

法律声明

□ 本课件包括：演示文稿，示例，代码，题库，视频和声音等，小象学院拥有完全知识产权的权利；只限于善意学习者在本课程使用，不得在课程范围外向任何第三方散播。任何其他人或机构不得盗版、复制、仿造其中的创意，我们将保留一切通过法律手段追究违反者的权利。

□ 课程详情请咨询

■ 微信公众号：小象

■ 新浪微博：ChinaHadoop



第二节课：RNN语言模型

□ RNN语言模型理论基础

□ RNN语言模型实践

□ 从RNN到LSTM

■ LSTM与长期记忆

■ 两句话解释BPTT

参考文献

□ RNN语言模型理论基础

- Recurrent neural network based language model. (Interspeech 2010, Tomas Mikolov et al.)
- [Unreasonable Effectiveness of Recurrent Neural Network](#)

□ RNN语言模型实践

- <https://github.com/sherjilozair/char-rnn-tensorflow>, 教科书级别的代码

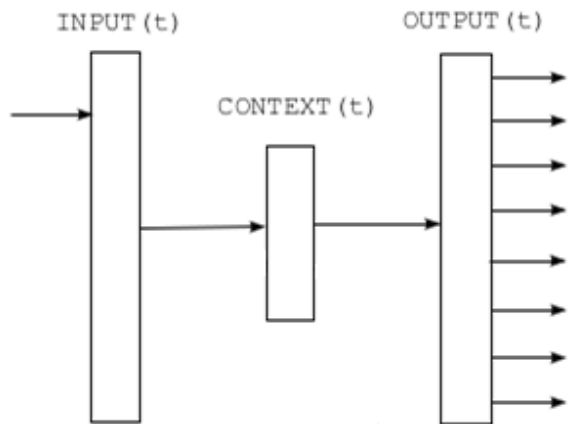
□ 从RNN到LSTM

- Understanding LSTM Networks –colah’s blog
- Deep Learning, Ian Goodfellow, Yoshua Bengio, Aaron Courville

RNN语言模型

RNNLM

CBOW/skip-gram的局限



□ CBOW/skip-gram 使用固定长度的语境.

- 语境的长度需要人工调节
- 只能利用固定长度的语境中的信息，不能利用更长范围内的词语之间的关联.

我们引入记忆

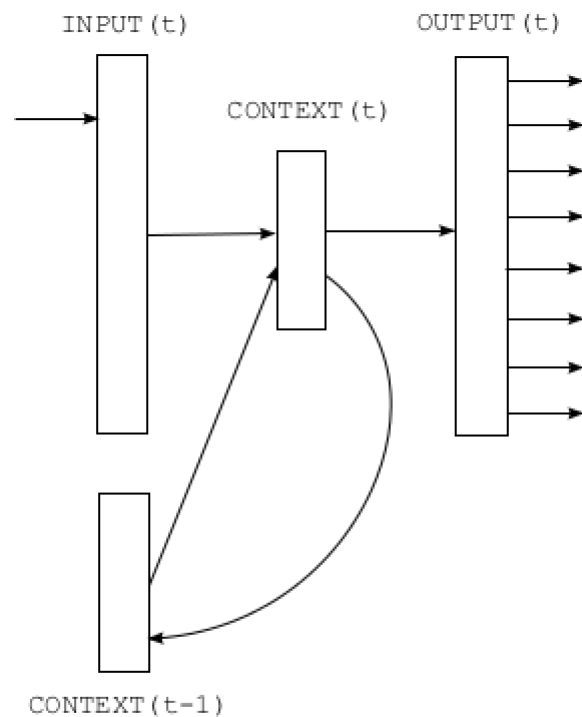
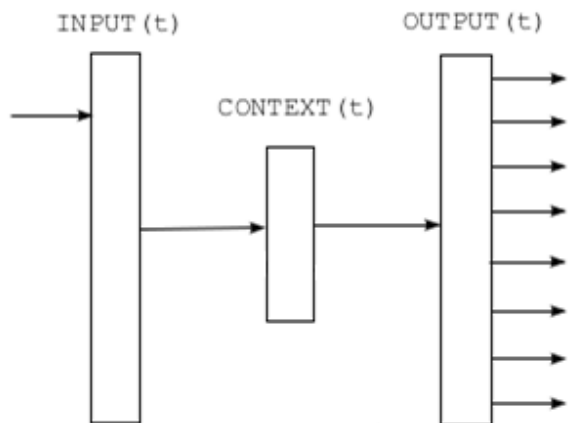


