

	<h2>Advanced Language Technologies</h2>
	<p>Information and Communication Technologies  Module "Knowledge Technologies"  <u>Jožef Stefan International Postgraduate School</u>  Winter 2011 / Spring 2012</p> <p>Lecture II.  Computer Corpora</p> <p><u>Tomaž Erjavec</u></p>

---

---

---

---

---

---

---

---

	<h2>Overview of the lecture</h2>
	<ol style="list-style-type: none"> <li>1. Background</li> <li>2. Corpus compilation and markup</li> <li>3. Morphosyntactic tagging</li> </ol>

---

---

---

---

---

---

---

---

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

---

---

---

---

---

---

---

---

	<h2>A corpus is:</h2>
	<ul style="list-style-type: none"> <li>• a large collection of texts</li> <li>• in digital format</li> <li>• language “as it is”</li> <li>• a sample of the language it is meant to represent</li> <li>• used for describing language (descriptive/empirical linguistics)</li> <li>• for computer scientists: a dataset</li> </ul>

---

---

---

---

---

---

---

---

	<h2>A more precise definition</h2>
	<ul style="list-style-type: none"> <li>• <b>Corpus</b> (plural <i>corpora</i>) is Latin for <i>body</i></li> <li>• Guidelines of the Expert Advisory Group on Language Engineering Standards, <b>EAGLES</b>: <ul style="list-style-type: none"> <li>• <b>Corpus</b> : <i>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.</i></li> <li>• <b>Computer corpus</b> : <i>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.</i></li> </ul> </li> </ul>

---

---

---

---

---

---

---

---

	<h2>Using corpora</h2>
	<ul style="list-style-type: none"> <li>• Applied linguistics: <ul style="list-style-type: none"> <li>• <i>Lexicography</i>: making dictionaries (first users of corpora)</li> <li>• <i>Translation studies</i>: translation equivalents with contexts translation memories, machine aided translations</li> <li>• <i>Language learning</i>: real-life examples, curriculum development</li> </ul> </li> <li>• Corpus linguistics: <ul style="list-style-type: none"> <li>• linguistics based not on introspection, but on observation of real data</li> </ul> </li> <li>• <i>Language technology</i>: <ul style="list-style-type: none"> <li>• testing set for developed methods</li> <li>• <i>training set</i> for inductive learning (<u>statistical Natural Language Processing</u>)</li> </ul> </li> </ul>

---

---

---

---

---

---

---

---

## 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 mark-up 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, JRC Acquis)
  - multilingual *comparable* corpora
- Annotation:
  - *plain text* corpora
  - *annotated* corpora

---

---

---

---

---

---

---

---

## Reference corpora

- Characteristics:
  - a sample of the "complete" language
  - large, expensive, detailed and explicit design criteria
  - typically of contemporary language
  - documented and annotated
  - legally clean, available (but usu. only via a concordancer)
- Criteria for including texts:
  - 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

---

---

---

---

---

---

---

---

## History of corpora

- First milestones:
  - Brown (1 million words) 1964; LOB (also 1M) 1974
- The spread of reference corpora:
  - Cobuild Bank of English (monitor, 100..200..M) 1980; BNC (100M) 1995; Czech CNC (100M) 1998; Croatian HNK (100M) 1999...
- EU corpus oriented projects in the '90: NERC, MULTEXT, MULTEXT-East,...
- Language resources brokers: LDC 1992, ELRA 1995
- Web as Corpus (2000..): ukWaC, itWaC, ... slWaC
- more, larger, for more languages, with diverse annotations: EUROPARL, JRC-ACQUIS, PDT, ...

---

---

---

---

---

---

---

---

## Slovene language corpora

- The „FIDA“ monolingual reference corpora (FF, IJS, DZS, Anebis):
- FIDA, 1998: 100M, ambiguous annotations
  - FidaPlus, 2006: 600M, unambiguous
  - Gigafida, 2012: 1000M, adds Web materials
- Freely available training sets:
- IJS, FF: JOS corpora (jos100k, jos1M)
  - SSJ project: ccFida, sj400k
- Parallel corpora:
- IJS: MULTEXT-East 1998-, SVEZ-IJS, 2004, JRC-ACQUIS, 2006
  - SVEZ: EuroKorpus
  - FF: TRANS, 2002
- Speech corpora:
- Laboratory for Digital Signal Processing, University of Maribor: SpeechDat, ONOMASTICA...
  - Laboratory of Artificial Perception, Systems and Cybernetics, University of Ljubljana: SQL, GOPOLIS,...

