



ML<sup>2</sup> Machine Learning  
for Language



# Inference Networks (InfNets) and Structured Prediction Energy Networks (SPENs)

Richard Yuanzhe Pang [yzpang.me](http://yzpang.me)

09/22/2020

Today's talk: mostly joint works with Lifu Tu, Kevin Gimpel (both currently at TTIC)

## → **Energy-based models**

Inference networks

Structured prediction energy networks

# Feed-forward models vs. energy-based models

- Feed-forward model **at inference time**
  - $\mathbf{y} = f(\mathbf{x})$  where  $f$  could be any function (e.g., a complicated neural network)
- Energy-based model **at inference time**

$$\mathbf{y} = \operatorname{argmin}_{\mathbf{y}'} E(\mathbf{x}, \mathbf{y}')$$

- Energy-based models

- Cons: energy functions may be hard to formulate/train
- Cons: inference may require dynamic programming or gradient descent
- Pros: can better capture x-y dependencies; for example, multiple  $\mathbf{y}$ 's can be compatible with a single  $\mathbf{x}$
- Pros: can inject expert knowledge to energy function; can produce parsimonious formulation => better generalization and better low-resource performance

- Feed-forward model **at training time**
  - One possibility:  $\min L(f(\mathbf{x}), \mathbf{y}_{\text{gold}})$  w.r.t.  $f$ 's params
- Energy-based model **at training time**
  - Goal: train  $E$ 's params; more complicated

# Two classes of learning methods for energy-based models

- **Contrastive:** push down on  $E(\mathbf{x}_i, \mathbf{y}_i)$ ; push up on other points  $E(\mathbf{x}_i, \mathbf{y}')$ 
  - Examples: contrastive divergence, metric learning, noise contrastive estimation, generative adversarial networks, denoising auto-encoder, masked auto-encoder
- **Architectural methods:** build  $E(\mathbf{x}, \mathbf{y})$  such that the volume of the low energy regions is limited or minimized through regularization
  - Examples: sparse coding, sparse auto-encoder, variational auto-encoders (VAEs), etc.
- **Today:** contrastive methods

Energy-based models

→ **Inference networks**

Structured prediction energy networks

# Inference networks (InfNets)

- Exact inference

$$\hat{\mathbf{y}} = \operatorname{argmin}_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} E_{\Theta}(\mathbf{x}, \mathbf{y})$$

- Gradient descent for inference

$$GD(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y} \in \mathcal{Y}_R(\mathbf{x})} E_{\Theta}(\mathbf{x}, \mathbf{y})$$

- Inference network for inference

- InfNet (can be a neural network):

$$A_{\Psi} : \mathcal{X} \rightarrow \mathcal{Y}_R$$

- Inference time:

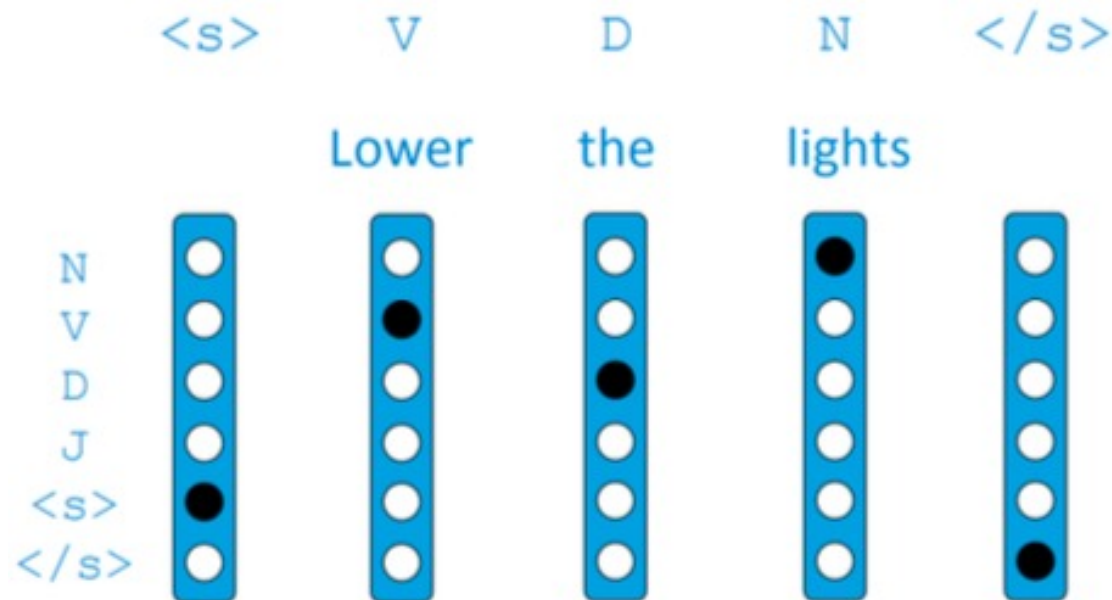
$$A_{\Psi}(\mathbf{x}) \approx \operatorname{argmin}_{\mathbf{y} \in \mathcal{Y}_R(\mathbf{x})} E_{\Theta}(\mathbf{x}, \mathbf{y})$$

