

# The Prosody of Emojis

Anonymous ACL submission

## Abstract

Prosodic features such as pitch, timing, and intonation are central to spoken communication, conveying emotion, intent, and discourse structure. In text-based settings, where these cues are absent, emojis act as visual surrogates that add affective and pragmatic nuance. This study examines how emojis influence prosodic realisation in speech and how listeners interpret prosodic cues to recover emoji meanings. Unlike previous work, we directly link prosody and emojis by analysing human speech data collected through a controlled elicited production task<sup>1</sup>. Using Bayesian multilevel modelling, we show that speakers systematically adapt their prosody based on emoji cues, and that listeners can recover intended meanings significantly above chance. Furthermore, our results reveal a clear hierarchy in prosodic shifts: greater semantic differences between emojis correspond to increased prosodic divergence. These findings suggest that emojis are meaningful carriers of prosodic intent that bridge the gap between digital text and spoken production.

## 1 Introduction

In spoken language, prosody, encompassing pitch, rhythm, and intonation, and other paralinguistic features, plays a crucial role in conveying linguistic nuances and socio-emotional cues (Cole, 2015; Hellbernd and Sammler, 2016; Ward and Levow, 2021). Acting as a dynamic carrier of information, prosody enriches spoken communication by modulating affect, signalling discourse structure, and shaping listener expectations. However, the shift toward computer-mediated communication presents challenges: text-based interaction strips away vital non-verbal elements like intonation, gesture, and facial expression, essential to holistic interpretation (Archer and Akert, 1977; Crystal, 2001).

To compensate for this loss, users of digital communication platforms increasingly rely on emojis,

<sup>1</sup>The data will be publicly released for research purposes.

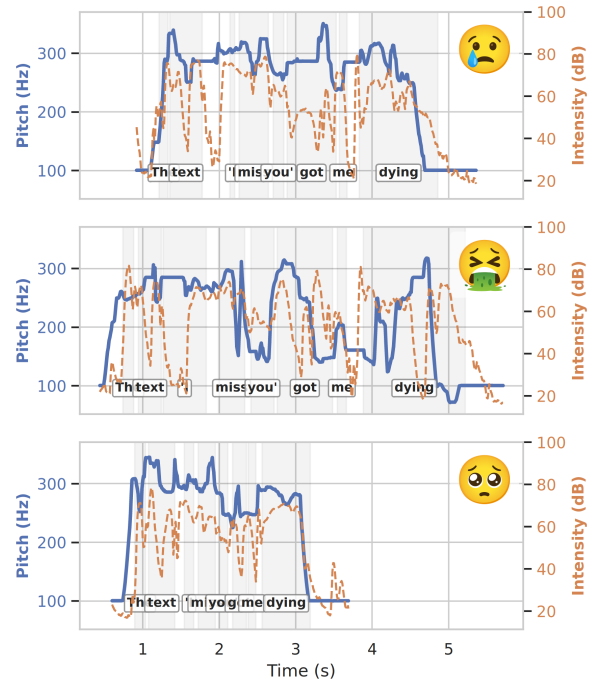


Figure 1: Prosodic variations for the sentence “The text ‘I miss you’ got me dying” across three emoji contexts. Each subplot displays the pitch (solid blue, Hz) and intensity (dashed orange, dB) contours.

which enrich digital communication by offering a visual medium for expressing emotions (Gülşen, 2016) and conveying semantic meanings (Na’aman et al., 2017), transcending textual constraints. Their interpretive flexibility, however, raises intriguing questions about their relationship to prosody. Prior research has shown that prosodic features cannot be fully recovered from text alone (Wolf et al., 2023). This leads to a compelling inquiry: can emojis, embedded in otherwise static written language, serve as mediators for prosodic intent, compensating for the absence of vocal cues? Specifically, we investigate the following research questions:

**RQ1:** Do emojis influence how prosodic features are expressed in spoken language?

057 **RQ2:** Can listeners interpret prosodic cues to  
058 recover intended emoji meanings in speech?

059 **RQ3:** How does inter-speaker similarity in  
060 prosodic realisation relate to listener interpretation  
061 of emoji meaning?

062 **RQ4:** Is there a correlation between the seman-  
063 tic distance of emojis and the magnitude of the  
064 prosodic shift in their vocal delivery?

065 To investigate these questions, we conducted  
066 four online experiments to collect a dataset of  
067 emoji-enriched speech and listener perception data  
068 (Figure 1). We employed a design that holds lexical  
069 content constant across emoji conditions to isolate  
070 the influence of emojis on prosody, thereby pre-  
071 venting confounding effects from varying syntax  
072 or vocabulary. While this involves reading specific  
073 sentences rather than recording spontaneous con-  
074 versations, we do not instruct participants on how  
075 to deliver the text or focus on the emojis. This  
076 approach captures a broad range of naturalistic ex-  
077 pressive strategies while ensuring that observed  
078 prosodic shifts reflect authentic interpretations of  
079 emoji meaning.

080 Our results demonstrate that emojis consistently  
081 invite prosodic modulation across diverse speakers,  
082 even with constant verbal content. This variation  
083 is often sufficient for listeners to recover intended  
084 emoji meanings significantly above chance. We  
085 find that prosodic convergence, where different  
086 speakers independently produce similar contours  
087 for the same emoji, significantly boosts listener suc-  
088 cess. This suggests that shared prosodic intuitions  
089 facilitate robust interpretation. Finally, our analysis  
090 reveals a hierarchy in prosodic shifts. The mag-  
091 nitude of prosodic change correlates with emoji  
092 semantic distance, as greater semantic dissimilar-  
093 ity between emojis leads to systematically larger  
094 prosodic divergence.

## 095 2 Background

### 096 2.1 Emoji Interpretation

097 Emojis constitute a pervasive form of digital par-  
098 alanguage: graphical cues that represent the non-  
099 verbal behaviours, such as facial expressions and  
100 gestures, intrinsic to face-to-face interaction (Al-  
101 shenqeei, 2016). From a psychological perspec-  
102 tive, emojis do not merely decorate text but ac-  
103 tively modulate affective processing. They have  
104 been shown to trigger cognitive responses in the  
105 receiver similar to those elicited by physical facial

expressions, thereby influencing social attributions  
of the sender, such as perceived warmth, trust, and  
sincerity, and promoting a sense of social presence  
(Boutet et al., 2021).

Beyond their social impact, emojis serve criti-  
cal pragmatic functions by acting as illocutionary  
force indicators. They allow users to clarify intent  
and manage face-threatening acts, such as soften-  
ing a critique or signalling irony, which are tasks  
typically reserved for vocal nuance (Li and Yang,  
2018; Holtgraves and Robinson, 2020). However,  
this interpretive process is fraught with ambiguity.  
While literal meanings are relatively stable across  
languages, figurative is closely intertwined with  
the sentiment of the context (Zhou et al., 2024a).  
Despite the presence of contextual cues intended to  
clarify intent, empirical evidence shows that misin-  
terpretation persists, suggesting that emojis remain  
a high-variance channel that context alone cannot  
resolve (Miller et al., 2017; Czkestopowska et al.,  
2022).

### 127 2.2 Prosody in Speech Processing

128 Prosody encompasses the suprasegmental features  
129 of speech that operate beyond individual phonemes,  
130 such as intonation, rhythm, pitch, and timing  
131 that structure discourse and convey speaker intent  
(Lehiste and Lass, 1976; Gerken and McGregor,  
1998). While speech processing systems often pri-  
132 oritise lexical cues, especially in cascade pipelines  
133 where prosodic nuance is lost during transcrip-  
134 tion (Zhou et al., 2024b; Tsiamas et al., 2024),  
135 research shows that these cues can independently  
136 support tasks like spoken question-answering, emo-  
137 tion recognition, and speaker profiling (Hirschberg,  
138 2002; Chi et al., 2025).

141 This importance extends to speech synthesis,  
142 where naturalness depends on accurate prosodic  
143 modelling. Recent efforts have focused on this by  
144 implementing expressive prosody transfer (Skerry-  
145 Ryan et al., 2018) and modelling human-like  
146 rhythm to improve listener engagement in Text-  
147 to-Speech (TTS) systems (Kane et al., 2024).

### 148 2.3 Prosody and Emojis

149 Given that emojis function as pragmatic substitutes  
150 for non-verbal cues, they present a unique oppor-  
151 tunity to recover the prosodic intent lost in text.  
152 Prior research suggests conceptual parallels, liken-  
153 ing emojis to written intonation or visual facial ex-  
154 pressions (James, 2017) that highlight information  
155 focus or compensate for a lack of vocal expressivity

(Wagner, 2016; Hu et al., 2019; Kaiser, 2021).

