

Improving Speech Processing through Social Signals: Automatic Speaker Segmentation of Political Debates using Role based Turn-Taking Patterns.

Fabio Valente
Idiap Research Institute
CH-1920 Martigny
Switzerland
fabio.valente@idiap.ch

Alessandro Vinciarelli
University of Glasgow
G12 8QQ Glasgow
United Kingdom
vincia@dcs.gla.ac.uk

ABSTRACT

Several recent works on social signals have addressed the problem of statistical modeling of social interaction in multi-party discussions showing that characteristics like turn-taking patterns can be modeled and predicted according to the role that each participant has in the discussion. Reversely this work investigates the use of social signals to improve conventional speech processing methods. In details we propose the use of turn-taking patterns induced by roles for improving speaker diarization, the task of determining 'Who spoke when' in an audio file. In detail, this work studies how to include this information as statistical prior on the speaker interactions for segmenting and clustering speakers in multi-party political debates. Experiments reveal that the proposed approach reduces the speaker error over the baseline by 25% when both the number of speakers and their roles are known and by 13% relative when the pattern information is estimated from the data. Furthermore we never verify a performance degradation in any recording. Experiments are also carried out to investigate the contribution of the first-order Markov assumption i. e. that the role of the speaker n is conditionally dependent on the role of the speaker $n - 1$.

Categories and Subject Descriptors:H.3.1[Content Analysis and Indexing]. **General Terms:** Algorithms. **Keywords:** Political Debates, Social Signals, Speaker Diarization, Turn-taking patterns, Speaker Roles.

1. INTRODUCTION

A large number of recent works have focused on the statistical modeling of social interaction in small groups discussions and in between those lot of attention has been devoted to the recognition of roles (both formal or informal) in multi-party discussions. Examples include the automatic recognition of roles in meetings recordings like CMU or AMIDA recordings [1],[2], the recognition of participant seniority (professor, phd or graduate student) in the ICSI

meeting data set [3] and the recognition of functional roles in the MSC corpus [4],[5].

Typically those works are based on the use of statistical classifiers trained on a set of automatically or semi-automatically derived audio features including the speaker turn durations, the overlap between speakers and the speaker turn statistics. They assume that the participants interactions and specifically the turn-taking patterns can be statistically modeled and provide enough information for recognizing the role of each speaker in the conversation. From an engineering point of view, the turn-taking can be automatically extracted using speaker diarization systems.

Speaker Diarization aims at inferring *who spoke when* in an audio stream and involves two simultaneous unsupervised tasks: (1) the estimation of the number of speakers, and (2) the association of speech segments to each speaker. Most of the recent efforts in the domain have addressed the problem using machine learning techniques, statistical methods (for a review see [6]) or signal processing techniques for enhancing the speech signal (in case of meeting recordings [7]) ignoring the fact that the data consists of instances of human conversations.

People interact in different ways depending on the context of the environment but "Their interactions involve behaviors associated with defined statuses and particular roles. These statuses and roles help to pattern our social interactions and provide predictability" [8].

The computational linguistic literature is rich on the analysis of human conversations; seminal works of [9],[10] show that conversations obey to predictable interactions pattern between participants and a speaker turn is related in predictable ways to the previous and next turn and follows a structure similar to a grammar [9]. The manner in which orderly conversation normally takes place i. e. the way speakers take turn in a conversation is typically referred as turn-taking.

This paper investigates whether the use of the statistical information derived from roles can reversely increase the performance of conventional audio processing systems like diarization. In details, this work discusses the use of turn-taking information induced by the roles that participants play in the discussion as prior information in the speaker diarization systems. Previous works have used participant interaction patterns to improve the diarization performance, e.g. [11], however this information was considered recording dependent and not induced by, or put in relation with, any

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to publish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SSPW'10, October 29, 2010, Firenze, Italy.

Copyright 2010 ACM 978-1-4503-0174-9/10/10 ...\$10.00.

