

Connectionist Temporal Classification

Waseem Gharbieh



Outline

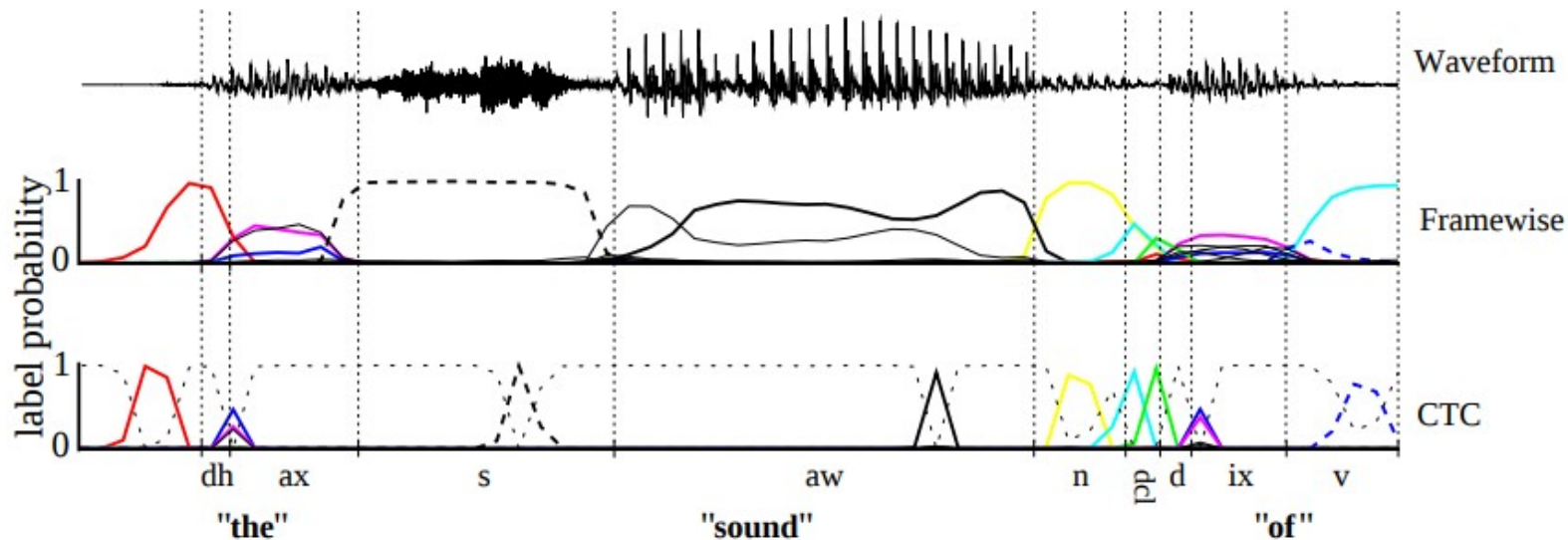
- Motivation
- Aligning Outputs And Labels
- CTC Loss
- Forward-Backward Algorithm
- Backprop
- Conclusions

Motivation

- Popular loss functions (such as MSE and XE) assume a one-to-one correspondence between the network's output and the target labels
- But what if there is no one-to-one correspondence? (as in speech recognition, or on-line handwriting recognition)
- Need a way to find the alignment between network outputs and target labels

Connectionist Temporal Classification (CTC)

- Given a sequence of N inputs and M labels, compute the loss between the N outputs and M labels

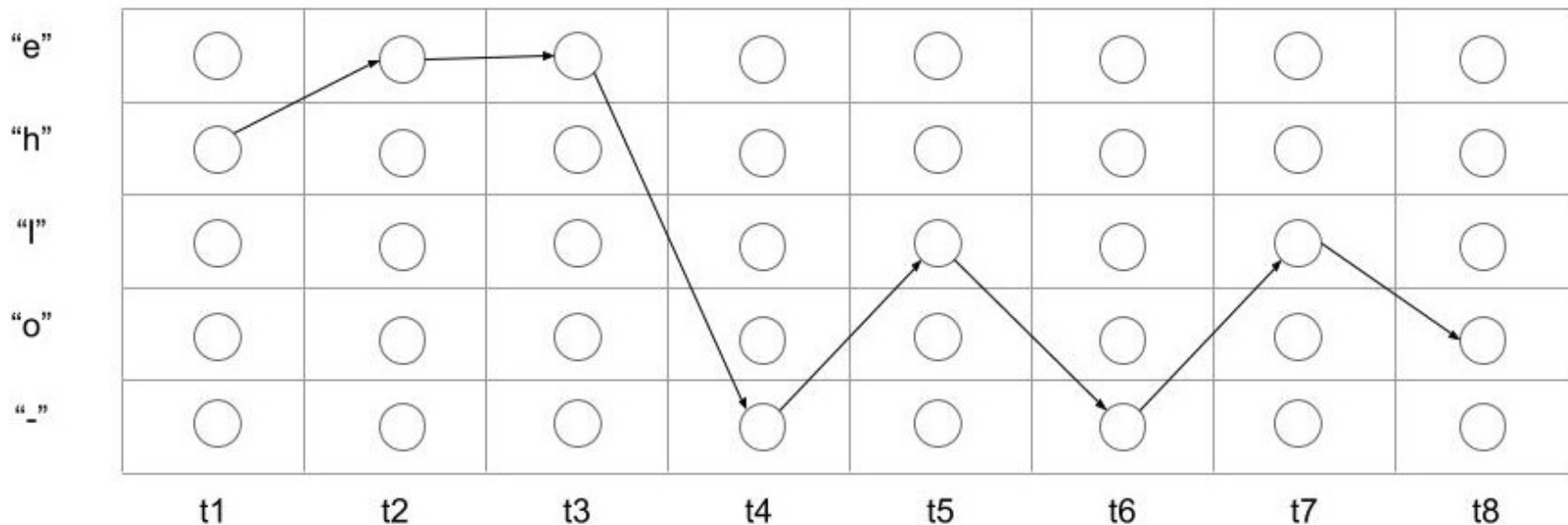


How To Align Outputs And Labels?