- The question is, how do we train the InfNet?

# Aside

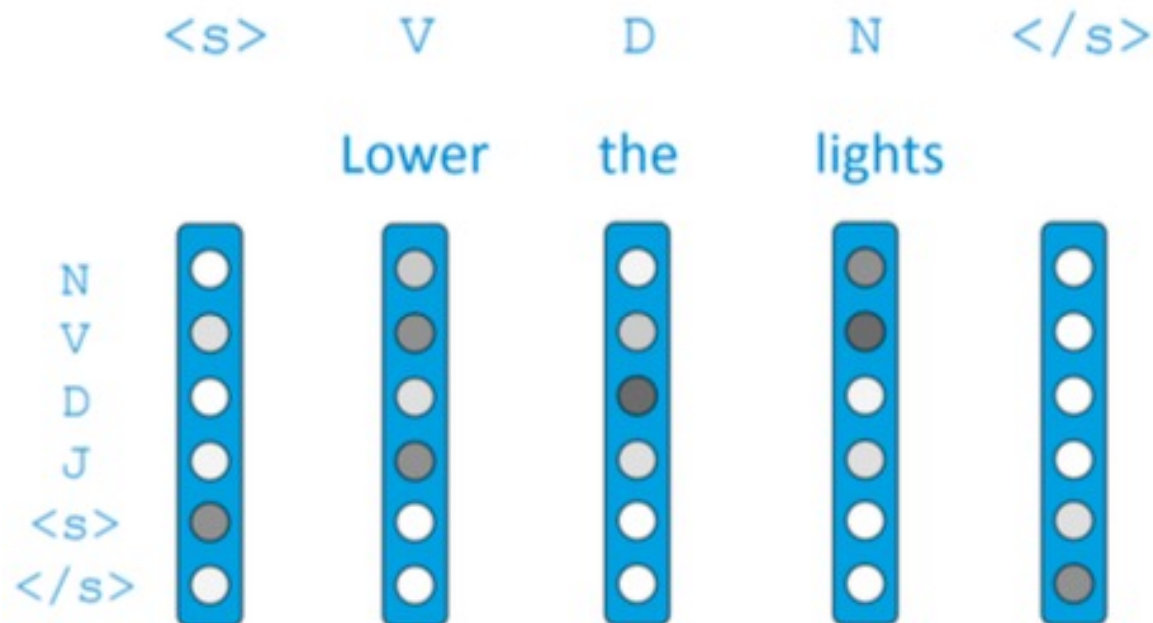
original “hard” space

- use 1 to represent V
- use 2 to represent D, etc.



“relaxed” space

- use a distribution to represent V
- e.g., (0, 1, 0, 0, 0, 0)
- or, (0.1, 0.8, 0.05, 0.05, 0, 0.1))



# Inference networks (InfNets)

- Exact inference

$$\hat{\mathbf{y}} = \operatorname{argmin}_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} E_{\Theta}(\mathbf{x}, \mathbf{y})$$

- Gradient descent for inference

$$GD(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y} \in \mathcal{Y}_R(\mathbf{x})} E_{\Theta}(\mathbf{x}, \mathbf{y})$$

