

# From Fiction to Fact: Fine-Grained Emotion Classification in COVID-19 Newspaper Discourse

Anonymous ACL submission

## Abstract

This study examines how a computational literary studies (CLS) emotion classification framework can be adapted to analyze newspaper discourse on COVID-19. We developed and tested single-layer and dual-layer BERT models to classify emotions at two levels: 9 primary emotion families (Level 1) and 87 subcategories (Level 2). Using 7,498 sentences from German newspapers, data sparsity directed our focus to the 10 most common Level-2 emotions. Our results revealed varied model performances across emotion categories. The single-layer model exhibited more consistent performance and a stronger correlation with emotion frequency. In contrast, the dual-layer model excelled at distinguishing specific emotions like interest, curiosity, and hope, although with greater variability. Both models struggled to recognize more complex emotions such as LOVE, DISGUST, and AMBIVALENCE. Our results underscore the complexities and potential of automated emotion detection in media discourse, highlighting the need for domain-specific classification methods.

## 1 Introduction

Emotions are a fundamental part of human cognition and life. Their expression is particularly diverse in language (Schwarz-Friesel, 2007), being omnipresent not only in literary texts (Anz, 2007) but also playing a significant role in seemingly neutral genres such as news reports. However, the systematic analysis and classification of emotions in language remains elusive.

The goals of the present paper are to (i) evaluate an emotion annotation framework developed for fictional texts and apply it to a different genre, i.e. news reporting, and (ii) to develop a classifier for

emotion annotation in news reporting and thus determine the extent to which our emotion annotation framework can be generalized to unseen data.

## 2 Analysis of Emotion in Language and Text

The study of emotions has a rich history across various disciplines, offering diverse conceptualizations. Aristotle defined 15 basic emotions (incl. desire, anger, fear, and joy), while later philosophers such as Descartes and Hume offered different taxonomies, ranging from two to six fundamental emotions (Süselbeck, 2019). In contemporary research, emotion classification mainly follows two approaches: structure-oriented and function-oriented (Schwarz-Friesel, 2017).

Structure-oriented classifications conceptualize emotions as innate and culture-independent and assume that certain emotions emerge from the architecture of the human brain (Damasio 1997, 2004). One example is Ekman's (1972, 1988, 1994) influential model of seven basic emotions—happiness, anger, sadness, fear, disgust, surprise, and contempt, whereas Plutchik's (1984) "wheel of emotions" refers to eight categories. Other notable structural-oriented frameworks list five (Oatley and Johnson Laird 1987), six (Argyle 1996), seven (Scherer (1993) or ten (Izard 1992) emotion types.

Function-oriented classifications, by contrast, differentiate emotions based on their referential targets and situational conditions. These include distinctions between target-oriented and non-directional emotions, environment-, body-, or pleasure-related emotions, and categorizations of relational, empathy, and target emotions (Mees, 1985; Holodyski, 2006). Furthermore, linguistic research has explored the pragmatics of the expression of emotion across communicative contexts, as well as semantic analyses on the linguistic representations

of emotions (Schwarz-Friesel, 2007). As the above discussion makes clear, there is no psychologically or linguistically motivated consensus as to how emotion in language should be analyzed. As a result, computational applications have pursued a variety of different emotion classifications.

A distinction can be made between the classification of positive versus negative statements and emotion detection, both of which have been explored with machine learning approaches (Ahman 2011; Perikos & Hatzilygeroudis 2016; Al-Baity et al. 2022; Machová et al. 2023; Maruf et al. 2024). Emotion analysis has in particular focused on social media texts due to their accessibility, processability, and abundance of data (Klinger et al., 2020; for an overview, see Acheampong et al., 2020; Liu, 2020; Peng et al., 2022).

Recent research has demonstrated the value of using BERT for emotion classification (Khurdula et al., 2024; Papadimitriou et al., 2024), which we will employ in our own analysis, alongside LSTMs. Moreover, while previous research has mostly focused on smaller emotion inventories (typically 10 or fewer), some recent research has adopted more fine-grained classifications, such as the 27 classes in Singh (2023), and 80 in Luca et al. (2024). We will follow this recent trend.

## 2.1 Analysis of Emotion in Computational Literary Studies

Emotions also play a fundamental role in literary texts. Since the 1990s, computational literary studies (CLS) have explored emotions in literature (Flüh, 2019; Winko, 2019), employing two main methodologies: lexicon-based (Strapparava & Valitutti, 2004; Taboada & Gillies, 2006; Bruggmann & Fabrikant, 2014; Lehmann, Mittelbach & Schmeier, 2017) and machine-learning approaches (Schmidt et al., 2018, Konle et al. 2022, 2023). Lexicon-based methods determine emotional content based on predefined dictionaries, but struggle with context-specific meanings. By contrast, machine-learning approaches generalize emotional content recognition from annotated data, offering greater adaptability to context-dependent usage (Schmidt et al., 2018).

CLS and CL share methodological roots but diverge in objectives; CLS typically focuses on fictional texts, whereas CL is applied to everyday and factual language. Both fields, however, face the limitations of lexicon-based sentiment analysis, which necessitates the development of large

language models (LLMs) fine-tuned for specific contexts and relying on extensive manual annotation (Borst et al., 2023).

## 2.2 A New Emotion Classification System

The discussion of previous research has shown that more fine-grained emotion classifications may be better suited for emotion analysis, but have so far been explored by only a limited number of studies. Our methodology is grounded in a text-centered, inductive approach to emotion classification. Starting with an exploratory annotation phase, we systematically identified and categorized emotion-bearing segments in fictional texts to develop a comprehensive tag set. This data-driven approach differs from traditional deductive frameworks that apply predetermined psychological or sociological categories. The resulting classification system captures the nuanced ways emotions manifest in fictional texts. In a next step, we examined (i) to what extent this emotion classification can be adapted to a different domain, i.e. news reporting, and (ii) to what extent this emotion classification can be generalized with the help of machine learning.

## 2.3 Emotions in Public Discourse

Newspapers play a critical role in shaping public discourse, oscillating between neutral reporting and emotionally charged narratives. However, even reporting aiming at neutrality and objectivity contains emotional language (Stenvall, 2014; Zappettini et al., 2021). This variety in journalistic styles makes newspapers an ideal test case to explore how media narratives influence societal perspectives and reactions through automated analysis (Schmitz, 2016; Storjohann & Cimander, 2022). In particular, we focus on reporting on the COVID-19 pandemic, in which emotional language played an important role (Lemor and Montpetit, 2024; Zhunis et al., 2022) and which had great practical relevance in its influence on public attitudes and compliance with rules and restrictions aimed at mitigating the spread and impact of the pandemic (Généreux et al., 2021).

