

端对端语音识别分享

分享人：黄佳恒
2019.12.16

思考

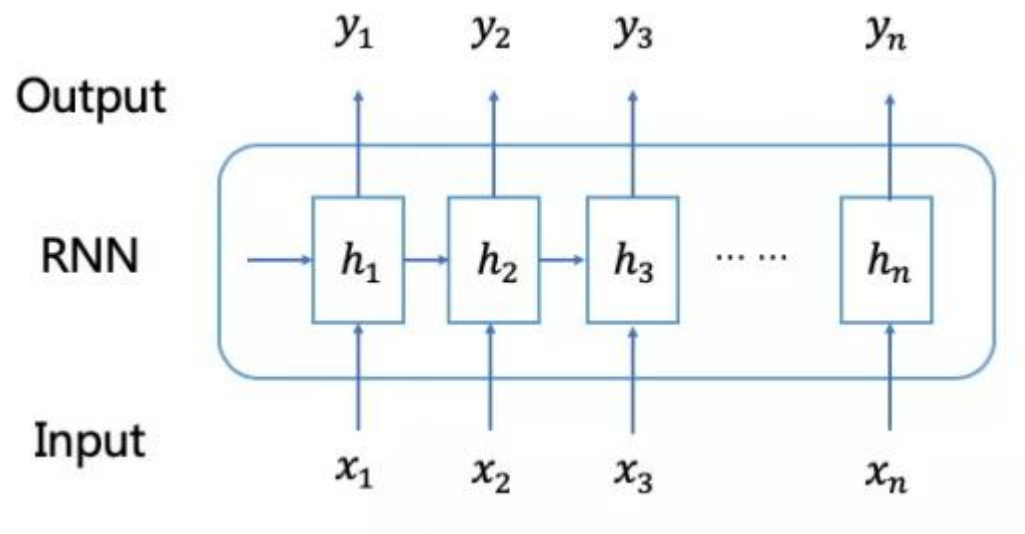
什么模型可以做语音识别？

RNN能做什么

1. 文本分类
2. 词性标注

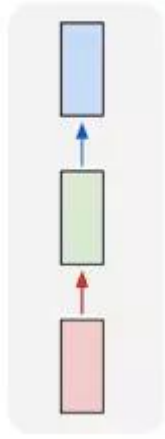
RNN要求：

输入序列和输出序列质检映射关系提前标注好

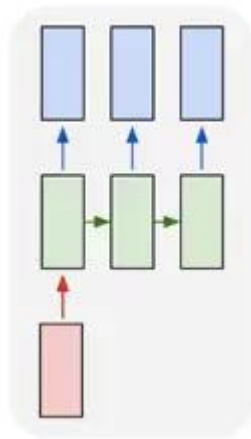


rnn

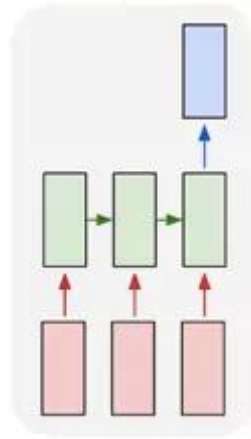
one to one



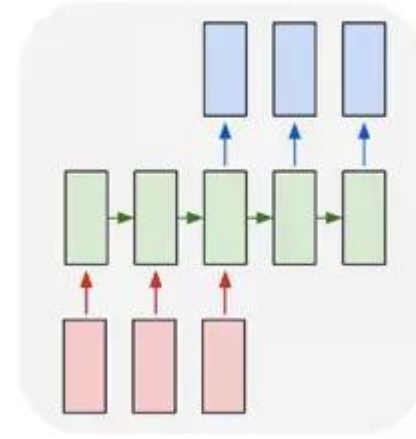
one to many



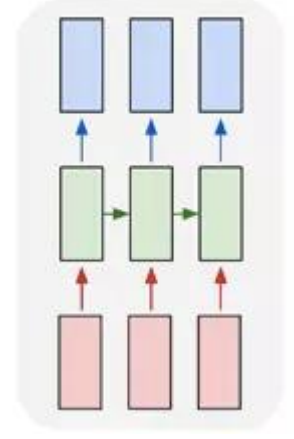
many to one



many to many



many to many



RNN不能做什么

天然难分割：音频、图像

为什么呢？

- 1) 音频的最小单位应该是什么？
- 2) 标注的量级大吗？

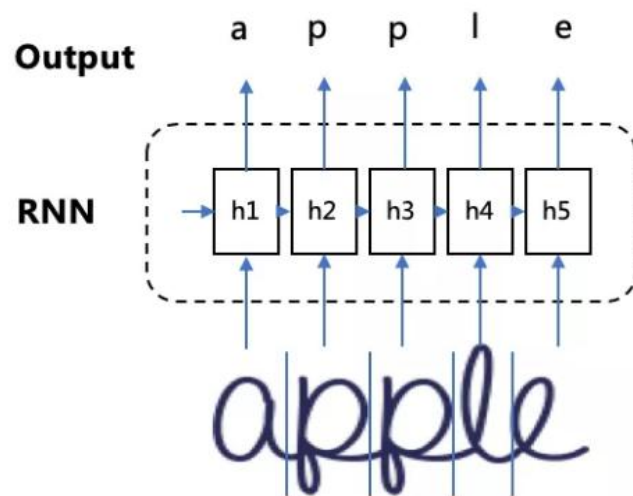
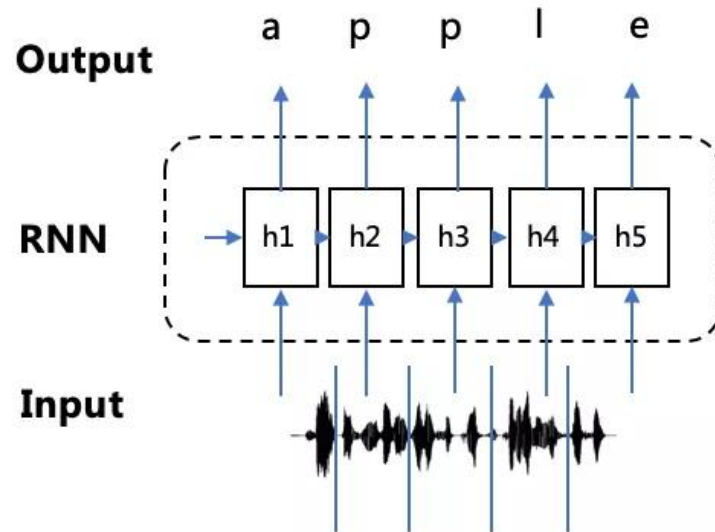
说话人语速有快有慢

xua

一个音素持续时间50-100ms之间

一般选取20、30、40、50ms

分割并标注映射关系的数据依赖是不切实际的



CTC

Connectionist Temporal Classification(连接时序分类器)

扩展了RNN的输出层，在输出序列和最终标签之间增加了多对一的空间映射，并在此基础上定义了CTC Loss函数

识别过程：<https://src.ailemon.me/blog/2019/20190718-1.gif>

t h e q u i c k b r o w n f o x



The quick brown fox

手写识别: 输入可以是 (x, y) 笔划的坐标或图像中的像素

j u m p s o v e r t h e l a z y d o g



语音识别: 输入可以是语谱图或其他基于频率的特征提取器

对齐

- 这种方法有两个问题。
- 通常，强制每个输入元素与某些输出对齐是没有意义的。例如，在语音识别中，输入元素可以是没有相应输出的静音段。
- 我们无法生成连续的多个字符的输出。比如对齐的 $[h, h, e, l, l, l, o]$ ，合并重复将产生“helo”而不是“hello”。

x_1	x_2	x_3	x_4	x_5	x_6	input (X)
c	c	a	a	a	t	alignment
c		a			t	output (Y)

引入空标记

.为了解决这些问题，CTC为允许的输出集引入了一个新的标记，这个标记有时被称为空标记。我们在这里将它称为 ϵ 。 ϵ 标记与任何内容都不对应，只会从输出中删除。

.CTC允许的对齐长度与输入的长度相同。我们允许在合并重复和删除 ϵ 标记后映射到Y的任何对齐：

.如果Y在一行序列中有两个相同的字符，那么有效的对齐必须在它们之间有一个 ϵ 。有了这个规则，我们就可以区分出“hello”的合并对齐和“helo”的合并对齐。

.让我们回到输出长度为6的输出[c, a, t]。以下是有效和无效对齐的更多示例。

h h e ϵ ϵ l l l ϵ l l o

h e ϵ l ϵ l o

h e l l o

h e l l o

First, merge repeat characters.

Then, remove any ϵ tokens.

The remaining characters are the output.

Valid Alignments

ϵ c c ϵ a t

c c a a t t

c a ϵ ϵ ϵ t

Invalid Alignments

c ϵ c ϵ a t

c c a a t

c ϵ ϵ ϵ | t t

corresponds to $Y = [c, c, a, t]$

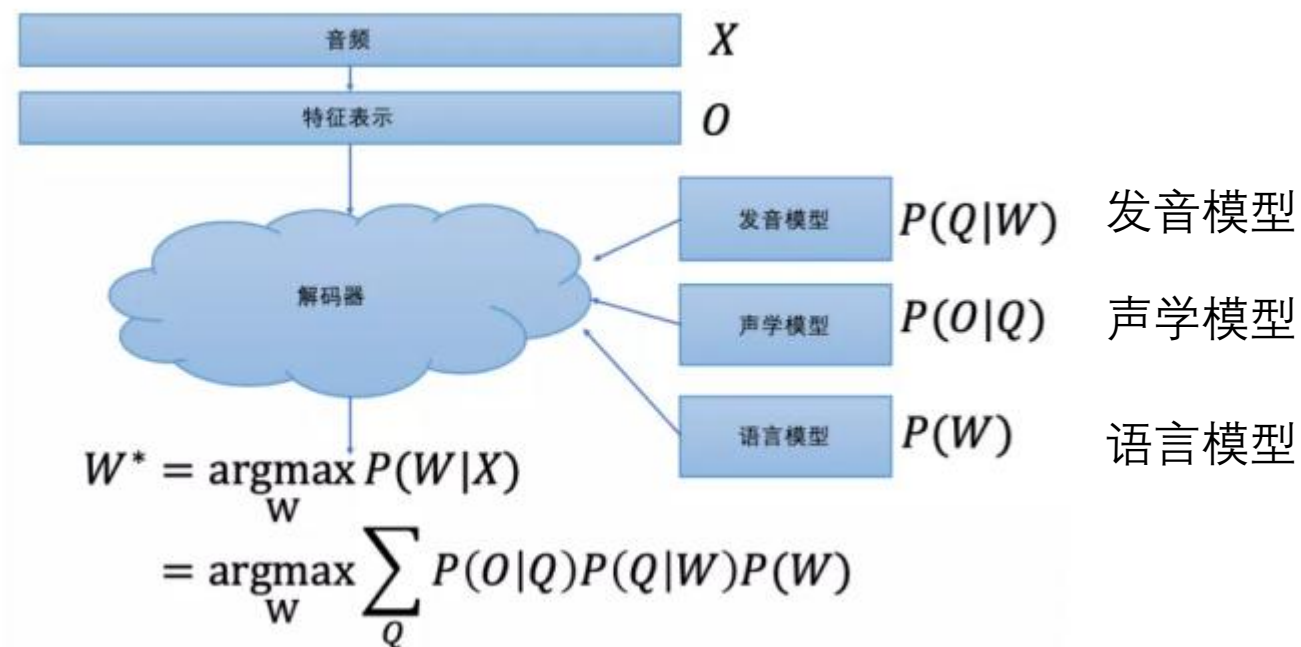
has length 5

missing the 'a'

