

# New Modality: Emoji Challenges in Prediction, Anticipation, and Retrieval

Spencer Cappallo , Stacey Svetlichnaya , Pierre Garrigues, Thomas Mensink ,  
and Cees G. M. Snoek , *Senior Member, IEEE*

**Abstract**—Over the past decade, emoji have emerged as a new and widespread form of digital communication, spanning diverse social networks and spoken languages. We propose treating these ideograms as a new modality in their own right, distinct in their semantic structure from both the text in which they are often embedded as well as the images which they resemble. As a new modality, emoji present rich novel possibilities for representation and interaction. In this paper, we explore the challenges that arise naturally from considering the emoji modality through the lens of multimedia research, specifically the ways in which emoji can be related to other common modalities such as text and images. To do so, we first present a large-scale data set of real-world emoji usage collected from Twitter. This data set contains examples of both text-emoji and image-emoji relationships within tweets. We present baseline results on the challenge of predicting emoji from both text and images, using state-of-the-art neural networks. Further, we offer a first consideration into the problem of how to account for new, unseen emoji—a relevant issue as the emoji vocabulary continues to expand on a yearly basis. Finally, we present results for multimedia retrieval using emoji as queries.

**Index Terms**—Content-based retrieval, image classification, machine learning, social computing.

## I. INTRODUCTION

EMOJI, small ideograms depicting objects, people, and scenes, have exploded in popularity. They are now available on all major mobile phone platforms and social media websites, as well as many other places. According to the *Oxford English Dictionary*, the term *emoji* is a Japanese coinage meaning ‘pictogram’, created by combining *e* (picture) with *moji* (letter or character). Emoji as we know them were first introduced as a set of 176 pictogram available to users of Japanese mobile phones. The available range of ideograms has expanded greatly over the previous years, with 1,144 single emoji characters defined in Unicode 10.0 and many more defined through combinations of two or more emoji characters. In this paper, we approach emoji

as a modality related to, but not contained within, text and images. We investigate the properties and challenges of relating these modalities to emoji, as well as the multimedia retrieval opportunities that emoji present.

The identification and benchmarking of novel modalities has a rich history in the multimedia community. When new modalities are identified, it is important to make first attempts to understand their relationship with already established information channels. One way in which to do this is to explore the cross-modal relationships between the modality and other modalities. When Lee *et al.* [19] identified nonverbal head nods as an information-rich and overlooked modality, they sought to provide understanding through prediction of them based on semantic understanding of the accompanying conversation transcript. Like emoji, new modalities are sometimes the result of a newly developed technology, as with 3D models [15] or the growth of microblogging [2]. Though ideograms are ancient, emoji are a modern technological evolution of that ancient idea. The march of technology sometimes facilitates new looks at old problems, such as the use of infrared imagery for facial recognition instead of natural images [43]. Often, the presentation of new tasks as research challenges can accelerate research progress, as it did with acoustic scenes [39] and video concepts [38]. We look to this history of multimedia challenge problems and identify emoji as an emerging modality worthy of a similar treatment. To facilitate further research on emoji, we propose three emoji challenge problems and present state-of-the-art neural network baselines for them, as well as a dataset for evaluation.

Despite their prevalence, research into emoji remains limited. The majority of prior research concerning emoji has focused on descriptive analysis, such as identifying how patterns of emoji usage shift among different demographics [5], [11], or has used them as a signal to indicate the emotional affect of accompanying media [16], [34]. The focus on sentiment is likely a result of there being a number of “face emoji” (e.g. 😊) which are designed to exhibit a particular emotion or reaction. These face emoji are by far the most visible emoji and among the most widely used [33], but the focus on them ignores the hundreds of other emoji which are worthy of study in their own right. Beyond these face emoji, the full set of emoji also contains a wide range of other objects, such as foods (🍕), signs (📺), and scenes (🏠) which may lack a strong sentimental signal [32]. Recently, Apple has introduced *Animoji* which allow users to animate select emoji with facial expressions, further broadening their range of

Manuscript received January 25, 2018; revised May 13, 2018 and June 29, 2018; accepted July 2, 2018. Date of publication August 1, 2018; date of current version January 24, 2019. This work was supported by the STW Story project. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Lei Zhang. (*Corresponding author: Spencer Cappallo.*)

S. Cappallo, T. Mensink, and C. G. M. Snoek are with the University of Amsterdam, 1012 WX Amsterdam, The Netherlands (e-mail: astucity@gmail.com; tmensink@uva.nl; cgmsnoek@uva.nl).

S. Svetlichnaya and P. Garrigues are with Yahoo Research, San Francisco, CA 94111 USA (e-mail: stacey.svet@gmail.com; pierre.garrigues@gmail.com).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TMM.2018.2862363



Fig. 1. Emoji prediction used for video summarization and query-by-emoji, adapted from our previous work [8]. **A.** The emoji summarization of the entire video presents a more complete representation of the video’s contents than a single screenshot might. **B.** Emoji can be used as a language-agnostic query language for media retrieval tasks. Here, emoji are used to retrieve photos from the MSCOCO dataset. Despite their limited vocabulary, emoji can be combined to compose more nuanced queries, such as *shoe+cat*. This results in a surprisingly flexible modality for both content description and retrieval.

emotional expression. While emoji can be powerful signals of emotion, focusing solely on the emotion-laden subset of emoji ignores the information conveyed and possibilities presented by the many other ideograms available.

In this work, we approach emoji as an information-rich modality in their own right. Though emoji are commonly embedded in text, we view them as distinct from text. Their visual nature allows for emoji to add richness of meaning and variety of semantics that is unavailable in pure text. When embedded in text, emoji sometimes simply replace a word, but more often they provide new information which was not contained in the text alone [1], [29]. Emoji can be used as a supplemental modality to clarify the intended sense of an ambiguous message [35], attach sentiment to a message [37], or subvert the original meaning of the text in ways a word could not [12], [30]. Emoji carry meaning on their own, and possess compositionality allowing for more nuanced semantics through multi-emoji phrases [22]. Many emoji are used in cases where the particular symbol resembles something else entirely, acting as a kind of visual pun. These qualities, along with a cross-language similarity of semantics [5], suggest that emoji, despite being unicode characters, are distinct from their frequent textual bedfellows.

Though emoji are represented by small pictures, they are distinct from standard images. As a form of symbology, the specifics of an emoji’s representation are often incidental to the underlying meaning of the ideogram. This is unlike images where the particulars of a given image are often more crucial than what it is representing generally. For example, a photo may be a photo of *your* dog, not just a photo representing the semantic notion of ‘dog’, while the dog emoji is unlikely to refer to one particular canine. This difference is further substantiated by the fact that emoji exist as nothing more than unicode characters. As characters, the details of their illustrations are left up to the platform supporting them, and significant variation for

a single emoji can exist between platforms [27], [42]. Furthermore, given the small size and illustrative nature of emoji, their low-level statistics will be very different from those of natural images. For these reasons, their behaviour and meaning is substantially different from that of images. Fig. 1 gives examples of video summary using emoji and query-by-emoji, which nicely demonstrate the way in which emoji as ideograms are related to but different from natural imagery.

Having established the view that emoji constitute a distinct modality from text or images, this paper seeks to explore the ramifications of this viewpoint through the lens of multimedia retrieval challenges. As a modality, we focus on the relationship between emoji and two other modalities, namely text and images.

This work makes the following contributions:

- We propose and support the treatment of emoji as a modality distinct from either text or images.
- We present a large scale dataset composed of real-world emoji usage on Twitter, containing both textual and text+image examples. We consider a wide range of over 1000 emoji, including the often overlooked long tail of emoji. To facilitate focus on the long tail of emoji usage, we present a balanced test set (in addition to the natural, unbalanced test set) which will give extra weight to those often overlooked long tail emoji. This dataset as well as the training splits are available for future researchers.
- We propose three challenge tasks for relating emoji to text and images, and present state-of-the-art, off-the-shelf baseline results on these. Namely, the tasks are emoji prediction from text and/or images, prediction of unanticipated emoji using their unicode description, and lastly multimedia retrieval using emoji as queries.

In the following section we give an overview of previous work on emoji. In Section III we present our dataset, and propose three

challenge tasks presented by the emoji modality. In Sections IV, V, and VI we present baseline results for each of these challenge tasks using state-of-the-art deep learning approaches. In Section VII, we conclude.

## II. RELATED WORK

Previous work on emoji in the scientific community has focused on using them as a source of sentiment annotation, or on descriptive analysis of emoji usage.

### A. *Emoji for Sentiment*

Much of the prior work has viewed emoji primarily as an indicator of sentiment. This is done either explicitly, through the direct consideration of sentiment, or implicitly, through the consideration of only popular emoji. The most popular emoji are disproportionately composed of sentiment-laden emoji. Face emojis, thumbs-up, and hearts have high incidence, while less emotional emoji such as symbols, objects, and flags, have much lower incidence. The result is that any work which considers only the most popular emoji may have an inherent bias toward emoji with heavy sentiment.

