

# The MuSe 2023 Multimodal Sentiment Analysis Challenge: Mimicked Emotions, Cross-Cultural Humour, and Personalisation

Lukas Christ  
EIHW, University of Augsburg  
Augsburg, Germany

Alexander Kathan  
EIHW, University of Augsburg  
Augsburg, Germany

Chris Gagne  
Hume AI  
New York, USA

Andreas König  
University of Passau  
Passau, Germany

Shahin Amiriparian  
EIHW, University of Augsburg  
Augsburg, Germany

Niklas Müller  
University of Passau  
Passau, Germany

Panagiotis Tzirakis  
Hume AI  
New York, USA

Alan Cowen  
Hume AI  
New York, USA

Björn W. Schuller  
GLAM, Imperial College London  
London, United Kingdom

Alice Baird  
Hume AI  
New York, USA

Steffen Klug  
University of Passau  
Passau, Germany

Eva-Maria Meßner  
University of Ulm  
Ulm, Germany

Erik Cambria  
Nanyang Technological University  
Singapore, Singapore

## ABSTRACT

The **Multimodal Sentiment Analysis Challenge (MuSe) 2023** is a set of shared tasks addressing three different contemporary multimodal affect and sentiment analysis problems: In the **Mimicked Emotions Sub-Challenge (MuSE-MIMIC)**, participants predict three continuous emotion targets. This sub-challenge utilises the **HUME-VIDMIMIC** dataset comprising of user-generated videos. For the **Cross-Cultural Humour Detection Sub-Challenge (MuSE-HUMOUR)**, an extension of the **Passau Spontaneous Football Coach Humour (PASSAU-SFCH)** dataset is provided. Participants predict the presence of spontaneous humour in a cross-cultural setting. The **Personalisation Sub-Challenge (MuSE-PERSONALISATION)** is based on the **Ulm-Trier Social Stress Test (Ulm-TSST)** dataset, featuring recordings of subjects in a stressed situation. Here, arousal and valence signals are to be predicted, whereas parts of the test labels are made available in order to facilitate personalisation. **MuSe 2023** seeks to bring together a broad audience from different research communities such as audio-visual emotion recognition, natural language processing, signal processing, and health informatics. In this baseline paper, we introduce the datasets, sub-challenges, and provided

feature sets. As a competitive baseline system, a Gated Recurrent Unit (GRU)-Recurrent Neural Network (RNN) is employed. On the respective sub-challenges' test datasets, it achieves a mean (across three continuous intensity targets) Pearson's Correlation Coefficient of .4727 for **MuSE-MIMIC**, an Area Under the Curve (AUC) value of .8310 for **MuSE-HUMOUR** and Concordance Correlation Coefficient (CCC) values of .7482 for arousal and .7827 for valence in the **MuSE-PERSONALISATION** sub-challenge.

## CCS CONCEPTS

• **Computing methodologies** → **Neural networks; Artificial intelligence; Computer vision; Natural language processing.**

## KEYWORDS

Multimodal Sentiment Analysis; Affective Computing; Humour Detection; Emotion Recognition; Multimodal Fusion; Challenge; Benchmark

## ACM Reference Format:

Lukas Christ, Shahin Amiriparian, Alice Baird, Alexander Kathan, Niklas Müller, Steffen Klug, Chris Gagne, Panagiotis Tzirakis, Eva-Maria Meßner, Andreas König, Alan Cowen, Erik Cambria, and Björn W. Schuller. 2023. The MuSe 2023 Multimodal Sentiment Analysis Challenge: Mimicked Emotions, Cross-Cultural Humour, and Personalisation. In *Proceedings of MuSe 2023 (4th Multimodal Sentiment Analysis Challenge and Workshop)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

In its 4th edition, the **Multimodal Sentiment Analysis (MuSe) Challenge** proposes three different tasks, namely categorical emotion

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 the author(s) 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](mailto:permissions@acm.org).

*4th Multimodal Sentiment Analysis Challenge and Workshop, October 2023, Ottawa, Canada*

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00  
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

prediction, cross-cultural humour detection, and personalised dimensional emotion regression. For the Emotional Mimicry Sub-Challenge (**MuSE-MIMIC**), emotional mimics are explored by introducing a first-of-its-kind, large-scale (557 subjects, 30+ hours), multimodal (audio, video, and text) dataset. The data were gathered in the wild, with subjects recording their own facial and vocal mimics to a wide range of ‘seed’ videos via their webcam. Subjects selected the emotions they perceived in each video out of 63 provided categories plus “no person can be seen or heard in this video” and rated each selected emotion on a 0–100 intensity scale. In this sub-challenge, participants will apply a multi-output regression to predict the intensities of three self-reported emotions from the subjects’ multimodal recorded responses related to decision-making emotional categories: Approval, Disappointment, and Uncertainty.

For the Cross-Cultural Humour Detection Sub-Challenge (**MuSE-HUMOUR**), participants will train their models to predict humour in German football press conference recordings. Different from the training data, the unseen test set consists of videos of English football press conferences, thus providing a cross-cultural evaluation setting. The Passau Spontaneous Football Coach Humour (**PASSAU-SFCH**) dataset [18] as featured in MuSe 2022 [17, 38, 63] is provided as the training set. For the test set, we extend PASSAU-SFCH with recordings of press conferences given by 7 different coaches from the English Premier League between September 2016 and September 2020. Both the training and the test set only contain segments in which the respective coach is speaking, accounting for more than 17 hours of data in total. While all provided videos are originally labelled according to the Humour Style Questionnaire (HSQ) proposed by Martin et al. [42], the prediction target in MuSE-HUMOUR is binary, i. e., presence or absence of humour.

In the further featured Personalisation Sub-Challenge (**MuSE-PERSONALISATION**), participants predict continuous estimations of valence and arousal using personalised approaches. Different from the usual speaker-independent challenge setup employed in recent years, participants will also be provided with labelled data of each subject from the test partition. Thus, MuSE-PERSONALISATION encourages the exploration of the adaptation of multimodal emotion recognition systems to individuals, taking their specific features into account. MuSE-PERSONALISATION utilises the Ulm-Trier Social Stress Test dataset (ULM-TSST) introduced in MuSe 2021 [54, 56]. ULM-TSST consists of recordings of 69 individuals undergoing the Trier Social Stress Test (TSST), a scenario designed to induce stress. Besides audio-visual recordings and their textual transcripts, ULM-TSST includes the subjects’ EDA, Electrocardiogram (ECG), Respiration (RESP), and heart rate (BPM) signals, each of them sampled at a rate of 1 kHz.

With their variety of data and prediction targets, the three sub-challenges target a broad audience, including but not limited to researchers interested in affective computing, multimodal representation learning, natural language processing, machine learning, and signal processing. MuSe is intended to serve as a common forum to compare different approaches to the proposed tasks, thereby leading to novel insights into the aptitude of different methods, modalities, and features.

In Section 2, the mentioned sub-challenges, their corresponding datasets, and the challenge protocol are outlined in more detail.

**Table 1: Statistics on the datasets of each sub-challenge. Given are the number of unique subjects(#), and the video durations in the format hh:mm:ss. Note that for MuSE-PERSONALISATION, partial recordings of the test subjects are also included in the train and development partition, as denoted by the round brackets.**

Partition	MuSE-HUMOUR		MuSE-MIMIC		MuSE-PERSONALISATION	
	#	Duration	#	Duration	#	Duration
Train	7	7:44:49	328	19:24:23	41 (+14)	3:39:56
Development	3	3:06:48	107	6:22:22	14 (+14)	1:20:10
Test	6	6:35:16	122	6:33:10	14	0:47:21
$\Sigma$	16	17:26:53	557	32:19:56	69	5:47:27

Next, Section 3 reports our pre-processing and feature extraction pipeline, as well as the experimental setup used to compute baseline results for each sub-challenge. The results are then presented in Section 4, before Section 5 concludes the paper.

## 2 THE THREE SUB-CHALLENGES

Following, we provide more details on each sub-challenge and dataset as well as the challenge protocol.

### 2.1 The MuSe-Mimic Sub-Challenge

The acquisition of data pertaining to human expressive behaviour remains a challenging task for researchers in the field of affective computing. To address this, researchers have employed various strategies, including the technique of mimicking human expressions in real-world settings [12]. The utilisation of this approach allows for specific emotions to be targeted with greater efficiency. Emotional mimicry has been observed in childhood development and typically occurs spontaneously during social interactions [32]. Furthermore, it has been demonstrated that humans possess a highly precise ability to mimic emotions, even down to specific features [33].