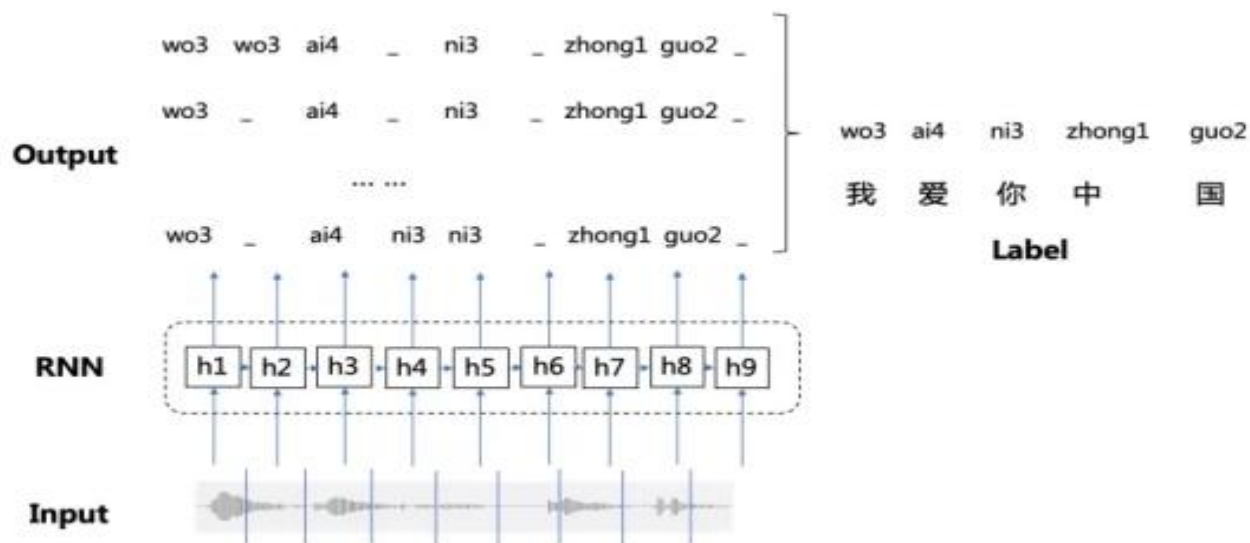
CTC Loss

1. X 是音频信号， O 是它的特征表示
一般使用MFCC、Fbank提取特征
2. Q 是 O 对应的发音序列，建模单元，
一般是音素
3. W 是音频信号 X 对应的文字

贝叶斯分解
引入隐变量（更符合常理）
求积分



CTC LOSS



- 训练一个时序分类器: $h(x) = \operatorname{argmax}_{z \in L^T} p(z|x)$
- Loss函数 $\mathcal{L}(S) = -\prod_{(x,z) \in S} p(x,z) =: -\sum_{(x,z) \in S} \ln p(x,z)$
 - $S \Rightarrow$ training data set from a fixed distribution $\mathcal{D}_{X \times Z}$
 - Input Space : $X = (\mathbb{R}^m)^+ \Rightarrow$ m 维实数向量序列集合
 - Target Space : $Z = L^+ \Rightarrow$ 由有限字符集 L 组成的序列集合
 - x 序列的长度大于或等于 z 序列的长度

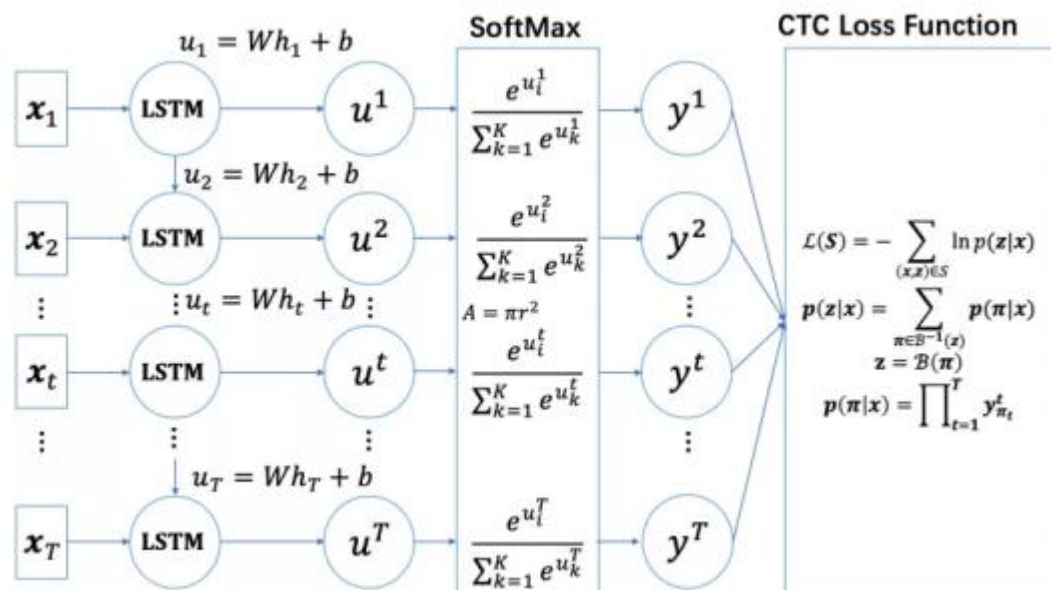
CTC LOSS

- 假设给定输入序列和模型参数，RNN每一时刻的输出之间是条件独立的，则：

- $p(\pi|x) = \prod_{t=1}^T y_{\pi_t}^t, \forall \pi \in L'^T$
- $p(z|x) = \sum_{\pi \in B^{-1}(z)} p(\pi|x)$, B^{-1} 是 z 全部路径集合的映射函数

- 我们得到CTC的Loss函数，如下：

- $\mathcal{L}(S) = -\ln \prod_{(x,z) \in S} p(z|x)$
- $= -\sum_{(x,z) \in S} \ln p(z|x)$
- $= -\sum_{(x,z) \in S} \ln \sum_{\pi \in B^{-1}(z)} p(\pi|x)$
- $= -\sum_{(x,z) \in S} \ln \sum_{\pi \in B^{-1}(z)} \prod_{t=1}^T y_{\pi_t}^t, \forall \pi \in L'^T$



动态规划

前向、后向到底解决了什么问题？

1. 计算时间复杂度优化
2. 方便求导的计算。因为梯度已经分解成了 α 和 β 的函数
复杂的损失函数分解成了 α 和 β 和 $y(t)$
3. 动态规划不是贪婪算法，没有损失最优解

LOSS损失函数计算

暴力计算 $P(z|x)$ 的复杂度非常高
使用Forward-Backward算法，利用动态规划求解

前向后向为甚什么会减少时间复杂度？

举个例子，斐波那契数列 0,1,1,2,3,5,8,13,... 有着一个相当简单的描述方式，它的每个数字都与前两个紧邻的数字相关。如果 $F(n)$ 是第 n 个数字，那么我们会得到 $F(n) = F(n-1) + F(n-2)$ 。这个在数学上称作*递归方程*或者*递推关系*。为了计算后面的项，它需要前面项的计算结果作为输入。

