



# Label-Aware Attention for Fine-Grained Emotion Detection

Yilin Wang, Ethan Wu, Ziqi Liu  
{yilinwan, yongyiw, ziqil2}@andrew.cmu.edu  
Carnegie Mellon University







# Label-Aware Attention for Fine-Grained Emotion Detection

Yilin Wang, Ethan Wu, Ziqi Liu  
 {yilinwan, yongyiw, ziqil2}@andrew.cmu.edu  
 Carnegie Mellon University



## Introduction

As an active NLP domain, emotion detection aims to determine the emotion(s) behind the plain text. GoEmotions dataset presents a multi-label classification challenge with more fine-grained emotion taxonomy than previous benchmarks. We propose a label-aware attention mechanism which significantly outperforms the baseline model.

## Task

GoEmotions consists of 58K English Reddit comments, each labelled with one or more of the 28 emotions, including neutral. The baseline classifier proposed by the author is trained by finetuning a BERT model on this dataset, which yields the macro F1-score 0.46.

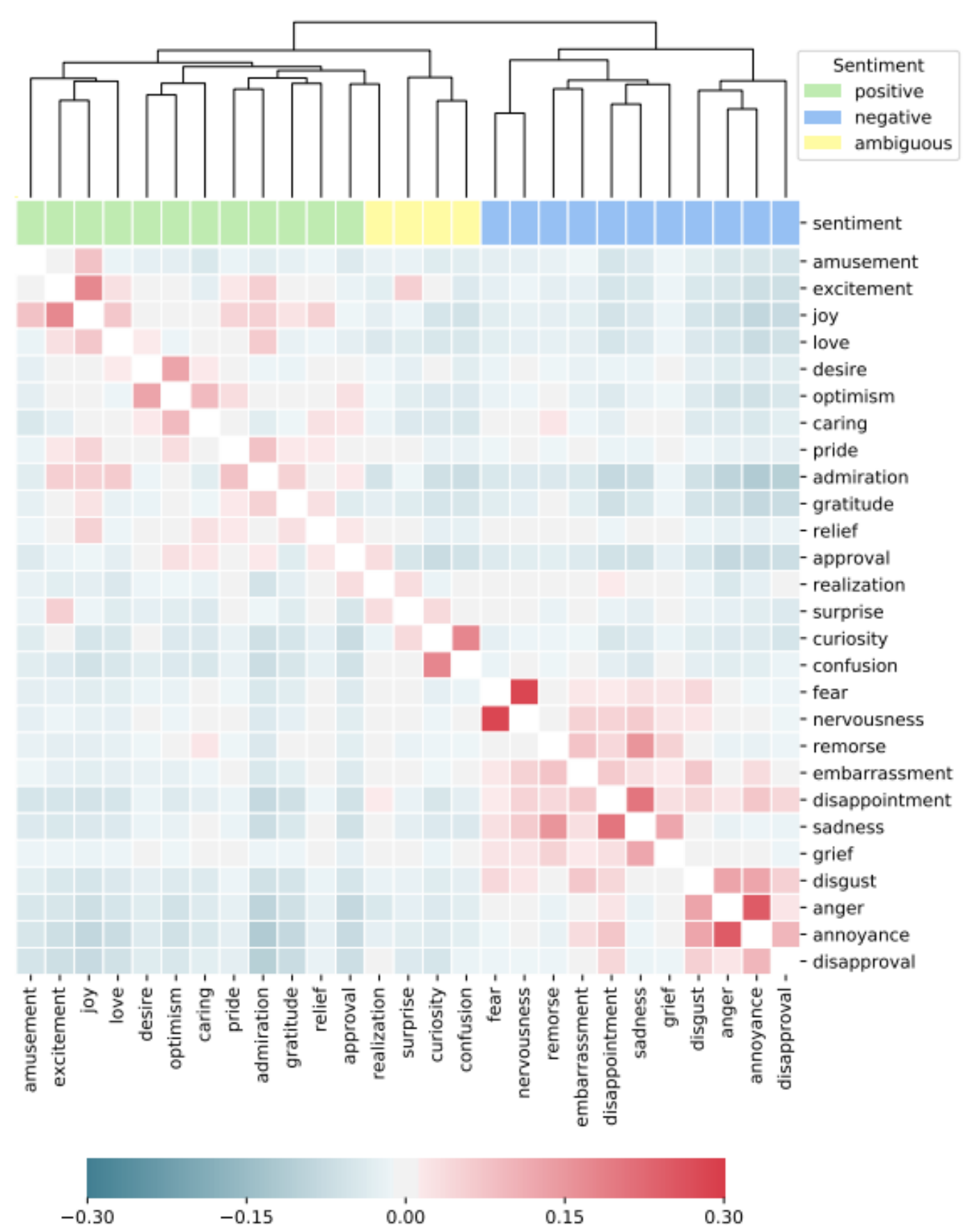


Figure 1: Label correlation

The challenge of the GoEmotions dataset comes from its fine-grained taxonomy. As illustrated in Figure 1, many emotions are in fact semantically and statistically correlated. After performing error analysis on the baseline model, we found that “neutral” is often confused with other emotions, and low-resourced emotions are often misclassified as high-resourced ones with similar sentiment.

Number of examples	58,009
Number of emotions	27 + neutral
Number of unique raters	82
Number of raters / example	3 or 5
Marked unclear or difficult to label	1.6%
Number of labels per example	1: 83% 2: 15% 3: 2% 4+: .2%
Number of examples w/ 2+ raters agreeing on at least 1 label	54,263 (94%)
Number of examples w/ 3+ raters agreeing on at least 1 label	17,763 (31%)

Table 1: Statistics of the GoEmotions dataset. We see that 83% of the samples have single label, and class “neutral” makes up 26% of all samples.

## Method

To mitigate the naively assumed independence between emotions (labels), we introduce an label-aware attention mechanism, where each emotion is enabled to focus on different aspects of a sentence, depending on the semantic representation of the emotion label.