To address the challenge of acquiring data related to human expressive behaviour, a multimodal dataset of mimics, called the HUME-VIDMIMIC dataset, has been made available for use in the MuSE-MIMIC sub-challenge. This dataset includes over 30 hours of video data, with a mean duration of 6.4 seconds, collected from 557 subjects located in the United States (cf. Table 1, for partitioning details). The subjects in the dataset range in age from 19 to 59 years and were recruited via Prolific, with reimbursement for their time.

In the dataset, each subject was instructed to mimic a seed video featuring someone expressing emotion, and then rate the intensity of the seed video, choosing from a selection of emotional classes. The emotional expressions of Approval, Disappointment, and Uncertainty are targeted for this sub-challenge. The labels of these three continuous intensity targets were determined through an agglomerative clustering approach applied to a filtered correlation matrix of all 63 available labels<sup>1</sup> (plus no person seen or heard in

<sup>1</sup>Admiration, Adoration, Aesthetic Appreciation, Amusement, Anger, Annoyance, Anxiety, Approval, Awe, Awkwardness, Boredom, Calmness, Concentration, Confidence, Confusion, Contemplation, Contempt, Contentment, Craving, Determination, Disappointment, Disapproval, Discomfort, Disgust, Distress, Doubt, Ecstasy, Embarrassment, Empathic Pain, Enthusiasm, Entrancement, Envy, Excitement, Fear, Frustration, Gratitude, Guilt, Hesitancy, Horror, Interest, Joy, Love, Nostalgia, Pain, Panic,

the video) with  $\rho \geq 0.3$ . For each sample, after removing from the 63 emotions the ones with a lower correlation and normalisation of the intensity per emotion, the target label was determined by taking the mean intensity rating for a seed-mimic pair, across the corresponding cluster of labels. The clustering is as follows: Approval unites Approval, Admiration, Adoration, Enthusiasm, Excitement, Joy, and Love; Uncertainty unites Uncertainty, Doubt, Romance, Sexual Desire, Shock, and Surprise (negative); and finally, Disappointment unites Disappointment, Anger, Annoyance, Contempt, Disapproval, Disgust, and Frustration.

The data preparation process involved a speaker-independent partitioning of the dataset into training, development, and test sets. An automatic transcription is provided for each sample, and the faces of individuals within the videos were detected at a frequency of 2 frames per second.

For the MuSE-MIMIC sub-challenge, the aim is to perform a multi-output regression from features extracted from the multimodal (audio, video, and text) data for the intensity of the 3 emotional intensity targets. Pearson’s correlation coefficient ( $\rho$ ) is used as the evaluation metric.

## 2.2 The MuSe-Cross-Cultural Humour Sub-Challenge

Humour – defined as an expression that establishes surprising or incongruent relationships or meaning with the intent to amuse [27]–constitutes one of the most complex phenomena in human social interaction with manifold potential positive or negative effects [13]. Thus, humour has been a focal research interest in affective computing and human-computer interaction, such as natural language interfaces [10]. As humour can be expressed both verbally and non-verbally, multimodal approaches are especially suited for detecting humour. Several datasets for multimodal humour recognition have been proposed [30, 44, 62], but typically rely on staged scenarios such as TED talks or TV series. Also, some studies equate audience laughter with (successful) humour which is a common, yet crucial flaw. In contrast, to the best of our knowledge, PASSAU-SFCH is the only database for predicting spontaneous and non-scripted displays of humour with a nuanced humour measurement.

Since humour is embedded in linguistic and contextual factors, a cross-cultural study can help to understand commonalities and differences in humour usage. Recently, studies are diving into the multimodal intricacies of humour in different countries, such as differences in displayed smiling behaviours of American and French persons [50] or amongst others, gesture and prosody for humour construction in German-Brazilian interactions [36]. However, to the best of our knowledge, automated multimodal cross-cultural humour detection, which sheds light on the transferability of humour, has not been done, yet. MuSE-HUMOUR is designed to elicit first insights for this topic.

Participants will train their models utilising PASSAU-SFCH’s press conference recordings of 10 different German Bundesliga football coaches. The English test set comprises press conference

recordings of 6 football coaches from the English Premier League. However, only one of them is a native English speaker, while the other five coaches come from 5 different countries (Argentina, France, Germany, Portugal, and Spain). Every coach in both the training/development and the test set is male. The training subjects are aged between 30 and 53 years, and the test subjects’ age span ranges from 47 to 57 years. We split the German coaches into a training and development partition for our baseline experiments, where the development partition is identical to that of 2022’s MuSE-HUMOUR challenge [2, 19]. Detailed statistics on the dataset can be found in Table 1.

We segment all videos such that the data only includes the parts in which the coach is actually speaking. Both training and test data include the segmented audio-visual recordings as well as manual transcripts with timestamps. All videos are initially annotated at a 2 Hz rate for the sentiment and direction (self-directed vs others-directed) of humorous utterances, following the HSQ [42]. We obtain binary labels denoting the presence or absence of humour as described in [19], leading to 2 s video frames, each of which is either considered humorous or not. In total, 438 % of the training data, 281% of the development data, and 617% of the test data are labelled as humorous in the gold standard.

Analogously to 2022’s MuSE-HUMOUR sub-challenge, we employ the AUC metric as the evaluation criterion.

## 2.3 The MuSe-Personalisation Sub-challenge

When working with real-world data, often significant differences can be observed between individuals (e. g., high variance in pitch within a person’s speech or variations in personality characteristics such as genders, age ranges, or cultural background) [1, 37, 61]. However, most approaches today tend to neglect these individual variations, resulting in models trained on a broad population which do not always generalise well to subjects not present in the training set [22]. To solve this problem, personalisation methods are needed that incorporate the characteristics of one individual, leading to personalised models which are capable of providing more accurate predictions [37]. MuSE-PERSONALISATION is aimed to serve as a benchmark for personalisation in multimodal affect analysis.

Participants of this sub-challenge are supplied with the multimodal ULM-TSST dataset. It consists of recordings of subjects participating in the TSST [35], which defines a stress-inducing, job interview-like scenario. Participants are asked to deliver a five-minute free speech presentation on why they are suited for a hypothetical job. ULM-TSST comprises recordings of these speeches by 69 subjects (age range 18-39), 49 of which are female, accounting for circa 6 hours of video data.

The data is partitioned into a training, development, and test set in a speaker-independent manner. The training dataset comprises 41, both development and test sets 14 videos. Note that the split is identical to the splits of Ulm-TSST employed for 2021’s and 2022’s Emotional Stress Sub-Challenge (MuSE-STRESS) challenges [19, 54]. In order to facilitate personalisation on the test subjects, we provide labelled parts of their videos as follows: we take the first 60 seconds of each test video and consider this our subject-specific training data. Moreover, the next 60 seconds are employed as subject-specific development set. The remaining parts of each test subject’s video

Pride, Realisation, Relief, Resentment, Romance, Sadness, Sarcasm, Satisfaction, Serenity, Sexual Desire, Shame, Shock, Surprise (negative), Surprise (positive), Sympathy, Tiredness, Triumph, and Uncertainty.

serve as the sub-challenge’s test data, i. e., their annotations are kept confidential. Table 1 lists key statistics of this partitioning of ULM-TSST.

ULM-TSST has been labelled with 2 Hz arousal and valence signals by three annotators. The gold standards are identical to those in 2022’s MuSE-STRESS challenge [19] which attracted considerable interest, e. g., [31, 40, 47, 64]. For the valence gold standard, we fuse these three annotations employing the Rater Aligned Annotation Weighting (RAAW) method [57]. RAAW combines temporal alignment utilising Canonical Time Warping (CTW) [67] with annotation fusion via Evaluator Weighted Estimator (EWE) [28]. Regarding arousal, for each video, the signals of the two annotators with the highest agreement are merged with the video subject’s downsampled and smoothed EDA signal, as EDA has been shown to be an objective indicator of arousal [15], in contrast to inevitably subjective annotations. Details on the gold standard creation can be found in [19], for extensive experiments on merging biosignals and annotation signals for arousal see [9].

We provide participants with the audiovisual recordings, manual transcripts and the ECG, RESP, and BPM signals.

