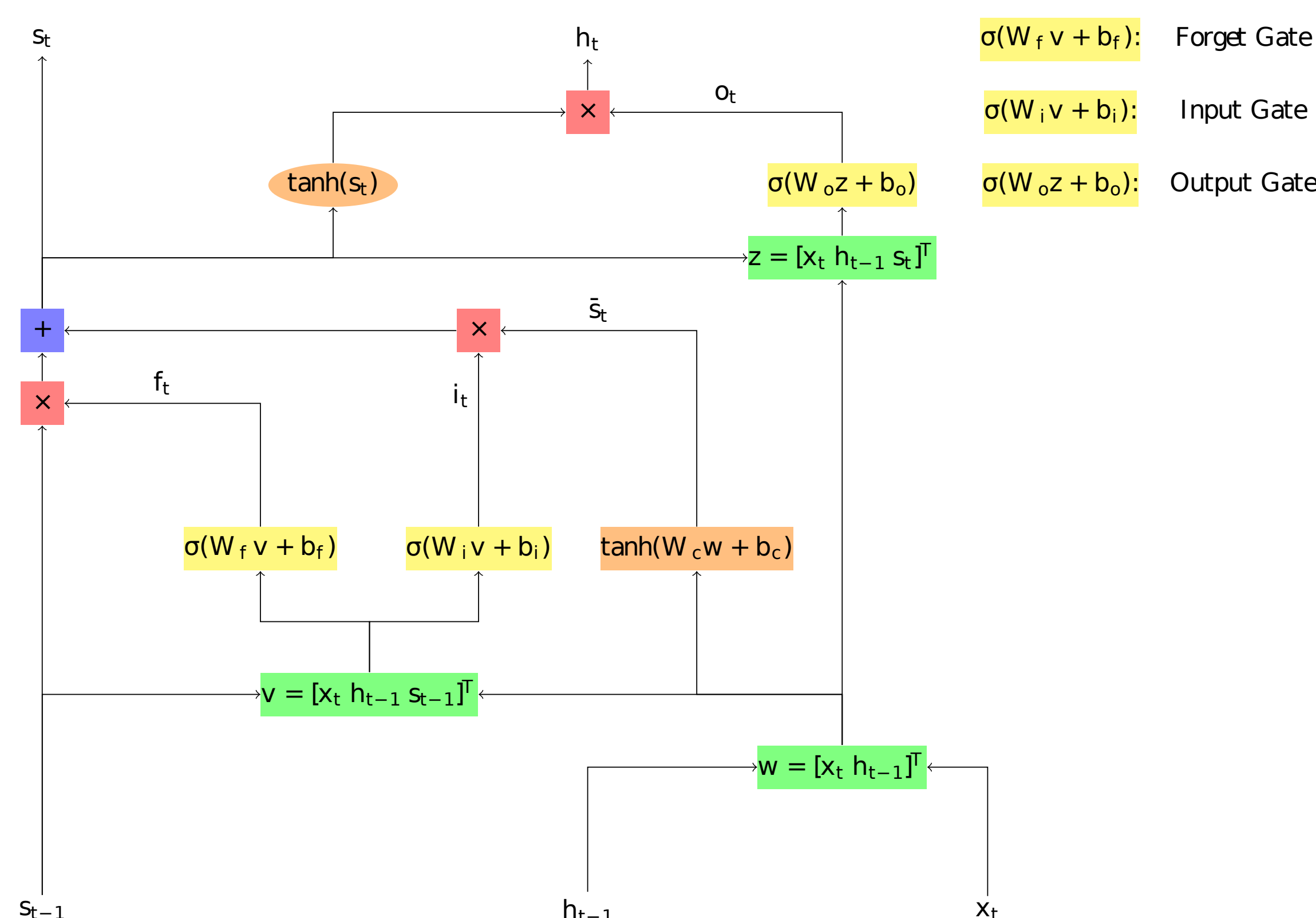


1. Goals

- Transcribe speech to text.
- Use an end to end system.
- Find suitable parameters.
- Explore regularization options.