- Collapse repeated letters together
- Define a new token called “Blank” (will be represented as “-”)
- Assuming an input of length 8, we can define a function such that:
 - $f(\text{“hell-loo”}) = f(\text{“hel-lo”}) = \text{“hello”}$
 - $f(\text{“hellllloo”}) = \text{“helo”}$
 - $f(\text{“cc-a--tt”}) = f(\text{“c-a-t”}) = \text{“cat”}$
 - $f(\text{“-c-a-t--”}) = f(\text{“-c-a-t-”}) = \text{“cat”}$
 - $f(\text{“c-aaa-at”}) = f(\text{“c-a-at”}) = \text{“caat”}$

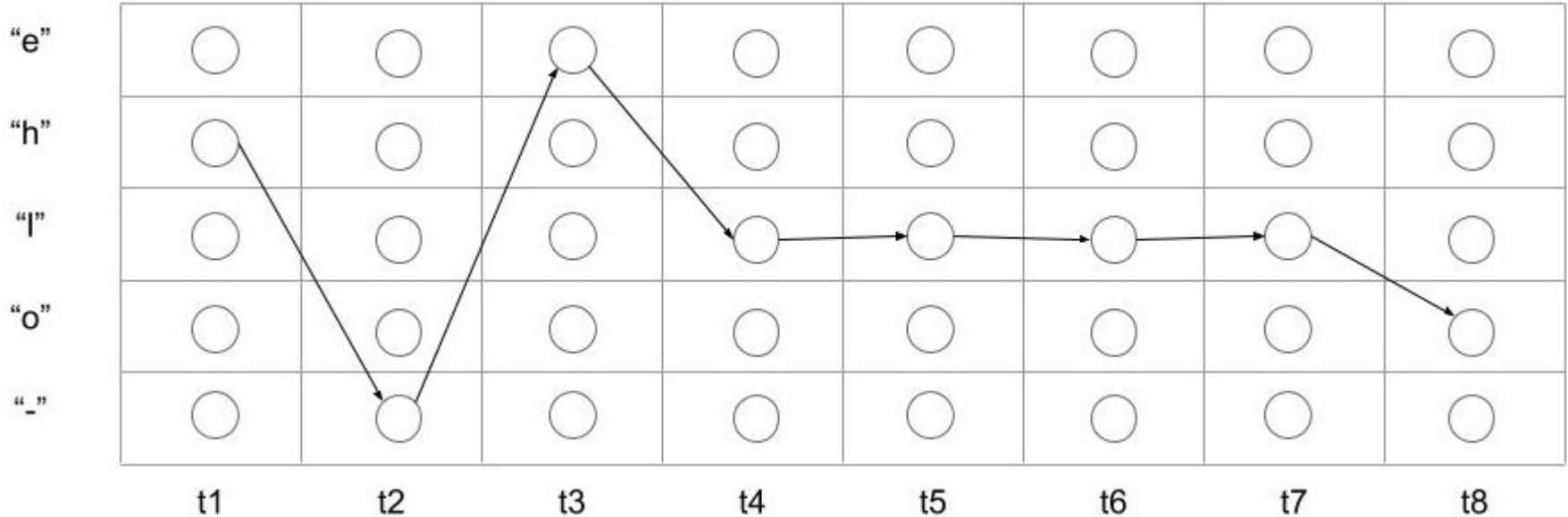
Some Examples

● $f(\text{"hee-l-lo"}) = \text{"hello"}$



Some Examples

● $f(\text{"h-ellllo"}) = \text{"helo"}$



There are 5^8 ways to go from t1 to t8, that's **390,625** possible paths!

CTC Loss

- CTC loss = $-\ln p(W)$
- So for the word “hello”, we would like to minimize $-\ln p(\text{“hello”})$
- Cannot be solved trivially, otherwise number of paths explode
- What can be done then? **Dynamic programming**

How Is This Done In Practice?

"a"	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
"h"	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
"u"	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
"e"	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
"_"	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
"l"	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
"_"	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
"l"	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
"_"	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
"o"	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
"_"	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	t1	t2	t3	t4	t5	t6	t7	t8

How Is This Done In Practice?

- $f(\text{"-hel-lo-"}) = \text{"hello"}$



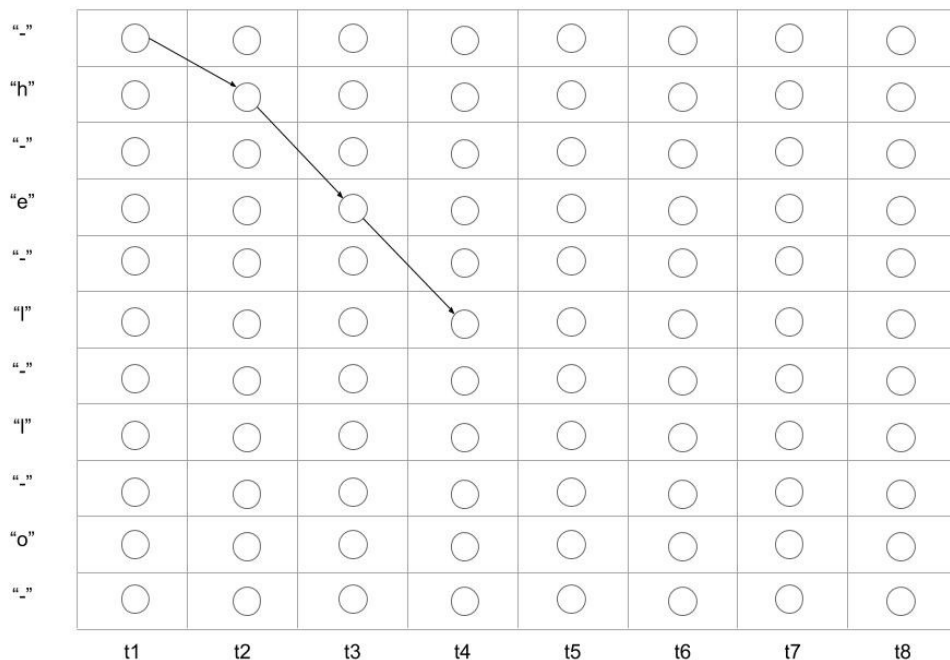
How Is This Done In Practice?

