

Shallow Processing Light Parsing & Named Entity Extraction

Günter Neumann
LT lab, DFKI

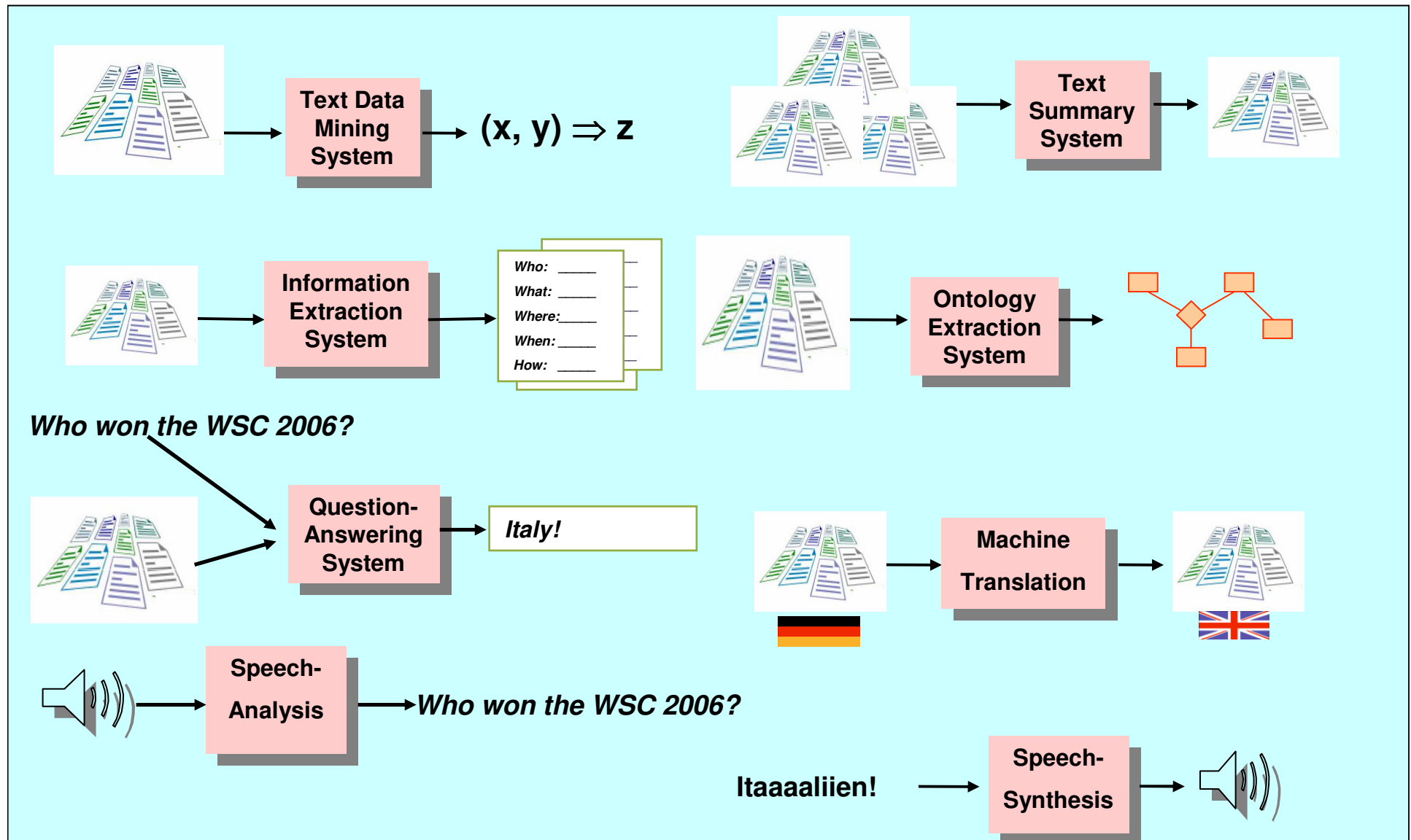
(includes modified slides from
Steven Bird, Gerd Dalemans, Karin Haenelt)

Human Language Technology

- *Human Language Technology LT* – covers
 - The design and implementation of algorithms, data and electronic devices for processing of natural language (text and speech), and
 - Their integration into real-world applications and products

- Language Technology defines the engineering part of computational linguistics

LT-methods cover many areas



Multi/cross-linguality is of great importance in all these areas!

LT Components
Lexical / Morphological Analysis

Tagging
Chunking

Syntactic Analysis

Grammatical Relation Finding
Named Entity Recognition
Word Sense Disambiguation

Semantic Analysis

Reference Resolution
Discourse Analysis

Text



Meaning

Applications

OCR
Spelling Error Correction
Grammar Checking
Information retrieval

Document Classification
Information Extraction
Summarization
Question Answering
Ontology Extraction and Refinement

Dialogue Systems
Machine Translation

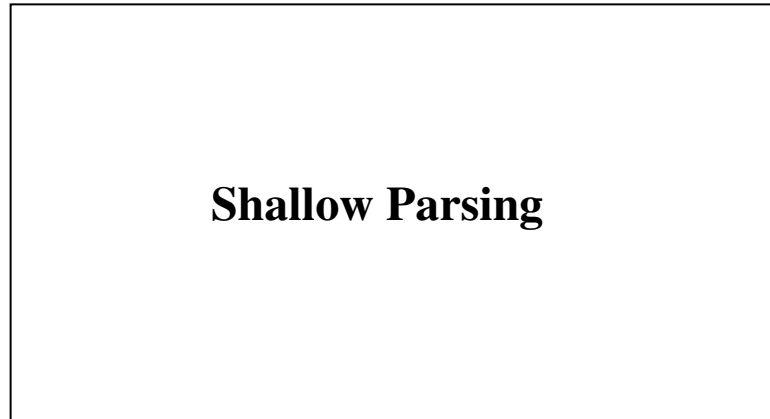
LT Components

Text

Applications

Lexical / Morphological Analysis

OCR



Spelling Error Correction

Grammar Checking

Information retrieval

Named Entity Recognition

Document Classification

Information Extraction

Word Sense Disambiguation

Summarization

Question Answering

Semantic Analysis

Ontology Extraction and Refinement

Reference Resolution

Discourse Analysis

Dialogue Systems

Machine Translation

Meaning

From POS tagging to IE - Classification-Based Perspective

- POS tagging

The/Det woman/NN will/MD give/VB Mary/NNP a/Det book/NN

- NP chunking

The/B-NP woman/I-NP will/B-VP give/I-VP Mary/B-NP a/B-NP book/I-NP

- Grammatical Relation Finding

[NP-SUBJ-1 the woman] [VP-1 will give] [NP-I-OBJ-1 Mary] [NP-OBJ-1 a book]]

- Semantic Tagging (as for Information Extraction)

[Giver the woman][will give][Givee Mary][Given a book]

- Semantic Tagging (as for Question Answering)

Who will give Mary a book?

[Giver ?][will give][Givee Mary][Given a book]

Parsing of unrestricted text

- Complexity of parsing of unrestricted text
 - Large sentences
 - Large data sources
 - Input texts are not simply sequences of word forms
 - Textual structure (e.g., enumeration, spacing, etc.)
 - Combined with structural annotation (e.g., XML tags)
 - Various text styles, e.g., newspaper text, scientific texts, blogs, email, ...
 - Demands high degree of flexibility and robustness

Motivations for Parsing

- Why parse sentences in the first place?
- Parsing is usually an intermediate stage
 - To uncover structures that are used by later stages of processing
- Full Parsing is a *sufficient* but not a *necessary* intermediate stage for many NLP tasks.
- Parsing often provides more information than we need.

Shallow Parsing Approaches

- Light (or “partial”) parsing
- Chunk parsing (a type of light parsing)
 - Introduction
 - Advantages
 - Implementations
- Divide-and-conquer parsing for German

Light Parsing

Goal: assign a *partial structure* to a sentence.

- Simpler solution space
- Local context
- Non-recursive
- Restricted (local) domain

Output from Light Parsing

- What kind of *partial structures* should light parsing construct?
- Different structures useful for different tasks:
 - Partial constituent structure
[NP I] [VP saw [NP a tall man in the park]].
 - Prosodic segments
[I saw] [a tall man] [in the park].
 - Content word groups
[I] [saw] [a tall man] [in the park].

Chunk Parsing

Goal: divide a sentence into a sequence of chunks.

- Chunks are non-overlapping regions of a text

[I] saw [a tall man] in [the park]

- Chunks are non-recursive
 - A chunk can not contain other chunks
- Chunks are non-exhaustive
 - Not all words are included in the chunks

Chunk Parsing Examples

- **Noun-phrase chunking:**
 - [I] saw [a tall man] in [the park].
- **Verb-phrase chunking:**
 - The man who [was in the park] [saw me].
- **Prosodic chunking:**
 - [I saw] [a tall man] [in the park].

Chunks and Constituency

Constituents: *[[a tall man] [in [the park]]]*.

Chunks: *[a tall man] in [the park]*.

- A constituent is part of some higher unit in the hierarchical syntactic parse
- Chunks are *not constituents*
 - Constituents are recursive
- But, chunks are typically sub-sequences of constituents
 - Chunks do not cross major constituent boundaries

1. [_{NP} [_{NP} G.K. Chesterton], [_{NP} [_{NP} author] of [_{NP} [_{NP} The Man] who was [_{NP} Thursday]]]]

2. [_{NP} G.K. Chesterton], [_{NP} author] of [_{NP} The Man] who was [_{NP} Thursday]

Chunk Parsing: Accuracy

Chunk parsing achieves high accuracy

- Small solution space
- Less word-order flexibility *within* chunks than *between* chunks
 - Fewer long-range dependencies
 - Less context dependence
- Better locality
- No need to resolve ambiguity
- Less error propagation

Chunk Parsing: Domain Specificity

Chunk parsing is less domain specific

- Dependencies on lexical/semantic information tend to occur at levels “higher” than chunks:
 - Attachment
 - Argument selection
 - Movement
- Fewer stylistic differences with chunks

Psycholinguistic Motivations

- **Chunks are processing units**
 - Humans tend to read texts one chunk at a time
 - Eye movement tracking studies
- **Chunks are phonologically marked**
 - Pauses
 - Stress patterns
- **Chunking might be a first step in full parsing**
 - Integration of shallow and deep parsing
 - Text zooming

Chunk Parsing: Efficiency

- Smaller solution space
- Relevant context is small and local
- Chunks are non-recursive
- Chunk parsing can be implemented with a finite state machine
 - Fast (linear)
 - Low memory requirement (no stacks)
- Chunk parsing can be applied to a very large text sources (e.g., the web)

Chunk Parsing Techniques

- Chunk parsers usually ignore lexical content
- Only need to look at part-of-speech tags
- Techniques for implementing chunk parsing
 - Regular expression matching
 - Chinking
 - Cascaded Finite state transducers

Regular Expression Matching

- Define a regular expression that matches the sequences of tags in a chunk
 - A simple noun phrase chunk regexp:
 - `<DT> ? <JJ> * <NN.??>`
- Chunk all matching subsequences:
 - In:
The /DT little /JJ cat /NN sat /VBD on /IN the /DT mat /NN
 - Out:
[The /DT little /JJ cat /NN] sat /VBD on /IN [the /DT mat /NN]
- If matching subsequences overlap, the first one gets priority
- Regular expressions can be cascaded

Chinking

- A *chink* is a subsequence of the text that is not a chunk.
- Define a regular expression that matches the sequences of tags in a chink.

