

# Multimodal Signal Processing

Human Interactions in Meetings

Edited by Steve Renals, Hervé Bourlard,  
Jean Carletta, and Andrei Popescu-Belis

CAMBRIDGE

CAMBRIDGE

more information - [www.cambridge.org/9781107022294](http://www.cambridge.org/9781107022294)



## Multimodal Signal Processing

### Human Interactions in Meetings

Bringing together experts in multimodal signal processing, this book provides a detailed introduction to the area, with a focus on the analysis, recognition, and interpretation of human communication. The technology described has powerful applications. For instance, automatic analysis of the outputs of cameras and microphones in a meeting can make sense of what is happening – who spoke, what they said, whether there was an active discussion, and who was dominant in it. These analyses are layered to move from basic interpretations of the signals to richer semantic information. The book covers the necessary analyses in a tutorial manner, going from basic ideas to recent research results. It includes chapters on advanced speech processing and computer vision technologies, language understanding, interaction modeling, and abstraction, as well as meeting support technology. This guide connects fundamental research with a wide range of prototype applications to support and analyze group interactions in meetings.

**Steve Renals** is Director of the Institute for Language, Computation, and Cognition (ILCC) and Professor of Speech Technology in the School of Informatics at the University of Edinburgh. He has over 150 publications in speech and language processing, is the co-editor-in-chief of *ACM Transactions on Speech and Language Processing*, and has led several large projects in the field. With Hervé Bourlard, he was the joint coordinator of the AMI and AMIDA European Integrated Projects, which form the basis for the book.

**Hervé Bourlard** is Director of the Idiap Research Institute in Switzerland, Professor at the Swiss Federal Institute of Technology at Lausanne (EPFL), and founding Director of the Swiss National Center of Competence in Research on Interactive Multimodal Information Management (NCCR IM2). He has over 250 publications, has initiated and coordinated numerous international research projects, and is the recipient of several scientific and entrepreneurship awards.

**Jean Carletta** is a Senior Research Fellow at the Human Communication Research Centre, University of Edinburgh. She was the scientific manager of the AMI and AMIDA Integrated Projects. A former Marshall Scholar, she has been on the editorial boards of *Computational Linguistics* and *Language Resources and Evaluation*.

**Andrei Popescu-Belis** is a Senior Researcher at the Idiap Research Institute in Switzerland. He is currently heading the Swiss Sinergia project COMTIS on machine translation, and has been a member of the technical committee of the IM2 NCCR since 2006.

This book provides a critical resource for understanding audio-visual social signalling between people. It will be invaluable for guiding future technical research on collaborative multimodal interaction, communication, and learning.

*Sharon Oviatt, Incaa Designs*

# Multimodal Signal Processing

## Human Interactions in Meetings

Edited by

**STEVE RENALS**

University of Edinburgh

**HERVÉ BOURLARD**

Idiap Research Institute

**JEAN CARLETTA**

University of Edinburgh

**ANDREI POPESCU-BELIS**

Idiap Research Institute



**CAMBRIDGE**  
UNIVERSITY PRESS

CAMBRIDGE UNIVERSITY PRESS  
Cambridge, New York, Melbourne, Madrid, Cape Town  
Singapore, São Paulo, Delhi, Mexico City

Cambridge University Press  
The Edinburgh Building, Cambridge CB2 8RU, UK

Published in the United States of America by Cambridge University Press, New York

[www.cambridge.org](http://www.cambridge.org)

Information on this title: [www.cambridge.org/9781107022294](http://www.cambridge.org/9781107022294)

© Cambridge University Press 2012

This publication is in copyright. Subject to statutory exception  
and to the provisions of relevant collective licensing agreements,  
no reproduction of any part may take place without the written  
permission of Cambridge University Press.

First published 2012

Printed in the United Kingdom at the University Press, Cambridge

*A catalog record for this publication is available from the British Library*

*Library of Congress Cataloging in Publication data*

Multimodal signal processing : human interactions in meetings / edited by Steve Renals . . . [et al.].

p. cm.

Includes bibliographical references and index.

ISBN 978-1-107-02229-4 (hardback)

1. Signal processing – Digital techniques. 2. Interactive multimedia. 3. Computer input-output  
equipment. 4. Computer conferencing – Technological innovations. I. Renals, Steve.

TK5102.9.M847 2012

621.382'2–dc23

2012000305

ISBN 978-1-107-02229-4 Hardback

Cambridge University Press has no responsibility for the persistence or  
accuracy of URLs for external or third-party internet websites referred to  
in this publication, and does not guarantee that any content on such  
websites is, or will remain, accurate or appropriate.

# Contents

*List of contributors*

*page xii*

<b>1</b>	<b>Multimodal signal processing for meetings: an introduction</b>	<b>1</b>
	Andrei Popescu-Belis and Jean Carletta	
1.1	Why meetings?	2
1.2	The need for meeting support technology	3
1.3	A brief history of research projects on meetings	3
1.3.1	Approaches to meeting and lecture analysis	4
1.3.2	Research on multimodal human interaction analysis	5
1.3.3	The AMI Consortium	6
1.3.4	Joint evaluation and dissemination activities	7
1.4	Outline of the book	8
1.5	Summary and further reading	9
1.6	Acknowledgments	10
<b>2</b>	<b>Data collection</b>	<b>11</b>
	Jean Carletta and Mike Lincoln	
2.1	The AMI Meeting Corpus design	12
2.1.1	The design team exercise	12
2.1.2	Ensuring generalizability	13
2.1.3	Including participants from outside the room	14
2.2	Multimodal recording	14
2.2.1	What was captured	15
2.2.2	Synchronization	16
2.2.3	Audio	18
2.2.4	Video	18
2.3	Transcription	19
2.4	Annotations	20
2.4.1	Video annotations	21
2.4.2	Language annotations	22
2.5	Handling multiple annotations	22
2.6	Public release	24
2.7	Summary and further reading	26

