# Learning to Parse Using a Tiny Corpus

Tao Lei, Yu Xin
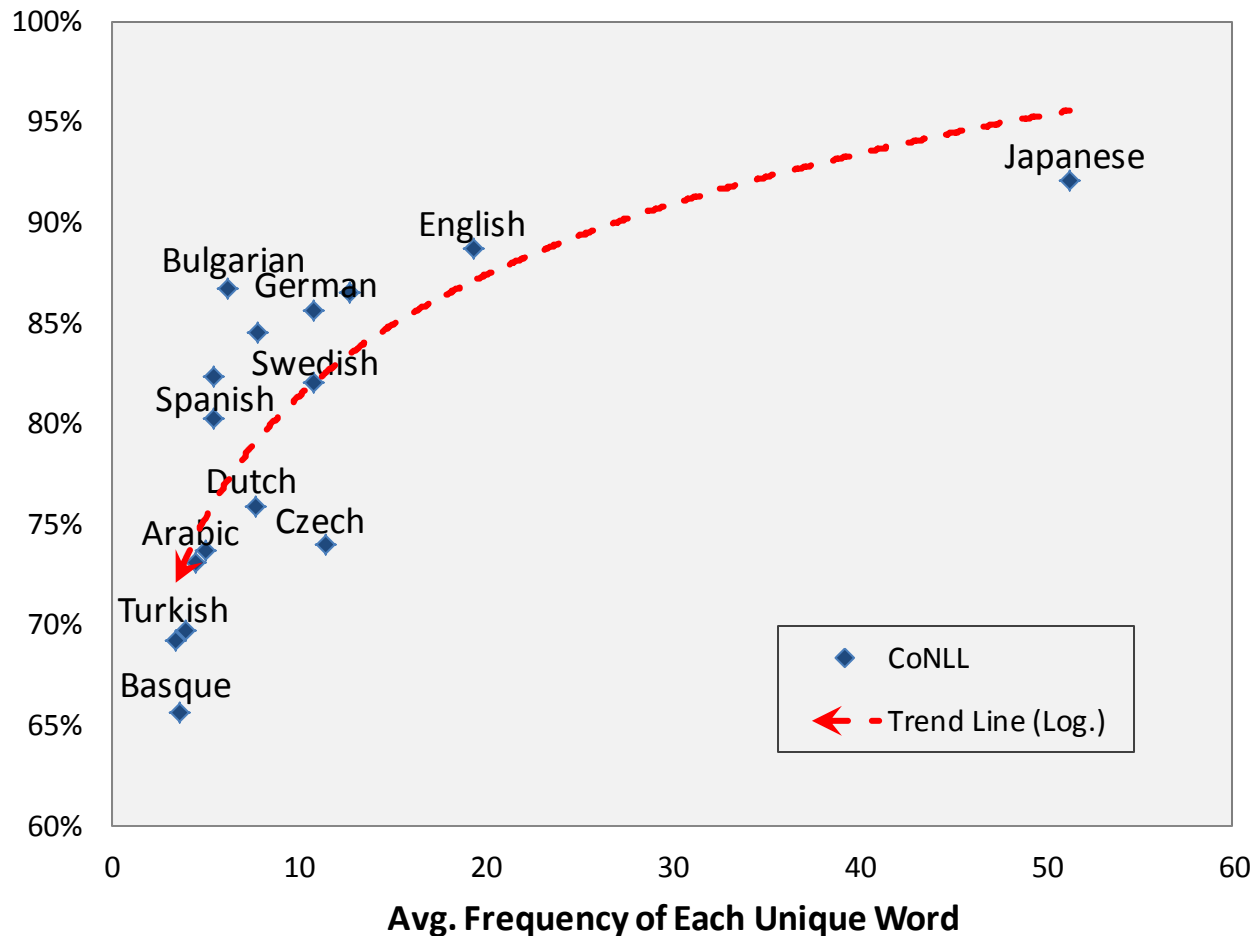
Regina Barzilay, Tommi Jaakkola

