# Mixed Signals: Understanding Model Disagreement in Multimodal Empathy Detection

**Anonymous ACL submission** 

#### Abstract

Multimodal models play a key role in empathy detection, but their performance can suffer when modalities provide conflicting cues. To understand these failures, we examine cases where unimodal and multimodal predictions diverge. Using fine-tuned models for text, audio, and video, along with a gated fusion model, we find that such disagreements often reflect underlying ambiguity, as evidenced by annotator uncertainty. Our analysis shows that dominant signals in one modality can mislead fusion when unsupported by others. We also observe that humans, like models, do not consistently benefit from multimodal input. These insights position disagreement as a useful diagnostic signal for identifying challenging examples and improving empathy system robustness.

#### 1 Introduction

004

007

013

015

017

021

034

040

041

Empathy recognition in human communication is a nuanced and multifaceted task, and a core component of socially intelligent communication systems (Fung et al., 2016). Empathy, commonly defined as the capacity to understand and share the emotional experiences of others, encompasses both cognitive perspective-taking and affective resonance (Baumeister and Vohs, 2007). In human interactions, language, speech, and visual cues jointly convey emotional intent (Holler and Levinson, 2019). For example, a seemingly neutral utterance might be perceived as warm or concerned when accompanied by a sympathetic tone or facial expression. For AI systems, effectively interpreting these multimodal signals requires not only accurate unimodal representations but also robust integration of potentially conflicting information across modalities. Despite recent advances in multimodal emotion recognition (Jabeen et al., 2021), empathy recognition remains particularly complex, as empathy often arises from subtle contextual cues that may not align across modalities (Hasan et al., 2023).



Figure 1: Given classifications provided by a single modality, we identify cases where integrating additional modalities leads to a different prediction. We analyze these flips to understand when and why they occur.

042

043

044

046

047

049

050

051

055

060

062

063

064

065

Our work investigates such complexity by examining when and why multimodal models misclassify empathy compared to their unimodal counterparts. We extend methods from prior work on dataset difficulty and human-model agreement (Swayamdipta et al., 2020; Saha et al., 2022), to the underexplored domain of empathy modeling. Through a combination of model-driven analysis and humansubject experiments, we pinpoint instances where unimodal and multimodal model predictions diverge, indicating conflicting cross-modal cues. By linking modality disagreement to human disagreement, we offer new insight into the limitations of current empathy modeling and highlight the value of disagreement-based analysis in socially grounded language tasks.

### 2 Related Work

**Empathy Modeling.** Early computational work on empathy focused on generating emotionally relevant textual responses (Rashkin et al., 2019; Li et al., 2019), but these approaches are inherently limited by the absence of non-verbal cues critical to empathic understanding. Recent datasets such as EMPATHICSTORIES++ (Shen et al., 2024), MEDIC (Zhu et al., 2023), EMMI (Galland et al.,

2024) and Chen et al. (2024) address this limitation 067 by incorporating speech, facial expressions, and 068 interaction context, enabling more comprehensive 069 modeling of empathy. These resources have motivated frameworks like PEGS (Zhang et al., 2024), which integrate text and visual stickers for affec-072 tive generation. Despite these advances, empathy 073 remains difficult to model due to its reliance on subtle, often conflicting signals across modalities. Prior work has largely focused on improving fusion strategies under the assumption that modalities 077 are complementary (Zadeh et al., 2017; Tsai et al., 078 2019), but has paid less attention to when fusion may fail or introduce noise.

> **Dataset Difficulty.** Complementary lines of work have investigated data difficulty and model disagreement as tools for understanding model behavior. Swayamdipta et al. (2020) propose the dataset cartography method to identify hard or ambiguous training samples; Saha et al. (2022) demonstrate that difficult instances are also harder for both humans and models to explain; Wang et al. (2023)'s Learning-From-Disagreement (LFD) framework underscores the importance of examining disagreements between models to gain deeper, actionable insights into their behaviors.

090

094

096

100

101

115

Yet, these methods remain underexplored in empathy modeling, where ambiguity is often intrinsic. Our work bridges this gap by using model-modality disagreement to identify inherently ambiguous empathy examples—cases where fusion misleads models, and where annotators also exhibit uncertainty.

