# High-order Low-rank Tensors for Semantic Role Labeling

Yuan Zhang, Tao Lei, Regina Barzilay

Lluís Màrquez, Alessandro Moschitti

**NAACL 2015**
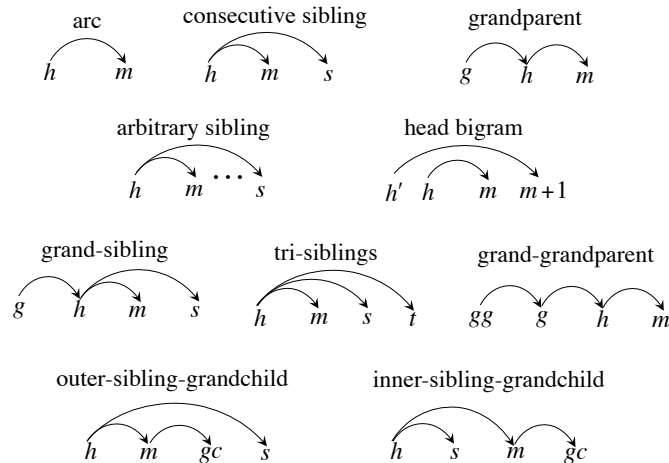
# Structural Prediction

- Traditional structural prediction requires huge feature engineering

Example: syntactic dependency parsing

*more than 10 groups of features*

*>100 feature templates*



- Problem: feature sparsity, hard to generalize to unseen data

# Structural Prediction

- Recent advance:

   learn low-dim. representations and their interactions
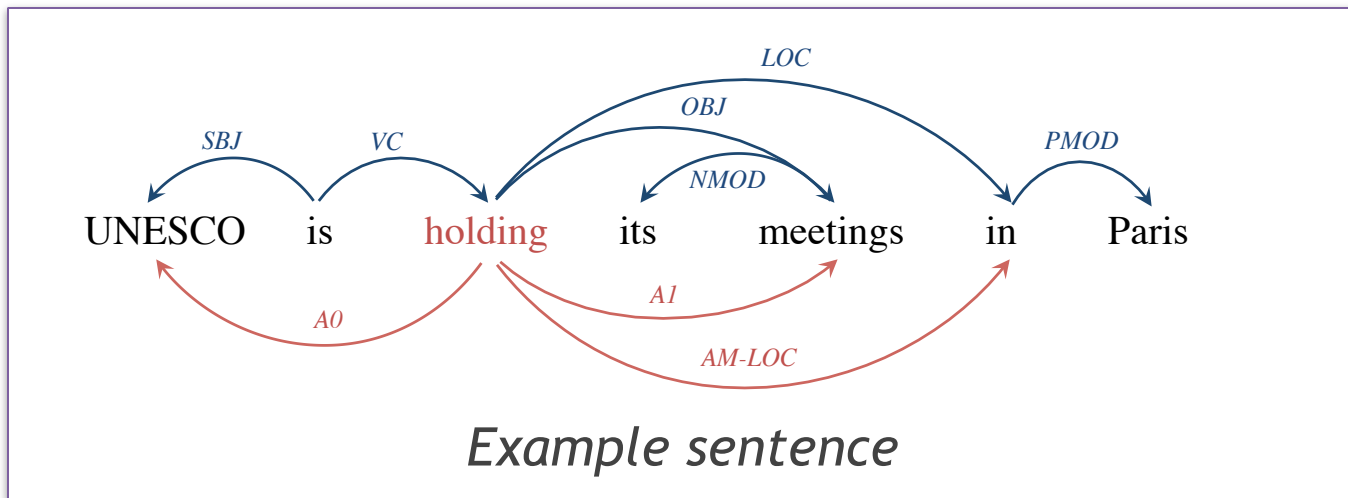   (compositions) to achieve better generalization

   ✦ Neural networks (Stenetorp 2013; Socher et al 2013;
      Chen and Manning 2014; Weiss et al 2015)

   ✦ Tensor factorization (Quattoni et al 2014; Lei et al 2014;
      Srikumar and Manning 2014)