CSAIL, MIT
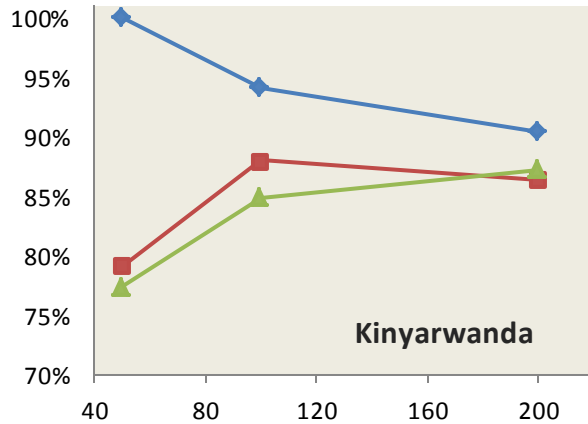
# Sparsity Problem

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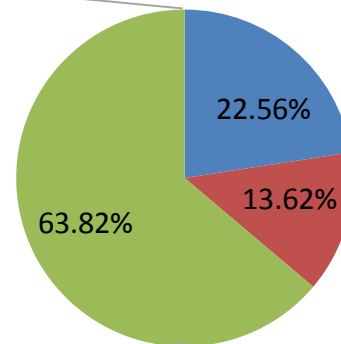  - *due to less frequent/unseen words and dependency arcs in data*

# Sparsity Problem

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



■ Seen Arc:    See dependency in the training data
■ Seen Words:  Only see the words in training
■ Unseen:      At least one word not in training

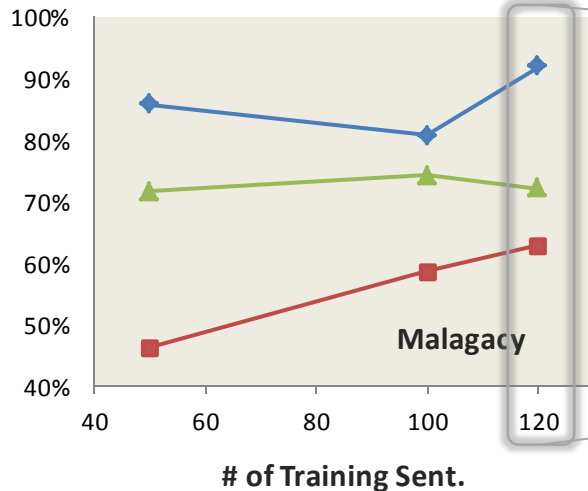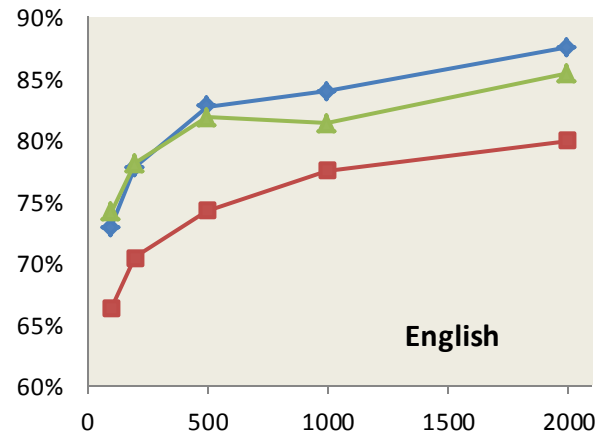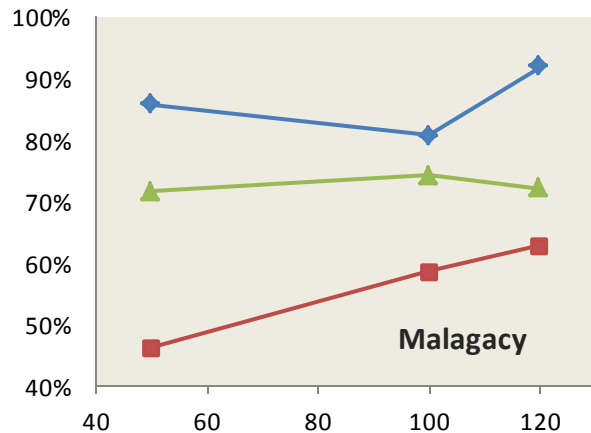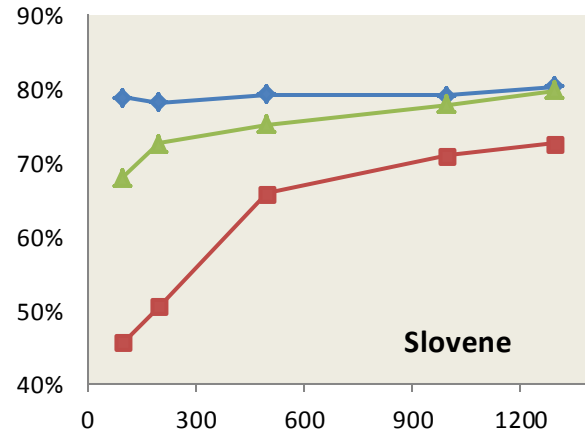a large portion of dependency arcs in test is unseen