Figure 1: *Simple recurrent neural network.*

RNN语言模型细节

一个RNN cell
是一个经典的神经网络

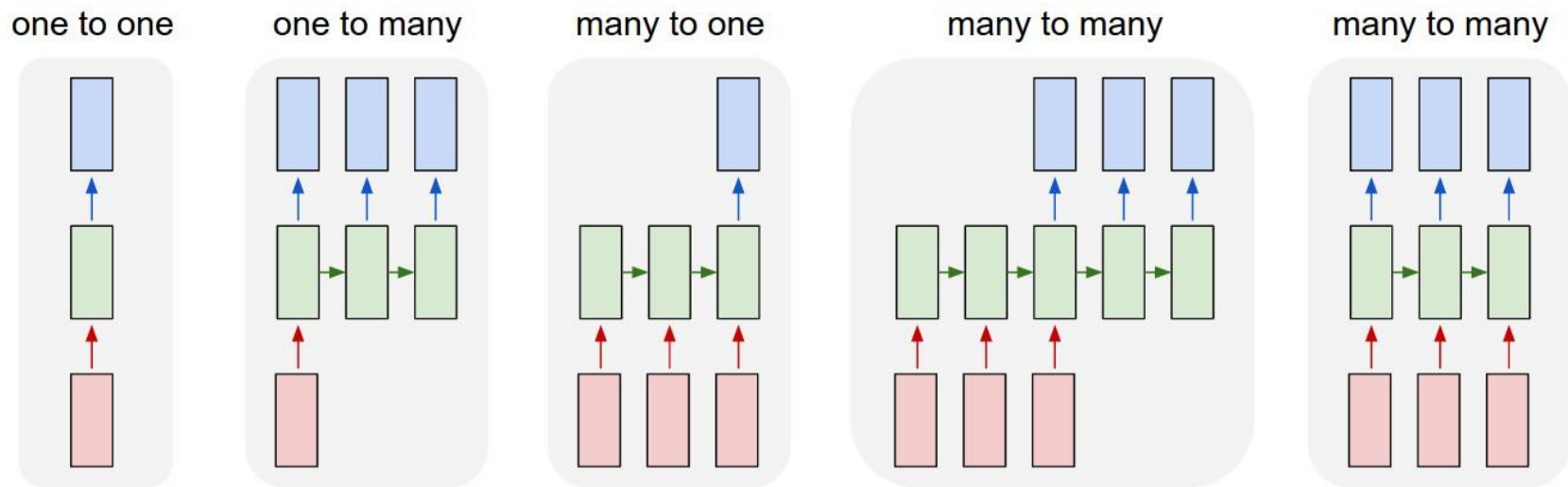
RNN语言模型细节

RNN语言模型细节

□ $P(W_{t+1} | W_t, W_{t-1}, \dots W_0)$

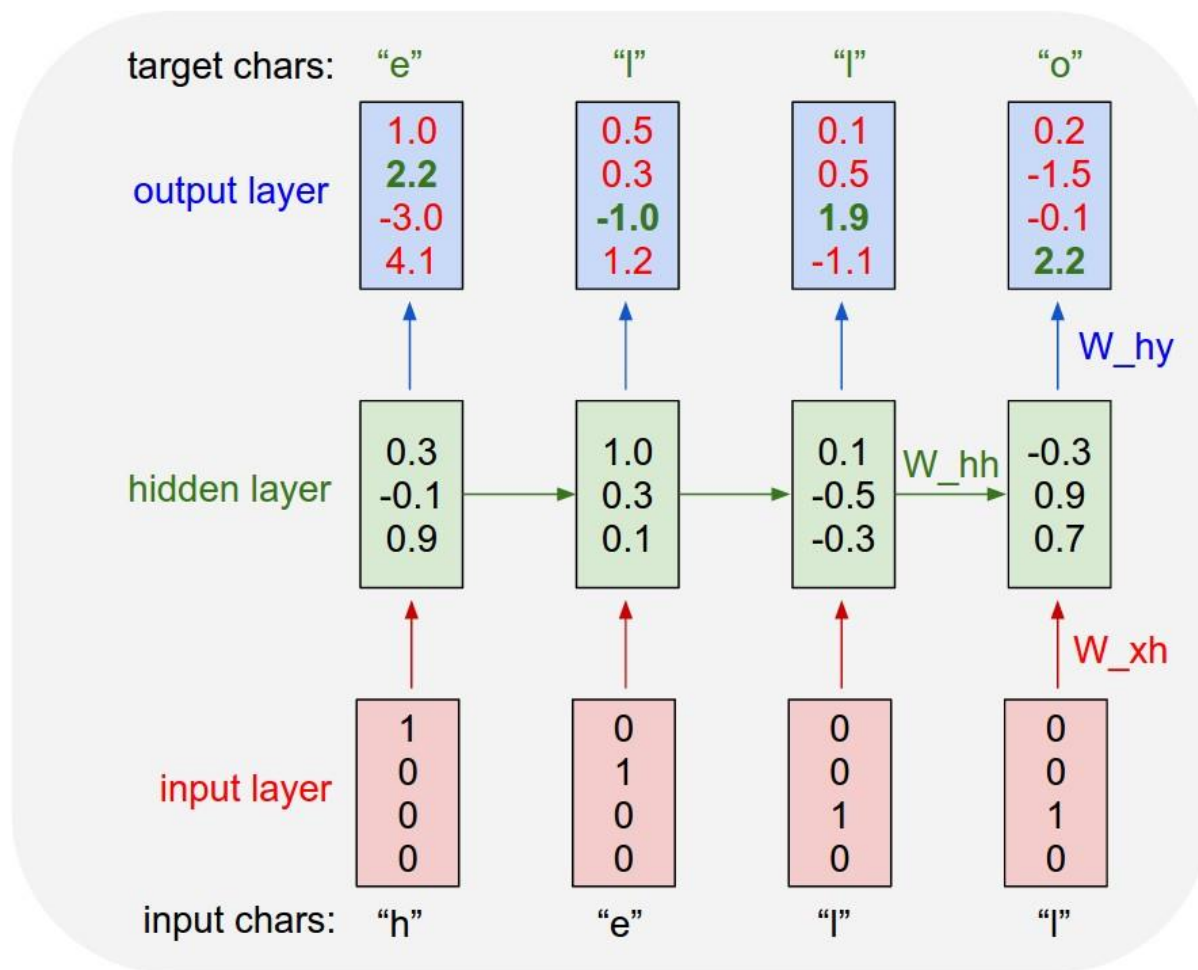
- 给定当前 (W_t) 以及更早 ($W_{t-1}, \dots W_0$) 的词语，预测句子中的下一个词语
- 在当前时间，有**两个**信息来源
 1. 当前词语的embedding
 2. 一个总结了之前所有单词的状态向量 S_{t-1}