– A simple chink regexp for finding NP chunks:

`(<VB.??> | <IN>)+`

- Chunk anything that is *not* a matching subsequence:

the/DT little/JJ cat/NN sat/VBD on /IN the /DT mat/NN

[the/DT little/JJ cat/NN] sat/VBD on /IN [the /DT mat/NN]

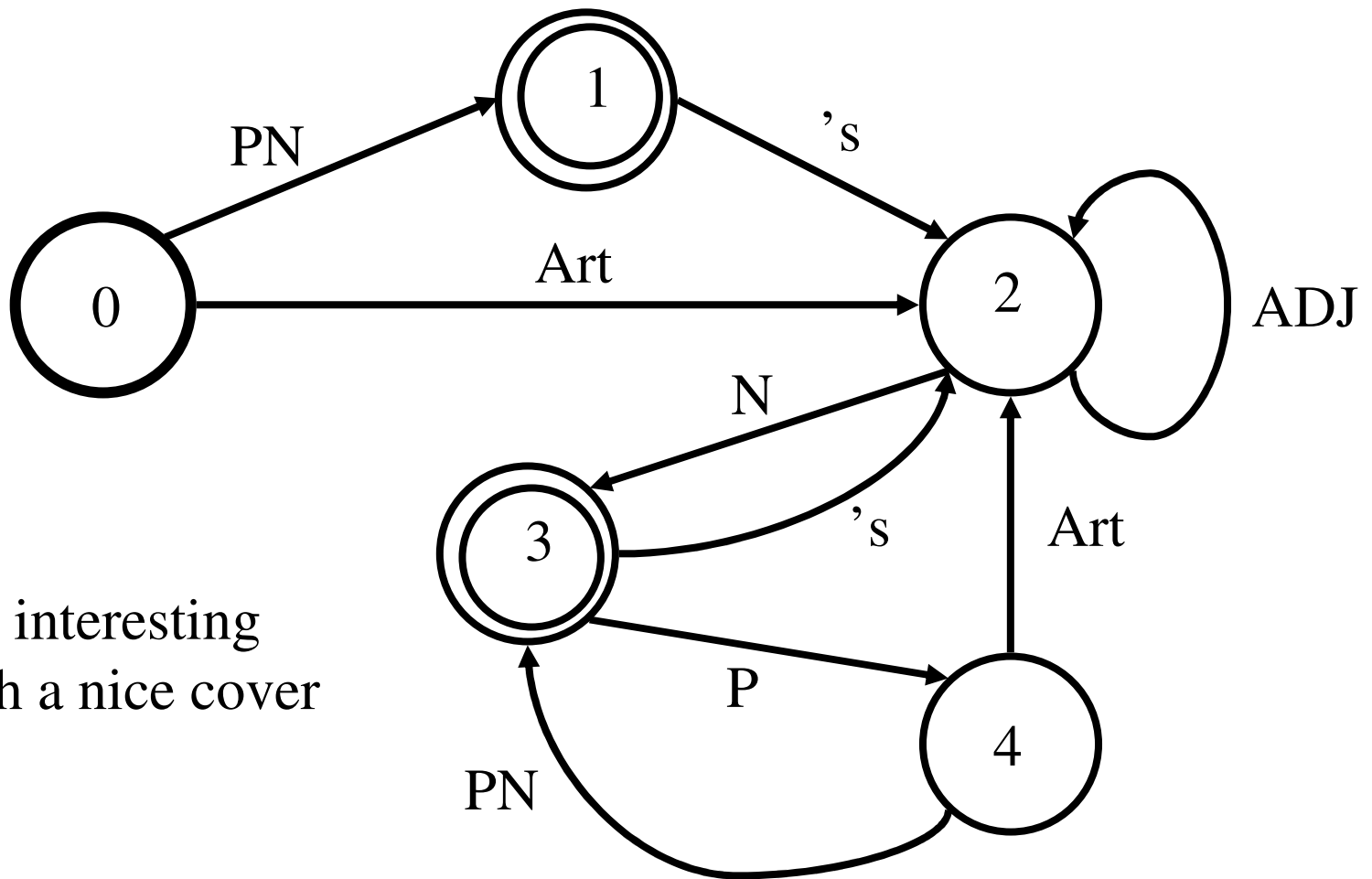
chunk

chink

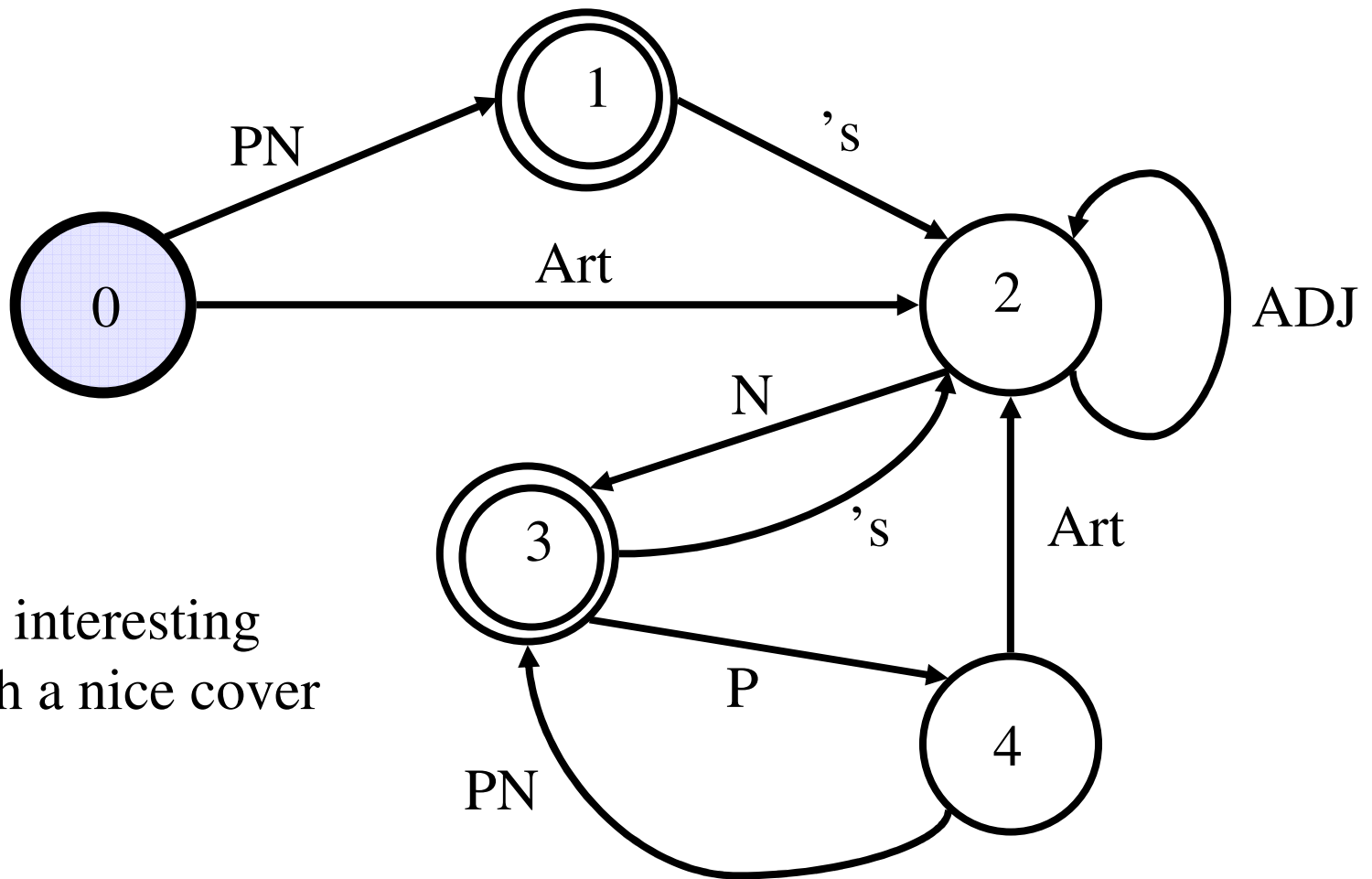
chunk

Finite State Approaches to Shallow Parsing

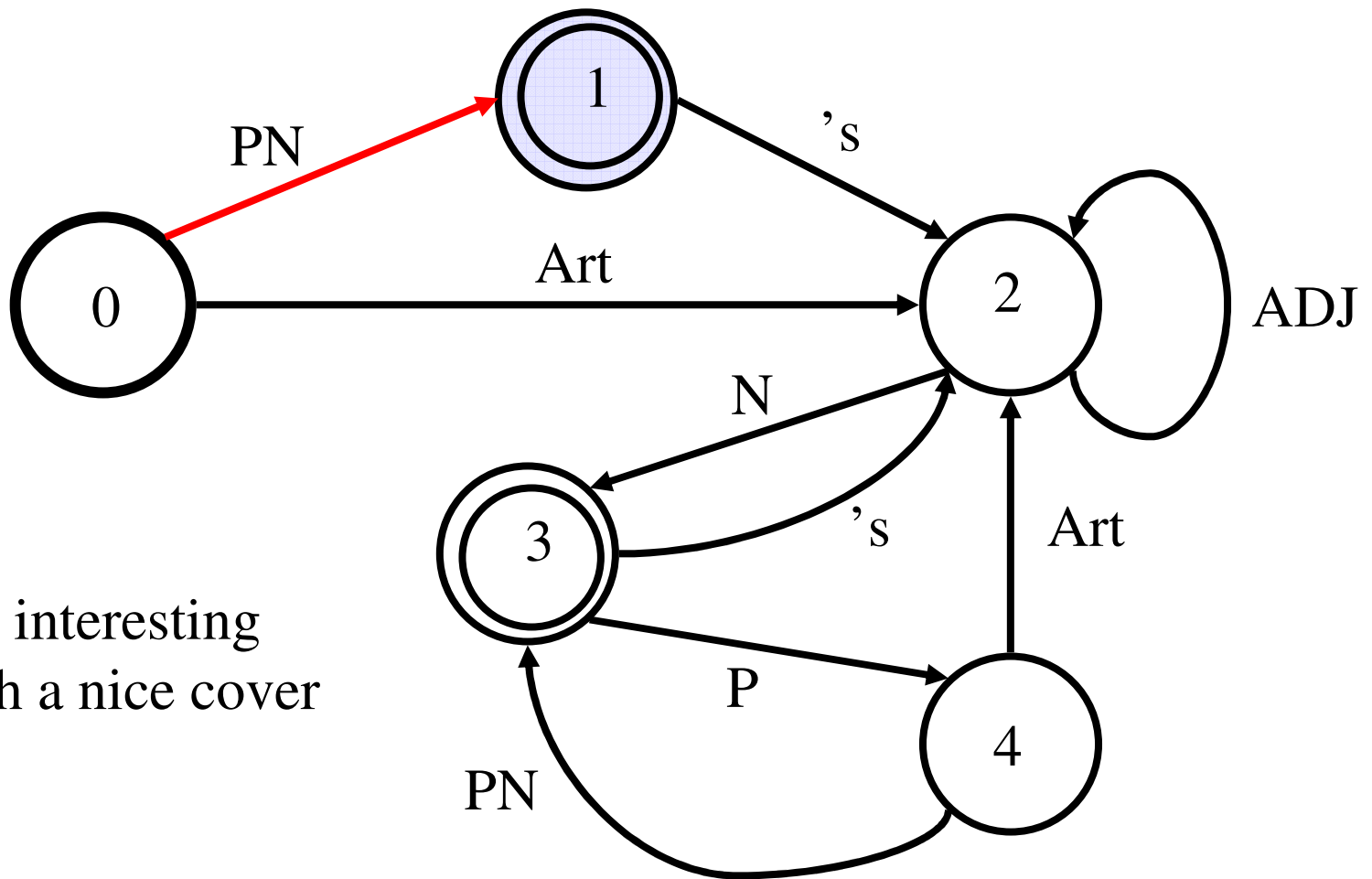
- Finite-state approximation of sentence structures (Abney 1995)
 - finite-state cascades: sequences of levels of regular expressions
 - recognition approximation: tail-recursion replaced by iteration