22.56%
13.62%
63.82%

# Sparsity Problem

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

# Sparsity Problem

- Reason: feature weights are simply zero when the features are not seen

I    eat    a    cake    with    frosting
PRON   VB    DT    NN    IN    NN

| eat⊕IN⊕frosting | VB⊕IN⊕NN | eat⊕IN⊕NN | ... |
|---|---|---|---|
| 0   + | 1.1   + | 0.9   = | 2.0 |

unseen in training data

I    eat    a    cake    with    frosting
PRON   VB    DT    NN    IN    NN

| cake⊕IN⊕frosting | NN⊕IN⊕NN | cake⊕IN⊕NN | ... |
|---|---|---|---|
| 0   + | 1.2   + | 0.5   = | 1.7 |

# 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
  - Propagate 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

# Motivating Example

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



*Dependency Arc*

{ eat⊕apple, eat⊕NN, VB⊕apple, VB⊕NN }

*Feature Strings*

|       | eat | apple | banana | VB | NN |
|-------|-----|-------|--------|----|----|
| eat   |     | 1     |        |    | 1  |
| apple |     |       |        |    |    |
| banana|     |       |        |    |    |
| VB    |     | 1     |        |    | 1  |
| NN    |     |       |        |    |    |

*Feature Matrix*

$\otimes$

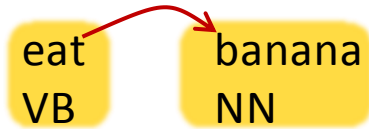|     |     |     |     |     |
|-----|-----|-----|-----|-----|
| ... | 2   | ... | ... | 4   |
| ... | 0   | 0   | ... | ... |
| ... | 0   | 0   | ... |     |
| ... | 1   | 0.9 | ... | 5   |
| ... | 0.1 | 0.1 | ... | ... |

*Parameter Matrix*

$= 12$

# Motivating Example

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

eat
VB

banana
NN

*Dependency Arc*

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

*Feature Strings*

**unseen entry**

|       | eat | apple | banana | VB | NN |
|-------|-----|-------|--------|----|----|
| eat   |     |       | 1      |    | 1  |
| apple |     |       |        |    |    |
| banana|     |       |        |    |    |
| VB    |     |       | 1      |    | 1  |
| NN    |     |       |        |    |    |

*Feature Matrix*

$\otimes$

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

*Parameter Matrix*

# Motivating Example

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



*Dependency Arc*

{ eat⊕banana, eat⊕NN, VB⊕banana, VB⊕NN }

*Feature Strings*

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

|  | 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*

10

# Motivating Example

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



*Dependency Arc*

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

*Feature Strings*



*Feature Matrix*

$\otimes$

*Parameter Matrix*

$= 11.9$

# Formulation (Simplified)

- Recall first-order decoding objective:

$$\widetilde{y}_i \;=\; \underset{y_i \in T(x_i)}{\operatorname{argmax}}\; S(y_i)$$

$$=\; \underset{y_i \in T(x_i)}{\operatorname{argmax}} \sum_{(h,c) \in y_i} s(h,c)$$

