

3.RNN& LSTM& GRU

1.RNN

循环神经网络(Recurrent Neural network, RNN) 是一种用于处理序列数据的神经网络。相比与一般的神经网络，其很适合用于处理序列变化的数据。

1.1RNN的基本结构

由于本质是处理序列数据（一般按时间顺序，也有可能按照文本顺序）。其基本结构如下：



现在看上去就比较清楚了，这个网络在t时刻接收到输入 x_t 之后，隐藏层的值（这个时刻下对应的状态值）是 S_t ，输出值是 O_t 。

$$S_t = f(U \cdot X_t + W \cdot S_{t-1})$$

$$O_t = g(V \cdot S_t)$$

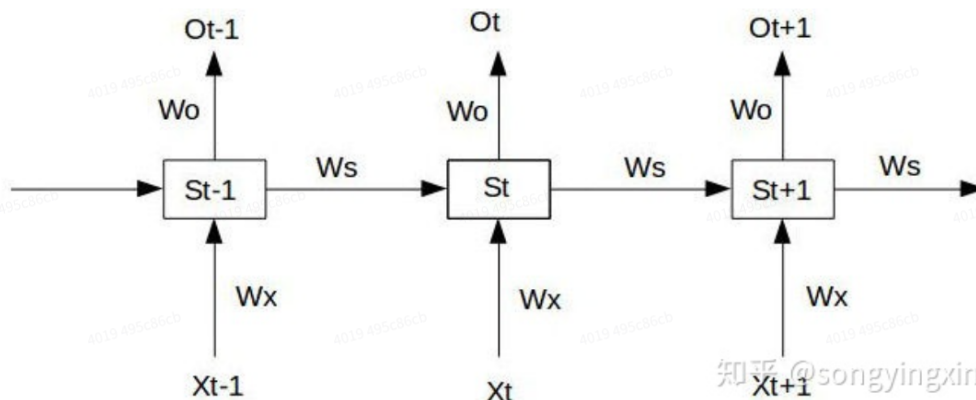
从公式中，可以看出， S_t 的值不仅仅取决于 x_t ，也取决于 S_{t-1} 。

同时，每个神经元都会接收上一个神经元的输出（其实这些神经元都是相同的，只有三个参数矩阵（ U, W, V ），每一层的参数相同，只是使用的状态不一样）。神经元的输出重新作为输入，因此将其称为循环神经网络。

1.2RNN的问题（梯度消失与梯度爆炸）

我们给定一个三个时间的RNN单元，如下：

我们假设 S_0 为给定值，且神经元中没有激活函数（便于分析），则前向过程如下：



我们假设最左端的输入 S_0 为给定值，且神经元中没有激活函数（便于分析），则前向过程如下：

$$\begin{aligned} S_1 &= W_x X_1 + W_s S_0 + b_1 & O_1 &= W_o S_1 + b_2 \\ S_2 &= W_x X_2 + W_s S_1 + b_1 & O_2 &= W_o S_2 + b_2 \\ S_3 &= W_x X_3 + W_s S_2 + b_1 & O_3 &= W_o S_3 + b_2 \end{aligned}$$

在 $t = 3$ 时刻，损失函数为 $L_3 = \frac{1}{2}(Y_3 - O_3)^2$ ，那么如果我们要训练RNN时，实际上就是对 W_x, W_s, W_o, b_1, b_2 求偏导，并不断调整它们以使得 L_3 尽可能达到最小（参见反向传播算法与梯度下降算法）。