RNN语言模型细节



来自大牛[Karpathy](#)的blog

RNN语言模型细节



RNN语言模型

代码实践

RNN语言模型实践

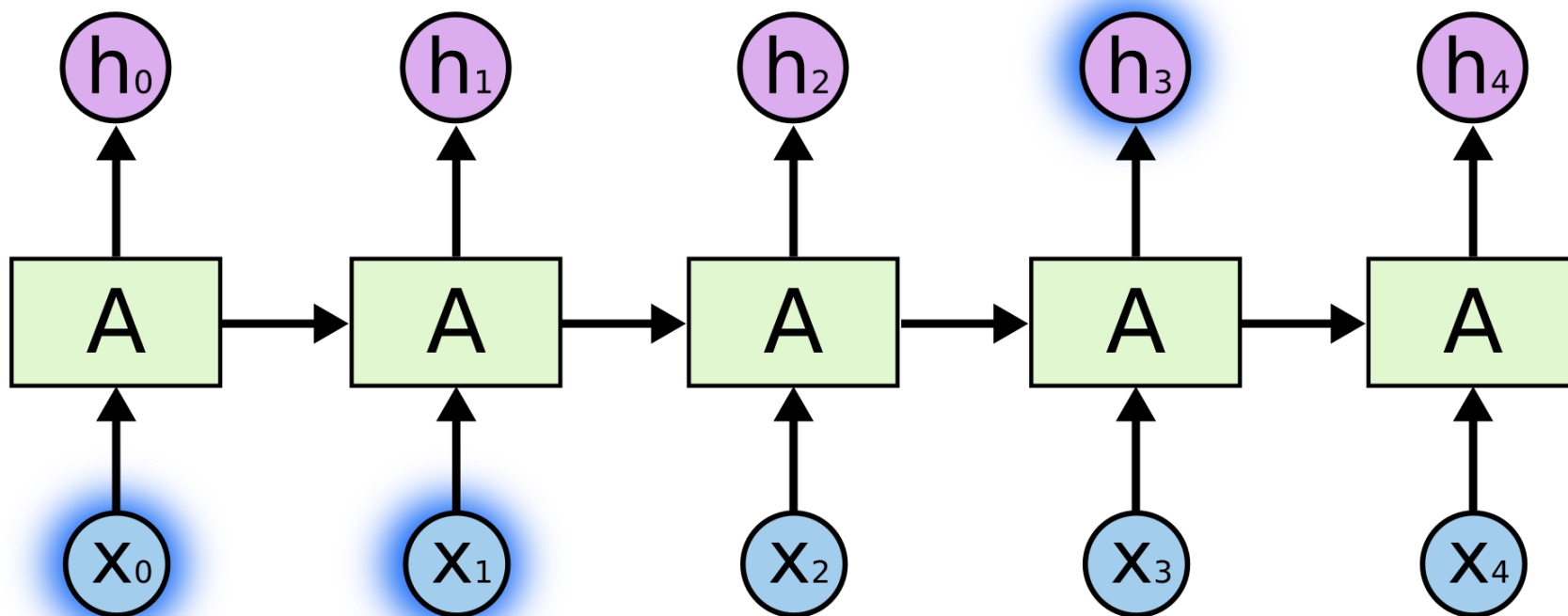
□ RNN 序列模型 第二版

- 使用 static_rnn
- 实践 cross-validation 交叉验证