 - interpretation approximation: embedding replaced by fixed levels
- Finite-state approximation of phrase structure grammars (Pereira/Wright 1997)
 - flattening of shift-reduce-recogniser
 - no interpretation structure (acceptor only)
 - used in speech recognition where syntactic parsing serves to rank hypotheses for acoustic sequences
- Finite-state approximation (Sproat 2002)
 - bounding of centre embedding
 - reduction of recognition capacity
 - flattening of interpretation structure



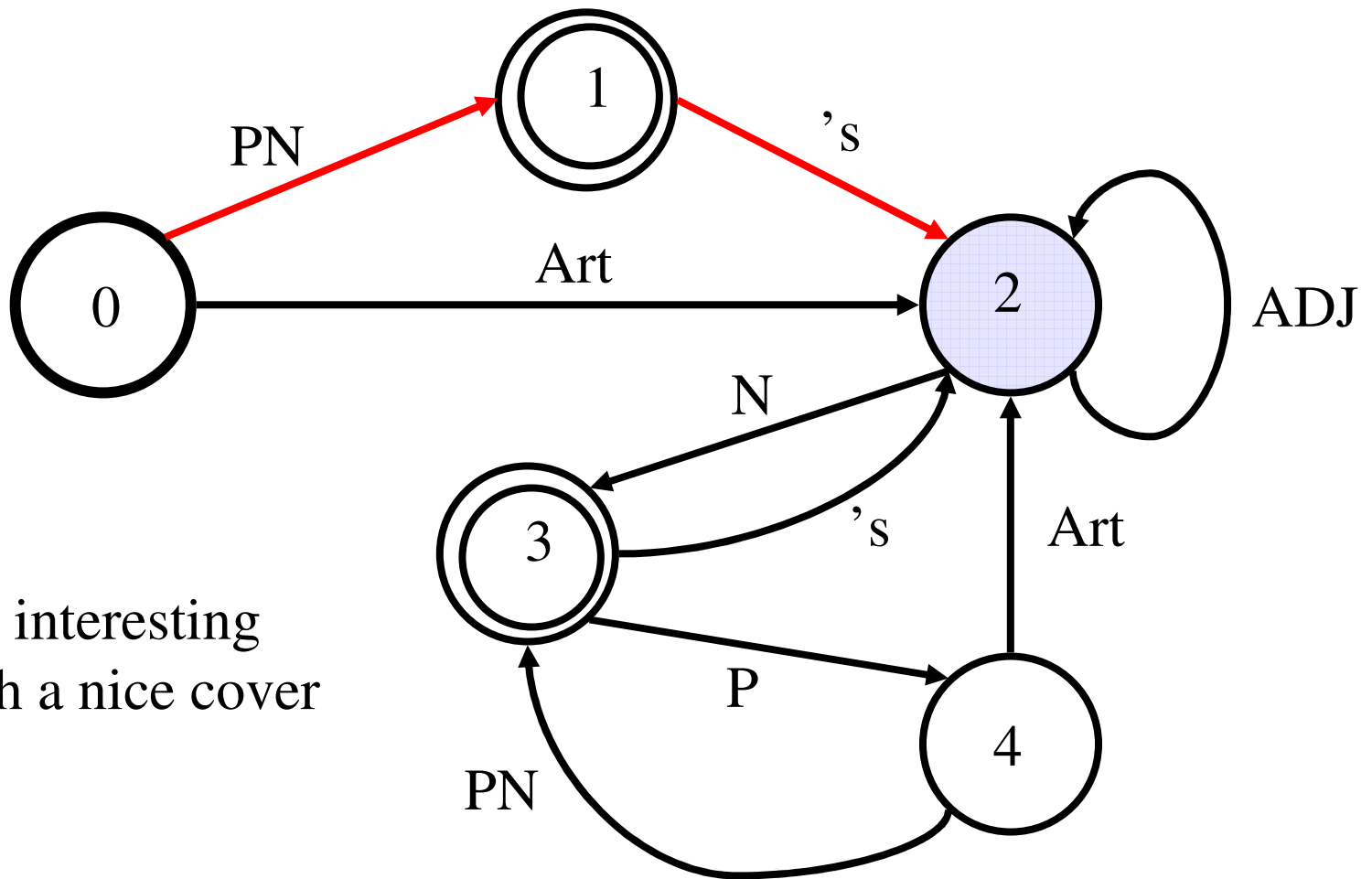
John's interesting
book with a nice cover



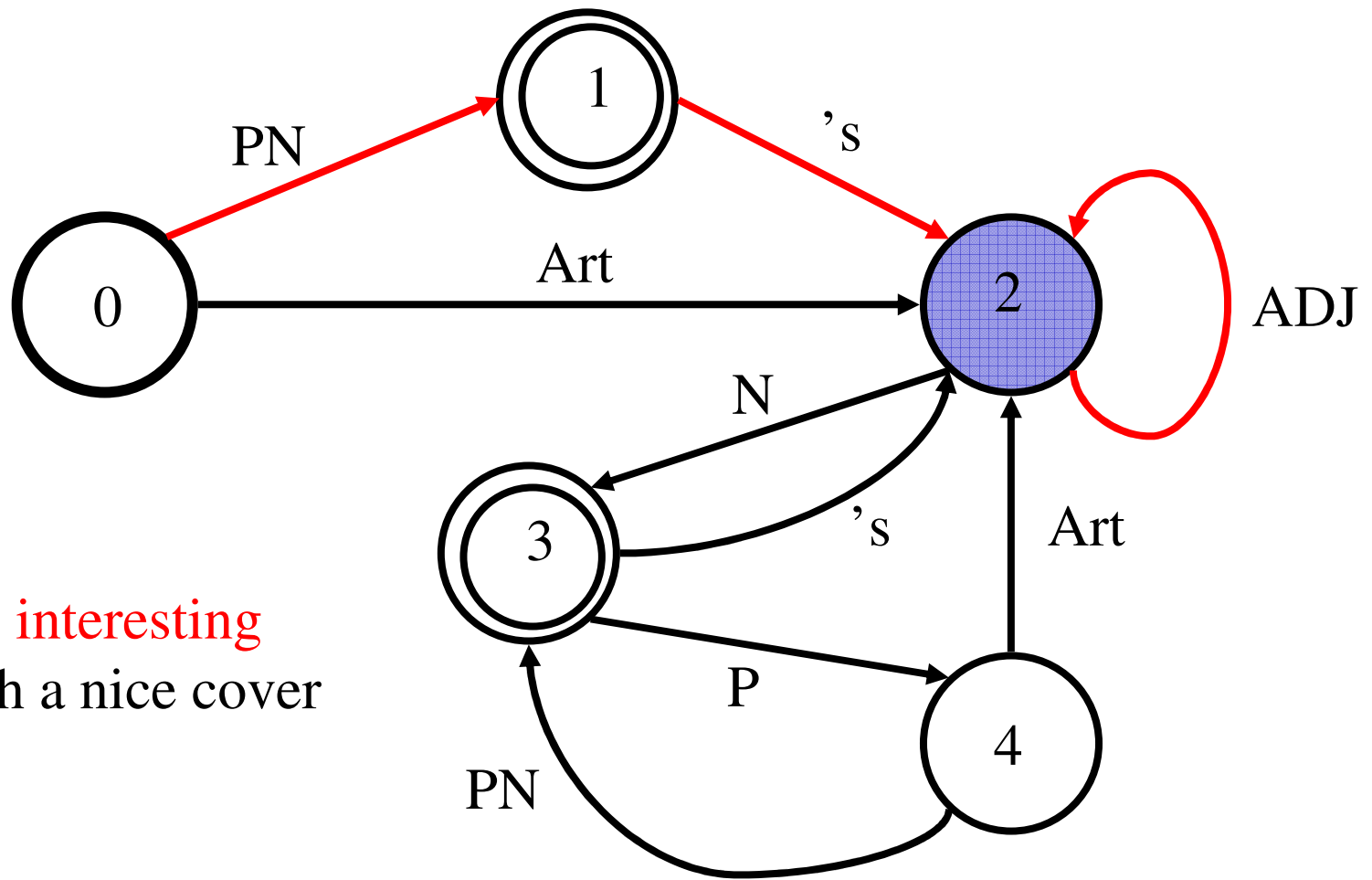
John's interesting
book with a nice cover



