

# Czech Broadcast News Corpus – Structural Metadata Annotation

## –Data Format Description–

The structural metadata (MDE) annotation for the Czech Broadcast News (BN) Corpus is stored in two different formats – QAn (Quick Annotator) and RTTM. Note that the formats slightly differ from those used for annotations from the earlier published Czech Broadcast Conversation Corpus because the BN annotations were created upon transcripts prepared based on different transcription guidelines. Character encoding in all files is ISO-8859-2.

The corpus uses the following file naming convention. All file names have the form “MMDDHHCN.xxx” where

MM is a number representing the month of the recorded broadcast (e.g., 02 for February)  
DD is a number representing the day of the recorded broadcast  
HH is a number representing the approximate beginning of the broadcast (in hours)  
C is a letter representing the radio station or the TV channel according to the following key:

- R ... CRo1 – Radiožurnál
- P ... CRo2 – Praha
- V ... CRo3 – Vltava
- C ... Ceska Televize
- M... Prima

N means “News”  
.x is either “.qan” or “.rttm”

## QAn (Quick Annotator) Format

This is the native format of the QAn ([Quick Annotator](#)) tool that was used to MDE annotate the corpus. It is based on the Transcriber’s XML (TRS, <http://trans.sourceforge.net/>) format, extended by some special tags representing structural metadata. The transcripts contain standard punctuation, but acceptable punctuation is limited to periods and question marks at the end of a sentence, and commas within a sentence. In addition to full words and punctuation, the transcripts contain the following special events:

**Table 1 Special events used in our “.qan” format**

Event	Description
<b>Speaker Noises:</b>	
<b>BREATH</b>	Audible breaths
<b>COUGH</b>	Audible cough
<b>LAUGH</b>	Laughter
<b>LIP-SMACK</b>	Lip smacks or tongue clicks
<b>SIGH</b>	Sigh

<b>Other Noises:</b>	
<b>NOISE</b>	Unspecified noise
<b>BACKGROUND_SPEECH</b>	Background or remote speech from other speakers
<b>MUSIC</b>	Music and jingles
<b>PAPER-RUSTLE</b>	Paper rustle
<b>Lexemes:</b>	
<b>ER</b>	Filled pause ‘er’
<b>UH</b>	Filled pause ‘uh’
<b>UM</b>	Filled pause ‘um’
<b>MH</b>	Interjection expressing agreement or disagreement
<b>UNINTELLIGIBLE</b>	Unintelligible word

Note that the list of special events also includes some special lexemes. This setting is used only to make them more visible in the transcript.

Word fragments and mispronounced words are marked by leading and trailing asterisks (e.g., *\*wor\** or *\*ord\** as fragments of *word*).

The transcripts also contain speaker information. For speakers whose identity was known, speaker IDs are simply their full names. For unknown speakers, we use IDs in the form “MMDDHHCN\_UGXX”, where the initial eight characters correspond to the transcript base name, U stands for Unknown. G for Gender (F or M), and XX is a two-digit number distinguishing multiple unknown speakers within a single transcript.

The MDE information is captured using the following symbols. The format uses two types of metadata tags: *SUs*, that are associated with interword boundaries, and *Labels*, that can span over one or more words. Thus, *Label* tags have the form of begin/end pairs, while *SUs* are only single tags.

The *SU* tags always start with “<mde:SU”. Further, the tags have the mandatory attribute “type”. The following types are used:

(a) *SU-external symbols:*

- "/." – Statement break *without* strong prosodic marking at boundary
- "//." – Statement break *with* strong prosodic marking at boundary
- "/?" – Question break *without* strong prosodic marking at boundary
- "//?" – Question break *with* strong prosodic marking at boundary
- “/-” – End of an incomplete (arbitrarily abandoned) *SU*
- “/~” – End of an incomplete *SU* interrupted by another speaker

(b) *SU-internal symbols:*

- "/", – Clausal break
- "/&"; – Coordination break

(c) *Interruption point symbol*<sup>1</sup>:

- " \* " – Interruption point within an edit disfluency (asterisk)

In addition to the mandatory attribute, the SU tags may also contain the optional attribute “previous”. This attribute indicates that the SU tag replaced a standard punctuation symbol (such as period or comma). This information is especially important for the annotation tool. If the annotator decides to delete an SU tag, the tool can display the original punctuation symbol again. For example, a tag with the “previous” attribute may be `<mde:SU type="/." prev=","/>`.

