

# Using Large Monolingual and Bilingual Corpora to Improve Coordination Disambiguation

Shane Bergsma

David Yarowsky

Kenneth Church

Department of Computer Science and HLTCOE, Johns Hopkins University

## Resolving Coordination Ambiguity

**Problem:** Which words are being linked by a conjunction?  
**Specifically:** Coordination in complex Noun Phrases (NPs)

[ **dairy** and **meat** ] production  
**asbestos** and [ **polyvinyl chloride** ]

Interpretation: yes: [ **dairy** production ] and [ **meat** production ]  
 no: [ **asbestos** chloride ] and [ **polyvinyl** chloride ]

Hard problem!  
 Treebank-trained  
 parsers fail on it

Why? Could help  
 Information  
 Retrieval

Could also help  
 syntactic Machine  
 Translation

Like PP-attachment, etc., **the specific lexical items matter**  
 • Both sequences above have part-of-speech tags: “NN and NN NN”  
 • We need to **look beyond labeled data for the key information**  
 e.g. paraphrase (on the web) diagnostic (Nakov & Hearst, 2005)

## Size Matters for Resolving Coordination

**WSJ portion of Penn Treebank**  
 • 1 MILLION words  
 (labeled data)

Marcus et al. (1993)  
 We use: Vadas & Curran (2007) annotations, giving syntax of all NPs in WSJ portion of Penn Treebank

**Bitexts**  
 • 1 BILLION words  
 (bilingual data)

Koehn (2005): Europarl  
 Callison-Burch et al. (2010): WMT 2010  
 We use: English-to-{Czech, Danish, German, Greek, Spanish, Finnish, French, Italian, Dutch, Portuguese, Swedish} bilingual data

**Web text (N-grams)**  
 • 1 TRILLION words  
 (unlabeled data)

Brants & Franz (2006): Google N-gram Corpus  
 We use Google N-grams V2 (Lin et al., 2010)

## Better Together: Mono and Bilingual Data

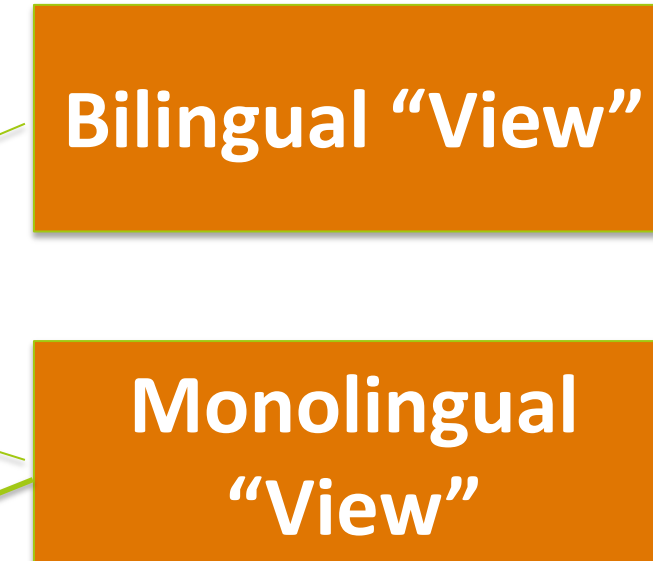
Phrase	Corpus Evidence (English and Foreign)	Pattern
<b>dairy</b>	English: production of <b>dairy</b> and <b>meat</b>	3 of 1 and 2
<b>1 and</b>	English: <b>dairy</b> production and <b>meat</b> production	1 3 and 2 3
<b>meat</b>	English: <b>meat</b> and <b>dairy</b> production	2 and 1 3
<b>2</b>	Spanish: producción láctea y cárnica (English: production <b>dairy</b> and <b>meat</b> )	3 1 ... 2
<b>3</b>	Finnish: maidon- ja lihantuotantoon (English: <b>dairy</b> - and <b>meat</b> production)	1- ... 2 3
	French: production de produits laitiers et de viande (English: production of products <b>dairy</b> and of <b>meat</b> )	3 ... 1 ... 2
<b>asbestos</b>	English: <b>polyvinyl chloride</b> and <b>asbestos</b>	2 3 and 1
<b>1 and</b>	English: <b>asbestos</b> , and <b>polyvinyl</b> chloride	1 , ... 2 3
<b>polyvinyl</b>	English: <b>asbestos</b> and <b>chloride</b>	1 and 3
<b>2</b>	Portuguese: o <b>amianto</b> e o <b>cloreto</b> de <b>polivinilo</b> (English: the <b>asbestos</b> and the <b>chloride</b> of <b>polyvinyl</b> )	1 ... 3 ... 2
<b>3</b>	Italian: l' <b>asbesto</b> e il <b>polivinilcloruro</b> (English: the <b>asbestos</b> and the <b>polyvinylchloride</b> )	1 ... 2 3

