

	<h2>Language Technologies</h2>
	<p>“New Media and eScience” MSc Programme Jožef Stefan International Postgraduate School</p> <p style="text-align: center;">Winter Semester, 2008/09</p> <p>Lecture III. Computer Corpora</p> <p><u>Tomaž Erjavec</u></p>

	<h2>Overview of the lecture</h2>
	<ol style="list-style-type: none">1. Background2. Corpus compilation and markup3. Morphosyntactic tagging

	<h2>Background</h2>
	<ul style="list-style-type: none">• What is a corpus?• Using corpora• Characteristics of a corpus• Typology of corpora• History• Slovene language corpora

A corpus is ~

- a large collection of texts
- in digital format
- language "as it is"
- a sample of the language it is meant to represent
- used for describing language (descriptive/empirical linguistics)

A more precise definition

- **Corpus** (plural *corpora*) is Latin for *body*
- Guidelines of the Expert Advisory Group on Language Engineering Standards, **EAGLES**:
 - **Corpus**: *A collection of pieces of language that are selected and ordered according to explicit linguistic criteria in order to be used as a sample of the language.*
 - **Computer corpus**: *a corpus which is encoded in a standardised and homogeneous way for open-ended retrieval tasks. Its constituent pieces of language are documented as to their origins and provenance.*
- For computer scientists: a dataset

Using corpora

- Applied linguistics:
 - *Lexicography*: making dictionaries (first users of corpora)
 - *Translation studies*: translation equivalents with contexts translation memories, machine aided translations
 - *Language learning*: real-life examples, curriculum development
- Corpus linguistics:
 - linguistics based not on introspection, but on observation of real data
- *Language technology*:
 - testing set for developed methods;
 - *training set* for inductive learning (statistical Natural Language Processing)

Characteristics of a (good) corpus

- *Quantity:*
the bigger, the better
- *Quality:*
the texts are authentic; the mark-up is validated
- *Simplicity:*
the computer representation is understandable, with the markup easily separated from the text
- *Documented:*
the corpus contains bibliographic and other meta-data

Typology of corpora I.

- Medium:
 - *written language*
 - *spoken language* (spoken, but in writing / transcription)
 - *speech corpora* (actual speech signal)
- Content:
 - *reference corpora* (representative), e.g. BNC
 - *sub-language corpora* (specialised), e.g. COLT
- Structure:
 - corpora with *integral* texts
 - corpora or of text *samples* (historical and legal reasons)
e.g. Brown

Typology of corpora II

- Time:
 - *static corpora*
 - *monitor corpora* (language change)
- Languages:
 - *monolingual corpora*
 - multilingual *parallel corpora* (e.g. Hansard, Europarl)
 - multilingual *comparable corpora*
- Annotation:
 - *plain text corpora*
 - *annotated corpora*

	<h2>Reference corpora</h2>
	<ul style="list-style-type: none"> ■ Characteristics: <ul style="list-style-type: none"> - a sample of the "complete" language - large, expensive, detailed and explicit design criteria - typically of contemporary language - documented and annotated - legally clean, available ■ Criteria for including texts: <ul style="list-style-type: none"> - representativeness: corpus includes "all" text types - balance: the sizes of text type samples are in proportion to their "importance" for the speakers of the language ■ methodology v.s. practical constraints

