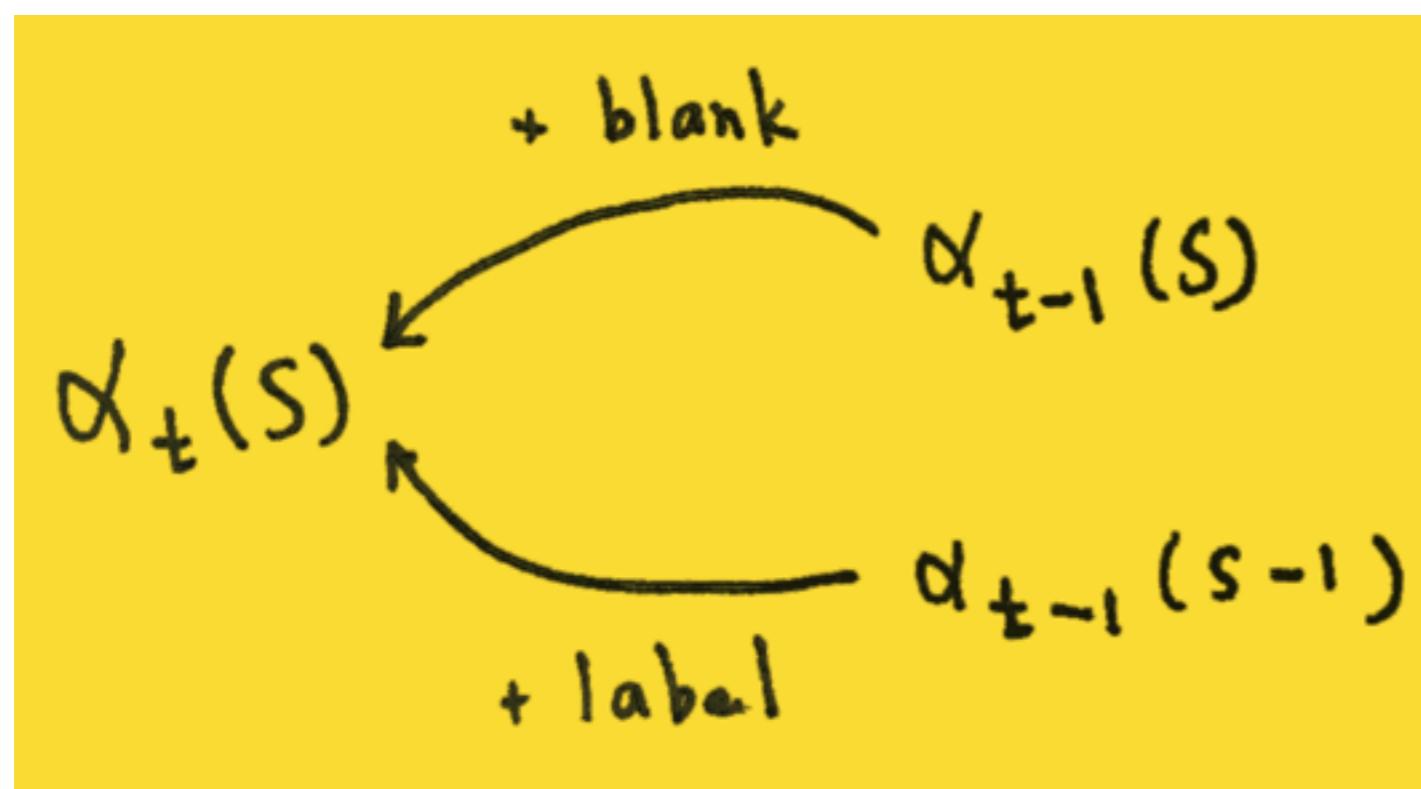
# Sum over paths probability, forward

first s symbol

$$\alpha_t(s) \equiv \sum_{\pi \in N^T : \pi_1 = s} \prod_{t'=1}^t y_{\pi_{t'}}^{t'}$$

timestep

$$B(\pi_{1:t}) = \ell_{1:S}$$



$$\alpha_1(1) = y_b^1$$

blank symbol

$$\alpha_1(2) = y_{\lambda_1}^1$$

first symbol of target

$$\alpha_1(s) = 0, \forall s > 2$$

only allow transition

- (1) between blank and non-blank label
- (2) distinct non-blank label

$$l_s = .$$

s-0 : .h.e.l.

s-1 : .h.e.l

s-2 : .h.e.

(s-2 blank to blank)

$$\bar{\alpha}_t(s) \equiv \alpha_{t-1}(s) + \alpha_{t-1}(s-1)$$

$$\alpha_t(s) = \begin{cases} \bar{\alpha}_t(s) \cdot y_{\lambda_s^t} & \text{if } \lambda_s^t = b \\ (\bar{\alpha}_t(s) + \alpha_{t-1}(s-2)) y_{\lambda_s^t} & \text{or } \lambda_{s-1}^t = \lambda_s^t \end{cases}$$