## 2.4 Challenge Protocol

In order to enter the challenge, participants need to hold an academic affiliation and complete the EULA available on the MuSe 2023 homepage<sup>2</sup>. The organisers do not take part in any sub-challenge as competitors. During the competition, participants submit their predictions for test labels on the CodaLab platform<sup>3</sup>. Up to 5 predictions are possible in each sub-challenge. All participating teams are encouraged to submit a paper describing their experiments and results. In order to officially win a sub-challenge, a paper, which must be accepted, is mandatory. Papers will undergo a double-blind peer-reviewing process.

## 3 BASELINE APPROACHES

We supply an extensive set of pre-extracted features in order to support participants in time-efficient model development. For each sub-challenge, seven feature sets (3 audio-based, 3 video-based, 1 text-based), as outlined in the following, are provided<sup>4</sup>.

### 3.1 Pre-processing

As described in Section 2, we partition every sub-challenge’s dataset into training, development, and test sets. In doing so, we take length, speaker independence, and label distributions into account (cf. Table 1). For MuSE-MIMIC – as can be seen in Table 1 –, a 60-20-20% split strategy is applied. There is no additional segmentation applied. Each sample contains a single mimic of a seed video, and labels were normalised per sample to range from [0:1].

The PASSAU-SFCH press conference recordings are segmented into clips in which the respective coach is speaking. We manually remove a few of these clips where the coach is using any language

other than English. Moreover, we discard clips whose audio quality is considerably impaired.

The recordings in the ULM-TSST dataset are cut to only show the job interview presentation defined by the TSST protocol. To preserve subjects’ privacy, we also delete parts of videos in which they disclose their names. Besides the segmentation of the test subjects’ videos described in Section 2.3, no further segmentation is conducted for MuSE-PERSONALISATION.

### 3.2 Audio

Before extracting audio features, we normalise all audio files to -3 decibels and convert them to mono, at 16 kHz, 16 bit. We then utilise the OPENSMILE [25] toolkit to compute handcrafted features. Moreover, we compute high-dimensional audio representations using both DEEPSPECTRUM [4] and a variant of WAV2VEC2.0 [7].

Both systems have proved valuable in audio-based Speech Emotion Recognition (SER) tasks [6, 11, 26].

**3.2.1 EGEMAPS.** Using the OPENSMILE toolkit [25]<sup>5</sup>, we extract 88 extended Geneva Minimalistic Acoustic Parameter Set (EGEMAPS) features [24] which have shown their suitability for sentiment analysis and SER tasks [8, 59]. For each sub-challenge, we apply the standard configuration and extract features with a window size of two seconds, and a hop size of 500 ms.

**3.2.2 DEEPSPECTRUM.** We apply DEEPSPECTRUM [4]<sup>6</sup> to obtain deep Convolutional Neural Network (CNN)-based representations from audio data. In particular, we first create a Mel-spectrogram (128 Mels, *viridis* colour mapping) from each audio file with a window size of one second and a hop size of 500 ms. We then forward the spectrogram representation into DENSENET121 and take the output of the last pooling layer as a 1024-dimensional feature vector. The efficacy of DEEPSPECTRUM has been shown for speech and audio recognition tasks [3, 5, 46].

**3.2.3 WAV2VEC2.0.** Self-supervised pretrained Transformer models have attracted considerable interest in computer audition recently [39]. A popular example of such a foundation model is WAV2VEC2.0 [7], which has frequently been employed for SER [45, 48]. As all sub-challenges are affect-related, we utilise a large version of WAV2VEC2.0 fine-tuned on the MSP-Podcast [41] dataset for continuous emotion recognition [60]<sup>7</sup>. We extract features for an audio signal by averaging over its representations in the final layer of this model, resulting in 1024-dimensional embeddings. Regarding MuSE-HUMOUR and MuSE-PERSONALISATION, we obtain 2 Hz features by sliding a 3 s window over each audio file, with a step size of 500 ms. Due to the comparably short segments in MuSE-MIMIC, a 2 s window size is applied there.

### 3.3 Video

We compute features of the visual modality based on the subjects’ faces. Therefore, we first automatically extract faces via Multi-task Cascaded Convolutional Networks (MTCNN). Subsequently, Facial Action Units (FAUs), FACENET512 and ViT representations are obtained for them.

<sup>2</sup><https://www.muse-challenge.org>

<sup>3</sup><https://codalab.lisn.upsaclay.fr/>. The link(s) to the respective challenges will be emailed to the participants.

<sup>4</sup>Note: Participants may employ other external resources (e. g., features, datasets, pretrained foundation models). The usage of additional resources, tools etc. is expected to be clearly stated in the corresponding paper.

<sup>5</sup><https://github.com/audeer/g/openmile>

<sup>6</sup><https://github.com/DeepSpectrum/DeepSpectrum>

<sup>7</sup><https://huggingface.co/audeer/wav2vec2-large-robust-12-ft-emotion-msp-dim>

**3.3.1 MTCNN.** We employ the MTCNN [65] face detection model<sup>8</sup>, to extract pictures of the subjects' faces.

In the videos of PASSAU-SFCH, typically several persons can be seen, whereas only the coach is relevant. In the first step, we attempt to filter out the respective coach's face per video automatically via clustering of face embeddings. Afterwards, we manually correct the thus obtained face sets and only keep the one corresponding to the coach. For ULM-TSST and HUME-VIDMIMIC, no such postprocessing is necessary. In both datasets' recordings, typically only one person is displayed per video. Subsequently, the extracted face images serve as the basis for feature extraction via PY-FEAT, FACENET512, and ViT.

**3.3.2 FAU.** The concept of FAUs [23] provides an interpretable way of encoding facial expressions by reference to the activation of certain facial muscles. As facial expressions contain important cues to a person's affective state, FAUs have received considerable attention from the AC community [66]. We compute automatic estimations of the activation of 20 different FAUs via the respective SVC model included in the PY-FEAT<sup>9</sup> library.

**3.3.3 FACENET512.** In order to compute high-dimensional face representations, we make use of the FACENET512 model [52] as implemented in the deepface library [53]. FACENET512 is trained on the task of face recognition and yields a 512-dimensional embedding for every face.

**3.3.4 Vision Transformer (ViT).** As an alternative vision-based strategy, we employ the DINO-trained ViT, which has been pre-trained on the ImageNet-1K dataset in a self-supervised manner using the self-distillation with no labels (DINO) method [14]. This model has demonstrated its efficacy for various image-based tasks, including emotion recognition from facial expressions [16, 58]. The model processes the extracted facial images and outputs a 384-dimensional embedding for each image. We do not conduct any further pretraining or fine-tuning.

## 3.4 Text: Transformers

Text-based features are obtained via pre-trained Transformer models in the fashion of BERT [21]. We compute sentence representations by taking the model's final layer's encoding for the CLS token, a special token referring to the whole input sentence. For MuSE-MIMIC, the English transcripts are encoded via a small pre-trained ELECTRA [20] model, yielding 256-dimensional embeddings. As MuSE-HUMOUR features a German training and development, but an English test set, we opt for the multilingual version of BERT [21]<sup>10</sup> here. This model has been pretrained on Wikipedia entries in 104 languages including German and English. It has been shown to generalise effectively from one language to another [49]. Regarding the German-only transcripts for MuSE-PERSONALISATION, we employ a German BERT model<sup>11</sup>. Both mentioned BERT variants account for 768-dimensional representations. Note that we do not fine-tune any of the models.

<sup>8</sup><https://github.com/ipazc/mtcnn>

<sup>9</sup><https://py-feat.org>

<sup>10</sup><https://huggingface.co/bert-base-multilingual-cased>

<sup>11</sup><https://huggingface.co/dbmdz/bert-base-german-cased>

## 3.5 Alignment

Each dataset includes audio, video, and transcripts. In order to ease the development of multimodal models utilising the provided features, we align the different features along with each other and the respective task's labelling scheme.