RNN语言模型

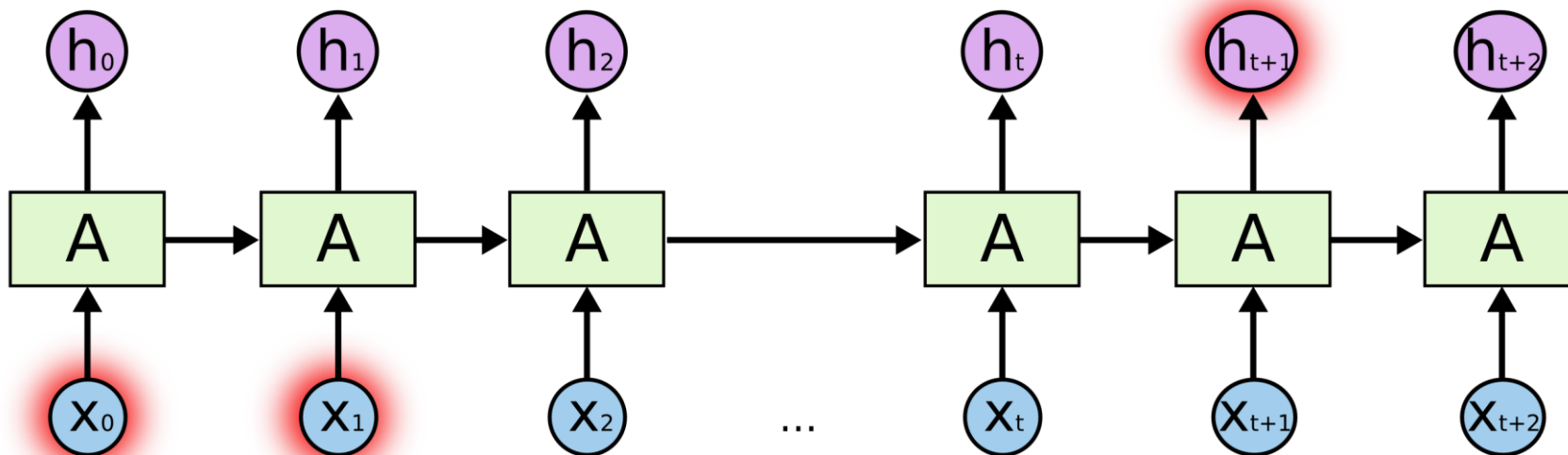
LSTM

相互关联的一些词距离很近时：



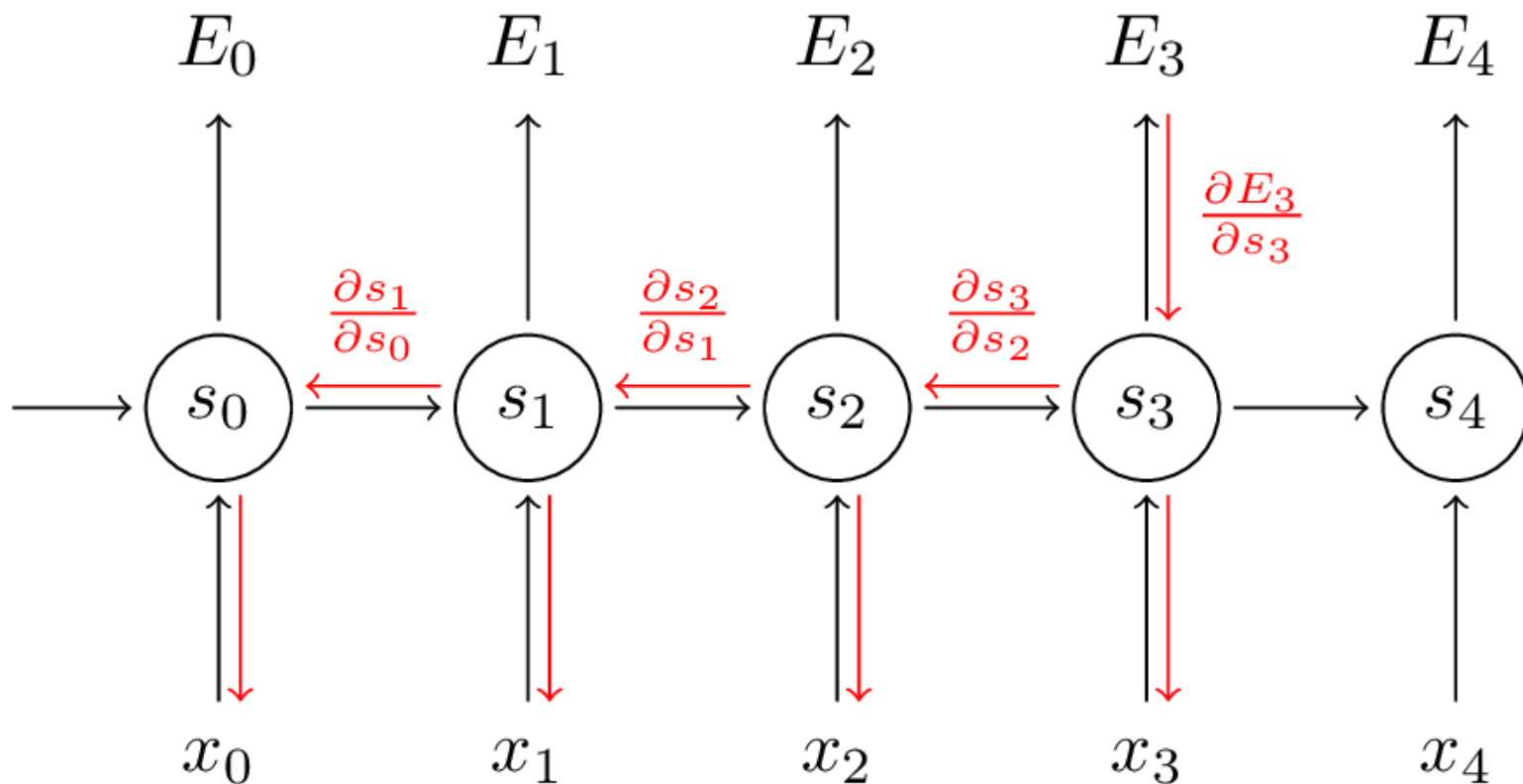
"the **clouds** are in the **sky**,"

