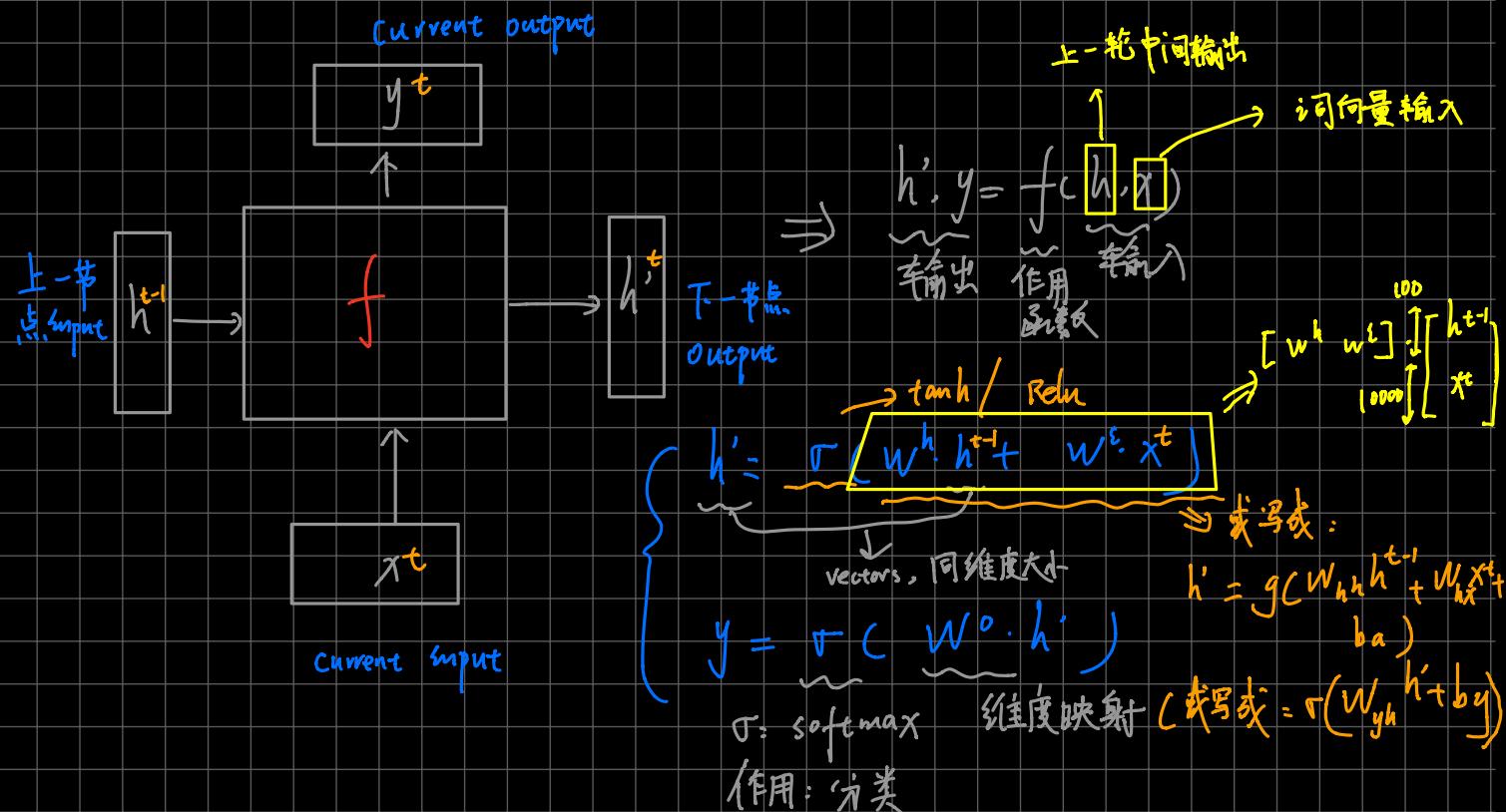
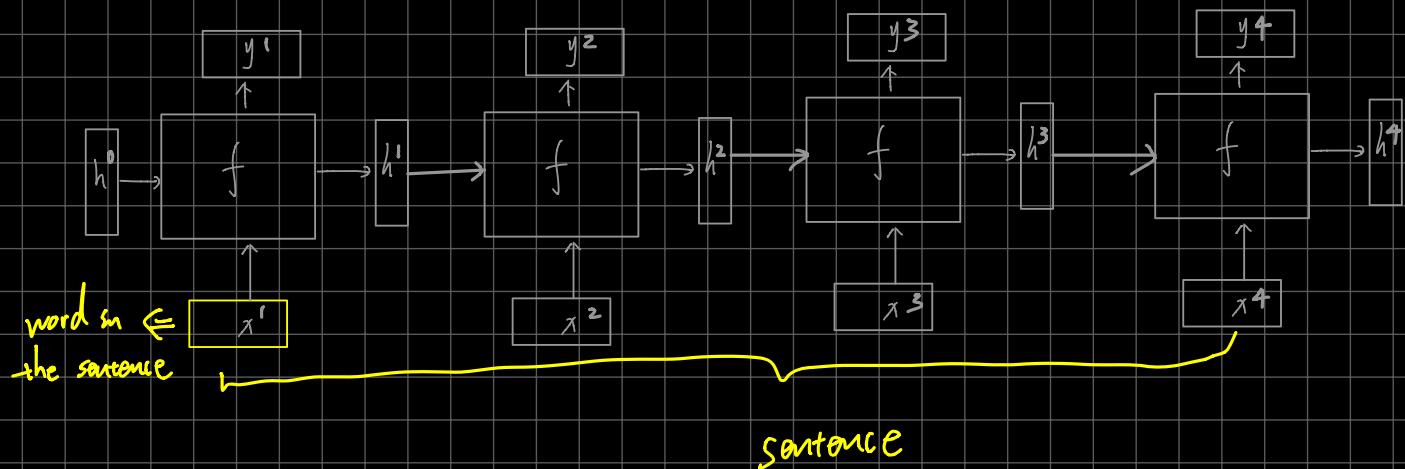


