

# Faster Parsing and Supertagging Model Estimation

Jonathan K. Kummerfeld<sup>a</sup> James R. Curran<sup>a</sup> Jessika Roesner<sup>b</sup>

School of Information Technologies<sup>a</sup>  
University of Sydney  
Australia

{jkum0593,james}@it.usyd.edu.au

Department of Computer Science<sup>b</sup>  
University of Texas at Austin  
USA

jessi@mail.utexas.edu

ALTW 2009



# Motivation – Parsing

Syntactic information is crucial for many tasks in NLP, such as QA and MT, but parsers are slow:

- State-of-the-art, usually  $< 1$  sentence / sec
- Fastest state-of-the-art,  $< 50$  sentences / sec

Far too slow to process the data available:

- $> 1,000,000,000,000$  words of English online
- More coming





# Tagging and Parsing

One claims he is pro – choice





# Part of Speech Tagging

One claims he is pro – choice  
*NN VBZ PRP VBZ JJ*





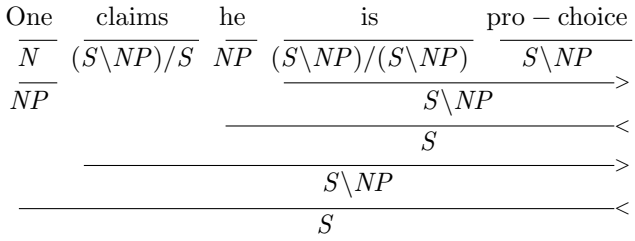
# Combinatory Categorical Grammar (CCG) – Supertagging

One	claims	he	is	pro – choice
$\overline{N}$	$(S \backslash NP) / S$	$\overline{NP}$	$\overline{(S \backslash NP) / (S \backslash NP)}$	$\overline{S \backslash NP}$





# Combinatory Categorical Grammar (CCG) – Parsing





# Supertagging Ambiguity

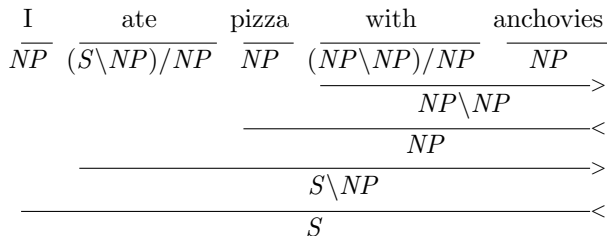
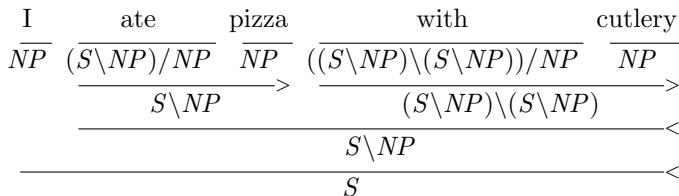
I ate pizza with cutlery

I ate pizza with anchovies





# Supertagging Ambiguity









# Motivation – Parsing

The key idea behind the speed of the fastest parsers today is to shift work from parsing to tagging:

For  $n$  words, each with  $k$  tags

- Tagging –  $O(nk)$
- Parsing –  $O(n^3k^2)$





# Outline

## Core Idea

- Provide fewer tags, but still include the tags the parser would have used anyway

## Implementation

- Perceptron Algorithms
- Parallelisation

## Results

- Modified rule usage
- Training data type and volume
- Algorithm comparison
- Feature extension





# Ideal World

$$\frac{\text{One}}{N} \quad \frac{\text{claims}}{(S \setminus NP) / S} \quad \frac{\text{he}}{NP} \quad \frac{\text{is}}{(S \setminus NP) / (S \setminus NP)} \quad \frac{\text{pro - choice}}{S \setminus NP}$$




# Current World – Problem

One      claims      he      is      pro – choice

$\frac{N}{N}$      $\frac{(S \setminus NP)/NP}{(S \setminus NP)/NP}$      $\frac{NP}{NP}$      $\frac{(S \setminus NP)/(S \setminus NP)}{(S \setminus NP)/(S \setminus NP)}$      $\frac{S \setminus NP}{S \setminus NP}$