## 3. Aims

This study bridges computational literary studies (CLS) and computational linguistics (CL) by

176 adapting an emotion classification system from  
177 literary texts to analyze newspaper articles about  
178 the COVID-19 pandemic. By adopting a detailed  
179 emotion classification system from CLS, we pur-  
180 sue two objectives:

181 (1) to test the automatic classification of an exten-  
182 sive category system originally developed for lit-  
183 erary studies (domain adaptation)

184 (2) to enable a more multifaceted analysis of emo-  
185 tions that reveals subtle emotional nuances and in-  
186 termediate states.

187 We aim to assess how well this framework cap-  
188 tures nuanced emotional expressions in public  
189 discourse, distinguishing relevant from irrelevant  
190 emotions. This domain adaptation seeks to de-  
191 velop a detailed understanding of emotional ex-  
192 pression in news media while evaluating general-  
193 ization potential to unseen data. Our assumption  
194 that human emotions extend beyond basic types to  
195 include complex feelings like nostalgia, envy, and  
196 pride underpins this approach, aiming for a more  
197 comprehensive view of human communication.

## 198 4. Data and Methods

199 In the present study, we expand on previous work  
200 that involved pilot annotations of emotions in 10  
201 literary texts to develop and operationalize a nu-  
202 anced tagset for emotion categories (AUTHOR  
203 A).

### 204 4.1 Emotion Classification Scheme

205 Our hierarchical classification spans two levels: 9  
206 Level-1 emotion families (in small capitals) and  
207 87 Level-2 subcategories (in italics). Within sub-  
208 categories, emotions are ordered by intensity (see  
209 Appendix 1).

210 The Level-1 categories consist of six basic  
211 emotions (LOVE, JOY/HAPPINESS, DISGUST, FEAR,  
212 GRIEF, and ANGER) and three additional categories  
213 (UNCATEGORIZED\_POS, UNCATEGORIZED\_NEG, AMBIV-  
214 ALENCE) to capture emotional states that extend be-  
215 yond the basic emotions.

216 UNCATEGORIZED\_POS and UNCATEGORIZED\_NEG  
217 include emotions that, while clearly positive or  
218 negative, could not be definitively assigned to any

219 basic emotion during initial annotation. AMBIVA-  
220 LENCE encompasses emotional states that simulta-  
221 neously exhibit both positive and negative quali-  
222 ties.<sup>1</sup>

## 223 4.2 Data

224 Our data is drawn from reporting on the COVID-  
225 19 pandemic and consists of a random sample of  
226 7,500 sentences drawn from 59 German newspa-  
227 pers sourced from Lexis Nexis, spanning from  
228 January 2020 to June 2022. All sentences con-  
229 tained at least one keyword from the semantic cat-  
230 egories of vaccination, COVID-19 names, or non-  
231 pharmaceutical interventions (e.g., lockdowns).  
232 Two sentences were removed because they con-  
233 sisted of one word only, leaving 7,498 sentences  
234 to be annotated.

235 Three human raters independently annotated  
236 each sentence for 87 Level-2 emotions, noting for  
237 each emotion whether it was present, absent, or  
238 potentially present. Annotators were provided  
239 with comprehensive guidelines outlining annota-  
240 tion criteria and examples for identifying the tar-  
241 get emotions in the dataset. The final dataset con-  
242 sisted of mean scores calculated from all three  
243 raters for Level-2 emotions. The frequency of the  
244 87 Level-2 emotions across the dataset was highly  
245 uneven and followed a Zipf distribution, with a  
246 few very frequent and many fairly infrequent  
247 emotion subcategories. For analysis, we focused  
248 on the 33 Level-2 emotions that occurred at least  
249 20 times (Sec. 5.1). However, to aggregate the  
250 emotions into their corresponding Level-1 catego-  
251 ries, the maximum scores of all 87 Level-2 emo-  
252 tions were used, regardless of their frequency  
253 (Sec. 5.2 and 5.3).

254 To illustrate our data, we present examples for  
255 the three most frequently occurring Level-2 emo-  
256 tions. The most frequent emotion, *concern*, was  
257 found in expressions of worry about ongoing  
258 COVID impacts, as shown in:

259 1. "Corona hat dafür gesorgt, dass es in den  
260 kommenden Jahren angespannt bleibt", sagte  
261 Wandrey. ["COVID has ensured that the situation  
262 will remain tense in the coming years," Wandrey  
263 said.]

264 2. Wenn sich Menschen mit einer  
265 Bedrohung wie der Corona-Pandemie

---

<sup>1</sup> We also include an *uncertain* category for emotion-  
ally charged text passages that are not captured by our  
classification system.

266 konfrontiert fühlen, gelangen sie häufig zu zwei  
267 Erkenntnissen, analysiert Perel: Die Welt, wie wir  
268 sie kennen, geht gerade verloren. [When people  
269 are confronted by a threat like the COVID pan-  
270 demic, they often realize two things, according to  
271 Perel: The world as we know it is currently being  
272 lost.]

273 The second most frequent emotion, *serious-*  
274 *ness*, was identified in statements on protective  
275 measures:

276 3. Der CSU-Politiker betonte aber, dass die  
277 bestehende Kontaktbeschränkung und das  
278 Distanzgebot weiterhin gelten. [However, the  
279 CSU politician emphasized that existing contact  
280 restrictions and distance requirements continue to  
281 apply.]

282 4. ZWIESEL Evangelische Gemeinde: Ab  
283 dem morgigen Sonntag ist in der Kreuzkirche das  
284 Tragen einer FFP2-Maske während des  
285 Gottesdienstes Pflicht. [ZWIESEL Protestant  
286 Church: From tomorrow, Sunday, FFP2 masks  
287 are compulsory during the church services in the  
288 Kreuzkirche.]

