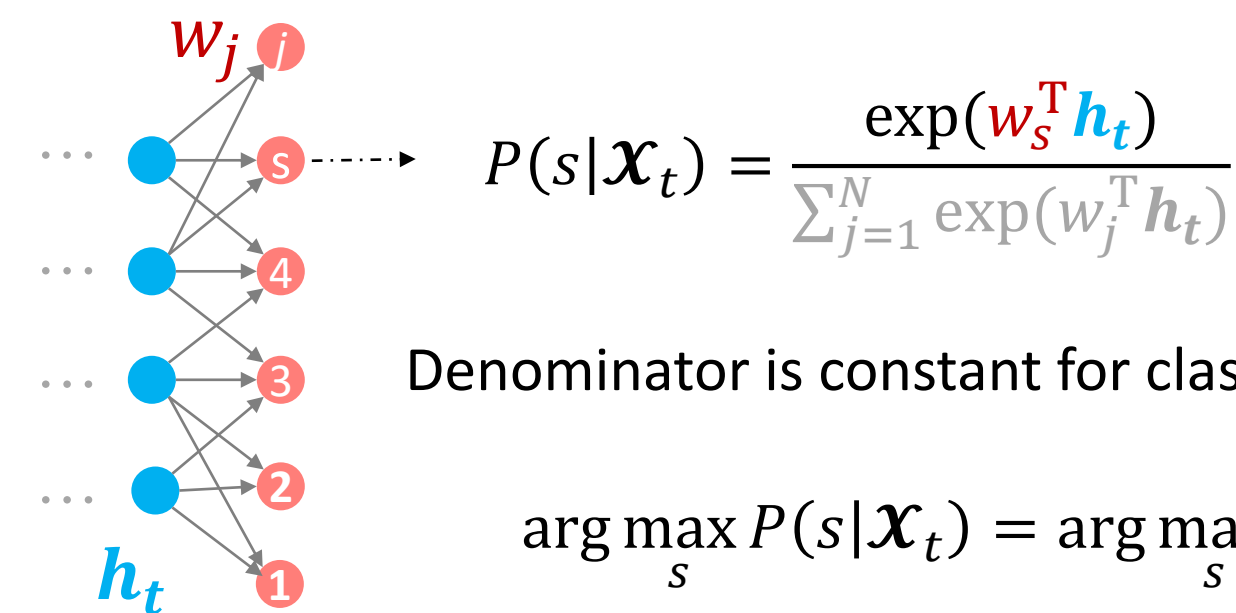


Shi-Xiong Austin Zhang, Rui Zhao, Chaojun Liu, Jinyu Li and Yifan Gong  
Microsoft, Redmond, WA, USA

