

# Ethics Sheet for Automatic Emotion Recognition and Sentiment Analysis

Saif M. Mohammad  
National Research Council Canada  
saif.mohammad@nrc-cnrc.gc.ca

*The importance and pervasiveness of emotions in our lives makes affective computing a tremendously important and vibrant line of work. Systems for automatic emotion recognition (AER) and sentiment analysis can be facilitators of enormous progress (e.g., in improving public health and commerce) but also enablers of great harm (e.g., for suppressing dissidents and manipulating voters). Thus, it is imperative that the affective computing community actively engage with the ethical ramifications of their creations. In this article, I have synthesized and organized information from AI Ethics and Emotion Recognition literature to present fifty ethical considerations relevant to AER. Notably, this ethics sheet fleshes out assumptions hidden in how AER is commonly framed, and in the choices often made regarding the data, method, and evaluation. Special attention is paid to the implications of AER on privacy and social groups. Along the way, key recommendations are made for responsible AER. The objective of the ethics sheet is to facilitate and encourage more thoughtfulness on why to automate, how to automate, and how to judge success **well before** the building of AER systems. Additionally, the ethics sheet acts as a useful introductory document on emotion recognition (complementing survey articles).*

## 1. Introduction

Emotions play a central role in our lives. Thus, affective computing, which deals with emotions and computation (often through AI systems) is a tremendously important and vibrant line of work. It is a sweeping interdisciplinary area of study exploring both fundamental research questions (such as *what are emotions?*) and commercial applications (such as *can machines detect consumer sentiment?*).

In her seminal book, *Affective Computing*, Dr. Rosalind Picard described Automatic Emotion Recognition (AER) as: “giving emotional abilities to computers” (Picard 2000). Such systems can be incredibly powerful: facilitators of enormous progress, but also enablers of great harm. In fact, some of the recent commercial and governmental uses of emotion recognition have garnered considerable criticism, including: infringing on one’s privacy, exploiting vulnerable sub-populations, and also allegations of downright

---

Submission received: 16 September 2021; accepted for publication: 6 December 2021.

<https://doi.org/10.1162/coli.a.00433>

© 2022 Crown in Right of Canada

pseudo-science (Wakefield 2021; ARTICLE19 2021; Woensel and Nevil 2019). Even putting aside high-profile controversies, emotion recognition impacts people and thus entails ethical considerations (big and small). Thus, it is imperative that the AER community actively engage with the ethical ramifications of their creations.

This article, which I refer to as an *Ethics Sheet for AER*, is a critical reflection of this broad field of study with the aim of facilitating more responsible emotion research and appropriate use of the technology. As described in Mohammad (2021a), an Ethics Sheet for an AI Task is a semi-standardized document that synthesizes and organizes information from AI Ethics and AI Task literature to present a comprehensive array of ethical considerations for that task. Thus, in some ways, an ethics sheet is similar to survey articles, except here the focus is on ethical considerations. It:

- Fleshes out assumptions hidden in how the task is framed, and in the choices often made regarding the data, method, and evaluation.
- Presents ethical considerations unique or especially relevant to the task.
- Presents how common ethical considerations manifest in the task.
- Presents relevant dimensions and choice points; along with tradeoffs.
- Lists common harm mitigation strategies.
- Communicates societal implications of AI systems to researchers, developers, and the broader community.

The sheet should flesh out various ethical considerations that apply at the level of the task. It should also flesh out ethical consideration of common theories, methodologies, resources, and practices used in building AI systems for the task. A good ethics sheet will question some of the assumptions that often go unsaid.

**Primary motivation for creating an Ethics Sheet for AER:** To provide a go-to point for a carefully compiled substantive engagement with the ethical issues relevant to emotion recognition, going beyond individual systems and datasets and drawing on knowledge from a large body of past work. The document will be useful to anyone who wants to build or use emotion recognition systems/algorithms for research or commercial purposes. Specifically, the main benefits can be summarized by the list below:

1. Encourages more thoughtfulness on why to automate, how to automate, and how to judge success **well before** the building of AER systems.
2. Helps us better navigate research and implementation choices.
3. Moves us toward consensus and standards.
4. Helps in developing better post-production documents such as datasheets and model cards.
5. Has citations and pointers; acts as a jumping off point for further reading.
6. Helps engage the various stakeholders of an AI task with each other. Helps stakeholders challenge assumptions made by researchers and developers. Helps develop harm mitigation strategies.