289 Finally, *disapproval* emerged as the third most  
290 frequently annotated Level-2 emotion, typically  
291 in sentences containing criticism towards regula-  
292 tions or social behavior during the pandemic:

293 5. Die geplanten Ausgangsbeschränkungen  
294 zwischen 21 und fünf Uhr seien bei einer Inzidenz  
295 von 100 "ein unverhältnismäßiger und  
296 epidemiologisch unbegründeter Eingriff in die  
297 Freiheit" der Bürger. [The planned curfews be-  
298 tween 9 p.m. and 5 a.m. at an incidence of 100 are  
299 "a disproportionate and epidemiologically un-  
300 founded encroachment on the freedom" of citi-  
301 zens.]

302 6. In Anbetracht steigender Covid-19-  
303 Zahlen in Japan mangle es Bach wohl an  
304 "normalem Menschenverstand", sagte etwa ein  
305 Regierungsberater für Infektionskrankheiten. [In  
306 view of rising COVID-19 numbers in Japan, Bach  
307 seems to lack "basic common sense," a govern-  
308 ment advisor for infectious diseases said.]

### 309 4.3 Model choice

310 We tested various approaches to generalize emo-  
311 tion annotations in our data, using LSTM and  
312 BERT models. LSTM models with FastText and  
313 Word2Vec embeddings tailored for German  
314 newspaper data showed poor performance in  
315 comparison to pre-trained BERT models, prompt-  
316 ing us to opt for the latter approach. BERT's deep

317 contextual embeddings and bidirectional pro-  
318 cessing offer nuanced semantic understanding,  
319 suitable for small datasets like ours (Delvin et al.,  
320 2019; Garí Soler & Apidianaki, 2021). We ex-  
321 plored two distinct BERT-based model configura-  
322 tions to investigate different approaches to emo-  
323 tion classification, a dual and a single-layer  
324 model. While the dual-layer BERT model aimed  
325 to capture both fine-grained Level-2 emotions and  
326 their superordinate Level-1 categories simultane-  
327 ously, the single-layer BERT model focused ex-  
328 clusively on Level-1 categories to examine  
329 whether a simpler architecture might achieve  
330 comparable results.

331 We split the data into training (75%), validation  
332 (15%), and testing (15%) sets, with performance  
333 evaluated using precision, recall, and Cohen's  
334 kappa, as well as Pearson correlations for emotion  
335 frequency relationships.

### 336 4.4 Dual-Layer BERT Model

337 For the dual-layer BERT model, we used the 'bert-  
338 base-german-cased' pre-trained model, which was  
339 trained on extensive German datasets, including  
340 Wikipedia (6 GB), OpenLegalData (2.4 GB), and  
341 news articles (3.6 GB), providing robust German  
342 language comprehension. The model's dual out-  
343 put layers served distinct purposes: one layer pre-  
344 dicted 33 Level-2 emotions, occurring at least 20  
345 times in the data, while the other layer focused on  
346 the nine Level-1 emotion families.

347 The dual-layer design treated all emotion vari-  
348 ables as statistically independent predictors, uti-  
349 lizing separate Dense layers with 33 and 10 neu-  
350 rons respectively for Level-1 and Level-2 emo-  
351 tions. Both layers employed the sigmoid activa-  
352 tion function, which operated independently for  
353 each neuron in the layer. This architecture allowed  
354 a single sentence to be associated with multiple  
355 emotions simultaneously, effectively capturing  
356 the complexity of emotional expression.

### 357 4.5 Single-Layer BERT Model

358 Our second configuration was a simpler model fo-  
359 cused on predicting only Level-1 emotions (emo-  
360 tion families). This model employed the same  
361 'bert-base-german-cased' pre-trained architecture  
362 and featured a single Dense layer with ten neurons  
363 corresponding to the Level-1 emotions. The  
364 model processed binarized input data, with Level-

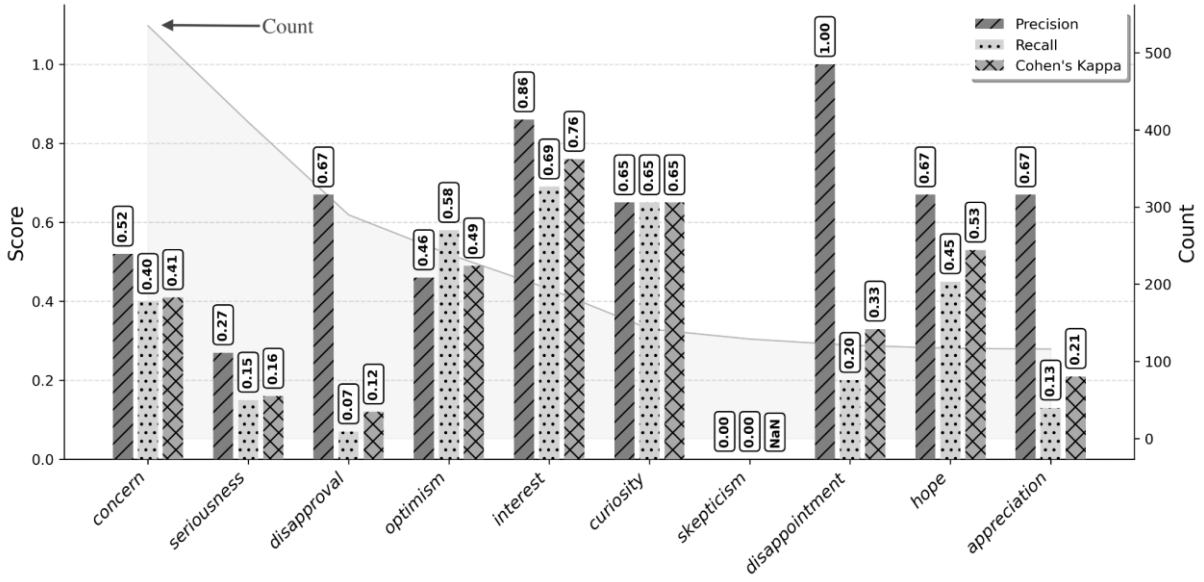


Figure 1: Dual-layer BERT model: precision, recall, Cohen’s kappa (bars), and count (line and gray shading in background) for the ten most frequent Level-2 emotions.

1 emotions represented by the maximum scores of  
 2 their associated Level-2 emotions. The use of the  
 3 sigmoid activation function for each neuron main-  
 4 tained the independent treatment of variables.