$$\lambda_s^t = b$$

$$\lambda_{s-1}^t = \lambda_s^t$$

$$l_s = |$$

s-0 : .h.e.l.l

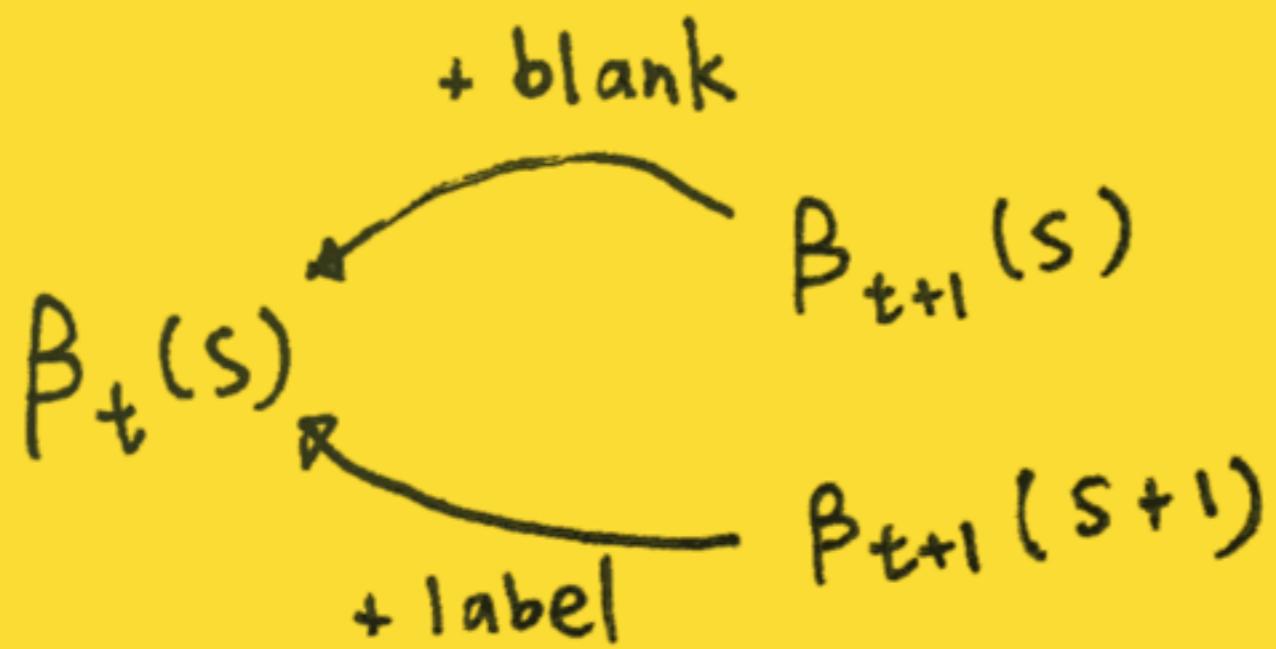
s-1 : .h.e.l.

s-2 : .h.e.l

(s-2 same non-blank label)

## Sum over paths probability, backward

$$\beta_t(s) \equiv \sum_{\pi \in N^T : \pi_t = t} \prod_{t'=t}^T y_{\pi_{t'}}^{s'}$$
$$B(\pi_{t,T}) = \lambda_{s=|\alpha|}$$



init

$$\beta_T(\lfloor \ell' \rfloor) = g_b^T$$

$$\beta_T(\lfloor \ell' \rfloor - 1) = g_{\ell_{\lfloor \ell' \rfloor}}^T$$

$$\beta_T(s) = 0, \quad \forall s < \lfloor \ell' \rfloor - 1$$

$$\bar{\beta}_t(s) \equiv \beta_{t+1}(s) + \beta_{t+1}(s+1)$$

$$\beta_t(s) = \begin{cases} \bar{\beta}_t(s) \cdot g_{\ell_s^t}^T \\ (\bar{\beta}_t(s) + \beta_{t+1}(s+2)) \cdot g_{\ell_s^t}^T \end{cases}$$

