

# The Correlation Between Emotion in Text and Speech Segments is Limited: A Cross-Modal Study

**David Lindevelt**  
Leiden University  
Leiden, The Netherlands  
david.lindevelt@gmail.com

**Suzan Verberne**  
Leiden University  
Leiden, The Netherlands  
s.verberne@liacs.leidenuniv.nl

**Joost Broekens**  
Leiden University  
Leiden, The Netherlands  
joost.broekens@gmail.com

## Abstract

Although expressive TTS systems aim to capture human-like emotion, little is known about how well emotional signals in text correspond to those in speech. In this short paper, we investigate how emotion (Valence, Arousal, Dominance) in text relates to emotion in speech. We use 8 large language models for identifying emotion in text and two audio models for emotion in speech, across three genres: Podcasts, Audiobooks and TED talks. Findings show that while language models perform well on emotion recognition from situational text, and the audio models perform well on speech, they show a strong correlation for Valence only. Further, the genre of the content significantly impacts the correlation: audiobooks exhibit higher text-audio correlation than TED talks. Finally, we show that more context for LLMs fails to improve this correlation between text and speech emotion prediction. Our results highlight that emotional signals in text do not correspond well to those in speech: emotion prediction from text alone is insufficient for emotional TTS.

## 1 Introduction

Text-to-speech (TTS) systems have made remarkable progress in recent years, with models achieving near-human-like naturalness in speech synthesis (Li et al., 2023; Lee et al., 2025; Manku et al., 2025). However, one of the remaining challenges in this field is accurate and controllable emotions in synthesized speech. This challenge stems from a fundamental training–inference mismatch. TTS systems learn emotional patterns from text–audio pairs, but at inference time they must rely only on text to generate emotionally appropriate speech. For this to work, the emotional information in the text must be related to the emotion expressed in the corresponding speech audio. We define the *emotion relation* as a correlation between segment-based emotion values from text and speech modalities.

We analyze how robust this relation is to changes in context, such as textual genre, and, whether more textual information improves this emotion relation.

The dimensional model of emotions (Russell and Mehrabian, 1977; Bradley and Lang, 1994), which includes three dimensions: Valence (positive – negative), Arousal (active – passive), and Dominance (dominant – submissive), is a commonly used, psychologically grounded representation for emotional states in both text and speech. Recent advances in large language models (LLMs) have shown promising results in assessing dimensional emotions from text (Broekens et al., 2023; Lei et al., 2023; Wu et al., 2025). Our research builds upon these developments. We investigate to what extent the acoustic emotional content relates to the textual emotional content in narrations, such as in audiobooks, podcasts and public lectures.

With this paper, we make three contributions: (1) We develop an analysis pipeline for extracting emotion values from text using state-of-the-art language models and assess their correlation with emotion values derived from corresponding audio narrations. (2) We evaluate the impact of both genre and contextual information (pre-text and post-text) on emotion prediction accuracy, simulating the contextual understanding that human narrators employ. (3) We show a weak relation between emotions identified in speech and text and thereby reveal an important limitation of current text to speech models that aim to express emotions in a natural way based only on text input.

## 2 Related work

Emotion prediction is an important topic in the fields of affective computing, speech synthesis and natural language processing. Good performance of speech emotion recognition (SER) can enhance human-machine interaction (Mohmad and Delhibabu, 2024).

Peng et al. (2021) employ an LSTM architecture operating on audio-only cochleagram features for a dimensional SER task. With the success of attention mechanisms researchers focused on transformer architecture to capture feature-rich representation of audio. Wagner et al. (2023) fine-tuned a wav2vec 2.0-based model (Baeovski et al., 2020) to predict emotions from speech audio.

Separate from SER, there is extensive research on sentiment analysis and emotion classification from text. Extracting emotions from text has been predominantly investigated with encoder-only (BERT-like) architectures (Siriwardhana et al., 2020; Acheampong et al., 2021). Recently, decoder-only language models have been explored (Broekens et al., 2023) for assessment of emotions in text, showing high correlation with human annotations. Wang et al. (2023) investigated emotional intelligence (EI) of LLMs using psychological tests and showed LLMs can achieve high EI scores.

Prior research showed that from text, Valence value predictions are the most feasible, while acoustic features proved to be better for Arousal prediction. To compensate for poor Valence prediction, audio-based prediction models incorporate text and audio modalities together in their architecture (Srinivasan et al., 2022; Atmaja and Akagi, 2020; Triantafyllopoulos et al., 2023).

In our analysis, we make use of decoder-only transformer models for predicting dimensional emotions from text and compare the output with speech-audio-only model predictions. We also examine the impact of text genre and additional textual context in prediction performance.

### 3 Experimental Methods

We first reproduce and extend the LLM-based emotion prediction capabilities reported in prior work (Broekens et al., 2023). We use a larger set of generative LLMs in the comparison. We compare different model sizes on the same dataset as presented in the paper, the ANET emotional situations dataset with human annotations. We slightly change the prompt for language models used in previous work and applied chain-of-thought prompting strategy with structured output when we predict emotion values using LLMs. The exact prompts used in the prediction process are adapted from prior work and provided in Appendix A.

Further, we reproduce the results of an open-source SOTA speech emotion recognition model

on the MSP Podcast test set and compare predicted emotions against human annotated emotions.

Then, we compare text and speech emotions for three data sets representing different genres: podcasts, audio books and TED talks. To do so, we use the same set of language models to predict emotions from text and compare the predicted emotions against the speech-based emotion predictions. For these comparisons, we quantify the relation between speech-based and text-based emotion predictions by computing the Pearson correlation coefficient for each emotion dimension (Arousal, Dominance, and Valence) for each dataset.

Finally, we test whether providing more context, in the form of more text before and after the target sentence, increases the performance of LLMs.