2. Long Short Term Memory



- Differentiable versions of computer memory chips.
- Gates i_t , f_t , o_t , state s_t and output h_t .

$$i_t = \sigma(W_i[x_t \parallel h_{t-1} \parallel s_{t-1}]^T + b_i), \quad (1)$$

$$f_t = \sigma(W_f[x_t \parallel h_{t-1} \parallel s_{t-1}]^T + b_f), \quad (2)$$

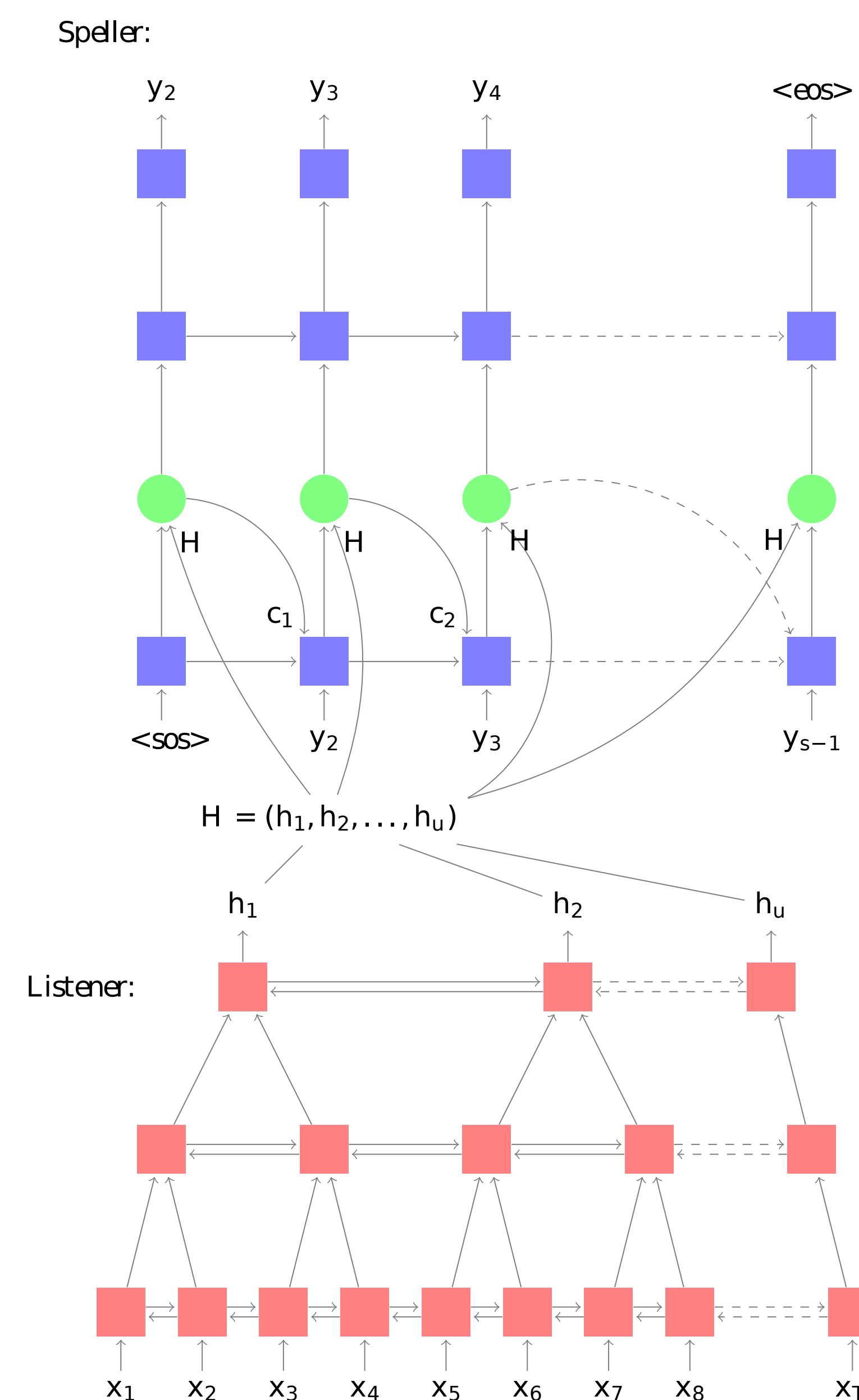
$$s_t = f_t s_{t-1} + i_t \tanh(W_s[x_t \parallel h_{t-1}]^T + b_s), \quad (3)$$

$$o_t = \sigma(W_o[x_t \parallel h_{t-1} \parallel s_t]^T + b_o), \quad (4)$$

$$h_t = o_t \tanh(s_t). \quad (5)$$

- Most important equation is the one for the state s_t , the state functions as cell memory.