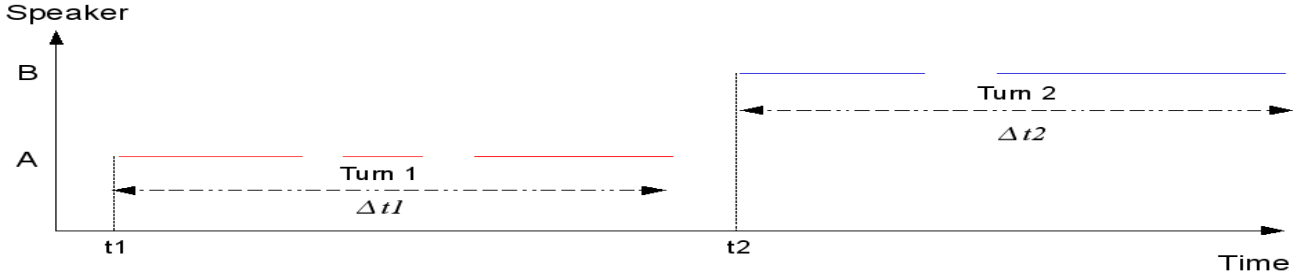


Figure 1: Speaker turn annotation: each turn is composed by a starting time t_i when a given speaker s_i grabs the floor of the discussion and a duration Δt_i which correspond at the time during which s_i holds the discussion floor. For instance the first turn is composed of three sentences from the same speaker separated by silence.

social phenomena. In this work, we make the following hypothesis: 1) the patterns are conditioned on the role that each speaker has in the conversation, 2) it can be estimated on an independent development data set.

Experiments are run on a database of political debates, described in section 2, annotated in terms of functional roles (moderator and guest) and in terms of agreeing/disagreeing groups. The choice of this dataset is motivated by previous studies on the relations between turn-taking patterns and roles (described in Section 3). Section 4 describes the baseline diarization system used for the experiments and presents its performance. Section 5 extends the diarization to include the prior information determined by the role that each speaker has in the conversation and presents results in three different case scenarios of progressively increasing difficulty i. e. considering the number of speakers and their roles known or unknown. Section 6 investigates the effectiveness of the first-order Markov assumption for modeling the turn-taking. Finally the paper is concluded in section 7 where future works and directions are also discussed.

2. DATA DESCRIPTION

The dataset used for this study consists of political debates [12] allowing the analysis of social phenomena like roles (functional and social), conflicts, agreement and disagreement. From a social interaction analysis point of view, political debates represent an excellent resource for their realism. In contrast with other benchmarks, political debates are real-world data. Debate participants do not act in a simulated social context, but participate in an event that has a major impact on their real life (for example, in terms of results at the elections). Thus, even if the debate format imposes some constraints, the participants are moved by real motivations leading to highly spontaneous social behavior.

Each debate revolves around a yes/no question like “Are you favorable to new laws on education?”. The participants state their answer (yes or no) at the beginning of the debate and do not change it during the discussion. Each debate involves a moderator and a variable number of guests (four or more). The dataset is annotated in terms of speaker times, i. e. who speaks when, and in terms of the role that each participant has in the discussion, i. e. moderator or guests. All debates include one moderator expected to ensure that all participants have at disposition the same amount of time for expressing their opinion. Furthermore, the moderator inter-

venes whenever the debate becomes too heated and people tend to interrupt one another or to talk together. The guests are labeled in terms of groups according to how they answer to the central question of the debate. Participants belonging to the same group agree with one another, while participants belonging to different groups disagree with one another.

More formally, for each debate the following triplets are available:

$$T = \{(t_1, \Delta t_1, s_1), \dots, (t_N, \Delta t_N, s_N)\} \quad (1)$$

where t_n is the beginning time of the n-th turn, Δt_n is its duration, s_n is the speaker associated with the turn and N is the total number of turns in the recording. The begin of the turn corresponds to the time at which the speaker s_n grabs the floor of the discussion and the length Δt_N corresponds to the time during which s_n holds the floor (see figure 1).