#### 3.1 Models

We select the wav2vec2-based model<sup>1</sup> as SOTA audio model for emotion prediction from speech audio. We also include a second WavLM-based model<sup>2</sup> to strengthen our analysis and evaluate generalizability. We look at how closely their predictions match with human annotations and their agreement across datasets.

We use 8 language models using the ollama<sup>3</sup> framework. These are llama3.3:70b, llama3.1:70b, phi3:14b, gemma2:9b, llama3.1:8b, qwen2:7b, openchat:7b, gemma:7b models. Each inference session follows default settings with the seed value set by experiment run. In total, 5 runs with different seed values are performed. We report the average of five different runs. The bold values in the tables represent statistically significant ( $p < 0.01$ ) values for at least 80% (4/5) of the runs.

#### 3.2 Datasets and pre-processing

We use four datasets for our analysis: 1. **ANET** (Bradley and Lang, 2007): 120 sentences describing emotional situations with dimensional emotion annotations (Pleasure, Arousal and Dominance) by human subjects. 2. **MSP Podcast** (Lotfian and Busso, 2017) version 1.8: 20,539 audio samples from podcasts with human emotion annotations on dimensions Valence, Arousal and Dominance. 3. **Libriheavy** (Kang et al., 2024): 20,539 samples from ‘small’ subset of dataset. It is an annotated

<sup>1</sup><https://huggingface.co/audeering/wav2vec2-large-robust-12-ft-emotion-msp-dim>

<sup>2</sup><https://huggingface.co/3loi/SER-Odyssey-Baseline-WavLM-Multi-Attributes>

<sup>3</sup><https://ollama.com>

version of Librilight (Kahn et al., 2020) dataset that consists of public domain audiobooks; it includes pre-text information for each segment as extended context. 4. **TEDLIUM** corpus release 2 (Rousseau et al., 2014): 89,986 samples from the train set of English speech recognition corpus from public TED talks.<sup>4</sup> For each sample, we keep preceding and following texts as extended context.

We present details on emotion annotations and emotionally balanced sample selection process for all datasets in Appendix B.

The speech-based audio-only models used in our experiments are trained on the train set of MSP Podcast dataset. We use the combined MSP Podcast test sets (test1 and test2) and selected emotionally balanced subset for our analysis. The MSP Podcast dataset does not include transcriptions. We use the OpenAI Whisper Turbo model (Radford et al., 2023) to transcribe the audio data. Emotion ratings in the MSP Podcast dataset range between 1 (very negative/calm/weak) and 7 (very positive/active/strong), for the dimensions Valence, Arousal, and Dominance. We apply min-max scaling to reduce their range to 0-1, using 1 and 7 as the defined minimum and maximum values.

## 4 Results

We describe our results along the lines of four research questions.

**How good are LLMs at extracting emotions from a given text?** Before comparing emotions extracted from text and speech, we first assess the capability of language models to predict dimensional emotions from text. Extending previous work (Broekens et al., 2023), we use 8 language models to predict dimensional emotions for the ANET dataset.

In this experiment, we evaluate both large and small models by comparing their predictions with ground truth text emotions annotated by human subjects. Table 1 reports average agreement over 5 different seeds. Our results show that all models achieve statistically significant and strong correlation with the human labels. Large models exhibit stronger correlations than smaller versions. This experiment confirms that decoder-only language models are able to reliably extract emotion dimensions in a given situational text, particularly when the text explicitly contains affective information, as is the case for the ANET dataset.

<sup>4</sup><https://www.ted.com>

Table 1: Correlation between text-based prediction and human ratings – ANET Dataset. Boldface indicates a significance level of  $p < 0.01$  in 80% of experiment runs across 5 seeds

	Arousal	Dominance	Valence
llama3.3:70b	<b>.929</b> <sub>(.003)</sub>	<b>.865</b> <sub>(.004)</sub>	<b>.966</b> <sub>(.002)</sub>
llama3.1:70b	<b>.918</b> <sub>(.008)</sub>	<b>.891</b> <sub>(.010)</sub>	<b>.967</b> <sub>(.009)</sub>
phi3:14b	<b>.822</b> <sub>(.043)</sub>	<b>.828</b> <sub>(.040)</sub>	<b>.870</b> <sub>(.053)</sub>
gemma2:9b	<b>.888</b> <sub>(.008)</sub>	<b>.666</b> <sub>(.026)</sub>	<b>.870</b> <sub>(.037)</sub>
llama3.1:8b	<b>.846</b> <sub>(.015)</sub>	<b>.508</b> <sub>(.048)</sub>	<b>.924</b> <sub>(.013)</sub>
qwen2:7b	<b>.553</b> <sub>(.017)</sub>	<b>.709</b> <sub>(.060)</sub>	<b>.950</b> <sub>(.017)</sub>
openchat:7b	<b>.740</b> <sub>(.032)</sub>	<b>.699</b> <sub>(.054)</sub>	<b>.914</b> <sub>(.046)</sub>
gemma:7b	<b>.803</b> <sub>(.038)</sub>	<b>.392</b> <sub>(.044)</sub>	<b>.901</b> <sub>(.017)</sub>

**How does text-based emotion prediction relate to speech-based emotions?** To quantify the relation between text-based and speech-based dimensional emotion, we use the MSP podcast dataset (Lotfian and Busso, 2017). First, we replicate the in-context performance evaluation of the open-source speech-based model as reported in Wagner et al. (2023). Our replication of their evaluation gave correlations of 0.772 (Arousal), 0.640 (Dominance), and 0.646 (Valence) on the MSP Podcast combined test set, validating the model’s performance for speech emotion prediction.