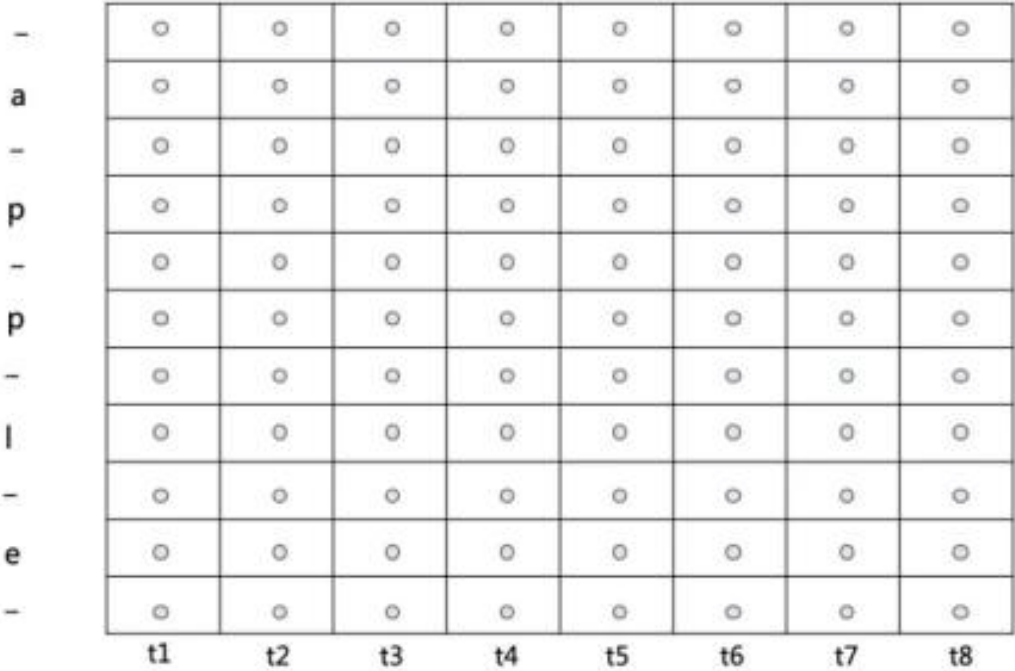
如果计算前 N 项斐波那契数列的和：

可以使用递归算法和迭代算法

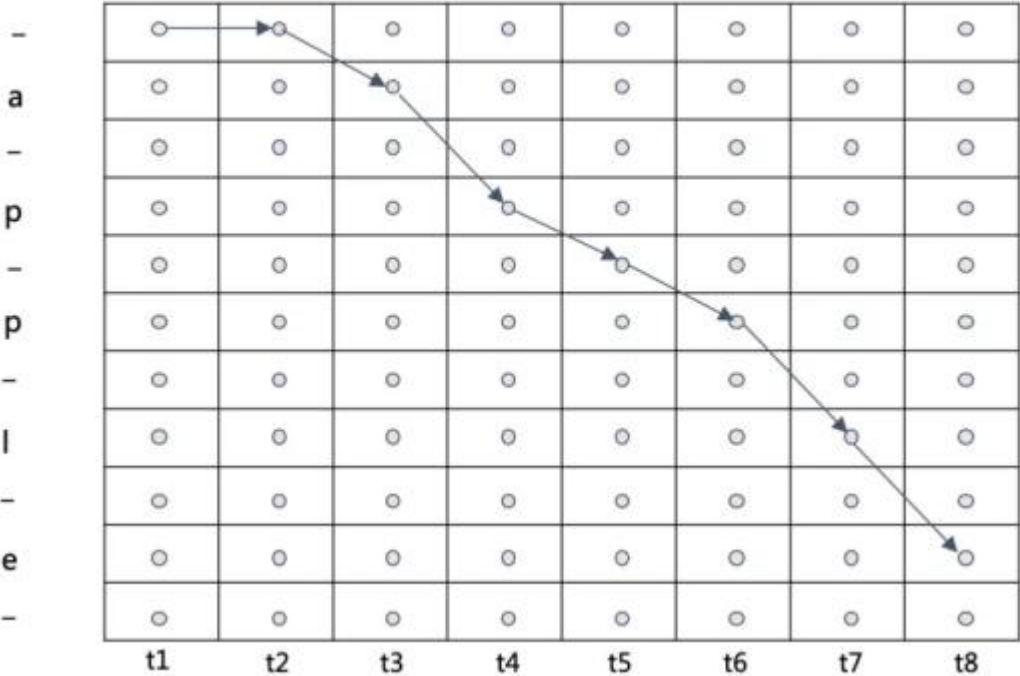
递归算法之所以速度慢，是因为它一遍又一遍地计算了相同的斐波那契数列

合法路径约束

以apple这个标签为例，按照时间序列 (T=8) 展开，它的路径搜索空间如下：

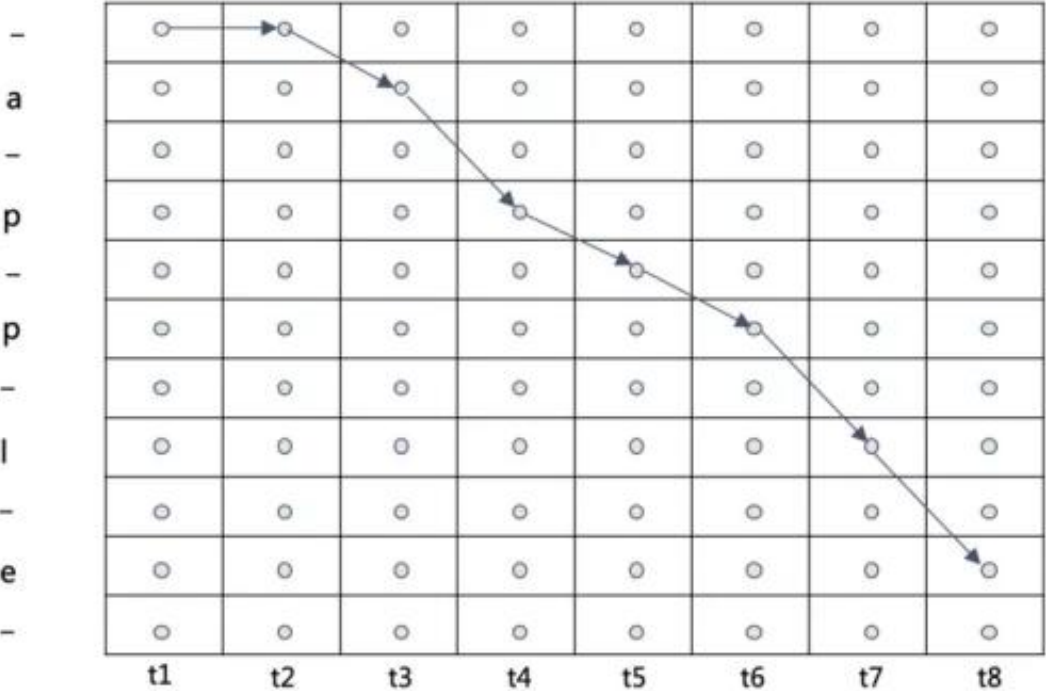


比如， $B(" _ _ a p _ p l e ") = "apple"$ ，对应的path如下：

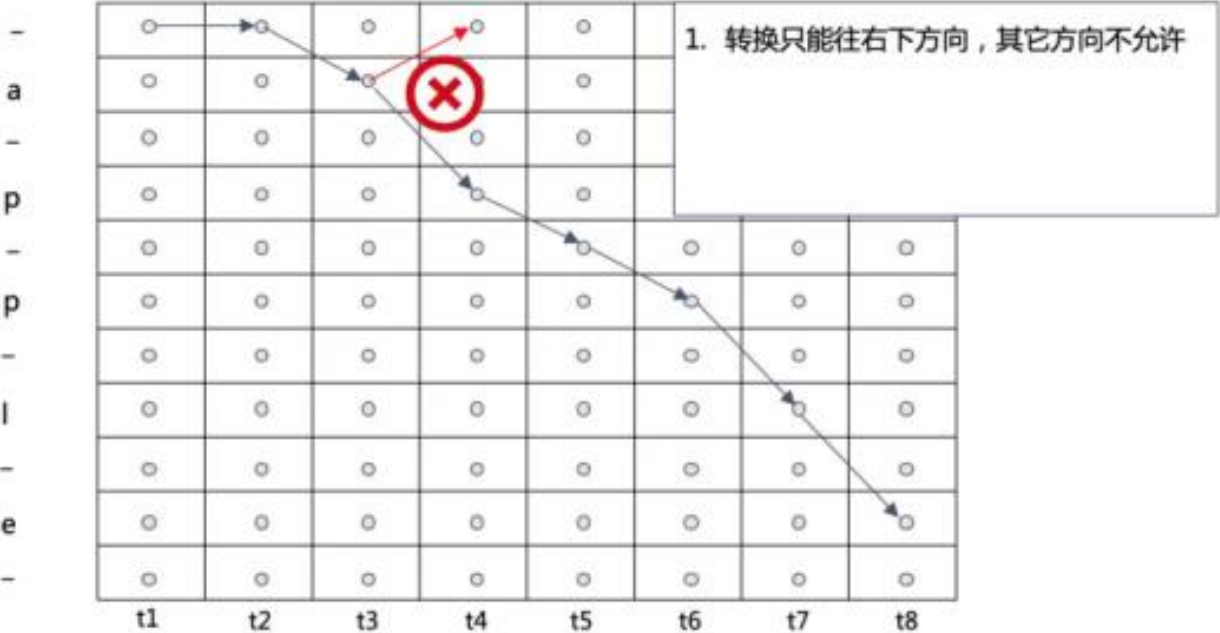


合法路径约束

$P(" _ _ a p _ p l e ") = y^1_{_} y^2_{_} y^3_a y^4_p y^5_{_} y^6_p y^7_l y^8_e$

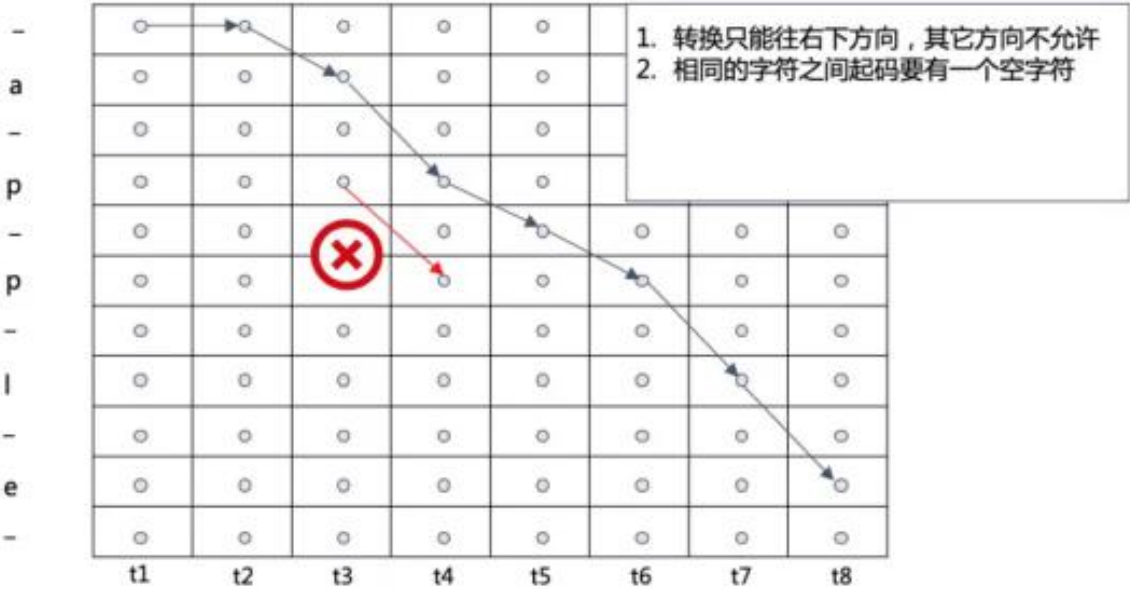


为了让所有的path都能在图中唯一、合法的表示，结点转换有如下约束：

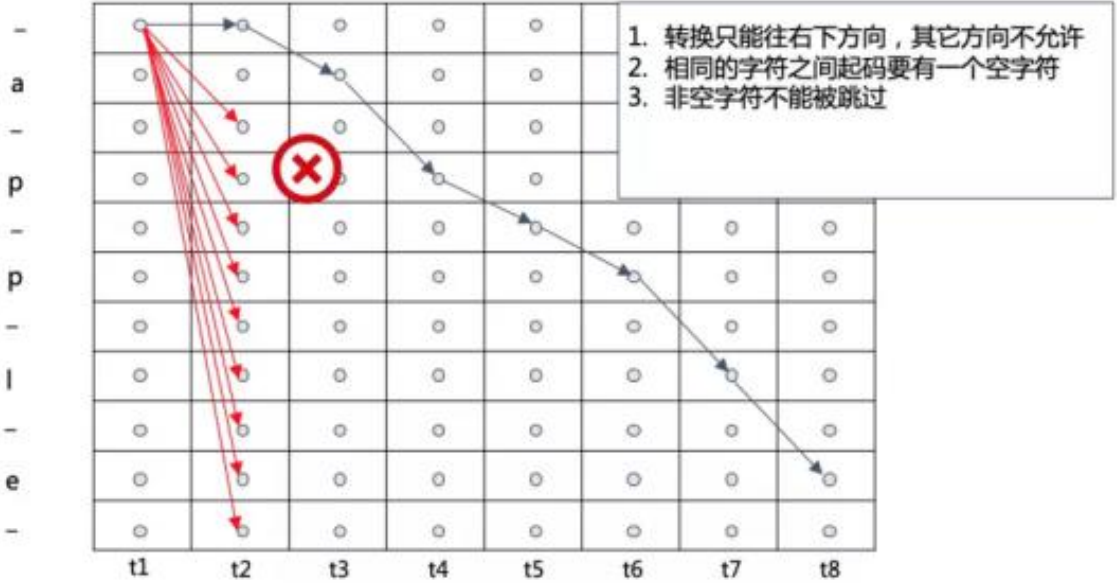


合法路径约束

为了让所有的path都能在图中唯一、合法的表示，结点转换有如下约束：

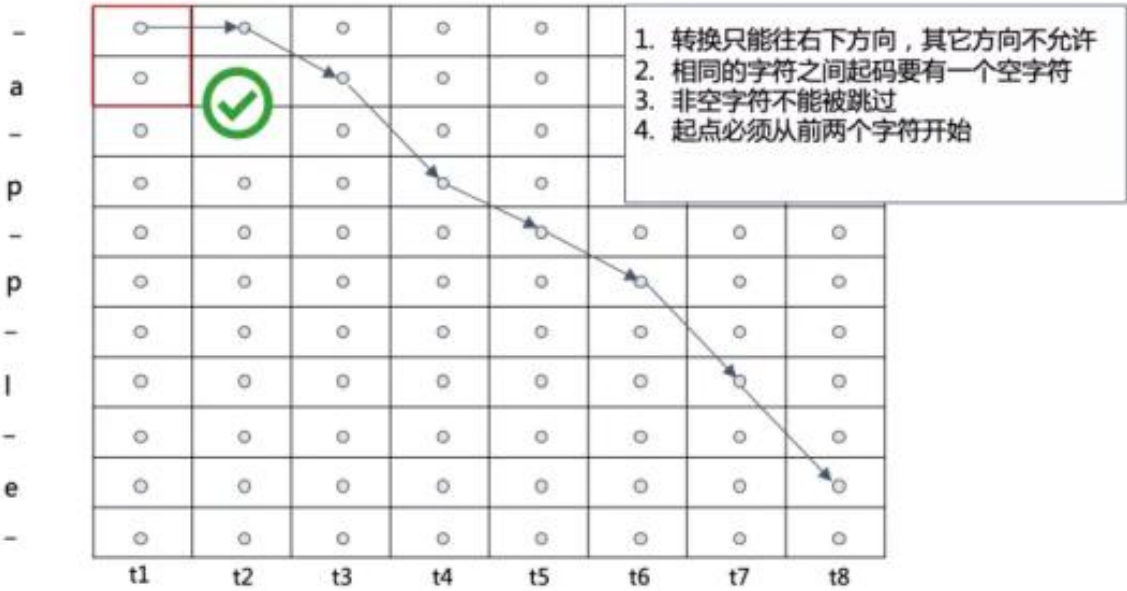


