

LSTM Back Propagation with code alongside samples

By Mohit Kumar

LSTM(Long Short Term Memory):Internals:Back Prop

#reverse

dh_next = np.zeros_like(h) #dh from the next character

dC_next = np.zeros_like(c)

dx=np.zeros_like(z).T

dy=yhat.copy()

dy=dy-np.reshape(symbols_out onehot,[-1,1])

dWy=np.dot(dy,h.T)

dBy=dy

dht=np.dot(dy.T,self.WOUT.T).T

for t in reversed(range(seq)):

z= zt[t].T

dht=dht+dh_next

dot=np.multiply(dht,self.tanh_array(ct[t])*self.dsigmoid(ot[t]))

dct=np.multiply(dht,ot[t]*self.dtanh(ct[t]))+dC_next

dcproj=np.multiply(dct,it[t]*(1-cproj[t]*cproj[t]))

dft=np.multiply(dct,coldt[t]*self.dsigmoid(ft[t]))

dit=np.multiply(dct,cproj[t]*self.dsigmoid(it[t]))

$$\frac{dE_t}{dB_{y_t}} = \frac{\partial E_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial pred_t}$$

$$\frac{dE_t}{dB_{y_t}} = \hat{y}_t - y_t$$

$$E_t = crossEntropyLoss(\hat{y}_t, y_t)$$

$$\hat{y}_t = softmax(pred_t)$$

$$pred_t = (h_t \cdot W_y) + b_y$$

$$h_t = o_t * tanh(c_t)$$

$$c_t = (f_t * c_{t-1}) + (i_t * \tilde{c}_t)$$

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

$$\tilde{c}_t = \tanh(W_c.[h_{t-1}, x_t] + b_c)$$

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

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

LSTM(Long Short Term Memory):Internals:Back Prop

#reverse

dh_next = np.zeros_like(h) #dh from the next character

dC_next = np.zeros_like(c)

dx=np.zeros_like(z).T

dy=yhat.copy()

dy=dy-np.reshape(symbols_out_onehot,[-1,1])

dWy=np.dot(dy,h.T)

dBy=dy

dht=np.dot(dy.T,self.WOUT.T).T

for t in reversed(range(seq)):

z= zt[t].T

dht=dht+dh_next

dot=np.multiply(dht,self.tanh_array(ct[t])*self.dsigmoid(ot[t]))

dct=np.multiply(dht,ot[t]*self.dtanh(ct[t]))+dC_next

dcproj=np.multiply(dct,it[t]*(1-cproj[t]*cproj[t]))

dft=np.multiply(dct,coltd[t]*self.dsigmoid(ft[t]))

dit=np.multiply(dct,cproj[t]*self.dsigmoid(it[t]))

$$\frac{dE_t}{dW_{y_t}} = \frac{\partial E_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial pred_t} \cdot \frac{\partial pred_t}{\partial W_{y_t}}$$