- Define score as matrix inner product:

$$s(h,c) = A \otimes \phi(h,c)$$

- Minimize the loss of training data:
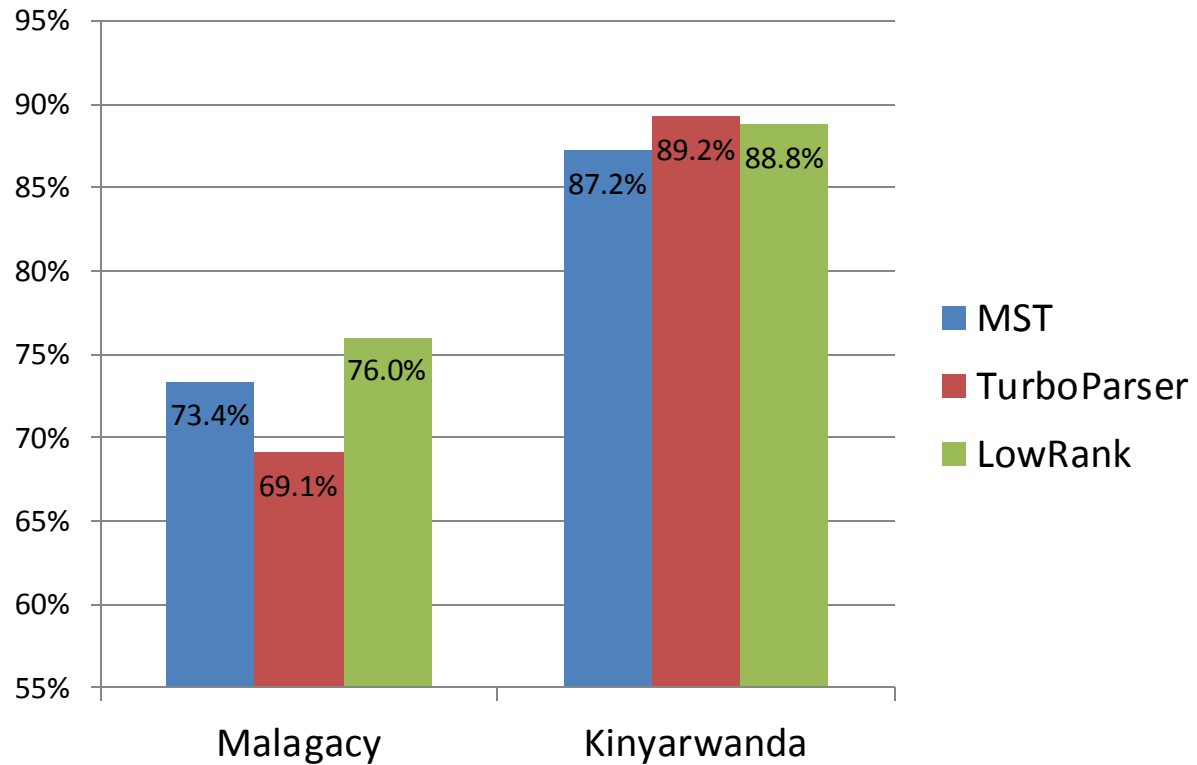
$$\mathcal{L}(D; A) = \frac{1}{N} \ell(x_i, \widehat{y}_i)$$

$$\text{s.t.} \quad \|A\|_* \leq C$$

<span style="color:red">Force A to be low-rank using nuclear norm constraint</span>