那么我们得到以下公式：

$$\begin{aligned} \frac{\partial L_3}{\partial W_0} &= \frac{\partial L_3}{\partial O_3} \frac{\partial O_3}{\partial W_0} \\ \frac{\partial L_3}{\partial W_x} &= \frac{\partial L_3}{\partial O_3} \frac{\partial O_3}{\partial S_3} \frac{\partial S_3}{\partial W_x} + \frac{\partial L_3}{\partial O_3} \frac{\partial O_3}{\partial S_3} \frac{\partial S_3}{\partial S_2} \frac{\partial S_2}{\partial W_x} + \frac{\partial L_3}{\partial O_3} \frac{\partial O_3}{\partial S_3} \frac{\partial S_3}{\partial S_2} \frac{\partial S_2}{\partial S_1} \frac{\partial S_1}{\partial W_x} \\ \frac{\partial L_3}{\partial W_s} &= \frac{\partial L_3}{\partial O_3} \frac{\partial O_3}{\partial S_3} \frac{\partial S_3}{\partial W_s} + \frac{\partial L_3}{\partial O_3} \frac{\partial O_3}{\partial S_3} \frac{\partial S_3}{\partial S_2} \frac{\partial S_2}{\partial W_s} + \frac{\partial L_3}{\partial O_3} \frac{\partial O_3}{\partial S_3} \frac{\partial S_3}{\partial S_2} \frac{\partial S_2}{\partial S_1} \frac{\partial S_1}{\partial W_s} \end{aligned}$$

将上述偏导公式与第三节中的公式比较，我们发现，随着神经网络层数的加深对 W_0 而言并没有什么影响（ W_0 直接作用与 O ，与之间状态无关），而对 W_x, W_s 会随着时间序列的拉长而产生梯度消失和梯度爆炸问题。

根据上述分析整理一下公式可得，对于任意时刻 t 对 W_x, W_s 求偏导的公式为：

$$\begin{aligned} \frac{\partial L_t}{\partial W_x} &= \sum_{k=0}^t \frac{\partial L_t}{\partial O_t} \frac{\partial O_t}{\partial S_t} \left(\prod_{j=k+1}^t \frac{\partial S_j}{\partial S_{j-1}} \right) \frac{\partial S_k}{\partial W_x} \\ \frac{\partial L_t}{\partial W_s} &= \sum_{k=0}^t \frac{\partial L_t}{\partial O_t} \frac{\partial O_t}{\partial S_t} \left(\prod_{j=k+1}^t \frac{\partial S_j}{\partial S_{j-1}} \right) \frac{\partial S_k}{\partial W_s} \end{aligned}$$

我们发现，导致梯度消失和爆炸的就在于 $\prod_{j=k+1}^t \frac{\partial S_j}{\partial S_{j-1}}$ ，而加上激活函数后的 S 的表达式为：

$$S_j = \tanh(W_x X_j + W_s S_{j-1} + b_1)$$

那么则有：

$$\prod_{j=k+1}^t \frac{\delta S_j}{\delta S_{j-1}} = \prod_{j=k+1}^t \tanh' W_s$$

而在这个公式中， \tanh 的导数总是小于1的，如果 W_s 也是一个大于0小于1的值，那么随着 t 的增大，上述公式的值越来越趋近于0，这就导致了梯度消失问题。那么如果 W_s 很大，上述公式会越来越趋向于无穷，这就产生了梯度爆炸。

面试回答：

在RNN的参数更新过程中，利用的是反向传播的链式法则，但是由于网络结构太深，梯度反向传播中的连乘效应随着网络层数的增加而不断加剧，就带来了梯度消失与梯度爆炸的问题。

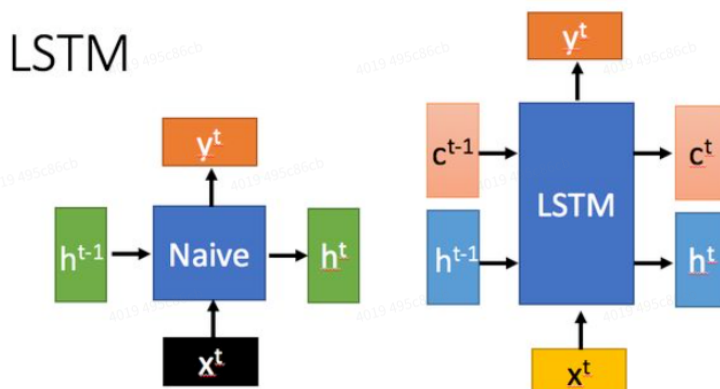
1.3为什么LSTM能解决梯度问题？

因为LSTM的门机制原因来解决了解决了RNN中的梯度消失问题。

因为门值的激活函数是sigmoid,使得三个门的输出要么接近0，要么接近1。当门为1时，梯度可以很好的在LSTM中进行传播，当门值为0的时候，说明上一时刻的信息对当前时刻无影响，也就没有必要进行梯度回传与更新。

