

Partial Parsing

Steven Abney

University of Tübingen

`abney@sfs.nphil.uni-tuebingen.de`

A tutorial presentation

Partial Parsing

- Standard parsers
 - Evaluate global parses, not partial parses
 - Do all-paths search (chart or no)
 - Why unrestricted text is difficult
 - Incompleteness of lexicon
 - Incompleteness of grammar
 - Incompleteness of semantics
 - Long sentences
 - Errors in input
 - Partial parsing
 - Produce forest
 - Speed
 - Reliability (precision)
 - Breadth
 - Robustness
 - Sacrifice depth of analysis
 - Levels
 - Breaking up “The Parsi
 - Fairly independent steps
 - Partial parsing is the : tagging
-
-

Overview

Chunks <i>Cass</i> <i>Chunks & dependencies</i> <i>Supertags</i> <i>Longest match</i> <i>Finite-state</i> <i>Chinks & chunks</i> <i>Ejerhed</i> <i>Church</i> <i>Fidditch</i> <i>Bourigault</i> <i>Voutilainen</i> <i>Chen & Chen</i> <i>Rooh</i>	MUC / IR <i>Futrelle</i> <i>BBN</i> <i>Seneff</i> <i>AutoSlog</i> <i>Fastus</i> <i>Copsy</i>	Phrase Spotting <i>Relative likelihood</i> <i>Alpha & beta</i>	Regre <i>Linear r</i> <i>Regress</i>
	HMMs <i>Generation</i> <i>Partial paths</i> <i>NP recognition</i>	Parameter Estimation <i>Smoothing</i> <i>Forward-backward</i>	Grammatic <i>Bayesian</i> <i>Finch</i> <i>Smith &</i>
		Finite-State grammars <i>HMMs are FSAs</i> <i>Composing FSAs</i>	Lingu <i>Function</i> <i>S-projec</i> <i>Chunks</i>

- Cascaded cheap analyzers
 1. Tag (Church tagger)
 2. First guess on NPs (Church NP-recognizer)
 3. Finite-state NP recognizer (correct some tagging and NP-boundary errors)
 4. Chunks
 5. Simplex clauses
 6. Clause repair
 7. Attachment
 - Each analyzer outputs a single ‘best’ answer
 - Local search, but no global search, within levels
 - Repair errors downstream
-
-

EOS

In_P

[South_{PN} Australia_{PN} beds_{NPI}]

of_P

[boulders_{NP}]

were_{Bed}

deposited_{Vbn}

by_P

[melting_{Vbg} icebergs_{NPI}]

in_P

[a_D gulf_N]

[that_{Wps}]

marked_{Vbd}

[the_D position_N]

of_P

[the_D Adelaide_{PN} geosyncline_N]

,

[and_D elongated_{Vbn}]

,

[sediment-filled_{Vbn} depression_N]

in_P

[the_D crust_N]

.

EOS

EOS

[PP In [NP South Australia beds

[PP of [NP boulders]]

[VP were deposited]

[PP by [NP melting icebergs]]

[PP in [NP a gulf]]

[WhNP that]

[VP marked]

[NP the position]

[PP of [NP the Adelaide geosyncline

,

[NP an elongated, sediment-filled

[PP in [NP the crust]]

.

EOS

EOS

[PP In [NP South Australia beds]]
[PP of [NP boulders]]
[VP were deposited]
[PP by [NP melting icebergs]]
[PP in [NP a gulf]]

[WhNP that]
[VP marked]
[NP the position]
[PP of [NP the Adelaide geosyncline]]

,
[NP an elongated, sediment-filled depression]
[PP in [NP the crust]]

.
EOS

[NoSubj

EOS

[PP In [NP South Australia
[PP of [NP boulders]]
Pred: [VP were deposited]
[PP by [NP melting icebergs]]
[PP in [NP a gulf]]

]

[SRC

Subj: [WhNP that]
Pred: [VP marked]
[NP the position]
[PP of [NP the Adelaide geosyncline]]

,

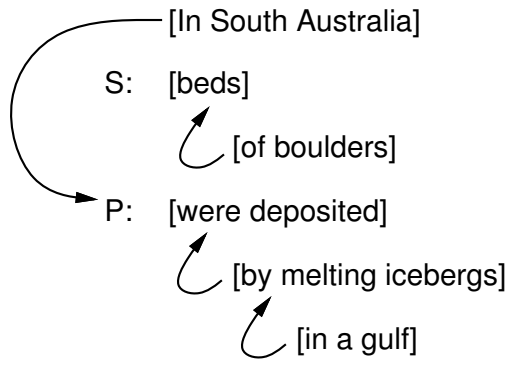
[NP an elongated, sediment-filled depression]
[PP in [NP the crust]]

]

.
EOS

[NoSubj

EOS

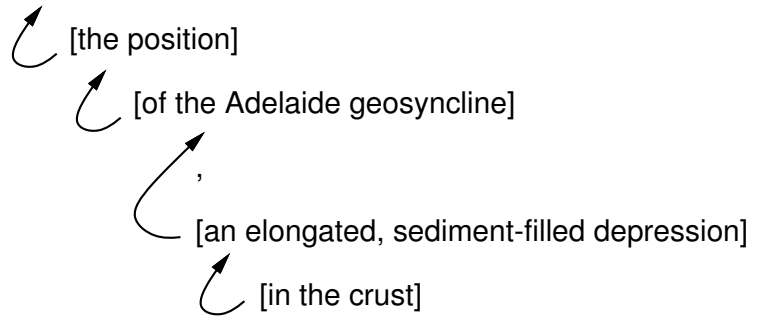


]

[SRC

S: [that]

P: [marked]

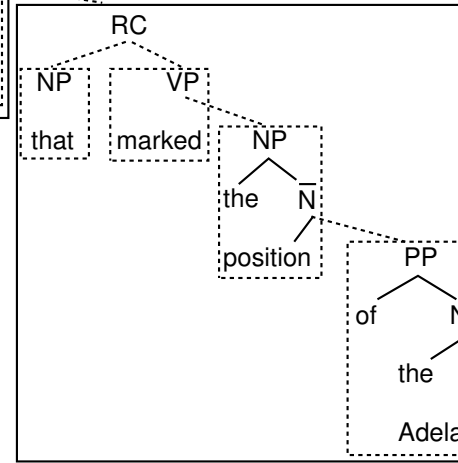
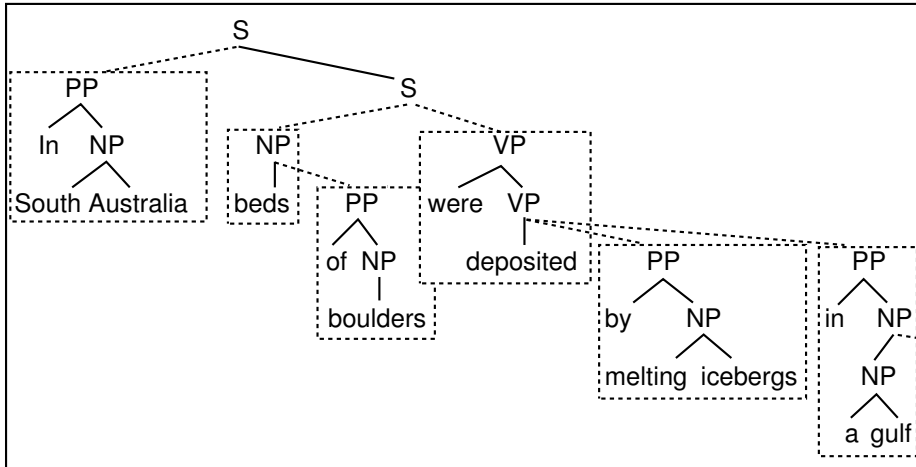


]

.

EOS

Chunks and Dependencies

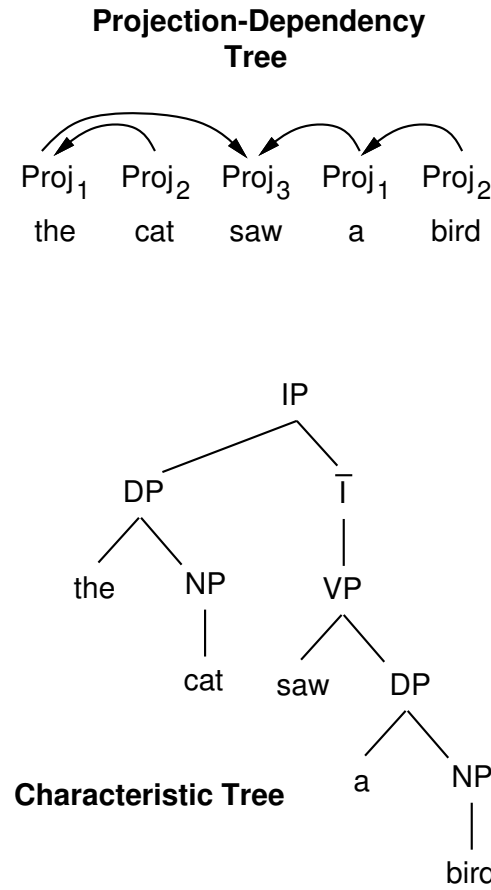
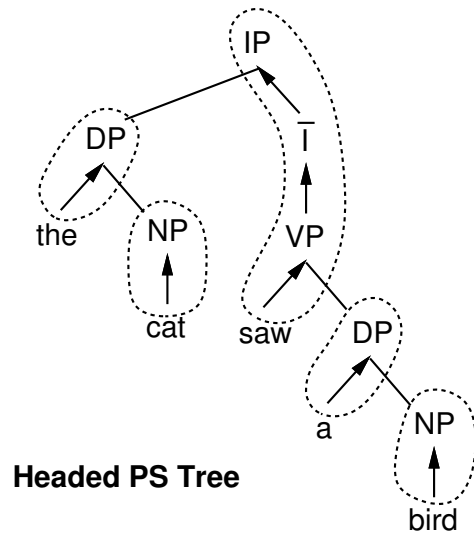


Chunks and Dependencies

- Factorization of the parsing problem
 - Dependencies: lexico-semantic, binary (head-dependent)
 - Chunks: syntactic category, finite-state sequences
 - Simplex clauses
 - Trapping all-ways ambiguities
 - E.g., no PP-attachment across clause boundary
 - (Chunks trap noun-modification ambiguities)
 - Instead of exponential global ambiguity, sequence of independent small sets
-
-

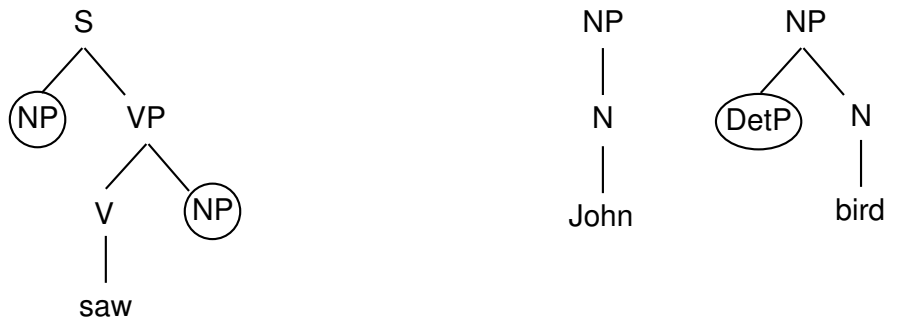
More generally

- Inspired by Gaifman [89]



Supertags

- Joshi & Srinivas [123]
- Instead of dependencies between projections, dependencies between elements



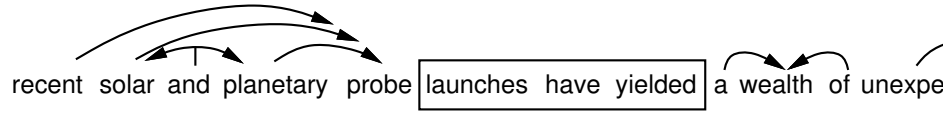
- The difference: dependencies can also represent adjunction, not just substitution
 - Parsing as tagging: elementary trees are ‘supertags’
 - Use standard tagging techniques (HMM’s)
 - Or take advantage of dependency information in supertags to identify relevant 2-grams
-
-

- Variant of dependency grammar
- Parsing as tagging
 - Syntactic category tag
 - Syntactic function tag
- Rules are rules for eliminating tags (“constraints”)

Vfin. . . → delete Ma
NomHead & . . . Vfin & ¬ NomHead. . . NomHead → keep only

- 1300 morphological rules, 120 syntactic rules
 - Ambiguous representation
-

recent	>N
solar	>N
and	CC
planetary	>N
probe	NH
launches	V
have	V
yielded	V
a	>N
wealth	NH
of	<N
unexpected	>N
data	NH



recent	>N
solar	>N
and	CC
planetary	>N
probe	>N
launches	NH
have	V
yielded	V
a	>N
wealth	NH
of	<N
unexpected	>N
data	NH



Creative Ambiguity

- Or, Lazy Disambiguation
- Or, Picking the Fights You Can Win
- D-theory [150]
- Unscoped quantificational formulae
- Ambiguity preservation in transfer in MT

Say which clause a PP belongs to
and where it's attached

Chunks

PP → (p | to) + (NP | vbg)

WhPP → (p | to) + WhNP

AdvP → (ql | preced | rb)* rb

AP → (AdvP | ql)* adj

Inf → to AdvP? VP-inf

VP → AdvP? (md | v-tns | hv-tns VPN? | be-tns (VPG | Vn)?)

VP-inf = AdvP? (vb | hv VPN? | be (VPG | Vn)?)

VPN = AdvP? (vbn | hvn | ben (VPG | Vn)?)

VPG = AdvP? (vbg | hvg | beg Vn?)

Vn = AdvP? (vbn | hvn | ben)

Other → any

Longest match heuristic

- Used in lexical analyzers for compilers
- Psychologically plausible

the emergency crews always dread is domestic violence

|-----|

|-----|

while she was mending the sock fell off her lap

|-----|

|-----|

Longest Match

- One automaton for each phrase category
- Start automata at position i (initially, $i = 0$)
- Take longest match

0 saw 1 horses 2 are 3 needed 4
└NP┘
└VP┘
└NP————┘

- Set $i := j$ and repeat

0 saw 1 horses 2 are 3 needed 4
└VP┘
└VP————┘

Effectiveness of longest match

- Take chunks out of the UPenn Treebank

NP → D N
NP → D Adj N
VP → V
VP → Hv Vbn
⋮

- At each point in string take longest matching pattern
 - Guess if multiple longest matches (of different category)
 - Punt one word if no match
- Performance: Precision .92
 Recall .88

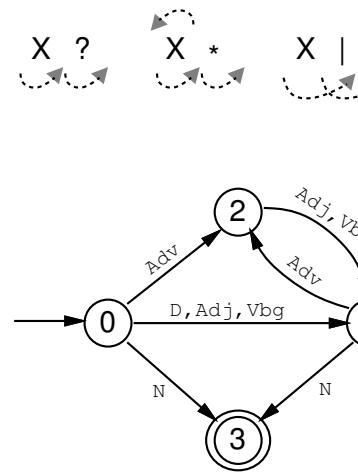
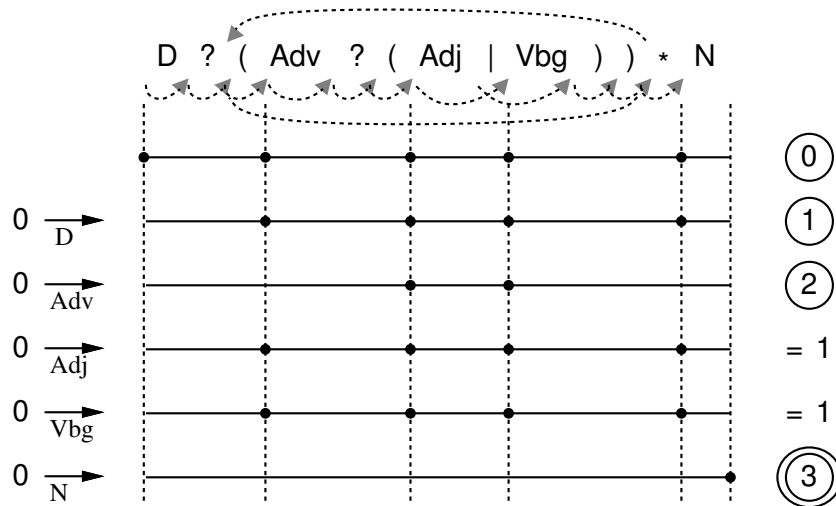


Finite-state techniques

- Hand-written grammar (regular expressions)

$NP \rightarrow Det? (Adj | Ing)* N$

- Compile into FSA



Clause

Extra-VPs → EOC+ pre NP mid VP post (VP post)+
Clause → EOC+ pre NP mid VP post
ObjRC → EOC* WhNP pre NP mid VP post
SubjRC → EOC* WhNP mid VP post
WhClause → EOC* (WhPP | wrb) pre NP mid VP post
VP-Conj → cc VP post
No-Subj → EOC+ pre VP post
No-VP → EOC+ post

pre = (X | Wh | PP-Conj)* ((, AdvP)? ,)?

mid = (X | EOC-Soft | NP)*

post = (X | NP)*

PP-Conj = PP (, N PP*)* cc NP

X = [^ Special]

Special = [EOC Wh NP VP]

EOC = [EOC-Hard EOC-Soft]

EOC-Hard = [: . eos]

EOC-Soft = [, cc cs that]

Wh = [WhNP WhPP wrb]

Bottom line

- Fast (once upon a time)
 - Pos: 4.2 ms/w
 - Cass: 15.0 ms/w
 - Total: 19.2 ms/w = 52 w/s
 - Accurate
 - ~5% error chunks
 - ~5% error subj & pred
 - BUT: Already in the tail
 - Only a few error types occur frequently
 - Only a few changes to the grammar will have much effect
 - The rest is sand
-
-

Parser speed

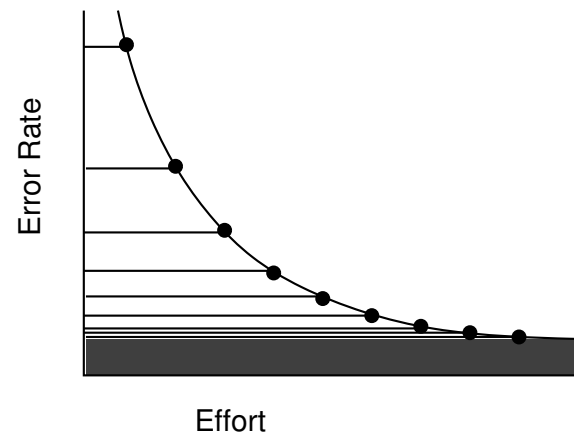
- Want a fast parser, get a fast machine
- Restricting search helps

Program	depth	sw	hardware	w/s	
Fidditch3	parse	C	SGI	5600	
Copsy	np	Pascal	BS2000	2700	
CG	dep		Sparc10	1550	± 250
Fidditch3	parse	C	Sun4	1200	
Pos	tag		Sun4	240	
Fidditch2	parse	Lisp	Sun4	62	
Cass	chunk	Lisp	Sun4	52	
Clarit	np	Lisp		50	
Fastus	chunk	Lisp	Sparc2	39	
Cass	chunk	Lisp	UX400S	32	
Scisor	skim			30	
Fidditch1	parse	Lisp	Sym-36xx	28	
McDonald	parse		MacII	14	± 6
Chupa	parse	Lisp	UX400S	1.1	
Traditional	parse			0.20	

Lies, damned lies, and statistics

- What would you get by guessing?
 - Tagging: always taking most-frequent tag \rightarrow 10% error
- Per-chunk error rate vs. per-sentence error rate
 - 5% chunk error
 - 10 chunks/sentence
 - $1 - (1 - .5)^{10} = 40\%$ sentence error

- Zipf's Law
 - A little effort goes a long way—*at first*
 - The down side—further significant error reduction requires horrendous effort



Chinks and chunks

- Venerable idea:
 - Function words are phrase delimiters (chinks)
 - Content words are phrase contents (chunks)
- Ross & Tukey [164]
 - Used for sorting KWIC index of statistical works

on the construction of Bose-Chaudhuri matrices
with the help of Abelian group characters

- fgroups
 - F+ C+
 - Used as low-level phrasal units in Bell Labs speech synthesizer
-

- Non-recursive (simplex) NP's and clauses
- Finite-state and stochastic methods
- Motivated in part by psycholinguistic studies
- Performance

	NP	Clause
Finite-state	3.3%	13%
Stochastic	1.4%	6.5%

- Application: text-to-speech (intonation)
-

Clause grammar

Clause → cc? NP ([cc p \$] NP)* adv? tns-v X* Punct?
| cc Adv? v X* Punct?
| cc? Comp+ X* Punct?
| cc? NP ([cc p \$] NP)* X* Punct?
| Verb X* Punct?
| cc? (Stray | NP)* X* Punct?

