# GUIDELINES
# FOR EPTIC COLLABORATORS

Drafted by Marta Kajzer-Wietrzny and Adriano Ferraresi (based on the previous version: EPTIC_01_2011_v3)

# Introduction

EPTIC, the European Parliament Translation and Interpreting Corpus, is an intermodal corpus created by transcribing European Parliament speeches and their interpretations, and by subsequently pairing the transcripts of interpreted speeches and their source texts with the corresponding translated versions and respective source texts. EPTIC's construction is made possible by the fact that, for each plenary session, the European Parliament published so-called "verbatim" reports of proceedings consisting of transcripts of the speeches and, until mid-2011, their translations into all EU official languages; despite being called "verbatim", these reports are often considerably edited (punctuation is added, context-related comments are removed, mistakes such as false starts, unfinished sentences or mispronunciations are corrected). The translations of the proceedings are then the result of an independently performed translation process based on the verbatim reports, without any reference to the interpreters' outputs.

An initial version of EPTIC (Bernardini et al. 2016) was based on the English <> Italian part of EPIC (Sandrelli and Bendazzoli 2005, Russo et al. 2010) and featured texts from the Parliament's part-session held in February 2004; these texts are no longer included in the current version of the corpus (Ferraresi and Bernardini 2019).

EPTIC now features texts from 2011 in 4 languages (English, French, Italian, Polish and Slovene). The corpus building procedure is a joint effort between the University of Bologna and teams from several other European universities, i.e. the University of Belgrade (Serbia), the Université Catholique de Louvain (Belgium), Adam Mickiewicz University in Poznań (Poland) and the University of Ljubljana (Slovenia), among others. Further enlargement is currently ongoing.

These guidelines are intended to help those involved in the process of transcription and coding of the texts, as well as their upload in SKEPTIC, an online platform designed to facilitate the collection and annotation procedure.

# Workflow

The following Section provides an overview of the complete workflow. Details on each step of the procedure are provided in the next Sections.

With the new SKEPTIC platform, corpus building is centred around individual speech Events. An Event coincides with a speech given by a speaker at the European Parliament. From the perspective of the teams collecting texts, an event can be conceived of as a hub which keeps together a set of inter-related files and texts: the video of the speech, the transcript of the speech in the original language and its interpretation(s), the corresponding verbatim report and its translation(s), as well as the alignment files and subtitles.

The recommended workflow within a single Event is as follows:

1. Download the multilingual video, the verbatim report and the translation of the speech from the EP website.
2. Prepare transcripts of the original speech and its interpretations in Notepad++ (on Windows computers) or BBEdit (on Mac computers), following the transcription procedures detailed below.
3. Log in to the SKEPTIC platform.
4. Add a new speaker and the relevant metadata (if the speaker is not in the list of speakers yet).
5. Add a new Event and the relevant metadata.
6. Add source texts (transcript of the original speech and corresponding verbatim report) to the Event, together with relevant metadata.
7. Add target texts (interpretations of the speech and corresponding translations) to the Event, together with relevant metadata.
8. Download sentence-split versions of the texts to be aligned from the relevant alignment section of the platform (texts are in XML format). A table with all the possible alignments is provided by the platform: you can decide whether to perform all (types of) alignment or only part of them depending on your needs.
9. Align the XML files in Intertext and generate an alignment file (the alignment file is in XML format, too).
10. Copy the content of the XML file in the relevant alignment section of the platform.
11. Download aligned corpus to use it e.g. in Sketch Engine
12. * Generate subtitles for spoken source and target texts in Aegisub (optional)
13. * Add subtitles in the platform (optional).

# Data collection

The speeches and the associated verbatim reports can be accessed via the European Parliament web page http://www.europarl.europa.eu/plenary/en/debates-video.html, which allows searches by parliamentary term and date(s) of the sittings. Clicking on the search results (bottom of page), which are organised by topic-based sections, opens a page containing the verbatim report of a given section (in HTML format), as well as links to video recordings (see Figures 1 and 2). The report is in the language in which the page is consulted: this means that for talks originally delivered in that language it constitutes a verbatim report of the original speech, while it is a translation for the others; the language of the page needs to be changed to access other language versions (a language code for the source is given in brackets if a speech is translated).

*An example of a verbatim report page with links to video recordings*

*An example of a video recording window*

The verbatim reports (sorted by date) can also be downloaded in PDF or ODT format from http://www.europarl.europa.eu/RegistreWeb/search/typedoc.htm?codeTypeDocu=PCRE. Obtaining the reports (and saving them as .txt documents, see Section 3) is a central step in the creation of the translation subcorpus components. For oral data, the reports can serve as a useful basis for a proper verbatim transcription, described in detail in Section 5.

