

A STUDY ON THE AMBIGUITY IN HUMAN ANNOTATION OF GERMAN ORAL HISTORY INTERVIEWS FOR PERCEIVED EMOTION RECOGNITION AND SENTIMENT ANALYSIS

Michael Gref¹, Nike Matthiesen², Sreenivasa Hikkal Venugopala^{1,3}, Shalaka Satheesh^{1,3},
Aswinkumar Vijayananth^{1,3}, Duc Bach Ha¹, Sven Behnke^{1,4}, Joachim Köhler¹



¹Fraunhofer Institute for Intelligent Analysis and Information Systems (IAIS)

²Haus der Geschichte der Bundesrepublik Deutschland Foundation (HdG)

³University of Applied Sciences Bonn-Rhein-Sieg

⁴Autonomous Intelligent Systems (AIS), Computer Science Institute VI, University of Bonn

Contact: michael.gref@iais.fraunhofer.de

Introduction

- Deep learning can help to make large audiovisual archives more accessible
- Prominent example in recent years: ASR for transcription of oral history interviews
- Beyond ASR: Sentiment analysis and emotion recognition can help capture, categorize and make different facets of interviews searchable
- However, humans perceive sentiments and emotions ambiguously and subjectively. Moreover, oral history interviews have multi-layered levels of complex facets of emotions
- To what degree humans and machines capture and assign these facets into commonly predefined classes? What are reasons for the ambiguity in human perception?

The HdG *Zeitzeugenportal* Corpus

Zeitzeugenportal (Portal of Oral History)

- <https://www.zeitzeugen-portal.de>
- German online service by *Haus der Geschichte* (House of the History) Foundation (HdG)
- Offers a central collection of contemporary German oral history interviews
- Currently hosts more than 8,000 clips from around 1,000 interviews

Project: Multi-Modal Mining for Oral History

- Goal: Investigation of complex search modalities for indexing oral history interviews
- Focus not only on *what* is being told but also on *how* it is being told
- Automated analysis can help better understand the role emotions play in historical remembering
- Recognition of perceived emotion

Table 1: Overview of HdG oral history data sets after annotation and split into speaker-independent subsets.

HdG Set	Videos	Segments	Hours
Training	104	1,863	6.35
Development	27	430	1.44
Test	33	471	1.74

Qualitative Survey of Annotators

- Narrative structure of oral history interviews has different levels. Accordingly, emotions become visible in different ways, such as those that arise during remembering or reported emotional situations.
- Annotators agreed: Ekman classes are insufficient to reflect the complexity of emotions in oral history interviews. Nuances of emotions do not fit into the six categories.
- Annotators intuitively combined multiple emotions to represent more complex emotions, such as hate (disgust + anger), despair/helplessness (fear + sadness), scorn (happiness + disgust), and overwhelm (happiness + surprise) in the annotation.
- Examples: *Overwhelm* was identified in some interviews about the Fall of the Berlin Wall. Disgust + anger occurred more frequently in narratives reporting oppression or persecution.

Text-Based Sentiment and Emotion Recognition

- BERT embeddings + DNN classifier
- Multi-staged training with German CMU-MOSEAS and HdG data set
- Input: Raw ASR transcripts
- Class weights are estimated using the compute-class-weights (from sklearn) to handle class imbalance

		Ground Truth						Precision
		Neutral	Happy	Sad	Anger	Disgust	Fear	Surprise
Prediction	Neutral	153	21	24	77.3%	198		
	Happy	30	66	16	58.9%	112		
	Sad	33	34	94	58.4%	161		
	Anger							
	Disgust							
	Fear							
Recall		70.8%	54.5%	70.1%	ACC	66.5%	771	

		Ground Truth						Precision
		Neutral	Happy	Sad	Anger	Disgust	Fear	Surprise
Prediction	Neutral	69	14	16	9	4	1	2
	Happy	24	46	11	8	3	2	3
	Sad	17	9	6	6	5	5	
	Anger	5	7	4	5	4	4	1
	Disgust	8	12	1	4	7	1	1
	Fear	19	16	27	4	10	12	5
Recall		47.9%	42.2%	9.0%	13.9%	20.6%	48.0%	7.7%
ACC		144	109	67	36	34	25	13
ACC		34.1%	34.2%	34.3%	34.4%	34.5%	34.6%	34.7%

Figure 4: Confusion matrix of the text-based sentiment analysis (left) and emotion recognition model (right) on HdG test.

