# An Emoji-aware Multitask Framework for Multimodal Sarcasm Detection

**Anonymous ACL submission**

## Abstract

Sarcasm is a case of implicit emotion and needs additional information like context and multi-modality for its better detection. But sometimes this additional information also fails to help in sarcasm detection. For example, the utterance "Oh yes, you've been so helpful. Thank you so much for all your help", said in a polite tone with a smiling face, can be understood easily as non-sarcastic because of its positive sentiment. But, if the above message is accompanied with a frustrated emoji 😫, the negative sentiment of emoji becomes evident and the intended sarcasm can be easily understood. Thus, in this paper, we propose the *SEEmoji* MUStARD, an extension of the multimodal MUStARD dataset. We annotate each utterance with relevant emoji, emoji's sentiment and emoji's emotion. We propose an emoji-aware multitask deep learning framework for multimodal sarcasm detection (i.e. primary task), and sentiment and emotion detection (i.e. secondary task) in a multimodal conversational scenario. Experimental results on the *SEEmoji* MUStARD show the efficacy of our proposed approach for sarcasm detection over the state-of-the-art.

## 1 Introduction

We know that sarcasm is implicit, we can also agree that sometimes just going through the utterance text is not enough to understand sarcasm. For example, the utterance (only text) "It's just a privilege to watch your mind at work" is positive in nature and if it is intended in a sarcastic manner, its next to impossible to understand it. If this utterance is multimodal in nature and is accompanied with a video of the facial expressions and the tone of the speaker, it can be easily understood that the utterance is sarcastic (Chauhan et al., 2020).

Emojis are a trending topic these days because they provide an expressive way to convey sentiment and emotion. They are also a convenient way of understanding the implicit sentiment and emotion of the utterance. As sarcasm is closely related with the understanding of implicit sentiment/emotion, we can hypothesize that emojis should help to understand if there is any intended sarcasm in the utterance or not.

Even though sarcasm is related with sentiment and emotion, sarcasm detection is very challenging and that is why everyone treats this task separately. But if we introduce emojis then somewhat sarcasm becomes easy to compare before. The main contributions and/or attributes of our proposed research are as follows: **a)** We propose the *SEEmoji* MUStARD, an extension of the multimodal MUStARD dataset (Chauhan et al., 2020). We manually annotate each utterance with relevant emoji, emoji's sentiment and emoji's emotion; **b)** We propose an emoji-aware multitask framework for multimodal sarcasm detection. In our multitask framework, sarcasm detection is treated as the primary task, whereas emotion and sentiment analysis are considered as auxiliary tasks; **c)** We propose a Gated Multimodal Attention mechanism for sarcasm detection; and **d)** We present the state-of-the-art systems for sarcasm detection in multimodal scenario.

## 2 Dataset

The MUStARD (Castro et al., 2019; Chauhan et al., 2020) dataset consists of conversational audio-visual utterances (total of 3.68 hours in length). The samples were gathered from four famous TV shows *viz.*, Buddies, The Big Bang Theory, The Golden Girls, and Sarcasmaholics Anonymous and annotated manually. This dataset has 690 samples, and each sample utterance (u) consists of its context (c) and multiple labels i.e., sarcasm ($S^r$), implicit sentiment ($I_s$), implicit emotion ($I_e$), explicit sentiment ($E_s$) and explicit emotion ($E_e$).

We have further annotated the MUStARD (Chauhan et al., 2020) dataset with extra information in the form of emojis ($E^m$), emoji's sentiment ($E_s^m$), and emoji's emotion ($E_e^m$). We use 25 dif-

ferent and most frequently used emojis on social media which represent different emotion as well as the sentiment. We take three sentiment values, *viz. positive, negative* or *neutral* for emoji's sentiment and nine emotion values, *viz.* anger (An), excited (Ex), fear (Fr), sad (Sd), surprised (Sp), frustrated (Fs), happy (Hp), neutral (Neu) and disgust (Dg). for emoji's emotion. We show some samples from the dataset and emojis in Table 1.

