

Integrating Plutchik’s Theory with Mixture of Experts for Enhancing Emotion Classification

Anonymous ACL submission

Abstract

Emotion significantly influences human behavior and decision-making processes. We propose a labeling methodology grounded in Plutchik’s Wheel of Emotions theory for emotion classification. Furthermore, we employ a Mixture of Experts (MoE) architecture to evaluate the efficacy of this labeling approach, by identifying the specific emotions that each expert learns to classify. Experimental results reveal that our methodology improves the performance of emotion classification.

1 Introduction

Emotion is essential in human life, having influence on our thoughts, behaviors, and communication. Recognizing the paramount importance of emotions, researchers have made significant efforts to analyze and understand them (Picard, 1997). A particularly important area of this research is emotion recognition in text, as it forms a substantial part of our daily interactions, including email and Social Network Service (SNS).

While sentiment analysis, categorizing text as positive, negative, or neutral, has advanced significantly, recognizing the full spectrum of emotions in text—such as joy, anger, sadness, and fear—remains a challenging task. Mao et al. (2023) report that RoBERTa large with HG-F24 achieved 84.7% accuracy on sentiment analysis of Amazon product reviews but only 40.9% accuracy in emotion detection using a Twitter (X) dataset.

Previous research utilizing deep learning technology has demonstrated significant promise in extracting emotions from text (Yu et al., 2018; Baziotis et al., 2018; Ying et al., 2019; Li and Xiao, 2023; Alhuzali and Ananiadou, 2021). Recently, Chen et al. (2023) conducted a study analyzing the role of emotions in controversial Reddit comments using language models. He et al. (2024) systematically measured the affective alignment of language models (LMs) by comparing LM-generated responses

to SNSs on two socio-political issues. However, these studies face challenges like sampling bias and subjective annotation. For instance, Chai et al. (2024) note that existing multilabel text classification models lack the ability to generalize complex concepts. Ahanin et al. (2023) argue that current methods overlook the sentiment polarity of words.

To tackle the problems in emotion annotation, we introduce a new labeling approach. Our primary objective is to enhance the expressiveness of emotion labels by applying Plutchik’s Wheel of Emotions and Diagram of Emotion Dyads. Furthermore, we employ a Mixture of Experts (MoE) framework for emotion classification, which identifies the specific emotions that each expert in the model is best at classifying. This approach seeks to validate the improved classification performance and specialization of experts in distinct emotional categories.

The key contributions of this research are listed as follows:

- We propose a new emotion labeling method based on Plutchik’s wheel of emotions theory.
- We leverage MoE that is trained on basic emotions and learns to classify composite emotions effectively.
- We conducted experiments to show the efficacy of the proposed method. The results demonstrate that our approach can effectively improve the performance of emotion classification tasks, especially for emotions that are typically harder to classify with traditional methods.

The structure of the paper is organized as follows. Section 2 provides a review of related work. Section 3 outlines our approach. Section 4 details the experimental design. Section 5 discusses the results, and Section 6 provides an in-depth analysis. The final section concludes with future research.

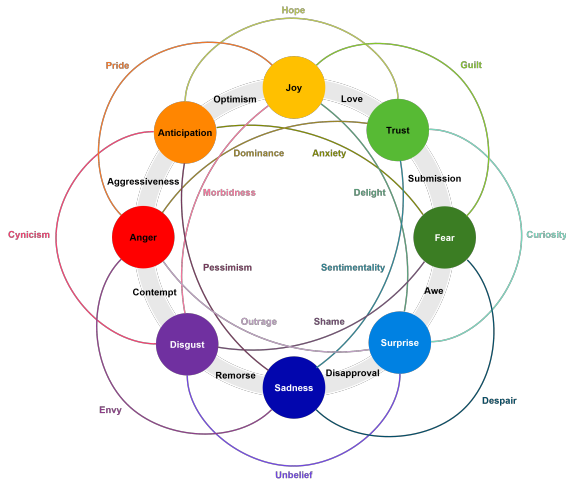


Figure 1: Plutchik’s Diagram of Emotion Dyads. Depicting the primary, secondary, and tertiary dyads formed by mixing the eight basic emotions.

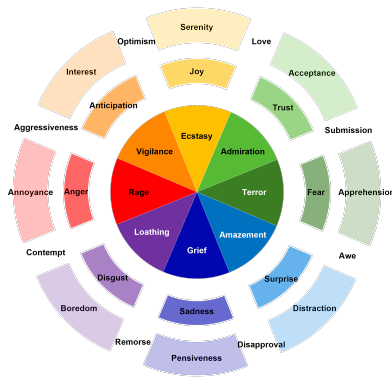


Figure 2: Plutchik’s Wheel of Emotions. The eight emotions are represented within the color spectrum, showing their mild and intense variations.

2 Related Work

2.1 Affective Computing

Emotions are physical and mental states induced by neurophysiological changes, often associated with specific thoughts, feelings, behavioral responses, and varying degrees of pleasure or displeasure (Damasio, 1998; Ekman and Davidson, 1994; Panksepp, 2004). They intertwine with mood, temperament, personality, disposition, and creativity (Averill, 1999). Recent research across psychology, medicine, history, sociology, and computer science highlights the complexity and importance of understanding emotions.