---

---

---

---

---

---

---

---

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

---

---

---

---

---

---

---

---

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

---

---

---

---

---

---

---

---

	<h2>Steps in the preparation of a corpus</h2>
	<ol style="list-style-type: none"> <li>1. Choosing the component texts and acquiring digital originals</li> <li>2. Up-translation to standard format</li> <li>3. Linguistic annotation</li> <li>4. Documentation</li> <li>5. Use</li> <li>6. Dissemination</li> </ol>

---

---

---

---

---

---

---

---

	<b>Getting the text</b>
	<ol style="list-style-type: none"> <li>1. Choosing the component texts: linguistic and non-linguistic criteria; availability; simplicity; size</li> <li>2. Copyright sensitivity of source (financial and privacy considerations); agreement with providers; usage, publication</li> <li>3. Acquiring digital originals OCR; digital originals; Web <ul style="list-style-type: none"> <li>▪ BootCat</li> </ul> </li> </ol>

---

---

---

---

---

---

---

---

	<b>Processing</b>
	<ol style="list-style-type: none"> <li>1. Conversion to common format consistency; character set encodings; structure <ul style="list-style-type: none"> <li>▪ Web as Corpus: Wacky tools</li> </ul> </li> <li>2. Documentation e.g. TEI header; Open Archives etc.</li> <li>3. Linguistic annotation language dependent methods; errors</li> </ol>

---

---

---

---

---

---

---

---

	<b>Use and dissemination</b>
	<ul style="list-style-type: none"> <li>• Using the corpus: <ul style="list-style-type: none"> <li>• concordancer (linguists) e.g. <u>Gigafida</u>, <u>SKE</u>, <u>iKorpus</u>, JOS, IMP</li> <li>• statistics extraction</li> <li>• development of new methods for analysis</li> </ul> </li> <li>• Dissemination: <ul style="list-style-type: none"> <li>• legalities (source copyright, corpus use agreement)</li> <li>• mode: concordancer or dataset</li> </ul> </li> </ul>

---

---

---

---

---

---

---

---

## Computer coding of corpora

- Encoding must ensure
  - durability
  - interchange between computer platforms
  - interchange between applications
- Basic standard: XML
  - companion standards: W3C Schema, ISO Relax NG, XSLT, XPath, XQuery, ...
- XML vocabulary of annotations of arbitrary texts: *Text Encoding Initiative*, TEI
- ISO TC 37 „Terminology and other language resources“: many standards for text encoding

---

---

---

---

---

---

---

---

## 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

```
<?xml version="1.0" encoding="utf-8"?>
<TEI xmlns="http://www.tei-c.org/ns/1.0" xml:lang="sl" xml:id="FPG_00008-1847">
  <teiHeader xml:lang="sl">
    <fileDesc>
      <titleStmt>
        <title>AHLib: Zschokke, Heinrich. "Čujte, čujte, kaj žganje dela!" (1847)</title>
        <principal>
          <name>Erich Prunč, Univerza Karl-Franzens v Gradcu</name>
        </principal>
        <respStmt>
          <name>Tomaž Erjavec, Institut "Jožef Stefan"</name>
          <resp>Računalniška obdelava</resp>
        </respStmt>
      </titleStmt>
      <editionStmt>
        <edition>1.0</edition>
      </editionStmt>
      <extent>124 pp</extent>
    </fileDesc>
  </teiHeader>
  ...

