

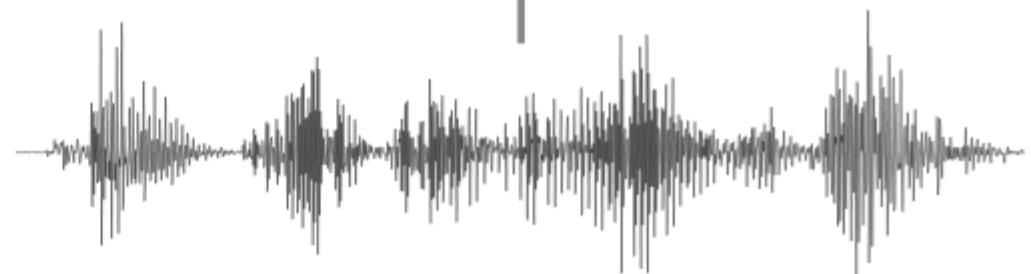
the quick brown fox



The quick brown fox

Handwriting recognition: The input can be (x, y) coordinates of a pen stroke or pixels in an image.

jumps over the lazy dog

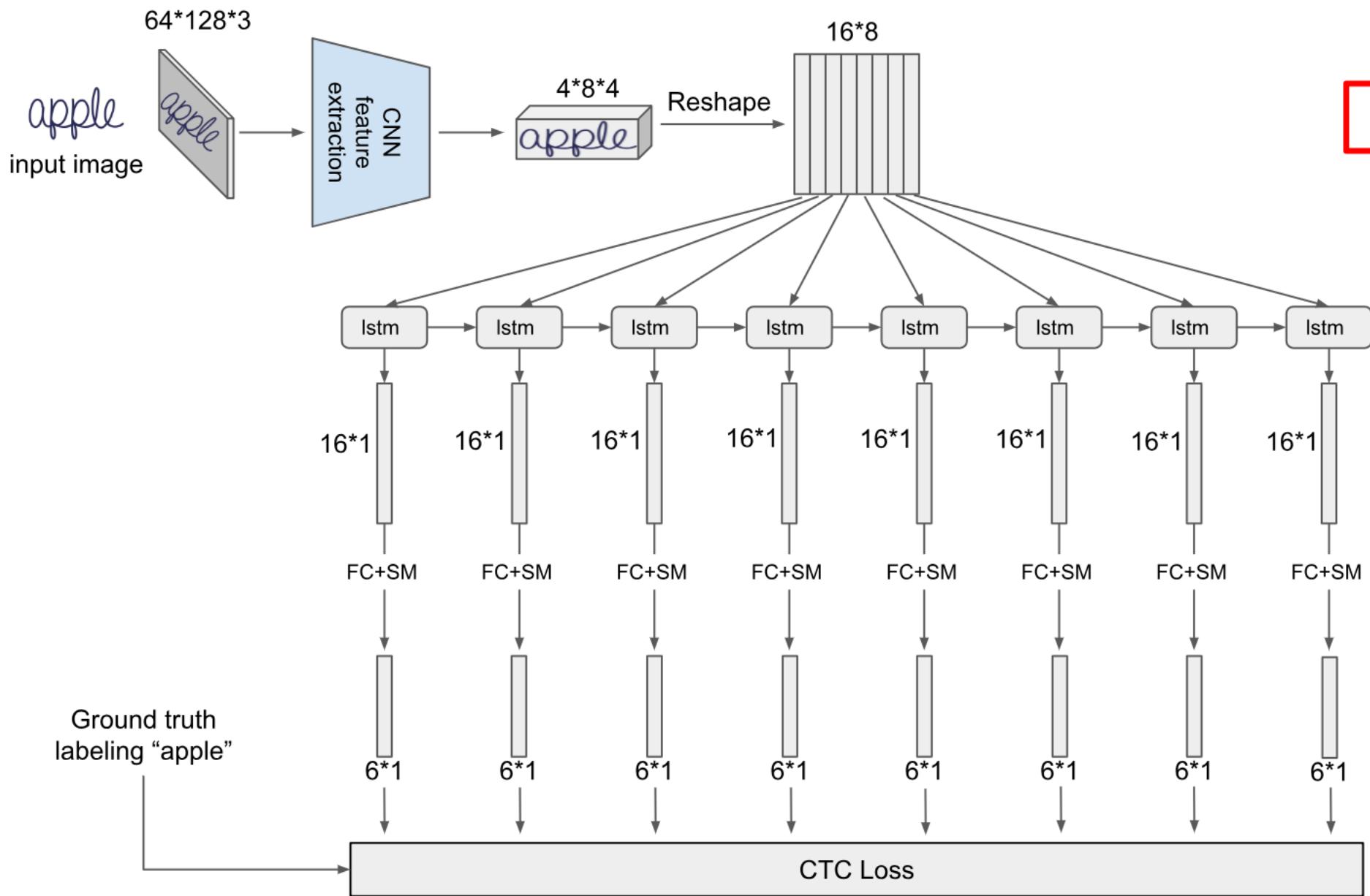


Speech recognition: The input can be a spectrogram or some other frequency based feature extractor.

Connectionist Temporal Classification Labelling Unsegmented Sequence Data with Recurrent Neural Networks

主讲人：杜臣

2018.07.14



[Link to CTC Loss paper](#)

1、CTC loss用来
处理输入特征和
label不对齐的训
练问题。

2、label长度需小
于输入长度。

blank“_”的作用

- 1、CTC引入了blank标记“_”，用来表示当前输出为空，以及用来分割连续的相同输出。
- 2、解码时先去除重复字符，再去除“_”

How CTC collapsing works

For an input,
like speech



Predict a
sequence of
tokens

h e l ε | o ε h e | | o

Merge repeats,
drop ε

h e l | o h e | o

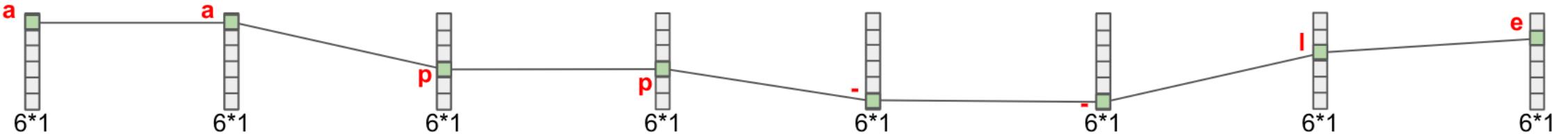
Final output

h e l | o h e | o



Path1: “ap-pl-ee” ————— B(“ap-pl-ee”) ————— Labeling: “apple”
 $p(\text{“ap-pl-ee”}) = y_a^1 \cdot y_p^2 \cdot y_l^3 \cdot y_p^4 \cdot y_l^5 \cdot y_e^6 \cdot y_e^7 \cdot y_e^8$

从RNN(LSTM)的输出中，能组成很多种路径，每一条路径的概率就是路径中每一时刻输出的概率乘。

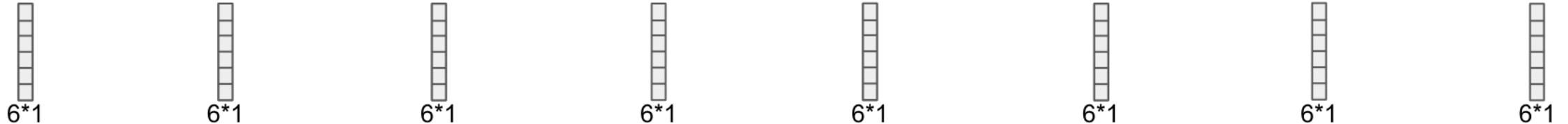


Path1: "ap-pl-ee" $\xrightarrow{B(\text{"ap-pl-ee"})}$ Labeling: "apple"

$$p(\text{"ap-pl-ee"}) = y_a^1 \cdot y_p^2 \cdot y_{-}^3 \cdot y_p^4 \cdot y_l^5 \cdot y_{-}^6 \cdot y_e^7 \cdot y_e^8$$

Path2: "aapp--le" $\xrightarrow{B(\text{"aapp--le"})}$ Labeling: "aple"

$$p(\text{"aapp--le"}) = y_a^1 \cdot y_a^2 \cdot y_p^3 \cdot y_p^4 \cdot y_{-}^5 \cdot y_{-}^6 \cdot y_l^7 \cdot y_e^8$$



Path1: "ap-pl-ee" $\xrightarrow{B(\text{"ap-pl-ee"})}$ Labeling: "apple"

$$p(\text{"ap-pl-ee"}) = y_a^1 \cdot y_p^2 \cdot y_{-}^3 \cdot y_p^4 \cdot y_l^5 \cdot y_{-}^6 \cdot y_e^7 \cdot y_e^8$$

Path2: "aapp--le" $\xrightarrow{B(\text{"aapp--le"})}$ Labeling: "aple"

$$p(\text{"aapp--le"}) = y_a^1 \cdot y_a^2 \cdot y_p^3 \cdot y_p^4 \cdot y_{-}^5 \cdot y_l^6 \cdot y_{-}^7 \cdot y_e^8$$

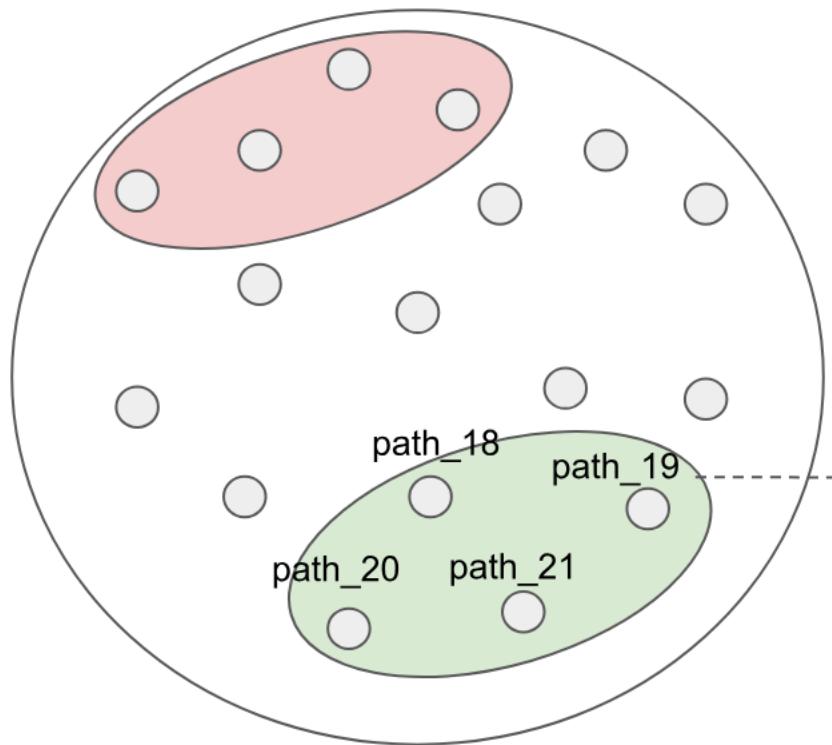
Path3: "aap--lzz" $\xrightarrow{B(\text{"aap--lzz"})}$ Labeling: "aplz"

...