- $f(\text{"-hel-lo-"}) = \text{"hello"}$



How Is This Done In Practice?

- $f(\text{"-hel-lo-"}) = \text{"hello"}$



How Is This Done In Practice?

- $f(\text{"-hel-lo-"}) = \text{"hello"}$



How Is This Done In Practice?

- $f(\text{"-hel-lo-"}) = \text{"hello"}$



How Is This Done In Practice?

- $f(\text{"-hel-lo-"}) = \text{"hello"}$



How Is This Done In Practice?

- $f(\text{"-hel-lo-"}) = \text{"hello"}$



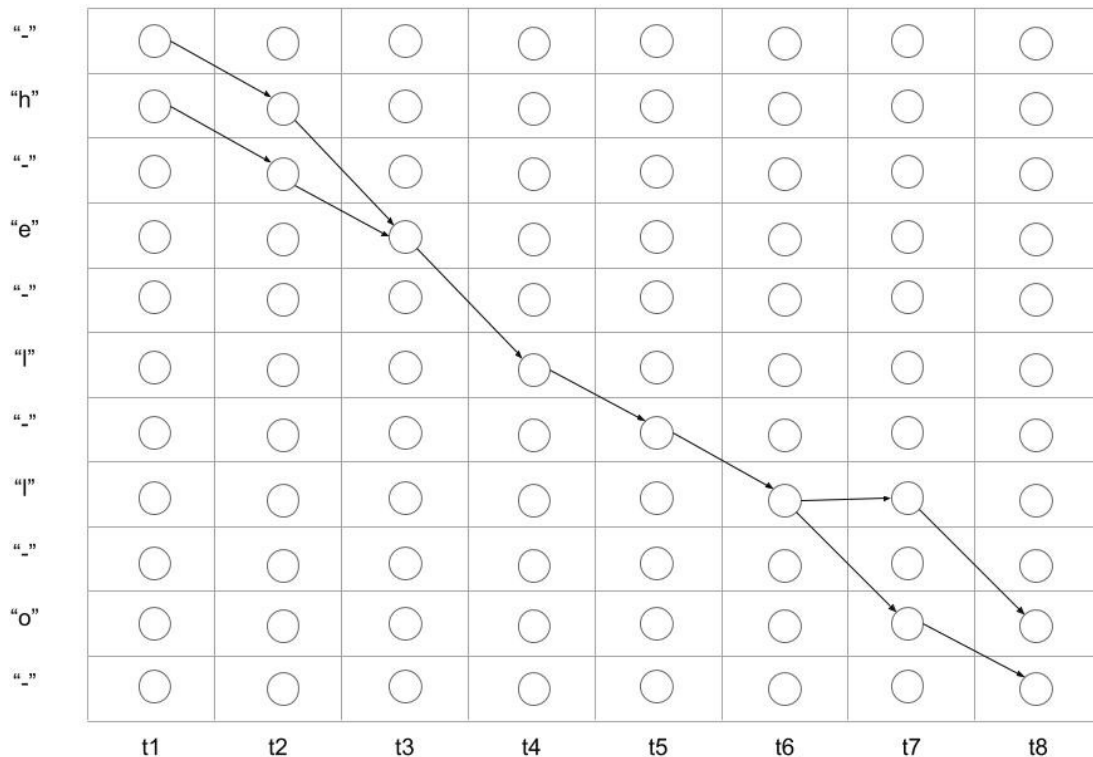
How Is This Done In Practice?