相互关联的一些词距离很远时：



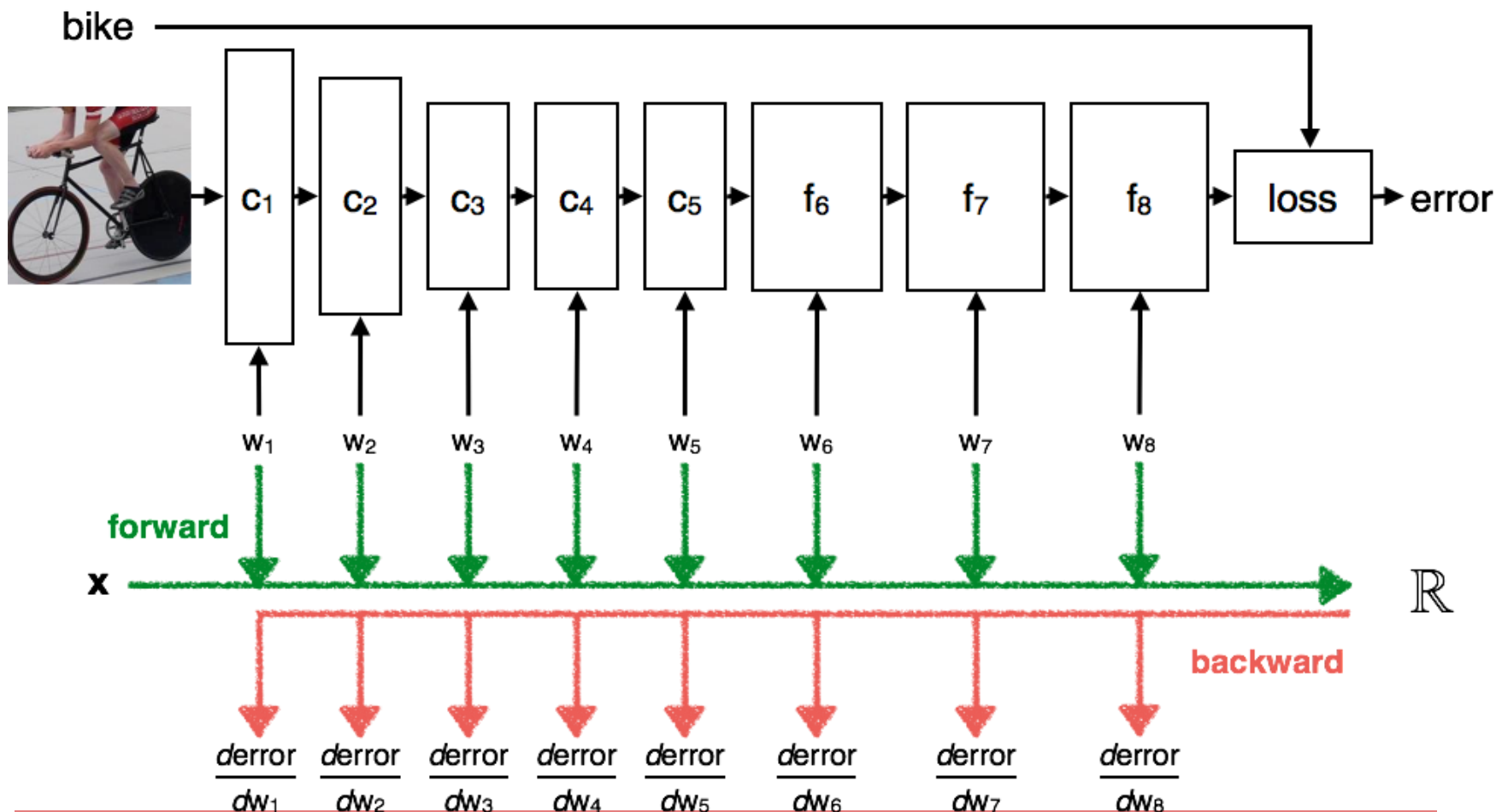
"I grew up in **France**...I ate a log of cheese.... I speak fluent ***French***."