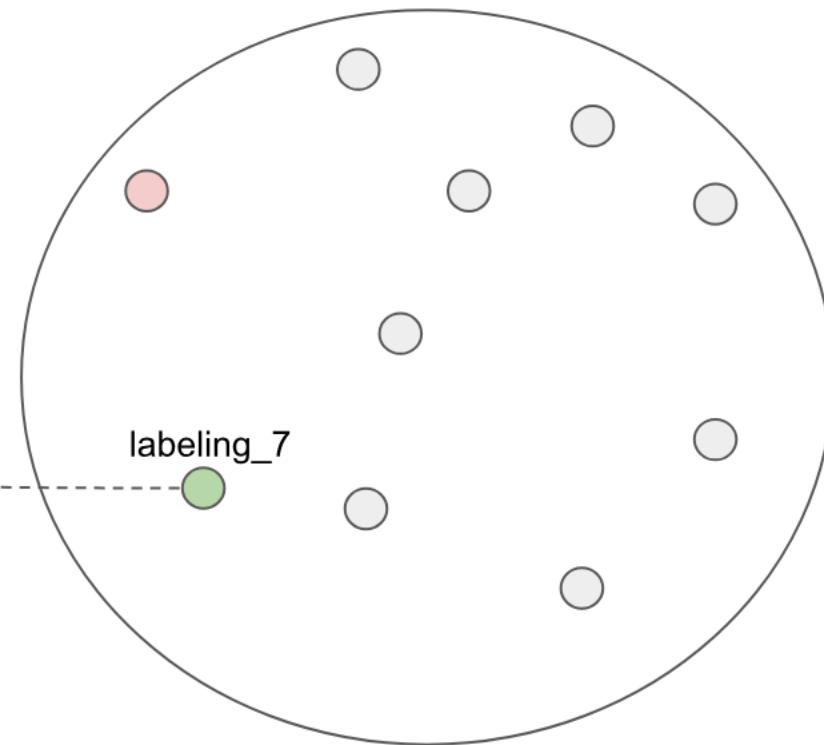
...

...

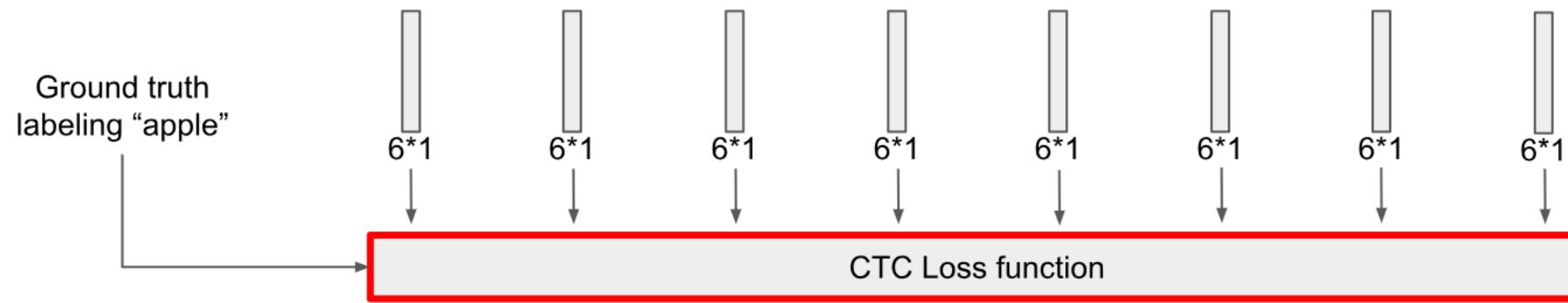
PathN: "--z--eee" $\xrightarrow{B(\text{"--z--eee"})}$ Labeling: "ze"



Paths that are correspond to certain label



$$\begin{aligned} p(\text{labeling}_7) &= \text{sum of probabilities of all corresponding paths} = \\ &= p(\text{path}_{18}) + p(\text{path}_{19}) + p(\text{path}_{20}) + p(\text{path}_{21}) \end{aligned}$$



$$\text{CTC Loss} = -\ln(p(\text{"apple"}))$$

It seems very simple, but there is one tricky problem.

In our example there are exist $6^8 = 1\ 679\ 616$ possible paths. For a larger dictionary size and for a larger number of lstm steps number of possible paths will be huge.

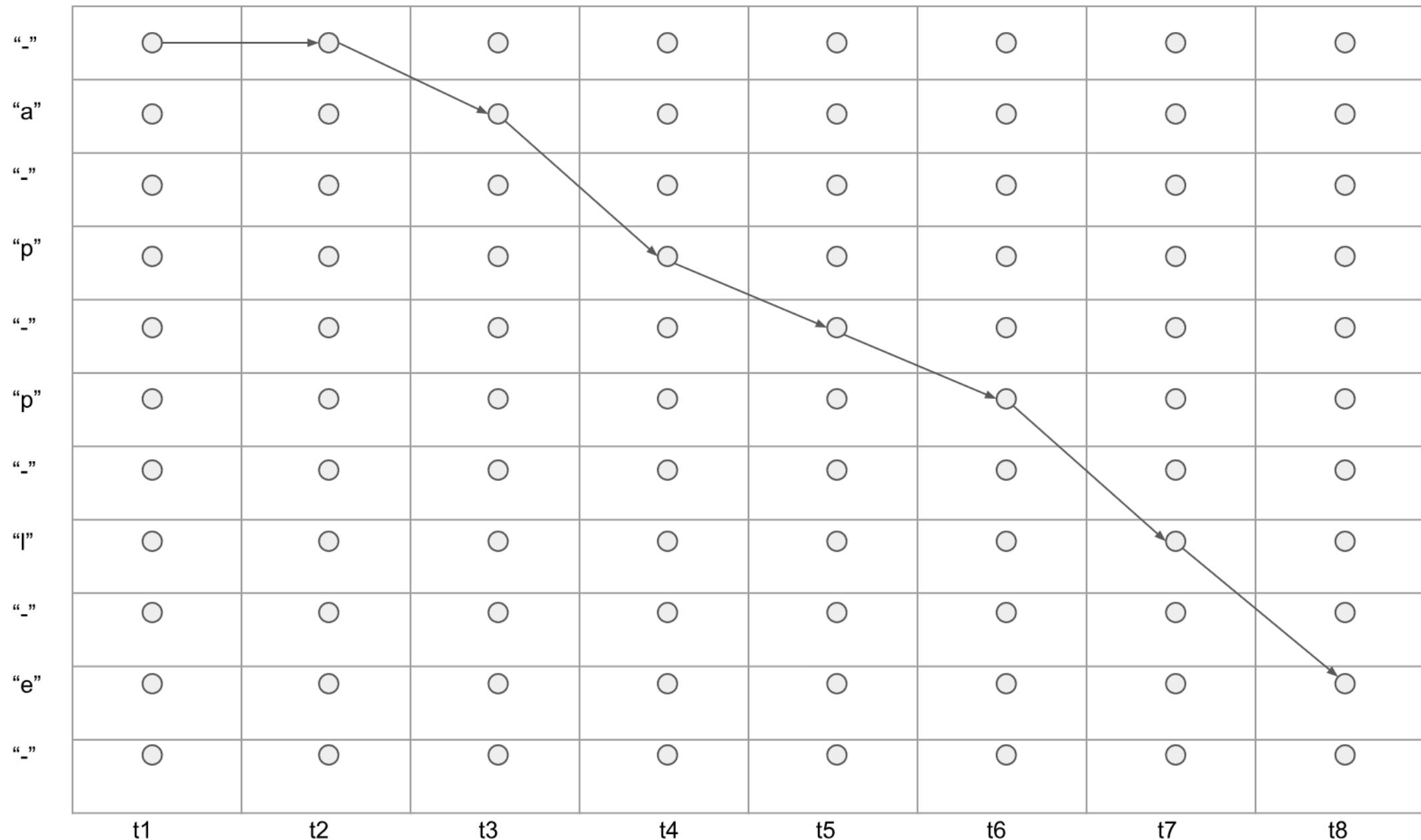
For a given labeling we can not compute the sum of all paths probabilities, because there are very many of these.

Fortunately there is an efficient way of calculation ground truth labeling probability. The problem can be solved with dynamic programming algorithm.