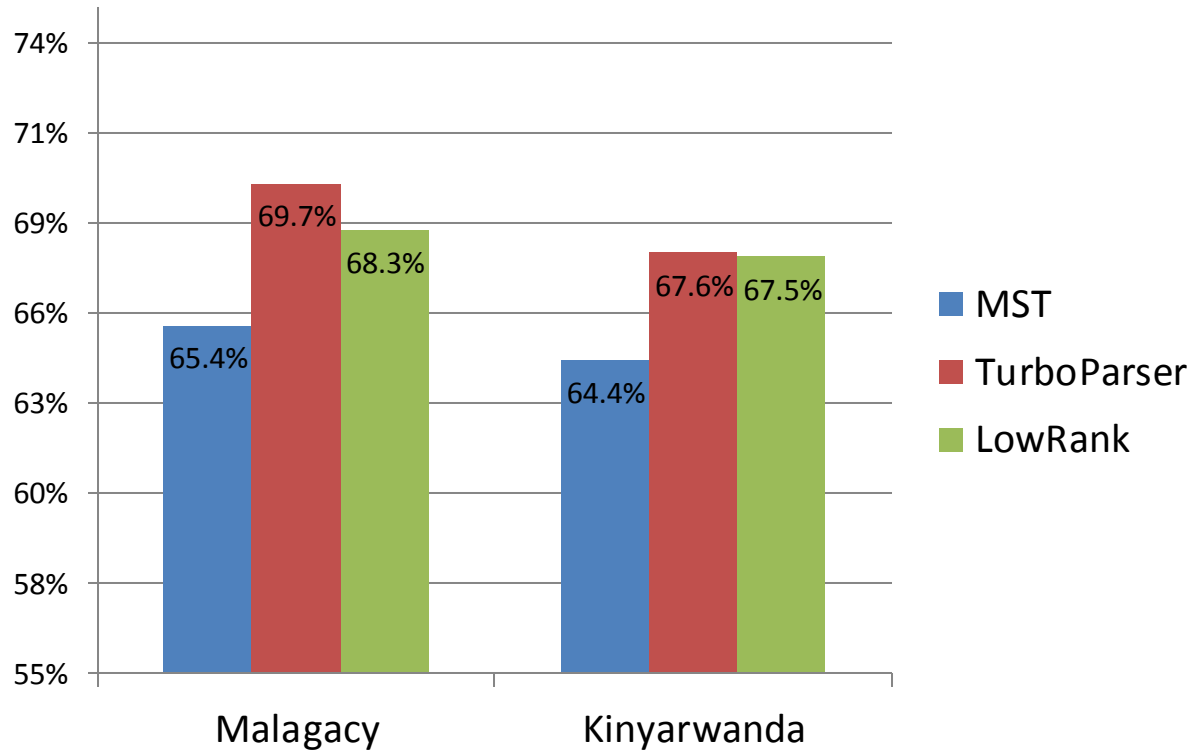  – online gradient descent algorithm available
  (Jaggi & Sulovsky, 2010) (Hazan, 2008)