Forward Propagation

I. Naive RNN:



△ 不论 input/output 序列有多长，我们只需要 1 个 f 函数。



II LSTM:

为何提出 LSTM:

解决 RNN 的 梯度消失问题, 在更长的序列中有更好的表现。

Gradient Vanish

然而 LSTM 并不能解决 梯度爆炸问题。 \Rightarrow 梯度剪裁算法

Gradient explosion

Gradient Clipping Algorithm.

梯度上限剪裁: 梯度的某个维度大于某个上限 \Rightarrow
剪裁为上限

Gradient clipping

梯度同比收缩: 梯度 L_2 范数 大于上限 \Rightarrow

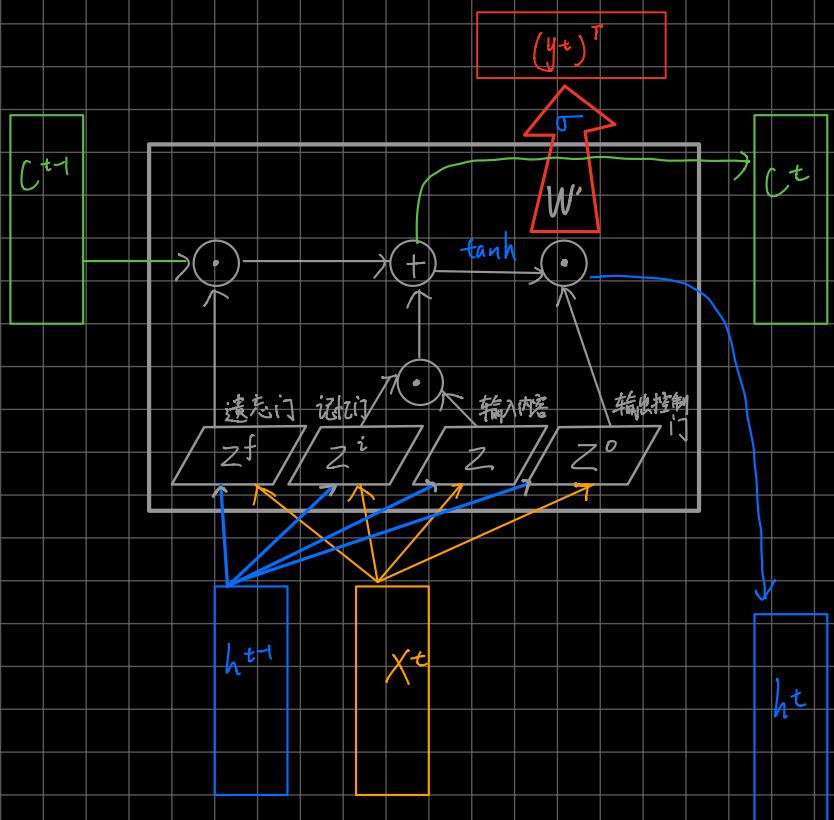
$$\text{梯度} = \frac{\text{梯度}}{\text{范数}}$$