# 3 Experiment 1: Identifying Complex Examples from Modality Disagreement

Disagreement between models trained on different 102 modalities can reveal challenging, nuanced, or am-103 biguous examples. Here, we identify and analyze 104 such cases of disagreement in empathy classifica-105 tion using a multimodal English empathy speech 106 dataset collected from Youtube (Chen et al., 2024) 107 (referred to as EMPSPEECH) consisting of 1,718 108 manually annotated English speech segments labeled as empathetic or neutral (Appendix B). 110

111Experimental Setup.Examples in EMP-112SPEECH are comprised of video segments113spanning three modalities: text (transcript), audio114(speech), and video.

We finetune two models per modality on the

Modality	Model	Accuracy	F1
Text	<b>RoBERTa</b> DeBERTa	<b>0.75</b> 0.69	<b>0.73</b> 0.68
Audio	HuBERT Wav2Vec2	<b>0.72</b> 0.68	<b>0.71</b> 0.63
Video	VideoMAE TimesFormer	<b>0.77</b> 0.64	<b>0.77</b> 0.62
Fusion (All Modalities)		0.76	0.72

Table 1: Performance of fine-tuned models acrossmodalities on the empathy classification task.

train set from EMPSPEECH: ROBERTA (Liu et al., 2019) and DEBERTA (He et al., 2021) for text, HU-BERT (Hsu et al., 2021) and WAV2VEC2 (Baevski et al., 2020) for audio, and VIDEOMAE (Tong et al., 2022) and TIMESFORMER (Bertasius et al., 2021) for video (Appendix A.1).<sup>1</sup> Then, we extract embeddings from each best-performing unimodal model (ROBERTA, HUBERT, and VIDEO-MAE, Table 1) to train a multimodal fusion model that projects all three modality embeddings into a shared latent space (Appendix A.2). 116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

**Results.** We evaluate all models (unimodal and multimodal) on the test split of EMPSPEECH to identify *disagreements*, or examples where two models with varying input modalities assign different labels, highlighting cases where different modalities may carry ambiguous, conflicting, or modality-specific signals.

Text disagrees most frequently with the other modalities (Table 2), suggesting that it does not frequently align with speech and visual cues. In contrast, audio and video are better aligned, likely due to a shared reliance on nonverbal expressive signaling such as prosody and facial expression. Our fusion model disagrees least with text, suggesting it relies more heavily on verbal content. This also may reveal an over-reliance on text-based signals by the original annotators.

Figure 2 visualizes disagreement regions between each unimodal model and the fusion model. We plot unimodal confidence (x-axis) against fusion confidence (y-axis) in the correct label; hence confidence greater than 0.5 resulted in a correct prediction. This yields four quadrants: green (multimodal correct, unimodal incorrect), red (multimodal incorrect, unimodal correct), blue (both correct), and yellow (both incorrect). Red and green quadrants are disagreement regions which we explore to identify complex examples.

<sup>&</sup>lt;sup>1</sup>The hidden layer dimensions of all models we consider are similar.



Figure 2: Comparing predictions between unimodal (text, audio, video) and multimodal models. We highlight regions where model predictions *agree* (blue and yellow quadrants) and disagree (red and green quadrants).

Modality	Text	Audio	Video
Text	_	0.338	0.318
Audio	0.338	_	0.253
Video	0.318	0.253	-
Full	0.214	0.383	0.331

Table 2: Pairwise disagreement rates between top unimodal models and between each unimodal model and the full model, measured as the proportion of test examples with differing predictions.

Feature	Red vs. Blue		Green vs. Blue	
	p-value Direction		p-value	Direction
valence arousal Mean Pitch dominance Min Pitch Jitter Max Intensity	0.0047 0.0065 0.0100 0.0108 0.0333 0.0347 0.1260	$\begin{array}{l} \mu_{blue} > \mu_{red} \\ \mu_{red} > \mu_{blue} \end{array}$	0.5166 0.0136 0.0001 0.0667 0.0001 0.0667 0.0023	$\begin{array}{l} \mu_{\rm green} > \mu_{\rm blue} \\ \mu_{\rm blue} > \mu_{\rm green} \\ \mu_{\rm green} > \mu_{\rm blue} \\ \mu_{\rm green} > \mu_{\rm blue} \end{array}$

Table 3: Statistically significant t-test results comparing red vs. blue and green vs. blue examples for audio features. See Appendix D for full table.