For all tasks, both the audio and the face-based features are extracted with a rate of 2 Hz by employing sliding windows of step size 500 ms for audio and sampling faces at a 2 Hz rate for video. In order to obtain sentence-wise timestamps for MuSE-HUMOUR and MuSE-PERSONALISATION from the manual transcripts, we employ a pipeline consisting of three steps: first, we utilise the Montreal Forced Aligner (MFA) [43] toolkit to generate word-level timestamps. We then automatically add punctuation to the transcript via the deepmultilingualpunctuation tool [29]. Finally, we use PySBD [51] to segment the transcript into sentences, such that sentence-wise timestamps can be inferred from the word-level timestamps. We then compute 2 Hz textual features by averaging the embeddings of all sentences whose timestamps overlap with the respective 500 ms windows. Since the videos in HUME-VIDMIMIC are typically only a few seconds long, we do not conduct a temporal alignment of their transcripts. The biosignals in ULM-TSST are originally sampled at a 1 KHz rate. In addition, we downsample them to 2 Hz, and, subsequently, smooth them via a Savitzky-Golay filter. We make both the 1 kHz and the downsampled biosignals available. Since ULM-TSST is also labelled at a 2 Hz rate, the features provided for MuSE-PERSONALISATION are already aligned to the annotation signals. As the labels in PASSAU-SFCH refer to windows of size 2 s, the features for MuSE-HUMOUR are not directly aligned to the annotations, but can easily be matched to them. In MuSE-MIMIC, the labels refer to whole videos, such that no alignment to the labels is necessary.

## 3.6 Baseline Training

Given that all tasks are based on video recordings and thus of sequential nature, we opt for a GRU-RNN followed by two feed-forward layers to serve as a competitive, yet simple baseline system. We conduct a hyperparameter search on this model for each individual sub-challenge and feature. More specifically, we optimise the GRU's hidden representations' size, the number of stacked GRU layers and the learning rate. Furthermore, we consider both unidirectional and bidirectional GRUs. Having trained all unimodal models, we also experiment with a simple late fusion approach to obtain multimodal baselines. More specifically, we average the best unimodal models' predictions, weighted by the models' performance on the development set. We make the code, hyperparameters, and best model checkpoints publicly available on GitHub<sup>12</sup>. Following, we outline task-specific training details.

**3.6.1 MuSE-MIMIC.** In order to predict the scores for the three clip-level emotion annotations in HUME-VIDMIMIC, we encode a video by feeding the video's features into the GRU model and taking its final hidden representation. The model is trained utilising Mean Squared Error (MSE) loss. Due to the large size of HUME-VIDMIMIC, we conduct the hyperparameter search on a subset of the data comprising 5 000 data points.

<sup>12</sup><https://github.com/EIHW/MuSe-2023>

**3.6.2 MUSE-HUMOUR.** Since PASSAU-SFCH’s labels refer to 2 s frames while the features are sampled every 500 ms, one training data point corresponds to a sequence of at most 4 steps. The late fusion is based on the predictions per 2 s frame. For each hyperparameter configuration, the training routine is run with 5 different, fixed random seeds. Due to the binary prediction target, we aim for the binary cross-entropy loss function.

**3.6.3 MUSE-PERSONALISATION.** Aiming to leverage both the full ULM-TSST dataset and the test-subject specific data, we employ a two-stage approach for the MUSE-PERSONALISATION baseline. Following [34], we first train a model on the whole dataset. Specifically, we do so via a hyperparameter search using 3 fixed random seeds. We then select the single best model checkpoint based on performance on the development set. In the second step, this model is duplicated for every test subject and further trained on the subjects’ data only, employing 5 fixed random seeds. For that, another hyperparameter search optimising the number of epochs and the learning rate is conducted. The final predictions are obtained by taking for every test subject the best among the 5 subject-specific models and having it predict the respective subjects’ test labels.

If, however, the best subject-specific model does not outperform the initial model on the respective subject’s development data, we take the initial model’s predictions for this subject instead. In both stages, the training data is segmented. Window size and step length of the segmentation are optimised within the hyperparameter search. Analogously to previous MuSe challenges [19, 54, 55], we employ CCC-Loss as the loss function. We limit our experiments to audio and video-based features and leave experimentation with both physiological signals and textual features to the participants.

## 4 BASELINE RESULTS

We train GRUs as described in the previous section. In the following, we present and discuss the baseline results.

### 4.1 MUSE-MIMIC

The results for MuSe-Mimic are given in Table 3.

The results of the experiment show that WAV2VEC2.0 features perform best among all unimodal approaches, with Mean Pearson’s values of .4317 and .4296 achieved on the development and test sets, respectively. ELECTRA features exhibit comparable performance, yielding Mean Pearson’s values of .4079 and .3855 on the development and test sets, respectively. In contrast, the other five audio- and video-based features display considerably lower performance, with FACENET512’s result on the test set (0.0292) indicating only a very weak correlation. Furthermore, all face-related features produce relatively low Pearson’s values, with FAUs performing the best among them, .1337 on the test set. With respect to the audio modality, WAV2VEC2.0 features stand out clearly, surpassing the  $\rho$  values of both EGEMAPS and DEEPSPECTRUM by a substantial margin. This could be due to the fact that the WAV2VEC2.0 model incorporates linguistic knowledge as it is pre-trained on an Automatic Speech Recognition (ASR) task [60]. Moreover, the relatively high Pearson values obtained with textual-only features (ELECTRA) indicate that the textual modality encodes essential information for the task. Therefore, it is likely that the model primarily exploits linguistic information embedded in the WAV2VEC2.0 embeddings,

**Table 2: Results for MUSE-MIMIC.** We report the Pearson’s correlation coefficient ( $\rho$ ) for the mean of the 3 emotional targets, *Approval*, *Disappointment*, *Uncertainty*. The reported score is the best among 5 fixed seeds, as well as the mean  $\rho$  over these seeds and the corresponding standard deviations.

Features	[Mean $\rho$ ]	
	Development	Test
<b>Audio</b>		
EGEMAPS	.0842 (.0739 $\pm$ .0067)	.0546 (.0462 $\pm$ .0070)
DEEPSPECTRUM	.0800 (.0748 $\pm$ .0043)	.0708 (.0734 $\pm$ .0072)
WAV2VEC2.0	.4317 (.4290 $\pm$ .0020)	.4296 (.4330) $\pm$ .0029)
<b>Video</b>		
FAU	.1280 (.1241 $\pm$ .0032)	.1337 (.1319 $\pm$ .0019)
ViT	.1202 (.1151 $\pm$ .0041)	.1068 (.1046 $\pm$ .0098)
FACENET512	.0669 (.0540 $\pm$ .0072)	.0292 (.0275 $\pm$ .0160)
<b>Text</b>		
ELECTRA	.4079 (.4027 $\pm$ .0028)	.3855 (.3902 $\pm$ .0037)
<b>Late Fusion</b>		
A+T	.4718 (.4695 $\pm$ .0022)	.4679 (.4657 $\pm$ .0025)
A+V	.4234 (.4131 $\pm$ .0094)	.4281 (.4209 $\pm$ .0079)
T+V	.4027 (.3983 $\pm$ .0049)	.3965 (.3869 $\pm$ .0075)
A+T+V	.4789 (.4761 $\pm$ .0024)	<b>.4727</b> (.4711 $\pm$ .0023)

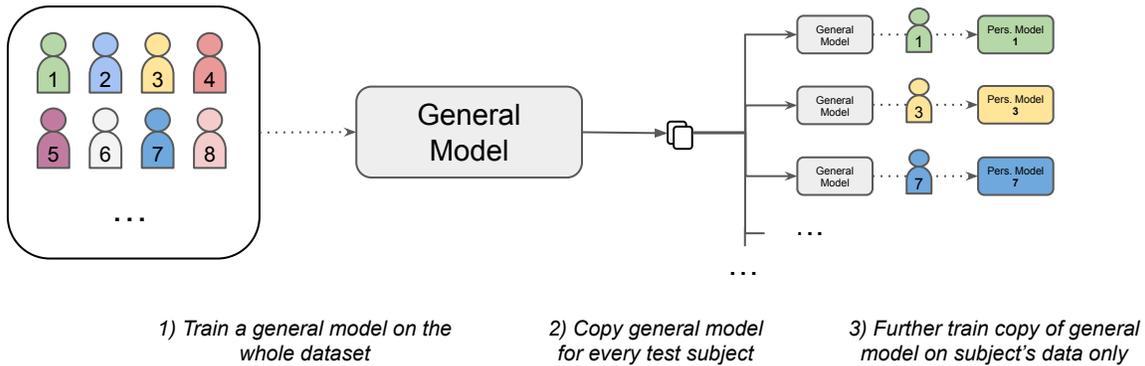
**Table 3: Class-wise Pearson’s correlation values ( $\rho$ ) for the best feature per modality and the late fusion of all of them.** Reported are the mean results across 5 fixed seeds and the respective standard deviations.