Despite extensive research, there is no universally accepted definition of emotion (Cabanac, 2002; Clore and Ortony, 2008). Emotions are categorized into various affects corresponding to specific situations (Barrett, 2006), and numerous theo-

ries have been proposed, each offering distinct perspectives on emotional experiences (James, 1884; Candland, 2003).

Ekman has significantly advanced our understanding of basic emotions through his research on facial expressions (Ekman, 1984). He identified six fundamental emotions: anger, disgust, fear, happiness, sadness, and surprise (Ekman, 1992a,b; Miller, 2016). Later, he expanded this list to include amusement, contempt, contentment, embarrassment, excitement, guilt, pride in achievement, relief, satisfaction, sensory pleasure, and shame, recognizing emotions not expressed solely through facial muscles (Ekman, 1999).

Our labeling method relies on Plutchik’s emotion theories (Plutchik, 2000, 1988), which define eight basic emotions, grouped as joy versus sadness; anger versus fear; trust versus disgust; and surprise versus anticipation. These basic emotions can combine to form complex emotions, as depicted in Figure 1; for instance, the complex emotion love is formed by joy and trust, while remorse is a mix of disgust and sadness. These complex emotions may arise from cultural conditioning or associations combined with the basic emotions. He further introduced twenty-four ‘Primary,’ ‘Secondary,’ and ‘Tertiary’ dyads, representing different emotion combinations, and noted that emotions can vary in intensity from mild to intense (Plutchik, 1991; Turner, 2000). As illustrated in Figure 2, annoyance, anger, and rage fall within the same category with different intensities.

2.2 Mixture of Expert

The Mixture of Experts (MoE) method divides complex problems into multiple sub-problems, using specialized models (i.e., experts) to address each sub-problem. MoE utilizes a gating network to combine the outputs of each expert model, selecting the most suitable expert for a given input. This approach is particularly useful for datasets with diverse characteristics, enhancing model performance and computational efficiency.

Eigen et al. (2013) introduced the idea of using multiple MoEs, each with its own gating network, as part of a deep model. This approach is more powerful since complex problems may contain many sub-problems, each requiring different experts. They also suggest that introducing sparsity could transform MoE into a tool for computational efficiency. Shazeer et al. (2017) proposed a new type of general-purpose neural network component:

a Sparsely-Gated Mixture-of-Experts Layer (MoE). This method uses Noisy top- k gating, which adds sparsity and noise to the Softmax Gate used in the MoE architecture (Jordan and Jacobs, 1994), selecting the top k values among the experts to produce the output. There are numerous other attempts to improve the gate network (Clark et al., 2022; Hazimeh et al., 2021; Zhou et al., 2022).

Lepikhin et al. (2020) replaced the Transformer Encoder’s FFN layer with MoE, distributing experts across devices. This had the drawback of slower speeds when computations concentrated on a single expert. Fedus et al. (2022) improved this by limiting each token to one expert ($k=1$) and restricting the number of tokens per expert. Jiang et al. (2024) used an MoE structure with Top- k Gating and SwiGLU as experts within the Mistral model’s Transformer block, improving performance across tasks and showing each expert specialized in specific tasks.

3 Method

This section describes our proposed method for emotion classification, utilizing the new labeling method based on Plutchik’s emotion theory and the implementation of the MoE structure in our model.

3.1 Plutchik Labeling

We redefine the dataset’s emotion labels for evaluation, based on the work of Plutchik (2000, 1988). Data labeled with our method are termed “Plutchik Labeling” and those without it as “Normal Labeling.” The Plutchik Labeling process follows the following rules:

- Labels corresponding to the eight basic emotions in Plutchik’s emotion theory were retained.
- Labels corresponding to primary, secondary, and tertiary dyads of the eight basic emotions were decomposed into their constituent emotions before labeling.
- Emotions that are combinations of opposite emotions were similarly decomposed into their constituent emotions before labeling.
- Mild and intense emotion labels were relabeled as the corresponding basic emotions.

While Plutchik’s emotion theory also hints at the existence of triads (Plutchik, 1991), these dataset

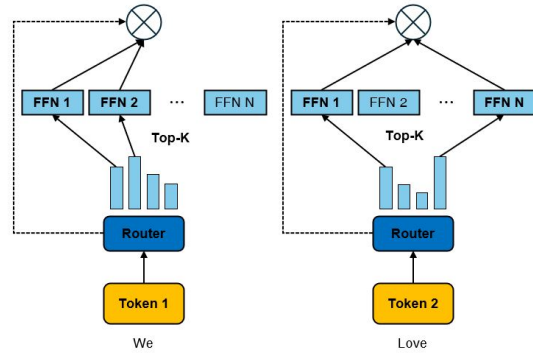


Figure 3: The Structure of Top- k MoE FFN.