CTC Loss calculation: dynamic programming
algorithm

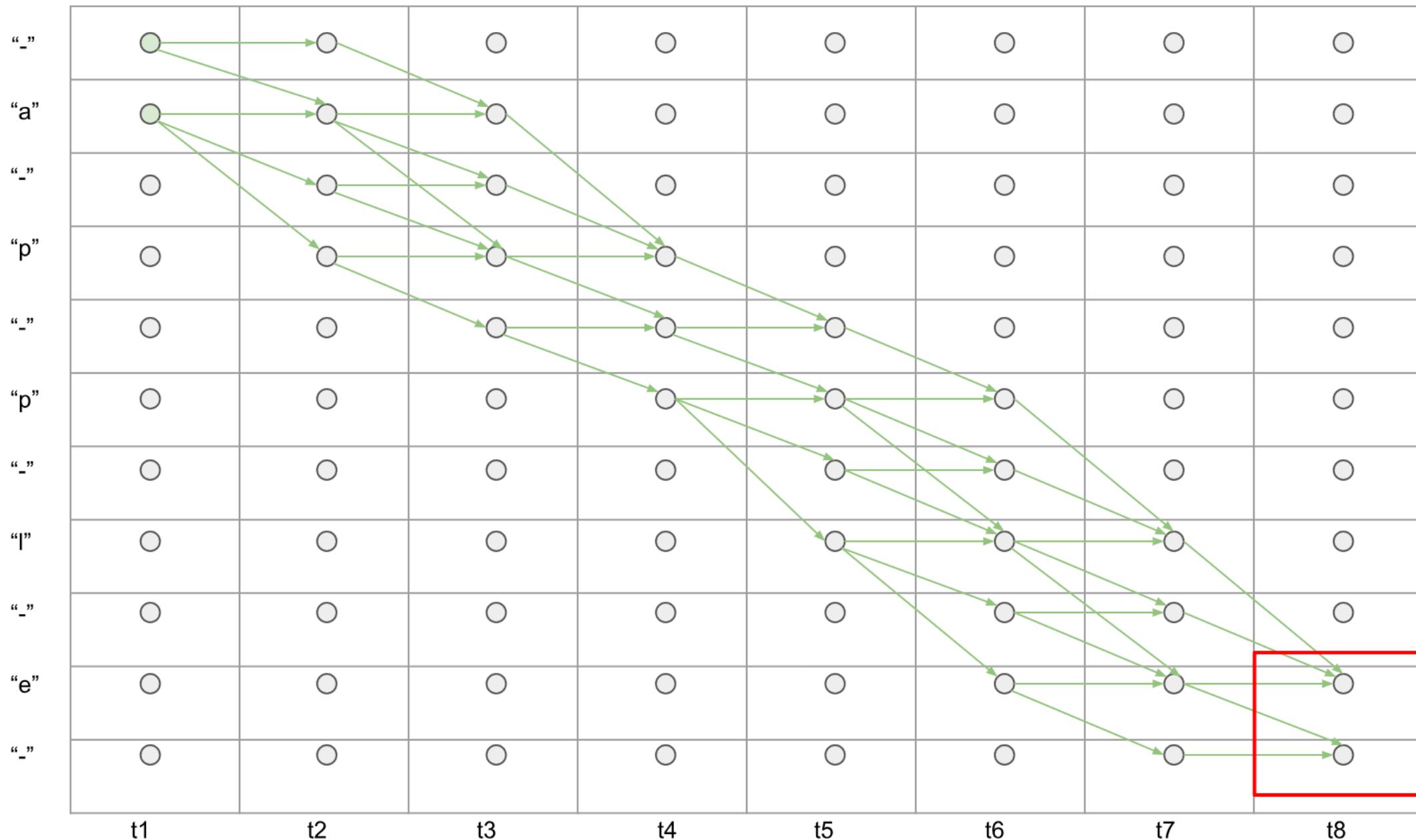
With this table we can model all paths that are correspond to ground truth labeling “apple”

For example: path “--ap-ple” can be mapped to labeling “apple”. i.e: $B(\text{“--ap-ple”}) = \text{“apple”}$

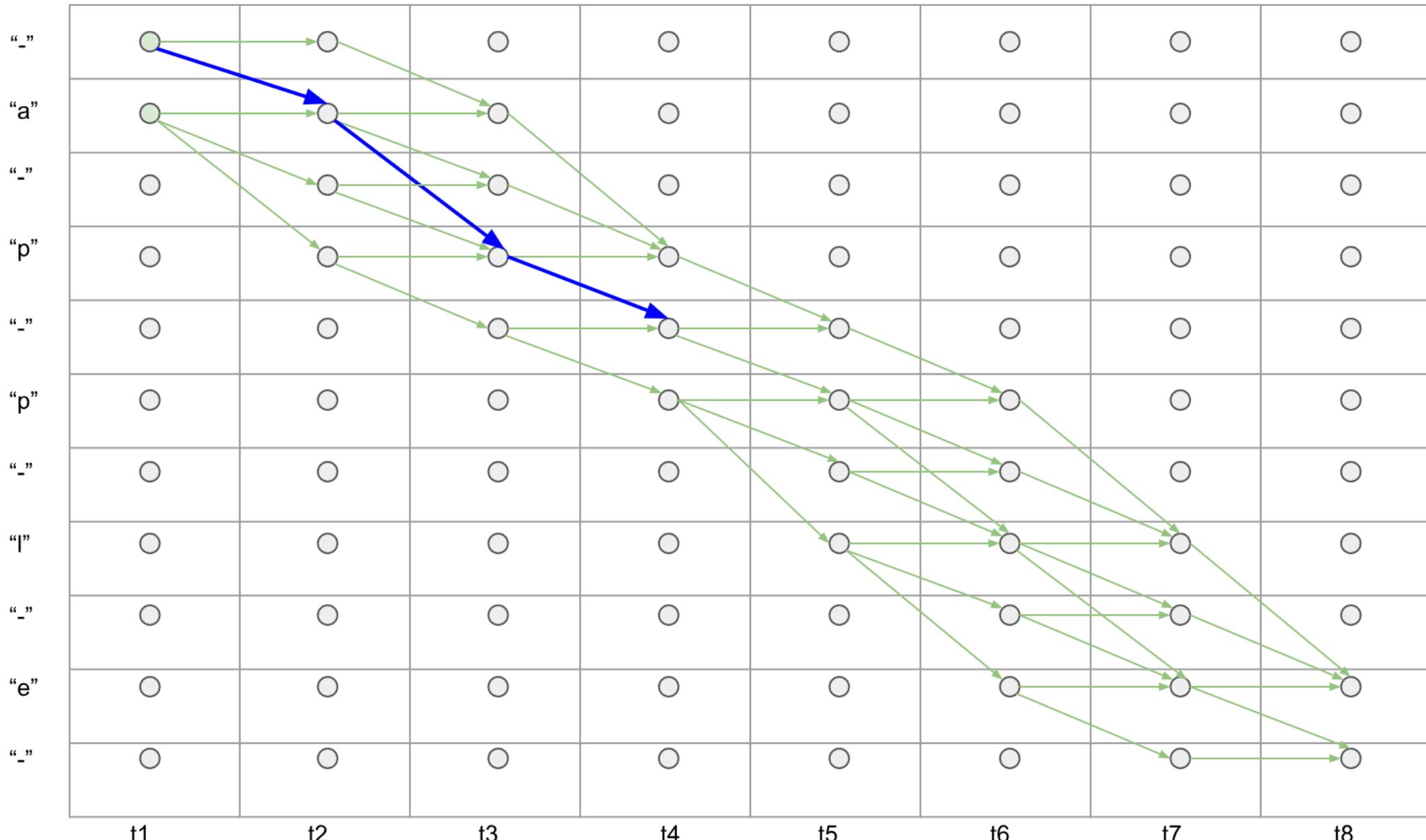


Initialization: paths can start only with this symbols.

Interesting note: all valid paths should end up with **this** nodes.



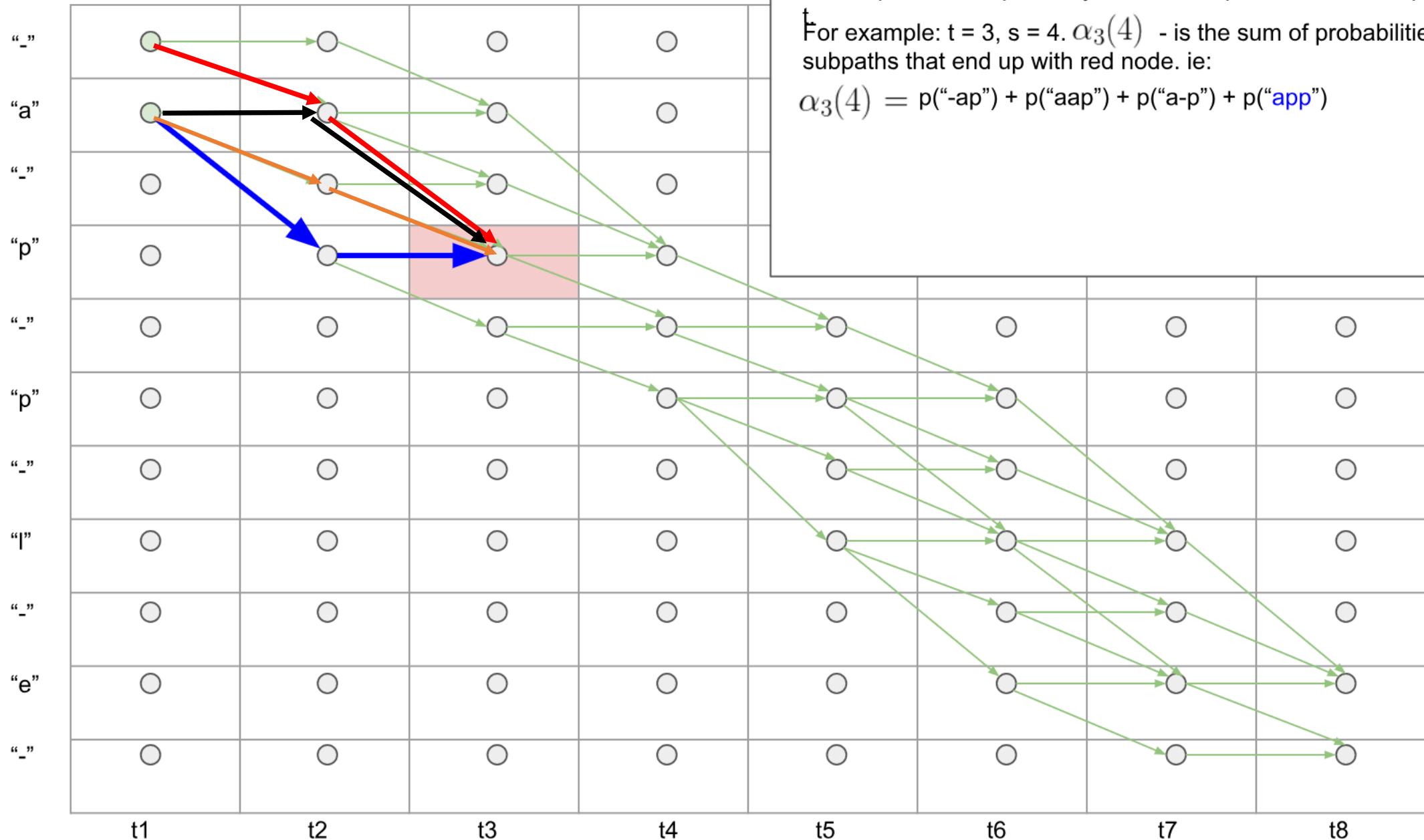
Let's consider subpath “-ap-” (blue color). Let me remind you that the probability of this subpath is simply the multiplication of corresponding network output probabilities - $p(“ - ap - ”) = y_-^1 \cdot y_a^2 \cdot y_p^3 \cdot y_-^4$