---

<b>3</b>	<b>Microphone arrays and beamforming</b>	<b>28</b>
	Iain McCowan	
3.1	Introduction	28
3.2	Foundations	29
3.2.1	Key terms	29
3.2.2	Key equations	29
3.2.3	Worked example	31
3.3	Design	33
3.3.1	Array geometry	33
3.3.2	Beamforming filters	34
3.4	Application to meetings	35
3.4.1	Array geometry	35
3.4.2	Beamforming filters	36
3.5	Summary and further reading	39
<b>4</b>	<b>Speaker diarization</b>	<b>40</b>
	Fabio Valente and Gerald Friedland	
4.1	Introduction	40
4.2	State of the art in speaker diarization	40
4.3	Information bottleneck diarization	42
4.3.1	Information bottleneck principle	43
4.3.2	IB-based speaker diarization	44
4.3.3	Extension to multiple features	44
4.3.4	Realignment	45
4.3.5	Experiments	46
4.4	Dialocalization	49
4.4.1	Features	50
4.4.2	Multimodal speaker diarization	50
4.4.3	Visual localization	51
4.4.4	Properties of the algorithm	52
4.5	Summary and further reading	55
4.6	Acknowledgments	55
<b>5</b>	<b>Speech recognition</b>	<b>56</b>
	Thomas Hain and Philip N. Garner	
5.1	General overview	56
5.1.1	Meetings are different	57
5.1.2	A brief history of meeting speech recognition	57
5.1.3	Outline	58
5.2	Meeting specifics	58
5.2.1	Data sources	59
5.2.2	Data analysis	60



5.2.3	Vocabulary	61
5.2.4	Language	61
5.2.5	Acoustic modeling	64
5.3	Transcribing the AMI Corpus	65
5.3.1	Meeting type and speaker variation	65
5.3.2	Close-talking performance	66
5.3.3	Distant microphones	67
5.4	The AMIDA system for meeting transcription	69
5.4.1	Front-end processing	70
5.4.2	Acoustic modeling	71
5.4.3	Offline processing	73
5.4.4	System overview	73
5.4.5	Results and conclusions	74
5.4.6	WebASR: meeting transcription on the Web	76
5.5	Online recognition	77
5.5.1	Overview	77
5.5.2	Architecture	78
5.5.3	Voice activity detection	79
5.5.4	Wrappers	80
5.6	Keyword spotting	80
5.6.1	Methods	80
5.6.2	Evaluation metrics and campaigns	81
5.7	Summary and further reading	82
<b>6</b>	<b>Sampling techniques for audio-visual tracking and head pose estimation</b>	<b>84</b>
	Jean-Marc Odobez and Oswald Lanz	
6.1	Introduction	84
6.2	State-space Bayesian tracking	86
6.2.1	The Kalman filter	86
6.2.2	Monte Carlo methods	87
6.3	Person tracking in rooms	88
6.3.1	Specific issues	88
6.3.2	Tracking in the image plane	89
6.3.3	Tracking in 3D space with calibrated cameras	90
6.3.4	Multi-object tracking inference	92
6.4	Head tracking and pose estimation	94
6.4.1	Head tracking	95
6.4.2	Joint head tracking and pose estimation	96
6.4.3	Head pose estimation in smart rooms	98
6.5	Audio-visual tracking	99
6.5.1	Audio-visual person tracking	99
6.5.2	Head pose tracking with audio information	101
6.6	Summary and further reading	102

<b>7</b>	<b>Video processing and recognition</b>	<b>103</b>
	Pavel Zemčik, Sébastien Marcel, and Jozef Mlích	
7.1	Object and face detection	103
7.1.1	Skin color detection	103
7.1.2	Face detection	105
7.1.3	Gaze and face expression detection	109
7.1.4	Object detection evaluation	113
7.2	Face recognition	115
7.2.1	Introduction to face recognition	115
7.2.2	Overview of face recognition techniques	116
7.2.3	Face verification	117
7.2.4	Face identification	117
7.2.5	Future research directions	118
7.3	Gesture recognition	118
7.3.1	Hand detection	119
7.3.2	Simple gestures	120
7.3.3	Compound gestures	121
7.4	Summary and further reading	123
7.5	Acknowledgments	124
<b>8</b>	<b>Language structure</b>	<b>125</b>
	Tilman Becker and Theresa Wilson	
8.1	Introduction	125
8.2	Dialogue acts	125
8.2.1	Dialogue act annotation schemes	126
8.2.2	Dialogue act segmentation	127
8.2.3	Dialogue act classification	129
8.2.4	Joint segmentation and classification	130
8.3	Structure of subjective language	131
8.3.1	Two schemes for annotating subjectivity in meetings	132
8.3.2	Experiments in subjectivity and sentiment recognition	134
8.3.3	Experiments in agreement and disagreement detection	136
8.4	Topic recognition	138
8.4.1	Topics in meetings	139
8.4.2	Evaluation metrics	139
8.4.3	Features and methods for topic segmentation	140
8.4.4	Topic labeling	142
8.5	Structure of decisions	142
8.5.1	Domain models and ontologies	143
8.5.2	Negotiation acts	143
8.5.3	Discourse model	144
8.5.4	Finding all decisions	145
8.5.5	Discourse memory	146

8.5.6	Decision summaries	147
8.6	Disfluencies	148
8.6.1	Classes of disfluencies	149
8.6.2	Statistical analysis	149
8.6.3	Hybrid disfluency detection	150
8.6.4	Detection modules	151
8.6.5	Hybrid combination	151
8.6.6	Results and discussion	152
8.7	Summary and further reading	153
<b>9</b>	<b>Multimodal analysis of small-group conversational dynamics</b>	<b>155</b>
	Daniel Gatica-Perez, Riëks op den Akker, and Dirk Heylen	
9.1	Introduction	155
9.2	Conversational dynamics phenomena: definitions	156
9.2.1	Conversational attention	156
9.2.2	Turn-taking and conversational floor	157
9.2.3	Addressing	158
9.3	Automatic analysis of small-group conversational dynamics	159
9.3.1	Visual attention	159
9.3.2	Turn-taking and conversational floor	161
9.3.3	Addressing	164
9.4	Towards social inference: dominance in small groups	165
9.4.1	Annotating dominance in meetings	165
9.4.2	Automatic dominance detection	166
9.5	Open issues	167
9.6	Summary and further reading	169
9.7	Acknowledgments	169
<b>10</b>	<b>Summarization</b>	<b>170</b>
	Thomas Kleinbauer and Gabriel Murray	
10.1	Introduction	170
10.2	Extractive summarization	171
10.2.1	Interpretation	172
10.2.2	Transformation	174
10.2.3	Generation	176
10.2.4	Focused extraction: decisions and action items	177
10.3	Abstractive summarization	178
10.3.1	Representation formalism	178
10.3.2	Interpretation	181
10.3.3	Transformation	183
10.3.4	Generation	184
10.3.5	Case studies	185
10.4	Evaluation	186