## Exploiting Unlabeled Text

Given some labeled examples of the two types of coordination:

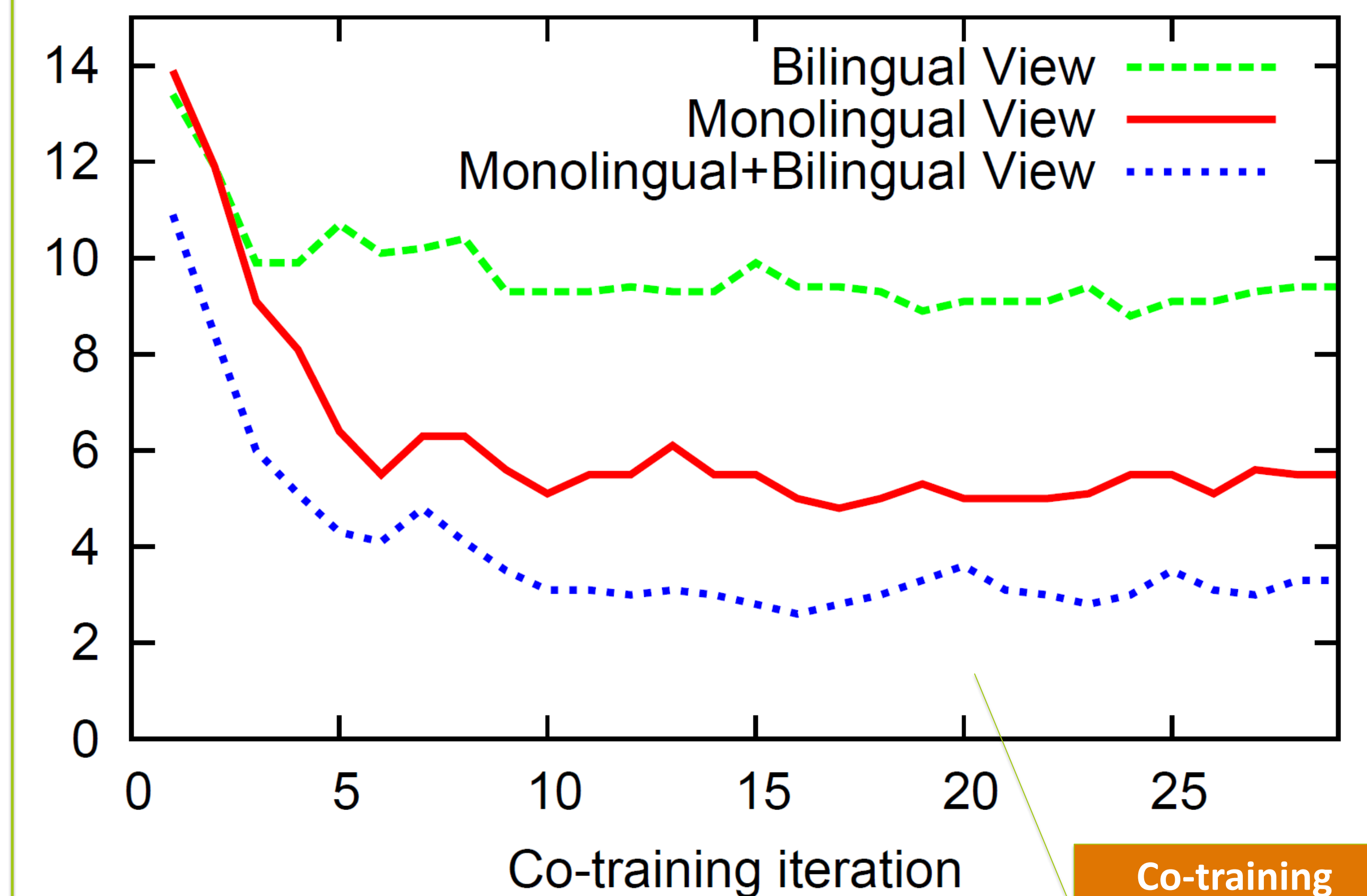
coal and steel money  
 insurrection and regime change  
 North and South Carolina  
 business and computer science  
 the Bosphorus and Dardanelles straits  
 pollution and transport safety  
 rocket and mortar attacks  
 the environment and air transport