- Inference network for inference

- InfNet (can be a neural network):

$$A_{\Psi} : \mathcal{X} \rightarrow \mathcal{Y}_R$$

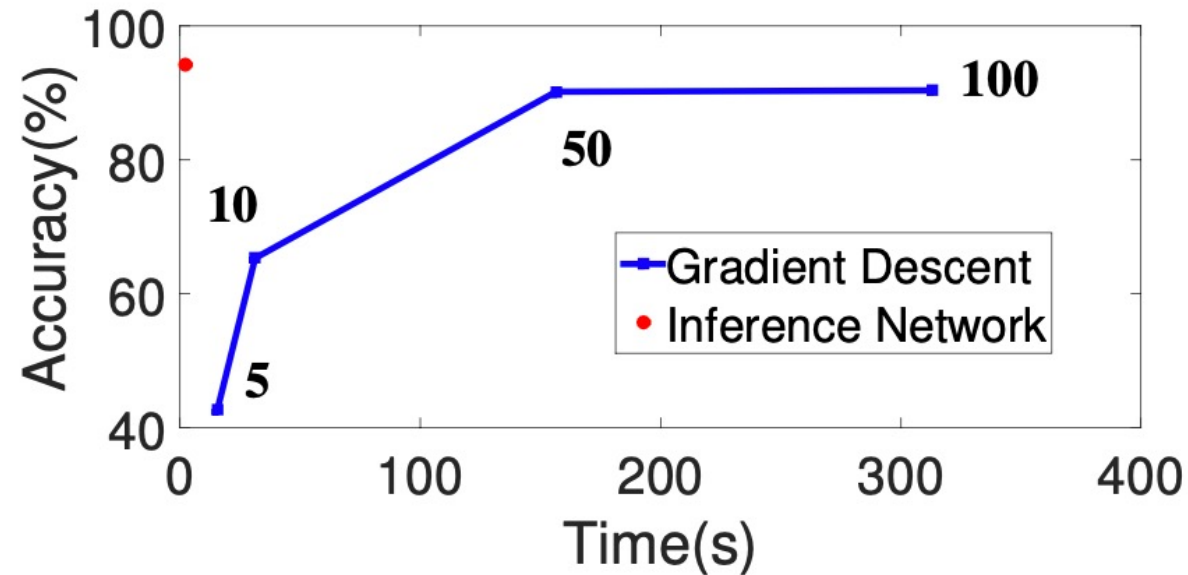
- Inference time:

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

- The question is, how do we train the InfNet?



## Inference network: fast and accurate



### Experimental Details

CCG Supertagging with 400 labels

Energy function: BLSTM-CRF

Inference network architecture: BLSTM

Gradient descent run for {5, 10, 50, 100} iterations

Energy-based models

Inference networks

→ **Structured prediction energy networks**

$$? = \operatorname{argmin}_y E(x, y')$$

# Structured prediction energy networks (SPENs)

- Original SPEN (Belanger and McCallum, 2016)

- Training:

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

- Inference:

$$GD(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y} \in \mathcal{Y}_R(\mathbf{x})} E_{\Theta}(\mathbf{x}, \mathbf{y})$$

# Structured prediction energy networks (SPENs)

- Original SPEN (Belanger and McCallum, 2016)

- Training:

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

- Approximate inference version

- Training:

$$\min_{\Theta} \max_{\Phi} \sum_{\langle \mathbf{x}_i, \mathbf{y}_i \rangle \in \mathcal{D}} [(\Delta(A_{\Phi}(\mathbf{x}_i), \mathbf{y}_i) - E_{\Theta}(\mathbf{x}_i, A_{\Phi}(\mathbf{x}_i)) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i))]_+$$

# Joint training of SPEN and InfNet

- Training:

$$\min_{\Theta} \max_{\Phi} \sum_{\langle \mathbf{x}_i, \mathbf{y}_i \rangle \in \mathcal{D}} [(\Delta(A_{\Phi}(\mathbf{x}_i), \mathbf{y}_i) - E_{\Theta}(\mathbf{x}_i, A_{\Phi}(\mathbf{x}_i)) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i))]_+$$

cost-augmented InfNet

test-time InfNet

- Inference:

$$A_{\Psi}(\mathbf{x})$$

Recall that InfNet is designed to approximate as follows

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

# Joint training of SPEN and InfNet

- Training:

$$\min_{\Theta} \max_{\Phi} \sum_{\langle \mathbf{x}_i, \mathbf{y}_i \rangle \in \mathcal{D}} [(\Delta(A_{\Phi}(\mathbf{x}_i), \mathbf{y}_i) - E_{\Theta}(\mathbf{x}_i, A_{\Phi}(\mathbf{x}_i)) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i))]_{+}$$