Furthermore each participant is labeled according to the role he or she has in the recording i. e. moderator m , or guest g . Guests are furthermore labeled in two groups $g1$ and $g2$ according to their agreement/disagreement. Let us designate with $\varphi(S) \rightarrow R$ the mapping between the K participants and the three roles $R = \{m, g1, g2\}$.

For performing this study the dataset is divided in two non-overlapping parts, a development dataset (composed of 25 debates for a total of 17 hours and 2600 speaker turns) and a test dataset (composed of 25 debates for a total of 15 hours and 2500 speaker turns).

3. TURN-TAKING PATTERNS AND ROLES

Previous works [13] on this dataset have shown that the sequence of speakers $S = \{s_1, \dots, s_n\}$ can be statistically modeled as a first-order Markov chain in which the probability of the participant s_n speaking after the participant s_{n-1} is regulated by their respective roles $\varphi(s_n)$ and $\varphi(s_{n-1})$.

Table 1 represents the conditional probability $p(\varphi(s_n)|\varphi(s_{n-1}))$ of a speaker role conditioned to the role of the previous speaker on the development dataset.

Those statistics are obtained disregarding overlapping speech regions (including backchannels). Although overlap regions are informative for both roles and agreement/disagreement detection, this work limits the statistics to non overlapping segments.

Table 1 can be interpreted in straightforward way: the moderator aims at sharing the available time in between the two groups and this is reflected in the fact that $p(g1|m)$ is

	Moderator	Group 1	Group 2
Moderator	0	0.51	0.49
Group 1	0.68	0.07	0.26
Group 2	0.67	0.25	0.08

Table 1: Transition matrix between roles estimated on the development data set.

approximately equal to $p(g2|m)$ as well as $p(m|g1)$ is approximately equal to $p(m|g2)$. On the other hand speakers with different opinions are more likely to take turn (on average) after a speaker they disagree with and this explains why $p(g2|g1)$ and $p(g1|g2)$ are considerably higher than $p(g1|g1)$ and $p(g2|g2)$. The probability $p(m|m)$ is equal to zero as there is only one moderator in each debate. In other words, the possible speaker sequences $S = \{s_1, \dots, s_N\}$ in a debate are not all equally probable and their probability can be simply estimated as:

$$\begin{aligned} p(S) &= p(s_1, \dots, s_n) = p(\varphi(s_1), \dots, \varphi(s_n)) = \\ &= p(\varphi(s_0)) \prod_{i=1}^N p(\varphi(s_n)|\varphi(s_{n-1})) \end{aligned} \quad (2)$$

where $p(\varphi(s_n)|\varphi(s_{n-1}))$ are elements of the matrix (1) and $p(\varphi(s_0))$ is the probability of the role associated with the speaker that opens the discussion.

Notably, the role that each participant has in a debate can be automatically estimated, finding the mapping φ^* between speaker and roles, that maximizes the probability of a given turn sequence (see [13],[2]) i. e.:

$$\varphi^* = \arg \max_{\varphi} p(\varphi(s_0)) \prod_{n=1}^N p(\varphi(s_n)|\varphi(s_{n-1})) \quad (3)$$

Results in [13] show that roles can be recognized from Eq. (3) when the speaker turns $S = \{s_1, \dots, s_n\}$ are the actual one (from manual data annotation) or are obtained using a speaker diarization system. Let us now describe the diarization system used in this study.

4. SPEAKER DIARIZATION SYSTEM

Speaker Diarization is the task that aims at inferring *who spoke when* in an audio stream. The system used here is a state-of-the-art system described in [14] and briefly summarized in the following.

Acoustic features consist of 19 MFCC coefficients extracted using a 30ms window shifted by 10ms. After speech/non-speech segmentation and rejection of non-speech regions, the acoustic features $X = \{x_1, \dots, x_T\}$ are uniformly segmented into chunks of 250ms. Then hierarchical agglomerative clustering is performed grouping together speech segments according to a distance inspired from information theory and the clustering stops when a criterion based on Normalized Mutual Information (NMI) is met (see [14] for details). This produces an estimate of the number of participants in the debate and a partition of the data in clusters, i. e., it associates each acoustic vector x_t to a speaker s .