Please note that, the motivation behind using emoji's sentiment and emotion information is to capture the relationship between sentiment and emotion of emojis and multimodal data.

| No. | Utterances | $E^m$ | $E_s^m$ | $E_e^m$ |
|-----|-----------|-------|---------|---------|
| 1 | *It's just a privilege to watch your mind at work.* | 🤢 | Neg | Dg |
| 2 | *To feed the cat Rose.* | 🙂 | Pos | Hp |
| 3 | *You're kidding, right?* | 😮 | Pos | Sp |
| **Emojis used for annotation** | | | | |
| 😠😡✨😂😭😱😳😔😒😤😶😀😥🧑😩😊😳😄❤️😃😆😄 😑😪😡🙏 | | | | |

Table 1: Samples From *SEEmoji* MUStARD

**Annotation Details:** We employ three graduate students highly proficient in the English language with prior experience in labeling emoji, sentiment and emotion. The guidelines for annotation, along with some examples, were explained to the annotators before starting the annotation process (c.f. Table 1). The annotators were given data without sarcasm labels and asked to annotate every utterance with one emoji and corresponding sentiment and emotion (only one emotion per utterance) of that emoji. A majority voting scheme was used for selecting the final emotion and sentiment. We achieve an overall Fleiss' (Fleiss, 1971) kappa score of 0.82, which is considered to be reliable.

## 3 Methodology

In this section, we describe our proposed methodology, where we aim to leverage the emoji information for solving the problem of multimodal sarcasm detection in a multitask framework[1]. We propose a multitask deep learning framework for sarcasm detection (primary task), and sentiment and emotions detection(secondary tasks) in a multimodal conversational scenario. We depict the overall architecture in Figure 1.

**Input Features:** The raw *utterance level* multimodal features are represented as text $T_u \in \mathbb{R}^{w \times 300}$ (*fastText* word embeddings (Joulin et al., 2016)) where w stands for number of words in an

---

[1] We shall make the datasets and codes available.

utterance, visual $V_u \in \mathbb{R}^{2048}$, acoustic $A_u \in \mathbb{R}^{283}$ and $E_u^m \in \mathbb{R}^{300}$ (*emoji2vec* emoji embeddings (Eisner et al., 2016)). Please note that we use same features for acoustic and visual modality and take average of the acoustic and visual features across the utterances for a fair comparison with the state-of-the-art systems. We show the detailed description of input features in appendix.

**Model description:** We first pass the $T_u$ through bi-directional Gated Recurrent Unit (Cho et al., 2014) ($BiGRU$) to learn the contextual relationship between the words, then pass though a dense layer (**BiGRU**+Dense as shown in Figure 1). Simultaneously, we pass $A_u$ and $V_u$ through the dense layer separately. Then we concatenate all the modalities together and pass through another dense layer to obtain multimodal representation (MR). Finally, we apply softmax layer to predict the $I_s$ & $E_s$ and sigmoid layer for $I_e$ & $E_e$. Similarly, we take *emoji embedding* as input and pass through a dense layer. Then, we apply softmax layer to predict the $E_s^m$ and $E_e^m$ and concatenate them (c.f. Figure 1) to obtain emoji representation (ER). The objective of ER is to enhance the information of the emoji's behaviour or nature.

As we know, emoji helps sarcasm but we do not know how to fused emoji with multimodal data. So, we take every possible combination where emoji can help. We first obtain *Emoji-aware Implicit Multimodal* Representation (EIMR) to capture the relationship between emoji and *implicit* behaviour of multimodal data by concatenating ER with MR, $I_s$ and $I_e$. Then, we obtain *Emoji-aware Explicit Multimodal* Representation (EEMR) to capture the relation between emoji and *explicit* behaviour of multimodal data by concatenating ER with MR, $E_s$ and $E_e$. Finally, we obtain *Emoji-aware Multimodal* Representation (EMR) to capture the relation between emoji and multimodal data, without including implicit/explicit sentiment and emotion information, by concatenating ER with MR.

