## Slide 1

# Natural Language Processing as Philology

David Smith
Department of Computer Science
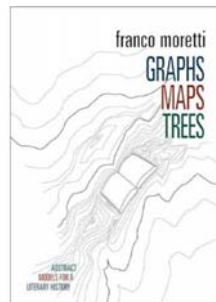UMass Amherst

## Slide 2



*Roman Jakobson*

Philology is the art of reading slowly.

## Slide 3



*Franco Moretti*

franco moretti
GRAPHS
MAPS
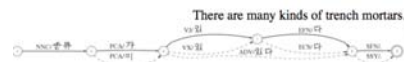TREES

ABSTRACT MODELS FOR A LITERARY HISTORY

## Slide 4

## Slide 5

# NLP Highlights

- Efficient algorithms for linguistic inference
  - Joint inference across many layers of language
- Adaptation to new languages and domains
- Inferring structure in large, noisy collections
  - Detecting text reuse and linkage
  - Inferring temporal sequence of events

## Slide 6

# Morphological Disambiguation

## Bare-Bones Dependency Structure

The computer knows you like your mother

## Syntax in Translation

Er wird in den Strassen wandern
*He will in the streets walk*

Google

He will walk in the streets

Er wird in den kleinen Strassen wandern
*He will in the small streets walk*

Google

He is in the small streets hike

## Who Did What To Whom?

Pierre  Vinken  ,  61  years  old  ,  will  join  the  board  as  a  nonexecutive  director

### PropBank join predicate

| ARG0 | ARG1 | ARG-PRD |
|------|------|---------|
| Vinken | board | director |

## SYNTAX AND PARSING

## Grammars and Trees

*Context-free grammars*

VP → V NP PP        0.0001

VP        0.0001

V   NP   PP

$\geq O(n^3)$

*Tree (substitution | insertion | adjoining) grammars*

VP        0.0001

VP   PP↓

V↓   NP↓        $\geq O(n^3), O(n^6)$

*Synchronous grammars*

S        S        0.0001
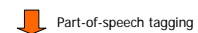
NP↓   VP        NP↓   VP

likes   NP↓     gusta   NP↓

$\geq O(n^6)$

Higher complexity of CCG, LFG, HPSG, Minimalist

## Dependency Parsing as Graph Inference

Raw sentence

He reckons the current account deficit will narrow to only 1.8 billion in September.

Part-of-speech tagging

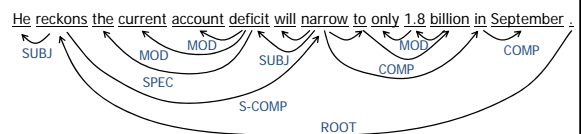POS-tagged sentence

He reckons the current account deficit will narrow to only 1.8 billion in September.
PRP  VBZ  DT  JJ  NN  NN  MD  VB  TO  RB  CD  CD  IN  NNP  .

Word dependency parsing

Word dependency parsed sentence

He reckons the current account deficit will narrow to only 1.8 billion in September .

SUBJ  MOD  MOD  SUBJ  MOD  COMP
SPEC  COMP
S-COMP
ROOT

slide adapted from Yuji Matsumoto

## How about structured outputs?

- Log-linear models great for n-way classification
- Also good for predicting sequences

find    preferred    tags

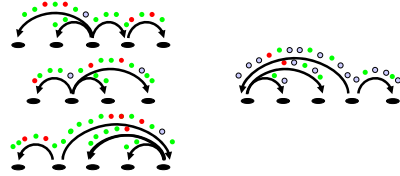but to allow fast dynamic programming, only use n-gram features

- Also good for dependency parsing

…find preferred links…

but to allow fast dynamic programming or MST parsing, only use single-edge features

---

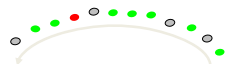## How about structured outputs?