Next, we compare text-based emotion predictions with speech-based emotion annotation and speech-based prediction for podcast transcripts. We find relatively high correlation for Valence (highest 0.630, llama3.1:70b) while performing substantially lower on Arousal (highest 0.384, llama3.3:70b) and Dominance (highest 0.117, gemma2:9b). Valence and Arousal correlations are statistically significant while Dominance shows no statistical significance. This is true for both the correlations between text-based prediction and ground-truth data, and text-based and speech-based predictions (see Table 4 in Appendix D). In other words, predicted text emotions correspond well with speech emotions in terms of Valence, but are less effective at capturing emotional cues related to Arousal and Dominance. This suggests that while both text-based and speech-based predictions are by themselves reliable, there seems to be weak relation and perhaps a genre dependency.

**Does text genre impact the relation between emotions predicted from text and speech?** To further investigate the effect of genre on the rela-

tion between text-based and speech-based emotion predictions, we repeated the analysis for audiobooks (Libriheavy) and TED talks (TEDLIUM). Since these datasets lack human annotations for text or speech emotions, and the speech-based prediction model worked well for the MSP podcast data, we treat predictions from SOTA audio model as "ground truth" for calculating the emotion relation. For comparison purposes we do the same for the podcast dataset. We find that the correlation for Valence is consistently higher than for Arousal and Dominance. Further, Dominance correlation is consistently low. Also, Valence and Arousal correlations between text-based and speech-based predictions are higher for audiobooks than the other two genres. Tables 5 and 7 in Appendices E and F, respectively, provide detailed results.

**Does additional context improve the relation between text and audio emotion predictions?** A potential reason for lower correlation on Arousal and Dominance is the lack of textual context. We test whether providing more context improves this relation. We use the same samples from both the audiobooks and TED talks, but we extend their context. For audiobooks, we combine pre-text with the main text, and for TEDLIUM, we include preceding and following sentences. LLMs are then prompted to extract emotions from the main text, considering the full context (see Appendix A).

Tables 2 and 3 show that while extended context may marginally enhance the relation in Arousal for the larger LLM models in the audiobook data, it does not significantly improve and may even slightly weaken the relation for Valence. This indicates that the benefit of additional context may depend on the emotion dimension and the nature of the dataset, and is not conclusive.

**Do our findings generalize to another audio model?** We find that the two audio models show high agreement with each other and show similar levels of agreement with human annotations, confirming our findings. Agreements measured with Pearson correlation coefficient and presented in Tables 9 (between the two models) and 10 (between models and human annotation) in Appendix G.

## 5 Conclusions

Language models excel on the emotion recognition task from situational text, but their predictions correlate poorly with speech-based emotions, partic-

Table 2: Correlation between *contextual* text and speech emotion prediction on Libriheavy

	Arousal	Dominance	Valence
llama3.3:70b	<b>.416</b> <sub>(.016)</sub>	.175 <sub>(.054)</sub>	<b>.700</b> <sub>(.006)</sub>
llama3.1:70b	<b>.427</b> <sub>(.021)</sub>	.015 <sub>(.025)</sub>	<b>.697</b> <sub>(.016)</sub>
phi3:14b	<b>.392</b> <sub>(.080)</sub>	.108 <sub>(.064)</sub>	<b>.591</b> <sub>(.066)</sub>
gemma2:9b	<b>.450</b> <sub>(.015)</sub>	<b>.297</b> <sub>(.026)</sub>	<b>.641</b> <sub>(.023)</sub>
llama3.1:8b	<b>.473</b> <sub>(.064)</sub>	.119 <sub>(.117)</sub>	<b>.642</b> <sub>(.040)</sub>
qwen2:7b	<b>.267</b> <sub>(.067)</sub>	.064 <sub>(.044)</sub>	<b>.679</b> <sub>(.016)</sub>
openchat:7b	<b>.320</b> <sub>(.099)</sub>	.006 <sub>(.037)</sub>	<b>.591</b> <sub>(.089)</sub>
gemma:7b	.230 <sub>(.052)</sub>	.144 <sub>(.100)</sub>	<b>.549</b> <sub>(.053)</sub>

Table 3: Correlation between *contextual* text and speech emotion prediction on TEDLIUM

	Arousal	Dominance	Valence
llama3.3:70b	.210 <sub>(.011)</sub>	.214 <sub>(.028)</sub>	<b>.627</b> <sub>(.007)</sub>
llama3.1:70b	.223 <sub>(.030)</sub>	.174 <sub>(.022)</sub>	<b>.624</b> <sub>(.017)</sub>
phi3:14b	.244 <sub>(.058)</sub>	.110 <sub>(.072)</sub>	<b>.538</b> <sub>(.077)</sub>
gemma2:9b	.135 <sub>(.029)</sub>	.116 <sub>(.058)</sub>	<b>.586</b> <sub>(.021)</sub>
llama3.1:8b	.208 <sub>(.041)</sub>	.217 <sub>(.037)</sub>	<b>.590</b> <sub>(.019)</sub>
qwen2:7b	<b>.340</b> <sub>(.041)</sub>	<b>.292</b> <sub>(.057)</sub>	<b>.647</b> <sub>(.011)</sub>
openchat:7b	<b>.237</b> <sub>(.058)</sub>	.181 <sub>(.054)</sub>	<b>.546</b> <sub>(.106)</sub>
gemma:7b	.155 <sub>(.043)</sub>	.048 <sub>(.063)</sub>	<b>.504</b> <sub>(.012)</sub>

ularly on the Arousal and Dominance in all tested genres. This shows that text alone is insufficient for emotional TTS. Second, genre impacts emotion relation, with audiobooks showing better text-speech correlation than TED talks and podcasts. Third, contrary to our expectation, increasing context did not improve the relation between text and speech emotions. Our results confirm the view of [Gunes and Schuller \(2013\)](#) that Activation (Arousal) is easily captured through acoustic features, Valence (positivity) is more reliably extracted from linguistic features, highlighting that emotional signals are not always shared across modalities.

Our findings motivate the development of text-based emotion predictors that match the emotion expressed in speech, but also show a need for improved TTS conditioning strategies that accurately control affective expression. Detailed understanding of the relation of textual and acoustic emotional content is essential for next generations of TTS. This will help natural-sounding speech synthesis but also enable better control over the subtle emotional nuances present in the source text, much like a skilled human narrator would have.