#### 4.6 Attempted Data Augmentation

5 Due to the uneven distribution of Level-2 emo-  
 6 tions (see Sec. 4.2) and the resulting data scarcity,  
 7 we explored options to augment the training data  
 8 synthetically. We opted for a strategy of translat-  
 9 ing all training sentences into another language  
 10 (both English and Russian) and back into German,  
 11 which was done with GPT-4o mini. To prevent  
 12 data leakage, sentences from the validation and  
 13 testing sets were excluded from the augmentation  
 14 process. The results showed slight precision and  
 15 recall differences post augmentation, but no clear  
 16 improvement, leading us to abandon this ap-  
 17 proach.

### 5. Results

#### 5.1 Dual-Layer BERT Model: Level 2

18 Although the dual-layer BERT model was trained  
 19 on 33 Level-2 emotions, robust metrics emerged  
 20 only for the 10 most frequent emotions due to data  
 21 sparsity in the remaining emotion classes (Fig. 1).  
 22 The model demonstrated the most consistent and  
 23 robust performance for *interest*, which showed the  
 24 highest level of agreement with human annotators

25 - a kappa of 0.76, along with strong precision and  
 26 recall. *Curiosity*, *hope*, and *optimism* formed a  
 27 second performance tier with balanced and mod-  
 28 erately high scores.

29 However, the model’s performance was notably  
 30 uneven. Emotions like *concern* and *seriousness*  
 31 showed moderate to low performance despite be-  
 32 ing the most frequent. *Disapproval*, *disappoint-*  
 33 *ment*, and *appreciation* revealed an interesting  
 34 pattern of high precision but low recall. This  
 35 asymmetry was particularly pronounced in *disap-*  
 36 *pointment*, which reached perfect precision while  
 37 capturing only a fifth of actual instances. The  
 38 most striking limitation was observed for *skepti-*  
 39 *cism*, which the model failed to identify entirely,  
 40 resulting in zero precision and recall and an unde-  
 41 fined Kappa score.

42 Overall, the model showed a conservative clas-  
 43 sification behavior with higher precision than re-  
 44 call values across emotions, and predominantly  
 45 low-to-moderate kappa scores (below 0.50). Inter-  
 46 estingly, correlation analysis revealed weak nega-  
 47 tive relationships between emotion frequency and  
 48 all performance metrics (precision:  $r = -0.25$ , re-  
 49 call:  $r = -0.03$ , kappa:  $r = -0.29$ ). However, none  
 50 of these relationships were statistically signifi-  
 51 cant. This indicates that emotion frequency was  
 52 not meaningfully associated with model perfor-  
 53 mance at Level 2.

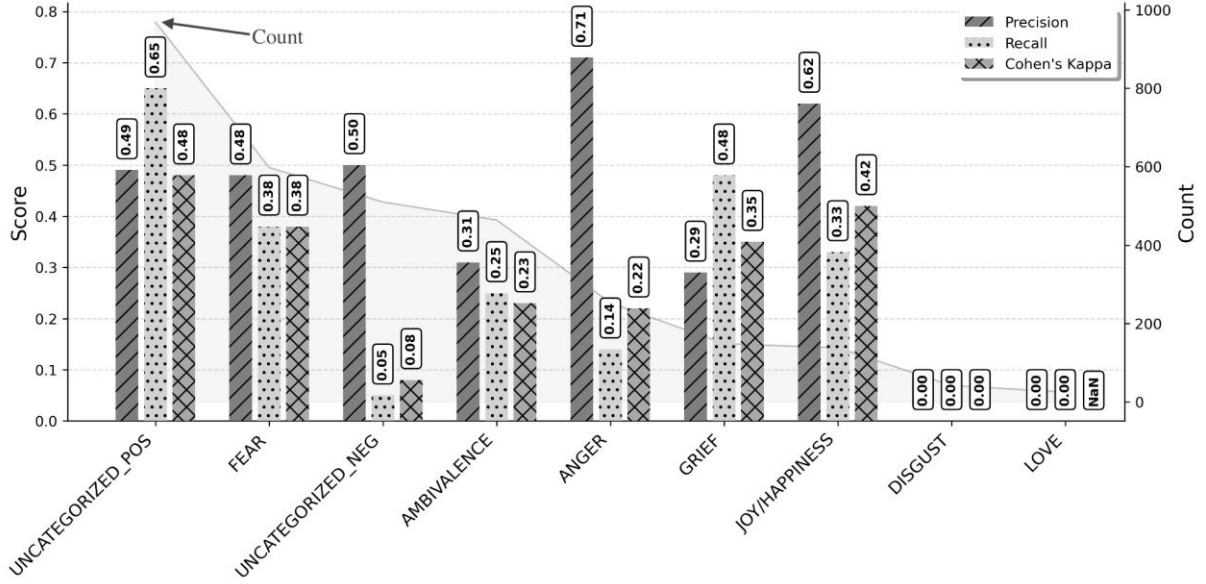


Figure 2: Dual-layer BERT model: precision, recall, Cohen’s kappa, and count for Level-1 emotions.

## 5.2 Dual-Layer BERT Model: Level 1

At Level 1, performance was variable across categories (Fig. 2). The most prevalent category, UNCATEGORIZED\_POS, achieved moderate performance with the highest kappa coefficient (0.65). FEAR, the second most frequent category, demonstrated comparable moderate performance levels. While ANGER had the highest precision (71%), it showed notably low recall (14%). This asymmetric pattern between precision and recall was also evident in UNCATEGORIZED\_NEG and JOY/HAPPINESS categories. Further, AMBIVALENCE consistently underperformed across all evaluation metrics. The model particularly struggled with DISGUST and LOVE categories, with LOVE showing zero recall and an undefined kappa value, indicating a complete absence of accurate predictions. In the Level-1 layer, the model generally prioritized accuracy over classification sensitivity, with UNCATEGORIZED\_POS and GRIEF being notable exceptions. Correlation analysis revealed moderate positive (but not statistically significant) associations between emotion frequency and performance metrics (precision:  $r = 0.45$ , recall:  $r = 0.62$ , kappa:  $r = 0.45$ ).

When compared to Level 2 (see Sec. 5.1), Level 1 demonstrated more uniform performance across categories. Both levels, however, exhibited conservative classification tendencies, manifesting in higher precision scores relative to recall. This effect was more pronounced in Level 2,

which displayed more substantial precision-recall disparities (e.g., an 80% gap for *disappointment*) compared to the Level 1 model (e.g., a 57% gap for ANGER).