As the diarization system classifies silence regions as non-speech, the actual turn-taking can be obtained bridging together consecutive speech segments from the same speaker separated by silence regions. For instance, the turn from the

	K known	K estimated
Speaker Error	6.2%	14.6%

Table 2: Speaker Error obtained in case of known and estimated number of speakers K . Results are reported on the test data set.

first speaker in figure 1 can be obtained bridging the silence regions that separates the three utterances spoken by the first speaker.

We refer this initial segmentation into speakers as T^* :

$$T^* = \{(t_1^*, \Delta t_1^*, s_1^*), \dots, (t_N^*, \Delta t_N^*, s_N^*)\} \quad (4)$$

After clustering, the speaker sequence is re-estimated using an ergodic Hidden Markov Model/Gaussian Mixture Model where each state represents a speaker. The emission probabilities are modeled as GMMs trained using acoustic vectors x_t assigned to speaker s . Each state enforces a minimum duration constraint. This step aims at refining the data partition obtained by the agglomerative clustering and improving the speaker segment boundaries [6].

The decoding is performed using a conventional Viterbi algorithm, i. e. the optimal speaker sequence $\mathbf{S}^* = (s_1, s_2, \dots, s_N)$ is obtained maximizing the following likelihood:

$$\mathbf{S}^* = \arg \max_{\mathbf{S}} \log p(X|S) \quad (5)$$

The emission probability $p(x_t|s_t)$ of the acoustic vector x_t conditioned to speakers s_t is:

$$\log p(x_t|s_t) = \log \sum_r w_{s_t}^r \mathcal{N}(x_t, \mu_{s_t}^r, \Sigma_{s_t}^r)$$

where $\mathcal{N}(\cdot)$ is the Gaussian pdf; $w_{s_t}^r$, $\mu_{s_t}^r$, $\Sigma_{s_t}^r$ are weights, means and covariance matrix corresponding to speaker model s_t . The output of the decoding step is a sequence of speakers with their associated speaking time.

Let us report the performance of this system on the 25 debates that compose the test data set. The most common metric for assessing diarization performances is the Diarization Error Rate [15] which is composed by speech/non-speech and speaker errors. As the same speech/non-speech segmentation is used across experiments, in the following only the speaker error is reported. Table 2 reports the speaker error in case of a-priori known number of speakers K and estimated number of speakers. When the number of speakers is estimated from data, the final speaker error is more than double compared to the case in which the number of speaker is known. This is mainly due to overlapping speech regions which can produce a number of spurious extra clusters that degrade the final system performance.

5. SPEAKER-TURNS BASED DIARIZATION

The decoding step only depends on the acoustic score $p(X|S)$ (see Eq. (5)) and completely neglects the fact that not all speaker sequences S have the same probability. In section 3, we discussed that the roles regulate the way speakers take turns and the probability of a given speaker sequence can be estimated using Eq. (2). It is thus straightforward to extend the objective function (see Eq. 5) in order

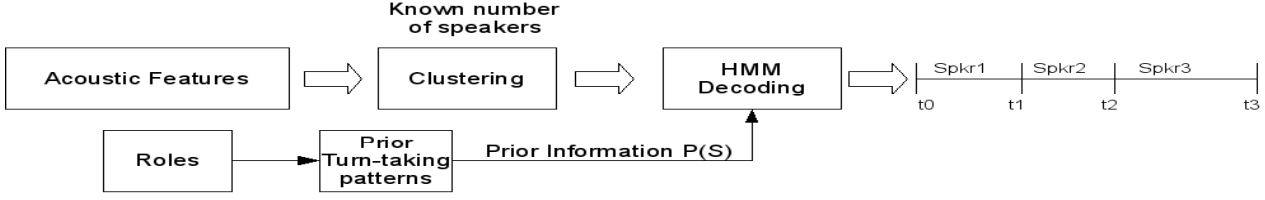


Figure 2: Schematic representation of the proposed system in case scenario 1 (known number of speakers and roles): the clustering stops when the known number of clusters is obtained; Speaker decoding is done combining the acoustic information with prior turn-taking information induced by participants role.

