

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/283470448>

Detecting Speech Interruptions for Automatic Conflict Detection

Article · February 2015

DOI: 10.1007/978-3-319-14081-0_18

CITATIONS

16

READS

303

2 authors, including:



Claude Montacie

Sorbonne Université

50 PUBLICATIONS 309 CITATIONS

SEE PROFILE

Detecting Speech Interruptions for Automatic Conflict Detection

Marie-José Caraty ^a, Claude Montacié ^b

^a LIPADE, Paris Descartes University, 45 rue des Saints Pères, 75006 Paris, France

^b STIH, Paris-Sorbonne University, 28 rue Serpente, 75006 Paris, France

Abstract This contribution is in the field of automatic detection of conflict in group discussions from voice analysis. A reliable detector of conflict would be useful for many applications, such as security in public places, the quality of customer services and the deployment of intelligent agents. Experiments were conducted on the SSPNet Conflict Corpus during the Interspeech'2013 Conflict Challenge. The audio clips, which were extracted from political debates, have been classified into two classes of conflict level (Low or High). In this study, we have used the turn-taking characteristics, such as interruptions, for improving the conflict detection. In a group discussion, overlapping speech (overlap) corresponds to interruption. Two overlap detectors have been developed using the SVM-classifier and audio features. The first detector aims at detecting whether interruptions occur in a speech segment. The second detector aims at detecting when interruptions occur in a speech segment and whether these interruptions are related to Low- or High-Level conflict. A multi-expert architecture has been defined to incorporate the knowledge that arises from the interruption detectors. The two-class conflict detector (Low or High Conflict) consists of an SVM-classifier that uses a composite feature set as input. This feature set is a concatenation of selected audio features and overlap detector-based features. Experiments provide an Unweighted Accuracy Recall (UAR) of 85.3% on the Test set. These results indicate an improvement of 4.5% compared to the Official Baseline System results. In conclusion, the interruptions in speech can be detected and can significantly improve the automatic conflict detection.

Keywords: Conflict detection; Political debate; Group discussion; Turn-taking characteristics; Overlapping speech; Multi-resolution detector; Audio features; Multi-expert architecture;

1 Introduction

Research in organization and management has investigated phenomena, such as the causes, effects and handling of interpersonal or intergroup conflicts. Various classifications of conflict handling, known as conflict-handling styles, have been proposed. The data on these social and psychological phenomena are collected from people who are involved in the conflict, witnesses of the conflict, or by extension, looking at a recording of the conflict escalation between the protagonists. A large quantity of data has been extracted from these recordings, such as the conversation, face and gesture interactions. The relations between these data and the measures of the conflict-handling styles have been studied and modeled. In this chapter, the conversational interactions during political debates have been studied to develop an automatic conflict detector from voice analysis. A reliable detector of conflict would be useful for many applications, such as security in public places, the quality of customer services and the deployment of intelligent agents. The development of such a system requires modeling of the conversational interactions as well as the search for specific interactions in relation to a given measure of conflict handling.

1.1 *Conflict-handling styles*

Conflict is commonly defined as an incompatibility between two or more opinions, principles, or interests. Understanding the causes, effects and handling of conflicts is important in organization and management. Many publications in this domain have investigated interpersonal or intergroup conflict and the ways in which people handle the conflict (Korabik et al. 1993, Sorenson et al. 1995, Thomas et al. 2008). Usually called the conflict-handling style in publications, this type of style refers to a behavioral response to a context-based conflict situation and to a personality-based

conflict-handling preference (Macintosh and Stevens 2008). A first classification of conflict style in one dimension, “cooperative-competitive”, has been shown to be insufficient to reflect the complexity of the subjects’ perception of conflict behavior. Conflict-handling styles were then classified with respect to two dimensions that have an evolution in their interpretation: “concern for production” and “concern for people” (Blake and Mouton 1964), “desire for own concerns” and “desire for others’ concerns” (Thomas and Kilmann 1974), “assertiveness” and “cooperativeness” (Thomas 1975), and “concern for self” and “concern for others” (Rahim 1983). According to the two dimensions, the number of conflict styles that have been investigated varies through publications from three styles (Oetzel et al. 2000) to seven styles (Euwema et al. 2003). However, the most widely used scheme classifies the conflict-handling styles into five styles (Rahim 1983), as follows: *avoid* (low concern for both self and others), *dominate* (high concern for self and low concern for others), *compromise* (moderate concern for both self and others), *oblige* (low concern for self and high concern for others) and *integrate* (high concern for self and others). For each style, the conflict-handling can be described by a strategy in attempting to reach its goal according to the style levels of concern for self and others. The definition of conflict is subject-dependent and context-dependent but also cultural-dependent. A model of Intercultural Conflict style (ICS) (Hammer 2005) was conceptualized in two dimensions: *Indirect* vs. *Direct* strategy for addressing disagreement and an *Emotional expressiveness* vs. *Restrained* strategy for addressing the affective dimension. Four basic styles were identified in Hammer’s ICS model: (1) *Accommodation* style: the subject has an indirect strategy and an emotional restraint; the conflict is handled with ambiguity in the language to prevent the conflict getting out of control and an emotional reserve to preserve interpersonal harmony. (2) *Dynamic* style: the subject has an indirect strategy and a high emotional expressiveness; the conflict is handled with repetition of one’s message, ambiguity in the language and an emotional expressiveness toward the disagreement and the other. (3) *Discussion* style: the subject has a direct strategy and an emotional restraint; the conflict is handled with precision in language and focus on the facts, (4) *Engagement* style: the subject has a direct strategy and a high

emotional expressiveness; the conflict is handled with an impartial view of each party toward a positive resolution and a degree of concern embedded in verbal and non-verbal emotional expressions. Various measurements were investigated for the perception of conflict-handling styles, with most relying on rating scales data. The management of differences exercise (MODE) (Thomas and Kilmann 1974) is a widely used paired-comparison in which the respondents must choose between two statements that describe the varying styles, each style being paired with each other style three times; the measure is based on the frequency of choice. Other measurements are mostly based on a Likert-type scale and a questionnaire for the respondents that is designed to capture the various styles (e.g., ROCI-II (Rahim 1983) and DUTCH (Euwema and van de Vliert 1990)). Reducing the biases that affect the rating scales in the way that the respondents use the scales, the Best-Worst Scaling (BWS) (Finn and Louvière 1992, Daly et al. 2010) is based on the maximal difference scaling, and the respondent must choose the best and the worst options from the subsets of a questionnaire.

1.2 Model of conversational interaction

Conversation is a social interaction between two or more people, where taking turns to talk is naturally observed. In the pioneering work of Sacks, Schegloff, and Jefferson (Sacks et al. 1974), an organizational model of turn-taking for conversation that is context-free, capable of context-sensitivity and having a cross-cultural validity was investigated. The constraints of their model were set in reference to the high cross-cultural flexibility of conversation accommodation, with a wide range of interaction in which there is a variety of persons and numbers of persons who are taking part. The authors proposed a model that relies on two components that are related to the Turn-Constructional Units (TCUs, the basic units of talk) and the turn-allocation at the end of each TCU for the next unit (the next speaker's TCU). TCUs end with points of possible completion (e.g., gap, query) called Transition-Relevant Places (TRPs), in which the turn transition could be relevant but is not necessary. Observed in any

conversation, 14 facts were listed. An excerpt of this list is the following: (a) mostly one party talks at a time, (b) the vast majority of turn-taking transitions is composed of transitions that have no/slight gap and no/slight overlap (c) the turn size varies, (d) overlapping speech is common, but brief, (e) two basic turn-allocation techniques are used: the “current selects next” technique when a current speaker can select a new speaker (e.g., addressing a question) and the “self-select” technique when a speaker can self-select in starting to talk, (f) repair mechanisms exist for addressing turn-taking violation; e.g., when overlapping speech occurs, one (or more) of the speakers will stop prematurely. A set of rules was edited for addressing turn transitions from TRP in such a way as to minimize the gap or overlap in the transitions. The turn transfer is defined according to the construction of the TCU, regardless of whether the “current speaker selects next” technique is used as well as the eventual application of “self-selection”. The rules are based on the purpose of no-gap-no-overlap transitions, for which ability is required in anticipating the precise moment at which a TCU is going to come to a completion point (i.e., a TRP). In related work (De Ruiter et al. 2006), the lexical and syntactic content of TCU was shown to be necessary for this anticipation, while the intonation contour was neither necessary nor sufficient for this projection. According to the turn-taking rule-set applied to a multi-party conversation, overlap is expected in the neighboring transition-relevant places: when a possible completion of the current TCU is wrongly projected by a party or when parties are competing in a self-selection mode for a next turn. In work that is related to turn-taking organization and that is beyond the ordinary conversation and is mostly unconstrained in terms of a role, a wide range of publications have studied the turn-taking practices and characteristics within various contexts of multi-party interactions. Distinctive features of turn-taking were found in institutional interactions in which turn-taking organization is more constrained and specialized according to the roles that are assigned to the group members (e.g., interviewer vs. interviewee, chair vs. participant). Studies on turn-taking management were investigated in institutional settings such as in a classroom (Mac Houl 1978, Mehan 1985, Lerner 1995), in courts (Atkinson and Drew 1979), in political interviews (Beattie 1982), in press conferences

(Schegloff 1987), in mediation (Garcia 1991), in professional meetings (Boden 1994), in talk shows in which interpersonal conflicts are expressed (Brinson and Winn 1997), in auctions (Heath and Luff 2007), in political debates (Valente and Vinciarelli 2010) and in political meetings that involve large groups of people in which everyone can contribute ideas, opinions, and proposals and in which opposition is also expressed (Mondada 2013). The role of the chair has been analyzed in various studies (Boden 1994, Svennevig 2008, Mondada 2012). Prediction of the speaker order in turn-taking was investigated in news, talk shows and meetings (Barzilay 2000, Vinciarelli 2009).

1.3 Guidelines and overview

In related work on conflict detection in conversational interactions, turn-taking patterns and overlaps between speakers are shown to be informative with respect to classification into the presence or absence of conflict. The total amount of overlap and the minimum pitch during overlap were found to be the features that correlated the most with conflict. A widely adopted classification of interruptions/overlaps is collaborative or competitive in reference to the “cooperative-competitive” dimension of the conflict style. While communication strategies are naturally collaborative, this preponderance is not the case for conflict dialogues, in which competitive strategies are the norm. The detection of competitive interruption is a difficult problem in relation to the search of the TRPs. Spectral content and intonation contour are not sufficient to locate these places. Furthermore, the perception of the conflict can be different in the case of the constrained organization of turn-taking, such as institutional interactions (interview, debate, meeting). Competitive strategies such as those of the moderator or the chairman appear to be natural in this context and are not perceived as conflicting.

