# O R PUS

## Parallel Corpora for Everyone
### (NLPL-Activity G)
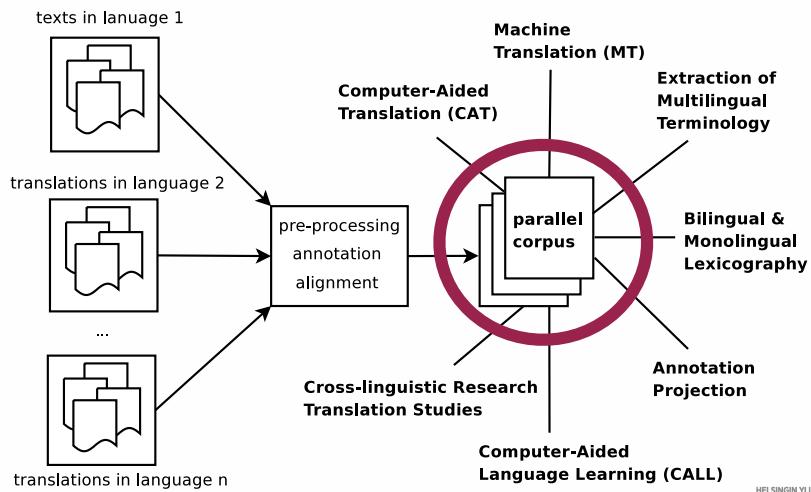
**Jörg Tiedemann**
**Department of Digital Humanities**
**University of Helsinki**

---

## Rosendal 2002

---

## The Amazing Utility of Parallel Corpora



texts in lanuage 1

translations in language 2

...

translations in language n

pre-processing
annotation
alignment

**parallel corpus**

**Machine Translation (MT)**

**Computer-Aided Translation (CAT)**

**Extraction of Multilingual Terminology**

**Bilingual & Monolingual Lexicography**

**Annotation Projection**

**Cross-linguistic Research Translation Studies**

**Computer-Aided Language Learning (CALL)**

---

## NoDaLiDa 2003 in Reykjavík



OPUS - an open source parallel corpus

http://logos.uio.no/opus/

Jörg Tiedemann
Department of Linguistics
Uppsala University
Box 527
SE-751 20 Uppsala, Sweden
joerg@stp.ling.uu.se

Lars Nygaard
Tekstlaboratoriet HF
University of Oslo
Postboks 1102 Blindern
0317 Oslo
lars.nygaard@ilf.uio.no

### 1 Introduction

Parallel corpora are useful in a wide variety of research areas, particularly in machine translation and lexicography. However, parallel corpora have been few, often unrepresentative, and not generally available. The aim of the OPUS project is to provide a public collection of parallel corpora which can be freely used and distributed. This makes it possible for everyone to

# Slide 1

Home / Query / WordAlign / Wiki    [books] [DGT] [DOGC] [ECB] [EMEA] [EUbooks] [EU] [Europarl] [GNOME] [GlobalVoices] [hren] [JRC] [KDE4/doc] [MBS] [MultiUN] [NCv9/v11] [OO/OO3] [subs/12/13/16]
[ParCor] [PHP] [SETIMES] [SPC] [Tatoeba] [TEP] [TedTalks] [TED] [Tanzil] [Ubuntu] [UN] [WikiSource] [Wikipedia] [WMT]

## O ᵣ PUS ... the open parallel corpus

OPUS is a growing collection of translated texts from the web. In the OPUS project we try to convert and align free online data, to add linguistic annotation, and to provide the community with a publicly available parallel corpus. OPUS is based on open source products and the corpus is also delivered as an open content package. We used several tools to compile the current collection. All pre-processing is done automatically. No manual corrections have been carried out.