**Gated Multimodal Attention:** We propose a Gated Multimodal Attention (GMA) mechanism. We first employ a gated architecture (Gated Multimodal Unit (GMU) (Arevalo et al., 2017)) to refine (or filter out the noise) an input representation (ER/EMR/EIMR/EEMR) *w.r.t.* all the participating input representations (ER, EMR, EIMR, and EEMR). After this, an attention mechanism is applied on the output of the gated architecture to decide which gated multimodal representation is
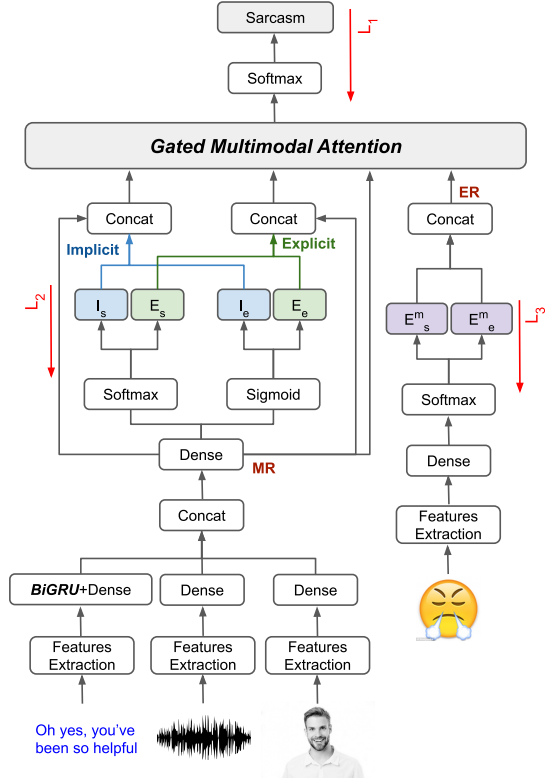
Figure 1: Overall architecture of the proposed emoji-aware multimodal sarcasm detection framework

contributing the most in sarcasm detection. This process is denoted by GMA.

Motivated by the residual skip connection (He et al., 2016), the outputs of GMA concatenated with the representations ER, EIMR, EEMR, and EMR. Finally, the concatenated representation is passed through an softmax layer for sarcasm detection. The gradients are updated based on three losses i.e., sarcasm ($loss_1$ or $L_1$), emoji's sentiment and emotion ($L_2$) and implicit/explicit sentiment and implicit/explicit emotion ($L_3$).

## 4 Experimental Results and Analysis

**Experimental Setup:** We evaluate our proposed model on the *SEEmoji* MUStARD. We perform our all experiments based on two setups i.e., Speaker Dependent and Speaker Independent. We do not take context and speaker information into consideration which is same as *utterance w/o context and w/o speaker* in (Castro et al., 2019; Chauhan et al., 2020). The detailed description of experimental setup is in appendix.

We implement our proposed model on the Python-based PyTorch deep learning library. As the evaluation metric, we employ precision (P), recall (R), and F1-score (F1) for implicit sentiment/emotion, explicit sentiment/emotion, emoji's sentiment, emoji's emotion and sarcasm detection. We use *Adam* as an optimizer, *Softmax* as a classifier for implicit/explicit sentiment, emoji' sentiment, emoji's emotion, and sarcasm detection, and the *categorical cross-entropy* as a loss function. For implicit/explicit emotion recognition, we use *Sigmoid* as an activation function and optimize the *binary cross-entropy* as the loss.

**Experimental Results:** In this section, we show the comparison between our proposed model and baselines i.e., Baseline-1 (Castro et al., 2019) and Baseline-2 (Chauhan et al., 2020) which also made use of the same dataset. We evaluate our proposed architecture with all the possible input combinations i.e., unimodal *(T, A, V)*, bimodal *(T+V, T+A, A+V)* and trimodal *(T+V+A)*. The results are shown in Table 2. For both the setups, we observe similar trend of performance improvement of our proposed model *(T+V+A)* over Baseline-1 (5.2 points ↑ and 7.0 points ↑ in F1-score) and Baseline-2 (4.1 points ↑ and 3.9 points ↑ in F1-score). Thus, we observe that emoji is helpful in improving the performance of sarcasm detection. For both the setups, we also observe that trimodal performs better than the unimodal and bimodal.