Our experiments relate to the classification of audio clips into two classes of conflict level (Low and High) during the Interspeech’2013 Conflict Challenge. The clips, which were extracted from political debates, have been

annotated into conflict levels, using crowdsourcing to model the perception of the people. For our design of the conflict detector, we categorized the overlapping speech into Low- and High-Level Conflict overlap. We made the assumption that these categories can be detected from acoustic cues. We focus our study on a multi-resolution framework for the detection of the overlaps and a multi-expert architecture to include knowledge about overlap in the automatic conflict detector.

This chapter is organized as follows. Section 2 presents the speech material that we used for the experiments on conflict detection; it describes and analyzes the statistical characteristics of the corpus while focusing on interruptions and the moderator’s role. Section 3 describes the Conflict challenge and the various audio feature sets that were used for our investigations. In Section 4, the multi-resolution framework of the overlap detectors is outlined, the relation of the types of overlap with the conflict level is introduced and assessed, and the results are discussed according to the official measure of the challenge in terms of the UAR. Section 5 describes the multi-expert architecture of the conflict detector. Various audio features that are related to the overlap detectors are presented. The results on the conflict detector task are discussed. Section 6 presents the study’s conclusions.

2 Speech Material

The SSPNet corpus (Kim et al. 2012a) is an international reference for social signal databases. In the context of political debates, this corpus allows investigations on conflict to occur during interactions between group members. SSPNet was used for our study in analyzing various turn-taking characteristics and testing models for conflict level detection.

2.1 *SSPNet corpus*

The “SSPNet Conflict Corpus” is a collection of 45 political debates in the French language that were televised in Switzerland. It represents approximately 12 hours of speech signals; 1,430 audio clips of 30 second duration were extracted from the corpus. A total of 157 individuals were speaking in the collection of debates (23 females and 134 males). In the various multi-party discussions of the debates, the roles of the group members were distinguished: a member of the group held the role of moderator, and the other members were participants who were taking part in the debate. Four moderators (1 female, 3 males) and 153 participants (22 females, 131 males) were counted in the database. The SSPNET corpus was distributed for the Interspeech 2013 Compare challenge. Data were split into the Train, Development and Test sets: 793 clips were in the Train set, 240 clips were in the Development set, and 397 were in the Test set. Metadata are available for the Train and Development sets.

The clips were annotated in terms of the conflict score in the range minus ten to plus ten by crowd-sourcing, to model the perceptions of the data consumers at a non-verbal level; metadata were taken to be Low-Level Conflict (LLC) when the score was lower than 0; otherwise, it was taken to be High-Level Conflict (HLC). Figure 1 shows the distribution of the clips of the Train set as a function of the Conflict Score Range (CSR). The clips are split into the two classes of level conflict (LLC and HLC); the dashed line shows the boundary between the LLC and HLC clips. LLC clips are predominantly represented in the database (63% for LLC vs. 37% for HLC).

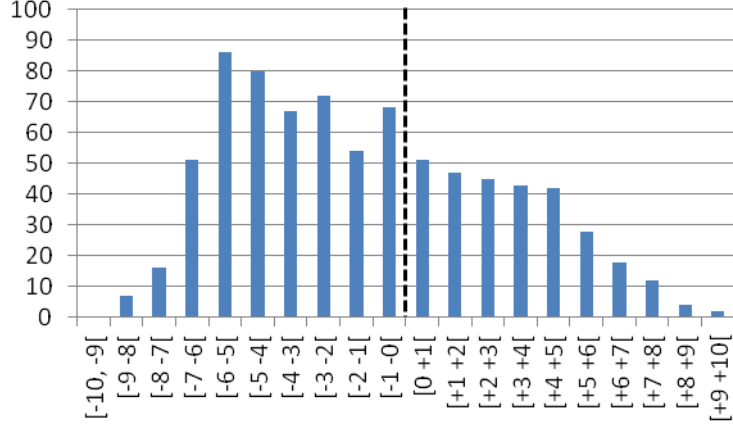


Fig. 1 Clip occurrence on the Train set as a function of the CSR

Segmentation metadata are available for each clip, indicating the diarization (“who spoke when”). From these metadata, we can compute the following statistics: (a) the overlap segment duration, (b) the clip overlap duration as the summation of each overlap segment duration of the clip, (c) the mean overlap duration of a clip as the ratio of the clip overlap duration to the number of overlaps occurring in the clip, (d) the percentage of overlap duration of the clip as the ratio of the clip overlap duration to the clip duration. Two clips of the Train set (#Train_0001 and #Train_0006) were chosen to depict the metadata for the two classes of conflict (LLC vs. HLC), to further describe the clip representation that we chose for conflict detection.

Table 1 gives an instance of diarization metadata for an LLC clip. From its conflict score (-7.2), Train_0001 clip was labeled as an LLC. There are five segments: in the fourth, a gap occurs (nobody is speaking); in the first, third and fifth segment, a lonely subject speaks (respectively, the speakers spk-50, spk-50 and spk-18), and in the second segment, two speakers are speaking at the same time (spk-47 and spk-18). The clip overlap duration is the lonely overlap segment duration (1.406 s), and the percentage of overlap duration of the clip is 4.7%.

Table 1 Speaker diarization metadata for the LLC clip #Train_0001

Start-time	End-time	Speaker-ID	Overlap.-Speaker-ID
0.0	13.009	spk-50	
13.009	14.415	spk-47	spk-18
14.415	23.686	spk-50	
23.686	24.709		
24.709	30.0	spk-18	

Conflict score: -7.2 – Conflict level: *Low* – Percentage of overlap duration: 4.7%

Table 2 gives an instance of diarization metadata for an HLC clip. From its conflict score (+7.3), Train_0006 clip was labeled as a HLC. This clip contains two overlaps, which are located at the second segment (spk-51 and spk-11) of 3.939 s duration and the fourth segment (spk-50 and spk-51) of 1.653 s duration. Its overlap number is higher than in #Train_0001 (2 vs. 1). The clip overlap duration is 5.592 s (the summation of the two overlap durations); the mean overlap duration of the clip is higher than in #Train_0001 (2.796 s vs. 1.406 s), and the percentage of overlap duration is higher than in #Train_0001 (18.6% vs. 4.7%).

Table 2 Speaker diarization metadata for the HLC clip #Train_0006

Start-time	End-time	Speaker-ID	Overlap.-Speaker_ID
0.0	14.939	spk-51	
14.939	18.878	spk-51	spk-11
18.878	28.347	spk-51	
28.347	30.0	spk-50	spk-51

Conflict score: +7.3 – Conflict level: *High* – Percentage of overlap duration: 18.6%

2.2 SSPNet Train set statistics

We analyzed the statistics of the SSPNet database Train set in focusing on the main characteristics of overlap segments; some statistics of the moderator were also investigated. The Train set includes 793 clips and has a total duration of 23,774 s (two clips duration are inferior to 30 s), with 82 speakers (one moderator and 81 participants).

We analyzed the 4,143 segments of 23,774 s duration that were obtained by the clip diarization given in the SSPNet database. These segments were split according to the number of speakers that occurred in the segment: (1) 34 segments of a total duration of 89.9 s, which correspond to gaps in which nobody is speaking, (2) 2,638 segments of a total duration of 20,083.5 s, in which a lonely subject is speaking (3) 1,471 segments of a total duration 3,600.6 s, in which two subjects are speaking. No segment was identified that had three or more speakers.

Figure 2 shows the histogram for each CSR of the average of the number of interruptions (i.e., the segments of overlapping speech) of the CSR clips. The horizontal dashed line represents the average of the number of interruptions of the Train set clips. Except for the CSR $([-1, 0[)$, all of the CSRs of LLC have a mean number of interruptions that are below the average value ($1.85 = 1,471/793$). The HLC clips have more interruptions than the LLC clips.

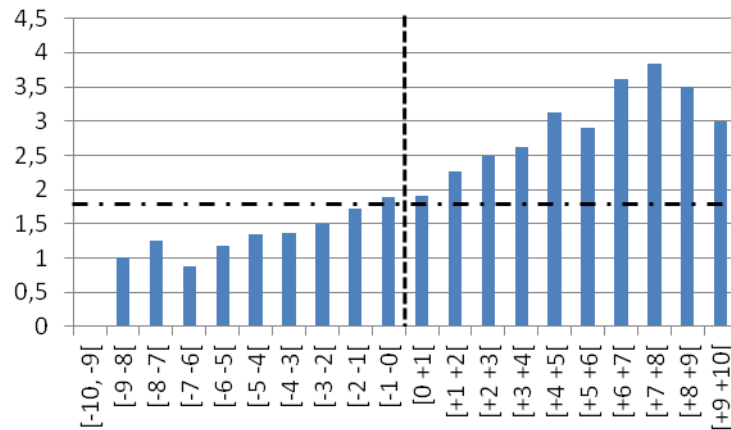


Fig. 2 Average of the number of interruptions as a function of the CSR

Figure 3 shows the histogram of the overlap mean duration for each CSR. The horizontal dashed line represents the average of the overlap duration in the Train set ($2.45 \text{ s} = 3600.6 / 1471$). HLC clips have a mean duration of overlap that is higher than the LLC clips.

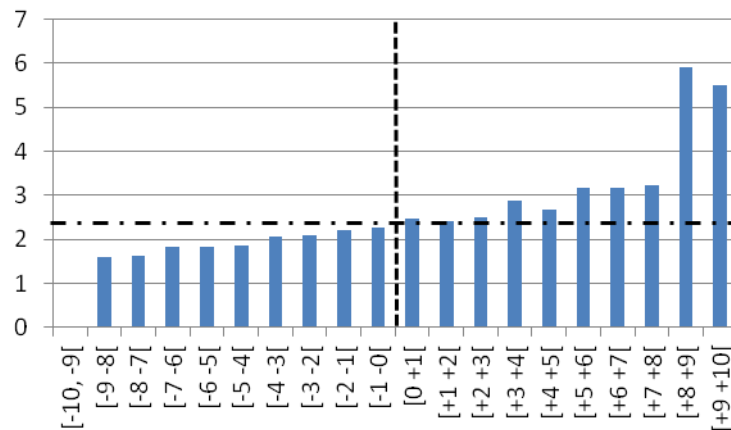


Fig. 3 Overlap mean duration (in s) as a function of the CSR

Figure 4 shows the histogram for each CSR of the percentage of overlap duration. The horizontal dashed line represents the mean percentage of the overlap duration of the Train set clips ($15.1\% = 3600.6 / 23,774$). The

conflict level is shown to be highly correlated to the percentage of overlap duration as in related work (Kim et al. 2012b).