Let me define variable $\alpha_t(s)$ - the total probability of all subpaths, whose prefix end up with symbol at s-th position in the sequence at time t

For example: $t = 3, s = 4$. $\alpha_3(4)$ - is the sum of probabilities of all subpaths that end up with red node. ie:

$$\alpha_3(4) = p("ap") + p("aap") + p("a-p") + p("app")$$

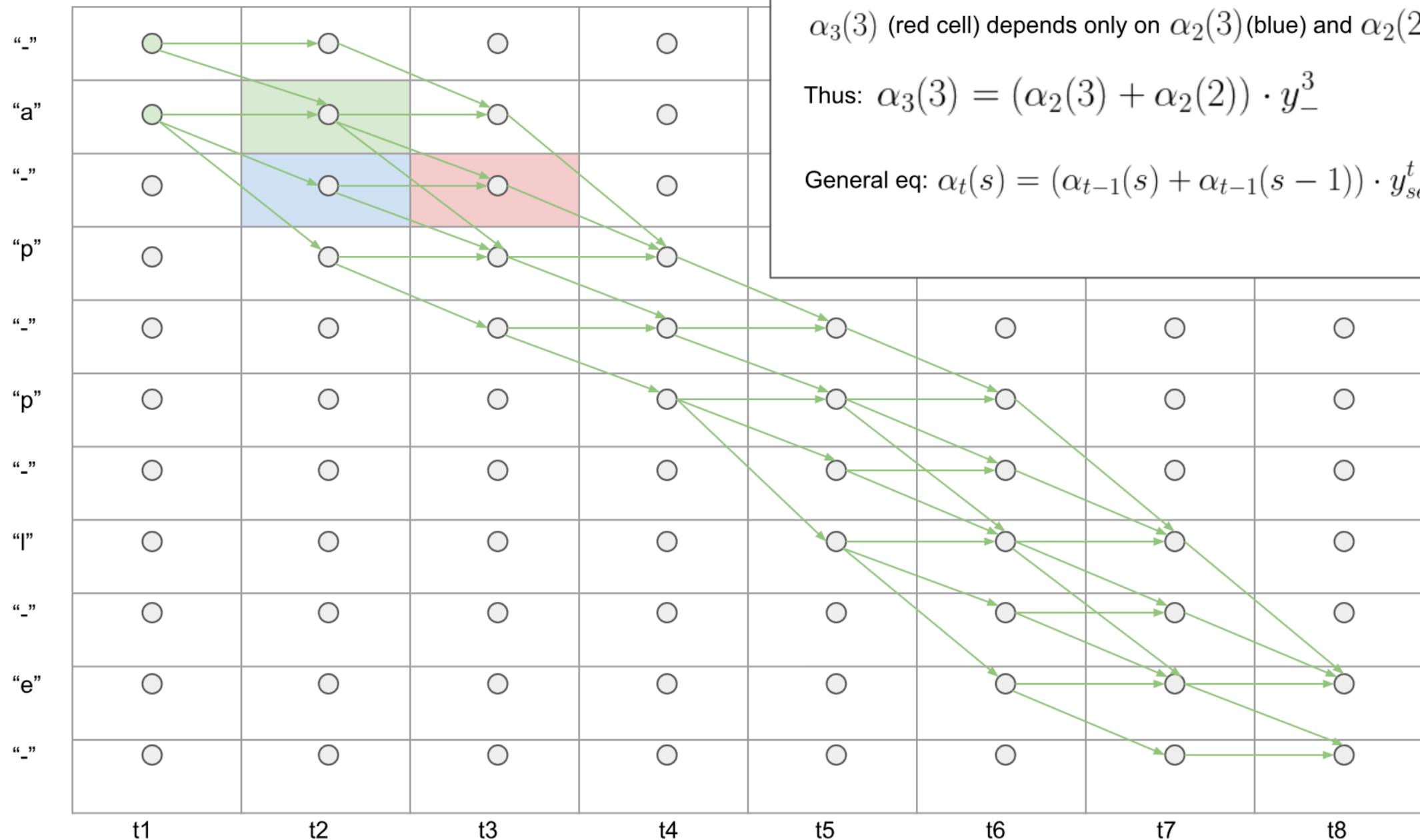


Case1. s-th symbol is blank. Example: s=3, 3-rd symbol is “-”

$\alpha_3(3)$ (red cell) depends only on $\alpha_2(3)$ (blue) and $\alpha_2(2)$ (green)

$$\text{Thus: } \alpha_3(3) = (\alpha_2(3) + \alpha_2(2)) \cdot y_-^3$$

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

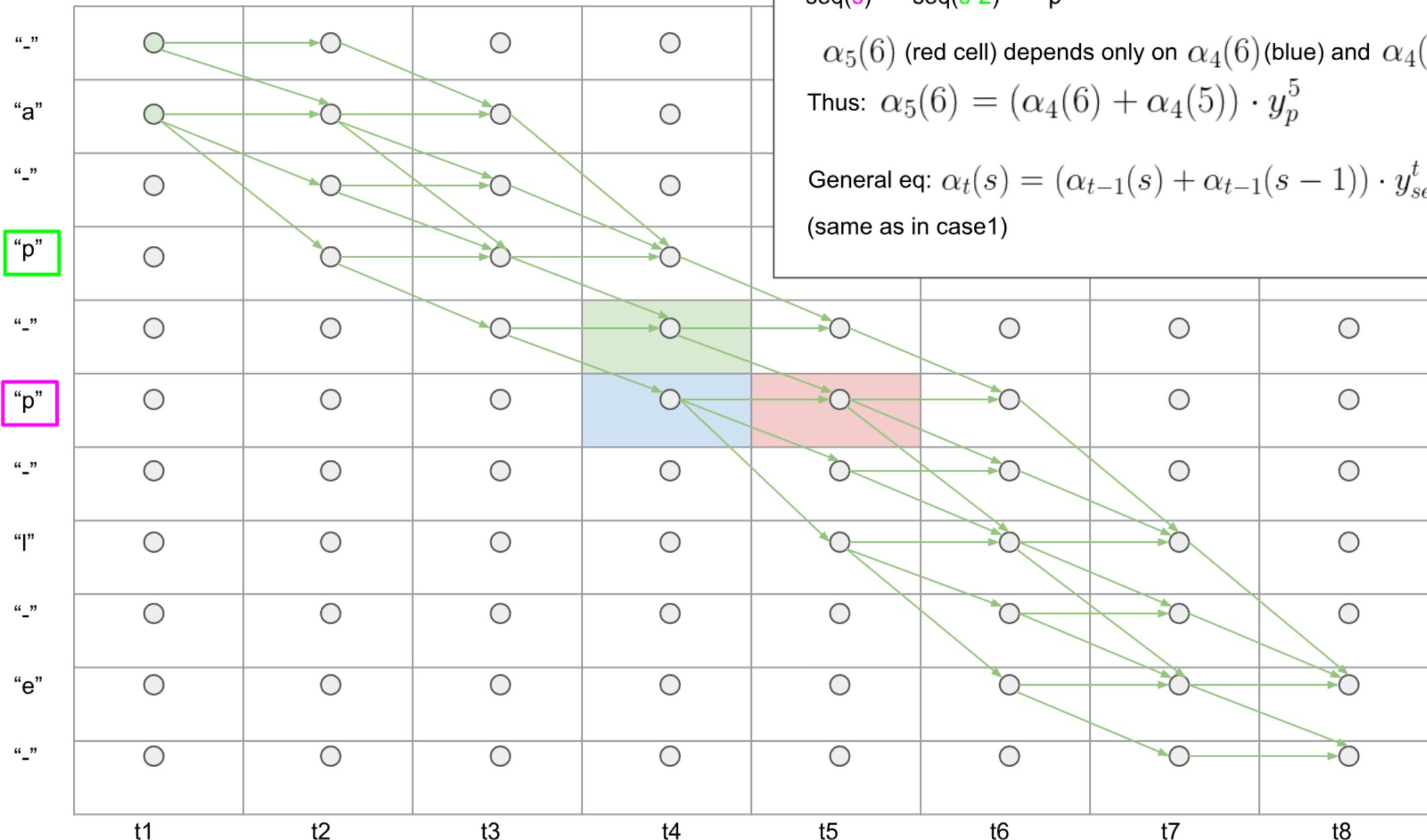


Case2. s-th symbol is equal to (s-2)-th symbol. Example: s=6,t=5.
 $\text{seq}(\text{s}) == \text{seq}(\text{s}-2) == "p"$