A notable distinction emerged in the relationship between emotion frequency and performance metrics across the two levels. While Level 2 exhibited weak negative correlations with frequency across all metrics, Level 1 showed moderate positive correlations. This pattern suggests that the increased granularity of emotion categories in Level 2 may have introduced complexities that counteracted potential benefits from higher frequency counts. However, it is important to note that none of these relationships achieved statistical significance.

## 5.3 Single-Layer BERT Model: Level 1

We finally turn to the single-layer BERT model, which exhibited varying performance across emotion categories (Fig. 3). It achieved balanced metrics for UNCATEGORIZED\_POS, with moderate performance for FEAR. However, challenges arose with categories like AMBIVALENCE which had notably low recall (10%) and kappa (0.15) values, indicating substantial difficulties in classification accuracy. DISGUST and LOVE yielded zero or undefined values across all performance metrics.

Precision consistently exceeded recall, indicating reliably positive predictions but challenges in comprehensive emotion identification. Moderate-to-low kappa scores ranged from 0.15 to 0.52.

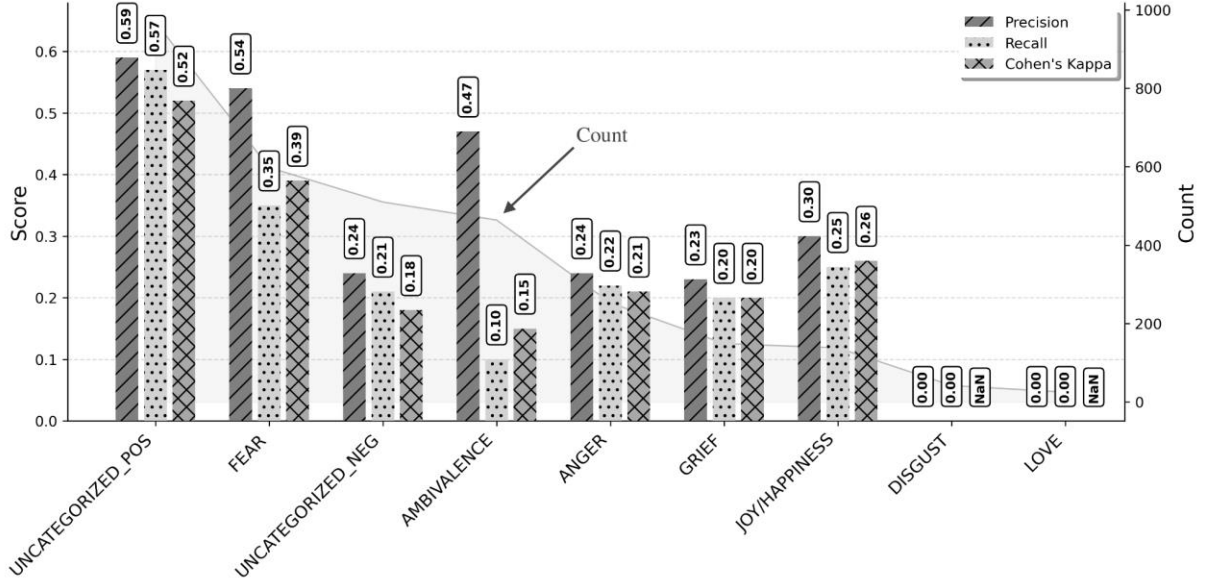


Figure 3: Single-layer BERT model: precision, recall, Cohen’s kappa, and count for Level-1 emotions.

Strong, statistically significant correlations were found between emotion frequency and both precision ( $r = 0.86$ ,  $p = 0.0$ ) and recall ( $r = 0.84$ ,  $p = 0.01$ ), suggesting frequency greatly influenced performance.

#### 5.4 Performance Comparison of Single-Layer and Dual-Layer BERT Models (Level 1)

Overall, the single- and dual-layer models revealed distinct classification patterns. Both effectively identified UNCATEGORIZED\_POS, but the dual-layer model prioritized sensitivity. In addition, both models struggled with rare emotions like DISGUST and LOVE. Yet, the dual-layer model showed advantages with mid-frequency emotions like ANGER and JOY/HAPPINESS, outperforming the single-layer model, particularly for GRIEF.

A key distinction emerged in precision-recall trade-offs. The dual-layer model exhibited more extreme variations, exemplified by ANGER (71% precision, 14% recall) and UNCATEGORIZED\_NEG (50% precision, 5% recall), while the single-layer model showed more moderate performance variations across categories. Although the single-layer model achieved slightly better agreement with human annotators for common emotions, the dual-layer model matched or exceeded these scores for medium-frequency emotions.

The single-layer model’s performance benefited more from higher frequency, shown by strong positive correlations with precision and recall, whereas the dual-layer model’s correlations were weaker and not statistically significant.

## 6. Discussion

This study set out to apply a fine-grained emotion classification system developed for literary texts to COVID-19 newspaper discourse. Overall, our findings indicate both promising advances and persistent challenges in adapting this framework to non-fiction texts.

### 6.1 Model Performance and Comparison

We observed that our dual-layer model excelled in distinguishing emotions such as ANGER and JOY/HAPPINESS, although it faced notable precision-recall trade-offs. Conversely, the single-layer model provided more consistent, balanced performance, effectively identifying UNCATEGORIZED\_POS and FEAR. These differences suggest that increased model complexity can enhance sensitivity for certain emotions while potentially compromising overall balance.

### 6.2 Emotion Frequency, Linguistic Markers, and Granularity

The relationship between emotion frequency and classification performance varied across model architectures. Notably, the single-layer model showed a statistically significant positive correlation between emotion frequency and performance, whereas the dual-layer model did not exhibit significant frequency-performance correlations. This discrepancy may indicate that in more granular classifications, the presence of clear linguistic markers is more influential than mere frequency.

Similar limitations have been observed by Demszyk et al. (2020), who analyzed 58,000 Reddit comments labeled with 27 emotions using a BERT-based model. They found that emotions with overt lexical markers (e.g., gratitude) were classified more successfully. Similarly, our findings suggest that in the dual-layer model, with its more granular classification, the presence of clear linguistic markers may be more important than mere frequency.

Comparing Level-1 and Level-2 classifications in the dual-layer model reveals insights into journalistic emotional granularity. Level-2 emotions showed a broader performance range and higher maximum kappa scores, but with greater variability. This result mirrors findings from Machová et al. (2023) suggesting that detecting multiple or weakly supported emotions remains a significant limitation in text-based emotion analysis.