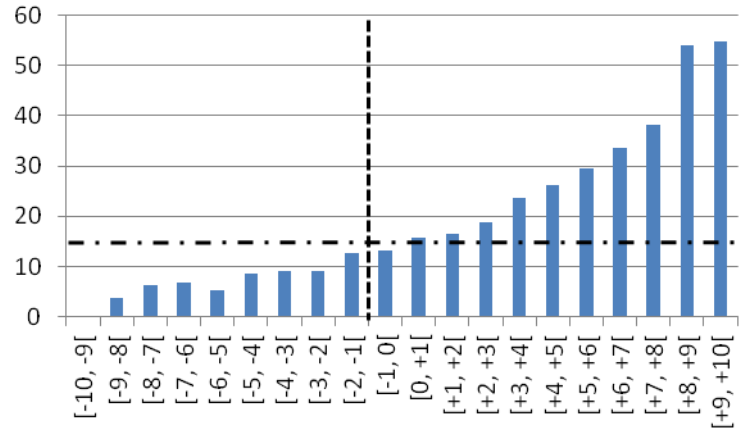


Fig. 4 Percentage of overlap duration as a function of the CSR

In the multi-party discussions of the debates, a member of the group held the role of moderator among the participant members who were taking part in the debate. We analyzed various statistics that were related to the moderator from the Train set.

In Table 3, the statistics on the speech duration of the moderator are given. The total speech duration (27,284.7 s) that we accounted for is different from the total segment duration (23,774 s) of the Train set; it was estimated as the duration of a segment in which a lonely subject is speaking plus twice the segment duration in which two subjects are speaking. The speech duration of the moderator was computed for the various classes of speech: Non-Overlap (Non-Ov) and Overlap (Ov) were split into the two conflict level classes, Low-Level Conflict Overlap (LLC-Ov) and High-Level Conflict Overlap (HLC-Ov). The speech duration of the moderator is given (in s). The percentage of the moderator speech duration was computed as the ratio between the speech duration of the moderator to the total speech duration. From the previous statistics, we note in Table 3 that the moderator speaks more during LLC-Ovs than during HLC-Ovs (39.3% vs. 20.4%) and Non-Ovs (39.3% vs. 15.9%).

Table 3 Statistics on the speech duration of the moderator

Moderator – spk-050	Train set	Non-Ov	Ov	LLC-Ov	HLC-Ov
Total speech duration (in s)	27,284.7	20,083.5	7,201.2	2,619.0 s	4,582.2
Speech duration of the moderator (in s)	5,149.7	3,183.7	1,966	1,029.8	936.2
Percentage of the moderator speech duration	18.9%	15.9%	27.3%	39.3%	20.4%

In Table 4, the modes of the interruptions that were related to the moderator were analyzed according to the following occurrences: “the moderator interrupted a participant” or “the moderator was interrupted by a participant”. Clips were extracted from the video into 30-second duration segments. The mode of interruption that was related to the moderator was defined from an overlap in which the moderator was speaking, by examining the previous segment: if in this segment the moderator was speaking, then the moderator was interrupted by a participant; otherwise, the moderator interrupted a participant. Taking off the first segment of each clip, an interruption occurs at the beginning of each overlap segment; the total number of interruptions in the Train set is 1,353 split into 604 interruptions in LLC-Ovs and 749 interruptions in HLC-Ovs. The number of interruptions by the moderator (respectively, the interruptions of the moderator) was computed for the overlaps and their two categories (LLC and HLC) as well as its percentage of occurrence. We note that the moderator interrupted the participants more often than the moderator was interrupted by the participants (47.7% vs. 14.6%). Moreover, the moderator interrupted the participants more in the LLC-Ovs than in the HLC-Ovs (59.1 vs. 38.4%).

Table 4 Statistics on interruption mode occurrences of the moderator

Moderator – spk-050	Ov	LLC-Ov	HLC-Ov
# Interruptions	1,353	604	749
# Interruptions by the moderator	645 (47.7%)	357 (59.1%)	288 (38.4%)

and occurrence percentage

# Interruptions of the moderator and occurrence percentage	198 (14.6%)	127 (21.0%)	71 (9.4%)
---	-------------	-------------	-----------

3.1 Conflict Challenge

The Conflict challenge was one of the shared tasks that was organized during the Interspeech'2013 Computational Paralinguistics Challenge (Schuller et al. 2013), which took place from January 15 to May 24, 2013. The task consisted of an automatic analysis of the group discussions, to retrieve the conflicts. The goal of this competition was to bridge the gap between research in automatic conflict detection and the low compatibility of the results. The task data were split into the Train, Development and Test sets. The speaker dependence between these sets was reduced to a minimum that was needed in the real-life settings. Metadata are available only for the Train and Development sets. The participants did not have access to the labels of the test data, and all of the machine learning algorithms were based only on the training and development data. However, each participant could upload the instance predictions, to receive the confusion matrix and the results from the test data set up to 5 times. The official measure of the competition is the UAR. An official system of conflict detection was also provided with the following characteristics: the WEKA data mining tool kit was used as a framework for the classification task (Hall et al. 2009), and the Support Vector Machine classifier (SVM) with linear Kernel and Sequential Minimal Optimization (SMO) was used for learning; the official set of features (6,373 features), which is referred to as the IS-2013 set, was a representation of the utterances, and the complexity parameter of the SVM-classifier was optimized by using UAR on the Development set.

3.1 Audio feature sets

In this section, we describe the audio feature sets that we used for analyzing speech segments. This speech representation (Vogt et al. 2005, Schuller et al. 2008) is a new paradigm for speech analysis. It contrasts with the standard paradigm for speech analysis (the sequence of observation vectors): regardless of its duration, a speech utterance is represented by a large set of features, which is termed an audio feature set. The feature set is based on several Low-Level Descriptors (LLDs) that are computed from short overlapping windows of the audio signal. These LLDs comprise the loudness, the harmonics-to-noise ratio, the zero-crossing-rate, the spectral and prosodic coefficients, the formant positions and bandwidths, the duration of voiced/unvoiced speech segments, and features derived from the long-term average spectrum such as band-energies, roll-off, and centroid as well as voice quality features such as jitter and shimmer. Various global statistical functions (functionals) are computed on these LLDs to obtain feature vectors of equal size for each speech utterance. The sequence of LLDs that are associated with speech utterances can have different lengths, depending on the duration; the use of functionals allows us to obtain one feature vector per speech utterance, with a constant number of elements. It avoids the use of the expensive procedures of time warping between sequences of different lengths such as dynamic programming algorithms. Some functionals aim at estimating the spatial variability (e.g., mean, standard deviation, quartiles 1-3) and others aim at the temporal variability (e.g., peaks, linear regression slope). The four audio feature sets that we used for our experiments include the set of features that are provided by the organizers of the Interspeech 2010 (IS-2010) Paralinguistic Challenge (Schuller et al. 2010), the set of features for the Interspeech 2011 (IS-2011) Speaker State Challenge (Schuller et al. 2011), the set of features for the Interspeech 2012 (IS-2012) Speaker Trait Challenge (Schuller et al. 2012) and the set of features for the Interspeech 2013 (IS-2013) Conflict Sub-Challenge (Schuller et al. 2013). All of the features were extracted using the open source openSMILE feature extraction tools (Eyben et al. 2010). The IS-2010 feature set consists of 1,582 audio features, which were computed from 38

LLDs and 21 functionals. The spectral features include loudness, Mel-frequency cepstral coefficients, Mel-frequency band-energy, and line spectral pair frequencies. The prosodic and voice quality features comprise the pitch frequency and envelope, jitter and shimmer. Functionals such as the mean, standard deviation, kurtosis, skewness, minimum and maximum value, relative position, linear regression coefficients, quartile and percentile coefficients were applied on the LLDs. The IS-2011 feature set consists of 4,368 audio features, which were computed from 59 LLDs and 39 functionals. Additional LLDs, such as the auditory spectrum-derived loudness measure, RASTA-style filtered auditory spectra and statistical spectral descriptors (such as flux, entropy, variance) have been introduced. Additional functionals, such as quadratic regression and linear predictive coefficients and peaks distances allowed a better estimation of the temporal variability. The IS-2012 feature set consists of 6,124 audio features, which were computed from 64 LLDs and 40 functionals. Few LLDs have been added, including the logarithmic harmonic-to-noise ratio, spectral harmonicity and psychoacoustic spectral sharpness. Functionals that are related to the local extrema, such as the statistics of inter-maxima distances, have been introduced. Useless functionals have been removed to limit the number of the audio features. The IS-2013 feature set consists of 6,373 audio features, computed from 59 LLDs and 48 functionals. A total of 724 audio features were removed from the IS-2012 feature set, and 972 were added. New functionals that were related to the local extrema, such as the modeling of inter-maxima, have been introduced.

Table 5 summarizes the main characteristics of the used feature sets. The first three feature sets were used for the detection of overlap, and the last feature set was the official feature set for the detection of conflict.

Table 5 Official feature sets of Interspeech Challenges

Feature set	IS-2010	IS-2011	IS-2012	IS-2013
# LLDs	38	59	64	59

# functional	21	39	40	48
# features	1,582	4,368	6,124	6,373

4 Interruption Detection

From the previous statistics analyzed in Section 2, the conflict level was shown to be highly correlated to the mean number of interruptions (cf. Figure 2), the mean duration of overlap (cf. Figure 3) and the percentage of overlap duration (cf. Figure 4). Detecting segments of overlap is a difficult problem without individual microphones (Yamamoto et al. 2005). The main problem is due to the non-stationary characteristics of the speech signal. An alternative approach is the use of a microphone array (Quinlan et al. 2007). In this case, the estimation of the number of signal sources allows the detection of segments that contain more than one source of speech. Another approach, which is applied for improving the speaker diarization system, is the speech segmentation by a three-class Hidden Markov Model (Boakye et al. 2008), with the three classes corresponding to Non-Speech, Speech, and Overlapping Speech. Mel-Frequency-Cepstrum Coefficients (MFCC), Root Mean Square (RMS) energy and Linear Predictive Coding (LPC) residual energy features have been used, and they provided a precision of 66% and a recall of 26%. In our approach, we have chosen to develop a multi-resolution framework to estimate the overlap duration percentage. This approach is based on the fusion of various overlap detectors, in which each detector is estimated on the segments of a fixed and chosen duration.