## 1. Introduction

- RNNs use the softmax activation function in the last layer



$$P(s|\mathcal{X}_t) = \frac{\exp(w_s^T h_t)}{\sum_{j=1}^N \exp(w_j^T h_t)}$$

Denominator is constant for classification

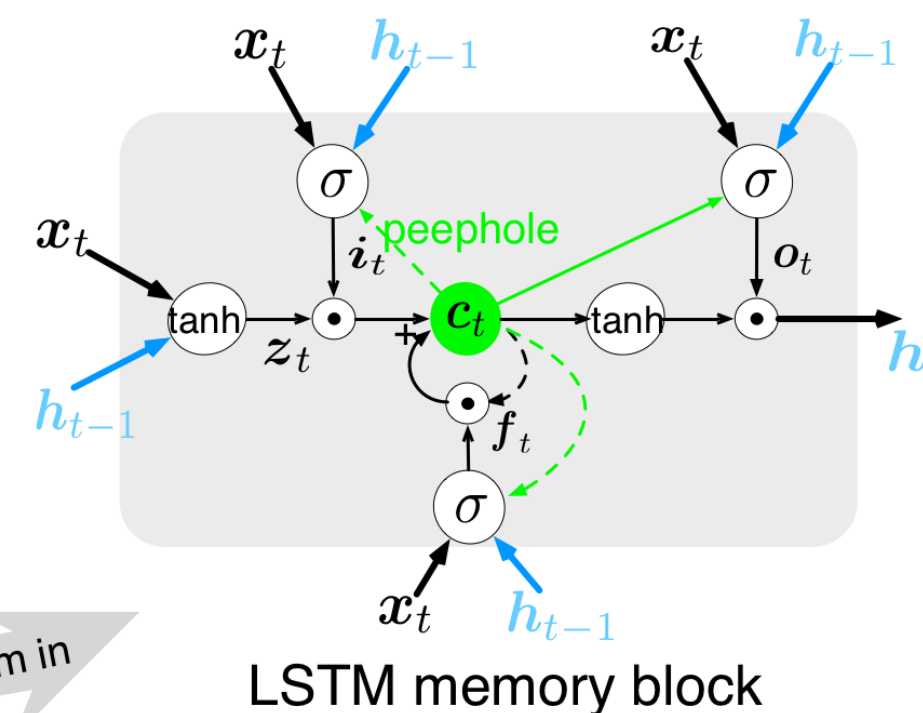
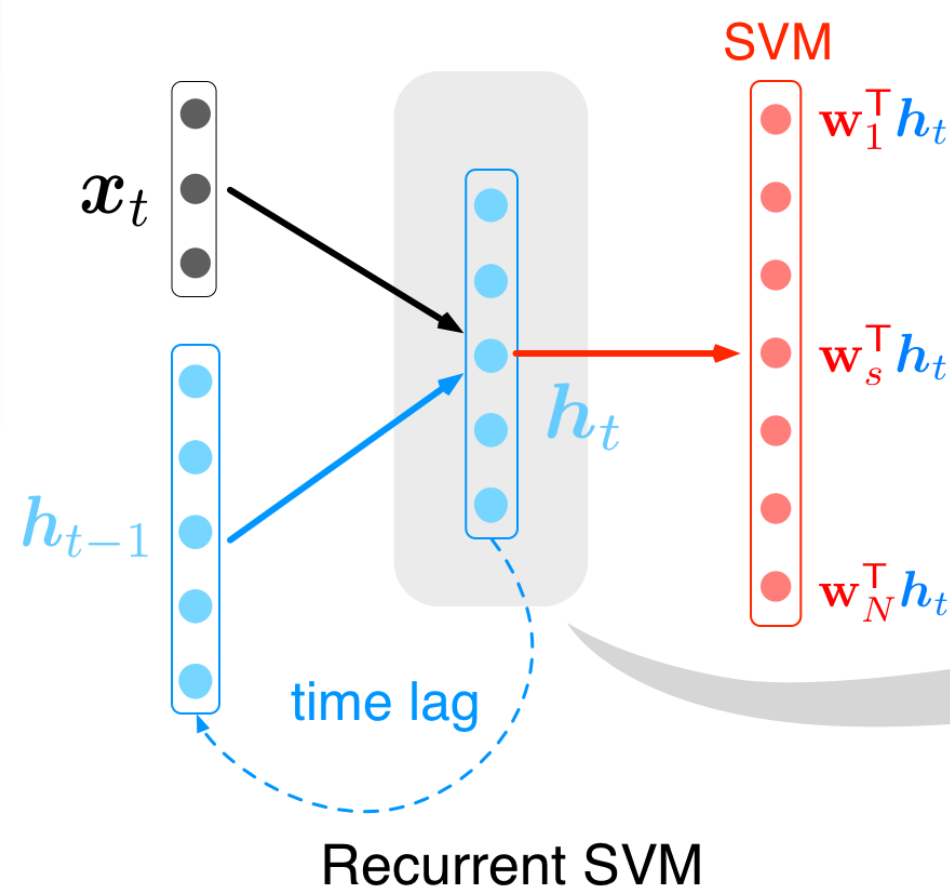
$$\arg \max_s P(s|\mathcal{X}_t) = \arg \max_s w_s^T h_t$$

- Multiclass SVM, same classification:  $\arg \max_s w_s^T \phi(\mathcal{X}_t)$

Replace softmax layer in RNN with SVM → Recurrent SVM

## 2. Recurrent SVM

Architecture: LSTM+SVM



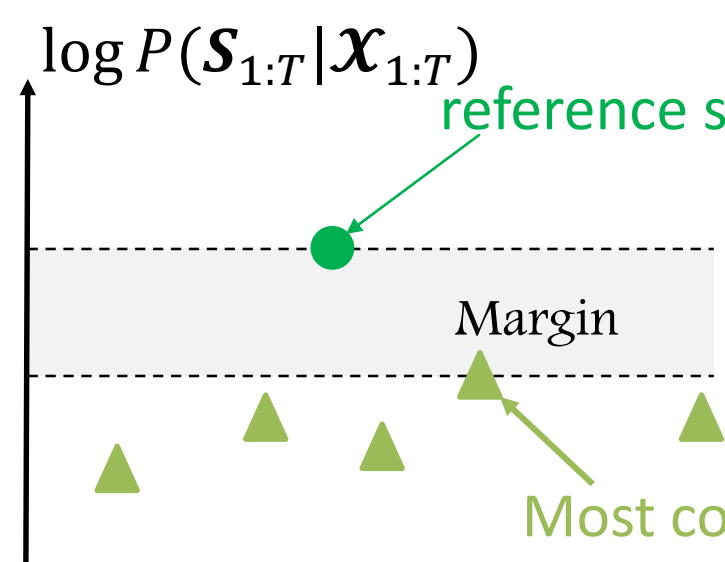
How to train them jointly?

**1<sup>st</sup>-step:** fixed LSTM, training SVM using the quadratic programming.

**2<sup>nd</sup>-step:** fixed SVM, training LSTM using the subgradient methods

## 3. Max-Margin Sequence Training

What is the Margin? Objective Function?