- Fine-tuning InfNet parameters:

$$\hat{\Psi} \leftarrow \hat{\Phi}$$

$$\hat{\Psi} = \operatorname{argmin}_{\Psi} E_{\Theta}(\mathbf{x}, A_{\Psi}(\mathbf{x}))$$

- Inference:

$$A_{\Psi}(\mathbf{x})$$

- Cost-augmented inference

$$\operatorname{argmin}_{\mathbf{y}'} (E_{\Theta}(\mathbf{x}, \mathbf{y}') - \Delta(\mathbf{y}', \mathbf{y}))$$

- Test-time inference

$$\operatorname{argmin}_{\mathbf{y}'} E_{\Theta}(\mathbf{x}, \mathbf{y}')$$

# Joint training of SPEN and (cost-augmented InfNet and test-time InfNet)

- Training:

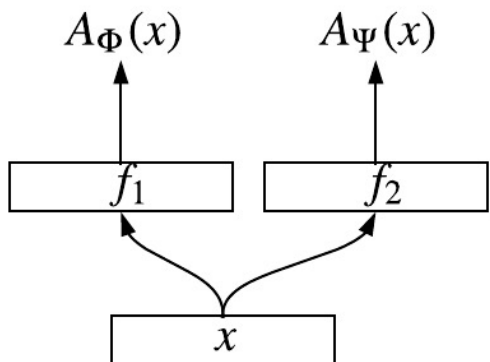
$$\min_{\Theta} \max_{\Phi, \Psi} \sum_{\langle \mathbf{x}_i, \mathbf{y}_i \rangle \in \mathcal{D}} [(\Delta(A_{\Phi}(\mathbf{x}_i), \mathbf{y}_i) - E_{\Theta}(\mathbf{x}_i, A_{\Phi}(\mathbf{x}_i)) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i))]_{+} + \lambda [-E_{\Theta}(\mathbf{x}_i, A_{\Psi}(\mathbf{x}_i)) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i)]_{+}$$

- Inference:

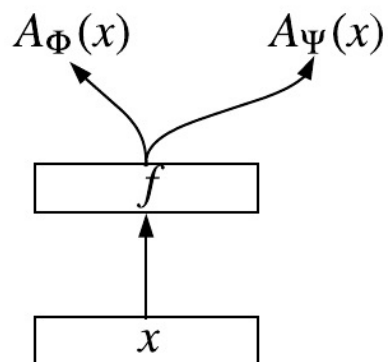
$$A_{\Psi}(\mathbf{x})$$

- An example energy function for sequence labeling

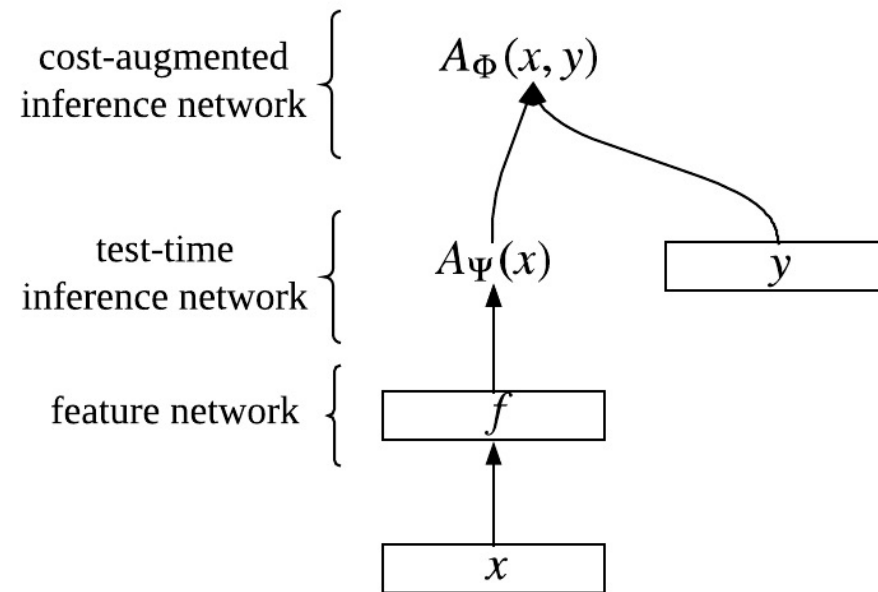
$$E_{\Theta}(\mathbf{x}, \mathbf{y}) = - \left( \sum_{t=1}^T \sum_{j=1}^L y_{t,j} (U_j^{\top} b(\mathbf{x}, t)) + \sum_{t=1}^T \mathbf{y}_{t-1}^{\top} W \mathbf{y}_t \right)$$