4.1 Clip segmentation and relabeling

The clips were segmented into consecutive audio segments. Three segment durations were chosen for the multi-resolution: one, two and five seconds. For a given duration of segment, two segment-based detectors

were designed: (1) the first detector is a two-class classifier that is referred to as an $\{N, O\}$ -detector; it classifies a segment into Non-Ov (N) or Ov (O), and (2) the second detector is a three-class classifier that is referred to as an $\{N, L, H\}$ -detector, which classifies a segment into Non-ov (N), LLC-Ov (L) or HLC-Ov (H). Then, for multi-resolution detection, six SVM-based overlap detectors have been developed: (1) three two-class SVM classifiers, which we called $\{N, O\}_1$, $\{N, O\}_2$ and $\{N, O\}_5$, for the three durations, and (2) three three-class SVM classifiers, which we called $\{N, L, H\}_1$, $\{N, L, H\}_2$ and $\{N, L, H\}_5$, for the three durations. These labels (N, O, H and L) were computed from the SSPNet corpus metadata using speaker segmentation and conflict metadata. The Train and Development sets were relabeled using the multi-resolution framework of overlap localization. For each clip, diarization and conflict information are now represented by 102 labels: 60 labels for $\{N, O\}_1$ and $\{N, L, H\}_1$, 30 labels for $\{N, O\}_2$ and $\{N, L, H\}_2$, and 12 labels for $\{N, O\}_5$ and $\{N, L, H\}_5$. These new labels will be used for the training and testing of the various overlap detectors.

In Figures 5 and 6, the row called *Time* gives the time in seconds in the range from 1 to 30 (i.e., the clip duration), and the row *Segmentation* is the representation of the diarization metadata of the clip: N-segments are colored in white, L-segments in grey and H-segments in black. The other rows contain the relabeling according to the various detectors. For the three rows $\{N, O\}_x$ ($x \in \{1, 2, 5\}$), a segment is labeled O when it contains a part of overlap and, otherwise, N. For the three rows $\{N, L, H\}_x$ ($x \in \{1, 2, 5\}$), overlap segments are labeled according to the conflict level of the clip: L for LLC-Ov and H for HLC-Ov.

Figure 5 gives an instance of metadata relabeling for the LLC clip #Train_0001. For this clip, an LLC-Ov occurs over 13.01 and 14.4 s. The relabeling is O for the segments 14 and 15 of $\{N, O\}_1$, the segments 7 and 8 of $\{N, O\}_2$ and the segment 3 of $\{N, O\}_5$. The relabeling is L for the segments 14 and 15 of $\{N, L, H\}_1$, the segments 7 and 8 of $\{N, L, H\}_2$ and the segment 3 of $\{N, L, H\}_5$.

Clip #Train_0001 - Conflict score -7.2 - Low-Level conflict																																
{N, L, H}_5	N					N					L					N					N					N						
{N, L, H}_2	N	N	N	N	N	N	N	L	L	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N		
{N, L, H}_1	N	N	N	N	N	N	N	N	N	N	N	N	L	L	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	
{N, O}_5	N					N					O					N					N					N						
{N, O}_2	N	N	N	N	N	N	N	N	N	N	O	O	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	
{N, O}_1	N	N	N	N	N	N	N	N	N	N	N	N	O	O	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N
Segmentation																																
Time (s)	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30		

Fig. 5 Train set relabeling for the Train_0001 clip of the Low-Level Conflict

Figure 6 gives an instance of metadata relabeling for the HLC clip #Train_0006. For this clip, HLC-Ovs occur over 14.9 and 18.9 s and over 28.3 and 30 s. The relabeling is O for the segments 15, 16, 17, 18, 19, 29 and 30 of {N, O}_1, for the segments 8, 9, 10 and 15 of {N, O}_2 and for the segments 3, 4 and 6 of {N, O}_5. The relabeling is H for the segments 15, 16, 17, 18, 19, 29 and 30 of {N, L, H}_1, for the segments 7, 8, 9, 10 and 15 of {N, L, H}_2 and for the segments 3, 4 and 6 of {N, L, H}_5.

Clip #Train_0006 - Conflict score 7.3 - High-Level conflict																															
{N, L, H}_5	N					N					H					H					N					H					
{N, L, H}_2	N	N	N	N	N	N	N	H	H	H	H	N	N	N	N	N	H														
{N, L, H}_1	N	N	N	N	N	N	N	N	N	N	N	H	H	H	H	N	N	N	N	N	N	N	N	N	N	N	H	H			
{N, O}_5	N					N					O					O					N					O					
{N, O}_2	N	N	N	N	N	N	N	N	N	N	O	O	O	O	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	O	
{N, O}_1	N	N	N	N	N	N	N	N	N	N	N	N	N	O	O	O	O	O	N	N	N	N	N	N	N	N	N	N	N	O	O
Segmentation																															
Time (s)	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	

Fig. 6 Train set relabeling for the Train_0006 clip of High-Level Conflict

4.2 Two-class {N, O} classifier

Using relabeling, three two-class SVMs ({N, O}_1, {N, O}_2, {N, O}_5) were estimated on the Train set. Each SVM classifies a segment of a given duration (1 s, 2 s and 5 s) into Overlap (O) or Non-Ov (N). To account for the imbalanced class distribution, the upper-represented category (N) was down-sampled by a given factor. A factor of four was applied for the {N, O}_1 detector, a factor of three for the {N, O}_2 detector and a factor of two for the {N, O}_5 detector. We investigated the effects of different feature sets on the accuracy of the overlap speech detection. Table 6

gives the accuracy rates (N-Acc. and O-Acc. in %) of the two-class classifiers on the two classes (N and O) of the Development set. Three audio feature sets were compared: IS-2010, IS-2011 and IS-2012 (cf. § 3.1 Section 3). For all two-class classifiers, IS-2010 was the best feature set, with a UAR (in %) of slightly over 80% for the three detectors.

Table 6 Accuracy rates of the detectors {N, O} on the Development set according to the feature sets

Detectors {N, O}	Feature set	N-Acc. (%)	O-Acc. (%)	UAR (%)
{N, O}_1	IS-2010	86.7	73.9	80.3
{N, O}_1	IS-2011	87.7	72.3	80.0
{N, O}_1	IS-2012	87.8	71.6	79.7
{N, O}_2	IS-2010	85.1	75.1	80.1
{N, O}_2	IS-2011	87.3	71.6	79.5
{N, O}_2	IS-2012	87.4	71.7	79.5
{N, O}_5	IS-2010	82.7	78.7	80.7
{N, O}_5	IS-2011	84.9	75.3	80.1
{N, O}_5	IS-2012	84.0	75.7	79.8

Table 7 gives the accuracy rates of the detectors {N, O}_1, {N, O}_2 and {N, O}_5 on the HLC and LLC clips of the Development set for the IS-2010 audio-feature set. The difference in the UAR between the HLC and LLC clips varies according to the detectors: 3.2% for {N, O}_1, 3.8% for {N, O}_2 and 2.9% for {N, O}_5. The best results are always obtained on the HLC clips. Several explanations are possible: an unbalanced class distribution (e.g., for the {N, O}_1 detector, 2,910 HLC-Ov instances vs. 1,833 LLC-Ov instances) or the heterogeneousness of the LLC-Ov class.

Table 7 Accuracy rates of the detectors {N, O} on the HLC and LLC clips of the Development set

Detectors {N, O}	Clip set	N-Acc. (%)	O-Acc.(%)	UAR (%)
{N, O}_1	HLC	78.6	80.8	78.7
{N, O}_1	LLC	93.0	58.3	75.5
{N, O}_2	HLC	74.1	83.3	78.7
{N, O}_2	LLC	92.9	56.9	74.9
{N, O}_5	HLC	69.9	87.2	78.6
{N, O}_5	LLC	90.6	60.9	75.7

4.3 Three-class {N,L,H} classifier

Previous studies presented different typologies of overlaps: overlap and backchannel with overlap (Gravano et al. 2011), competitive and collaborative overlaps (Oertel et al. 2012). A backchannel indicates that the speaker producing them follows and understands the other speaker. They are generally words, onomatopoeias or other sounds produced in the background (Clancy et al. 1996). Collaborative or competitive interruptions are manifested by speech overlap, but only overlap from a competitive interruption can it be related to a conflict (Kurtié et al. 2012). In competitive overlaps, the incoming speaker attempts to forcefully take over the turn. In collaborative overlaps, the incoming speaker assists the current speaker in his or her speech. We chose to build classes of LLC-Ovs and HLC-Ovs by making the hypothesis that they would be separable acoustically and useful for conflict detection. This choice is supported by the observation that some of the LLC-Ovs of the Train set were backchannel with overlaps or/and collaborative overlaps.

Using relabeling, three three-class SVM classifiers ({N, L, H}_1, {N, L, H}_2 and {N, L, H}_5) were estimated on the Train set. Each SVM classifies a segment of a given duration (1 s, 2 s and 5 s) into an H, L or N. To account

for the imbalanced class distribution, the upper-represented category (N) was down-sampled by a given factor. A height factor was applied for the {N, L, H}_1 detector, which was a factor of six for the {N, L, H}_2 detector and a factor of three for the {N, L, H}_5 detector. We investigated the effects of different feature sets on the accuracy rate of the overlap speech detection. Table 8 gives the accuracy rates of three-class classifiers on the Development set. Three audio feature sets were compared: IS-2010, IS-2011 and IS-2012. IS-2010 was the best feature set for {N, L, H}_1, having a UAR of 61.1%. IS-2011 was the best feature set for {N, L, H}_2, with a UAR of 61.3%. IS-2010 was the best feature set for {N, L, H}_5, with a UAR of 63.5%. The LLC-Ovs are more difficult to detect than the HLC-Ovs. Furthermore, the detection rate of the LLC-Ovs appears to decrease with the duration of the analyzed segment: 44.7% for {N, L, H}_5 (5 s), 35.9% for {N, L, H}_2 (2 s) and 31.7% for {N, L, H}_1 (1 s). A possible explanation would be that the detector {N, L, H}_5 allows a better estimation of the overlap durations than the other detectors and, consequently, a better discrimination of the LLC- and HLC-Ovs. Indeed, the duration of the LLC-Ovs is lower on average than the HLC-Ovs (1.98 s vs. 2.75 s). For further experiments, we chose to use only the best three-class classifier: {N, L, H}_5 with the IS-2010 audio feature set.