```

---

---

---

---

---

---

---

---

## 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="#Vcips-sma">Bil</w>
<w lemma="biti" ana="#Vcip3s-n">je</w>
<w lemma="jasen" ana="#Afpmsn">jasen</w><pc>,</pc>
<w lemma="mrzel" ana="#Afpmsn">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><pc>.</pc>
</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">
...
```

---

---

---

---

---

---

---

---



## Methods for linguistic markup

- *hand annotation*: documentation, first steps generic (XML, spreadsheet) or specialised editors
- *semi-automatic*: morphosyntactic and other linguistic annotation  
cyclic approach: machine, hand, validate, correct, machine, ...
- *machine, with hand-written rules*: tokenisation regular expression
- *machine, with inductive models*: "supervised learning"; HMMs, decision trees, inductive logic programming,...
- *machine, with inductively built models from un-annotated data*: "unsupervised learning"; clustering techniques
- overview of the field

---

---

---

---

---

---

---

---

## III. Morphosyntactic tagging

- 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:
  - *saw* → singular noun „I brought a saw“
  - *saw* → past tense of verb „I saw a tree“
- Purposes and applications (examples):
  - pre-processing step for further analyses:
    - lemmatisation
    - syntactic structure, etc.
  - text indexing, e.g. nouns are more useful than verbs
  - pronunciation in speech processing

---

---

---

---

---

---

---

---

## Steps in tagging

- for each word token in text the tagger needs to know all its possible tags (ambiguity class)  
→ a morphological lexicon
- given the context in which the word appears in, the tagger must decide in the correct tag:
  - he saw/V a man carrying a saw/N
- so, tagging performs 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 ambiguity class of a word and their frequencies in the training corpus
  - tag n-grams

---

---

---

---

---

---

---

---

## 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"> <li>• <u>TnT</u> (Trigrams 'n Tags)</li> <li>• induces lexicon and tag trigrams from the training corpus</li> <li>• has heuristics to tag unknown words</li> <li>• has no problem with large tagsets</li> <li>• fast in training and tagging</li> <li>• freely available for non-commercial use</li> <li>• but only as a Linux executable</li> <li>• OS alternative: <u>hunpos</u></li> </ul>

---

---

---

---

---

---

---

---

	<h2>Yet another Tagger</h2>
	<p>For a while, trying out new approaches to tagging was in fashion</p> <ul style="list-style-type: none"> <li>• Maximum Entropy taggers</li> <li>• Support Vector Machine taggers</li> <li>• Memory based taggers</li> <li>• ...</li> </ul>

---

---

---

---

---

---

---

---

	<h2>Tagsets</h2>
	<ul style="list-style-type: none"> <li>• A tagset is a set of part-of-speech tags</li> <li>• Classical 8 classes (Thrax, 100 BC): noun, verb, article, participle, pronoun, preposition, adverb, conjunction</li> <li>• But tagset typically use more tags than that!</li> <li>• Criteria: <ul style="list-style-type: none"> <li>• specifiability: degree to which humans use the tagset uniformly on the same text</li> <li>• accuracy: evaluation of output on tagged text</li> <li>• suitability for intended application</li> </ul> </li> </ul>

---

---

---

---

---

---

---

---

## Tagsets for English

- 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.
  - 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
  - ...

---

---

---

---

---

---

---

---

## Morphosyntactic tagsets

- For inflectionally rich languages (e.g. Slavic), tagsets contain much more information than just PoS
- Slovene, Czech, etc. > 1,000 different morphosyntactic tags
  - gender, number, case, animacy, definiteness, ...
- Efforts to standardise tagsets across languages:
  - Eagles
  - MULTEXT
  - MULTEXT-East

---

---

---

---

---

---

---

---

## MULTEXT-East

- EU project in '90s: development of language resources for Central and East-European languages
- Several later releases, V4 in 2010 (17 languages)
- Development of morphosyntactic specifications, lexica and annotated corpus
- Parallel annotated corpus: Orwell's 1984
- Web site: <http://nl.ijs.si/ME/>

---

---

---

---

---

---

---

---



## jos100k encoding