## Limitations

For Libriheavy and TEDLIUM, we rely on an open-source audio-only model trained on podcasts, potentially introducing bias. Additionally, we use English-only data to conduct the analysis, future work should include human annotations, multilingual data and consider alternative metrics to account for imbalanced emotion distributions, particularly for extreme values that may skew Pearson correlation results. Further, binned emotion distributions for MSP Podcast, Libriheavy audiobooks, and TEDLIUM (see Appendix C) reveal that audiobooks show more neutral emotions, while TEDLIUM emotions in higher bins similarly to MSP Podcast; all show limited coverage of edge emotions. This may influence the correlations.

## References

- Francisca Adoma Acheampong, Henry Nunoo-Mensah, and Wenyu Chen. 2021. Transformer models for text-based emotion detection: a review of bert-based approaches. *Artificial Intelligence Review*, 54(8):5789–5829.
- Bagus Tris Atmaja and Masato Akagi. 2020. Improving valence prediction in dimensional speech emotion recognition using linguistic information. In *2020 23rd Conference of the Oriental COCOSDA International Committee for the Co-ordination and Standardisation of Speech Databases and Assessment Techniques (O-COCOSDA)*, pages 166–171. IEEE.
- 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*, 33:12449–12460.
- Margaret M Bradley and Peter J Lang. 1994. Measuring emotion: the self-assessment manikin and the semantic differential. *Journal of behavior therapy and experimental psychiatry*, 25(1):49–59.
- M.M. Bradley and P.J. Lang. 2007. Affective norms for english text (anet): Affective ratings of text and instruction manual. *Technical Report. D-1, University of Florida, Gainesville, FL*.
- Joost Broekens, Bernhard Hilpert, Suzan Verberne, Kim Baraka, Patrick Gebhard, and Aske Plaat. 2023. Fine-grained affective processing capabilities emerging from large language models. In *2023 11th international conference on affective computing and intelligent interaction (ACII)*, pages 1–8. IEEE.
- Hatice Gunes and Björn Schuller. 2013. Categorical and dimensional affect analysis in continuous input: Current trends and future directions. *Image and Vision Computing*, 31(2):120–136.
- Jacob Kahn, Morgane Riviere, Weiyi Zheng, Evgeny Kharitonov, Qiantong Xu, Pierre-Emmanuel Mazaré, Julien Karadayi, Vitaliy Liptchinsky, Ronan Collobert, Christian Fuegen, and 1 others. 2020. Libri-light: A benchmark for asr with limited or no supervision. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7669–7673. IEEE.
- Wei Kang, Xiaoyu Yang, Zengwei Yao, Fangjun Kuang, Yifan Yang, Liyong Guo, Long Lin, and Daniel Povey. 2024. Libriheavy: A 50,000 hours asr corpus with punctuation casing and context. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 10991–10995. IEEE.
- Sang-Hoon Lee, Ha-Yeong Choi, Seung-Bin Kim, and Seong-Whan Lee. 2025. Hierspeech++: Bridging the gap between semantic and acoustic representation of speech by hierarchical variational inference for zero-shot speech synthesis. *IEEE Transactions on Neural Networks and Learning Systems*.
- Shanglin Lei, Guanting Dong, Xiaoping Wang, Keheng Wang, Runqi Qiao, and Sirui Wang. 2023. Instructerc: Reforming emotion recognition in conversation with multi-task retrieval-augmented large language models. *arXiv preprint arXiv:2309.11911*.
- Yinghao Aaron Li, Cong Han, Vinay Raghavan, Gavin Mischler, and Nima Mesgarani. 2023. Styletts 2: Towards human-level text-to-speech through style diffusion and adversarial training with large speech language models. *Advances in Neural Information Processing Systems*, 36:19594–19621.
- Reza Lotfian and Carlos Busso. 2017. Building naturalistic emotionally balanced speech corpus by retrieving emotional speech from existing podcast recordings. *IEEE Transactions on Affective Computing*, 10(4):471–483.
- Ruskin Raj Manku, Yuzhi Tang, Xingjian Shi, Mu Li, and Alex Smola. 2025. Emergenttts-eval: Evaluating tts models on complex prosodic, expressiveness, and linguistic challenges using model-as-a-judge. *arXiv preprint arXiv:2505.23009*.
- GH Mohmad and Radhakrishnan Delhibabu. 2024. Speech databases, speech features and classifiers in speech emotion recognition: A review. *IEEE Access*.
- Zhichao Peng, Jianwu Dang, Masashi Unoki, and Masato Akagi. 2021. Multi-resolution modulation-filtered cochleagram feature for lstm-based dimensional emotion recognition from speech. *Neural Networks*, 140:261–273.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR.

- Anthony Rousseau, Paul Deléglise, Yannick Esteve, and 1 others. 2014. Enhancing the ted-lium corpus with selected data for language modeling and more ted talks. In *LREC*, pages 3935–3939.
- James A Russell and Albert Mehrabian. 1977. Evidence for a three-factor theory of emotions. *Journal of research in Personality*, 11(3):273–294.
- Shamane Siriwardhana, Andrew Reis, Rivindu Weerasekera, and Suranga Nanayakkara. 2020. Jointly fine-tuning "bert-like" self supervised models to improve multimodal speech emotion recognition. *arXiv preprint arXiv:2008.06682*.
- Sundararajan Srinivasan, Zhaocheng Huang, and Katrin Kirchhoff. 2022. Representation learning through cross-modal conditional teacher-student training for speech emotion recognition. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6442–6446. IEEE.
- Andreas Triantafyllopoulos, Uwe Reichel, Shuo Liu, Stephan Huber, Florian Eyben, and Björn W Schuller. 2023. Multistage linguistic conditioning of convolutional layers for speech emotion recognition. *Frontiers in Computer Science*, 5:1072479.
- 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.
- Xuena Wang, Xueting Li, Zi Yin, Yue Wu, and Jia Liu. 2023. Emotional intelligence of large language models. *Journal of Pacific Rim Psychology*, 17:18344909231213958.
- Zehui Wu, Ziwei Gong, Lin Ai, Pengyuan Shi, Kaan Donbekci, and Julia Hirschberg. 2025. Beyond silent letters: Amplifying llms in emotion recognition with vocal nuances. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 2202–2218.