Furthermore, interruption point tags may also receive the attribute “auto”. If the tag looks like `<mde:SU type="*" auto="1"/>`, it indicates that the Interruption point tag was inserted automatically by the annotation tool at the right edge of the preceding Delreg.<sup>2</sup>

### Labels

Label-type tags start with “<mde:Label”. All Label tags have two mandatory attributes – “type” and “extent”. The attribute extent may have two values – “begin” and “end”, to indicate tag pairs. The attributes “type” may have the following values:

- “A/P” – Aside/Paraphrased
- “Backchannel” – Backchannel uttered by other speaker than
- “Correction” – Correction of previous Delreg
- “DM” – Discourse marker
- “DR” – Discourse marker of subtype “Discourse response”
- “Delreg” – Deletable region
- “EET” – Explicit editing term
- “FP” – Filled pause

An example of a Label tag is “<mde:Label type="DM" extent="begin" />”. The following text serves as an example of a metadata annotated speech transcript in the “.qan” format:

```
<Event desc="BREATH" type="noise" extent="instantaneous"/> to
<mde:Label type="Delreg" extent="begin" /> bylo <mde:Label type="Delreg" extent="end" />
<mde:SU type="*" auto="1"/> <mde:Label type="Correction" extent="begin" /> bylo
<mde:Label type="Correction" extent="end" /> velkou zkouškou vládnoucí strany <mde:SU
type="/," prev=","/> protože
<Event desc="BREATH" type="noise" extent="instantaneous"/>
prakticky od roku devadesát šest v České republice neexistuje většinová vláda
<mde:SU type="//." prev=","/>
```

<sup>1</sup> Note that Interruption Points are included with SU tags only for the sake of format simplicity because same as SU boundary symbols, they are associated with interword boundaries. Strictly speaking, they should be in a separate group, but we did not want to introduce another group that would only include a single tag category.

<sup>2</sup> Note that unlike English, Czech MDE does not use automatic interruption points before fillers.

## RTTM Format

The RTTM format also provides information about structural metadata that enrich standard speech transcripts. The format described herein is based on the RTTM-format-v13 used for MDE in the EARS project. The original RTTM format could not be used in the exact form employed in the EARS project because annotation modifications<sup>3</sup> introduced in the Czech MDE annotation project had to be reflected. Note that the published RTTM files only contain description of those regions of data that have MDE annotation (i.e., sections with overlapping speech are not present in RTTMs).

The format uses object-oriented representation of the rich text data. There are four general object categories to be represented. They are STT objects, MDE objects, source (speaker) objects, and structural objects. Each of these general categories may be represented by one or more types and subtypes, as shown in Table 1. Note that the object subtypes that are generally allowed but do not appear in this corpus are marked with asterisks.

**Table 2 Rich Text object types and subtypes**

Type	Subtypes
<b>Structural types:</b>	
<b>SEGMENT</b>	<NA>
<b>STT types:</b>	
<b>LEXEME</b>	<b>lex, fp, frag, interjection, un-lex<sup>4</sup>, and other*</b>
<b>NON-LEX</b>	<b>laugh, breath, lip-smack, cough, sigh, and other*</b>
<b>NON-SPEECH</b>	<b>noise, music, background_speech, paper-rustle and other*</b>
<b>MDE types:</b>	
<b>FILLER</b>	<b>discourse_marker, discourse_response<sup>5</sup>, explicit_editing_term, backchannel*, and other*</b>
<b>EDIT</b>	<NA>
<b>CORRECTION</b>	<NA>
<b>IP</b>	<b>edit, filler*, edit&amp;filler*, and other*</b>
<b>SU</b>	<b>/, //, /?, //?, /~, /-<sup>6</sup></b>
<b>CB</b>	<b>coordinating, clausal, and other*</b>
<b>A/P</b>	(none)

<sup>3</sup> The modifications are described in the document “Structural Metadata Annotation for Czech: An Overview” that is also included in the corpus documentation.

<sup>4</sup> Un-lex is used to tag unintelligible words.

<sup>5</sup> By definition, discourse\_response is a subtype of discourse\_marker. They are listed on the same level herein only because the RTTM format does not allow to define “subsubtypes”.

<sup>6</sup> Since there are more SU subtypes in Czech MDE than in the original standard, we rather use a symbolic instead of word representation of SU subtypes for the sake of simplicity.

<b>SPEAKER</b>	(none)
<b>Source information:</b>	
<b>SPKR-INFO</b>	<b>adult_male, adult_female, child*, and unknown*</b>

Except for the static speaker information object [**SPKR-INFO**], each object exhibits a temporal extent with a beginning time and duration. (The duration of interruption points [**IP**] and clausal boundaries [**CB**] is zero by definition.)