Despite these theoretical links, empirical analysis of how emojis influence the acoustic realisation of speech is limited. While some work has proposed emojis as graphical surrogates for prosody based on surveys (Alnuzaili et al., 2024), or used emoji-enriched speech to increase the variance of TTS systems (Tuttösí et al., 2025), research has yet to characterise the systematic prosodic adaptations human speakers make in response to different emoji semantics. In contrast, our study examines how speakers naturally modulate their acoustic prosody when encountering emoji-enriched utterances. By analysing real human speech, we provide empirical evidence that emoji semantics directly influence prosodic expression, bridging the gap between symbolic digital cues and spoken production.

### 3 Methods

#### 3.1 Data Collection

To investigate our research questions, we conducted four online experiments<sup>2</sup>, with each experiment employing a unique pool of participants. To isolate the influence of emojis on prosody, we employed a design that holds lexical content constant across conditions. While this involves elicited speech rather than spontaneous conversation, this approach is necessary to generate the parallel data required to compare how the same lexical string is transformed by different emoji contexts, a phenomenon rarely captured in naturalistic corpora.

The experimental framework is built upon a base corpus of 230 unique utterances, each of which originally contained at least one emoji. These were manually selected and edited from publicly available user-generated textual content to cover a broad range of topics (e.g., music, gaming, film) and diverse emoji usage. All selected utterances were specifically screened to ensure they remained natural when spoken aloud.

#### Experiment 1: Utterance-Emoji Stimulus

The objective of Experiment 1 was to construct a collection of utterances with identical lexical content but varying emoji usage. Participants were presented with the base utterances (with original emojis removed) and instructed to create up to five variations per sentence by inserting different emo-

<sup>2</sup>Appendix G for full instructions and Appendix B for data summaries.

jis to alter the emotional tone or interpretation. To ensure variety, they were advised against using emojis with highly similar meanings for the same utterance. Each utterance was annotated independently by two annotators, resulting in 1,595 unique emoji-augmented stimuli.

#### Experiment 2: Emoji-Influenced Prosodic Data

We collected a corpus of emoji-influenced speech by asking participants to read a subset of the stimuli annotated in Experiment 1. Each speaker recorded approximately 25 utterances, ensuring they produced at least two distinct variations of each base sentence (including a non-emoji control version). To minimise bias in prosodic expression, participants were not informed of the focus on emojis. Following the recordings, participants summarised the perceived meaning of each emoji in one word to verify their engagement and interpretation. The final dataset comprised 3,153 valid recordings from 129 speakers. To account for potential sources of variance in our analysis, we also collected demographic data, including age, gender, and self-reported general emoji usage (measured on a scale of 1–5). Additionally, participants reported their recording hardware (headset, mobile device, or built-in computer microphone).

#### Experiment 3: Judging Prosodic Variation Within Speakers

To test listeners’ perception of prosodic changes within the same speaker, we sampled 645 recording pairs which had identical lexical content. These pairs represented two distinct contrast types: *Dual Emoji*, comparing two emoji renditions (e.g., *Text* + 😞 vs. *Text* + 😊) and *Single Emoji*, comparing an emoji rendition against an emoji-free baseline. For each pair, listeners: judged whether the recordings differed in prosodic realisation (RQ1), and matched each recording to its corresponding emoji (RQ2). Listeners were instructed to focus on vocal expressiveness and ignore recording artefacts. Each pair received judgments from an average of three annotators. Consistent with Experiment 2, we collected demographic data from these listeners, including age, gender, and general emoji usage (scale 1-5).

#### Experiment 4: Identifying Shared Prosody Across Speakers

Experiment 4 investigated whether inter-speaker similarity in prosodic realisation, or prosodic convergence, facilitates the interpretation of emoji

Predictor	Est. ( $\beta$ )	95% CI
(Intercept)	-0.25	[-0.66, 0.15]
Age	-0.04	[-0.16, 0.07]
Gender (Female)	0.16	[-0.12, 0.44]
Mic (Headset)	0.10	[-0.26, 0.46]
Mic (Mobile)	-0.15	[-0.42, 0.12]
Emoji Use	<b>0.18</b>	<b>[0.06, 0.30]</b>
Contrast Type (Dual)	<b>0.29</b>	<b>[0.03, 0.55]</b>

Table 1: Model 1 results for perceived prosodic differentiation. Bold indicates 95% Credible Intervals that strictly exclude zero.

meaning (RQ3). To ensure sufficient expressive range, we excluded speakers who demonstrated low prosodic variation in Experiment 3 (see Section 4.2). The final stimulus set comprised 49 inter-speaker audio pairs, each consisting of recordings of the same sentence and emoji produced by two different speakers. Participants performed two tasks: judging the prosodic similarity of the recordings, and matching each recording to either the target emoji or a random distractor. Crucially, while both recordings in a pair shared the same target emoji, participants were not informed of this fact. Each pair received judgments from an average of four annotators. As in Experiments 2 and 3, we collected demographic data from these listeners, including age, gender, and general emoji usage.

### 3.2 Statistical Analysis

To evaluate our research questions, we implemented Bayesian linear mixed-effects regression models (Appendix A). This framework allows us to directly quantify the strength of evidence for our hypotheses by estimating the full probability distribution of our parameters. We modelled the influence of emoji cues on both vocal production and listener perception, including the accuracy of recovering intended meanings. By including random effects for speakers, listeners, and specific utterances, we ensured that our results account for individual expressive variability and the inherent prosodic differences between sentences.

## 4 Prosody in Emoji-Enriched Speech

### 4.1 Speaker Variation in Prosodic Expression

Our first inquiry examines whether emojis elicit prosodic variation (RQ1). We analysed data from Experiment 3, where listeners compared pairs of

recordings produced by the same speaker. To quantify this effect, we modelled listener judgments using a Bayesian multilevel logistic regression:

$$\text{Model 1: PerceivedVariation} \sim \text{Age} + \text{Gender} + \text{Mic} + \text{EmojiUse} + \text{ContrastType} + (1 \mid \text{Utterance}) + (1 \mid \text{Listener})$$

In this model, PerceivedVariation is the binary response variable (whether the two recordings sounded different or not), while ContrastType (Single vs. Dual Emoji comparisons), self-reported EmojiUse, and speaker demographics serve as the fixed effects.

The results, detailed in Table 1, provide a nuanced answer to RQ1. The influence of emojis appears to be a capacity that varies significantly between individuals. While most speakers exhibit moderate variation, there is significant spread, with clear tails of highly expressive and highly monotonous speakers (see Appendix B.1). Model 1 helps explain this variance through EmojiUse ( $\beta = 0.18, CI[0.06, 0.30]$ ). The positive effect suggests that for frequent emoji users, the parallel between emoji semantics and prosody is stronger, leading to more distinct acoustic signals. Furthermore, ContrastType systematically influenced detectability. The positive coefficient for Dual Emoji trials ( $\beta = 0.29, CI[0.03, 0.55]$ ) indicates that listeners were significantly more likely to detect a difference when comparing two distinct emojis than when comparing an emoji utterance against a neutral baseline. This implies that the prosodic modifications speakers make are not binary. Instead, the signal is graded, becoming most acoustically distinct when speakers must actively differentiate between conflicting states. These findings verify the existence of the parallel proposed in RQ1: emojis can and do influence prosodic expression. However, this influence is not uniform. Instead, emojis function as an expressive resource that is modulated by the speaker’s digital fluency and the communicative necessity of the contrast.