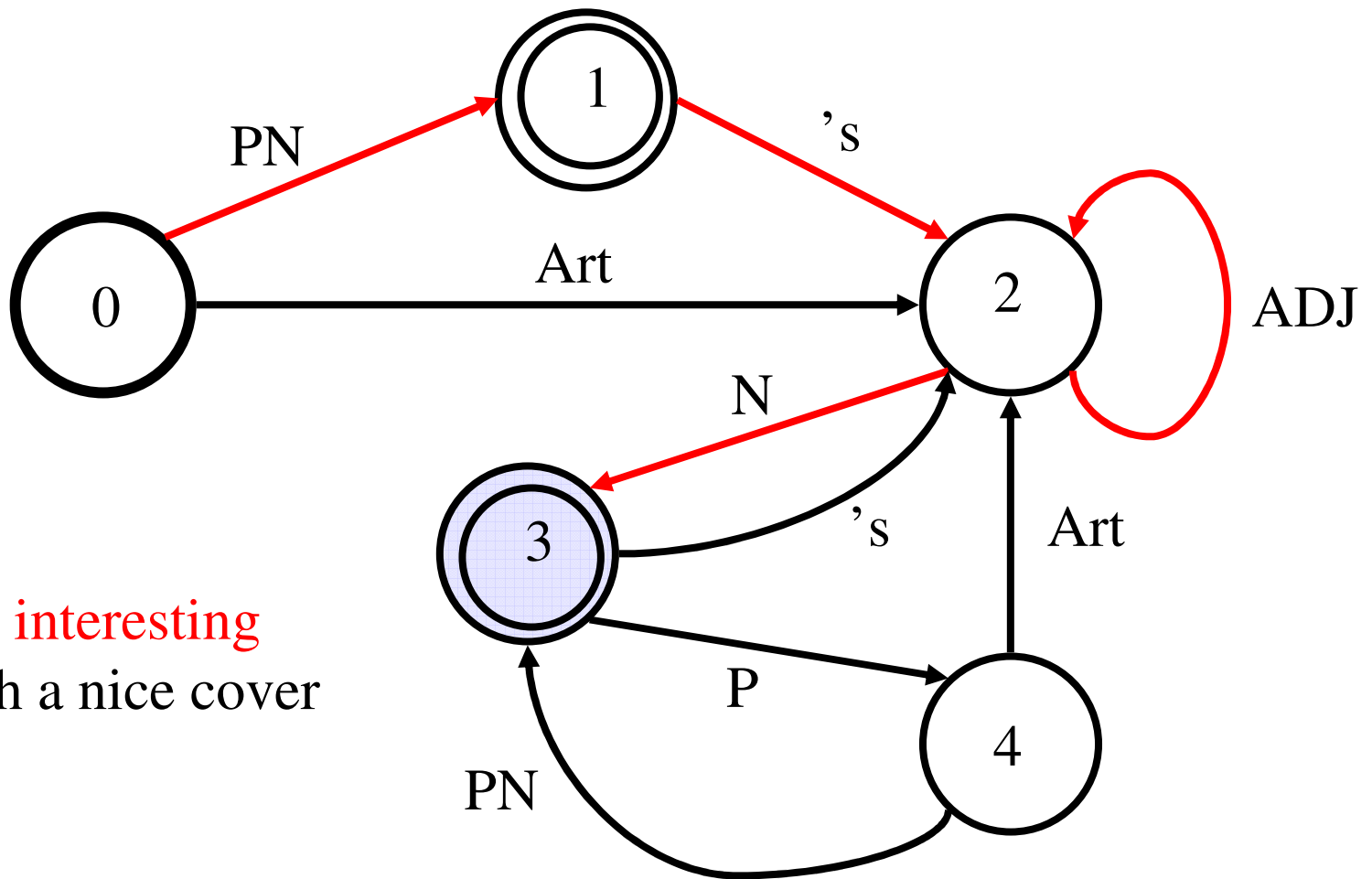
John's interesting
book with a nice cover



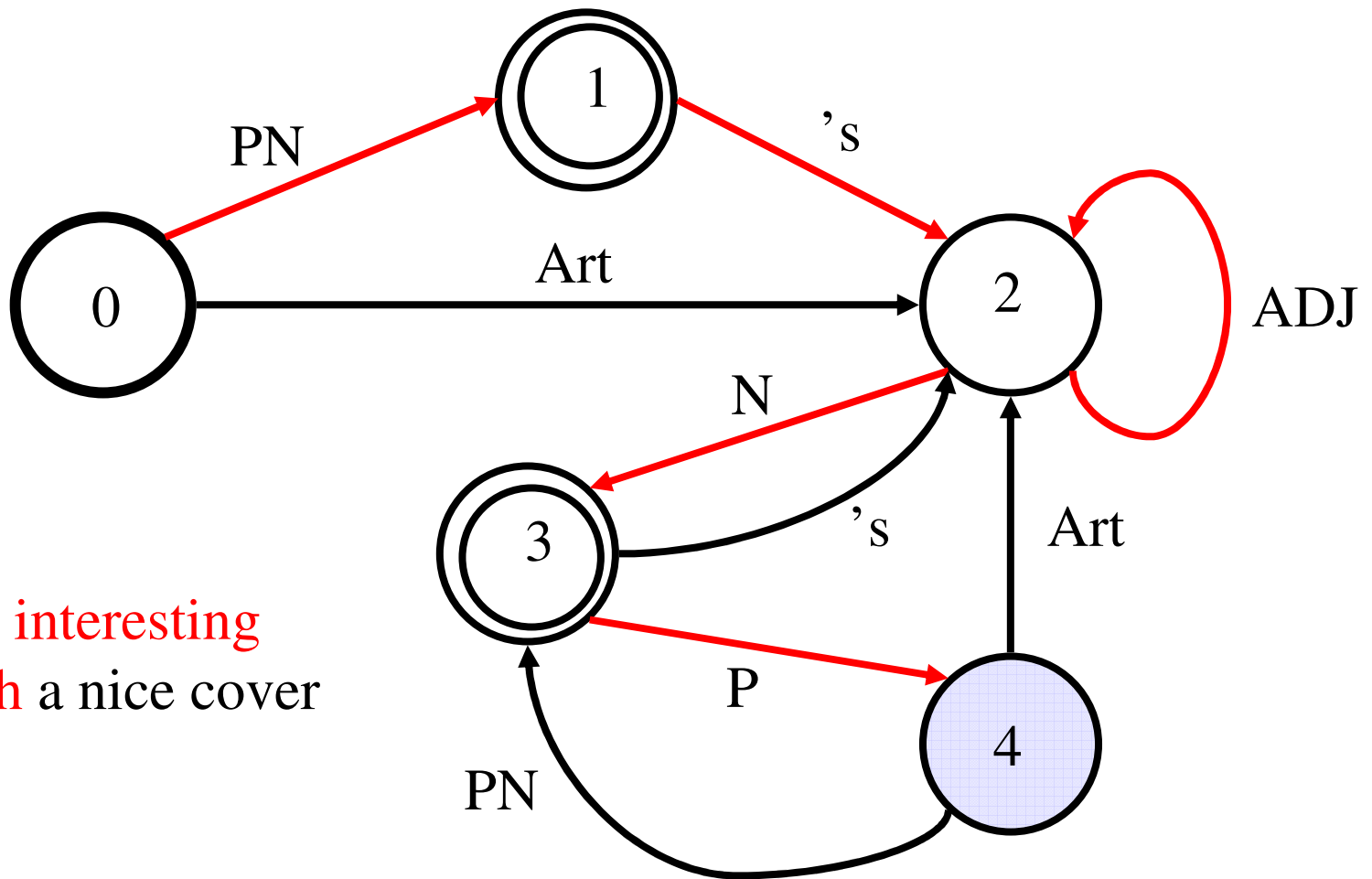
John's interesting
book with a nice cover



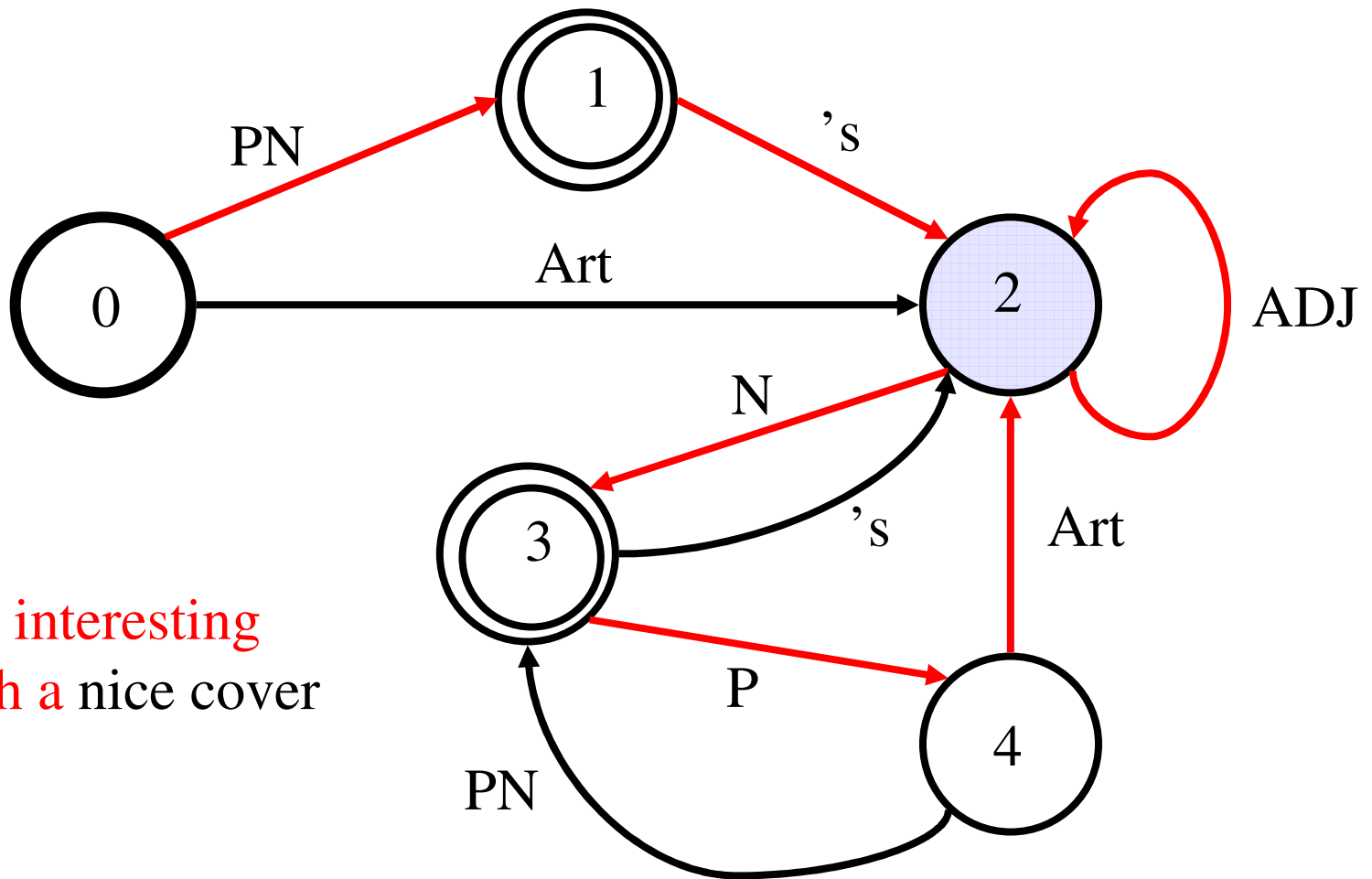
John's interesting
book with a nice cover