2.LSTM的结构详解

2.1粗对比



c change slowly $\Rightarrow c^t$ is c^{t-1} added by something

h change faster $\Rightarrow h^t$ and h^{t-1} can be very different

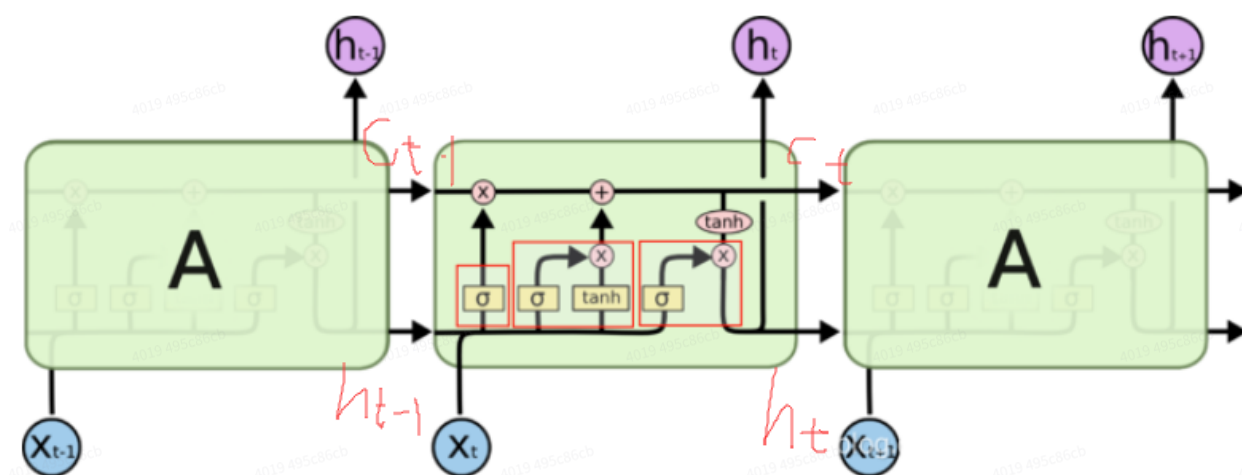
相比RNN只有一个传递状态 h_t ,LSTM有两个传输状态,一个 C_t (cell state)和一个 h_t (hidden state)。Tips,RNN中的 h_t 相对于LSTM中的 C_t 。

其中对于传递下去的 c_t 改变的很慢,通常输出的 C_t 是上一个状态传过来的 $C(t-1)$ 加上一些数值。

h_t 在不同节点下往往会有很大的区别。

2.2LSTM详解

为了解决梯度消失和爆炸以及更好的预测和分类序列数据等问题,rnn逐渐转变为lstm。



$$i^{(t)} = \sigma(W^{(i)}x^{(t)} + U^{(i)}h^{(t-1)})$$

(Input gate)

$$f^{(t)} = \sigma(W^{(f)}x^{(t)} + U^{(f)}h^{(t-1)})$$

(Forget gate)

$$o^{(t)} = \sigma(W^{(o)}x^{(t)} + U^{(o)}h^{(t-1)})$$

(Output/Exposure gate)

$$\tilde{c}^{(t)} = \tanh(W^{(c)}x^{(t)} + U^{(c)}h^{(t-1)})$$

(New memory cell)

$$c^{(t)} = f^{(t)} \circ c^{(t-1)} + i^{(t)} \circ \tilde{c}^{(t)}$$

(Final memory cell)

$$h^{(t)} = o^{(t)} \circ \tanh(c^{(t)})$$

<https://chehongshu.blog.csdn.net>

看的不是特别懂，下面就来逐一分析。

2.2.1遗忘门（第一个框）

这个阶段主要是对上一个节点传进来的输入进行选择性的忘记。简单来说就是会“忘记不重要的，记住重要的”。具体来说是通过计算得到的 $f(t)$ 表示forget来作为忘记门控，来控制上一个状态的 $C(t-1)$ 哪些需要留，哪些需要忘。（第二个公式）