X = [^ Comp Punct]

Comp = [cs to wdt wrb wps wpo wp\$ wql]

Punct = [, . - :]

Adv = [rb rbr]

Verb = [tns-v vbg vbn beg hvg]

Stray = [Adv rp ql neg nr jj jjr p]

Example

[the jury further said in term-end presentments]
[that the City Executive Committee ,]
[which had over-all charge of the election ,]
[deserves the praise and thanks of the City of Atlanta for the manner]
[which the election was conducted .]

Church [57, 58]

- Stochastic tagger, followed by nonrecursive NP recognizer
- Between any pair of tags, we can insert one of:

[]] [-

- Must keep track of whether inside or outside of NP

[the [corrosion weight loss [

- Computation:

B:	[-	-	-]	[...
I:	0	1	1	1	1	0	...
T:	\$D	DN	NN	NN	NP	PD	...

- Choose the sequence of brackets with the highest probability
-

Probabilities

B:	[-	-	-]	[...
I:	0	1	1	1	1	0	...
T:	\$D	DN	NN	NN	NP	PD	...

- Estimate by counting in parsed corpus

$$\hat{\Pr}(B|T) = \frac{f(B, T)}{f(T)}$$

 - Including inside/outside constraint

$$\begin{aligned} & *[[\Pr(B = b|[T, I = 1) \\ & *]] \Pr(B = b|[T, I = 0) \\ & *]] \Pr(B = b|[T, I = 0) \end{aligned}$$

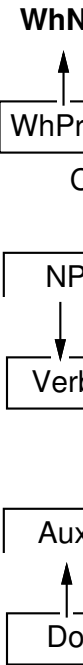
$$\Pr(B|T, I)$$

 - Choices at different positions independent

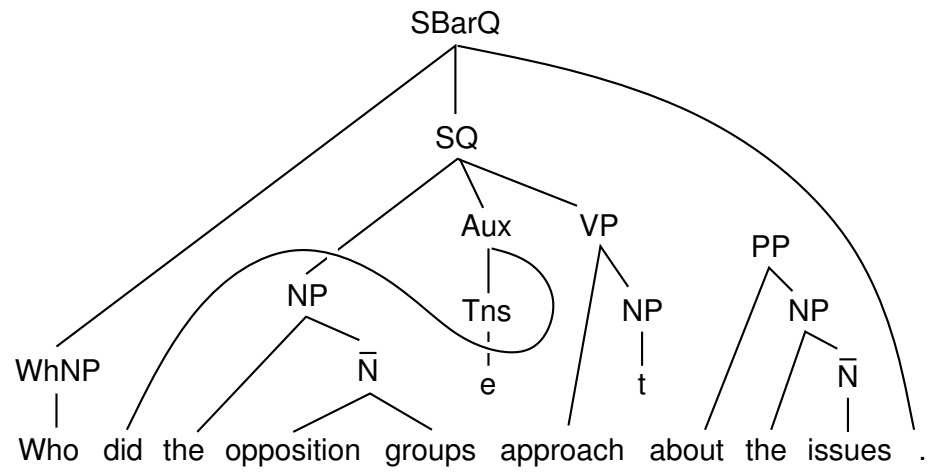
$$\Pr(\mathbf{B}|\mathbf{T}, \mathbf{I}) = \prod_i \Pr(\mathbf{B}_i|\mathbf{T}_i, \mathbf{I}_i)$$
-

- Industrial-strength version of Marcus Parser

Create	Recognizing leading edge of new node
Attach	Recognizing material belong to current node
Drop (Close)	Recognizing leading edge of material following node
Switch	Subject-aux inversion
Insert	Recognizing empty category
Attention-shift	Recognizing leading edge of NP in lookahead
Punt	Avoid an attachment decision

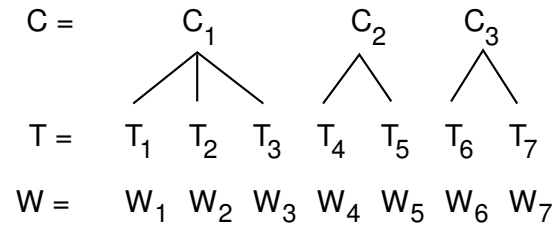


Fidditch tree



- Extraction of likely multi-word terms for automatic indexing
 - Phrase boundaries
 - Chunks: things that can't be chunks
 - E.g., Verbs, Pron, Conj, Prep (except *de*, *a*, Det
 - *un [traitement de texte] est installe sur le [disque dur de la stati*
 - Parsing/extraction
 - Rules for extracting smaller potential terms
 - E.g. $N_1 \text{ Adj } P \text{ D } N_2 \text{ P } N_3 \rightarrow N_1 \text{ Adj}, N_2 \text{ P } N_3$
 - *disque dur, station de travail*
 - 800 such rules, manually built and tested
-
-

- Building sequence of chunks on tags



- Best chunk

$$\begin{aligned} C^* &= \operatorname{argmax}_C \Pr(C|W) \\ &\hat{=} \operatorname{argmax}_C \Pr(C|T) \\ &\hat{=} \operatorname{argmax}_C \prod_i \Pr(C_i|C_1, \dots, C_{i-1}, T) \\ &\hat{=} \operatorname{argmax}_C \prod_i \Pr(C_i|C_{i-1}, T) \\ &\hat{=} \operatorname{argmax}_C \prod_i \Pr(C_i|C_{i-1}) \Pr(C_i|T) \quad ! \end{aligned}$$

- Probabilities estimated from parsed corpus (Susanne)
-

Internal probability and contextual probability

- Church and DeRose also say

$$\Pr(\mathbf{T}_i | \mathbf{T}_{i-1}, \mathbf{W}_i) \hat{=} \Pr(\mathbf{T}_i | \mathbf{T}_{i-1})$$

- Doesn't necessarily hurt performance

- But:

D = throw of die

E = 1 if D is even, 0 otherwise

L = 1 if $D \leq 3$, 0 otherwise

$$\Pr(D = 2 | E = 1, L = 1)$$

$$\Pr(D = 2 | E = 1) \Pr(D = 2 | L = 1)$$

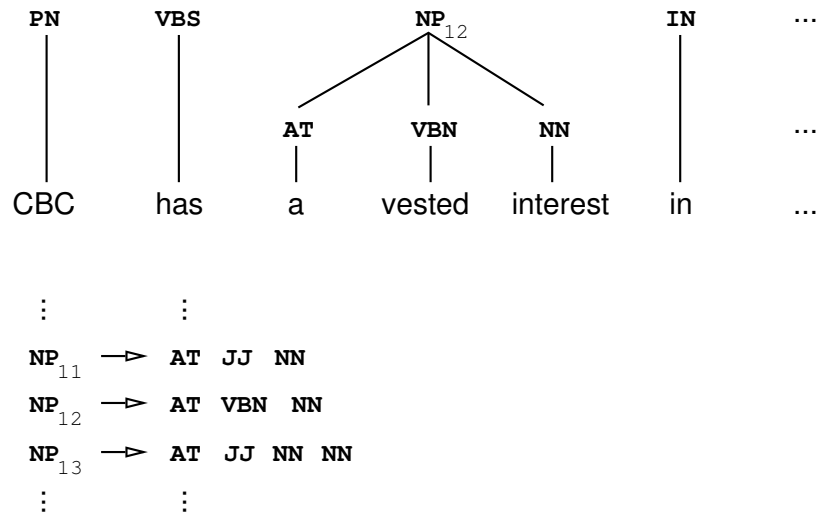
- Combining information sources: multivariate regression

- Alternative: HMM

$$\Pr(\mathbf{T} | \mathbf{W}) \propto \Pr(\mathbf{T}, \mathbf{W})$$

$$\hat{=} \prod_i \Pr(\mathbf{T}_i | \mathbf{T}_{i-1}) \Pr(\mathbf{W}_i | \mathbf{T}_i)$$

- Modified Hidden Markov Model



- Generation probabilities $\Pr(x_i|x_{i-1})$ $\Pr(w|t)$
- Choose the structure by which the words were most likely generated

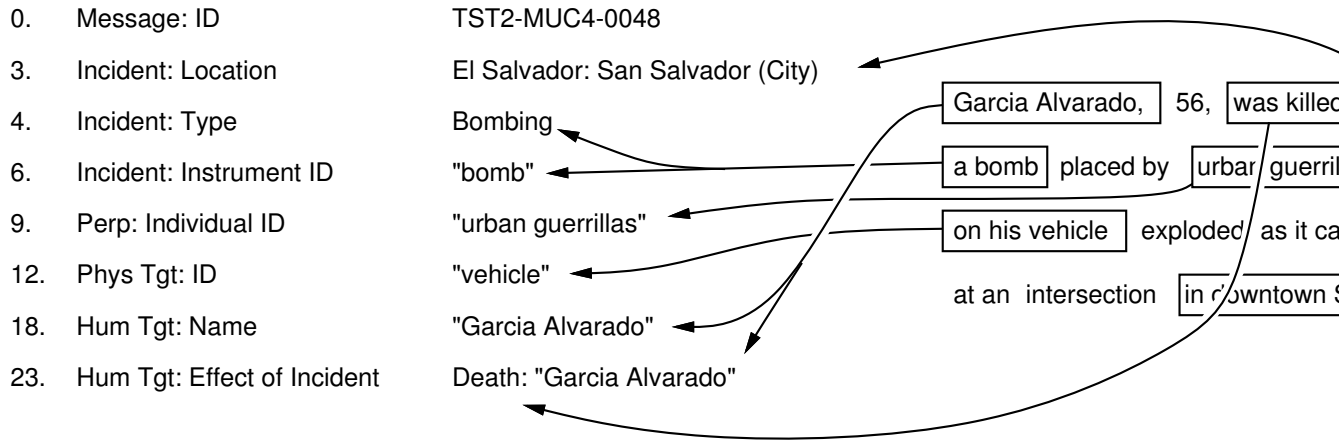
Threads

- Determinism
 - Local evaluation of pieces
 - Dependency grammar
DG \leftrightarrow CFG \leftrightarrow chunks
 - Levels/cascade
 - Specialized grammars
 - Creative ambiguity
 - Longest match
 - Likelihood
 - HMM's
 - Regression
 - Induction (bootstrapping, C)
 - Linguistic/psycholinguistic
-
-

MUC

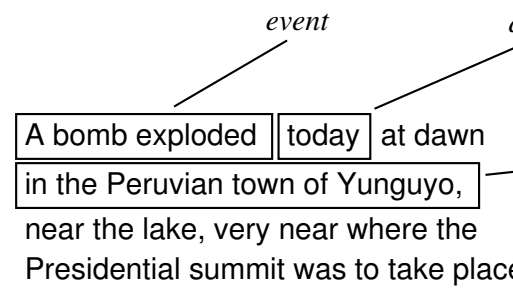
- Message Understanding Conference
 - Task: data extraction from news reports
 - Filter out irrelevant texts
 - Tokenize and clean
 - Trigger on tokens
 - Fill semantic frames
 - Merge frames to fill data templates
-
-

Example



- Partial parsing for handling unrestricted text
- Message Understanding doesn't require complete parse

- Data extraction
- Message routing
- Message prioritization



- Questions

- Effectiveness of fragment recognition?
- How to interpret fragments?

- Interpretation

- Identify headword to get
of phrase
- Make attachment if cla
requirement

Tokenize and clean

- Issues
 - Spelling errors
 - Foreign words / foreign names
 - Punctuation
 - Formulae
 - Graphics / Formatting
 - Sentence, paragraph boundaries
 - Requirements
 - Fast
 - Highly reliable (snowball)
 - When in doubt, pass on ambiguity
 - Shades into partial parsing
-
-

- Examples

7.3	sodium chloride
36,768	CO ₂
2,6-diaminohexanoic acid	3.4×10^{-8}
³ H	

Cells were suspended in a medium containing $3.05 \times 10^{-2} \mu\text{M}$ L-[*methyl*-³H]-methionine, α -methyloaspartate and AIBU⁸.

- Deterministic subgrammars

- Hand-correction

Examples

- Date/time expressions

24.10.94	10:06 p.m.
10/24/94	2000 GMT
Tues. the 24th Oct., 1994	two-thirty
Thu, 06 Oct 1994 11:47:55 EDT	

- Names

- Person: *John T. Smith, Juan Mercedes Garcia de Mendoza, Kim*
- Place: *the Orontes River; Mt. Pinatubo; Paris, TX*
- Organization: *IBM; AT&T; Mt. Sinai Publishing Co., Inc.*
- Titles: *Green County Sheriff's Deputy Gordon Caldwell*

- Bibliographic conventions

Smyth (1990)
Fig. 2
... as is probable.⁶
NEW ORLEANS, 19 Jun 93 (API) –

- State of the art: write little grammars by hand

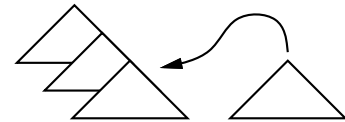
PLUM (BBN) [17]

- Uses de Marcken parser to get fragments
- Semantic frames tied to words

$\text{bomb}_V (\text{subj} [1], \text{obj} [2])$ $\left[\begin{array}{l} \text{bombing} \\ \text{ti-perp-of} \langle \text{person} \rangle \llbracket 1 \rrbracket \\ \text{object-of} \langle \text{any} \rangle \llbracket 2 \rrbracket \end{array} \right]$

- Frame of fragment is gotten from head
- Assemble fragments deterministically via *attachment*

- Try leftward attachments first
- Try low attachments before high
- Take first attachment satisfying slot constraints

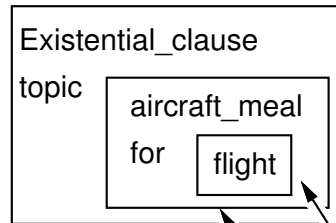


- Start with standard full-sentence parser
 - Parse fails: no $S[0, n]$
 - Consider $X[i, j]$ for X “major” and $i = 0$
 - Take longest match (maximize j)
 - Set $i = j$, repeat
 - If no $X[i, j]$, take next word, set $i = i + 1$, repeat
 - Use discourse processor to integrate fragments
 - Bottom line: good, but not as good as full-sentence parser
-
-

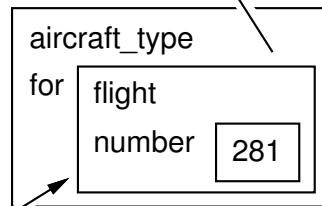
Using Discourse Processor

(*what are the meals*) and (*aircraft for flight two eighty one*) and
also for (*flight two oh one*)

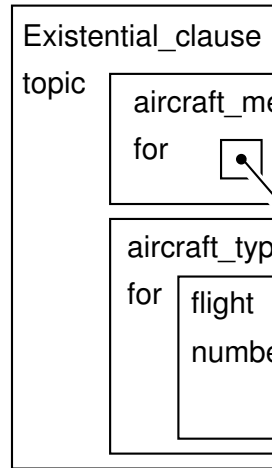
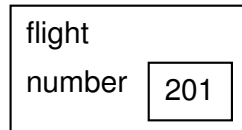
what are the meals



aircraft for flight two eighty one



flight two oh one



Big problem for frame-based systems

- Building lexicon of frames
 - Frames provide robustness: assemble any way they fit
 - Acquiring new frames from corpora
 - To name a few at random: [16, 34, 40, 44, 54, 60, 77, 95, 103, 128, 135, 176, 199]
 - UMass: AutoSlog
-
-

- Input: examples of correct slot files

The ARCE battalion command has reported that about 50 peasants of various ages have been kidnapped by terrorists of the Farabundo Marti National Liberation Front in San Miguel department.

[perp-indiv-id “terrorists”]

- Parse sentence, look at region around given word

actor: peasants

verb: kidnapped [PASSIVE]

prep: by

pobj: terrorists of FMNL

- Propose pattern

verb = kidnapped [PASSIVE]

actor = ANY

PP_{by} = $\left[\begin{array}{l} \text{ORGANIZATION} \\ \text{TERRORIST} \\ \text{PROPER-NAME} \\ \text{HUMAN} \end{array} \right.$

- Automatic evaluation of precision/recall possible
-

The inspiration for FASTUS was threefold. First, we were struck by the performance that the group at the University of Massachusetts got out of a fast system. It was clear they were not doing anything like the depth of prepositional syntactic analysis, or pragmatics that was being done by the systems at SRI, Bell Electric, or New York University. They were not doing a lot of processing, but they were doing the *right* processing.