### 4.2 Listener Interpretation of Prosodic Variation

Beyond mere production, we investigated whether these variations effectively communicate symbolic intent and allow listeners to recover intended emoji meanings (RQ2). We tested this using a second model:

$$\text{Model 2: IdSuccess} \sim \text{SpeakerExpressivity} + \text{ContrastType} + (\text{SpeakerEmojiUse} \times$$

Predictor	Est. ( $\beta$ )	95% CI
<i>Model 2: Intra-speaker Identification (RQ2)</i>		
(Intercept)	-0.16	[-0.45, 0.11]
Speaker Expressivity	<b>0.24</b>	<b>[0.10, 0.38]</b>
Contrast Type (Dual)	0.03	[-0.27, 0.32]
Speaker Emoji Use	0.05	[-0.09, 0.19]
Listener Emoji Use	-0.15	[-0.31, 0.01]
Speaker $\times$ Listener Use	0.09	[-0.05, 0.23]
<i>Model 3: Inter-speaker Identification (RQ3)</i>		
(Intercept)	-1.79	[-3.38, -0.17]
Perceived Similarity (Yes)	<b>2.42</b>	<b>[1.66, 3.22]</b>

Table 2: Model 2 (intra-speaker) and Model 3 (inter-speaker) results for emoji identification success. Bold indicates 95% Credible Intervals that strictly exclude zero.

$$\text{ListenerEmojiUse}) + (1 \mid \text{Utterance}) + (1 \mid \text{Listener})$$

In this model,  $\text{IdSuccess}$  denotes the correct assignment of emoji meanings to both recordings in the pair. Key predictors included  $\text{SpeakerExpressivity}$  (a proxy for expressive bandwidth derived from Experiment 3),  $\text{ContrastType}$  (Dual vs. Single), and the interaction between speaker and listener  $\text{EmojiUse}$  to test for shared digital fluency effects. Model 2 (Table 2) revealed a baseline identification probability of  $P(\text{IdSuccess}) \approx 0.46$  (Intercept =  $-0.16$ ). Although the intercept coefficient includes zero (indicating a probability near 50%), this value sits significantly above the 25% chance level<sup>3</sup>, confirming that the task was generally performable.

The primary driver of communicative success was  $\text{Speaker Expressivity}$ : for every standard deviation increase in a speaker’s prosodic differentiation, the log-odds of correct emoji identification increased reliably ( $\beta = 0.24$ ,  $CI[0.10, 0.38]$ ). Interestingly, while  $\text{Contrast Type}$  was a major factor in the perception of differences (Model 1), it did not reliably impact identification success in Model 2 ( $\beta = 0.03$ ,  $CI[-0.27, 0.32]$ ). This suggests that once a prosodic difference is detected, the intended meaning is equally recoverable whether comparing two different emojis or an emoji against a neutral baseline. Additionally, the lack of a reliable interaction between speaker and listener emoji familiarity indicates that this communication is mutually intelligible regardless of digital background, rely-

<sup>3</sup>Chance performance is 25% based on the four possible emoji assignment combinations available to the listener (Log-odds  $\approx -1.1$ ).

ing instead on universal expressive mechanisms. These findings address RQ2 by demonstrating that while listeners can interpret prosodic cues to recover emoji meanings significantly above chance, the success of this communication is not universal; rather, it is a high-variance channel that depends largely on the individual speaker’s ability to encode pragmatic intent into vocal variation.

### 4.3 Prosodic Convergence and Emoji Interpretation

RQ3 shifts focus to the inter-speaker dynamics explored in Experiment 4. Specifically, we investigated whether prosodic convergence, defined as instances where different speakers independently produce similar contours for the same emoji, facilitates higher listener accuracy. To test this, we modelled the relationship between the listener’s perception of similarity and their ability to correctly identify the emoji using a third Bayesian model:

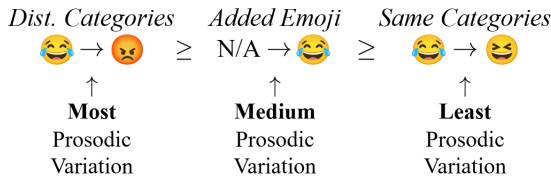
$$\text{Model 3: IdSuccess} \sim \text{PerceivedSimilarity} + (1 \mid \text{Utterance}) + (1 \mid \text{Listener}) + (1 \mid \text{Speaker}_A) + (1 \mid \text{Speaker}_B)$$

The fixed effect  $\text{PerceivedSimilarity}$  is a binary predictor indicating whether the listener judged the recordings from the two speakers as sounding prosodically similar. Random effects account for variance across utterances, listeners, and the two unique speakers involved in each pair ( $\text{Speaker}_A$  and  $\text{Speaker}_B$ ). The regression results, presented as Model 3 in Table 2, confirm that perceived similarity is a robust predictor of successful emoji recovery. The model found a significant positive effect for similarity ( $\beta = 2.42$ , 95% CI [1.66, 3.22]). This result suggests that listeners are approximately 11 times more likely to achieve discriminative success when speakers converge on a shared prosodic signature. Even after controlling for random variation, the evidence for this effect remained extreme with a posterior probability of 1.0. These findings provide a direct answer to RQ3 by demonstrating that inter-speaker similarity in prosodic realisation is fundamentally tied to listener interpretation success. When speakers independently converge on a shared prosodic strategy, listener recovery of intended meaning improves significantly. However, communicative success is notably more fragile when prosodic strategies diverge. In such cases, listeners often perceive the recordings as sounding different and prefer to assign distinct emoji mean-

ings to each, reflecting the increased uncertainty caused by a lack of prosodic agreement (further details in Appendix B.3). Thus, while emoji vocalisation reflects a stable and shared paralinguistic code that can ensure semantic intent remains recoverable across different voices, its effectiveness is largely contingent on whether speakers align their prosodic intuitions.

## 5 Prosodic Distance and Emoji Semantic Variations

This section investigates how emoji semantics influence prosodic variation. To address RQ4, we test the hypothesis that prosodic differences are greater when associated emojis belong to distinct semantic categories than when they share a category. Additionally, we examine whether the inclusion of an emoji in a previously emoji-free sentence induces greater prosodic change than switching between emojis within the same semantic group.



### 5.1 Emoji Semantic Categories

Accurately categorising emojis poses a challenge because standard groupings, such as Unicode emoji groups and subgroups<sup>4</sup>, do not always reflect the communicative or semantic similarities perceived by users (e.g., 😄 and 😊 belong to the “face-affection” subgroup). To address this, we derived emoji embeddings from participant intents (Experiment 1) and interpretations (Experiment 2). Despite individual variation, these embeddings showed strong correlations between intents and interpretations, indicating a substantial shared understanding of the communicative roles of emojis (see Appendix E for detailed analysis).

To identify meaningful semantic groups, we performed unsupervised clustering on the interpretation embeddings using HDBSCAN (Campello et al., 2013) with a minimum cluster size of 3 (see Appendix F for further discussion). Based on these clusters, we established three experimental conditions for our prosodic analysis: recording pairs where emojis belong to different semantic clusters

<sup>4</sup><https://unicode.org/Public/emoji/16.0/emoji-test.txt>

(Diff); recording pairs where emojis belong to the same semantic cluster. (Same); pairs consisting of an emoji-free recording matched with an emoji-enriched version of the same utterance (Added).

All relevant recording pairs were extracted and categorised accordingly, excluding speakers with minimal prosodic differentiation ( $\leq 30\%$ , as identified in Section 4.1).

### 5.2 Methodology

For each recording, we extracted multiple prosodic features to capture variation across several dimensions. Pitch contours were estimated using the FCNF0++ model (Morrison et al., 2023), and intensity was calculated via the Praat Parselmouth interface (Jadoul et al., 2018). To account for varying recording lengths, we employed Dynamic Time Warping (DTW) to align pitch and intensity contours. From these paths, we derived a Temporal Warping metric defined as the percentage of non-diagonal steps in the alignment to quantify rhythmic divergence and local temporal distortion. Additionally, we calculated the absolute difference in speech rate (characters per second) as a measure of global articulation speed via the Stopes toolkit (Seamless Communication et al., 2023).

Beyond raw acoustic features, we extracted prosodic embeddings from the 7th layer of the Masked Prosody Model (Wallbridge et al., 2025). These frame-wise, 256-dimensional representations are explicitly trained to model pitch, energy, and voice activity. Distances between these sequences were computed using multivariate DTW to capture holistic prosodic divergence across these latent dimensions.

