

# Coding Instructions for Topic Segmentation of the ICSI Meeting Corpus

(version 1.2)

Weiqun Xu\*, Jean Carletta and Jonathan Kilgour

School of Informatics, University of Edinburgh  
Email: {wxu, jeanc, jonathan}@inf.ed.ac.uk

June 14, 2005

## ABSTRACT

This paper is a manual that instructs annotator how to work on topic segmentation in the ICSI meeting corpus. After introducing some general ideas of the task and the tool, it explains how to do the job step by step. Some advice is further given for external users.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>The Task</b>	<b>2</b>
2.1	Segmentation . . . . .	2
2.1.1	Sub-topics . . . . .	2
2.1.2	Hints for finding segment boundaries . . . . .	2
2.2	Topic Description . . . . .	3
2.2.1	Standard “Topic” Descriptions . . . . .	3
<b>3</b>	<b>The coding tool: ICSI topic segmenter</b>	<b>4</b>

## 1 Introduction

The ICSI Meeting Corpus is a collection of recordings of meetings from university research groups at ICSI, in California. Many people study these meetings because one of the big technology challenges at the moment is to build a meeting browser — that is, something that can be used to find out what happened at a meeting. A big part of meeting browsing is knowing what the people in a meeting were talking about — the topic — and when they changed topics. It’s easier to make a machine understand topics and topic changes if someone tells

---

\*Contact person. All kinds of comment and feedback are warmly welcome.

it about topics and topic changes on some examples. Your job is to listen to some meeting recordings, divide the meetings up by topic into “segments”, and briefly describe what the topic is for each one. We expect this to take 4-5 hours for new coders (but 2-3 hours for experienced coders) for every hour of recorded meeting.

## **2 The Task**

### **2.1 Segmentation**

You might be expecting us to tell you exactly how many topic segments each meeting should have, but the truth is, it varies considerably. We have in mind that a typical one-hour meeting might have something in the order of six to ten segments, but there could be meetings that discuss one topic extensively, and others that handle a very large number of topics, all briefly. For this reason, you should divide the meeting into segments in the way that you find most natural.

Everything that is said during a meeting should end up in some segment, but since some material in meetings isn’t really about a topic we provide some standard kinds of segments to make this easier.

#### **2.1.1 Sub-topics**

Sometimes you won’t be sure whether to mark part of a meeting as one segment or two, because there are really two segments but they are related to each other. For instance, if a group were talking about what they liked about Edinburgh, they might talk first about the free museums and then they might talk about the architecture. If they talked about these two things completely separately, without saying anything about why the two go together, then they would be separate topics. However, if they made clear that there were some connection, for instance, by saying that they were talking about what they liked about Edinburgh and then introducing these themes, the overall segment would be about “good things about Edinburgh”, with the sub-segments about “free museums” and “architecture”. You can mark sub-segments wherever you like, to cover part of the material in a segment.

If you feel that there is a clear subtopic happening, then mark it and describe it. Otherwise, don’t bother to subdivide. In theory every time someone talks they’re saying something different from the last person and therefore it should be possible to mark a new subtopic, but we don’t want this level of detail. A sub-segment should be something that the group is discussing, not just something one person threw into the discussion. We would expect some sub-topics in the meeting corpus, and possibly (but rarely) some sub-sub-topics, but if you find highly nested structures (with lots of detailed sub-sub-topics), you should consider whether you might be subdividing topics too finely.

#### **2.1.2 Hints for finding segment boundaries**

There are some clues that should help you find segment boundaries.

The first is that people in meetings quite often announce topic changes. For example,

*Ok, so the s the the next thing we had on the agenda was something about alignment [BMR019-740]*

You should be careful that the meeting actually moves on to the next topic at this point — sometimes someone will intervene with more material on the previous topic — but otherwise such utterances are pretty good indicators of a boundary. Note that as well as signalling a topic shift, these utterances often give some clear idea about the topic which you can put into topic description.

The second clue is that if the group discusses the agenda at the beginning, it can be useful to look for where the agenda items appear, even if groups don't always follow the agenda that they set. Again, the agenda can be a useful source of topic descriptions.

The last set of clues are words like “anyway” and “so”, which can be used to indicate a topic shift, like in the example above. When you listen to the recordings, these indicators can sound quite distinctive from other ways of using the same words.

## 2.2 Topic Description

As well as saying where the discussion of a topic starts and ends, you need to give a short English description of the topic. This can usually be based on a few keywords from the discussion, and needs to be detailed enough that someone could figure out later on what the topic was, but not so detailed that writing the descriptions is a major part of the work. We have in mind short phrases for these, of perhaps six or seven words. Groups will quite often discuss a topic and then return to it later in the same meeting. When this happens, you should use exactly the same description. The software lets you do this without typing it in again. We'll know which topic a subtopic goes with, so it's OK if it is necessary to read both the topic and subtopic label to understand what the subtopic is (e.g., “why we like Doris” could have “her hairstyle” for a subtopic rather than “hairstyle as a reason for liking Doris”).