为了让所有的path都能在图中唯一、合法的表示，结点转换有如下约束：

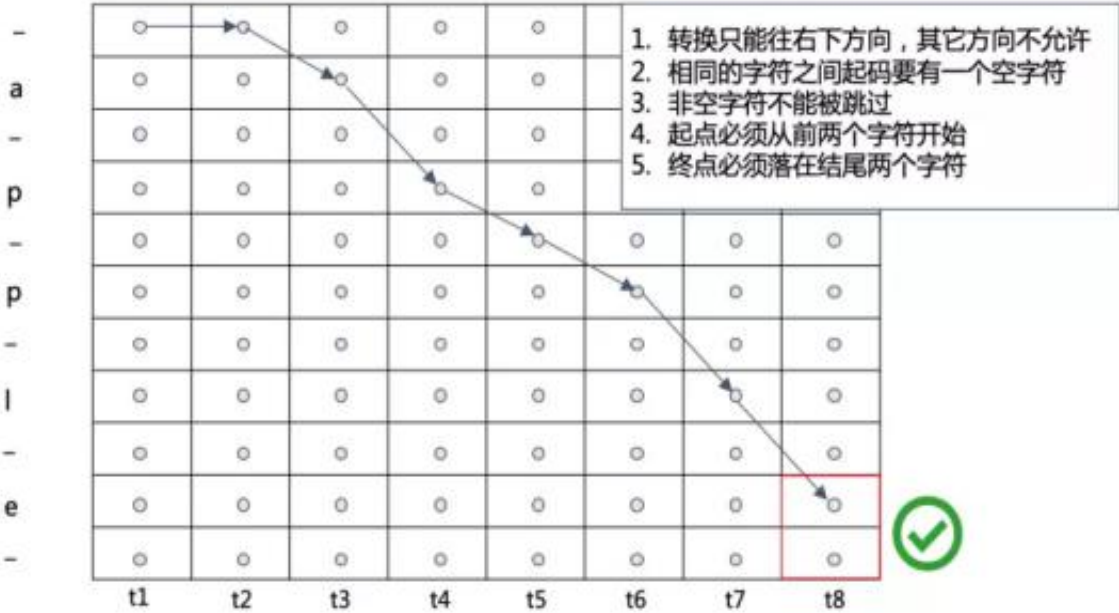


合法路径约束

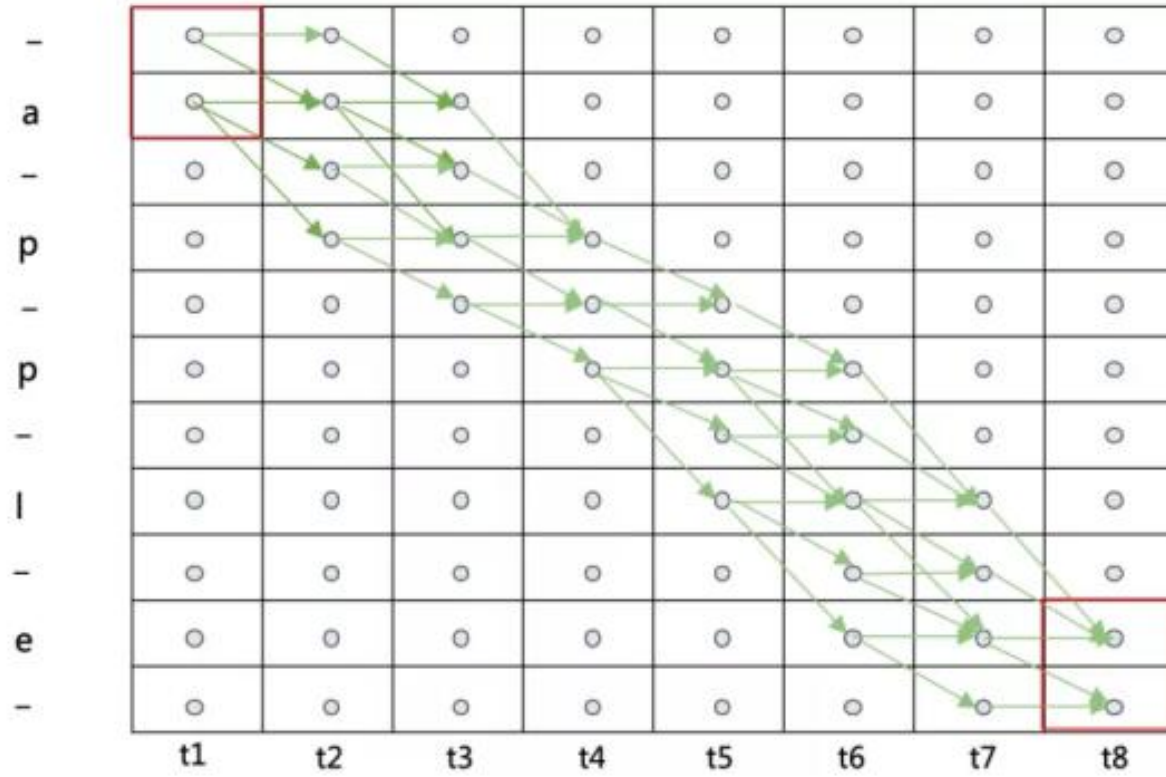
为了让所有的path都能在图中唯一、合法的表示，结点转换有如下约束：



为了让所有的path都能在图中唯一、合法的表示，结点转换有如下约束：

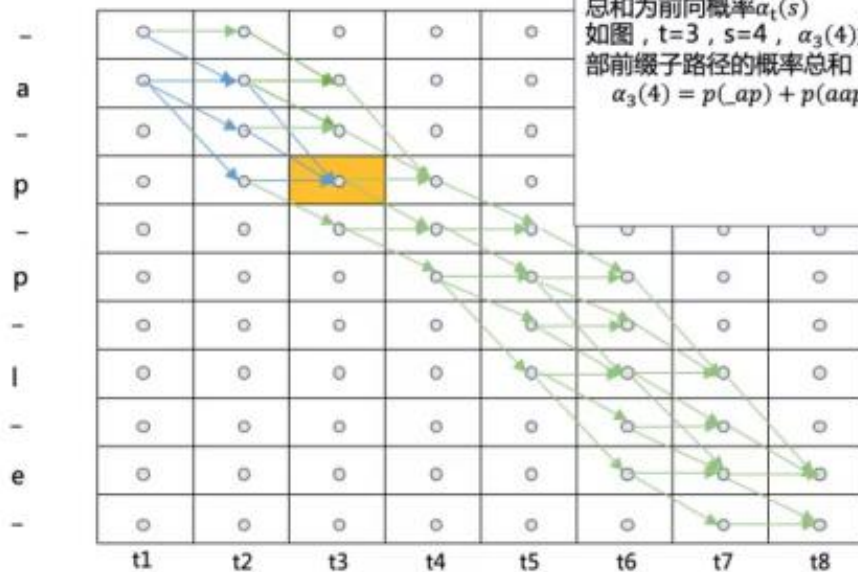


合法路径约束



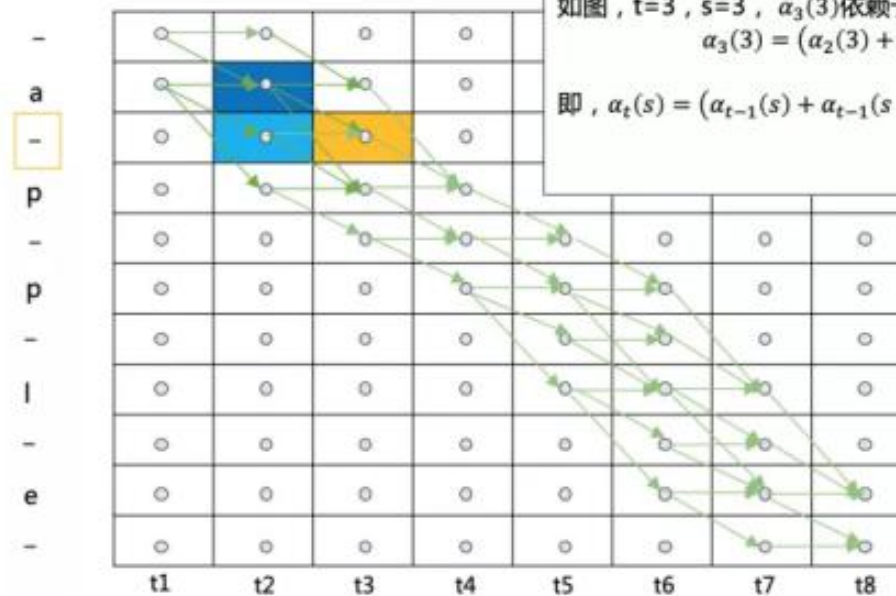
前向后向

将某节点的全部路径分为前后两部分



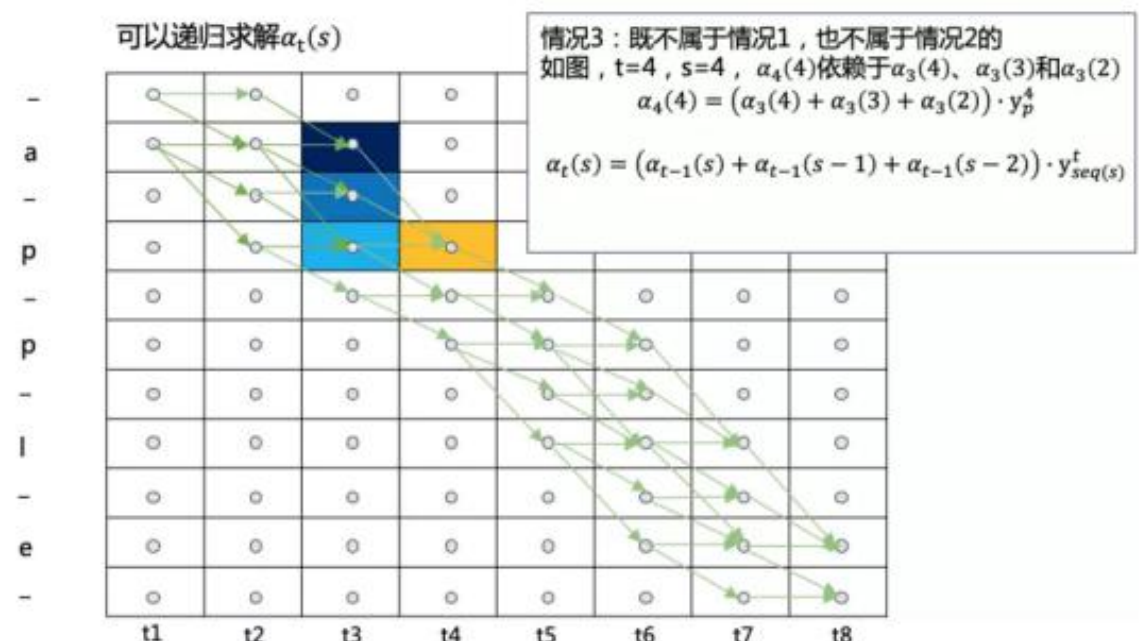
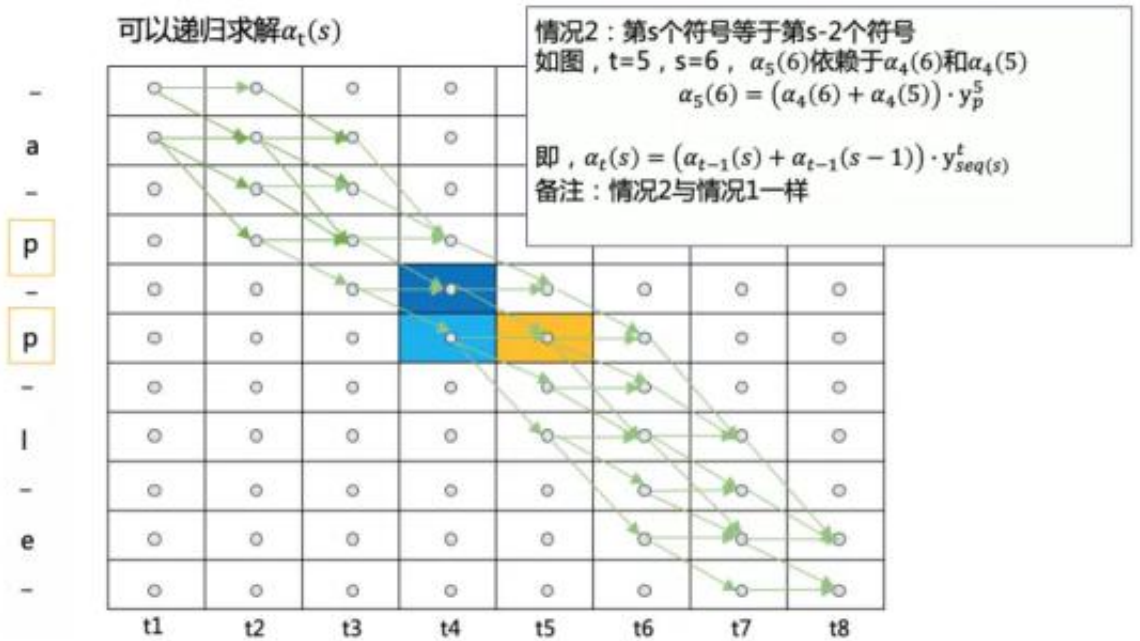
定义在时刻 t 经过节点 s 的全部前缀子路径的概率总和为前向概率 $\alpha_t(s)$
 如图, $t=3, s=4$, $\alpha_3(4)$ 为所有经过该节点的全部前缀子路径的概率总和:
 $\alpha_3(4) = p(_ap) + p(aap) + p(a_p) + p(app)$

可以递归求解 $\alpha_t(s)$



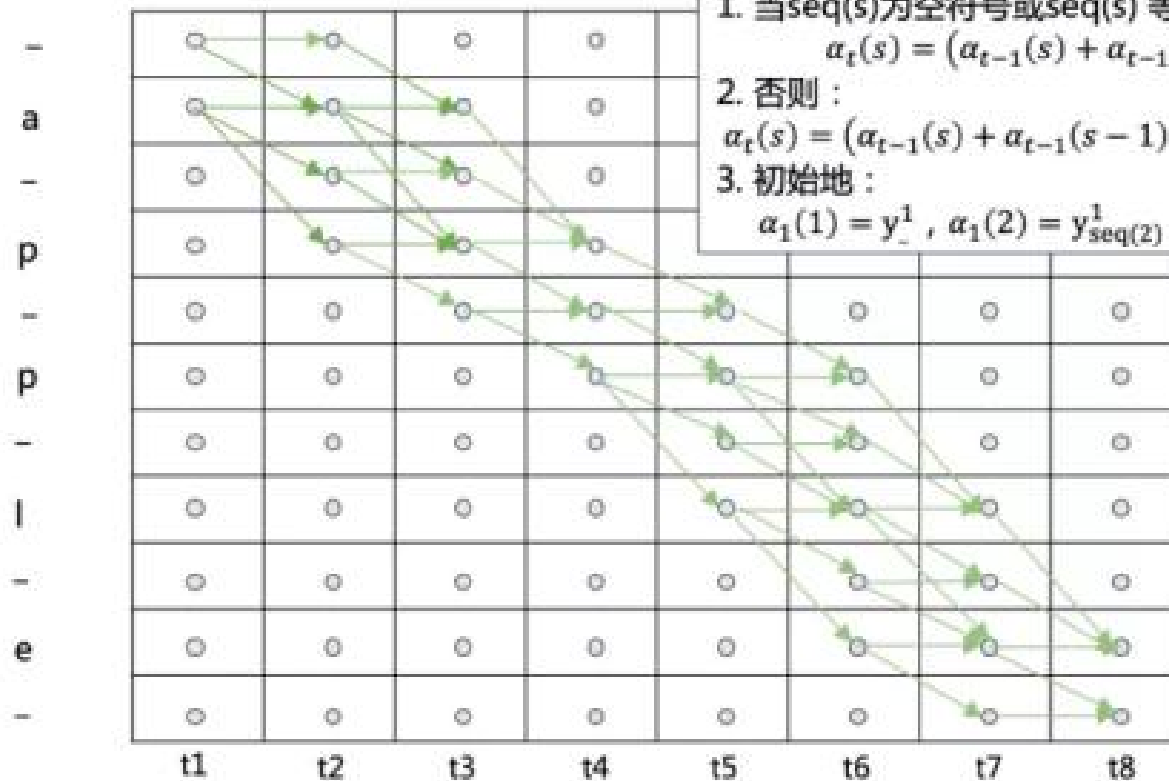
情况1: 第 s 个符号为空符号blank