To determine if the magnitude of the prosodic shift varies systematically according to emoji semantic relationships (RQ4), we employed Bayesian linear mixed-effects regression. We fitted separate models for each acoustic feature using the following specification:

$$\text{Model 4: ProsodicDistance} \sim \text{SemanticRelation} + (1 + \text{SemanticRelation} | \text{Speaker}) + (1 | \text{Utterance})$$

The dependent variable, ProsodicDistance, represents the calculated acoustic difference between the recordings. To satisfy normality assumptions and facilitate direct comparison across features, all distance metrics were log-transformed and Z-standardised prior to fitting. The fixed effect, SemanticRelation, categorises the pairing

Hypothesis	Est. ( $\beta$ )	95% CI	ER
<i>Embeddings</i>			
Distinct > Added	<b>0.10</b>	<b>[0.05, 0.14]</b>	<u><math>\geq 100</math></u>
Same < Added	0.05	[-0.03, 0.13]	0.2
Distinct > Same	0.05	[-0.02, 0.12]	<u>6.2</u>
<i>Intensity</i>			
Distinct > Added	<b>0.06</b>	<b>[0.02, 0.11]</b>	<u>63.0</u>
Same < Added	-0.01	[-0.09, 0.08]	1.2
Distinct > Same	0.07	[-0.01, 0.15]	<u>12.1</u>
<i>Character Speech Rate</i>			
Distinct > Added	0.09	[0.00, 0.19]	<u>21.2</u>
Same < Added	-0.01	[-0.20, 0.18]	1.1
Distinct > Same	0.10	[-0.07, 0.28]	<u>5.1</u>
<i>Temporal Warping</i>			
Distinct > Added	0.04	[-0.04, 0.12]	<u>3.3</u>
Same < Added	-0.08	[-0.22, 0.06]	<u>5.1</u>
Distinct > Same	0.12	[-0.01, 0.25]	<u>12.9</u>
<i>Pitch</i>			
Distinct > Added	0.02	[-0.05, 0.09]	2.2
Same < Added	-0.01	[-0.13, 0.11]	1.3
Distinct > Same	0.03	[-0.08, 0.15]	2.2

Table 3: Model 4 results for prosodic shifts across semantic conditions relative to the ‘Added’ baseline. ER = Evidence Ratio. Bold indicates 95% credible intervals that strictly exclude zero (certainty of non-zero magnitude). Underline indicates ER > 3 (substantial evidence for directionality).

as *Added Emoji* (Baseline), *Same Category* (Congruent), or *Distinct Category* (Divergent). The random effects structure included random intercepts for speakers and utterances, as well as by-speaker random slopes for *SemanticRelation* to capture individual variability in prosodic strategies.

### 5.3 Results

Table 3 summarises the posterior estimates and evidence ratios for each acoustic feature. Given the high inter-speaker variability inherent in prosodic realisation, we provide both 95% Credible Intervals (CI) and Evidence Ratios (ER) to distinguish between certainty of effect magnitude and probability of effect direction. Following Jeffreys’ scale (Jeffreys, 1998), we interpret an ER > 3 as “substantial” and > 10 as “strong” evidence for a directional shift, allowing for the identification of consistent trends in cases where high variance causes the 95% CI to overlap zero.

We observed a clear hierarchy where prosodic divergence was largest for Distinct Category pairings, while Same Category pairings clustered near the

Feature	-	-	-
Embeddings	6.03	6.14	<b>6.18</b>
Intensity	<b>0.58</b>	<b>0.58</b>	0.53
Speech Rate	1.28	<b>1.68</b>	0.58
Temp. Warp	4.38	<b>4.48</b>	4.45
Pitch	1.49	<b>2.67</b>	2.60

Table 4: Prosodic distances between three recordings of the sentence “The text ‘I miss you’ got me dying”, spoken by a single participant with different emoji cues. Bolded values indicate the highest distance within each row.

Added Emoji baseline. This “conflict effect” was supported by very strong directional evidence in the Prosodic Embeddings ( $\beta = 0.10$ , ER > 100), driven by a robust increase in Intensity ( $\beta = 0.06$ , ER = 63.0) and a high directional probability for increased Speech Rate (ER = 21.2) relative to the baseline. Crucially, speakers also differentiated conflict from congruence: Distinct pairs exhibited greater Temporal Warping (ER = 12.9) and Intensity (ER = 12.1) than Same pairs, indicating moderate to strong directional evidence for prosodic differentiation. Conversely, synonymous emojis functioned as prosodic stabilisers, showing negligible deviation from the baseline (ER  $\approx$  1.2) and even a potential reduction in temporal warping ( $\beta = -0.08$ , ER = 5.1). Notably, while raw Pitch Distance showed the expected directionality (*Distinct* > *Added*), the evidence was inconclusive due to high variance (ER  $\approx$  2.2). This stands in contrast to Intensity, which provided a robust univariate signal, and the Prosodic Embeddings, which showed very strong differentiation. This discrepancy likely stems from the acoustic fragility of raw pitch tracking: features like creaky voice or irregular phonation, which are common in expressive speech, introduce noise that obscures the DTW signal.

## 6 Discussion

### 6.1 Prosodic Differentiation Across Emojis

Figure 1 illustrates prosodic modulation (RQ1) for the sentence “The text ‘I miss you’ got me dying” across three emoji contexts.

The and versions exhibit comparable high, sustained pitch contours, though is more compressed in duration. In contrast, the version diverges with a rising pitch on “you” and a sarcastic profile reinforced by a loudness spike on “got”.

569 These observations align with the log-transformed  
570 distances in Table 4. The semantically similar 🙄  
571 – 🙄 pair shows minimal prosodic distance, while  
572 the 🙄 – 🙄 pair exhibits the highest divergence  
573 in temporal alignment and pitch contours. These  
574 results align with the broader hierarchy identified  
575 in Model 4 regarding the overall magnitude of di-  
576 vergence. While the general model indicates that  
577 distinct categories drive larger shifts globally, this  
578 specific case illustrates the divergence primarily  
579 through pitch and timing.

## 580 6.2 Semantic Overlap and Expressive 581 Limitation

582 While many speakers modulate prosody in re-  
583 sponse to emojis, others produce similar patterns  
584 across different emoji conditions. These overlaps  
585 raise questions about whether they stem from mean-  
586 ingful semantic similarities or individual expressive  
587 limitations, a distinction that is particularly chal-  
588 lenging when lexical content implies a dominant  
589 tone.

590 Semantic redundancy often explains this lack of  
591 contrast. For instance, in utterances like “*Happy*  
592 *birthday, have a great day*” or “*He started salsa*  
593 *dancing*,” the inherent positivity of the text aligns  
594 closely with emojis such as 😄, 😍, or 🦄. In  
595 these cases, the lack of prosodic differentiation  
596 likely reflects shared affective intent or semantic  
597 alignment rather than a failure to encode the cue.  
598 Conversely, other examples suggest expressive con-  
599 straints. The sentence “*Welp, this is gonna hurt*  
600 *to watch*” was read in an emotionally flat manner  
601 for both 😬 and 😞 despite their divergent tones  
602 of nervous tension and sadness. Such instances  
603 may reflect a difficulty or unwillingness to encode  
604 nuanced emotional distinctions in prosody, or a  
605 mismatch between the emoji and the speaker’s per-  
606 sonal interpretation.

607 Together, these examples highlight a central chal-  
608 lenge for RQ1: similar prosody across conditions  
609 may reflect semantic overlap, speaker limitations,  
610 or the deliberate downplaying of emoji cues. Dis-  
611 entangling these factors is methodologically diffi-  
612 cult and would require detailed manual inspection  
613 alongside speaker-level analysis.

## 614 6.3 Prosodic Divergence and Interpretation 615 Failure

616 Finally, we consider a case where prosodic strate-  
617 gies fail to support accurate listener interpretation  
618 (RQ2/RQ3) involving the sentence “*I’m not sur-*

619 *prised* 😲”, where 😲 is the intended emoji and  
620 🙄 the random alternative in Experiment 4.

621 The first speaker produced a lowering pitch con-  
622 tour with a subtle lengthening on “surprised”, giv-  
623 ing the utterance a flat, resigned tone signalling  
624 emotional detachment rather than shock. Post-  
625 hoc interpretation revealed the speaker viewed the  
626 emoji as representing “shame,” indicating a mis-  
627 alignment between their internal representation of  
628 😲 and its conventional use. The second speaker  
629 used a faster tempo and sharper accents that were  
630 perceived as sarcastic rather than alarmed. Listen-  
631 ers in both cases favoured the 🙄 emoji, likely due  
632 to the mismatch between these prosodic cues and  
633 the semantics of the target emoji.

634 This example highlights several key challenges.  
635 Prosodic strategies vary widely between speakers  
636 and may not always align with shared semantic  
637 expectations. Furthermore, listener interpretation  
638 depends not only on prosodic cues but also on cul-  
639 turally or contextually anchored understandings of  
640 emoji meaning. Even when prosody is used expres-  
641 sively, its communicative intent can remain opaque  
642 without a clear shared mapping between prosodic  
643 form and emoji semantics. These failures reinforce  
644 our Model 3 findings, demonstrating that success-  
645 ful communication is contingent on speakers and  
646 listeners aligning their paralinguistic intuitions.