7. Acts as a useful introductory document on emotion recognition (complements survey articles).

Note that even though this sheet is focused on AER, many of the ethical considerations apply broadly to natural language tasks in general. Thus, it can serve as a useful template to build ethics sheets for other tasks.

**Target audience:** The primary audience for this sheet are researchers, engineers, developers, and educators from various fields (especially NLP, ML, AI, data science, public health, psychology, and digital humanities) who build, make use of, or teach about AER technologies; however, much of the discussion should be accessible to various other stakeholders of AER as well, including policy/decision makers, and those who are impacted by AER. I hope also that this sheet will act as a springboard for the creation of a sheet where non-technical stakeholders are the primary audience.

**Process:** My own research interests are at the intersection of emotions and language—to understand how we use language to express our feelings. I created this sheet to gather and organize my thoughts around responsible emotion recognition research, and I am hopeful that it is of use to others as well. Discussions with various scholars from computer science, psychology, linguistics, neuroscience, and social sciences (and their comments on earlier drafts) have helped shape this sheet. An earlier draft of this material was also posted as a blog post with an explicit invitation for feedback. Valuable insights from the community were then incorporated into this document. That said, it should be noted that I do not speak for the AER community. There is no “objective” or “correct” ethics sheet. This sheet should be taken as one perspective among many in the community. I welcome dissenting views and encourage further discussion. These can lead to periodically revised or new ethics sheets. As stated in Mohammad (2021a):

*Multiple ethics sheets can be created (by different teams and approaches) to reflect multiple perspectives, viewpoints, and what is important to different groups of people. We should be wary of the world with single authoritative ethics sheets per task and no dissenting voices.*

The rest of the article is organized as follows: Section 2 is a preface to the ethics sheet, Section 3 presents the ethics sheet for AER (50 considerations), and this is followed by summarizing thoughts in Section 4. The Appendix compiles a list of succinct recommendations for responsible AER (drawn from the discussions on ethical considerations in Section 3).

## 2. Preface for the Ethics Sheet on AER

Let us consider a few rapid-fire questions to set the context. A good ethics sheet makes us question our assumptions. So let us start at the top:

**Q1. Should we be building AI systems for automatic emotion recognition? Is it ethical to do so?**

**A.** This is a good question. This sheet will not explicitly answer the question, but it will help in clarifying and thinking about it. This sheet will sometimes suggest that certain applications in certain contexts are good or bad ideas, but largely it will discuss what are the various considerations to be taken into account: whether to build or use a particular

system, how to build or use a particular system, what is more appropriate for a given context, how to assess success, and so forth.

The above question is also somewhat under-specified. We first need to clarify...

## **Q2. What does automatic emotion recognition mean?**

**A.** Emotion recognition can mean many things, and it has many forms. (This sheet will get into that.) Emotion recognition can be deployed in many contexts. For example, many will consider automated insurance premium decisions based on inferred emotions to be inappropriate. However, studying how people use language to express gratitude, sadness, and so on, is considered okay in many contexts. A human-computer interaction system benefits from being able to identify which utterances can convey anger, joy, sadness, hate, and so forth. (Not having such capabilities will lead to offensive, unempathetic, and inappropriate interactions.) Many other contexts are described in the sheet.

## **Q3. Can machines infer one's true emotional state ever?**

**A.** No. (This sheet will get into that.)

## **Q4. Can machines infer some small aspect of people's emotions (or emotions that they are trying to convey) in some contexts, to the extent that it is \*useful\*?**

**A.** In my view, yes. In a limited way, this is analogous to machine translation or Web search. The machine does not understand language, nor does it understand what the user really wants, nor the social, cultural, or embodied context, but it is able to produce a somewhat useful translation or search result with some likelihood; and it produces some amount of inappropriate and harmful results with some likelihood. However, unlike machine translation or search, emotions are much more personal, private, and complex. People cannot fully determine each other's emotions. People cannot fully determine their own emotional state. But we make do with our limitations and infer emotions as best we can to function socially. We also have moral and ethical failures. We cause harm because of our limitations, and we harbor stereotypes and biases.