- In this work,

   we extend our tensor factorization method to SRL
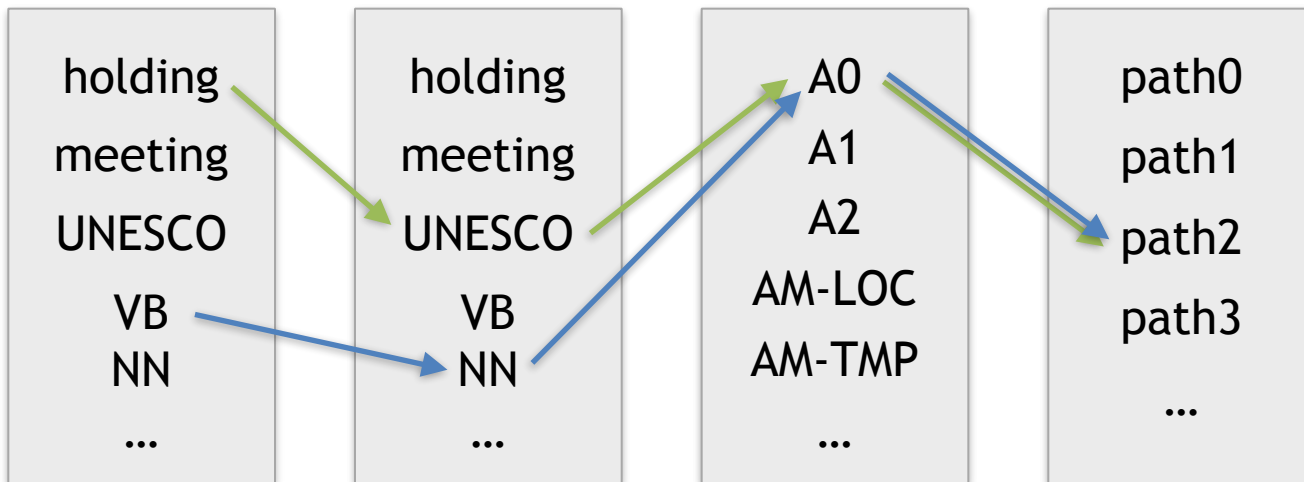
CSAIL

# Feature Construction in SRL

- Features defined over tuples ( *pred*, *arg*, *role*, *path* )

  - ✦ *pred* — predicate       holding
  - ✦ *arg* — argument       UNESCO
  - ✦ *role* — role label       *A0*
  - ✦ *path* — syntactic path



*Example sentence*

# Feature Construction in SRL

- Features defined over tuples ( *pred*, *arg*, *role*, *path* )

Selecting 1 up to 4 of them to construct features:

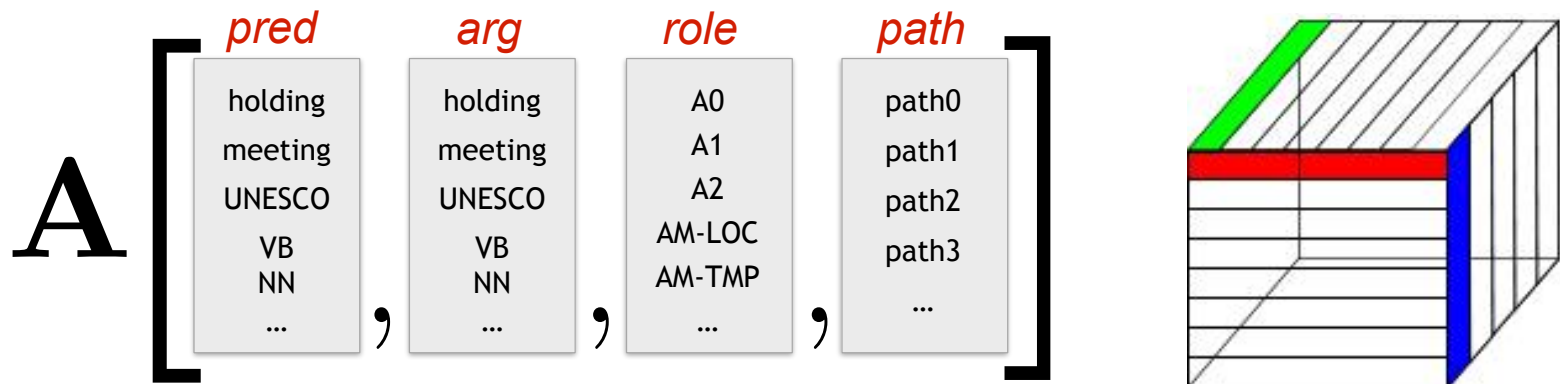| holding | holding | A0 | path0 |
| meeting | meeting | A1 | path1 |
| UNESCO | UNESCO | A2 | path2 |
| VB | VB | AM-LOC | path3 |
| NN | NN | AM-TMP | ... |
| ... | ... | ... | |