- $f(\text{"-hel-lo-"}) = \text{"hello"}$



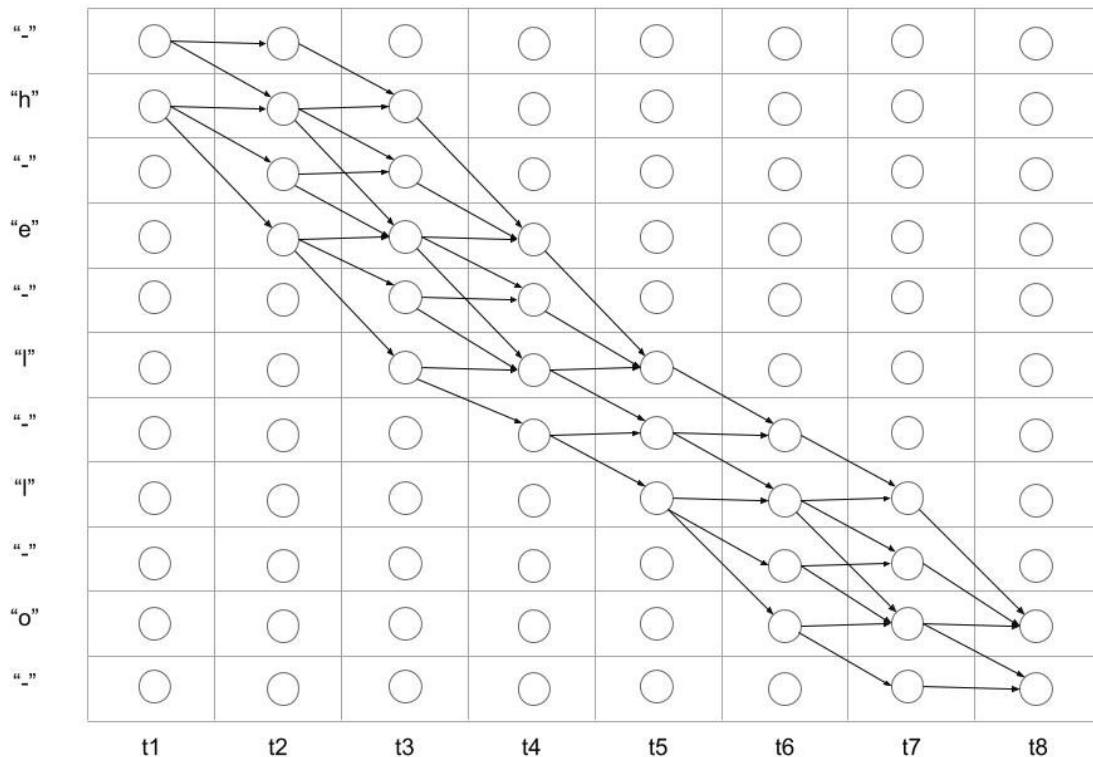
Alternative Paths

- $f(\text{"-hel-lo-"}) = \text{"hello"}$
- $f(\text{"-hel-llo"}) = \text{"hello"}$
- $f(\text{"h-el-llo"}) = \text{"hello"}$
- $f(\text{"h-el-lo-"}) = \text{"hello"}$



All Possible Paths

- Captures variations in pronunciation
- From f("--hel-lo")
- To f("hel-lo--")
- Note f("hel---lo")



CTC Loss Components

$$\alpha_t(s)$$

- Forward Variable: Calculates the total probability from the first timestep till timestep \mathbf{t} and token \mathbf{s}

$$\beta_t(s)$$

- Backward Variable: Calculates the total probability from timestep \mathbf{t} and token \mathbf{s} till last timestep

Forward Calculation Example 1

- Let's calculate $\alpha_3(4)$

- There are 4 paths

$$p(" - he") = y_{-}^1 \cdot y_h^2 \cdot y_e^3$$

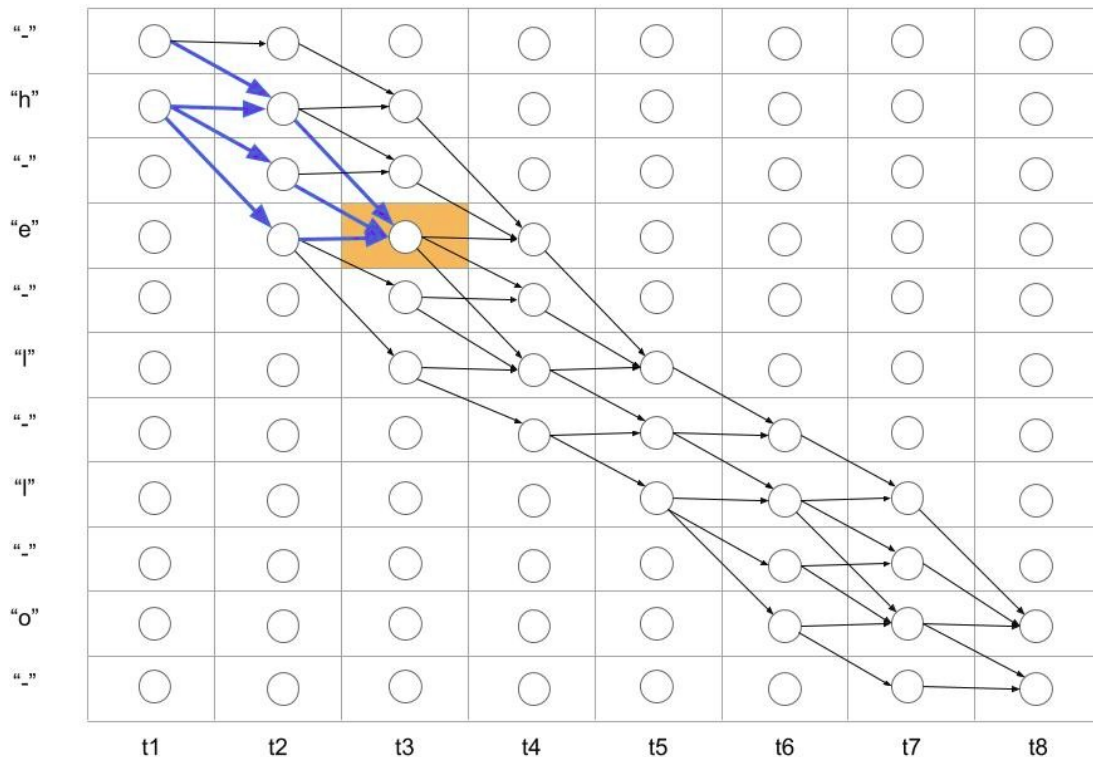
$$p("hhe") = y_h^1 \cdot y_h^2 \cdot y_e^3$$