The second source of inspiration was Pereira's work on finite-state approaches to grammars, especially the speed of the implementation.

Speed was the third source. It was simply too embarrassing to have to report at the MUC-3 conference that it took TACITUS 36 hours to process 100 sentences. FASTUS has brought that time down to 11 minutes.

Fastus

- Triggering: single keywords from patterns plus known proper names
 - Phrase recognition
 - Noun groups
 - Verb groups
 - P, Conj, RelPro, *ago*, *that*
 - Keep only longest match (nested, not overlapping)
 - Patterns
 - killing of ⟨HumanTarget⟩
 - ⟨GovtOfficial⟩ accused ⟨PerpOrg⟩
 - bomb was placed by ⟨Perp⟩ on ⟨PhysicalTarget⟩
 - Merge compatible incidents
-
-

Fastus example

Noun Group: Salvadoran President-elect

Name: Alfredo Cristiani

Verb Group: condemned

Noun Group: the terrorist

Verb Group: killing

Prep: of

Noun Group: Attorney General

Name: Roberto Garcia Alvarado

Conj: and

Verb Group: accused

Noun Group: the Farabundo Marti National Liberation Front (FMLN)

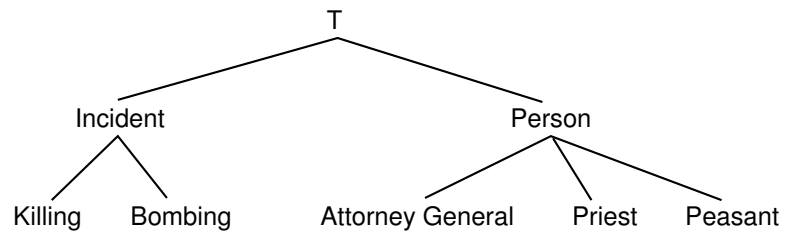
Prep: of

Noun Group: the crime

Fastus merging

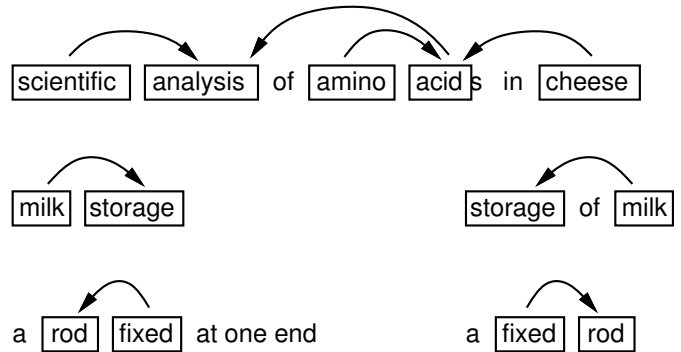
- Lots of frame scraps
- Merge if all slot-fillers compatible

$$\begin{bmatrix} \text{Incident} & \text{Killing} \\ \text{Perp} & - \\ \text{Confid} & - \\ \text{HumTarg} & \text{"Alvarado"} \end{bmatrix} + \begin{bmatrix} \text{Incident} & \text{Incident} \\ \text{Perp} & \text{FMLN} \\ \text{Confid} & \text{Suspected} \\ \text{HumTarg} & - \end{bmatrix} \Rightarrow \begin{bmatrix} \text{Incident} & \text{Killing} \\ \text{Perp} & \text{FMLN} \\ \text{Confid} & \text{Suspected} \\ \text{HumTarg} & \text{"Alvarado"} \end{bmatrix}$$



Schwarz: Copsy [169]

- Dependency parsing of noun phrases to improve precision in IR



- Recognition rules must be

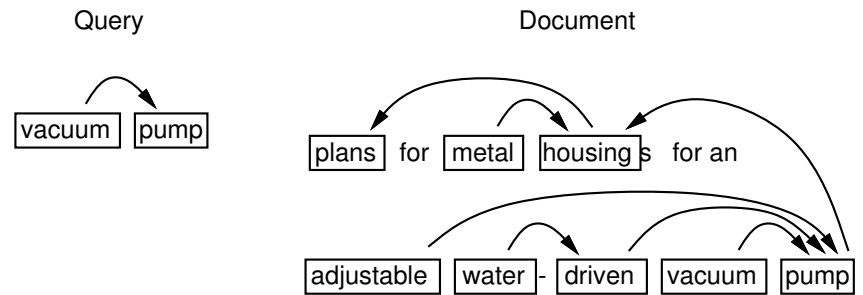
- Relevant
- Highly accurate
- Cheap to apply

- Normalization

- Dependencies
 - Development labor-intensive
 - 200 proposed rules tested
 - 15,000 matching sentences
 - final rules
-

Copsy: matching

- Index only words, not phrases
- Presearch: boolean OR of words in query
- Parse query, match against parsed documents in initial return set



- Fast enough to parse documents at search time (19 Kb/s)
 - Only 10% space overhead, however
-

More threads

- Interpretation
 - Dependencies \leftrightarrow Slots
 - “class = head class” is consequence
 - Merging if slot-fillers are compatible
 - Applications
 - Bootstrapping (collocations, alignment, ...)
 - MUC (Data extraction)
 - Terminology extraction
 - IR
 - Language models, spoken language understanding
-
-

Generation via Hidden Markov Model [160]

- Finite set of states s_i
 - Finite set of output symbols w_i
 - Random variables
State at time t Q_t
 - Random variables
Observation at time t O_t
 - Transition probabilities
 $\Pr(Q_{t+1} = s_j | Q_t = s_i)$ a_{ij}
 - Emission probabilities
 $\Pr(O_t = w | Q_t = s_i)$ $b_i(w)$
 - Initial probabilities
 $\Pr(Q_1 = s_i)$ π_i
-
-

Example: Tagger

- States are tags $\{\$, N, \text{Pron}, V, D\}$
- Output symbols are words $\{I, \text{see}, a, \text{bird}, .\}$

- Transition matrix

	\$	N	Pron	V	D
\$	0	.2	.5	0	.2
N	.3	.3	0	.4	0
Pron	.2	.1	0	.6	.1
V	.4	.2	.2	0	.2
D	0	1	0	0	0

- Emission matrix

	I	see	a	bird	.
\$	0	0	0	0	1
N	.1	.1	.1	.7	0
Pron	1	0	0	0	0
V	0	.9	0	.1	0
D	0	0	1	0	0

- Initial matrix

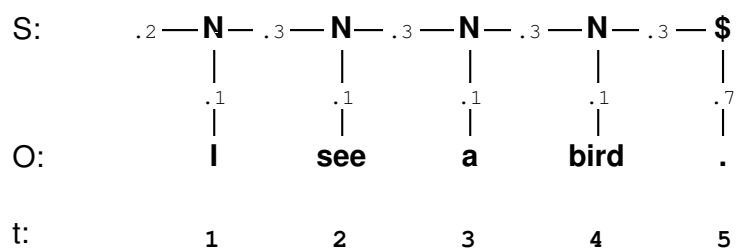
\$	N	Pron	V	D
0	.2	.5	0	.3

Probability of Generating a Structure

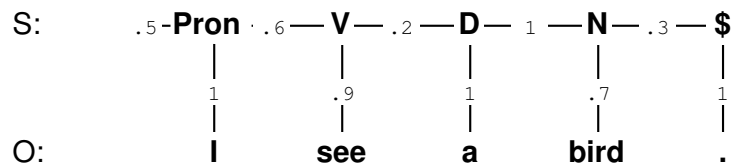
	\$	N	Pron	V	D
\$	0	.2	.5	0	.2
N	.3	.3	0	.4	0
Pron	.2	.1	0	.6	.1
V	.4	.2	.2	0	.2
D	0	1	0	0	0

	I	see	a	bird	.
\$	0	0	0	0	1
N	.1	.1	.1	.7	0
Pron	1	0	0	0	0
V	0	.9	0	.1	0
D	0	0	1	0	0

\$	0
----	---



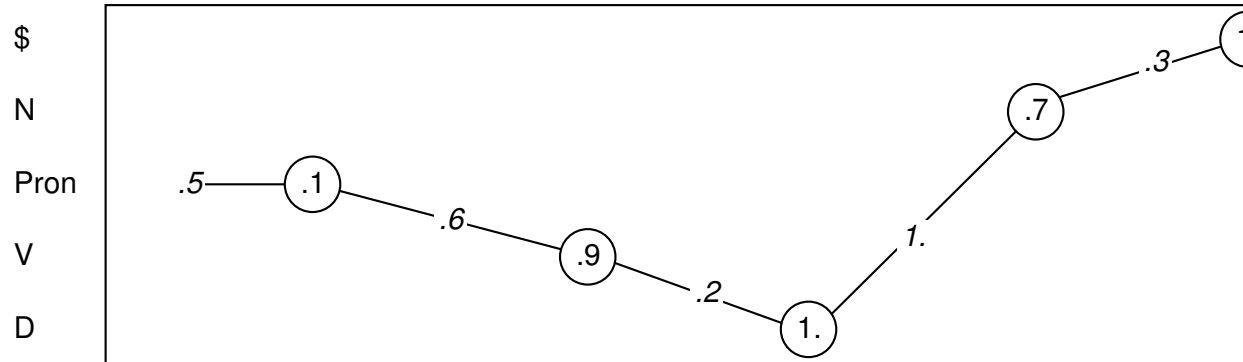
e^{-16}



$e^{-4.5}$

State sequence as path

π_{Pron} $b_{Pron}(l)$ $a_{Pron,V}$ $b_V(\text{see})$ $a_{V,D}$ $b_D(a)$ $a_{D,N}$ $b_N(\text{bird})$ $a_{N,\$}$ $b_{\$}$



Paths

State sequence (path)	\mathbf{q}	=	(q_1, \dots, q_T)
Observation sequence	\mathbf{o}	=	(o_1, \dots, o_T)
Probability	$\Pr(\mathbf{q}, \mathbf{o})$	=	$\Pr(Q_1 = q_1, \dots, Q_T = q_T, O_1 = o_1, \dots)$
Likelihood of path	$L(\mathbf{q})$	=	$\Pr(\mathbf{q}, \mathbf{o})$

'Best' = Maximum Likelihood

- We want

$$\mathbf{q}^* = \operatorname{argmax}_{\mathbf{q}} \Pr(\mathbf{q}|\mathbf{o})$$

- By definition

$$\Pr(\mathbf{q}|\mathbf{o}) = \frac{\Pr(\mathbf{q}, \mathbf{o})}{\Pr(\mathbf{o})}$$

- Since $\Pr(\mathbf{o})$ is constant

$$\Pr(\mathbf{q}|\mathbf{o}) \propto \Pr(\mathbf{q}, \mathbf{o})$$

- Therefore

$$\operatorname{argmax}_{\mathbf{q}} \Pr(\mathbf{q}|\mathbf{o}) = \operatorname{argmax}_{\mathbf{q}} \Pr(\mathbf{q}, \mathbf{o})$$

- Substituting

$$\mathbf{q}^* = \operatorname{argmax}_{\mathbf{q}} L(\mathbf{q})$$

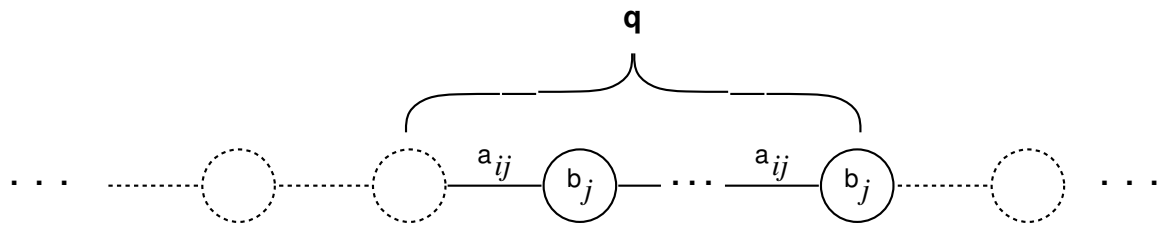
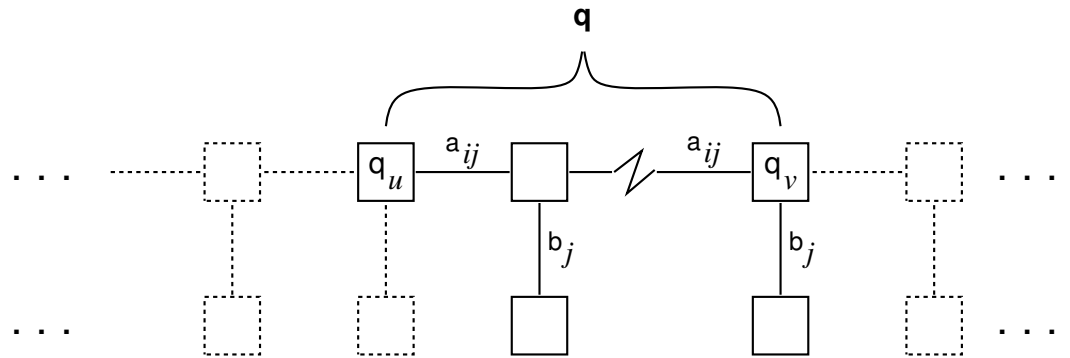
- That is, \mathbf{q}^* is the maximum-likelihood state sequence
-
-

Partial paths

Partial path from u to v

$$\mathbf{q} : u, v$$

Context-independent likelihood $\lambda(\mathbf{q}) = \prod_{t=u}^{v-1} a_{ij} b_j(\mathbf{o}_{t+1})$ where $q_t = s_i, o_t$

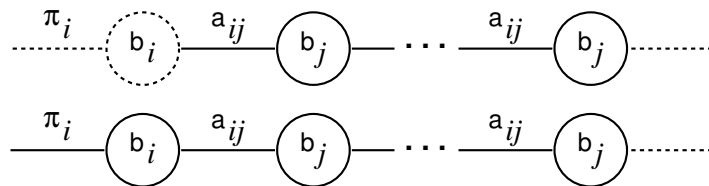


Partial paths

- Special case: initial

$\mathbf{q} : 1, t \quad \lambda(\mathbf{q}) :$

$$\lambda^0(\mathbf{q}) = \pi_i b_i(\mathbf{o}_1) \lambda(\mathbf{q})$$



- Relation to likelihood

if $\mathbf{q} : 1, T$ then $L(\mathbf{q}) = \lambda^0(\mathbf{q})$

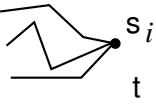
Computing q^*

- Most-likely partial sequence

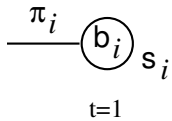
$$q_t^*(i) = \operatorname{argmax}_{\mathbf{q}:1,t|\mathbf{q}_t=s_i} \lambda^0(\mathbf{q})$$

- Likelihood thereof

$$\delta_t(i) = \max_{\mathbf{q}:1,t|\mathbf{q}_t=s_i} \lambda^0(\mathbf{q})$$



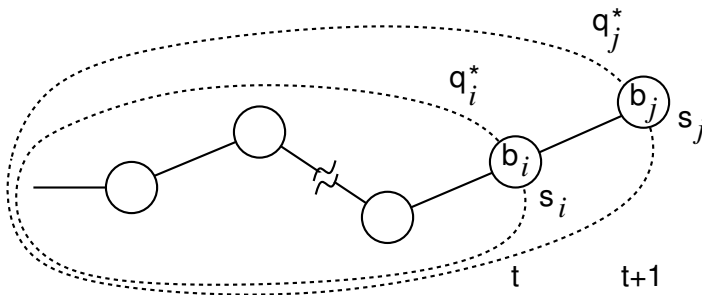
- Time 1



$$\delta_1(i) = \pi_i b_i(\mathbf{o}_1)$$

$$q_1^*(i) = \langle s_i \rangle$$

- Time $t + 1$



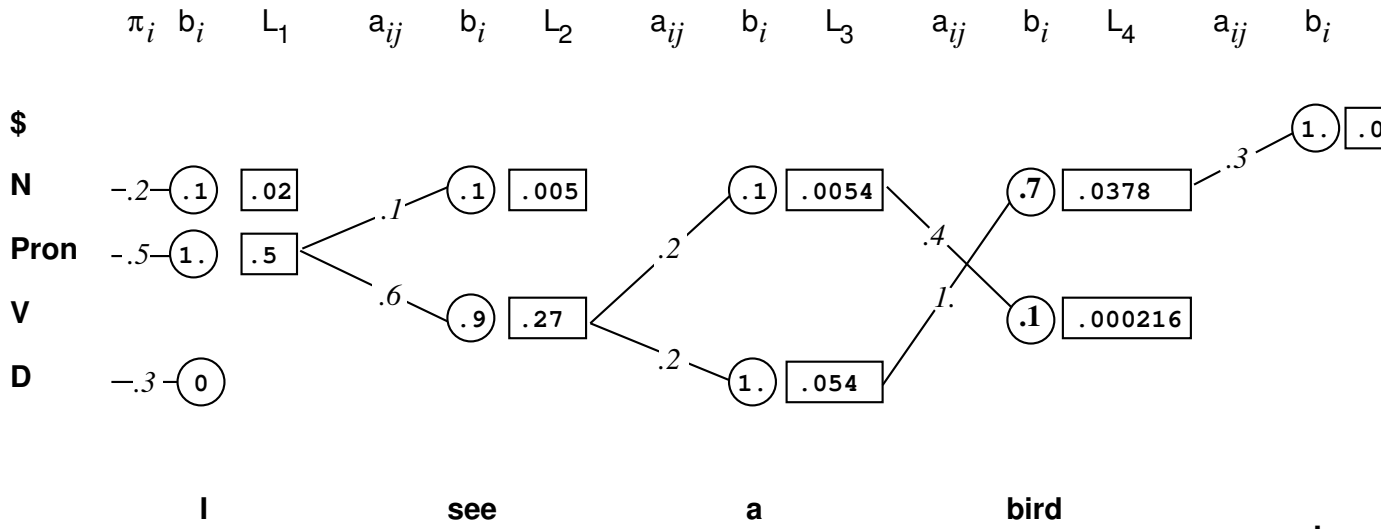
$$\delta_{t+1}(j) = \max_i \delta_t(i) a_{ij} b_j(\mathbf{o}_{t+1})$$

$$i^* = \operatorname{argmax}_i \delta_t(i) a_{ij} b_j(\mathbf{o}_{t+1})$$

$$q_{t+1}^*(j) = q_t^*(i^*) \wedge \langle s_j \rangle$$

Computing q^*

- Recursive definitions for $q_t^*(i)$, $\delta_t(i)$
- Fill in array by increasing values of variable of recursion (t)