Features	[Class-wise $\rho$ ]		
	Approval	Disappointment	Uncertainty
WAV2VEC2.0 (A)	.5139 $\pm$ .0032	.4813 $\pm$ .0023	.3045 $\pm$ .0036
ELECTRA (T)	.4590 $\pm$ .0055	.4227 $\pm$ .0054	.2889 $\pm$ .0034
FAU (V)	.1786 $\pm$ .0030	.1407 $\pm$ .0041	.0763 $\pm$ .0064
A+T+V	<b>.5536</b> $\pm$ .0029	<b>.5139</b> $\pm$ .0032	<b>.3395</b> $\pm$ .0015

thus outperforming the audio-only features DEEPSPECTRUM and EGEMAPS.

Table 3 provides further insights into the performance of the best models with respect to individual emotion classes. The analysis reveals that, for all three labels, the highest mean class-wise correlation is achieved for Approval, while Uncertainty consistently proves to be the most challenging emotion target. For instance, in the case of WAV2VEC2.0, a  $\rho$  value of .5139 is observed for Approval, while the corresponding value for Uncertainty is only .3045. This trend is also reflected in the late fusion results (A+T+V), which outperform all unimodal class-wise results with mean correlations of .5536, .5139, and .3395 for Approval, Disappointment, and Uncertainty, respectively.

Moreover, as evident from Table 2 and Table 3, that performance is enhanced by the multimodal late fusion approach. Notably, the



**Figure 1: Training method employed in MuSe-Personalisation.** Note that the IDs in the figure do not correspond to actual IDs in the ULM-TSST dataset.

combination of audio and text yields promising results, improving upon the Wav2Vec2.0 outcomes by approximately .04 on both the development and test sets. Finally, the fusion of predictions obtained with the best features per modality, namely Wav2Vec2.0, FAUs, and ELECTRA, leads to a Pearson’s of .4789 and .4727 for development and test, respectively. Therefore, although the textual modality seems to play a critical role in the task at hand, multimodal approaches are expected to outperform text-only methods, indicating that all three modalities contain valuable information that can contribute to the task.

## 4.2 MuSe-HUMOUR

Table 4 reports the results for MuSe-HUMOUR.

All features lead to results above chance, i.e., 0.5 AUC. With AUC-Scores of .8435 and .7940 on the development and test data, respectively, Wav2Vec2.0 is the best-performing feature overall. Analogously to the results for MuSe-MIMIC, one reason might be that Wav2Vec2.0, pretrained on an ASR task, also encodes linguistic information already [60] and can thus not be regarded as a strictly unimodal feature extractor. It clearly outperforms eGEMAPS and DEEPSPECTRUM, which are, by construction, based on audio signals only. eGEMAPS and DEEPSPECTRUM generalise well to the test data, with the mean AUC for eGEMAPS decreasing from .6836 to .6554 and the mean AUC of DEEPSPECTRUM even improving slightly from .6936 on the development data to .7019 on the test data set. In contrast, notable generalisation issues can be observed for the video modality, in particular for the FAU and FACENET512 features. The mean AUC-Score obtained with FAU features drops from .7702 on the development data to .5948 when evaluating on the cross-cultural test data. Considering only the best seed, the textual features (.7572 AUC) outperform all other features but Wav2Vec2.0 on the test set. However, their performance is rather unstable. Their standard deviations of .0717 on the development data and .0830 on the test data are higher than those of all other features.

As for the simple late fusion approach, all possible modality combinations yield improvements over the respective unimodal results, demonstrating the multimodal nature of humour expression.

**Table 4: MuSe-HUMOUR baseline results.** Each line refers to experiments conducted with 5 fixed seeds and reports the best AUC-Score among them, together with the mean AUC-Scores and their standard deviations across the 5 seeds.

Features	[AUC]	
	Development	Test
<b>Audio</b>		
eGEMAPS	.7235 (.6836 ± .0254)	.6672 (.6554 ± .0169)
DEEPSPECTRUM	.6969 (.6936 ± .0022)	.7012 (.7019 ± .0025)
Wav2Vec2.0	.8435 (.8332 ± .0082)	.7940 (.7929 ± .0113)
<b>Video</b>		
FAU	.7879 (.7702 ± .0087)	.6398 (.5948 ± .0546)
ViT	.8277 (.7890 ± .0257)	.7457 (.7478 ± .0093)
FACENET512	.7342 (.6608 ± .0671)	.5350 (.5446 ± .0188)
<b>Text</b>		
BERT	.8105 (.7635 ± .0717)	.7572 (.7108 ± .0830)
<b>Late Fusion</b>		
A+T	.8791 (.8600 ± .0218)	.8218 (.8067 ± .0149)
A+V	.8656 (.8362 ± .0205)	.8222 (.8125 ± .0112)
T+V	.8428 (.8125 ± .0246)	.7907 (.7780 ± .0259)
A+T+V	.8759 (.8504 ± .0209)	<b>.8310</b> (.8244 ± .0168)

The combination of the best audio, text, and video models accounts for the best AUC-Score on the test set overall, namely .8310.

## 4.3 MuSe-Personalisation

Table 5 reports the results obtained for MuSe-Personalisation. In every experiment besides Valence prediction with Wav2Vec2.0, the result on the test partition is considerably lower than the CCC value on the development set. For both modalities and prediction targets, at least one feature set exists that leads to a CCC value larger than 0.5 on the test partition. Regarding the video modality, FACENET512 clearly trumps the other two feature sets for both

arousal and valence, with its arousal CCC value of .5959 on the test data exceeding the second best video feature’s (FAUs) by more than 0.200. Moreover, FACENET512 is the best-performing feature for valence, surpassing all other unimodal results on both the development and the test data. As for audio, there is no such clearly superior feature set. While EGEMAPS accounts for the best results on the development data of both arousal and valence, its result on the test sets are outperformed by DEEPSPECTRUM for arousal and WAV2VEC2.0 for valence. In general, audio features can be said to achieve better results for predicting arousal than face-based feature sets. The CCC values of both EGEMAPS and DEEPSPECTRUM on the test set for arousal, namely .7395 and .7482, outperform the best visual-based result, i. e., the CCC value of .5959 obtained using FACENET512. For valence, such a tendency does not exist. Here, only one feature set per modality yields a CCC value slightly above 0.5 on the test data, namely WAV2VEC2.0 accounting for a CCC of .5232 and FACENET512 leading to a CCC of .5654.

The late fusion results demonstrate the complementarity of the audio and video modality. In particular, the fusion of EGEMAPS and FACENET512 for valence considerably improves upon their individual performance on the test set, resulting in a CCC value of .7827. This effect is less extreme regarding arousal, where the late fusion result of .7450 on the test set only slightly exceeds the respective result achieved utilising EGEMAPS only, i. e., .7395. Overall, the best combined score of .7639 on the test set is obtained using the late fusion approach. The differences between one feature’s performance on the development and on the test data can be large. This issue might be addressed by a more sophisticated usage of the provided data. In particular, participants are not bound to use the first 60 seconds as training and the second minute as development data.

**Table 5: Results for MUSE-PERSONALISATION. Reported are the CCC values for valence, and physiological arousal after subject-specific training, i. e., based on the best personalised model among five seeds for each subject. Combined refers to the mean of the respective feature’s valence and arousal CCC values.**

Features	Arousal [CCC]		Valence [CCC]		Combined [CCC]
	Dev.	Test	Dev.	Test	Test
<b>Audio</b>					
EGEMAPS	.9073	.7395	.5892	.3944	.5670
DEEPSPECTRUM	.8064	<b>.7482</b>	.3536	.2836	.5159
WAV2VEC2.0	.7421	.5325	.5142	.5232	.5279
<b>Video</b>					
FAU	.6382	.3766	.1468	.1076	.2421
ViT	.2691	.0001	.6050	.4490	.2246
FACENET512	.8260	.5959	.6491	.5654	.5807
<b>Late Fusion</b>					
A + V	.9145	.7450	.8559	<b>.7827</b>	<b>.7639</b>

## 5 CONCLUSIONS

We introduced MuSe 2023 – the 4th Multimodal Sentiment Analysis challenge. MuSe 2023 comprises three sub-challenges:

For the MUSE-MIMIC Sub-Challenge, the novel HUMVIDMIMIC data set is made available, consisting of recordings of people mimicking emotional videos. Participants predict the degree of Approval, Disappointment, and Uncertainty for each video.