 如图, $t=3, s=3$, $\alpha_3(3)$ 依赖于 $\alpha_2(3)$ 和 $\alpha_2(2)$
 $\alpha_3(3) = (\alpha_2(3) + \alpha_2(2)) \cdot y_{seq}^3$
 即, $\alpha_t(s) = (\alpha_{t-1}(s) + \alpha_{t-1}(s-1)) \cdot y_{seq}^t$

前向后向



前向后向

可以递归求解 $\alpha_t(s)$



所以，可以基于动态规划算法计算 $\alpha_t(s)$ ：

1. 当 $\text{seq}(s)$ 为空符号或 $\text{seq}(s)$ 等于 $\text{seq}(s-2)$ 时：

$$\alpha_t(s) = (\alpha_{t-1}(s) + \alpha_{t-1}(s-1)) \cdot y_{\text{seq}(s)}^t$$

2. 否则：

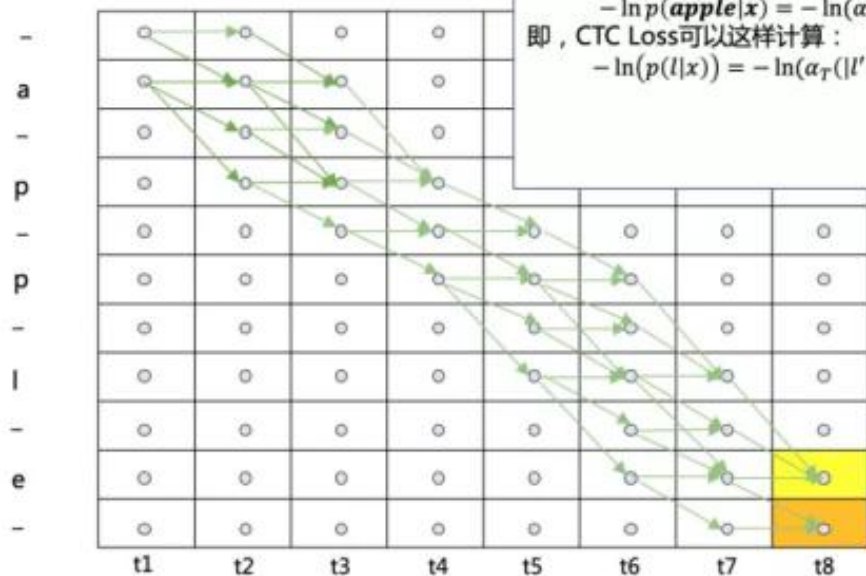
$$\alpha_t(s) = (\alpha_{t-1}(s) + \alpha_{t-1}(s-1) + \alpha_{t-1}(s-2)) \cdot y_{\text{seq}(s)}^t$$

3. 初始地：

$$\alpha_1(1) = y_-^1, \alpha_1(2) = y_{\text{seq}(2)}^1, \alpha_1(s) = 0, \forall s > 2$$