Original Emot.	Augmented Emot.
Love	Joy, Trust
Optimism	Anticipation, Joy
Pessimism	Anticipation, Sadness

Table 1: Rules for relabeling compound emotions as the corresponding basic emotions in SemEval-2018.

did not provide sufficient detail on these emotions. Therefore, our study does not consider the triads, higher-order combinations, or the intensity of emotions.

3.2 Mixture of Emotion Expert

We aim to apply MoE to each model to determine whether each expert can be trained as a specialist in individual emotions. As mentioned earlier, there are several methods for gating that connect inputs to specific experts. Following the approach in Jiang et al. (2024), we selected the k most relevant experts for each token. The reason for experimenting with multiple values of k instead of fixing it is to account for complex emotions such as love and optimism, which are described as mixtures of several basic emotions according to Plutchik (2000, 1988). This consideration is crucial when tokens contain complex emotions. For the implementation of MoE, we refer to Mixtral (Jiang et al., 2024).

To compare how well the model understands emotions when MoE is applied, we used the existing FFN network as experts. To observe the performance changes with minimal parameter modifications, we replaced the FFN in the last transformer block of each model with an MoE structure.

4 Experiment

This section details the experimental design for evaluating the effectiveness of the proposed method

Original Emot.	Augmented Emot.
Admiration	Trust
Annoyance	Anger
Confusion	Anticipation, Surprise
Curiosity	Surprise, Trust
Disappointment	Sadness, Surprise
Disapproval	Sadness, Surprise
Excitement	Fear, Joy
Grief	Sadness
Love	Joy, Trust
Optimism	Anticipation, Joy
Pride	Anger, Joy
Remorse	Disgust, Sadness

Table 2: Rules for relabeling compound, mild, and intense emotions as the corresponding basic emotions in GoEmotions.

Emotion	train	valid	test
Anger	2544	315	1101
Anticipation	978	124	425
Disgust	2602	319	1099
Fear	1242	121	485
Joy	2477	400	1442
Love	700	132	516
Optimism	1984	307	1143
Pessimism	795	100	375
Sadness	2008	265	960
Surprise	361	35	170
Trust	357	43	153

Table 3: Emotion distribution across train, validation, and test sets for SemEval-2018 with Normal labeling.

Emotion	train	valid	test
Anger	2544	315	1101
Anticipation	3216	453	1688
Disgust	2602	319	1099
Fear	1242	121	485
Joy	2991	454	1669
Sadness	2266	292	1049
Surprise	361	35	170
Trust	975	161	621

Table 4: Emotion distribution across train, validation, and test sets for SemEval-2018 with Plutchik labeling.

Emotion	train	valid	test
Admiration	4130	488	504
Anger	1567	195	198
Annoyance	2470	303	320
Confusion	1368	152	153
Curiosity	2191	248	284
Disappointment	1269	163	151
Disapproval	2022	292	267
Disgust	793	97	123
Excitement	853	96	103
Fear	596	90	78
Grief	77	13	6
Joy	1452	172	161
Love	2086	252	238
Optimism	1581	209	186
Pride	111	15	16
Remorse	545	68	56
Sadness	1326	143	156
Surprise	1060	129	141

Table 5: Emotion distribution across train, validation, and test sets for GoEmotions with Normal labeling.

in multi-label emotion classification.

4.1 Experimental Setup

Our experiments utilize two transformer-based models, Llama-2(Touvron et al., 2023) and Mistral(Jiang et al., 2023), each with 7 billion parameters, chosen for their effectiveness across various domains. Their unmodified versions served as baselines for comparison. The models were accessed and utilized through the Hugging Face API. We fine-tuned the models using Q-LoRA(Dettmers et al., 2024). For all experiments, we used the same hyperparameters except for the k value. Performance was evaluated by averaging the results over five runs for each setting. Detailed hyperparameter configurations are provided in Section A.1.

Emotion	train	valid	test
Anger	3877	464	504
Anticipation	2944	360	336
Disgust	1334	164	179
Fear	1448	186	181
Joy	5801	707	669
Sadness	4928	643	607
Surprise	7472	944	951
Trust	8125	956	994

Table 6: Emotion distribution across train, validation, and test sets for GoEmotions with Plutchik labeling.

4.2 Labeling for Building Datasets

We chose the evaluation datasets based on the following criteria: 1) inclusion of all 8 basic emotions from Plutchik’s wheel, or 2) inclusion of emotions corresponding to Plutchik’s ‘Primary’, ‘Secondary’, and ‘Tertiary’ dyads, which, when decomposed, satisfy criterion 1. As a result, we selected SemEval-2018 (Mohammad et al., 2018) and GoEmotions (Demszky et al., 2020).

SemEval-2018 includes tweets, each labeled with one or more of 11 emotions or marked as Neutral. GoEmotions consists of 58K Reddit comments from 2005 to 2019, labeled with one or more of 27 emotions or Neutral. The rules for applying Plutchik labeling in these datasets are detailed in Table 1 and 2.