If machines are to be a part of this world and interact with people in any useful and respectful way, then they must have at least some limited emotion recognition capabilities; and thereby will also cause some amount of harm. Thus, if we use them, it is important that we are aware of the limitations; design systems that protect and empower those without power; deploy them in the contexts they are designed for; use them to assist human decision making; and work to mitigate the harms they will perpetrate. We need to hold AER systems to high standards, not just because it is a nice aspirational goal, but because machines impact people at scale (in ways that individuals rarely can) and emotions define who we are (in ways that other attributes rarely do). I hope this sheet is useful in that regard.

## **3. Main Sheet (version 1.0)**

This ethics sheet for automatic emotion recognition has four sections: Modalities and Scope, Task, Applications, and Ethical Considerations. The first three are brief and set the context. The fourth presents various ethical considerations of AER as a numbered list, organized in thematic groups.

### 3.1 Modalities and Scope

**Modalities:** Work on AER has used a number of modalities (sources of input), including:

- Facial expressions, gait, proprioceptive data (movement of body), gestures
- Skin and blood conductance, blood flow, respiration, infrared emanations
- Force of touch, haptic data (from sensors of force)
- Speech, language (especially written text, emoticons, emojis)

All of these modalities come with benefits, potential harms, and ethical considerations.

**Scope:** This sheet will focus on AER from written text and AER in NLP, but several of the listed considerations apply to AER in general (regardless of modality, and regardless of field such as NLP or Computer Vision).

### 3.2 Task

Automatic emotion recognition from one's utterances (written or spoken) is a broad umbrella term used to refer to a number of related tasks such as those listed below (note that each of these framings has ethical considerations and may be more or less appropriate for a given context):

1. Inferring emotions felt by the speaker (e.g., given Sara's tweet, what is Sara feeling?); inferring emotions of the speaker as perceived by the reader/listener (e.g., what does Li think Sara is feeling?); inferring emotions that the speaker is attempting to convey (e.g., what emotion is Sara trying to convey?). These may be correlated, but they can be different depending on the particular instance. The first framing "inferring emotions felt by the speaker" is fairly common in scientific literature, but also perhaps most often misused/misinterpreted. More on this in the Ethical Considerations section.
2. Inferring the intensity of the emotions discussed above.
3. Inferring patterns of speaker's emotions over long periods of time, across many utterances; including the inference of moods, emotion dynamics, and emotional arcs (e.g., tracking character emotion arcs in novels and tracking impact of health interventions on a patient's well-being).
4. Inferring speaker's emotions/attitudes/sentiment toward a target product, movie, person, idea, policy, entity, and so on (e.g., does Sara like the new phone?).
5. Inferring emotions evoked in the reader/listener (e.g., what feelings arise in Li on reading Sara's tweet?). This may be different among different readers because of their past experiences, personalities, and world-views: For example, the same text may evoke different feelings among people with opposing views on an issue.
6. Inferring emotions of people mentioned in the text (e.g., given a tweet that mentions Moe, what emotional state of Moe is conveyed in the tweet?).

7. Inferring emotionality of language used in text (regardless of whose emotions) (e.g., is the tweet about happy things, angry feelings, ...?).
8. Inferring how language is used to convey emotions such as joy, sadness, loneliness, hate, and so forth.
9. Inferring the emotional impact of sarcasm, metaphor, idiomatic expression, dehumanizing utterance, hate speech, and so on.

**Note 1:** The term Sentiment Analysis is commonly used to refer to the task described in bullet 4, especially in the context of product reviews (sentiment is commonly labeled as positive negative, or neutral) (Turney 2002; Pang, Lee, and Vaithyanathan 2002). On the other hand, determining the predilection of a person toward a policy, party, issue, or similar, is usually referred to as Stance Detection, and involves classes such as favor and against (Mohammad et al. 2016; Mohammad, Sobhani, and Kiritchenko 2017).

