

Lifu Tu Kevin Gimpel
Toyota Technological Institute at Chicago

ABSTRACT

- structured prediction is challenging due to **exponentially-large** output spaces
- How to **speed up the inference time**?
- Structured prediction energy networks (SPENs; Belanger & McCallum 2016): use neural networks to define structured energy functions
- Belanger & McCallum used **gradient descent for inference** with SPENs
- We replace this use of gradient descent with **a neural network trained to approximate structured argmax inference**.
- We develop large-margin training criteria for joint training of the **structured energy function** and **inference network**
 - On **multi-label classification**, high accuracy and **10-60x speed-ups** compared to Belanger et al. (2017)
 - Sequence labeling**: accuracies comparable to exact inference with **faster inference** speeds
 - Improved accuracy by augmenting energy with **“tag language model”**
- We also show how inference networks can replace dynamic programming at test time for conditional random fields (see paper for details)

INFERENCE NETWORKS

We define an inference network $A_\Psi(x)$ with the goal that

$$A_\Psi(x) \approx \operatorname{argmin}_{y \in \mathcal{Y}_R(x)} E_\Theta(x, y)$$

SPENs TRAINING

SPENs are trained with the following SSVM loss:

$$\min_{\Theta} \sum_{(x_i, y_i) \in \mathcal{D}} \left[\max_{y \in \mathcal{Y}_R(x)} (\Delta(y, y_i) - E_\Theta(x_i, y) + E_\Theta(x_i, y_i)) \right]_+$$

Here $[\cdot]_+ = \max(0, \cdot)$, $\Delta(y, y_i)$ is the error function between a prediction and the ground truth

However, the **“cost-augmented” inference** step is expensive.

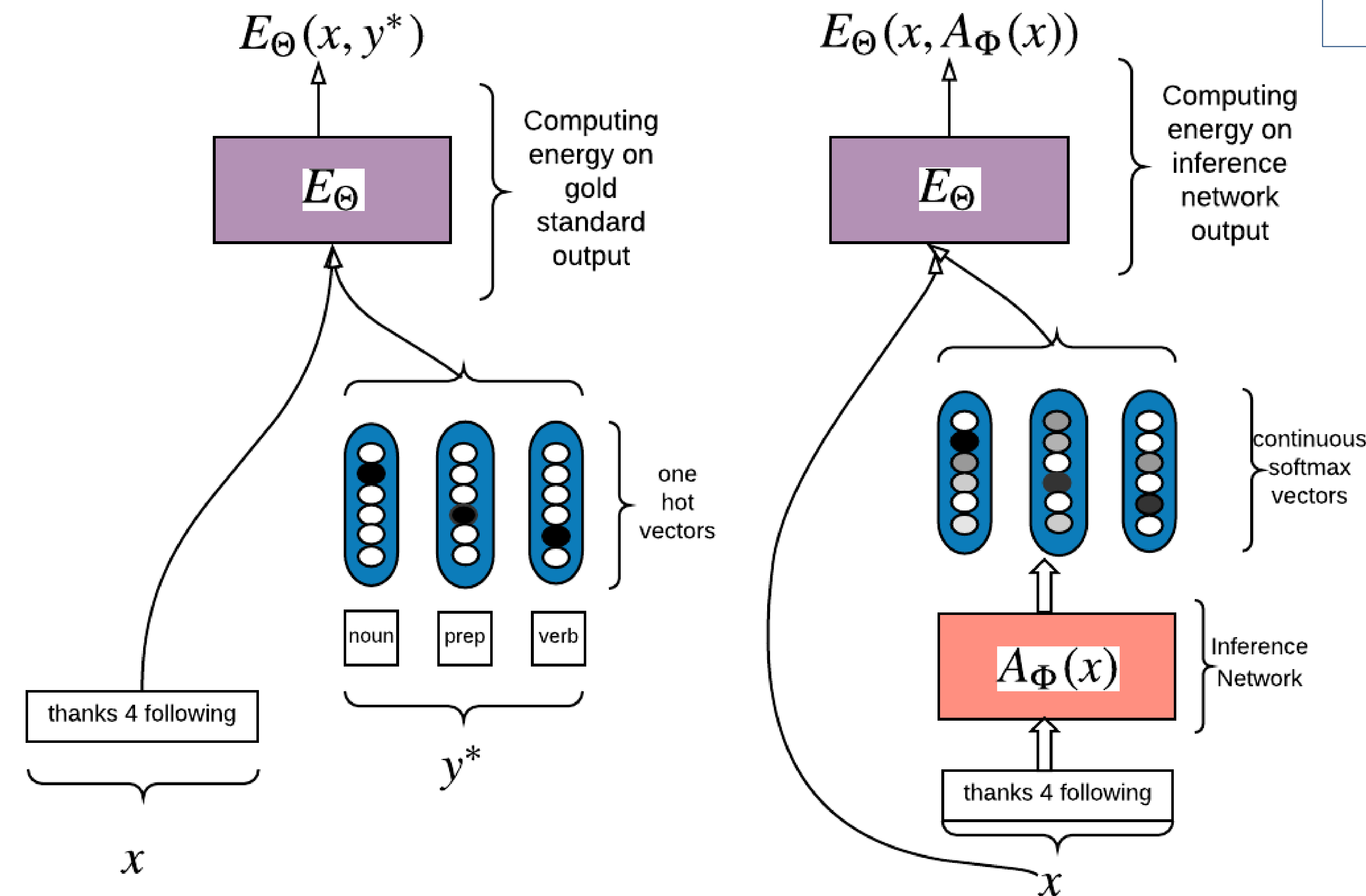
- In SPENs, this step uses **gradient step**. \mathcal{Y}_R : relaxed output space
- However, hard to do exact optimization in the inner loop with gradient descent