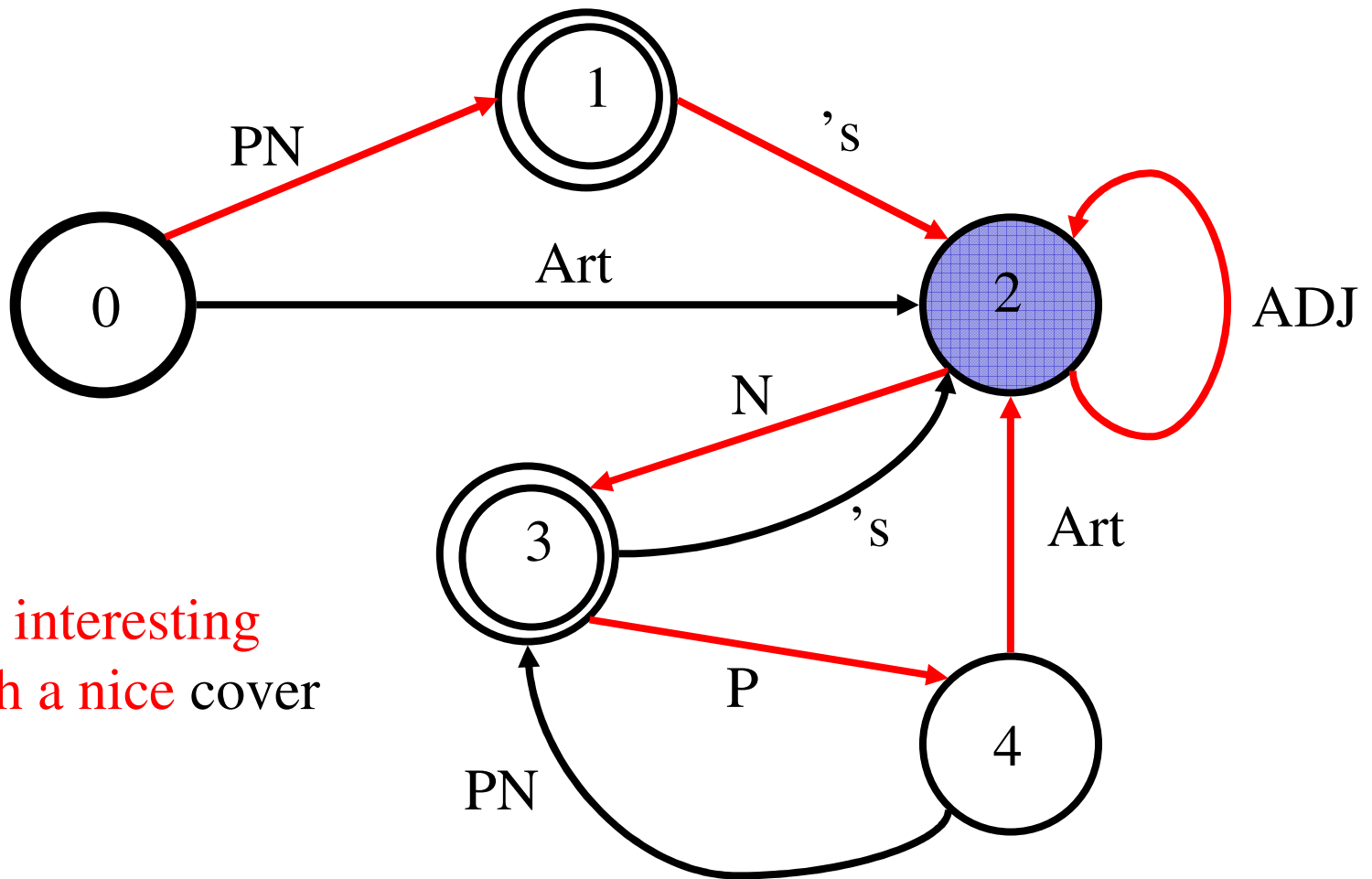
John's interesting
 book with a nice cover



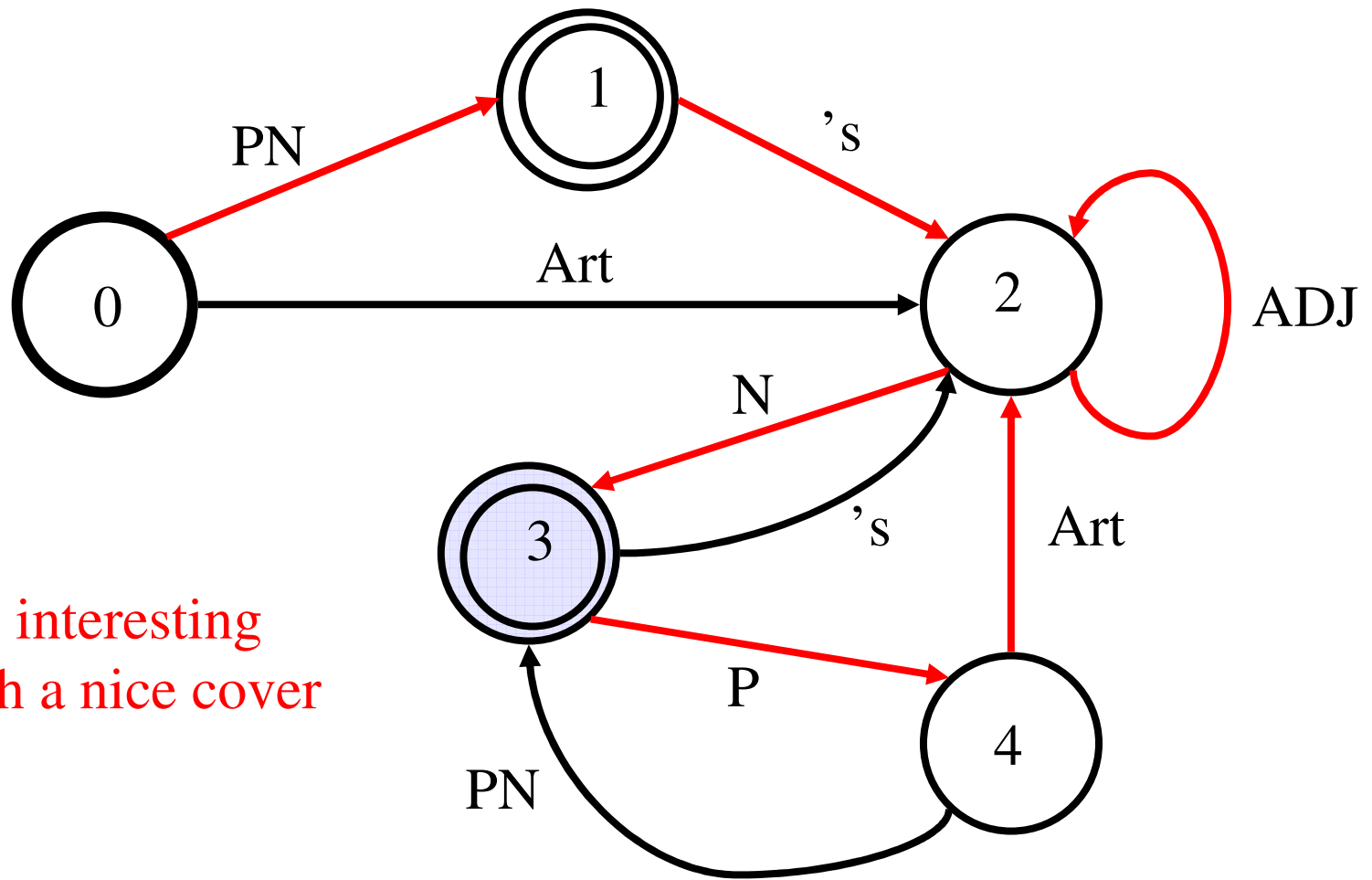
John's interesting
book with a nice cover



John's interesting
book with a nice cover



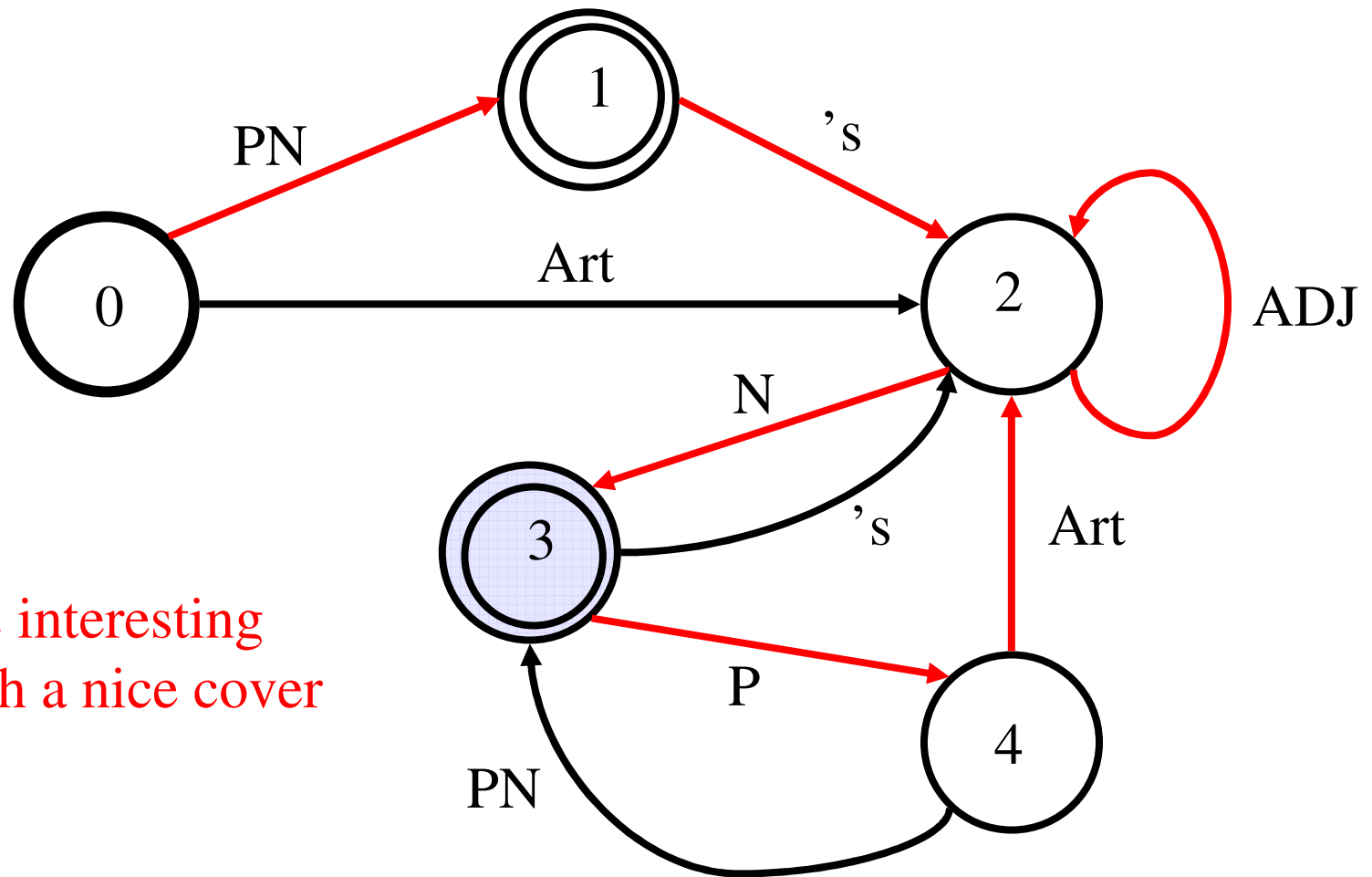
John's interesting
book with a nice cover