The OPUS collection is growing! Check this page from time to time to see new data arriving ...
Contributions are very welcome! Please contact <jorg.tiedemann@lingfil.uu.se >

**Latest News**

- 2016-01-08: New version: OpenSubtitles2016
- 2015-10-15: New versions of TED2013, NCv9
- 2014-10-24: New: JRC-Acquis
- 2014-10-20: NCv9, TED talks, DGT, WMT
- 2014-08-21: New: Ubuntu, GNOME
- 2014-07-30: New: Translated Books
- 2014-07-27: New: DOGC, Tanzil
- 2014-05-07: Parallel coref corpus ParCor

Search & download resources: -- select --   -- select --   all

**Search & Browse**

- OPUS multilingual search interface
- Europarl v7 search interface
- Europarl v3 search interface
- OpenSubtitles search interface
- EUconst search interface
- Word Alignment Database

**Tools & Info**

- OPUS Wiki
- Tools for tagging and parsing
- Downloads (tools and models)
- Other annotation and corpus tools
- Experimental visualization tool for monolingual and parallel treebanks (demo)
- Uplug at bitbucket
- A reliable Language Identifier

**Sub-corpora (downloads & infos):**

- Books - A collection of translated literature (DOGC2014-07-17.tar.gz - 236 MB)
- DGT - A collection of EU Translation Memories provided by the JRC
- DOGC - Documents from the Catalan Goverment (DOGC2014-07-17.tar.gz - 702 MB)
- ECB - European Central Bank corpus
- EMEA - European Medicines Agency documents (EMEA0.3.tar.gz - 5.0 GB)
- The EU bookshop corpus (EUbookshop/EUbookshop0.2.tar.gz - 33 GB)
- EUconst - The European constitution (EUconst0.1.tar.gz - 67 MB)
- EUROPARL v7 - European Parliament Proceedings (Europarlv7.tar.gz - 8.4 GB)
- EUROPARL - European Parliament Proceedings (Europarl3.tar.gz - 3.6 GB)
- GNOME - GNOME localization files (GNOME2014-08-20.tar.gz - 9 GB)
- Global Voices - News stories in various languages (GlobalVoices2015.tar.gz - 1.1 GB)
- The Croatian - English WaC corpus (hrenWaC1.tar.gz - 48 MB)
- JRC-Acquis- legislative EU texts
- KDE4 - KDE4 localization files (v.2) (KDE4.tar.gz - 1.4 GB)
- KDEdoc - the KDE manual corpus ( KDEdoc.tar.gz - 35 MB)
- MBS - Belgisch Staatsblad corpus
- MultiUN - Translated UN documents

---

# Slide 2

## O ᵣ PUS

### Highlights

- over 200 language and language variants
- 29 million documents, 3.2 billion sentences, 28 billion tokens
- several domains (legislation, medical, subtitles ...)
- 12,572 language pairs, 10.8 billion translation units
- available in several formats (OPUS XML, TMX, Moses)

### Tools & online services

- tools for conversion, annotation, alignment
- online search interfaces
- word alignment database + concordances

**http://opus.nlpl.eu**

---

# Slide 3

## What Is Included?

O ᵣ PUS

- Copyright-free Books
- DGT translation memories
- DOGC (Catalan Government)
- European Central Bank
- European Medicines Agency
- EU Bookshop
- EuroParl
- GNOME, KDE, OpenOffice, PHP, Ubuntu localisation files
- Global Voices News
- Croatian-English WaC

- JRC-Acquis
- Belgisch Stadsblad
- United Nations corpora
- News Commentary, WMT sets
- OpenSubtitles
- SETIMES
- Tatoeba
- TedTalks
- Tanzil Quran Translations
- Wikipedia, WikiSource

---

# Slide 4

## How To Find Resources

link to corpus website    select source language    select target language    select size    UD parsed    word alignment    bilingual dictionaries    alternative alignments

Search & download resources:    en (English)    fr (French)    >10M

Language resources: click on [ tmx | moses | xces | lang-id ] to download the data! (raw = untokenized, ud = parsed with universal dependencies, alg = word alignments and phrase tables)