ENERGIES FOR SEQUENCE LABELING

$$E_\Theta(x, y) = - \left(\sum_t U_{y_t}^\top f(x, t) + \sum_t W_{y_{t-1}, y_t} \right) \leftarrow \text{CRF}$$

$$E_\Theta(x, y) = - \left(\sum_t \sum_{i=1}^L y_{t,i} (U_i^\top f(x, t)) + \sum_t y_{t-1}^\top W y_t \right) \leftarrow \text{continuous label space}$$

$$E^{\text{TLM}}(y) = - \sum_{t=1}^{|y|+1} \log(y_t^\top \text{TLM}(\langle y_0, \dots, y_{t-1} \rangle)) \leftarrow \text{General}$$



ENERGIES FOR MULTI-LABEL CLASSIFICATION

$$E^{\text{loc}}(x, y) = \sum_{i=1}^L y_i b_i^\top F(x) \quad E^{\text{lab}}(y) = c_2^\top g(C_1 y)$$

$$E_\Theta(x, y) = E^{\text{loc}}(x, y) + E^{\text{lab}}(y)$$

JOINT ADVERSARIAL TRAINING OF SPENs AND INFERENCE NETWORKS

In SSVM loss, replacing expensive **“cost-augmented” inference** step with $A_\Phi(x)$, then the new training objective is to minimize:

$$\min_{\Theta} \max_{\Phi} \sum_{(x_i, y_i) \in \mathcal{D}} [\Delta(A_\Phi(x_i), y_i) - E_\Theta(x_i, A_\Phi(x_i)) + E_\Theta(x_i, y_i)]_+$$

We optimize the objective with the following two steps:

The objective for the cost-augmented inference network is:

$$\hat{\Phi} \leftarrow \operatorname{argmax}_{\Phi} [\Delta(A_\Phi(x_i), y_i) - E_\Theta(x_i, A_\Phi(x_i)) + E_\Theta(x_i, y_i)]_+ + \text{Reg.}$$

The objective for the energy function is:

$$\hat{\Theta} \leftarrow \operatorname{argmin}_{\Theta} [\Delta(A_\Phi(x_i), y_i) - E_\Theta(x_i, A_\Phi(x_i)) + E_\Theta(x_i, y_i)]_+ + \lambda \|\Theta\|_2^2$$