- Create features for **counts of patterns** in:
  - a) **Bitexts** (Foreign Patterns)
  - b) **Web text** (English Patterns)
- Also create **binary features** for words/tags
- Train a logistic regression classifier to classify the coordination type
- Iteratively **co-train** using the Monoclassifier to label new examples for the bilingual classifier, and vice versa



## Results on Europarl Examples

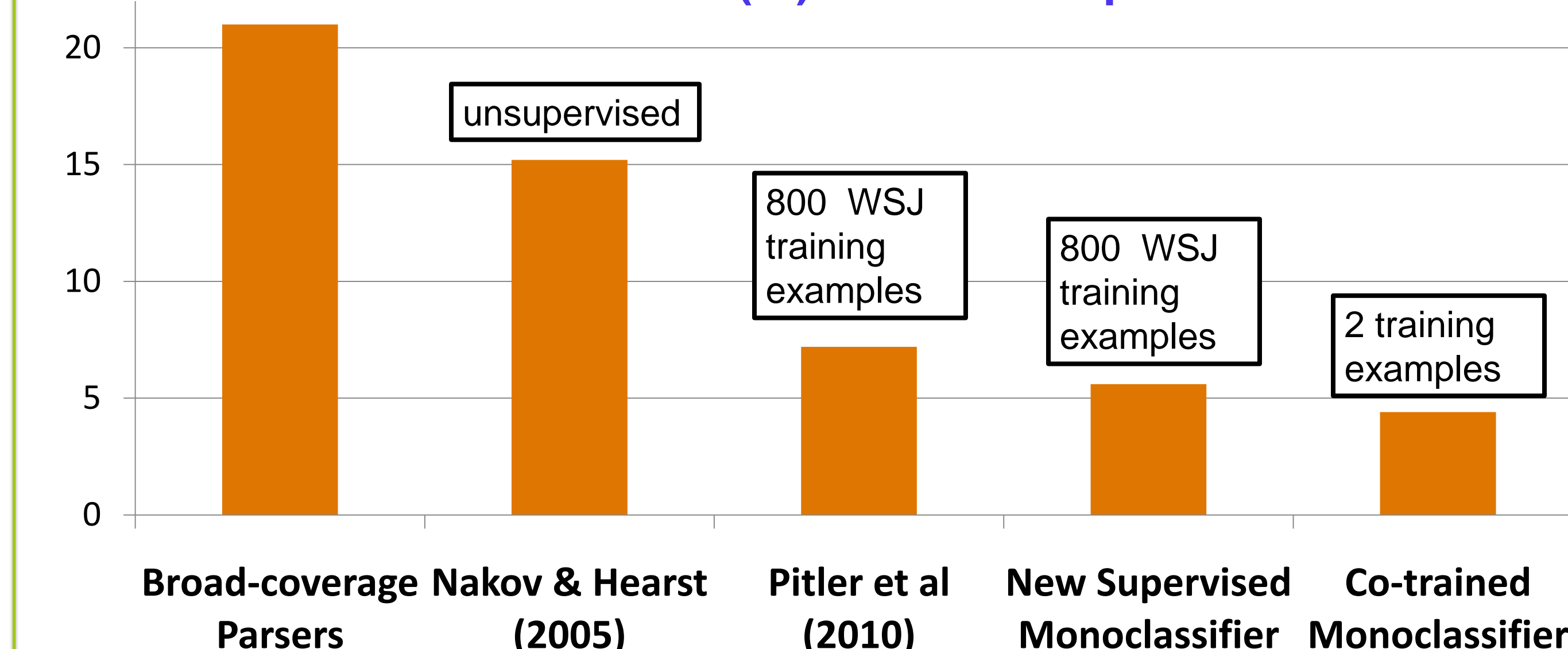
Error rate (%) of Co-trained Classifiers on Europarl test data



Co-training Works!

## Results on WSJ Examples

Error rate (%) on WSJ Corpus



• **Co-trained classifier using only monolingual features advances the state-of-the-art on WSJ data**

Check out our new annotated data and evaluation scripts at:  
<http://www.clsp.jhu.edu/sbergsma/coordNP.ACL11.zip>