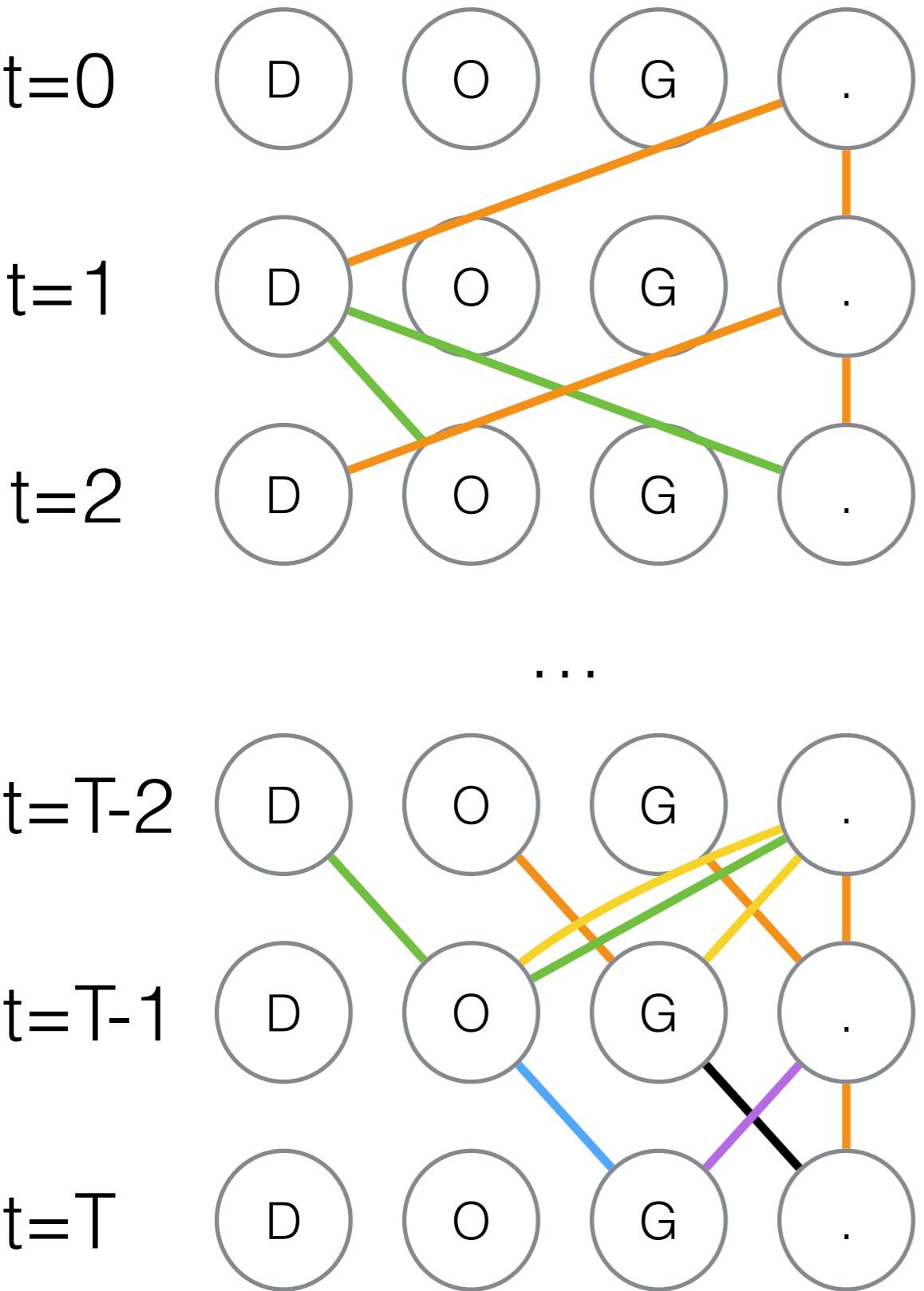
if  $\ell_s^t = b$   
or  $\ell_{s+2}^t = \ell_s^t$

- only allow transition  
 (1) between blank and non-blank label  
 (2) distinct non-blank label

$\underline{l}_s = .$   
 $s+0 : .l.l.o.$   
 $s+1 : l.l.o.$   
 $s+2 : .l.o.$   
 (  $s+2$  blank to blank )

$\underline{l}_s = l$   
 $s+0 : l.l.o.$   
 $s+1 : .l.o.$   
 $s+2 : l.o.$   
 (  $s+2$  same non-blank label )

# Sum over paths probability



Target : dog -> .d.o.g.

s=1 —— .

s=2 —— .d

s=3 —— .d.

