

# Learning to Parse Using a Tiny Corpus

Tao Lei, Yu Xin

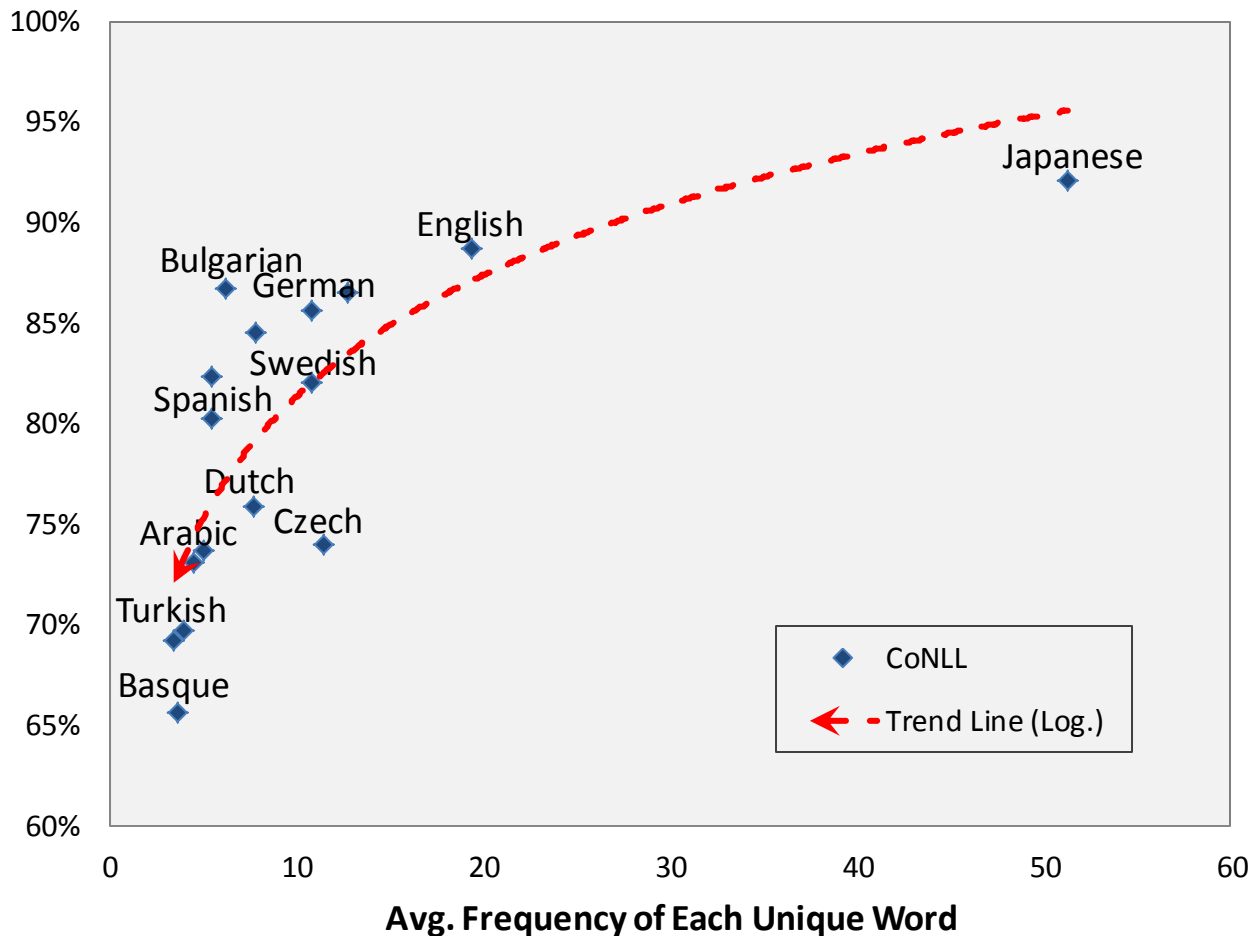
Regina Barzilay, Tommi Jaakkola

CSAIL, MIT



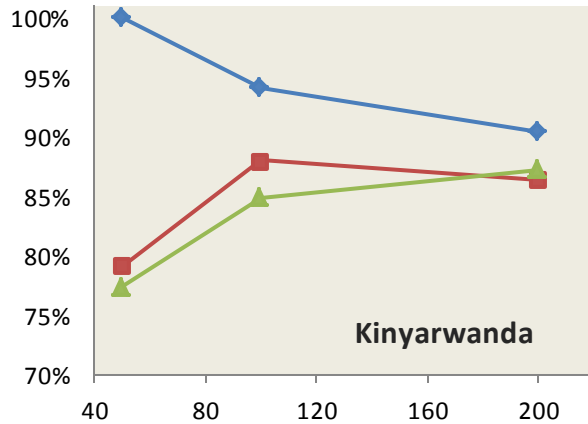