## A LLM Prompts

Here is the prompt structure used in extracting emotions from a given text using LLMs.:

```
[{
  "role": "system",
  "content": "Valence, Arousal and Dominance are three affective dimensions that you can use to identify the sentiment in sentences. Assume that these dimensions can take values between 0 and 1, with 0 being low, and 1 being high. Remember that dominance assesses the extent to which the main person in the situation experiences the amount of control it can assert over the situation. Assess according to these dimensions the sentiment in the sentences I will give you after. Be precise, and output the values (up until two digits after the decimal point) in a JSON format please. Use 'Valence', 'Arousal' and 'Dominance' as keys in JSON format. Just acknowledge you got it."
},
{
  "role": "assistant",
  "content": "I'll assess the sentiment in the given sentence according to the Valence, Arousal, and Dominance dimensions and output the values in a JSON format. Please go ahead and provide the sentence. I'll output the values in a JSON format once all the assessments are complete."
},
{
  "role": "user",
  "content": "[SENTENCE]"
}]
```

Prompt structure for contextual emotion extraction:

```
[{
  "role": "system",
  "content": "Valence, Arousal and Dominance are three affective dimensions that you can use to identify the sentiment in sentences. Assume that these dimensions can take values between 0 and 1, with 0 being low, and 1 being high. Remember that dominance assesses the extent to which the main person in the situation experiences the amount of control it can assert over the situation. You will be provided with three parts of text; 1. Pre-text: The text that comes before the main sentence. 2. Main sentence: The sentence to be evaluated. 3. Post-text: The text that comes after the main sentence. Consider the entire context (pre-text, main sentence, and post-text) when evaluating the emotional dimensions, but focus your assessment primarily on the main sentence. The surrounding context should inform your understanding and interpretation of the main sentence's emotional content. Assess according to these dimensions the sentiment in the main sentence. Be precise, and output the values (up until two digits after the decimal point) in a JSON format please. Use 'Valence', 'Arousal' and 'Dominance' as keys in JSON format. Just acknowledge you got it."
},
{
  "role": "assistant",
  "content": "I'll assess the sentiment in the given sentence according to the Valence, Arousal, and Dominance dimensions and output the values in a JSON format. I will focus on assessing main sentence while using surrounding text as contextual cue for my assessment. Please go ahead and provide the previous text, main sentence and post text. I'll output the values in a JSON format once all the assessments are complete."
},
{
  "role": "user",
  "content": "Pre-Text: [PRE-TEXT]"; Main Sentence: [SENTENCE]; Post-Text: [POST-TEXT]"
}]
```

## B Details on data processing

To calculate correlations, we sample all except the ANET datasets and force the sample to a) cover the VAD space as best as possible and b) keep the number of segments close to ANET dataset. We use the speech-based model to predict emotion values from audiobooks and TED talks and then sampled based on predicted speech emotions. We sampled a subset of the MSP Podcast dataset based on ground truth emotion values. For the sample selection process we aim to cover emotional space as equally as possible to not bias the statistical correlation analyses. We divide the 0-1 emotion value range into 5 bins. For each bin we randomly select maximum 3 samples. If the bin does not have enough available samples to select we select the remaining as whole. We reduce the sampling criteria for MSP-Podcast as maximum

2 samples to keep sample sizes closer to each other across datasets. MSP Podcast, audiobooks, and TEDLIUM datasets yield 162, 115, 148 samples respectively.

Libriheavy and TEDLIUM datasets lack human annotated emotion labels. After validating the performance of SOTA speech-only Wav2Vec2 model, we use it as proxy and its prediction for emotion annotations. For a fair comparison we also use its prediction of MSP Podcast dataset in our analysis.

### C Coverage of Emotions in Audio Datasets

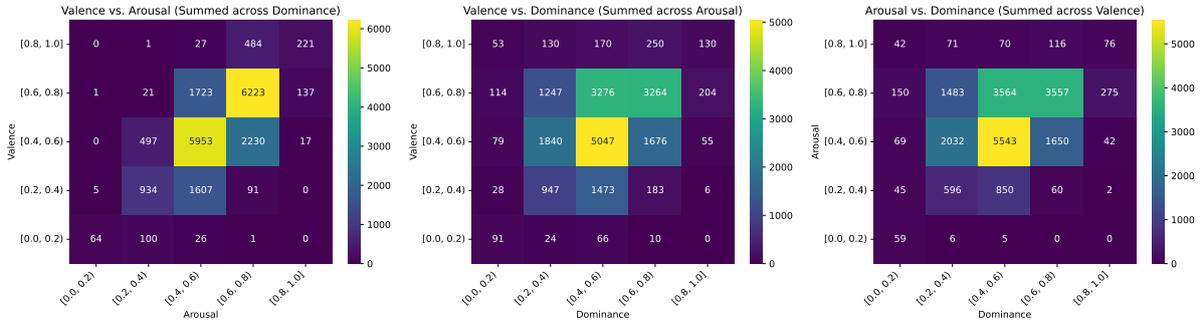


Figure 1: MSP Podcast dataset emotion coverage (Ground Truth Human Annotations).

