

# How Good Is GPT’s “Emojinal Intelligence”? Investigating Emoji Patterns in LLM-Generated Social Media Text

Michael Dunn and Kenneth Hopkinson

Air Force Institute of Technology, Wright-Patterson Air Force Base, USA

[michael.dunn@afit.edu](mailto:michael.dunn@afit.edu)

[kenneth.hopkinson@afit.edu](mailto:kenneth.hopkinson@afit.edu)

**Abstract:** Recent advancement in Large Language Models (LLMs) has opened the prospect of generating text for social media content that mimics human writing. The misuse of these tools presents urgent dilemmas, motivating the need to better understand the structure and patterns of LLM-generated content. Human communication on the Internet has developed relevant linguistic adaptations, including the use of emoji to augment traditional text. This study investigates the ability of one LLM, OpenAI’s GPT-3.5, to replicate human emoji usage in social media contexts. Drawing upon a dataset of nearly three thousand US English human-written tweets, we employed GPT-3.5-Turbo to generate social-media-style content and analyzed the use of emoji in the resulting text. We compared the patterns of emoji usage between the LLM-generated and the human-written datasets, particularly frequency, types of emoji commonly used, and emoji sequences (n-grams). Our results revealed notable differences in all categories. While human-written tweets were more likely to use faces, hearts, and repetitive sequences of emoji, LLM-created content had a broader variety of emoji, with a preference for literal representations of the text’s subject matter, producing diverse and unique emoji combinations.

**Keywords:** Large language models, Natural language processing, Computational linguistics, Social media, Emoji

---

## 1. Introduction

The recent advances in natural language generation by Large Language Models (LLMs) has enabled these models to produce text that closely resembles human-authored content. These LLMs, such as OpenAI’s GPT-3.5 and GPT-4, are notable for both the quality and flexibility of their output, able to generate written compositions in a wide variety of styles and on nearly any topic as prompted by the user. This capability has and will continue to enable bad actors to use LLMs to speed up and potentially even automate detrimental content on social media, including spam, disinformation, and trolling. To combat such misuse, researchers and national security professionals need a better understanding of how LLM-generated content may appear in the online information environment.

Human language adapts remarkably well to changes in technology and communication media, and Internet-based communications platforms have been no different (McCulloch, 2019a). One language adaptation driven by the Internet, and especially social media, has been the invention and widespread adoption of emoji. These pictograms enable Internet users to add extra dimensions to their informal written communications ordinarily not possible with traditional text alone. Despite being quite a recent innovation, the use of emoji has already developed distinct and pervasive usage patterns. These are not rules to be followed but patterns that have emerged organically from humans experimenting, observing, sharing, and modeling.

The purpose of this work is to investigate the abilities of a current state-of-the-art LLM (GPT-3.5 Turbo) to mimic human emoji usage in social media content. We conduct an experiment to generate social-media-style content with the LLM, then compare the patterns of emoji use in the LLM-generated content to those of human social media users.

## 2. Background

Emoji, small pictograms that can be integrated into online text communication, were first created in Japan in 1999, and then introduced to the broader international community when they were made available on Apple’s iPhone in 2011. Since text is limited in its ability to express tone, inflection, body language, and other embodied aspects of verbal communication, emoji have become an extremely popular method to encode additional meaning into informal digital communication. For a more thorough discussion of the history of emoji and their various linguistic properties, see McCulloch (2019).

Linguistics researchers have identified a number of patterns in the use of emoji in online communications. Gawne and McCulloch (2019) framed emoji usage as a form of digital gesture – a sort of linguistic property, but distinct from having grammar as a separate language. Their analysis demonstrated that the patterns of emoji usage are tightly analogous to those identified in the linguistics subfield of Gesture Studies. The gesture functions most relevant to our study here are illustrative, beat, and illocutionary.

Illustrative gestures are used to refer to physical objects or people in the real world, such as making the form of an item with the hands or indicating size with outstretched arms. Emoji used in this way depict graphically a thing or person in the real world – foods, sports gear, transportation vehicles, people in occupational clothing or engaging in an activity, etc. – and are generally used in a straightforward and literal manner. The majority of emoji options fall into this category, though it is not the most popular in actual usage. A property of this function is that there is no prescribed form. A sequence of birthday-related emoji will mean the same thing regardless of order or the specific order used (e.g. 🎂🎉🎂🎉, 🎂🎉🎂🎉, 🎂🎉🎂🎉, or 🎂🎉🎂🎉).