$\alpha_5(6)$ (red cell) depends only on $\alpha_4(6)$ (blue) and $\alpha_4(5)$ (green)

Thus: $\alpha_5(6) = (\alpha_4(6) + \alpha_4(5)) \cdot y_p^5$

General eq: $\alpha_t(s) = (\alpha_{t-1}(s) + \alpha_{t-1}(s - 1)) \cdot y_{\text{seq}(s)}^t$
 (same as in case1)

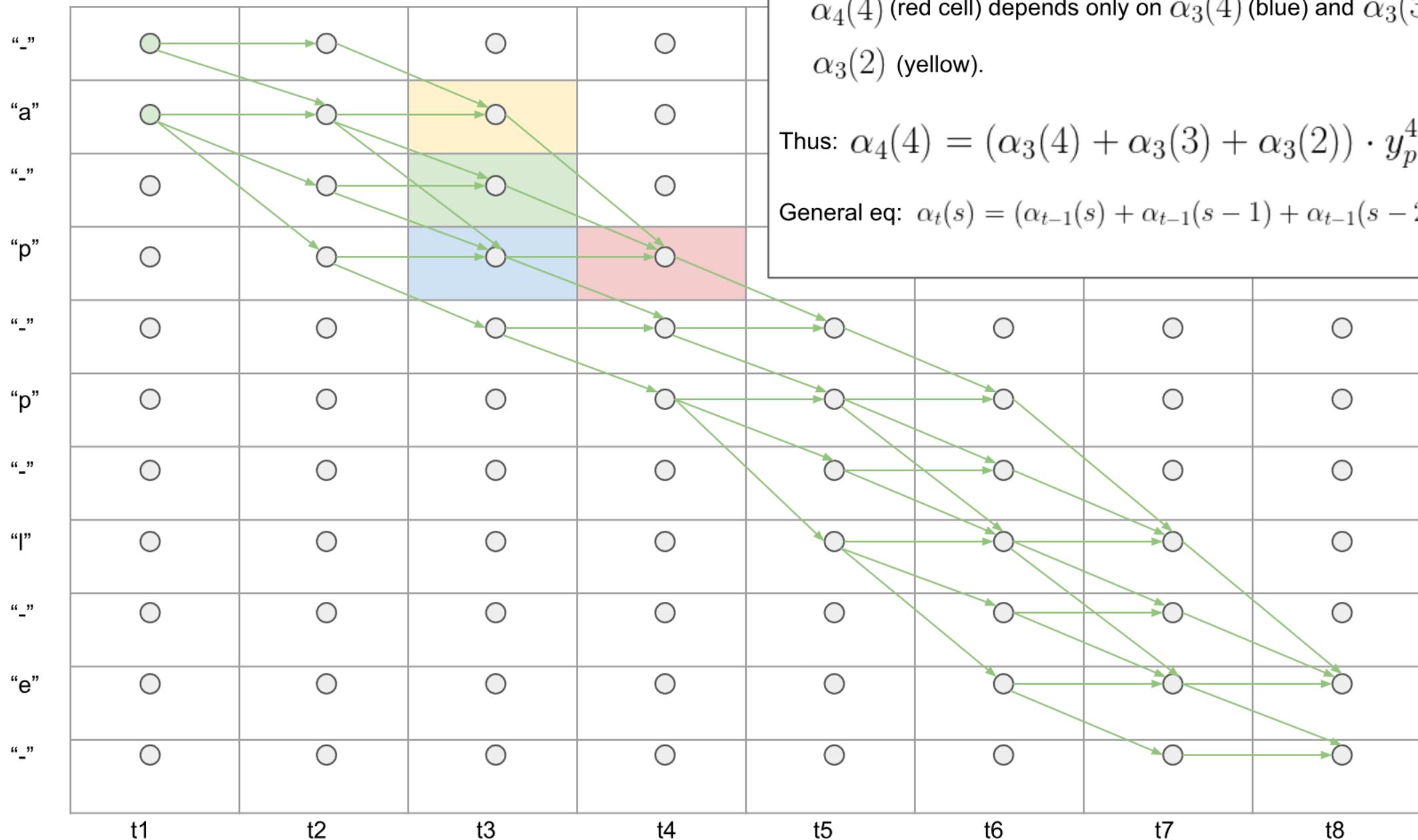


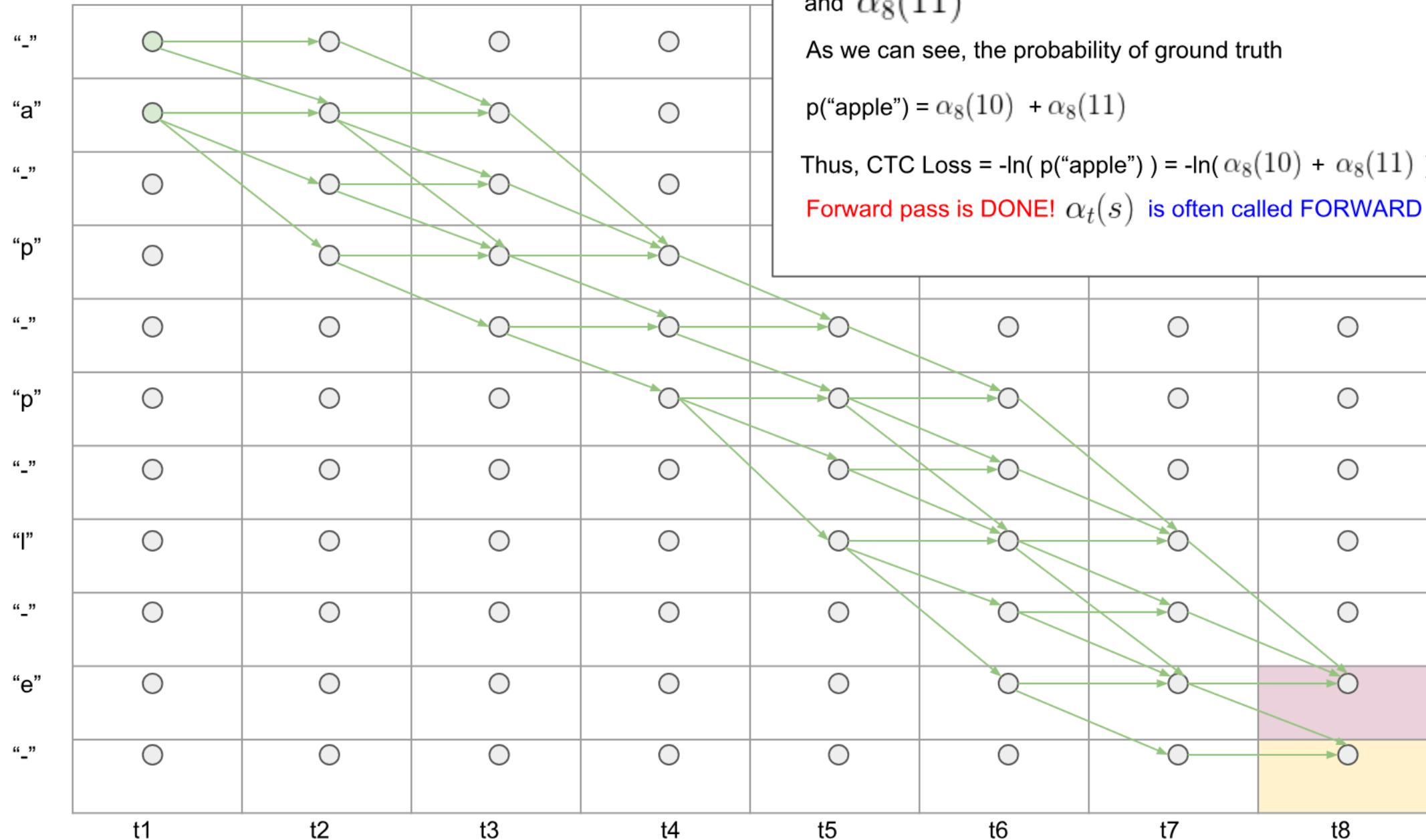
Case3. Otherwise. Example: $s=4, t=4$.

$\alpha_4(4)$ (red cell) depends only on $\alpha_3(4)$ (blue) and $\alpha_3(3)$ (green) and $\alpha_3(2)$ (yellow).

$$\text{Thus: } \alpha_4(4) = (\alpha_3(4) + \alpha_3(3) + \alpha_3(2)) \cdot y_p^4$$

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





With this dynamic programming algorithm we can calculate $\alpha_8(10)$ and $\alpha_8(11)$

As we can see, the probability of ground truth

$$p(\text{"apple"}) = \alpha_8(10) + \alpha_8(11)$$

$$\text{Thus, CTC Loss} = -\ln(p(\text{"apple})) = -\ln(\alpha_8(10) + \alpha_8(11))$$

Forward pass is DONE! $\alpha_t(s)$ is often called FORWARD VARIABLE.