Each combination defines a feature

Needs to learn the corresponding feature weights (i.e. parameters)

# A Tensor View of the Parameters

- Parameters of feature combinations indexed by a 4-way tensor:
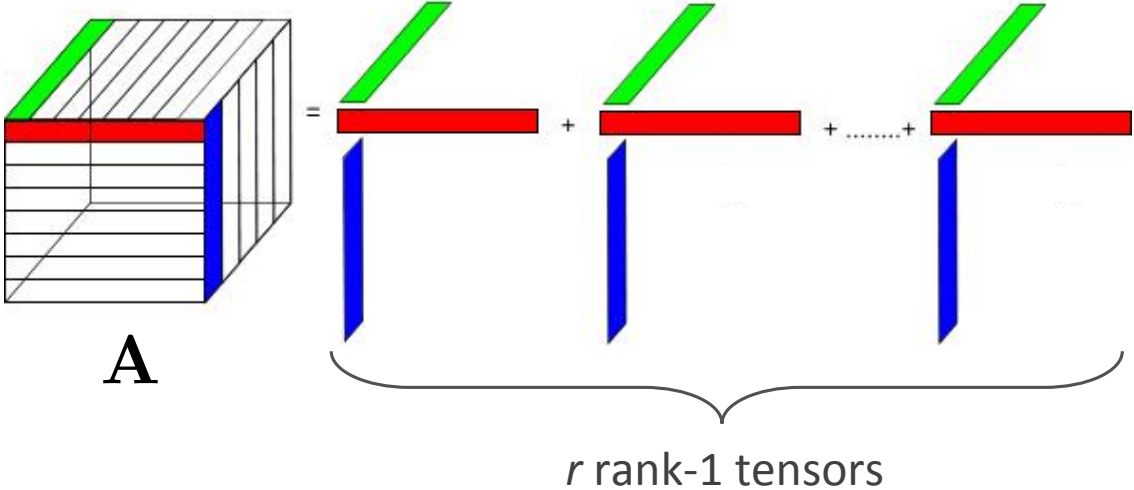


*Entries of A stores the feature weights*

# Avoid Explosion via Low-rank

- Learn a low-rank factorization of A, optimized for parsing

$$A = \sum_{i=1}^{r} P(i) \otimes Q(i) \otimes R(i) \otimes S(i)$$



A

*r* rank-1 tensors

*\* here we use 3-way tensor for better visualization*

# Online Learning

- Adopt standard mar-margin framework

$$\forall \mathbf{z}_{\text{sem}} \in \mathcal{Z}(\hat{\mathbf{x}}, \hat{\mathbf{y}_{\text{syn}}}) :$$
$$S_{sem}(\hat{\mathbf{x}}, \hat{\mathbf{y}_{\text{syn}}}, \hat{\mathbf{z}_{\text{sem}}}) \geq S_{sem}(\hat{\mathbf{x}}, \hat{\mathbf{y}_{\text{syn}}}, \mathbf{z}_{\text{sem}}) + cost(\hat{\mathbf{z}_{\text{sem}}}, \mathbf{z}_{\text{sem}})$$

*score of gold structure*     *score of pred. structure*     *margin*

Optimize parameters to satisfy this as much as possible

- Jointly update all parameter matrices via a new modified version of passive-aggressive algorithm

$$\Delta \boldsymbol{\theta} = \max \left\{ C, \ \frac{loss(\boldsymbol{\theta})}{\|g\boldsymbol{\theta}\|^2} \right\} \ g\boldsymbol{\theta}$$

CSAIL

# Tensor Initialization

- Performance can be impacted by initial values of P,Q,R,S

- Basic initialization steps:

(i) learn a traditional model, obtain sparse subset of parameter values

(ii) store the values as a sparse tensor *T*

(iii) find a low-rank approximation of T

$$\min_{P,Q,R,S} \|T - \sum_i P(i) \otimes Q(i) \otimes R(i) \otimes S(i)\|_2^2$$

