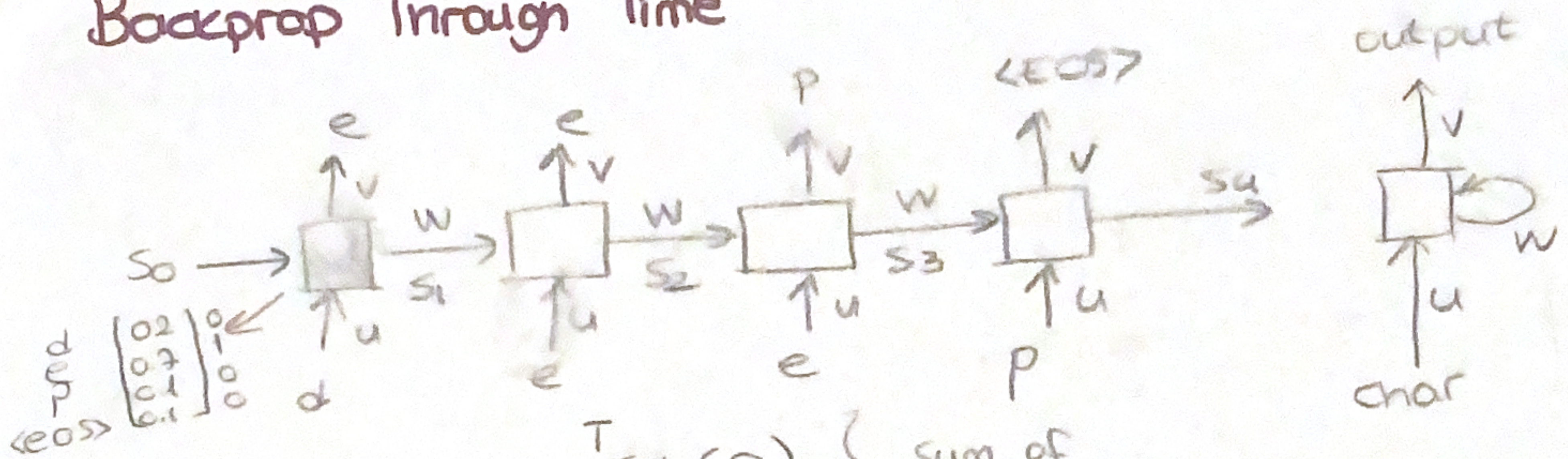


Backprop Through Time



Total Loss $\rightarrow \sum_{t=1}^T d_t(\theta)$ } Sum of loss over all steps

$d_t(\theta) = -\log(y_{t,c})$
 (prob produced for true class at time step i (0.7 in above example, log will decrease if it gets closer to 1))

For backprop we compute gradients for w, u, v

For v : $\frac{\partial d(\theta)}{\partial v} = \frac{\partial d_1(\theta)}{\partial v} + \frac{\partial d_2(\theta)}{\partial v} + \dots$ } standard backprop

\downarrow
 this is not being affected by states
 v is a matrix

s_4 depends on s_3 & w
 s_3 depends on s_2 & w

For w : $\frac{\partial d(\theta)}{\partial w}$

$\frac{\partial d_4(\theta)}{\partial w} = \frac{\partial d_4(\theta)}{\partial s_4} \cdot \frac{\partial s_4}{\partial w}$

we can't calculate this
 $s_4 = G(w, s_3, b)$

Total derivative not a pain.

Explicit $\rightarrow \frac{\partial^* s_4}{\partial w}$ } treats all inputs as constant

Implicit \rightarrow Sum all internal paths from s_4 to w .

$\frac{\partial s_4}{\partial w} = \frac{\partial^* s_4}{\partial w} + \frac{\partial s_4}{\partial s_3} \cdot \frac{\partial s_3}{\partial w} = \frac{\partial^* s_4}{\partial w} + \frac{\partial s_4}{\partial s_3} \left[\frac{\partial^* s_3}{\partial w} + \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial w} \right]$

\dots

$\frac{\partial s_2}{\partial s_1} \cdot \frac{\partial s_1}{\partial w}$