(a) separated networks



(b) shared feature networks



(c) stacked networks with  $y$  as extra input to  $A_\Phi$

- Training:

$$\min_{\Theta} \max_{\Phi, \Psi} \sum_{\langle \mathbf{x}_i, \mathbf{y}_i \rangle \in \mathcal{D}} [(\Delta(A_\Phi(\mathbf{x}_i), \mathbf{y}_i) - E_\Theta(\mathbf{x}_i, A_\Phi(\mathbf{x}_i)) + E_\Theta(\mathbf{x}_i, \mathbf{y}_i))]_+ + \lambda [-E_\Theta(\mathbf{x}_i, A_\Psi(\mathbf{x}_i)) + E_\Theta(\mathbf{x}_i, \mathbf{y}_i)]_+$$

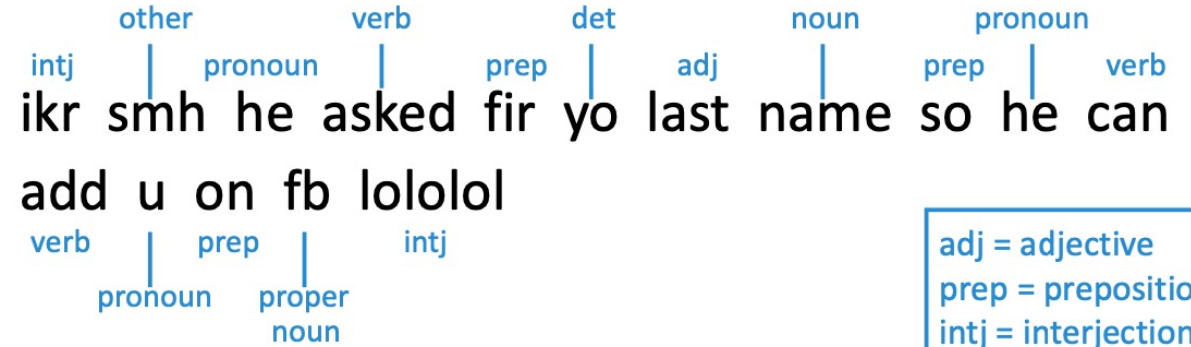
- Inference:

$$A_\Psi(\mathbf{x})$$

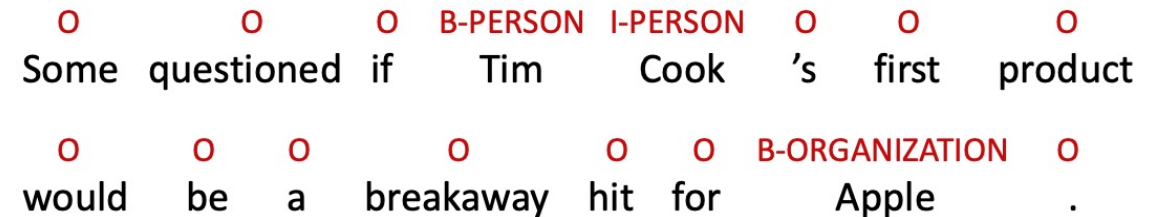


# Experimental setup

1. Part-of-speech (POS) tagging
  2. Named entity recognition (NER) labeling
  3. Constituency parsing (skip)
- Only experimented on **sequence labeling** tasks. Tu and Gimpel (2018) has some multi-label classification results.
  - Note that although these are NLP experiments, the approach can be applied to structured prediction tasks in general.