# Tensor Initialization

- Performance can be impacted by initial values of P,Q,R,S

- Basic initialization steps:

> (i) learn a traditional model, obtain sparse subset of parameter values

> (ii) store the values as a sparse tensor *T*

> (iii) find a low-rank approximation of T
>
> $$\min_{P,Q,R,S} \| T - \sum_i P(i) \otimes Q(i) \otimes R(i) \otimes S(i) \|^2$$

In our previous work (Lei et al 2014), we use SVD initialization, which doesn't apply here

CSAIL

# Iterative Power Method for Initialization

- Approximately find one component — P(i), Q(i), R(i) and S(i) using an iterative algorithm, one by one

1: Randomly initialize four unit vectors $p$, $q$, $r$ and $s$
2: $T' = T - \sum_j P(j) \otimes Q(j) \otimes R(j) \otimes S(j)$
3: **repeat**
4:    $p = \langle T', -, q, r, s \rangle$ and normalize it
5:    $q = \langle T', p, -, r, s \rangle$ and normalize it
6:    $r = \langle T', p, q, -, s \rangle$ and normalize it
7:    $s = \langle T', p, q, r, - \rangle$
8:    $norm = \|s\|_2^2$
9: **until** $norm$ converges
10: $P(i) = p$ and $Q(i) = q$
11: $R(i) = r$ and $S(i) = s$
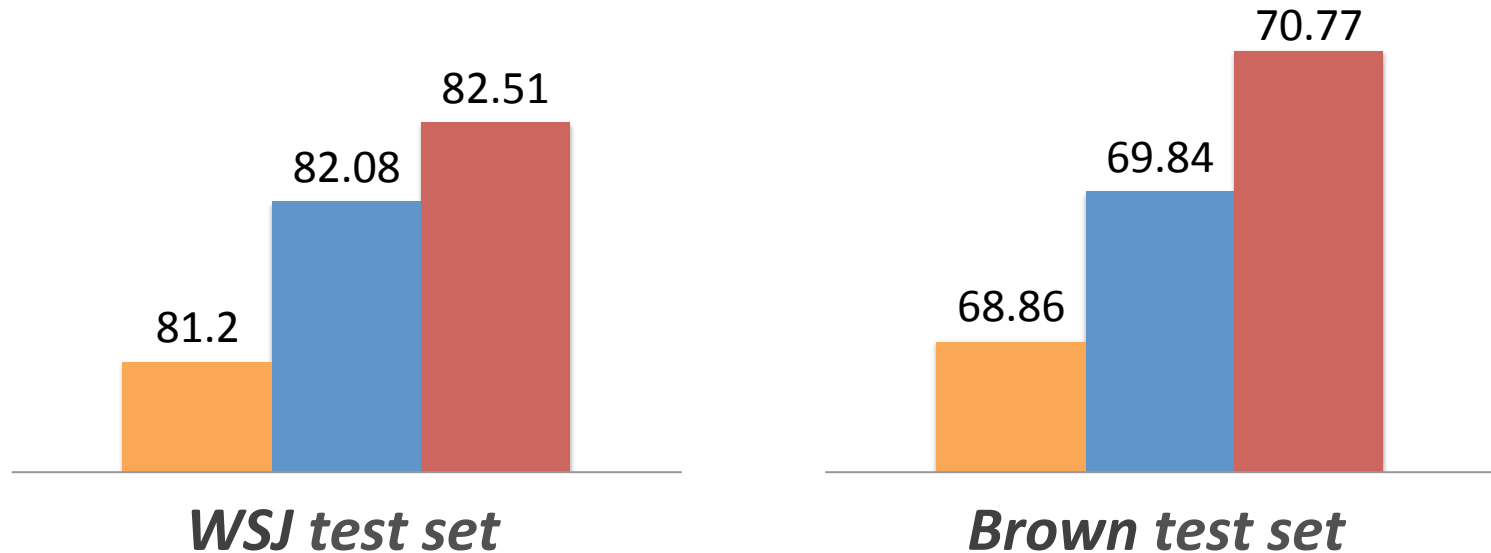
Optimize one vector while fixing the other three

# Experimental Setup

- Decoding:  weighted bipartite assignment (Lluís et al. 2013)