$$\frac{dE_t}{dW_{y_t}} = (\hat{y}_t - y_t) \cdot h_t$$

$$E_t = crossEntropyLoss(\hat{y}_t, y_t)$$

$$\hat{y}_t = softmax(pred_t)$$

$$pred_t = (h_t \cdot W_y) + b_y$$

$$h_t = o_t * tanh(c_t)$$

$$c_t = (f_t * c_{t-1}) + (i_t * \tilde{c}_t)$$

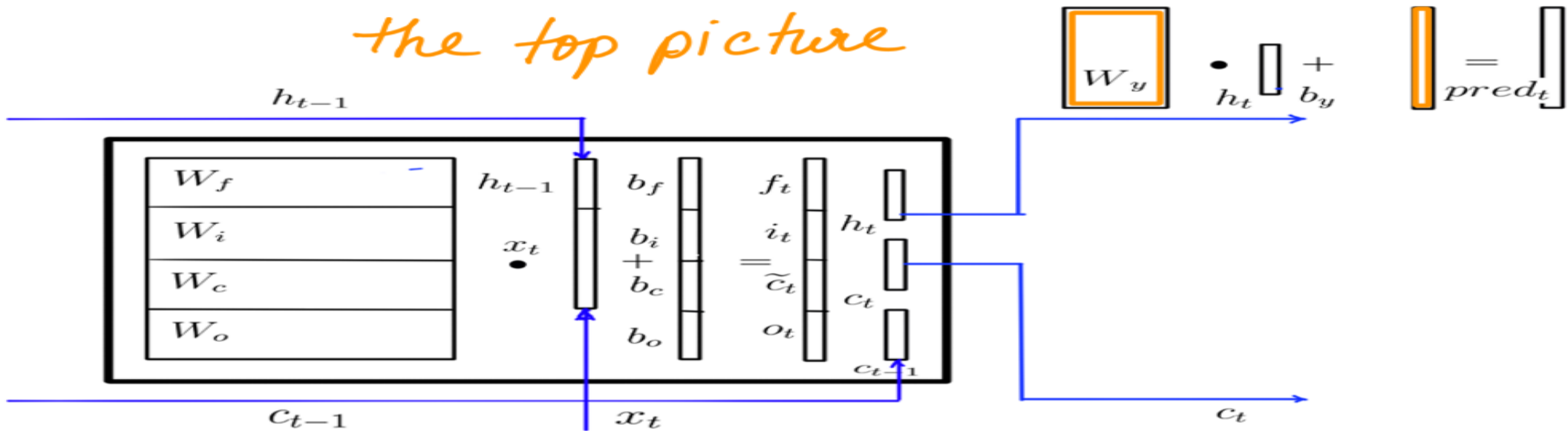
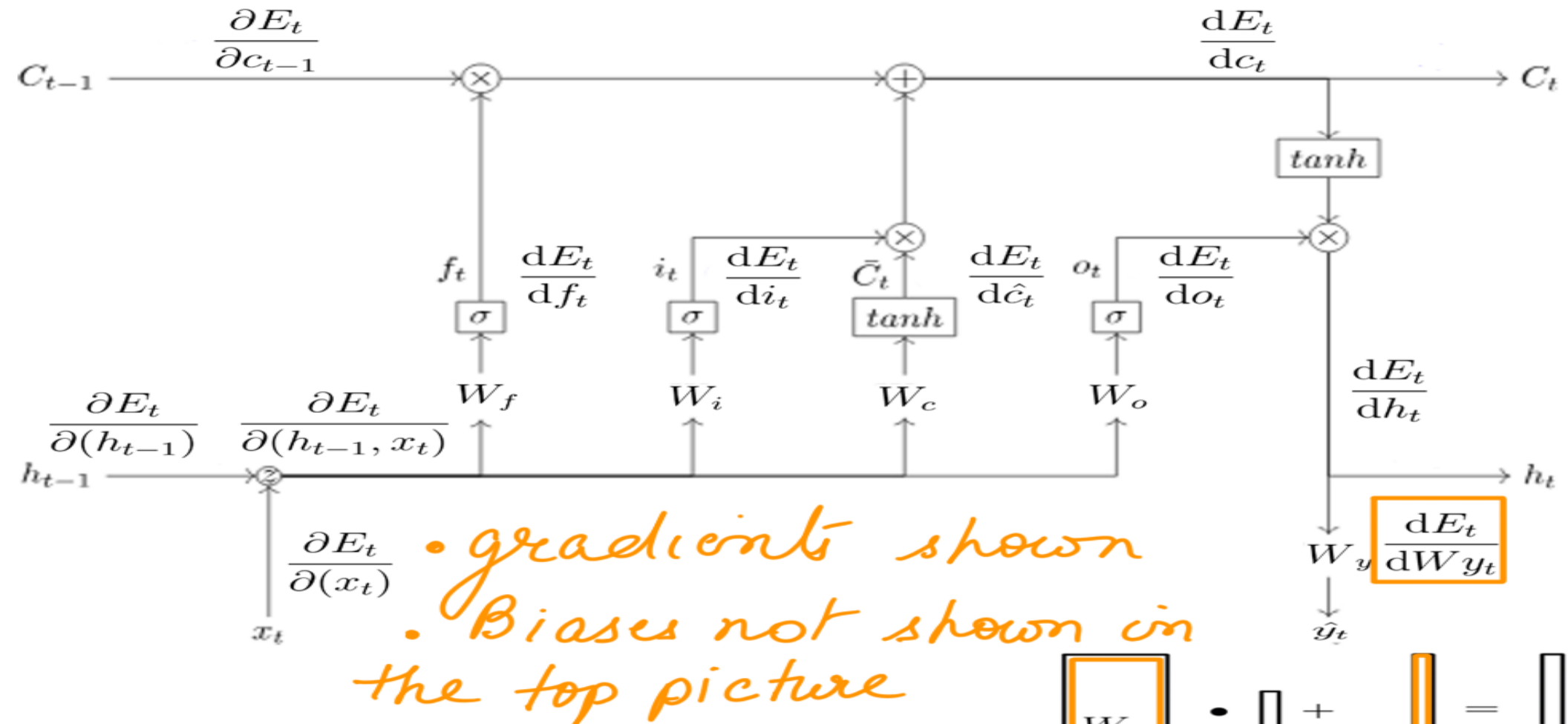
$$o_t = \sigma(W_o.[h_{t-1}, x_t] + b_o)$$