Beat gestures are repetitive and rhythmic, used to emphasize important points of what the speaker is saying. A politician pounding a fist on the podium or a sports fan pumping an arm in the air while cheering are quintessential beat gestures. When watching a conversation between two people too far away to hear, it is very often possible to tell who is speaking, as they will move their hands and arms while speaking in a way that they do not while listening. Emoji usage patterns in online written communication are highly similar to beat gestures in verbal communication. Some are explicitly tied to physical gesture, such as repeated raised hands (🙌🙌) or clapping between/after words (WHAT 👏 IS 👏 GOING 👏 ON 👏) (Tatman, 2017). Most beat emoji, however, are repetition sequences, such as tears of joy (😂😂), heart eyes (😍😍😍), or loudly sobbing (😭😭😭😭). In fact, this is the most common type of emoji sequence. In a corpus of English-language users of the SwiftKey keyboard app, the top 9 emoji bigrams (two-character sequences) and the top 22 emoji trigrams and quadrigrams (three- and four-character sequences, respectively) were all pure repetition (McCulloch and Gawne, 2018).

Gestures that do not relate directly to the semantic meaning of the spoken or written words but rather “indicate the type of speech act” (Gawne and McCulloch, 2019) are described by some linguists as expressing the illocutionary force of the communication. That is, they give some indication of the speaker’s intent, such as conviction, indifference, or sarcasm. Emoji used for this function are most often based on facial expressions or body language. For example, a smiling face (😊) might indicate pro-social intent, a smirk (😏) could indicate sarcasm, or a face throwing a kiss (😘) affection or flirtatiousness. As with gestures (McNeill, 2005, cited in Gawne and McCulloch, 2019), all of these functions that emoji can serve are non-exclusive, i.e. even a single usage can fall into more than one category. Emoji with illocutionary meaning can easily be repeated for emphasis as a beat gesture.

### 3. Methodology

The purpose of this experiment was to study emoji usage in LLM-generated social media text as compared to the patterns identified in the literature for human users. Since the emoji used are highly dependent on the intended meaning of the content, it was important that the subject matter of the LLM-generated text appropriately represented real human-written text on social media. To this end, all of the LLM-generated content was designed to mirror one-to-one the content of a representative dataset of human-written text. The dataset used is the UMass Global English on Twitter<sup>1</sup> Dataset, which contains over 10,000 randomly-sampled tweets, annotated with (among other things) English, non-English, and country (Blodgett et al., 2017; University of Massachusetts Amherst, 2017). Since the patterns of human emoji usage described in (McCulloch and Gawne, 2018; Medlock and McCulloch, 2016) are based on a corpus from users with their language set to US English, we focus our analysis initially on tweets annotated as English with country specified as US. This subset consists of 2,735 tweets.

**Table 1: Examples of original tweets, topic, and a new tweet generated by GPT 3.5**

Original tweet, from the UMass Global English on Twitter Dataset (Blodgett et al., 2017)	Topic, as summarized by GPT 3.5	LLM-generated tweet, from GPT 3.5 based on topic description
@[username] It's all about being real and true to yourself. The opposite of real is fake.	Authenticity and self-acceptance are key.	Just remember, being true to yourself and accepting who you are is the ultimate key to happiness and success. #authenticity #selfacceptance 🧡🌟

<sup>1</sup> In July 2023, Twitter changed the name of its platform to X. Since both the dataset and most of the literature come from before this name change, for consistency this paper refers to said platform as Twitter.

Everything I don't want to happen is happening. Smh I hate my life right now ! 😞	Negative feelings about current life circumstances.	Feeling super overwhelmed and discouraged by my current life circumstances right now 😞 Just trying to stay positive and remember that tough times don't last forever 🙌 #perspective #staystrong
I just wanna sleep all day tomorrow!	Desire for extended sleep and relaxation.	Current mood: just want to crawl back into bed and sleep for like, three days straight 😞 🌙 #needmoreest #relaxationgoals

Each of the tweets in this US English subset is then summarized by individually feeding through GPT 3.5 (specifically GPT-3.5-turbo-0125), using OpenAI’s API (OpenAI, 2023), with a prompt to “describe the overall topic of this tweet in eight words or fewer.”<sup>2</sup> The LLM-created content was then generated by prompting GPT 3.5 to write a tweet about each of those topics in turn (see Table 1 for examples). Only one artificial tweet was generated from the topic of each authentic tweet. Our dataset for analysis then is a set of 2,735 LLM-generated artificial English-language tweets. Using AntConc (version 4.2.4) corpus linguistics toolkit and regular expressions, we extracted list of the most common emoji n-grams for n=1, 2, 3, and 4. In addition to comparing LLM-generated tweets against the human patterns observed in the literature, we also look at emoji usage in the original dataset using the same extraction of common emoji n-grams.