## 647 7 Conclusion

648 This work provides the first systematic evidence  
649 linking emoji semantics to prosodic expression and  
650 interpretation in speech. We establish that emo-  
651 jis act as mediators of prosodic intent, allowing  
652 speakers to encode pragmatic nuance into other-  
653 wise static text. Our results reveal a hierarchy in  
654 this process: the magnitude of prosodic divergence  
655 is systematically governed by the semantic distance  
656 between emojis, with intensity and timing serving  
657 as primary acoustic carriers of emoji-driven intent.  
658 Finally, we demonstrate that these modulations en-  
659 able listeners to recover intended meanings with  
660 accuracy significantly above chance.

661 Beyond these findings, this study highlights emo-  
662 jis as a structured paralinguistic layer bridging the  
663 gap between digital text and spoken production.  
664 Future work could leverage these emoji–prosody  
665 mappings to enhance multimodal embeddings or  
666 drive acoustic emoji prediction, positioning emojis  
667 as vital tools for modelling the nuances of human  
668 communication.

## 669 Limitations

670 While our empirical analysis establishes a robust  
671 baseline for emoji-driven prosody, this study repre-  
672 sents an initial exploration rather than an exhaus-  
673 tive account. First, we prioritised lexical control  
674 via elicited reading over spontaneous conversation.  
675 Although participants received no instructions to  
676 focus on emojis, the resulting prosody may be more  
677 deliberate than in casual speech. Second, untrained  
678 speakers exhibited significant variation in expres-  
679 sivity. While our Bayesian models accounted for  
680 individual baselines, we observed that the clarity  
681 of prosodic cues is graded and varies by speaker,  
682 rather than being uniform across the population.  
683 Furthermore, our curated stimuli may introduce  
684 selection bias, and we did not control for cross-  
685 platform emoji rendering. Finally, as this work  
686 is restricted to British English, findings may not  
687 generalise to languages with distinct prosodic or  
688 pragmatic norms.

## 689 Ethical Considerations

690 This research was approved by a departmen-  
691 tal ethics board (details are omitted to preserve  
692 anonymity), with all participants providing explicit  
693 informed consent for the use and release of their  
694 recordings and annotations for research purposes.  
695 Recruitment was conducted via Prolific at a com-  
696 pensation rate of £9/hour, adhering to fair pay and  
697 living wage standards. To maintain privacy, data  
698 was anonymised at the source using Prolific IDs,  
699 ensuring no direct identifiers were collected. To  
700 verify data integrity and safety, we utilized an Auto-  
701 matic Speech Recognition (ASR) model (Whisper)  
702 to ensure recordings matched the target stimuli and  
703 contained no personally identifying information.  
704 This automated process was supplemented by man-  
705 ual audits of a subset of recordings and annotations  
706 to verify the absence of harmful content.

## 707 References

708 Ehab Saleh Alnuzaili, Muhammad Waqar Amin, Sami  
709 Saad Alghamdi, Nazir Ahmed Malik, Abdulbasit  
710 A. Alhaj, and Asad Ali. 2024. Emojis as graphic  
711 equivalents of prosodic features in natural speech:  
712 evidence from computer-mediated discourse of what-  
713 sapp and facebook. *Cogent Arts & Humanities*,  
714 11(1):2391646.

715 Hamza Alshenqeeti. 2016. Are emojis creating a new  
716 or old visual language for new generations? a socio-

semiotic study. *Advances in language and Literary  
Studies*, 7(6). 717  
718

Dane Archer and Robin M Akert. 1977. Words and  
everything else: Verbal and nonverbal cues in so-  
cial interpretation. *Journal of personality and social  
psychology*, 35(6):443. 719  
720  
721  
722

Dale J Barr, Roger Levy, Christoph Scheepers, and  
Harry J Tily. 2013. Random effects structure for  
confirmatory hypothesis testing: Keep it maximal.  
*Journal of memory and language*, 68(3):255–278. 723  
724  
725  
726

Isabelle Boutet, Megan LeBlanc, Justin A Chamber-  
land, and Charles A Collin. 2021. Emojis influence  
emotional communication, social attributions, and  
information processing. *Computers in Human Be-  
havior*, 119:106722. 727  
728  
729  
730  
731

Paul-Christian Bürkner. 2021. [Bayesian item response  
modeling in R with brms and Stan](#). *Journal of Statis-  
tical Software*, 100(5):1–54. 732  
733  
734

Ricardo JGB Campello, Davoud Moulavi, and Jörg  
Sander. 2013. Density-based clustering based on  
hierarchical density estimates. In *Pacific-Asia confer-  
ence on knowledge discovery and data mining*, pages  
160–172. Springer. 735  
736  
737  
738  
739

Jie Chi, Maureen de Seyssel, and Natalie Schluter. 2025.  
[The role of prosody in spoken question answering](#).  
In *Findings of the Association for Computational  
Linguistics: NAACL 2025*, pages 8468–8479, Al-  
buquerque, New Mexico. Association for Computa-  
tional Linguistics. 740  
741  
742  
743  
744  
745

Jennifer Cole. 2015. Prosody in context: A review.  
*Language, Cognition and Neuroscience*, 30(1-2):1–  
31. 746  
747  
748

David Crystal. 2001. *Language and the Internet*. Cam-  
bridge University Press. 749  
750

Justyna Czkestopchowska, Kristina Gligorić, Maxime  
Peyrard, Yann Mentha, Michał Bień, Andrea Grütter,  
Anita Auer, Aris Xanthos, and Robert West. 2022.  
On the context-free ambiguity of emoji. In *Proceed-  
ings of the International AAI Conference on Web  
and Social Media*, volume 16, pages 1388–1392. 751  
752  
753  
754  
755  
756

LouAnn Gerken and Karla McGregor. 1998. An  
overview of prosody and its role in normal and disor-  
dered child language. *American Journal of Speech-  
Language Pathology*, 7(2):38–48. 757  
758  
759  
760

Tüge T Gülşen. 2016. You tell me in emojis. In *Com-  
putational and cognitive approaches to narratology*,  
pages 354–375. IGI Global. 761  
762  
763

Nele Hellbernd and Daniela Sammler. 2016. Prosody  
conveys speaker’s intentions: Acoustic cues for  
speech act perception. *Journal of Memory and Lan-  
guage*, 88:70–86. 764  
765  
766  
767

Julia Hirschberg. 2002. Communication and prosody:  
Functional aspects of prosody. *Speech Communica-  
tion*, 36(1-2):31–43. 768  
769  
770