NP-Recognizer as HMM

- States

[] I - #

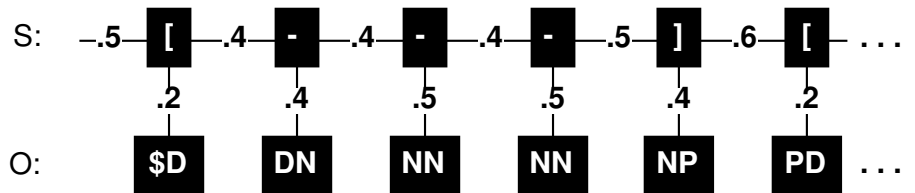
- π, a, b

[]	I	-	#
.5				.5

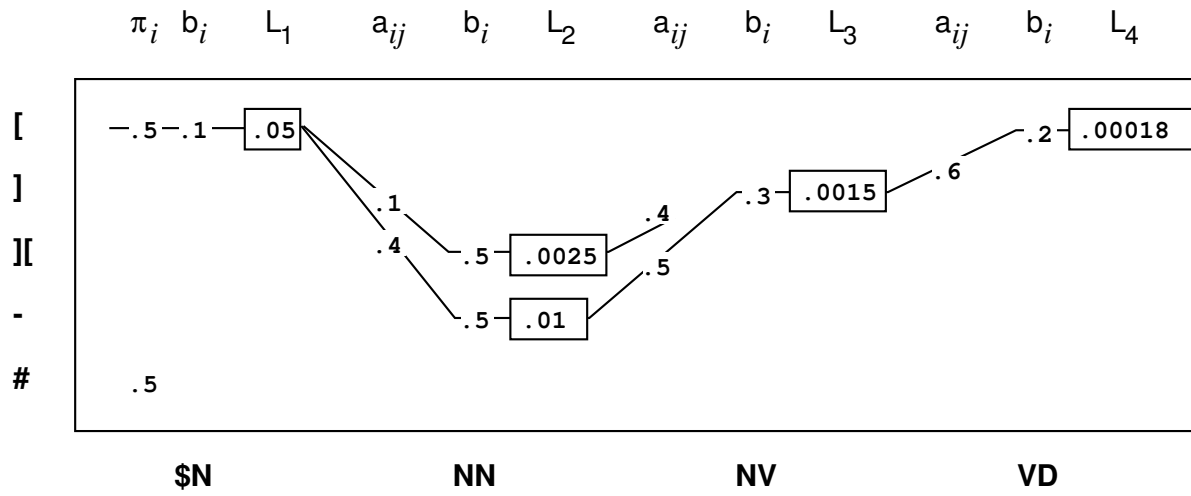
[]	I	-	#
.6	.5	.1	.4	
	.4		.6	.4
	.5	.1	.4	
.2				.8

	\$N	\$D	\$P	N\$	NN	ND	NV	NP	DN	V\$	VN	VD	VV	VP	P\$	PN	PD	PV
[.1	.2									.15	.2				.15	.2	
]				.25			.3	.4										.05
I					.5	.5												
-					.5			.1	.4									
#		.05	.3							.1			.1	.4	.05			

- $L(S)$



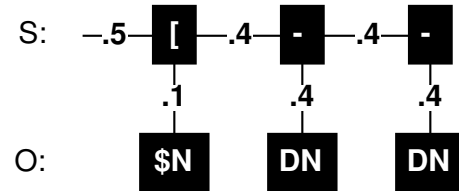
Viterbi with brackets



$\left[\begin{array}{cc} N & N \\ \text{computer} & \text{science} \end{array} \right] V \left[\begin{array}{c} D \\ \text{a} \dots \end{array} \right]$

Matching up pairs

- HMM does not guarantee that tag-pairs match up



- Define $L'(\mathbf{q}, \mathbf{o}) = \begin{cases} \alpha L(\mathbf{q}, \mathbf{o}) & \text{if } \mathbf{o} \text{ has matching tag-pairs} \\ 0 & \text{otherwise} \end{cases}$

– α is normalization constant to guarantee that

$$\sum_{\mathbf{q}, \mathbf{o}} L'(\mathbf{q}, \mathbf{o}) = 1$$

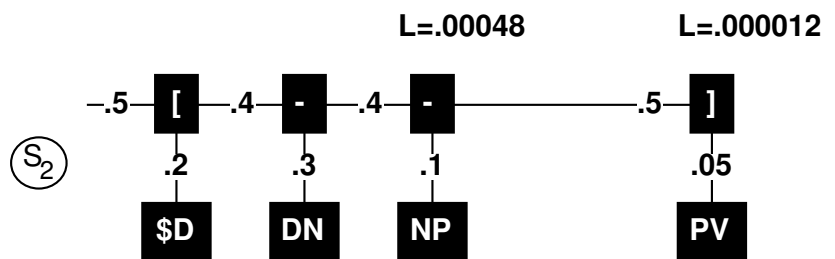
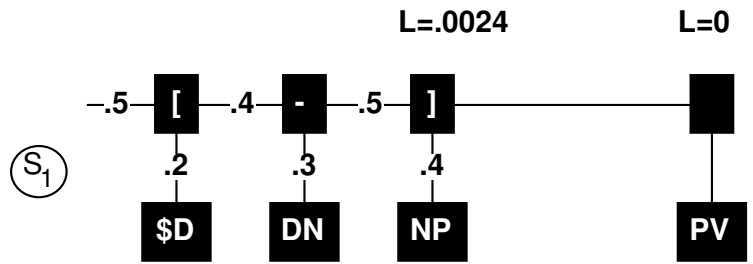


Phrase spotting

- Identifying individual phrases reliably
 - E.g. for terminology extraction
 - Aim: high precision, high recall, on individual phrases
Don't care about getting complete, consistent parse for sentences
 - Issues
 - Can't ignore context of candidate phrase
 - Can't directly compare $\lambda(\mathbf{q})$ and $\lambda(\mathbf{q}')$
 - How do we compute $\Pr(\mathbf{q}|\mathbf{o})$ for partial paths?
-
-

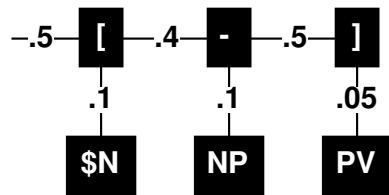
Phrase spotting

1. Can't just ignore context



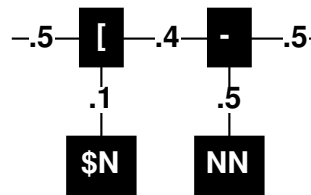
Phrase spotting

2. Can't just compare likelihoods

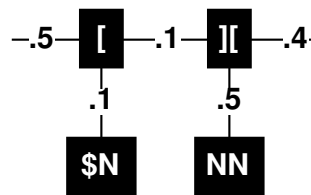


$L = .00005$
 $\Pr(S|O) = 1$

$L = .0015$
 $\Pr(S|O) = .833$



$L = .0003$
 $\Pr(S|O) = .167$



Relative likelihood

- The likelihood of being right, given the input

$$\begin{aligned}\Pr(\mathbf{q}|\mathbf{o}) &= \frac{\Pr(\mathbf{q}, \mathbf{o})}{\Pr(\mathbf{o})} \\ &= \frac{\Pr(\mathbf{q}, \mathbf{o})}{\sum_{\mathbf{q}'} \Pr(\mathbf{q}', \mathbf{o})} \\ &= \frac{L(\mathbf{q})}{\sum_{\mathbf{q}'} L(\mathbf{q}')}\end{aligned}$$

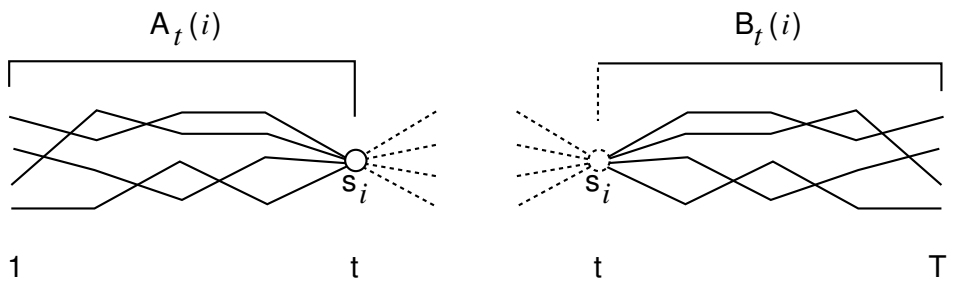
- For complete state-sequences, most-likely path is most-reliable

$$\operatorname{argmax}_{\mathbf{q}} \Pr(\mathbf{q}, \mathbf{o}) = \operatorname{argmax}_{\mathbf{q}} \Pr(\mathbf{q}|\mathbf{o})$$

- Not so for partial paths
-

Partial paths

- Prefix and suffix paths



$$A_t(i) = \{\mathbf{q} : 1, t | \mathbf{q}_t = s_i\}$$

$$B_t(i) = \{\mathbf{q} : t, T | \mathbf{q}_t = s_i\}$$

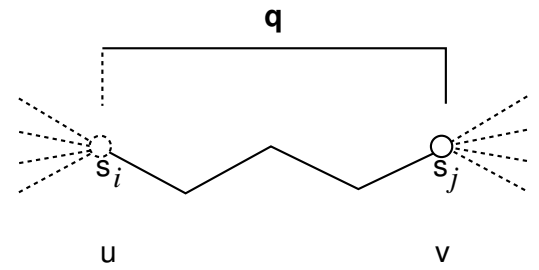
$$\alpha_t(i) = \sum_{\mathbf{q} \in A_t(i)} \lambda^0(\mathbf{q})$$

$$\beta_t(i) = \sum_{\mathbf{q} \in B_t(i)} \lambda(\mathbf{q})$$

Partial paths

- Partial-path likelihood

$$L(\mathbf{q}) = \Pr(\mathbf{q}, \mathbf{o}) = \alpha_u(i)\lambda(\mathbf{q})\beta_v(j)$$

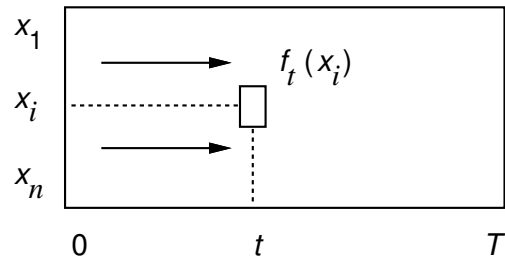


- Relative likelihood

$$\Pr(\mathbf{q}|\mathbf{o}) = \frac{L(\mathbf{q})}{\sum_{\mathbf{q}':u,v} L(\mathbf{q}')}$$

Dynamic Programming

- $f_t(x_i)$ only requires values for $f_u(x_j)$ for $u < t$
- t is variable of recursion
- Fill in array by increasing t



- Example: $\delta_t(i)$
-

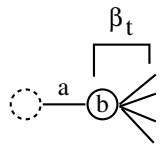
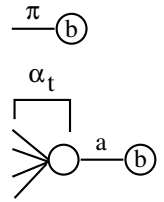
Recursive definitions for α , β

$$\alpha_1(i) = \pi_i b_i(\mathbf{o}_1)$$

$$\alpha_{t+1}(j) = \sum_i \alpha_t(i) a_{ij} b_j(\mathbf{o}_{t+1})$$

$$\beta_T(i) = 1$$

$$\beta_{t-1}(i) = \sum_j a_{ij} b_j(\mathbf{o}_t) \beta_t(j)$$



Remaining problems for phrase-spotting

- Dependent on global analysis
 - Search is linear-time, but can be moderately expensive if large numbers
 - Poor enough models of ‘garbage’ can damage estimates of $\Pr(\mathbf{q}|\mathbf{o})$ for \mathbf{q}
 - Can’t always reliably segment text into sentences
 - Integrating multiple information sources
-
-

Another problem: relative likelihood is not precision

- Some misspellings are undetectable at word level

combing appositive NPs
we had a rather milk winter

- Don't want to assume all words are misspelled (search)
 - Would like to detect problem by low relative likelihood
 - But if there's only one analysis, relative likelihood = 1, no matter how analysis
 - Precision is corpus-global measure of relative likelihood
E.g., of all the times we've seen "D Adv N N \$", how often has it been an
 - Have to estimate precision directly: it is neither likelihood nor relative likelihood
-
-

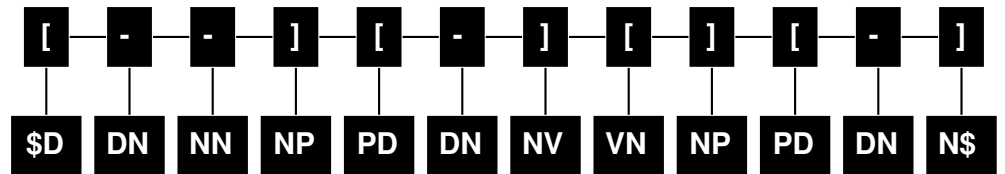
Parameter estimation

- With parsed corpus: count

$$a_{ij} = \Pr(Q_{t+1} = s_j | Q_t = s_i) \hat{=} \frac{f(Q_t = s_i, Q_{t+1} = s_j)}{f(Q_t = s_i)}$$

$$b_i(w) = \Pr(\mathbf{o}_t = w | Q_t = s_i) \hat{=} \frac{f(Q_t = s_i, \mathbf{o}_t = w)}{f(Q_t = s_i)}$$

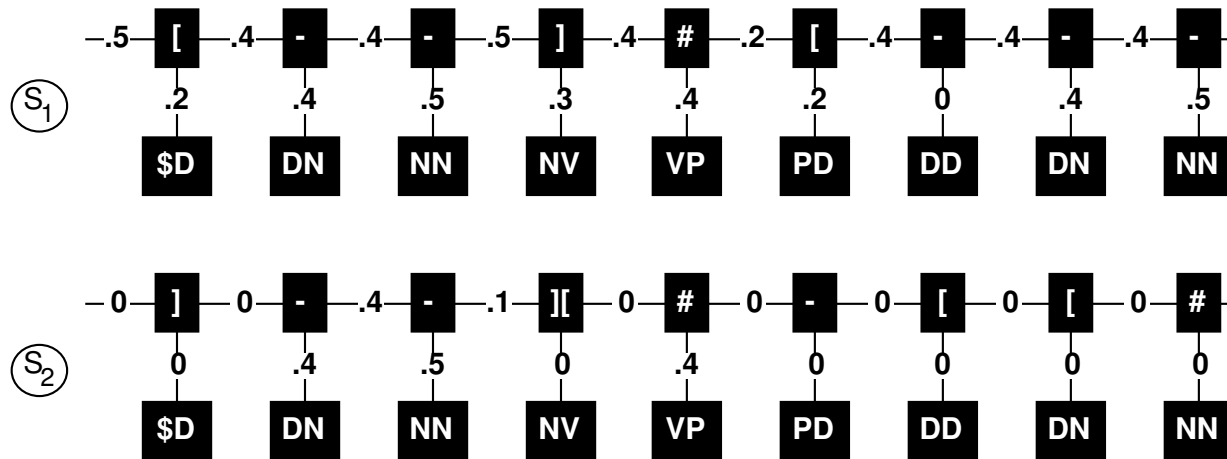
- Corpus is one giant observation sequence



$$\hat{a}_{[-]} = \frac{f([-])}{f([)} = \frac{3}{4} \quad \hat{b}_{[}(PD) = \frac{f([, PD)}{f([)} = \frac{2}{4}$$

Why zeros are a problem

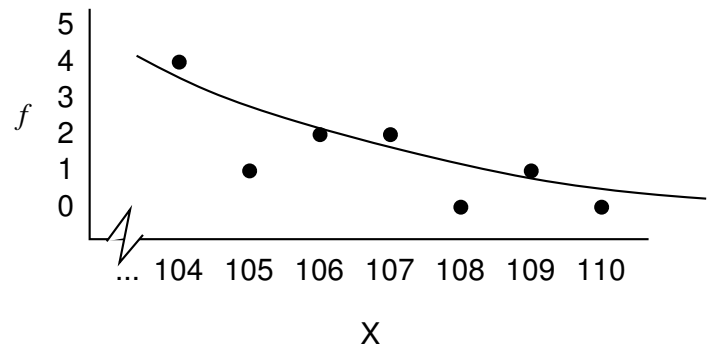
- Two structures with same likelihood: $L = 0$



- But if we replace '0' with '.01':
 $L(S_1) = e^{-24}$
 $L(S_2) = e^{-70}$

Smoothing

- Choosing a good value to replace the zeros
- From choosing a smooth curve:



Good-Turing [59]

f	n_f	$f \cdot n_f$
9	22,280	200,520
8	27,710	221,680
7	35,709	249,963
6	48,190	289,140
5	68,379	341,895
4	105,668	422,672
3	188,933	566,799
2	449,721	899,442
1	2,018,046	2,018,046
0	74,671,100,000	0

$$\bar{f} \cdot n_f = (f + 1) \cdot n_{f+1}$$

$$\bar{f} = \frac{(f + 1) \cdot n_{f+1}}{n_f}$$

Cat-cal

- Categorize and calibrate
- Some of the events with 0 counts in training have > 0 counts in test
- Group by count

$$G_e = \{e' | f(e') = f(e)\}$$

- Re-estimate counts for groups from cross-validation corpus
- Re-estimate individual counts as group count times probability of choosing of group

$$\bar{f}(e) = \bar{f}(G_e) \cdot \Pr(e|G_e)$$

Cat-cal

	Corpus 1	Corpus2	$\bar{f}(G_i)$	$\Pr(e G_e)$	\bar{f}	
G_2	[-	2	3	3	1	3
	[]	1	2		.3	1.2
G_1	- -	1	2	4	.3	1.2
	-]	1	0		.3	1.2
G_0] [-	0	0		.2	.4
	# #	0	0		.2	.4
] #	0	0	2	.2	.4
	#]	0	0		.2	.4
] []	0	0		.2	.4
] []	0	0		.2	.4

Without Parsed Corpus

- Probability of transition from s_i to s_j at t to $t + 1$

$$\Pr(Q_t = s_i, q_{t+1} = s_j | \mathbf{o}) = \Pr(\mathbf{q} | \mathbf{o}) \text{ for } \mathbf{q} : t, t + 1, \mathbf{q}_t = s_i, \mathbf{q}_{t+1} = s_j$$