# Sparsity Problem

- Data sparsity makes parsing harder
  - due to *less frequent/unseen words* and *dependency arcs* in data

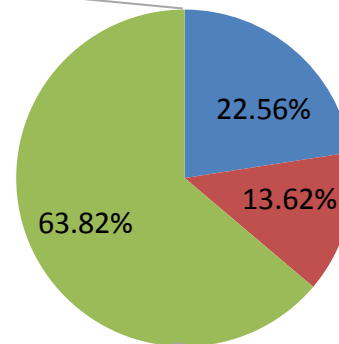
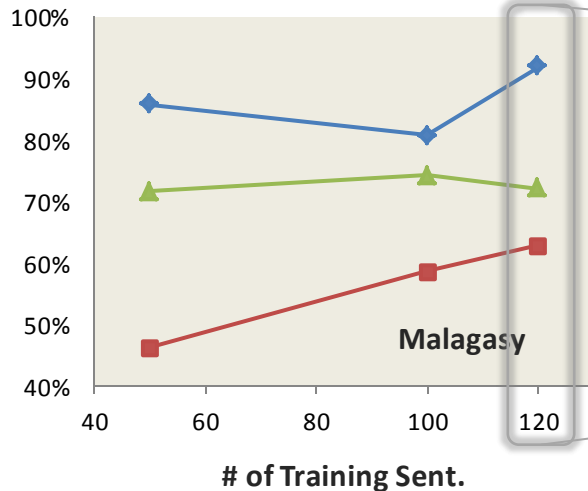