### 6.3 Implications for Automated Emotion Detection

Both models encountered specific challenges, revealing key patterns in pandemic reporting. Both models struggled with AMBIVALENCE and *skepticism*, despite their potential relevance to COVID-19 news analysis. This finding underscores the challenge of detecting context-dependent emotions, which are often conveyed through subtle linguistic cues. This result aligns with Machová et al.'s (2023) findings, which highlighted the inherent challenge of modeling emotions that involve mixed or contradictory feelings as well as complex emotional expressions that depend on context or cultural understanding, including sarcasm, irony, idioms, and metaphors.

### 6.4 Generalization of Emotion Detection

Regarding generalization capability, both models exhibited conservative classification tendencies, resulting in higher precision but lower recall scores. This result suggests that while positive predictions generalize well to unseen data, many valid emotional instances may be overlooked.

In our study, 77 out of 87 original Level-2 categories were removed due to data sparsity, underscoring the challenges involved in fine-grained emotion classification. It is plausible to expect that a substantially larger dataset would have yielded more relevant data points, and thus better training data, for rare emotion classes. However, a Zipfian distribution of emotion categories may

be inherent to the expression of emotion in language, and partly also context-dependent – the most frequent Level-2 emotions in our COVID-pandemic discourse data, namely *concern*, *seriousness*, *disapproval*, are unlikely to be the most frequent emotions in, say, romantic novels. Ultimately, an important obstacle to the automated detection of emotion in language may be its context-dependent nature. This conclusion is further supported by the lack of a clear relationship between frequency and model performance across emotion categories – more training data does not necessarily imply better model performance, perhaps due to the expression of some emotions being more context-dependent than that of others.

Taken together, these findings have important implications for the analysis of news discourse. Tasks requiring high-precision classification of specific emotions (e.g., *interest* in vaccine development) seem to benefit from a dual-layer architecture, despite its performance variability. Conversely, applications demanding stable and balanced performance across emotion categories (e.g., monitoring broad societal sentiment during vaccination campaigns—UNCATEGORIZED\_POS, FEAR, UNCATEGORIZED\_NEG) may find the single-layer model more suitable, albeit at the expense of granularity.

## 7. Conclusion

This study represents a pioneering effort to apply an emotion classification scheme from computational literary studies to newspaper discourse covering the COVID-19 pandemic. By utilizing dual-layer and single-layer BERT models we demonstrated both the potential and challenges of automated emotion identification in journalistic discourse. Our findings suggest that emotion classification is complex, with varying performance across different emotional categories and model architectures.

While our current study focused on German newspaper articles, we plan to extend this methodology to social media data (such as Sailunaz et al. 2018) and other linguistic contexts, thus refining our understanding of emotional communication across different communicative domains. The proposed methodological innovations and insights provide a foundation for more advanced computational approaches to emotion analysis in text.



## 8. Limitations

This study faced several methodological and practical constraints that should be considered when interpreting the results. First, significant class imbalance in the dataset posed a major challenge, necessitating the exclusion of 77 out of 87 Level-2 emotion categories due to insufficient representation. This substantial reduction in emotional granularity, while methodologically necessary, limited our ability to fully evaluate the effectiveness of fine-grained emotion classification in COVID-related newspaper articles. Targeted data collection strategies, such as selectively sampling articles likely to contain underrepresented emotions, can improve the representation of these categories.

The annotation process presented several interconnected challenges. The task of emotion identification requires high emotional literacy and a solid understanding of emotional nuances from annotators. Despite providing comprehensive guidelines, the inherent complexity of emotion recognition, particularly in journalistic text where emotions may be subtly expressed or implied, likely contributed to annotation inconsistencies. This challenge was particularly apparent in the classification of ambivalent emotions. This limitation could be addressed by developing more precise, domain-specific annotation guidelines, and implementing a multi-stage annotation process with cross-validation among annotators.

In addition, the granularity of our classification scheme, while theoretically comprehensive, proved challenging to implement in practice. The attempt to distinguish between 87 different emotional categories may have been overly ambitious for the newspaper domain, where emotional expression tends to be more restrained and less varied than in literary texts or social media posts. This suggests that a more domain-appropriate classification system might be necessary for analyzing journalistic content.

Finally, our decision to conduct annotation at the sentence level, while practical for implementation, may have limited our ability to capture emotional content that develops across multiple sentences or requires broader context for proper interpretation. Emotions in news articles often emerge through extended narrative development and contextual framing that may be better captured within paragraphs rather than sentences.

These limitations point to several potential improvements for future research: collecting more

balanced datasets, developing more robust and domain-specific annotator training guidelines along with a multi-phase annotation approach, and potentially adjusting the granularity of emotion classification to better match the journalistic context. Additionally, exploring paragraph-level annotation might provide more insight into how emotions are conveyed in Covid-related news articles through extended context and narrative development.

## Acknowledgements

ChatGPTo3-mini was used while editing this paper. Further acknowledgements will be added once the paper is accepted.