$$\tilde{c}_t = \tanh(W_c.[h_{t-1}, x_t] + b_c)$$

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

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

LSTM(Long Short Term Memory):Internals:Back Prop



LSTM(Long Short Term Memory):Internals:Back Prop

#reverse

dh_next = np.zeros_like(h) #dh from the next character

dC_next = np.zeros_like(c)

dx=np.zeros_like(z).T

dy=yhat.copy()

dy=dy-np.reshape(symbols_out_onehot,[-1,1])

dWy=np.dot(dy,h.T)

dBy=dy

dht=np.dot(dy.T,self.WOUT.T).T

for t in reversed(range(seq)):

z= zt[t].T

dht=dht+dh_next

dot=np.multiply(dht,self.tanh_array(ct[t])*self.dsigmoid(ot[t]))

dct=np.multiply(dht,ot[t]*self.dtanh(ct[t]))+dC_next

dcproj=np.multiply(dct,it[t]*(1-cproj[t]*cproj[t]))

dft=np.multiply(dct,coldt[t]*self.dsigmoid(ft[t]))

dit=np.multiply(dct,cproj[t]*self.dsigmoid(it[t]))

$$\frac{dE_t}{dh_t} = \frac{\partial E_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial pred_t} \cdot \frac{\partial pred_t}{\partial h_t}$$