Annotation of Perceived Emotions and Sentiment in Oral History

- Annotation of perceived emotion and sentiment by three employees at the HdG, who have an academic background in history
- The annotators received pre-segmented interview videos based on ASR, i.e., video-stream including audio and a raw ASR transcript

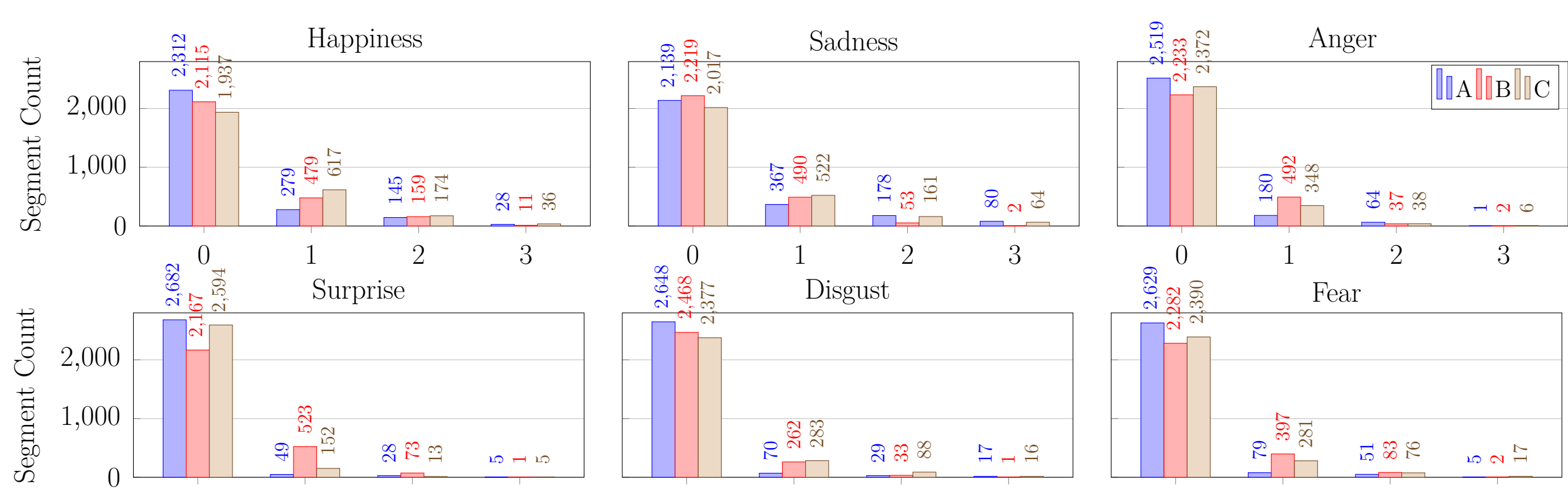


Figure 1: Histograms of the annotation scores for each emotion. Each color bar represents a different annotator.