$$p("h - e") = y_h^1 \cdot y_{-}^2 \cdot y_e^3$$

$$p("hee") = y_h^1 \cdot y_e^2 \cdot y_e^3$$

- Final probability is

$$p(" - he") + p("hhe") + p("h - e") + p("hee")$$



Forward Calculation Example 1

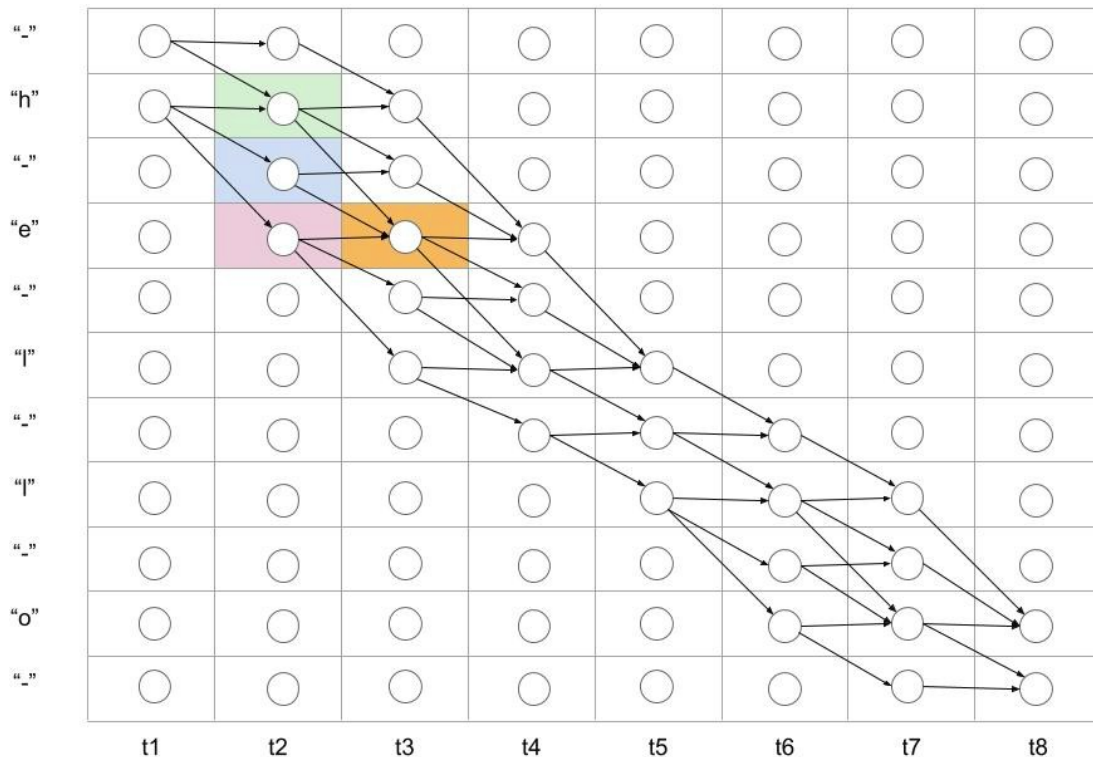
$$p("-he") + p("hhe") + p("h-e") + p("hee")$$

- This expression can also be calculated using dynamic programming

$$\alpha_3(4) = (\alpha_2(4) + \alpha_2(3) + \alpha_2(2)) \cdot y_e^3$$

- Or in general

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



Forward Calculation Example 2

- The top component can be calculated using dynamic

programming

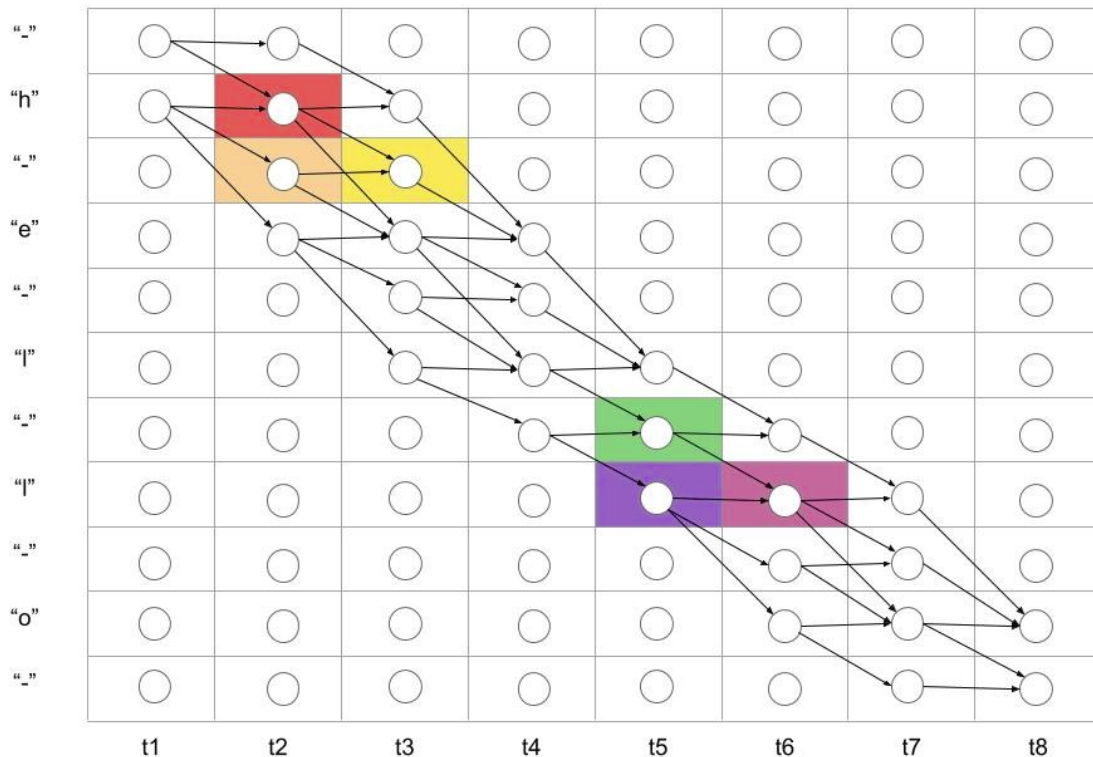
$$\alpha_3(3) = (\alpha_2(3) + \alpha_2(2)) \cdot y_3^3$$

- Or in general (for both)