	<h2>History</h2>		
	<table style="width: 100%; border: none;"> <tr> <td style="width: 50%; vertical-align: top;"> <p>History of Computational linguistics:</p> <ul style="list-style-type: none"> ■ 1950 -- 1960: empiricism weak computers: frequency lists ■ 1970 -- 1980: cognitive modeling (generative approaches, artificial intelligence) deep analysis / "basic science": computational linguistics ■ 1990 -- ...: empiricist revival, also combined approaches quantity / usefulness: language technologies ■ 2000 -- ...: The Web </td> <td style="width: 50%; vertical-align: top;"> <p>History of computer corpora:</p> <ul style="list-style-type: none"> ■ First milestones: <u>Brown</u> (1 million words) 1964; <u>LOB</u> (also 1M) 1974 ■ The spread of reference corpora: Cobuild Bank of English (monitor, 100..200..M) 1980; <u>BNC</u> (100M) 1995; Czech <u>CNC</u> (100M) 1998; Croatian <u>HNK</u> (100M) 1999... ■ Slovene language reference corpora: <u>FIDA</u> (100M), <u>Nova Beseda</u> (100M...) 1998; <u>FIDA+</u> (600M) 2006. ■ EU corpus oriented projects in the '90: <u>NERC</u>, <u>MULTEXT-East</u>,... ■ Language resources brokers: <u>LDC</u> 1992, <u>ELRA</u> 1995 ■ Web as Corpus (2000) ■ more, larger, for more languages, with diverse annotations </td> </tr> </table>	<p>History of Computational linguistics:</p> <ul style="list-style-type: none"> ■ 1950 -- 1960: empiricism weak computers: frequency lists ■ 1970 -- 1980: cognitive modeling (generative approaches, artificial intelligence) deep analysis / "basic science": computational linguistics ■ 1990 -- ...: empiricist revival, also combined approaches quantity / usefulness: language technologies ■ 2000 -- ...: The Web 	<p>History of computer corpora:</p> <ul style="list-style-type: none"> ■ First milestones: <u>Brown</u> (1 million words) 1964; <u>LOB</u> (also 1M) 1974 ■ The spread of reference corpora: Cobuild Bank of English (monitor, 100..200..M) 1980; <u>BNC</u> (100M) 1995; Czech <u>CNC</u> (100M) 1998; Croatian <u>HNK</u> (100M) 1999... ■ Slovene language reference corpora: <u>FIDA</u> (100M), <u>Nova Beseda</u> (100M...) 1998; <u>FIDA+</u> (600M) 2006. ■ EU corpus oriented projects in the '90: <u>NERC</u>, <u>MULTEXT-East</u>,... ■ Language resources brokers: <u>LDC</u> 1992, <u>ELRA</u> 1995 ■ Web as Corpus (2000) ■ more, larger, for more languages, with diverse annotations
<p>History of Computational linguistics:</p> <ul style="list-style-type: none"> ■ 1950 -- 1960: empiricism weak computers: frequency lists ■ 1970 -- 1980: cognitive modeling (generative approaches, artificial intelligence) deep analysis / "basic science": computational linguistics ■ 1990 -- ...: empiricist revival, also combined approaches quantity / usefulness: language technologies ■ 2000 -- ...: The Web 	<p>History of computer corpora:</p> <ul style="list-style-type: none"> ■ First milestones: <u>Brown</u> (1 million words) 1964; <u>LOB</u> (also 1M) 1974 ■ The spread of reference corpora: Cobuild Bank of English (monitor, 100..200..M) 1980; <u>BNC</u> (100M) 1995; Czech <u>CNC</u> (100M) 1998; Croatian <u>HNK</u> (100M) 1999... ■ Slovene language reference corpora: <u>FIDA</u> (100M), <u>Nova Beseda</u> (100M...) 1998; <u>FIDA+</u> (600M) 2006. ■ EU corpus oriented projects in the '90: <u>NERC</u>, <u>MULTEXT-East</u>,... ■ Language resources brokers: <u>LDC</u> 1992, <u>ELRA</u> 1995 ■ Web as Corpus (2000) ■ more, larger, for more languages, with diverse annotations 		

	<h2>Slovene language corpora</h2>
	<p>Monolingual reference corpora:</p> <ul style="list-style-type: none"> ■ ZRC SAZU: <u>Beseda</u>, 1998; <u>Nova beseda</u>, 2000- ■ DZS, Amebis, FF, IJS: <u>FIDA</u>, 1998, <u>FidaPlus</u>, 2006 <p>Parallel corpora:</p> <ul style="list-style-type: none"> ■ IJS: <u>MULTEXT-East</u> 1998--, <u>IJS-FLAN</u> 1999--, <u>SVEZ-IJS</u>, 2004, <u>JRC-ACQUIS</u>, 2006 ■ SVEZ: EuroKorpus ■ FF: <u>TRANS</u>, 2002 ■ UP: Turist Corpus, 2008 <p>Speech corpora:</p> <ul style="list-style-type: none"> ■ Laboratory for Digital Signal Processing, University of Maribor: <u>SpeechDat</u>, <u>ONOMASTICA</u>... ■ Laboratory of Artificial Perception, Systems and Cybernetics, University of Ljubljana: <u>SQEL</u>, <u>GOPOLIS</u>,...