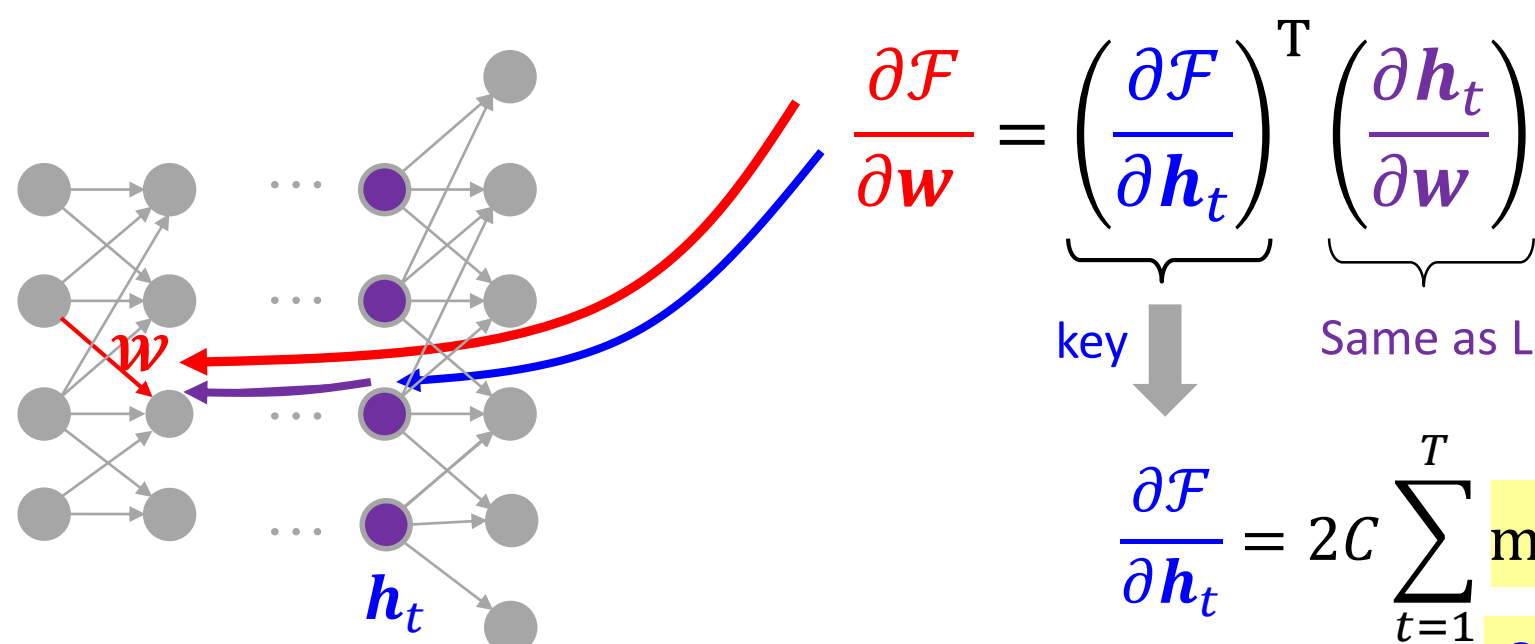
$$\min_{\bar{S}} \left\{ \log \frac{P(S|\mathcal{X})}{P(\bar{S}|\mathcal{X})} \right\} = \mathcal{F}(w^{SVM}, w^{LSTM})$$

How to train the last layer  $w^{SVM}$ ?

**1<sup>st</sup>-step:** fixed  $w^{LSTM}$ ,  $\mathcal{F}(w^{SVM}, w^{LSTM})$  is convex, training  $w_{SVM}$  ⇔ the structured SVM

How to train the previous layers  $w^{LSTM}$ ?

**2<sup>nd</sup>-step:** fixed  $w^{SVM}$ ,  $\mathcal{F}(w^{SVM}, w^{LSTM})$  is non-differentiable, training  $w^{LSTM}$  requires subgradients



$$\frac{\partial \mathcal{F}}{\partial w} = \left( \frac{\partial \mathcal{F}}{\partial h_t} \right)^T \left( \frac{\partial h_t}{\partial w} \right)$$

key ↓ Same as LSTM, Same BPTT ☺

$$\frac{\partial \mathcal{F}}{\partial h_t} = 2C \sum_{t=1}^T \max\{\text{Score}_{\bar{S}} - \text{Score}_S, 0\} (w_{\bar{S}_t}^{SVM} - w_{S_t}^{SVM})$$

Only the support vectors have gradients! ☺

## 4. Experiments & Conclusion

- Training data:** 60 hours of US-English Windows Phone Short Message Dictation
- Testing data:** 3 hours of data from same task

Model	WER (%)
6-layer DNN (MMI training)	21.1
4-layer LSTM (MMI training)	20.8
Recurrent SVM (Max Margin training)	19.8

- 4.8% WERR. More results in the paper.

Model	WER (%)
Recurrent SVM (only train last layer)	20.2
Recurrent SVM (+ previous layer)	19.8

- 65% gains are from updating the LSTM layers

### References

[1] Schmidhuber, Jürgen, et al. Evolino for Recurrent Support Vector Machines. Euro. Symp. on Artificial Neural Networks, 2006.