Sometimes you won't be sure what words to use to describe a topic, especially if there are clear keywords in the text. It's better to provide some description than none, even if you have to just describe the segment in terms of how it contributes to the meeting (such as “clarification of technical issue affecting recognizer performance”).

### 2.2.1 Standard “Topic” Descriptions

We provide five special, standardized topic descriptions to make coding easier and to ensure some consistency across the coded meetings: “digit task” for a peculiar task the participants do involving reading digits aloud; “opening” and “closing” for the special kinds of things that happen at the beginning and ends of meetings; “chitchat” for social chatter; and “complex” for material that is just too difficult to code.

#### **Digit task**

Because the researchers involved in the meetings are themselves studying automatic meeting processing, they sometimes do things during their meetings particular to aid their research. Near the beginning or end of most of the meetings, they perform a task that involves reading a series of digits. This is to

help their speech recognizer, and isn't really a part of the meeting business. If they do this task during the meeting you are coding, place these utterances in a segment of their own labelled "digit task".

### **Opening**

The "opening" topic description can be used for all of the little things that groups tend to do at the beginning of a meeting: take attendance, review or set the agenda, say how long the meeting is, and so on. It's useful for us to know where the opening begins and ends, but the opening doesn't contain particular topics, or at least not technical ones, just discussion about how to conduct the meeting. We don't need opening material to be further segmented, except that if the group performs the digit task during the meeting opening, the digit task should be in its own topic segment or sub-segment. It doesn't matter whether you place the digit task within the opening as a sub-topic or next to it, so don't fuss about which structure to use.

### **Closing**

Just as meetings can have openings that aren't about a real topic, they can also be closed in a similar way, for instance, with time dedicated to setting the next meeting date and reviewing who will do what in preparation for the next meeting (or deciding this if it hasn't been done during the meeting itself). Mark these as "closing", and again, do not segment further, except for the digit task.

### **Chitchat**

Sometimes during a meeting the participants just chat aimlessly, usually about social matters. This especially happens after the microphones have been switched on but before the beginning of the meeting "proper", and again at the end. It can also happen in the middle of a meeting, for instance, when a projector breaks, or simply if someone drags the group off-topic. When this happens, divide the meeting so that these areas form their own segments, but label them with the special topic description, "chitchat".

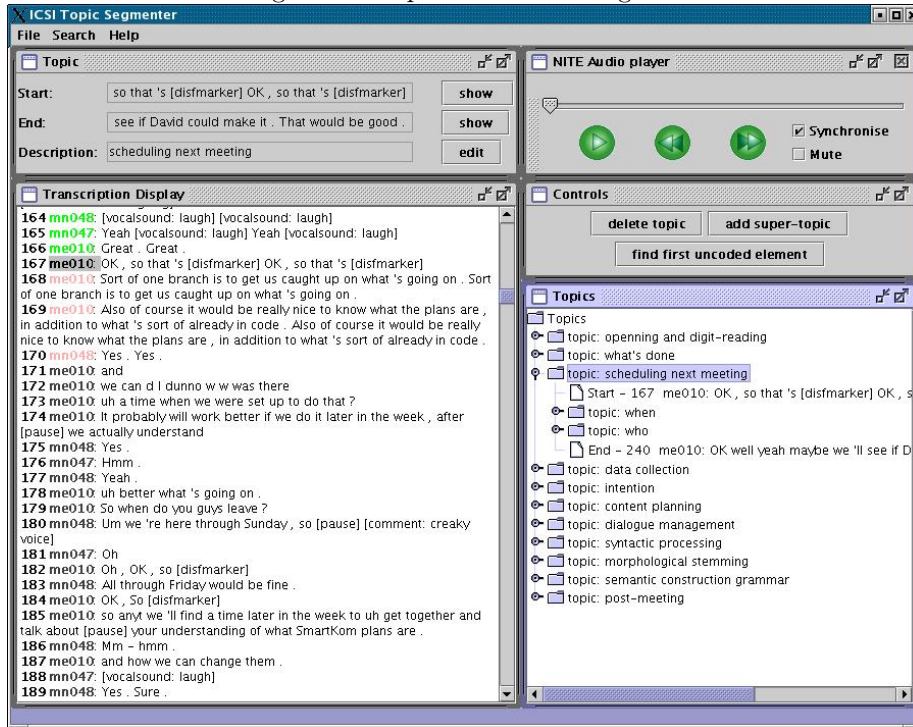
It can also happen that while the group is having a proper discussion of some topic, there will be one or two quick utterances, like jokes, that you might be tempted to code as "chitchat", but the group doesn't really get pulled off the topic they are discussing. Don't bother to segment around these cases — just leave the utterances within the wider segment.