771	Thomas Holtgraves and Caleb Robinson. 2020. Emoji can facilitate recognition of conveyed indirect meaning. <i>PLoS one</i> , 15(4):e0232361.	
772		
773		
774	Jiaxiong Hu, Qian Yao Xu, Limin Paul Fu, and Yingqing Xu. 2019. Emojilization: An automated method for speech to emoji-labeled text. In <i>Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems</i> , pages 1–6.	
775		
776		
777		
778		
779	Lawrence Hubert and Phipps Arabie. 1985. Comparing partitions. <i>Journal of classification</i> , 2:193–218.	
780		
781	Yannick Jadoul, Bill Thompson, and Bart de Boer. 2018. <a href="#">Introducing Parselmouth: A Python interface to Praat</a> . <i>Journal of Phonetics</i> , 71:1–15.	
782		
783		
784	Allan James. 2017. Prosody and paralinguistics in speech and the social media: The vocal and graphic realization of affective meaning. <i>Linguistica</i> , 57(1):137–149.	
785		
786		
787		
788	Harold Jeffreys. 1998. <i>The theory of probability</i> . OUP Oxford.	
789		
790	Elsi Kaiser. 2021. Focus marking with emoji: On the relation between information structure and expressive meaning. In <i>Colloque de syntaxe et sémantique à Paris (CSSP), Paris, France</i> .	
791		
792		
793		
794	Joseph Kane, Michael N Johnstone, and Patryk Szweczyk. 2024. Voice synthesis improvement by machine learning of natural prosody. <i>Sensors</i> , 24(5):1624.	
795		
796		
797		
798	Ilse Lehiste and Norman J Lass. 1976. Suprasegmental features of speech. <i>Contemporary issues in experimental phonetics</i> , 225:239.	
799		
800		
801	Li Li and Yue Yang. 2018. Pragmatic functions of emoji in internet-based communication—a corpus-based study. <i>Asian-Pacific Journal of Second and Foreign Language Education</i> , 3:1–12.	
802		
803		
804		
805	Leland McInnes, John Healy, Nathaniel Saul, and Lukas Grossberger. 2018. Umap: Uniform manifold approximation and projection. <i>The Journal of Open Source Software</i> , 3(29):861.	
806		
807		
808		
809	Hannah Miller, Daniel Kluver, Jacob Thebault-Spieker, Loren Terveen, and Brent Hecht. 2017. Understanding emoji ambiguity in context: The role of text in emoji-related miscommunication. In <i>Eleventh international AAAI conference on web and social media</i> .	
810		
811		
812		
813		
814	Max Morrison, Caedon Hsieh, Nathan Pruyne, and Bryan Pardo. 2023. Cross-domain neural pitch and periodicity estimation. In <i>arXiv preprint arXiv:2301.12258</i> .	
815		
816		
817		
818	Noa Na’aman, Hannah Provenza, and Orion Montoya. 2017. <a href="#">Varying linguistic purposes of emoji in (Twitter) context</a> . In <i>Proceedings of ACL 2017, Student Research Workshop</i> , pages 136–141, Vancouver, Canada. Association for Computational Linguistics.	
819		
820		
821		
822		
	R Core Team. 2025. <i>R: A Language and Environment for Statistical Computing</i> . R Foundation for Statistical Computing, Vienna, Austria.	823 824 825
	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-BERT: Sentence embeddings using Siamese BERT-networks</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	826 827 828 829 830 831 832 833
	Seamless Communication et al. 2023. Seamless: Multilingual expressive and streaming speech translation. <i>arXiv preprint arXiv:2312.05187</i> .	834 835 836
	RJ Skerry-Ryan, Eric Battenberg, Ying Xiao, Yuxuan Wang, Daisy Stanton, Joel Shor, Ron Weiss, Rob Clark, and Rif A Saurous. 2018. Towards end-to-end prosody transfer for expressive speech synthesis with tacotron. In <i>international conference on machine learning</i> , pages 4693–4702. PMLR.	837 838 839 840 841 842
	Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. MpNet: Masked and permuted pre-training for language understanding. <i>Advances in neural information processing systems</i> , 33:16857–16867.	843 844 845 846 847
	Ioannis Tsiamas, Matthias Sperber, Andrew Finch, and Sarthak Garg. 2024. <a href="#">Speech is more than words: Do speech-to-text translation systems leverage prosody?</a> In <i>Proceedings of the Ninth Conference on Machine Translation</i> , pages 1235–1257, Miami, Florida, USA. Association for Computational Linguistics.	848 849 850 851 852 853
	Paige Tuttösi, Shivam Mehta, Zachary Syvenky, Bermet Burkanova, Gustav Eje Henter, and Angelica Lim. 2025. Emojivoice: Towards long-term controllable expressivity in robot speech. <i>arXiv preprint arXiv:2506.15085</i> .	854 855 856 857 858
	Nguyen Xuan Vinh, Julien Epps, and James Bailey. 2010. <a href="#">Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance</a> . <i>Journal of Machine Learning Research</i> , 11(95):2837–2854.	859 860 861 862 863
	Michael Wagner. 2016. How to be kind with prosody. In <i>Speech prosody</i> , 1:250–1253.	864 865
	Sarenne Wallbridge, Christoph Minixhofer, Catherine Lai, and Peter Bell. 2025. Prosodic structure beyond lexical content: a study in self-supervised learning. In <i>proceedings of Interspeech 2025</i> .	866 867 868 869
	Nigel Ward and Gina-Anne Levow. 2021. <a href="#">Prosody: Models, methods, and applications</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: Tutorial Abstracts</i> , pages 26–28, Online. Association for Computational Linguistics.	870 871 872 873 874 875 876

877 Lukas Wolf, Tiago Pimentel, Evelina Fedorenko, Ryan  
 878 Cotterell, Alex Warstadt, Ethan Wilcox, and Tamar  
 879 Regev. 2023. [Quantifying the redundancy between  
 880 prosody and text](#). In *Proceedings of the 2023 Confer-  
 881 ence on Empirical Methods in Natural Language  
 882 Processing*, pages 9765–9784, Singapore. Associa-  
 883 tion for Computational Linguistics.

884 Giulio Zhou, Sydelle De Souza, Ella Markham,  
 885 Oghenetekevwe Kwakpovwe, and Sumin Zhao.  
 886 2024a. [Semantics and sentiment: Cross-lingual varia-  
 887 tions in emoji use](#). In *Proceedings of the 2024 Confer-  
 888 ence on Empirical Methods in Natural Language Pro-  
 889 cessing*, pages 18698–18712, Miami, Florida, USA.  
 890 Association for Computational Linguistics.

891 Giulio Zhou, Tsz Kin Lam, Alexandra Birch, and Barry  
 892 Haddow. 2024b. [Prosody in cascade and direct  
 893 speech-to-text translation: a case study on Korean  
 894 wh-phrases](#). In *Findings of the Association for Com-  
 895 putational Linguistics: EACL 2024*, pages 674–683,  
 896 St. Julian’s, Malta. Association for Computational  
 897 Linguistics.

## 898 A Model Specifications

899 All models were implemented in R ([R Core  
 900 Team, 2025](#), version 4.5.2) using the brms pack-  
 901 age ([Bürkner, 2021](#)), which interfaces with the Stan  
 902 probabilistic programming platform. Binary  
 903 outcomes for interpretation and perception trials  
 904 (RQ1, RQ2, and RQ3) were modelled using a  
 905 Bernoulli distribution with a logit link function.  
 906 Continuous prosodic distance measures (RQ4)  
 907 were modelled using Gaussian distributions. We  
 908 applied weakly informative priors to the fixed ef-  
 909 fects ( $Normal(0, 0.5)$ ) to provide regularisation  
 910 and ensure stable posterior convergence. This  
 911 prior choice reflects a conservative assumption that  
 912 large effect sizes are relatively unlikely. Centred  
 913 around a null effect of zero, these priors help pre-  
 914 vent the models from over-fitting to noise in the  
 915 high-variance speech data. Following the princi-  
 916 ple of the maximal random effects structure ([Barr  
 917 et al., 2013](#)), we included random intercepts for  
 918 speakers, listeners, and utterances to account for  
 919 the non-independence of observations. Where ex-  
 920 perimental design permitted, specifically in Model  
 921 4, we included random slopes for semantic con-  
 922 ditions by speaker. This accounts for individual  
 923 variability in how speakers adapt their prosodic  
 924 strategies to different emoji cues, ensuring that our  
 925 population-level estimates are robust to speaker-  
 926 level idiosyncratic behaviour.

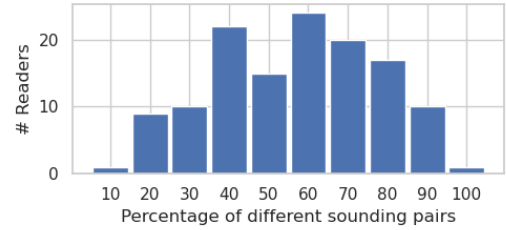


Figure 2: Readers (Experiment 2) grouped by the proportion of recordings judged as prosodically different (Experiment 3). Each bin represents a 10% range (e.g., 0%–9%), from least to most prosodically varied speakers.

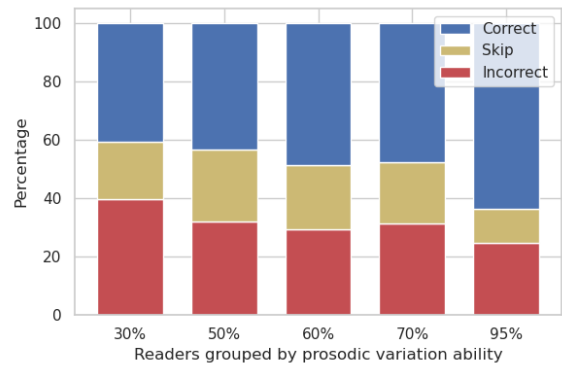


Figure 3: Listener (Experiment 3) accuracy in emoji assignment for prosodically distinct pairs, grouped by speaker variation. Bars show correct, incorrect, and skipped responses from least to most prosodically varied speakers.

## 927 B Data statistics

### 928 B.1 Speaker Expressivity

929 The distribution in Figure 2 categorises speakers  
 930 by the proportion of their recording pairs judged as  
 931 prosodically distinct. While 20 readers exhibited  
 932 minimal differentiation (0–30%), the majority (81  
 933 speakers) showed moderate variation (31–70%),  
 934 and a notable subset (28 speakers) consistently pro-  
 935 duced distinct renditions (71–100%). These re-  
 936 sults demonstrate that while individual variability  
 937 is present, the prevalence of distinct recordings sup-  
 938 ports the conclusion that emoji cues systematically  
 939 influence how prosody is realised.