The MUSE-HUMOUR Sub-Challenge is based on an extension of the PASSAU-SFCH dataset [18], aiming at the recognition of spontaneous humour in press conferences. A cross-cultural setting is put forward for this task, where participants train their model on German recordings, but have to predict humorous utterances in English data.

In the MUSE-PERSONALISATION Sub-Challenge, participants are encouraged to develop methods for tailoring models to specific individuals. MUSE-PERSONALISATION employs the ULM-TSST dataset already featured in previous iterations of MuSe [19, 54]. While all sub-challenges feature rather simple scenarios involving one person, we believe that the variety of prediction targets and modalities will also foster progress in empathetic dialogue systems and conversational sentiment analysis.

We utilised publicly available software to pre-compute a wide selection of features for the audio, video and text modalities. Moreover, on the basis of these features, we trained a simple GRU model to obtain the following official baseline results on the respective test sets: for MUSE-MIMIC, a mean  $\rho$  value of .4727 was obtained via the late fusion of the audio, video, and text features. Similarly, for MUSE-HUMOUR a late fusion of all three modalities accounts for the baseline AUC value of .8310. Regarding MUSE-PERSONALISATION, DEEPSPECTRUM features lead to the best arousal CCC value of .7482, while a late fusion of audio and video features yielded the best results for both valence (.7827 CCC) and the mean of arousal and valence, namely .7639 CCC.

We made the code, data sets and features publicly available. The results of our fairly simple baseline systems provide first insights into the suitability of the different modalities and features for the proposed sub-challenges. We expect that these baseline results can be improved considerably via more sophisticated models and methods. MuSe 2023 is intended to be a common platform for developing and evaluating such novel multimodal approaches.

For future efforts, beyond the mere optimisation of performances targeted by this challenge, many more will need to be faced, including pressing aspects such as dependability, explainability, fairness, ‘green’ efficient processing, to name but some of the most urgent ones.

## 6 ACKNOWLEDGMENTS

This project has received funding from the Deutsche Forschungsgemeinschaft (DFG) under grant agreement No. 461420398, and the DFG’s Reinhart Koselleck project No. 442218748 (AUDI0NOMOUS).

## REFERENCES

- [1] Miray Akyunus, Tülin Gençöz, and B Türküler Aka. 2021. Age and sex differences in basic personality traits and interpersonal problems across young adulthood. *Current Psychology* 40 (2021), 2518–2527.
- [2] Shahin Amiriparian, Lukas Christ, Andreas König, Eva-Maria Meßner, Alan Cowen, Erik Cambria, and Björn W Schuller. 2022. MuSe 2022 Challenge: Multimodal Humour, Emotional Reactions, and Stress. In *Proceedings of the 30th ACM International Conference on Multimedia*. Association for Computing Machinery, New York, NY, USA, 7389–7391.
- [3] Shahin Amiriparian, Nicholas Cummins, Sandra Ottl, Maurice Gerczuk, and Björn Schuller. 2017. Sentiment Analysis Using Image-based Deep Spectrum Features. In *Proceedings 2nd International Workshop on Automatic Sentiment Analysis in the Wild (WASA 2017) held in conjunction with the 7th biannual Conference on Affective Computing and Intelligent Interaction (ACII 2017)*. AAAC, IEEE, San Antonio, TX, 26–29.
- [4] Shahin Amiriparian, Maurice Gerczuk, Sandra Ottl, Nicholas Cummins, Michael Freitag, Sergey Pugachevskiy, and Björn Schuller. 2017. Snore Sound Classification Using Image-based Deep Spectrum Features. In *Proceedings INTERSPEECH 2017, 18th Annual Conference of the International Speech Communication Association*. ISCA, ISCA, Stockholm, Sweden, 3512–3516.
- [5] Shahin Amiriparian, Maurice Gerczuk, Lukas Stappen, Alice Baird, Lukas Koebe, Sandra Ottl, and Björn Schuller. 2020. Towards Cross-Modal Pre-Training and Learning Tempo-Spatial Characteristics for Audio Recognition with Convolutional and Recurrent Neural Networks. *EURASIP Journal on Audio, Speech, and Music Processing* 2020, 19 (2020), 1–11.
- [6] Shahin Amiriparian, Tobias Hübner, Vincent Karas, Maurice Gerczuk, Sandra Ottl, and Björn W. Schuller. 2022. DeepSpectrumLite: A Power-Efficient Transfer Learning Framework for Embedded Speech and Audio Processing From Decentralized Data. *Frontiers in Artificial Intelligence* 5 (2022), 10 pages. <https://doi.org/10.3389/frai.2022.856232>
- [7] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems* 33 (2020), 12449–12460.
- [8] Alice Baird, Shahin Amiriparian, and Björn Schuller. 2019. Can deep generative audio be emotional? Towards an approach for personalised emotional audio generation. In *2019 IEEE 21st International Workshop on Multimedia Signal Processing (MMSp)*. IEEE, IEEE, Kuala Lumpur, Malaysia, 1–5.
- [9] Alice Baird, Lukas Stappen, Lukas Christ, Lea Schumann, Eva-Maria Meßner, and Björn W Schuller. 2021. A Physiologically-adapted Gold Standard for Arousal During a Stress Induced Scenario. In *Proceedings of the 2nd Multimodal Sentiment Analysis Challenge, co-located with the 29th ACM International Conference on Multimedia (ACMMM)*. ACM, Association for Computing Machinery, Changu, China, 69–73.
- [10] Kim Binsted et al. 1995. Using humour to make natural language interfaces more friendly. In *Proceedings of the ai, alife and entertainment workshop, intern. Joint conf. On artificial intelligence*.
- [11] Björn W. Schuller and Anton Batliner and Christian Bergler and Cecilia Mascolo and Jing Han and Iulia Lefter and Heysem Kaya and Shahin Amiriparian and Alice Baird and Lukas Stappen and Sandra Ottl and Maurice Gerczuk and Panagiotis Tzirakis and Chloë Brown and Jagmohan Chauhan and Andreas Grammenos and Apinan Hasthanasombat and Dimitris Spathis and Tong Xia and Pietro Cicuta and Leon J. M. Rothkrantz and Joeri Zwerts and Jelle Treep and Casper Kaandorp. 2021. The INTERSPEECH 2021 Computational Paralinguistics Challenge: COVID-19 Cough, COVID-19 Speech, Escalation & Primates. In *Proceedings INTERSPEECH 2021, 22nd Annual Conference of the International Speech Communication Association*. ISCA, ISCA, Brno, Czechia, 431–435.
- [12] Jeffrey A Brooks, Panagiotis Tzirakis, Alice Baird, Lauren Kim, Michael Opara, Xia Fang, Dacher Keltner, Maria Monroy, Rebecca Corona, Jacob Metrick, et al. 2023. Deep learning reveals what vocal bursts express in different cultures. *Nature Human Behaviour* 7, 2 (2023), 240–250.