John's interesting
book with a nice cover

Pattern-matching

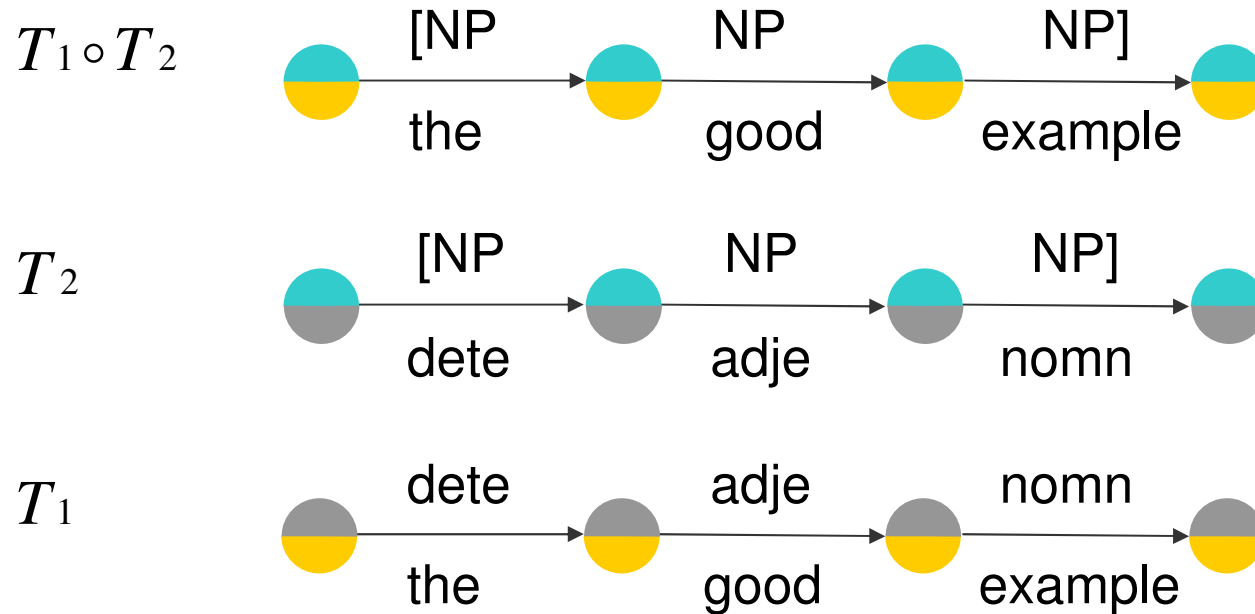
PN 's (ADJ)* N P Art (ADJ)* N



John's interesting
book with a nice cover

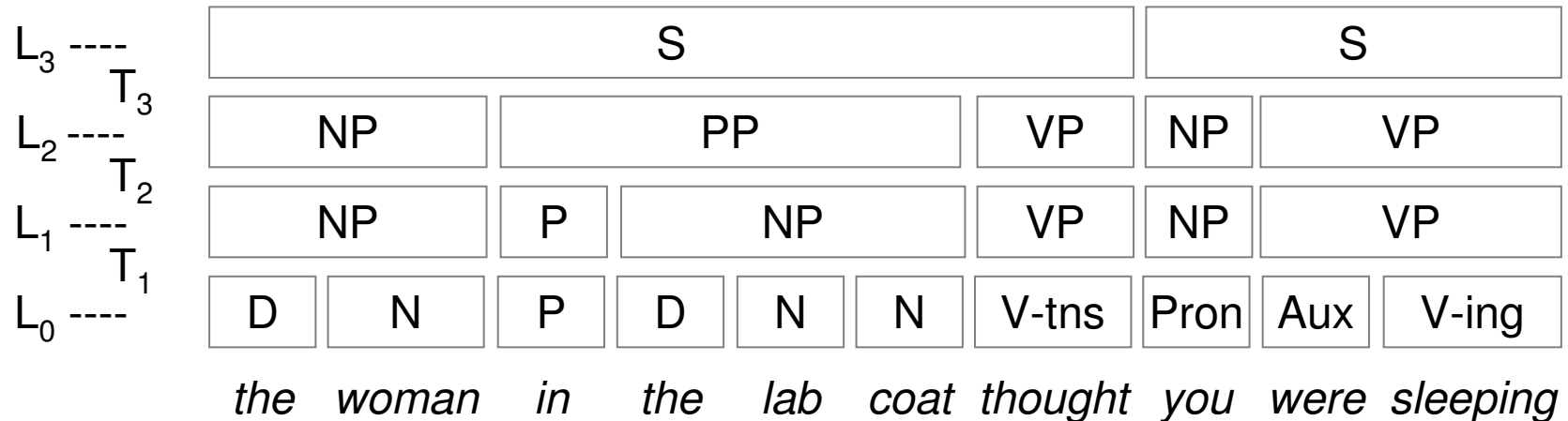
Syntactic Structure: Finite State Cascades

- functionally equivalent to composition of transducers,
 - but without intermediate structure output
 - the individual transducers are considerably smaller than a composed transducer



Syntactic Structure: Finite-State Cascades (Abney)

Finite-State Cascade



Regular-Expression Grammar

$$L_1: \left\{ \begin{array}{l} NP \rightarrow D? N^* N \\ VP \rightarrow V - tns \mid Aux V - ing \end{array} \right\}$$

$$L_2: \{ PP \rightarrow P NP \}$$

$$L_3: \{ S PP^* NP PP^* VP PP^* \}$$

NOTE:
No recursion allowed

Syntactic Structure:

Finite-State Cascades (Abney)

- cascade consists of a sequence of levels
- phrases at one level are built on phrases at the previous level
- no recursion:
 - phrases never contain same level or higher level phrases
- two levels of special importance
 - chunks: non-recursive cores (NX, VX) of major phrases (NP, VP)
 - simplex clauses: embedded clauses as siblings
- patterns:
 - reliable indicators of bits of syntactic structure

Syntactic Structure:

Finite-State Cascades (Abney)

- each transduction is defined by a set of patterns
 - category
 - regular expression
- regular expression is translated into a finite-state automaton
- level transducer
 - union of pattern automata
 - deterministic recognizer
 - each final state is associated with a unique pattern
- heuristics
 - longest match (resolution of ambiguities)
- external control process
 - if the recognizer blocks without reaching a final state, a single input element is punted to the output and recognition resumes at the following word

Syntactic Structure:

Finite-State Cascades (Abney)

- patterns: reliable indicators of bits of syntactic structure
- parsing
 - easy-first parsing (easy calls first)
 - proceeds by growing islands of certainty into larger and larger phrases
 - no systematic parse tree from bottom to top
 - recognition of recognizable structures
 - containment of ambiguity
 - prepositional phrases and the like are left unattached
 - noun-noun modifications not resolved

Syntactic Structure: Finite-State Cascades (Abney)

- **extended patterns**
 - include actions
 - after a phrase with pattern p has been recognised an internal transducer for pattern p is used
 - to flesh out the phrase with features and internal structure
 - insert brackets (non-deterministic, not left-to-right)
- **features represented as bit vectors**
 - unification by bit operations
 - phrases are not rejected in case of unification failures

An alternative FST cascade for German (free word order), Neumann et al.

Most partial parsing approaches following a bottom-up strategy:

Major steps

lexical processing

including morphological analysis, POS-tagging, Named Entity recognition

phrase recognition

general nominal & prepositional phrases, verb groups

clause recognition via domain-specific templates

templates triggered by domain-specific predicates attached to relevant verbs;
expressing domain-specific selectional restrictions for possible argument fillers

Bottom-up chunk parsing

perform clause recognition after phrase recognition is completed

However a bottom-up strategy showed to be problematic in case of German free text processing

Crucial properties of German

1. highly ambiguous morphology (e.g., case for nouns, tense for verbs)
2. free word/phrase order
3. splitting of verb groups into separated parts into which arbitrary phrases and clauses can be spliced in (e.g., *Der Termin **findet** morgen **statt**. The date **takes place** tomorrow.*)

Main problem in case of a bottom-up parsing approach:
Even recognition of simple sentence structure depends heavily on performance of phrase recognition.

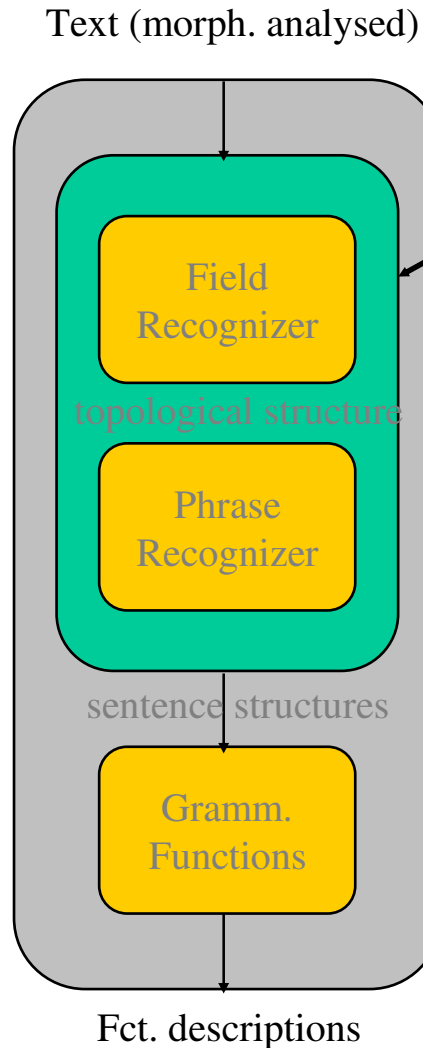
NP ist gängige Praxis.

[NP Die vom Bundesgerichtshof und den Wettbewerbern als Verstoß gegen das Kartellverbot gezeigte zentrale TV-Vermarktung] ist gängige Praxis.