…find preferred links…

but to allow fast dynamic programming or MST parsing, only use single-edge features

---

## Edge-Factored Parsers (McDonald et al. 2005)
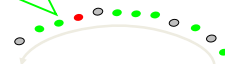
- Is this a good edge?

yes, lots of green …

Byl   jasný   studený   dubnový   den   a   hodiny   odbíjely   třináctou

"It was a bright cold day in April and the clocks were striking thirteen"

---

## Edge-Factored Parsers (McDonald et al. 2005)

- Is this a good edge?

jasný ← den
("bright day")

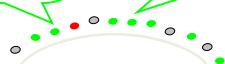Byl   jasný   studený   dubnový   den   a   hodiny   odbíjely   třináctou

"It was a bright cold day in April and the clocks were striking thirteen"

---

## Edge-Factored Parsers (McDonald et al. 2005)

- Is this a good edge?

jasný ← den
("bright day")

jasný ← N
("bright NOUN")

Byl   jasný   studený   dubnový   den   a   hodiny   odbíjely   třináctou
V      A        A          A         N    J     N        V         C

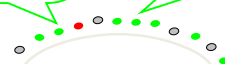"It was a bright cold day in April and the clocks were striking thirteen"

---

## Edge-Factored Parsers (McDonald et al. 2005)

- Is this a good edge?

jasný ← den
("bright day")

jasný ← N
("bright NOUN")

A ← N

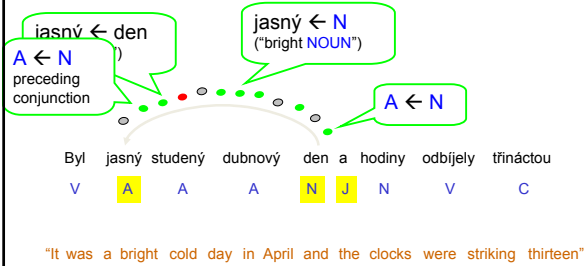Byl   jasný   studený   dubnový   den   a   hodiny   odbíjely   třináctou
V      A        A          A         N    J     N        V         C

"It was a bright cold day in April and the clocks were striking thirteen"
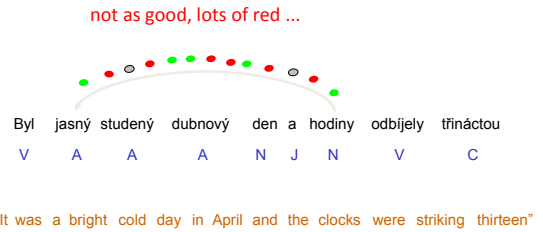
# Edge-Factored Parsers (McDonald et al. 2005)

- Is this a good edge?

jasný ← den

A ← N
preceding
conjunction

jasný ← N
("bright NOUN")

A ← N

| Byl | jasný | studený | dubnový | den | a | hodiny | odbíjely | třináctou |
|-----|-------|---------|---------|-----|---|--------|----------|-----------|
| V | A | A | A | N | J | N | V | C |

"It was a bright cold day in April and the clocks were striking thirteen"

---

# Edge-Factored Parsers (McDonald et al. 2005)

- How about this competing edge?

not as good, lots of red ...

| Byl | jasný | studený | dubnový | den | a | hodiny | odbíjely | třináctou |
|-----|-------|---------|---------|-----|---|--------|----------|-----------|
| V | A | A | A | N | J | N | V | C |

"It was a bright cold day in April and the clocks were striking thirteen"

---

# Edge-Factored Parsers (McDonald et al. 2005)

- How about this competing edge?

jasný ← hodiny
("bright clocks")

... undertrained ...

| Byl | jasný | studený | dubnový | den | a | hodiny | odbíjely | třináctou |
|-----|-------|---------|---------|-----|---|--------|----------|-----------|
| V | A | A | A | N | J | N | V | C |

"It was a bright cold day in April and the clocks were striking thirteen"

---

# Edge-Factored Parsers (McDonald et al. 2005)

- How about this competing edge?

jasný ← hodiny
("bright clocks")

... undertrained ...

jasn ← hodi
("bright clock,"
stems only)

| Byl | jasný | studený | dubnový | den | a | hodiny | odbíjely | třináctou |
|-----|-------|---------|---------|-----|---|--------|----------|-----------|
| V | A | A | A | N | J | N | V | C |
| byl | jasn | stud | dubn | den | a | hodi | odbí | třin |