- Probability of being in s_i at t

$$\Pr(Q_t = s_i | \mathbf{o}) = \Pr(\mathbf{q} | \mathbf{o}) \text{ for } \mathbf{q} : t, t, \mathbf{q}_t = s_i$$

Without parsed corpus

- Use relative likelihood of transitions/emissions
- Suppose $\Pr(s_i \rightarrow_t s_j | \mathbf{o}) = .25$
 - Then if the Markov process generates \mathbf{o} 100 times, we expect it to see s_i
 - Equivalently, we take $\Pr(s_i \rightarrow_t s_j | \mathbf{o})$ as a fractional count
- Sum across time positions

$$f(s_i \rightarrow s_j | \mathbf{o}) = \sum_t \Pr(s_i \rightarrow_t s_j | \mathbf{o})$$

- Use same re-estimation formulae as for parsed corpus

$$a_{ij} = \Pr(Q_{t+1} = s_j | Q_t = s_i) \hat{=} \frac{f(Q_t = s_i, Q_{t+1} = s_j)}{f(Q_t = s_i)}$$
$$b_i(w) = \Pr(\mathbf{o}_t = w | Q_t = s_i) \hat{=} \frac{f(Q_t = s_i, \mathbf{o}_t = w)}{f(Q_t = s_i)}$$

Iteration

- To compute $\Pr(\mathbf{s}_i \rightarrow \mathbf{s}_j | \mathbf{o})$, etc., we need initial guess

$$M_0 = (a_0, b_0, \pi_0)$$

- Iterate using fractional counts to get M_{i+1} from M_i
- Likelihood of model

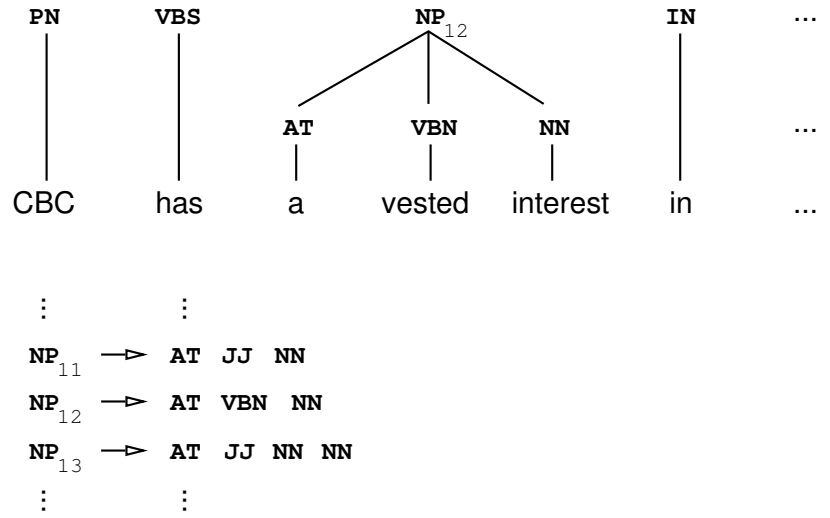
$$L(M) = \Pr(\mathbf{o}, M) = \sum_{\mathbf{q}} \Pr(\mathbf{q}, \mathbf{o}, M)$$

- It can be shown that

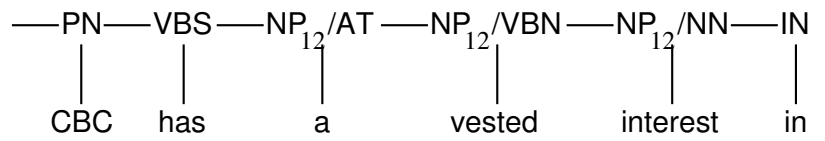
$$L(M_{i+1}) \geq L(M_i)$$

- But:
 - Local maximum
 - Overtraining
-
-

Root



Can be mapped to a standard HMM:



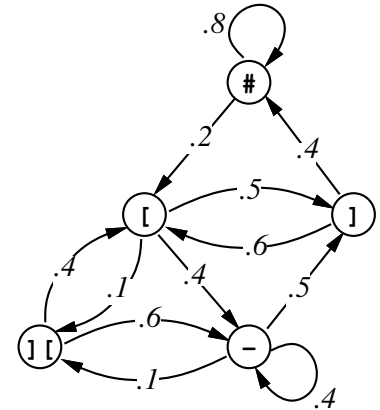
Extensions

- Could also ‘tie’ states
 - E.g. set $b_{\text{NP}_{12}/\text{AT}} = b_{\text{AT}}$
 - Estimate

$$\hat{b}_{\text{NP}_{12}/\text{AT}}(w) = \hat{b}_{\text{AT}}(w) = \frac{f(\text{NP}_{12}/\text{AT}, w) + f(\text{AT}, w)}{\sum_{w'} [f(\text{NP}_{12}/\text{AT}, w') + f(\text{AT}, w')]}$$

- Generalizing to categories other than NP
 - Leads to: finite-state chunks
-
-

An HMM is a (stochastic) FSA



	[]	l	-	#
[.5	.1	.4	
]	.6				.4
l		.4	.6		
-		.5	.1	.4	
#	.2				.8

Composing FSA's

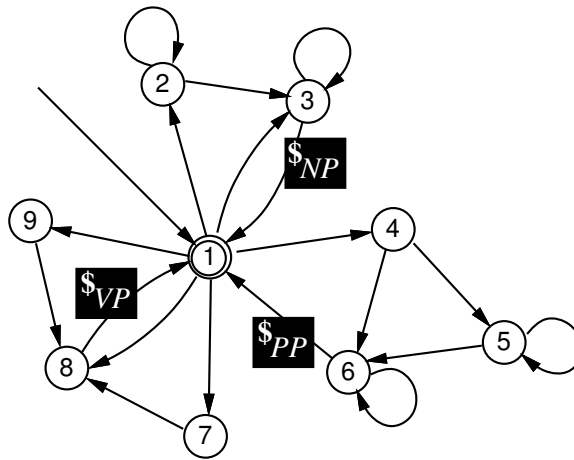
NP = D? Adj* N+ $\$_{NP}$

PP = P NP $\$_{PP}$

VP = (V |Hv Vbn |Be Vbg) $\$_{VP}$

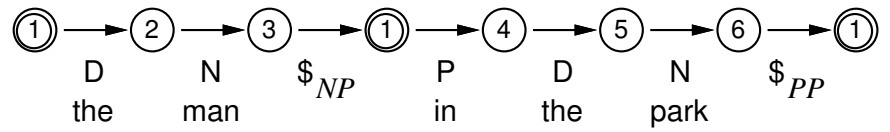
Chunk = NP |PP |VP

S → Chunk+

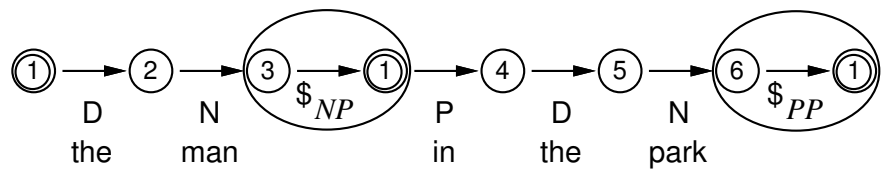


Matching

- Works great if the \$'s are in the input

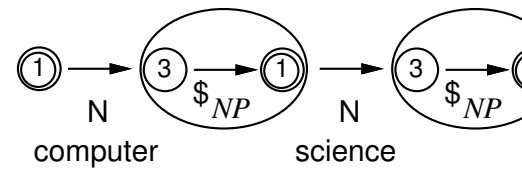
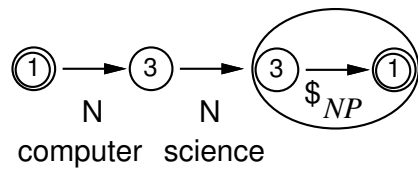


- Fold \$'s into surrounding states



Result

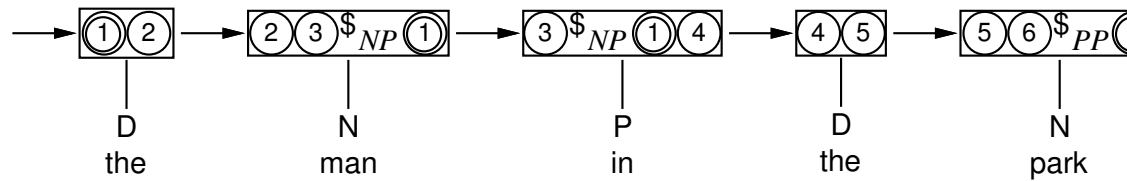
- Add new states $3\$_{NP}1$, $6\$_{NP}1$, $8\$_{NP}1$
- Tie transitions to transitions from original state 1
- Now non-deterministic



- Parse is uniquely recoverable from state-sequence
-
-

Final step

- FSA scans on arcs, HMM emits on states
- Turn state-pairs into states



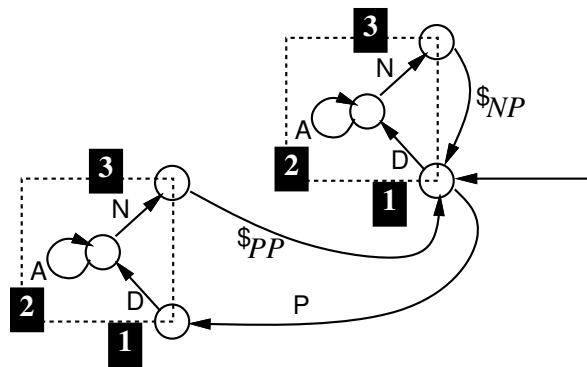
- Transition from $[ij]$ to $[jk]$ corresponds to transition from j to k in the un
 - Initial probability of $[1i]$ represents probability of transition from initial sta
-
-

Cascaded FSA's

- More of the same medicine

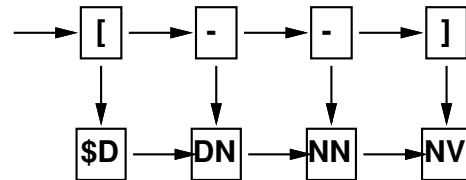
Clause \rightarrow PP* NP PP* VP NP? PP* . $\$_{Clause}$

- Insert a copy of the PP regex at each place there's a PP
- Build a large FSA from the resulting regex
- Tie corresponding transitions in different copies of sub-regex

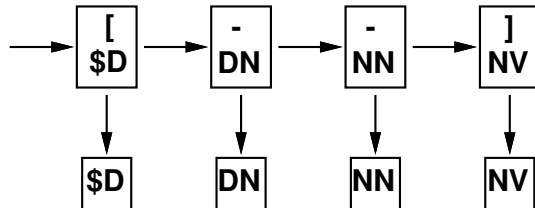


Multiple information sources

- Suppose choice of bracket depends on preceding bracket and preceding tag



- Remember, we cannot do: $\Pr(\mathbf{o}_{t+1} | \mathbf{q}_{t+1}, \mathbf{o}_t) = \Pr(\mathbf{o}_{t+1} | \mathbf{q}_{t+1}) \Pr(\mathbf{o}_t | \mathbf{q}_{t+1})$
- We must estimate the entire distribution $\Pr(\mathbf{o}_{t+1}, \mathbf{q}_{t+1}, \mathbf{o}_t)$
- In effect, we must fold together all information sources into single state

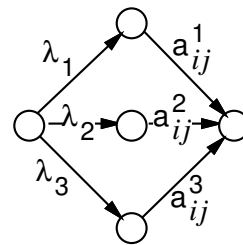


Also for phrase spotting

- Integrate multiple info sources in estimating a_{ij} , $b_i(w)$
 - Folding info sources together leads to state-space explosion, sparse data
 - Combine information from features of state to estimate transition/emission
 - Integrate multiple info sources in estimating precision of phrase-spotting p
 - Longest match vs. longer-same-cat vs. longer-other-cat vs. overlapping
 - Collocation score
 - Tagging score
 - Phrase type
 - Etc.
-
-

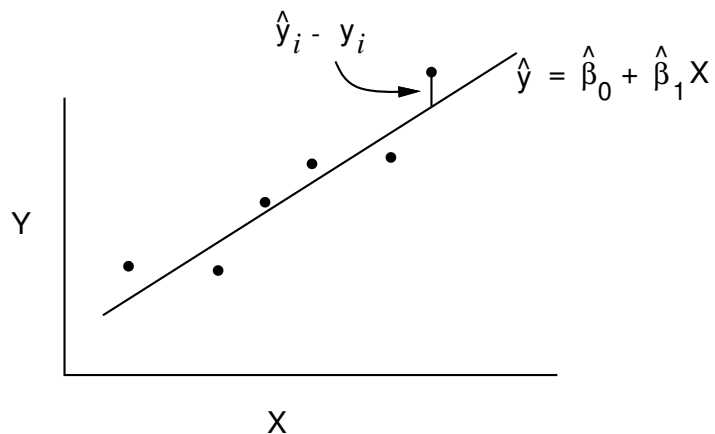
Linear interpolation

- Separately train submodels M_1, M_2, \dots
- E.g., M_1 is an HMM that only looks at previous bracket, and M_2 looks at tag
- Combine into single model
 - Hold a_{ij}^k fixed
 - Train λ_k
 - Transition probability in combined HMM is $\sum_k \lambda_k k_{ij} = \sum_k \Pr(M_k) \Pr$



Regression

- “Regression analysis is the part of statistics that deals with investigation of the relationship between two or more variables related in a nondeterministic fashion” [68]
- For example: linear regression



$$Y = \beta_0 + \beta_1 X + \epsilon$$
$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

- Estimating β_0, β_1 : minimize squared error $\sum(\hat{y} - y)^2$
 - Minimum can be determined analytically from observed pairs (x_i, y_i)
 - For given value x , we have point estimate \hat{y} and probability distribution \hat{p}
-

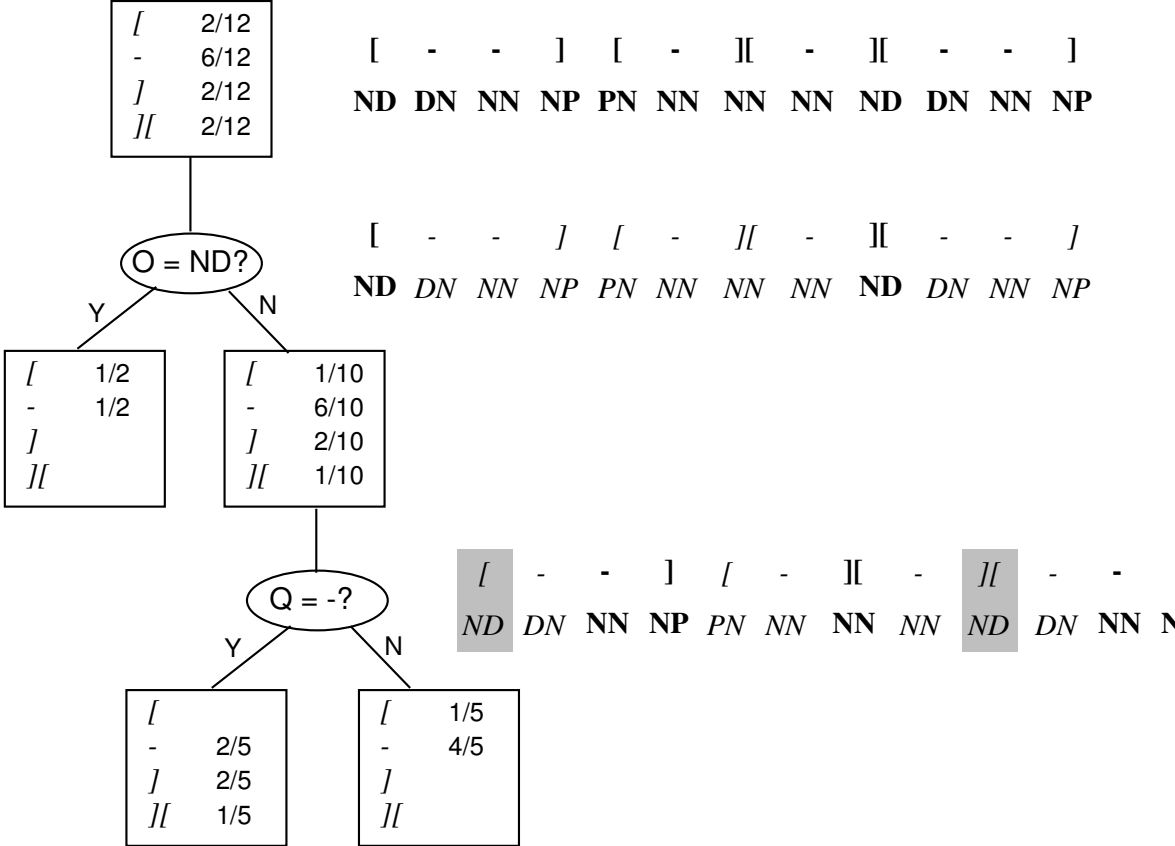
Multivariate regression

- Combining info from multiple variables

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_n X_n + \epsilon$$

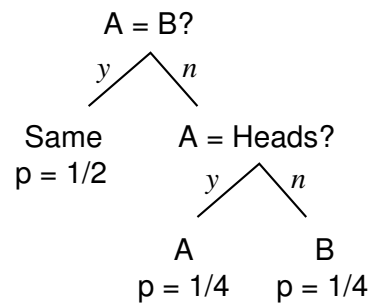
- X_i are *predictor variables*
 - Estimate β_i by minimizing squared error
 - To do so, need observations $(x_{1i}, \dots, x_{ni}, y_i)$
 - For given values $\langle x_1, \dots, x_n \rangle$ of predictor variables, we have point estimate \hat{y} for Y
 - Only useful if relationship is approximately linear (though polynomial generalizations exist)
-
-

Regression trees [38, 20]



How good is a split?