Table 8 Accuracy rates of the detectors {N, L, H} on the Development set according to the feature sets

Detectors	Feature set	N-Acc. (%)	L-Acc. (%)	H-Acc. (%)	UAR (%)
{N, L, H}_1	IS-2010	78.0	31.7	73.5	61.1
{N, L, H}_1	IS-2011	79.9	32.7	70.5	61.0
{N, L, H}_1	IS-2012	79.4	31.4	71.4	60.7
{N, L, H}_2	IS-2010	79.5	32.6	71.2	61.2
{N, L, H}_2	IS-2011	78.1	35.9	70.0	61.3
{N, L, H}_2	IS-2012	80.5	31.5	68.0	60.0
{N, L, H}_5	IS-2010	77.5	44.7	68.3	63.5

{N, L, H}_5	IS-2011	76.4	40.0	67.7	61.4
{N, L, H}_5	IS-2012	80.8	38.2	67.4	62.1

The detector {N, L, H}_5 classifies the segments of the Development set into the three classes (N, L and H). In table 9, we investigate the speech duration of the moderator for the segments labeled with the three classes. As previously quoted in Section 2 for similar experiments that are related to the moderator on the Train set, the total speech duration was estimated by the duration of the segment in which a lonely subject is speaking plus twice the segment duration in which two subjects are speaking. The speech duration of the moderator (in s) was computed for the various classes of speech (N, L and H). The percentage of the moderator speech duration was computed as the ratio of the speech duration of the moderator to the total speech duration. We note that the moderator spoke more during the Ls than the Hs (27.1% vs. 11.0%) and the Ns (27.1% vs. 18.1%). These results confirm those obtained in Table 3 (cf. Section 2) from the Train set. An interpretation of the result similarity could be a possible relation between the interruptions caused by the moderator and a Low Level of Conflict. We have to notice that the moderator (female) of the Train set is different than the moderator (male) of the Development set.

Table 9 Statistics on the speech duration of the moderator (Partition of the Development clips by the {N, L, H}_5 detector)

Detector {N, L, H}_5	N	L	H
Total speech duration (in s)	4,284.5	1,561.2	2,350.1
Speech duration of the moderator (in s)	775.3	423.4	258.2
Percentage of the moderator speech duration (in %)	18.1%	27.1%	11.0%

4.4 Audio characteristics of overlaps

Previous studies (Smolenski et al 2011, Shokouhi et al. 2013) have shown that the audio characteristics of overlapping speech are different from speech in which a lonely speaker occurs. We looked for the discriminating cues (1) between Ov and Non-Ov and (2) more specifically between HLC-Ov and LLC-Ov. For these investigations, we chose to study the segments that had a 5-second duration in the Train set for the best accuracy results of the 5-second-based {N, O} and {N, L, H} detectors (see, respectively, Tables 7 and 8 in Section 4). The 38 Low-Level Descriptors (LLD) of the IS-2010 feature set have been used as audio characteristics. The relevance of the LLD was analyzed with respect to the classes Non-Ov/Ov, which are referred to as {N, O}, and the HLC-Ov/LLC-Ov, which are referred to as {H, L}. For each LLD, the relevance is given by the information gain (Raubert et al. 1993), which is computed on the segments of 5 s duration with the following formula: $H(\text{class}) - H(\text{class}/\text{LLD})$, where H is the Shannon entropy. Four steps were defined to compute the entropy: (1) filtering of the IS-2010 features according to a given LLD, (2) clustering of the segments of the Train set using the filtered features, (3) computation of the contingency table from the class and the cluster associated with each segment and (4) estimation of the entropy from the table of contingency. Table 10 gives the information gain computed on the Train set of the five best-ranked LLDs (over 38 LLDs) in discriminating LLC-Ovs and HLC-Ovs. The most relevant LLDs are the logarithmic powers of Mel-frequency bands and, more precisely, the high-frequency bands and the normalized loudness. These results show that various acoustic differences exist between the two types of overlaps.

Table 10 Information gain of the five best-ranked LLDs of the IS-2010 audio feature set in discriminating LLC-Ovs and HLC-Ovs

Low-Level Descriptors (LLD)	Inf. Gain	Rank (/38)
log power [3934 Hz -5649 Hz]	0,130	1
log power [2682 Hz -3934 Hz]	0,119	2

log power [1768 Hz -2682 Hz]	0,107	3
normalized loudness	0.102	4
log power [5649 Hz -8000 Hz]	0,102	5

Table 11 gives the information gain that is computed on the Train set of the five best-ranked LLDs (over 38 LLDs) in discriminating Ovs and Non-Ovs. According to the information gain rank, the most relevant LLDs are the fundamental frequency, the logarithmic powers, especially in low-frequency bands, the jitter and the first Mel Frequency Cepstral coefficient. The usual representation techniques and algorithms are designed and interpreted for speech signals in which a lonely subject is speaking. In the case of overlapping speech in which two or more subjects are speaking, the usual algorithms are not adapted (e.g., the pitch algorithm); the computation of one fundamental frequency has no sense, and its computation was shown to be the most discriminant cue for detecting Ov/Non-Ov. For a speech representation such as the logarithmic power in the Mel frequency bands, the low frequency bands in which the first two formants of the speaker occur were also shown to be discriminant. Last, the jitter related to the pitch and the first Mel Frequency Cepstrum Coefficient related to the energy of the segment were also shown to be relevant for the discrimination Ov/Non-Ov.

Table 11 Information gain of the five best-ranked LLDs of the IS-2010 audio feature set in discriminating Ovs and Non-Ovs

Low-Level Descriptors (LLD)	Inf. Gain	Rank (/38)
fundamental frequency (F0)	0,141	1
log power [614 Hz -1101 Hz]	0,129	2
log power [0 Hz -259 Hz]	0,127	3
jitter (DDP)	0,124	4
first mel frequency cepstral coef.	0,121	5

5 Conflict Detector

Overlap Detectors have been developed and assessed, to incorporate their knowledge in an improved Conflict Detector (Conflict/Non-Conflict). Incorporating prior knowledge (Krupka and Tishby 2007, Li et al. 2008) in classification systems allowed an increase in the performance in many applications of pattern recognition (e.g., biomedical image, pathological voice). Various methods have been developed for Neural Network systems (Chen et al. 2000) and SVM-classifiers (Decoste and Scholkopf 2002, Lauer and Bloch 2008). As defined by Schölkopf and Smola (2001), the methods developed include prior knowledge in an SVM-classifier and can be divided into three categories: (1) the kernel methods with selection of the most appropriate kernel or the creation of a new kernel, (2) the optimization methods with the addition of constraints, and (3) the sample methods with data generation or modification of data representations,. We have chosen the last category by developing an SVM-based detector, using as input a composite feature set. This feature set is a concatenation of selected audio features and posterior-based features that are computed from the posterior probabilities of the Overlap Detectors. The architecture characteristics of this classification system are close to those used in a mixture of experts (Jordan and Jacobs 1994). These approaches have theoretical advantages, such as a reduction in the hypothesis space and learning consistency. As described in Figure 7, the multi-expert architecture scheme of the Conflict Detector has consisted of a set of Overlap Detectors (e.g., X, Y) and a Conflict/Non-Conflict Detector (C). A specialized audio-feature set (e.g., X-Feat. set) was associated with each Overlap Detector (e.g., X), to represent the utterances. A Conflict audio-Feature set (Cf-Feat. set) was associated with the Conflict Detector. This feature set consisted of the selection of the relevant features (Feat. Select.) that were extracted from the Overlap Feature set (Ov-Feat. set) and the IS-2013 feature set (cf. § 3.1 Section 3). A set of Functionals (Funct.) was applied to the posterior probabilities of the Overlap Detectors (e.g., X-Post and Y-Post) to obtain the Ov-Feat. set.

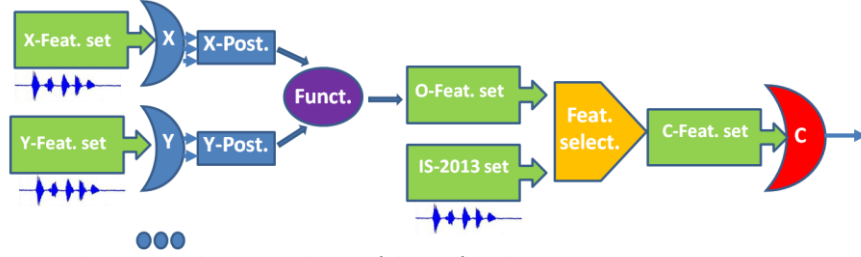


Fig. 7 Multi-expert architecture scheme of the conflict detector

We chose the Overlap Detectors giving the best UAR on the Development set (cf. Table 6 §4.2. and Table 8 §4.3 in Section 4): three two-class (Non-Ov/Ov) SVM-based detectors ($\{N, O\}_1$, $\{N, O\}_2$, $\{N, O\}_5$) and one three-class (Non-Ov/LLC-Ov/HLC-Ov) SVM-based detector ($\{N, L, H\}_5$).