12

# Results

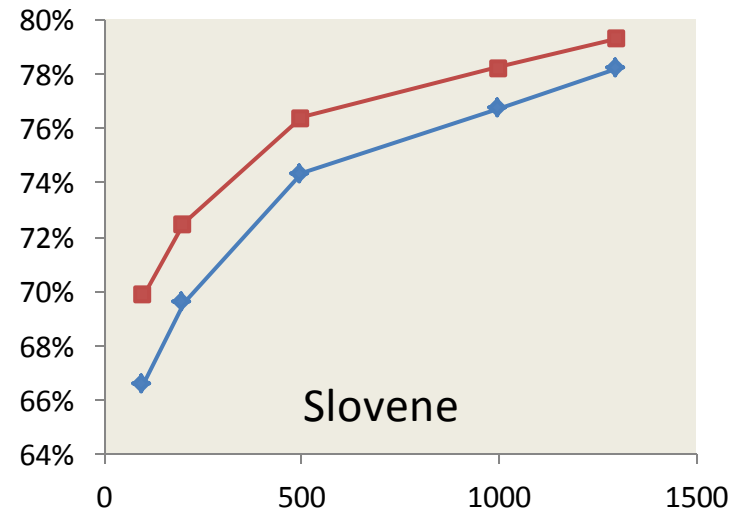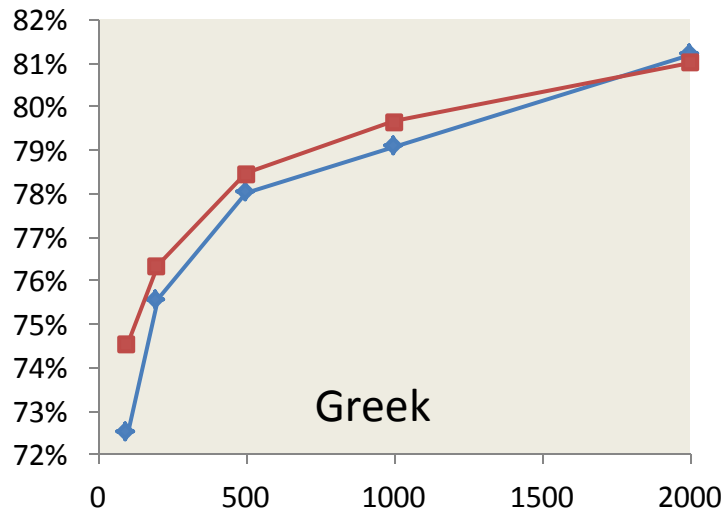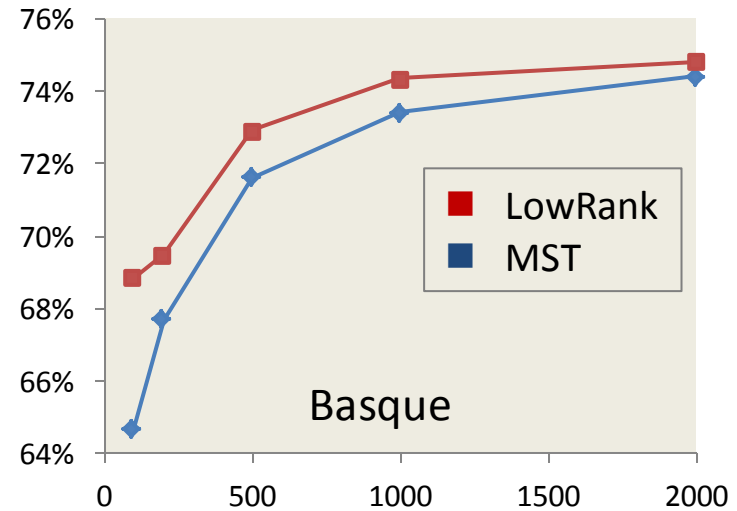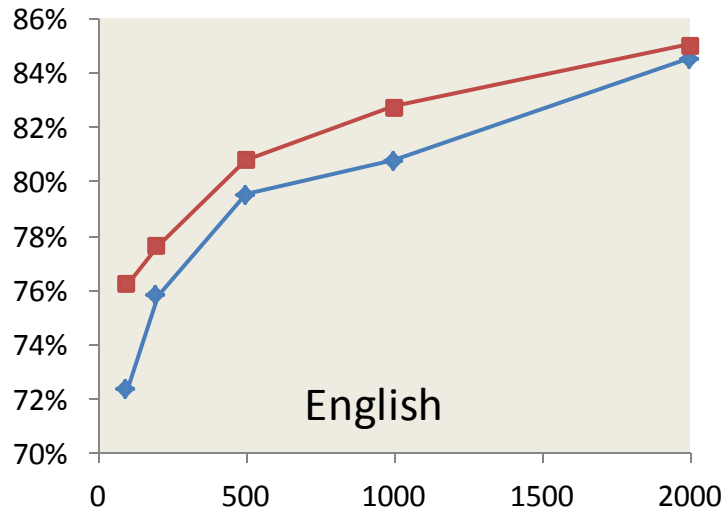- Models trained with gold POS tags

# Results

- Models trained with auto POS tags by TurboTagger

# Results

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

# Results

- Adding unsupervised word embeddings to English

|  | MST | LowRank | LowRank+wv |
|---|---|---|---|
| **100** | 72.4% | 76.3% | 76.6% (+0.3%) |
| **200** | 75.8% | 77.7% | 78.0% (+0.3%) |
| **500** | 79.5% | 80.8% | 81.4% (+0.6%) |
| **1000** | 80.8% | 82.8% | 82.8% (+0.0%) |
| **2000** | 84.5% | 85.1% | 85.8% (+0.7%) |