

Look Who's Talking

Visual Identification of the Active Speaker in Multi-party Human-robot Interaction

Kalin Stefanov
KTH Royal Institute of
Technology
Lindstedtsvägen 24
Stockholm, Sweden
kalins@kth.se

Akihiro Sugimoto
National Institute of
Informatics
2-1-2 Hitotsubashi, Chiyoda
Tokyo, Japan
sugimoto@nii.ac.jp

Jonas Beskow
KTH Royal Institute of
Technology
Lindstedtsvägen 24
Stockholm, Sweden
beskow@kth.se

ABSTRACT

This paper presents analysis of a previously recorded multi-modal interaction dataset. The primary purpose of that dataset is to explore patterns in the focus of visual attention of humans under three different conditions - two humans involved in task-based interaction with a robot; the same two humans involved in task-based interaction where the robot is replaced by a third human, and a free three-party human interaction. The paper presents a data-driven methodology for automatic visual identification of the active speaker based on facial action units (AUs). The paper also presents an evaluation of the proposed methodology on 12 different interactions with an approximate length of 4 hours. The methodology will be implemented on a robot and used to generate natural focus of visual attention behavior during multi-party human-robot interactions.

CCS Concepts

•Human-centered computing → Human computer interaction (HCI); •Computing methodologies → Activity recognition and understanding; Machine learning;

Keywords

Multi-modal interaction, Human-robot interaction, Active speaker identification

1. INTRODUCTION

Successful multi-party human-robot interaction requires keeping track of the active speaker. In this paper we describe a methodology of detecting who is speaking based solely on visual features. For the purpose of efficient real world human-robot interaction, we have two main requirements. The first one is that we should be able to make decisions in *real-time* (possibly with a short lag), which in

practice means that the system should not require any future information beyond a limited fixed lookahead window. In a spoken human-robot interaction system, in practice, it is sufficient if the system can classify each detected utterance as coming from a particular speaker by the time it is recognized by the speech recognition module. The second requirement is that the classification should be *independent* of the speaker.

The basic problem of identifying the source speaker in a video is a recurring one in the area of multi-modal interfaces, and different applications place different requirements on the solutions. An information theoretical approach exploiting mutual correlations to associate an audio source with regions of a video stream was demonstrated by [5], while [10] showed that audiovisual correlation may be used to automatically find the correct temporal synchronization between audio and a talking face. A general pattern recognition framework was used by [3]. The above approaches are all evaluated on small amounts of data, and usability in real world scenarios has not been demonstrated. Real world applications where visual speaker identification has been mostly explored, include speaker diarization ([2],[8] and [6]) and video conference management [14]. [7] applied visual activity (the amount of movement) and focus of visual attention as features to determine who is the current speaker on real meeting room corpus data, however the results were lower than audio-only diarization. Most methods primarily use visual features from the lower half face, often directly calculated from pixel values, but body movement has also been shown to increase recognition rates, [12] and [1]. We believe that the challenge of identifying the active speaker in more dynamic and cluttered environments remains. We would like to derive an approach that can handle such type of interactions. For example, we do not want to impose limitations such as specific hardware arrangement or participants' location in the environment.

Since one main requirement of our application is that it should be speaker independent, it is convenient to use derived features in the form of speaker independent facial parameters (AUs) as input features, in contrast to the methods referenced above, which typically use pixel-based feature representations of some form. In addition, AU features are efficiently calculated which is important for the real-time applicability of the system.

The rest of the paper is structured as follows; Section 2 briefly introduces the dataset used for evaluation of the proposed methodology, and Section 3 provides a description of

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. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ASSP4MI'16, November 16 2016, Tokyo, Japan

© 2016 ACM. ISBN 978-1-4503-4557-6/16/11...\$15.00

DOI: <http://dx.doi.org/10.1145/3005467.3005470>

the methodology. We then describe two main experiments conducted and the results obtained in Section 4. Finally, we conclude and discuss directions for future work in Section 5.

2. DATASET

The analysis presented in this paper is based on a newly recorded multi-modal, multi-party dataset [11]. The main purpose of that dataset is to explore patterns in the focus of visual attention of humans under three different conditions: two humans involved in task-based interaction with a robot; the same two humans involved in task-based interaction where the robot is replaced by a third human, and a free three-party human interaction. The dataset contains two parts: 6 sessions, each of which is with duration of approximately 30 minutes, and 9 sessions, each of which is with duration of approximately 40 minutes. Both parts of the dataset are rich in modalities and recorded data streams. They include the streams of three Kinect v2 devices (color, depth, infrared, body and face data), three high quality audio streams, three high resolution GoPro video streams, touch data for the task-based interactions and the system state of the robot. In addition, the second part of the dataset introduces the data streams from three Tobii Pro Glasses 2 eye trackers. The language of all interactions is English and all data streams are spatially and temporally aligned. All interactions in the dataset occur around a round table and the participants are seated. Figure 1 illustrates the spatial configuration of the setup.