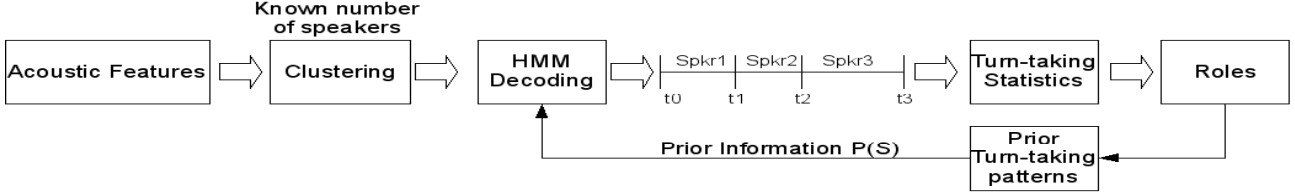


Figure 3: Schematic representation of the proposed system in case scenario 2 (known number of speakers and unknown roles): the clustering stops when the known number of clusters is obtained; turn-taking statistics obtained from the diarization output are used to recognize speaker roles. Roles are then used to compute the prior probability of a speaker sequences $P(S)$ which is used then in the diarization system.

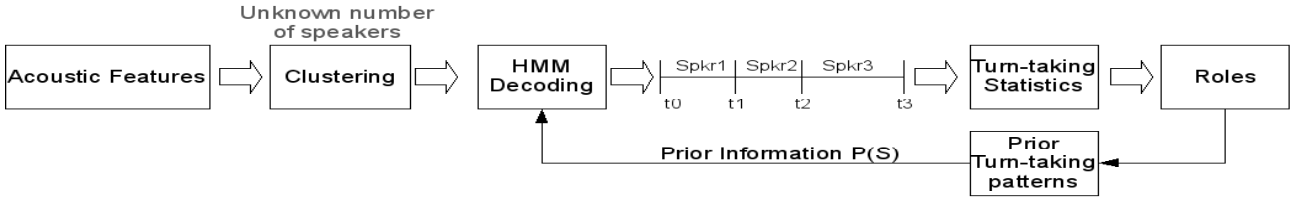


Figure 4: Schematic representation of the proposed system in case scenario 3 (unknown number of speakers and unknown roles): the clustering stops when the NMI criterion is met; turn-taking statistics obtained from the diarization output are used to recognize speaker roles. Roles are then used to compute the prior probability of a speaker sequences $P(S)$ which is used then in the diarization system.

to include this type of information i. e.:

$$\mathbf{S}^* = \arg \max_{\mathbf{S}} \log p(X|S)p(S) = \arg \max_{\mathbf{S}} \log p(X|S)p(\varphi(S)) \quad (6)$$

In other words, the optimal speaker sequence (and the associated speaker times) can be obtained combining the evidence from the acoustic score $p(X|S)$ together with the prior probability of a given sequence $p(S)$. This is somehow similar to what is done in Automatic Speech Recognition (ASR) where sentences (i. e. word sequences) are recognized combining acoustic information together with linguistic information captured in the language model. Looking at Eq. (6), it is possible to notice that while the acoustic score $p(X|S)$ is modeled using a probability density function, i. e. a GMM, $p(S)$ is a probability; as in ASR, we introduce a factor λ tuned on the development data set to scale $P(S)$ at the same order of magnitude of $p(X|S)$:

$$\mathbf{S}^* = \arg \max_{\mathbf{S}} [\log p(X|S)p(\varphi(S))^\lambda] \quad (7)$$

Eq (7) can be solved using a Viterbi decoder that includes the prior probability of different speaker sequences.