### 940 B.2 Listener Accuracy - Intra Speaker

941 Figure 3 illustrates how a speaker’s expressiveness  
 942 directly impacts communicative success. A Spear-  
 943 man’s rank correlation revealed a significant posi-  
 944 tive relationship between speaker variation and  
 945 correct emoji identification ( $\rho = 0.90, p = .037$ )

Response	Overall	Correct	Partial	Incorrect
Same	34.74%	66.67%	6.06%	27.27%
Different	65.26%	14.75%	74.59%	10.66%

Table 5: Results from the inter-speaker experiment showing how listeners judge prosodic similarity between two recordings and how these judgments affect their accuracy in assigning each recording to the correct emoji.

and a significant negative correlation for incorrect assignments ( $\rho = -0.90, p = .037$ ). For the most expressive speakers (95% bin), accuracy reached 62.50%, significantly exceeding both the incorrect rate (25.52%) and a random-choice baseline. Conversely, in the lowest variation group (30%), correct (41.09%) and incorrect (39.11%) rates were nearly equal, indicating that listeners were reduced to random guessing. While skip rates dropped from 20% to 11.98% in the highest bin, this trend was not statistically significant ( $\rho = -0.10, p = .87$ ).

### B.3 Listener Accuracy - Inter Speaker

Results in Table 5 show that 34.74% of inter-speaker pairs were judged as sounding similar, while 65.26% were judged as different. This suggests that speakers often employ flexible prosodic strategies for the same emoji. In this experiment, listeners were presented with two recordings and two emoji options, allowing for four possible assignment combinations. We define Correct as the proper identification of the target emoji for both recordings, while Partial correctness occurs when a listener assigns the target emoji to only one of the two recordings. Incorrect reflects a failure to identify the target emoji in either instance. When listeners perceived prosodic similarity, they achieved full correctness 66.67% of the time, indicating strong interpretive agreement when prosody converges. For pairs judged as different, the high partial correctness rate (74.59%) reflects increased uncertainty; listeners often successfully “recovered” the emoji for one speaker while failing for the other. However, even when prosody differed, the full correctness rate (14.75%) still exceeded the incorrect rate (10.66%), suggesting that divergent realisations can still convey intended meaning, albeit with less reliability.

Predictor	Est. ( $\beta$ )	95% CI
(Intercept)	-1.867	[-2.279, -1.492]
Speaker Expressivity	-0.182	[-0.368, 0.003]
Contrast Type (Single)	0.196	[-0.203, 0.591]
Speaker Emoji Use	-0.045	[-0.233, 0.135]
Listener Emoji Use	0.151	[-0.184, 0.481]
Speaker $\times$ Listener Use	-0.063	[-0.257, 0.126]

Table 6: Bayesian multilevel logistic regression results predicting the probability of a “semantic tie” (assigning the same emoji to both recordings)

## C Analysis of Perceptual Overlap - RQ2

To further investigate the mechanics of communicative success, we analysed the likelihood of “semantic ties.” These represent instances where a listener assigned the same emoji to both recordings within a single trial, modelled as a “Skip” in Figure 3. A high rate of semantic ties indicates perceptual overlap, where the speaker’s prosodic modulation is insufficient to signal a change in symbolic intent. We modelled the probability of such ties using a Bayesian multilevel logistic regression:

$$\begin{aligned} \text{Model 5: SemanticTie} \sim & \text{SpeakerExpressivity} + \text{ContrastType} + \\ & (\text{SpeakerEmojiUse} \times \text{ListenerEmojiUse}) + \\ & (1 \mid \text{Utterance}) + (1 \mid \text{Listener}) \end{aligned}$$

The results in Table 6 show a strongly negative intercept ( $\beta = -1.87, CI[-2.28, -1.49]$ ), indicating that the baseline probability of a semantic tie is very low ( $P \approx 0.13$ ). This suggests that listeners are naturally predisposed to seek out prosodic distinctions rather than perceiving the recordings as identical. The primary driver of further reducing this perceptual ambiguity was Speaker Expressivity ( $\beta = -0.18, CI[-0.37, 0.003]$ ). While the credible interval narrowly overlaps zero, the negative trend suggests that higher expressive bandwidth effectively helps listeners to recognise distinct meanings. By producing more distinctive prosodic signals, highly expressive speakers successfully disambiguate the potential semantic tie, making it easier for listeners to identify that two different meanings were intended. This finding reinforces the conclusion that the interpretability of emoji-enriched speech is directly contingent on the speaker’s ability to vocalise pragmatic contrasts.

	Ambiguity	Correct	Partial	Incorrect
?	0.1577	25.0	75.0	0.0
😄	0.1707	83.3	16.7	0.0
👉	0.2131	0.0	100.0	0.0
😞	0.2443	66.7	33.3	0.0
😱	0.3750	0.0	20.0	80.0
😞	0.4630	100.0	0.0	0.0
😄	0.5158	66.7	0.0	33.3
😞	0.5410	22.2	44.4	33.3
👍	0.6593	0.0	50.0	50.0

Table 7: Examples illustrating emoji ambiguity and Experiment 4 listener interpretation accuracy (in percentages), sorted from low to high ambiguity.

## D Listener Accuracy and Emoji Ambiguity

Table 7 presents a subset of emojis from Experiment 4 alongside listener interpretation accuracy and ambiguity scores derived from the semantic variation metrics reported by Czkestopchowska et al. (2022). These scores reflect cross-participant variability in emoji interpretation when presented without context, where higher values indicate greater semantic ambiguity. Although the sample size precludes statistical generalisation, the examples illustrate a qualitative trend: emojis with lower ambiguity scores tend to yield fewer incorrect interpretations, whereas highly ambiguous emojis result in greater listener disagreement.

Specific emojis revealed distinct patterns within the data. Emojis such as 👉 and 😞 demonstrated clear prosodic convergence alongside high accuracy. In contrast, emojis like 😞 and 😄 were reliably recognised despite divergent prosody, suggesting that inherent semantic context may at times override prosodic mismatch. Conversely, for 😱, 😞, and 😞, prosodic similarity did not consistently ensure accurate interpretation.

Qualitative comparison with existing ambiguity scores suggests that while semantic ambiguity does not dictate the degree of prosodic convergence, emojis with lower inherent ambiguity tend to result in fewer incorrect interpretations. These tentative patterns suggest that higher semantic ambiguity may reduce the recoverability of emoji intent from prosody alone, as listeners lack a stable semantic “anchor” to which they can map vocal variations.

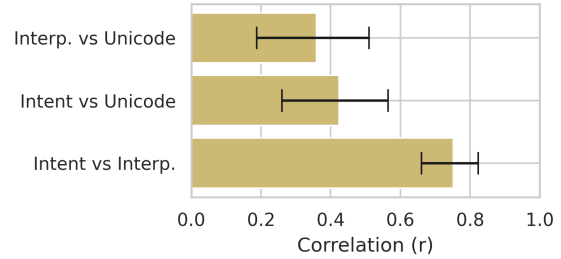


Figure 4: Mean Pearson correlations between emoji meaning spaces (Intent vs Interpretation vs Unicode Descriptions). Error bars indicate 95% confidence intervals computed using Fisher z-transformation.

## E Emoji Meaning Analysis

To analyse the semantic consistency of emoji usage, we compared three meaning representations for each emoji: the *intent* (from Experiment 1), the *interpretation* (from Experiment 2), and the official Unicode *description*.

For each emoji, we aggregated the most frequent annotations and merged them into a representative string per meaning type. We then generated sentence embeddings using the sentence-transformers/all-mpnet-base-v2 model (Reimers and Gurevych, 2019; Song et al., 2020). This yielded one embedding per emoji for each semantic category.

We applied UMAP (McInnes et al., 2018) to reduce the embeddings to two dimensions for visualisation. Emojis were plotted directly in 2D space using their glyphs, after normalising representations (e.g., removing skin tone and gender modifiers) to avoid redundancy and improve clarity (Figure 5).