For a fair comparison, we excluded data for emotions not covered by Plutchik’s 8 basic emotions or their dyads, as well as Neutral, in all experiments. The final datasets are detailed in Tables 3, 4, 5, and 6. We fine-tuned the classification models using the train sets and evaluated their performance with the test sets.

5 Result

5.1 Main Result

Table 7 and 8 present the F1-scores of our proposed methods on two dataset. Table 7 shows the performance for different k values when applying MoE in Normal Labeling. For SemEval-2018, the macro-F1 indicates the model exceeds baseline performance at $k=2$, achieving the highest performance. In GoEmotions, the Mistral model surpasses the baseline across all k values, peaking at $k=4$, while the Llama2 model underperforms at all k values. The micro-F1 shows the highest performance at $k=4$ in all cases.

Overall, SemEval-2018 shows a consistent trend in macro-F1 changes with varying k values, unlike GoEmotions. This inconsistency, shown in Table 5, is due to significant label imbalance in GoEmotions. Elbayad et al. (2023) and Fedus et al. (2022) explain that MoE models tend to overfit on low-resource data, suggesting that the experts in the MoE model failed to learn effectively for certain emotions due to extreme imbalance. Additionally, the ‘grief’ and ‘pride’ have significantly fewer test samples, leading to high variance in performance metrics. Thus, performance comparisons using macro-F1 in GoEmotions may not be accurate.

Table 8 presents the performance using MoE

Top- k	SemEval-2018		GoEmotions	
	miF1	maF1	miF1	maF1
baseline	70.7	56.4	64.2	58.7
1	70.6	56.4	63.5	58.5
2	70.8	57.0	63.8	58.0
3	70.7	56.1	63.8	58.0
4	70.8	55.9	64.3	58.7
baseline	70.3	55.4	63.7	58.2
1	70.5	55.4	63.8	58.9
2	70.3	55.5	64.1	58.9
3	69.6	54.7	64.0	59.2
4	70.7	54.6	64.2	59.3

Table 7: F1 scores of the models with Normal Labeling. Upper: Llama2, Lower: Mistral

Top- k	SemEval-2018		GoEmotions	
	miF1	maF1	miF1	maF1
baseline	74.9	68.0	75.6	70.9
1	61.2	57.8	75.7	71.3
2	74.7	68.0	75.6	70.8
3	75.0	68.4	75.8	71.1
4	74.6	67.4	75.7	71.0
baseline	74.4	67.1	75.01	70.4
1	60.6	56.2	74.5	69.8
2	74.7	67.0	74.9	70.3
3	74.9	67.6	74.6	70.1
4	74.6	67.0	75.1	70.7

Table 8: F1 scores of the models with Plutchik Labeling. Top: Llama2, Bottom: Mistral.

with Plutchik Labeling varying the k values. With SemEval-2018, the highest macro-F1 was obtained at $k=3$, outperforming the baseline model. In GoEmotions, the Mistral model achieved the highest score at $k=4$, while the Llama2 model exceeded the baseline at $k=1$. The highest micro-F1 was generally obtained at $k=3$, except for the Mistral model on GoEmotions, which showed different patterns.

Plutchik Labeling resulted in more stable and superior performance than Normal Labeling, especially in GoEmotions, mitigating severe label imbalance. The MoE-trained model consistently outperformed the baseline model across various k values.

Figure 4 illustrates the changes in macro-F1 performance across both datasets with varying k values. When applying Plutchik Labeling, there is a significant improvement in performance compared to Normal Labeling in both the baseline and all MoE configurations. Especially, in SemEval-2018,

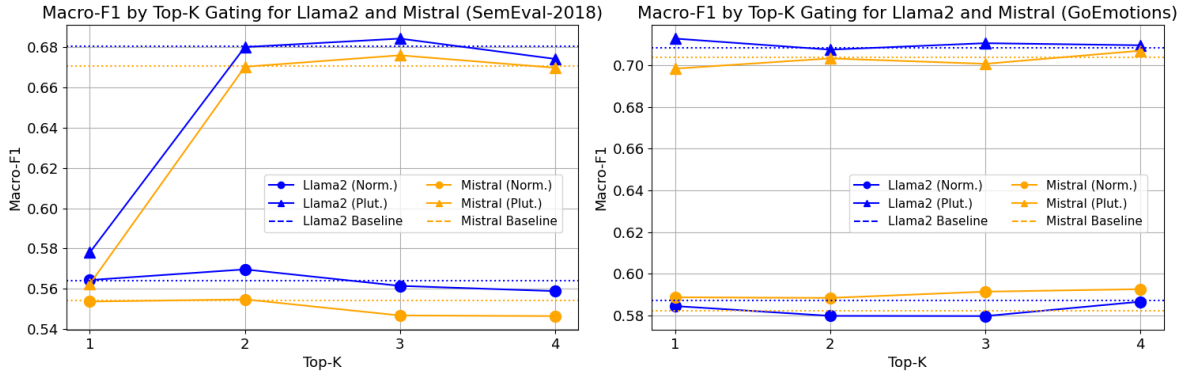


Figure 4: The macro-F1 scores of the MoE model across each datasets, k values, and labeling methods.

when k is set to 1, the performance improvement with Plutchik Labeling is less pronounced compared to the baseline and other k values. This suggests that selecting at least two or more experts in SemEval-2018 allows for better interpretation of emotions.