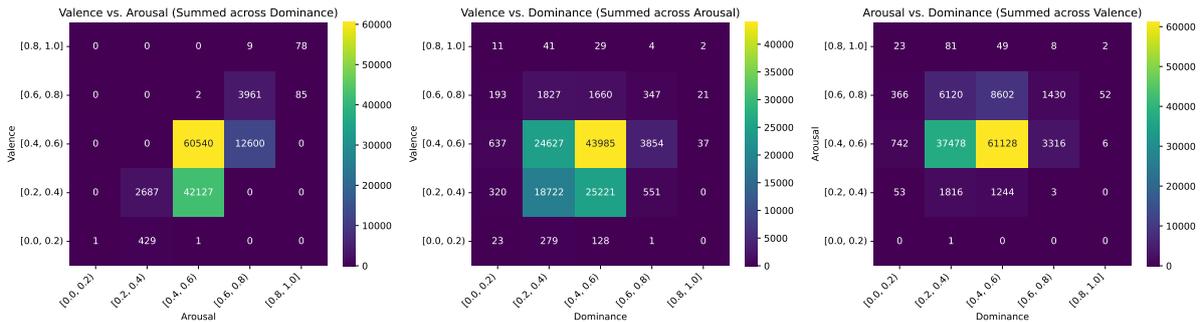


Figure 2: Libriheavy dataset emotion coverage (Audio Model Predictions).

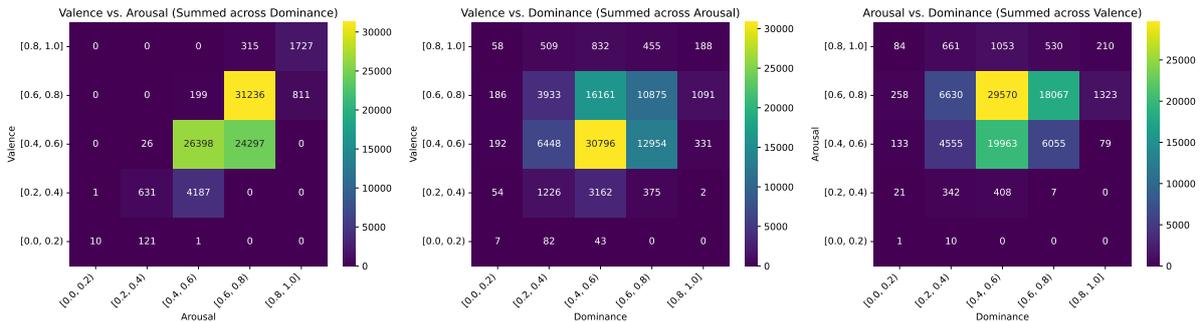


Figure 3: TEDLIUM dataset emotion coverage (Audio Model Predictions).

## D Results on the MSP Podcast dataset

Table 4: Correlation between text-based emotion prediction and speech-based emotion on MSP Podcast

Model	Wav2Vec2			WavLM			Human Annotations		
	A	D	V	A	D	V	A	D	V
llama3.3:70b	<b>.384</b> <sub>(.012)</sub>	.083 <sub>(.015)</sub>	<b>.618</b> <sub>(.014)</sub>	<b>.351</b> <sub>(.017)</sub>	.007 <sub>(.013)</sub>	<b>.589</b> <sub>(.014)</sub>	<b>.298</b> <sub>(.010)</sub>	.041 <sub>(.011)</sub>	<b>.604</b> <sub>(.002)</sub>
llama3.1:70b	<b>.344</b> <sub>(.026)</sub>	.092 <sub>(.042)</sub>	<b>.630</b> <sub>(.022)</sub>	<b>.313</b> <sub>(.030)</sub>	.018 <sub>(.051)</sub>	<b>.610</b> <sub>(.021)</sub>	<b>.265</b> <sub>(.035)</sub>	.088 <sub>(.040)</sub>	<b>.584</b> <sub>(.037)</sub>
phi3:14b	<b>.295</b> <sub>(.032)</sub>	-.041 <sub>(.049)</sub>	<b>.506</b> <sub>(.088)</sub>	<b>.269</b> <sub>(.023)</sub>	-.075 <sub>(.035)</sub>	<b>.471</b> <sub>(.091)</sub>	<b>.223</b> <sub>(.036)</sub>	-.011 <sub>(.037)</sub>	<b>.478</b> <sub>(.057)</sub>
gemma2:9b	<b>.381</b> <sub>(.027)</sub>	.117 <sub>(.028)</sub>	<b>.483</b> <sub>(.024)</sub>	<b>.360</b> <sub>(.021)</sub>	.092 <sub>(.025)</sub>	<b>.460</b> <sub>(.028)</sub>	<b>.281</b> <sub>(.034)</sub>	.074 <sub>(.023)</sub>	<b>.467</b> <sub>(.022)</sub>
llama3.1:8b	<b>.323</b> <sub>(.050)</sub>	.075 <sub>(.069)</sub>	<b>.545</b> <sub>(.026)</sub>	<b>.298</b> <sub>(.059)</sub>	.066 <sub>(.065)</sub>	<b>.514</b> <sub>(.016)</sub>	<b>.252</b> <sub>(.038)</sub>	.018 <sub>(.035)</sub>	<b>.545</b> <sub>(.025)</sub>
qwen2:7b	.210 <sub>(.079)</sub>	.097 <sub>(.060)</sub>	<b>.589</b> <sub>(.027)</sub>	.156 <sub>(.074)</sub>	.063 <sub>(.076)</sub>	<b>.541</b> <sub>(.023)</sub>	.161 <sub>(.063)</sub>	.025 <sub>(.058)</sub>	<b>.509</b> <sub>(.029)</sub>
openchat:7b	<b>.296</b> <sub>(.066)</sub>	.014 <sub>(.041)</sub>	<b>.511</b> <sub>(.090)</sub>	<b>.252</b> <sub>(.069)</sub>	-.032 <sub>(.052)</sub>	<b>.479</b> <sub>(.100)</sub>	.197 <sub>(.075)</sub>	.037 <sub>(.038)</sub>	<b>.495</b> <sub>(.107)</sub>
gemma:7b	<b>.358</b> <sub>(.060)</sub>	-.039 <sub>(.069)</sub>	<b>.417</b> <sub>(.017)</sub>	<b>.369</b> <sub>(.062)</sub>	-.061 <sub>(.056)</sub>	<b>.402</b> <sub>(.009)</sub>	<b>.306</b> <sub>(.031)</sub>	-.036 <sub>(.066)</sub>	<b>.419</b> <sub>(.034)</sub>

Note: Values show mean correlation (standard deviation) across 5 seeds.  
**Bold values** indicate statistical significance ( $p < 0.01$ ) in  $\geq 80\%$  of runs.