"It was a bright cold day in April and the clocks were striking thirteen"

---

# Edge-Factored Parsers (McDonald et al. 2005)

- How about this competing edge?

jasný ← hodiny
("bright clocks")

... undertrained ...

jasn ← hodi
("bright clock,"
stems only)

$A_{plural}$ ← $N_{singular}$

| Byl | jasný | studený | dubnový | den | a | hodiny | odbíjely | třináctou |
|-----|-------|---------|---------|-----|---|--------|----------|-----------|
| V | A | A | A | N | J | N | V | C |
| byl | jasn | stud | dubn | den | a | hodi | odbí | třin |

"It was a bright cold day in April and the clocks were striking thirteen"

---

# Edge-Factored Parsers (McDonald et al. 2005)

- How about this competing edge?

jasný ← hodiny

A ← N
where N follows
a conjunction

jasn ← hodi
("bright clock,"
stems only)

$A_{plural}$ ← $N_{singular}$

| Byl | jasný | studený | dubnový | den | a | hodiny | odbíjely | třináctou |
|-----|-------|---------|---------|-----|---|--------|----------|-----------|
| V | A | A | A | N | J | N | V | C |
| byl | jasn | stud | dubn | den | a | hodi | odbí | třin |

"It was a bright cold day in April and the clocks were striking thirteen"

## Slide 1

# Edge-Factored Parsers (McDonald et al. 2005)

- Which edge is better?
  - "bright day" or "bright clocks"?

| Byl | jasný | studený | dubnový | den | a | hodiny | odbíjely | třináctou |
|-----|-------|---------|---------|-----|---|--------|----------|-----------|
| V | A | A | A | N | J | N | V | C |
| byl | jasn | stud | dubn | den | a | hodi | odbí | třin |

"It was a bright cold day in April and the clocks were striking thirteen"

## Slide 2

# Edge-Factored Parsers (McDonald et al. 2005)

our current weight vector

- Which edge is better?
- Score of an edge e = θ ⋅ **features**(e)
- Standard algos ➔ valid parse with max <u>total</u> score

| Byl | jasný | studený | dubnový | den | a | hodiny | odbíjely | třináctou |
|-----|-------|---------|---------|-----|---|--------|----------|-----------|
| V | A | A | A | N | J | N | V | C |
| byl | jasn | stud | dubn | den | a | hodi | odbí | třin |

"It was a bright cold day in April and the clocks were striking thirteen"

## Slide 3

# Edge-Factored Parsers (McDonald et al. 2005)

our current weight vector

- Which edge is better?
- Score of an edge e = θ ⋅ **features**(e)
- Standard algos ➔ valid parse with max <u>total</u> score

can't have both
(one parent per word)

can't have both
(no crossing links)

Can't have all three
(no cycles)

Thus, an edge may lose (or win) because of a consensus of <u>other</u> edges.

## Slide 4

# Finding Highest-Scoring Parse

- Convert to context-free grammar (CFG)
- Then use dynamic programming

The cat in the hat wore a stovepipe. ROOT

let's vertically stretch
this graph drawing

wore ← ROOT
cat
The in stovepipe
hat a
the

each subtree is a linguistic constituent
(here a <u>noun phrase</u>)

## Slide 5

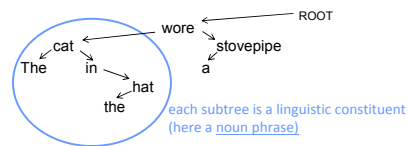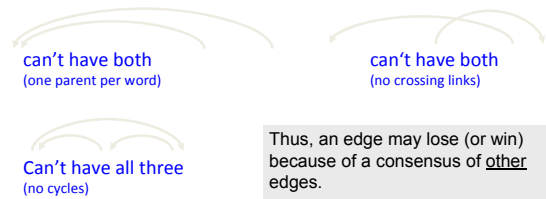# Finding Highest-Scoring Parse