5.2 Underperforming Emotions

To evaluate the improving emotion classification performance by Plutchik Labeling, we investigated whether Plutchik Labeling could enhance the classification of emotions that were poorly classified under Normal Labeling. Specifically, we focused on Underperforming emotions, defined as those with an F1-score below 0.6 using the Normal Labeling dataset. This indicates that these emotions were challenging for the model to classify accurately.

Table 9¹ presents the F1-scores for Underperforming Emotions in SemEval-2018. When applying Plutchik Labeling, Pessimism is decomposed into Anticipation and Sadness, resulting in the removal of the Pessimism label. For Basic Emotions, both Anticipation and Trust showed significant improvement in classification performance due to data augmentation. However, in the case of Surprise, the transition from Normal Labeling to Plutchik Labeling did not benefit from data augmentation.

Table 10¹ presents the F1-scores for each Underperforming Emotions in GoEmotions. Basic emotions such as anger, disgust, and surprise, identified as Underperforming Emotions, demonstrated substantial improvement with the Plutchik Labeling.

¹AN: Anger, ANO: Annoyance, ANT: Anticipation, CO: Confusion, CUR: Curiosity, DIS: Disappointment, DAP: Disapproval, DIG: Disgust, EXC: Excitement, GRF: Grief, LO: Love, OPT: Optimism, PES: Pessimism, PRI: Pride, REM: Remorse, SUR: Surprise, TRU: Trust

Weak Emot.	Llama2		Mistral	
	Norm.	Plut.	Norm.	Plut.
ANT	24.0	66.8	24.3	69.4
PES	33.1	-	32.6	-
SUR	28.3	27.9	25.7	24.2
TRU	12.8	57.8	11.2	58.3
maF1	24.6	42.7	23.4	50.6

Table 9: F1-scores of Underperforming Emotions in SemEval-2018.

Most of the other Underperforming Emotions in GoEmotions are either complex emotions or mild or intense emotions, making direct comparisons with Plutchik Labeling challenging.

By comparing the macro-F1 scores of Underperforming Emotions between Normal Labeling and Plutchik Labeling in Tables 9 and 10, we observe a significant overall improvement in classification performance for Underperforming Emotions across both datasets. This enhancement indicates that our proposed method can effectively improve emotion classification tasks, especially for emotions that are typically harder to classify accurately. This demonstrates the potential of Plutchik Labeling to enhance the robustness and accuracy of emotion classification systems.

5.3 Complex Emotions

To assess if our MoE approach better classifies complex emotions, we compared the F1-scores of complex emotions between the baseline and MoE models under Normal Labeling.

Table 11¹ presents the classification performance of complex emotions in SemEval-2018, comparing the baseline with the Top-2 MoE models. The MoE approach resulted in a substantial improvement in macro-F1, notably increasing the performance for

Weak Emot.	Llama2		Mistral	
	Norm.	Plut.	Norm.	Plut.
AN	57.0	66.4	51.2	65.0
ANO	45.3	-	45.2	-
CO	57.7	-	58.0	-
DIS	32.0	-	35.6	-
DAP	57.9	-	57.5	-
DIG	48.9	56.8	46.1	56.8
EXC	47.8	-	50.0	-
GRF	29.5	-	29.4	-
PRI	43.9	-	42.2	-
SUR	60.8	77.5	58.3	76.5
maF1	48.1	66.9	47.4	66.1

Table 10: F1-scores of Underperforming Emotions in GoEmotions.

Comp Emot.	llama2		mistral	
	baseline	$k=2$	baseline	$k=2$
LO	62.4	61.8	59.0	60.8
OPT	70.7	71.7	71.0	72.4
PES	33.1	37.7	32.6	37.3
maF1	55.4	57.1	54.2	56.8

Table 11: F1-scores of complex emotions in SemEval-2018.

Pessimism, previously categorized as a Weak Emotion.

Table 12¹ shows the complex emotion classification performance of the baseline and Top-4 MoE models on GoEmotions. Based on macro-F1, Llama2 exhibited a slight increase in overall classification performance, while Mistral showed a slight decrease. Specifically, Llama2’s performance decreased for emotions such as confusion and pride, whereas Mistral saw decreases for confusion, curiosity, disappointment, disapproval, and pride. Pride, with insufficient data representation, poses a challenge for performance improvement due to data imbalance. According to Plutchik (1991), confusion, curiosity, disappointment, and disapproval share elements with surprise. Clore and Ortony (2013) explains that emotions like surprise are a neutral cognitive state that can be positive or negative, focusing neither on affect nor evaluation. Due to these characteristics, the MoE model likely struggled with classifying surprise and related complex emotions. Analysis of recall values revealed that applying MoE decreased recall for confusion, curiosity, disappointment, and disapproval, while recall for surprise increased.