△ 在 RNN 及其变种 (LSTM, GRU) 网络中常出现
梯度爆炸问题 \Rightarrow 梯度随步长累加。

△ 组成 RNN 单元 (LSTM, GRU...) 的组合 (z^f, z^i, z^o 等)

最好理解为对探索的一种约束, 而非一种工程意义上
的设计 (e.g 遗忘门, 记忆门等)。

\Rightarrow 因为 RNN 单元的作用由权重来决定, 权重是通过端到端的
学习来调整的, 因此 RNN 单元类型不同只能说是对探索的
限制不同, 并非是具有特定的工程意义。



門控碼狀，即逐元素乘積。

$$C^t = z^f \odot C^{t-1} + z^i \odot z^o \odot h^t$$

$$h^t = \tanh(C^t) \odot z^o$$

$$y^t = \sigma(W^T \cdot h^t)$$

Input ↓

門控狀態 $(0,1)$

$$z^f = \sigma(W^f \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix})$$

$$z^i = \sigma(W^i \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix})$$

輸入內容 $\in [-1,1]$

$$z^o = \sigma(W^o \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix})$$

C^t : cell state

相当于RNN中 h^t , 主
要用于保存先前结点数据, 改变很慢, 有利于长期记忆

h^t : hidden state

在不同节点下有很
大区别, 有利于
短期记忆.

LSTM 内部处理三阶段:

I. 遗忘阶段 \Rightarrow 通过遗忘门 Z^f 控制 y , 忘记不重要的

II. 选择记忆阶段 \Rightarrow 通过选择记忆门 Z^i 控制, 记住重要的

III. 车输出阶段 \Rightarrow 通过输出控制门 Z^o 控制, 决定哪些会被当成车输出.

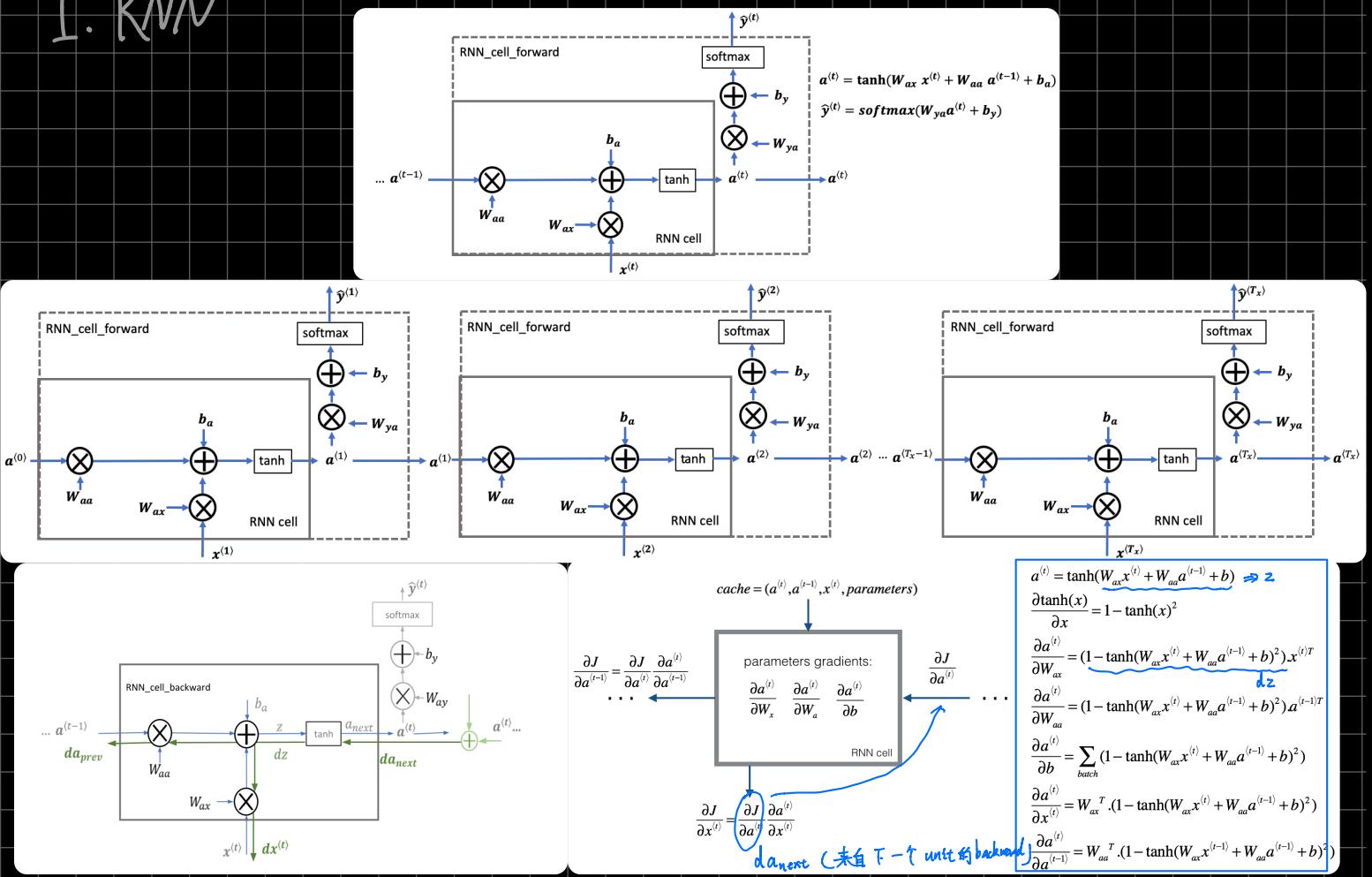


LSTM 相比 RNN 引入参数变多, 输入、输出多了传输状态

使训练难度增大许多 \Rightarrow CRU 模型, 参数更少, 更容易训练

Backward Propagation

I. RNN



$$\frac{\partial J}{\partial W_{ax}} = \frac{\partial J}{\partial a^{(t)}} \cdot \frac{\partial a^{(t)}}{\partial W_{ax}}$$

da_{next} \cdot (\dots)

$a^{(t)} = \tanh(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b_a) \quad (-)$

$\frac{\partial \tanh(x)}{\partial x} = 1 - \tanh^2(x) \quad (-)$

$\frac{\partial J}{\partial W_{ax}} = (da_{next} * (1 - \tanh^2(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b_a)))x^{(t)T} \quad (1)$