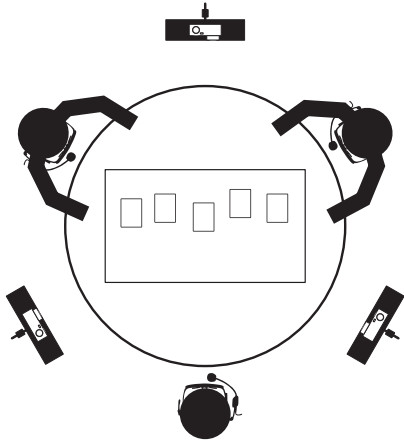


Figure 1: Spatial configuration of the setup and the location of different sensors used in the dataset.

3. METHODOLOGY

Our goal is to model two distinct behaviors during face-to-face interactions - *speaking* and *not speaking* and build an efficient and accurate classifier which can distinguish the current state of the participant given features extracted in real time. The information used for classification will vary both in space and time, therefore we use a spatio-temporal approach to modeling. One such approach is the Hidden Markov Model (HMM), [9] and [13], which is a double stochastic process governed by:

- an underlying Markov chain with a finite number of states;
- a set of random functions, each associated with one state.

In discrete time instants, the process is in one of the states and generates an observation symbol according to random function corresponding to that state. Each transition between the states has a pair of probabilities, defined as follows:

- transition probability, which provides the probability for undergoing transition;
- output probability, which defines the conditional probability of emitting an output symbol from a finite alphabet when given the state.

For the classification problem, the goal is to classify the unknown class of an observation sequence O into one of C classes. If we denote C different models by λ_c , $1 \leq c \leq C$, then the observation sequence is classified into class c^* , where $c^* = \text{argmax}_{c \in C} P(O|\lambda_c)$.

The generalized topology of an HMM is a fully connected structure, known as an *ergodic* model, where any state can be reached from any other state. We employ an 8 state fully connected HMM, as depicted in Figure 2, to model 2 different classes - speaking and not speaking.

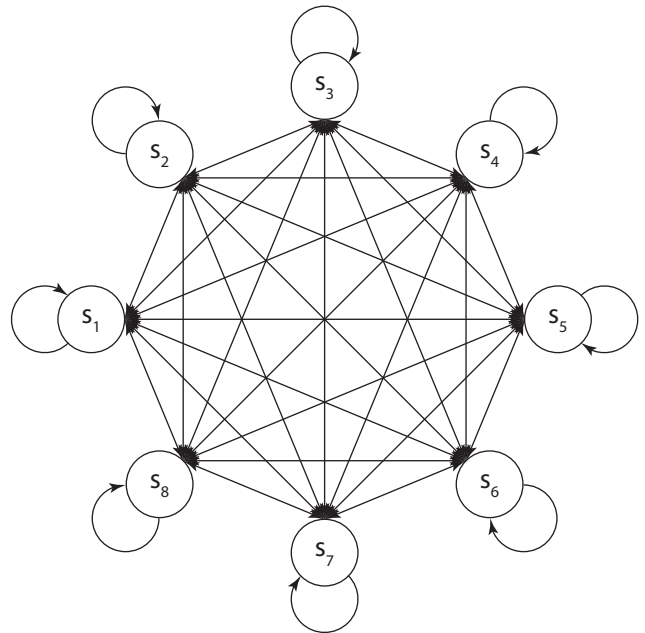


Figure 2: An 8 state fully connected HMM.

The requirement to be able to identify the active speaker in real time imposes limitation on the features that we can use to build the models. We use features that can be efficiently extracted during the interaction, specifically, the results presented in this paper are based on facial animation units (AUs), which are automatically calculated by the Kinect v2 SDK. The exact features from the Facial Action Coding System (FACS) [4] are,

- JawOpen, JawSlideRight;
- LipPucker, LipStretcherRight, LipStretcherLeft, LipCornerPullerLeft, LipCornerPullerRight, LipCornerDepressorLeft, LipCornerDepressorRight, LowerlipDepressorLeft, LowerlipDepressorRight;
- LeftcheekPuff, RightcheekPuff;
- LefteyeClosed, RighteyeClosed, RighteyebrowLowerer, LefteyebrowLowerer.