Comp Emot.	llama2		mistral	
	baseline	$k=4$	baseline	$k=4$
CO	57.7	57.2	58.0	57.3
CUR	67.4	67.6	68.2	67.0
DIS	32.0	33.7	35.6	30.4
DAP	57.9	58.6	57.5	56.6
EXC	47.8	50.7	50.0	54.7
LO	83.3	83.9	84.2	85.6
OPT	68.7	70.3	69.8	69.9
PRI	43.9	38.2	42.2	41.9
REM	70.6	71.9	71.6	72.8
maF1	58.8	59.1	59.7	59.6

Table 12: F1-scores of complex emotions in GoEmotions.

6 Analysis

To clarify the relationships between emotions, an analysis was conducted comparing the predominant selections made by experts for each emotion. By tracking the output values of the Gate Layer in a MoE (Mixture of Experts) model, we identified which Experts were primarily selected for each emotion. Our approach involved selecting Experts for each token and aggregating the selection proportions of the Top- k Experts per token for each input. The value of k corresponds to the Top- k used in the MoE, and the sum of the selection proportions of the Top- k Experts per token equals 1. Subsequently, inputs were grouped by their labels (emotions), and the aggregate Expert selection proportions for each label were computed and standardized. Using the compiled frequency of Expert selections for each emotion, we plotted the emotion-emotion correlations to examine the relationships between different emotions.

The Figure 5a reveals that ‘joy’, ‘love’, and ‘optimism’ exhibit strong correlations, indicating that positive emotions are closely interconnected. In contrast, ‘anger’, ‘sadness’, and ‘disgust’ show strong positive correlations with each other, as well as with ‘fear’ and ‘pessimism’, forming a cluster of negative emotions. Additionally, ‘optimism’ and ‘pessimism’, as well as ‘love’ tends to show high correlations with ‘joy’ and ‘trust’, ‘optimism’ with ‘joy’, and ‘pessimism’ with ‘anticipation’ and ‘sadness’, allowing us to understand the similarities between complex emotions and their component basic emotions.

In Figure 5b, ‘joy’, ‘love’, ‘optimism’, and ‘admiration’ exhibit strong positive correlations, in-

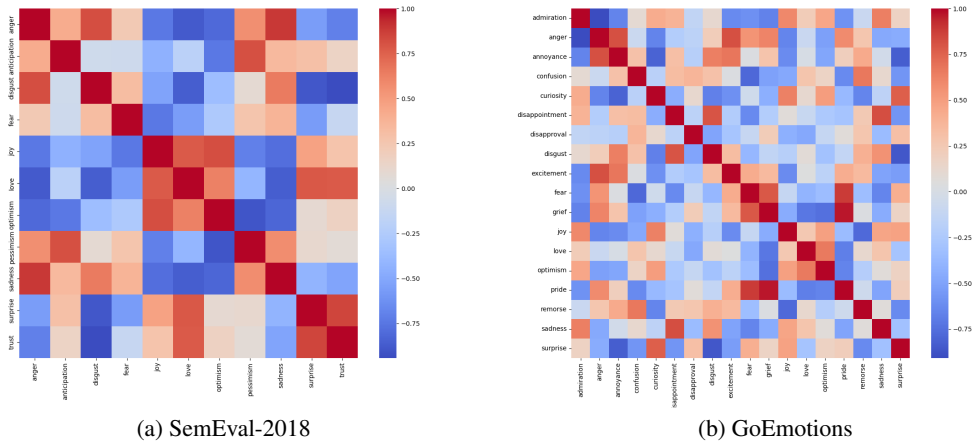


Figure 5: (a): Emotion correlations in Normal Labeling with Top-2 Gating. (b): Emotion correlations from in Normal Labeling with with Top-4 Gating.

dicating their close interrelation as positive emotions. Conversely, ‘anger’, ‘annoyance’, ‘excitement’, ‘fear’, ‘grief’, and ‘pride’ form a group of negative emotions, with ‘admiration’ and ‘anger’ showing a strong negative correlation, highlighting their opposing nature. Furthermore, the complex emotions ‘disappointment’ and ‘curiosity’ show high correlations with ‘sadness’ and ‘surprise’, respectively, while ‘anger’ correlates strongly with ‘annoyance’ and ‘sadness’ with ‘grief’. These patterns reveal the similarities between complex emotions and their component emotions, as well as the relationships between basic emotions and their mild or intense counterparts.

Overall, while the tendency to choose Experts for each emotion does not perfectly align with Plutchik’s emotion theory, the results show a significant degree of similarity. This suggests that our approach is valid for emotion analysis. These findings contribute to understanding the interrelations of emotions and can enhance the development of emotion prediction models.

7 Conclusion

Our approach is based on Plutchik’s emotion theory and the MoE architecture to enhance the performance of multi-label emotion classification tasks. The proposed methodologies were evaluated against baseline models, demonstrating significant improvements in classification performance. Notably, our approach excelled in accurately identifying emotions that were challenging to classify with traditional methods and showed superior performance in recognizing complex emotions.