- Convert to context-free grammar (CFG)
- Then use dynamic programming
  - CKY algorithm for CFG parsing is $O(n^3)$
  - Unfortunately, $O(n^5)$ in this case
    - to score "cat ← wore" link, not enough to know this is NP
    - must know it's rooted at "cat"
    - so expand nonterminal set by $O(n)$: {$NP_{the}$, $NP_{cat}$, $NP_{hat}$, ...}
      - so CKY's "grammar constant" is no longer constant ☺

cat stovepipe
The in a
hat
the

each subtree is a linguistic constituent
(here a <u>noun phrase</u>)

## Slide 6

# Finding Highest-Scoring Parse

- Convert to context-free grammar (CFG)
- Then use dynamic programming
  - CKY algorithm for CFG parsing is $O(n^3)$
  - Unfortunately, $O(n^5)$ in this case
  - Solution: Use a different decomposition (Eisner 1996)
    - Back to $O(n^3)$

wore ← ROOT
cat
The in stovepipe
hat a
the

each subtree is a linguistic constituent
(here a <u>noun phrase</u>)

## Finding Highest-Scoring Parse

- Convert to context-free grammar (CFG)
- Then use dynamic programming
  - CKY algorithm for CFG parsing is $O(n^3)$
  - Unfortunately, $O(n^5)$ in this case
  - Solution: Use a different decomposition (Eisner 1996)
    - Back to $O(n^3)$
- Can play usual tricks for dynamic programming parsing
  - Further refining the constituents or spans
    - Allow prob. model to keep track of even more internal information
  - A*, best-first, coarse-to-fine } require "outside" probabilities of constituents, spans, or links
  - Training by EM etc.

---

# Hard Constraints on Valid Trees

our current weight vector

- Score of an edge e = θ ( **features**(e)
- Standard algos ➔ valid parse with max <u>total</u> score

can't have both
(one parent per word)

can't have both
(no crossing links)

Can't have all three
(no cycles)

Thus, an edge may lose (or win) because of a consensus of <u>other</u> edges.

---

# Non-Projective Parses

ROOT    I    'll    give    a    talk    tomorrow    on    bootstrapping

subtree rooted at "talk"
is a discontiguous noun phrase

can't have both
(no crossing links)

The "projectivity" restriction.
Do we really want it?

---

# Non-Projective Parses

ROOT    I    'll    give    a    talk    tomorrow    on    bootstrapping

occasional non-projectivity in English

ROOT    ista    meam    norit    gloria    canitiem
        $that_{NOM}$    $my_{ACC}$    may-know    $glory_{NOM}$    $going\text{-}gray_{ACC}$

That glory may-know my going-gray
(i.e., it shall last till I go gray)

frequent non-projectivity in Latin, etc.

---

## Finding highest-scoring <u>non-projective</u> tree

- Consider the sentence "John saw Mary" (left).
- The Chu-Liu-Edmonds algorithm finds the maximum-weight spanning tree (right) – may be non-projective.
- Can be found in time $O(n^2)$.

root    9    10    saw    30    9    20    30    John    0    Mary    11    3

root    10    saw    30    John    30    Mary

Every node selects best parent
If cycles, contract them and repeat

---

### Summing over all non-projective trees

## Finding highest-scoring <u>non-projective</u> tree

- Consider the sentence "John saw Mary" (left).
- The Chu-Liu-Edmonds algorithm finds the maximum-weight spanning tree (right) – may be non-projective.
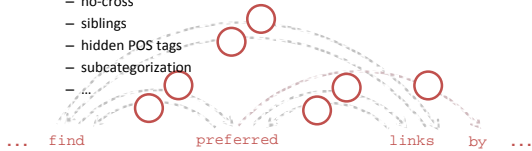- Can be found in time $O(n^2)$.

- How about total weight Z of all trees?
- How about outside probabilities or gradients?
- Can be found in time $O(n^3)$ by matrix determinants and inverses (Smith & Smith, 2007).