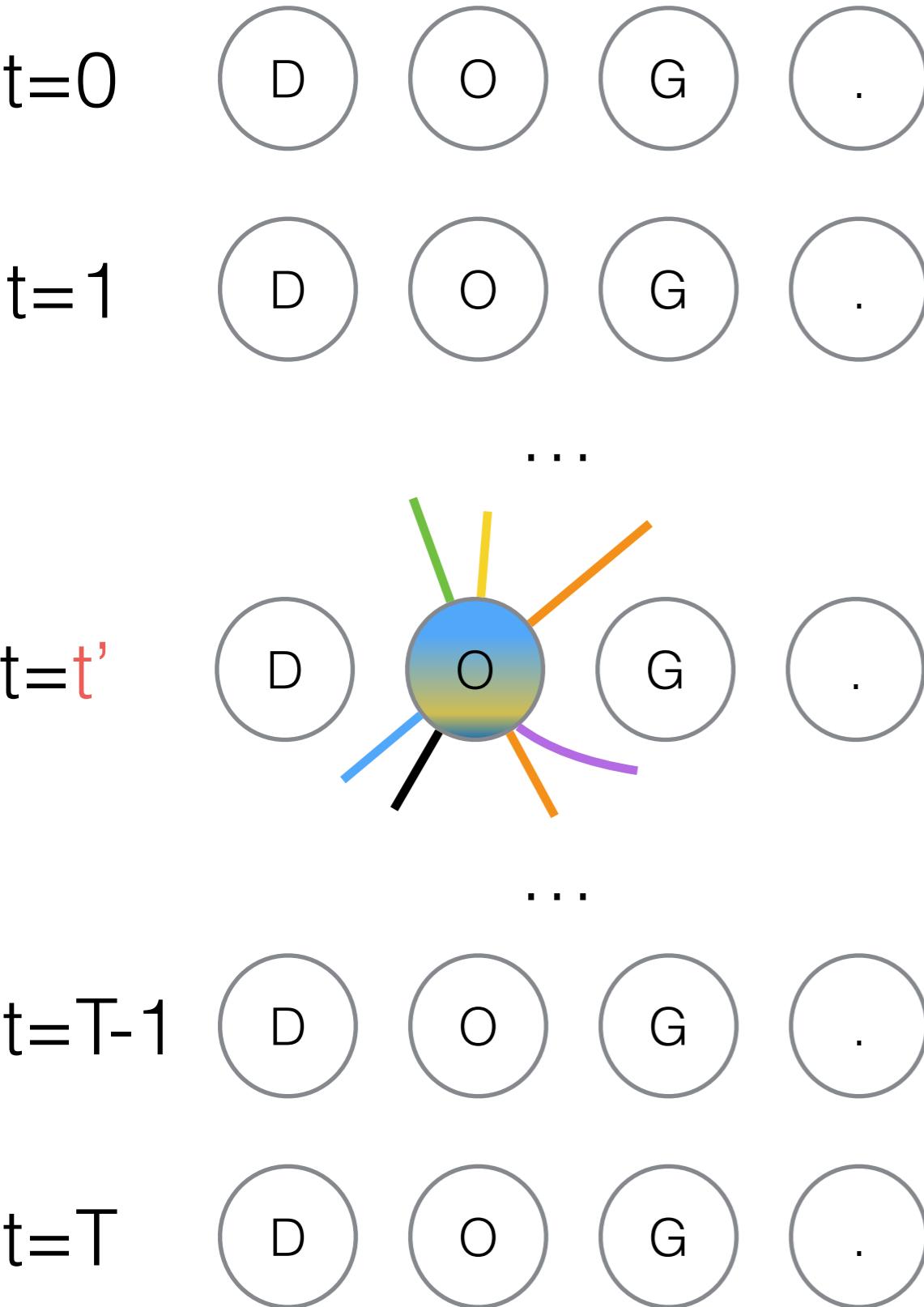
s=4 —— .d.o

s=5 —— .d.o.

s=6 —— .d.o.g

s=7 —— .d.o.g.

# Sum over paths probability



backward-forward : π  
 the probability of all the paths corresponding to  $\lambda$   
 that go through the symbol at time t

$$\alpha_t(s) \cdot \beta_t(s) = \sum_{\pi \in \mathcal{B}^t(\lambda)} y_{\lambda_s^t}^t \frac{\prod_{t=1}^T y_{\pi_t}^t}{\pi_t = \lambda_s^t}$$

$$\Rightarrow \frac{\alpha_t(s) \cdot \beta_t(s)}{y_{\lambda_s^t}^t} = \sum_{\substack{\pi \in \mathcal{B}^t(\lambda) : \\ \pi_t = \lambda_s^t}} p(\pi | x)$$

$$p(\lambda | x) = \sum_{s=1}^{|\lambda'|} \frac{\alpha_t(s) \cdot \beta_t(s)}{y_{\lambda_s^t}^t}$$

$$\frac{\partial O}{\partial y_k^t} = - \frac{\partial \ln(P(z|x))}{\partial y_k^t}$$

$$= \frac{-1}{P(z|x)} \cdot \frac{\partial P(z|x)}{\partial y_k^t}$$

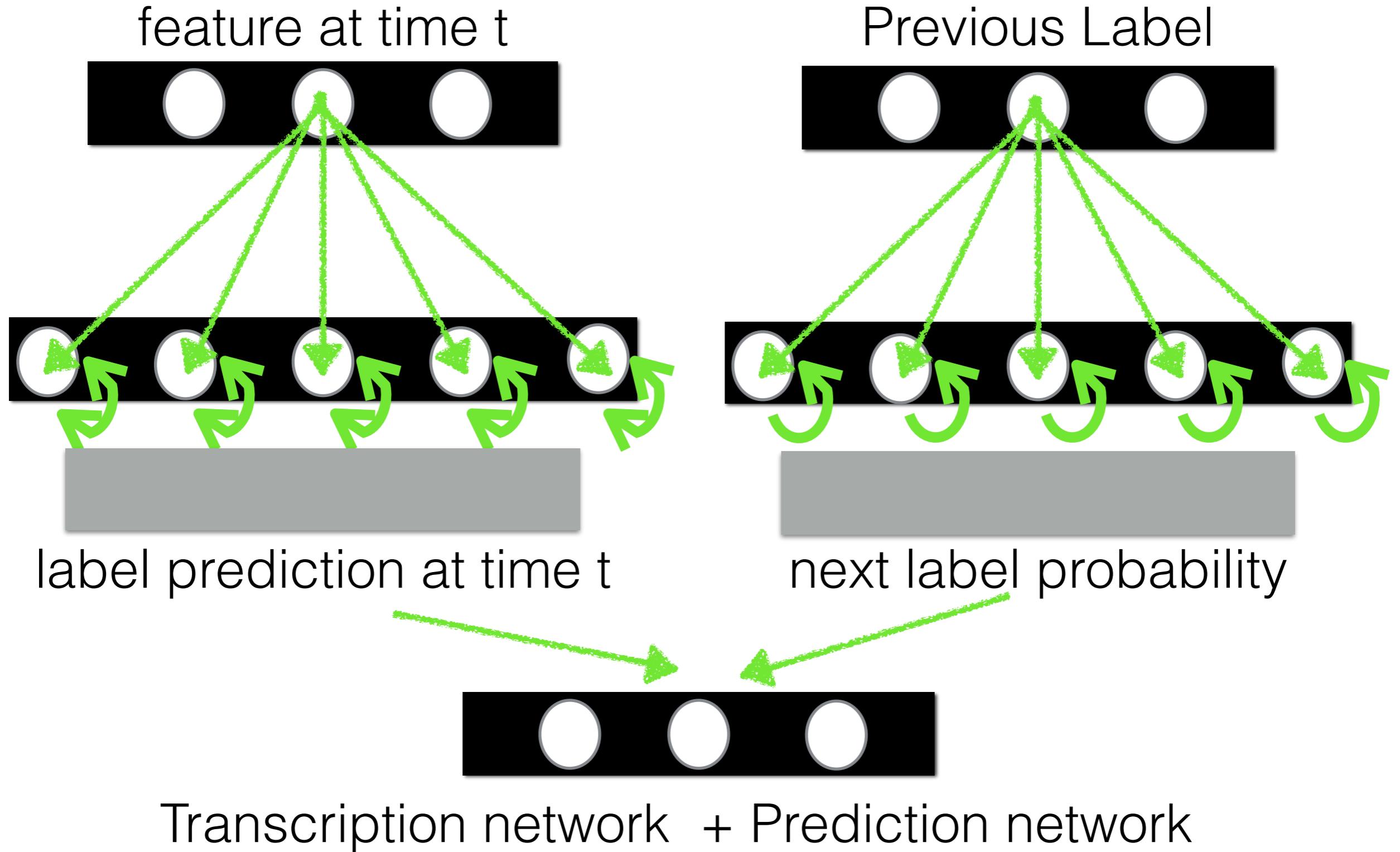
$$= \frac{-1}{\sum_{s=1}^{|z'|} \frac{\alpha_t(s) \cdot \beta_t(s)}{y_{z'_s}^t}} \cdot \frac{1}{y_k^t} \sum_{s \in lab(z,k)} \alpha_t(s) \cdot \beta_t(s)$$