For training and testing the models we use 51 dimensional feature vector with 17 raw AU feature values, 17 first order AU differences (the difference between the feature value at time step t and at time step $t-1$) and 17 second order AU differences (the differences between the first order differences). The observations are modeled as multivariate Gaussian distribution with the full covariance matrix. We have tested the diagonal covariance matrix which results in significantly lower performance (assuming feature independence in this case is not justifiable). The number of hidden states is motivated by experiments on the data and was chosen to be 8 with accuracy and efficiency in mind. We remark that the number of AUs is 17 because this is the total number detected by Kinect. Furthermore, we remark that we do not assume only one active speaker (or speaker at all), the models are meant to be used on all visible faces when speech is detected, therefore we also allow overlaps.

4. EXPERIMENTS

We have conducted two main experiments to test our modeling approach - speaker dependent and speaker independent experiments. The acoustic signal is passed through an automatic voice activity detector (VAD) which produces intervals of speech and no speech. The silence threshold for the VAD is fixed to 200ms. We then use a window with a fixed length to extract the observation sequences from the synchronized Kinect data for both classes. The presented results are based on three different window lengths - 200ms, 500ms and 1s. The data for the AUs in the first three sessions of the dataset is missing, therefore the results presented next are based on 12 interactions (not all 15).

4.1 Speaker Dependent

The speaker dependent experiment is designed as follows. We use a 10-fold cross validation procedure to build 10x2 different subject specific models, 2 models (positive and negative) per subject with the same topology as depicted in Figure 2. We train the models using 90% of the subject's data, for speaking and not speaking classes, and use the rest of the available data for that subject to test the models. Figure 3 illustrates the results of this experiment for 12 interactions, 3 subjects per interaction.

The confusion matrices for window length of 1s, included in Table 1, show that in the speaker dependent case the negative (not speaking) class is modeled slightly better than the positive (speaking) class by the 8 state HMM. Specifically, 77.3% of the negative instances are correctly classified and 75.6% of the positive instances are correctly classified.

The class distribution is skewed, with the negative (not speaking) class having more instances. Specifically, the negative class for window length of 1s is 67.6%. The average accuracy of the proposed methodology in this case is 76.8%

which is a significant improvement compared to choosing the dominant class every time.

4.2 Speaker Independent

We designed the speaker independent experiment in the following way. First we built 36x2 subject specific models for each of the two classes. Then given all data for subject n we made classification with the rest 35x2 models (excluding the two models for the current subject). The final class assignment for the observation sequence is calculated by a majority vote between all 35x2 models (the majority voter takes the votes after $\text{argmax } P(O|positive), P(O|negative)$ for 35x2 models). Figure 4 illustrates the results of this experiment for 12 interactions, 3 subjects per interaction, for window length of 1s.

The confusion matrices for window length of 1s, included in Table 2, show that in the speaker independent case the positive (speaking) class is modeled better than the negative (not speaking) class by the 8 state HMM. Specifically, 43.6% of the negative instances are correctly classified and 74.8% of the positive instances are correctly classified.

The class distribution is skewed, with the negative (not speaking) class having more instances. Specifically, the negative class for window length of 1s is 67.6%. The average accuracy of the proposed methodology in this case is 53.7% which is significantly lower compared to choosing the dominant class every time.

5. CONCLUSIONS

We have presented a methodology that attempts to identify the active speaker during different types of multi-party interactions: collaborative task-based interactions with and without a robot and free three-party human interactions. This methodology involves efficient and automatic extraction of visual features, facial action units, and Hidden Markov Models spatio-temporal modeling. The evaluation performed shows that the not speaking class is more complex than the speaking counterpart when expressed in the selected feature space.

A desired property of the proposed approach is speaker independence. The conducted experiments show that the proposed methodology can exhibit an independence property for one of the classes modeled - the performance for the speaking class is the same for both experiments. However, the performance in the speaker independent case for the not speaking class considerably degraded from the one in the speaker dependent case.

Future work will involve the definition of different topology for the not speaking class (currently we use the same 8 state fully connected HMM for both classes), and investigation of other approaches to spatio-temporal data modeling. Furthermore, we plan to implement an additional design of the speaker independent experiment, where only two models are built from all data for $n-1$ subjects and tested on the data for the left out subject. We will also investigate additional features from the acoustic channel, which can be efficiently extracted during the interactions in order to increase the overall system accuracy. The introduction of new features will be accompanied with analysis of their individual contribution to the task of successful active speaker identification.