## E Results on Libriheavy Dataset

Table 5: Correlation between text-based and speech-based emotion prediction on Libriheavy dataset

Model	Wav2Vec2			WavLM		
	A	D	V	A	D	V
llama3.3:70b	<b>.328</b> <sub>(.011)</sub>	.172 <sub>(.021)</sub>	<b>.783</b> <sub>(.010)</sub>	<b>.399</b> <sub>(.010)</sub>	.094 <sub>(.022)</sub>	<b>.664</b> <sub>(.011)</sub>
llama3.1:70b	<b>.374</b> <sub>(.039)</sub>	.072 <sub>(.044)</sub>	<b>.793</b> <sub>(.018)</sub>	<b>.422</b> <sub>(.031)</sub>	.005 <sub>(.040)</sub>	<b>.701</b> <sub>(.024)</sub>
phi3:14b	<b>.350</b> <sub>(.048)</sub>	.104 <sub>(.063)</sub>	<b>.671</b> <sub>(.091)</sub>	<b>.353</b> <sub>(.041)</sub>	.040 <sub>(.063)</sub>	<b>.615</b> <sub>(.093)</sub>
gemma2:9b	<b>.437</b> <sub>(.010)</sub>	.198 <sub>(.030)</sub>	<b>.713</b> <sub>(.013)</sub>	<b>.494</b> <sub>(.011)</sub>	.147 <sub>(.028)</sub>	<b>.636</b> <sub>(.027)</sub>
llama3.1:8b	<b>.418</b> <sub>(.050)</sub>	.201 <sub>(.108)</sub>	<b>.704</b> <sub>(.019)</sub>	<b>.484</b> <sub>(.042)</sub>	.180 <sub>(.093)</sub>	<b>.628</b> <sub>(.029)</sub>
qwen2:7b	<b>.345</b> <sub>(.023)</sub>	.195 <sub>(.051)</sub>	<b>.728</b> <sub>(.019)</sub>	<b>.303</b> <sub>(.021)</sub>	.102 <sub>(.039)</sub>	<b>.626</b> <sub>(.032)</sub>
openchat:7b	<b>.332</b> <sub>(.065)</sub>	.057 <sub>(.083)</sub>	<b>.720</b> <sub>(.101)</sub>	<b>.331</b> <sub>(.071)</sub>	-.032 <sub>(.088)</sub>	<b>.634</b> <sub>(.102)</sub>
gemma:7b	.250 <sub>(.054)</sub>	.031 <sub>(.075)</sub>	<b>.619</b> <sub>(.007)</sub>	<b>.254</b> <sub>(.038)</sub>	-.003 <sub>(.058)</sub>	<b>.577</b> <sub>(.006)</sub>

Note: Values show mean correlation (standard deviation) across 5 seeds.  
**Bold values** indicate statistical significance ( $p < 0.01$ ) in  $\geq 80\%$  of runs.

Table 6: Correlation between *contextual* text-based emotion prediction and speech-based emotion prediction on Libriheavy

Model	Wav2Vec2			WavLM		
	A	D	V	A	D	V
llama3.3:70b	<b>.416</b> (.016)	.175 (.054)	<b>.700</b> (.006)	<b>.468</b> (.018)	.060 (.047)	<b>.618</b> (.007)
llama3.1:70b	<b>.427</b> (.021)	.015 (.025)	<b>.697</b> (.016)	<b>.484</b> (.017)	-.070 (.029)	<b>.626</b> (.023)
phi3:14b	<b>.392</b> (.080)	.108 (.064)	<b>.591</b> (.066)	<b>.392</b> (.076)	.023 (.060)	<b>.519</b> (.064)
gemma2:9b	<b>.450</b> (.015)	<b>.297</b> (.026)	<b>.641</b> (.023)	<b>.511</b> (.015)	<b>.266</b> (.027)	<b>.599</b> (.021)
llama3.1:8b	<b>.473</b> (.064)	.119 (.117)	<b>.642</b> (.040)	<b>.502</b> (.056)	.067 (.100)	<b>.576</b> (.027)
qwen2:7b	<b>.267</b> (.067)	.064 (.044)	<b>.679</b> (.016)	.251 (.053)	-.004 (.039)	<b>.563</b> (.014)
openchat:7b	<b>.320</b> (.099)	.006 (.037)	<b>.591</b> (.089)	<b>.303</b> (.088)	-.076 (.031)	<b>.496</b> (.097)
gemma:7b	.230 (.052)	.144 (.100)	<b>.549</b> (.053)	.195 (.050)	.102 (.090)	<b>.458</b> (.063)

*Note: Values show mean correlation (standard deviation) across 5 seeds.  
**Bold values** indicate statistical significance ( $p < 0.01$ ) in  $\geq 80\%$  of runs.*

## F Results on TEDLIUM Dataset

Table 7: Correlation between text-based emotion prediction and speech-based emotion prediction on TEDLIUM dataset