**Ablation Study:** To understand the effect of *Emoji* and proposed *GMA*, we perform an ablation study on our proposed model. The results are shown in Table 3. For both the setups, we observe that proposed model outperformed *Proposed w/o Emoji* (2.9 points ↑ and 2.7 points ↑ in F1-score) and *proposed w/o GMA* (2.0 points ↑ and 2.1 points ↑ in F1-score).

**Impact of Emoji:** Empirically, we have shown that emoji helps sarcasm (C.f. Table 3). We take some examples from the dataset (c.f. Table 4), which are sarcastic, to show the effect of emojis. Each example has positive implicit/explicit sentiment. The predictions made by the model, proposed w/o emoji, are incorrect for sarcasm but correct for implicit and explicit sentiment.

Now, when emojis are used, the model, correctly predicts all the examples as sarcastic. We observe that emoji's sentiment is playing an important role for sarcasm detection. The sentiment displayed by the emojis is negative. This helps the model to understand the contrast between the sentiments displayed by the utterance and the emoji. Thus, it correctly interprets that the utterance is sarcastic.

| | Speaker Dependent | | | | | | | | | Speaker Independent | | | | | | | | |
| | Proposed | | | Baseline-1 (2019) | | | Baseline-2 (2020) | | | Proposed | | | Baseline-1 (2019) | | | Baseline-2 (2020) | | |
| Labels | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| T | 69.9 | 69.7 | 69.6 | 65.1 | 64.6 | 64.6 | - | - | - | 63.1 | 61.2 | 62.0 | 60.9 | 59.6 | 59.8 | - | - | - |
| A | 69.1 | 68.0 | 67.4 | 65.9 | 64.6 | 64.6 | - | - | - | 67.2 | 67.4 | 67.3 | 65.1 | 62.6 | 62.7 | - | - | - |
| V | 75.1 | 74.2 | 74.0 | 68.1 | 67.4 | 67.4 | - | - | - | 65.4 | 65.7 | 65.5 | 54.9 | 53.4 | 53.6 | - | - | - |
| T+V | 75.1 | 74.8 | 74.7 | 72.0 | 71.6 | 71.6 | 72.7 | 71.9 | 71.6 | 66.2 | 66.6 | 66.2 | 62.2 | 61.5 | 61.7 | 65.5 | 65.5 | 65.7 |
| T+A | 70.6 | 70.3 | 70.1 | 66.6 | 66.2 | 66.2 | 62.2 | 61.1 | 59.6 | 69.5 | 66.0 | 65.9 | 64.7 | 62.9 | 63.1 | 59.1 | 60.0 | 50.3 |
| A+V | 76.1 | 75.7 | 75.6 | 66.2 | 65.7 | 65.7 | 72.7 | 71.9 | 71.8 | 68.9 | 69.1 | 68.2 | 64.1 | 61.8 | 61.9 | 65.6 | 63.8 | 63.9 |
| T+V+A | 77.9 | 76.9 | 76.7 | 71.9 | 71.4 | 71.5 | 73.4 | 72.7 | 72.6 | 70.0 | 69.7 | 69.8 | 64.3 | 62.6 | 62.8 | 69.5 | 66.0 | 65.9 |

Table 2: Comparative analysis between our proposed model, and Baseline-1 and Baseline-2

| Setup | Speaker Dependent | | | Speaker Independent | | |
| | P | R | F1 | P | R | F1 |
|---|---|---|---|---|---|---|
| Proposed w/o Emoji | 74.2 | 73.4 | 73.8 | 67.8 | 66.5 | 67.1 |
| Proposed w/o GMA | 74.4 | 74.9 | 74.7 | 67.7 | 67.9 | 67.7 |
| Proposed | 77.9 | 76.9 | 76.7 | 70.0 | 69.7 | 69.8 |

Table 3: Ablation study

| | Utterances | W/o Emoji | W/ Emoji | |
| | | $S^r$ | $E^m$ | $S^r$ |
|---|---|---|---|---|
| 1 | Oh, I'm so glad you asked it like that. You. | NS | 💩 | S |
| 2 | We can? Ok I am trying that. | NS | 😠 | S |
| 3 | Wow you look just like your son, Mrs. Tribbiani | NS | 😃 | S |

Table 4: Comparison between w/ Emoji and w/o Emoji

**Impact of GMA:** Empirically, we have shown the effectiveness of the GMA (C.f. Table 3). We show the heatmap of an utterance "Oh, I'm so glad you asked it like that. You." (c.f. Table 4) and we have already shown that emoji 💩 help to predict this utterance as sarcastic. To prove this, we show the attention heatmap for this utterance in Figure 2. We see that ER contributing more than others which means emoji is more evident for this utterance to predict correctly. Thus, this also proves our hypothesis that emoji help sarcasm.