```
<s xml:id="F0020003.557.2">
  <w xml:id="F0020003.557.2.1" lemma="ta" msd="Zk-sei">To</w><S/>
  <w xml:id="F0020003.557.2.2" lemma="biti" msd="Gp-ste-n">je</w><S/>
  <term type="sloWNet" sortKey="kraj" key="ENG20-08114200-n">
    <w xml:id="F0020003.557.2.3" lemma="turističen,
      msd="Ppnmein">turističen</w><S/>
    <w xml:id="F0020003.557.2.4" lemma="kraj" msd="Somei">kraj</w>
  </term>
  <c xml:id="F0020003.557.2.5">.</c><S/>
</s>
<linkGrp type="syntax" targFunc="head argument" corresp="#F0020003.557.2">
  <link type="ena" targets="#F0020003.557.2.2 #F0020003.557.2.1"/>
  <link type="modra" targets="#F0020003.557.2 #F0020003.557.2.2"/>
  <link type="dol" targets="#F0020003.557.2.4 #F0020003.557.2.3"/>
  <link type="dol" targets="#F0020003.557.2.2 #F0020003.557.2.4"/>
  <link type="modra" targets="#F0020003.557.2 #F0020003.557.2.5"/>
</linkGrp>
```

---

---

---

---

---

---

---

---

---

---

## Processing Historical Language

- Interesting for diachronic linguistics and better access to digital libraries
- Problems:
  - difficult to obtain good transcriptions
  - great variation in spelling
  - no resources for tool training
- Historical slv:
  - Late standardisation (XIX ≠ XX)
  - Before 1850: f fh s sh z zh → s š z ž c č
  - No corpora/lexica of historical Slovene

---

---

---

---

---

---

---

---

---

---

## Background



- **AHLib** (2004–08)  
Deutsch-slowenische/kroatische Übersetzung 1848–1918
  - Scans + correction + (lemmatisation) of ger→slv books
  - AAS & Karl-Franzens University, Graz (prof. Erich Prunč)
  - JSI: correction & lemmatisation environment
- **EU IP IMPACT** (ext. 2010–2011)
  - Better OCR for historical texts
  - NUK: GTD transcriptions
  - JSI: (semi)manual lexicon construction
- **Google award** (2011+2012)  
Developing language models for historical Slovene
  - ZRC SAZU: transcriptions of old texts
  - JSI: annotating a corpus of XIX<sup>th</sup> century Slovene

42

---

---

---

---

---

---

---

---

---

---

	<h2>Producing the goo300k corpus</h2>
	<ul style="list-style-type: none"> <li>• Representative &amp; balanced, sampled</li> <li>• Corpus element: unbroken &amp; contiguous text from 1 page</li> <li>• Sampled by decade &amp; text</li> <li>• Target size: 1,000 pages (~300,000 tokens)</li> <li>• Encoded in TEI P5</li> <li>• Automatically annotated</li> <li>• Tool for manual annotation: IMPACT INL Cobalt</li> <li>• Annotator training &amp; management: May</li> <li>• Manual correction: June–November</li> <li>• Fixing bugs &amp; packaging: December - April</li> </ul>

---

---

---

---

---

---

---

---

	<h2>Annotation tool</h2>
	<p>Approach:</p> <ul style="list-style-type: none"> <li>• Modernise, then process as contemporary language</li> <li>• Language independent (trainable) modules</li> </ul> <p>Steps:</p> <ol style="list-style-type: none"> <li>1. <b>T</b>okenisation (miToken)</li> <li>2. <b>T</b>ranscription (Vaam)</li> <li>3. <b>T</b>agging (TnT)</li> <li>4. <b>L</b>emmatisation (CLOG)</li> </ol> <p>= ToTrTaLe</p> <ul style="list-style-type: none"> <li>■ Pipeline in Perl</li> <li>■ TEI P5 I/O</li> </ul>

---

---

---

---

---

---

---

---

	<h2>Extracted lexicon</h2>
	<ul style="list-style-type: none"> <li>■ Also encoded in TEI</li> <li>■ Lemma oriented</li> <li>■ Useful for enabling full-text searching in DL</li> <li>■ Also for humans: look up of extinct words</li> </ul>

---

---

---

---

---

---

---

---

## Conclusions

- What is a corpus
- How to make it
- How to annotate it
- Case studies: MULTEXT-East, JOS, IMP

---

---

---

---

---

---

---

---