## Local factors for parsing

– So what factors shall we multiply to define parse probability?
  - Unary factors to evaluate each link in isolation
  - Global TREE factor to <u>require</u> that the links form a legal tree
    – this is a "hard constraint": factor is either 0 or 1
  - Second-order effects: factors on 2 variables
    – grandparent
    – no-cross
    – siblings
    – hidden POS tags
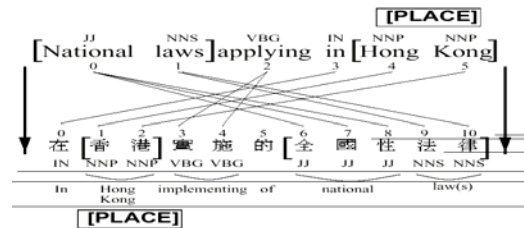    – subcategorization
    – ...

... find       preferred       links   by ...

## Future Opportunities

- Efficiently modeling more hidden structure
  - POS tags, link roles, secondary links (DAG-shaped parses)
- Beyond dependencies
  - Constituency parsing, traces, lattice parsing
- Beyond parsing
  - Alignment, translation
  - Bipartite matching and network flow
  - Joint decoding of parsing and other tasks (IE, MT, reasoning ...)
- Modeling sentence processing
  - BP is a *parallel, anytime* process

## ADAPTATION

## Projecting Hidden Structure

## Projection



- Train with bitext
- Parse one side
- Align words
- Project dependencies
- Many to one links?
- Invalid trees?
- Hwa et al.: fix-up rules
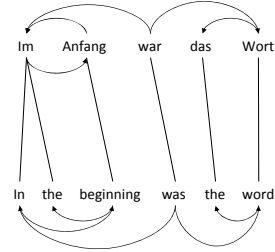- Ganchev et al.: trust only some links

## Divergent Projection



Auf   diese   Frage   habe   ich   leider   keine   Antwort   bekommen

NULL

I   did   not   unfortunately   receive   an   answer   to   this   question

head-swapping        monotonic        null        siblings

## Free Translation

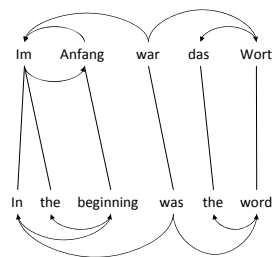## What's Wrong with Projection?



- Hwa et al. Chinese data:
  - 38.1% F1 after projection
  - Only 26.3% with automatic English parses
  - Cf. 35.9% for attach right!
  - 52.4% after fix-up rules
- Only 1-to-1 alignments:
  - 68% precision
  - 11% recall

## Projection



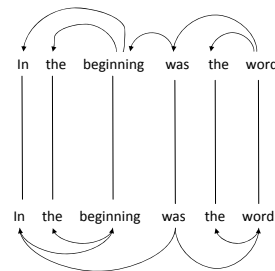- Different languages
- Similar meaning
- Divergent syntax

## Adaptation



- Same sentence
- Divergent syntax

## A Lack of Coordination

## Prepositions and Auxiliaries

## Adaptation Recipe

- Acquire (a few) trees in target domain
- Run source-domain parser on training set
- Train parser with features for:
  - Target tree alone
  - Source and target trees together
- Parse test set with:
  - Source-domain parser
  - Target-domain parser

## Why?

- Why not just modify source treebank?
- Source parser could be a black box
  - Or rule based
- Vastly shorter training times with a small target treebank
  - Linguists can quickly explore alternatives
  - Don't need dozens of rules
- Other benefits of stacking
- And sometimes, divergence is very large

## MODEL STRUCTURE

## What We're Modeling



| | This paper |
|---|---|
| Generative | $p(t,a,w \mid t',w')$ |
| Conditional | $p(t \mid t',a,w,w')$ |
| | Ongoing work |
| | $p(t,t',a \mid w,w')$ |

$$s(t,t',a,w,w') = \sum_i \theta_i f_i(t,w)$$
$$+ \sum_j \theta_j g_j(t,t',a,w,w')$$
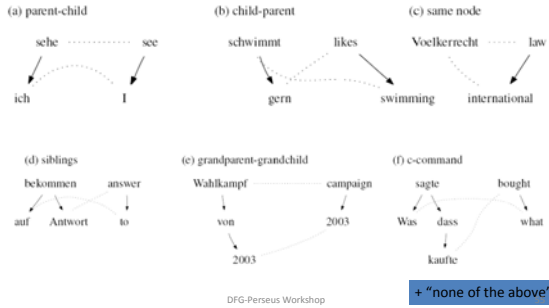
## Stacking