NP ist gängige Praxis.

[NP Central television marketing censured by the German Federal High Court and the guards against unfair competition as an infringement of anti-cartel legislation] is common practice.

In order to overcome these problems we propose the following two phase divide-and-conquer strategy



Divide-and-conquer strategy

1. Recognize verb groups and topological structure (*fields*) of sentence domain-independently;

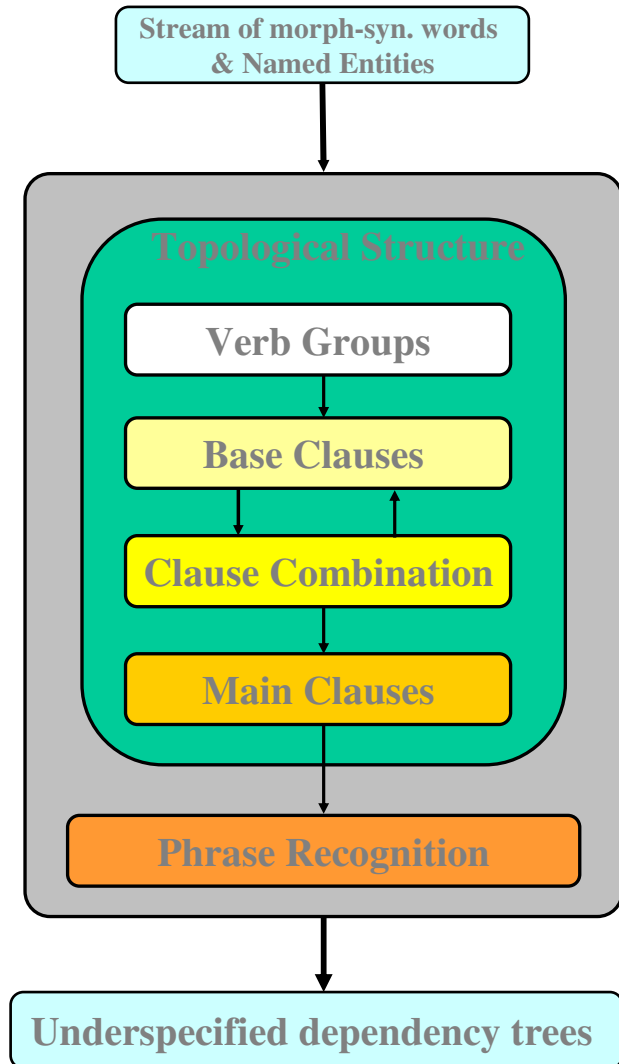
FrontField LeftVerb MiddleField RightVerb RestField

2. Apply general as well as domain-dependent phrasal grammars to the identified fields of the main and sub-clauses

[CoordS [CSent Diese Angaben konnte der Bundesgrenzschutz aber nicht bestätigen], [CSent Kinkel sprach von Horrorzahlen, [Relcl denen er keinen Glauben schenke]]].

This information couldn't be verified by the Border Police, Kinkel spoke of horrible figures that he didn't believe.

The divide-and-conquer parser is realized by means of a cascade of finite state grammars



Weil die Siemens GmbH, die vom Export lebt, Verluste erlitt, mußte sie Aktien verkaufen.

Because the Siemens Corp which strongly depends on exports suffered from losses they had to sell some shares.

Weil die Siemens GmbH, die vom Export **Verb-FIN**, Verluste **Verb-FIN**, **Modv-FIN** sie Aktien **FV-Inf**.

Weil die Siemens GmbH, **Rel-Clause** Verluste **Verb-FIN**, **Modv-FIN** sie Aktien **FV-Inf**.

Subconj-Clause,
Modv-FIN sie Aktien **FV-Inf**.

Clause

Semantic Analysis

Selected Approaches (1)

- *Chunk linking and chunk attachment (Abney)*
 - Interpretation steps in partial parsing
 - linking of hitherto unconnected structures (attachment of modifiers, prepositional phrases, determination of subject and object)
 - interpretation basis: case frames, corpus examples
- *Finite state filtering (Grefenstette, 1999)*
 - layered finite-state parser
 - groups adjacent syntactically related units
 - extracts non-adjacent n-ary grammatical relations.
 - high level specifications of regular expressions or describing the patterns to be extracted.

Semantic Analysis

Selected Approaches (2)

- *head-modifier-pairs*
 - mass data parsing with identifying pairs like [H: extraction, M: information]
 - used in information retrieval for enriching the document index and improving retrieval efficiency (Strzalkowski/Lin/Ge/Perez-Carballo, Jose (1999)).
- *fact extraction in fixed domains*
 - information patterns in highly standardized text types (weather forecasts, stock market reports)
 - example: biography
 - [A-Z][a-z]*“, “[A-Z][a-z]*“, *[0-9]{4}“ in “[A-Z][a-z]*“, † „[0-9]{4}“ in “[A-Z][a-z]*
 - *Buonarroti, Michelangelo, *1475 in Caprese , † 1564 in Roma*

Semantic Analysis

Selected Approaches (3)

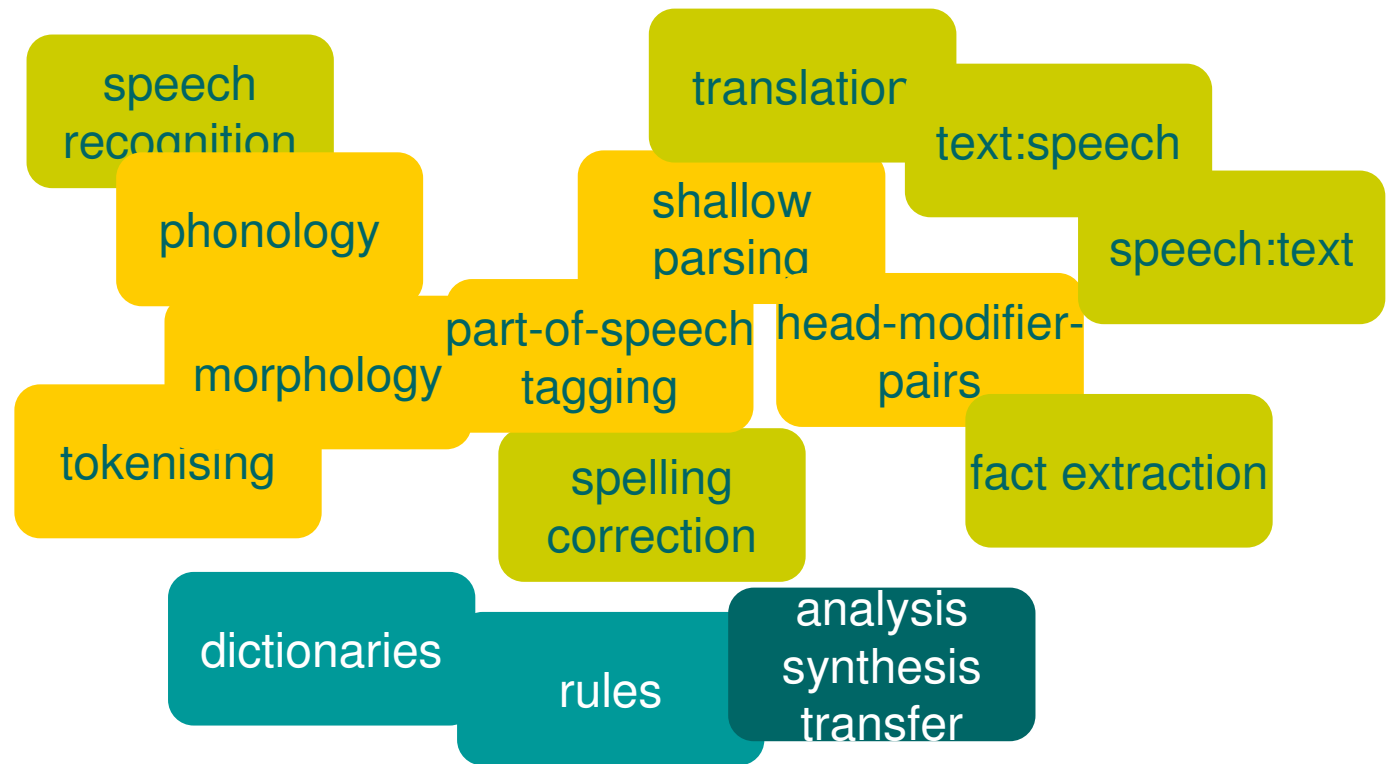
- *message understanding/information extraction*
 - filling in relational database templates from newswire texts
 - approach of FASTUS ¹⁾: cascade of five transducers
 - recognition of names,
 - fixed form expressions,
 - basic noun and verb groups
 - patterns of events
 - <company> <form><joint venture> with <company>
 - "Bridgestone Sports Co. said Friday it has set up a joint venture in Taiwan with a local concern and a Japanese trading house to produce golf clubs to be shipped to Japan."

| | |
|----------------|---|
| Relationship | TIE-UP |
| Entities | Bridgestone Sports Co. a local concern a Japanese trading house |
| JV Company | - |
| Capitalization | - |

¹⁾ Hobbs/Appelt/Bear/Israel/Kehler/Martin/Meyers/Kameyama/Stickel/Tyson (1997)

Application Perspective

- Introduction
- Morphology
- Syntax
- Semantics
- Summary



References

- Abney, Steven (1996). Tagging and Partial Parsing. In: Ken Church, Steve Young, and Gerrit Bloothoof (eds.), *Corpus-Based Methods in Language and Speech*. Kluwer Academic Publishers, Dordrecht. <http://www.vinartus.net/spa/95a.pdf>
- Abney, Steven (1996a) *Cascaded Finite-State Parsing*. Viewgraphs for a talk given at Xerox Research Centre, Grenoble, France. <http://www.vinartus.net/spa/96a.pdf>
- Abney, Steven (1995). Partial Parsing via Finite-State Cascades. In: *Journal of Natural Language Engineering*, 2(4): 337-344. <http://www.vinartus.net/spa/97a.pdf>
- Barton Jr., G. Edward; Berwick, Robert, C. und Eric Sven Ristad (1987). *Computational Complexity and Natural Language*. MIT Press.
- Beesley Kenneth R. und Lauri Karttunen (2003). *Finite-State Morphology*. Distributed for the Center for the Study of Language and Information. (CSLI- Studies in Computational Linguistics)
- Bod, Rens (1998). *Beyond Grammar. An Experienced-Based Theory of Language*. CSLI Lecture Notes, 88, Standford, California: Center for the Study of Information and Language
- Grefenstette, Gregory (1999). Light Parsing as Finite State Filtering. In: Kornai 1999, S. 86-94. earlier version in: Workshop on Extended finite state models of language, Budapest, Hungary, Aug 11--12, 1996. ECAI'96. <http://citeseer.nj.nec.com/grefenstette96light.html>
- Hobbs, Jerry; Doug Appelt, John Bear, David Israel, Andy Kehler, David Martin, Karen Meyers, Megumi Kameyama, Mark Stickel, Mabry Tyson (1997). *Breaking the Text Barrier*. FASTUS Presentation slides. SRI International. <http://www.ai.sri.com/~israel/Generic-FASTUS-talk.pdf>
- Jurafsky, Daniel und James H. Martin (2000): *Speech and Language Processing. An Introduction to Natural Language Processing, Computational Linguistics and Speech Recognition*. New Jersey: Prentice Hall.
- Kornai, András (ed.) (1999). *Extended Finite State Models of Language*. (Studies in Natural Language Processing). Cambridge: Cambridge University Press.
- Koskenniemi, Kimmo (1983). *Two-level morphology: a general computational model for word-form recognition and production*. Publication 11, University of Helsinki. Helsinki: Department of Genral Linguistics