## References

- Al-Baity HH, Alshahrani HJ, Nour MK, Yafaz A, Alghushairy O, Alsini R, Othman M. Computational Linguistics Based Emotion Detection and Classification Model on Social Networking Data. *Applied Sciences*. 2022; 12(19):9680. <https://doi.org/10.3390/app12199680>
- Acheampong, F. A., Wenyu, C., & Nunoo-Mensah, H. (2020). Text-based emotion detection: Advances, challenges, and opportunities. *Engineering Reports*, 2(7), e12189.
- Ahmad, K., Workshop on Emotion, M., & EMOT. (2011). *Affective Computing and Sentiment Analysis: Emotion, Metaphor and Terminology*. Springer Science+Business Media B.V. <https://doi.org/10.1007/978-94-007-1757-2>.
- Anz, T. (2007). Kulturtechniken der Emotionalisierung: Beobachtungen, Reflexionen und Vorschläge zur literaturwissenschaftlichen Gefühlsforschung. In K. Eibl, K. Mellmann, & R. Zymner (Eds.), *Im Rücken der Kulturen* (pp. 207–239). Paderborn: Schöningh.
- Argyle, M. (1996). *Körpersprache und Kommunikation*. Paderborn: Junfermann.
- Borst, J., Klähn, J., & Burghardt, M. (2023). Death of the dictionary? – The rise of zero-shot sentiment classification. In *Proceedings of the Computational Humanities Research Conference 2023* (pp. 303–319).
- Bruggmann, A., & Fabrikant, S. I. (2014). Spatializing a digital text archive about history. In K. Janowicz, B. Adams, G. McKenzie, & T. Kauppinen (Eds.), *Workshop on Geographic Information Observatories 2014: Proceedings* (GIO 2014 / GIScience: 8, pp. 6–14). Aachen: CEUR Workshop Proceedings.
- Damásio, A. R. (1997). *Descartes' Irrtum. Fühlen, Denken und das menschliche Gehirn*. München: List.
- Damásio, A. R. (2004). Emotions and feelings. A neurobiological perspective. In A. S. R. Manstead, N. Frijda & A. Fisher (Eds.), *Feelings and emotions. The Amsterdam Symposium* (pp. 49–57). Cambridge: Cambridge University Press.
- Demszky, D., Movshovitz-Attias, D., Ko, J., Cowen, A., Nemade, G., & Ravi, S. (2020). GoEmotions: A dataset of fine-grained emotions. *arXiv preprint arXiv:2005.00547*.
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)* (pp. 4171–4186). Minneapolis, Minnesota: Association for Computational Linguistics.
- Ekman, P. (1972). Universals and cultural differences in facial expression of emotion. In J. K. Cole (Eds.), *Nebraska Symposium on Motivation* (pp. 207–283). Lincoln: University of Nebraska Press.
- Ekman, P. (1988). *Gesichtsausdruck und Gefühl. 20 Jahre Forschung von Paul Ekman*. Paderborn: Junfermann.
- Ekman, P. (1994). The nature of emotions. Fundamental questions. New York: Oxford University Press.
- Flüh, M. (2019). „Sentimentanalyse“. In: *forTEXT. Literatur digital erforschen*. URL: <https://fortext.net/routinen/methoden/sentimentanalyse> [Access: 13. January 2025].
- Garí Soler, A., & Apidianaki, M. (2021). Let's play mono-poly: BERT can reveal words' polysemy level and partitionability into senses. *Transactions of the Association for Computational Linguistics*, 9, pp. 825–844. [https://doi.org/10.1162/tacl\\_a\\_00400](https://doi.org/10.1162/tacl_a_00400)
- Généreux, M., David, M. D., O'Sullivan, T., Carignan, M. È., Blouin-Genest, G., Champagne-Poirier, O., ... & Roy, M. (2021). Communication strategies and media discourses in the age of COVID-19: an urgent need for action. *Health Promotion International*, 36(4), 1178–1185.
- Holodyski, M. (2006). *Emotionen – Entwicklung und Regulation*. Heidelberg: Springer-Medizin-Verlag.
- Izard, C. E. (1992). Basic emotions, relations among emotions, and emotion-cognition relation. *Psychological Review*, 99, 561–565.
- Khurdula, H. V., Pagutharivu, A., & Yoo, J. S. (2024, March). The Future of Feelings: Leveraging Bi-LSTM, BERT with Attention, Palm V2 & Gemini Pro for Advanced Text-Based Emotion Detection. In *South-eastCon 2024* (pp. 275–278). IEEE.
- Kim, E., & Klinger, R. (2019). A survey on sentiment and emotion analysis for computational literary studies. *Zeitschrift für digitale Geisteswissenschaften*. [https://doi.org/10.17175/2019\\_008](https://doi.org/10.17175/2019_008)
- Klinger, R., Kim, E., & Padó, S. (2020). Emotion analysis for literary studies. In N. Reiter, A. Pichler, & J. Kuhn (Eds.), *Reflektierte algorithmische Textanalyse* (pp. 237–269). Berlin/Boston: De Gruyter.