## 3.1 Modality-Based Feature Analysis

155

156

158

159

160

161

162

163

164

165

167

168

169

170

171

To better understand examples in disagreement regions, we extract and analyze modality-based human interpretable features.

Audio. We extract 12 prosodic and paralinguistic features from audio signals: 9 low-level acoustic features using PRAAT (Boersma and Weenink, 1992–2022) and PARSELMOUTH (Jadoul et al., 2018), and 3 high-level affective dimensions—valence, arousal, and dominance—using a finetuned WAV2VEC2 (Wagner et al., 2023). We compare feature distributions using t-tests for examples in disagreement quadrants (red and green) compared to those in the blue quadrant, signifying a non-ambiguous, easy examples. Blue examples have several significantly elevated pitch-related values than red examples (Table 3), suggesting that

AU	p (R vs B)	Dir	p (G vs B)	Direction
AU04	<b>0.0106</b>	red > blue	0.3682	green > blue
AU12	<b>0.0174</b>	blue > red	0.8977	green > blue
AU05	0.1837	blue > red	<b>&lt;0.0001</b>	<b>blue</b> > <b>green</b>

Table 4: Statistically significant t-test results comparing AU activation rates between red vs. blue and green vs. blue. See Appendix D for full table

172

173

174

175

176

177

178

180

181

182

183

185

186

187

188

189

190

191

192

193

194

195

197

198

199

stronger prosodic fluctuations are frequently corroborated by other modalities. Examples in the green quadrant show significantly higher *Max Intensity* than in blue, potentially reflecting the role of volume-based emphasis in aiding unimodal predictions. Furthermore, affective dimensions such as valence, arousal, and dominance are significantly lower in red examples, reinforcing the idea that red examples are not simply noisy, but structurally ambiguous: they express strong unimodal signals that are complicated when analyzed alongside other modalities.

Video. We examine facial action unit (AU) activations (Baltrušaitis et al., 2016) from video. AU04 (Brow Lowerer), AU12 (Lip Corner Puller), and AU05 (Upper Lid Raiser) show significant differences across example types, revealing how specific facial expressions contribute to perceptual ambiguity (Table 4). AU04 is more active in red examples than blue, indicating that despite its visually strong presence, its signal conflicts with other modalities. In contrast, AU12, associated with positive affect, and AU05, which is linked to attentiveness (Friesen and Ekman, 1978), both show greater activation in blue examples than in red and green, respectively, suggesting that these expressions may serve as clearer cues that are more consistently interpreted across modalities. Our findings indicate that fine-grained facial signals may contribute to perceptual complexity in the visual stream.



Figure 3: UMAP projections of text-only embeddings for empathetic (left) and neutral (right) examples, colored by modality disagreement class. Red and green points tend to cluster near the decision boundary, indicating high ambiguity.

**Text.** Finally, visualizing UMAP (Sainburg et al., 2021) projections of text embeddings (Figure 3) reveals that examples in disagreement regions (red and green) tend to cluster along the boundary between consistently correct (blue) and consistently incorrect (yellow) examples. Rather than forming distinct or isolated groups, disagreement examples appear in transition zones within the embedding space—areas where semantic cues are less definitive, supporting our hypothesis that red and green examples are inherently ambiguous and difficult and illustrating that modality disagreement is a reliable signal of challenging examples in empathy detection.

204

210

212

213 214

215

216

217

218

221

# 4 Experiment 2: Characterizing Complex Examples

We further assess whether model disagreements stem from data ambiguity by conducting a human annotation study to understand whether examples where models disagree are similarly challenging for human annotators.

Annotation Setup. We sample 204 examples evenly split across the four quadrants of each Figure 2 modality plot. For each example, annotators provide a binary judgment (empathetic or neutral) from a unimodal signal, then a judgment from the full multimodal version (see Appendix C for instructions), allowing us to track how human predictions shift with additional modality signals and understand the cognitive burden of multimodal integration. All examples were annotated by one author and one external annotator.