| corpus | doc's | sent's | en tokens | fr tokens | XCES/XML | raw | TMX | Moses | mono | raw | ud | alg | dic | freq | | other files |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| EUbookshop | 16947 | 10.8M | 406.8M | 431.8M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | | en-fr | | en fr | [query] [sample] | |
| MultiUN | 87480 | 14.2M | 373.8M | 454.6M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | | en-fr | | en fr | [query] [sample] | |
| OpenSubtitles2018 | 55650 | 45.2M | 363.4M | 338.0M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | | | dic | en fr | [query] [sample] | [alt] |
| OpenSubtitles2016 | 44253 | 37.3M | 299.0M | 276.0M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | | en | | | | |
| DGT | 26879 | 3.1M | 72.8M | 68.7M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | en fr | en-fr | dic | en fr | [query] [sample] | |
| Europarl | 9428 | 2.1M | 59.9M | 65.7M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | en fr | en-fr | dic | en fr | [query] [sample] | |
| JRC-Acquis | 12056 | 0.8M | 34.2M | 36.4M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | | | | en fr | [query] [sample] | |
| Wikipedia | 2 | 0.8M | 23.0M | 17.8M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | | | | en fr | [query] [sample] | |
| EMEA | 1933 | 1.1M | 12.0M | 14.8M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | | en-fr | | en fr | [query] [sample] | |
| GlobalVoices | 14501 | 0.3M | 7.0M | 7.4M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | | en-fr | | en fr | [query] [sample] | |
| ECB | 1 | 0.2M | 5.7M | 6.5M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | | en-fr | dic | en fr | [query] [sample] | |
| News-Commentary11 | 7398 | 0.2M | 6.7M | 5.2M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | | en-fr | dic | en fr | [query] [sample] | |
| GNOME | 2293 | 0.9M | 5.6M | 5.3M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | | | | en fr | [sample] | |
| News-Commentary | 1 | 0.2M | 4.7M | 5.4M | [ xces en fr ] | [ en fr ] | [ tmx ] | [ moses ] | en fr | en fr | | | | | [sample] | |
| **total** | 278822 | 117.2M | 1.7G | 1.7G | 117.2M | | 89.M | 108.M | | | | | | | | |

data in XML (tokenized)    untokenized XML    bilingual TMX    aligned plain text    monolingual plain text    word frequencies    search interface    alignment sample

# On-Line Search



# On-Line Search



# A Multilingual Word-Alignment Database

# File Structure in OPUS

On Taito: **/proj/nlpl/data/OPUS/corpus**

```
Europarl
Europarl/xml
Europarl/xml/en
Europarl/xml/en/ep-10-07-05-004.xml.gz
Europarl/xml/fr
Europarl/xml/fr/ep-10-07-05-004.xml.gz
...
Europarl/xml/en-fr.xml.gz
...
Europarl/raw
Europarl/raw/en
Europarl/raw/en/ep-10-07-05-004.xml.gz
Europarl/raw/fr
Europarl/raw/fr/ep-10-07-05-004.xml.gz
...
Europarl/parsed
```

# Internal XML Format

```
<?xml version="1.0" encoding="utf-8"?>
<document>
  <CHAPTER ID="0">
    <P id="1"></P>
    <SPEAKER ID="1" LANGUAGE="DE" NAME="Rübig">
      <P id="2">
        <s id="1">Madam President, I saw a few boats landing at
                  Parliament this week and notified the security
                  service.</s>
        <s id="2">Not only were there language difficulties; the
                  telephone line was so poor that it was almost
                  impossible to communicate.</s>
        <s id="3">I would be most obliged if the number on which
                  the security service can be reached could also
                  be clearly displayed in the House, so that if
                  anyone wants to report an incident, they can do
                  so quickly and efficiently.</s>
      </P>
      <P id="3"></P>
    </SPEAKER>
    ...
```

# Sentence Alignment

```
<?xml version="1.0" encoding="utf-8"?>
<!DOCTYPE cesAlign PUBLIC "-//CES//DTD XML cesAlign//EN" "">
<cesAlign version="1.0">

  <linkGrp targType="s" fromDoc="en/ep-00-01-17.xml.gz"
                          toDoc="fr/ep-00-01-17.xml.gz">

    <link xtargets="1;1" />
    <link xtargets="2;2" />
    <link xtargets="3;3 4" />
    ...
```

aligns sentence 3 with
sentences 3 and 4

# Sentence Alignment

```
<?xml version="1.0" encoding="utf-8"?>
<!DOCTYPE cesAlign PUBLIC "-//CES//DTD XML cesAlign//EN" "">
<cesAlign version="1.0">

 <linkGrp targType="s" toDoc="fr/2005/CES_AC71_2005_5SUMMARY.xml.gz"
                       fromDoc="en/2005/CES_AC71_2005_5SUMMARY.xml.gz">

   <link certainty="0"        xtargets=";1 2" id="SL1" />
   <link certainty="0.612088" xtargets="1;3"  id="SL2" />
   <link certainty="0.173077" xtargets="2;4"  id="SL3" />
   <link certainty="1.65065"  xtargets="3;5"  id="SL4" />
   <link certainty="1.63824"  xtargets="4;6"  id="SL5" />
   <link certainty="-0.3"     xtargets=";7"   id="SL6" />
```