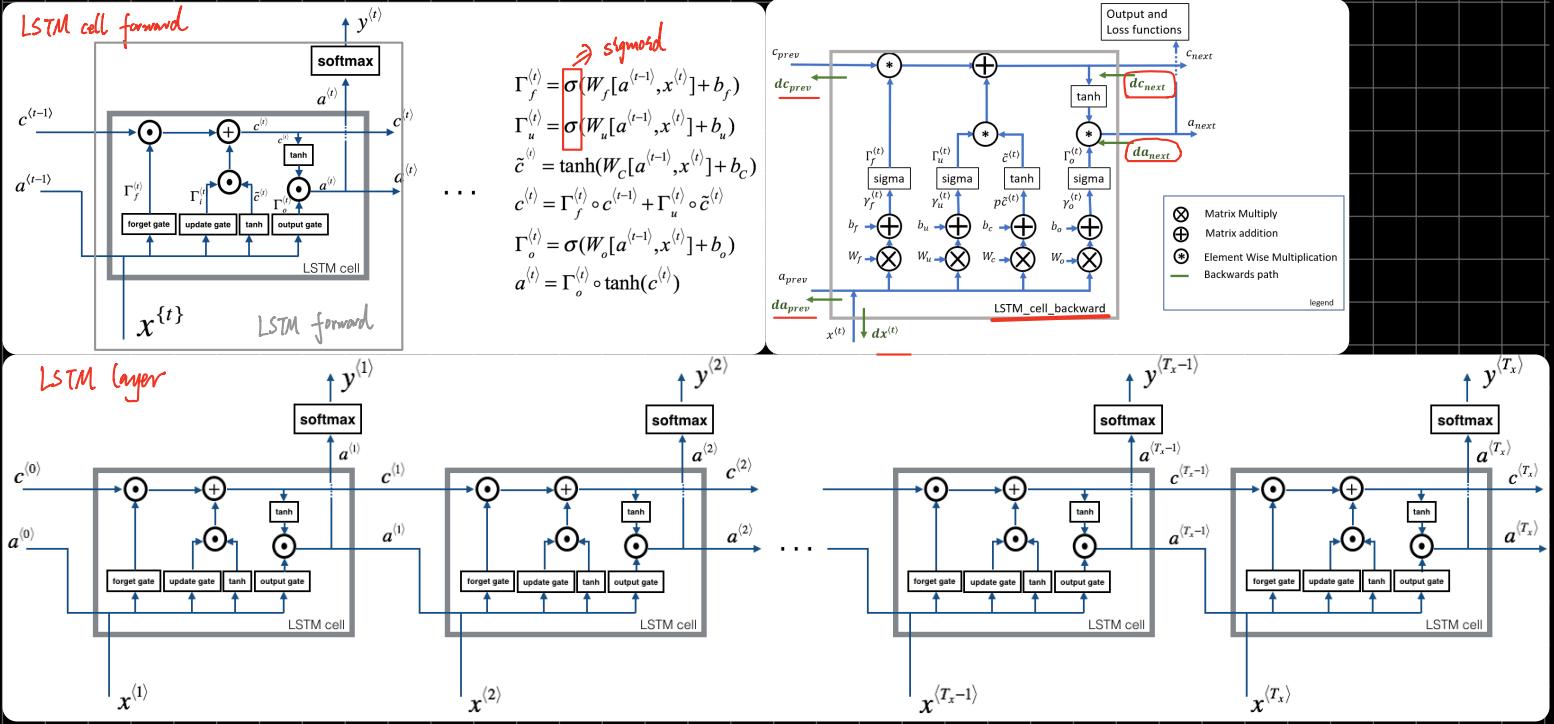
$\frac{\partial J}{\partial W_{aa}} = (da_{next} * (1 - \tanh^2(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b_a)))a^{(t-1)T} \quad (2)$

$\frac{\partial J}{\partial b_a} = \sum_{batch} (da_{next} * (1 - \tanh^2(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b_a))) \quad (3)$

$\frac{\partial J}{\partial x^{(t)}} = W_{ax}^T (da_{next} * (1 - \tanh^2(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b_a))) \quad (4)$

$da_{prev} = W_{aa}^T (da_{next} * (1 - \tanh^2(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b_a))) \quad (5)$

II. LSTM



LSTM Backward: 分 0. I II. III 共 3 步计算，输入为 da_{next} , dc_{next} , 输出为 $[da_{prev}^T, dc_{prev}^T, dx_t^T]$

3.2.2 gate derivatives

Note the location of the gate derivatives ($\gamma_{\cdot}^{(t)}$) between the dense layer and the activation function (see graphic above). This is convenient for computing parameter derivatives in the next step.