In the most general case, the number of speakers as well as

their roles are unknown. To incrementally study the integration of prior information $p(S)$, three different case scenarios are proposed. The development data set is used to estimate the probabilities $p(\varphi(s_n)|\varphi(s_{n-1}))$ and the scaling factor λ as well as the decoder insertion penalty, while performances are reported on the evaluation data set.

5.1 Case 1

The number of participants K (thus speakers) in the debate is known as well as the mapping speakers-role $\varphi(\cdot)$. The entire process is schematically depicted in Figure 2.

Those assumptions significantly simplify the problem. The clustering stops whenever the number of clusters is equal to the actual number of participants in the recording and the mapping speaker-role is obtained from the manual reference thus the prior $P(S)$ can be directly estimated from Eq. (2). Table 4 (first line) reports the speaker error obtained with conventional decoding and with role-based decoding. The inclusion of the prior information reduces the speaker error from 6.2% to 4.6% i. e. a relative improvement of 25%. The improvement is verified on all the recordings from the data set.

Table 3: Speaker Error obtained in the three case scenarios using first order Markov assumption and independence assumption. In brackets the relative improvement is reported w. r. t. no prior information.

Prior	$P(\varphi(s_n))$	$P(\varphi(s_n) \varphi(s_{n-1}))$	$P(\varphi(s_n) \varphi(s_{n-1}, \varphi(s_{n-2}))$
Case 1 - Sp. Err.	5.8 (+6%)	4.6 (+25%)	4.6 (+25%)
Case 2 - Sp. Err.	5.9 (+6%)	4.9 (+20%)	4.9 (+20%)
Case 3 - Sp. Err.	13.9 (+4%)	12.7 (+13%)	12.6 (+13%)

5.2 Case 2

The number of participants K in the debate is known but the mapping speakers-role $\varphi^*(.)$ is estimated from the segmentation T^* . The entire process is schematically depicted in Figure 3.

As before, the clustering stops whenever the number of clusters is equal to the actual number of participants in the recording producing an initial solution T^* . The mapping speakers-role $\varphi^*(.)$ is estimated from the segmentation T^* using the following maximization:

$$\varphi^* = \arg \max_{\varphi} p(\varphi(s_0^*)) \prod_{n=1}^N p(\varphi(s_n^*)|\varphi(s_{n-1}^*)). \quad (8)$$

The optimization (8) is performed exhaustively searching the space of possible mappings speakers-roles, i.e., $\varphi(\{s_k\}) \rightarrow \{m, g1, g2\}$ and selecting the one that maximize the probability of the speaker sequence s^* , i.e., Eq. (8). The size of the search space is reduced making the assumptions that 1) there is always one moderator 2) there are always two groups of opponents and a group is composed of at least one speaker.

Approximatively 70% of the speaker time is correctly labeled in terms of roles. Table 4 (second line) reports the associated speaker errors. The use of prior information reduces the error from 6.2% to 4.9% i.e. a relative improvement of 20%.

5.3 Case 3

The number of participants K and the mapping speakers-role $\varphi^*(.)$ are both unknown and estimated from data. The entire process is schematically depicted in Figure 4.

The clustering stops whenever the Normalized Mutual Information (NMI) criterion is met [14] producing an estimated number of speakers. This number is typically larger than the actual number because of overlapping speech regions which produce spurious extra clusters. The mapping speakers-role $\varphi^*(.)$ is then estimated from the segmentation T^* as before. Also spurious extra clusters are mapped into a role according to Eq. 8. Table 4 (third line) reports the speaker error with conventional decoding and with the proposed decoding. The use of prior information reduces the speaker error from 14.6% to 12.7% i.e. a relative improvement of +13%.

Figure 5 plots the speaker error with and without prior information for the 25 recordings that compose the test data set in Case 3. The proposed approach reduces the speaker error on 23 out of 25 debates in Case 3. The error does not decrease in two recordings with high speaker error. In Case 1 and Case 2 (not plotted), the improvements are verified on all the 25 recordings. We do not verify a degradation in performance in any recording.