alignment scores

insertion of a target language sentence

# UD-Parsed Corpora (UDPipe)

```
<?xml version="1.0" encoding="utf-8"?>

<document>
  <CHAPTER ID="002">
    <P id="1">
      <s id="1">
        <w xpos="NOUN" head="1.2" feats="Number=Plur" upos="NOUN"
lemma="document" id="1.1" deprel="nsubj">Documents</w>
        <w xpos="VERB" head="0" feats="Mood=Ind|Tense=Past|
VerbForm=Fin" upos="VERB" misc="SpaceAfter=No" lemma="receive"
id="1.2" deprel="root">received</w>
        <w xpos="PUNCT" head="1.2" upos="PUNCT" lemma=":"
id="1.3" deprel="punct">:</w>
        <w xpos="VERB" head="1.2" feats="Mood=Imp|VerbForm=Fin"
upos="VERB" lemma="see" id="1.4" deprel="parataxis">see</w>
        <w xpos="PROPN" head="1.4" feats="Number=Plur"
upos="PROPN" misc="SpaceAfter=No" lemma="Minutes" id="1.5"
deprel="obj">Minutes</w>
      </s>
    </P>
  </CHAPTER>
</document>
```

## Downloadable Packages

On Taito: **/proj/nlpl/data/OPUS/download**

```
Europarl/en.tar.gz
Europarl/en.raw.tar.gz
Europarl/en-pt.xml.gz
Europarl/en-pt.txt.zip
Europarl/en-pt.tmx.gz
```

On Abel: **/projects/nlpl/data/OPUS/download**

```
Europarl/en.raw.tar.gz
Europarl/en-pt.xml.gz
```

---

## Word Alignment and Phrase Tables

On Taito: **/proj/nlpl/data/OPUS/wordalign**

```
wordalign/Europarl/en-fr/model/aligned.grow.gz
wordalign/Europarl/en-fr/model/aligned.srctotgt.gz
wordalign/Europarl/en-fr/model/aligned.grow-diag-final-and.gz
wordalign/Europarl/en-fr/model/aligned.intersect.gz
wordalign/Europarl/en-fr/model/aligned.tgttosrc.gz
wordalign/Europarl/en-fr/model/phrase-table-filtered.gz
wordalign/Europarl/en-fr/c.true.fr.gz
wordalign/Europarl/en-fr/c.true.en.gz
wordalign/Europarl/en-fr/ids.gz
```

```
en/0/1089124/4995691.xml.gz    fr/0/1089124/4588599.xml.gz    1         1
en/0/1089124/4995691.xml.gz    fr/0/1089124/4588599.xml.gz    2 3       2
en/0/1089124/4995691.xml.gz    fr/0/1089124/4588599.xml.gz    5 6       3
en/0/1089124/4995691.xml.gz    fr/0/1089124/4588599.xml.gz    7 8 9     4
...
```

---

## OPUS Tools

On Taito: **module load nlpl-opus**

```
opus-read corpusname/lang-pair
opus-read -d corpusname lang-pair
opus-read -d corpusname -s srclang -t trglang

# print alignments with alignment certainty > LinkThr=0
opus-read -c 0 align-file.xml

# alignments with max 2 source sentences and 3 target sent's
opus-read -S 2 -T 3 align-file.xml
```

---

## Links

Main webpage: http://opus.nlpl.eu

Documentation: http://opus.nlpl.eu/trac (wiki)

Search interfaces:

- http://opus.nlpl.eu/bin/opuscqp.pl
- http://opus.nlpl.eu/cwb/Europarl7/frames-cqp.html
- http://opus.nlpl.eu/cwb/OpenSubtitles/frames-cqp.html

Word alignment dictionary: http://opus.lingfil.uu.se/lex.php

Location in NLPL: $NLPL-HOME/data/OPUS