	<h2>II. Compilation and markup of corpora</h2>
	<ul style="list-style-type: none"> •Steps in the preparation of a corpus •What annotation can be added to the text •Computer coding of corpora •Markup Methods

	<h2>Before making your own corpus</h2>
	<p>check if an appropriate corpus is already available</p> <ul style="list-style-type: none"> ■ google ■ corpora@lists.uib.no ■ LDC, ELRA

	<h2>Steps in the preparation of a corpus</h2>
	<ol style="list-style-type: none"> 1. Choosing the component texts and acquiring digital originals 2. Up-translation to standard format 3. Linguistic annotation 4. Documentation 5. Use and Dissemination

Getting the text

1. Choosing the component texts:
linguistic and non-linguistic criteria;
availability; simplicity; size
2. Copyright
sensitivity of source (financial and
privacy considerations); agreement
with providers; usage, publication
3. Acquiring digital originals
OCR; digital originals; Web
 - BootCat

Processing

1. Conversion to common format
consistency; character set encodings;
structure
 - Web as Corpus: Wacky tools
2. Documentation
e.g. TEI header; Open Archives etc.
3. Linguistic annotation
language dependent methods; errors

Use and dissemination

- Using the corpus:
 - concordancer (linguists)
e.g. NovaBeseda, FidaPLUS, SKE,
iKorpus
 - statistics extraction
 - development of new methods for analysis
- Dissemination:
 - legalities (source copyright, corpus use
agreement)
 - mode: concordancer or dataset

Computer coding of corpora

- Encoding must ensure
 - durability
 - interchange between computer platforms
 - interchange between applications
- Basic standard: *Extended Markup Language, XML*
 - a number of companion standards and technologies: XSLT, XML Schema, ISO Relax NG, XPath, XQuery, ...
- The vocabulary of annotations for corpora and other language resources are defined by the *Text Encoding Initiative, TEI*
- XML/TEI used much wider than just for corpora:
 - annotation of dictionaries: *English-Slovene, Japanese-Slovene* (from jaSlo)
 - for annotating *text-critical editions*

Corpus annotation

- Annotation = interpretation
- Documentation about the corpus (*example*)
 - Document structure (*example*)
 - Basic linguistic markup: sentences, words (*example*), punctuation, abbreviations (*example*)
 - Lemmas and morphosyntactic descriptions (*example*)
 - Syntax (*example*)
 - Alignment (*example*)
 - Terms, semantics, anaphora, pragmatics, intonation,...

Example: TEI header

```
<teiHeader id="ecmr.H" type="text" lang="sl-en" creator="ET" status="update"
date.created="1999-04-13" date.updated="1999-06-22" >
<fileDesc>
<titleStmt>
<title lang="sl">Ekonomsko ogledalo; 13 &scaron;tevilik 98/99</title>
<title lang="en">Slovenian Economic Mirror; 13 issues, 98/99</title>
<respstmt>
<name>Andrej Skubic, FF</name>
<resp lang="sl">Zagotovitev digitalnega originala, poravnava</resp>
<resp lang="en">Provision of digital original, alignment</resp>
<name>Tomaž&scaron; Erjavec, IJS</name>
<resp lang="sl">Tokenizacija, pretvorba v TEI</resp>
<resp lang="en">Tokenisation, conversion to TEI</resp>
</respstmt>
</titleStmt> ...
```

Example: text structure

```
<quote id="Osl.1.8.18" rend="center;it">
<lg id="Osl.1.8.18.1">
  <l id="Osl.1.8.18.1.1">Tam pod kostanjevim drevesom</l>
  <l id="Osl.1.8.18.1.2">izdala si me,</l>
  <l id="Osl.1.8.18.1.3">izdal sem te,</l>
  <l id="Osl.1.8.18.1.4">ne da bi trenila z očesom.</l>
</lg>
</quote>
<p id="Osl.1.8.19">
  <s id="Osl.1.8.19.1">Trije možje se niso niti ganili.</s>
  <s id="Osl.1.8.19.2">Toda ko je <name>Winston</name>
  znova pogledal v Rutherfordov propadli obraz, je opazil, da so
  njegove oči polne solz.</s> ...
```

Example: morphosyntactic tagging

```
<s id="Osl.1.2.2.1">
<w lemma="biti" ana="Vcps-sma">Bil</w>
<w lemma="biti" ana="Vcip3s-n">je</w>
<w lemma="jasen" ana="Afpmsnn">jasen</w><c>,</c>
<w lemma="mrzel" ana="Afpmsnn">mrzel</w>
<w lemma="aprilski" ana="Aopmsn">aprilski</w>
<w lemma="dan" ana="Ncmsn">dan</w>
<w lemma="in" ana="Ccs">in</w>
<w lemma="ura" ana="Ncfpn">ure</w>
<w lemma="biti" ana="Vcip3p-n">so</w>
<w lemma="biti" ana="Vmpps-pfa">bile</w>
<w lemma="trinajst" ana="Mcnpl">trinajst</w><c>.</c>
</s>
```