**Note 2:** Many AER systems focus only on the emotionality of the language used (bullet 7), even though their stated goal might be one of the other bullets. This may be appropriate in restricted contexts such as customer reviews or personal diary blog posts, but not always. (More on this in §3.4.1 *Ethical Considerations: Task Design.*)

**Note 3:** There also exist tasks that focus not directly on emotions, but on associated phenomena, such as: whose emotions, who/what evoked the emotion, what types of human need was met or not met resulting in the emotion, and so on.

See these surveys for more details: Mohammad (2021b) examines emotions, sentiment, stance, and so forth; Zhang, Wang, and Liu (2018) focuses on sentiment analysis tasks; Soleymani et al. (2017) surveys multi-modal techniques for sentiment analysis.

### 3.3 Applications

The potential benefits of AER are substantial. Below is a sample of some existing applications. (Note that this is not an endorsement of these applications. All of the applications come with potential harms and ethical considerations. Use of AER by the military, for intelligence, and for education are especially controversial.)

- Public Health: Assist public health research projects, including those on loneliness (Guntuku et al. 2019; Kiritchenko et al. 2020), depression (De Choudhury et al. 2013; Resnik et al. 2015), suicidality prediction (MacAvaney et al. 2021), bipolar disorder (Karam et al. 2014), stress (Eichstaedt et al. 2015), and well-being (Schwartz et al. 2013).
- Commerce/Business: Track sentiment and emotions toward one's products, track reviews, blog posts, YouTube videos and comments; develop virtual assistants, writing assistants; help advertise products that one is more likely to be interested in.
- Government Policy and Public Health Policy: Track and document views of the broader public on a range of issues that impact policy (tracking amount of support and opposition, identify underlying issues and pain points, etc.). Governments and health organizations around the world are also interested in tracking how effective their messaging has been in response to crises such as pandemics and climate change.

- **Art and Literature:** Improve our understanding of what makes a compelling story, how do different types of characters interact, what are the emotional arcs of stories, what is the emotional signature of different genres, what makes well-rounded characters, why does art evoke emotions, how do the lyrics and music impact us emotionally, and so on. Can machines generate art (generate paintings, stories, music, etc.)?
- **Social Sciences, Neuroscience, Psychology:** Help answer questions about people. What makes people thrive? What makes us happy? What can our language tell us about our well-being? What can language tell us about how we construct emotions in our minds? How do we express emotions? How different are people in terms of what different emotion words mean to them and how they use emotional words?
- **Military, Policing, and Intelligence:** Track how sets of people or countries feel about a government or other entities (controversial); track misinformation on social media.

### 3.4 Ethical Considerations

The usual approach to building an AER system is to **design the task** (identify the process to be automated, the emotions of interest, etc.), compile appropriate **data** (label some of the data for emotions—a process referred to as human annotation), train ML models that capture patterns of emotional expression from the data—**the method**, and **evaluate** the models by examining their predictions on a held-out test set. There are ethical considerations associated with each step of this development process. Considerations for privacy and social groups are especially pertinent for AER and cut across task, design, data, and evaluation.

This section describes fifty considerations grouped under the themes: *Task Design, Data, Method, Impact and Evaluation*, and *Implications for Privacy and Social Groups*. First I present an outline of the considerations along with a summary for each grouping. This is followed by five sub-sections (§3.4.1 through §3.4.5) that present, in detail, the ethical considerations associated with the five groups.

#### I. TASK DESIGN

**Summary:** This section discusses various ethical considerations associated with the choices involved in the framing of the emotion task and the implications of automating the chosen task. Some important considerations include: Whether it is even possible to determine one's internal mental state. Whether it is ethical to determine such a private state. And, who is often left out in the design of existing AER systems. I discuss how it is important to consider which formulation of emotions is appropriate for a specific task/project, while avoiding careless endorsement of theories that suggest a mapping of external appearances to inner mental states.