Results. Annotator *disagreement*, measured with
Cohen's Kappa (Cohen, 1960), can signal complex
phenomena in examples (Jiang and de Marneffe,
2022; Pavlick and Kwiatkowski, 2019) such as uncertainty in meaning leading to discrepancies in

Quadrant	Unimodal Judgment	Multimodal Judgment	Δ
Red	0.301	0.164	-0.137
Blue	0.379	0.646	0.267
Yellow	0.225	0.329	0.104
Green	0.482	0.218	-0.264

Table 5: Cohen's Kappa between internal and external annotators, computed separately for each quadrant and prediction round.

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

266

267

268

269

270

271

272

273

274

275

276

277

278

reasoning. In disagreement regions (red and green), we see a *decrease* in annotator agreement between unimodal and multimodal judgments (Table 5), indicating that humans diverge when weighing signals across modalities. In contrast, annotator agreement *improves* on examples where unimodal and multimodal model predictions are in agreement, supporting our these examples are relatively unambiguous and reliably interpreted once the full context is available (Table 5). Collectively, these results corroborate our hypothesis that modality disagreement can serve as a valuable signal for identifying ambiguous, challenging, or complex instances that are also difficult for human annotators.

# 5 Discussion and Conclusion

We have demonstrated how disagreement, both between modalities and between humans and models, can serve as a diagnostic lens to understand the complexity of multimodal empathy detection, challenging the assumption that more signal from other modalities reliably yields better performance. Our analysis reveals that disagreement between unimodal and multimodal models is often not arbitrary, but instead marks the presence of subtle, ambiguous, or context-sensitive cues that challenge fusion models and human annotators alike. Our findings emphasize the necessity of high-quality, contextsensitive annotation in socially complex tasks like empathy detection, where model errors may reflect genuine human uncertainty. This framework provides a scalable method for identifying ambiguity and enhancing model reliability, especially in recognizing complex emotional states that involve inherent disagreement and uncertainty. Our work lays the foundation for several directions of future work, such as creating adversarial test sets to evaluate empathy detection systems in realistic scenarios or the identification of challenging examples for human annotation in an active learning setup to improve model robustness.

## 6 Limitations

279

281

282

293

294

297

301

303

310

315

316

317

319

324

325

We acknowledge several limitations in our study. Our analyses are based on a limited dataset and a small number of human annotators. Given that empathy is inherently subjective, annotations may vary due to individual interpretations, potentially introducing biases rather than reflecting universal properties of the data. Additionally, we rely on a single dataset, and future work should investigate whether the patterns we observe hold across other datasets and domains.

Our data is also derived from U.S.-based, English-language television and interview content. As such, the generalizability of our findings to multilingual or culturally diverse settings may be limited. Future research should investigate these patterns in varied cultural and linguistic environments to better assess the broader applicability of our conclusions.

## 7 Ethics Statement

We use a publicly available dataset and strictly use open-source models for analysis.

All annotations were conducted by an author and an individual affiliated with the research team. No participants were recruited via crowdsourcing or external platforms, and no monetary compensation was provided, as the annotators were contributing in a research capacity. We provide detailed information on what we ask the annotators to annotate and how we plan to use the data. The annotators willingly agreed to participate with full knowledge of the task. No sensitive or identifying information were collected from annotators.

We note that empathy expression may vary across cultures, and our findings may not generalize to non-English or non-Western contexts. We encourage future work to explore these questions in more diverse settings.

We will release all code and experimental resources at https://anonymous.4open.sc ience/r/multimodal-empathy-disag reement-F48B to support reproducibility.