$$\frac{dE_t}{dh_t} = (\hat{y}_t - y_t \cdot W_y) + dh_next$$

$$E_t = crossEntropyLoss(\hat{y}_t, y_t)$$

$$\hat{y}_t = softmax(pred_t)$$

$$pred_t = (h_t \cdot W_y) + b_y$$

$$h_t = o_t * tanh(c_t)$$

$$c_t = (f_t * c_{t-1}) + (i_t * \tilde{c}_t)$$

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

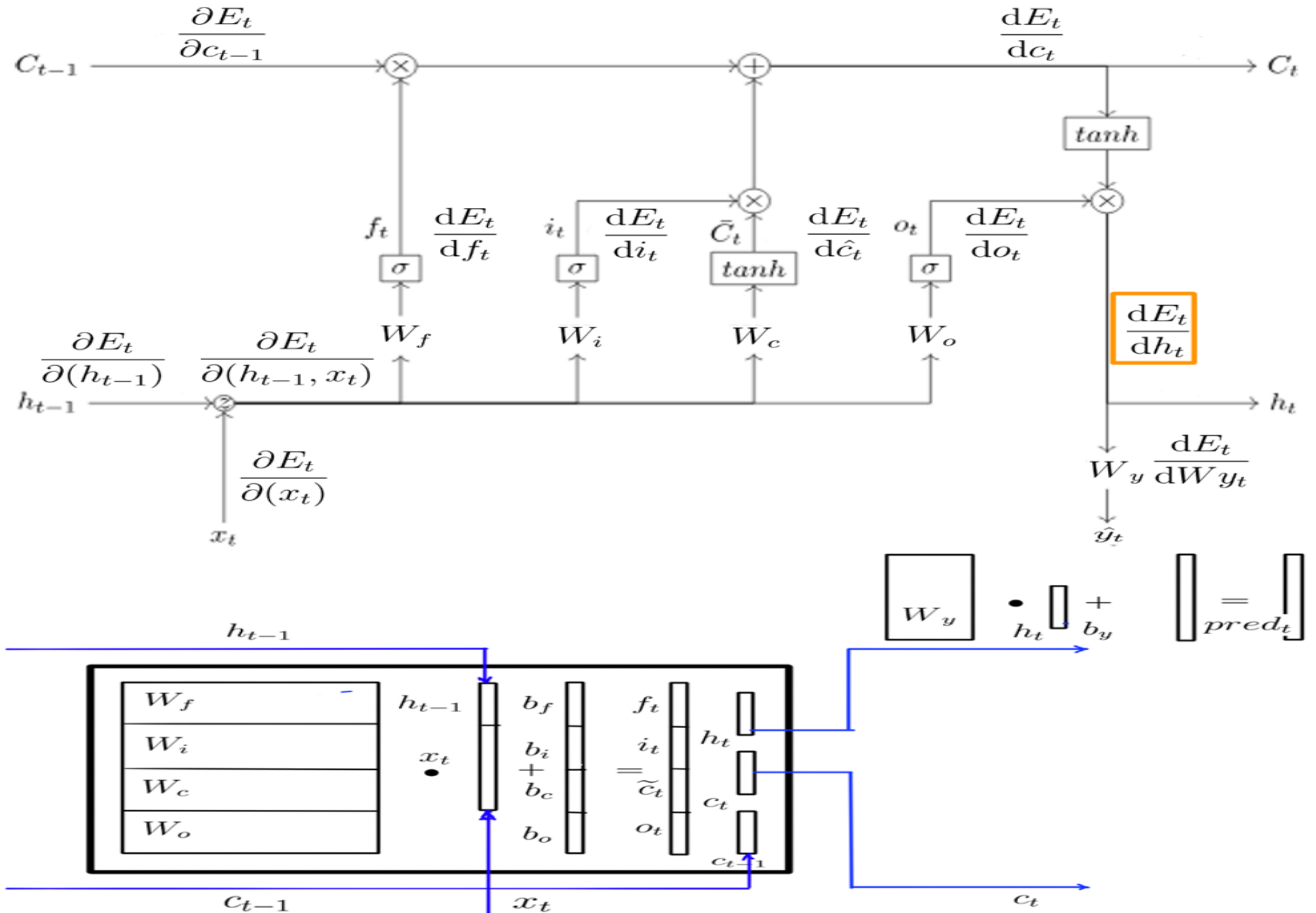
$$\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

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

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

- Wout is Wy.
- Recurrent element
- zeros to begin with
- Done sequence number of times.

LSTM(Long Short Term Memory):Internals:Back Prop



LSTM(Long Short Term Memory):Internals:Back Prop

#reverse

```
dh_next = np.zeros_like(h) #dh from the next character
dC_next = np.zeros_like(c)
dx=np.zeros_like(z).T
dy=yhat.copy()
dy=dy-np.reshape(symbols_out_onehot,[-1,1])
dWy=np.dot(dy,h.T)
dBy=dy
dht=np.dot(dy.T,self.WOUT.T).T
for t in reversed(range(seq)):
    z= zt[t].T
    dht=dht+dh_next
    dot=np.multiply(dht,self.tanh_array(ct[t])*self.dsigmoid(ot[t]))
    dct=np.multiply(dht,ot[t]*self.dtanh(ct[t]))+dC_next
    dcproj=np.multiply(dct,it[t]*(1-cproj[t]*cproj[t]))
    dft=np.multiply(dct,coldt[t]*self.dsigmoid(ft[t]))
    dit=np.multiply(dct,cproj[t]*self.dsigmoid(it[t]))
```