# Sparsity Problem

- Prediction is worse when the arc is **not seen** in the training data



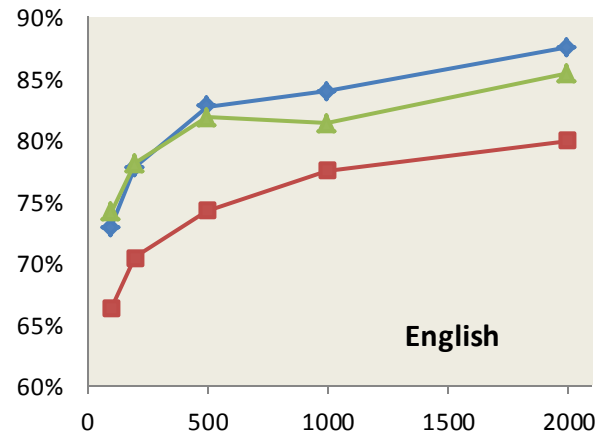
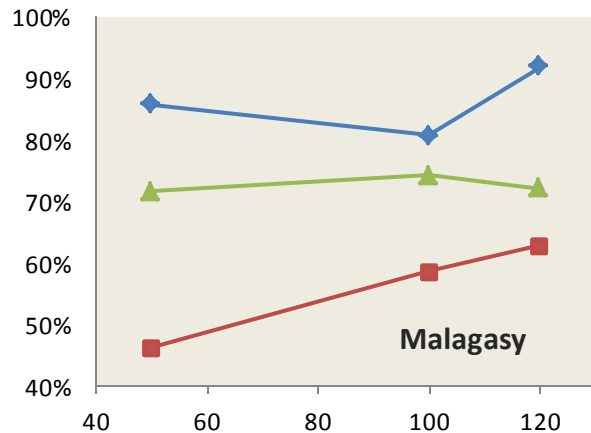
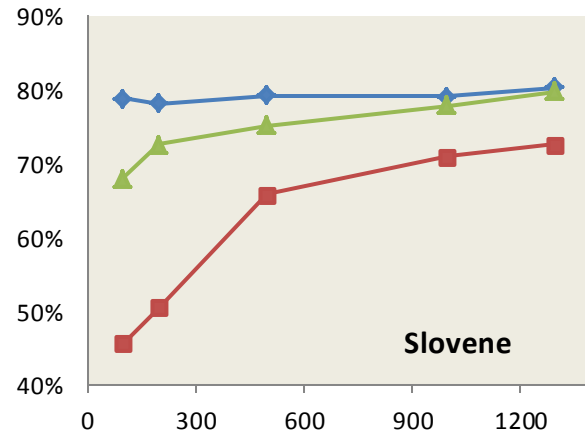
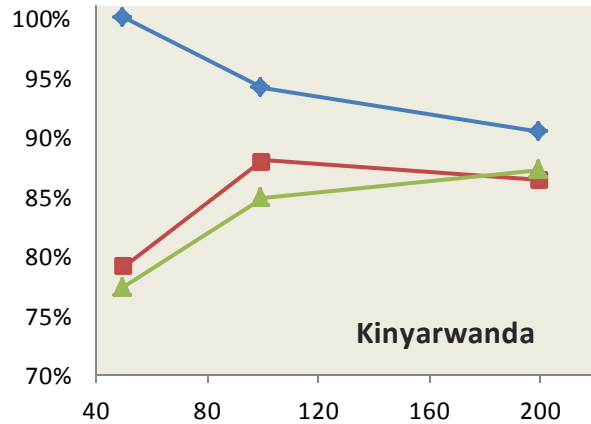
- Seen Arc: See dependency arc in the training data
- Seen Words: See the words in training but arc unseen
- Unseen: At least one word not in training



a large portion of dependency arcs in test is unseen

# Sparsity Problem

- Prediction is worse when the arc is **not seen** in the training data



# Sparsity Problem

- Feature weights are zeros when the features are not seen



# Opportunity and Challenge

To deal with sparsity problem, we will

- Make the model flexible to add **various rich features**
  - For example, words, coarse-to-fine POS tags and word embeddings
  - Adjust complexity based on how much training data it has
- Model **interactions** between feature weights
  - E.g. **propagating weights** from seen features to unseen features

# Motivating Example: Matrix Completion

- Learn a matrix (or high-order tensor) that has a lot of **unseen entries**
  - Example: image



**Input**

image with missing values



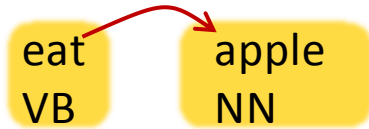
**Output**

re-constructed image

- Example: Netflix problem
  - Users give only a few movie ratings. Predict unseen ratings

# Motivating Example

- In our case: learn a parameter matrix (or tensor) with unseen weights



*Dependency Arc*

$\{ \text{eat} \oplus \text{apple}, \text{eat} \oplus \text{NN}, \text{VB} \oplus \text{apple}, \text{VB} \oplus \text{NN} \}$

*Feature Strings*

	eat	apple	banana	VB	NN
eat		1			1
apple					
banana					
VB		1			1
NN					

*Feature Matrix*



...	<b>2</b>	...	...	<b>4</b>
...	0	0	...	...
...	0	0	...	...
...	<b>1</b>	0.9	...	<b>5</b>
...	0.1	0.1	...	...