2.2.2输入门（第二个框）

这个阶段将这个阶段的输入有选择性的进行“记忆”。主要是会对输入 $x(t)$ 进行选择记忆。哪些重要则着重记录下来，哪些不重要，则少记一些。而选择的门控信号则是由 i 代表(information)来进行控制。（第一个公式）

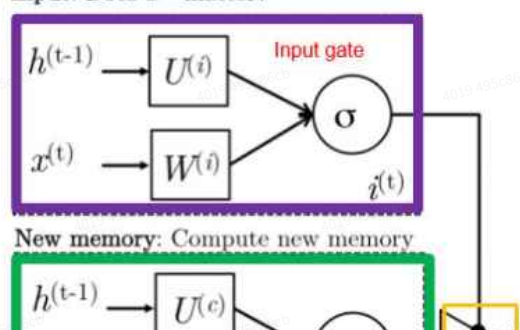
2.2.3输出门（第三个框）

这个阶段将决定哪些将会被当成当前状态的输出。

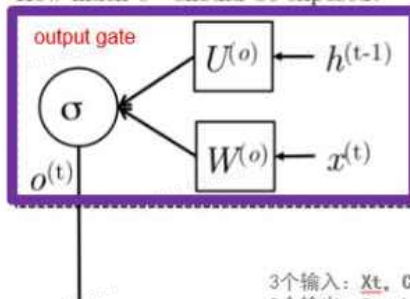
主要是通过 $o(t)$ 来进行控制的(第三个公式)。并且还对上一阶段得到 c 的进行了放缩（通过一个 \tanh 激活函数进行变化）。（第四个公式）

2.3LSTM各部分展开图

Input: Does $x^{(t)}$ matter?

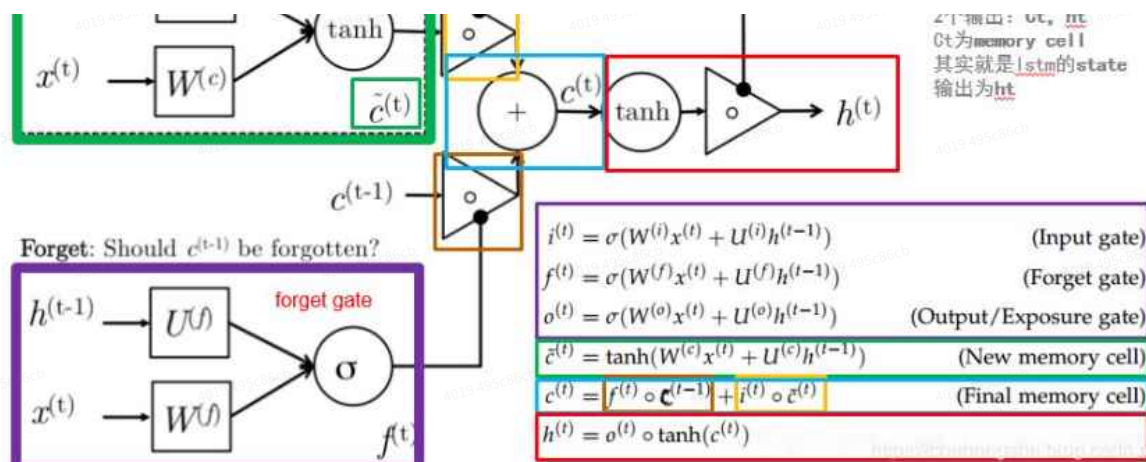


Output/Exposure:
How much $c^{(t)}$ should be exposed?



3个输入: x_t, c_{t-1}, h_{t-1}

2个输出: z_t, o_t



1. 首先输入为三个值，一个是此刻的输入值 x_t ，另一个是上一时刻的状态值 $c(t-1)$ ，最后一个是上一个单元的输出生 $h(t-1)$ 。
2. 最终输出为两个值，一个是此刻产生的状态值 c_t 和输出 h_t 。
3. 首先是输入值 x 和上一个单元的输出生 h ，分别两个输入都有对应的权重，在经过sigmoid激活作用下得到0-1的值，也就是三个门值。（得到紫色框中的三个门值）

$$i^{(t)} = \sigma(W^{(i)}x^{(t)} + U^{(i)}h^{(t-1)}) \quad (\text{Input gate})$$

$$f^{(t)} = \sigma(W^{(f)}x^{(t)} + U^{(f)}h^{(t-1)}) \quad (\text{Forget gate})$$

$$o^{(t)} = \sigma(W^{(o)}x^{(t)} + U^{(o)}h^{(t-1)}) \quad (\text{Output/Exposure gate})$$