### 321 References

 Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. In Advances in Neural Information Processing Systems (NeurIPS). Tadas Baltrušaitis, Peter Robinson, and Louis-Philippe Morency. 2016. Openface: An open source facial behavior analysis toolkit. In 2016 IEEE Winter Conference on Applications of Computer Vision (WACV), pages 1–10. 327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

346

347

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

366

367

368

369

370

371

372

373

374

375

376

377

378

381

- Roy F. Baumeister and Kathleen D. Vohs. 2007. *Encyclopedia of Social Psychology*, volume 1. Sage.
- Gedas Bertasius, Heng Wang, and Lorenzo Torresani. 2021. Is space-time attention all you need for video understanding? In *Proceedings of the 38th International Conference on Machine Learning (ICML).*
- Paul Boersma and David Weenink. 1992–2022. *Praat: doing phonetics by computer [Computer program]*. Version 6.2.14, retrieved 24 May 2022.
- Run Chen, Haozhe Chen, Anushka Kulkarni, Eleanor Lin, Linda Pang, Divya Tadimeti, Jun Shin, and Julia Hirschberg. 2024. Detecting empathy in speech. In *Proceedings of Interspeech 2024*, Dublin, Ireland. ISCA.
- J. Cohen. 1960. A Coefficient of Agreement for Nominal Scales. *Educational and Psychological Measurement*, 20(1):37.
- Wallace V. Friesen and Paul Ekman. 1978. Facial Action Coding System: A technique for the measurement of facial movement. Consulting Psychologists Press.
- Pascale Fung, Dario Bertero, Yan Wan, Anik Dey, Ricky Ho Yin Chan, Farhad Bin Siddique, Yang Yang, Chien-Sheng Wu, and Ruixi Lin. 2016. Towards empathetic human-robot interactions. *arXiv preprint arXiv:1605.04072*.
- Lucie Galland, Catherine Pelachaud, and Florian Pecune. 2024. Emmi: Empathic multimodal motivational interviews dataset: Analyses and annotations. *arXiv preprint arXiv:2406.16478*.
- Md Rakibul Hasan, Md Zakir Hossain, Shreya Ghosh, Aneesh Krishna, and Tom Gedeon. 2023. Empathy detection from text, audiovisual, audio or physiological signals: A systematic review of task formulations and machine learning methods. *arXiv preprint arXiv:2311.00721*.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: Decoding-enhanced bert with disentangled attention. *Proceedings of the International Conference on Learning Representations (ICLR).*
- Judith Holler and Stephen C. Levinson. 2019. Multimodal language processing in human communication. *Trends in Cognitive Sciences*, 23(8):639–652.
- Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. In Proceedings of the Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL).

- 386 394 395 396 397 400 401

- 402 403
- 404 405
- 406 407
- 408 409
- 410 411
- 412
- 413 414
- 415 416 417

418 419

420 421

422

423 424

425 426

427 428

429 430

435

436 437

- Summaira Jabeen, Xi Li, Amin Muhammad Shoib, Songyuan Li, and Abdul Jabbar. 2021. Recent advances and trends in multimodal deep learning: A review. arXiv preprint arXiv:2105.11087.
- Yannick Jadoul, Bart de Boer, and Peter Thompson. 2018. Introducing parselmouth: A python interface to praat. Journal of Phonetics, 71:1–15.
- Nan-Jiang Jiang and Marie-Catherine de Marneffe. 2022. Investigating reasons for disagreement in natural language inference. Transactions of the Association for Computational Linguistics, 10:1357–1374.
- Qintong Li, Hongshen Chen, Zhaochun Ren, Pengjie Ren, Zhaopeng Tu, and Zhumin Chen. 2019. Empdg: Multiresolution interactive empathetic dialogue generation. arXiv preprint arXiv:1911.08698.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach.
- Ellie Pavlick and Tom Kwiatkowski. 2019. Inherent disagreements in human textual inferences. Transactions of the Association for Computational Linguistics, 7:677-694.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards empathetic opendomain conversation models: A new benchmark and dataset. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 5370-5381, Florence, Italy. Association for Computational Linguistics.
- Swarnadeep Saha, Peter Hase, Nazneen Rajani, and Mohit Bansal. 2022. Are hard examples also harder to explain? a study with human and model-generated explanations. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 2121-2131, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Tim Sainburg, Leland McInnes, and Timothy O Gentner. 2021. Parametric umap embeddings for representation and semisupervised learning. Neural Computation, 33(11):2881-2907.
- Jocelyn Shen, Yubin Kim, Mohit Hulse, Wazeer Zulfikar, Sharifa Alghowinem, Cynthia Breazeal, and Hae Park. 2024. EmpathicStories++: A multimodal dataset for empathy towards personal experiences. In Findings of the Association for Computational Linguistics: ACL 2024, pages 4525-4536, Bangkok, Thailand. Association for Computational Linguistics.
- Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A. Smith, and Yejin Choi. 2020. Dataset cartography: Mapping and diagnosing datasets with training dynamics. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP),

pages 9275-9293, Online. Association for Computational Linguistics.