### **Complex**

It can be difficult to understand what is happening in these meetings for several reasons. The subject matter is pretty esoteric, and during informal meetings among people who know each other well, they tend to raise whatever issues they want to discuss whenever they spot an opportunity, rather than rigidly following a set agenda. Also, sometimes topics are so intertwined with each other it can be difficult to tell where one ends and another begins, especially if both need to be discussed before some decision can be made.

As long as you can find the boundaries between topics and pick out some keywords to put in the brief topic description, what you do will be OK. Segmenting the topics but failing to describe them is better than nothing. If you really can't make enough sense of some part of the meeting to segment it, then put that part in the special topic description, "complex".

Figure 1: Snapshot of the coding tool



### 3 The coding tool: ICSI topic segmenter

We've written software specifically for this task that will allow you to view a transcription, play the meeting recording, and segment the meeting and add the topic annotation. (The software comes with on-line help, so you can look at that as well as reading this document.) The coding tool will work on uncoded, coded, or partly coded meetings, so you can stop and restart at any time (but remember to save your work!) or just review meetings you coded earlier.

Opening a meeting usually results in five windows on a common desktop (see figure 1):

**TOPIC:** which shows some information about the current topic segment, including its START and END and its topic DESCRIPTION. You may quickly move to the corresponding utterances by clicking the SHOW button. You can also change the description by clicking the EDIT button. When you click the EDIT button, a DESCRIBE TOPIC window will pop up, in which you can either *Choose an existing topic* (predefined or previously coded) from a drop down list or *add a new one* by inputting some free description.

**NITE AUDIO PLAYER:** which plays the audio of the opened meeting. <sup>1</sup>

There are three buttons, two check boxes, and one progress bar.

<sup>1</sup>If there are problems with the sound card or the coding tool can not locate the corresponding audio file, the player may not work properly or even not appear. Contact your coordinator if this happens.

**Progress bar** – which shows audio play progress. You may also slide forward or backward to your desired part.

**Buttons** ▷ = play or || = pause, ▷▷ = fast forward, ◀◀ = fast rewind.

#### **Check boxes**

**Synchronise** when checked (default), audio is synchronised with text highlighted in green. This is very helpful when you browse the meeting for the first time, but might become annoying when you just want to skim the text later. This feature can be disabled by unchecking the box.

**Mute** turn off the sound while it's playing.

**TRANSCRIPTION DISPLAY:** which displays the transcription of the opened meeting. Every utterance is preceded by an automatically generated line number and a speaker ID label. The transcription window divides up the meeting by who said what. Beyond that, the way it divides something one person said into numbered lines is fairly meaningless. Most of the time, you will be able to say that an entire utterance belongs to one topic segment or another, placing the boundary between utterances.

If you want to add a topic segment that uses complete numbered lines, left click on the speaker label of the first utterance in the segment and then right click on the speaker label of the last utterance in the segment, which has to be after the first segment and not already in a topic segment. When there is a valid right click, the interface will pop up a window for you to describe its topic. You can do that either by choosing one from the drop-down list or by entering a new description. You can choose to split a numbered line (if you have to) by right clicking on a word instead of a speaker ID to end the topic segment. In the current version of the tool you can't start a topic by left clicking on a word, so you'll have to code the topic segment that covers the first half of the line before you code the one that covers the second half.

If you want to start play the audio from an arbitrary line, select the line, then press **Ctrl** and right click at the same time.

**CONTROLS:** which provides additional ways of changing the topic coding and moving around the transcript.

**Delete Topic:** which deletes the topic (i.e., a segment and its topic description) after you click a topic in the TOPICS window. (Note: any sub-structure will not be deleted, instead all the children of the deleted topic will be upgraded one level.) The interface is slow to adjust after you delete a topic, so be patient.

**Add Super-Topic:** which adds a super/parent topic. First, select multiple consecutive topics in the TOPICS window. (You can do this by selecting the first topic in the series with the left mouse button, holding down the *shift* key, and selecting the last topic in the series.) Then click this button to make the selected topics sub-topics of some new parent topic. Of course, you need to give a description of the super-topic.

**Find First Uncoded Element:** which helps you quickly jump to the first uncoded utterance.

**TOPICS:** which displays all the coded topics in the meeting. A topic node is represented by its description with its start and end utterances and possible sub-topic(s). You can do nothing in this window but select one or more topic nodes.

[TIPS FOR CODING] When you start on a new meeting, we suggest that you begin by listening, using the transcription to skip forward once you have a sense of what's going on. If you're in a room with other people please use headphones. You can code any part of the meeting at any time, but it's best to work in an orderly fashion, from beginning to end. If subtopics feature heavily in the meeting you're coding, you may find it easier to segment the entire meeting and then return to add the subtopics afterwards. You should keep a list of the meetings that you are finished with, as opposed to those that you've partially done, or those that you've segmented but wish to check — the software won't help you do this. If you feel it helpful to take some notes during coding, feel free to use a piece of scratch paper or two. Also, remember to keep track of the hours you spend.