Several works look at the effect that including emoji can have on the perception of accompanying text. Some find that the inclusion of emoji increases the perceived level of sentiment attached to a message [29], [32], [37]. Similarly, the work from [36] finds that emoji correlate to a more positive perception for messages in a dating app than messages that don't contain emoji. These works demonstrate that emoji can be a useful supplementary signal for sentiment within text messages, but these works focus primarily on face emoji designed specifically for the communication of emotion. In contrast, Riordan [35] investigates the affect of non-face emoji. They found that even non-face emoji can increase perceived emotion, and also can improve clarity of text that is otherwise ambiguous. Some text phrases are ambiguous when considered alone, but the inclusion of another modality (emoji) can help readers to pinpoint the intended sense (e.g. "I took the shot" vs "I took the shot 🎯").

A notable work of sentiment analysis of emoji is [32], which annotated a collection of tweets with sentiment and presented sentiment rankings for 751 emoji (the most frequent in their data). Their work demonstrated that while some emoji have very strong positive sentiment scores, others were very neutral, being rarely associated with strong positive or negative sentiment. Similarly, they observed that some emoji are used frequently to denote both strong positive and negative sentiment. These observations suggest that treating emoji as merely a straightforward signal of sentiment is misguided, and that there is a more nuanced richness and variety to emoji meaning.

Lastly, some works consider emoji, particularly face emoji, as a pure sentiment signal. The approach by [34] incorporates emoji as an input source for evaluating the sentiment of social media messages mentioning particular brands. Going a step further, Guthier *et al.* [16] assumes emoji to be a reliable ground truth for sentiment. They construct a dataset for sentiment predic-

tion and use a set of emoji to automatically annotate the dataset. Given the broad ambiguity of usage and the sentiment gap between emoji and text explored in other works, such an approach may yield noisy annotation.

### B. *Analysis of Emoji Usage*

Numerous works have helped to glean insight into the properties and trends of real-world emoji usage. Several have looked at the manner in which emoji usage varies between different countries and cultures [5], [21], [23]. Meanwhile Chen *et al.* [11] analyzes differences in emoji usage patterns between genders. While there are differences between how specific communities may use emoji, the data makes clear that emoji usage is on the rise globally [21], [46]. This further supports our viewpoint that emoji are their own modality, as they are not tied to any one particular culture or language and share semantic commonalities which are orthogonal to the community that uses them.

Several works look at the problem of ambiguity in the perceived meaning of emoji [27], [28], [42]. In general, they find a degree of ambiguity with emoji, and that the choice of illustration used by a particular platform (e.g. iOS or Android) can increase this confusion. Notably, Miller *et al.* [28] observes that the inclusion of an additional input modality (in the form of textual context) improves the distinctiveness of meaning substantially. This observation is well in line with what has been known in the multimedia community for years: that a multimodal approach can improve prediction. Ambiguity between the message intent from the author of an emoji-containing message and its interpretation by readers has also been investigated [7]. The ambiguity and breadth of possible meaning for a given emoji helps to make emoji a challenging modality for algorithmic understanding, worthy of pursuing and with a high ceiling for perfection.

The relationship among emoji themselves has been studied in [6], [33], [45]. The work of [33] gives a thorough analysis of emoji usage, and proposes a model for analyzing the relatedness of pairs of emoji. Similarly, Barbieri *et al.* [6] looks at the problem of trying to identify text tokens which are most closely related to a given emoji. The authors do this by learning a shared embedding space using a skip-gram model [25], and identifying those text tokens closest to the emoji within this mutual semantic space. While both [33] and [6] learn models that could be applied to emoji prediction, they both focus instead on descriptive analysis of emoji usage.

Along similar lines, there has been some recent work on identifying the different ways in which emoji can be used in combination with text. [1], [12], [29] use emoji either as a straightforward replacement for text, or as a supplementary contribution which alters or enhances the meaning of the text. The work of [12] constructs a dataset of 4100 tweets that have been annotated to indicate whether the emoji contain redundant information (already contained in the text) or not. Among their collection of annotated tweets, they found that the non-redundant class was the largest class of emoji. This result supports our proposition that emoji are distinct from, though entwined with, any text that accompanies them.

While works such as [1], [6], [33] tackle the problem of understanding emoji usage through building models on top of real world usage data, there has also been work which tries to build an emoji understanding in a more hand-crafted fashion. For example, Wijeratne *et al.* [44] acquires a structured understanding of emoji usage through combining several user-defined databases of emoji meaning. Their later work then uses this data to learn a model for sentiment analysis which performs comparably to models trained directly on real world usage data [45]. This kind of structured, pre-defined understanding of emoji is similar to the no-example approach explored in our previous work [8] and further explored in this work. This work, however, targets emoji as a rich, informative modality rather than only a means for performing sentiment analysis.

[22] is an early investigation into the compositionality of emoji. They find that emoji can be combined to create new composed meanings, a finding which lends support to the notion of composing queries from multiple emojis that is discussed in this work.

Much of the analysis of these works support our philosophy of treating emoji as a modality in their own right. In contrast to these works and to complement them, rather than trying to provide descriptive analysis of emoji usage, we focus on how the emoji can be used with and related to other modalities.

### C. Cross-modal Emoji Prediction

A few recent works have investigated the problem of emoji prediction, which is closer to our position of emoji-as-modality.

Our previous work was the first to look at the problem of emoji prediction, and approached it from a zero-shot perspective due to a lack of an established dataset [8]. Following on from the work, a query-by-emoji video search engine was also proposed [9]. These works reported quantitative results only on related tasks in other modalities, and presented only qualitative results for the emoji modality. We instead present results on a large scale, real-world emoji dataset, with proposed tasks and state-of-the-art supervised baselines.

Felbo *et al.* [14] train a model to predict emoji based on input text. Rather than using the model directly for the task of emoji prediction, they use this model as a form of pre-training for learning a sentiment prediction network. Additionally, their emoji model is intentionally limited to 64 emoji chosen for having a high degree of sentiment. Our aim is to treat emoji as an end goal rather than an intermediary, and to consider the full breadth of emoji available including rare emoji or emoji with little or no sentiment attached to them.

Barbieri *et al.* [4] looked at the problem of emoji prediction based on an input text. Their setting is most similar to the one considered in this paper. However, they focus strictly on text, while we also consider images. Further, Barbieri *et al.* restrict their labels to only the top 20 most frequent emoji within their dataset. Along similar lines, Li *et al.* [20] uses a convolutional network to predict 100 common emoji based on a corresponding text from weibo or another social media network. Both of these papers consider only the most common emoji. There are thousands of emoji, and the longtail of the available emoji present a valuable and difficult prediction task. We consider the full range

of emoji present in our dataset, and look at the problems involved with tackling this longtail. We further distinguish our work by also considering the problem of newly introduced emoji, which is important as the set of available ideograms is growing every year.

El Ali *et al.* [13] is, to the best of our knowledge, the only previous work that considers supervised prediction of emoji from images. Their work looks at the problem of translating images of faces into corresponding face emojis. We take a broader approach both on the image and annotation sides, seeking to instead predict any sort of relevant emoji based on a wide variety of images.

## III. NEW MODALITY

There is no guarantee that a simple explanation of what an emoji depicts will encompass its full semantic burden. Emoji are inherently representational, so by definition some overlap in semantics is expected, but that overlap may be incomplete in terms of real-world usage. For example, the emoji for *cactus* 🌵 is not used only to represent a cactus, but is also widely used to signify a negative sentiment due to its resemblance to a certain hand gesture. This discrepancy between the intended semantics and the actual semantics leads us to propose learning the semantics directly from real-world usage in a large dataset collected from Twitter.

Motivated by our view that emoji constitute a separate modality, in this section we outline our methodological approach to establishing baseline analysis and results for the emoji modality. We begin by establishing three emoji challenge tasks, and subsequently propose a large dataset of real-world emoji usage as a testbed for exploring these challenges. We further propose evaluation criteria to quantify and compare performance on these challenges and dataset. An overview of how these three tasks differ in their objectives and the information available to them is provided in Fig. 2.

### A. Emoji Challenges

1) *Emoji Prediction - How to predict emoji?:* There are thousands of emoji, and new ones are added every year. As they develop into an ever richer information signal, it is useful to understand how emoji are related to other modalities. The most straightforward way to go about this is to look at how well we can predict emoji given another, related input. Since emoji can be flexible in their usage, the question becomes: Given some input text and/or image, can we predict the relevant emoji that would accompany that input? This work seeks to present strong first baselines for the problem.

We propose an Emoji Prediction challenge where the objective is to predict relevant emoji from alternative input modalities. Using real-world training examples correlating text and images to emoji annotations, models seek to predict relevant emoji when presented with test examples.