Videos with multilingual audio tracks can be downloaded in MP4 format from this address: http://www.europarl.europa.eu/plenary/en/debate-details.html?date=20110119&detailBy=date , where the string of numbers written in red corresponds to the date of the debate (and can be changed, e.g. if one wants to collect speeches/texts from a different day/month in

which Parliamentary sessions were held). Please download videos in MP4 format and only resort to WMV when the MP4 file is unavailable or damaged.
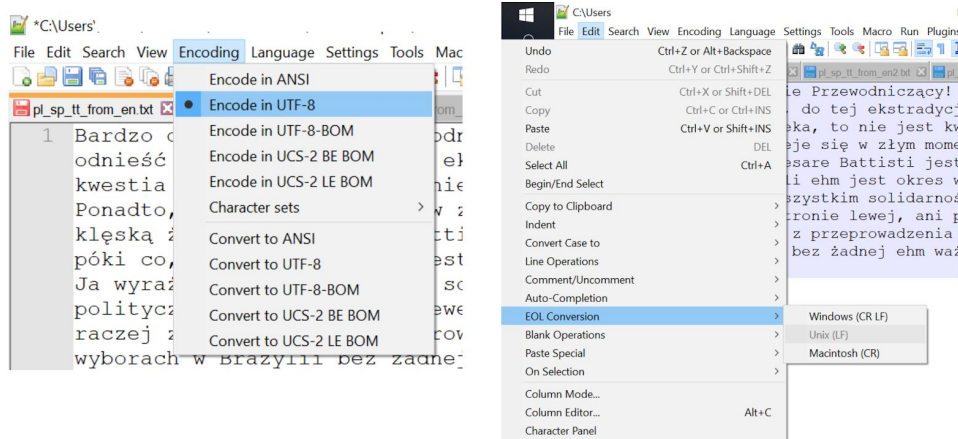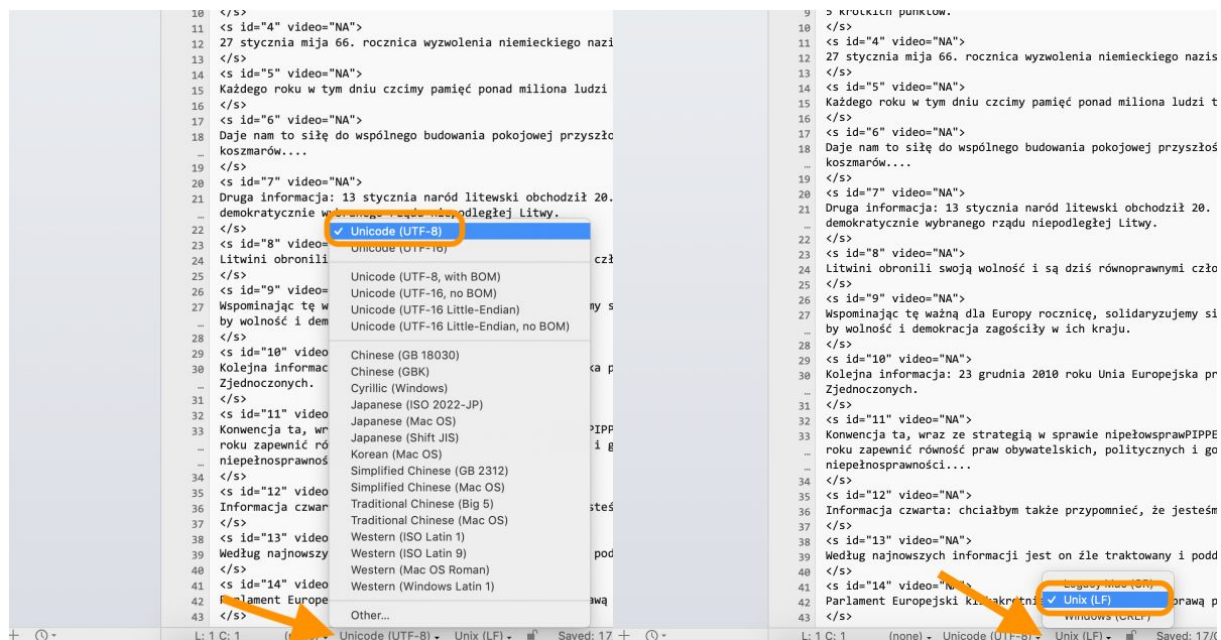


*An example of video download options.*

To make the corpus useful for linguistic research, **only texts reaching a minimum length of 60 words are to be included**. Due to being very short and repetitive, speeches by Presidents, Vice-presidents or Chairmen when moderating the debate should be excluded too.