These objects are represented individually, one object per record, using a flat record format with object attributes stored in white-space separated fields. The format is shown in Table 2.

**Table 3 Object record format for RTTM objects**

<b>Field 1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>
type	file	chnl	Tbeg	tdur	ortho	stype	name	conf

where

*file* is the waveform file base name (i.e., without path names or extensions).

*chnl* is the waveform channel (e.g., “1” or “2”).

*tbeg* is the beginning time of the object, in seconds, measured from the start time of the file.<sup>7</sup>

If there is no beginning time, use *tbeg* = “<NA>”.

*tdur* is the duration of the object, in seconds.<sup>4</sup> If there is no duration, use *tdur* = “<NA>”.

*stype* is the subtype of the object. If there is no subtype, use *stype* = “<NA>”.

*ortho* is the orthographic rendering (spelling) of the object for STT object types. If there is no orthographic representation, use *ortho* = “<NA>”.

*name* is the name of the speaker. *name* must uniquely specify the speaker within the scope of the *file*. If *name* is not applicable or if no claim is being made as to the identity of the speaker, use *name* = “<NA>”.

*conf* is the confidence (probability) that the object information is correct. If *conf* is not available, use *conf* = “<NA>”.

This format, when specialized for the various object types, results in the different field patterns shown in table 3.

**Table 4 Format specialization for specific object types**

<b>Field 1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>
<i>Type</i>	<i>file</i>	<i>chnl</i>	<i>tbeg</i>	<i>tdur</i>	<i>ortho</i>	<i>stype</i>	<i>name</i>	<i>conf</i>
<b>SEGMENT</b>	file	chnl	tbeg	tdur	<NA>	eval or <NA>	name or <NA>	conf or <NA>
<b>LEXEME NON-LEX</b>	file	chnl	tbeg	tdur	ortho or <NA>	stype	name	conf or <NA>

<sup>7</sup> If *tbeg* and *tdur* are “fake” times that serve only to synchronize events in time and that do not represent actual times, then these times are tagged with a trailing asterisk (e.g., *tbeg* = **12.34\*** rather than **12.34**).

<b>NON-SPEECH</b>	file	chnl	tbeg	tdur	<NA>	stype	<NA>	conf or <NA>
<b>FILLER EDIT CORRECTION SU</b>	file	chnl	tbeg	tdur	<NA>	stype	name	conf or <NA>
<b>IP CB</b>	file	chnl	tbeg	<NA>	<NA>	stype	name	conf or <NA>
<b>A/P SPEAKER</b>	file	chnl	tbeg	tdur	<NA>	<NA>	name	conf or <NA>
<b>SPKR-INFO</b>	file	chnl	<NA>	<NA>	<NA>	stype	name	conf or <NA>

The following table shows mapping between the QAn and the RTTM format for the events that are defined in both formats.

**Table 5 Mapping between QAn and RTTM annotation**

<b>QAn Type</b>	<b>RTTM Type</b>	<b>RTTM Subtype</b>
<b>/.</b>	<b>SU</b>	<b>/.</b>
<b>//.</b>	<b>SU</b>	<b>//.</b>
<b>/?</b>	<b>SU</b>	<b>/?</b>
<b>//?</b>	<b>SU</b>	<b>//?</b>
<b>/~</b>	<b>SU</b>	<b>/~</b>
<b>/-</b>	<b>SU</b>	<b>/-</b>
<b>/,</b>	<b>CB</b>	<b>clausal</b>
<b>/&amp;</b>	<b>CB</b>	<b>coordinating</b>
<b>*</b>	<b>IP</b>	<b>edit</b>
<b>A/P</b>	<b>A/P</b>	<b>&lt;NA&gt;</b>
<b>Backchannel</b>	<b>FILLER</b>	<b>backchannel</b>
<b>Correction</b>	<b>CORRECTION</b>	<b>&lt;NA&gt;</b>
<b>DM</b>	<b>FILLER</b>	<b>discourse_marker</b>
<b>DR</b>	<b>FILLER</b>	<b>discourse_response</b>
<b>Delreg</b>	<b>EDIT</b>	<b>&lt;NA&gt;</b>
<b>EET</b>	<b>FILLER</b>	<b>explicit_editing_term</b>
<b>FP</b>	<b>LEXEME</b>	<b>fp</b>