Additionally, analyzing expert selection tendencies based on emotion correlations showed that our model’s behavior aligns closely with Plutchik’s emotion theory. This alignment enhances classification accuracy and provides a theoretically grounded understanding of emotional interactions. Our research presents a robust framework for multi-label emotion classification, integrating psychological theories and advanced machine learning techniques in emotion recognition tasks.

Limitation

This study acknowledges several limitations. First, utilizing Plutchik’s emotion theory requires the dataset to include all eight basic emotions defined by the theory, posing a challenge for datasets lacking these emotions. Furthermore, excluding emotions not covered by Plutchik’s emotion theory can be inefficient, making careful selection of datasets crucial. Future research could improve the labeling method by incorporating additional emotion models, such as the OCC model (Clore and Ortony, 2013).

Second, during the application of MoE, we encountered a known issue where tokens clustered around specific experts. This imbalance suggests the model may not fully leverage all experts. We plan to design a more sophisticated MoE structure to address this in the near future.

References

Zahra Ahanin, Maizatul Akmar Ismail, Narinderjit Singh Sawaran Singh, and Ammar AL-Ashmori.

491	2023. Hybrid feature extraction for multi-label emotion classification in english text messages . <i>Sustainability</i> , 15(16).	Antonio R Damasio. 1998. Emotion in the perspective of an integrated nervous system published on the world wide web on 27 january 1998.1. <i>Brain Research Reviews</i> , 26(2):83–86.	544 545 546 547
494	Hassan Alhuzali and Sophia Ananiadou. 2021. SpanEmo: Casting multi-label emotion classification as span-prediction . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 1573–1584, Online. Association for Computational Linguistics.	Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 2020. GoEmotions: A dataset of fine-grained emotions . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4040–4054, Online. Association for Computational Linguistics.	548 549 550 551 552 553 554
501	J R Averill. 1999. Individual differences in emotional creativity: structure and correlates. <i>J. Pers.</i> , 67(2):331–371.	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms . <i>Advances in Neural Information Processing Systems</i> , 36.	555 556 557 558
504	Lisa Feldman Barrett. 2006. Solving the emotion paradox: Categorization and the experience of emotion . <i>Personality and Social Psychology Review</i> , 10(1):20–46. PMID: 16430327.	David Eigen, Marc’ Aurelio Ranzato, and Ilya Sutskever. 2013. Learning factored representations in a deep mixture of experts . <i>arXiv preprint arXiv:1312.4314</i> .	559 560 561
508	Christos Baziotis, Athanasiou Nikolaos, Alexandra Chronopoulou, Athanasia Kolovou, Georgios Paraskevopoulos, Nikolaos Ellinas, Shrikanth Narayanan, and Alexandros Potamianos. 2018. Ntue-slp at semeval-2018 task 1: Predicting affective content in tweets with deep attentive rnns and transfer learning . In <i>Proceedings of The 12th International Workshop on Semantic Evaluation</i> . Association for Computational Linguistics.	Paul Ekman. 1984. Expression and the nature of emotion .	562 563
513		Paul Ekman. 1992a. Are there basic emotions? <i>Psychological review</i> , 99(3):550–553.	564 565
514		Paul Ekman. 1992b. An argument for basic emotions . <i>Cognition & Emotion</i> , 6:169–200.	566 567
515		Paul Ekman. 1999. <i>Basic Emotions</i> . John Wiley Sons, Ltd.	568 569
516		Paul Ekman and Richard J. Davidson, editors. 1994. <i>The Nature of Emotion: Fundamental Questions</i> . Oxford University Press USA.	570 571 572
517	Michel Cabanac. 2002. What is emotion? <i>Behavioural Processes</i> , 60(2):69–83.	Maha Elbayad, Anna Sun, and Shruti Bhosale. 2023. Fixing moe over-fitting on low-resource languages in multilingual machine translation . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 14237–14253.	573 574 575 576 577
518		William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity . <i>Journal of Machine Learning Research</i> , 23(120):1–39.	578 579 580 581
519	D. Candland. 2003. <i>Emotion</i> . Core books in psychology. Authors Choice Press.	Hussein Hazimeh, Zhe Zhao, Aakanksha Chowdhery, Maheswaran Sathiamoorthy, Yihua Chen, Rahul Mazumder, Lichan Hong, and Ed Chi. 2021. Dselectk: Differentiable selection in the mixture of experts with applications to multi-task learning . <i>Advances in Neural Information Processing Systems</i> , 34:29335–29347.	582 583 584 585 586 587 588
521	Yuyang Chai, Zhuang Li, Jiahui Liu, Lei Chen, Fei Li, Donghong Ji, and Chong Teng. 2024. Compositional generalization for multi-label text classification: A data-augmentation approach . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(16):17727–17735.	Zihao He, Siyi Guo, Ashwin Rao, and Kristina Lerman. 2024. Whose emotions and moral sentiments do language models reflect? <i>arXiv preprint arXiv:2402.11114</i> .	589 590 591 592
522		William James. 1884. II.—WHAT IS AN EMOTION ? <i>Mind</i> , os-IX(34):188–205.	593 594
523			
524			
525			
526			
527	Kai Chen, Zihao He, Rong-Ching Chang, Jonathan May, and Kristina Lerman. 2023. Anger breeds controversy: Analyzing controversy and emotions on reddit . In <i>Social, Cultural, and Behavioral Modeling</i> , pages 44–53, Cham. Springer Nature Switzerland.		
528			
529			
530			
531			
532	Aidan Clark, Diego de Las Casas, Aurelia Guy, Arthur Mensch, Michela Paganini, Jordan Hoffmann, Bogdan Damoc, Blake Hechtman, Trevor Cai, Sebastian Borgeaud, et al. 2022. Unified scaling laws for routed language models . In <i>International conference on machine learning</i> , pages 4057–4086. PMLR.		
533			
534			
535			
536			
537			
538	Gerald Clore and Andrew Ortony. 2008. Handbook of emotions . <i>Appraisal theories: How cognition shapes affect into emotion</i> , pages 628–642.		
539			
540			
541	Gerald L Clore and Andrew Ortony. 2013. Psychological construction in the OCC model of emotion . <i>Emot. Rev.</i> , 5(4):335–343.		
542			
543			