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

- Note that

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

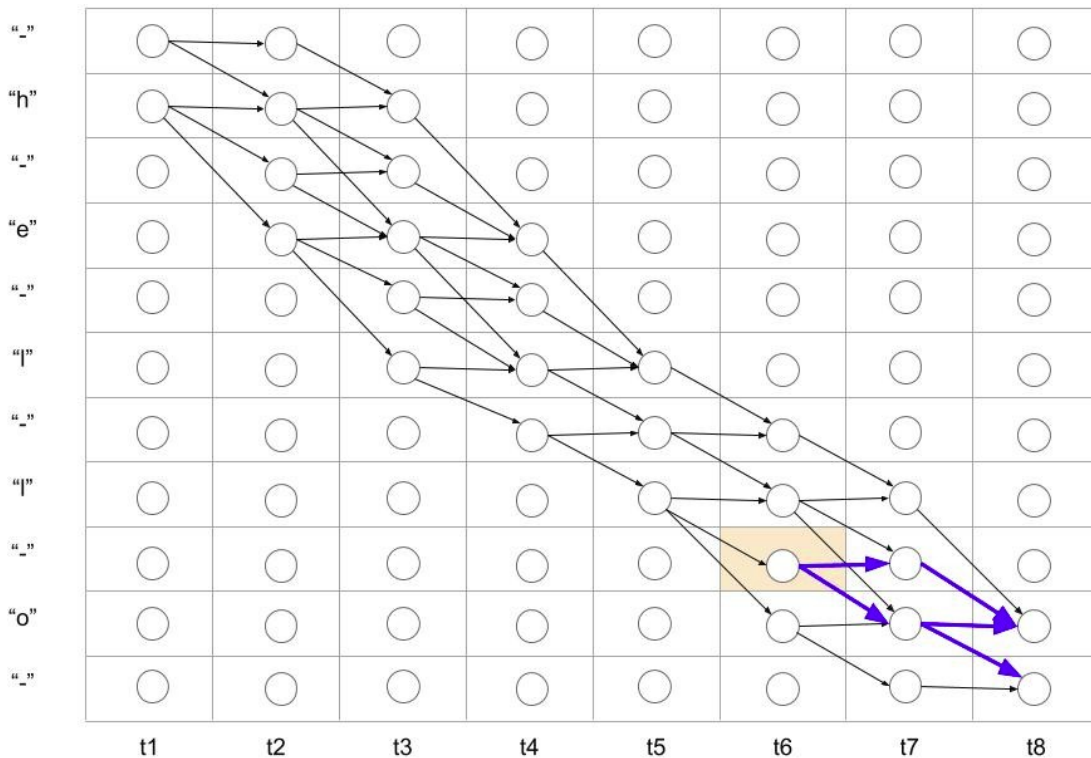


Backward Calculation Example

- This phase is just the opposite of the forward phase

$$\beta_6(9) = p(" -- o") + p("- oo") + p("- o -")$$

- The same logic here applies as the forward calculation



Complete Path Calculation

- The forward pass gives

$$\alpha_3(2) = p(\text{"--h"}) + p(\text{"-hh"}) + p(\text{"hhh"})$$

$$\alpha_3(2) = y_{-}^1 y_{-}^2 y_h^3 + y_{-}^1 y_h^2 y_h^3 + y_h^1 y_h^2 y_h^3$$

The backward pass results in

$$\beta_3(2) = p(\text{"hel-lo"}) = y_h^3 y_e^4 y_l^5 y_o^6 y_l^7 y_o^8$$

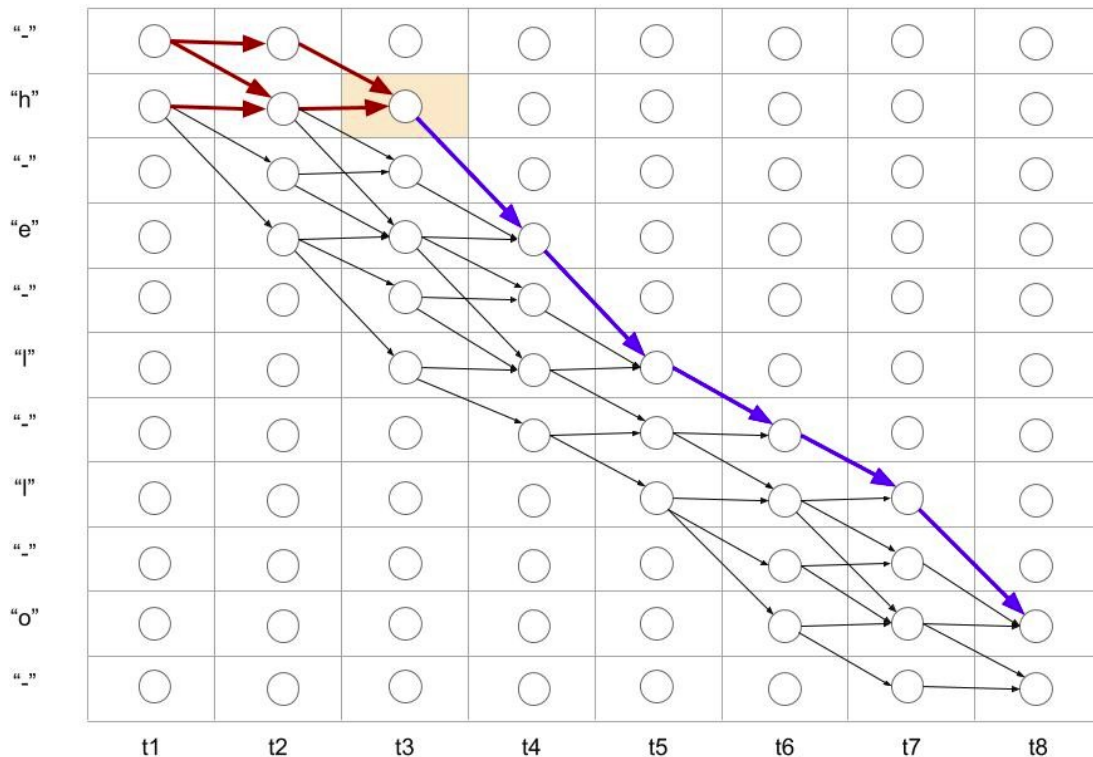
The final result

$$\begin{aligned} \alpha_3(2) \cdot \beta_3(2) &= y_{-}^1 \cdot y_{-}^2 \cdot y_h^3 \cdot y_h^3 \cdot y_e^4 \cdot y_l^5 \cdot y_o^6 \cdot y_l^7 \cdot y_o^8 \\ &\quad + y_{-}^1 \cdot y_h^2 \cdot y_h^3 \cdot y_h^3 \cdot y_e^4 \cdot y_l^5 \cdot y_o^6 \cdot y_l^7 \cdot y_o^8 \\ &\quad + y_h^1 \cdot y_h^2 \cdot y_h^3 \cdot y_h^3 \cdot y_e^4 \cdot y_l^5 \cdot y_o^6 \cdot y_l^7 \cdot y_o^8 \end{aligned}$$