Model	Wav2Vec2			WavLM		
	A	D	V	A	D	V
llama3.3:70b	<b>.281</b> <sub>(.005)</sub>	<b>.274</b> <sub>(.035)</sub>	<b>.685</b> <sub>(.013)</sub>	<b>.306</b> <sub>(.008)</sub>	.190 <sub>(.029)</sub>	<b>.698</b> <sub>(.016)</sub>
llama3.1:70b	<b>.301</b> <sub>(.024)</sub>	.218 <sub>(.020)</sub>	<b>.645</b> <sub>(.029)</sub>	<b>.317</b> <sub>(.037)</sub>	.128 <sub>(.025)</sub>	<b>.638</b> <sub>(.044)</sub>
phi3:14b	<b>.243</b> <sub>(.027)</sub>	.025 <sub>(.029)</sub>	<b>.599</b> <sub>(.090)</sub>	<b>.271</b> <sub>(.042)</sub>	-.028 <sub>(.048)</sub>	<b>.626</b> <sub>(.086)</sub>
gemma2:9b	<b>.263</b> <sub>(.018)</sub>	<b>.270</b> <sub>(.058)</sub>	<b>.560</b> <sub>(.020)</sub>	<b>.289</b> <sub>(.021)</sub>	<b>.250</b> <sub>(.059)</sub>	<b>.556</b> <sub>(.005)</sub>
llama3.1:8b	.179 <sub>(.047)</sub>	.138 <sub>(.056)</sub>	<b>.616</b> <sub>(.040)</sub>	.210 <sub>(.062)</sub>	.123 <sub>(.051)</sub>	<b>.612</b> <sub>(.029)</sub>
qwen2:7b	<b>.270</b> <sub>(.052)</sub>	<b>.292</b> <sub>(.057)</sub>	<b>.662</b> <sub>(.042)</sub>	<b>.246</b> <sub>(.055)</sub>	.247 <sub>(.052)</sub>	<b>.672</b> <sub>(.036)</sub>
openchat:7b	.203 <sub>(.118)</sub>	.128 <sub>(.077)</sub>	<b>.609</b> <sub>(.109)</sub>	.170 <sub>(.108)</sub>	.060 <sub>(.086)</sub>	<b>.615</b> <sub>(.103)</sub>
gemma:7b	<b>.281</b> <sub>(.022)</sub>	-.052 <sub>(.041)</sub>	<b>.472</b> <sub>(.009)</sub>	<b>.279</b> <sub>(.032)</sub>	-.110 <sub>(.044)</sub>	<b>.528</b> <sub>(.007)</sub>

Note: Values show mean correlation (standard deviation) across 5 seeds.  
**Bold values** indicate statistical significance ( $p < 0.01$ ) in  $\geq 80\%$  of runs.

Table 8: Correlation between *contextual* text-based emotion prediction and speech-based emotion prediction on TEDLIUM

Model	Wav2Vec2			WavLM		
	A	D	V	A	D	V
llama3.3:70b	.210 <sub>(.011)</sub>	.214 <sub>(.028)</sub>	<b>.627</b> <sub>(.007)</sub>	<b>.241</b> <sub>(.023)</sub>	.125 <sub>(.024)</sub>	<b>.654</b> <sub>(.005)</sub>
llama3.1:70b	.223 <sub>(.030)</sub>	.174 <sub>(.022)</sub>	<b>.624</b> <sub>(.017)</sub>	<b>.233</b> <sub>(.027)</sub>	.078 <sub>(.027)</sub>	<b>.664</b> <sub>(.016)</sub>
phi3:14b	.244 <sub>(.058)</sub>	.110 <sub>(.072)</sub>	<b>.538</b> <sub>(.077)</sub>	.244 <sub>(.055)</sub>	.045 <sub>(.068)</sub>	<b>.561</b> <sub>(.081)</sub>
gemma2:9b	.135 <sub>(.029)</sub>	.116 <sub>(.058)</sub>	<b>.586</b> <sub>(.021)</sub>	.164 <sub>(.029)</sub>	.046 <sub>(.063)</sub>	<b>.578</b> <sub>(.024)</sub>
llama3.1:8b	.208 <sub>(.041)</sub>	.217 <sub>(.037)</sub>	<b>.590</b> <sub>(.019)</sub>	<b>.242</b> <sub>(.038)</sub>	.220 <sub>(.043)</sub>	<b>.629</b> <sub>(.037)</sub>
qwen2:7b	<b>.340</b> <sub>(.041)</sub>	<b>.292</b> <sub>(.057)</sub>	<b>.647</b> <sub>(.011)</sub>	<b>.263</b> <sub>(.054)</sub>	.216 <sub>(.045)</sub>	<b>.653</b> <sub>(.011)</sub>
openchat:7b	<b>.237</b> <sub>(.058)</sub>	.181 <sub>(.054)</sub>	<b>.546</b> <sub>(.106)</sub>	.214 <sub>(.043)</sub>	.104 <sub>(.062)</sub>	<b>.569</b> <sub>(.107)</sub>
gemma:7b	.155 <sub>(.043)</sub>	.048 <sub>(.063)</sub>	<b>.504</b> <sub>(.012)</sub>	.151 <sub>(.033)</sub>	.007 <sub>(.061)</sub>	<b>.508</b> <sub>(.009)</sub>

Note: Values show mean correlation (standard deviation) across 5 seeds.  
**Bold values** indicate statistical significance ( $p < 0.01$ ) in  $\geq 80\%$  of runs.

## G Audio Model Comparison

Table 9: Correlation between emotion predictions of the two audio models Wav2Vec2 and WavLM across datasets

	Arousal	Dominance	Valence
MSP Podcast	<b>.928</b>	<b>.915</b>	<b>.823</b>
Libriheavy	<b>.934</b>	<b>.919</b>	<b>.869</b>
TEDLIUM	<b>.935</b>	<b>.931</b>	<b>.868</b>

Table 10: Correlation between audio models on MSP Podcast dataset human annotations

	Arousal	Dominance	Valence
Wav2Vec2	<b>.772</b>	<b>.640</b>	<b>.646</b>
WavLM	<b>.765</b>	<b>.677</b>	<b>.577</b>