BPTT Gradient: 爆炸或者消失

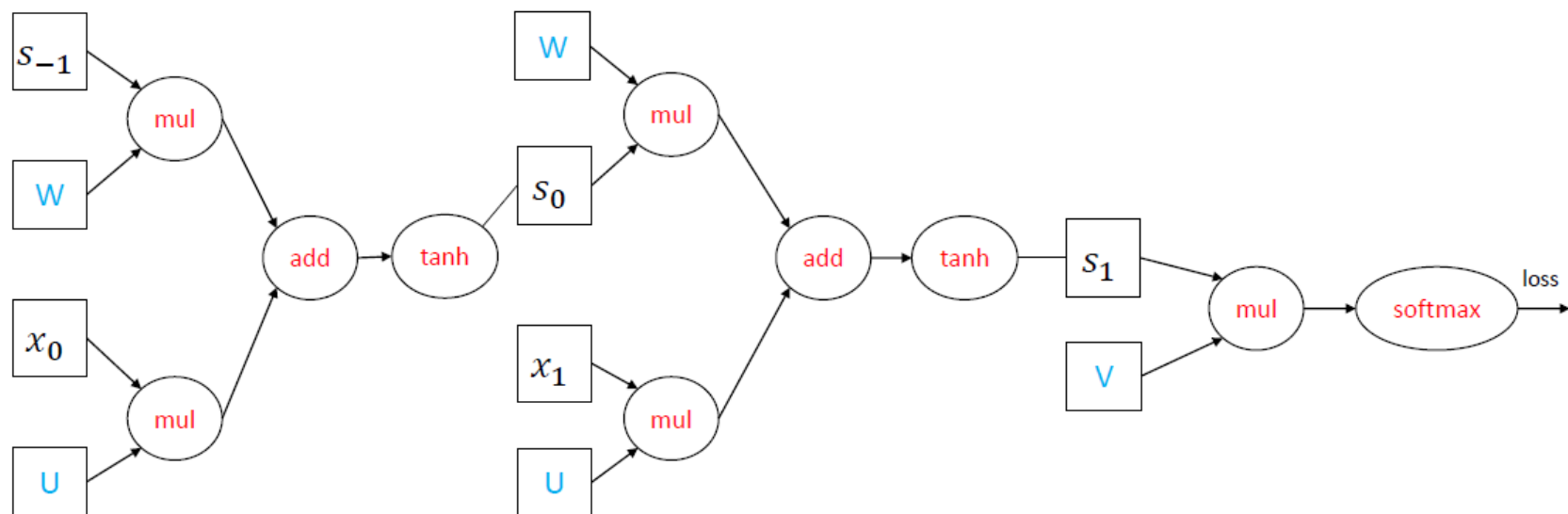


BPTT Gradient: 爆炸或者消失

问题来了，**CNN**的**depth**和**RNN**的**depth**有什么区别吗？



BPTT Gradient: 爆炸或者消失



$$\begin{aligned} s_0 W + \dots &\rightarrow s_1 \\ s_1 W + \dots &\rightarrow s_2 \\ s_2 W + \dots &\rightarrow s_3 \\ &\dots \\ s_{t-1} W + \dots &\rightarrow s_t \end{aligned}$$

在每一个时间点，同一个参数 W 和状态向量相乘
在长度为 $T+1$ 的序列里面， W 和状态向量相乘 T 次

BPTT Gradient: 爆炸或者消失

and lacking inputs \mathbf{x} . As described in Sec. 8.2.5, this recurrence relation essentially describes the power method. It may be simplified to

$$\mathbf{h}^{(t)} = (\mathbf{W}^t)^\top \mathbf{h}^{(0)}, \quad (10.37)$$

and if \mathbf{W} admits an eigendecomposition of the form

$$\mathbf{W} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top, \quad (10.38)$$

with orthogonal \mathbf{Q} , the recurrence may be simplified further to

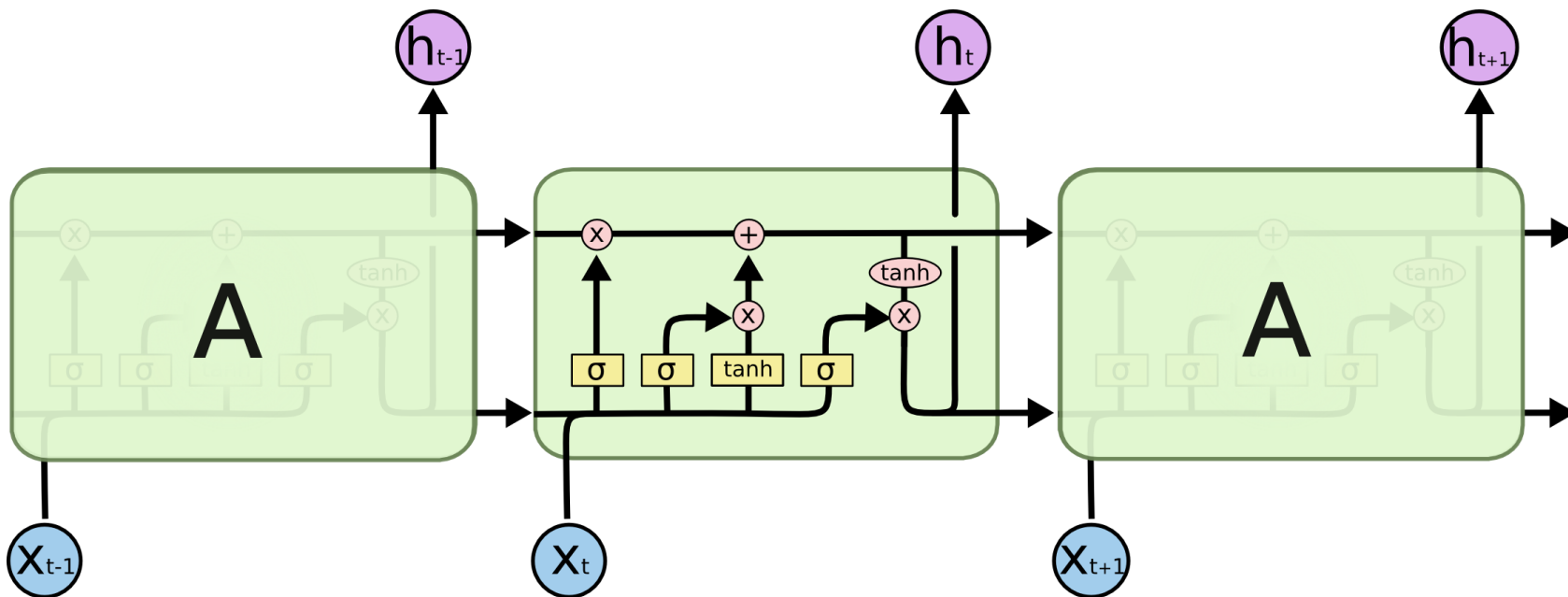
$$\mathbf{h}^{(t)} = \mathbf{Q}^\top \mathbf{\Lambda}^t \mathbf{Q} \mathbf{h}^{(0)}. \quad (10.39)$$