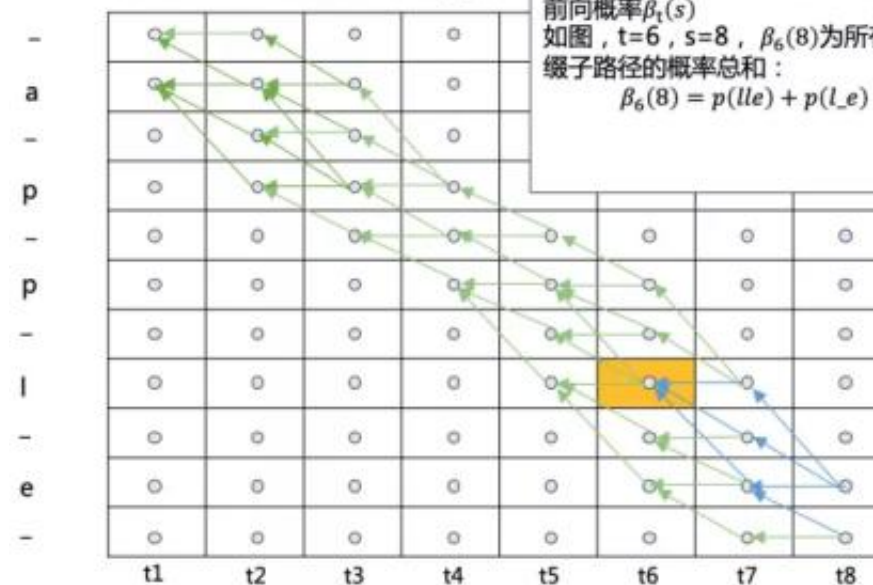
前向后向

利用 $\alpha_t(s)$ 计算CTC Loss



比如, $T=8$, 标签apple的CTC Loss计算的 $\alpha_t(s)$:
 $-\ln p(\mathbf{apple}|\mathbf{x}) = -\ln(\alpha_8(10) + \alpha_8(11))$
即, CTC Loss可以这样计算:
 $-\ln(p(l|x)) = -\ln(\alpha_T(|l'|) + \alpha_T(|l'| - 1))$

类似的, 可以定义反向概率 $\beta_t(s)$



定义在时刻 t 经过节点 s 的全部后缀子路径的概率总和为前向概率 $\beta_t(s)$
如图, $t=6, s=8$, $\beta_6(8)$ 为所有经过该节点的全部后缀子路径的概率总和:
 $\beta_6(8) = p(lle) + p(l_e) + p(lee) + p(le_)$

前向后向

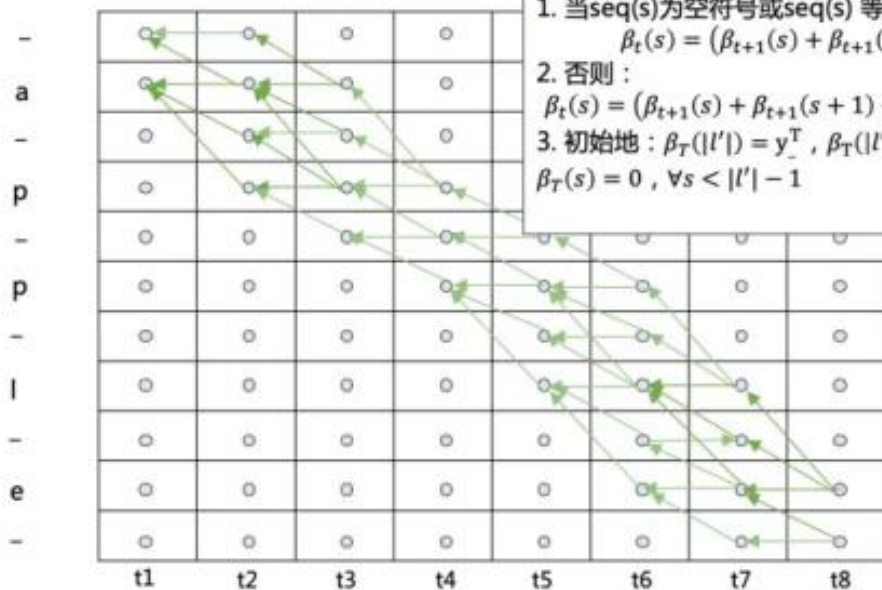
类似地，可以递归求解 $\beta_t(s)$

类似地，可以基于动态规划算法计算 $\beta_t(s)$ ：

1. 当 $\text{seq}(s)$ 为空符号或 $\text{seq}(s)$ 等于 $\text{seq}(s+2)$ 时：

$$\beta_t(s) = (\beta_{t+1}(s) + \beta_{t+1}(s+1)) \cdot y_{\text{seq}(s)}^t$$
2. 否则：

$$\beta_t(s) = (\beta_{t+1}(s) + \beta_{t+1}(s+1) + \beta_{t+1}(s+2)) \cdot y_{\text{seq}(s)}^t$$
3. 初始地： $\beta_T(|l'|) = y_{-}^T$ ， $\beta_T(|l'| - 1) = y_{\text{seq}(|l'|-1)}^T$ ，
 $\beta_T(s) = 0, \forall s < |l'| - 1$

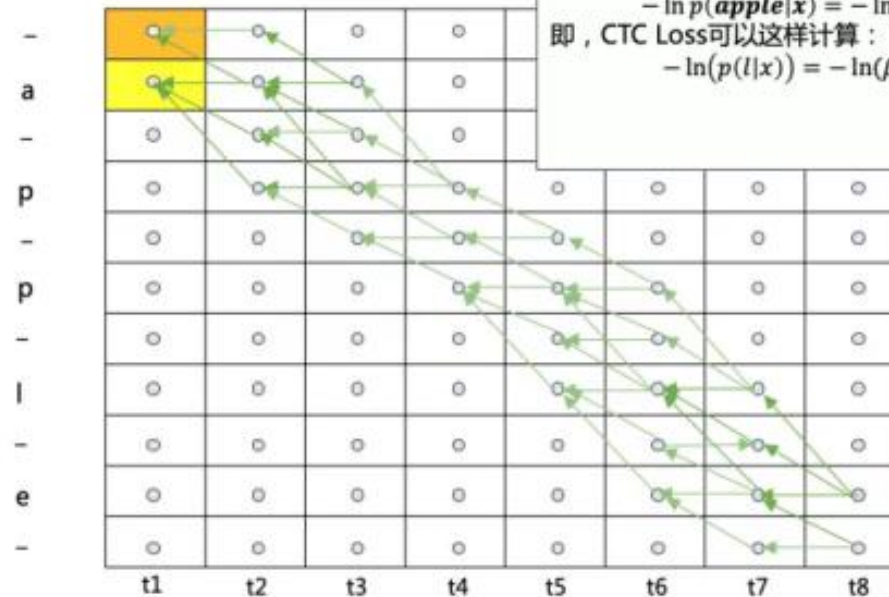


可以利用 $\beta_t(s)$ 计算CTC Loss

比如， $T=8$ ，标签apple的CTC Loss计算的 $\beta_t(s)$ ：

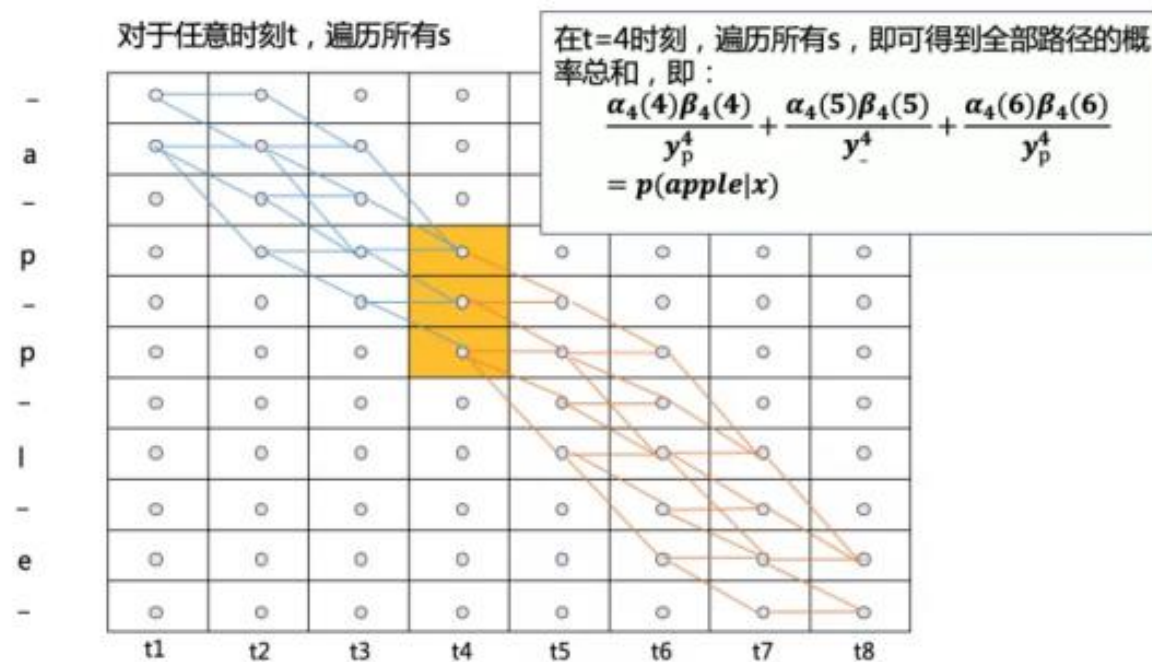
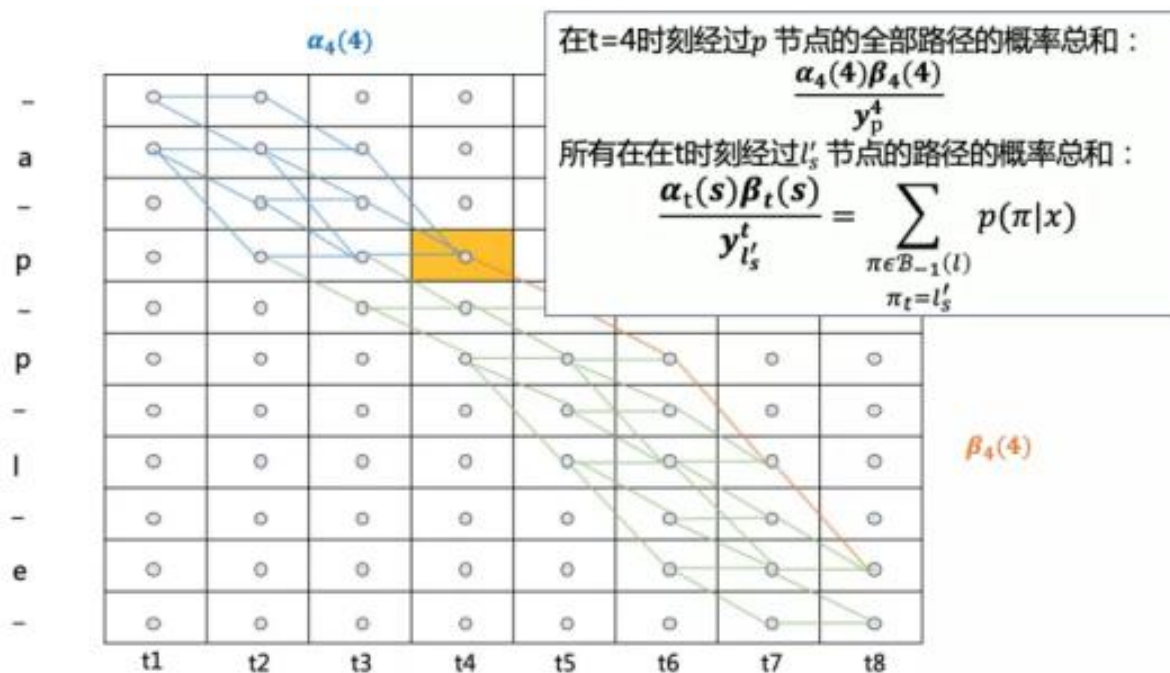
$$-\ln p(\text{apple}|x) = -\ln(\beta_1(1) + \beta_1(2))$$