$$Z_t \stackrel{\text{def}}{=} \sum_{s=1}^{|l'|} \frac{\hat{\alpha}_t(s) \hat{\beta}_t(s)}{y_{l'_s}^t}$$

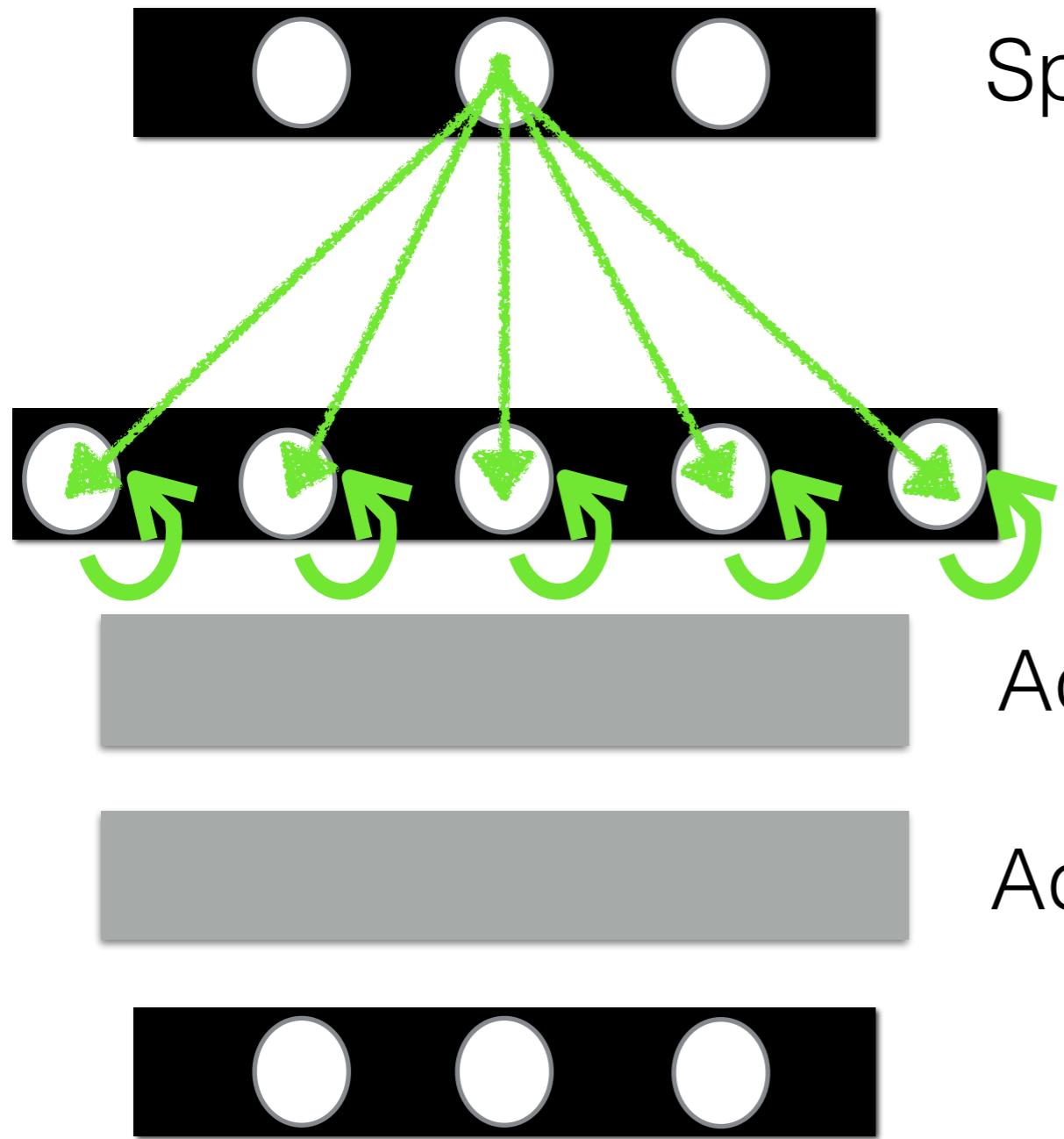
last we derivates w.r.t. unnormalised outputs, get error signal