来自 Deep Learning, Ian Goodfellow

BPTT Gradient: 爆炸或者消失

- 在较长的时间跨度下，不稳定的gradient是学习长距离关联的主要障碍
- LSTMs 被设计用来描述长距离关联（long-term dependency）.
 - Long
 - Short Term
 - Memory networks

LSTM



LSTM是一种特殊的RNN CELL

LSTM: Output Gate

中国 男足 踢得 真烂

气得我 想 大吃 一顿

假定我们用**256**维的向量表示RNN的状态
(256在RNN应用中算不上特别高维)

LSTM : Output Gate

中国 男足 踢得 真烂

气得我 想 大吃 一顿

我们需要使用所有的**256**个维度去表示
“中国男足踢得真烂”这个简单的意思吗？

LSTM : Output Gate

□ 中国男足踢得真烂。气得我想大吃一顿。最近有没有什么新开的好吃的餐馆？哎呀我的手机好慢，你帮我去xx app 上查一下？

1. 256个维度都用来既记录第一句，又记录第二句，还记录第三句
2. 256个维度有100个记录第一句，78个记录第二句，78个记录第三句 (better 😊)

LSTM : Output Gate

我们要求只有一部分维度用来
记忆/总结目前为止的输入

中国 男足 踢得 真烂

气得我 想 大吃 一顿

LSTM : Output Gate

中国 男足 踢得 真烂

气得我 想 大吃 一顿

LSTM : Output Gate

中国 男足 踢得 真烂

气得我 想 大吃 一顿

LSTM : Input Forget Gate

中国 男足 踢得 真烂

气得我 想 大吃 一顿

让我们切换到下一个词语（气得我...）的视角

LSTM : Input Forget Gate

中国 男足 踢得 真烂

气得我 想 大吃 一顿

让我们切换到下一个词语（气得我...）的视角

LSTM : Input Forget Gate

中国 男足 踢得 真烂

气得我 想 大吃 一顿

- 256个维度有100个记录“中国 男足 踢得 真烂”这句话
- 进一步，我们假设其中25个维度侧重表示“真烂”这个词

LSTM : Input Forget Gate

中国 男足 踢得 真烂

气得我想 大吃 一顿

- 进一步，我们假设其中25个维度侧重表示“真烂”这个词
- 考虑到当前的词语是“气得我”，和语境里面的“真烂”这个词关系最密切

LSTM : Input Forget Gate

中国 男足 踢得 真烂

气得我想 大吃 一顿