Table 4: Speaker Error obtained in the three case scenarios without and with use of prior information. In brackets the relative improvement is reported w. r. t. no prior information.

	decoding no prior	decoding with prior
Case 1 - Sp. Err.	6.2	4.6 (+25%)
Case 2 - Sp. Err.	6.2	4.9 (+20%)
Case 3 - Sp. Err.	14.6	12.7 (+13%)

6. DECODING WITH TRIGRAMS

Finally, to quantify the contribution of the first-order Markov assumption, the previous experiments are repeated replacing the conditional distributions $P(\varphi(s_n)|\varphi(s_{n-1}))$ with $P(\varphi(s_n))$ (i.e. a unigram distribution) and with $P(\varphi(s_n)|\varphi(s_{n-1}, \varphi(s_{n-2}))$ (i.e. a trigram distribution) during the decoding step. Statistics are estimated on the development data set. Scale factor and insertion penalty are estimated accordingly on the development data set; the results are reported in Table 3.

The improvements w. r. t. the baseline are significantly smaller when the unigram prior $P(\varphi(s_n))$ is used showing that most of the gain is obtained conditioning the role of each speaker to the role of the previous one. On the other hand, trigrams $P(\varphi(s_n)|\varphi(s_{n-1}, \varphi(s_{n-2}))$ does not provide any improvement w. r. t. bigrams in Case 1 and Case 2 and they perform slightly better in Case 3 suggesting that bigrams capture already most of the information induced by the speaker roles. Furthermore higher order interaction patterns like trigrams seem to be effective only at higher speaker error rates.

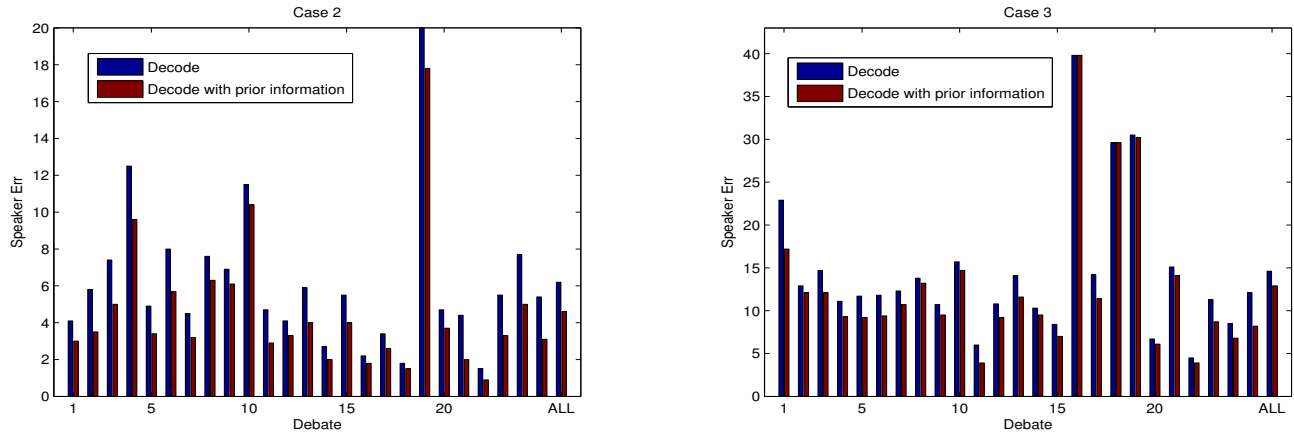
7. DISCUSSIONS

A large body of recent works has focused on the recognition of roles in multi-party discussions. Turn-taking patterns, i.e. the tendency of participants to interact or to react to certain persons rather than others, represents a powerful cue for inferring the role that each speaker has in a discussion[2],[13],[5].