$$\begin{aligned} d\gamma_f^{(t)} &= da_{next} * \tanh(c_{next}) * \Gamma_o^{(t)} * (1 - \Gamma_o^{(t)}) \\ d\gamma_u^{(t)} &= (dc_{next} * \Gamma_u^{(t)} + \Gamma_o^{(t)} * (1 - \tanh^2(c_{next})) * \tilde{c}^{(t)} * da_{next}) * (1 - (\tilde{c}^{(t)})^2) \\ d\gamma_c^{(t)} &= (dc_{next} * \tilde{c}^{(t)} + \Gamma_o^{(t)} * (1 - \tanh^2(c_{next})) * \tilde{c}^{(t)} * da_{next}) * \Gamma_c^{(t)} * (1 - \Gamma_c^{(t)}) \\ d\gamma_o^{(t)} &= (dc_{next} * c_{prev} + \Gamma_o^{(t)} * (1 - \tanh^2(c_{next})) * c_{prev} * da_{next}) * \Gamma_f^{(t)} * (1 - \Gamma_f^{(t)}) \end{aligned}$$

3.2.3 parameter derivatives

$$\begin{aligned} \text{e.g.: } dW_f &= \frac{\partial J}{\partial W_f} = \frac{\partial J}{\partial a_{next}} \cdot \frac{\partial a_{next}}{\partial W_f} = \frac{\partial J}{\partial a_{next}} \cdot \frac{\partial \gamma_f^{(t)}}{\partial W_f} \cdot \frac{\partial \gamma_f^{(t)}}{\partial a_{next}} = \frac{\partial J}{\partial a_{next}} \cdot \frac{\partial \gamma_f^{(t)}}{\partial W_f} \cdot X \cdot [a^{(t-1)} \cdot x_t]^T \\ dW_u &= d\gamma_u^{(t)} \left[\begin{array}{c|c} a_{prev} & x_t \end{array} \right]^T \\ dW_c &= d\gamma_c^{(t)} \left[\begin{array}{c|c} a_{prev} & x_t \end{array} \right]^T \\ dW_o &= d\gamma_o^{(t)} \left[\begin{array}{c|c} a_{prev} & x_t \end{array} \right]^T \end{aligned}$$

To calculate db_f, db_u, db_c, db_o you just need to sum across all 'm' examples (axis=1) on $d\gamma_f^{(t)}, d\gamma_u^{(t)}, d\gamma_c^{(t)}, d\gamma_o^{(t)}$ respectively. Note that you should have the keepdims = True option.

$$\begin{cases} (7) \\ (8) \\ (9) \\ (10) \end{cases} \quad \sum \uparrow \Rightarrow \text{II}$$

$$db_f = \sum_{batch} d\gamma_f^{(t)} \quad (15)$$

$$db_u = \sum_{batch} d\gamma_u^{(t)} \quad (16)$$

$$db_c = \sum_{batch} d\gamma_c^{(t)} \quad (17)$$

$$db_o = \sum_{batch} d\gamma_o^{(t)} \quad (18)$$

Finally, you will compute the derivative with respect to the previous hidden state, previous memory state, and input.

$$III. \text{ 最终要backward } da_{prev} = W_f^T d\gamma_f^{(t)} + W_u^T d\gamma_u^{(t)} + W_c^T d\gamma_c^{(t)} + W_o^T d\gamma_o^{(t)} \quad (19)$$

Here, to account for concatenation, the weights for equations 19 are the first n_a, (i.e. $W_f = W_f[:, :n_a]$ etc...)

$$dc_{prev} = dc_{next} * \Gamma_f^{(t)} + \Gamma_o^{(t)} * (1 - \tanh^2(c_{next})) * \Gamma_f^{(t)} * da_{next} \quad (20)$$

$$dx^{(t)} = W_f^T d\gamma_f^{(t)} + W_u^T d\gamma_u^{(t)} + W_c^T d\gamma_c^{(t)} + W_o^T d\gamma_o^{(t)} \quad (21)$$

where the weights for equation 21 are from n_a to the end, (i.e. $W_f = W_f[:, n_a:]$ etc...)

△ Note: 在 RNN, LSTM 的 Backward 过程中，计算出的 gradients 均需保存于 cache 之中。等到所有该层的 units 均通过 Backward 计算之后，再通过： I. Gradient Clipping (此对 RNN 非常重要) II. update parameters.

实例：① 恐龙名字生成 project

② Jazz 音乐生成 project