- 考虑到当前的词语是“**气得我**”，和语境里面的“**真烂**”这个词关系最密切
- 我们把历史状态里面和当前最密切的维度提取出来

LSTM : Input Forget Gate

根据当下的单词，我们检查哪些维度包含的信息和当下单词密切相关？

LSTM : Input Forget Gate

但是，我看到的是 o_t 处理过的二手信息——一手信息啥样的？

LSTM : Input Forget Gate

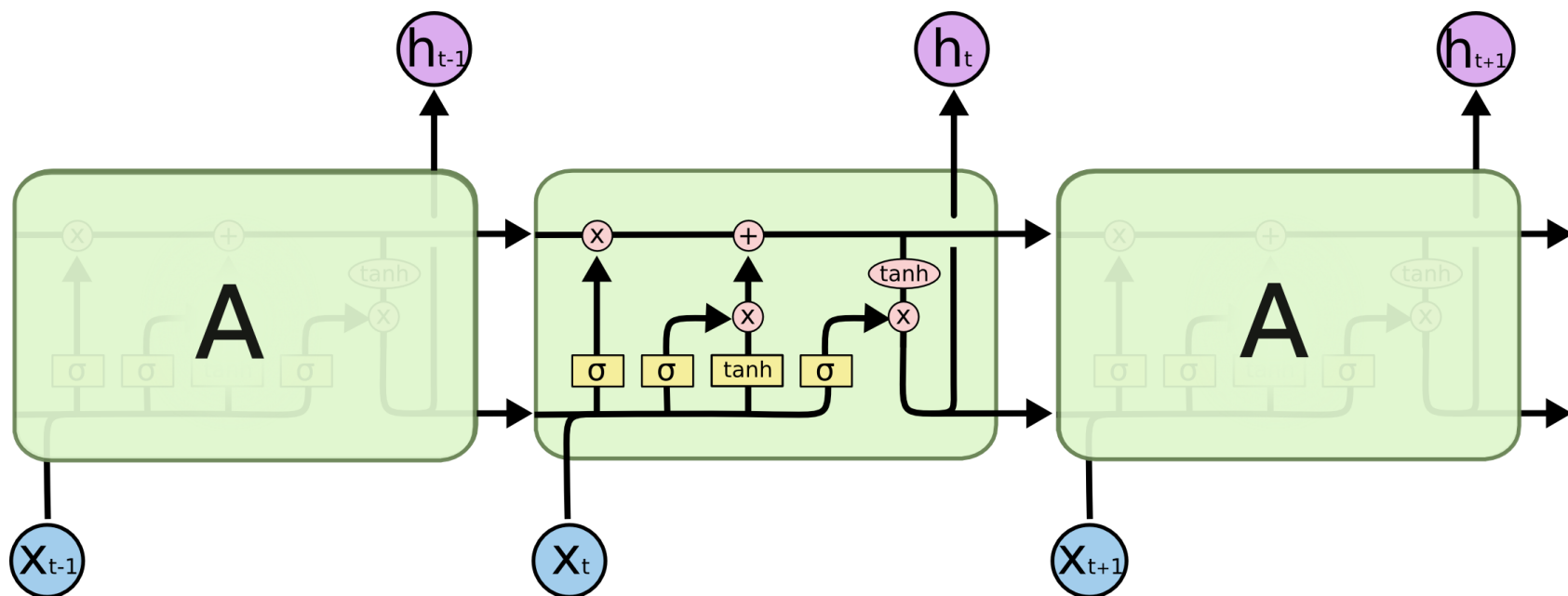
根据当前词语，
挑选的一手信息

LSTM : Input Forget Gate

两类信息如何结合起来？

LSTM : Input Forget Gate

LSTM 小结



LSTM 有两套state: C & h

两句话解释BPTT

□ Back Propagation Through Time

- 无论多么复杂的深度学习模型都是DAG
- 对于任何一个node X , 想要计算 $\frac{\partial L}{\partial X}$, 只需要计算所有children的 $\frac{\partial L}{\partial y}$, $y \in Children(X)$: $\frac{\partial L}{\partial X} = \sum_y W_y \frac{\partial L}{\partial y}$

上述理解是否合理？

□ 上述理解：

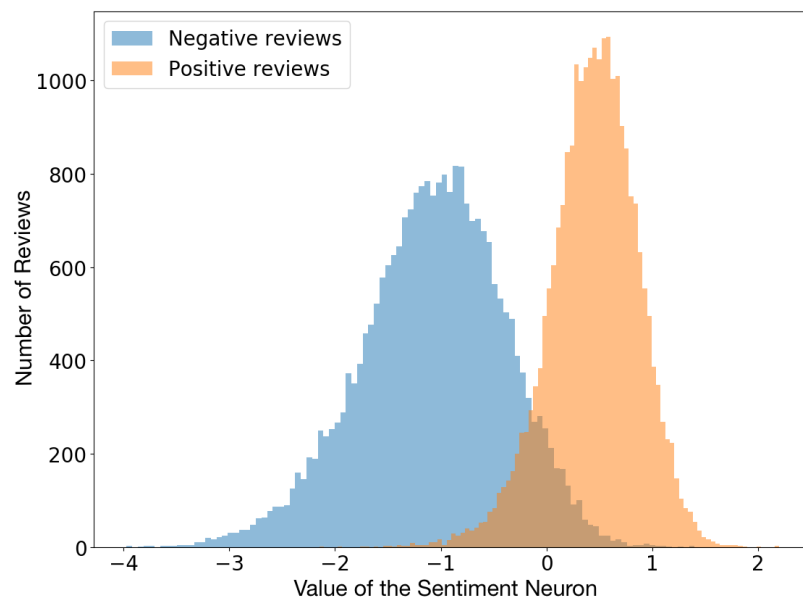
- 256个维度有100个记录第一句，78个记录第二句，78个记录第三句
- 256个维度有100个记录“中国男足踢得真烂”这句话；进一步，我们假设其中25个维度侧重表示“真烂”这个词
- mLSTM有4,096个维度
- 发现4096个维度里面有一个neuron node，可以相当准确地预测评论的褒贬

一个支持上述理解的工作

□ OpenAI: Unsupervised Sentiment Neuron

- 在82million Amazon 评论数据集上训练RNNLM
- 使用一个LSTM的变种 (multiplicative LSTM)

1. mLSTM有4,096 个维度
2. 这4096个维度里面有一个neuron node, 可以相当准确地预测评论的褒贬



LSTM代码演示

疑问

□ 问题答疑：<http://www.xxwenda.com/>

■ 可邀请老师或者其他回答问题

联系我们

小象学院：互联网新技术在线教育领航者

- 微信公众号：大数据分析挖掘
- 新浪微博：ChinaHadoop