# File format

All texts (both verbatim reports/translations and transcripts of speeches/interpretations) should be converted to text-only format, edited in Notepad++ (for Windows; http://notepad-plus-plus.org/) or BBEdit (for Mac; https://www.barebones.com/products/bbedit/ ) and saved in TXT format.

Character encoding should be set to UTF-8 (Notepad++: Menu "Encoding>UTF-8"; BBedit: check bar at the bottom of the interface) and end-of-line to UNIX (Notepad++: Menu "Edit>EOL conversion>UNIX"; BBedit: click on menus at the bottom of the interface).

# Transcription of source (original) and interpreted speeches

Both the texts belonging to the interpretation subcorpora and those making up the translation subcorpora are modified following some editing conventions. Orality traits that were not retained in the verbatim reports are reinserted in the oral subcorpora (sources and targets of interpreting) based on the video recordings. Common changes include reinserting omissions as well as truncated words, repetitions, disfluencies, pauses, etc. While many of the editing conventions described below are relevant only for the interpreting subcorpus, being closely linked to orality (see points 1-6), others apply both to the interpreting and the translation subcorpora (7-12). The applicable conventions should be followed closely.

**Punctuation**: It should follow the syntactic structure of the speech and the speaker's intonation.

**Truncated words**: If the speaker does not pronounce the entire word, this should be indicated with a dash (*propo-*). If the speaker makes a pause while pronouncing a word, the word should first be written in its correct form, and then put between slashes with an underscore where the speaker stops (*proposal /pro_posal/*).

**Disfluencies**: The word is first written in its correct form and then reported between slashes as effectively pronounced by the speaker (*proposal /preposal/*), indicating, if necessary, mistakes regarding stress (*proposal /pròposal/*) or word merging (*evitare eventuali /eventuare/*). The same convention is used when the speaker insists on one part of the word, extending one particular sound (*proposal /prooposal/*). Pronunciation markers due to non-nativeness should not be indicated between slashes, for they do not constitute a mistake but point to the origin of the speaker.

**Pauses**: There are two kinds of pauses, filled and empty. Filled pauses should be indicated with "ehm" (coughing does not count as a filled pause). A single "ehm" should always be inserted, regardless of the pause duration or repetition. A capitalized "Ehm" should be inserted at the beginning of the sentence, for example "Ehm the European Parliament already accepted this resolution". Empty pauses should be indicated with "…". Empty pauses are defined as those that last an amount of time which sounds unnatural, and that interrupt the discourse flow in an unnatural way. When "…" is used at the end of a text, it means that the speaker is interrupted and does not finish a meaningful sentence.

**Unclear bits**: Use "#" when you cannot hear a portion of the speech because of technical problems, or you do not understand what the speaker says, as well as if the speaker switches to a different language (one that is not the language of the subcorpus at issue). Only one # should be used regardless of the actual number of words that are unclear or spoken in another language. For example, if a English MEP starts a speech saying "Grazie singnor presidente" and then continues in English: "Ladies and gentlemen" it should be transcribed in the following way: # Ladies and gentlemen

**Calques and "made-up" words**: Non-conventional words which can be understood only based on the source text should be indicated in braces (e.g. when the French interpreter uses the word "manganelle" for the Italian word "manganello", {manganelle} should be written). Similar interventions should be explained in the "comments" field.

**Applause**: If the speech is interrupted by or ends with applauses, [applause] should be inserted.

**Numbers**: Numbers are written as figures, following the conventions of individual languages. Numbers in the millions and billions are not to be written as figures (*12 million, rather than 12,000,000*). Fractions are written as words. If numbers are not pronounced in their conventional way, that should be indicated between slashes (*2010 /twenty ten/*). Ordinals are written following the conventions of each language (*er/ère/ème* for French, *st/nd/rd/th* for English, *o/a* for Italian). Ordinals included in the names of important people or historical figures are written in Roman numerals (*John Paul II*). The symbol "%" is always written as a word in the relevant language.

**Character issues**: Inverted commas and apostrophes should be standardised using only those set by default in Notepad++/BBEdit (" / ').