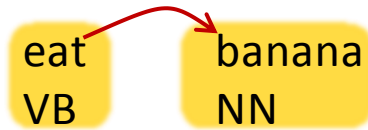
*Parameter Matrix*

**= 12**



# Motivating Example

- In our case: learn a parameter matrix (or tensor) with unseen weights



*Dependency Arc*

{ eat $\oplus$ banana, eat $\oplus$ NN, VB $\oplus$ banana, VB $\oplus$ NN }

*Feature Strings*

	eat	apple	banana	VB	NN
eat			1		1
apple					
banana					
VB			1		1
NN					

*Feature Matrix*

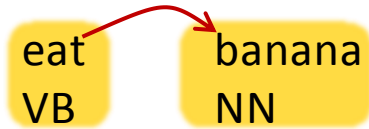


...	2	<b>??</b>	...	4
...	0	0	...	...
...	0	0	...	...
...	1	<b>0.9</b>	...	<b>5</b>
...	0.1	0.1	...	...

*Parameter Matrix*

# Motivating Example

- In our case: learn a parameter matrix (or tensor) with unseen weights



*Dependency Arc*

{ eat $\oplus$ banana, eat $\oplus$ NN, VB $\oplus$ banana, VB $\oplus$ NN }

*Feature Strings*

	eat	apple	banana	VB	NN
eat			1		1
apple					
banana					
VB			1		1
NN					

*Feature Matrix*



...	2	??	...	4
...	0	0	...	...
...	0	0	...	...
...	1	0.9	...	5
...	0.1	0.1	...	...

*Parameter Matrix*

similar columns because  
"apple" and "banana" have  
similar syntactic behavior

# Motivating Example

- In our case: learn a parameter matrix (or tensor) with unseen weights



*Dependency Arc*

{ eat $\oplus$ banana, eat $\oplus$ NN, VB $\oplus$ banana, VB $\oplus$ NN }

*Feature Strings*

	eat	apple	banana	VB	NN
eat			1		1
apple					
banana					
VB			1		1
NN					

*Feature Matrix*



...	2	<b>2</b>	...	<b>4</b>
...	0	0	...	...
...	0	0	...	...
...	1	<b>0.9</b>	...	<b>5</b>
...	0.1	0.1	...	...

*Parameter Matrix*

= 11.9

# Formulation

- Recall first-order decoding objective:

$$\begin{aligned}\tilde{y}_i &= \operatorname{argmax}_{y_i \in T(x_i)} S(y_i) \\ &= \operatorname{argmax}_{y_i \in T(x_i)} \sum_{(h,c) \in y_i} s(h,c)\end{aligned}$$

- Define score as matrix (tensor) inner product:

$$s(h,c) = \boldsymbol{\theta} \otimes \phi(h,c)$$

$$s(h,c) = w_h^T \mathbf{A} w_c + \boldsymbol{\eta}^T d(h,c)$$

Model Parameters:  $\boldsymbol{\theta} = \{\mathbf{A}, \boldsymbol{\eta}\}$

Feature Matrix/Vector:  $w_h w_c^T$   $d(h,c)$

# Formulation

- Minimize the loss of training data:

$$\min_{A, \eta} \mathcal{L}(D; A, \eta) = \frac{1}{N} \ell(x_i, \hat{y}_i)$$

$$\text{s. t. } \|A\|_* + \lambda \|\eta\| \leq C$$

Force  $A$  to be low-rank using  
nuclear norm constraint

- online learning algorithm available

[\(Jaggi & Sulovsky, 2010\)](#) [\(Hazan, 2008\)](#)

Initialize  $Z^{(1)} := v_0 v_0^T$  for arbitrary unit vector  $v$ .

For  $k = 1$  to  $T$  do

    Compute  $v_k := \text{ApproxEV} \left( -\nabla f(Z^{(k)}), \frac{C_f}{k^2} \right)$ .