$$= [p(\text{"--hel-lo"}) + p(\text{"-hhel-lo"}) + p(\text{"hhhel-lo"})] \cdot y_h^3$$

going through **n** at $t=3$

$$\frac{\alpha_3(2) \cdot \beta_3(2)}{y_h^3}$$



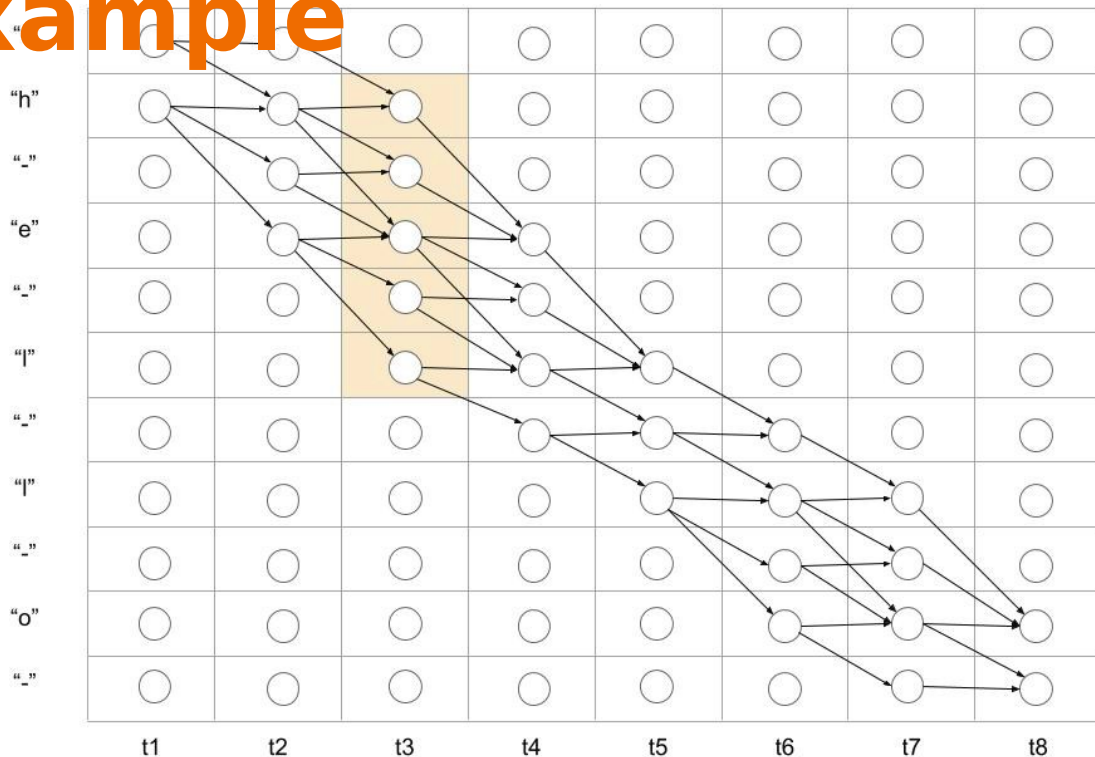
Calculating Loss At A Particular Timestep Example

- $p(\text{"hello"})$ at $t=3$ is the sum of the probabilities of all paths through all the symbols

$$p(\text{"hello"}) = \sum_{s=2}^6 \frac{\alpha_3(s) \cdot \beta_3(s)}{y_{seq(s)}^3}$$

- In general, the probability of the ground truth label is

$$p(\text{"hello"}) = \sum_{s=1}^{|seq|} \frac{\alpha_t(s) \cdot \beta_t(s)}{y'_{seq(s)}}$$



How To Do Backprop?

$$\frac{\partial(-\ln p(\text{"hello"}))}{\partial y_k^t} = \frac{-1}{p(\text{"hello"})} \frac{\partial p(\text{"hello"})}{\partial y_k^t}$$

- Since $p(\text{"hello"}) = \sum_{s=1}^{|\text{seq}|} \frac{\alpha_t(s) \cdot \beta_t(s)}{y_{\text{seq}(s)}^t}$

$$\frac{\partial p(\text{"hello"})}{\partial y_k^t} = \frac{-1}{y_k^{t^2}} \sum_{s:\text{seq}(s)=k} \alpha_t(s) \cdot \beta_t(s)$$