**Capitalisation**: To decide whether to use capital letters, refer to the Interinstitutional style guides (http://bookshop.europa.eu/en/interinstitutional-style-guide-2011-pbOA3110655/),

official websites or monolingual dictionaries: for English the Oxford Learner's Dictionary (http://www.oxfordlearnersdictionaries.com/), for French the CNRTL (http://www.cnrtl.fr/definition/), for Italian Vocabolario Treccani (http://www.treccani.it/vocabolario/). To distinguish between the European Commission and various parliamentary commissions, use capital letters only for the former. The same applies to the word "commissioner" – a "Commissioner" is a member of the European Commission, a "commissioner" is a member of a parliamentary commission.

**Acronyms and names not referring to people**: Acronyms and names of programmes, institutions, regulations, action plans, etc. should be written following the conventions set by the Interinstitutional style guides (see the link above), or relevant official websites. Names of programmes, institutions, regulations, action plans, etc. are to be put in capital letters only if the speaker uses their official names.

**Titles**: Titles are written as follows (for English, French, Italian and Polish respectively): *Mr, Mrs, Ms, sir*; *madam, monsieur, madame, mademoiselle*; *signor, signore, signora, pan, pani.*

Table 1. Editing conventions

| Text Feature | Example | Editing convention |
|---|---|---|
| **Punctuation** | | *Based on syntax and intonation* |
| **Truncated words** | propo | *propo-* |
| **Mispronounciations** | Pro posal | *proposal /pro_posal/* |
| | preposal | *proposal /preposal/* |
| | pròposal | *proposal /pròposal/* |
| | prooposal | *proposal /prooposal/* |
| | eventuare | *evitare eventuali /eventuare/* |
| **Pauses** | filled | *ehm* |
| | empty | *…* |
| **Interrupted speech** | | *…* |
| **Numbers** | numbers | *in figures* |
| | 326,000,000 | *326 million/millions/milioni* |
| | 326,000,000,000 | *326 billion/milliards/miliardi* |
| | 34% | *34 percent/pour cent/per cento* |
| | years | *in figures* |

|  | fractions | *in words* |
|---|---|---|
|  | ordinals | *figures+st/nd/rd/th/er/ère/ème/o/ a or Roman numerals* |
| **Unclear bits** |  | *#* |
| **Language switching** |  | *#* |
| **Calques/Made-up words** | manganelle | *{manganelle}* |
| **Applause** |  | *[applause]* |
| **Capitalisation** |  | *Refer to Interinstitutional style guides, official websites or monolingual dictionaries* |
|  | European Commission | *Commission/Commission/Commiss ione* |
|  | parliamentary commission | *commission/commission/commissi one* |
|  | member of the European Commission | *Commissioner/Commissaire/Comm issario* |
|  | member of a parliamentary commission | *commissioner/commissaire/commi ssario* |
| **Acronyms and non-personal names** |  | *Refer to Interinstitutional style guides or official websites* |
| **Titles** |  | *Mr/Mrs/Ms/sir/madam monsieur/madame/mademoiselle signor/signore/signora/pan/pani* |

# Uploading data: Starting your work with the SKEPTIC platform

The SKEPTIC platform can be accessed at https://skeptic.dipintra.it. Only registered users can upload new data to SKEPTIC. In order to obtain logging credentials for your team of contributors please contact us by sending an e-mail to the address eptic@dipintra.it stating the number of accounts needed and the name of the team member coordinating the collection procedure.

# Icons

| | |
|---|---|
| | view |
| | edit |
| | delete |
| | add (e.g. text, event, speaker) |
| | work on event |
| | download |
| | add alignment |
| | view alignment (i.e. the alignment has been already provided) |
| | text is missing |
| | not applicable |

# How to start

Once you log in, you will see the following view.

## SKEPTIC

Welcome to SKEPTIC

Before adding new data to the corpus, check if your speaker's metadata has already been entered in the database of speakers. You can do that by clicking the **Speakers** icon on the left.

# How to add new Speakers

If the required speaker is not visible on the list, you can proceed to add the new metadata by clicking the **Add Speaker** icon in the left bottom corner.

💬 Events                                    👤 Speakers

## Speakers

| Full name | Gender | Political group | Political function | Country | | | |
|-----------|--------|-----------------|--------------------|---------|---|---|---|
| Tavares Rui | M | GUE-NGL | MEP | Portugal | 👁 | ✏ | 🗑 |
| Jerzy Buzek | M | N/A | President of the EP | Poland | 👁 | ✏ | 🗑 |

Add Speaker 📄

In the next window you need to fill in speaker's details. **All metadata in this and other sections need to be entered in English**. While entering the name of the speaker, make sure you paste exactly the same name as stated on the European Parliament website, even it the name is really long. Write "N/A" in the sections that do not apply to your speaker (e.g. a guest speaker will not belong to any political groups at the EP).

