# Unleashing the killer corpus: experiences in creating the multi-everything AMI Meeting Corpus
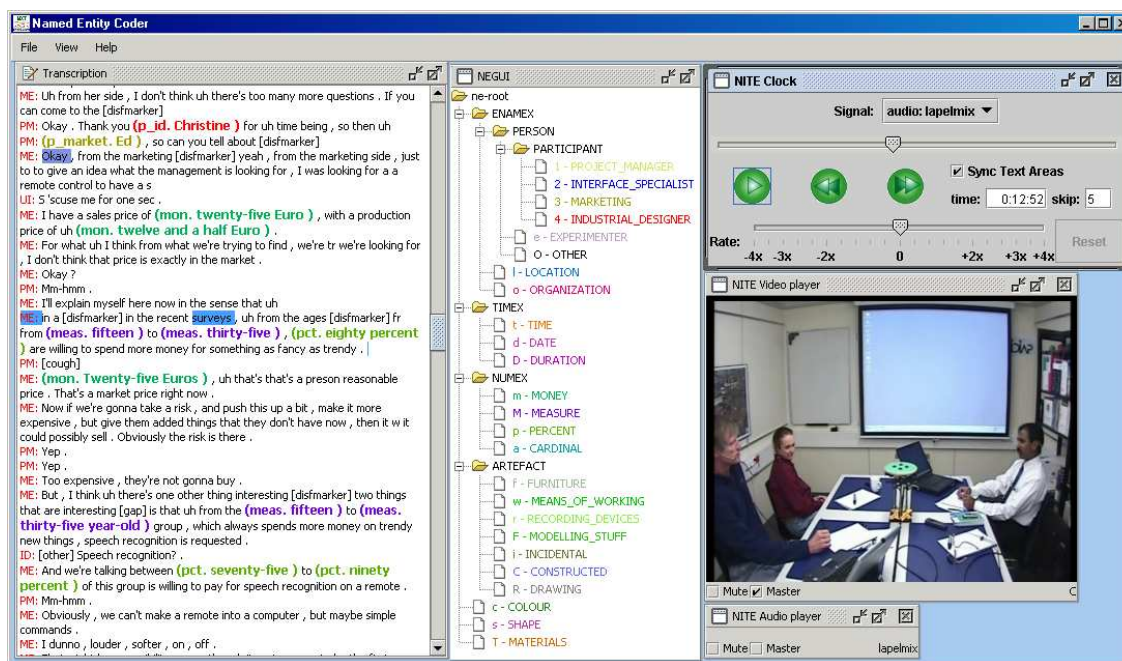
## Jean Carletta

University of Edinburgh
2 Buccleuch Place
J.Carletta@ed.ac.uk

### Abstract

Creating the AMI Meeting Corpus was an ambitious endeavour, probably more ambitious than the people who first thought of it realize even now. It contains 100 hours of meetings captured using a whole host of synchronized recording devices, and is designed to support work in speech and video processing, language engineering, corpus linguistics, and organizational psychology. It has been transcribed orthographically, with annotated subsets for everything from named entities, dialogue acts, and summaries to simple gaze and head movement. In this keynote speech, I describe the data and talk about what it was like to manage the process of creating it, distributed over six sites. If this is "killer" data, that presupposes a platform that it will "sell"; in this case, that is the NITE XML Toolkit, which is the only viable way to create, store, and browse annotations for the same base data, at least where creation is on this large a scale or there are structural relationships among the annotations. I'll end with details of the data's imminent public release, but on the way I'll take in what happens when you mix engineers and psychologists, sing the praises of Wikis, and argue that this data has really been in the making since the 1950's.

Three rooms, 100 hours of meeting data.



NXT's named entity coder — one of the many tools used to annotate and view AMI data.