2) *Emoji Anticipation - What to do about new emoji?:* A large real-world dataset provides the opportunity for learning how to use emoji in a natural way that reflects their true semantics. However, new emoji are added to the unicode specification every year, and will be deployed to users before their real world

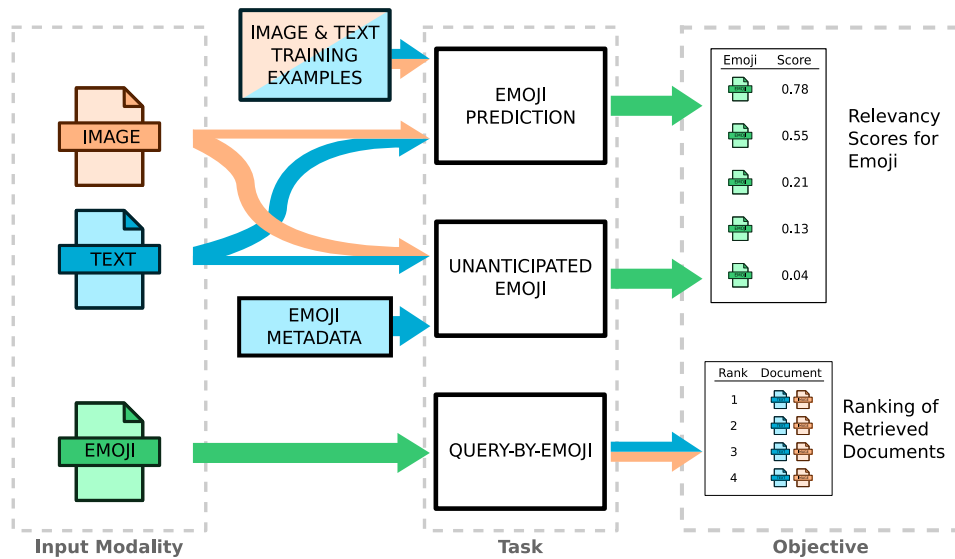




















Fig. 2. Overview of our three proposed tasks. Emoji Prediction and Unanticipated Emoji both seek to score emoji based on other input modalities. Their difference is that Emoji Prediction has the benefit of emoji-annotated training examples to learn from, while Unanticipated Emoji simulates the setting of newly released emoji where there is no training data available. Instead, Emoji Anticipation must use textual metadata describing the emoji to relate them to the input modalities. Query-by-Emoji seeks to retrieve relevant multi-modal documents using queries composed with emoji.

TABLE I  
UNICODE-PROVIDED EMOJI NAMES AND KEYWORDS, ALONG WITH THE REPRESENTATIONS FOR THAT PARTICULAR EMOJI ON THREE MAJOR PLATFORMS

Name	Keywords	Apple	Google	Twitter
dog	pet			
person in steamy room	sauna, steam room, hamam, steambath			
person climbing	climber			
sad but relieved face	disappointed, face, relieved, whew			
dizzy face	dizzy, face			
face with steam from nose	face, triumph, won			

The Name and Keywords can be Used During the Emoji Anticipation task, Though they Might not Align well with Popular Usage.

usage can be known. A similar challenge is also present in the related phenomenon of message stickers — small illustrations that can be sent in lieu of text. Stickers share some similarity in function to emoji, but are platform specific and can be released without major oversight, meaning the likelihood of significant training data is small. Any system that seeks to understand or suggest emoji (or stickers) to users should be prepared to deal with the challenge of new, previously unseen emoji.

In the Emoji Anticipation challenge, real world training data of emoji usage is no longer available. This simulates the situation when a new crop of emoji have been announced, but have not yet been deployed onto common platforms. Systems seeking to understand and predict these emoji must therefore exploit alternative knowledge sources. We present the problem as a zero-shot cross-modal problem, where we have only textual metadata regarding the emoji and must then try to determine its relevancy to images or text. An example of the information available is presented in Table I. This task shares some resemblance to that of zero-shot image classification [3], [31] or zero example video retrieval [10], [18]. Generally, in zero-shot classification the model has a disjoint set of seen and unseen

classes, and attempts to leverage the knowledge of seen classes as well as external information to classify the unseen classes. Our setting differs from this, as we test our model in a setting where it has seen no direct examples of the target modality whatsoever.

3) *Query-by-Emoji - Can we query with emoji?*: Not only can emoji be predicted for a given input modality, but they can also be used as queries to retrieve other modalities. Emoji have some unique advantages for retrieval tasks. The limited nature of emoji (1000+ ideograms as opposed to 100,000+ words) allows for a greater level of certainty regarding the possible query space. Furthermore, emoji are not tied to any particular natural language, and most emoji are pan-cultural. This means that emoji can be deployed as a query language in situations where a spoken language might fail. For example, with children who haven't yet learned to read, or perhaps even high intelligence animals such as apes. Further, the square form factor of emoji works naturally with touch screen interfaces. Many of these advantages are shared by any ideogram scheme, but emoji have the additional benefit of exceptional cultural penetration. Because emoji are already adopted and used daily by millions,

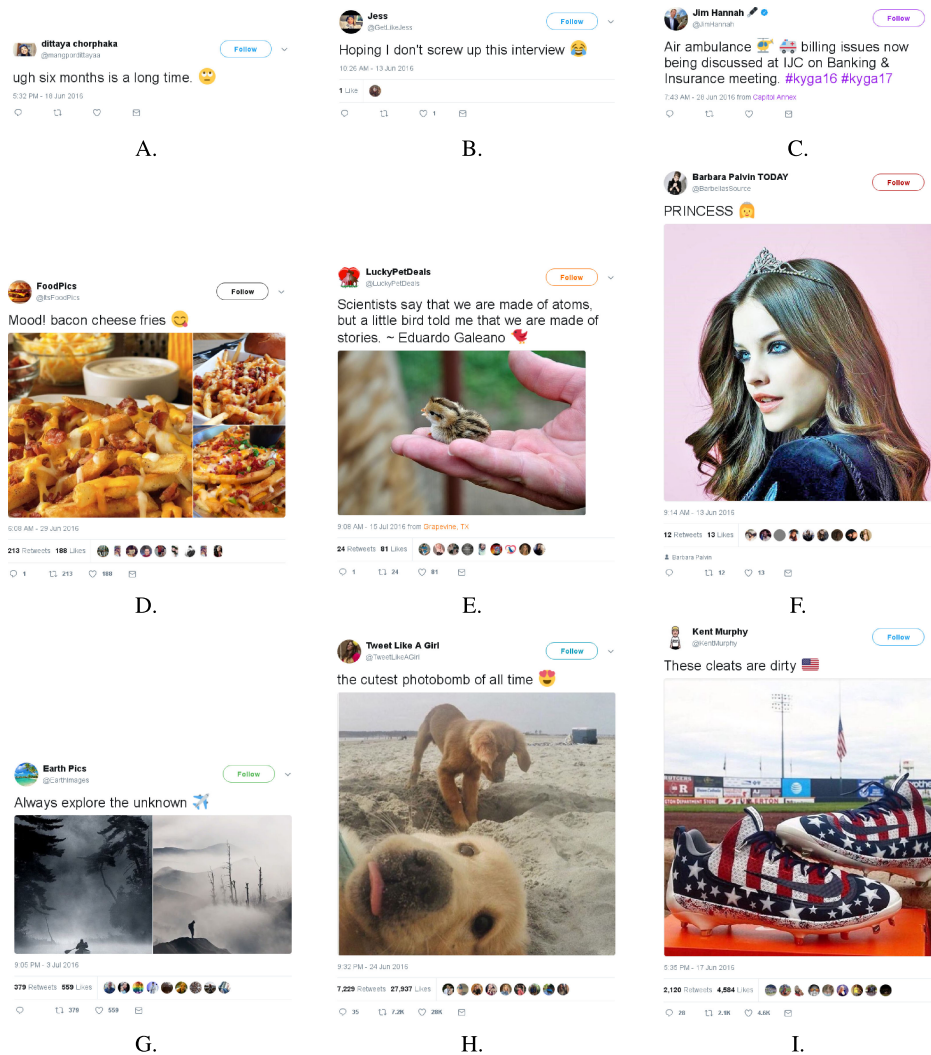


Fig. 3. Example tweets from the proposed Twemoji dataset. Emoji are removed and used as ground truth annotation. The top row gives examples of text-only tweets, while the bottom rows contain both the text and image modalities. We see the interactions between the three modalities (text, images, and emoji) can vary. For example, F has a strong alignment between all three, while the correlation between the emoji and the tweet is more obvious in the image than the text. Sometimes emoji re-confirm content, as in E, and sometimes they express a sentiment as in D. G gives an example where the emoji modify the content semantics—the airplane emoji adds a suggestion of travel that is not strictly present in either the text or image modalities. Emoji are intertwined with their related modalities, but are definitely not subsumed by them.

the cognitive burden to learn what emoji are available to use as queries is significantly decreased. Indeed, platforms such as Microsoft Bing and Instagram have already begun allowing the inclusion of emoji in their search systems, highlighting the need for a benchmark assessment within the multimedia community for this emerging problem.

In the Query-by-Emoji challenge, we aim to quantify performance on the task of multimedia retrieval given an emoji query. Samples in the test set should be ranked by the model for a given emoji query, and performance will be evaluated based on whether those documents are considered relevant to that emoji or not.