6. REFERENCES

- [1] S. Alexanderson, J. Beskow, and D. House. Automatic Speech/Non-Speech Classification Using Gestures in Dialogue. In *Swedish Language Technology Conference*, 2014.
- [2] X. Anguera, S. Bozonnet, N. Evans, C. Fredouille, G. Friedland, and O. Vinyals. Speaker Diarization: A Review of Recent Research. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2012.
- [3] P. Besson and M. Kunt. Hypothesis Testing for Evaluating a Multimodal Pattern Recognition Framework Applied to Speaker Detection. *Journal of NeuroEngineering and Rehabilitation*, 2008.
- [4] P. Ekman and W. Friesen. *Facial Action Coding System: A Technique for the Measurement of Facial Movement*. 1978.
- [5] J. W. Fisher, T. Darrell, W. T. Freeman, and P. Viola. Learning Joint Statistical Models for Audio-Visual Fusion and Segregation. In *Advances in Neural Information Processing Systems*, 2000.
- [6] G. Friedland, C. Yeo, and H. Hung. Visual Speaker Localization Aided by Acoustic Models. In *ACM International Conference on Multimedia*, 2009.
- [7] H. Hung and S. O. Ba. Speech/Non-Speech Detection in Meetings from Automatically Extracted Low Resolution Visual Features. Technical report, 2009.
- [8] H. J. Nock, G. Iyengar, and C. Neti. Speaker Localisation Using Audio-Visual Synchrony: An Empirical Study. In *Image and Video Retrieval*, 2003.
- [9] L. R. Rabiner. A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition. In *Proceedings of IEEE*, 1988.
- [10] M. Slaney and M. Covell. Facesync: A Linear Operator for Measuring Synchronization of Video Facial Images and Audio Tracks. In *Advances in Neural Information Processing Systems*, 2000.
- [11] K. Stefanov and J. Beskow. A Multi-party Multi-modal Dataset for Focus of Visual Attention in Human-human and Human-robot Interaction. In *Language Resources and Evaluation Conference*, 2016.
- [12] H. Vajaria, S. Sarkar, and R. Kasturi. Exploring Co-Occurrence Between Speech and Body Movement for Audio-Guided Video Localization. *IEEE Transactions on Circuits and Systems for Video Technology*, 2008.
- [13] J. Yamato, J. Ohya, and K. Ishii. Recognizing Human Action in Time-Sequential Images Using Hidden Markov Model. In *Computer Vision and Pattern Recognition*, 1992.
- [14] C. Zhang, P. Yin, Y. Rui, R. Cutler, P. Viola, X. Sun, N. Pinto, and Z. Zhang. Boosting-Based Multimodal Speaker Detection for Distributed Meetings. *IEEE Transactions on Multimedia*, 2008.

	neg	pos		neg	pos		neg	pos		neg	pos		neg	pos		neg	pos	
neg	81.7	18.3	(464)	75.4	24.6	(371)	81.3	18.7	(650)	76.4	23.6	(480)	85.1	14.9	(536)	82.5	17.5	(608)
pos	17.3	82.7	(51)	6.4	93.6	(21)	39.7	60.3	(33)	44.6	55.4	(82)	61.7	38.3	(92)	35.6	64.4	(37)
	Subject 1			Subject 2			Subject 3			Subject 4			Subject 5			Subject 6		
neg	83.1	16.9	(446)	81.9	18.1	(484)	75.6	24.4	(403)	72.4	27.6	(497)	85.4	14.6	(962)	74.5	25.5	(762)
pos	14.5	85.5	(35)	52.2	47.8	(96)	35.0	65.0	(68)	18.6	81.4	(114)	31.1	68.9	(74)	10.4	89.6	(34)
	Subject 7			Subject 8			Subject 9			Subject 10			Subject 11			Subject 12		
neg	61.2	38.8	(472)	76.5	23.5	(675)	73.5	26.5	(767)	88.3	11.7	(416)	76.7	23.3	(447)	75.0	25.0	(411)
pos	23.3	76.7	(106)	31.3	68.7	(110)	36.0	64.0	(90)	17.9	82.1	(77)	20.1	79.9	(61)	24.2	75.8	(70)
	Subject 13			Subject 14			Subject 15			Subject 16			Subject 17			Subject 18		
neg	75.5	24.5	(497)	79.4	20.6	(557)	70.3	29.7	(479)	86.4	13.6	(317)	82.8	17.2	(288)	63.4	36.6	(178)
pos	22.4	77.6	(95)	17.7	82.3	(59)	25.0	75.0	(79)	19.1	80.9	(30)	5.2	94.8	(9)	12.3	87.7	(27)
	Subject 19			Subject 20			Subject 21			Subject 22			Subject 23			Subject 24		
neg	78.9	21.1	(557)	80.8	19.2	(709)	73.1	26.9	(529)	88.6	11.4	(340)	81.2	18.8	(341)	78.6	21.4	(487)
pos	18.3	81.7	(98)	19.8	80.2	(59)	34.8	65.2	(136)	24.0	76.0	(69)	11.2	88.8	(39)	31.8	68.2	(41)
	Subject 25			Subject 26			Subject 27			Subject 28			Subject 29			Subject 30		
neg	79.5	20.5	(601)	76.6	23.4	(672)	67.0	33.0	(300)	87.2	12.8	(1267)	77.0	23.0	(483)	59.4	40.6	(729)
pos	28.9	71.1	(125)	28.5	71.5	(82)	25.1	74.9	(177)	21.9	78.1	(70)	34.1	65.9	(349)	22.4	77.6	(113)
	Subject 31			Subject 32			Subject 33			Subject 34			Subject 35			Subject 36		