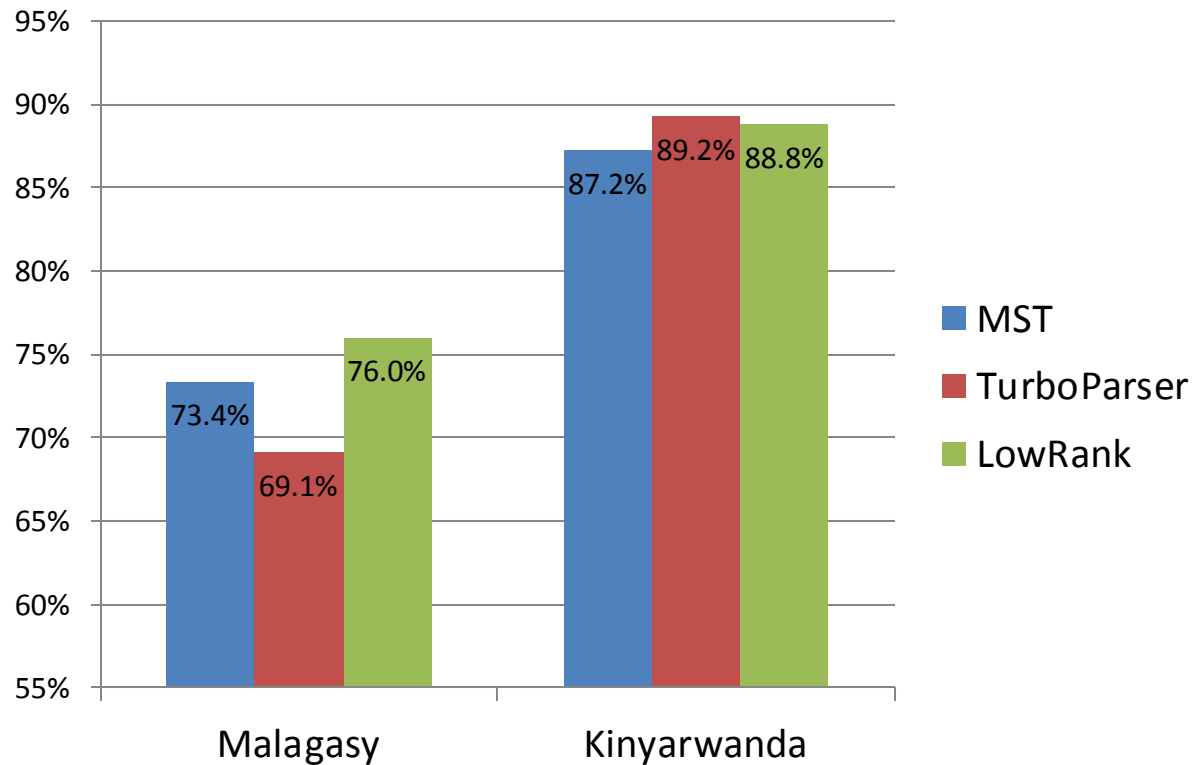
    Set  $\alpha_k := \frac{2}{k}$ .

    Set  $Z^{(k+1)} := Z^{(k)} + \alpha_k (v_k v_k^T - Z^{(k)})$ .