$$\frac{\partial O^{ML}(\{(x, z)\}, \mathcal{N}_w)}{\partial u_k^t} = y_k^t - \frac{1}{y_k^t Z_t} \sum_{s \in lab(z, k)} \hat{\alpha}_t(s) \hat{\beta}_t(s)$$

# A. Graves. Sequence Transduction with Recurrent Neural Networks



# A. Graves. Towards End-to-end Speech Recognition with Recurrent Neural Networks ( **Google Deepmind** )



Spectrogram feature

$$CTC(\mathbf{x}) = -\log \Pr(\mathbf{y}^* | \mathbf{x})$$

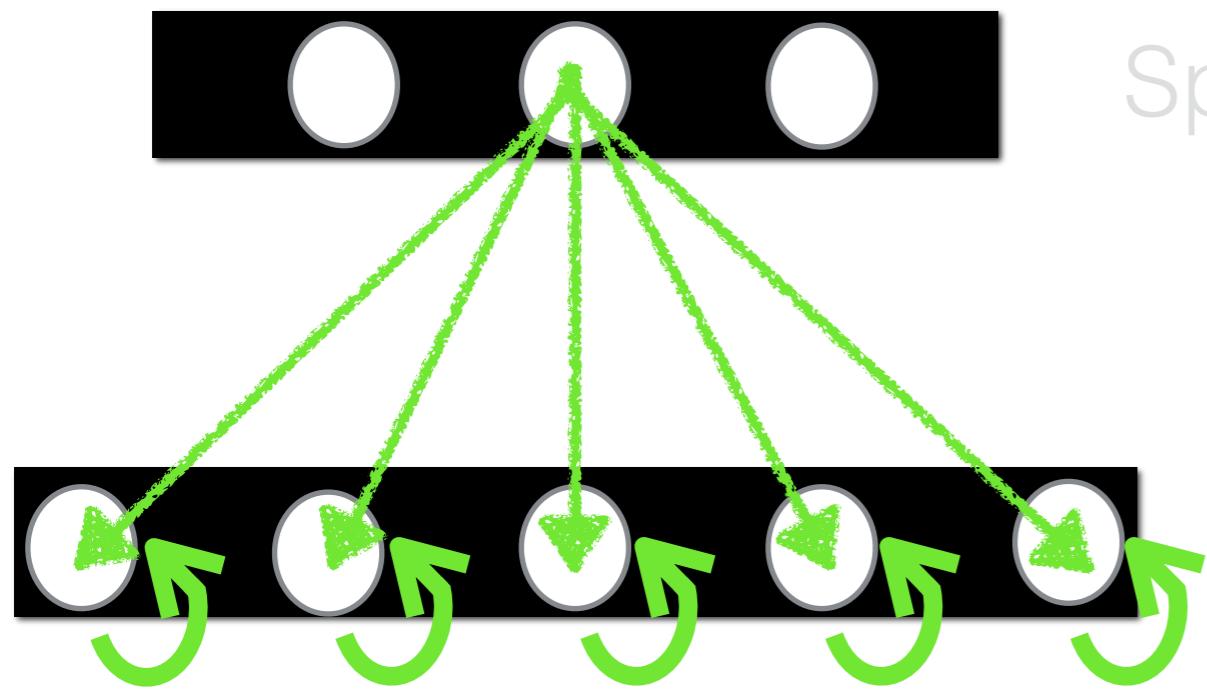
Additional CTC Output Layer

Additional WER Output Layer

$$\mathcal{L}(\mathbf{x}) = \sum_{\mathbf{y}} \Pr(\mathbf{y} | \mathbf{x}) \mathcal{L}(\mathbf{x}, \mathbf{y})$$

Minimise Objective Function -> Minimise Loss Function

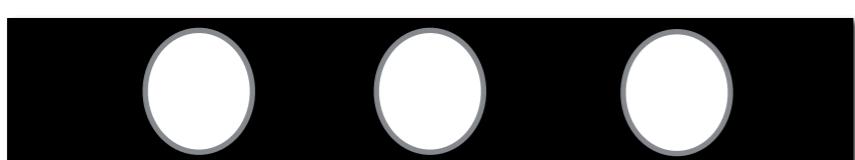
# Andrew Y. Ng. Deep Speech: Scaling up end-to-end speech recognition ( **Baidu Research**)