总结：前向递推公式

This gives us the following rules for initialisation

$$\alpha_1(1) = y_b^1$$

$$\alpha_1(2) = y_{l_1}^1$$

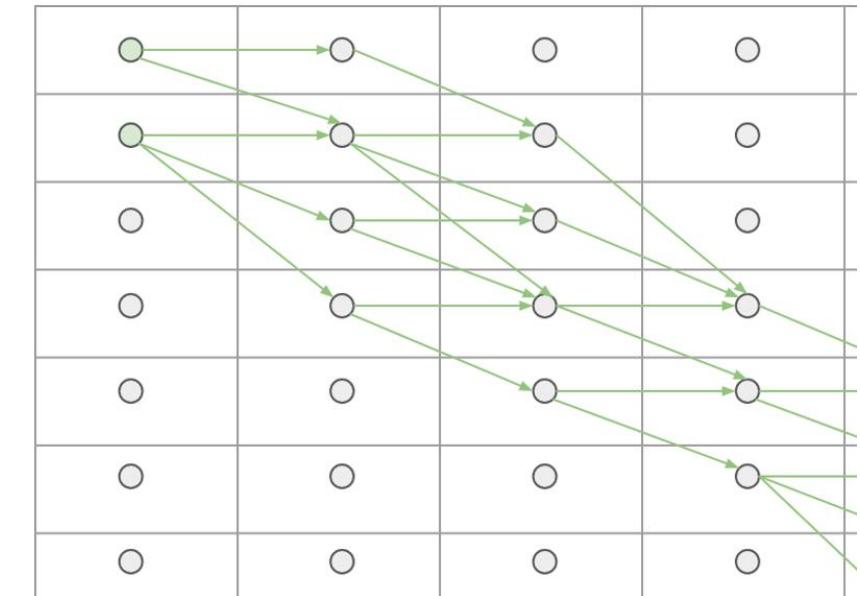
$$\alpha_1(s) = 0, \forall s > 2$$

and recursion

$$\alpha_t(s) = \begin{cases} \bar{\alpha}_t(s)y_{l'_s}^t & \text{if } l'_s = b \text{ or } l'_{s-2} = l'_s \\ (\bar{\alpha}_t(s) + \alpha_{t-1}(s-2))y_{l'_s}^t & \text{otherwise} \end{cases} \quad (6)$$

where

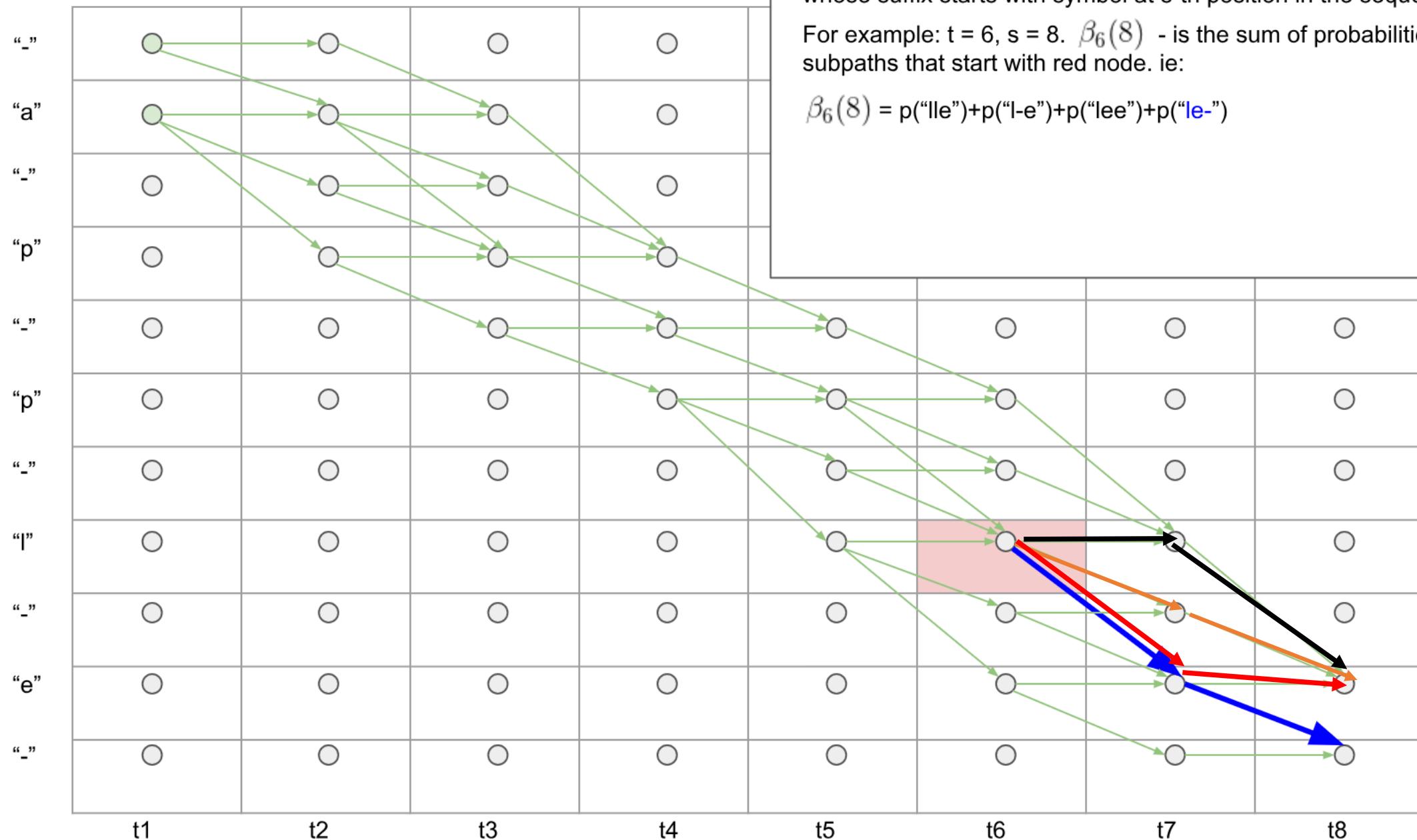
$$\bar{\alpha}_t(s) \stackrel{\text{def}}{=} \alpha_{t-1}(s) + \alpha_{t-1}(s-1). \quad (7)$$



Let me define variable $\beta_t(s)$ - the total probability of all subpaths, whose suffix starts with symbol at s-th position in the sequence at time t.

For example: $t = 6, s = 8$. $\beta_6(8)$ - is the sum of probabilities of all subpaths that start with red node. ie:

$$\beta_6(8) = p("lle") + p("l-e") + p("lee") + p("le-")$$



总结：后向递推公式

Similarly, the backward variables $\beta_t(s)$ are defined as the total probability of $\mathbf{l}_{s:|\mathbf{l}'|}$ at time t .

$$\beta_t(s) \stackrel{\text{def}}{=} \sum_{\substack{\pi \in N^T: \\ \mathcal{B}(\pi_{t:T}) = \mathbf{l}_{s:|\mathbf{l}'|}}} \prod_{t'=t}^T y_{\pi_{t'}}^{t'} \quad (9)$$

$$\beta_T(|\mathbf{l}'|) = y_b^T$$

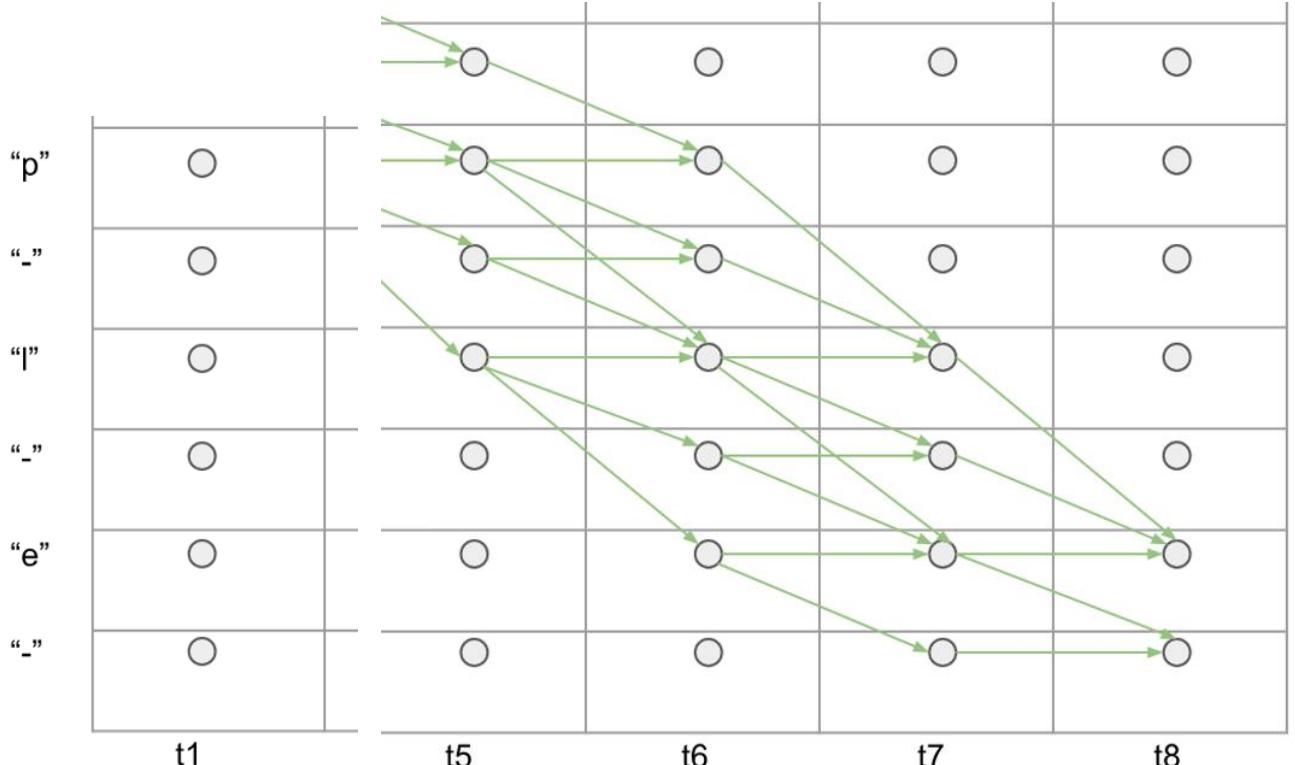
$$\beta_T(|\mathbf{l}'| - 1) = y_{\mathbf{l}'_{|\mathbf{l}'|}}^T$$