Example: alignment

```
<linkGrp id="Osl.1" type="body" targtype="s"
domains="Oen Osl">
  <link xtargets="Osl.1.2.2.1 ; Oen.1.1.1.1">
  <link xtargets="Osl.1.2.2.2 ; Oen.1.1.1.2">
  <link xtargets="Osl.1.2.3.1 ; Oen.1.1.2.1">
  <link xtargets="Osl.1.2.3.2 ; Oen.1.1.2.2">
  ...
  <link xtargets="Osl.1.2.6.5 ; Oen.1.1.5.5">
  <link xtargets="Osl.1.2.6.6 ; Oen.1.1.5.6 Oen.1.1.5.7">
  <link xtargets="Osl.1.2.6.7 ; Oen.1.1.5.8">
  ...
```

	<h2>Methods for linguistic markup</h2>
	<ul style="list-style-type: none"> ■ <i>hand annotation</i>: documentation, first steps generic (XML, spreadsheet) editors or specialised editors ■ <i>semi-automatic</i>: morphosyntactic and other linguistic annotation cyclic approach: machine, hand, validate, correct, machine, ... ■ <i>machine, with hand-written rules</i>: tokenisation regular expression ■ <i>machine, with inductively built models from annotated data</i>: "supervised learning"; HMMs, decision trees, inductive logic programming,... ■ <i>machine, with inductively built models from un-annotated data</i>: "unsupervised learning"; clustering techniques ■ <u>overview of the field</u>

	<h2>III. Morphosyntactic tagging</h2>
	<ul style="list-style-type: none"> ■ Better known as part-of-speech (PoS) tagging ■ Tagging is the task of labeling each word in a sequence of words with its appropriate part-of-speech ■ Words are often ambiguous with respect to their POS: <ul style="list-style-type: none"> - <i>saw</i> → singular noun - <i>saw</i> → past tense of verb <i>see</i> ■ Purposes and applications (examples): <ul style="list-style-type: none"> - pre-processing step for further analyses: <ul style="list-style-type: none"> ■ lemmatisation ■ syntactic structure, etc. - text indexing, e.g. nouns are more useful than verbs - pronunciation in speech processing

	<h2>Steps in tagging</h2>
	<ul style="list-style-type: none"> ■ for each word token in text the tagger needs to know all its possible tags --> a morphological lexicon ■ given the context in which the word appears in, the tagger must decide in the correct tag: <ul style="list-style-type: none"> - he saw/V a man carrying a saw/N ■ so, tagging performs a limited syntactic disambiguation

Example: Penn Treebank

Under/IN the/DT proposal/NN ,/, Delmed/NNP
would/MD issue/VB about/IN 123.5/CD
million/CD additional/JJ Delmed/NNP
common/JJ shares/NNS to/TO
Fresenius/NNP at/IN an/DT average/JJ
price/NN of/IN about/IN 65/CD cents/NNS
a/DT share/NN ,/, though/IN under/IN
no/DT circumstances/NNS more/JJR than/IN
75/CD cents/NNS a/DT share/NN ./.

PoS taggers

- Most taggers induce the language model from a hand-annotated corpus
- Typically, two resources are induced:
 - lexicon, giving the possible tags of a word and their frequencies in the training corpus
 - how the tag of a word depends on its local context

Tagging with Markov Models

- Sequence of tags in a text is regarded a Markov chain
- Limited horizon: A word's tag only depends on the previous tag: $p(x_{i+1} = \tilde{t} \mid x_1, \dots, x_i) = p(x_{i+1} = \tilde{t} \mid x_i)$
- Time invariant: This dependency does not change over time: $p(x_{i+1} = \tilde{t} \mid x_i) = p(x_2 = \tilde{t} \mid x_1)$
- Task: Find the most probable tag sequence for a sequence of words
- Maximum likelihood estimate of tag t^* following \tilde{t} :
 $p(t^* \mid \tilde{t}) = f(\tilde{t}, t^*) / f(\tilde{t})$
- Optimal tags for a sentence:
 $t'_{1,n} = \arg \max p(t_{1,n} \mid w_{1,n}) = \prod p(w_i \mid t_i) p(t_i \mid t_{i-1})$