## Add Speaker

**Full Name**

[                                                                                    ]

**Gender**

○ Female  ○ Male  ○ Unspecified

**Political Group**

[                                                                                    ]

**Political Function**

[                                                                                    ]

**Country**

[                                                                                    ]

[ 🖫 Save ]    [ ☒ Cancel ]

Viewing/Editing metadata of an existing speaker. [THIS FUNCTIONALITY IS NOT IMPLEMENTED YET].

# How to add new Events

Each speech delivered at the European Parliament is catalogued as an Event and transcripts, verbatim reports, interpretations and translations in other language versions are associated with it . You can view all Events by selecting **Event** in the upper left corner of the screen.

## Events

Add Event 📄

| Speaker | Date | Topic | Specific topic | Source Language | Created | Modified | Owner | | | |
|---------|------|-------|----------------|-----------------|---------|----------|-------|---|---|---|
| Stefan Füle | 2011-01-17 | Politics | Statement by the President of the European Parliament on the situation in Tunisia | English | 2019-02-21 11:05 | 2019-02-21 11:05 | Marta K | 🔧 | 📝 | 🗑 |
| Jerzy Buzek | 2011-02-11 | Opening of the session | XXX | Polish | 2019-02-18 14:11 | 2019-02-18 14:52 | Marta K | 🔧 | 📝 | 🗑 |
| Tavares Rui | 2011-01-17 | Politics | Order of business | English | 2018-10-04 17:02 | 2019-02-20 15:39 | Adriano F | 🔧 | 📝 | 🗑 |

If you are about to add interpretations and translations to an already existing event, click on the **"Work on Event" icon** next to the **Event** of your choice. A new page will open providing details on that event: scroll down the page to add missing language versions.

If your event is not on the list, you should add it by clicking **Add Event**.

# Adding metadata of a new event

When creating a new event, you will see a page similar to the one below. What you should do here:
- Add the speaker from a drop-down list (if the speaker's name is not there you first need to add the speaker in the Speakers section);
- Specify the "Source language", i.e. the one in which the speech was originally delivered;
- Indicate if the speaker is a native speaker of the Source language (e.g. if the language is an official language used in the country s/he represents).
- Indicate the date of the speech;
- Indicate the Topic of the speech, choosing among one of the following: Agriculture and Fisheries, Economics and Finance, Environment, Health, Justice, Politics, Procedure and Formalities, Science and Technology, Society and Culture, Transport.
- Indicate the "Specific topic" of the speech. This must correspond to the title of the plenary session section in which the text belongs: these titles can be found on the European Parliament page from which the verbatim reports and/or the videos were downloaded.
- Indicate mode of delivery: impromptu (i.e. spoken without notes), mixed (i.e. spoken with some notes), read (i.e. read out in its entirety).
- After you've entered all the data click "Save".

## Add Event

**Speaker**

> Tavares Rui

**Source Language**

> English

**Speaker is native**

○ Yes  ○ No  ◉ N/A

**Date (YYYY-MM-DD)**

[Choose year... ▼] [Choose month... ▼] [Choose day... ▼]

**Topic**

**Topic Specific**

**Delivery**

○ Impromptu  ○ Mixed  ○ Read

**Notes**

[🖫 Save]  [☒ Cancel]

## Editing metadata of an existing event

In case you need to edit/correct the details of any Event, you should go to the main Events section, and click on the yellow icon next to the Event of interest.

### Events

Add Event 📄

| Speaker | Date | Topic | Specific topic | Source Language | Created | Modified | Owner | | | |
|---------|------|-------|----------------|-----------------|---------|----------|-------|---|---|---|
| Stefan Füle | 2011-01-17 | Politics | Statement by the President of the European Parliament on the situation in Tunisia | English | 2019-02-21 11:05 | 2019-02-21 11:05 | Marta K | 🔧 | 📝 | 🗑 |
| Jerzy Buzek | 2011-02-11 | Opening of the session | XXX | Polish | 2019-02-18 14:11 | 2019-02-18 14:52 | Marta K | 🔧 | 📝 | 🗑 |
| Tavares Rui | 2011-01-17 | Politics | Order of business | English | 2018-10-04 17:02 | 2019-02-20 15:39 | Adriano F | 🔧 | 📝 | 🗑 |

## Adding multilingual videos [TBA]

This function is unavailable at the moment, please download videos with multilingual audio tracks and store them in your computer. If a multilingual video is unavailable, download monolingual videos for all the languages for which you have a transcription. The name of the video  should correspond to the Event ID (which can be identified in each event's page as shown in the screenshot below) + the language code of the speech + the date + the speaker's surname, e.g. **4_EN_20110117_fule**.