We evaluated the alignment between the three emoji meaning spaces by computing pairwise Pearson correlation coefficients between corresponding emoji embeddings. To assess the reliability of these correlations across the dataset, we applied the Fisher z-transformation, which stabilises the variance of correlations and enables us to compute statistically sound confidence intervals for the mean. Figure 4 shows the mean correlation and associated 95% confidence intervals across all emojis with at least five annotations. We observe the strongest alignment between intent and interpretation embeddings (mean  $r = 0.752$ , CI: [0.660, 0.823]), reflecting high agreement between what users intend and how others interpret emoji use. In contrast, the unicode description space shows significantly

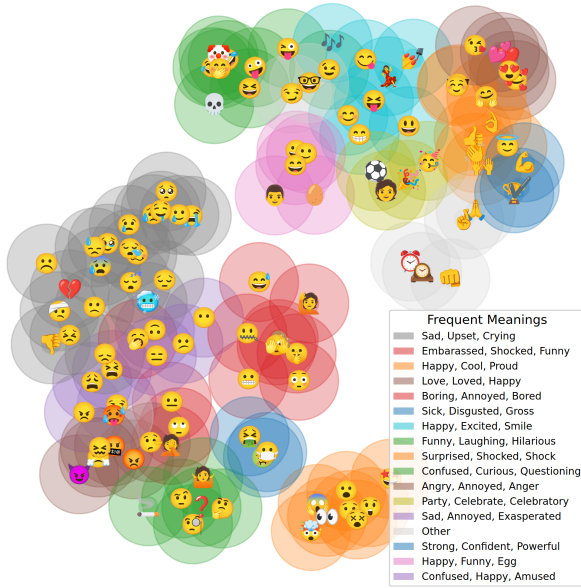


Figure 5: Emoji embeddings reduced with UMAP. Clusters computed with HDBSCAN (on 49 components). Labels are top 3 keywords among the words used to describe the emojis.

weaker alignment with both intent ( $r = 0.424$ , CI: [0.260, 0.564]) and interpretation ( $r = 0.359$ , CI: [0.188, 0.510]), reinforcing prior findings that Unicode definitions often fail to capture how emojis are actually used in communicative context.

## F Emoji Clustering

To explore emergent semantic groupings among emojis, we performed unsupervised clustering on the embeddings derived from Experiments 1 and 2, as well as from the emoji Unicode descriptions. We first applied Principal Component Analysis (PCA) to determine the intrinsic dimensionality of the embedding space. Based on the explained variance, we selected 49 dimensions to preserve approximately 90% of the variance.

Clustering was performed using HDBSCAN (Campello et al., 2013, Hierarchical Density-Based Spatial Clustering of Applications with Noise), a density-based algorithm that can identify an optimal number of clusters based on the density distribution of the data. HDBSCAN does not require the number of clusters to be predefined and is robust to noise and outliers. Emojis that did not belong to any cluster were labelled as noise. We set `min_cluster_size` to 3 to allow for finer-grained clusters, and used Euclidean distance as the similarity metric.

Figure 5 shows the two-dimensional projections

Cluster Comparison	ARI	NMI
Int. vs Inter.	0.249	0.613
Int. vs Desc.	0.077	0.431
Inter. vs Desc.	0.062	0.443

Table 8: Adjusted Rand Index (ARI) and Normalized Mutual Information (NMI) between clusterings based on different semantic spaces.

of the resulting embeddings and clusters. The Figure shows labels for each cluster derived from the most frequent emoji interpretations within the group. These derived clusters serve as our semantic categories in subsequent analyses.

To assess the alignment between emoji clusters derived from different semantic spaces (Intent, Interpretation, and Description), we computed the *Adjusted Rand Index* (Hubert and Arabie, 1985, ARI) and *Normalized Mutual Information* (Vinh et al., 2010, NMI) between each pair of clusterings. These metrics quantify the similarity of cluster assignments while accounting for chance (ARI) and shared information (NMI).

The ARI scores suggest moderate agreement between *Intent* and *Interpretation* clusters (ARI = 0.249), whereas clusterings involving *Description* show weak alignment (ARI < 0.1). NMI results are consistent with this pattern: the Intent and Interpretation spaces share more information (NMI = 0.613) than either does with the Description space (NMI  $\approx$  0.43). These findings indicate that user-generated intentions and interpretations exhibit a shared semantic structure that is not well captured by official Unicode descriptions.

## G Instructions and Trial Samples

### G.1 Experiment 1

Participants were given the following instructions:“

1. Select ONE emoji and place it in any suitable position within the provided sentence. Only use emojis and spaces, avoiding inserting any other elements.
2. Describe how the chosen emoji modifies the sentence. Keep it short! One word should be enough.
3. Repeat until you have filled out at least two forms with a suitable emoji and its corresponding descriptions.

1154 *Avoid choosing multiple emojis that are too similar.*

1155 *For instance, if you select a happy emoji, con-*  
1156 *sider avoiding another happy emoji to add diversity*  
1157 *to your selections.”*

1158 Figure 6a shows an example of a trial page for  
1159 Experiment 1.

## 1160 G.2 Experiment 2

1161 Participants were asked to complete the Recording  
1162 Tasks first, and then provide their interpretation  
1163 of the encountered emojis. Here, the instructions  
1164 given to the participants.

### 1165 G.2.1 Recording Task

1166 *“Click on the mic icon to start recording audio, then*  
1167 *click again to stop.*

1168 *You will be able to playback your recording. You*  
1169 *can confirm by pressing ‘Continue’ or try again*  
1170 *by pressing ‘Record again’ (it will start recording*  
1171 *immediately).*

1172 *Try to record in a quiet environment.”*

1173 Figure 6b shows the recording interface.

### 1174 G.2.2 Annotation Task

1175 *“In this second (and last) part, you will be asked to*  
1176 *describe how you interpreted the emoji used in the*  
1177 *previous section.*

1178 *Keep it short! One word should be enough.*

1179 *Leave blank if there wasn’t any emoji.”*

## 1180 G.3 Experiment 3

1181 Figure 6c shows an example of the interface for  
1182 Experiment 3. The same interface was presented  
1183 for experiments with Inter and Intra-Speaker Data.

### 1184 G.3.1 Inter-Speaker Data

1185 *“In this study, you will listen to pairs of recorded*  
1186 *speech. Each pair will contain the same words,*  
1187 *but they may differ in how they are spoken. For*  
1188 *example, in rhythm, intonation, intensity, use of*  
1189 *pauses and so on. Your task is to determine whether*  
1190 *there are intentional variations between the two*  
1191 *recordings. Please disregard any differences that*  
1192 *do not seem deliberate. For example, small lexical*  
1193 *variations or stutter.*

1194 *You will be provided with transcriptions of the*  
1195 *recordings, which may contain different emojis or*  
1196 *none at all. Your task is to match each recording to*  
1197 *the transcription that best corresponds to it.*

1198 *Try to assign each recording to a different tran-*  
1199 *scription first!”*

### 1200 G.3.2 Intra-Speaker Data

1201 *“In this study, you will listen to pairs of recorded*  
1202 *speech from two different speakers. Each pair will*  
1203 *contain the same words, but they may differ in*  
1204 *how they are spoken, for example, in rhythm, in-*  
1205 *tonation, intensity, use of pauses and so on. Your*  
1206 *task is to determine whether there are intentional*  
1207 *variations between the two recordings. Please dis-*  
1208 *regard any differences that do not seem deliberate.*  
1209 *For example, small lexical variations or stutter.*

1210 *You will hear two recordings of the same sen-*  
1211 *tence, spoken by different speakers.*

1212 *Your task is to listen carefully to each recording*  
1213 *and choose the emoji that best matches its tone.*

1214 *Each speaker produced their recording based on*  
1215 *an emoji prompt.*

1216 *You can assign the same emoji to both record-*  
1217 *ings, or a different one to each, depending on what*  
1218 *you perceive.”*

**Add ONE emoji to each sentence (at least two)  
and describe how the emojis modify their meaning  
(preferably in ONE word)**

1. \_\_\_\_\_ Marlon is such a clown. Why did he do that? \_\_\_\_\_ *Description*
2. \_\_\_\_\_ Marlon is such a clown. Why did he do that? \_\_\_\_\_ *Description*
3. \_\_\_\_\_ Marlon is such a clown. Why did he do that? \_\_\_\_\_ *Description*
4. \_\_\_\_\_ Marlon is such a clown. Why did he do that? \_\_\_\_\_ *Description*
5. \_\_\_\_\_ Marlon is such a clown. Why did he do that? \_\_\_\_\_ *Description*

(a) Experiment 1 - Emoji Annotation Task


Looks awesome! I look forward to this! 😊





(b) Experiment 2 - Recording Task

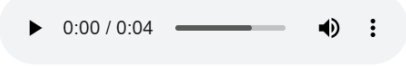
**Listen to the recordings and assign them the best matching transcriptions**



[1/20]

A -  0:00 / 0:04

This league makes no sense  This league makes no sense 

---

B -  0:00 / 0:04

This league makes no sense  This league makes no sense 

---

Do the two recordings sound the same?

(c) Experiment 3 and 4 - Listening Task

Figure 6: Example of trials' main page for the online experiments.