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

- Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. 2022. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pretraining. In Proceedings of the 33rd Conference on Neural Information Processing Systems (NeurIPS).
- Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J. Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2019. Multimodal transformer for unaligned multimodal language sequences. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 6558-6569, Florence, Italy. Association for Computational Linguistics.
- Johannes Wagner, Andreas Triantafyllopoulos, Hagen Wierstorf, Maximilian Schmitt, Felix Burkhardt, Florian Eyben, and Björn W. Schuller. 2023. Dawn of the transformer era in speech emotion recognition: Closing the valence gap. IEEE Transactions on Pattern Analysis and Machine Intelligence, 45(9):10745-10759.
- Junpeng Wang, Liang Wang, Yan Zheng, Chin-Chia Michael Yeh, Shubham Jain, and Wei Zhang. 2023. Learning-From-Disagreement: A Model Comparison and Visual Analytics Framework . IEEE Transactions on Visualization & Computer Graphics, 29(09):3809-3825.
- Amir Zadeh, Minghai Chen, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2017. Tensor fusion network for multimodal sentiment analysis. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pages 1103-1114, Copenhagen, Denmark. Association for Computational Linguistics.
- Yiqun Zhang, Fanheng Kong, Peidong Wang, Shuang Sun, SWangLing SWangLing, Shi Feng, Daling Wang, Yifei Zhang, and Kaisong Song. 2024. STICKERCONV: Generating multimodal empathetic responses from scratch. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 7707-7733, Bangkok, Thailand. Association for Computational Linguistics.
- Zhou'an Zhu, Xin Li, Jicai Pan, Yufei Xiao, Yanan Chang, Feiyi Zheng, and Shangfei Wang. 2023. Medic: A multimodal dataset for empathic dialogue in counseling. arXiv preprint arXiv:2305.14221.

583

584

585

586

# A Model Training Details

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

504

505

507

508

510

511

512

513

514

515

516

517

518

519

522

524

526

527

Data was split into train, test and validation sets using random sampling, with an 80-10-10 split. We run fine-tuning and inference for all open-source models on an A100 GPU in Google Colab.

#### A.1 Unimodal Model Training Details

Each model is trained on a binary empathy classification task using precomputed 768-dimensional embeddings. We freeze all but the final two transformer layers and train for 15 epochs with a learning rate of 5e-6 and batch size of 8.

#### A.2 Fusion Model Details

Each unimodal model representation is independently gated and passed through an attention mechanism that computes modality-specific weights.
The weighted embeddings are aggregated and classified using a three-layer feedforward network with max pooling. The fusion model is trained for 10 epochs using a learning rate of 1e-4 and includes modality dropout during training.

# **B** Dataset Details

We use a multimodal empathy dataset (Chen et al., 2024) consisting of 346 English-language videos totaling approximately 53 hours, collected from YouTube between 2020 and 2022 using keywords like "empathy" and "empathetic training." The dataset includes empathy training sessions, therapy roleplays, interviews, TED Talks, and TV/movie scenes, comprising both acted (62%) and spontaneous (38%) speech. Each video was labeled by at least three expert annotators as either empathetic or neutral, with final labels determined by majority vote. Metadata such as speaker gender, topic, and emotional context was manually annotated, covering themes like therapy, parenting, workplace dynamics, and social relationships. From this collection, a subset of 65 videos was transcribed, diarized, and manually re-aligned using Praat to ensure accurate speaker segmentation and time alignment. This process resulted in 1,718 annotated segments with speaker labels, timestamps, transcripts, and empathy stage annotations, enabling fine-grained analysis of empathy in naturalistic and semi-scripted settings.