10.4.1	Intrinsic evaluation	187
10.4.2	Extrinsic evaluation	188
10.5	Conclusion and discussion	191
10.6	Further reading	192
<b>11</b>	<b>User requirements for meeting support technology</b>	<b>193</b>
	Denis Lalanne and Andrei Popescu-Belis	
11.1	Models for the software development process	193
11.1.1	Definitions. The waterfall model	193
11.1.2	Limits of the waterfall model	194
11.1.3	An iterative process for meeting support technology	195
11.2	Determining user requirements: two approaches	196
11.2.1	Analysis of current practices for meeting archiving and access	197
11.2.2	A practice-centric study of access to past meeting information	198
11.2.3	Elicitation of requirements from potential users	199
11.3	Query analysis	200
11.4	From requirements to specifications	201
11.5	Summary and further reading	203
11.6	Acknowledgments	203
<b>12</b>	<b>Meeting browsers and meeting assistants</b>	<b>204</b>
	Steve Whittaker, Simon Tucker, and Denis Lalanne	
12.1	Introduction	204
12.2	Meeting browsers	205
12.2.1	Categorization of meeting browsers	205
12.2.2	Meeting browsers from the AMI and IM2 Consortia	208
12.2.3	Conference recording and browsing	212
12.3	Meeting assistants: real-time meeting support	212
12.3.1	Improving user engagement in meetings	213
12.3.2	Suggesting relevant documents during meetings	214
12.4	Summary and perspectives	216
<b>13</b>	<b>Evaluation of meeting support technology</b>	<b>218</b>
	Simon Tucker and Andrei Popescu-Belis	
13.1	Approaches to evaluation: methods, experiments, campaigns	218
13.2	Technology-centric evaluation	220
13.2.1	Target tasks and meeting data	221
13.2.2	Observations and user comments	221
13.2.3	Implications	222
13.3	Task-centric evaluation: the BET method and its results	223
13.3.1	Defining the task: the Browser Evaluation Test	223
13.3.2	Applying the BET: evaluation results	225

---

13.3.3	Discussion of the BET	228
13.4	User-centric approaches	229
13.5	A software process perspective on achievements	230
13.6	Summary and further reading	231
<b>14</b>	<b>Conclusion and perspectives</b>	<b>232</b>
	Hervé Bourlard and Steve Renals	
14.1	Goals and achievements	232
14.2	Perspectives	236
	<i>References</i>	238
	<i>Index</i>	271

# Contributors

Tilman Becker, DFKI Saarbrücken, Germany

Hervé Bourlard, Idiap Research Institute, Martigny, Switzerland

Jean Carletta, University of Edinburgh, UK

Gerald Friedland, International Computer Science Institute, Berkeley, CA, USA

Philip N. Garner, Idiap Research Institute, Martigny, Switzerland

Daniel Gatica-Perez, Idiap Research Institute, Martigny, Switzerland

Thomas Hain, University of Sheffield, UK

Dirk Heylen, University of Twente, the Netherlands

Thomas Kleinbauer, Monash University, Australia

Denis Lalanne, University of Fribourg, Switzerland

Oswald Lanz, FBK-IRST, Trento, Italy

Mike Lincoln, University of Edinburgh, UK

Sébastien Marcel, Idiap Research Institute, Martigny, Switzerland

Iain McCowan, Dev-Audio Pty Ltd, Southport, QLD, Australia

Jozef Mlích, Brno Institute of Technology, Czech Republic

Gabriel Murray, University of British Columbia, Vancouver, BC, Canada

Jean-Marc Odobez, Idiap Research Institute, Martigny, Switzerland

Rieks op den Akker, University of Twente, the Netherlands

Andrei Popescu-Belis, Idiap Research Institute, Martigny, Switzerland

Steve Renals, University of Edinburgh, UK

Simon Tucker, University of Sheffield, UK

Fabio Valente, Idiap Research Institute, Martigny, Switzerland

Steve Whittaker, University of California at Santa Cruz, CA, USA

Theresa Wilson, Johns Hopkins University, Baltimore, MD, USA

Pavel Zemčík, Brno Institute of Technology, Czech Republic

# 1 Multimodal signal processing for meetings: an introduction

---

Andrei Popescu-Belis and Jean Carletta

This book is an introduction to multimodal signal processing. In it, we use the goal of building applications that can understand meetings as a way to focus and motivate the processing we describe. Multimodal signal processing takes the outputs of capture devices running at the same time – primarily cameras and microphones, but also electronic whiteboards and pens – and automatically analyzes them to make sense of what is happening in the space being recorded. For instance, these analyses might indicate who spoke, what was said, whether there was an active discussion, and who was dominant in it. These analyses require the capture of multimodal data using a range of signals, followed by a low-level automatic annotation of them, gradually layering up annotation until information that relates to user requirements is extracted.

Multimodal signal processing can be done in real time, that is, fast enough to build applications that influence the group while they are together, or offline – not always but often at higher quality – for later review of what went on. It can also be done for groups that are all together in one space, typically an instrumented meeting room, or for groups that are in different spaces but use technology such as videoconferencing to communicate. The book thus introduces automatic approaches to capturing, processing, and ultimately understanding human interaction in meetings, and describes the state of the art for all technologies involved.

Multimodal signal processing raises the possibility of a wide range of applications that help groups improve their interactions and hence their effectiveness between or during meetings. However, developing applications has required improvements in the technological state of the art in many arenas.

The first arena comprises core technologies like audio and visual processing and recognition that tell us basic facts such as who was present and what words were said. On top of this information comes processing that begins to make sense of a meeting in human terms. Part of this is simply combining different sources of information into a record of who said what, when, and to whom, but it is often also useful, for instance, to apply models of group dynamics from the behavioral and social sciences in order to reveal how a group interacts, or to abstract and summarize the meeting content overall. Finding ways to integrate the varying analyses required for a particular meeting support application has been a major new challenge.

*Multimodal Signal Processing: Human Interactions in Meetings*, ed. Steve Renals, Hervé Bourlard, Jean Carletta, and Andrei Popescu-Belis. Published by Cambridge University Press. © Cambridge University Press 2012.

Finally, moving from components that model and analyze multimodal human-to-human communication scenes to real-world applications has required careful user requirements capture, as well as interface and systems design. Even deciding how to evaluate such systems breaks new ground, whether it is done intrinsically (that is, in terms of the accuracy of the information the system presents) or from a user-centric point of view.

## 1.1 Why meetings?

The research described in this book could be applied to just about any setting where humans interact face-to-face in groups. However, it is impossible to design reasonable end-user applications without focusing on a specific kind of human interaction. Meetings provide a good focus for several reasons.

First, they are ubiquitous. Meetings pervade nearly every aspect of our communal lives, whether it is in work, in the running of community groups, or simply in arranging our private affairs. Meetings may not be the only way in which humans interact, but they are a frequent and understandable one, with obvious practical relevance.