Figure 2: The heatmaps represent attention weights of a particular utterance across ER, EMR, EIMR and EEMR.

**Error Analysis:** We perform error analysis for our proposed model. We take some samples which are incorrectly predicted by our proposed model and analyze our model's shortcomings. We take two utterances i) *"Yes you can. You're thinking about time, you can't go back in time."* and given label is not sarcastic (NS) with emoji 😐 and ii) *"I thought if I littered, that crying Indian might come by and save us."* and given label is sarcastic with

emoji 😐. For both utterances, the implicit/explicit sentiment of the utterances is positive and the emoji is 😐 (expressionless). Even though, the information for sentiment and emoji types are same for both but one utterance is non-sarcastic while the second utterance is sarcastic. With this, the model fails to learn the subtle difference between the utterances as the emojis do not provide any additional distinguishable information to the model about the utterances during training.

## 5 Conclusion and Future Work

In this paper, we have created *SEEmoji* MUStARD by manually annotating an existing MUStARD dataset with emoji, emoji's sentiment and emotion labels. In our multitask framework, sarcasm is treated as the primary task, whereas emotion and sentiment analysis are considered as secondary tasks. We have proposed a Gated Multimodal Attention based emoji-aware-multitask learning framework for sarcasm prediction. Empirical results of our proposed model, on the newly annotated dataset, achieve state-of-the-art performance over the existing methods.

During the annotation, we found that the dataset is very small for a complex architecture to learn a complex problem like sarcasm. We think that increasing the size of the dataset by annotating more samples should be helpful to gain improvement in performance.

## 6 Ethical Declaration

The dataset used in this paper is freely available and we extend the dataset by annotating (Emoji, Emoji's sentiment, and Emoji's emotion) the dataset, and has been used only for the purpose of academic research. The annotation for extending the dataset was done by human experts, who are the regular employee of our research group. There are no other issues to declare.

## References

John Arevalo, Thamar Solorio, Manuel Montes-y-Gómez, and Fabio A González. 2017. Gated multi-modal units for information fusion. *arXiv preprint arXiv:1702.01992*.

Santiago Castro, Devamanyu Hazarika, Verónica Pérez-Rosas, Roger Zimmermann, Rada Mihalcea, and Soujanya Poria. 2019. Towards multimodal sarcasm detection (an _obviously_ perfect paper). *arXiv preprint arXiv:1906.01815*.

Dushyant Singh Chauhan, Dhanush S R, Asif Ekbal, and Pushpak Bhattacharyya. 2020. Sentiment and emotion help sarcasm? a multi-task learning framework for multi-modal sarcasm, sentiment and emotion analysis. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4351–4360, Online. Association for Computational Linguistics.

KyungHyun Cho, Bart van Merrienboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder-decoder approaches. *CoRR*, abs/1409.1259.

Ben Eisner, Tim Rocktäschel, Isabelle Augenstein, Matko Bošnjak, and Sebastian Riedel. 2016. emoji2vec: Learning emoji representations from their description. *arXiv preprint arXiv:1609.08359*.

Joseph L. Fleiss. 1971. Measuring nominal scale agreement among many rater. *Psychological Bulletin*, 76:378–382.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.

Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Hérve Jégou, and Tomas Mikolov. 2016. Fasttext.zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651*.

## A Dataset

### A.1 The role of emoji's sentiment and emoji's emotion

We know that sentiment and emotion of multimodal data helps in better sarcasm detection. To compliment the sentiment and emotion information, which plays a significant role in sarcasm detection, of multimodal data, we also use emoji's sentiment and emotion information. The idea is to capture the relation between sentiment & emotion of emojis and multimodal data, and the combined effect they have in better understanding of sarcasm.