4. 和3差不多，依然还是 输入值 x 和上一个单元的输出生 h ，两个值有对应的权重和3中的描述一模一样，唯一的区别在于有一个tanh激活函数，最后相当于得到此时输入得到的当前state，也就是new memory。

这里可以理解为输入其实是近似的 x 和 h 的concatenate操作，经过正常的神经网络的权重，最后经过tanh激活函数得到此时输入的当前的state， x 相当于此刻的输入， h 为前面历史的输入，合在一起就是整个序列的信息，也就是此时的new memory。

$$\tilde{c}^{(t)} = \tanh(W^{(c)}x^{(t)} + U^{(c)}h^{(t-1)}) \quad (\text{New memory cell})$$

5. 最后输出的state，也就是final memory的计算利用了input gate和forget gate，output gate只与输出有关。

final memory的计算自然而然和上一步算得此时的记忆state相关并且和上一个输出的final memory相关，故为忘记门和 C_{t-1} 的乘积加上上一步算出来的此时单元的C和输入门的乘积为最终的state。

忘记门和 C_{t-1} 的乘积：选择遗忘哪一些之前的信息

单元的C和输入门的乘积：选择保留当前状态的那些信息

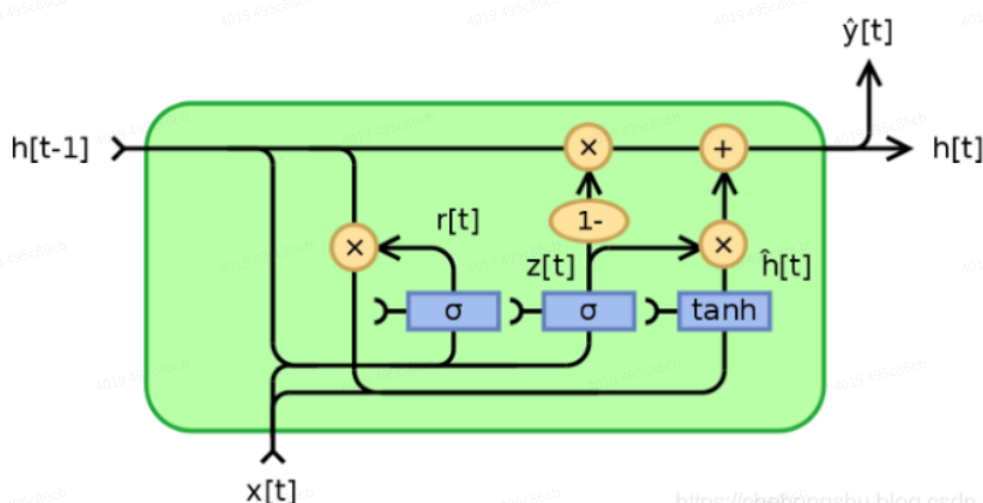
$$c^{(t)} = f^{(t)} \circ c^{(t-1)} + i^{(t)} \circ \tilde{c}^{(t)} \quad (\text{Final memory cell})$$