## B. Dataset

To facilitate research on these challenges, it is necessary to use a dataset with sufficient examples of the relationship between emoji and other modalities. Existing works on emoji have ei-

ther forgone the use of an annotated emoji dataset or have used datasets comprised of only a small subset of available emoji. Both of these settings are artificial and fail to adequately represent the challenge and promise of emoji. Instead, we target the full range of potential emoji, including their very long tail, and seek to learn their real-world usage rather than place any prior assumptions on them. We construct our dataset, which we call Twemoji, from the popular microblogging platform Twitter, and also identify two valuable subsets of the dataset. The dataset and details of the splits discussed below are publicly available.<sup>1</sup>

To generate a representative emoji dataset, we collected 25M tweets via the Twitter streaming API during the summer of 2016, filtering these to 15M unique English language tweets that contain at least one emoji. Fig. 3 gives some examples of tweets in our dataset. Emoji are common on Twitter, appearing in roughly

<sup>1</sup>Twemoji Dataset, DOI: 10.21942/uva.5822100

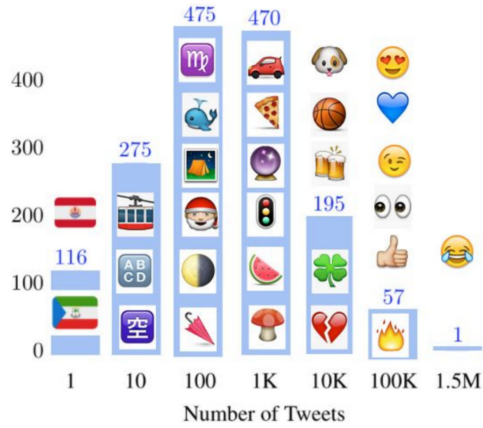


Fig. 4. Emoji Usage Histogram. The bars show the count of emoji appearing in at least N tweets—e.g. 275 different emoji each appear in 10–100 tweets. In each column, some examples of the emoji in that rarity bracket are displayed.

TABLE II  
TWEMOJI DATASET AND SUBSET STATISTICS

	Full	Balanced	Images
# Train Samples	13M	13M	917K
# Validation Samples	1M	1M	80K
# Test Samples	1M	10K	80K
# Emoji Present	1242	1242	1082

Full is the entire collection, balanced has a class-balanced test set but uses the same training and validation sets, and images is composed from those tweets with attached images.

1% of the tweets posted during our collection period. However, the usage frequency is heavily skewed (see Fig. 4). 🤔 is the most commonly used emoji, and it appears in 1.57M tweets. The top emoji (appearing in 100K+ tweets) are mostly facial expressions, hearts, and a few hand gestures (🙄, 👍, ❤️). Most emoji in the dataset have only hundreds (👉, 🍷) and thousands (🐦, 🍷) of examples. Flags and symbols compose the bulk of the rarer emoji.

A fraction of the tweets also contain images, which allow us to present results for the relationship between not only text and emoji but also images and emoji. We therefore present three selections of this dataset: Full, comprised of all tweets in the collection; Balanced, which has a test set constructed with a flattened distribution across emoji; and, Images, which is comprised of those tweets in the collection containing both emoji and images. We present statistics for the three subsets in Table II, and describe their composition below.

1) *Twemoji (Full)*: The Twitter data set is split randomly into training, validation, and test sets containing 13M, 1M, and 1M tweets, respectively. Input and annotation pairs are created by removing the emoji from the tweets’ text to use as annotation. This approach means that the data set is multi-label, though the preponderance of tweets have only one correct annotation. Fig. 5 shows the number of tweets with a given emoji annotation count. Noting that the  $y$ -axis is plotted on a log scale, we see that there are almost an order of magnitude more tweets with one emoji than with two emoji, and the numbers continue to drop.

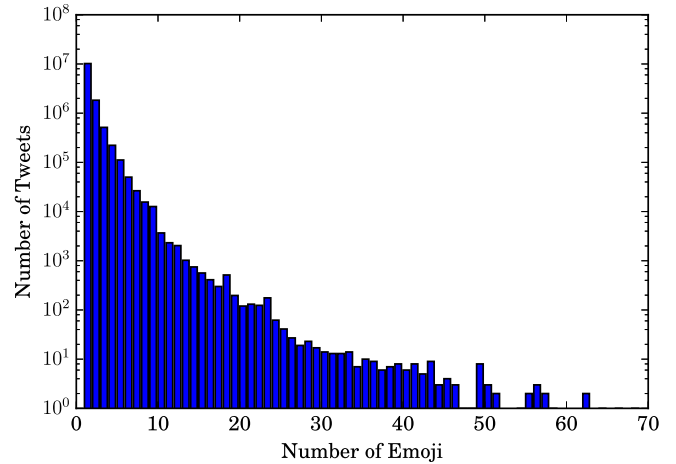


Fig. 5. Frequency of tweets containing multiple, distinct emoji in the Twemoji-Full training set, plotted on a log scale. We see that a few tweets contain many emoji, but the majority of tweets contain only one or two different emoji.

Emoji relevance annotations are treated as binary, so multiple occurrences of the same emoji in a single tweet are only counted once. A few tweets contain very many emoji. These are perhaps tweets where emoji are being used as a visual language.

The use of emoji as annotation assumes that the majority of emoji provide only supplementary information, and are not operating merely as one-to-one replacements for text tokens (e.g., “in 🚗 going to 🏠 to meet new 🤔” is no longer parseable text without the emoji, while for “awesome day 🤔” the message remains complete without the emoji).

2) *Twemoji Balanced*: While there is undoubtedly a natural imbalance of emoji popularity, we assume that current emoji interfaces may be a contributing factor to the distribution skew of emoji usage. The difficulty in navigating to a desired emoji, compounded with users being unfamiliar with rarer emoji, means that the heavy skew of the distribution could be a self-fulfilling prophecy and an undesired one. Further, it is not clear that the skew of commonly used emoji says anything about their relevance for new tasks like summarization using emoji. We therefore target the case when all emoji are used equally often. Targeting an equal balance ensures that commonly overlooked emoji will still be suggested, and can help eliminate undesired dataset biases. To evaluate this, we test on a more balanced, randomly selected subset of the test set in addition to the full, unbalanced test set. Use of this balanced test set also helps present a more complete picture regarding algorithm performance, by giving extra weight to the more difficult-to-predict long tail of emoji.

The balanced subset is selected such that no single emoji annotation applies to more than 10 examples. To train toward this objective while still leveraging the breadth of the available data, we construct our mini batches so that each emoji has an equal chance of being selected. Namely, the likelihood of selecting a particular sample  $x_i$  is

$$p(x_i) = \frac{\text{cnt}(y_i)^{-1}}{\sum \text{cnt}(y_i)^{-1}} \quad (1)$$

where  $\text{cnt}(y_i)$  returns the total count of samples with the same emoji annotation  $y_i$ . While over time this assures that every emoji equally contributes to the model updates, the model will still gain a more nuanced understanding of the more common emoji due to the diversity of the training samples.

3) *Twemoji Images*: Not all of the images contained in the tweets were still available on the internet, but those that were were downloaded. From these, we constructed a subset of the dataset for which both image and text inputs were available. Due to the prevalence of image-sharing on Twitter and the internet as a whole, a large number of tweets contain the exact same image as other tweets. We use the image-bearing tweets in the full Twemoji test set as our test set. We allow duplicate images between the train and test sets, but only when the emoji annotation of the test set differs from that in the training set. This results in a training set of 900K images, and validation and test sets of 80K images.

### C. Evaluation Protocols

1) *Emoji Prediction*: Performance in the Emoji Prediction challenge is reported in both Top- $k$  accuracy and mean samplewise Average Precision (msAP). Top- $k$  accuracy corresponds directly to the scenario in which a system is suggesting some  $k$  emoji that the user may wish to include during message composition, and the system should try to ensure that at least some of these emoji are relevant. As our dataset is multi-label, we calculate Top- $k$  accuracy by considering a prediction as correct if *any* predicted class in the top  $k$  is annotated as relevant, and a prediction as false if there are none. This means that an emoji ranking for a given input may score a relevant emoji as very unlikely, but still be marked as correct if a different, relevant emoji is correctly predicted in the top  $k$ . For  $N$  samples, where each input  $x_i$  has a corresponding binary vector  $y_i$  indicating emoji relevancy, the top- $k$  accuracy is calculated with

$$\text{ind}_k(x_i, y_i) = \begin{cases} 1 & \sum_{j \in \text{top}_k(p(y_i|x_i))} y_i^j > 0 \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

$$\text{Top-}k = \frac{\sum_{i=0}^N \text{ind}_k(x_i, y_i)}{N} \quad (3)$$

where  $\text{top}_k(p(y_i|x_i))$  yields the indices of the  $k$  highest scoring class predictions, and  $y_i^j$  corresponds to the value of the  $j$ th element of  $y_i$ .