#### C Annotation Instructions

We employed two annotators, one of the paper's authors and an non-author, both fluent English speakers based in the United States. No additional demographic information was collected, as the annotation was conducted internally for research purposes.

Annotators were asked to provide two judgments per example, labeling each as either empathetic or neutral (Figure 4). A excerpt describing empathy (drawn from the Encyclopedia of Social Psychology, Volume 1, (Baumeister and Vohs, 2007)) was provided to ensure a consistent conceptual foundation for annotation:

Empathy is often defined as understanding another person's experience by imagining oneself in that other person's situation: One understands the other person's experience as if it were being experienced by the self, but without the self actually experiencing it. There are three commonly studied components of emotional empathy. The first is feeling the same emotion as another person (sometimes attributed to emotional contagion, e.g., unconsciously "catching" someone else's tears and feeling sad oneself). The second component, personal distress, refers to one's own feelings of distress in response to perceiving another's plight. The third emotional component, feeling compassion for another person, is the one most frequently associated with the study of empathy. Cognitive empathy refers to the extent to which we perceive or have evidence that we have successfully guessed someone else's thoughts and feelings.

Annotators were given an annotation flag indicating which modality to use for the first pass; for instance, if the flag was text, only the transcript was to be used to make the first prediction. After submitting the first judgment, annotators were then given access to the full video, including all available audio, visual, and textual information. They were then asked to provide a second prediction.

#### **D** Full Feature Comparisons

Tables 6, 7 and 8 provide additional results from the t-tests comparing examples across different confidence quadrants. Table 6 provides an internal comparison between the disagreement quadrants. Table 7 presents the full version of the audio feature comparisons summarized in Table 3. Table 8 expands on the facial feature comparisons shown in Table 4.

## **E** Feature Distributions

Figures 5 and 6 visualize the distributions of key features across confidence quadrants. Figure 5 presents the distribution of selected audio features (e.g., pitch, intensity) for red, green, and blue examples, highlighting acoustic patterns associated



Figure 4: Annotation interface

Feature	t-stat	p-value	Mean Comparison
Mean Pitch	2.453	0.0159	$\mu_{ m red} > \mu_{ m green}$
Max Intensity	-2.124	0.0366	$\mu_{ ext{green}} > \mu_{ ext{red}}$
Max Pitch	2.016	0.0465	$\mu_{ m red} > \mu_{ m green}$
Min Pitch	2.007	0.0475	$\mu_{ m red} > \mu_{ m green}$
valence	-1.908	0.0593	$\mu_{ m green} > \mu_{ m red}$
arousal	1.827	0.0705	$\mu_{ m red} > \mu_{ m green}$
speaking_rate	1.773	0.0807	$\mu_{ m red} > \mu_{ m green}$
dominance	1.712	0.0899	$\mu_{ m red} > \mu_{ m green}$
Shimmer	0.773	0.4416	$\mu_{ m red} > \mu_{ m green}$
Jitter	0.622	0.5355	$\mu_{ m red} > \mu_{ m green}$
Mean Intensity	0.544	0.5886	$\mu_{ m red} > \mu_{ m green}$
HNR	0.508	0.6129	$\mu_{ m red} > \mu_{ m green}$
Min Intensity	-0.429	0.6685	$\mu_{ m green} > \mu_{ m red}$

Table 6: T-test results comparing audio features between red and green examples. Statistically significant results are bolded.

with model disagreement. Figure 6 shows acti-587 vation rates for facial Action Units (AUs) in red, 588 green, and blue examples, illustrating how specific 589 facial expressions vary across agreement condi-590 tions. These visualizations complement the statis-591 592 tical comparisons reported in Tables 7 and 8, providing a more interpretable view of the underlying 593 feature dynamics. 594