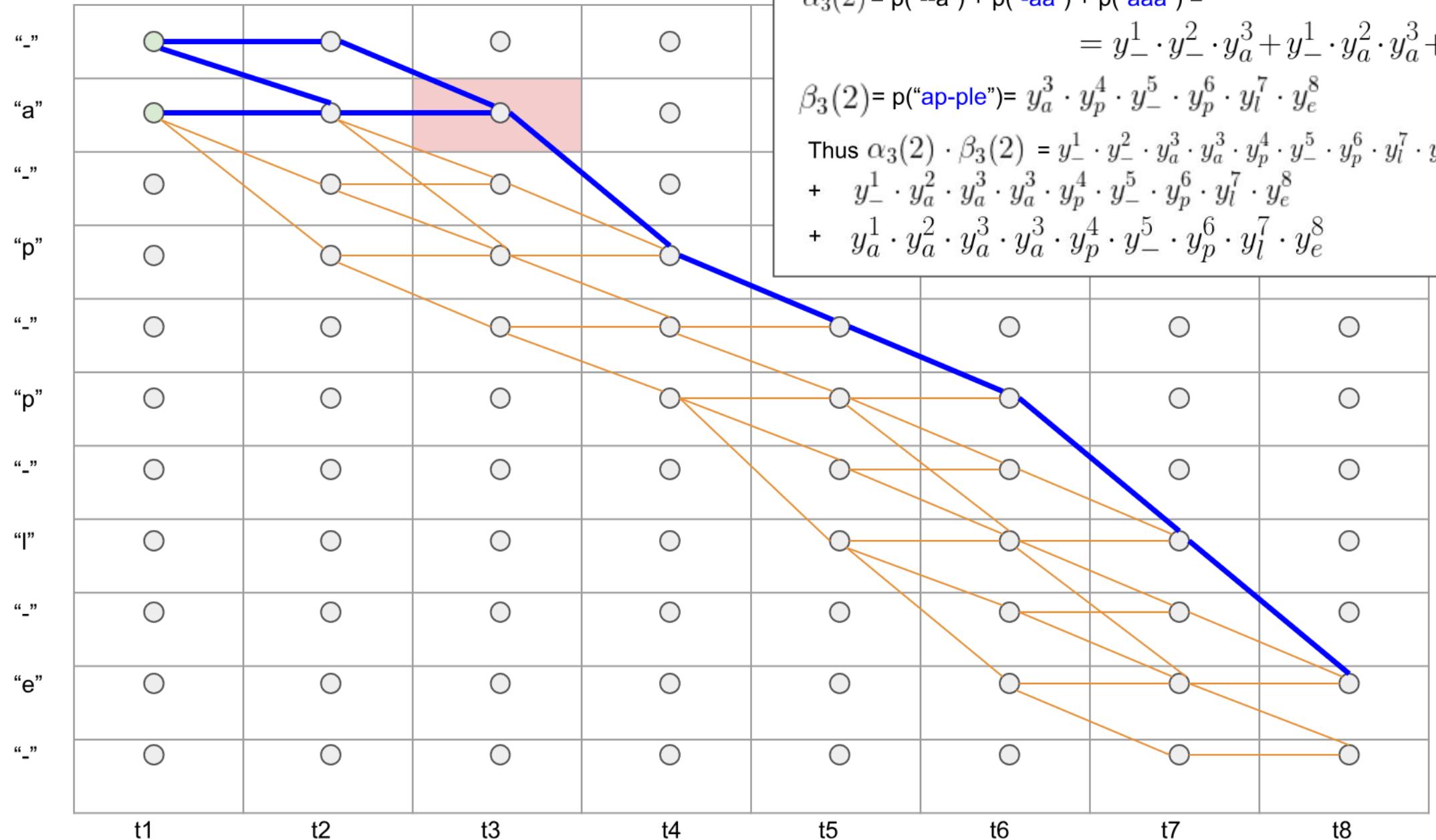
$$\beta_T(s) = 0, \forall s < |\mathbf{l}'| - 1$$

$$\beta_t(s) = \begin{cases} \bar{\beta}_t(s)y_{\mathbf{l}'_s}^t & \text{if } \mathbf{l}'_s = b \text{ or } \mathbf{l}'_{s+2} = \mathbf{l}'_s \\ (\bar{\beta}_t(s) + \beta_{t+1}(s+2))y_{\mathbf{l}'_s}^t & \text{otherwise} \end{cases} \quad (10)$$

where

$$\bar{\beta}_t(s) \stackrel{\text{def}}{=} \beta_{t+1}(s) + \beta_{t+1}(s+1). \quad (11)$$





Let's consider $t=3, s=2$.

$$\alpha_3(2) = p("aa") + p("-aa") + p("aaa") =$$

$$= y_-^1 \cdot y_-^2 \cdot y_a^3 + y_-^1 \cdot y_a^2 \cdot y_a^3 + y_a^1 \cdot y_a^2 \cdot y_a^3$$

$$\beta_3(2) = p("apple") = y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8$$

$$\text{Thus } \alpha_3(2) \cdot \beta_3(2) = y_-^1 \cdot y_-^2 \cdot y_a^3 \cdot y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8 +$$

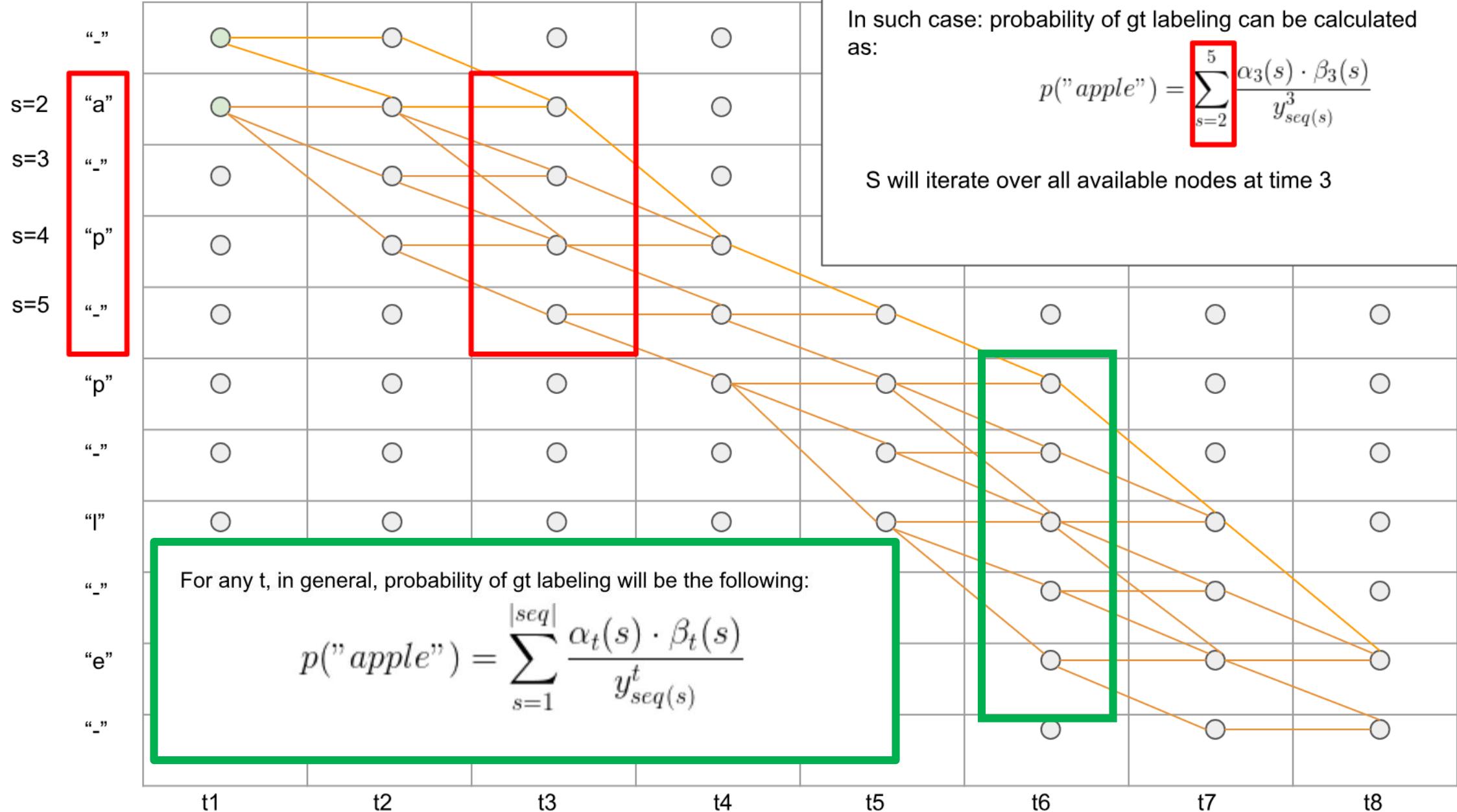
$$+ y_-^1 \cdot y_a^2 \cdot y_a^3 \cdot y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8$$

$$+ y_a^1 \cdot y_a^2 \cdot y_a^3 \cdot y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8$$

$$\begin{aligned}
 \text{Thus } \alpha_3(2) \cdot \beta_3(2) &= y_-^1 \cdot y_-^2 \cdot y_a^3 \cdot \boxed{y_a^3} \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8 + \\
 &+ y_-^1 \cdot y_a^2 \cdot y_a^3 \cdot \boxed{y_a^3} \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8 \\
 &+ y_a^1 \cdot y_a^2 \cdot y_a^3 \cdot \boxed{y_a^3} \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8 = \\
 &= (p("~-aap-ple") + p("-aap-ple") + p("aaap-ple")) * y_a^3
 \end{aligned}$$

$$\alpha_3(2) \cdot \beta_3(2) = (p("~-aap-ple") + p("-aap-ple") + p("aaap-ple")) * y_a^3$$