### Event 4: Stefan Füle

| | |
|---|---|
| **Speaker** | Stefan Füle |
| **Source language** | English |
| **Speaker is native** | N/A |
| **Date** | 2011-01-17 |
| **Topic** | Politics |
| **Specific topic** | Statement by the President of the European Parliament on the situation in Tunisia |
| **Delivery** | read |
| **Video** | 0 |
| **Notes** | |

Event ID: 4                     Owner: Marta K

# How to add texts

Texts of transcripts, verbatim reports, interpretations and translations are associated with Events and any newly created event needs to be "filled" with texts.

You can start by adding the transcript of the source speech. To do that click on the "Add text" icon in the Spoken Source Text area.



You will be prompted to select a local file from your computer and save it.



The text is then automatically imported and segmented into sentences. Your task is to:
- Check that segmentation into sentences is correct. This is especially important since all alignments are based on this segmented version of the text
- Add the Text Duration, which should be calculated in seconds starting from the moment the speaker starts speaking to the moment when s/he stops speaking.

# Edit Text

### Polish spoken source

**Plain Text**

Mam uwagi wstępne. 5 krótkich punktów. 27 stycznia mija 66. rocznica wyzwolenia niemieckiego nazistowskiego obozu koncentracyjnego i zagłady Auschwitz-Birkenau. Każdego roku w tym dniu czcimy pamięć ponad miliona ludzi tam zamordowanych. Daje nam to siłę do wspólnego budowania pokojowej przyszłości zjednoczonej Europy, która nie będzie znała w przyszłości takich koszmarów. ... Druga informacja: 13 stycznia naród litewski

**Sentence Split Text**

```
<?xml version="1.0"?>
<xml>
 <text>
  <p id="1">
   <s id="19:1"></s>
```

Once all the information has been entered and checked:
- change text status to submitted;
- click Save.

```
  <p id="1">
   <s id="23:1">Merci monsieur le Président, ehm il s'agissait de mettre... donc ehm une d'introduire une demande de retrait d'un débat d'urgence sur l'extra
```

**Notes**

**Video**

-1

**Subtitled Text**

**Duration**

100

**Interpreter**

**Status**

○ Draft  ● Submitted  ○ Approved  ○ Published

[ 🖫 Save ]   [ ⊠ Cancel ]

Perform these operations for all the texts related to the Event that you have collected, e.g. Spoken Source Texts, Written Source Texts, Spoken Target Texts, Written Target Texts. It might be the case that  before adding the Target texts, the Target language has to be added first. To add another language version to an event scroll down the Event view to the very end and select the target language.



# How to align texts

Before proceeding to alignment, you first need to obtain the sentence-segmented XML versions of the speeches/texts. To obtain these texts:
- go to Events page and click on "Work on event";
- in the Event view click on "Toggle all alignments" or on "Alignments" under the text(s) you want to align.

- The list that opens shows all possible alignments. Click on the icon corresponding to the alignment you want to perform. E.g. the screenshot highlights the icon corresponding to the alignment between a Polish Spoken (source) text and an English Spoken (target) text. Notice that you may want to add all possible alignments or only part of them.

## Source Texts



- Download the XML-segmented files from the page that opens after you've clicked on the "Add" icon.

## Add Alignment



- Remember the order in which these files are displayed in the interface: this will have to be the same order in which you load them in Intertext.

## Aligning texts in Intertext

Download and install Intertext Editor (https://wanthalf.saga.cz/intertext)

In Intertext choose Alignment > New alignment.

Create new alignment from custom text or XML files

Now name the files according to the following conventions:

- Document name: EVENT ID (as provided in SKEPTIC platform: see "Section Adding multilingual videos");
- Version 1 name: name of the XML file you downloaded from the left-hand side of the "Add alignment" view;
- Version 2 name: name of the XML file you downloaded from the right-hand side of the "Add alignment" view.



After clicking OK you will be prompted to import the left-side text (Version 1 above). **NB:** If you have already aligned this file to a different file belonging in the same Event, you will be asked if you want to use the version of the file stored locally. Click yes and you will

automatically proceed to selecting right-side text (Version 2), even if no additional window with this message pops up.