# Current World – Solution

One	claims	he	is	pro – choice
$\frac{N/N}{N}$	$\frac{(S\backslash NP)/NP}{N}$	$\frac{NP}{NP}$	$\frac{(S\backslash NP)/(S\backslash NP)}{(S\backslash NP)/NP}$	$\frac{S\backslash NP}{(S\backslash NP)\backslash(S\backslash NP)}$
$(S/S)/(S/S)$			$(S\backslash NP)/(S\backslash NP)$	$(S\backslash NP)/S$
				$N$
				$(S\backslash NP)/PP$
				$(S\backslash NP)/NP$
				$N/N$
				$(S\backslash NP)/(S\backslash NP)$





# Adaptive Supertagging

One	claims	he	is	pro – choice
$\overline{N/N}$	$\overline{N}$	$\overline{NP}$	$\overline{(S\backslash NP)/(S\backslash NP)}$	$\overline{S\backslash NP}$
			$(S\backslash NP)/NP$	$(S\backslash NP)/PP$
			$(S\backslash NP)/(S\backslash NP)$	$(S\backslash NP)/NP$
				$N/N$

How do we teach the supertagger to produce these tags?  
Use the parser!





# Outline

## Core Idea

- Provide fewer tags, but still include the tags the parser would have used anyway

## Implementation

- **Perceptron Algorithms**
- **Parallelisation**

## Results

- Modified rule usage
- Training data type and volume
- Algorithm comparison
- Feature extension







# Implementation

Component	Initial System	Additions
Statistical Feature Extraction	3 Types Single thread	+9 Types Parallel
Parameter Estimation	BFGS, GIS Single thread	AP, MIRA Parallel





# Implementation – Extra Constraint

Added a constraint that only allows Backward Composition to occur if both children are type raised





# Implementation – AP and MIRA

Algorithm	Training Time (sec)		
	40k	80k	440k
GIS	7,200	14,000	*
BFGS	6,300	13,000	*
AP	76	160	950
MIRA	96	200	1,200

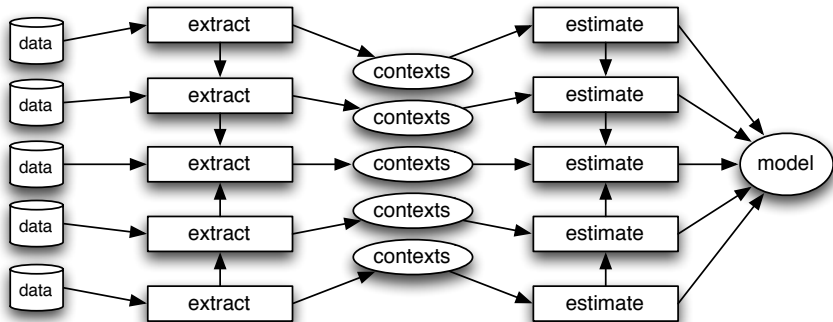




# Implementation – Initial System



# Implementation – Parallelised



# Implementation – Parallelised Weight Estimation

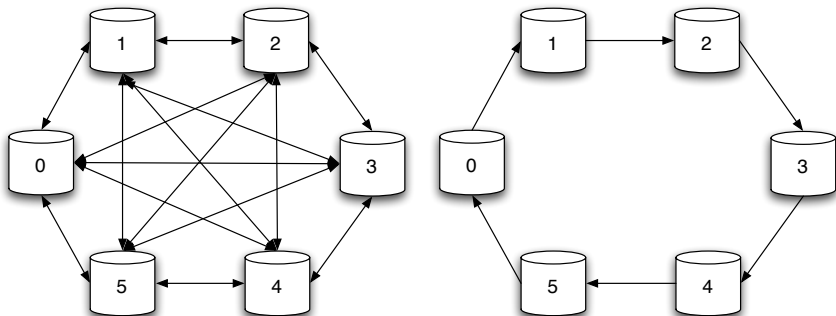


Figure: Information flow for parallel model estimation





# Outline

## Core Idea

- Provide fewer tags, but still include the tags the parser would have used anyway