End for

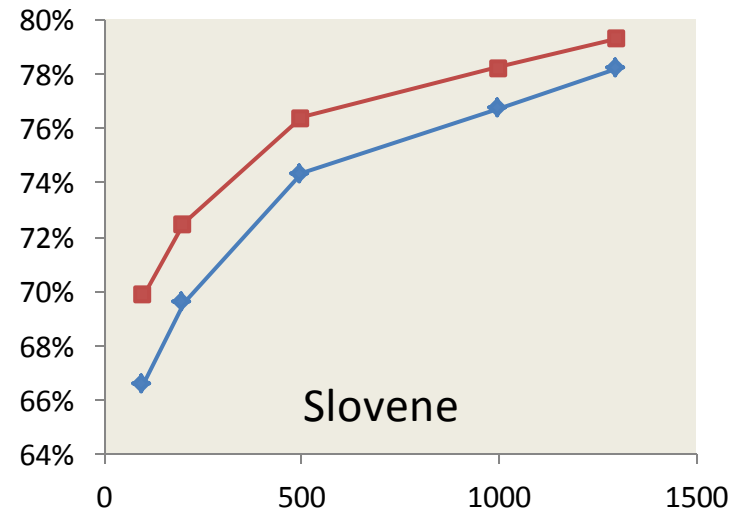
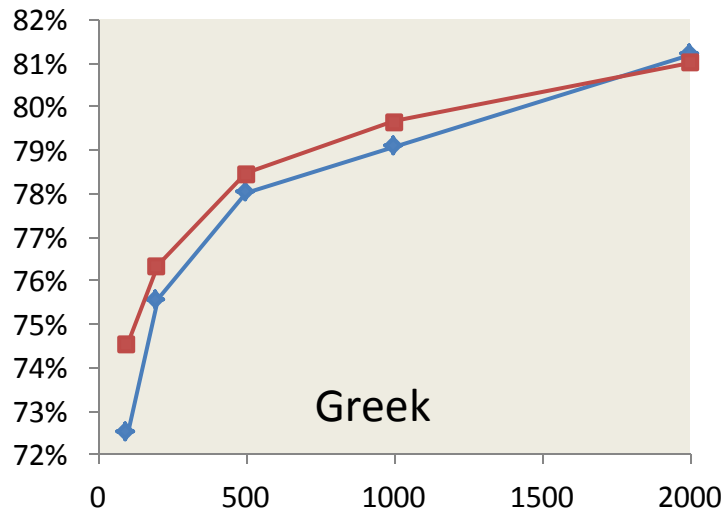
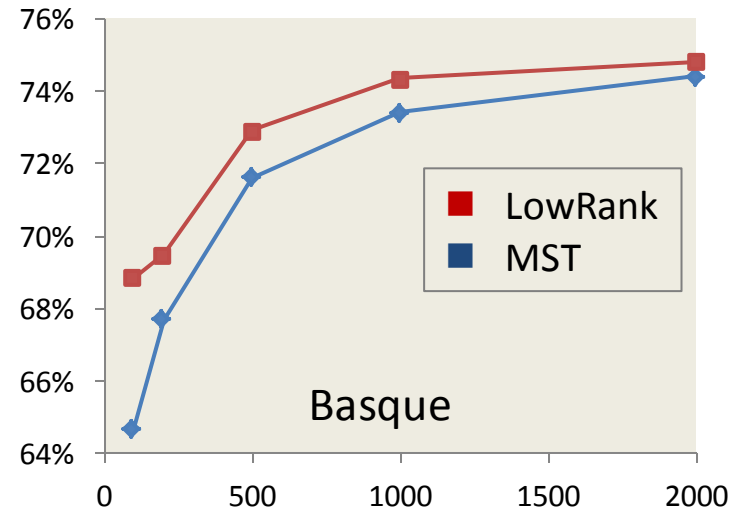
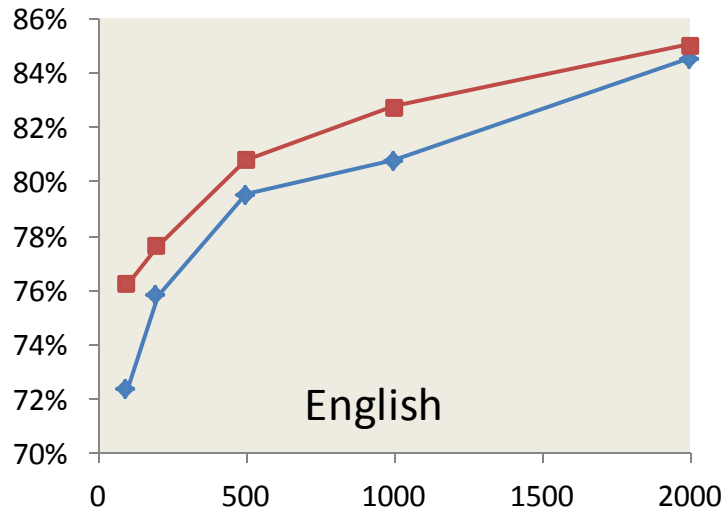
# Results