…

Model 2 has features for when to trust Model 1
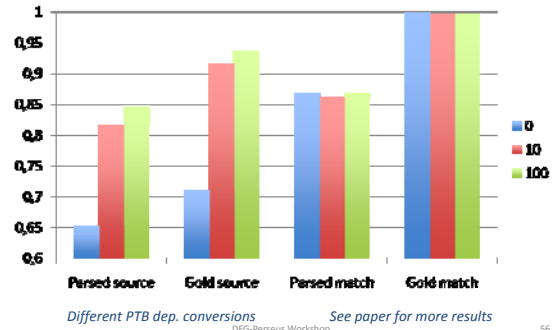


Model 2

Input

Model 1

## Quasi-Synchronous Grammar

- Generative or conditional monolingual model of target language or tree
- Condition target trees on source structure
- Applications to
  - Alignment (D. Smith & Eisner '06)
  - Question Answering (Wang, N. Smith, Mitamura '07)
  - Paraphrase (Das & N. Smith '09)
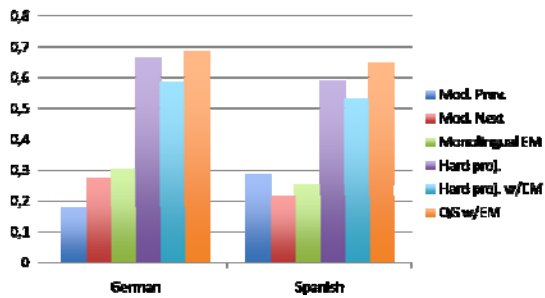  - Translation (Gimpel & N. Smith '09)
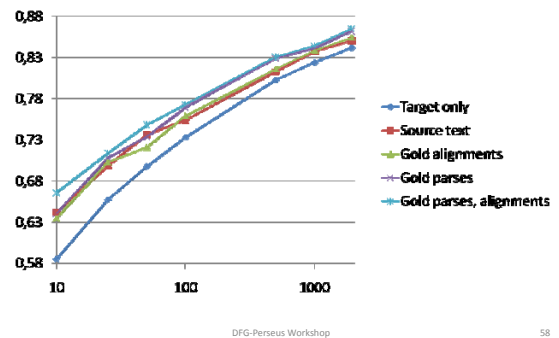
## Dependency Relations



(a) parent-child  
sehe ........ see  
ich → I

(b) child-parent  
schwimmt likes  
gern → swimming

(c) same node  
Voelkerrecht ........ law  
→ international

(d) siblings  
bekommen answer  
auf Antwort → to

(e) grandparent-grandchild  
Wahlkampf campaign  
von → 2003

(f) c-command  
sagte bought  
Was dass → what  
kaufte

+ "none of the above"

## Adaptation Results



Parsed source | Gold source | Parsed match | Gold match

*Different PTB dep. conversions*    *See paper for more results*

## Unsupervised Projection



German    Spanish

- Mod. Prnc.
- Mod. Next
- Monolingual EM
- Hard proj.
- Hard proj. w/CM
- QS w/EM

## Supervised Projection



- Target only
- Source text
- Gold alignments
- Gold parses
- Gold parses, alignments

## MINING A MILLION BOOKS

## Mining Information from Books



Several other derivations are given; among the rest that of the learned Bochart seems to have been most generally adopted: according to him, the Phœnicians called the ifland Barat-Anac, that is, the country of tin or lead; which name might by the Romans have been formed into Britannia, or Britannicæ infulæ.

Bochart  
Phoenicians  
Barat-Anac  
Britannia  
Britannicæ infulæ

- Modern OCR, several errors/page
- Names are worse:
  - In one study, 35% names incorrectly transcribed
- Errors propagate to later steps
- Train language models and name extractors on noisy corpus

## Learning by Reading



Correct *Orozco* occurs 149x in doc.
*Tlachquiauhco* occurs once in doc.

Inducing language- and corpus-specific features

## Reuse & Quotation

## UMass Book Search