- [13] Arnie Cann, Amanda J Watson, and Elisabeth A Bridgewater. 2014. Assessing humor at work: The humor climate questionnaire. *Humor* 27, 2 (2014), 307–323.
- [14] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. 2021. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*. 9650–9660.
- [15] Delphine Caruelle, Anders Gustafsson, Poja Shams, and Line Lervik-Olsen. 2019. The use of electrodermal activity (EDA) measurement to understand consumer emotions—a literature review and a call for action. *Journal of Business Research* 104 (2019), 146–160.
- [16] Aayushi Chaudhari, Chintan Bhatt, Achyut Krishna, and Pier Luigi Mazzeo. 2022. ViTFER: Facial Emotion Recognition with Vision Transformers. *Applied System Innovation* 5, 4 (2022), 80.
- [17] Chengxin Chen and Pengyuan Zhang. 2022. Integrating Cross-Modal Interactions via Latent Representation Shift for Multi-Modal Humor Detection. In *Proceedings of the 3rd International on Multimodal Sentiment Analysis Workshop and Challenge (Lisboa, Portugal) (MuSe' 22)*. Association for Computing Machinery, New York, NY, USA, 23–28. <https://doi.org/10.1145/3551876.3554805>
- [18] Lukas Christ, Shahin Amiriparian, Alexander Kathan, Niklas Müller, Andreas König, and Björn W Schuller. 2022. Multimodal Prediction of Spontaneous Humour: A Novel Dataset and First Results. *arXiv preprint arXiv:2209.14272* (2022).
- [19] Christ, Lukas and Amiriparian, Shahin and Baird, Alice and Tzirakis, Panagiotis and Kathan, Alexander and Müller, Niklas and Stappen, Lukas and Meßner, Eva-Maria and König, Andreas and Cowen, Alan and Cambria, Erik and Schuller, Björn W. 2022. The MuSe 2022 Multimodal Sentiment Analysis Challenge: Humor, Emotional Reactions, and Stress. In *Proceedings of the 3rd International on Multimodal Sentiment Analysis Workshop and Challenge (Lisboa, Portugal)*. Association for Computing Machinery, New York, NY, USA, 5–14.
- [20] Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators. In *ICLR*. <https://openreview.net/pdf?id=r1xMH1BtvB>
- [21] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 4171–4186.
- [22] George Doddington, Walter Liggett, Alvin Martin, Mark Przybocki, and Douglas Reynolds. 1998. *Sheep, goats, lambs and wolves: A statistical analysis of speaker performance in the NIST 1998 speaker recognition evaluation*. Technical Report. National Inst of Standards and Technology Gaithersburg Md.
- [23] Paul Ekman and Wallace V Friesen. 1978. Facial action coding system. *Environmental Psychology & Nonverbal Behavior* (1978).
- [24] Florian Eyben, Klaus R Scherer, Björn W Schuller, Johan Sundberg, Elisabeth André, Carlos Busso, Laurence Y Devillers, Julien Epps, Petri Laukka, Shrikanth S Narayanan, et al. 2015. The Geneva minimalist acoustic parameter set (GeMAPS) for voice research and affective computing. *IEEE Transactions on Affective Computing* 7, 2 (2015), 190–202.
- [25] Florian Eyben, Martin Wöllmer, and Björn Schuller. 2010. Opensmile: the munich versatile and fast open-source audio feature extractor. In *Proceedings of the 18th ACM International Conference on Multimedia*. Association for Computing Machinery, Firenze, Italy, 1459–1462.
- [26] Maurice Gerczuk, Shahin Amiriparian, Sandra Ottl, and Björn Schuller. 2022. EmoNet: A Transfer Learning Framework for Multi-Corpus Speech Emotion Recognition. *IEEE Transactions on Affective Computing* 13 (2022).
- [27] Panagiotis Gkorezis, Eugenia Petridou, and Panteleimon Xanthiakos. 2014. Leader positive humor and organizational cynicism: LMX as a mediator. *Leadership & Organization Development Journal* 35 (2014), 305–315.
- [28] Michael Grimm and Kristian Kroschel. 2005. Evaluation of natural emotions using self assessment manikins. In *IEEE Workshop on Automatic Speech Recognition and Understanding, 2005*. IEEE, IEEE, Cancun, Mexico, 381–385.
- [29] Oliver Guhr, Anne-Kathrin Schumann, Frank Bahrmann, and Hans Joachim Böhme. 2021. FullStop: Multilingual Deep Models for Punctuation Prediction. In *Proceedings of the Swiss Text Analytics Conference 2021*. CEUR Workshop Proceedings, Winterthur, Switzerland. [http://ceur-ws.org/Vol-2957/sepp\\_paper4.pdf](http://ceur-ws.org/Vol-2957/sepp_paper4.pdf)
- [30] Md Kamrul Hasan, Wasifur Rahman, AmirAli Bagher Zadeh, Jianyuan Zhong, Md Iftekhar Tanveer, Louis-Philippe Morency, and Mohammed (Ehsan) Hoque. 2019. UR-FUNNY: A Multimodal Language Dataset for Understanding Humor. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 2046–2056. <https://doi.org/10.18653/v1/D19-1211>
- [31] Yu He, Licai Sun, Zheng Lian, Bin Liu, Jianhua Tao, Meng Wang, and Yuan Cheng. 2022. Multimodal Temporal Attention in Sentiment Analysis. In *Proceedings of the 3rd International on Multimodal Sentiment Analysis Workshop and Challenge (Lisboa, Portugal) (MuSe' 22)*. Association for Computing Machinery, New York, NY, USA, 61–66. <https://doi.org/10.1145/3551876.3554811>
- [32] Ursula Hess and Agneta Fischer. 2013. Emotional mimicry as social regulation. *Personality and social psychology review* 17, 2 (2013), 142–157.
- [33] Ursula Hess and Agneta Fischer. 2014. Emotional mimicry: Why and when we mimic emotions. *Social and personality psychology compass* 8, 2 (2014), 45–57.
- [34] Alexander Kathan, Shahin Amiriparian, Lukas Christ, Andreas Triantafyllopoulos, Niklas Müller, Andreas König, and Björn W. Schuller. 2022. A Personalised Approach to Audiovisual Humour Recognition and Its Individual-Level Fairness. In *Proceedings of the 3rd International on Multimodal Sentiment Analysis Workshop and Challenge (Lisboa, Portugal) (MuSe' 22)*. Association for Computing Machinery, New York, NY, USA, 29–36. <https://doi.org/10.1145/3551876.3554800>
- [35] Clemens Kirschbaum, Karl-Martin Pirke, and Dirk H Hellhammer. 1993. The 'Trier Social Stress Test' – a tool for investigating psychobiological stress responses in a laboratory setting. *Neuropsychobiology* 28, 1-2 (1993), 76–81.
- [36] Anna Ladilova and Ulrike Schröder. 2022. Humor in intercultural interaction: A source for misunderstanding or a common ground builder? A multimodal analysis. *Intercultural Pragmatics* 19, 1 (2022), 71–101.
- [37] Jialin Li, Alia Waleed, and Hanan Salam. 2023. A Survey on Personalized Affective Computing in Human-Machine Interaction. *arXiv preprint arXiv:2304.00377* (2023).