As this approach may be overly optimistic on multi-label samples, we also report the mean samplewise Average Precision. This measures the performance of the algorithm across the entire ranking of emoji for a given input. It evaluates how accurately ranked the emoji are for a given image and/or text input.

$$\text{msAP} = \frac{1}{N} \sum_i \frac{\sum_j^C \text{Prec}(j) \times y_i^j}{\sum y_i} \quad (4)$$

where  $\text{Prec}(j)$  gives the precision of the prediction at rank  $j$ ,  $C$  is the total number of emoji, and  $y_i^j$  gives the value of  $y_i$  at the index  $j$ .

2) *Emoji Anticipation*: Emoji Anticipation differs from Emoji Prediction in its absence of training data, but the test

set and goal of the challenge is shared with Emoji Prediction. For this reason, results are again reported in both Top- $k$  accuracy and msAP.

3) *Query-by-Emoji*: Query-by-Emoji turns the problem on its head: given a query emoji, the goal is to retrieve a ranked list of documents considered relevant due to their text or image content. As this corresponds to a more classical retrieval problem, we report results in mean Average Precision (mAP) across all single emoji queries

$$\text{mAP} = \frac{1}{Q} \sum_i \frac{\sum_k^N \text{Prec}(j) \times y_k^i}{\sum_j y_k^i} \quad (5)$$

where  $Q$  is the number of single emoji queries,  $N$  is the number of samples, and  $y_k^i$  corresponds to the relevancy of query  $i$  to the  $k$ th ranked sample.

## IV. EMOJI PREDICTION

### A. Baselines

1) *Text-to-Emoji*: Our baseline text model consists of a bi-directional LSTM, which processes the text both in standard order and reverse order, on top of a word embedding layer [25]. LSTMs use their memory to help emphasize relevant information [17], but there is still a degradation of information propagation. The bi-directional nature of the LSTM helps to combat this effect and ensure that information from the beginning of the sentence isn't lost in the representation.

Words are placed in a vector embedding space, passed through our bi-directional LSTM layers, and the resultant representations are combined and fed to a softmax layer that attempts to predict relevant emoji. A softmax is used despite the multi-label nature of some of the data, because the majority of tweets contain only one label and those that do not may bear some relationship between labels. Text from the Twemoji dataset is tokenized and used to train the model. The validation set is used to determine after how many epochs to stop training (to avoid overfitting).

2) *Image-to-Emoji*: Similar to the approach for text-based prediction, we can also train a model for image-to-emoji prediction using our data. We use a CNN to represent images accompanying tweets. It is a GoogLeNet architecture trained to predict 13K ImageNet classes [24], [40]. We use the representation yielded at the penultimate layer for our image input. We train a single softmax layer on top of this representation with emoji prediction as the objective, with the weights prior to this softmax frozen. An end-to-end convolutional model could also be trained with sufficient training data, but it would be difficult to amass the requisite number of training samples, particularly for the longtail of the emoji usage distribution.

3) *Fusion*: For the combination of both text and image modalities, a late fusion approach is used. As both the text-based neural network and the image-based convolutional network output emoji confidence scores in a softmax layer, their format is directly comparable. Given confidence scores  $p_{txt}(y|x_{txt})$  predicting the likelihood of a given emoji  $y$  for some text  $x_{txt}$  and the corresponding scores  $p_{img}(y|x_{img})$  for some image  $x_{img}$ ,



TABLE III  
RESULTS FOR TEXT-BASED EMOJI PREDICTION

Training	Dataset	Top-1	Top-5	Top-10	Top-100	msAP
Balanced	Full	13.0	30.0	41.0	84.0	19.4
Full	Full	<b>23.7</b>	<b>47.5</b>	<b>59.9</b>	<b>94.9</b>	<b>32.6</b>
Balanced	Balanced	<b>35.1</b>	<b>48.3</b>	<b>54.7</b>	<b>87.7</b>	<b>35.1</b>
Full	Balanced	26.4	40.8	48.8	74.5	24.5

We report the model trained with the balanced test set as a goal, as well as trained toward the skewed distribution. the balanced sampling scheme outperforms the raw distribution on the balanced testset



Fig. 6. Examples of the hardest emoji to predict (red), the easiest (green), and those in between. Ambiguous faces are difficult to predict, while emoji tied concretely to an event, object, or place tend to be the easiest.

we give a combined prediction:

$$p(y|x_{txt}, x_{vis}) = \alpha p_{txt}(y|x_{txt}) + (1 - \alpha)p_{img}(y|x_{img}) \quad (6)$$

where  $\alpha$  is a modality weighting parameter in the range  $[0, 1]$  which is determined through validation.

## B. Results

1) *Text-to-Emoji*: The results for prediction on the Twemoji test sets are shown in Table III. Fig. 6 gives examples of those emoji the baseline models find difficult or easy to predict. We see that some of the most difficult emoji to predict include ambiguous face emoji where no clear emotion is displayed. Among the easiest emoji to predict are flag emoji and emoji tied closely to particular events, such as Christmas or birthdays. We also see less obvious emoji such as 📺 included. This is likely due to the resemblance of 📺 to a recording symbol on a video camera, as it is often used in conjunction with tweets containing links to video. It is likely this co-occurrence that makes it a particularly easy emoji to predict. Such usage underscores the necessity of using real world emoji usage where possible, as the unicode name for 📺 is merely ‘Large Red Circle’ which gives little to relate it to video.

It is worth noting that the numbers here reflect accuracy on predicting the emoji that *were* used, which are not necessarily all the emoji which could have been used. It is likely that some emoji were predicted which could be argued as relevant but which happened to not be the particular emoji the Twitter user selected. While the results should be considered indicative,

the annotations used cannot be considered absolute due to the subjectivity of emoji.

In Table III, we report the performance of the proposed model aimed at the balanced test set, as well as the same model trained on the unbalanced distribution. We note that the balanced sampling model performs much stronger on the balanced dataset. This is expected, as we targeted a balanced distribution during training, due to the assumption that some amount of the data bias was due to intrinsic bias in input interfaces. While we target a balanced distribution, the model can also be trained without balanced sampling to learn the skewed distribution. The model, when trained without balanced sampling, achieves higher performance on the full, unbalanced test set. From a practical standpoint, this is a far less interesting result due to the heavy skew in data. While this greatly improves the performance on the raw test set, the performance on the balanced subset diminishes significantly. We restrict all further discussion to only models that have been trained with a balanced sampling regime.

2) *Image-to-Emoji*: As described previously, we train a model to predict emoji based on CNN representations of images. In the top section of Table IV, we present the results of the image-trained model on the available image-bearing test set. We also present results for testing the text-trained model on this subset. We see that the image modality is competitive to the text modality for the prediction of emoji. This suggests that the emoji may often be as related to the images as they are to the text content. Overall, the performance of the models is broadly similar to those on the full Twemoji dataset, which is encouraging. It suggests that the relationship between the input data and the annotation is not too dissimilar to the whole set in this subset.








Table V gives some qualitative examples of results for emoji prediction on image and text inputs, along with the ground truth emoji annotation. Example C captures the food aspect of the image which is missed in the text modality, but neither are able to predict the true emoji. This is an example where the information contained in the emoji modality is mostly orthogonal to that in the text or image. We see in example F that the text-based prediction is led astray by the mention of food while the image-based method focuses on the emotional reaction expected from cuddling animals. The correct emoji, 🐾, appears in the top 100 results for the image-based baseline, while it is in the 400 s for the text modality. Some examples are easily handled by both the text and image modalities, such as A – this may be due to a strong association between the 🔥 emoji and sneaker enthusiasts. Example B is an interesting one, because both the image and the text contained the context of artwork, but the image was able to retrieve the artwork’s content and associate it with the correct emoji 🎨 while that content was not available in the text.

3) *Fusion*: In the bottom of Table IV, we provide scores for a fusion of both the image and text modalities. We see a significant improvement across most metrics through the fusion of both modalities, which tells us that they have complementary information. Though this could be an artifact of the representations used in either modality, it is reasonable to assume that the semantics of the emoji are not strictly tied to either modality, which is evidence that emoji should be considered as a modal-

TABLE IV  
RESULTS OF THE CNN-BASED IMAGE-INPUT MODEL AND THE BI-DIRECTIONAL LSTM TEXT-INPUT MODEL ON TWEMOJI-IMAGES,  
AS WELL AS THE FUSION OF THE TWO

	Model	Top-1	Top-5	Top-10	Top-100	msAP
Single Modality	Image only	14.7	33.0	44.0	86.4	17.0
	LSTM (Text Input)	17.7	33.5	43.4	81.3	22.3
Fusion	Image + LSTM ( $\alpha = 0.6$ )	<b>20.6</b>	<b>40.3</b>	<b>51.5</b>	<b>89.3</b>	<b>27.0</b>