Many emoji are encoded as multiple symbols that are then combined and rendered as one icon. This includes appending skin-tone markers for emoji depicting body parts and gender markers for depictions of persons (The Unicode Consortium, n.d.). Additionally, a special non-printable character called a Zero Width Joiner (ZWJ) can be used to combine emoji in different ways (“Zero Width Joiner Emoji,” n.d.): some combinations create a person with a depicted profession, optionally with gender and skin tone markers; some create family groupings. For this analysis, all skin-tone and gender modified emoji are counted together with their unmodified base emoji. ZWJ groupings are treated as n-gram sequences of their component emoji, as ZWJ emoji groupings are relatively new and not uniformly supported.

#### 4. Results

Across the 2,735 tweets generated by the LLM, there were 4,434 displayable Unicode characters associated with emoji<sup>3</sup> – an average of 1.6 emoji characters per generated tweet. While there are some overlaps between the list of top emoji in the LLM-generated tweets and the human-written tweets, they are largely very different distributions. Not only does the LLM’s dataset contain vastly more emoji for the same number of tweets, but most of the top-20 common emoji (Table 2) are greatly overrepresented compared to human-written content. For instance, the most common emoji in the LLM set is raising hands (🙌), appearing in fully 6% of tweets; by contrast, this emoji appears in less than 0.4% of tweets in the human set. The flexed muscle arm (💪) is even more skewed. It was in only 6 tweets in the human set, but 133 LLM-generated tweets, more than 22 times as frequent! Eleven of the top 20 most common emoji in the human set of tweets are faces, including the top three (see Table 4); only four of the top 20 in the LLM set are faces, the highest ranked being only the 10th most frequent (😄). The results from the human-written tweets in this dataset are fairly consistent with emoji usage on Twitter overall. Of the top 20 most common emoji in the original human dataset 10 are also in the top 20 most popular emoji used globally on Twitter as of 2023, according to Emojitracker (Emojipedia, 2023).

**Table 2: Most common emoji in LLM-generated tweets**

Rank	Emoji	No.of tweets	Total no.of emoji
1	🙌	165	165
2	❤️	138	138
3	👉	133	133
4	🎉	128	129

<sup>2</sup>We found in the methodology design process that if asked to describe the topic of a tweet without any restrictions, GPT 3.5 often would produce a summary as long or longer than the original tweet itself, and sometimes even repeat whole chunks of the input verbatim.

<sup>3</sup>This total includes 347 skin-tone and gender modifiers, and 64 “Regional Indicator Symbol Letters” that are used on compatible platforms in 2-character groupings to encode a single national or international flag (e.g. **u s** would be displayed as 🇺🇸).

Rank	Emoji	No.of tweets	Total no.of emoji
5	🔥	126	147
6	☀️	100	103
7	🎵	94	94
8	👛	93	99
8	✨	93	94
10	😄	77	81
11	😬	72	74
12	😍	61	61
13	😏	56	56
14	❤️	54	54
15	🚗	51	52
16	✳️	42	42
16	👤	42	42
18	🎆	41	50
19	💕	40	40
20	😏	39	39

Notes: N=2,735; ranked by the total number of tweets each emoji appears in; total no. of emoji includes some sequences of repeated emoji (see Table 3)

Table 3: Most common emoji n-gram sequences in LLM-generated tweets

Bigrams		Trigrams		Quadrigrams	
Emoji	No.	Emoji	No.	Emoji	No.
📱✍️	12	🔥🔥🔥	9	😄😄😄😄*	6
🔥🔥	11	😄😄😄*	7	📱👛😄\$*	3
🎆🍰	11	😄😄😄*	6		
🍔🍷	10	📱✍️👛	4		
🔥🎵	8	📱👛😄	4		
📱👛	7	☀️😄👛*	3		
😄😄*	7	🍔🍷🍷	2		
😄😄*	7	📱✍️✨	2		
👛😄*	7	💕👤👤	2		
🍷🍷	7	🍷👤🍷	2		
☀️😄	6				
😄😄*	6				
💕✨	6				
😍👤	6				
🚗👤	6				

Notes: Ranked by total number of tweets each sequence appears in; only n-grams appearing in at least 2 tweets are shown; \* indicates ZWJ groupings.