- We want to reduce uncertainty about dependent variable
- Uncertainty = entropy
- 1 bit = the uncertainty in one equally-likely two-way guess
- E.g. flip two coins: Same, A, B



Entropy

- Point entropy η – number of 2-way choices to reach given result

$$\begin{aligned}\eta(\text{Same}) &= 1 \\ \eta(\text{A}) &= 2 \\ \eta(\text{B}) &= 2\end{aligned}$$

- Probability p of ending up at result

$$\begin{aligned}p(\text{Same}) &= 1/2 \\ p(\text{A}) &= 1/4 \\ p(\text{B}) &= 1/4\end{aligned}$$

- Entropy is average number of 2-way choices = weighted average of η

$$\begin{aligned}&= p(\text{Same})\eta(\text{Same}) + p(\text{A})\eta(\text{A}) + p(\text{B})\eta(\text{B}) \\ &= \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{4} \cdot 2 \\ &= 1.5\end{aligned}$$

Entropy

- In binary-branching tree of uniform depth η containing N leaves

$$\begin{aligned} N &= 2^\eta & \eta &= \log_2 N \\ p &= \frac{1}{N} & \text{i.e., } N &= \frac{1}{p} \\ & & \eta &= \log_2 \frac{1}{p} \end{aligned}$$

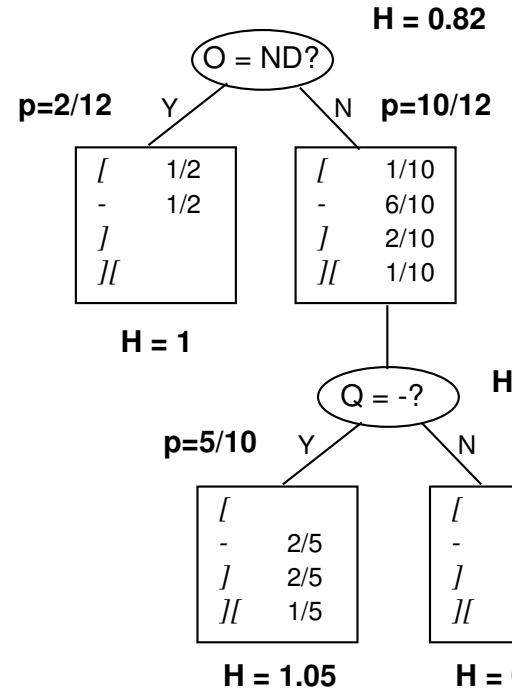
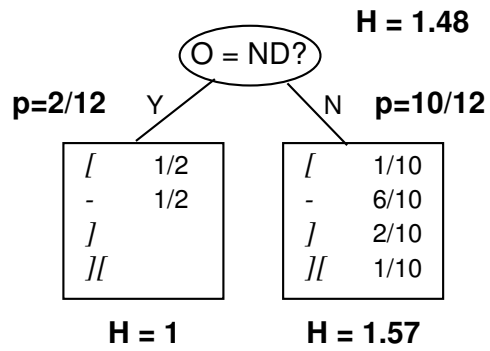
- The same relation can be used generally

$$\eta_i = \log_2 \frac{1}{p_i} \quad H = \sum_i p_i \eta_i$$

- Entropy is maximized when all choices are equally likely (maximum uncertainty)
 - The more skewed the distribution, the lower the entropy, the lower the uncertainty
-
-

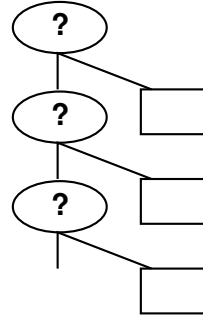
Goodness of split

- Goodness of split is reduction in uncertainty: $1.48 - 0.82 = 0.66$



Decision lists

- Binary decision tree in which one daughter of every node is a leaf



- Alternative to greedy algorithm (Yarowsky [198])
 - Discriminator: question + answer (Y/N)
 - Evaluate each discriminator independently on *all* data
 - Goodness of discriminator is inverse to uncertainty of resulting leaf dist.
 - Sort discriminators by goodness to create decision list
-
-

Transformation-based regression (Brill [41])

- Initial assignment rules

E.g., assign most frequent bracket to tag-pairs

- Error-correction rules $Y \rightarrow Y' / X_1 = x_1, \dots, X_n = x_n$
- Predictor variables: X_1, \dots, X_n and Y
- Dependent variable: $Y' = Y$ at $t + 1$

- Iterate

- Evaluate all potential rules
- Choose best (greedy)
- Apply, creating a new corpus

- Evaluation

- Reduction in error rate
- Errors in corpus after applying rules
- Corpus errors before applying rules

- Like decision lists, trains on all data
 - Only gives point estimate, not distribution
-

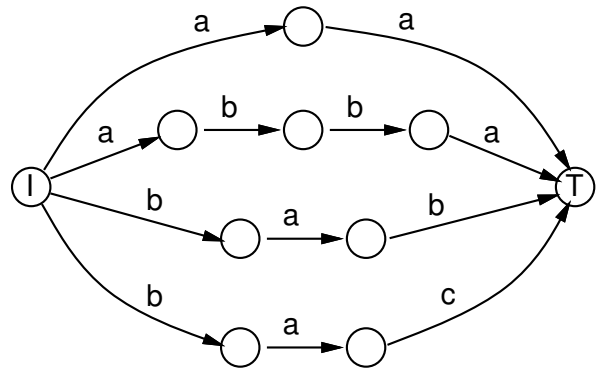
Summary

- User identifies relevant attributes (predictor variables)
 - Automatic search through space of discriminators (boolean combinations of predictor variables)
 - Point estimate and probability distribution
 - State = set of values for predictor variables
 - Discriminator = set of states
-
-

Grammatical inference: Regular grammars

- Canonical grammar exactly generates training corpus

aa
abba
bab
bac



- Prior and posterior
 - Canonical grammar has perfect fit to data
Highest conditional probability $\Pr(\mathbf{o}|G)$
 - Canonical grammar generally is overly complex
Low prior probability $\Pr(G)$
 - Likelihood is posterior probability $\Pr(\mathbf{o}, G) = \Pr(\mathbf{o}|G) \Pr(G)$
 - Search for maximum-likelihood grammar
 - Operation on grammar: merge two states into one
 - Greedy search
 - Consider each pair of states
 - Compute posterior probability if we merge this pair
 - Choose best pair, merge, iterate
 - Quit if no pair improves likelihood
-
-

Context-free grammars

- Canonical grammar: one production for each sentence

$S \rightarrow \text{sentence}_1$

$S \rightarrow \text{sentence}_2$

\vdots

- Operators

- Merge nonterminals

- Structuring

Substitute (new) nonterminal X everywhere for sequence Y_1, \dots, Y_n

Add new rule $X \rightarrow Y_1, \dots, Y_n$

Infering partial grammars: collocations

- Chuch, Gale, Hanks & Hindle [60]

- Use MI to induce \sim selectional restrictions

drink —: \langle Qty \rangle beer, tea, Pepsi, champagne, liquid, . . .

- Preprocess with Fidditch to find head-head pairs

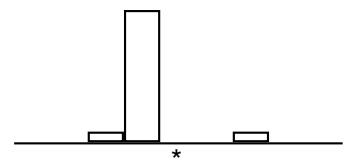
- Smadja [177, 176]

- Use strength of association \sim MI

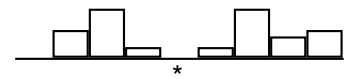
- Also use entropy of positional distribution

doctor:

honorary



nurse



- Postprocess with Cass
-
-

- Word distribution vectors

	a__	aardvark__	...	__zoologic	__zygote
fish	216	0	...	0	2
habitat	1	5	...	0	0

- Measures of vector (dis)similarity

Manhattan, Euclidean, dot product, cosine, correlation, rank correlation, divergen

- Cluster words using one of the distance metrics to form parts of speech
- Compute distribution vectors for part-of-speech sequences
- Cluster part-of-speech sequences to form phrase classes

E.g. 'NP': **C8** (*it*), **C8 C3** (*her status*), **C1 C91 C3** (*the following section*), ...

- Special role for function words
- Identify function words by high frequency
 - Another way: bursty → content word (Gale, p.c.; Jones & Sinclair [122])
- Cluster function words

F0: *a, an, her, his, ...*

F1: *he, I, she, then, ...*

F2: *are, be, had, has, ...*

- Form chunks & chunks

F0 C C C F7 F0 C
a tiny bird sat in the tree

- Collect content-word contexts

tiny: F0 _ C C F7

bird: F0 C _ C F7

- Cluster contexts to form content-word categories

F0 C24 C51 C40 F7 F0 C24 C51
a tiny bird sat in a hollow tree

- Build chunk & chunk grammar

FP0 → F0 C24 C51 C40 F7

FP1 → F7 F0

FP2 → F0 C24 C51 F\$

- Generalize using substitution operator

CP1 → C24 C51

American structuralists

- Two measures of phrasehood
 - Substitution (distributional similarity)
 - Cohesiveness
- Substitution

he ~ the man $\left\{ \begin{array}{l} \text{—} \textit{laughed} \\ \text{—} \textit{saw him} \\ \textit{he saw } \text{—} \end{array} \right.$

- Also used by Brill to induce trees
 - Current information-theoretic instantiation:
 - Substitution = divergence
 - Cohesiveness = mutual information
-
-

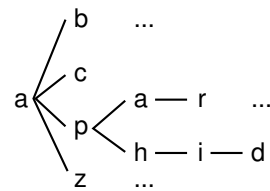
Zellig Harris [99]

- American structuralist
 - Sought objective, operational definitions for linguistic concepts
 - Phoneme, morpheme, word, phrase
- “From phoneme to morpheme” [99]

– Look at number possible continuations for a word prefix

ap— { a(rtr
e(rtu
h(id,
:
}

– Within morpheme, number of possible continuations decreases because



– Jumps back up at boundary

Example

	h__	he__	hes__	hesc__	hescl__	hescle__	hescleve__	hescleve__	
a	and	al	afraid	alm	ad	an			
b		built	bad						
c		came	clever						
d		dgo	dead						
e	elp	ehee	ells	entered	ever		er		
f		ft	orit			ft			
x		xagon	eroxel						
y	ype	y	yodeled	ythed	yde				
z		zoomed	zoomed						
	6	26	26	9	6	[7]	1	1	

Harris

- Do it backwards, too

Agreement itdisturb__smethatheleft

Cranberry words cran__berry

Ambiguous prefix hed__eparatelyneedsit

- Only practical way of getting utterances is elicitation
-
-

- Chomsky: “We can be fairly certain that there will be no operational criteria for any but the most elementary [linguistic] notions”
- Seeks operational definition for *phrase* nonetheless
- Phrase = sequence of word-categories co-occurring more frequently than exp
- “Bond”

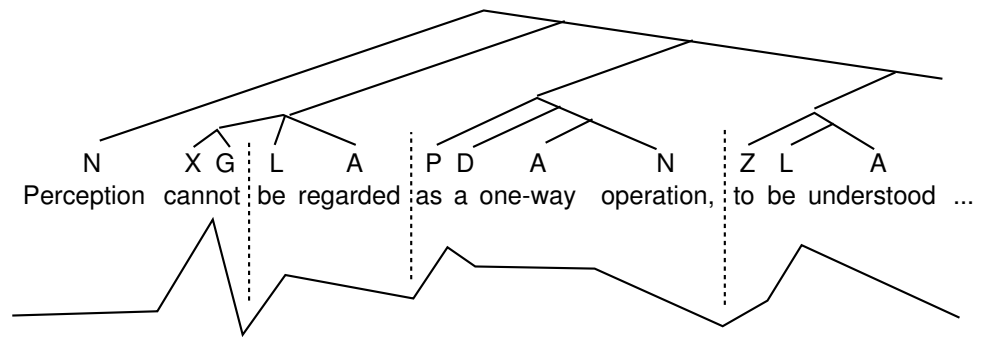
$$B_F(i) = \frac{\Pr(t_{i+1}|t_1, \dots, t_i)}{\Pr(t_{i+1})} \quad B_B(i) = \frac{\Pr(t_{i-1}|t_i, \dots, t_n)}{\Pr(t_{i-1})}$$
$$B(i) = \frac{1}{2}[B_F(i) + B_B(i)]$$

Note: $\log B_F(i) = I(t_{i+1}; t_1, \dots, t_i)$

- Phrase boundaries at minima in B
-
-

Stolz

- Estimates: hand-counted all cat-sequences in a 68,000-word corpus
- Test: 13 sentences from *Scientific American*
- Hand-parsed, differences arbitrated among three judges
- Example



- Sequences of categories

$$B(i) = \log \frac{\Pr(t_1, \dots, t_i | t_{i+1}, \dots, t_n)}{\Pr(t_{i+1})} = \log \frac{\Pr(t_{i+1}, \dots, t_n | t_1, \dots, t_i)}{\Pr(t_1, \dots, t_i)}$$

- Estimate as product of n-gram MI for windows around i
- Find minimum in window, truncate sentence, repeat

$t_1 \ t_2 \ t_3 \ | \ t_4 \ t_5 \ \dots$
 $t_1 \ t_2 \ t_3 \ || \ t_4 \ t_5 \ | \ t_6 \ t_7 \ t_8 \ \dots$

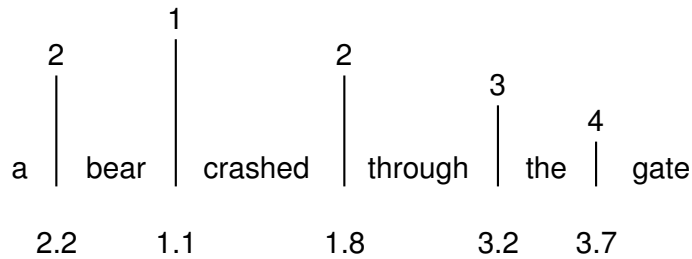
- Alternative beginning and end of sentence
 - Recurse to find constituents inside these
-

Bottom line

- Works OK for lowlevel phrases
 - Important that one use categories, not words
 - Else lexical association pulls phrases apart
a strong interest__in
 - Function words predict following function words better than following c
of__the wilderness
 - Result
an interest__in pictures__of__the Tetons
 - Less good at higher levels of structure: here lexical associations are needed
-
-

Operational definitions of phrases

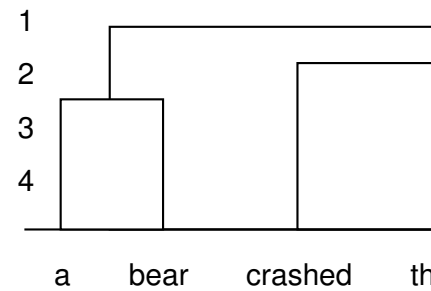
- Performance structures
- Naive parsing [96]
 - Subjects divide sentence, redivide



- Take average prominence of boundary across subjects

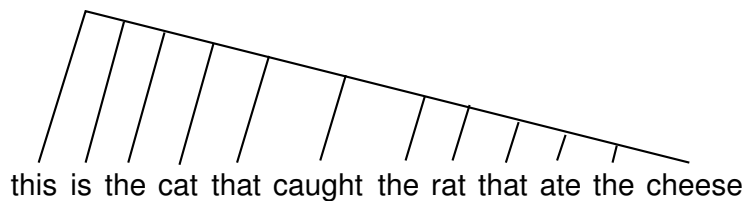
- Also: transitional error probabilities, pausing, sentence comprehension

– Dendograms (performance)



Performance structures

- Differ from traditional phrase structures
 - Flatter, no deep right branching
 - Chunk boundaries stable, higher-level boundaries less syntactically predicted
- Prosodic phrases differ from traditional phrases in the same way



- Selkirk: ϕ -phrases [172]
 - Gee & Grosjean [92]: use ϕ -phrases to predict performance structures
 - Bachenko & Fitzpatrick [18] turn it around and use Gee & Grosjean algorithm for intonation for text-to-speech
-
-

Linguistics

- The levels sentence, clause, phrase, word are traditional
- Quirk et al. [159] have VP stop at verb

[*NP* The weather] [*VP* has been] [*AdjP* remarkably warm]

- Postmodifiers of nouns often assumed Chomsky-adjoined

[*NP* [*NP* the man] [*PP* in the park]]

- Bloch 1946 [31] defines phrases prosodically: “pause-groups”

a little dog , with a big bone
*a little , dog with a big , bone

Function Words

- Suzuki (1824)

- *si*: noun, verb, adjective – “[*si*] denotes something”

- *zi*: particles – “[*zi*] denotes nothing; it only attaches ‘voice of the heart’

- Aristotle

- Words without meaning: complementizers, conjunctions, etc.