Second, what happens in meetings (or, as often, what does not) is actually important. For many people, meetings are the milestones by which they pace their work. In truly collaborative decision-making, the meeting is where a group's goals and work take shape. Even in groups where the real decision-making takes place behind the scenes, in the absence of written documents the meeting itself is where a group's joint intention is most fully and most clearly expressed. Being able to understand what happens in meetings is bound to be useful, whether the goal is to reveal the content of the meeting or simply to identify where a group's process could be improved.

Third, because of changes in modern society, meetings present an obvious opportunity. Many organizations operate globally. There are few jobs for life. In the face of staff churn and business fragmentation, it is increasingly difficult for organizations simply to keep and access the institutional memory they need in order to make good decisions. Adequately documenting everything in writing is expensive, if not impossible. This makes it economically important to get better control of the information locked in meetings, starting from adequate options to record, analyze, and access some of the media related to them.

Finally, a great many meetings take place in settings where there is already, or is developing, a sense that the benefits of recording outweigh privacy considerations. Many organizations already record and archive at least their key meetings routinely, even without decent tools for sifting later through what they have stored. This is not just a matter of the technology for recording being cheap enough (although of course this is a factor), but of the organizations hoping to function better thanks to the recordings. This in itself brings benefits for an organization's members, but there can be more personal benefits too. Meetings may be ubiquitous, but we cannot always be at all of the ones that affect us. Being able to glean their content efficiently is likely to help.



## 1.2 The need for meeting support technology

Like other business processes, meetings are going digital. Increasingly, people are using computer technology alone and in conjunction with broadband networks to support their meeting objectives. E-mail is used to pass around files for people to read prior to a meeting. Collaborative workspaces in corporate networks and on the Internet offer geographically distributed collaborators a virtual repository for documents related to a project or a meeting. Electronic meeting support systems, such as interactive network-connected white boards and videoconferencing appliances, are available for the benefit of those who share the same room as well as those who are in remote locations.

Meetings play a crucial role in the generation of ideas, documents, relationships, and actions within an organization. Traditionally, depending on the type of meeting, either everyone will take whatever style of notes they please, or one person will create official written minutes of the meeting. Whatever the form of written record, it will be subjective and incomplete. Even with the best minutes, business questions often appear later, which can only possibly be resolved by going back to what actually happened. The technology now exists to capture the entire meeting process, keeping the text and graphics generated during a meeting together with the audio and video signals.

If only people could use the multimedia recordings of meetings to find out or remember what they need to know about the outcome of a meeting, then using these recordings would become an attractive adjunct (or even, alternative) to note taking. This can only happen once it is possible to recognize, structure, index, and summarize meeting recordings automatically so that they can be searched efficiently. One of the long-term goals of meeting support technology is to make it possible to capture and analyze what a group of people is doing together in a room-sized space using portable equipment, and to put together a wide range of applications supporting the group, using configurable componentry or web services for tasks like recognizing the speech, summarizing, and analyzing the group's interaction. This will enable companies to make use of archives of meetings, for instance, for audit purposes or to promote better cohesion in globalized businesses. Different configurations of the same underlying components will also help people who work away from the office to participate more fully in meetings. These possibilities indicate that we are at the point of a big technological breakthrough.

## 1.3 A brief history of research projects on meetings

The ideas presented in this book stem for a large part, though not exclusively, from the contributions made by the members of the AMI Consortium. This network of research and development teams was formed in the year 2003 building upon previous collaborations. However, several other large initiatives focused as well on multimodal signal processing and its application to meeting analysis and access, and were either precursors or contemporaries of AMI.

### 1.3.1 Approaches to meeting and lecture analysis

The understanding of human communication has long been a theoretical goal of artificial intelligence, but started having also practical value for information access through the 1990s, as more and more audio-visual recordings were available in digital formats. During the 1990s, separate advances in the audio and video analysis of recordings led to the first implemented systems for interaction capture, analysis, and retrieval. The early Filochat system (Whittaker *et al.*, 1994b) took advantage of handwritten notes to provide access to recordings of conversations, while BBN's Rough'n'Ready system (Kubala *et al.*, 1999) enhanced audio recordings with structured information from speech transcription supplemented with speaker and topic identification. Video indexing of conferences was also considered in early work by Kazman *et al.* (1996). Multi-channel audio recording and transcription of business or research meetings was applied on a considerably larger scale in the Meeting Recorder project at ICSI, Berkeley (Morgan *et al.*, 2001, 2003), which produced a landmark corpus that was reused in many subsequent projects.

Around the year 2000, it became apparent that technologies for meeting support needed to address a significant subset of the modalities actually used for human communication, not just one. This in turn required appropriate capture devices, which needed to be placed in instrumented meeting rooms, due to constraints on their position, size, and connection to recording devices, as exemplified by the MIT Intelligent Room with its multiple sensors (Coen, 1999). The technology seemed mature enough, however, for corporate research centers to engage in the design of such rooms and accompanying software, with potential end-user applications seeming not far from reach.

For instance, Classroom 2000 (Abowd, 1999) was an instrumented classroom intended to capture and render all aspects of the teaching activities that constitute a lecture. The Microsoft Distributed Meetings system (Cutler *et al.*, 2002) supported live broadcast of audio and video meeting data, along with recording and subsequent browsing. Experiments with lectures in this setting, for example for distance learning, indicated the importance of video editing based on multimodal cues (Rui *et al.*, 2003). Instrumented meeting or conference rooms were also developed by Ricoh Corporation, along with a browser for audio-visual recordings (Lee *et al.*, 2002), and by Fuji Xerox at FXPAL, where the semi-automatic production of meeting minutes, including summaries, was investigated (Chiu *et al.*, 2001).

However, even if companies were eager to turn meeting support technology into products, it became clear that in order to provide intelligent access to multimedia recordings of human interaction a finer-grained level of content analysis and abstraction was required, which could simply not be achieved with the knowledge available around the year 2000. Technology for remote audio-visual conferencing has been embedded into a host of successful products,<sup>1</sup> but without analyzing the conveyed signals and generally with highly limited recording or browsing capabilities.

<sup>1</sup> To name but a few: HP's Halo (now owned by Polycom) or CISCO's WebEx for the corporate market, and Skype, iChat, or Adobe Connect as consumer products.