$$\frac{dE_t}{do_t} = \frac{\partial E_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial o_t}$$

$$\frac{dE_t}{do_t} = \frac{\partial E_t}{\partial h_t} \cdot \tanh(c_t) \cdot \text{dsigmoid}(o_t)$$

```
def dsigmoid(self,f):
    return f*(1-f)
```

$$E_t = \text{crossEntropyLoss}(\hat{y}_t, y_t)$$

$$\hat{y}_t = \text{softmax}(\text{pred}_t)$$

$$\text{pred}_t = (h_t \cdot W_y) + b_y$$

$$h_t = o_t * \tanh(c_t)$$

$$c_t = (f_t * c_{t-1}) + (i_t * \tilde{c}_t)$$

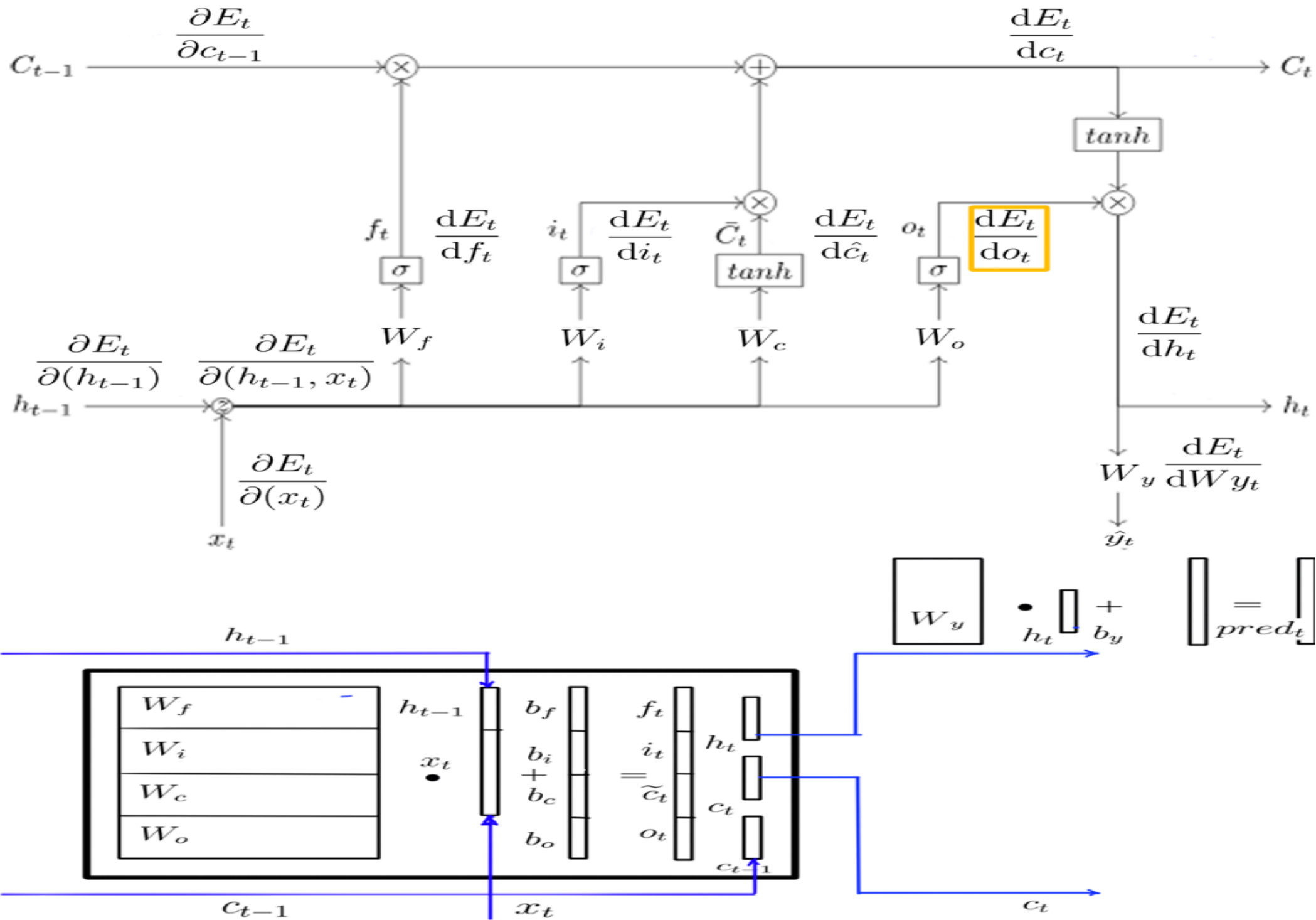
$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