References

- Kunze, Jürgen (2001). *Computerlinguistik. Voraussetzungen, Grundlagen, Werkzeuge*. Vorlesungsskript. Humboldt Universität zu Berlin. http://www2.rz.hu-berlin.de/compling/Lehrstuhl/Skripte/Computerlinguistik_1/index.html
- Manning, Christopher D.; Schütze, Hinrich (1999). *Foundations of Statistical Natural Language Processing*. Cambridge, Mass., London: The MIT Press. <http://www.sultry.arts.usyd.edu.au/fsnlp>
- Mohri, Mehryar (1997). Finite State Transducers in Language and Speech Processing. In: *Computational Linguistics*, 23, 2, 1997, S. 269-311. <http://citeseer.nj.nec.com/mohri97finitestate.html>
- Mohri, Mehryar (1996). On some Applications of finite-state automata theory to natural language processing. In: *Journal of Natural Language Engineering*, 2, S. 1-20.
- Mohri, Mehryar und Michael Riley (2002). *Weighted Finite-State Transducers in Speech Recognition (Tutorial)*. Teil 1: <http://www.research.att.com/~mohri/postscript/icslp.ps>, Teil 2: <http://www.research.att.com/~mohri/postscript/icslp-tut2.ps>
- G. Neumann, C. Braun and J. Piskorski (2000) [A Divide-and-Conquer Strategy for Shallow Parsing of German Free Texts](#) Proceedings of ANLP-2000, Seattle, Washington, pages 239-246
- Partee, Barbara; ter Meulen, Alice and Robert E. Wall (1993). *Mathematical Methods in Linguistics*. Dordrecht: Kluwer Academic Publishers.
- Pereira, Fernando C. N. and Rebecca N. Wright (1997). Finite-State Approximation of Phrase-Structure Grammars. In: Roche/Schabes 1997.
- Roche, Emmanuel und Yves Schabes (Eds.) (1997). *Finite-State Language Processing*. Cambridge (Mass.) und London: MIT Press.
- Sproat, Richard (2002). *The Linguistic Significance of Finite-State Techniques*. February 18, 2002. <http://www.research.att.com/~rws>
- Strzalkowski, Tomek; Lin, Fang; Ge, Jin Wang; Perez-Carballo, Jose (1999). Evaluating Natural Language Processing Techniques in Information Retrieval. In: Strzalkowski, Tomek (Ed.): *Natural Language Information Retrieval*, Kluwer Academic Publishers, Holland : 113-145
- Woods, W.A. (1970). Transition Network Grammar for Natural Language Analysis. In: *Communications of the ACM* 13: 591-602.

Named Entity Extraction

Machine Learning for
Named Entity
Extraction

The who, where, when & how much in a sentence

- The task: identify lexical and phrasal information in text which express references to named entities NE, e.g.,
 - person names
 - company/organization names
 - locations
 - dates×
 - percentages
 - monetary amounts
- Determination of an NE
 - Specific type according to some taxonomy
 - Canonical representation (template structure)

Example of NE-annotated text

Delimit the named entities in a text and tag them with NE types:

```
<ENAMEX TYPE=„LOCATION“>Italy</ENAMEX>'s business world was rocked by the announcement <TIMEX TYPE=„DATE“>last Thursday</TIMEX> that Mr. <ENAMEX TYPE=„PERSON“>Verdi</ENAMEX> would leave his job as vice-president of <ENAMEX TYPE=„ORGANIZATION“>Music Masters of Milan, Inc</ENAMEX> to become operations director of <ENAMEX TYPE=„ORGANIZATION“>Arthur Andersen</ENAMEX>.
```

- „Milan“ is part of organization name
- „Arthur Andersen“ is a company
- „Italy“ is sentence-initial ⇒ capitalization useless

Difficulties of Automatic NEE

- Potential set of NE is too numerous to include in dictionaries/Gazetteers
- Names changing constantly
- Names appear in many variant forms
- Subsequent occurrences of names might be abbreviated

- list search/matching does not perform well
 - context based pattern matching needed

Difficulties for Pattern Matching Approach

Whether a phrase is a named entity, and what name class it has, depends on

– Internal structure:

„Mr. Brandon“

– Context:

„The new company, SafeTek, will make air bags.“

– Feiyu Xu, researcher at DFKI, Saarbrücken

NE is an interesting problem

- Productivity of name creation requires lexicon free pattern recognition
- NE ambiguity requires resolution methods
- Fine-grained NE classification needs fine-grained decision making methods
 - Taxonomy learning
- Multi-linguality
 - A text might contain NE expressions from different languages
 - New pilot challenge in [ACE'2007](#)
 - Extract all NEs mentioned in a Mandarin/Arabic text
 - Translate them to English

NE Co-reference

*Norman Augustine ist im Grunde seines Herzens ein friedlicher Mensch. "Ich könnte niemals auf irgend etwas schießen", versichert der 57jährige Chef des US-Rüstungskonzerns **Martin Marietta Corp. (MM)**. ... Die Idee zu diesem Milliardendeal stammt eigentlich von GE-Chef John F. Welch jr. Er schlug Augustine bei einem Treffen am 8. Oktober den Zusammenschluss beider Unternehmen vor. Aber Augustine zeigte wenig Interesse, **Martin Marietta** von einem zehnfach grösseren Partner schlucken zu lassen.*

- **Martin Marietta** can be a person name or a reference to a company
- If **MM** is not part of an abbreviation lexicon, how do we recognize it?
 - Also by taking into account NE reference resolution.

References for NEE

- Alfonseca, Enrique; Manandhar, S. 2002. An Unsupervised Method for General Named Entity Recognition and Automated Concept Discovery. In *Proc. International Conference on General WordNet*.
- D. Bikel, S. Miller, Richard Schwartz and Ralph Weischedel: "Nymble: a High-Performance Learning Name Finder" ANLP 1997.
- Chen, H. H.; Lee, J. C. 1996. Identification and Classification of Proper Nouns in Chinese Texts. In *Proc. International Conference on Computational Linguistics*.
- Fleischman, Michael; Hovy. E. 2002. Fine Grained Classification of Named Entities. In *Proc. Conference on Computational Linguistics*.
- Evans, Richard. 2003. A Framework for Named Entity Recognition in the Open Domain. In *Proc. Recent Advances in Natural Language Processing*.
- Huang, Fei. 2005. *Multilingual Named Entity Extraction and Translation from Text and Speech*. Ph.D. Thesis. Pittsburgh: Carnegie Mellon University.
- M. Collins, Y. Singer "Unsupervised Models for Named Entity Classification", EMNLP 1999.
- A. McCallum and W. Li, "Early Results for Named Entity Recognition with Conditional Random Fields, Features Induction and Web-Enhanced Lexicons", CoNLL 2003.
- Nadeau, David, Satoshi Sekine (2006). A survey of named entity recognition and classification, Special issue of *Linguisticæ Investigationes* 30:1 (2007)
 - I have a copy in my office.
- Neumann (2007): web course page
 - <http://www.dfki.de/%7Eneumann/meta-ner/SoftWareProject.html>
- http://en.wikipedia.org/wiki/Named_entity_recognition
- E. Riloff, R. Jones: "Learning Dictionaries for Information Extraction by Multi-Level Bootstrapping", AAAI 1999.
- Yangarber, Lin, Grishman, Coling 2002
- Lin, Yangarber, Grishman, ICML 2003

Chunk parsing exercise

- **Goal:**
 - Write a program that uses regular expressions to recognize noun groups/chunks.
- **Starting point:**
 - <http://www.cnts.ua.ac.be/conll2000/chunking/>
 - Check CoNLL 2002 homepage and learn about goals and data format
- **Two different data sets are issued**
 - Download training and test data
 - Evaluation script (in Perl)

Chunk parsing exercise

- Form of annotation:
 - Each line a word + its annotation
 - POS
 - NP/PP/VP chunks
- Sequence of words are annotated according to the IOB1 standard
 - I-XP: words inside a chunk
 - B-XP: beginning of a XP chunk
 - O: for all elements outside any chunk
 - XP stands for NP, PP, VP, ADJP, SBAR

Example of an chunk annotated sentence in CoNLL format

He PRP B-NP
reckons VBZ B-VP
the DT B-NP
current JJ I-NP
account NN I-NP
deficit NN I-NP
will MD B-VP
narrow VB I-VP
to TO B-PP
only RB B-NP
I-NP
1.8 CD I-NP
billion CD I-NP
in IN B-PP
September NNP B-NP
.. O

In order to use standard RE packages transform data into a string

- **Input string**
“The/**Det** woman/**NN** will/**MD** give/**VB** Mary/**NNP** a/**Det** book/**NN**”
- **Output string**
“The/**B-NP** woman/**I-NP** will/**B-VP** give/**I-VP** Mary/**B-NP** a/**B-NP** book/**I-NP**”
- **Alternative:**
 - In: “DET NN MD VB NNP Det NN”
 - Out: “B-NP I-NP ... “
- **Note:** Gets easier input, probably more difficult output

In order to simplify the problem

- You might implement separate programs for recognizing each chunk type individually

Structure of program: read in data

- Read in each sentence from file (all tokens between `\newline`) using `readline()` function
- Two possibilities
 - Fetch two first elements (which are separated by `\tab`)
 - Just fetch POS tag (second element)
- Append each such token to a string (initialized by „“)

Structure of program: define regular expressions (e.g., Python)

```
import re
# pattern for html tags of form <TAG> or </TAG> or both ending with digit
# if pattern is found bind it to pattern variable ?P<tag>

def apply_re(f):
    expr = re.compile(r"(?P<tag>(<(/?)?[a-zA-Z]+(\d)?>))")
    file = open(f, 'r')
    for line in file.readlines():
        # searches only first match
        #res = expr.search(line)
        # searches all matches and binds it to a list
        res = expr.findall(line)
        if res == None:
            print "No match!"
        else:
            for tg in res:
                print tg
    return file.close()
```

Write result to output in CoNLL format

- Transform each string to CoNLL format
- Append it to some output file
- Call Perl script `conlleval.txt` to evaluate your results