5.1 Posterior probabilities

Logistic regression models (Hosmer and Lemeshow 2000) were used to obtain the posterior probabilities from the four Overlap Detectors ($\{N, O\}_1$, $\{N, O\}_2$, $\{N, O\}_5$ and $\{N, L, H\}_5$). These posterior probabilities of the Overlap Detectors provide information about the uncertainty of belonging to one class: for example, the probability of 60% of a segment to be an overlap and 40% to be a non-overlap. There are various strategies for computing these probabilities, such as Platt's method (Platt 2000), isotonic regression (Zadrozny and Elkan 2002) and Bayesian methods (Sollich 2002). These probabilities are useful to integrate expert classifiers such as overlap classifiers in a global decision process. This approach is a flexible architecture for making decisions without global optimization. The method of computation of the posterior probabilities depends on the chosen set of clips. The goal is to obtain a consistent computation of the posterior probabilities from the different corpora (Train, Development and Test sets). For the Train and Development sets, the posterior probabilities have been computed by performing cross-predictions on the union of these two sets. This process consists of splitting the data set into s disjoint folds and predicting class posterior probabilities of each instance of a fold from a model trained on the $s-1$ other

fold. Sixteen folds have been chosen that have participant independence between two folds. For the Test set, the posterior probabilities have been computed from a model trained on the union of the Train and Development sets. A total of 120 posterior probabilities were computed for each clip: 60 for the $\{N, O\}_1$ detector, 30 for $\{N, O\}_2$, 12 for $\{N, O\}_5$ and 18 for $\{N, L, H\}_5$.

Figure 8 and Figure 9 give an instance of the posterior probabilities from the four Overlap Detectors ($\{N, O\}_1$, $\{N, O\}_2$, $\{N, O\}_5$ and $\{N, L, H\}_5$) respectively for the LLC clip #Train_0001 and the HLC clip #Train_0006. The row called *Time* gives the time from 1 to 30 s (clip duration), the row *Segmentation* is the representation of the diarization metadata of the clip: N-segments are colored in white, L-segments in grey and H-segments in black. The other rows contain the posterior probabilities presented as a percentage. For the three rows $\{N, O\}_x$ ($x \in \{1, 2, 5\}$), a segment of posterior probabilities that was higher than 50 was detected as O; otherwise, it was detected as N. The posterior probabilities that were associated with the $\{N, L, H\}_5$ detector are presented in the three other rows $\{N, L, H\}_5$ (N), $\{N, L, H\}_5$ (L) and $\{N, L, H\}_5$ (H) for, respectively, Non-Overlap (N), Low-Level-Conflict (L) and High-Level-Conflict (H). For a given segment, the higher probability (in bold) corresponds to the class that was detected._

In Figure 8, the class O was detected for the segments 14 and 15 of $\{N, O\}_1$, the segments 8 and 9 of $\{N, O\}_2$, and the segments 3 and 4 of $\{N, O\}_5$. Class H was detected for the segment 3 for $\{N, L, H\}_5$. There are three wrong detections: the class O instead of N for the segment 9 of $\{N, O\}_2$ and the segment 4 of $\{N, O\}_5$, and the class H was detected instead of L for segment 3 for $\{N, L, H\}_5$.

Clip #Train_0001 - Conflict level -7.2 - Low-Level conflict																														
{N, L, H}_5 (N)	77				73				11				54				99				97									
{N, L, H}_5 (H)	00				01				51				21				00				02									
{N, L, H}_5 (L)	03				06				38				23				01				01									
{N, O}_5 (O)	01				08				60				80				01				01									
{N, O}_2 (O)	10	07	02	08	00	06	36	77	65	10	00	01	06	13	03															
{N, O}_1 (O)	03	08	01	01	04	01	03	04	06	01	04	10	10	76	55	40	19	24	03	29	02	01	03	02	35	03	14	04	02	03
Segmentation																														
Time (s)	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30

Fig. 8 Overlap posterior probabilities as percentages for the Train_0001 clip with Low-Level conflict

In Figure 9, the class O was detected for the segments 10, 11, 16, 17, 18, 19, 21, 22, 23, 28 and 30 of {N, O}_1, for the segments 5 and 8 through 15 of {N, O}_2, and for the segments 2, 4, 5 and 6 of {N, O}_5. Class H was detected for the segments 2, 4, 5 and 6 of {N, L, H}_5. There are 9 wrong detections: the class O instead of N for the segments 10, 11 and 28 of {N, O}_1, the segments 4, 13 and 14 of {N, O}_2, the segment 2 of {N, O}_5, the class N instead of O for the segment 29 of {N, O}_1 and the class H instead of N for the segment 2 of {N, L, H}_5. We note that there was no wrong decision for segments 21, 22 and 23 of {N, O}_1, for the segments 11 and 12 of {N, O}_2, and for segment 5 of {N, O}_5 and {N, L, H}_5; after listening, an overlap occurs effectively from 20.3 to 22.2 s but was not labeled in the metadata.

Clip #Train_0006 - Conflict level 7.3 - High level conflict																														
{N, L, H}_5 (N)	98				09				72				01				03				02									
{N, L, H}_5 (H)	01				84				26				99				94				89									
{N, L, H}_5 (L)	01				07				02				00				03				09									
{N, O}_5 (O)	03				62				12				99				95				99									
{N, O}_2 (O)	03	11	15	05	54	20	20	98	99	99	98	68	54	75	97															
{N, O}_1 (O)	01	38	35	05	02	07	04	13	12	89	68	17	19	06	02	99	99	99	99	16	87	99	58	19	05	46	21	82	29	99
Segmentation																														
Time (s)	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30

Fig. 9 Overlap posterior probabilities as percentages for the Train_0006 clip with High-Level conflict

5.2 Overlap feature sets

One hundred and twenty posterior probabilities are computed for each clip. These values depend on the time and represent the temporal shape

of a conflict in terms of the overlap. There are specific temporal shapes for conflict escalation (Kim et al. 2012c), but the 797 clips of the Train set are insufficient to model these temporal shapes. We have chosen to apply statistical functionals to the posterior probabilities; the purpose was to obtain an Overlap feature set that is related to the percentage of overlap duration. Three functionals have been chosen: mean, correlation, and covariance. The mean functional was applied to the posterior probabilities of $\{N, O\}_1$, $\{N, O\}_2$, $\{N, O\}_5$ for the class O and to the posteriors of $\{N, L, H\}_5$ for the classes N, L and H. The correlation functional was applied between the posterior probabilities of the class O for all combinations of $\{N, O\}_1$, $\{N, O\}_2$ and $\{N, O\}_5$. Table 10 gives a list of the ten features that were computed by the mean and correlation functionals.

Table 10 List of the features computed by the mean and correlation functionals

Mean and correlation functionals	Feature name
Mean (post ($\{N, O\}_1$ (O)))	O1
Mean (post ($\{N, O\}_2$ (O)))	O2
Mean (post ($\{N, O\}_5$ (O)))	O5
Correlation (post ($\{N, O\}_1$ (O)), post ($\{N, O\}_2$ (O)))	O12
Correlation (post ($\{N, O\}_1$ (O)), post ($\{N, O\}_5$ (O)))	O15
Correlation (post ($\{N, O\}_2$ (O)), post ($\{N, O\}_5$ (O)))	O25
Correlation (post ($\{N, O\}_1$ (O)), post ($\{N, O\}_2$ (O)), post ($\{N, O\}_5$ (O)))	O125
Mean (post ($\{N, L, H\}_5$ (N)))	N5
Mean (post ($\{N, L, H\}_5$ (L)))	L5
Mean (post ($\{N, L, H\}_5$ (H)))	H5

Functional Covariance is a functional of a functional. It was applied to the mean and correlation functionals. The interest of this functional is to reveal the co-factors. Two Overlap feature sets have been defined. The first feature set, called Ov-1, consisted of 28 features; it was computed by the

covariance functional applied to the features that are related to the {N, O} detectors (O1, O12, O15, O2, O25, and O125). The second feature set, called Ov-2, consisted of 55 features; it was computed by the covariance functional applied to the features that are related to the {N, O} and {N, L, H} detectors (O1, O12, O15, O2, O25, O125, N5, L5, and H5). These two feature sets will allow a contrastive test to measure the contribution of the {N, L, H}_5 detector in the detection of conflict. The method of information gain was used to analyze the feature relevance of the Ov-2 set in comparison with the IS-2013 set. Table 11 gives the information gain computed on the Train set and the rank on 6,428 features (55 features from the Ov-2 set and 6,373 features from the IS-2013 set) of the most relevant features for the conflict detection. The best feature is the Cov_O125_O125 feature (which is equal to O125 multiplied by O125). The 12th rank of the Cov_H5_O1 feature shows that the {N, L, H}_5 detector is relevant for the detection of conflict. A total of 36 out of 55 features of the Ov-2 set have better information gain than those of the IS-2013 set. These results show the interest of the Overlap feature sets for the detection of conflict.

Table 11 Information gain of the fifteen best-ranked LLDs of the audio feature set, including the Ov-2 feature set and the IS-2013 feature set

Features	Information Gain	Rank (/6,428)
Cov_O125_O125	0.43862	1
Cov_O12_O125	0.43758	2
Cov_O1_O12	0.43586	3
Cov_O1_O125	0.43177	4
Cov_O15_O125	0.42914	5
Cov_O12_O12	0.42858	6
Cov_O25_O125	0.41965	7
Cov_O12_O15	0.41957	8

Cov_O12_O25	0.41431	9
Cov_O15_O15	0.41429	10
Cov_O15_O25	0.41325	11
Cov_H5_O1	0.40915	12
Cov_H5_O15	0.40849	13
Cov_O1_O25	0.40705	14
Cov_H5_O12	0.40509	15

5.3 Conflict feature sets

From two initial feature sets (Ov-1 and IS-2013; and Ov-2 and IS-2013), two Conflict feature sets (Cf-1 and Cf-2) were selected by a backward selection algorithm when maximizing UAR on the Development set for the conflict detection task. Table 12 gives the characteristics of the Cf-1 and Cf-2 sets of the Conflict detector using these feature sets. The Cf-1 feature set consists of 315 features (15 features from the Ov-1 set and 300 features from the IS2013 set). The Cf-2 feature set consists of 335 features (45 features from Ov-2 set and 290 features from the IS2013-set).

Table 12 Characteristics of the Conflict feature sets

Feature set	Selected feature set	# of selected features	# selected feat. from Ov features	# of selected feat. from IS-2013
Ov-1 and IS-2013 (6,428 features)	Cf-1	315	15	300
Ov-2 and IS-2013 (6,401 features)	Cf-2	335	45	290

Table 13 gives the accuracy (UAR in %) of the conflict detection on the Development set using the various feature sets (IS-2013, Cf-1, Cf-2). The results show an improvement of 8.3% using the Cf-1 set and 9.2% using

the Cf-2 set on the Development set compared to the baseline results that use the IS-2013 set (UAR of 79.1%). These results show also that the majority of the features of the Cf-2 set are relevant and not redundant. It confirms that the two types of detectors ($\{N, L, H\}$ and $\{N, O\}$) are relevant for the detection of conflict.