Table 1: Confusion matrices for the speaker dependent experiment. The observation sequence window length is 1s, *neg* refers to the *not speaking* class and *pos* refers to the *speaking* class. The number in brackets represents instance count.

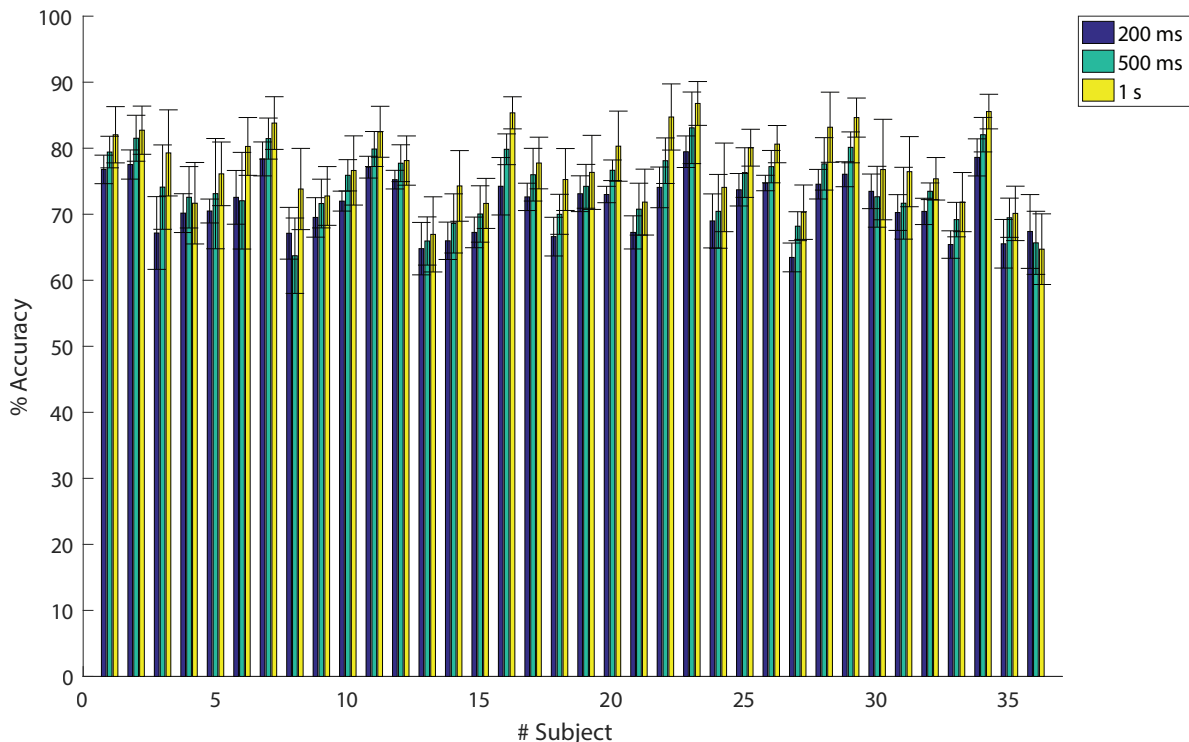


Figure 3: 10-fold cross validated results for the speaker dependent experiment. The accuracy is calculated according to $accuracy = \frac{tp+tn}{tp+tn+fp+fn}$, for 10 different train/test iterations and averaged to produce the bar graph for each subject for each fixed window length. The error bar illustrates the standard deviation.