6. 输出门只与输出相关，最终的输出h为输出门乘以 $\tanh(c)$

$$h^{(t)} = o^{(t)} \circ \tanh(c^{(t)})$$

3.GRU

因为LSTM的训练比较慢，而GRU在其上稍微修改，速度可以快很多，而精度基本不变，所以GRU也十分流行



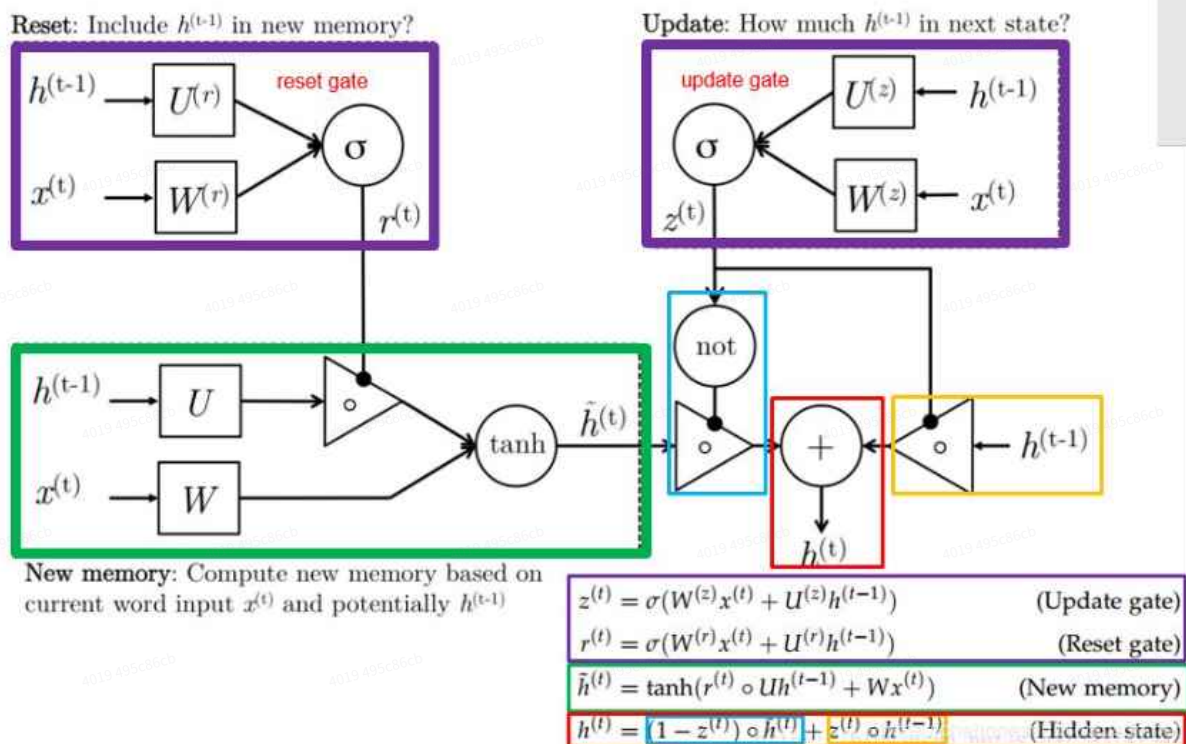
$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

换个图看看：



1. 这里GRU只有两个gate，一个是reset gate，一个是update gate，
update gate的作用类似于input gate和forget gate，
(1-z)相当于input gate，z相当于forget gate。
2. 输入为两个值，输出也为一个值，输入为输入此时时刻值x和上一个时刻的输出ht-1，输出这个时刻的输出值ht。
3. 首先依然是利用xt和ht-1经过权重相乘通过sigmoid，得到两个0-1的值，即两个门值。

$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

4. 接下来这里有一些不同，并且经常容易搞混淆。对于LSTM来说依然还是xt与ht-1分别权重相乘相加，之后经过tanh函数为此时的new memory。

而GRU为在这个计算过程中，在 h_{t-1} 与权重乘积之后和reset gate相乘，之后最终得到new memory，这里的reset gate的作用为让这个new memory包括之前的 h_{t-1} 的信息的多少。

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

5. 接下来和lstm得到final memory其实一样，只是GRU只有两个输入，一个输出，其实这里h即输出也是state，就是说GRU的输出和state是一个值，所以4步骤得到的是new_h，这步骤得到的是final_h，通过update gate得到。

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

4.GRU与LSTM之间的比较

4.1 结构上

1. lstm为三个输入 x_t ， h_{t-1} ， c_{t-1} ，两个输出。gru为两个输入 x_t ， h_{t-1} ，一个输出 h_t ，输出即state。
2. lstm有三个门，输入输出忘记门。gru有两个门，reset，update 门。
3. update 类似于 input gate和forget gate。

4.2 功能上

4. GRU参数更少，训练速度更快，相比之下需要的数据量更少。
5. 如果有足够的数据，LSTM的效果可能好于GRU。