- Words with meaning: nouns, verbs, adjectives

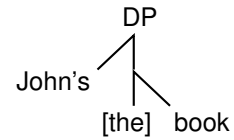
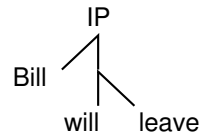
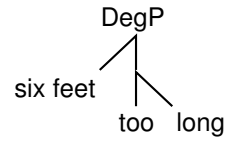
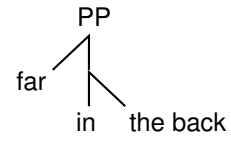
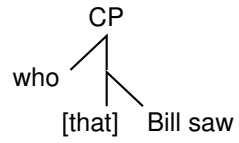
- Psychology

- Some aphasias selectively affect function words or content words

- Slips of the tongue interchange F-F, C-C, but not F-C

Uniform syntactic treatment

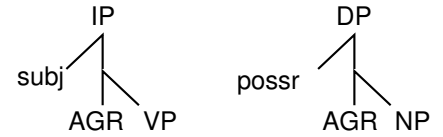
- Function words have subjects and complements [5]



Motivation: Agreement in noun phrase

- English: Tensed verb is *first* verb, not e.g. head:

leaves
was leaving
has been leaving



- Yup'ik: noun phrase has AGR, too

angute-m	kiputaa- \emptyset	“the man bought it”
angute-t	kiputaa-t	“the men bought it”
angute-m	kuiga- \emptyset	“the man’s river”
angute-t	kuiga-t	“the men’s river”

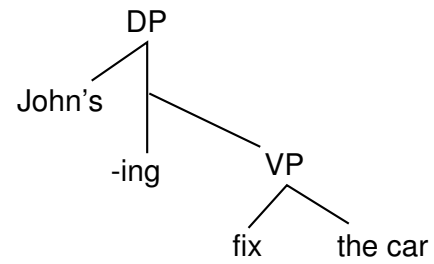
- Turkish

el	“hand”
senin el-in	“your hand”
onun el-i	“his hand”

Motivation: Gerund

- The Poss-Ing gerund is a gryphon

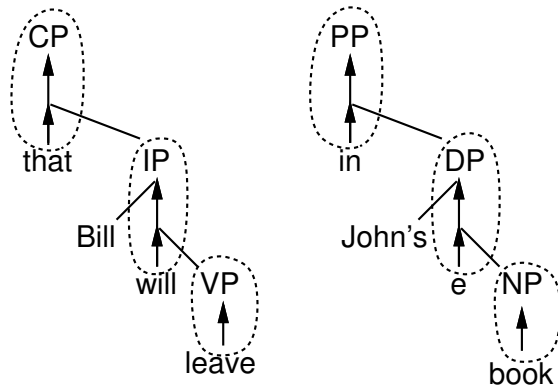
[*NP* John's [*VP* fixing the car]]



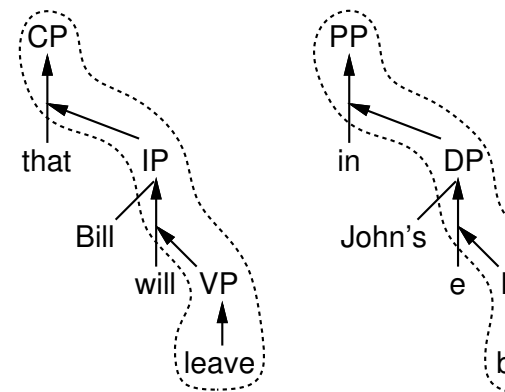
S-projection

- Whether to “count” function words as heads

fine grain (c-projection)

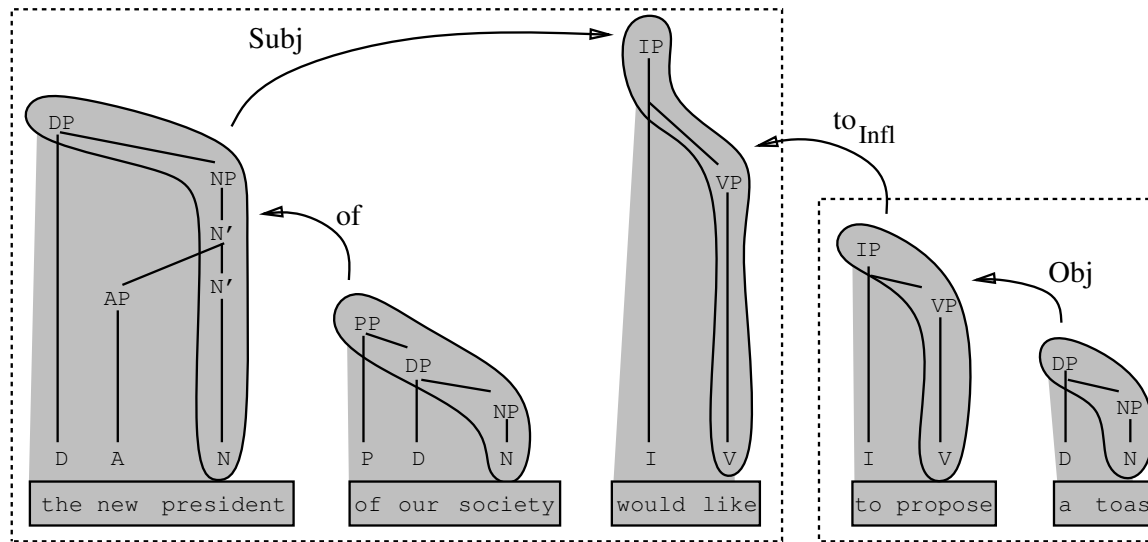


coarse grain (s-projection)



Chunks and clauses

- Chunk: connected piece of tree covered by an s-projection
- Clause: chunks dominated by same clausal node



Syntactic usefulness of chunks

- No chunk within a chunk [7]

* [a proud [of his son] man]	[a man] [proud] [of his
* [a [so tall] man]	[so tall] [a man]
* [a [six feet] tall man]	[six feet] [tall], [a six-foot]
* [was [every three weeks] fixing] his bike	[was frequently fixing]

- More precisely, F-C selection must be in same chunk
-
-

General [2, 3, 4, 35, 36, 50, 61, 62, 81, 82, 84, 116, 117, 118, 129, 143, 144, 148, 200]

Tagging [10, 19, 28, 56, 57, 66, 90, 91, 124, 125, 126, 131, 138, 153, 163, 168, 188]

HMMs [21, 22, 23, 24, 25, 49, 64, 67, 78, 115, 119, 155, 157, 160, 161]

Search [156]

The Inside-Outside Algorithm [85, 86, 136, 137]

Regression [20, 30, 29, 38, 41, 42, 45, 46, 154, 162]

Partial Parsing [6, 7, 8, 9, 11, 37, 43, 47, 48, 51, 52, 53, 57, 58, 112, 65, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 100, 101, 102, 103, 104, 107, 110, 113, 114, 120, 121, 127, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 149, 152, 163, 164, 165, 166, 169, 178, 182, 186, 190, 191, 192, 194, 195, 196, 197, 198, 199]

Grammatical Inference, Acquisition [1, 12, 13, 14, 15, 16, 32, 33, 39, 40, 55, 58, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 111, 130, 167, 175, 179, 181, 184, 187, 189, 199]

Mutual Information Parsing [98, 99, 146, 185]

Prosody and Performance Structures [18, 26, 27, 31, 63, 92, 96, 97, 105, 106, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 193]

References

- [1] AAI. *Proceedings on Machine Learning of Natural Language and Ontology*. GMD (DFKI), 1991. Spring Symposium. DFKI Publication D-91-09.

- [2] AAI. *Workshop on Statistically-Based NLP Techniques*, July 1992. Workshop meeting.
- [3] AAI. *Fall Symposium on Probability and Natural Language Processing*, 1993.
- [4] Jan Aarts and T. van den Heuvel. Computational tools for the syntactic analysis. *Linguistics*, 23:303–335, 1985.
- [5] Steven Abney. *The English Noun Phrase in its Sentential Aspect*. PhD thesis, MIT, MA, 1987.
- [6] Steven Abney. Rapid incremental parsing with repair. In *Proceedings of the 6th Conference: Electronic Text Research*, pages 1–9, Waterloo, Ontario, October 1990. Waterloo.
- [7] Steven Abney. Syntactic affixation and performance structures. In D. Bouchard, editors, *Views on Phrase Structure*. Kluwer Academic Publishers, 1990.
- [8] Steven Abney. Parsing by chunks. In Robert Berwick, Steven Abney, and Carol Fisiak, editors, *Principle-Based Parsing*. Kluwer Academic Publishers, 1991.
- [9] Steven Abney. Chunks and dependencies: Bringing processing evidence to bear. *Computational Linguistics and the Foundations of Linguistic Theory*. CSLI, To appear.
- [10] Eric Akkerman, Hetty Voog–van Zutphen, and Willem Meijs. *A Computerized Level Tagging. ASCOT Report 2*. Language and Computers: Studies in Practical Linguistics 1. Edited by Jan Aarts and Willem Meijs. Rodopi, Amsterdam, 1988.
- [11] H. Alshavi. Processing dictionary definitions with phrasal pattern hierarchies. *Linguistics*, 13:195–202, 1987.
- [12] A. Andreewsky, C. F. Fluhr, and F. Debili. Computational learning of semantic structures for the generation and automatical analysis of content. *Information Processing*, 1990.

- [13] Angluin and Smith. Inductive inference: Theory and methods. *ACM Computing Surveys*, 15:319–347, 1983.
- [14] D. Angluin. Inductive inference of formal languages from positive data. *Information Processing Letters*, 45:117–135, 1980.
- [15] D. Angluin. Learning regular sets from queries and counterexamples. *Information Processing Letters*, 75:87–106, 1987.
- [16] Peter Anick and James Pustejovsky. An application of lexical semantics to knowledge-based parsing. In *COLING 90, vol. 2*, pages 7–12, 1990.
- [17] Damaris Ayuso et al. Bbn: Description of the PLUM system as used for MUC-4. In *Proceedings, Fourth Message Understanding Conference (MUC-4)*, pages 169–176, San Francisco, CA, 1991. Morgan Kaufmann.
- [18] Joan Bachenko and Elizabeth Fitzpatrick. A computational grammar of discourse phrasing in English. *Computational Linguistics*, 16(3):155–170, 1990.
- [19] L. R. Bahl and R. Mercer. Part-of-speech assignment by a statistical decision tree. In *International Symposium on Information Theory*, Ronneby, Sweden, 1976.
- [20] Lalit R. Bahl, Peter F. Brown, Peter V. de Souza, and Robert L. Mercer. A tree-based language model for natural language speech recognition. *IEEE Trans. on Acoustics, Speech, and Signal Processing*, pages 507–514, 1991.
- [21] L.R. Bahl, F. Jelinek, and R.L. Mercer. A maximum likelihood approach to continuous speech recognition. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, PAMI-13(3):358–369, 1991.
- [22] L.E. Baum. An inequality and associated maximization technique in statistical estimation of probabilistic functions of markov processes. *Inequalities*, 3:1–8, 1972.

- [23] L.E. Baum and J.A. Egon. An inequality with applications to statistical estimation functions of a markov process and to a model for ecology. *Bull. Amer. Meterol.* 1967.
- [24] L.E. Baum, T. Petrie, G. Sopules, and N. Weiss. A maximization technique statistical analysis of probabilistic functions of markov chains. *Annals of Math* 41:164–171, 1970.
- [25] L.E. Baum and G.R. Sell. Growth functions for transformations on manifolds 27(2):211–227, 1968.
- [26] John Bear and Patti Price. Prosody, syntax and parsing. In *28th Annual Meeting for Computational Linguistics*, pages 17–22, 1990.
- [27] Mary Beckman and Janet Pierrehumbert. Intonational structure in japanese and e *Yearbook*, 3:255–310, 1986.
- [28] J. Benello, A. Mackie, and J. Anderson. Syntactic category disambiguation with *Computer Speech and Language*, 3(3), 1989.
- [29] Ezra Black, F. Jelinek, J. Lafferty, R. Mercer, and S. Roukos. Decision tree mod labeling of text with parts-of-speech. In *Darpa Workshop on Speech and Natur* Mateo, CA, 1992. Morgan Kaufman.
- [30] Ezra Black, Fred Jelinek, John Lafferty, David M. Magerman, Robert Mercer, and Towards history-based grammars: Using richer models for probabilistic parsing *Meeting of the Association for Computational Linguistics*, Columbus, Ohio, 1993
- [31] Bernard Bloch. Studies in colloquial Japanese II: Syntax. *Language*, 22:200–248,
- [32] Blumer, Ehrenfeucht, Haussler, and Warmuth. Learnability and Vapnik-Chervon *Journal of the ACM*, 36(4), October 1989.

- [33] B. Boguraev, R. Byrd, J. Klavans, and M. Neff. From structural analysis of le semantics in a lexical knowledge base. In Zernik, editor, *Proceedings of the First Lexical Acquisition Workshop*. IJCAI, Detroit, 1989.
- [34] Branimir Boguraev and Ted Briscoe, editors. *Computational Lexicography for Natural Language Processing*. Longman, New York, 1989.
- [35] T. Booth. Probabilistic representation of formal languages. In *Tenth Annual IEEE Symposium on Switching and Automata Theory*, October 1969.
- [36] T.L. Booth and R.A. Thompson. Applying probability measures to abstract languages. *Trans. Comput.*, C-22:442–450, 1973.
- [37] Didier Bourigault. Surface grammatical analysis for the extraction of terminology. In *COLING-92, Vol. III*, pages 977–981, 1992.
- [38] L. Breiman, J. Friedman, R. Olshen, and C. Stone. *Classification and Regression Trees*. Wadsworth and Brooks, 1984.
- [39] Michael Brent. Semantic classification of verbs from their syntactic contexts. ms.
- [40] Michael R. Brent. Automatic acquisition of subcategorization frames from untagged corpora. In *Proceedings of the 29th Annual Meeting of the Association for Computational Linguistics*, pages 209–214, 1991.
- [41] Eric Brill. *Transformation-Based Learning*. PhD thesis, Univ. of Pennsylvania, 1990.
- [42] Eric Brill. Some advances in transformation-based part of speech tagging. In *Proceedings of the 32nd Annual Meeting of the Association for Computational Linguistics*, 94, 1994.
- [43] Ted Briscoe. Noun phrases are regular: a reply to Professor Sampson. In W. Meijs, editor, *Linguistics and Beyond*. Rodopi, 1987.

- [44] Ted Briscoe, Ann Copestake, and Bran Boguraev. Enjoy the paper: Lexical semantics. In *COLING-90, vol. 2*, pages 42–47, 1990.
- [45] P. Brown, V. Della Pietra, S. Della Pietra, and R. Mercer. Class-based n -gram models for natural language. IBM internal research report, IBM, Yorktown Heights, New York 10590, 1990.
- [46] P. Brown, V. Della Pietra, P. deSouza, J. Lai, and R. Mercer. Class-based n -gram models for natural language. *Computational Linguistics*, 18(4):467–480, 1992.
- [47] Jaime G. Carbonell and Philip J. Hayes. Recovery strategies for parsing extragrammatical sentences. *American Journal of Computational Linguistics*, 9(3–4):123–146, 1983.
- [48] Jaime G. Carbonell and Philip J. Hayes. Robust parsing using multiple construction rules. In Leonard Bolc, editor, *Natural Language Parsing Systems*, pages 1–32. Springer-Verlag, Heidelberg, West Germany, 1987.
- [49] R. L. Cave and L. P. Neuwirth. Hidden Markov Models for English. In J. D. Acosta, editor, *Hidden Markov Models for Speech*. IDA-CRD, Princeton, NJ, October 1980.
- [50] Eugene Charniak. (*Statistical NLP*). MIT Press, 1993.
- [51] Y. Chiaramella, B. Defude, M. Bruandet, and D. Kerkouba. Iota: A full text information retrieval system. In *Proc. of ACM ICRDIR*, pages 207–213, 1986.
- [52] M. Chitrao and R. Grishman. Statistical parsing of messages. In *Proceedings of the Conference on Artificial Intelligence and Natural Language Processing*. Morgan Kaufman: New York, 1990.
- [53] M. Chodorow and J. Klavans. Locating syntactic patterns in text corpora. IBM report, IBM, Yorktown Heights, New York 10598, 1990.
- [54] Y. Choueka. Looking for needles in a haystack or locating interesting collocations in large textual databases. In *Proceedings of the RIAO-88, 609-623*. Cambridge, MA, 1988.

- [55] Y. Choueka, S.T. Klein, and E. Neuwitz. Automatic retrieval of frequent idiomatic expressions in a large corpus. *ALLC Journal*, 4(1):34–38, 1983.
- [56] Y. Choueka and S. Luisgnan. Disambiguation by short contexts. *Computers and* 19(3):147–157, 1985.
- [57] Kenneth Church. A stochastic parts program and noun phrase parser for unre. *Proceedings of the Second Conference on Applied Natural Language Processing*, 1988.
- [58] Kenneth Church. Stochastic parts program and noun phrase parser for unre. *ICASSP 89*, pages 695–698, 1989.
- [59] Kenneth Church and William Gale. A comparison of the Enhanced Good-Tu. Estimation methods for estimating probabilities of English Bigrams. *Computational Linguistics*, 5, 1991.
- [60] Kenneth Church, William Gale, Patrick Hanks, and Donald Hindle. Parsing, word typical predicate-argument relations. In *International Workshop on Parsing Te* 389–98, 1989.
- [61] Kenneth Church and Robert Mercer. Introduction to the special issue on comput using large corpora. *Computational Linguistics*, 19(1):1–24, 1993.
- [62] Computational linguistics **19**(1–2): Special issue on using large corpora, 1993.
- [63] W. Cooper and J. Paccia-Cooper. *Syntax and speech*. Harvard University Press, 1980.
- [64] M. Cravero, L. Fissore, R. Pieraccini, and C. Scagliola. Syntax driven recogniti words by markov models. In *ICASSP 84*, 1984. ¶Stochastic Parsing¶.

- [65] Carl G. de Marcken. Parsing the LOB corpus. In *ACL 28*, pages 243–251, 1990.
- [66] S. DeRose. Grammatical category disambiguation by statistical optimization. *Linguistics*, 14(1), 1988.
- [67] A.-M. Deroualt. Context-dependent phonetic Markov models for large vocabulary recognition. *Proc. IEEE ICASSP*, 1:360–363, 1987.
- [68] Jay L. Devore. *Probability and Statistics for Engineering and the Sciences*. Brooks/Cole, Monterey, CA, 2nd edition edition, 1987.
- [69] M. Dillon and A. Gray. Fasit: A fully automatic syntactically based indexing system. *Journal of the ACM*, 34(2):99–108, 1983.
- [70] Eva Ejerhed. Finding clauses in unrestricted text by finitary and stochastic methods. In *Proceedings of the 2nd Conference on Applied Natural Language Processing.*, Austin, Texas, 1983.
- [71] Eva Ejerhed and Kenneth Church. Finite state parsing. In Fred Karlsson, editor, *the Seventh Scandinavian Conference of Linguistics*, pages 410–432, Hallituskatu 10, Helsinki 10, Finland, 1983. University of Helsinki, Department of General Linguistics.
- [72] D. Evans, K. Ginther-Webster, M. Hart, R. Lefferts, and I. Monarch. Automatic selective nlp and first-order thesauri. In *Proc. of RIAO 91 (Barcelona)*, pages 62–71, 1991.
- [73] David Evans. Concept management in text via natural-language processing: The clarit approach. In *Text-Based Intelligent Systems: AAAI Spring Symposium*. AAAI, 1991.
- [74] David A. Evans, Steve K. Henderson, Robert G. Lefferts, and Ira A. Monarch. The clarit project. Technical Report CMU-LCL-91-2, Laboratory for Computational Linguistics, Carnegie-Mellon University, November 1991.
- [75] J. Fagan. *Experiments in Automatic Phrase Indexing of Document Retrieval: Syntactic and Non-Syntactic Methods*. PhD thesis, Cornell University, Ithaca, NY, 1983.

- [76] J. Fagan. The effectiveness of a non-syntactic approach to automatic phrase index retrieval. *JASIS*, 40(2):115–132, 1989.
- [77] Jean Fargues and Adeline Perrin. Synthesizing a large concept hierarchy from free text. In *COLING 90, vol. 2*, pages 112–117, 1990.
- [78] J. D. Ferguson, editor. *Hidden Markov Models for Speech*. IDA-CRD, Princeton, 1980.
- [79] Steven Paul Finch. *Finding Structure in Language*. PhD thesis, University of Edinburgh, 1989.
- [80] K. S. Fu. *Syntactic pattern recognition and applications*. Prentice-Hall, Englewood Cliffs, 1982.
- [81] K.S. Fu. *Syntactic Methods in Pattern Recognition*. Springer-Verlag, New York, 1988.
- [82] K.S. Fu. *Syntactic Methods in Pattern Recognition*. Academic Press, New York, 1988.
- [83] K.S. Fu and T.L. Booth. Grammatical inference: Introduction and survey. *IEEE Transactions on Man and Cybernetics*, 5, 1975. In two parts.
- [84] K.S. Fu and T. Huang. Stochastic grammars and languages. *Int'l. J. of Computer Science*, 1(2):135–170, 1972.
- [85] T. Fujisaki. A stochastic approach to sentence parsing. In *ACL 22*, 1984. Also appeared in *ACL '84*.
- [86] T. Fujisaki, F. Jelinek, J. Cocke, and E. Black. Probabilistic parsing method for natural language. In *Proceedings of the International Workshop on Parsing Technologies*, pages 1–10, 1989.
- [87] Robert P. Futrelle et al. Preprocessing and lexicon design for parsing technical text. In *International Workshop on Parsing Technologies*, pages 31–40, 1991.