即，CTC Loss可以这样计算：

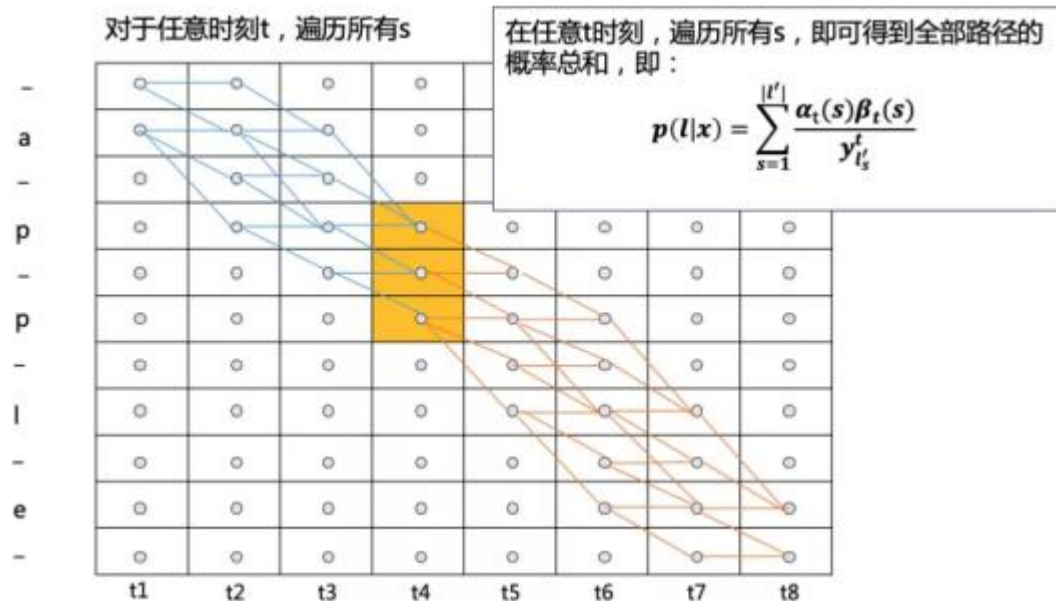
$$-\ln(p(l|x)) = -\ln(\beta_1(1) + \beta_1(2))$$


前向后向

把前向概率和后向概率结合起来也可以计算CTC Loss函数



前向后向



前向、后向计算LOSS

- 对于时序长度为 T 的输入序列 x 和输出序列 z ，前向概率：

- $\alpha_t(s) = \sum_{\pi_t=l'_s} \sum_{\pi \in \mathcal{B}^{-1}(z)} p(\pi_{1:t}|x)$

- $\alpha_1(1) = y_-^1, \alpha_1(2) = y_{l'_2}^1, \alpha_1(s) = 0, \forall s > 2$

- $\alpha_t(s) = 0, \forall s < |l'| - 2(T - t) - 1$

- $\alpha_t(s) = \begin{cases} (\alpha_{t-1}(s) + \alpha_{t-1}(s-1))y_{l'_s}^t & \text{if } l'_s = \mathbf{b} \text{ or } l'_{s-2} = l'_s \\ (\alpha_{t-1}(s) + \alpha_{t-1}(s-1) + \alpha_{t-1}(s-2))y_{l'_s}^t & \text{otherwise} \end{cases}$

- 利用 $\alpha_t(s)$ 计算CTC Loss：

- $-\ln(p(l|x)) = -\ln(\alpha_T(|l'|) + \alpha_T(|l'| - 1))$

- 对于时序长度为 T 的输入序列 x 和输出序列 z ，后向概率：

- $\beta_t(s) = \sum_{\pi_t=l'_s} \sum_{\pi \in \mathcal{B}^{-1}(z)} p(\pi_{t:T}|x)$

- $\beta_T(|l'|) = y_-^T, \beta_T(|l'| - 1) = y_{l'_{|l'|-1}}^T, \beta_T(s) = 0 \forall s > |l'| - 1$

- $\beta_t(s) = 0 \forall s > 2t \text{ and } \forall s > |l'|$

- $\beta_t(s) = \begin{cases} (\beta_{t+1}(s) + \beta_{t+1}(s+1))y_{l'_s}^t & \text{if } l'_s = \mathbf{b} \text{ or } l'_{s+2} = l'_s \\ (\beta_{t+1}(s) + \beta_{t+1}(s+1) + \beta_{t-1}(s+2))y_{l'_s}^t & \text{otherwise} \end{cases}$

- 利用 $\beta_t(s)$ 计算CTC Loss：

- $-\ln(p(l|x)) = -\ln(\beta_1(1) + \beta_1(2))$

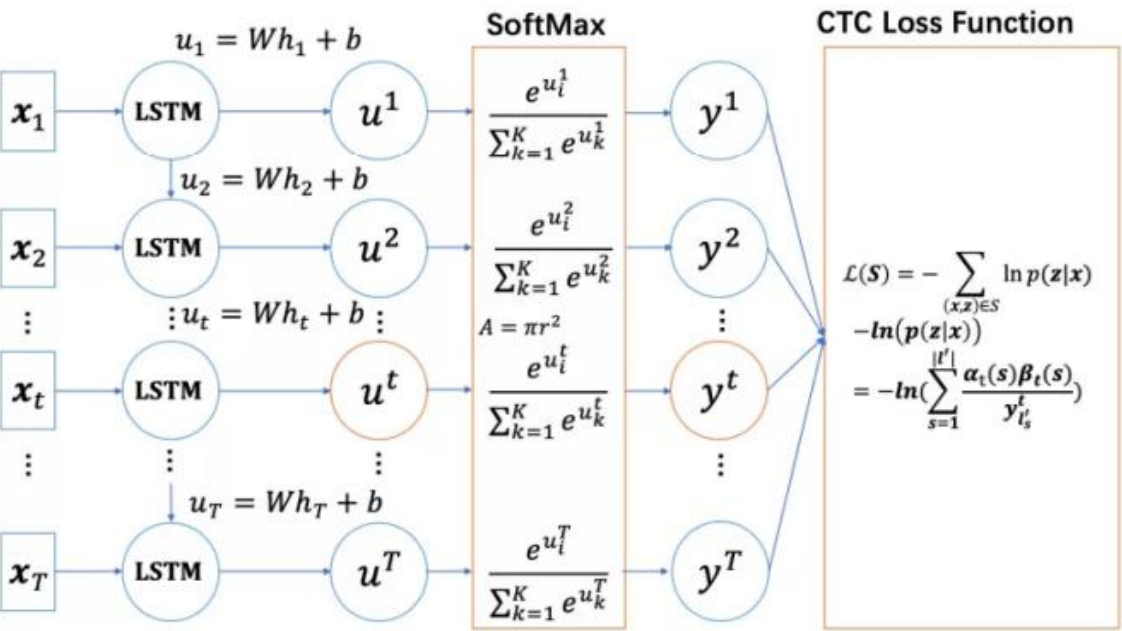
前向后向计算 LOSS

根据任意时刻的前向概率和后向概率计算CTC Loss函数，得到以下结论：

- 对于任意时刻 t ，利用前向概率和后向概率计算CTC Loss：

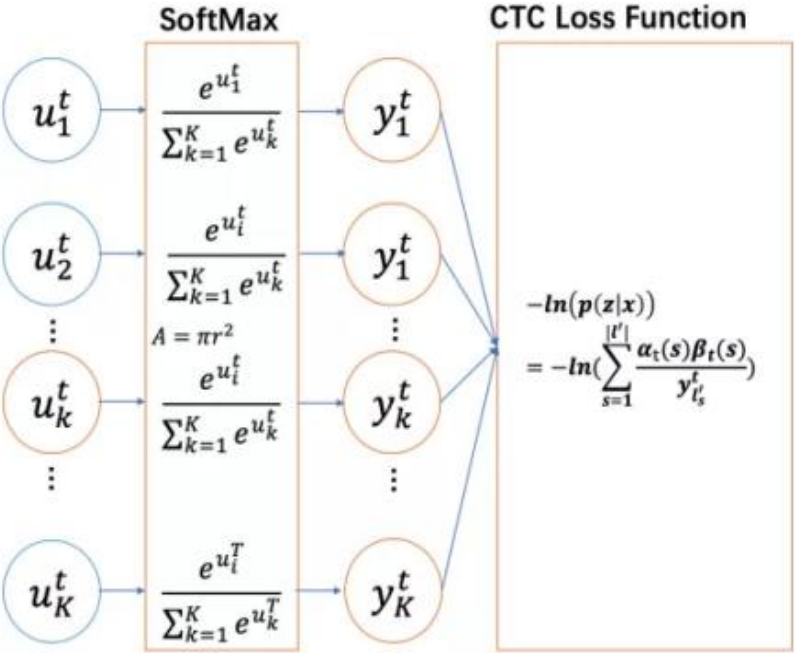
$$p(l|x) = \sum_{s=1}^{|l'|} \frac{\alpha_t(s)\beta_t(s)}{y_{l'_s}^t}$$
$$-\ln(p(l|x)) = -\ln\left(\sum_{s=1}^{|l'|} \frac{\alpha_t(s)\beta_t(s)}{y_{l'_s}^t}\right)$$

CTC Loss求导



$$p(l|x) = \sum_{s=1}^{|I'|} \frac{\alpha_t(s)\beta_t(s)}{y_{l'_s}^t}$$
$$\frac{\partial p(l|x)}{\partial y_{k'}^t} = \frac{1}{y_{k'}^t} \sum_{s \in \text{inab}(l, k')} \alpha_t(s)\beta_t(s) \quad (1)$$
$$\frac{\partial \ln(p(l|x))}{\partial y_{k'}^t} = \frac{1}{p(l|x)y_{k'}^t} \sum_{s \in \text{inab}(l, k')} \alpha_t(s)\beta_t(s) \quad (2)$$
$$\frac{\partial y_{k'}^t}{\partial u_k^t} = y_{k'}^t \delta_{kk'} - y_{k'}^t y_k^t \quad (3)$$
$$-\frac{\partial \ln(p(l|x))}{\partial u_k^t} = -\sum_{k'=1}^K \frac{\partial \ln(p(l|x))}{\partial y_{k'}^t} \frac{\partial y_{k'}^t}{\partial u_k^t} \quad (4)$$