Spectrogram feature

Additional CTC Output Layer

Additional WER Output Layer

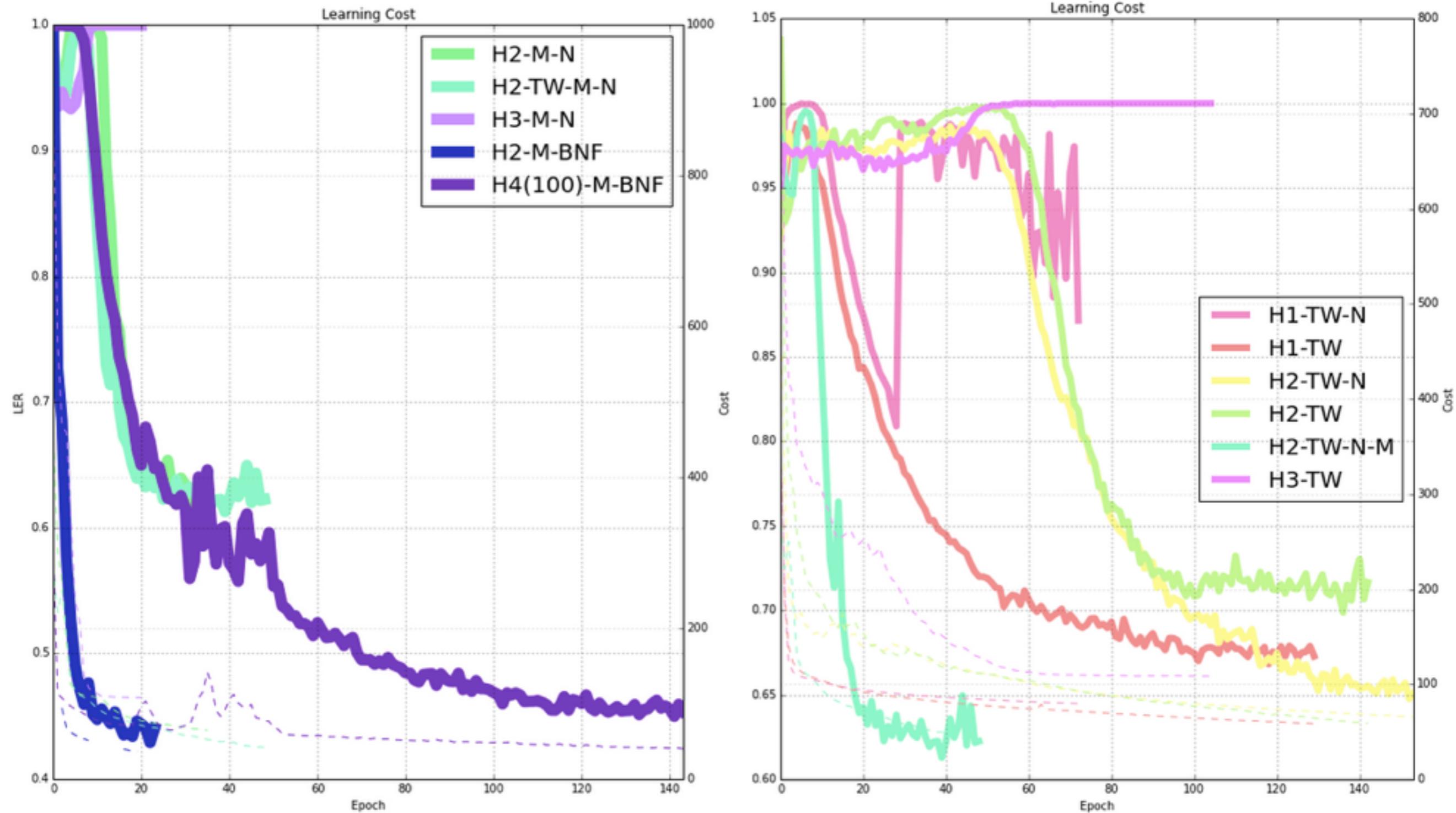


RNN Output : arther n tickets for the game  
Decode Target: are there any tickets for the gam

Language Model Transcription

$$Q(L) = \log(P(L|x)) + a * \log(P_{LM}(L)) + b * \text{word\_count}(L)$$

# Low-resources language experiment



Target: ? ka didto sigi ra og tindog  
Output: h bai to giy ngtndog

test 0.44,  
test 0.71

# Reference

1. Boulard and Morgan (1994) Connectionist speech recognition A hybrid approach.
2. A. Grave. Supervised Sequence Labelling with Recurrent Neural Networks
3. A. Grave. Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Network.
4. A. Graves. Sequence Transduction with Recurrent Neural Networks
5. A. Graves. Towards End-to-end Speech Recognition with Recurrent Neural Networks ( Google Deepmind )
6. Andrew Y. Ng. Deep Speech: Scaling up end-to-end speech recognition ( Baidu Research)
7. standford CTC: <https://github.com/amaas/stanford-ctc>
8. R. Pascanu , On the difficulty of training recurrent neural networks
9. L. Besacier. Automatic Speech Recognition for Under-Resourced Languages: A Survey
- 10.A. Gibiansky. <http://andrew.gibiansky.com/blog/machine-learning/speech-recognition-neural-networks/>
11. LSTM Mathematic formalism ( mandarin )  
<http://blog.csdn.net/u010754290/article/details/47167979>
11. Assumption of Markov Model  
<http://jedlik.phy.bme.hu/~gerjanos/HMM/node5.html>
12. Story about ANN-HMM Hybrid ( mandarin )  
HMM : <http://www.taodocs.com/p-5260781.html>
13. Generative model v.s. discriminative model  
<http://stackoverflow.com/questions/879432/what-is-the-difference-between-a-generative-and-discriminative-algorithm>