TABLE V  
EXAMPLES OF TEXT-TO-EMOJI AND IMAGE-TO-EMOJI PREDICTION RESULTS ON THE TWEMOJI-IMAGES TEST SET

	Image	Text	Image-only	Text-only	True Emoji
A		rt U : nah this neymar x jordan collab is pure heat	🔥 🌑 🌑 🤩 🌟	🔥 🌑 ❤️ 📷 🤩	🔥
B		rt U : one of the short poetry i have done , #watercolor #art	🌟 🤩 🎧 🎨 🤩	😂 🤩 🖋️ 🏠 🌲	🌟
C		thank you	👉 🤩 🤩 🤩 🙌	😘 🤩 ❤️ 🤩 ❤️	📺 📺
D		turned my ghetto concrete workshop room into my own cool little space	😂 🤩 🤩 🤩 🌑	🔥 🤩 🎹 📺 🗣️ 🙌	😎
E		no one will ever understand what it's like to have a best friend like this so lucky i am U	❤️ ❤️ 🤩 🎉 🤩	❤️ 🤩 🎵 ❤️ 📢	❤️ 🤩 🤩
F		not food	🤩 ❤️ 🤩 🤩 🤩	👉 🤩 🍷 🍔 🍷	🐝
G		im that weird girl that likes to hold snakes	🤩 🤩 ❤️ 🤩 🤩	🤩 📺 🤩 100 🏪	🤩

We display the top 5 highest scoring emoji for a given sample. Sometimes images or text capture important predictive content that isn't present in the other modality, and sometimes both modalities fail to yield the expected emoji. Often, most of the suggested emoji seem reasonable from a subjective standpoint, which suggests that perfection on the evaluation metrics is not required for useful models

ity in their own right. In Fig. 7, we show the per-sample mAP (ranking emoji given an image + text input) performance as a function of the fusion weighting parameter  $\alpha$ . We see that the curve hits its peak near the center, with a skew toward the text input. This suggests a slightly stronger correlation between the emoji modality and text than between emoji and images.

In Fig. 8, we report the per-class difference in the msAP metric. This difference is calculated by subtracting the image-based performance from the text-based performance. A value of 0.0 would therefore mean that both methods performed identically well (or poorly), a positive value indicates that the text-based model performed better, and a negative indicates that the image-based model performed better. A strong bias toward the text-based approach is observed across almost all emoji. It is impossible to say whether this reflects the strength of cross-modal affinities, but it does tell us that the model we use for relating text to emoji is stronger than that for images.

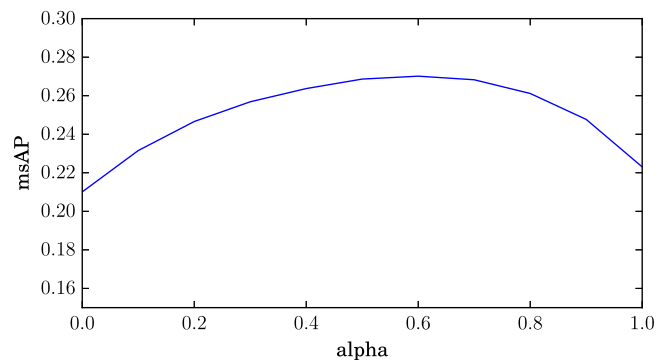


Fig. 7. Effect of modality-weighting parameter  $\alpha$  on the prediction of Twemoji-Images, measured in mean samplewise Average Precision.  $\alpha = 1.0$  corresponds to using only the text predictions, while  $\alpha = 0.0$  uses only image predictions. Peak performance occurs near  $\alpha = 0.6$ . Improvement through the combination of both modalities tells us that the modalities have complementary information for the prediction of emoji.

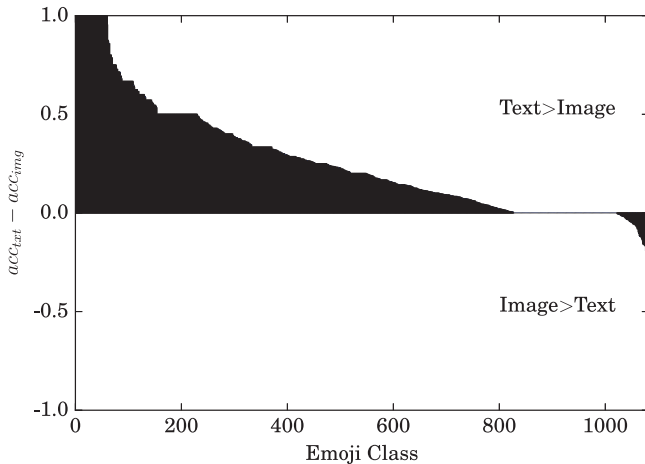


Fig. 8. Per-class performance difference between text and image modalities. This graph shows the difference in Top-5 accuracy between using only the text input modality to predict emoji and using only the image input modality. For roughly 80% of the emoji, text outperforms images for our dataset and baselines.

## V. EMOJI ANTICIPATION

### A. Baselines

1) *Text- and/or Image-to-Emoji*: Word embeddings have been used for the task of zero-shot image classification as a means to transfer knowledge from one class to another [31]. To place an emoji within this embedding space without the need for training examples, a short textual description of the emoji can be used as its representation.

We utilize a word2vec representation [26] that is pre-trained on a corpus of millions of lines of text accompanying Flickr photos [41]. Input modalities are then embedded in this shared space, where relationships between items are evaluated by their similarity in the space. Text terms are placed directly in the space through vocabulary look-up, as the embedding is originally trained on text. In the case of images, the names of the 15 highest scoring visual concepts are used, weighted by their confidence scores. We use 13K visual concept scores that come from the same GoogLeNet-style CNN used to extract high level features in the supervised setting.

To place the emoji modality within this mutual vector space, we use text terms extracted from the unicode-specified emoji title and descriptions. Emoji are unicode characters, and the details of their illustration are left to the implementation of the platform which incorporates them. However, when new emoji are accepted into the unicode specification, they are presented with a title and description. As there are generally only a few terms in the unicode metadata (see Table I for examples), we take the averaged word2vec vector representation of the words in this specification as a vector representative of that emoji within our space.

For emoji prediction using a fusion of text and image inputs, we use a simple weighted late fusion approach in the manner described in the previous section. Because we don't have any validation (or training) data in the unfamiliar emoji setting, the weighting parameter  $\alpha$  cannot be experimentally determined.

TABLE VI  
EMOJI ANTICIPATION RESULTS, REPORTED ON TWEMOJI-IMAGES

Model	Top-1	Top-5	Top-10	Top-100	msAP
Random	0.0	0.4	0.9	8.1	0.5
Zero-shot Text	1.1	2.5	3.9	20.9	1.9
Zero-shot Images	1.3	3.0	4.3	21.4	2.1
Fusion ( $\alpha = 0.5$ )	<b>1.5</b>	<b>3.8</b>	<b>5.7</b>	<b>23.8</b>	<b>2.5</b>

Emoji are predicted without any direct supervision data, analogous to what must be done when new emoji are released. We see improvement across all metrics when a fusion of the input modalities is used.

Instead, we assign  $\alpha = 0.5$ , giving both text and visual modalities equal priority in our model.

### B. Results





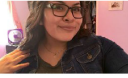










In Table VI we give results for emoji prediction on the Twemoji-Images dataset using only the text modality, only the image modality, and the fusion of the two (using  $\alpha = 0.5$ ). We observe that, as would be expected, the overall scores are much lower than the supervised approaches in the previous section. Though the results are small, they are significantly above random. The top-1 accuracy of random guesses on the Twemoji-Images test set is on the order of 0.08% compared with 1.5% for the fusion of the zero-shot results.

A surprising result is that the image modality actually outperforms the text modality in most of the metrics. Because the semantic space is learned on textual data, one might expect the text modality to be the most reliably embedded modality within the shared space, but that does not seem to be the case. Perhaps this is a result of many distracting terms in the textual data, which supervised approaches learn to filter out. Meanwhile, the limited vocabulary of the CNN concepts are likely to be a strong signal. Nonetheless, the fusion of the two modalities improves performance across all metrics.

The names of emoji may be reasonable, but might not capture unexpected uses. For example, *fireworks* 🎆 could be used for 'north star' or 'sun' based solely on its particular illustration here – usages that would be unlikely to be captured based on the title alone. Similarly, *ghost* 👻 has an especially friendly illustration, with the spectre appearing to wave hello. Such usage based on the visual appearance can easily diverge relative to the drier, more descriptive title.