After clicking OK you will be prompted to import the left side text (Version 2 above).  Again, if you have already aligned this file to a different file belonging in the same Event, you will be asked if you want to use the version of the file stored locally. Click yes.

Select OK and both files will be aligned automatically.



This automatic alignment needs to be corrected manually.



You can right click on the segment that needs to be relocated and chose "Move text up" or "Move text down". The result is visible below.

If two segments correspond to one, they should be merged by choosing the command "Merge segments (move both up/ bove both down)".

**Important note**: Please make sure that you merge/insert/delete/move **segments** and NOT elements (see screenshot below). You should also ensure that your edits do not affect the number of blue triangles/arrows in the left margin. These correspond to the sentences as identified by the SKEPTIC segmenter, and adding or removing segments makes the produced alignment unusable.



Before finishing your work, make sure all segments are "confirmed" (marked with a green tick).

| en_sp_tt_16 | S |
|---|---|
| ⟫ Thank you.<br>▸ I've spoken many times ... during ehm parliamentary debates on ... ehm limiting the emission of ... ehm emissions of greenhouse ehm gases that is so much supported by the European Commission and the European Parliament. | ✔ |
| ▸ But I tried to draw your attention to the fact that it has ehm significant consequences on ehm many ... industries. | ✔ |
| ▸ And this leads to ehm ... the de- relocation of many factories to countries outside the- of the European Union. | ✔ |
| ▸ Ehm these remarks ehm unfortunately, are not being taken under consideration by the European Commission and we learn that ehm yet another big factory from Germany or Poland plan to relocate {delocate} the production to the former USSR, or even to South ehm America. | ✔ |
| ▸ And I think that the European Commission should take some action, ehm because it doesn't have any vision with regard to that. | ✔ |

Once your alignment is correct, you can export it from "Alignment > Export".



File names will be determined automatically, so you do not need to specify them. You only need to select the folder where the export files will be saved.

By default, Intertext exports three files: text 1 (the filename of which corresponds to the ID of the Event + ID of Version 1), text 2 (the filename of which corresponds to the ID of the Event + ID of Version 2), and the actual alignment file (the filename of which corresponds to the ID of the Event + ID of Version 1 + ID of Version 2), as in the following example:



The only file that is needed is the XML alignment file (the one with the longest name). The other two files are just copies of the files that were imported into Intertext.

Open the XML alignment file (in Notepad++ or BBEdit) and copy its contents (CTRL/CMD + C).



Go back to the SKEPTIC platform and in the window from which you downloaded the texts, paste the alignment information in the relevant box "Alignment File".

## Add Alignment

**Polish Spoken Source Text -> English Spoken Target Text**

⬇ Download Polish Spoken Source Text      ⬇ Download English Spoken Target Text

**Alignment File**

```
<?xml version='1.0' encoding='utf-8'?>
<linkGrp toDoc='2.en_sp_tt_16.xml' fromDoc='2.pl_sp_st_15.xml'>
<link type='2-1' xtargets='16:1 16:2;15:1' status='man'/>
<link type='1-1' xtargets='16:3;15:2' status='man'/>
<link type='1-1' xtargets='16:4;15:3' status='man'/>
```

Click Save and you will be taken to the Event view, where the status of the alignment should now be green.



# How to batch download aligned texts from SKEPTIC

Log in to the SKEPTIC platform and select export corpus.



In the next step specify parameters of the texts that you want to download. To download a parallel corpus select the parameters similar to the example shown below.

**Export corpus**

There are two text fields, which allow to specify the parameters of the texts.
In the example the downloaded corpus will constitute of:
English written source texts (Text 1 on the left) and French written target texts (Text 2 on the right). There are a number of options that help further narrow down the selection of texts. Depending on your needs, you can select texts that have the draft status, those that have been submitted, have already been approved or published. Further options allow to specify the preferred text length, duration of the recording, interpreter gender, nativeness of the interpreter, interpreter nickname (if specified in the corpus) and even choose texts that contain specific words. If these are not relevant to your search select 'unspecified'.

Having specified text parameters, you can further narrow down the selection of texts for your corpus to the events of your interest.

First, specify if the downloaded corpora should be aligned. Then you can add additional details regarding the original speaker and the event at which the speech was delivered.



In the example presented above, the downloaded texts will be aligned and will include all written source texts produced by native speakers of English and delivered impromptu.
When all the details have been specified, choose "Submit" to download a parallel corpus.

In the next steps SKEPTIC will inform you how many files match your specification and will list available files.