### 1.3.2 Research on multimodal human interaction analysis

The need for advanced multimodal signal processing for content abstraction and access has been addressed in the past decade by several consortia doing mainly fundamental research. Only such collaborative undertakings could address the full complexity of human interaction in meetings, which had long been known to psychologists (e.g., Bales, 1950, McGrath, 1984). Moreover, only such consortia appeared to have the means to collect large amounts of data in normalized settings and to provide reference annotations in several modalities, as needed for training powerful machine learning algorithms. The public nature of most of the funding involved in such initiatives ensured the public availability of the data.

Two projects at Carnegie Mellon University (CMU) were among the first to receive public funding to study multimodal capture, indexing, and retrieval, with a focus on meetings. The target of the Informedia project was first the cross-modal analysis of speech, language, and images for digital video libraries (1994–1999), and then the automatic summarization of information across multimedia documents (1999–2003) (Wactlar *et al.*, 1996, 2000). In parallel, CMU's Interactive Systems Laboratory initiated a project on meeting record creation and access (Waibel *et al.*, 2001a). This was directly concerned with recording and browsing meetings based on audio and video information, emphasizing the role of speech transcription and summarization for information access (Burger *et al.*, 2002).

In Europe, the FAME project (Facilitating Agent for Multicultural Exchange, 2002–2005) developed the prototype of a system that made use of multimodal information streams from an instrumented room (Rogina and Schaaf, 2002) to facilitate cross-cultural human–human conversation. A second prototype, the FAME Interactive Space (Metze *et al.*, 2006), provided access to recordings of lectures via a table top interface that accepted voice commands from a user. The M4 European project (MultiModal Meeting Manager, 2002–2005), introduced a framework for the integration of multimodal data streams and for the detection of group actions (McCowan *et al.*, 2003, 2005b), and proposed solutions for multimodal tracking of the focus of attention of meeting participants, multimodal summarization, and multimodal information retrieval. The M4 Consortium achieved a complete system for multimodal recording, structuring, browsing, and querying an archive of meetings.

In Switzerland, the IM2 National Center of Competence in Research is a large long-term initiative (2002–2013) in the field of Interactive Multimodal Information Management. While the range of topics studied within IM2 is quite large, the main application in the first two phases (2002–2009) has focused on multimodal meeting processing and access, often in synergy with the AMI Consortium. The IM2 achievements in multimodal signal processing (see for instance Thiran *et al.*, 2010) are currently being ported, via user-oriented experiments, to various collaborative settings.

Two recent joint projects were to a certain extent parallel to the AMI and AMIDA projects. The CHIL European project (Computers in the Human Interaction Loop, 2004–2007) has explored the use of computers to enhance human communication in smart environments, especially within lectures and post-lecture discussions, following

several innovations from the CMU/ISL and FAME projects mentioned above (Waibel and Stiefelhagen, 2009). The US CALO project (Cognitive Assistant that Learns and Organizes, 2003–2008) has developed, among other things, a meeting assistant focused on advanced analysis of spoken meeting recordings, along with related documents, including emails (Tür *et al.*, 2010). Its major goal was to learn to detect high-level aspects of human interaction which could serve to create summaries based on action items.

It must be noted that projects in multimodal signal processing for meetings appear to belong mainly to three lineages: one descending from CMU/ISL with the FAME and CHIL projects (with emphasis on lectures, video processing and event detection), another one from ICSI MR to CALO (with emphasis on language and semantic analysis), and finally the lineage from M4 and IM2 to AMI and AMIDA (with a wider and balanced approach). Of course, collaborations between these three lineages have ensured that knowledge and data have moved freely from one to another.

### 1.3.3 The AMI Consortium

The technologies and applications presented in this book are closely connected to the research achievements of the AMI Consortium, a group of institutions that have advanced multimodal signal processing and meeting support technology. The AMI Consortium was constituted around 2003, building on existing European and international expertise, and on previous collaborations. The consortium was funded by the European Union through two successive integrated projects: Augmented Multiparty Interaction (AMI, 2003–2006) and Augmented Multiparty Interaction with Distance Access (AMIDA, 2006–2009). As a result, the consortium was highly active for more than seven years, which represents a particularly long-term multi-disciplinary research effort, surpassed only by certain national initiatives such as the Swiss IM2 NCCR (twelve years). This book presents only a selection of what the AMI Consortium has achieved, but also includes relevant advances made by the wider research community.

The AMI Consortium has included both academic partners (universities and not-for-profit research institutes) and non-academic ones (companies or technology transfer organizations). Although the partnership has varied over the years, the academic partners were the Idiap Research Institute, the University of Edinburgh, the German Research Center for AI (DFKI), the International Computer Science Institute (ICSI, Berkeley), the Netherlands Organization for Applied Scientific Research (TNO), Brno University of Technology, Munich University of Technology, Sheffield University, the University of Twente, and the Australian CSIRO eHealth Research Center. The primary non-academic partners were Philips and Noldus Information Technology. Interested companies who were not project partners were able to interact with the AMI Consortium through the AMI Community of Interest and in focused “mini-project” collaborations. These interactions allowed industry to influence the research and development work based on market needs and to prepare to use AMI technology within existing or future products and services.

### 1.3.4 Joint evaluation and dissemination activities

In many fields, the existence of a shared task – with standardized data sets and evaluation metrics – has served as a driving force to ensure progress of the technology. Shared tasks offer an accurate comparison of methods at a given time. They also provide training and test data, thus lowering the entry cost for new institutions interested in solving the task. Shared tasks and standardized evaluation began in 1988 for automatic speech recognition, and since then, the approach has spread more widely.

For multimodal signal processing applied to meetings or lectures, two initiatives have promoted shared tasks: the Rich Transcription (RT) evaluations and the Classification of Events Activities and Relationships (CLEAR) ones. In both series, the US National Institute for Standard Technology (NIST) has played a pivotal role in gathering normalized data that was considered by participants to be representative of the addressed research questions. Along with external data from the AMI and CHIL consortia, NIST has also produced original data in its own instrumented meeting rooms, starting from the Smart Spaces Laboratory (Stanford *et al.*, 2003).

The NIST annual RT evaluations started as early as 2001 for broadcast news and telephone conversations, and meetings were targeted starting 2004. Following increasing interest, the most visible results were produced in the 2005–2007 campaigns, the latter one being organized and published jointly with CLEAR (Stiefelhagen *et al.*, 2008); a smaller workshop was further held in 2009. The goal of the RT evaluations was to compare the performance of systems submitted by participants on meetings of varying styles recorded using multiple microphones. The systems were mainly for automatic speech recognition (producing text from speech, including punctuation and capitalization) and for speaker diarization (determining who spoke when). RT differed from other campaigns for speech recognition, such as broadcast news, in its emphasis on multiple, simultaneous speakers and on non-intrusive capture devices, but did not target higher-level information extraction capabilities on meeting signals, such as those developed by AMI or CALO.