Table 13 UAR in the Conflict detection task on the Development set according to the Conflict feature sets

Feature set	IS-2013	Cf-1	Cf-2
# features	6,373	315	335
UAR (Devel. Set)	79.1%	87.4%	88.3%

5.4 Conflict detectors

Two Conflict Detectors have been developed. Figure 10 resumes the architecture characteristics of the first Conflict Detector, called the Simple Overlap-based Conflict Detector (SO-Conflict Detector). This detector was based on a set of Overlap Detectors (1, 2 and 3) that correspond to the three multi-resolution-based $\{N, O\}$ Detectors and a Conflict Detector (4). The IS-2010 feature set (1,582 features) was used for the Overlap Detectors. The Cf-1 feature set (315 features) was associated with the Conflict Detector. The Cf-1 feature set was obtained by a backward selection algorithm from the Ov-1 feature set (28 features) and the IS-2103 set (6,373 features).

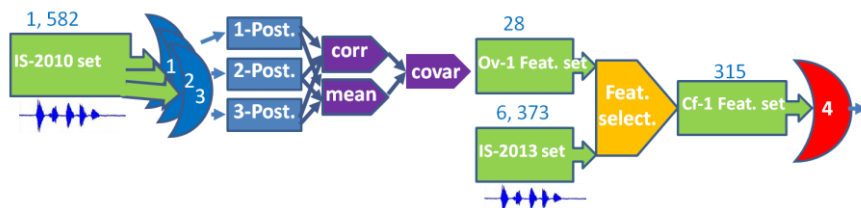


Fig. 10 Architecture scheme of the SO-Conflict Detector

Figure 11 resumes the architecture characteristics of the second Conflict Detector, called the Advanced Overlap-based Conflict Detector (AO-Conflict Detector). This detector was based on a set of Overlap Detectors (1, 2, 3 and 4) and a Conflict Detector (5). The IS-2010 audio-feature set (1,582 features) was used for the Overlap Detectors. The Cf-2 feature set (335 features) was associated with the Conflict Detector. The Cf-2 feature set was obtained by a backward selection algorithm from the Ov-2 feature set (55 features) and the IS-2103 set (6,373 features).

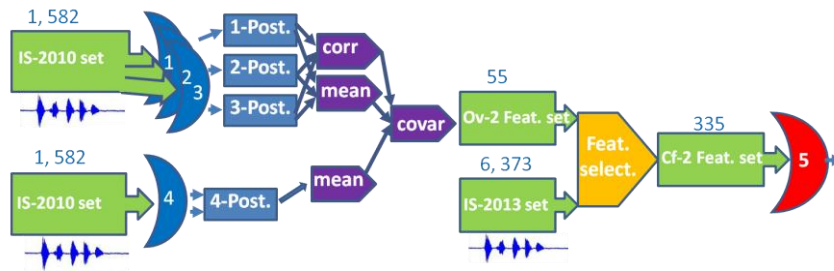


Fig. 11 Architecture scheme of the AO-Conflict Detector

5.5 Conflict detection on the Test set

The Test set of the Interspeech'2013 Conflict Challenge (Schuller et al. 2013) consisted of 397 clips with no information or metadata available. Table 14 gives the results obtained on the Test set during the Conflict Challenge. Experiments gave a UAR of 83.4% for the SO-Conflict Detector and a UAR of 85.3% for the AO-Conflict detector. These results show an improvement of 2.6% (SO-Conflict Detector) and 4.5% (AO-Conflict Detector) on the Test set compared to the baseline results with the IS-2013 set (UAR of 80.8%) for the Conflict Detection task. These results confirm also that the two types of overlap detectors ($\{N, L, H\}$ and $\{N, O\}$) are relevant for the detection of conflict. The other results are those obtained by the other participants. In (Grèzes et al. 2013), a UAR of 83.1% was obtained on the Test set using a unique feature: the percentage of overlap predicted by an SVM-based regression model. In (Rasanen et Pohjalainen 2013), a UAR of 83.9% was obtained on the Test set using 349 relevant features

selected from the IS-2013 feature set. Feature relevance was computed by a random process. We notice that the two better results were obtained by a similar number of features (335 vs. 349).

Table 14 Assessment on the Test set

Conflict Detector	# features	UAR (%) on Test set
SO-Conflict Detector	315	83.4%
AO-Conflict Detector	335	85.3%
IS-2013 Baseline system	6,373	80.8%
(Grèzes et al. 2013)	1	83.1%
(Räsänen and Pohjalainen 2013)	349	83.9%

6 Conclusions

This article presents and assesses a detection system of conflict in group discussions from voice analysis. The system was based on a multi-expert architecture and detected two states (Conflict/Non-Conflict). The analysis of the Train set of the SSPNet database has demonstrated that the conflict level was highly correlated with the mean number of interruptions, the mean duration of overlap and the percentage of overlap duration. The multi-expert architecture enabled knowledge regarding overlaps to be used in the Conflict Detector.

The concept of LLC-Ovs and HLC-Ovs has been introduced and investigated. Two types of Overlap Detectors have been developed: the first type aims at detecting whether a speech segment contains overlap, and the second type aims at detecting whether a speech segment contains an LLC-Ov or HLC-Ov. The accuracy of the detectors shows that the LLC-Ovs and HLC-Ovs can be modeled. The high-frequency Mel bands and the normalized loudness are shown to be the audio characteristics that are relevant to discriminating these two types of overlap. A multi-resolution

framework has been developed for the Overlap Detectors, to improve the robustness of the detection. Three segment durations have been chosen (1 s, 2 s and 5 s). The experiments have shown that these detectors were not redundant.

A composite set of 335 features, which consist of audio-based features and overlap detector-based features, has been defined for the Conflict Detection task of the Conflict Interspeech'2013 challenge. The performance obtained for the Test set gave a UAR of 85.3%. These results show an improvement of 4.5% compared to the results of the baseline System of the Conflict challenge (UAR of 80.8%).

These experiments have shown the capability of a multi-expert architecture to integrate a piece of conflict knowledge. Other knowledge that is related to the turn-taking patterns, such as the modeling of the moderator role (Vinciarelli 2007), or that is related to the non-verbal interactions, such as the movements of the body, the head and the arms, could be integrated into the Conflict Detector.

Acknowledgments Many thanks to Björn Schuller (TUM, Germany), Stefan Steidl (FAU Erlangen-Nuremberg, Germany), and Anton Batliner (TUM, Germany) for the organization of the Interspeech'2013 Conflict Challenge and special thanks to Alessandro Vinciarelli (University of Glasgow, UK) for the SSPNet Conflict Corpus.

References

- Atkinson, J.M., Drew, P., (1979). *Order in Court: The Organisation of Verbal Interaction in Judicial Settings*. Atlantic Highlands, NJ: Humanities Press.
- Barzilay, R., Collins, M., Hirschberg, J., Whittaker, S. (2000). The rules behind the roles: identifying speaker roles in radio broadcasts, paper presented at 17th National Conference on Artificial Intelligence, Austin, USA, 30 July– 3 August, 679–684.
- Beattie, G.W., (1982). Turn-taking and interruption in political interviews: Margaret Thatcher and Jim Callaghan compared and contrasted. *Semiotica*, 39(1-2), 93–114.
- Blake, R.R. and Mouton, J.S., 1964. *The Managerial Grid*. Houston: Gulf Publishing.
- Boakye, K., Trueba-Hornero, B., Vinyals, O. and Friedland, G., (2008). Overlapped Speech Detection for Improved Diarization in Multi-Party Meetings, paper presented at ICASSP Conference, Las Vegas, USA, 31 March-4 April, 4353–4356.
- Boden, D., (1994). *The Business of Talk. Organizations in Action*. London: Polity Press.
- Brinson, S.L., and Winn, J.E., (1997). Talk shows' representations of interpersonal conflicts. *Journal of Broadcasting and Electronic Media*, 41(1), 25–39.

- Chen, Z., Feng, T.J. and Houkes, Z., (2000). Incorporating a priori knowledge into initialized weights for neural classifier. paper presented at International Joint Conference on Neural Networks (IJCNN), Como, Italy, 24-27 July, 291-296.
- Clancy, P.M., S.A. Thompson, R. Suzuki and H. Tao, (1996). The conversational use of reactive tokens in English, Japanese and Mandarin. *Journal of Pragmatics*, 26, 355-387.
- Daly, T.M., Lee, J.A., Soutar, G. N. and Rasmi, S., (2010). Conflict-handling style measurement: a best-worst scaling application. *International Journal of Conflict Management*, 21(3), 281-308.
- De Ruiter, J.P., Mitterer, H. and Enfield, N. J., (2006). Projecting the end of a speaker's turn: A cognitive cornerstone of conversation. *Language*, 82(3), 515-535.
- Decoste, D. and Scholkopf, B., (2002). Training invariant support vector machines. *Machine Learning*, 46(1-3), 161-190.
- Euwema, M.C. and Van de Vliert, E. (1990). Behaviour and the escalation in hierarchical conflicts. *Toegepaste Sociale Psychologie*. 4, 28 - 41.
- Euwema, M.C., Van de Vliert, E., and Bakker, A.B., (2003). Substantive and relational effectiveness of organizational conflict behavior, *International Journal for Conflict Management*, 14, 119 -139.
- Eyben, F., Wöllmer, M. and Schuller, B., (2010). openSMILE The Munich Versatile and Fast Open-Source Audio Feature Extractor. paper presented at the ACM Multimedia Conference (MM), Florence, Italy, 25-29 October, 1459-1462.
- Finn, A. and Louviere, J.J., (1992). Determining the appropriate response to evidence of public concern: the case of food safety. *Journal of Public Policy and Marketing*, 11, 19-25.
- Garcia, A., (1991). Dispute resolution without disputing: how the interactional organization of mediation hearings minimizes argumentative talk. *American Sociological Review*, 56, 818-835.
- Gravano, A. and Hirschberg, J., (2011). Turn-taking cues in task oriented dialogue, *Computer Speech and Language*, 25(3), 601-634.
- Grèzes, F., Richards, J. and Rosenberg A., (2013). Let Me Finish: Automatic Conflict Detection Using Speaker Overlap. paper presented at the Interspeech Conference, Lyon, France, 25-29 August, 5 pages.
- Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P. and Witten, I., (2009). The WEKA Data Mining Software: An Update. *SIGKDD Explorations*, 11, 10-18.
- Hammer, M.R., (2005). The Intercultural Conflict Style Inventory: A conceptual framework and measure of intercultural conflict resolution approaches. *International Journal of Intercultural Relations*, 29(6), 675-695.
- Heath, C., Luff, P., (2007). Ordering competition: the interactional accomplishment of the sale of art and antiques at auction. *British Journal of Sociology* 58, 63-85
- Hosmer, D.W. and Lemeshow, S., (2000). *Applied Logistic Regression*. 2nd edition, New York: Wiley.
- Jordan, M.I. and Jacobs, R.A., (1994). Hierarchical mixtures of experts and the EM algorithm. *Neural Computation*, 6, 181-214.
- Karatzoglou, A., Smola, A., Hornik, K. and Zeileis, A., (2004). Kernlab: An S4 package for kernel methods. *Journal of Statistical Software*, 11(9), 1-20. <http://www.jstatsoft.org/v11/i09/>
- Kim, S., Filippone, M., Valente, F. and Vinciarelli, A., (2012a). Predicting the conflict level in television political debates: an approach based on crowdsourcing, nonverbal communication and gaussian processes. paper presented at the ACM Conference on Multimedia. Nara, Japan, 793-796.
- Kim, S., Valente, F. and Vinciarelli, A., (2012b). Automatic Detection of Conflicts in Spoken Conversations: Ratings and Analysis of Broadcast Political Debates. paper presented at ICASSP, Kyoto, Japan, 25-30 March, 5089-5092.