## Named Entity Recognition



**B = "begin"**  
**I = "inside"**  
**O = "outside"**

# Our hypotheses

1. **(InfNet + SPEN) > MLE**
2. (Cost-augmented InfNet + test-time InfNet + SPEN) > (InfNet + SPEN)
3. Tag language models help

	POS				NER				NER+ F1 (%)
	acc (%)	T	I	speed	F1 (%)	T	I	speed	
BiLSTM	88.8	166K	166K	–	84.9	239K	239K	–	89.3
margin-rescaled	89.4	333K	166K	–	85.2	479K	239K	–	89.5
perceptron	88.6	333K	166K	–	84.4	479K	239K	–	89.0

# Our hypotheses

1. (InfNet + SPEN) > MLE
2. **(Cost-augmented InfNet + test-time InfNet + SPEN) > (InfNet + SPEN)**
3. Tag language models help

	POS				NER				NER+ F1 (%)
	acc (%)	T	I	speed	F1 (%)	T	I	speed	
BiLSTM	88.8	166K	166K	–	84.9	239K	239K	–	89.3

margin-rescaled	89.4	333K	166K	–	85.2	479K	239K	–	89.5
perceptron	88.6	333K	166K	–	84.4	479K	239K	–	89.0

## **SPENs with inference networks, compound objective, CE, no zero truncation (this paper):**

separated	89.7	500K	166K	66	85.0	719K	239K	32	89.8
shared	89.8	339K	166K	78	85.6	485K	239K	38	90.1
<b>stacked</b>	<b>89.8</b>	<b>335K</b>	<b>166K</b>	<b>92</b>	<b>85.6</b>	<b>481K</b>	<b>239K</b>	<b>46</b>	<b>90.1</b>

# Our hypotheses

1. (InfNet + SPEN) > MLE
2. (Cost-augmented InfNet + test-time InfNet + SPEN) > (InfNet + SPEN)
3. **Tag language models help**

# Joint training of SPEN and (cost-augmented InfNet and test-time InfNet)

- Training:  $\min_{\Theta} \max_{\Phi, \Psi} \sum_{\langle \mathbf{x}_i, \mathbf{y}_i \rangle \in \mathcal{D}} [(\Delta(A_{\Phi}(\mathbf{x}_i), \mathbf{y}_i) - E_{\Theta}(\mathbf{x}_i, A_{\Phi}(\mathbf{x}_i)) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i))]_{+} + \lambda [-E_{\Theta}(\mathbf{x}_i, A_{\Psi}(\mathbf{x}_i)) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i)]_{+}$
- Inference:  $A_{\Psi}(\mathbf{x})$
- An example energy function for sequence labeling

$$E_{\Theta}(\mathbf{x}, \mathbf{y}) = - \left( \sum_{t=1}^T \sum_{j=1}^L y_{t,j} (U_j^{\top} b(\mathbf{x}, t)) + \sum_{t=1}^T \mathbf{y}_{t-1}^{\top} W \mathbf{y}_t \right)$$

$$\bar{\mathbf{y}}_t = h(\mathbf{y}_0, \dots, \mathbf{y}_{t-1})$$

$$\bar{\mathbf{y}}_t = h(\mathbf{x}_0, \dots, \mathbf{x}_{t-1}, \mathbf{y}_0, \dots, \mathbf{y}_{t-1})$$

$$E^{\text{TLM}}(\mathbf{y}) = - \sum_{t=1}^{T+1} \log(\mathbf{y}_t^{\top} \bar{\mathbf{y}}_t)$$

# Our hypotheses

1. (InfNet + SPEN) > MLE
2. (Cost-augmented InfNet + test-time InfNet + SPEN) > (InfNet + SPEN)
3. **Tag language models help**

	NER	NER+	NER++
margin-rescaled	85.2	89.5	90.2
compound, stacked, CE, no truncation	85.6	90.1	90.8
+ global energy GE(c)	<b>86.3</b>	<b>90.4</b>	<b>91.0</b>

Table 4: NER test F1 scores with global energy terms.

# Conclusions and thoughts

1. SPENs are powerful but learning and inference are hard
2. Inference networks can make it easier and more efficient to use SPENs
3. Separating inference networks for the two inference problems (cost-augmented and test-time inference) improves accuracy and leads to complementary functionality
4. Adding global energy terms leads to further improvements
5. Next step: move to generation tasks and model other types of data (not only sequence labeling)