The performance of this baseline approach can likely be improved by focusing on improving the quality of the mapping of the three modalities to the mutual space. The embedding of emoji, for example, could likely be improved by manually specifying additional relevant text terms. The terms contained in the unicode specification focus on being descriptive about the emoji, focusing on *what* it is, rather than how it might be used. Though difficult to experimentally evaluate in an objective manner, adding some extra terms based on postulated usage to the emoji representation could be one way to boost performance without significant extra effort. For example, ▶ has the title “black right-pointing triangle”, which is a description of what the emoji is but says little about how it might be used. Adding potentially related

TABLE VII  
TOP RANKED DOCUMENTS FOR THREE EMOJI QUERIES

Query:			
1	 you can't imagine how much i miss you #facetimemenash	 rt U : glasses ... no glasses ... glasses	 graduation part N : my favorite fish in the sea
2	 rt U : so sad #orlando #rip	 rt U : the bigger the better when it comes to eyewear ! by U . london	 rt U : it's a fishy kinda day ... fish platter and salmon & smoked UNKNOWN fish cakes
3	 rt U : this is so sad #prayforturkey	 rt U : glasses or no glasses	 N days to go ! just keep swimming swimming swimming UNKNOWN
4	 my heart goes out to the families and friends who lost their loved ones terrible and sad news ! #istanbul	 glasses	 rt U : i found dory

We see a correspondence between the baseline's prediction of certain emoji and current events, with relationships between *finding dory* and the tropical fish emoji, as well as sad current events and the pensive face emoji. Non-relevant results, like those for eyeglasses, may appear subjectively to be relevant but there is clearly a nuance in the usage of the eyeglass emoji that is being overlooked

TABLE VIII  
QUERY-BY-EMOJI RESULTS FOR BOTH SUPERVISED AND ZERO-SHOT BASELINES

Method	Twemoji (Full)	Twemoji (Balanced)	Twemoji (Images)
Random	0.1	0.3	0.2
LSTM (Text)	19.3	35.5	20.2
CNN (Image)	–	–	<b>22.0</b>
Fusion	–	–	21.2
Zero-shot Text	0.5	2.0	1.5
Zero-shot Images	–	–	0.8
Zero-shot Fusion	–	–	<b>1.3</b>

Results are reported in percentage mAP. In the supervised setting, we find the images to slightly outperform the text, but in the zero-shot setting the performance is reversed.

terms such as *next* or *play* or *therefore* might capture probable usage semantics that are absent in a pure description of the emoji itself. Indeed, due to the particular illustration of this emoji, the term *black* in the description is actually misleading as there is nothing black about the right-pointing triangle in this rendering.

## VI. QUERY-BY-EMOJI


### A. Baselines

The baselines in previous sections give normalized scores across possible emoji given the input modalities. By calculating these normalized scores for all documents, we are able to rank the documents in order of predicted relevancy to a given emoji query. In this way, we can then perform retrieval per-emoji

across these documents. All results in this section are therefore produced by applying the baseline models described in the previous sections to all documents within the test database, and performing retrieval based on per-emoji class scores.

### B. Results

Table VIII gives results for the Query-by-Emoji task. Surprisingly, we see that retrieving tweets using only the supervised image understanding slightly outperforms both text-only and the fusion of the two. This result is markedly different from the emoji prediction task where text outperformed images. This could possibly be the result of a very strong correlation within high probability image-emoji pairs.

In Table VII, some qualitative query-by-emoji results are shown. We observe strong signals for correlations with current events that occurred during the data collection period of the dataset. Tragic events occurred during this period in both Orlando and Turkey, and the model picked up a strong relationship between the “pensive face”  and these topics. Similarly, the movie *Finding Dory* was released during this time, and we see it present in the high-ranked predictions for the tropical fish. The exploitation and mapping of these emoji-event relationships presents interesting avenues for future research.

For the eyeglasses emoji, the top-ranked results from our baseline model did not contain the eyeglasses emoji. The top four results all contain glasses in the image and a mention of ‘glasses’ or ‘eyewear’ in the text, but the authors opted for alternative emoji during composition. While these results undoubtedly have a level of subjective relevance, the authors

clearly felt that other emoji were called for. Perhaps the eyeglass emoji is considered too redundant when the content is already contained in both the text and images. Learning to identify and exploit these subtle distinctions is an open problem for future, improved models.

## VII. CONCLUSION

In this paper, we have approached emoji as a modality distinct from text and images. There is sufficient motivation for doing so, and considerable future opportunities for research and applications with the emoji modality. We have proposed a large scale dataset of real-world emoji usage, containing the semantic relationships between emoji and text as well as emoji and images. We have defined three challenge tasks with evaluation on this dataset, and provided baseline results for all three. We have looked at the problem of predicting emoji from text and/or images, both with the use of ample training data and in the absence of any. We have also looked at the problem of using emoji as queries for cross-modal retrieval. Emoji are everywhere, and are becoming more pervasive. They already possess a distinct semantic space that can be utilized as a strong information signal as well as a novel means of interaction with data, through both query-by-emoji as well as emoji summarization of content. Furthermore, their semantic richness will only increase as new emoji continue to be introduced. It is our hope that this work and the challenge tasks defined within will spur further research and understanding of emoji within the multimedia community.

## REFERENCES

- [1] W. Ai *et al.*, "Untangling emoji popularity through semantic embeddings," in *Proc. 11th Int. AAAI Conf. Web Social Media*, 2017, pp. 2–11.
- [2] L. M. Aiello *et al.*, "Sensing trending topics in twitter," *IEEE Trans. Multimedia*, vol. 15, no. 6, pp. 1268–1282, Oct. 2013.
- [3] Z. Akata, F. Perronnin, Z. Harchaoui, and C. Schmid, "Label-embedding for image classification," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 7, pp. 1425–1438, Jul. 2016.
- [4] F. Barbieri, M. Ballesteros, and H. Saggion, "Are emojis predictable? in *Proc. 15th Conf. Eur. Chapter Assoc. Comput. Linguistics: vol. 2, Short Papers*, 2017, pp. 105–111.
- [5] F. Barbieri, G. Kruszewski, F. Ronzano, and H. Saggion, "How cosmopolitan are emojis?: Exploring emojis usage and meaning over different languages with distributional semantics," in *Proc. ACM Conf. Multimedia*, 2016, pp. 531–535.
- [6] F. Barbieri, F. Ronzano, and H. Saggion, "What does this emoji mean? a vector space skip-gram model for twitter emojis," in *Proc. Lang. Resources Eval. Conf.*, 2016.
- [7] J. Berengueres and D. Castro, "Sentiment perception of readers and writers in emoji use," 2017, arXiv:1710.00888.
- [8] S. Cappallo, T. Mensink, and C. G. M. Snoek, "Image2emoji: Zero-shot emoji prediction for visual media," in *Proc. ACM Conf. Multimedia*, 2015, pp. 1311–1314.
- [9] S. Cappallo, T. Mensink, and C. G. M. Snoek, "Query-by-emoji video search," in *Proc. ACM Conf. Multimedia*, 2015, pp. 735–736.
- [10] J. Chen, Y. Cui, G. Ye, D. Liu, and S.-F. Chang, "Event-driven semantic concept discovery by exploiting weakly tagged internet images," in *Proc. Int. Conf. Multimedia Retrieval*, 2014, pp. 1–8.
- [11] Z. Chen *et al.*, "Through a gender lens: An empirical study of emoji usage over large-scale android users," 2017, arXiv:1705.05546.
- [12] G. Donato and P. Paggio, "Investigating redundancy in emoji use: Study on a twitter based corpus," in *Proc. 8th Workshop Comput. Approaches Subjectivity, Sentiment Social Media Anal.*, 2017, pp. 118–126.
- [13] A. El Ali, T. Wallbaum, M. Wasmann, W. Heuten, and S. C. Boll, "Face2emoji: Using facial emotional expressions to filter emojis," in *Proc. SIGCHI Conf. Human Factors Comput. Syst.*, 2017, pp. 1577–1584.
- [14] B. Felbo, A. Mislove, A. Søgaard, I. Rahwan, and S. Lehmann, "Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm," in *Proc. Conf. Empirical Methods Natural Language Process.*, 2017, pp. 1615–1625.
- [15] Y. Gao and Q. Dai, "View-based 3d object retrieval: Challenges and approaches," *IEEE MultiMedia*, vol. 21, no. 3, pp. 52–57, Jul./Sep. 2014.
- [16] B. Guthrie, K. Ho, and A. El Saddik, "Language-independent data set annotation for machine learning-based sentiment analysis," in *Proc. IEEE Int. Conf. Syst., Man, Cybern.*, 2017, pp. 2105–2110.
- [17] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.
- [18] L. Jiang, T. Mitamura, S.-I. Yu, and A. G. Hauptmann, "Zero-example event search using multimodal pseudo relevance feedback," in *Proc. Int. Conf. Multimedia Retrieval*, 2014, p. 297.
- [19] J. Lee and S. C. Marsella, "Predicting speaker head nods and the effects of affective information," *IEEE Trans. Multimedia*, vol. 12, no. 6, pp. 552–562, Oct. 2010.
- [20] X. Li, R. Yan, and M. Zhang, "Joint emoji classification and embedding learning," in *Proc. Asia-Pacific Web Web-Age Inf. Manage. Joint Conf. Web Big Data*, 2017, pp. 48–63.
- [21] N. Ljubešić and D. Fišer, "A global analysis of emoji usage," in *Proc. 10th Web as Corpus Workshop*, 2016, pp. 82–89.
- [22] R. P. López and F. Cap, "Did you ever read about frogs drinking coffee? investigating the compositionality of multi-emoji expressions," in *Proc. 8th Workshop Comput. Approaches Subjectivity, Sentiment Social Media Anal.*, 2017, pp. 113–117.
- [23] X. Lu *et al.*, "Learning from the ubiquitous language: An empirical analysis of emoji usage of smartphone users," in *Proc. ACM Int. Joint Conf. Pervasive Ubiquitous Comput.*, 2016, pp. 770–780.
- [24] P. Mettes, D. C. Koelma, and C. G. M. Snoek, "The imagenet shuffle: Reorganized pre-training for video event detection," in *Proc. Int. Conf. Multimedia Retrieval*, 2016, pp. 175–182.
- [25] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," 2013, arXiv:1301.3781.
- [26] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Proc. 26th Adv. Int. Conf. Neural Inf. Process. Syst.*, 2013, pp. 3111–3119.
- [27] H. Miller *et al.*, "blissfully happy or ready to fight: Varying interpretations of emoji," in *Proc. 10th Int. Conf. Web Social Media*, 2016, pp. 259–268.
- [28] H. J. Miller, D. Kluver, J. Thebault-Spieker, L. G. Terveen, and B. J. Hecht, "Understanding emoji ambiguity in context: The role of text in emoji-related miscommunication," in *Proc. 11th Int. Conf. Web Social Media*, 2017, pp. 152–161.
- [29] N. Na'aman, H. Provenza, and O. Montoya, "Varying linguistic purposes of emoji in (twitter) context," in *Proc. Assoc. Comput. Linguistics*, 2017, pp. 136–141.
- [30] K. Njenga, "Social media information security threats: Anthropomorphic emoji analysis on social engineering," in *Proc. IT Convergence Security*, 2017, pp. 185–192.
- [31] M. Norouzi *et al.*, "Zero-shot learning by convex combination of semantic embeddings," in *Proc. Int. Conf. Learn. Represent.*, 2014.
- [32] P. K. Novak, J. Smailović, B. Sluban, and I. Mozetič, "Sentiment of emojis," *PLoS one*, vol. 10, no. 12, 2015, Art. no. e0144296.
- [33] H. Pohl, C. Domin, and M. Rohs, "Beyond just text: Semantic emoji similarity modeling to support expressive communication," *ACM Trans. Comput.-Human Interact.*, vol. 24, no. 1, 2017, Art. no. 6.
- [34] M. Rathan, V. R. Hulipalled, K. Venugopal, and L. Patnaik, "Consumer insight mining: Aspect based twitter opinion mining of mobile phone reviews," *Appl. Soft Comput.*, vol. 68, pp. 765–773, 2017.
- [35] M. A. Riordan, "The communicative role of non-face emojis: Affect and disambiguation," *Comput. Human Behav.*, vol. 76, pp. 75–86, 2017.
- [36] D. Rodrigues, D. Lopes, M. Prada, D. Thompson, and M. V. Garrido, "A frown emoji can be worth a thousand words: Perceptions of emoji use in text messages exchanged between romantic partners," *Telematics Inf.*, vol. 34, pp. 1532–1543, 2017.
- [37] M. Shiha and S. Ayvaz, "The effects of emoji in sentiment analysis," *IJCEE*, vol. 9, pp. 360–369, 2017.
- [38] C. G. M. Snoek, M. Worrington, J. C. Van Gemert, J.-M. Geusebroek, and A. W. M. Smeulders, "The challenge problem for automated detection of 101 semantic concepts in multimedia," in *Proc. 14th ACM Int. Conf. Multimedia*, 2006, pp. 421–430.