The differences in multi-character sequences of emoji between LLM-generated text and human-composed is even more stark. Linguistics researchers studying patterns of emoji use online (by humans) have found that when multiple emoji are used in a sequence, it is overwhelmingly most common in repeated patterns (McCulloch and Gawne, 2018; Medlock and McCulloch, 2016). McCulloch and Gawne (2018) found that of the top 200 most common n-grams (n=2, 3, and 4) in their SwiftKey US English corpus, roughly half were strict repetition of identical emoji characters. At the top this trend is even stronger: nine of the top ten emoji bigrams, and all of the top ten trigrams and quadrigrams were strict repetition. Even among non-identical sequences, most show some form of repetition, such as different types of hearts (e.g. ❤️💖), or different faces expressing the same sentiment (e.g. 😊😄). This is true of the UMass corpus (Blodgett et al., 2017) as well, with 15 of the top 25 most common emoji bigrams being straight repetition, and all but one of the emoji trigrams and quadrigrams appearing more than once being pure repetition. The LLM-generated tweets show a very different pattern of emoji use. Only one of the top 50 most common emoji bigrams and trigrams are identical repetitions. In fact, in the entire dataset of 2,735 LLM-generated tweets, only three repetition sequences appear in more than one tweet: 🔥🔥🔥 nine times, 🔥🔥 twice, and 😄😄 twice. There are zero sequences of four repeated emoji that appear more than once in the entire LLM-generated corpus, and only one sequence of four or more repeated emoji at all, appearing just once (★★★★★, describing a movie). What we see instead are mostly small heterogeneous groupings, thematically matching the text, but not matching the repetition patterns prevalent in real human usage. The LLM had a strong tendency to use pairs of different emoji to reflect the topic literally, as with illustrative gestures. Use of emoji as beat or illocutionary features, both of which are much more common than illustrative uses among human writers, was noticeably lacking among the LLM-generated text.

**Table 4: Most common emoji the original set of human-composed tweets**

Rank	Emoji	No. of tweets	Total no. of emoji
1	😄	63	115
2	😍	35	49
3	🤔	28	49
3	👉	28	33
5	❤️	27	37
6	😏	20	24
7	😞	19	21
8	😞	17	22
8	💖	17	19
8	😊	17	19
11	😏	13	13
12	👉	11	13
13	😏	10	13
13	😎	10	12
13	👉	10	11
16	💙	9	15
16	💯	9	13
16	😊	9	11
19	🔥	8	12
19	🎉	8	9

Notes: N=2,735; ranked by the total number of tweets each emoji appears in; total no. of emoji includes sequences of repeated emoji (see Table 5).

Table 5: Most common emoji n-gram sequences in the original set of human-composed tweets

Bigrams		Trigrams		Quadrigrams	
Emoji	No.	Emoji	No.	Emoji	No.
😂😂	21	😂😂😂	13	😂😂😂😂	7
😂😂	8	😂😂😂	7	😂😂😂😂	2
😂😂	7	😂😂😂	4	❤❤❤❤	2
❤❤	4	❤❤❤	3	❤❤❤❤	2
😂😂	3	😂😂😂	2		
😂😂	2	❤❤❤	2		
❤❤	2	🏀🏀🏀	2		
👏👏	2	100 100 100	2		
🏀🏀	2	🔥🔥🔥	2		
100 100	2	😂😂😂	2		
🔥🔥	2	😂😂😂	2		
😂😂	2	😂😂😂	2		
😂😂	2				
😂😂	2				
😂😂	2				
😂😂	2				
😂😂	2				
😂❤	2				
😂😂	2				
😂😂	2				
😂😂	2				
😂❤	2				
👏👏	2				

Notes: Ranked by total number of tweets each sequence appears in; only n-grams appearing in at least 2 tweets are shown.