- Kim, S., Yella, S.H., and Valente, F.A., (2012c). Automatic Detection of Conflict Escalation in Spoken Conversations. paper presented at Interspeech Conference, Portland, USA, OR, 9-13 September, 4 pages.
- Korabik, K., Baril, G.L. and Watson, C., (1993). Managers' conflict management style and leadership effectiveness: The moderating effects of gender. *Sex Roles*, 29(5-6), 405-418.
- Krupka, E. and Tishby, N., (2007). Incorporating Prior Knowledge on Features into Learning. *Journal of Machine Learning Research*, 227-234.
- Kurtié, E., Brown, G.J. and Wells, B., (2012). Resources for turn competition in overlapping talk. *Speech Communication*, 55, 1-23, doi:10.1016/j.specom.2012.10.002.
- Lauer, D.F. and Bloch, G., (2008). Incorporating prior knowledge in support vector machines for classification: a review. *Neurocomputing*, 71(7-9), 1578-1594.
- Lerner, G.H., (1995). Turn design and the organization of participation in instructional activities. *Discourse Processes*, 19 (1), 111-131.
- Li, Y., de Ridder, D., Duin, R.P.W. and Reinders, M.J.T., (2008). Integration of prior knowledge of measurement noise in kernel density classification. *Pattern Recognition*, 41, 320-330.
- Macintosh, G. and Stevens, C.J., (2008). Personality, motives and conflict strategies in everyday service encounters. *International Journal of Conflict Management*, 19(2), 112-131.
- Mac Houl, A., (1978). The organization of turns at formal talk in the classroom. *Language in Society*, 7, 183-213.
- Mehan, H., The structure of classroom discourse (1985). In: Dijk, T.A. (Ed.), *Handbook of Discourse Analysis*, 3:120-131. New York: Academic Press.
- Mondada, L., (2012). The dynamics of embodied participation and language choice in multilingual meetings. *Language in Society*, 41, 1-23.
- Mondada, L., (2013). Embodied and spatial resources for turn-taking in institutional multi-party interactions: Participatory democracy debates. *Journal of Pragmatics*, 46(1), 39-68.
- Oertel, C., Włodarczak, M., Tarasov, A., Campbell, N. and Wagner, P., (2012). Context cues for classification of competitive and collaborative overlaps. paper presented at Speech Prosody Conference, Shanghai, China, 22-25 May, 4 pages.
- Oetzel, J.G., Ting-Toomey, S., Yokochi, Y., Masumoto, T., and Takai, J. (2000). A typology of facework behaviors in conflicts with best friends and relative strangers. *Communication Quarterly*, 4, 397-419.
- Platt, J.C., (2000). Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In: P. J. Bartlett, B. Schölkopf, D. Schuurmans, A. J. Smola (Eds.), *Advances in Large Margin Classifiers*, Cambridge: MIT Press, 61-74.
- Quinlan A. and Asano, F., (2007). Detection of overlapping speech in meeting recordings using the modified exponential fitting test. paper presented at the European Signal Processing Conference, Poznan, Poland, 3-7 September, 2360-2364.
- Rahim, M.A., (1983). A measure of styles of handling interpersonal conflict. *The Academy of Management Journal*, 26(2), 368-376.
- Räsänen, O. and Pohjalainen, J., (2013). Random Subset Feature Selection in Automatic Recognition of Developmental Disorders, Affective States, and Level of Conflict from Speech. paper presented at the Interspeech Conference, Lyon, France, 25-29 August, 5 pages.
- Rauber, T.W., Steiger-Garcia, A.S., (1993). Feature selection of categorical attributes based on contingency table analysis. paper presented at the Portuguese Conference on Pattern Recognition, Porto, Portugal.
- Sacks, H., Schegloff, E. A. and Jefferson, G., (1974). A simplest systematics for the organization of turn-taking for conversation. *Language*, 50(4), 696-735.
- Schegloff, E.A., (1987). Between macro and micro: contexts and other connections. In: Alexander, J., Giesen, B., Munch, R., Smelser, N. (Eds.), *The Micro-Macro Link*. Berkeley: University of California Press, 207-234.
- Schölkopf, B.A.J. and Smola, A.J., (2001). *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond.*, Cambridge, MA, USA: MIT Press.

- Schuller, B., Wimmer, M., Moesenlechner, L., Kern, C., Arsic, D. and Rigoll, G., 2008. "Brute-forcing Hierarchical Functional for Paralinguistics: A Waste of Feature Space?". In: Proceedings of ICASSP, pp. 4501–4504.
- Schuller, B., Steidl, S., Batliner, A., Burkhardt, F., Devillers, L., Müller, C. and Narayanan, S., (2010). The Interspeech 2010 paralinguistic challenge. paper presented at the Interspeech Conference, Makuhari, Japan, 26-30 September, 2794–2797.
- Schuller, B., Batliner, A., Steidl, S., Schiel, F. and Krajewski, J., (2011). The Interspeech 2011 Speaker State Challenge". paper presented at the Interspeech Conference, Florence, Italy, 28-31 August, 4 pages.
- Schuller, B., Steidl, S., Batliner, A., Noth, E., Vinciarelli, A., Burkhardt, F., van Son, R., Weninger, F., Eyben, F., Bocklet, T., Mohammadi, G. and Weiss, B., (2012). "The Interspeech 2012 Speaker Trait Challenge. paper presented at the Interspeech Conference. Portland, OR, USA, 9-13 September, 4 pages.
- Schuller, B., Steidl, S., Batliner, A., Vinciarelli, A., Scherer, K., Ringeval, F., Chetouani, M., Weninger, F., Eyben, F., Marchi, E., Mortillaro, M., Salamin, H., Polychroniou, A., Valente, F. and Kim S., (2013). The Interspeech 2013 Computational Paralinguistics Challenge: Social Signals, Conflict, Emotion Autism. paper presented at the Interspeech Conference, Lyon, France, 25-29 August, 5 pages.
- Shokouhi, N., Sathyanarayana, A., Sadjadi, S.O. and Hansen J.H.L., (2013). Overlapped-speech detection with applications to driver assessment for in-vehicle active safety systems. paper presented at ICASSP Conference, Vancouver, Canada, 26-31 May, 2834–2838.
- Smolenski, B. and Ramachandran, R., (2011). Usable speech processing: A filterless approach in the presence of interference. *Circuits and Systems Magazine*, IEEE, 11(2), 8–22.
- Sollich, P., (2002). Bayesian methods for support vector machines: evidence and predictive class probabilities. *Machine Learning*, 46, 21–52.
- Sorenson, P.S., Hawkins, K. and Sorenson, R.L., (1995). Gender, psychological type and conflict style preference. *Management Communication Quarterly*, 9(1), 115–126.
- Svennevig, J., (2008). Exploring leadership conversations. *Management Communication Quarterly*, 21, 529–536.
- Thomas, K.W., and Kilmann, R.H. (1974). *Conflict MODE instrument* Tuxedo, New York: XICOM.
- Thomas, K.W., (1975). Conflict and conflict management. In Dunnette, M. (Ed.), *Handbook of industrial and organizational psychology*. Chicago: Rand McNally
- Thomas, K.W., Thomas, G.F. and Schaubhut, N., (2008). Conflict styles of men and women at six organization levels. *International Journal of Conflict Management*, 19(2), 148–166.
- Valente, F. and Vinciarelli, A. (2010). Improving Speech Processing through Social Signals: Automatic Speaker Segmentation of Political Debates using Role based Turn-Taking Patterns. paper presented at the International Workshop on Social Signal Processing, Firenze, Italy, 25-29 October, 29-34.
- Vinciarelli, A., 2007. Speakers role recognition in multiparty audio recordings using social network analysis and duration distribution modeling. *Multimedia, IEEE Transactions on*, 9(6), pp. 1215–1226.
- Vinciarelli, A., (2009). Capturing order in social interactions. *Signal Processing Magazine*, IEEE, 26(5), 133–152.
- Vogt, T. and André, E., (2005). Comparing feature sets for acted and spontaneous speech in view of automatic emotion recognition. paper presented at the ICME Conference, Amsterdam, The Netherlands, 6-8 July, 474–477.
- Yamamoto, K., Asano, F., Yamada, T. and Kitawaki, N., (2005). Detection of Overlapping Speech in Meetings Using Support Vector Regression, paper presented at the International Workshop on Acoustic Echo and Noise Control (IWAENC), Eindhoven, The Netherlands, 12-15 September, 2158-2165.

Zadrozny, B. and Elkan, C., (2002). Transforming classifier scores into accurate multiclass probability estimates. paper presented at the International Conference on Knowledge Discovery and Data Mining, Edmonton, Canada, 23-25 July, 694–699.