## B Input Features

**Text Features:** Let us assume, in an utterance, there are $n_t$ number of words $w_{1:n_t} = w_1, ..., w_{n_t}$, where $w_j \in \mathbb{R}^{300}$. Each word, $w_j$, is represented as a vector using *fastText* word embeddings.

**Visual Features:** Let us assume that the number of visual frames for an utterance be $n_v$. We take the average of all frames to extract the utterance level information for the visual modality where $V_u \in \mathbb{R}^{2048}$.

**Acoustic Features:** Given $n_a$ number of frames for the acoustic *w.r.t.* an utterance, we take the average of all the frames to extract the utterance level information where $A_u \in \mathbb{R}^{283}$.

**Emoji:** There is one emoji (say $E^m$) associated with each utterance. The pre-trained emoji embeddings are obtained using *emoji2vec* where $E_u^m \in \mathbb{R}^{300}$.

Please note that we take average of the acoustic and visual features across the utterances for a fair comparison with the state-of-the-art.

## C Experimental Setup

We evaluate our proposed model on the *SEEmoji* MUStARD. We perform our all experiments based on two setups i.e., Speaker Dependent Setup and Speaker Independent Setup. We do not take context and speaker information into consideration (*utterance w/o context and w/o speaker*). We perform *grid search* to obtain the optimal hyper-parameters (c.f. Table 5). Though our aim is to use a generic hyper-parameter configuration for all our experiments. There are two setups which are as follows;

**Speaker Dependent Setup:** In this setup, five-fold cross-validation was performed for the experiments, where each fold takes samples randomly in a stratified manner from all the TV shows.

**Speaker Independent Setup:** In this experiment, samples from three TV shows (i.e., The Golden Girls, Big Bang Theory, and Sarcasmaholics Anonymous) were taken in the training set while samples from the fourth TV show (i.e., Friends) were taken in the test set. Following this step, we were able to reduce the effect of the speaker in the model.

### C.1 Computational Budget

We use GPUs[2] for all experiments. Our model only take approx 1.5GB GPU memory. It takes 2-3 seconds per epoch approximately.

---

[2]GPU: 1080Ti with 32GB, RAM: 256GB

| Parameters | Speaker Dependent | Speaker Independent |
|---|---|---|
| Bi-GRU | $2 \times 300$N | |
| Dense layer | 300N, D=0.3 | |
| Activations | *ReLu* | |
| Optimizer | *Adam (lr=0.001)* | |
| Outputs | *Softmax* $(I_s, E_s, E_s^m, E_e^m, S^r)$ & *Sigmoid* $(E_s, E_e)$ | |
| Loss | *Categorical cross-entropy* $(I_s, E_s, E_s^m, E_e^m, S^r)$ *Binary cross-entropy* $(E_s, E_e)$ | |
| Epochs | *200* | |
| Batch | 64 | 16 |

Table 5: Model configurations

## D   Description of Emojis

We use 25 frequently used emojis on social medial and the detailed description of emojis are as follows;

### D.1   Anger

: A yellow face with a frowning mouth and eyes and eyebrows scrunched downward in anger.
: A red face with an angry expression: frowning mouth with eyes and eyebrows scrunched downward. Bears the same expression as  Angry Face on most platforms and may convey more intense degrees of anger, e.g., hate or rage.

### D.2   Excited

: The glittering flashes of sparkles. Generally depicted as a cluster of three, yellow four-point stars, with one large sparkle and two small ones to its left or right.Commonly used to indicate various positive sentiments, including love, happiness, beauty, gratitude, and excitement. May also be used to convey newness or cleanliness.
: A yellow face smiling with open hands, as if giving a hug. May be used to offer thanks and support, show love and care, or express warm, positive feelings more generally. Due to its hand gesture, often used to represent jazz hands, indicating such feelings as excitement, enthusiasm, or a sense of flourish or accomplishment.

### D.3   Fear

: A face with small, open eyes, open frown, raised eyebrows, and a pale blue forehead, as if experiencing a cold flash.
: A yellow face screaming in fear, depicted by wide, white eyes, a long, open mouth, hands pressed on cheeks, and a pale blue forehead, as if it has lost its color. Its expression evokes Edvard Munch's iconic painting The Scream.