Speaker diarization represents a solution for their automatic extraction. This work discusses the use of turn-taking patterns as a priori information in diarization systems. In contrary to related works [11], the patterns are explicitly put in relation with the roles that each speaker has in the discussions and they are estimated on an independent development data set. Experiments are carried out on political debates labeled in terms of speaker time, participant roles and participant agreeing/disagreeing groups.

Results show that whenever the number of participants in the debate as well as their roles are known the speaker error is reduced by 25%; whenever the second one is not available the improvement is 20%. In the most general case,

Figure 5: Speaker error obtained using realignment with and without prior information for the 25 recordings that compose the test data set for Case 2 (left picture) and for Case 3 (right picture). The speaker error is reduced on all the debates in Case 2 and on 23 debates out of 25 in Case 3. We do not verify a degradation in performance in any recording.



i.e., unknown number of speakers and unknown roles, the use of prior information reduces the error by 13% relative.

Future works will investigate a number of issues related to the turn-taking patterns estimation and their integration into the diarization system. For instance, higher order dependencies between speakers could be considered and the prior information could also be integrated into the clustering step as compared to the decoding step. This work focuses on political debates however the approach could be extended to other types of data like meetings (spontaneous or scripted) and broadcast conversations in which annotations in terms of roles, formal or informal, are available with the central question on how those statistics generalize across data.

8. ACKNOWLEDGMENTS

This work has been supported by the Swiss National Science Foundation under the NCCR IM2 grant and by the EU Network of Excellence SSPNet.

9. REFERENCES

- [1] Banerjee S. and Rudnick A., “Using simple speech-based features to detect the state of a meeting and the roles of the meeting participants,” *Proceedings of ICSLP*, 2004.
- [2] Salamin H. et al., “Automatic role recognition in multiparty recordings: Using social affiliation networks for feature extraction,” *IEEE Transactions on Multimedia*, 2010.
- [3] Laskowski K. et al., “Modeling vocal interaction for text-independent participant characterization in multi-party conversation,” *Proceedings of the 9th ISCA/ACL SIGdial Workshop on Discourse and Dialogue*, 2008.
- [4] Zancaro M. et al., “Automatic detection of group functional roles in face to face interactions,” *Proceedings of ICMI*, 2006.
- [5] Dong W. et al., “Using the influence model to recognize functional roles in meetings,” *Proceedings of ICMI*, 2007.
- [6] Tranter S.E. and Reynolds D.A., “An overview of automatic speaker diarization systems,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14(5), 2006.
- [7] Xavier Anguera Jose M. Pardo and Check Wooters, “Speaker diarization for multiple-distant-microphone meetings using several sources of information,” *IEEE Transactions on Computers*, vol. 56, no. 9, September 2007.
- [8] Tischler H., *Introduction to Sociology*, Harcourt Brace College Publisher, 1990.
- [9] Sacks H., Schegloff D., and Jefferson G., *A simple systematic for the organization of turn-taking for conversation*, Number 5. 1974.
- [10] Bettie A. and Geoffrey, *Talk: an analysis of speech and non-verbal behaviour in conversation.*, Milton Keynes: Open University Press, 1983.
- [11] Han K.J. and Narayanan S.S., “Improved speaker diarization of meeting speech with recurrent selection of representative speech segments and participant interaction pattern modeling,” in *Proceedings of ICSLP*, 2009.
- [12] Vinciarelli A. et al., “Canal9: A database of political debates for analysis of social interactions,” in *Proceedings of the International Conference on Affective Computing and Intelligent Interaction*, September 2009, pp. 1–4.
- [13] Vinciarelli A., “Capturing order in social interactions,” *IEEE Signal Processing Magazine*, September 2009.
- [14] Vijayasenan D. et al., “An information theoretic approach to speaker diarization of meeting data,” *IEEE Transactions on Audio Speech and Language Processing*, vol. 17, no. 7, pp. 1382–1393, September 2009.
- [15] “<http://www.itl.nist.gov/iad/mig/tests/rt/>,” .