把公式(2)、(3)代入(4)得



CTC Loss求导

$$p(l|x) = \sum_{s=1}^{|l'|} \frac{\alpha_t(s)\beta_t(s)}{y_{l'_t}^t} \quad (1)$$

$$\frac{\partial p(l|x)}{\partial y_{k'}^t} = \frac{1}{y_{k'}^t{}^2} \sum_{s \in \text{lab}(l, k')} \alpha_t(s)\beta_t(s)$$

$$\frac{\partial \ln(p(l|x))}{\partial y_{k'}^t} = \frac{1}{p(l|x)y_{k'}^t{}^2} \sum_{s \in \text{lab}(l, k')} \alpha_t(s)\beta_t(s) \quad (2)$$

$$\frac{\partial y_{k'}^t}{\partial u_k^t} = y_{k'}^t \delta_{kk'} - y_{k'}^t y_k^t \quad (3)$$

$$-\frac{\partial \ln(p(l|x))}{\partial u_k^t} = -\sum_{k'=1}^K \frac{\partial \ln(p(l|x))}{\partial y_{k'}^t} \frac{\partial y_{k'}^t}{\partial u_k^t} \quad (4)$$

把公式(2)、(3)代入(4)得

$$\begin{aligned} -\frac{\partial \ln(p(l|x))}{\partial u_k^t} &= -\sum_{k'=1}^K \frac{\partial \ln(p(l|x))}{\partial y_{k'}^t} \frac{\partial y_{k'}^t}{\partial u_k^t} \\ &= \sum_{k'=1}^K \frac{y_k^t - \delta_{kk'}}{p(l|x)y_{k'}^t} \sum_{s \in \text{lab}(l, k')} \alpha_t(s)\beta_t(s) \\ &= -\sum_{k'=1}^K \frac{y_{k'}^t \delta_{kk'} - y_{k'}^t y_k^t}{p(l|x)y_{k'}^t{}^2} \sum_{s \in \text{lab}(l, k')} \alpha_t(s)\beta_t(s) \\ &= \sum_{k'=1}^K \frac{y_k^t}{p(l|x)y_{k'}^t} \sum_{s \in \text{lab}(l, k')} \alpha_t(s)\beta_t(s) - \sum_{k'=1}^K \frac{\delta_{kk'}}{p(l|x)y_{k'}^t} \sum_{s \in \text{lab}(l, k')} \alpha_t(s)\beta_t(s) \\ &= \frac{y_k^t}{p(l|x)} \sum_{k'=1}^K \frac{1}{y_{k'}^t} \sum_{s \in \text{lab}(l, k')} \alpha_t(s)\beta_t(s) - \sum_{k'=1}^K \frac{\delta_{kk'}}{p(l|x)y_{k'}^t} \sum_{s \in \text{lab}(l, k')} \alpha_t(s)\beta_t(s) \\ &= \frac{y_k^t}{p(l|x)} p(l|x) - \sum_{k'=1}^K \frac{\delta_{kk'}}{p(l|x)y_{k'}^t} \sum_{s \in \text{lab}(l, k')} \alpha_t(s)\beta_t(s) \\ &= y_k^t - \frac{1}{p(l|x)y_k^t} \sum_{s \in \text{lab}(l, k')} \alpha_t(s)\beta_t(s) \end{aligned}$$

BeamSearch

Beam Search的过程非常简单，每一步搜索选取概率最大的W个节点进行扩展，W也称为Beam Width，其核心还是计算每一步扩展节点的概率

b	0.3	0.3	0.1
a	0.2	0.3	0.3
—	0.5	0.4	0.6
	t1	t2	t3

横轴表示时间，纵轴表示每一步输出层的概率， $T=3$ ，字符集为{a, b}

穷举搜索

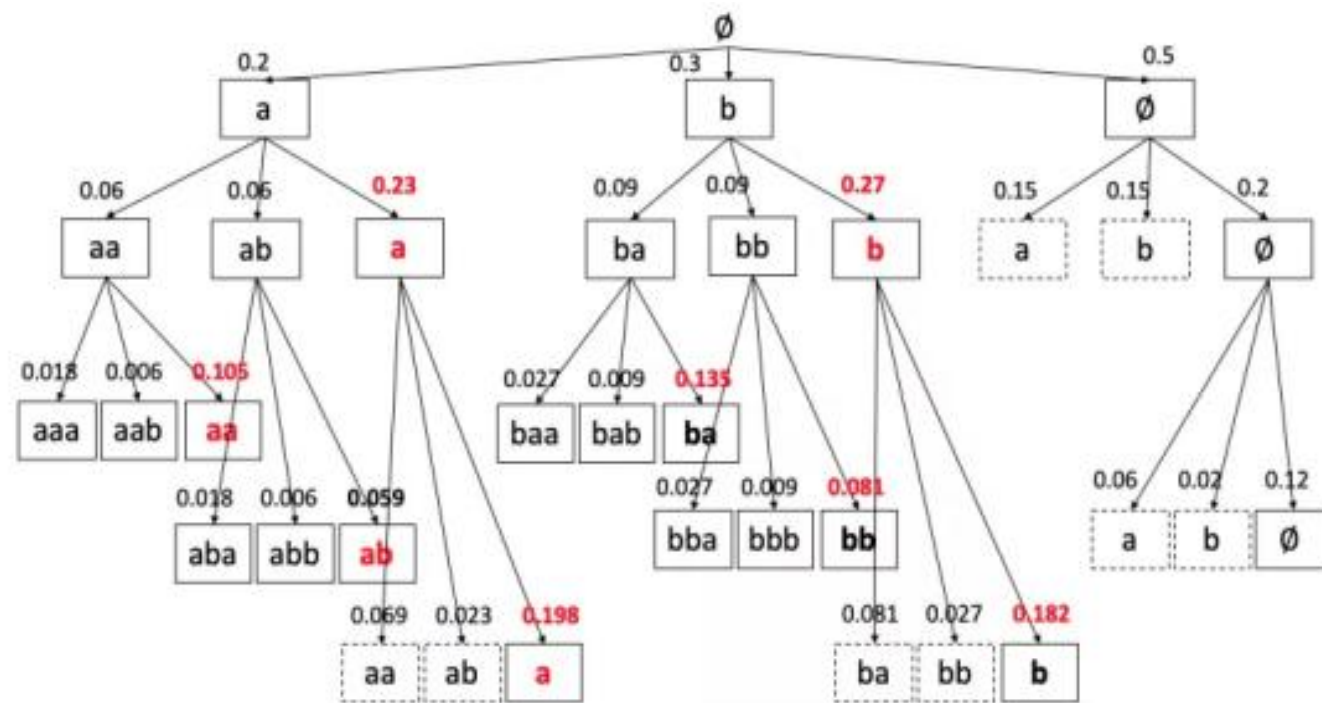
Beam Search的过程非常简单，每一步搜索选取概率最大的W个节点进行扩展，W也称为Beam Width，其核心还是计算每一步扩展节点的概率

b	0.3	0.3	0.1
a	0.2	0.3	0.3
—	0.5	0.4	0.6
	t1	t2	t3

横轴表示时间，纵轴表示每一步输出层的概率， $T=3$ ，字符集为{a, b}

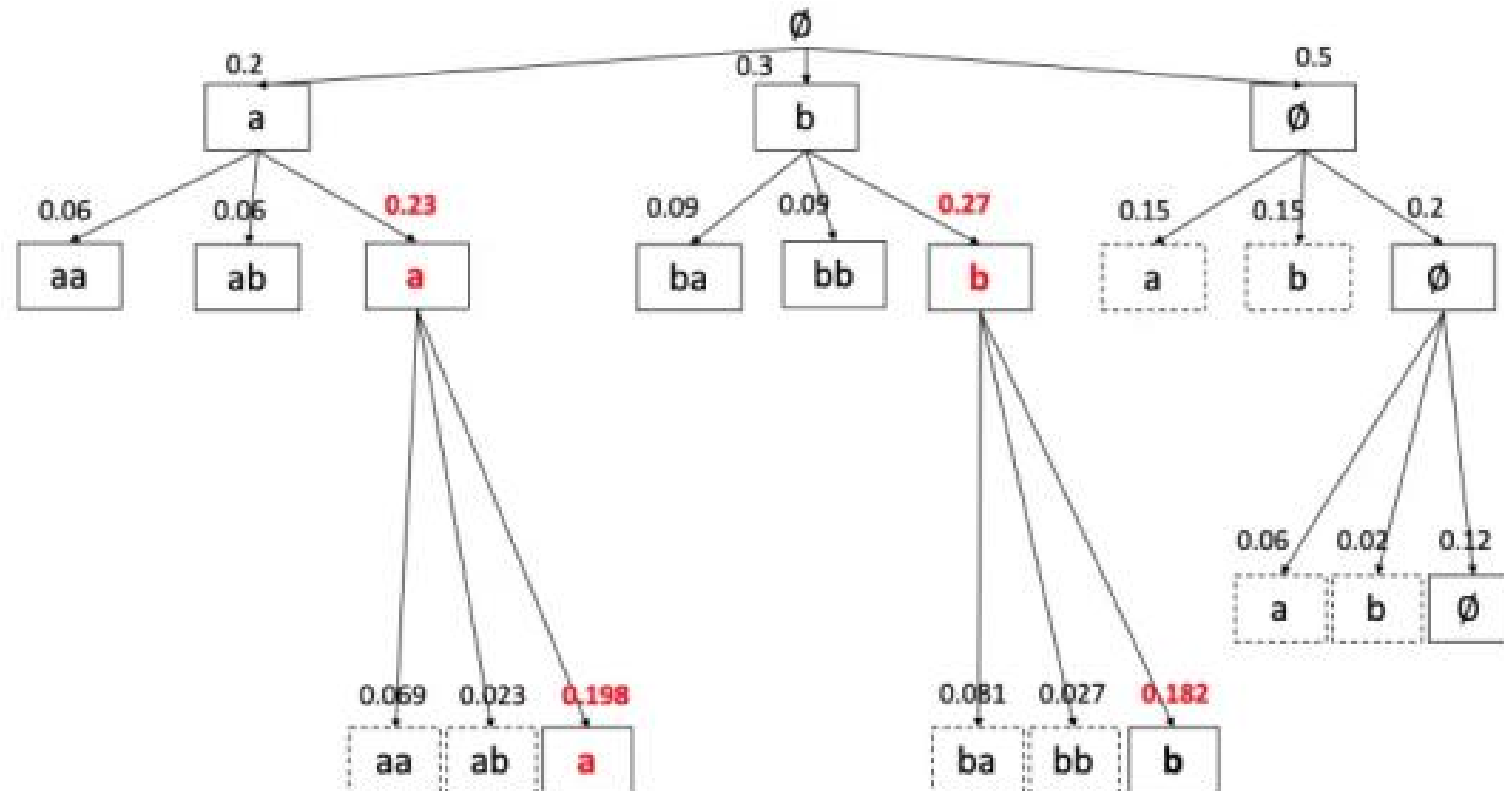
穷举搜索

如果对它的搜索空间进行穷举搜索，则每一步都展开进行搜索



BamSearch

穷举搜索每一步都要扩展全部节点，能保证最终找到最优解（上图中例子最优解 $l^*=b$ ， $p(l^*)=0.164$ ），但搜索复杂度太高，而Beam Search的思路很简单，每一步只选取扩展概率最大的 W 个节点进行扩展



感谢聆听