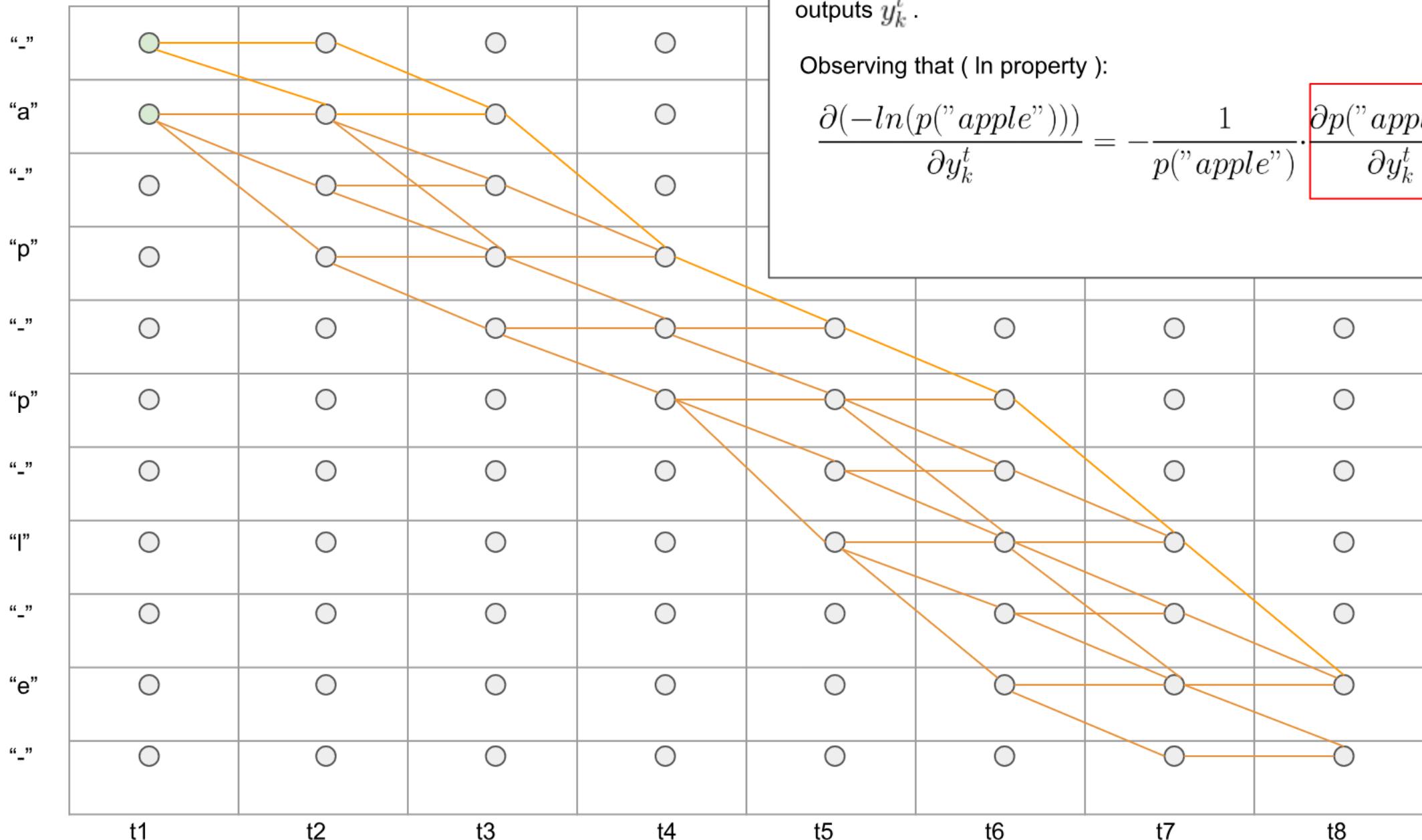
$\frac{\alpha_3(2) \cdot \beta_3(2)}{y_a^3}$ = sum of probabilities of all paths that go through second sequence symbol (symbol "a") at time 3



To do backprop, we should differentiate loss with respect to all network outputs y_k^t .

Observing that (In property):

$$\frac{\partial(-\ln(p("apple")))}{\partial y_k^t} = -\frac{1}{p("apple")}. \boxed{\frac{\partial p("apple")}{\partial y_k^t}}$$

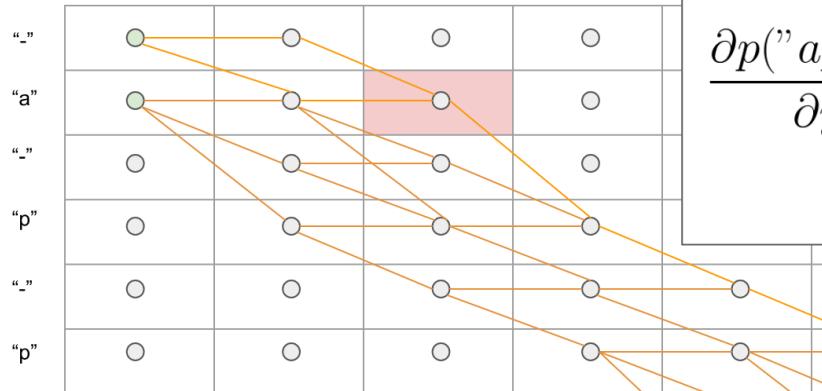


$$\frac{\partial p(\text{"apple"})}{\partial y_k^t}$$

To differentiate this with respect to y_k^t we need only consider those paths going through symbol k at time t.

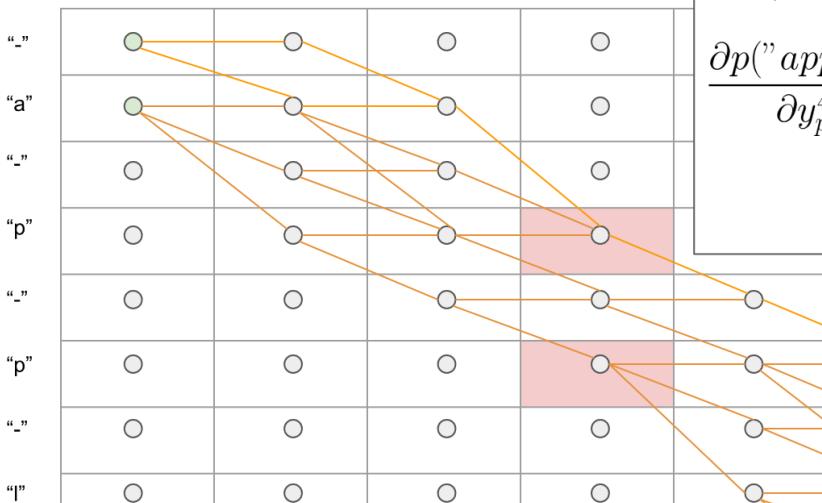
$$\frac{\partial p(\text{"apple"})}{\partial y_k^t} = -\frac{1}{y_k^{t2}} \cdot \sum_{s:seq(s)=k} \alpha_t(s) \cdot \beta_t(s)$$

此处应为“+”



Example1:

$$\frac{\partial p(\text{"apple"})}{\partial y_a^3} = -\frac{1}{y_a^{32}} \cdot \alpha_3(2) \cdot \beta_3(2)$$



Example2:

$$\frac{\partial p(\text{"apple"})}{\partial y_p^4} = -\frac{1}{y_p^{42}} \cdot (\alpha_4(4)\beta_4(4) + \alpha_4(6)\beta_4(6))$$

Conclusion

- Dynamic programming
- Matrix α (forward variables) is used to compute loss
- Matrix β (backward variables) is used to compute gradients

CTC解码

Best path decoding

按照上面思路，需要找到序列 I 的所有路径的概率，一种简化的方式是：找到路径中概率最大的，然后其对应的序列 I 就是最优序列，这个方法被称为「Best path decoding」。

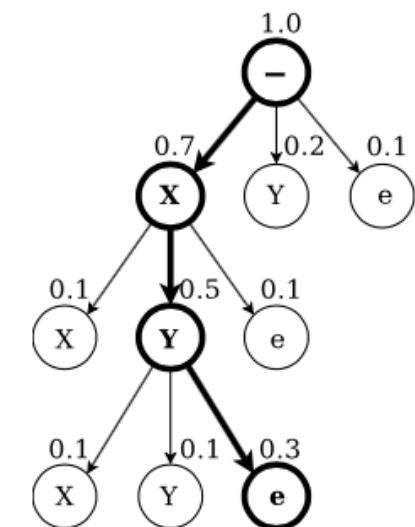
这个方法相当于是假设，最优序列的最优路径也是全局最优的（最优表示概率最大），形式化表示为：

$$h(\mathbf{x}) \approx \mathcal{B}(\pi^*)$$

where $\pi^* = \arg \max_{\pi \in N^t} p(\pi | \mathbf{x}).$

prefix search decoding

接着介绍第二种方法，「prefix search decoding」是一种剪枝的算法，剪枝主要在两个方面，一是同路径不重复计算，二是不可能状态不再搜索，下图中第一层的Y不搜索就是因为同层的X和下层的Y概率都比他高。



References

- [1] Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks
- [2] <https://blog.csdn.net/luodongri/article/details/77005948>
- [3] <https://distill.pub/2017/ctc/>
- [4] https://docs.google.com/presentation/d/12gYcPft9_4cxk2AD6Z6ZIJNa3wvZCW1ms31nhq51vMk/pub?start=false&loop=false&delayms=3000&slide=id.g24e9f0de4f_0_21796
- [5] <https://www.youtube.com/watch?v=eYIL4TMAeRI&t=549s>
- [6] <https://www.youtube.com/watch?v=c86gfVGcvh4>