3. Listen Attend and Spell



- Mel feature inputs.
- Listener consists of one bidirectional and one or more pyramidal bidirectional LSTM layer.
 - pyramidal layers concatenate features over the feature dimension in time:

$$h_t^n = \text{BLSTM}(h_{t-1}^n, [h_{2t}^{n-1}, h_{2t+1}^{n-1}]). \quad (6)$$

- Speller consists of an attend and spell cell.
 - The attend and spell function evaluates:

$$s_i = \text{RNN}(s_{i-1}, y_{i-1}, c_{i-1}), \quad (7)$$

$$c_i = \text{AttentionContext}(s_i, \mathbf{H}), \quad (8)$$

$$P(y_i | \mathbf{x}, y_{<i}) = \text{CharacterDistribution}(s_i, c_i). \quad (9)$$

- The attention context function computes:

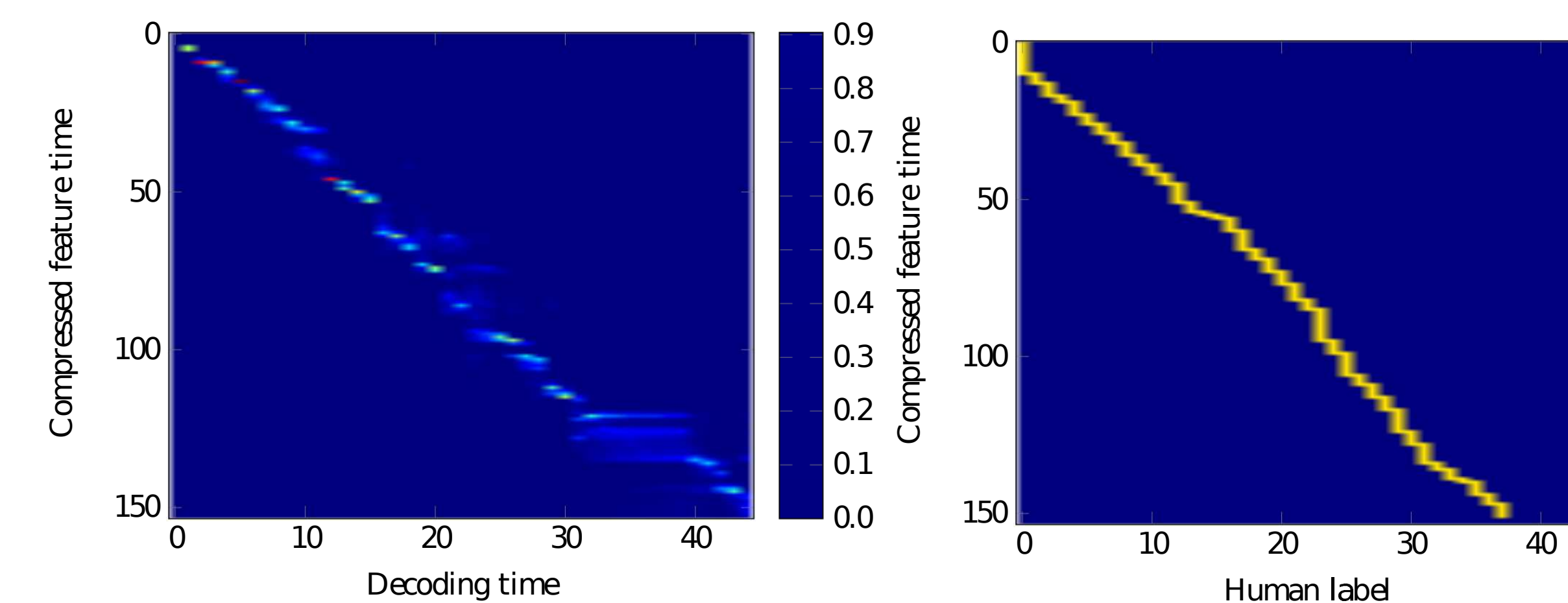
$$e_{i,u} = \phi(s_i)^T \psi(h_u), \quad (10)$$