- Models trained with gold POS tags
- Evaluated without “#” tag at the end of sentences



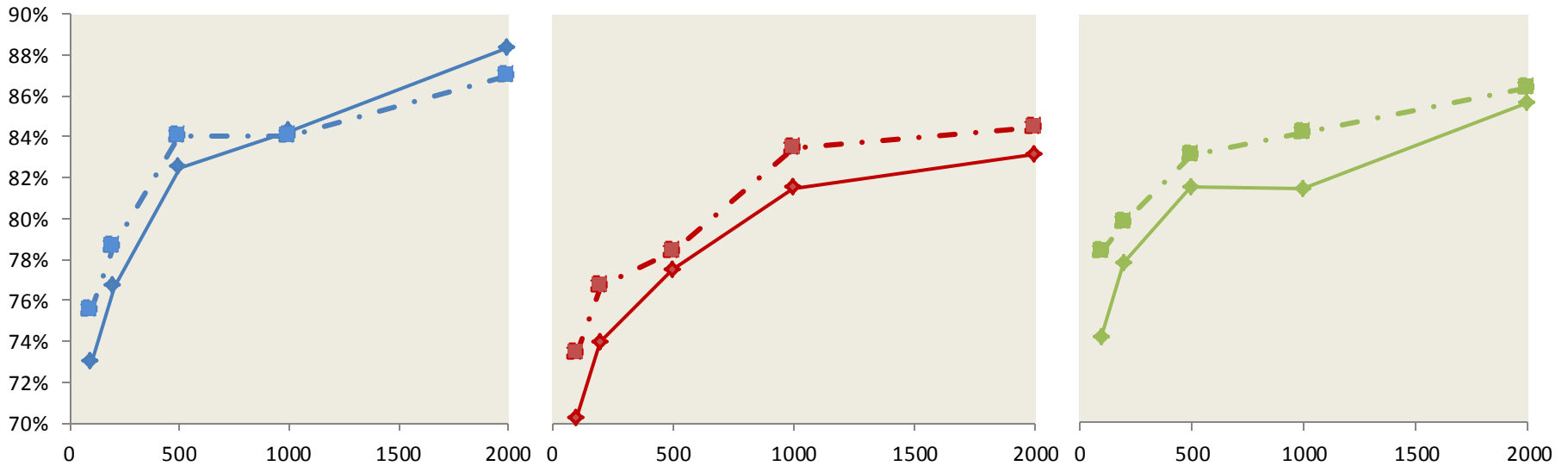
# Results

- Results on CoNLL shared task (up to 2000 sentences)



# Results

- MST parser: solid lines
- Low-rank parser: dotted lines



- Seen Arc: See dependency in the training data
- Seen Words: See the words in training but arc unseen
- Unseen: At least one word not in training



# Results

- Adding unsupervised word embeddings to English

	MST	MST+label	LowRank	LowRank+vv
100	72.5%	72.4%	76.3%	<b>76.6% (+0.3%)</b>
200	75.8%	75.8%	77.7%	<b>78.0% (+0.3%)</b>
500	79.4%	79.5%	80.8%	<b>81.4% (+0.6%)</b>
1000	80.9%	80.8%	<b>82.8%</b>	<b>82.8% (+0.0%)</b>
2000	83.7%	84.5%	85.1%	<b>85.8% (+0.7%)</b>

- **MST+label:** MST parser trained with labeled dependencies
- Other models are trained with only unlabeled dependencies.

# Thanks

- Current implementation available at:  
<http://people.csail.mit.edu/taolei/muri/lowrankparser.zip>