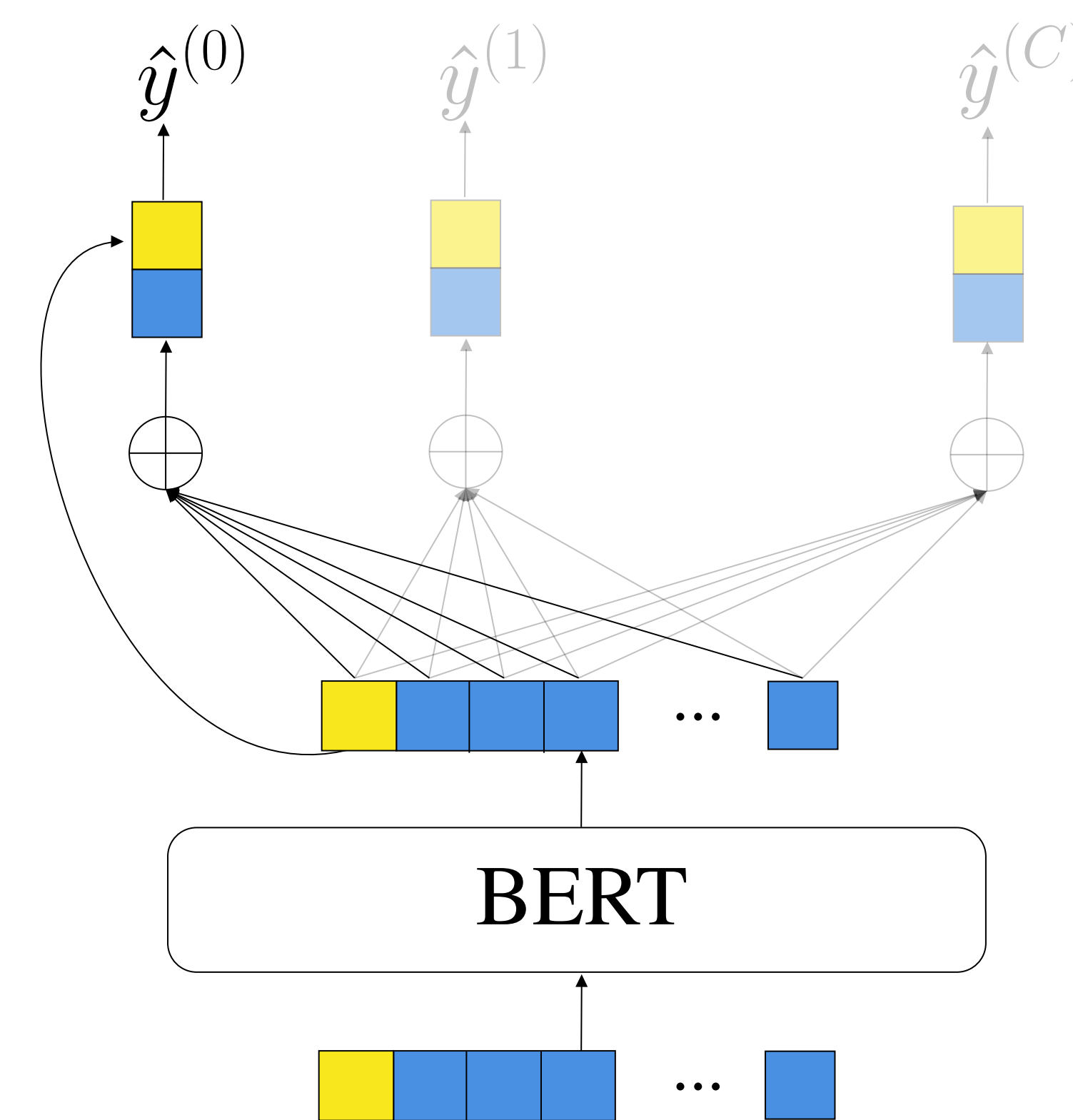


Figure 2: Visualization of the label-aware attention mechanism. The yellow block corresponds to the [CLS] token or its vector representation.

More formally, for an emotion  $i$  we create a trainable representation vector  $\mathbf{e}^{(i)} \in \mathbb{R}^{d'}$  and use it to query the  $d$ -dimensional final-layer hidden states  $\{\mathbf{h}_j\}_{j=1}^L$  from the pretrained language model, where  $L$  is the input sequence length. A softmax layer is applied to the bilinear product to compute the attention weights, which linearly combine contextualized embeddings:

$$\alpha_j^{(i)} = \frac{\exp(\mathbf{e}^{(i)T} \mathbf{W}_{\text{attn}} \mathbf{h}_j)}{\sum_{k=1}^L \exp(\mathbf{e}^{(i)T} \mathbf{W}_{\text{attn}} \mathbf{h}_k)} \quad \tilde{\mathbf{h}}^{(i)} = \sum_{k=1}^L \alpha_k^{(i)} \mathbf{h}_k$$

where  $\mathbf{W}_{\text{attn}} \in \mathbb{R}^{d' \times d}$  is a trainable transformation matrix and  $\tilde{\mathbf{h}}^{(i)}$  is a label-aware re embedding for emotion  $i$ . To predict the likelihood of emotion  $i$ , we concatenate the sentence embedding  $\mathbf{h}_0$  (which corresponds to the [CLS] token) and the label-aware embedding, and then pass the result to a dense layer with sigmoid activation:

$$\hat{y}^{(i)} = \sigma(\mathbf{w}_{\text{out}}^T [\mathbf{h}_0 \| \tilde{\mathbf{h}}^{(i)}] + b_{\text{out}})$$

## Results

Emotions	Baseline (F1)	Ours (F1)
Admiration	0.65	<b>0.68</b>
Amusement	<b>0.80</b>	<b>0.80</b>
Anger	0.47	<b>0.48</b>
Annoyance	0.34	<b>0.36</b>
Approval	0.36	<b>0.40</b>
Caring	0.39	<b>0.44</b>
Confusion	0.37	<b>0.42</b>
Curiosity	0.54	<b>0.55</b>
Desire	<b>0.49</b>	0.48
Disappointment	0.28	<b>0.33</b>
Disapproval	0.39	<b>0.41</b>
Disgust	0.45	<b>0.46</b>
Embarrassment	0.43	<b>0.45</b>
Excitement	0.34	<b>0.41</b>
Fear	0.60	<b>0.68</b>
Gratitude	0.86	<b>0.90</b>
Grief	0.00	<b>0.35</b>
Joy	0.51	<b>0.59</b>
Love	0.78	<b>0.80</b>
Nervousness	<b>0.35</b>	0.34
Neutral	<b>0.68</b>	0.66
Optimism	0.51	<b>0.55</b>
Pride	0.36	<b>0.50</b>
Realization	0.21	<b>0.25</b>
Relief	0.15	<b>0.35</b>
Remorse	0.66	<b>0.67</b>
Sadness	0.49	<b>0.53</b>
Surprise	0.50	<b>0.55</b>
Macro Avg.	0.46	<b>0.51</b>
Std	0.19	<b>0.01</b>

Table 2: F1 scores across all emotions in GoEmotions dataset, in comparison with baseline. Our model achieves significant improvement on Macro-average F1 and consistent improvements across all emotion groups. Moreover, the model shows much lower variance under different seeds, an indication of robust design.