- [38] Jia Li, Ziyang Zhang, Junjie Lang, Yueqi Jiang, Liuwei An, Peng Zou, Yangyang Xu, Sheng Gao, Jie Lin, Chunxiao Fan, Xiao Sun, and Meng Wang. 2022. Hybrid Multimodal Feature Extraction, Mining and Fusion for Sentiment Analysis. In *Proceedings of the 3rd International on Multimodal Sentiment Analysis Workshop and Challenge* (Lisboa, Portugal) (*MuSe' 22*). Association for Computing Machinery, New York, NY, USA, 81–88. <https://doi.org/10.1145/3551876.3554809>
- [39] Shuo Liu, Adria Mallol-Ragolta, Emilia Parada-Cabaleiro, Kun Qian, Xin Jing, Alexander Kathan, Bin Hu, and Björn W. Schuller. 2022. Audio self-supervised learning: A survey. *Patterns* 3, 12 (2022), 100616. <https://doi.org/10.1016/j.patter.2022.100616>
- [40] Yiping Liu, Wei Sun, Xing Zhang, and Yebao Qin. 2022. Improving Dimensional Emotion Recognition via Feature-Wise Fusion. In *Proceedings of the 3rd International on Multimodal Sentiment Analysis Workshop and Challenge* (Lisboa, Portugal) (*MuSe' 22*). Association for Computing Machinery, New York, NY, USA, 55–60. <https://doi.org/10.1145/3551876.3554804>
- [41] R. Lotfian and C. Busso. 2019. Building Naturalistic Emotionally Balanced Speech Corpus by Retrieving Emotional Speech From Existing Podcast Recordings. *IEEE Transactions on Affective Computing* 10, 4 (October-December 2019), 471–483. <https://doi.org/10.1109/TAFFC.2017.2736999>
- [42] Rod A Martin, Patricia Publik-Doris, Gwen Larsen, Jeanette Gray, and Kelly Weir. 2003. Individual differences in uses of humor and their relation to psychological well-being: Development of the Humor Styles Questionnaire. *Journal of research in personality* 37, 1 (2003), 48–75.
- [43] Michael McAuliffe, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger. 2017. Montreal Forced Aligner: Trainable Text-Speech Alignment Using Kaldi. In *Proceedings of INTERSPEECH*, Vol. 2017. International Speech Communication Association (ISCA), Stockholm, Sweden, 498–502.
- [44] Anirudh Mittal, Pranav Jeevan, Prerak Gandhi, Diptesh Kanojia, and Pushpak Bhattacharyya. 2021. "So You Think You're Funny?": Rating the Humour Quotient in Standup Comedy. arXiv:arXiv preprint arXiv:2110.12765
- [45] Edmilson Morais, Ron Hoory, Weizhong Zhu, Itai Gat, Matheus Damasceno, and Hagai Aronowitz. 2022. Speech emotion recognition using self-supervised features. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 6922–6926.
- [46] Sandra Ottl, Shahin Amiriparian, Maurice Gerczuk, Vincent Karas, and Björn Schuller. 2020. Group-level Speech Emotion Recognition Utilising Deep Spectrum Features. In *Proceedings of the 8th ICMI 2020 EmotiW – Emotion Recognition In The Wild Challenge (EmotiW 2020)*, 22nd ACM International Conference on Multimodal Interaction (ICMI 2020). ACM, ACM, Utrecht, The Netherlands, 821–826.
- [47] Ho-min Park, Ilho Yun, Ajit Kumar, Ankit Kumar Singh, Bong Jun Choi, Dhananjay Singh, and Wesley De Neve. 2022. Towards Multimodal Prediction of Time-Continuous Emotion Using Pose Feature Engineering and a Transformer Encoder. In *Proceedings of the 3rd International on Multimodal Sentiment Analysis Workshop and Challenge* (Lisboa, Portugal) (*MuSe' 22*). Association for Computing Machinery, New York, NY, USA, 47–54. <https://doi.org/10.1145/3551876.3554807>
- [48] Leonardo Pepino, Pablo Riera, and Luciana Ferrer. 2021. Emotion Recognition from Speech Using wav2vec 2.0 Embeddings. In *Proc. Interspeech 2021*. ISCA, ISCA, Brno, Czechia, 3400–3404. <https://doi.org/10.21437/Interspeech.2021-703>
- [49] Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How Multilingual is Multilingual BERT?. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 4996–5001. <https://doi.org/10.18653/v1/P19-1493>
- [50] Béatrice Priego-Valverde, Brigitte Bigi, Salvatore Attardo, Lucy Pickering, and Elisa Gironzetti. 2018. Is smiling during humor so obvious? a cross-cultural comparison of smiling behavior in humorous sequences in american english and french interactions. *Intercultural Pragmatics* 15, 4 (2018), 563–591.
- [51] Nipun Sadvilkar and Mark Neumann. 2020. PySBD: Pragmatic Sentence Boundary Disambiguation. In *Proceedings of Second Workshop for NLP Open Source Software (NLP-OSS)*. Association for Computational Linguistics, Online, 110–114. <https://www.aclweb.org/anthology/2020.nlposs-1.15>
- [52] Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. FaceNet: A unified embedding for face recognition and clustering. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE. <https://doi.org/10.1109/cvpr.2015.7298682>
- [53] Sefik Ilkin Serengil and Alper Ozpinar. 2020. LightFace: A Hybrid Deep Face Recognition Framework. In *2020 Innovations in Intelligent Systems and Applications Conference (ASYU)*. IEEE, 23–27. <https://doi.org/10.1109/ASYU50717.2020.9259802>
- [54] Lukas Stappen, Alice Baird, Lukas Christ, Lea Schumann, Benjamin Sertolli, Eva-Maria Messner, Erik Cambria, Guoying Zhao, and Björn W Schuller. 2021. The MuSe 2021 multimodal sentiment analysis challenge: sentiment, emotion, physiological-emotion, and stress. In *Proceedings of the 2nd on Multimodal Sentiment Analysis Challenge*. Association for Computing Machinery, New York, NY, USA, 5–14.
- [55] Lukas Stappen, Alice Baird, Georgios Rizos, Panagiotis Tzirakis, Xinchun Du, Felix Hafner, Lea Schumann, Adria Mallol-Ragolta, Bjoern W. Schuller, Julia Lefter, Erik Cambria, and Ioannis Kompatsiaris. 2020. MuSe 2020 Challenge and Workshop: Multimodal Sentiment Analysis, Emotion-Target Engagement and Trustworthiness Detection in Real-Life Media. In *Proceedings of the 1st International on Multimodal Sentiment Analysis in Real-Life Media Challenge and Workshop*. ACM, Association for Computing Machinery, New York, NY, USA, 35–44.
- [56] Lukas Stappen, Eva-Maria Meßner, Erik Cambria, Guoying Zhao, and Björn W. Schuller. 2021. MuSe 2021 Challenge: Multimodal Emotion, Sentiment, Physiological-Emotion, and Stress Detection. In *29th ACM International Conference on Multimedia (ACMMM)* (Virtual Event, China). ACM, Association for Computing Machinery, New York, NY, USA.
- [57] Lukas Stappen, Lea Schumann, Benjamin Sertolli, Alice Baird, Benjamin Weigel, Erik Cambria, and Björn W Schuller. 2021. MuSe-Toolbox: The Multimodal Sentiment Analysis Continuous Annotation Fusion and Discrete Class Transformation Toolbox. In *Proceedings of the 2nd Multimodal Sentiment Analysis Challenge, co-located with the 29th ACM International Conference on Multimedia (ACMMM)*. ACM, Association for Computing Machinery, Changu, China, 75–82.
- [58] Lorenzo Vaiani, Moreno La Quatra, Luca Cagliero, and Paolo Garza. 2022. ViPER: Video-Based Perceiver for Emotion Recognition. In *Proceedings of the 3rd International on Multimodal Sentiment Analysis Workshop and Challenge* (Lisboa, Portugal) (*MuSe' 22*). Association for Computing Machinery, New York, NY, USA, 67–73. <https://doi.org/10.1145/3551876.3554806>
- [59] Bogdan Vlasenko, RaviShankar Prasad, and Mathew Magimai.-Doss. 2021. Fusion of Acoustic and Linguistic Information using Supervised Autoencoder for Improved Emotion Recognition. In *Proceedings of the 2nd on Multimodal Sentiment Analysis Challenge*. Association for Computing Machinery, New York, NY, USA, 51–59.
- [60] J. Wagner, A. Triantafyllopoulos, H. Wierstorf, M. Schmitt, F. Burkhardt, F. Eyyen, and B. W. Schuller. 2023. Dawn of the Transformer Era in Speech Emotion Recognition: Closing the Valence Gap. *IEEE Transactions on Pattern Analysis & Machine Intelligence* 01 (2023), 1–13. <https://doi.org/10.1109/TPAMI.2023.3263585>
- [61] Yanna J Weisberg, Colin G DeYoung, and Jacob B Hirsh. 2011. Gender differences in personality across the ten aspects of the Big Five. *Frontiers in psychology* 2 (2011), 178.
- [62] Jiaming Wu, Hongfei Lin, Liang Yang, and Bo Xu. 2021. MUMOR: A Multimodal Dataset for Humor Detection in Conversations. In *CCF International Conference on Natural Language Processing and Chinese Computing*. Springer, Springer, Qingdao, China, 619–627.
- [63] Haojie Xu, Weifeng Liu, Jiangwei Liu, Mingzheng Li, Yu Feng, Yasi Peng, Yunwei Shi, Xiao Sun, and Meng Wang. 2022. Hybrid Multimodal Fusion for Humor Detection. In *Proceedings of the 3rd International on Multimodal Sentiment Analysis Workshop and Challenge* (Lisboa, Portugal) (*MuSe' 22*). Association for Computing Machinery, New York, NY, USA, 15–21. <https://doi.org/10.1145/3551876.3554802>
- [64] Sarthak Yadav, Tilak Purohit, Zohreh Mostaani, Bogdan Vlasenko, and Mathew Magimai.-Doss. 2022. Comparing Biosignal and Acoustic Feature Representation for Continuous Emotion Recognition. In *Proceedings of the 3rd International on Multimodal Sentiment Analysis Workshop and Challenge* (Lisboa, Portugal) (*MuSe' 22*). Association for Computing Machinery, New York, NY, USA, 37–45. <https://doi.org/10.1145/3551876.3554812>
- [65] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. 2016. Joint Face Detection and Alignment Using Multitask Cascaded Convolutional Networks. *IEEE Signal Processing Letters* 23 (04 2016).
- [66] Ruicong Zhi, Mengyi Liu, and Dezheng Zhang. 2020. A comprehensive survey on automatic facial action unit analysis. *The Visual Computer* 36 (2020), 1067–1093.
- [67] Feng Zhou and Fernando De la Torre. 2015. Generalized canonical time warping. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 38, 2 (2015), 279–294.