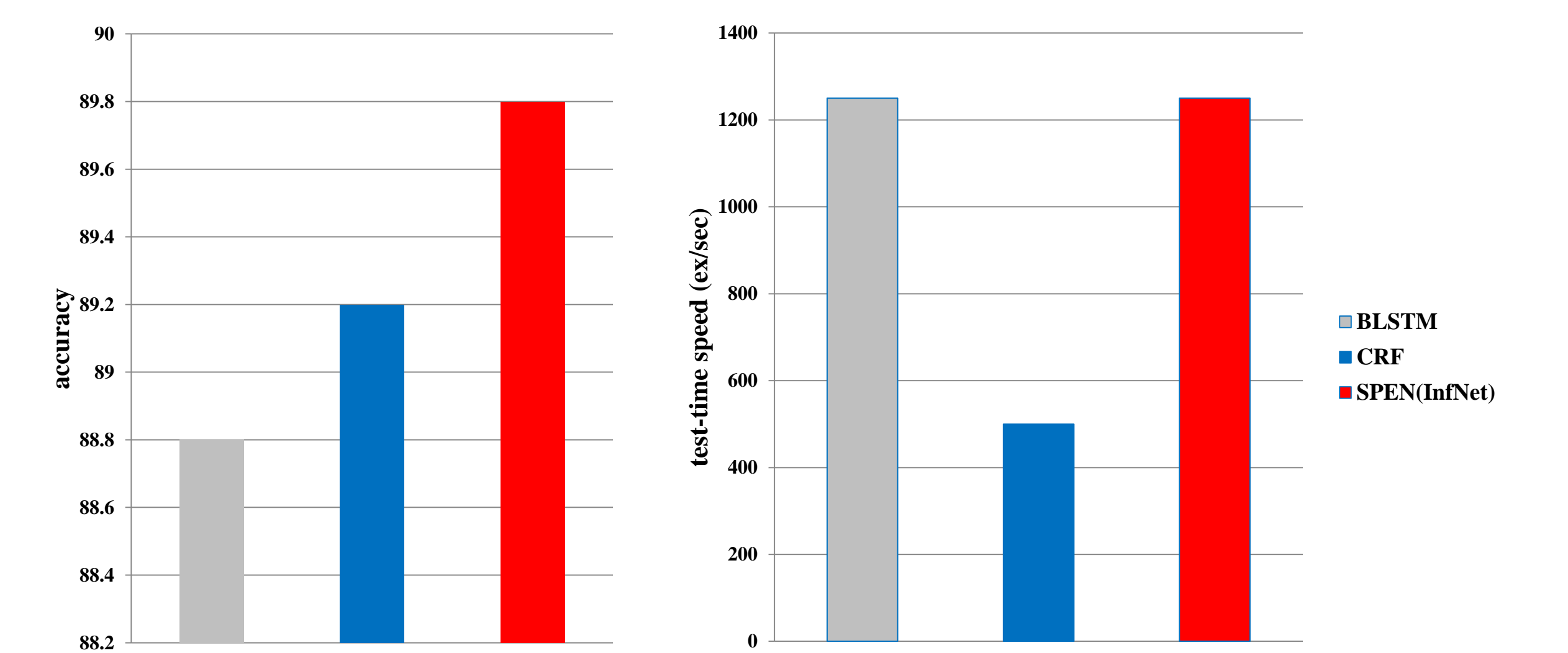
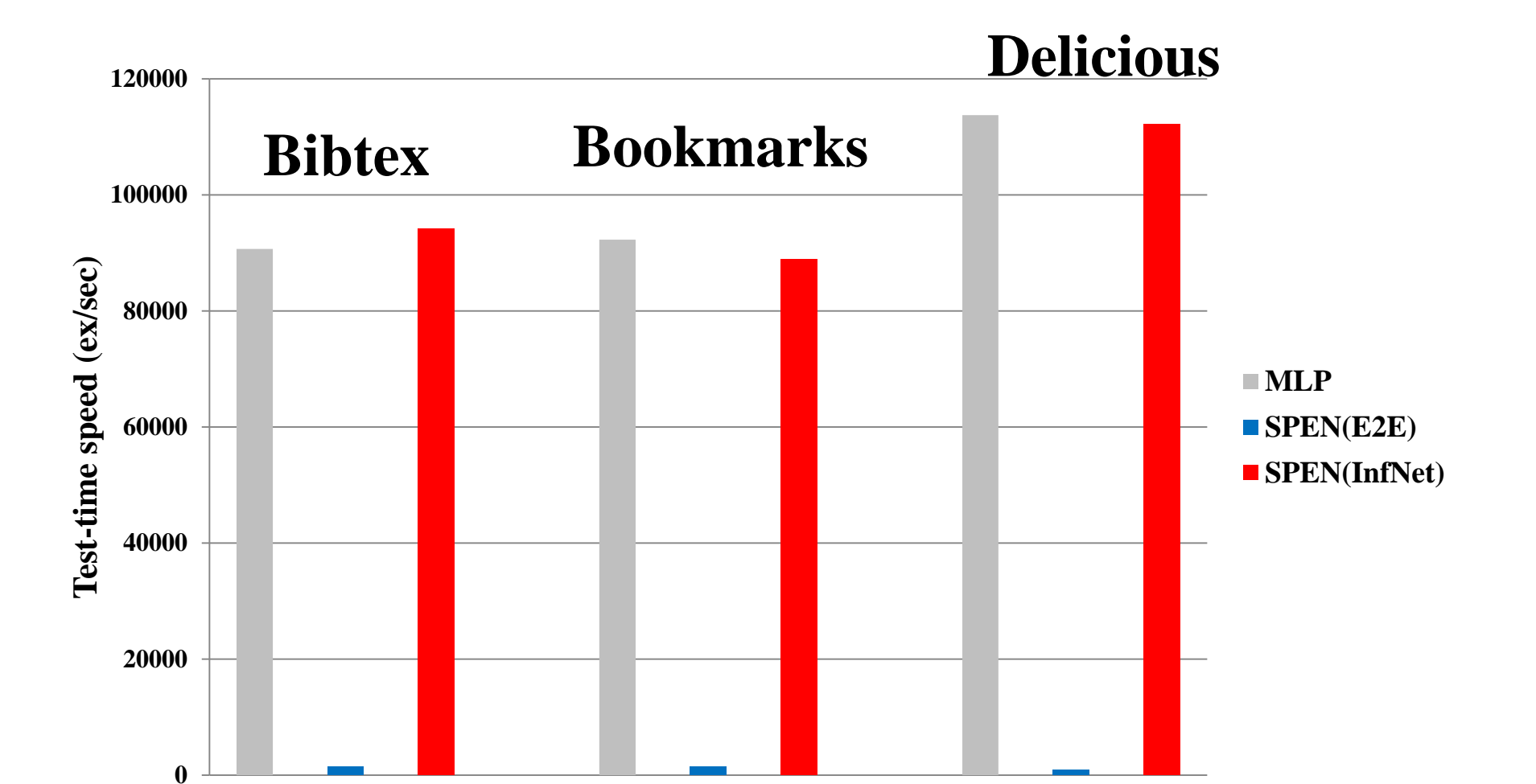