## Implementation

- Perceptron Algorithms
- Parallelisation

## Results

- **Modified rule usage**
- **Training data type and volume**
- **Algorithm comparison**
- **Feature extension**





# Extra Constraint on Rule Application

Parser	F-score (%)	Speed (sent / sec)
C&C 1.02	83.22	31.7
Modified	83.41	47.8







# Plan

- Acquire a large set of unannotated data – Wikipedia
- Parse the corpus
- Retrain the supertagger, using the parsed sentences

## Variations

- Amount of data
- Estimation algorithms
- Feature set



# Training Data Type and Volume

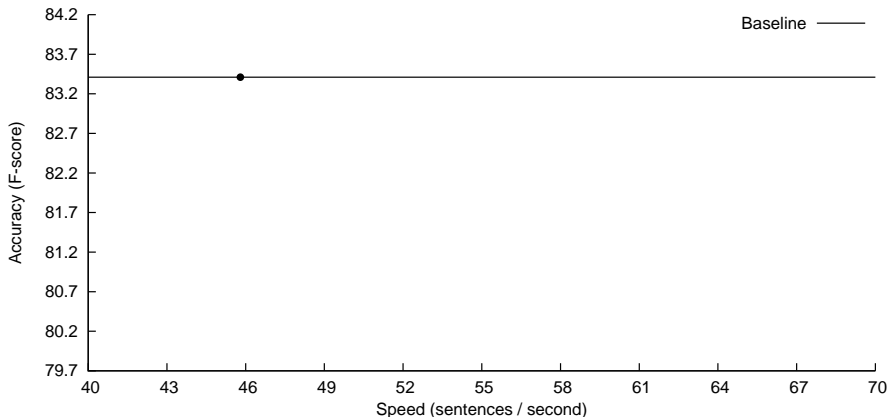


Figure: Evaluation on the Wall Street Journal



# Training Data Type and Volume

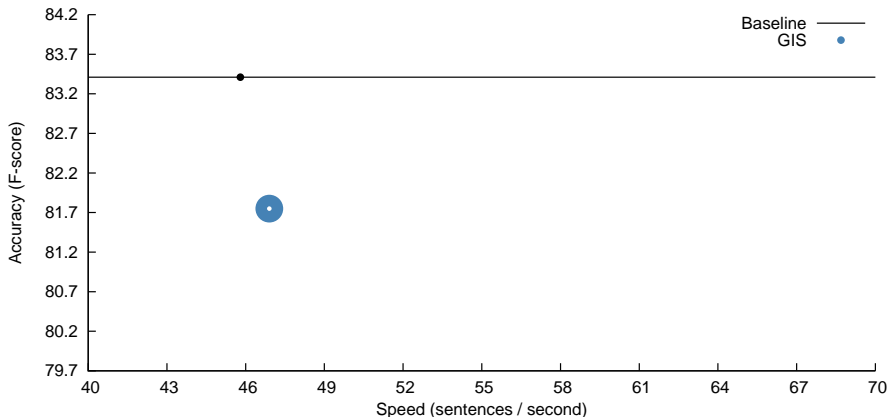


Figure: Evaluation on the Wall Street Journal



