

## An object-oriented representation of speech for EARS:

There are four general object categories to be represented. They are STT objects, MDE objects, source (speaker) objects, and structural objects.<sup>1</sup> Each of these general categories may be represented by one or more types and subtypes, as shown in table 1.

**Table 1 Rich Text object types and subtypes**

Type	Subtypes
<b>Structural types:</b>	
<b>SEGMENT</b>	<b>eval</b> , or (none)
<b>NO_SCORE</b>	(none)
<b>NO_RT_METADATA</b>	(none)
<b>STT types:</b>	
<b>LEXEME</b>	<b>lex</b> , <b>fp</b> , <b>frag</b> , <b>un-lex</b> <sup>2</sup> , <b>for-lex</b> , <b>alpha</b> <sup>3</sup> , <b>acronym</b> <sup>3</sup> , <b>interjection</b> <sup>3</sup> , <b>propername</b> <sup>3</sup> , and <b>other</b>
<b>NON-LEX</b>	<b>laugh</b> , <b>breath</b> , <b>lip-smack</b> , <b>cough</b> , <b>sneeze</b> , and <b>other</b>
<b>NON-SPEECH</b>	<b>noise</b> , <b>music</b> , and <b>other</b>
<b>MDE types:</b>	
<b>FILLER</b>	<b>filled_pause</b> , <b>discourse_marker</b> , <b>explicit_editing_term</b> , and <b>other</b>
<b>EDIT</b>	<b>repetition</b> , <b>restart</b> , <b>revision</b> , <b>simple</b> , <b>complex</b> , and <b>other</b>
<b>IP</b>	<b>edit</b> , <b>filler</b> , <b>edit&amp;filler</b> , and <b>other</b>
<b>SU</b>	<b>statement</b> , <b>backchannel</b> , <b>question</b> , <b>incomplete</b> , <b>unannotated</b> , and <b>other</b>
<b>CB</b>	<b>coordinating</b> , <b>clausal</b> , and <b>other</b>
<b>A/P</b>	(none)
<b>SPEAKER</b>	(none)
<b>Source information:</b>	
<b>SPKR-INFO</b>	<b>adult_male</b> , <b>adult_female</b> , <b>child</b> , and <b>unknown</b>

The STT, MDE and Source information objects are potential research target. And, except for the static speaker information object [**SPKR-INFO**], each object exhibits a temporal extent with a beginning time and a duration. (The duration of interruption points [**IP**] and clausal boundaries [**CB**] is zero by definition.)

<sup>1</sup> Structural objects are important because they are produced by LDC to provide a modicum of temporal organization in the annotation and identify non-evaluable regions.

<sup>2</sup> Un-lex is also used to tag words that are infected with or affected by laughter.

<sup>3</sup> This subtype is an optional addition to the previous set of lexeme subtypes which is provided to supplement the interpretation of some lexemes.

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 2 Object record format for EARS objects**

Field 1	2	3	4	5	6	7	8	9
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>4</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 3 Format specialization for specific object types**

Field 1	2	3	4	5	6	7	8	9
<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>NOSCORE</b>	file	chnl	tbeg	tdur	<NA>	<NA>	<NA>	<NA>
<b>NO_RT_METADATA</b>	file	chnl	tbeg	tdur	<NA>	<NA>	<NA>	<NA>
<b>LEXEME NON-LEX</b>	file	chnl	tbeg	tdur	ortho or <NA>	stype	name	conf or <NA>
<b>NON-SPEECH</b>	file	chnl	tbeg	tdur	<NA>	stype	<NA>	conf or <NA>
<b>FILLER EDIT SU</b>	file	chnl	tbeg	tdur	<NA>	stype	name	conf or <NA>
<b>IP</b>	file	chnl	tbeg	<NA>	<NA>	stype	name	conf

<sup>4</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 should be tagged with a trailing asterisk (e.g., *tbeg* = **12.34\*** rather than **12.34**).

<b>CB</b>								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>

## Transforming EARS objects to a sequence of events:

For purposes of training recognition algorithms it may be desirable to transform the object data into a chronological sequence of object boundary events. This may be done by representing separately the beginning and end of every object and then sorting these event records into chronological order. The record format for these temporally sequenced records is identical to the record format for the objects themselves, except for three additional fields used to identify the type of boundary, the object that the boundary derives from, and the boundary time. The format is shown in table 4.

**Table 4 Event record format for temporally sequenced events**

Field 1	2	3	{Fields 4...12}
eventType	objectID	eventTime	{object fields 1...9}

where

*eventType* is the type of event being represented, namely *beg*, *end*, and *obj* (where *obj* is an object that is represented by a single point in time). For speaker information, which is not a function of time, use *eventType* = "<NA>".

*objectID* is an (arbitrary) ID that is unique to the object that this event derives from. For speaker information, which is not a function of time, use *eventType* = "<NA>".

*eventTime* is the time of the event, in seconds, measured from the start time of the file.<sup>5</sup> For speaker information, which is not a function of time, use *eventType* = "<NA>".

This format, when specialized for the various event types, results in the different field patterns shown in table 5.

**Table 5 Format specialization for specific event types**

Field 1	2	3	4	{Fields 5...12}
<i>eventType</i>	<i>objectID</i>	<i>eventTime</i>	<i>objectType</i>	{object fields 2...9}
<b>beg end</b>	objectID	eventTime	SEGMENT LEXEME NON-LEX NON-SPEECH FILLER EDIT SU A/P SPEAKER	{object fields 2...9}
<b>obj</b>	objectID	eventTime	IP CB	{object fields 2...9}

<sup>5</sup> If *eventTime* is a "fake" time that serves only to synchronize events in time and that does not represent an actual time, then this time should be tagged with a trailing asterisk (e.g., *eventTime* = **12.34\*** rather than **12.34**).

<NA>	<NA>	<NA>	<b>SPKR-INFO</b>	{object fields 2...9}
------	------	------	------------------	-----------------------