$$\alpha_{i,u} = \frac{\exp(e_{i,u})}{\sum_u \exp(e_{i,u})}, \quad (11)$$

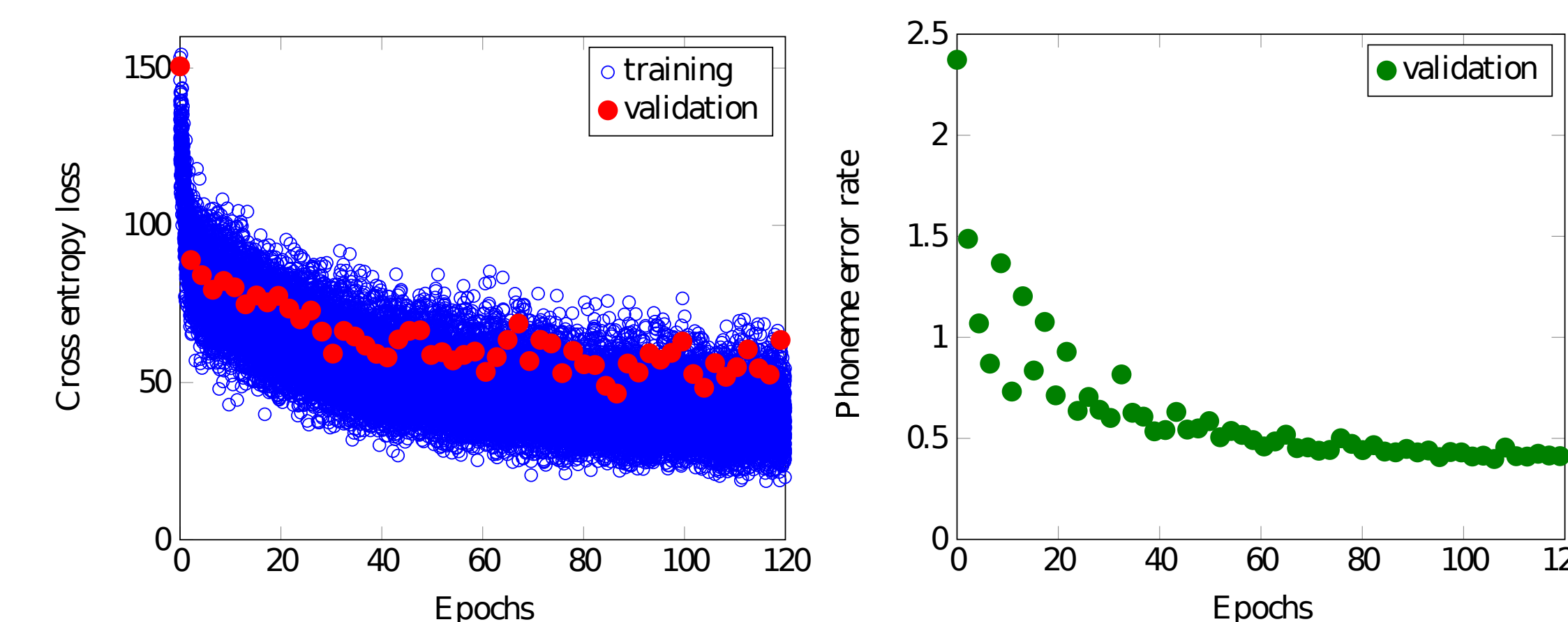
$$c_i = \sum_u \alpha_{i,u} h_u. \quad (12)$$

- ϕ and ψ are MLPs.

4. Results



- Network alignment and human alignment resemble each other.
- Result obtained using greedy decoding test set error rate 0.55.



- Improved results were found using a larger network, beam search and dropout, test set error rate 0.45.
- Utterance fm1d0_sx295:
 - Human labeling:
 - $\langle \text{sos} \rangle$ sil ih f sil k eh r l sil k ah m z sil t ah m aa r ah hh
 - ae v er rey n jh f er m iy dx iy ng ih sil t uw sil $\langle \text{eos} \rangle$
 - Network labeling:
 - $\langle \text{sos} \rangle$ sil hh ih f sil k ih r ow sil k ah m sil sil t ah m aa aa hh
 - hh v v er ey n n sil f f er m iy iy iy sil sil t uw sil sil $\langle \text{eos} \rangle$
 - Error ratio 0.36 for this utterance.

5. Literature

- [1] W. Chan, N. Jaitly, Q.V. Le, and O. Vinyals, "Listen attend and Spell." *arXiv preprint*, pages, 2015.
- [2] A. Graves. "Supervised Sequence Labelling with Recurrent Neural Networks." *Springer, Berlin Heidelberg*, 2012.