# Training Data Type and Volume

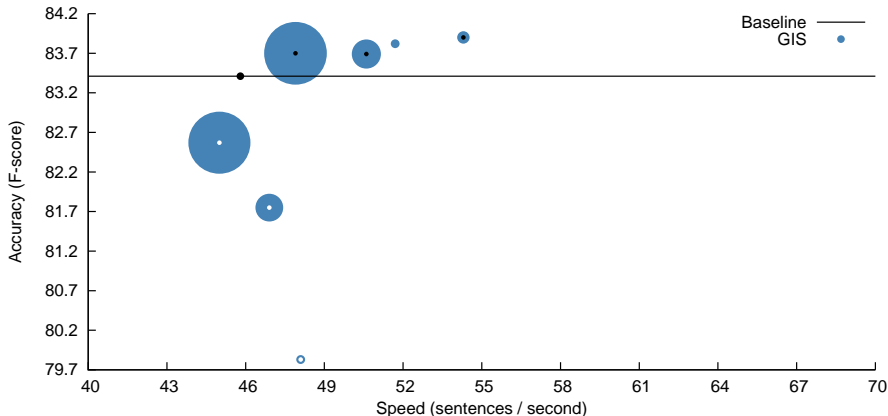


Figure: Evaluation on the Wall Street Journal



# Training Data Type and Volume

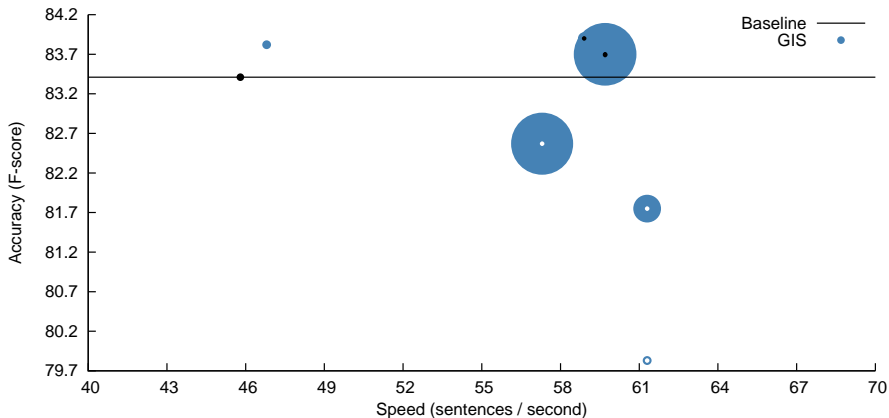


Figure: Evaluation on Wikipedia



# Algorithm Comparison

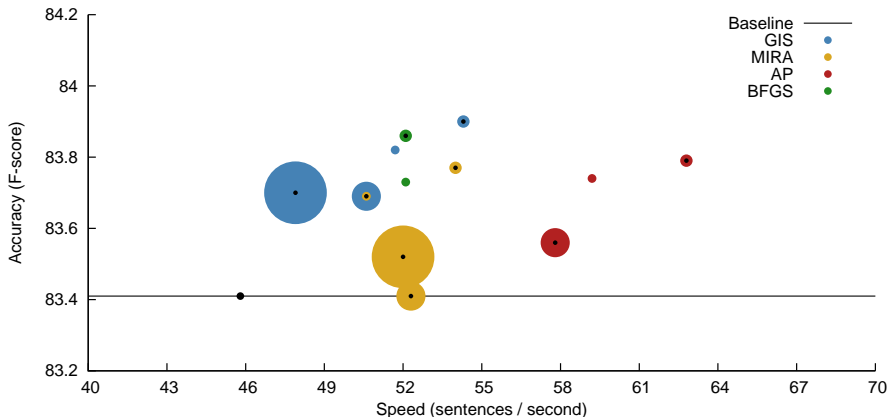


Figure: Evaluation on the Wall Street Journal



# Algorithm Comparison

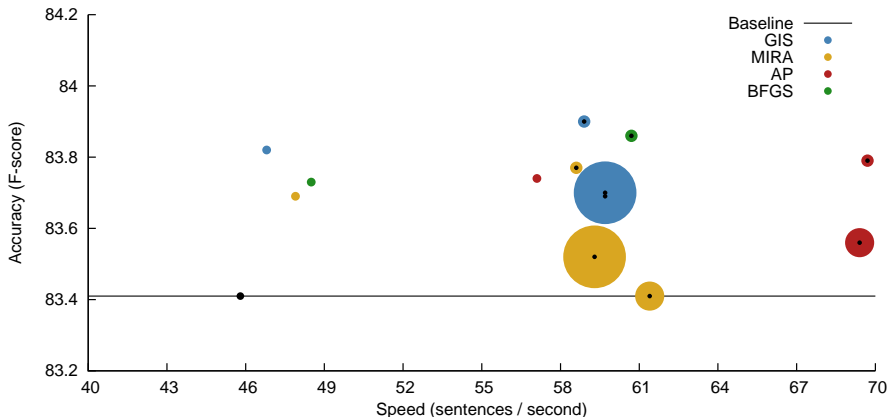


Figure: Evaluation on Wikipedia



# Feature Extension

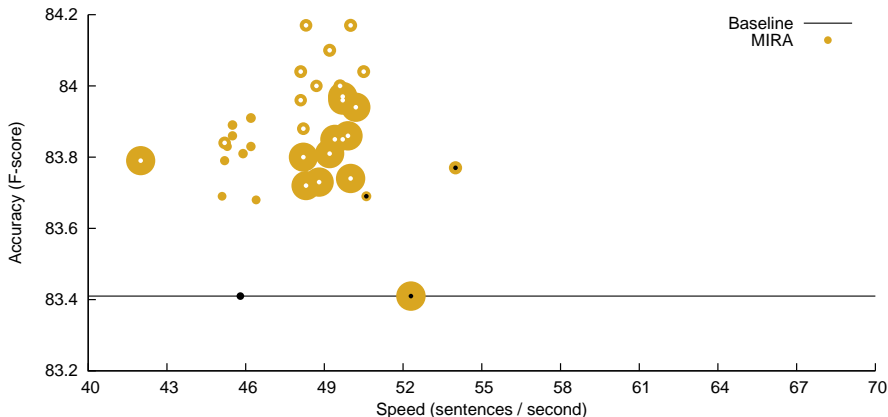


Figure: Evaluation on the Wall Street Journal





# Feature Extension

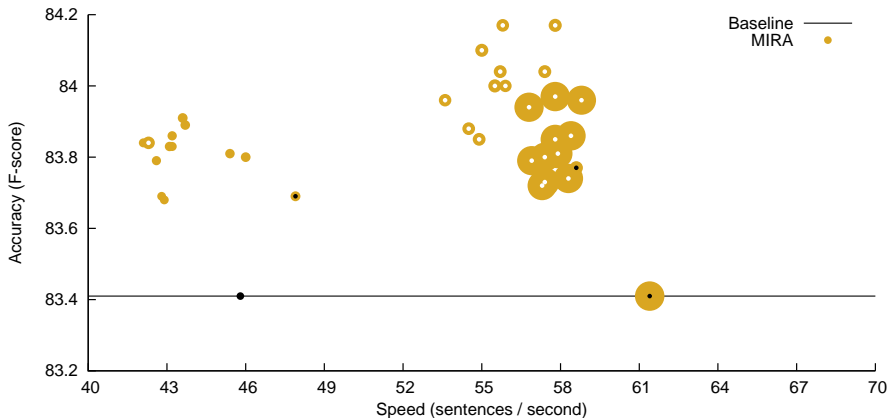


Figure: Evaluation on Wikipedia





# Future Work

- Other domains
- Expanded training sets
- Co-training
- Online learning





# Conclusion

Improved training:

- Enabled access to more text
- Constructed an effective source of more text

Improved parsing speed:

- Added an extra constraint on rule usage
- Trained models that are adapted to the parser

Improved parsing accuracy:

- Constructed statistical models using more evidence
- Expanded the set of statistical features





# Conclusion

Metric	Initial	Final	Ratio
<b>Training</b>			
Sentences	40k	80k	2
Time (secs)	6,300	160	1/40
<b>Accuracy</b>			
F-score (%)	83.22	83.79	n/a
<b>Speed</b>			
WSJ (sents / sec)	31.7	62.8	2.0
Wikipedia (sents / sec)	30.8	69.7	2.3





# Acknowledgements

- Johns Hopkins University, CLSP Summer Workshop
- Capital Markets Cooperative Research Centre Limited