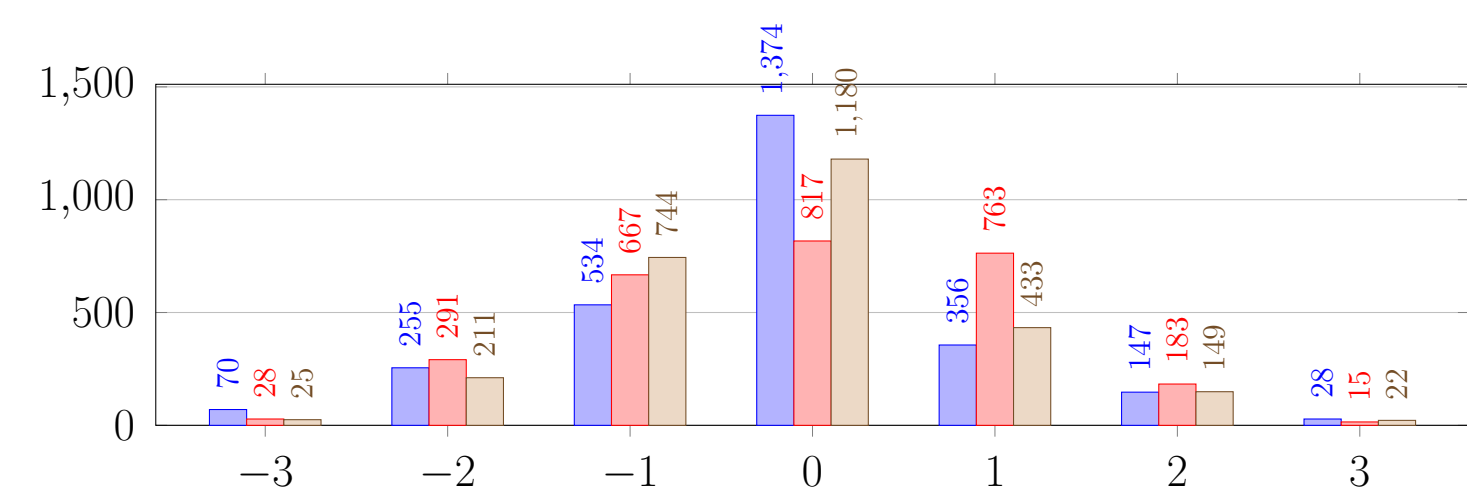


Figure 2: Sentiment annotation score histograms.

Correlation Analysis of Annotation Pairs

Table 2: Spearman rank-order correlation coefficients between the annotated labels of two transcribers.

Class	Transcriber Pairs			Avg.
	A-B	A-C	B-C	
sentiment	0.66	0.61	0.61	0.63
happy	0.52	0.52	0.60	0.55
sad	0.45	0.52	0.44	0.47
anger	0.29	0.35	0.36	0.33
surprise	0.14	0.26	0.19	0.20
disgust	0.31	0.32	0.38	0.34
fear	0.36	0.38	0.41	0.38