The CLEAR evaluations were sponsored by the US VACE program (Video Analysis and Content Extraction) with support from CHIL and an infrastructure provided by NIST. The CLEAR 2006 and 2007 evaluations (Stiefelhagen and Garofolo, 2007, Stiefelhagen *et al.*, 2008) targeted mainly the problems of person and face tracking, head pose estimation, and acoustic event detection using signals from several capture devices (cameras, microphones) in instrumented meeting rooms. Several conditions were tested for each track, although some of them remained experimental only. The CLEAR evaluations used data from CHIL and AMI, as well as NIST and VACE (Chen *et al.*, 2005), some of it being shared with RT.

Beyond the established scientific events and scholarly journals which disseminate work on meeting analysis and access, the community has also created a new dedicated forum, the Machine Learning for Multimodal Interaction (MLMI) workshops, initiated more specifically by the AMI and IM2 consortia. Many of the research results gathered in this book were originally presented at MLMI

workshops.<sup>2</sup> Due to converging interests and complementarity, joint events between MLMI and the International Conference on Multimodal Interfaces (ICMI) were organized in 2009 and 2010. Following their success, the two series merged their advisory boards and decided to hold annual conferences under the name of International Conference on Multimodal Interaction.

## 1.4 Outline of the book

In order to design tools with the potential to unlock the business value contained in meetings, researchers in several related fields must collaborate. There are many places to find information about components like speech recognition that are the building blocks for the new technology. However, understanding the global picture requires a basic understanding of work from a wide range of disciplines, and help for developing that understanding is much harder to find. One particular challenge is in how to use what organizational and social psychologists know about human groups to determine user requirements and methods of testing technologies that users cannot really imagine yet. Another is in joining work on individual communication modalities like speech and gesture into a truly multimodal analysis of human interaction. While this book does not pretend to offer a fully integrated approach, the longevity of the collaborations between its authors has enabled many new connections and the feeling that it was possible to understand and achieve more by working together. One of the goals of this book is to pass on that understanding, making it easier for new researchers to move from their single disciplines into a rewarding and exciting area.

The book begins with something that underpins everything that follows: the data. Chapter 2 presents a hardware and software infrastructure for meeting data collection and annotation, initially designed for the comprehensive recording of four-person meetings held in instrumented meeting rooms. The rooms were used to record the AMI Meeting Corpus (Carletta, 2007), which consists of 100 hours of meeting recordings, along with manually produced transcriptions and other manual annotations that describe the behavior of meeting participants at a number of levels.

After Chapter 2, the book contains two unequal parts: Chapters 3–10 and Chapters 11–13. The first part explains the range of technological components that make up multimodal signal processing. Each chapter takes one kind of analysis that an application might need and describes what it does, how it works (and how well), and what the main issues are for using it. The advances in audio, visual, and multimodal signal processing are primarily concerned with the development of algorithms that can automatically answer, using the raw audio-video streams, questions such as the following ones: What has been said during the meeting? Who has spoken when? Who and where are the persons in the meeting? How do people behave in meetings? What is the essence of what has been said? In general, the order of the chapters reflects a progress towards

<sup>2</sup> The workshop proceedings were published as revised selected papers in Springer's Lecture Notes in Computer Science series, numbers 3361 (Martigny, 2004), 3869 (Edinburgh, 2005), 4299 (Bethesda, MD, 2006), 4892 (Brno, 2007) and 5237 (Utrecht, 2008).

more and more content abstraction, building up higher and higher levels of information from raw audio and video signals.

Chapters 3 to 5 build up towards an understanding of what was said in a meeting, primarily (but not entirely) based on audio signals, from microphone arrays (Chapter 3) to speaker diarization (determining who spoke when, Chapter 4) and automatic meeting transcription (Chapter 5). Chapters 6 and 7 move to focus more substantively on video processing as a source of information, again building upwards from the raw signals. Chapter 6 deals with tracking individual people, and especially their heads, as they move through a space. Chapter 7 then builds on this work to discuss methods for finding people and faces in recordings, recognizing faces, and interpreting head and hand gestures.

The remaining chapters in the first part of the book develop more of what a layperson would consider an understanding of a meeting. Chapter 8 describes analyses that begin to make sense of the words that were said, such as removing disfluencies, identifying questions, statements, and suggestions, or identifying subjective statements, such as positive opinions. Chapter 9 is more social in nature, and covers the analysis of conversational dynamics, in particular in terms of which speakers are being most dominant conversationally, and the different roles that they take in the meeting. Finally, Chapter 10 addresses a higher-level but very important task: that of creating useful summaries of meetings.

The second part of the book (Chapters 11–13) considers how to design, build, and test applications that use multimodal signal processing to analyze meetings. It takes the reader from the methods for identifying user needs for meeting support technology and their results (Chapter 11), through a range of meeting browsing applications that draw on underlying components from the first part (Chapter 12), to the methods for evaluating them (Chapter 13). The focus is particularly on meeting browsers, the most mature of the new technologies, which allow users to find information from past meetings, but the material also covers applications that support groups as they meet.

Finally, the conclusion (Chapter 14) abstracts from the lessons learned in analyzing meetings, and adopts a critical perspective to show what interesting and scientific challenges are still left ahead of us, and their potential impact in other application domains, such as social signal processing.

## 1.5 Summary and further reading

Multimodal signal processing has now had a decade of investment, including the promotion of shared tasks that allow the results from different techniques to be compared. It has benefited immensely from hardware advances that make synchronized recordings of audio and video signals relatively cheap to make and store. There are now many different automatic analyses available as components for systems that will do new, useful, and interesting things with these recordings. Although meeting support technology is only one of the many possibilities, the emergence of corporate meeting archives and the business value locked in them make it an obvious choice.



We conclude this introduction (and, indeed, every chapter of this book) with suggestions for further reading. These include mostly books at comparable levels of generality; more focused articles on specific topics are indicated in the respective chapters, while the names of relevant periodicals and conference series can simply be found by browsing the bibliography at the end of the book.