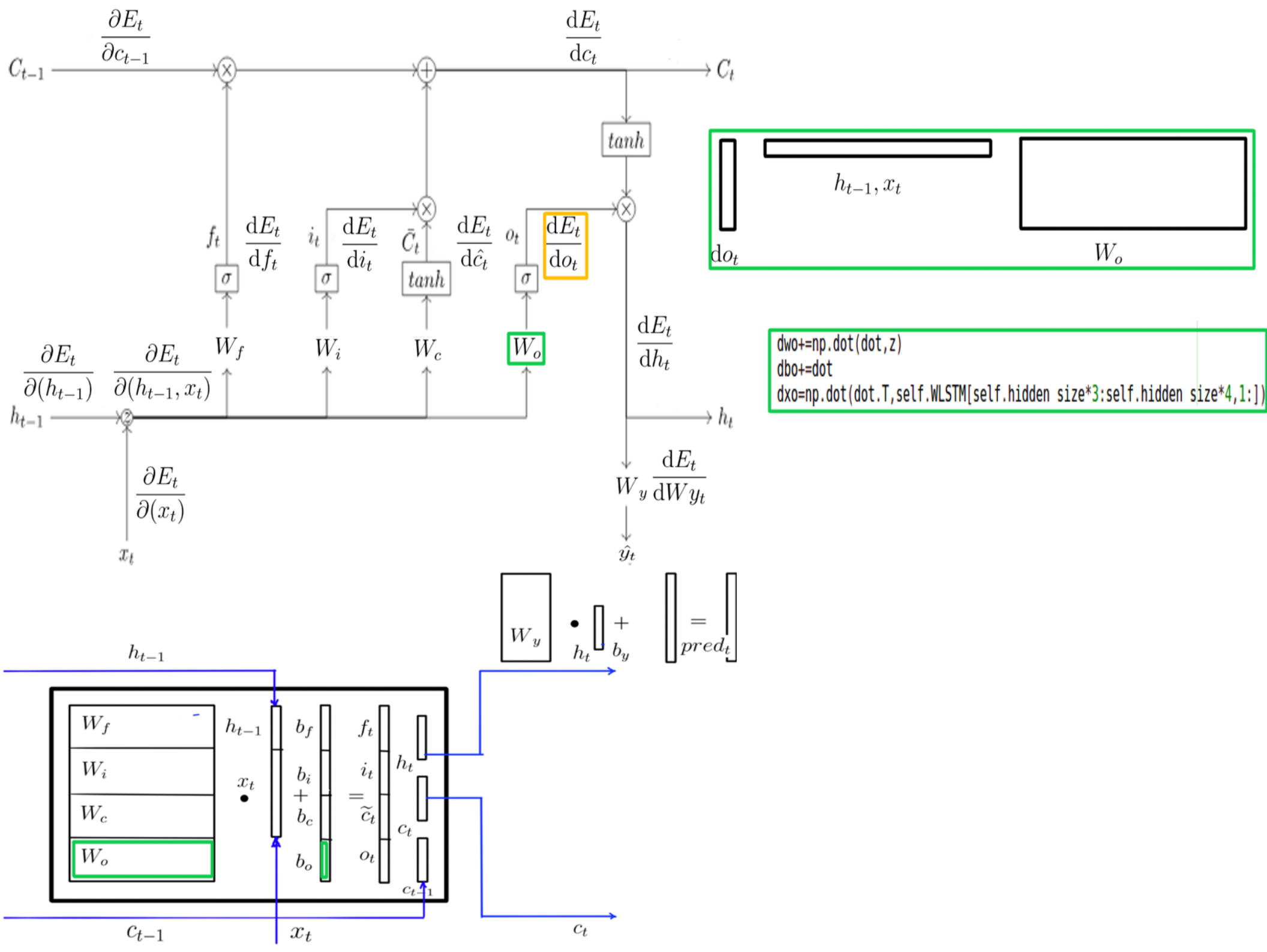
$$\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