Choose "Export Parallel Corpus" to download the aligned version of the files (in .xls format). If you do not need to investigate speech disfluencies, truncated words and filled or empty pauses in spoken texts select the relevant option. If you are not interested in parallel corpora, you can also download each monolingual corpus separately (in .xml format).

# How to upload aligned texts to SKETCH ENGINE

Go to https://www.sketchengine.eu and log in with your (institutional) credentials.
In the start page, select "My corpora"



Then create a New corpus.

Provide a name for your corpus and select "Multilingual corpus". Then upload the .xls file downloaded from SKEPTIC and click "Next".



Confirm/change corpus details and click "Next".

SETTINGS    type to search 🔍                                    Get more space ⊕

Each language in the source file will be processed into a separate monolingual
corpus and aligned with the corresponding corpus in the other language(s). Below
you can change the corpus names and/or the automatically detected languages

Corpus name (French)       EPTIC, EN-FR, French

Corpus language (French)   French                              ▼

Corpus name (English)      EPTIC, EN-FR, English

Corpus language (English)  English                             ▼

BACK    NEXT

Your corpus will be compiled:
Go

CREATE CORPUS    type to search 🔍              Get more space ⊕    🔗    ❓

COMPILATION

**EPTIC, EN-FR, French** (French)              Compiling…

**EPTIC, EN-FR, English** (English)            Compiling…

Estimated time: 100,000 words - less than a minute, 10,000,000 words - a minute or two

LEAVE

Corpus building will continue in background.

Go to Corpus Dashboard and explore your corpus.

Choose Parallel Concordance to carry out a parallel search.



*Adding subtitles [TBA]

* Downloading Subtitling Templates [TBA]

* Subtitling in Aegisub [TBA]

* Uploading Subtitled videos [TBA]

# References

Bernardini S., Ferraresi A., Miličević M. (2016). "From EPIC to EPTIC: exploring simplification in interpreting and translation from an intermodal perspective". *Target* 28, 61-86.

Ferraresi, A., Bernardini, S. (2019). "Building EPTIC: A many-sided, multi-purpose corpus of EU Parliament proceedings". In I. Doval and M. T. Sánchez Nieto (eds.) *Parallel Corpora for Contrastive and Translation Studies. New resources and applications*. Amsterdam and Philadelphia: John Benjamins. 123-139.

Sandrelli A., Bendazzoli C. (2005). "Lexical patterns in simultaneous interpreting: a preliminary investigation of EPIC (European Parliament Interpreting Corpus)". *Proceedings from the Corpus Linguistics Conference Series 1*; University of Birmingham, Birmingham.

Russo M., Bendazzoli C., Sandrelli A., Spinolo N. (2010). "The European Parliament Interpreting Corpus (EPIC): implementation and developments". Paper presented at the *international conference  Emerging Topics in Translation and Interpreting/Nuovi percorsi in traduzione e interpretazione*, 16-18 giugno 2010, SSLMIT, Università di Trieste.

RECOMMENDED SOFTWARE:
Intertext Editor: https://wanthalf.saga.cz/intertext
Notepad++: https://notepad-plus-plus.org
BBEdit: https://www.barebones.com/products/bbedit/ [Select "Free Download"]

# APPENDIX

Summary of metadata and annotation files to be submitted online

| **Event** |
| --- |
| <ul><li>speaker name = e.g. Tavares Rui [selected from a drop-down list]</li><li>source language = [language of the original speech]</li><li>speaker is native = Y/N [native = the speaker speaks in his country's official language]</li><li>date (YYYY-MM-DD)</li><li>topic = e.g. politics [selected from a drop-down list]</li><li>specific topic = e.g. order of business [copied from the EP website]</li><li>delivery = impromptu/ mixed/ read</li></ul> |

| **Speaker** | | | |
|---|---|---|---|
| ● speaker name = e.g. Tavares Rui [copied from the EP website] <br> ● gender = M/ F <br> ● country = e.g. Portugal [selected from a drop-down list] <br> ● political function = e.g MEP [selected from a drop-down list] <br> ● political group = e.g. GUE-NGL [selected from a drop-down list] | | | |
| **Spoken Sources** | **Written Sources** | **Spoken Targets** | **Written Targets** |
| ● duration= e.g. 79 [in seconds] <br> ● alignment file(s) | ● alignment file(s) | ● interpreter gender = M/F <br> ● interpreter native = Y/N <br> ● interpreter id = [optional] <br> ● duration = e.g. 67 [in seconds] <br> ● alignment file(s) | ● alignment file(s) |