The books by [Thiran \*et al.\* \(2010\)](#) and [Waibel and Stiefelbogen \(2009\)](#) draw on some of the same core technologies as the present book, but cover certain additional aspects not dealt with here, such as human–computer interaction (HCI), speech synthesis, or multimodal fusion. The second book is a collection of papers summarizing achievements from the CHIL project, each of them with a close focus on specific research results. Books like those by [Cassell \*et al.\* \(2000\)](#) and by [Stock and Zancanaro \(2005\)](#) are in the same general area of multimodal interaction, but focus on presenting, not obtaining, information from multimodal data. An overview of machine learning algorithms for processing monomodal communication signals similar to those analyzed in this book is provided by [Camastra and Vinciarelli \(2008\)](#). There are many books about multimodal HCI, such as those by [Wahlster \(2006\)](#), from the SmartKom project, or by [Grifoni \(2009\)](#), which include spoken and multimodal dialogue interfaces and mobile devices. The proceedings of the MLMI conferences series of work mentioned in Section 1.3.4 represent additional collections of in-depth research articles (e.g., [Popescu-Belis and Stiefelbogen, 2008](#)).

## 1.6 Acknowledgments

Most of the contributors to this book, though not all, have been connected to some extent to the AMI Consortium. The editors and authors are grateful for the significant support of the European Union, through the Sixth Framework Programme for Research in its Information Society Technology (IST) thematic priority, as well as the support of the Swiss National Science Foundation through its NCCR division.

More specifically, the following grants have supported the research presented here, as well as the preparation of the book itself: the AMI EU integrated project (FP6, no. IST-2002-506811), the AMIDA integrated project of the EU (FP6, no. IST-033812), and the IM2 NCCR of the Swiss SNSF. Unless otherwise stated, the research work described in this book was funded by these sources. Additional funding sources are acknowledged at the end of each chapter.

The editors would like to thank the staff at Cambridge University Press, in particular Dr. Philip Meyler and Ms. Mia Balashova, their copy-editor Mr. Jon Billam, as well as Dr. Pierre Ferrez from Idiap, for their help with the production of this book.



## 2 Data collection

---

Jean Carletta and Mike Lincoln

One of the largest and most important parts of the original AMI project was the collection of a multimodal corpus that could be used to underpin the project research. The AMI Meeting Corpus<sup>1</sup> contains 100 hours of synchronized recordings collected using special instrumented meeting rooms. As well as the base recordings, the corpus has been transcribed orthographically, and large portions of it have been annotated for everything from named entities, dialogue acts, and summaries to simple gaze and head movement behaviors. The AMIDA Corpus<sup>2</sup> adds around 10 hours of recordings in which one person uses desktop videoconferencing to participate from a separate, “remote” location.

Many researchers think of these corpora simply as providing the training and test material for speech recognition or for one of the many language, video, or multimodal behaviors that they have been used to model. However, providing material for machine learning was only one of our concerns. In designing the corpus, we wished to ensure that the data was coherent, realistic, useful for some actual end applications of commercial importance, and equipped with high-quality annotations. That is, we set out to provide a data resource that might bias the research towards the basic technologies that would result in useful software components. In addition, we set out to create a resource that would be used not just by computationally oriented researchers, but by other disciplines as well. For instance, corpus linguists need naturalistic data for studying many different aspects of human communication. Organizational psychologists want to know how well-functioning work groups behave and how to measure their effectiveness. Where it does not distort the research, it makes sense to reuse our data for these purposes. Moreover, interacting with these disciplines helps us in our own technological goals. For instance, within the project, developing a basic understanding of how addressing works in face-to-face meetings has helped us know what help people who are connecting by phone or laptop might need in order to participate fully. If we can measure group effectiveness, then we can test whether our technologies actually make groups more effective. It was important to us to make the corpora useful in all of these ways.

<sup>1</sup> <http://corpus.amiproject.org>

<sup>2</sup> <http://corpus.amidaproject.org>

## 2.1 The AMI Meeting Corpus design

The simplest approach to corpus building is simply to collect whatever material is easy to obtain but still varies enough to represent the range of source types and phenomena one wishes to cover. For our needs, this approach was insufficient, primarily because this makes it very difficult to assess group effectiveness. Although one can ask group members what they think, using, for instance, questionnaires, the results can reflect how well they get on more than how well the group is actually doing. If all the groups are of the same type (for instance, all primary health care teams or all design meetings) then domain experts might be brought in to rank the groups, but this is difficult to organize and highly subjective. Moreover, real work groups are never truly comparable – even if they are doing the same thing, some will always face more difficult conditions than others. Even understanding why they behave as they do can be difficult if they have a long history of working with each other or if they are not entirely sure what they are meant to achieve.

Our solution to the problem of needing to know how well teams do under different conditions is the standard one employed by psychologists: experimental control. That is, rather than simply observe real workplace teams, we have people who don't really work together play different roles in a team where we control the task.

### 2.1.1 The design team exercise

About 70% of the AMI Meeting Corpus makes use of an exercise in which four participants spend half a day playing different roles in a fictional design team. Design teams are a good target for our meeting archive technologies: they are widespread in industry, have functional meetings with clear goals that are relatively measurable, and rely heavily on information from past meetings. We first tell the participants something about the task for which they have been hired – the design of a new kind of remote control – and give them separate crash courses in what it means to be a marketing expert, project manager, industrial designer, or user interface designer. We then leave them to do their work, recording what they do.

Although the company and team are fictional, we make the exercise as realistic as we can while still completing it in a day. So, for instance, although there are four set meeting times for the group to achieve steps towards their goal, the team members do their “ordinary” work, too – gathering information, working out ideas, and preparing presentations for the meetings – in a completely normal working environment. They receive automated emails from other people in the company throughout the day, some of which give them relevant information or advice, and some of which change their goals.

Because we give the groups a number of constraints that their design should meet – such as that it should be in the company's colors – we can measure how well they do by checking their designs against the full set of constraints. Meeting the constraints requires the participants to pool the information they have been given, including new information

as it comes in. There is an obvious optimal solution that met all of the constraints, but few of the groups have found it.

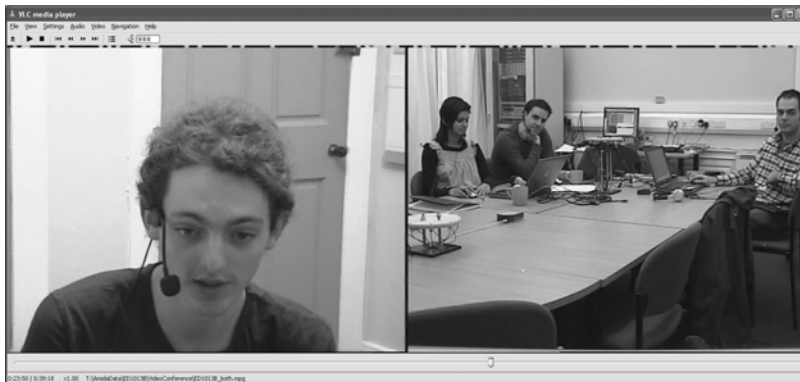
Our experimental control raises an interesting possibility: not only does it allow us to compare different groups, we can potentially compare different groups under different conditions. For instance, one way of showing that our technologies help meeting participants would be to run new groups that have access to them.