- Dataset: CoNLL-2009 joint syntactic and semantic parsing

- Features:

    a traditional set of 14 templates (Johansson, 2009)
    + our tensor component

- Baselines:

    best systems participated CoNLL-2009 and their
    improved versions
    (Che et al., 2009; Zhao et al., 2009; Bjorkelund et al., 2010; Roth and Woodsend, 2014)

    All explored much richer feature sets,  language-specific tuning and system combination
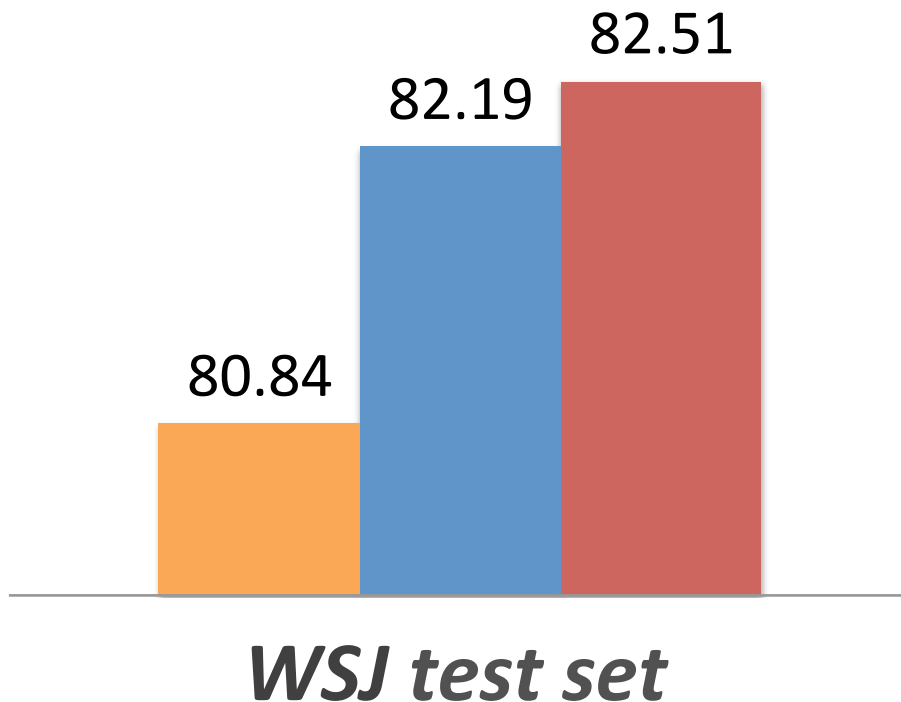
CSAIL

# Result on English



Legend: ■ CoNLL 2nd ■ CoNLL 1st ■ Our system

WSJ test set: CoNLL 2nd = 81.2, CoNLL 1st = 82.08, Our system = 82.51

Brown test set: CoNLL 2nd = 68.86, CoNLL 1st = 69.84, Our system = 70.77

outperforms best single system (w/o reranking) with statistical significance

# 3-way vs. 4-way tensor

3-way tensor by merging "role" and "path" into one mode



82.51

82.19

80.84

- basic features
- +3-way tensor
- +4-way tensor

***WSJ test set***

# Random vs. PM Initialization

82.51

81.63

80.84

***WSJ test set***

- basic features

- random init.

- power method init

CSAIL

# Overall Improvement

| Dataset | w/ tensor | w/o tensor |
|---------|-----------|------------|
| English | 82.51 | 80.84 |
| Catalan | 74.67 | 71.86 |
| Chinese | 69.16 | 68.43 |
| German | 76.94 | 74.03 |
| Spanish | 75.58 | 72.85 |
| **Average** | 75.77 | 73.60 |

Adding tensor component leads to > 2% absolute gain in F-score

CSAIL

# Thank you!

- RBG dependency parser
  https://github.com/taolei87/RBGParser

- Semantic role labeling parser
  https://github.com/taolei87/SRLParser

# Overall Improvement

| Dataset | w/ tensor | w/o tensor | CoNLL-1 (Zhao et al) |
|---------|-----------|------------|----------------------|
| English | **82.51** | 80.84 | 82.08 |
| Catalan | 74.67 | 71.86 | **76.78** |
| Chinese | **69.16** | 68.43 | 68.52 |
| German | **76.94** | 74.03 | 74.65 |
| Spanish | 75.58 | 72.85 | **77.33** |
| **Average** | 75.77 | 73.60 | 75.84 |