- [39] D. Stowell, D. Giannoulis, E. Benetos, M. Lagrange, and M. D. Plumbley, "Detection and classification of acoustic scenes and events," *IEEE Trans. Multimedia*, vol. 17, no. 10, pp. 1733–1746, Oct. 2015.
- [40] C. Szegedy *et al.*, "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015, pp. 1–6.
- [41] B. Thomee *et al.*, "The new data and new challenges in multimedia research," 2015, arXiv:1503.01817.
- [42] G. W. Tigwell and D. R. Flatla, "Oh that's what you meant!: Reducing emoji misunderstanding," in *Proc. 18th Int. Conf. Human-Comput. Interact. With Mobile Devices Services Adjunct*, 2016, pp. 859–866.
- [43] S. Wang *et al.*, "A natural visible and infrared facial expression database for expression recognition and emotion inference," *IEEE Trans. Multimedia*, vol. 12, no. 7, pp. 682–691, Nov. 2010.
- [44] S. Wijeratne, L. Balasuriya, A. Sheth, and D. Doran, "Emojinet: Building a machine readable sense inventory for emoji," in *Proc. 8th Int. Conf. Social Inf.*, 2016, pp. 527–541.
- [45] S. Wijeratne, L. Balasuriya, A. Sheth, and D. Doran, "A semantics-based measure of emoji similarity," in *Proc. Int. Conf. Web Intell.*, 2017, pp. 646–653.
- [46] R. Zhou, J. Hentschel, and N. Kumar, "Goodbye text, hello emoji: Mobile communication on wechat in china," in *Proc. CHI Conf. Human Factors Comput. Syst.*, 2017, pp. 748–759.



**Spencer Cappallo** received the B.Sc. degree in physics (with honors) from the University of St. Andrews, St. Andrews, U.K., in 2006 and received the M.Sc. degree with distinction in scientific computing from the University of Nottingham, Nottingham, U.K., in 2013. He received the Ph.D. degree in computer vision and machine learning from the University of Amsterdam, Amsterdam, The Netherlands, in 2018. His research interests include multimedia and reinforcement learning, with a particular interest in video understanding, weak supervision, and learning to predict subjective qualities of data.



develop a visual similarity search engine with LookFlow, which Yahoo acquired in 2013.

**Stacey Svetlichnaya** received the B.S. degree in 2011 and the M.S. degree in 2012 in symbolic systems from Stanford University, Stanford, CA, USA. She is a Senior Research Engineer with the Oath AI team, San Francisco, CA, USA (previously Yahoo Vision & Machine Learning). Her recent work includes image aesthetic quality and style classification, object recognition, and photo caption generation. She has worked extensively on Flickr image search and data pipelines, as well as automating content discovery and recommendation. Prior to Flickr, she helped de-



Flickr, Tumblr, or Yahoo Mail. Prior to Yahoo, he cofounded IQ Engines, a startup that developed image recognition APIs and was acquired by Yahoo in 2013.

**Pierre Garrigues** received the Engineering degree from Ecole Polytechnique, Palaiseau, France, in 2003, and the Ph.D. degree from the Department of Electrical Engineering and Computer Sciences, University of California at Berkeley, Berkeley, CA, USA, in 2009. He is a Senior Director of Research with Yahoo. He leads a team of scientists and engineers focused on machine learning and artificial intelligence, with primary applications in computer vision and natural language processing. His work has been deployed in multiple Yahoo products such as



His research interests include machine learning and computer vision, with special interests in learning representations for image classification, zero-example prediction, and for 3-D understanding. Prof. Mensink is a recipient of the ACM Multimedia Best Paper award (2014), a prestigious NWO VENI Grant (2015), and the ACM ICMR Best Paper Award (2016).

**Thomas Mensink** received the M.Sc. degree (*cum laude*) in artificial intelligence from the University of Amsterdam, Amsterdam, The Netherlands, in 2007. He received the Ph.D. degree in computer science in 2012 from the University of Grenoble, Grenoble, France, while working both at the LEAR team of INRIA Grenoble (now THOTH) and at Xerox Research Centre Europe (now Naver Labs). Since 2012, he has been with the University of Amsterdam, first as a Postdoctoral Researcher (2012–2017) and is currently an Assistant Professor in 3-D Deep Learning.



of Amsterdam, as well as visiting Scientist with Carnegie Mellon University and UC Berkeley, the Head of R&D, University spin-off Euvision Technologies, and Managing Principal Engineer with Qualcomm Research Europe. He is the Lead Researcher of the award-winning MediaMill Semantic Video Search Engine, which is the most consistent top performer in the yearly NIST TRECVID evaluations. He has published more than 200 refereed journal and conference papers, and frequently serves as an area chair of the major conferences in multimedia and computer vision. His research interests focus on video and image recognition.

Prof. Snoek was general chair of ACM Multimedia 2016 in Amsterdam, founder of the VideOlympics 2007–2009, and a member of the editorial board for *ACM Transactions on Multimedia*. He is the recipient of an NWO Veni Award, a Fulbright Junior Scholarship, an NWO Vidi Award, and the Netherlands Prize for ICT Research. Several of his Ph.D. students and Post-docs have won awards, including the IEEE TRANSACTIONS ON MULTIMEDIA Prize Paper Award, the SIGMM Best Ph.D. Thesis Award, the Best Paper Award of ACM Multimedia, an NWO Veni Award, and the Best Paper Award of ACM Multimedia Retrieval. Five of his former mentees serve as Assistant and Associate Professors.

**Cees G. M. Snoek** received the M.Sc. degree in business information systems in 2000 and the Ph.D. degree in computer science in 2005 both from the University of Amsterdam, Amsterdam, The Netherlands. He is a Full Professor in computer science with the University of Amsterdam, where he heads the Intelligent Sensory Information Systems Lab. He is also the Director of the QUVA Lab, the joint research lab of Qualcomm and the University of Amsterdam on deep learning and computer vision. He was previously an Assistant and Associate Professor with the University