### D.4   Sad

: A yellow face with raised eyebrows and a slight frown, shedding a single, blue tear from one eye down its cheek. May convey a moderate degree of sadness or pain,
: A pensive, remorseful face. Saddened by life. Quietly considering where things all went wrong. Depicted as a yellow face with sad, closed eyes, furrowed eyebrows, and a slight, flat mouth. May convey a variety of sad emotions, including feeling disappointed, hurt, or lonely.  Less intense than other sad emojis like  Loudly Crying Face and more introspective.
: A yellow face with an open mouth wailing and streams of heavy tears flowing from closed eyes. May convey inconsolable grief but also other intense feelings, such as uncontrollable laughter, pride or overwhelming joy.
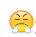
### D.5   Surprised

: A yellow face with small, open eyes, raised eyebrows, and a small, open mouth, as if it has been hushed by concern or correction.  Meaning widely varies, but its expression is commonly taken as surprise, embarrassment, or mild excitement.
: A yellow face with small, open eyes and a large, round mouth, slack with surprise or shock, as if saying Wow! or Oh my! May convey such feelings as awe or disbelief, often milder or more ironic in tone than  Face Screaming in Fear.

### D.6   Frustrated

: A hand shown pressing against the head of a person, commonly written as facepalm. Used to display frustration or embarrassment at the ineptitude of a person or situation.May be used in a similar context to the acronym SMH (shaking my head), or in relation to the Picard Facepalm meme.
: A yellow face with closed eyes, furrowed eyebrows, broad frown, and two puffs of steam blowing out of its nose, as if in a huff or fuming. May convey various negative emotions, including irritation, anger, and contempt. May also convey feelings of pride, dominance, and empowerment.

### D.7   Happy

: A yellow face with a big grin and scrunched, X-shaped eyes, tilted on its side as if rolling on the floor laughing (the internet acronym ROFL). Sheds two tears and tilts right on most platforms. Often

6

conveys hysterical laughter more intense than 😂 Face With Tears of Joy.

😘: A yellow face winking with puckered lips blowing a kiss, depicted as a small, red heart. May represent a kiss goodbye or good night and convey feelings of love and affection more generally.

😊: A yellow face with smiling eyes and a broad, closed smile turning up to rosy cheeks. Often expresses genuine happiness and warm, positive feelings.

😂: Emoji Meaning A yellow face with a big grin, uplifted eyebrows, and smiling eyes, each shedding a tear from laughing so hard.

❤️: A classic red love heart emoji, used for expressions of love and romance. This is the most popular heart emoji A similar emoji exists for the heart suit in a deck of playing cards.

😍: A yellow face with an open smile, sometimes showing teeth, and red, cartoon-styled hearts for eyes. Often conveys enthusiastic feelings of love, infatuation, and adoration, e.g., I love/am in love with this person or thing.

😀: A yellow face with smiling eyes and full-toothed grin, as if saying Cheese! for the camera. Teeth may be smoothed-over or crosshatched. Often expresses a radiant, gratified happiness. Tone varies, including warm, silly, amused, or proud.

## D.8 Neutral

😐: A yellow face with simple, open eyes and a flat, closed mouth. Intended to depict a neutral sentiment but often used to convey mild irritation and concern or a deadpan sense of humor.

😑: A yellow face with flat, closed eyes and mouth. May convey a sense of frustration or annoyance more intense than suggested by 😐 Neutral Face, as if taking a moment to collect itself.

## D.9 Disgust

🤢: A sickly-green face with concerned eyes and puffed, often red cheeks, as if holding back vomit. May represent physical illness or general disgust.

🤮: A yellow face with scrunched, X-shaped eyes spewing bright-green vomit. May represent physical illness or disgust, more intensely than 🤢 Nauseated Face.

💩: A swirl of brown poop, shaped like soft-serve ice cream with large, excited eyes and a big, friendly smile. May be used to represent feces and other bathroom topics as well as stand in for their many related slang terms. It also enjoys a wide range of idiosyncratic applications, such as conveying a sense of whimsy or silliness, given its fun, happy expression.