### 2.1.2 Ensuring generalizability

Using a role-play exercise does have some drawbacks. The first is that because all the groups talk about the same thing, the vocabulary size of meetings involving our role play is roughly half that of other spontaneous speech. This makes it easier to build components like speech recognizers that work well on the corpus. That sounds like an advantage, but it is not – it just makes things more difficult when it comes to deploying real systems. The second is that because the groups aren't necessarily motivated in the same way, and to the same degree, as real workplace teams, we can't be sure that real teams act enough like our role-playing groups to make the data pertinent. Our answer to these drawbacks is to collect most of our data using our role-playing method, but to collect some real meetings as well. As stated above, 70% of the AMI Meeting Corpus consists of our experimental role play, and the other 30% contains whatever data we found we could obtain in our instrumented rooms, but with priority given to series of recurring meetings and to meetings in which the participants were discussing things that really mattered to them. There is an inevitable bias towards scientific discussions in areas related to our own research. The additional data helps to ensure that the material we train our recognizers on is not too idiosyncratic, as well as giving us test material that will tell us how much the domain matters.

There are also other ways we have tried to safeguard the corpus so that technologies built using it will port well. All rooms have their own individual acoustic and visual properties, and training and testing components on data from only one room may limit their use on data from other environments. In addition, the layout of the room can have a significant effect on the participants' interaction – for instance, sometimes whiteboards are behind projector screens, which makes them unlikely to be used during meetings that also contain presentations. In order to avoid such issues we capture data from three different meeting rooms. Each contains essentially the same capture devices, but each is significantly different in its physical properties.

Finally we must consider the participants themselves. Working in multiple languages greatly increases the cost of data collection and component development, and so we reluctantly limited ourselves to English, in favor of being able to consider more components and more applications. Even so, there is one type of language breakthrough that we felt we could make: ensuring that our participants had a broad range of accents and language backgrounds, so that our technologies would be more robust than those developed on the more usual mono-cultural data. The case for our approach is particularly strong when developing meeting applications: in a globalized era, people from different



**Fig. 2.1** Sample image from the videoconference system.

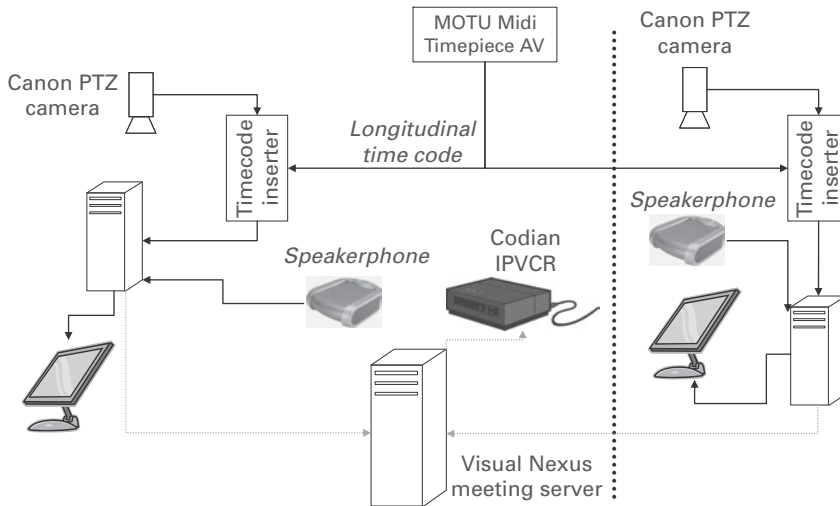
language backgrounds often do meet together, and simply because of its dominance “international English” is often their choice of language.

### 2.1.3 Including participants from outside the room

The final thing to note about our design is the inclusion of an additional ten hours of data in which one person participates “remotely” in the design role-play from a different location. This data comes separately as the AMIDA Meeting Corpus. In our arrangement, the meeting participants are free to move the cameras that they use to communicate so that they focus on anything they wish; example images from the conference are shown in Figure 2.1. Any presentations shown using the overhead projector in the main room were also viewable on the remote participant’s laptop, and conversely, they could also choose to present slides in the main room. Although we were recording high-quality audio and video in both rooms, the audio and video feeds the participants used to communicate were the lower-quality ones that are typical of desktop conferencing systems. A diagram of the videoconference hardware is shown in Figure 2.2. Although in theory the remote participant could have been anywhere, to make data capture easier, they were in a room that was visually and acoustically isolated from the main one, but close enough to run cables to the rest of the recording equipment.

## 2.2 Multimodal recording

Although we tend to think of communication as mostly about the words we speak, in order to understand what groups do together during meetings, even human analysts need to work from multimodal recordings. Participants speak, gesture, present, and write simultaneously, and to analyze these modalities individually is to miss out on important information which may radically alter one’s view of the meeting. For instance, a literal transcription of the speech from a meeting may show a recommendation which seems to be unopposed, but which the project manager dismisses out of hand with a



**Fig. 2.2** Videoconference equipment.

shake of their head. To allow complete understanding of the interaction, each modality must be captured with sufficient quality to allow individual analysis. In addition, all modalities must be captured simultaneously and in a synchronized manner, so that when they are combined they give a complete representation of the meeting.

### 2.2.1 What was captured

For deployed technology, there is a limit to the data that it makes sense to capture and analyze. However, for our purpose – creating a corpus that will result in a suite of component technologies and end-user applications – it makes sense to capture everything we can as completely as possible. This both allows us to develop the full range of applications, from those with poor audio only to those that make use of the richest range of signals, and to determine where the trade-offs lie in the cost and inconvenience of different methods for data capture against what we can get out of the signals. For each meeting in the AMI Meeting Corpus, we have collected:

- Audio recordings, including both far-field recordings from microphones placed around the room and recordings from close talking microphones for each participant.
- Video recordings, including both wide-angle views of the entire meeting room and close-up views of each participant.
- Images of any slides projected from the participants' laptops, given as a JPEG image for each slide change with accompanying transcription, derived using optical character recognition (OCR). We produce the images by capturing the video temporarily, and then extracting still images for each slide change, judged as occurring whenever enough pixels change on the image and then the image stays static enough to indicate that the new slide was not just “flipped through.”