- Konle, L., Jannidis, F., Kröncke, M., & Winko, S. (2022). Emotions and literary periods. In *DH Conference Abstracts*. Digital Humanities.
- Konle, L., Kröncke, M., Winko, S., & Jannidis, F. (2023). Connecting the dots: Variables of literary history and emotions in German-language poetry. *Journal of Computational Literary Studies*, 2\*(1), 1–22. <https://doi.org/10.48694/jcls.3604>
- Landweer, H., & Renz, U. (2008). Zur Geschichte philosophischer Emotionstheorien. In H. Landweer & U. Renz (Eds.), *Handbuch klassische Emotionstheorien* (pp. 1–18). Berlin, New York: de Gruyter.
- Lehmann, J., Mittelbach, M., & Schmeier, S. (2017). *Quantifizierung von Emotionswörtern in Texten*. DARIAH-DE Working Papers Nr. 24. Göttingen: DARIAH-DE. <https://doi.org/urn:nbn:de:gbv:7-dariah-2017-4-5>
- Lemor, A., & Montpetit, É. (2024). Exploring the role of uncertainty, emotions, and scientific discourse during the COVID-19 pandemic. *Policy and Society*, puae010.
- Liu, B. (2020). *Sentiment analysis: Mining opinions, sentiments, and emotions*. Cambridge University Press.
- Luca, M., Lopez, G., Longa, A., & Kaul, J. (2024). How are You Really Doing? Dig into the Wheel of Emotions with Large Language Models. In *2024 Artificial Intelligence for Business (AIxB)* (pp. 72–75). IEEE.
- Machová, K., Szabóová, M., Paralič, J., & Mičko, J. (2023). Detection of emotion by text analysis using machine learning. *Frontiers in Psychology*, 14. <https://doi.org/10.3389/fpsyg.2023.1190326>
- Maruf, A. A., Khanam, F., Haque, M. M., Jiyad, Z. M., Mridha, M. F., & Aung, Z. (2024). Challenges and opportunities of text-based emotion detection: A survey. *IEEE Access*, 12, pp. 18416–18450. <https://doi.org/10.1109/ACCESS.2024.3356357>
- Oatley, K., & Johnson-Laird, P. (1987). Towards a cognitive theory of emotions. *Cognition and Emotion*, 1, 29–50.
- Papadimitriou, O., Kanavos, A., Vonitsanos, G., Maragoudakis, M., Karkazis, P., & Mylonas, P. (2024, November). Enhancing Emotion Classification with a Hybrid BERT and CNN Architecture. In *2024 19th International Workshop on Semantic and Social Media Adaptation & Personalization (SMAP)* (pp. 156–161). IEEE.
- Peng, S., Cao, L., Zhou, Y., Ouyang, Z., Yang, A., Li, X., ... & Yu, S. (2022). A survey on deep learning for textual emotion analysis in social networks. *Digital Communications and Networks*, 8(5), 745–762.
- Perikos, I., & Hatzilygeroudis, I. (2016). Recognizing emotions in text using ensemble of classifiers. *Engineering Applications of Artificial Intelligence*, 51, 191–201. <https://doi.org/10.1016/j.engappai.2016.02.002>
- Plutchik, R. (1984). Emotions. A general psycho-evolutionary theory. In K. Scherer & P. Ekman (Eds.), *Approaches to emotion* (S. 197–200). Hillsdale: Erlbaum.
- Sailunaz, K., Özeyer, T., Rokne, J., & Alhajj, R. (2018). Text-based analysis of emotion by considering tweets. In T. Özeyer & R. Alhajj (Eds.), *Machine learning techniques for online social networks* (pp. 219–236). Cham: Springer.
- Scherer, K. (1993). On the nature and function of emotion: A component process approach. *Theorien und aktuelle Probleme der Emotionspsychologie*.
- Schmidt, T., Burghardt, M., & Wolff, C. (2018). Herausforderungen für Sentiment Analysis bei literarischen Texten. In M. Burghardt & C. Müller-Birn (Eds.), *INF-DH 2018* (pp. 1–9). Bonn: Gesellschaft für Informatik e.V. <https://doi.org/10.18420/infdh2018-16>
- Schmitz, U. (2016). 92. Kulturwissenschaftliche Orientierung in der Medienlinguistik. In L. Jäger, W. Holly, P. Krapp, S. Weber & S. Heekeren (Eds.), *Sprache - Kultur - Kommunikation / Language - Culture - Communication: Ein internationales Handbuch zu Linguistik als Kulturwissenschaft / An International Handbook of Linguistics as a Cultural Discipline* (pp. 901–908). Berlin, Boston: De Gruyter Mouton. <https://doi.org/10.1515/9783110224504-093>
- Schwarz-Friesel, M. (2007). *Sprache und Emotionen*. UTB.
- Singh, G., Brahma, D., Rai, P., & Modi, A. (2023). Text-based fine-grained emotion prediction. *IEEE Transactions on Affective Computing*.
- Stenvall, M. (2014). Presenting and representing emotions in news agency reports: On journalists' stance on affect vis-à-vis objectivity and factuality. *Critical Discourse Studies*, 11(4), 461–481.
- Storjohann, P., & Cimander, L. (2022). Approaching the Covid-19 discourse through neologisms in public communication. In M. Jakosz & M. Kałasznik (Eds.), *Corona-Pandemie: Diverse Zugänge zu einem aktuellen Superdiskurs* (pp. 25–51). Göttingen: V&R.
- Strapparava, C., & Valitutti, A. (2004). WordNet-Affect: An affective extension of WordNet. In M. T. Lino, M. F. Xavier, F. Ferreira, R. Costa, & R. Silva

909 (Eds.), *Proceedings of the 4th International Confer-*  
910 *ence on Language Resources and Evaluation (LREC:*  
911 *4, Vol. 4, pp. 1083–1086). Paris.*

912 Süselbeck, J. (2019). Sprache und emotionales  
913 Gedächtnis. Zur Konstruktion von Gefühlen und  
914 Erinnerungen in der Literatur und den Medien. In H.  
915 Kappelhoff, J. H. Bakels, H. Lehmann & C. Schmitt  
916 (Eds.), *Emotionen. Ein interdisziplinäres Handbuch*  
917 (pp. 282–295). Stuttgart: Metzler.

918 Taboada, M., Gillies, M. A., & McFetridge, P. (2006).  
919 Sentiment classification techniques for tracking liter-  
920 ary reputation. In *LREC workshop: Towards compu-*  
921 *tational models of literary analysis* (LREC: 5, pp. 36–  
922 43). Paris.

923 Winko, S. (2019). Literaturwissenschaftliche  
924 Emotionsforschung. In H. Kappelhoff, J.-H. Bakels,  
925 H. Lehmann, & C. Schmitt (Eds.), *Emotionen. Ein in-*  
926 *terdisziplinäres Handbuch* (pp. 397–407). Stuttgart:  
927 Metzler.

928 Zappettini, F., Ponton, D. M., & Larina, T. V. (2021).  
929 Emotionalisation of contemporary media discourse. A  
930 research agenda. *Russian Journal of Linguistics*,  
931 25(3), 587-595.

932 Zhunis, A., Lima, G., Song, H., Han, J., & Cha, M.  
933 (2022, April). Emotion bubbles: Emotional composi-  
934 tion of online discourse before and after the COVID-  
935 19 outbreak. In *Proceedings of the ACM Web Confer-*  
936 *ence 2022* (pp. 2603-2613).

937

938

## Appendix A: Emotion Classification Scheme

Category system for annotating emotions with nine supercategories and 87 subcategories. Within the subcategories, the emotion types are ordered from intense to less intense; positive emotions are red, negative emotions are blue, and emotions with an ambivalent valence are orange

**LOVE:** *affection, kindness, trust, intimacy, devotion, worship*

**JOY/HAPPINESS:** *contentment, pleasure, amusement, humor, (joyful) anticipation, enthusiasm, delight*

**DISGUST:** *weariness, reluctance, aversion, dislike, contempt*

**FEAR:** *concern, hesitancy, nervousness, creepiness, dread, terror, horror, consternation, panic*

**GRIEF:** *dejection, loneliness, sorrow, melancholy, despair, suffering*

**ANGER:** *disappointment, annoyance, indignation, resentment, rage, bitterness, hate*

**UNCATEGORIZED\_POS:** *curiosity, appreciation, admiration, hope, pride, self-confidence, material desire, relief, interest, serenity, empathy, friendliness, gratitude, optimism, schadenfreude, helpfulness*

**UNCATEGORIZED\_NEG:** *longing, masochism, confusion, aggression, nostalgia, impatience, disapproval, skepticism, greed/desire, perplexity, shame, remorse, jealousy, boredom, madness, compassion*

**AMBIVALENCE:** *courage, seriousness, astonishment, disregard, defiance, love-hate, being deeply moved, impulsiveness, reverence, humility, sadism, mockery, emotional coldness, vehemence*

**uncertain:** *[annotation]*