Inter-Class Correlation Analysis

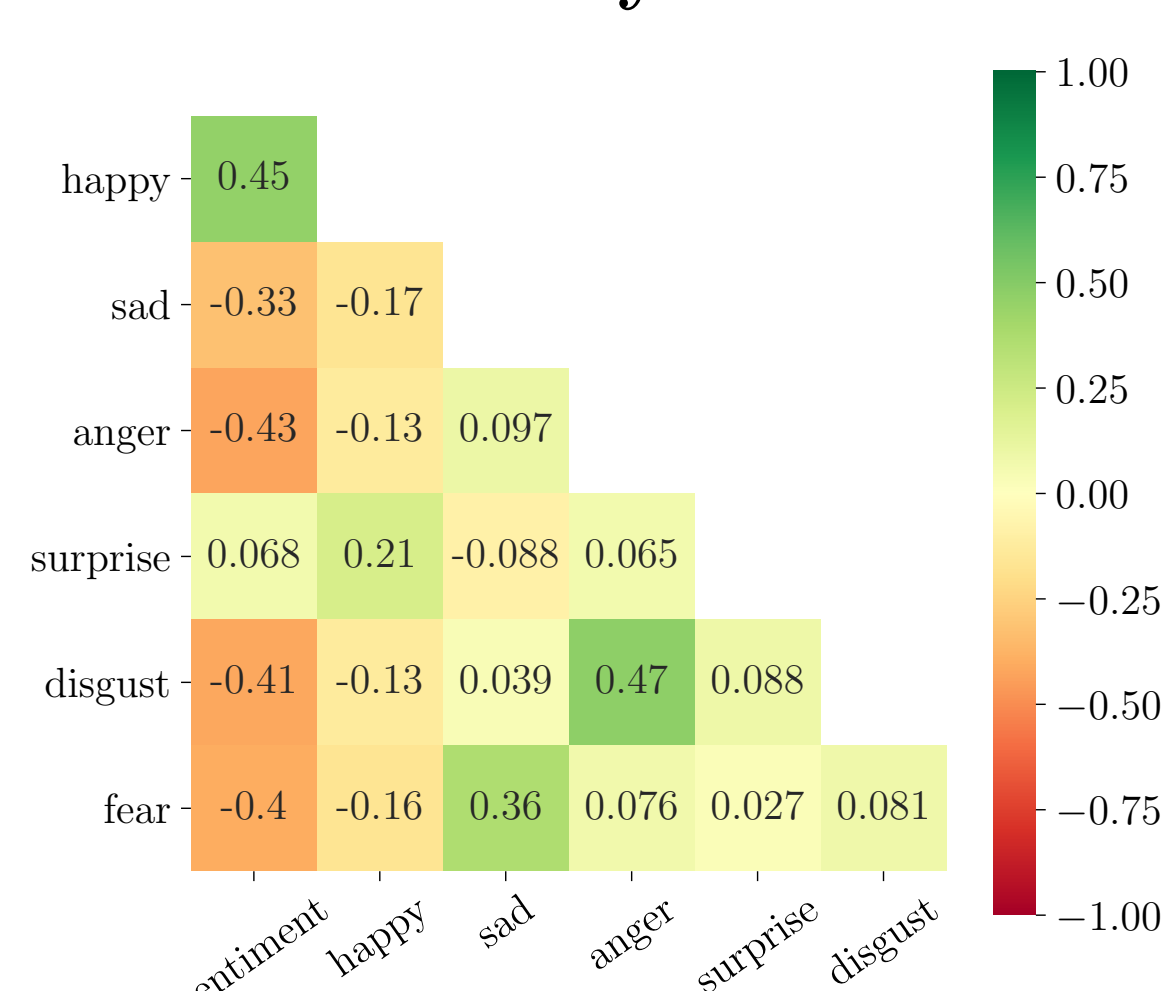


Figure 3: Spearman correlation between the annotated average scores of classes.

Speech Emotion Recognition

- Hybrid model: SVM + embeddings from a pretrained VGG-19 model (ImageNet)
- Input: Grayscale log-melspectrograms
- Combined training set: HdG, CMU-MOSEAS (German part), CMU-MOSEI, Berlin EmoDB

		Ground Truth			Precision
		Neutral	Happy	Sad	
Prediction	Neutral	358	277	134	46.6%
	Happy	1025	1261	678	42.5%
	Sad	1655	1501	1987	38.6%
Recall		41.8%	41.5%	71.0%	ACC
ACC		3038	3039	2799	8876

		Ground Truth			Precision
		Neutral	Happy	Sad	
Prediction	Neutral	17	15	4	47.2%
	Happy	67	73	48	38.8%
	Sad	84	24	24	18.2%
Recall		10.1%	65.2%	31.6%	ACC
ACC		168	112	76	356

		Ground Truth		Precision
		Happy	Sad	
Prediction	Happy	1953	1026	65.6%
	Sad	1086	1773	62.0%
Recall		64.3%	63.3%	ACC
ACC		3039	2799	5838

		Ground Truth		Precision
		Happy	Sad	
Prediction	Happy	106	58	64.6%
	Sad	6	18	75.0%
Recall		94.6%	23.7%	ACC
ACC		112	76	188

a) 3 class classification experiment

b) 2 class classification experiment

Figure 5: Different classification experiments for SER models. For each experiment, the training results are on the left. The results on the HdG test set are on the right.

Facial Emotion Recognition

- Frame Attention Networks (FAN) model (ResNet-18 feature embedding module and a subsequent frame attention network)
- Pre-training with Microsoft FER+ and AFEW data sets; fine-tuning with HdG

		Ground Truth			Precision
		Happy	Sad	Anger	
Prediction	Happy	94	12	88.7%	
	Sad	17	63	78.8%	
	Anger	4	30	42.9%	
Recall		84.7%	84.0%	ACC	
ACC		111	75	186	

		Ground Truth			Precision
		Happy	Sad	Anger	
Prediction	Happy	84	6	46	61.8%
	Sad	4	30	36	42.9%
	Anger	23	39	84	57.5%
Recall		75.7%	40.0%	50.6%	ACC
ACC		111	75	166	352

		Ground Truth			Precision
		Happy	Sad	Anger	
Prediction	Happy	92	7	3	90.2%
	Sad	19	65	10	69.1%
	Anger	3	3	1	25.0%
Recall		82.0%	86.7%	7.1%	ACC
ACC		111	75	14	200

		Ground Truth			Precision
		Happy	Sad	Anger	
Prediction	Happy	66	2	1	65.3%
	Sad	3	25	29	43.9%
	Anger	42	48	13	50.5%
Recall		39.5%	33.3%	0.0%	ACC
ACC		111	75	14	166

Figure 6: Confusion matrices from results of the FER experiments.

Summary and Conclusion

- Comparing the annotations of three persons using Ekman classes commonly used in emotion recognition revealed substantial differences in human perception for oral history interviews.
- Annotators combine different emotion classes in the annotation to represent the complex emotions of oral history interviews.
- Perceptual ambiguity and other challenges, such as class imbalance and lack of training data, currently limit the opportunities of emotion recognition for oral history archives. The systems can only differentiate fundamental emotions in our interviews, such as happiness and sadness.
- Nonetheless, our work uncovers promising observations and possibilities for further research.