- [88] Salton. G. and C. Buckley. A comparison between statistically and syntactically phrases. Report TR89-1027, Cornell University, Dept. of Computer Science, Ithaca, NY, 1989.
- [89] Haim Gaifman. Dependency systems and phrase-structure systems. *Information Processing Letters*, 15:304–337, 1965.
- [90] R. Garside. The CLAWS word-tagging system. In Garside R., F. Leech, and G. Sampson, editors. *The Computational Analysis of English*. Longman, 1987.
- [91] R. Garside, F. Leech, and G. Sampson, editors. *The Computational Analysis of English*. Longman, 1987.
- [92] James Paul Gee and Francois Grosjean. Performance structures: A psycholinguistic appraisal. *Cognitive Psychology*, 15:411–458, 1983.
- [93] Lila Gleitman. The structural sources of verb meanings. *Language Acquisition*, 1(1):1–66, 1984.
- [94] Ralph Grishman, Lynette Hirschman, and Ngo Thanh Nhan. Discovery procedure for selectional patterns: Initial experiments. *Computational Linguistics*, 12(3), 1986.
- [95] Ralph Grishman and John Sterling. Acquisition of selectional patterns. In *Computational Linguistics*, 18:658–664, 1992.
- [96] F. Grosjean, L. Grosjean, and H. Lane. The patterns of silence: Performance structures in speech production. *Cognitive Psychology*, 11:58–81, 1979.
- [97] Michael Halliday. *Intonation and Grammar in British English*. Mouton, The Hague, 1969.
- [98] Zellig Harris. From morpheme to utterance. *Language*, 22, 1946.
- [99] Zellig Harris. From phoneme to morpheme. *Language*, 31, 1955.

- [100] Donald Hindle. Deterministic parsing of syntactic non-fluencies. In *ACL 21 (MIT)*, 1983.
- [101] Donald Hindle. User manual for Fidditch. Technical Memorandum #7590-142 Laboratory, 1983.
- [102] Donald Hindle. Acquiring disambiguation rules from text. In *Proceedings of Meeting of the Association of Computational Linguistics*, Vancouver, British Columbia, 1983.
- [103] Donald Hindle. Noun classification from predicate-argument structures. In *28th Annual Meeting of the Association of Computational Linguistics, Pittsburgh*, 268–275, 1990.
- [104] Donald Hindle. A parser for text corpora. In A. Zampolli, editor, *Computational Linguistics and the Lexicon*. Oxford University Press, New York, 1994.
- [105] Julia Hirschberg. Prosody bibliography. E-mail, 1986.
- [106] Julia Hirschberg. Using text analysis to predict intonational boundaries. Manuscript, MIT Speech Recognition Laboratories, 1991.
- [107] Jerry R. Hobbs, Douglas Appelt, Mabry Tyson, and Megumi Kameyama. Fast extraction of information from text. In *ARPA Workshop on Human Language Technology*, 1993. Defense Advanced Research Projects Agency (DARPA), Morgan Kaufmann, CA, 1993.
- [108] Jerry R. Hobbs et al. SRI International: Description of the FASTUS system used in the *Proceedings, Fourth Message Understanding Conference (MUC-4)*, pages 268–275, 1992. Morgan Kaufmann.
- [109] James Jay Horning. *A Study of Grammatical Inference*. PhD thesis, Stanford (California), 1969.

- [110] Kuang hua Chen and Hsin-Hsi Chen. Extracting noun phrases from large-scale approach and its automatic evaluation. In *Proceedings of ACL*, 1994. Available Archive.
- [111] Institute of Electrical Engineers and Institute of Mathematics, University of Essex. *Inference: Theory, Applications and Alternatives*, Colchester, UK, 1993. IEE Publication no. 1993/092.
- [112] J. of computational linguistics **9**(3-4): Special issue on dealing with ill-formed text.
- [113] Ajay N. Jain. Parsing complex sentences with structured connectionist networks. *Journal of Computational Linguistics*, 3:110-120, 1990.
- [114] Ajay N. Jain. *PARSEC: A Connectionist Learning Architecture for Parsing Spoken Language*. Ph.D. thesis, CMU, Pittsburgh, PA, 1991. Available as Technical Report CMU-CS-91-107.
- [115] F. Jelinek. A fast sequential decoding algorithm using a stack. *IBM Journal of Research and Development*, 13:675-685, November 1969.
- [116] F. Jelinek. Self-organized language modeling for speech recognition. *W & L*, pages 1-10, 1986.
- [117] F. Jelinek. Self-organized language modeling for speech recognition. In *IBM European Conference on Advances in Speech Recognition (Oberlech, Austria)*, 1986.
- [118] F. Jelinek, J. D. Lafferty, and R. L. Mercer. Basic methods of probabilistic context-free parsing. Research Report RC 16374 (#72684), IBM, Yorktown Heights, New York 10598, 1983.
- [119] F. Jelinek and R. Mercer. Interpolated estimation of Markov source parameters from continuous speech. In E.S Gelsema and Kanal L.N., editors, *Pattern Recognition in Practice*, pages 381-397. North Holland Publishing Co., 1980.
- [120] K. Jensen and G.E. Heidorn. The fitted parse: 100English. Computer Science Research Report RC9729 (#42958), IBM Research Division (San Jose), 1982.

- [121] K. Jensen, G.E. Heidorn, L.A. Miller, and Y. Ravin. Parse fitting and prose fixing of ill-formedness. *Computational Linguistics*, 9(3–4):147–161, 1983.
- [122] S. Jones and J. Sinclair. English lexical collocations: A study in Computational Linguistics. *de Lexicologie*, 24:15–49, 1974.
- [123] Aravind K. Joshi and B. Srinivas. Disambiguation of super parts of speech (or super parsing). In *COLING-94*, 1994.
- [124] Gunnel Kællgren. Tagging pilys 47. Technical report, Institute of Linguistics, University of Stockholm, 1982.
- [125] Gunnel Kællgren. Making maximal use of surface criteria in large-scale parsing: unpublished ms., 1990.
- [126] Fred Karlsson. Morphological tagging of Finnish. In *Computational Morphosyntax No. 13*, pages 115–136. University of Helsinki, Department of General Linguistics, 1989.
- [127] Fred Karlsson. Parsing and constraint grammar. unpublished ms., Research Unit for Linguistics, Helsinki, Finland, 1989.
- [128] Judith Klavans. Complex: a computational lexicon for natural language systems. 1988.
- [129] Judith L Klavans. Bibliography on corpus analysis and tagging. presented at the role of large text corpora in building natural language systems at the 13th international conference on computational linguistics (coling), 1990.
- [130] Julian Kupiec. Training stochastic grammars from unlabelled text corpora. Ms., 1989.
- [131] Julian Kupiec. Augmenting a hidden Markov model for phrase-dependent word recognition. *Proceedings of the DARPA Speech and Natural Language Workshop*. Morgan Kaufmann, 1989.

- [132] Julian Kupiec. An algorithm for finding noun phrase correspondences in bilingual text. *ACL*, pages 17–22, 1993.
- [133] Stan Kwasny and N. Sondheimer. Relaxation techniques for parsing ill-formed sentences. *Journal of Computational Linguistics*, 7(2):99–108, 1981.
- [134] Francois-Michel Lang. Parsing incomplete sentences. In *Proceedings of COLING '88*, 1988.
- [135] Francois-Michel Lang and Lynette Hirschman. Improved portability and parsing of natural language. In *Proceedings of the Second Conference on Natural Language Processing*, Austin, TX, 1988. ACL.
- [136] K. Lari and S.J. Young. The estimation of stochastic context-free grammars using the forward algorithm. *Computer Speech and Language*, 4:35–56, 1990.
- [137] K. Lari and S.J. Young. Applications of stochastic context-free grammars using the forward algorithm. *Computer Speech and Language*, 5:237–257, 1991.
- [138] Leech, Garside, and Atwell. The automatic grammatical tagging of the LOB corpus. *Journal of Computational Linguistics*, 7:13–33, 1983.
- [139] Wendy Lehnert et al. University of Massachusetts: MUC-4 test results and analysis. *Fourth Message Understanding Conference (MUC-4)*, pages 151–158, San Francisco, 1988. Morgan Kaufmann.
- [140] L. Lesmo and P. Torasso. Interpreting syntactically ill-formed sentences. In *COLING '88*, 1988.
- [141] W.J.M. Levelt. Hierarchical chunking in sentence processing. *Perception & Psychophysics*, 103, 1970.
- [142] D. D. Lewis and W. B. Croft. Term clustering of syntactic phrases. In *SIGIR '90*, 1990.

- [143] Mark Liberman. Tutorial: Statistical methods in nl processing. In *EACL-93*, 1993.
- [144] Mark Liberman and Mitch Marcus. (statistical nlp). *CACM*, 1994?
- [145] David D. MacDonald. An efficient chart-based algorithm for partial parsing of u
In *Proceedings of the 3rd Conference on Applied Natural Language Processing*, 1993.
- [146] D. Magerman and M. Marcus. Parsing a natural language using mutual informa
Proceedings of AAAI-90, 1990.
- [147] Christopher D. Manning. Automatic acquisition of a large subcategorization dic
pora. In *31st Annual Meeting of the Association for Computational Linguistics*
1993.
- [148] Mitch Marcus. (statistical nlp). Tutorial, ACL 1992, 1992.
- [149] Mitchell Marcus. Building non-normative systems: the search for robustness: an o
20, page 152, 1982.
- [150] Mitchell P. Marcus, Donald Hindle, and Margaret M. Fleck. D-theory: Talking ab
trees. Manuscript, Bell Laboratories.
- [151] James G. Martin. Rhythmic (hierarchical) versus serial structure in speech and
Psychological Review, 79(6):487–509, 1972.
- [152] Chris S. Mellish. Some chart-based techniques for parsing ill-formed input. In *Pr
'89*, 1989.
- [153] Meterer, Schwartz, and Weischedel. Studies in part of speech labelling. In *Procedin
Speech and Natural Language Workshop*. Morgan Kaufmann, 1991.
- [154] Frederick Mosteller and John W. Tukey. *Data Analysis and Regression*. Addison-V
Company, Reading MA, 1977.

- [155] Douglas B. Paul. Speech recognition using Hidden Markov Models. *Lincoln Language* 3(1):41–62, 1990.
- [156] Judea Pearl. *Heuristics: Intelligent Search Strategies for Computer Problem Solving*. Wesley Publishing Company, Reading, MA, 1984.
- [157] Joseph Picone. Continuous speech recognition using Hidden Markov Models. *IEEE Transactions on Speech and Audio Processing*, pages 26–41, July 1990.
- [158] James Pustejovsky, Sabine Bergler, and Peter Anick. Lexical semantic techniques for natural language processing. Ms., Brandeis, 1992.
- [159] R. Quirk, S. Greenbaum, G. Leech, and J. Svartik. *A Comprehensive Grammar of the English Language*. Longman: London, 1985.
- [160] L. R. Rabiner. A tutorial on Hidden Markov Models and selected applications in speech processing. *Proceedings of the IEEE*, 77(2):257–285, February 1989.
- [161] L.R. Rabiner and B.H. Juang. An introduction to Hidden Markov Models. *IEEE Transactions on Speech and Audio Processing*, page 4ff, January 1986.
- [162] Lance A. Ramshaw. Exploring the statistical derivation of transformational rules for part-of-speech tagging. In *Proceedings of the ACL Balancing Act Workshop*, 1990.
- [163] Mats Rooth. Unitary stochastic part-of-speech and phrase tagging. Manuscript, University of Stuttgart, 1994.
- [164] Ian C. Ross and John W. Tukey. Introduction to these volumes. In *Index of Statistical Probability*, pages iv–x. The R & D Press, Los Altos, CA, 1975.
- [165] G. Ruge, C. Schwarz, and A. Warner. Effectiveness and efficiency in natural language processing for large amounts of text. *JASIS*, 42(6):450–456, 1991.

- [166] Gerard Salton and Maria Smith. On the application of syntactic methodologies in text analysis. In *Proceedings of the 12th Annual International ACM/SIGIR Conference and Development in Information Retrieval*, pages 137–150, 1989.
- [167] E. Sanchis, F. Casacuberta, I. Galiano, and E. Segarra. Learning structural models from text units through grammatical inference. In *IEEE ICASSP, Vol. 1*, pages 189–192, 1990.
- [168] Beatrice Santorini. Part-of-speech tagging guidelines for the Penn Treebank Project. Report MS-CIS-90-47/LINC LAB 178, Dept. of Computer and Information Science, University of Pennsylvania, Philadelphia, PA, July 1990.
- [169] Christoph Schwarz. Automatic syntactic analysis of free text. *JASIS*, 41(6):408–414, 1990.
- [170] Elisabeth O. Selkirk. On prosodic structure and its relation to syntactic structure. In J. Ohala, editor, *Nordic Prosody II*. Tapir, Trondheim, 1978.
- [171] Elisabeth O. Selkirk. Prosodic domains in phonology: Sanskrit revisited. In M. A. Jones and J. Kean, editors, *Juncture*, pages 107–129. Anma Libri, Saratoga, CA, 1980.
- [172] Elisabeth O. Selkirk. On the nature of phonological representations. In T. Myer and J. Anderson, editors, *The Cognitive Representation of Speech*. North-Holland Publishing Company, Amsterdam, 1981.
- [173] Elisabeth O. Selkirk. *Phonology and Syntax: The Relation between Sound and Meaning*. MIT Press, Cambridge, MA, 1984.
- [174] Stephanie Seneff. A relaxation method for understanding spontaneous speech utterances. In *Proceedings, Speech and Natural Language Workshop*, San Mateo, CA, 1992. DARPA, Addison-Wesley, Harman Publishers.
- [175] *Extraction of Hierarchical Structure for Machine Learning of Natural Language*. IBM Research and Language Technology (ITK), 1992. Proceedings 92/1, ISBN 90-74029-02-7.

- [176] Frank Smadja. *Extracting Collocations from Text. An Application: Language* thesis, Columbia University, New York, NY, 1991.
- [177] Frank Smadja and Kathy McKeown. Automatically extracting and representing language generation. In *Proceedings of the 28th Annual Meeting of the Association of Linguistics*, pages 252–259, 1990.
- [178] A. Smeaton. Using parsing of natural language as part of document retrieval. CSC/88/R1, University of Glasgow, 1988.
- [179] A.R. Smith et al. Application of a sequential pattern learning system to connection. In *ICASSP '85*, 1985.
- [180] Tony C. Smith and Ian H. Witten. Language inference from function words. Manuscript of Calgary and University of Waikato, January 1993.
- [181] P. Smyth and R.M. Goodman. An information theoretic approach to rule induction. *IEEE Transactions on Knowledge and Data Engineering*, 4(4):301–317, August 1992.
- [182] David Stallard and Robert Bobrow. Fragment processing in the DELPHI system. *Speech and Natural Language Workshop*, pages 305–310, San Mateo, CA, 1992. Research Projects Agency (DARPA), Morgan Kaufmann Publishers.
- [183] Andreas Stolcke and Stephen Omohundro. Inducing probabilistic grammars by merging. From cmp-lg archive, 1994.
- [184] Andreas Stolcke and Stephen M. Omohundro. Best-first model merging for hidden Markov model induction. Technical report, International Computer Science Institute, Berkeley, January 1994.
- [185] W. Stolz. A probabilistic procedure for grouping words into phrases. *Language* 42:235, 1965.

- [186] Tomek Strzalkowski. Ttp: A fast and robust parser for natural language. In *C* 198–204, 1992.
- [187] Stan Szpakowicz. Semi-automatic acquisition of conceptual structure from technical texts. *Man-Machine Studies*, 33:385–397, 1990.
- [188] E. Tzoukermann and B. Merialdo. Some statistical approaches for tagging unpublished ms., IBM, T. J. Watson Research Center, Yorktown Heights, New York, 1990.
- [189] L. G. Valiant. A theory of the learnable. In *Proceedings of the ACM Symposium on Computing*, pages 436–445, Washington, D.C., 1984. ACM Press.
- [190] Jacques Vergnes. ?? In *COLING '90*, 1990.
- [191] Atro Voutilainen. NPtool, a detector of English noun phrases. In *Proceedings of Very Large Corpora*, pages 48–57, 1993.
- [192] Atro Voutilainen, Juha Heikkilä, and Arto Anttila. Constraint grammar of english oriented introduction. Technical Report Publication No. 21, University of Helsinki General Linguistics, Helsinki, 1992.
- [193] Michelle Q. Wang and Julia Hirschberg. Predicting intonational phrasing from text. *J. Acoust. Soc. Am.*, 87:1037–1047, 1990.
- [194] Weischedel and Black. Responding intelligently to unparsable inputs. *Amer. J. of CL*, 6(2):97–109, 1980.
- [195] Ralph Weischedel et al. Partial parsing: A report on work in progress. In *DARPA Speech and Natural Language Workshop*, pages 204–209, Asilomar, CA, 1980.
- [196] R.M. Weischedel and N.K. Sondheimer. Meta-rules as a basis for processing ill-formed sentences. *Amer. J. of CL*, 9:161–177, 1983.

- [197] Yorick Wilks, Louise Guthrie, Joe Guthrie, and Jim Cowie. Combining weak n-gram and rule-based methods for large scale text processing. In Paul S. Jacobs, editor, *Text-Based Intelligent Systems: Theory and Practice in Information Extraction and Retrieval*, pages 35–58. Lawrence Erlbaum Associates, Hillsdale, New Jersey, 1992.
- [198] David Yarowsky. Decision lists for lexical ambiguity resolution. Manuscript, University of Pennsylvania, 1994.
- [199] Uri Zernik and Paul Jacobs. Tagging for learning: Collecting thematic relations. In *COLING '90 vol. 1*, pages 34–39, 1990.
- [200] Wu Zhibiao. A survey of statistical-based approaches to nlp. ms., 1993.