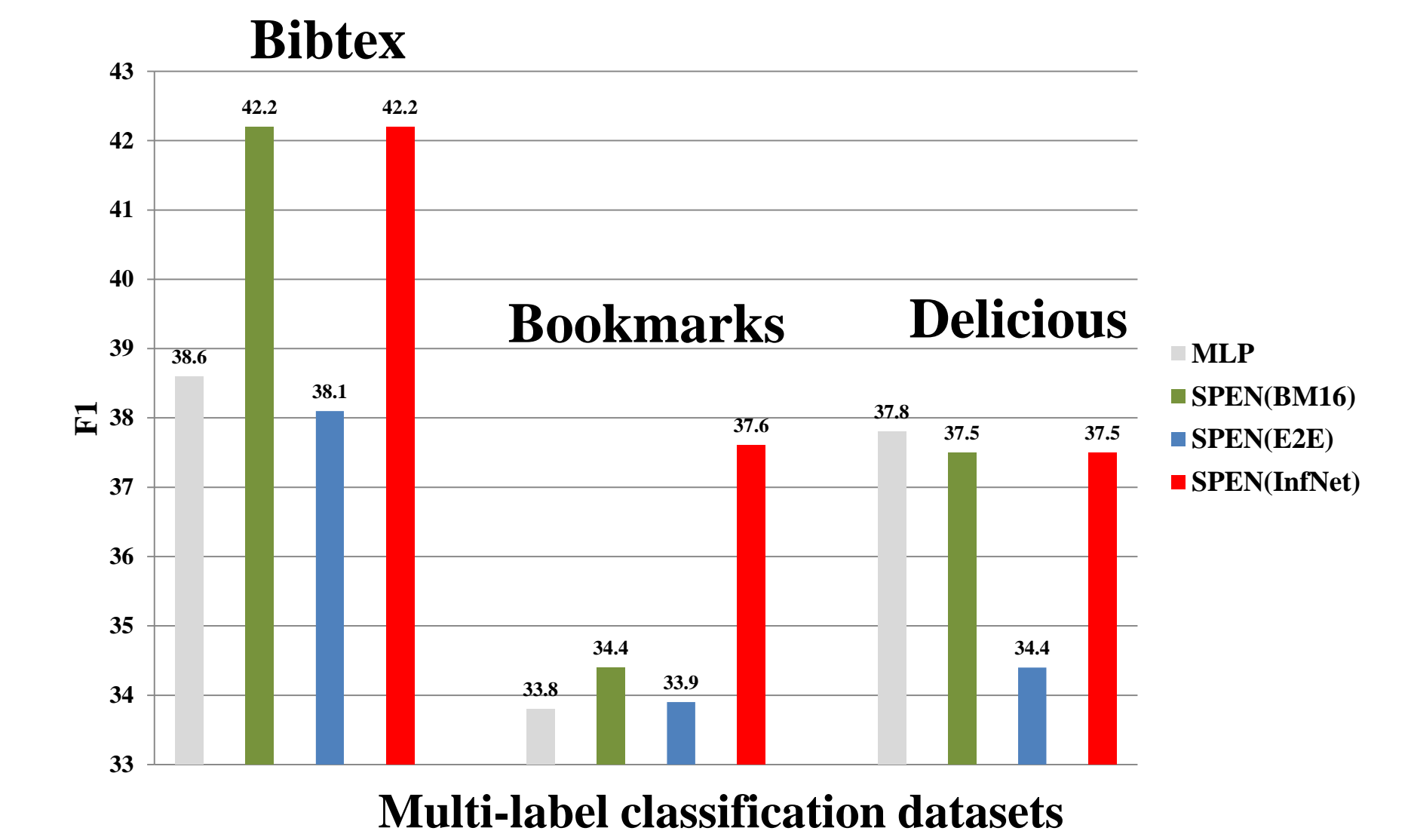
Training iterates between updating Φ and Θ