595	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .	646
596		647
597		648
598		649
599		650
		651
600	Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. <i>arXiv preprint arXiv:2401.04088</i> .	
601		
602		
603		
604		
605	Michael I. Jordan and Robert A. Jacobs. 1994. Hierarchical mixtures of experts and the em algorithm . <i>Neural Computation</i> , 6(2):181–214.	
606		
607		
608	Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. 2020. Gshard: Scaling giant models with conditional computation and automatic sharding. In <i>International Conference on Learning Representations</i> .	
609		
610		
611		
612		
613		
614	Jinfen Li and Lu Xiao. 2023. Multi-emotion recognition using multi-emobert and emotion analysis in fake news . page 128–135.	
615		
616		
617	Rui Mao, Qian Liu, Kai He, Wei Li, and Erik Cambria. 2023. The biases of pre-trained language models: An empirical study on prompt-based sentiment analysis and emotion detection . <i>IEEE Transactions on Affective Computing</i> , 14(3):1743–1753.	
618		
619		
620		
621		
622	Harold L. Miller. 2016. <i>The SAGE Encyclopedia of Theory in Psychology</i> .	
623		
624	Saif Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. SemEval-2018 task 1: Affect in tweets . In <i>Proceedings of the 12th International Workshop on Semantic Evaluation</i> , pages 1–17, New Orleans, Louisiana. Association for Computational Linguistics.	
625		
626		
627		
628		
629		
630	J. Panksepp. 2004. <i>Affective Neuroscience: The Foundations of Human and Animal Emotions</i> . Series in Affective Science. Oxford University Press.	
631		
632		
633	Rosalind W. Picard. 1997. <i>Affective computing</i> .	
634	R. Plutchik. 1991. <i>The Emotions</i> . University Press of America.	
635		
636	Robert Plutchik. 1988. <i>The Nature of Emotions: Clinical Implications</i> , pages 1–20. Springer US, Boston, MA.	
637		
638		
639	Robert Plutchik. 2000. Emotions in the practice of psychotherapy: Clinical implications of affect theories .	
640		
641	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. <i>arXiv preprint arXiv:1701.06538</i> .	
642		
643		
644		
645		
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	
	J. Turner. 2000. <i>On the Origins of Human Emotions: A Sociological Inquiry into the Evolution of Human Affect</i> . Stanford University Press.	652
		653
		654
	Wenhao Ying, Rong Xiang, and Qin Lu. 2019. Improving multi-label emotion classification by integrating both general and domain-specific knowledge. In <i>Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)</i> , pages 316–321.	655
		656
		657
		658
		659
	Jianfei Yu, Luís Marujo, Jing Jiang, Pradeep Karuturi, and William Brendel. 2018. Improving multi-label emotion classification via sentiment classification with dual attention transfer network . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1097–1102, Brussels, Belgium. Association for Computational Linguistics.	660
		661
		662
		663
		664
		665
		666
		667
	Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M Dai, Quoc V Le, James Laudon, et al. 2022. Mixture-of-experts with expert choice routing. <i>Advances in Neural Information Processing Systems</i> , 35:7103–7114.	668
		669
		670
		671
		672

Hyperparameter	Value
epoch	10
gradient_accumulation_steps	4
learning_rate	1e-4
warmup_ratio	0.1
max_grad_norm	0.3
weight_decay	0.001
batch_Size	8
quant_type	nf4
lora_r	8
lora_alpha	8
lora_dropout	0.1
num_expert	8

Table 13: Hyperparameter Settings for our experiments.

A Appendix

A.1 Hyperparameters

Table 13 shows the hyperparameter values applied to the models used in our experiments. Except for the K value, all hyperparameters were kept constant across all experiments. Each condition was tested five times.