Backprop Example 1

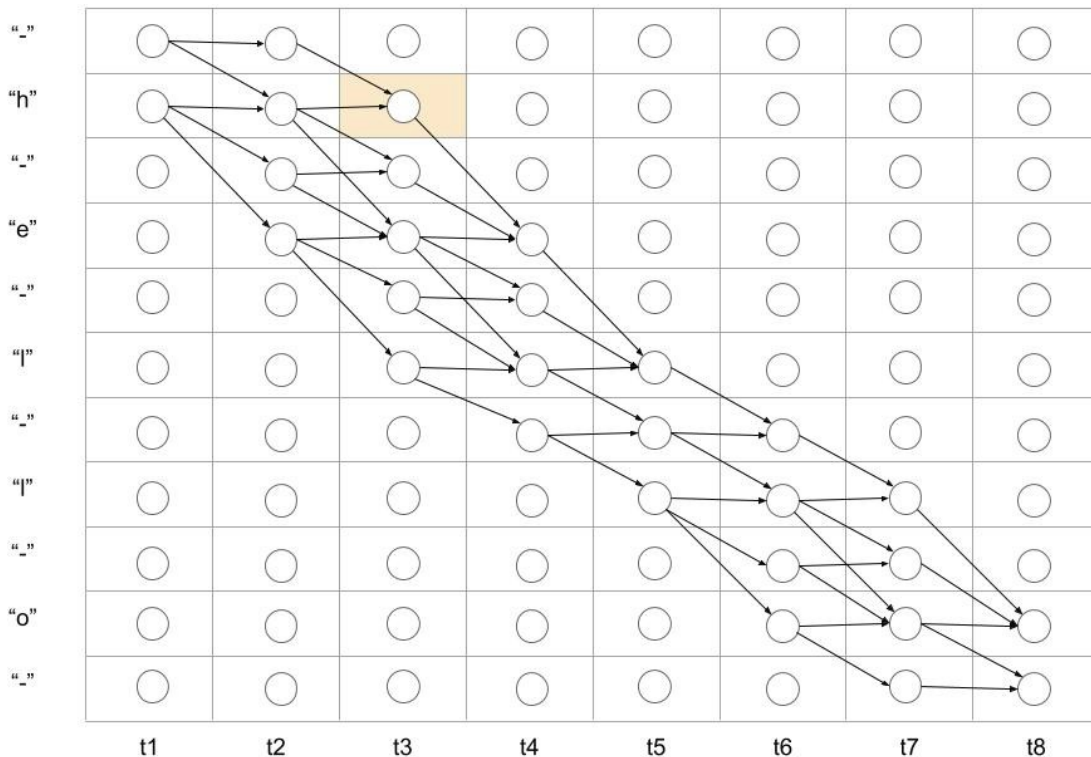
- Given

$$\frac{\partial p(\text{"hello"})}{\partial y_k^t} = \frac{-1}{y_k^t{}^2} \sum_{s: \text{seq}(s)=k} \alpha_t(s) \cdot \beta_t(s)$$

- For $t=3$, and $k=\mathbf{h}$

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

Since **h** occurs only at $s=2$



Backprop Example 2

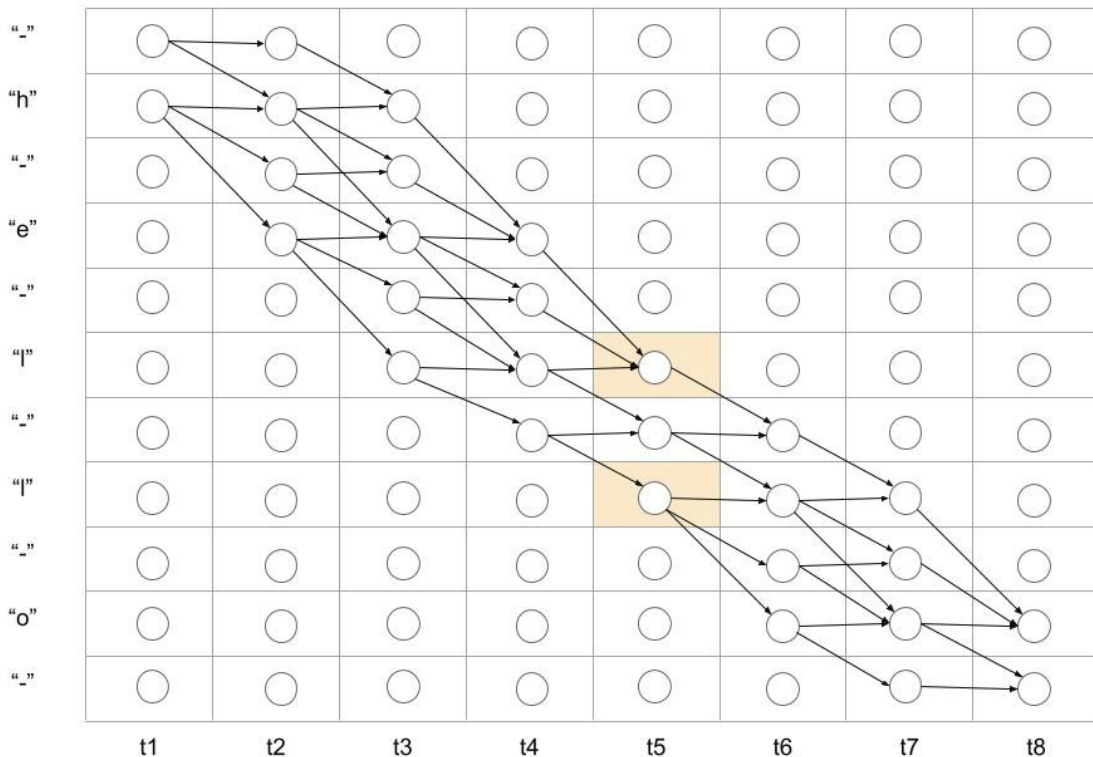
- Given

$$\frac{\partial p(\text{"hello"})}{\partial y_k^t} = \frac{-1}{y_k^{t^2}} \sum_{s: \text{seq}(s)=k} \alpha_t(s) \cdot \beta_t(s)$$

- For $t=5$, and $k=\text{I}$

$$\frac{\partial p(\text{"hello"})}{\partial y_l^5} = \frac{-1}{y_l^{5^2}} \cdot (\alpha_5(6) \cdot \beta_5(6) + \alpha_5(8) \cdot \beta_5(8))$$

Since **I** occurs at $s=6$ and $s=8$



Conclusions

- CTC Loss allows training models on sequences whose number of inputs is different than the number of labels
- It makes use of dynamic programming to calculate path probabilities efficiently
- CTC treats every timestep independently

Th-aa-nk -Yo-uu--!

References

1. Graves, Alex, et al. "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks." Proceedings of the 23rd international conference on Machine learning. ACM, 2006.
2. https://en.wikipedia.org/wiki/Connectionist_temporal_classification
3. <https://www.youtube.com/watch?v=c86gfVGcvh4>
4. <https://www.youtube.com/watch?v=eYIL4TMAeRI>