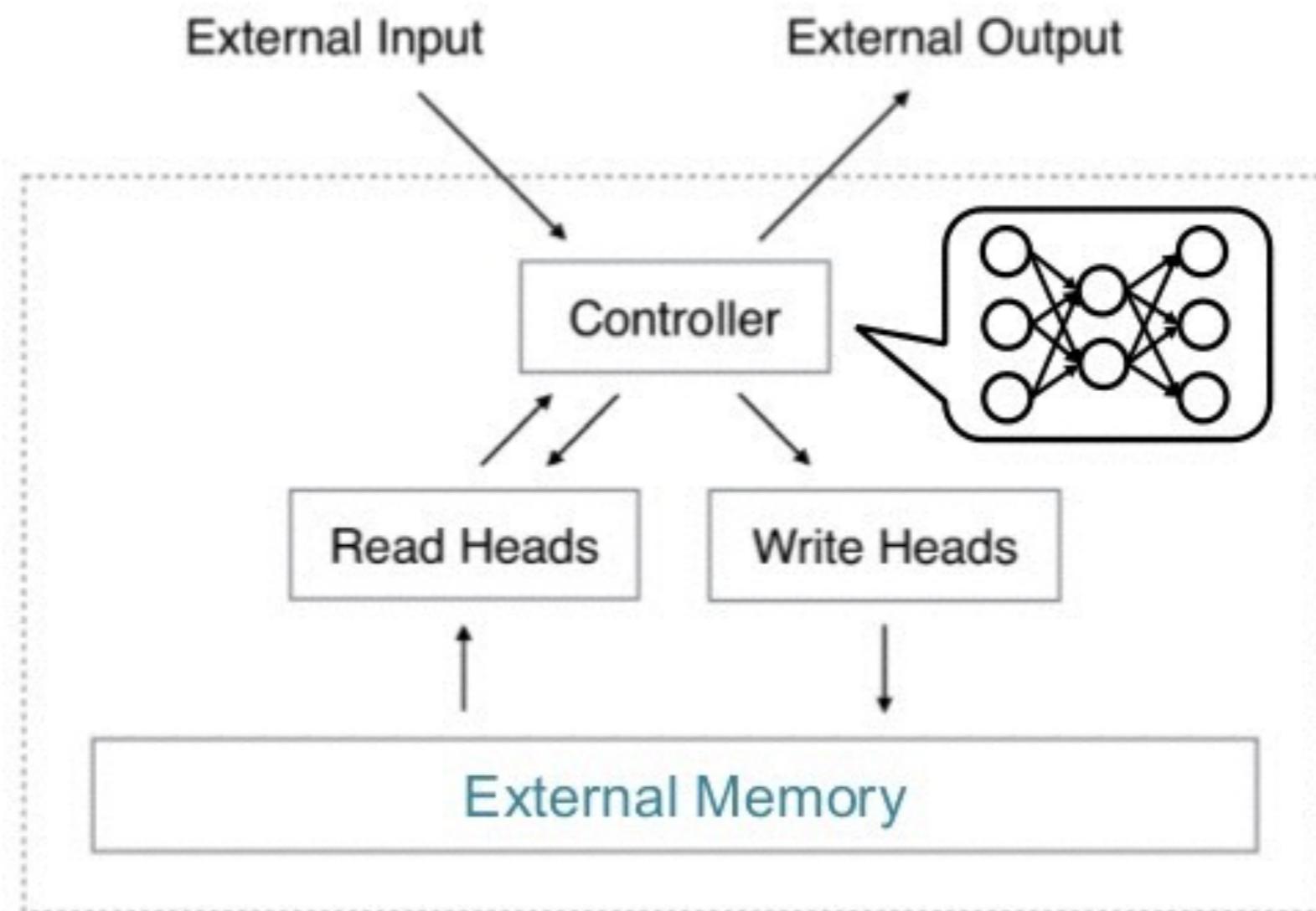
## Memory Network :

### A. Grave. Neural Turing Machine

#### Neural Turing Machine

- "Neural Turing Machine" is NN which has the capability of coupling to the *external memories*.

(Controller is NN with parameters for coupling to external memories)



## Appendix A. Generative v.s. discriminative

A **generative** model learns the **joint** probability  $p(x,y)$ .  
A **discriminative** model learns the **conditional** probability  $p(y|x)$ .

data input : (1,0), (1,0), (2,0), (2,1)

Generative

	y=0	y=1
x=1	1/2	0
x=2	1/4	1/4

Discriminative

	y=0	y=1
x=1	1	0
x=2	1/2	1/2

According to Andrew Y. Ng. On Discriminative vs. Generative classifier: A comparison of logistic regression and naive Bayes

The overall gist is that **discriminative models generally outperform generative models in classification tasks.**