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

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

LSTM(Long Short Term Memory):Internals:Back Prop





LSTM(Long Short Term Memory):Internals:Back Prop

#reverse

dh_next = np.zeros_like(h) #dh from the next char

dC_next = np.zeros_like(c)

dx=np.zeros_like(z).T

dy=yhat.copy()

dy=dy-np.reshape(symbols_out_onehot,[-1,1])

dWy=np.dot(dy,h.T)

dBy=dy

dht=np.dot(dy.T,self.WOUT.T).T

for t in reversed(range(seq)):

z= zt[t].T

dht=dht+dh_next

dot=np.multiply(dht,self.tanh_array(ct[t])*self.dsigmoid(ot[t]))

dct=np.multiply(dht,ot[t]*self.dtanh(ct[t]))+dC_next

dcproj=np.multiply(dct,it[t]*(1-cproj[t]*cproj[t]))

dft=np.multiply(dct,coldt[t]*self.dsigmoid(ft[t]))

dit=np.multiply(dct,cproj[t]*self.dsigmoid(it[t]))

$$\frac{dE_t}{dc_t} = \frac{\partial E_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial c_t}$$

$$\frac{dE_t}{dc_t} = \frac{\partial E_t}{\partial h_t} \cdot o_t \cdot dtanh(c_t) + dc_next$$

• Because of recurrent cell state c_{t-1} .

```
def dtanh(self,f):
    tanhf=np.tanh(f)
    return 1 - tanhf * tanhf
```