	<h2>Most popular Markov model tagger</h2>
	<ul style="list-style-type: none"> ■ TnT (Trigrams 'n Tags) ■ induces lexicon and tag trigrams from the training corpus ■ has heuristics to tag unknown words ■ has no problem with large tagsets ■ fast in training and tagging ■ freely available for non-commercial use: http://www.coli.uni-saarland.de/~thorsten/tnt/ ■ but only as a Linux executable

	<h2>TreeTagger</h2>
	<ul style="list-style-type: none"> ■ uses decision trees ■ relatively fast ■ comes with lots of models for various languages ■ executables freely available http://www.ims.uni-stuttgart.de/projekte/corplex/TreeTagger/

	<h2>Transformation-based Tagging (TbT)</h2>
	<ul style="list-style-type: none"> ■ Basic idea: transform an imperfect tagging into one with fewer errors by changing wrong tags ■ Features that trigger changes can be conditioned on words and on more context and are user specified ■ Components: <ul style="list-style-type: none"> - specification of transformations - learning algorithm: constructs a ranked list of transformations ■ A transformation consists of two parts: <ul style="list-style-type: none"> - triggering environment + rewrite rule ■ Examples: <ul style="list-style-type: none"> - if previous tag is TO and current tag is NN then change it to VB - if one of previous two words is <i>n?</i> and current tag is VBP then change it to VB - if next tag is JJ and current tag is JJR then change it to RBR - if one of previous three tags is MD and current tag is VBP then change it to VB

	Yet another Tagger
	<p>For a while, trying out new approaches to tagging was in fashion</p> <ul style="list-style-type: none"> ■ Maximum Entropy taggers ■ Support Vector Machine taggers ■ Memory based taggers ■ ...

	Tagsets
	<ul style="list-style-type: none"> ■ A tagset is a set of part-of-speech tags ■ Classical 8 classes (Thrax, 100 BC): noun, verb, article, participle, pronoun, preposition, adverb, conjunction ■ But all tagset use more tags than that! ■ Criteria: <ul style="list-style-type: none"> - specificity: degree to which humans use the tagset uniformly on the same text - accuracy: evaluation of output on tagged text - suitability for intended application

	Tagsets for English
	<ul style="list-style-type: none"> ■ For English, there exist several tagsets: Brown, CLAWS, Penn, ... ■ English tagsets include PoS + some other morphological (inflectional) properties: 30-80 tags ■ Penn Treebank Tagset for English: 37 tags, e.g. <ul style="list-style-type: none"> - JJ adjective, positive - JJR adjective, comparative - JJS adjective, superlative - NN non-plural common noun - NNS plural common noun - NNP non-plural proper name - NNPS plural proper name - IN preposition - ...

	<h2>Morphosyntactic tagsets</h2>
	<ul style="list-style-type: none"> ■ For inflectionally rich languages (such as Slavic languages), tagsets contain much more information than just PoS ■ Slovene, Czech, etc. > 1000 different morphosyntactic tags <ul style="list-style-type: none"> – gendar, number, case, animacy, definiteness, ... ■ Efforts to standardise tagsets across languages: <ul style="list-style-type: none"> – Eagles – MULTEXT – MULTEXT-East

	<h2>MULTEXT-East</h2>
	<ul style="list-style-type: none"> ■ EU project in '90s: development of language resources for Central and East-European languages ■ also development of morphosyntactic specifications, lexica and annotated corpus ■ Parallel annotated corpus: <u>Orwell's 1984</u> ■ Several later releases, V3 in 2004 ■ Project Web site: http://nl.ijs.si/ME/

	<h2>MULTEXT-East morphosyntactic specifications</h2>
	<ul style="list-style-type: none"> ■ Specify <ul style="list-style-type: none"> – what morphosyntactic features particular languages distinguish, – what their names and values are, – how they can be mapped to tags (morphosyntactic descriptions, MSDs) ■ e.g. that <i>Ncms</i> is: <ul style="list-style-type: none"> – a valid for Slovene – is equivalent to <i>PoS:Noun, Type:common, Gender:masculine, Number:singular</i> ■ http://nl.ijs.si/ME/V3/msd/html/

JOS morphosyntactic specifications

- only for Slovene
- based on MULTTEXT-East but changed some features and lexical assignments
- also moved to 100% XML/TEI encoding
- bi-lingual (Slovene and English)
- also made annotated corpora:
 - jos100k (hand validated)
 - jos1M (partially hand validated)
- <http://nl.ijs.si/jos/> - still work in progress!