## 5. Conclusion

This experiment quantitatively analyzed how emoji are used in LLM-generated social media style content and compared that to the typical use by human social media users. As we demonstrated in section 4, tweets written by this LLM, OpenAI’s GPT-3.5-Turbo, display markedly different patterns of emoji usage than humans. The most common emoji usage by humans consists largely of faces, hearts, and hands (Emojipedia, 2023), primarily to convey tone and intent as a substitute for body language. Repetition is typical, and by far the most common types of emoji sequences found in the wild are straight repetition. LLM-produced content, on the other hand, includes a much wider variety of emoji, with fewer faces and more literal representations of the text’s subject matter. GPT 3.5 seems to favor a few emoji in particular, notably the raised hands (👏) and flexed bicep (💪). Some of this overuse is likely due to the LLM’s deliberate tuning by its developers. For instance, many uses of the flexed bicep emoji are in response to prompted topics expressing negativity, to which GPT 3.5 added an unprompted positive outlook (see Table 1 for one example).

In stark contrast to humans' proclivity for repetition using only a small subset of the most popular emoji, the LLM's emoji groupings were diverse and varied, resulting in many unique combinations rather than the same predictable few. This can result in the LLM-generated "tweets" looking a lot more like generic corporate brand accounts than average users (one example from the LLM content: "Don't forget to tune in to the big game tonight at 6:30 PM! Who's ready for some intense action on the field? 🏈 ⚽ #sports #excited #gametime" [A football and a soccer ball? Which 'big game' are we watching exactly?]).

These results have implications for the open question of detecting LLM-generated content on social media. While emoji use alone cannot flag an online post as artificial, patterns of emoji use over time could be used as one element contributing to detection of accounts posting large amounts of LLM-generated content.

## References

- Blodgett, S.L., Wei, J.T.Z. and O'Connor, B. (2017) 'A Dataset and Classifier for Recognizing Social Media English', *3rd Workshop on Noisy User-Generated Text, W-NUT 2017 - Proceedings of the Workshop*, pp. 56–61. doi:10.18653/V1/W17-4408.
- Emojipedia (2023) *Emojitracker*. Available at: <https://emojitracker.com/> (Accessed: 13 March 2024).
- Gawne, L. and McCulloch, G. (2019) 'Emoji as Digital Gestures', *Language@Internet*, 17, p. article 2. Available at: <https://www.languageatinternet.org/articles/2019/gawne/index.html> (Accessed: 11 February 2024).
- McCulloch, G. (2019a) *Because Internet: Understanding the New Rules of Language*. New York: Riverhead Books.
- McCulloch, G. (2019b) 'Emoji and Other Internet Gestures', in *Because Internet: Understanding the New Rules of Language*. New York: Riverhead, pp. 155–196.
- McCulloch, G. and Gawne, L. (2018) 'Emoji Grammar as Beat Gestures', in *Proceedings of the 1st International Workshop on Emoji Understanding and Applications in Social Media (Emoji2018)*. Available at: <http://ceur-ws.org>.
- McNeill, D. (2005) *Gesture and Thought*. Chicago: University of Chicago Press.
- Medlock, B. and McCulloch, G. (no date) *The Linguistic Secrets Found in Billions of Emoji - SXSW 2016 presentation*. Available at: <https://www.slideshare.net/SwiftKey/the-linguistic-secrets-found-in-billions-of-emoji-sxsw-2016-presentation-59956212> (Accessed: 11 February 2024).
- OpenAI (2023) *Text generation models, OpenAI developer platform documentation*. Available at: <https://platform.openai.com/docs/guides/text-generation> (Accessed: 6 December 2023).
- Tatman, R. (2017) *Where 🙌 do 🙌 the 🙌 claps 🙌 go 🙌 when 🙌 you 🙌 write 🙌 like 🙌 this 🙌?, Making Noise & Hearing Things*. Available at: <https://makingnoiseandhearingthings.com/2017/07/13/where-🙌-do-🙌-the-🙌-claps-🙌-go-🙌-when-🙌-you-🙌-write-🙌-like-🙌-this-🙌/> (Accessed: 3 March 2024).
- The Unicode Consortium (no date) *Full Emoji Modifier Sequences, v15.1*.
- University of Massachusetts Amherst (no date) *Twitter Language Identification: Data and Software, SLANG Lab at UMass*. Available at: <https://slanglab.cs.umass.edu/TwitterLangID/> (Accessed: 4 March 2024).
- Zero Width Joiner Emoji (no date). Available at: <https://emojipedia.org/zero-width-joiner> (Accessed: 4 March 2024).