The final inference network with retuning (Ψ is initialized with Θ):

$$\hat{\Psi} \leftarrow \operatorname{argmin}_{\Psi} \sum_{x \in \mathcal{X}} E_\Theta(x, A_\Psi(x)) + \text{Reg.}$$

TEST-TIME INFERENCE

EXPERIMENTS



Twitter part-of-speech tagging

	test accuracy(%)
-TLM	89.6
+TLM	90.2

Twitter POS test accuracies when adding tag language model(TLM)

#	tweet (target word in bold)	-TLM	+TLM
1	... that's a t-17, technically . does that count as top-25 ?	determiner	pronoun
2	... lol you know im down like 4 flats on a cadillac ... lol ...	adjective	preposition
3	... them who he is : he wants her to like him for his pers ...	preposition	verb
4	Cut my hair , gag and bore me	noun	verb
5	I wonder when Nic Cage is going to film " Another Something Something Las Vegas " .	noun	verb
6	... they had their fun , we hd ours ! :) lmaooo	proper noun	verb
7	lmao I am not a sheep who listens to it cos everyone else does ...	verb	preposition
8	Noo its not cuss you have swag andd you wont look dumb ! ...	noun	coord. conj.

Examples of improvements when using tag language model

References

- Belanger, David and McCallum, Andrew. Structured Prediction Energy Networks. ICML2016
- LeCun, Yann and Chopra, Sumit and Hadsell, Raia and Ranzato, A. A Tutorial on Energy-Based Learning. MIT Press 2006
- Belanger, D. Yang, B., McCallum, A. "End-to-End Learning for Structured Prediction Energy Networks." Arxiv Preprint