Feature	p (Red vs Blue)	Direction	p (Green vs Blue)	Direction
valence	0.0047	$\mu_{ m blue} > \mu_{ m red}$	0.5166	$\mu_{\rm green} > \mu_{\rm blue}$
arousal	0.0065	$\mu_{ m blue} > \mu_{ m red}$	0.0136	$\mu_{ m blue} > \mu_{ m green}$
Mean Pitch	0.0100	$\mu_{ ext{blue}} > \mu_{ ext{red}}$	0.0001	$\mu_{ m blue} > \mu_{ m green}$
dominance	0.0108	$\mu_{ m blue} > \mu_{ m red}$	0.0667	$\mu_{ m blue} > \mu_{ m green}$
Min Pitch	0.0333	$\mu_{ ext{blue}} > \mu_{ ext{red}}$	0.0001	$\mu_{ m blue} > \mu_{ m green}$
Jitter	0.0347	$\mu_{ m red} > \mu_{ m blue}$	0.0667	$\mu_{ m green} > \mu_{ m blue}$
Max Intensity	0.1260	$\mu_{ m red} > \mu_{ m blue}$	0.0023	$\mu_{ ext{green}} > \mu_{ ext{blue}}$
Mean Intensity	0.1599	$\mu_{ m red} > \mu_{ m blue}$	0.5329	$\mu_{ m blue} > \mu_{ m green}$
HNR	0.2217	$\mu_{ m blue} > \mu_{ m red}$	0.2055	$\mu_{ m blue} > \mu_{ m green}$
speaking_rate	0.2723	$\mu_{ m blue} > \mu_{ m red}$	0.9991	$\mu_{\rm green} > \mu_{\rm blue}$
Shimmer	0.4122	$\mu_{ m red} > \mu_{ m blue}$	0.1541	$\mu_{\rm blue} > \mu_{\rm green}$
Max Pitch	0.6845	$\mu_{ m red} > \mu_{ m blue}$	0.2647	$\mu_{\rm blue} > \mu_{\rm green}$
Min Intensity	0.7999	$\mu_{ m blue} > \mu_{ m red}$	0.1571	$\mu_{\rm blue} > \mu_{\rm green}$

Table 7: T-test results comparing audio features between red vs. blue and green vs. blue examples. Statistically significant p-values are bolded.

AU	p (Red vs Blue)	Direction	p (Green vs Blue)	Direction
AU04: Brow Lowerer	0.0106	red > blue	0.3682	green > blue
AU12: Lip Corner Puller	0.0174	blue > red	0.8977	green > blue
AU05: Upper Lid Raiser	0.1837	blue > red	<0.0001	blue > green
AU17: Chin Raiser	0.2256	red > blue	0.9802	blue $>$ green
AU10: Upper Lip Raiser	0.2275	blue $>$ red	0.6700	green > blue
AU45: Blink	0.3200	blue > red	0.7462	green > blue
AU07: Lid Tightener	0.3252	blue > red	0.9318	blue > green
AU14: Dimpler	0.4593	red > blue	0.0652	green > blue
AU20: Lip Stretcher	0.5701	blue > red	0.7907	blue > green
AU09: Nose Wrinkler	0.6211	blue > red	0.7639	green > blue
AU25: Lips Part	0.6227	blue > red	0.7492	blue > green
AU01: Inner Brow Raiser	0.6529	blue > red	0.4674	green > blue
AU23: Lip Tightener	0.6630	red > blue	0.3474	green > blue
AU28: Lip Suck	0.6735	red > blue	0.9846	green > blue
AU26: Jaw Drop	0.6851	red > blue	0.4596	blue > green
AU06: Cheek Raiser	0.7097	blue > red	0.3201	green > blue
AU15: Lip Corner Depressor	0.9528	red > blue	0.4834	green > blue
AU02: Outer Brow Raiser	0.9647	blue > red	0.6677	green > blue

Table 8: T-test results comparing AU activation rates between red vs. blue and green vs. blue. Bolded p-values are statistically significant.



Figure 5: Distribution of audio features for red, green and blue examples across the confidence quadrants. Red examples are those correctly classified by the unimodal audio model but misclassified by the multimodal model; green examples represent the reverse. Blue examples represent those correctly classified by both the unimodal audio model and the multimodal model. Significant differences appear in pitch and intensity-based features.



Figure 6: AU activation rates for red, green, and blue examples. Red bars indicate examples where the unimodal visual model predicted correctly but the multimodal model did not (Red: Unimodal > 0.5, Multimodal < 0.5). Green bars show the reverse. Blue bars indicate examples where both the unimodal and multimodal models correctly predicted the label.