$$E_t = crossEntropyLoss(\hat{y}_t, y_t)$$

$$\hat{y}_t = softmax(pred_t)$$

$$pred_t = (h_t \cdot W_y) + b_y$$

$$h_t = o_t * tanh(c_t)$$

$$c_t = (f_t * c_{t-1}) + (i_t * \tilde{c}_t)$$

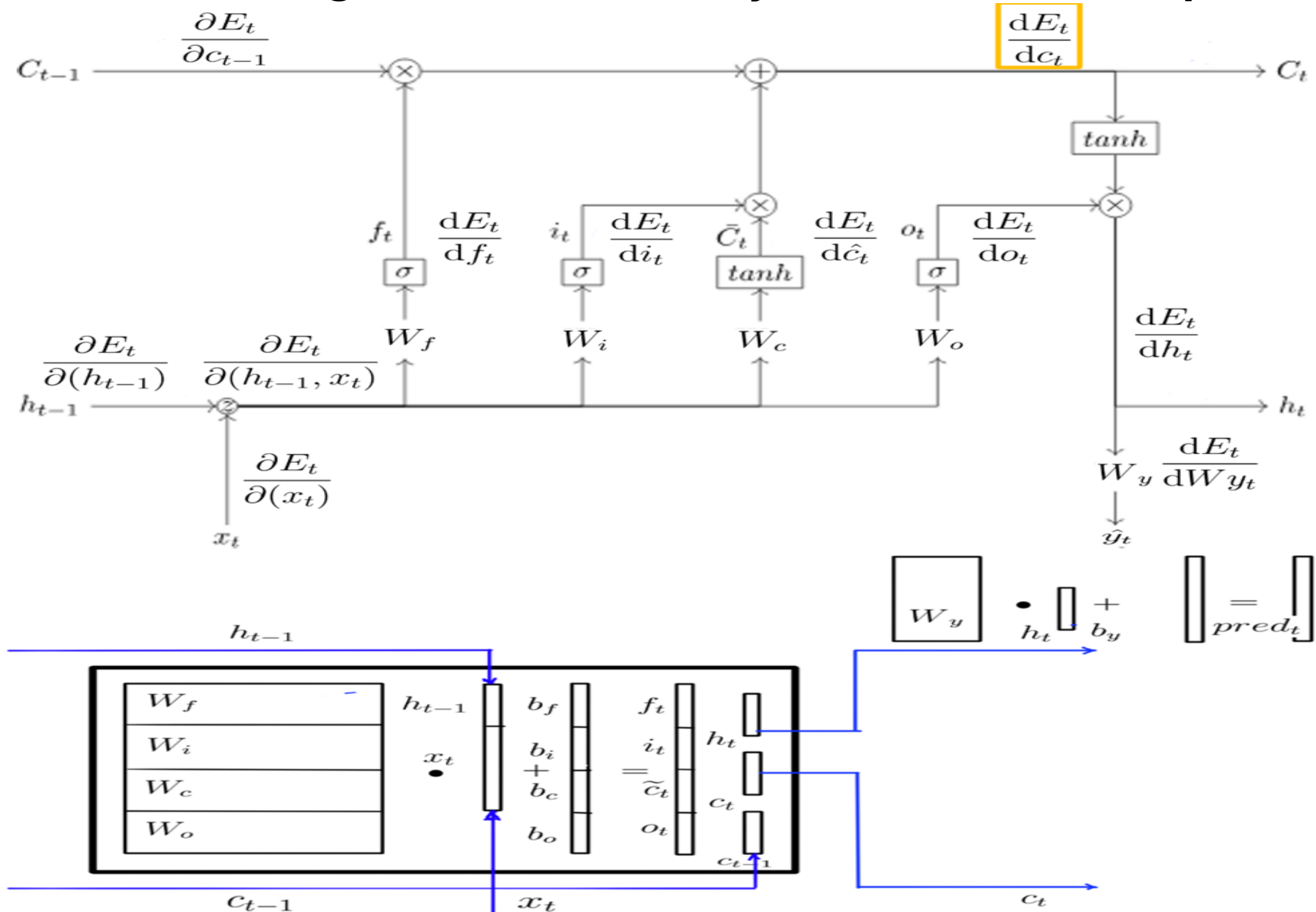
$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

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

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

LSTM(Long Short Term Memory):Internals:Back Prop



LSTM(Long Short Term Memory):Internals:Back Prop

#reverse

dh_next = np.zeros_like(h) #dh from the next character

dC_next = np.zeros_like(c)

dx=np.zeros_like(z).T

dy=yhat.copy()

dy=dy-np.reshape(symbols_out_onehot,[-1,1])

dWy=np.dot(dy,h.T)

dBy=dy

dht=np.dot(dy.T,self.WOUT.T).T

for t in reversed(range(seq)):

z= zt[t].T

dht=dht+dh_next

dot=np.multiply(dht,self.tanh_array(ct[t])*self.dsigmoid(ot[t]))

dct=np.multiply(dht,ot[t]*self.dtanh(ct[t]))+dC_next

dcproj=np.multiply(dct,it[t]*[1-cproj[t]*cproj[t]])

dft=np.multiply(dct,coldt[t]*self.dsigmoid(ft[t]))

dit=np.multiply(dct,cproj[t]*self.dsigmoid(it[t]))

$$\frac{dE_t}{d\hat{c}_t} = \frac{\partial E_t}{\partial c_t} \cdot \frac{\partial c_t}{\partial \hat{c}_t}$$

$$\frac{dE_t}{d\hat{c}_t} = \frac{\partial E_t}{\partial c_t} \cdot i_t \cdot dtanh(\hat{c}_t)$$

• Previously calculated
dct

• inlined dtanh

• \hat{c} is cproj.

$$E_t = crossEntropyLoss(\hat{y}_t, y_t)$$

$$\hat{y}_t = softmax(pred_t)$$

$$pred_t = (h_t \cdot W_y) + b_y$$

$$h_t = o_t * tanh(c_t)$$

$$c_t = (f_t * c_{t-1}) + (i_t * \tilde{c}_t)$$

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

$$\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

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

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

LSTM(Long Short Term Memory):Internals:Back Prop