##### A. Theoretical Foundations

1. Emotion Task and Framing
2. Emotion Model and Choice of Emotions
3. Meaning and Extra-Linguistic Information
4. Wellness and Emotion
5. Aggregate Level vs. Individual Level

##### B. Implications of Automation

6. Why Automate (Who Benefits; Will this Shift Power)
7. Embracing Neurodiversity
8. Participatory/Emancipatory Design

9. Applications, Dual use, Misuse
10. Disclosure of Automation

## II. DATA

Summary: This section has three themes: implications of using datasets of different kinds, the tension between human variability and machine normativeness, and the considerations regarding the people who have produced the data. Notably, I discuss how on the one hand is the tremendous variability in human mental representation and expression of emotions, and on the other hand, is the inherent bias of modern machine learning approaches to ignore variability. Thus, through their behavior (e.g., by recognizing some forms of emotion expression and not others), AI systems convey to the user what is “normal”; implicitly invalidating other forms of emotion expression.

### C. Why This Data

11. Types of Data
12. Dimensions of Data

### D. Human Variability vs. Machine Normativeness

13. Variability of Expression and Mental Representation
14. Norms of Emotion Expression
15. Norms of Attitudes
16. One “Right” Label or Many Appropriate Labels
17. Label Aggregation
18. Historical Data (Who is Missing and What are the Biases)
19. Training-Deployment Differences

### E. The People Behind the Data

20. Platform Terms of Service
21. Anonymization and Ability to Delete One’s Information
22. Warnings and Recourse
23. Crowdsourcing

## III. METHOD

Summary: This section discusses the ethical implications of doing AER using a given method. It presents the types of methods and their tradeoffs, as well as considerations of who is left out, spurious correlations, and the role of context. Special attention is paid to green AI and the fine line between emotion management and manipulation.

### F. Why This Method

24. Types of Methods and their Tradeoffs
25. Who is Left Out by this Method
26. Spurious Correlations
27. Context is Everything
28. Individual Emotion Dynamics
29. Historical Behavior is not always indicative of Future Behavior
30. Emotion Management, Manipulation
31. Green AI

## IV. IMPACT AND EVALUATION

Summary: This section discusses ethical considerations associated with the impact of AER systems using both traditional metrics as well as through a number of other criteria beyond metrics. Notably, this latter subsection discusses interpretability, visualizations, building safeguards, and contestability, because even when systems work as designed, there will be some negative consequences. Recognizing and planning for such outcomes is part of responsible development.

### G. Metrics

32. Reliability / Accuracy
33. Demographic Biases
34. Sensitive Applications
35. Testing (on Diverse Datasets, on Diverse Metrics)



#### H. Beyond Metrics

36. Interpretability, Explainability
37. Visualization
38. Safeguards and Guard Rails
39. Harms even when the System Works as Designed
40. Contestability and Recourse
41. Be wary of Ethics Washing

### V. IMPLICATIONS FOR PRIVACY, SOCIAL GROUPS

Summary: This section presents ethical implications of AER for privacy and for social groups. These issues cut across Task Design, Data, Method, and Impact. I discuss both individual and group privacy. The latter becomes especially important in the context of soft-biometrics determined through AER that are not intended to be able to identify individuals, but rather identify groups of people with similar characteristics. I discuss the need for work that does not treat people as a homogeneous group (ignoring sub-group differences) but rather explores disaggregation and intersectionality, while minimizing reification and essentialization of social constructs.

#### I. Implications for Privacy

42. Privacy and Personal Control
43. Group Privacy and Soft Biometrics
44. Mass Surveillance vs. Right to Privacy, Expression, Protest
45. Right Against Self-Incrimination
46. Right to Non-Discrimination

#### J. Implications for Social Groups

47. Disaggregation
48. Intersectionality
49. Reification and Essentialization
50. Attributing People to Social Groups

One can read these various sections in one go, or simply use it as a reference when needed (jumping to sections of interest).

#### 3.4.1 Task Design (Ten Considerations)

##### A. Theoretical Foundations

*Domain naivete is not a virtue.*

Study the theoretical foundations for the task from relevant research fields such as psychology, linguistics, and sociology, to inform the task formulation.

**#1. Emotion Task and Framing:** Carefully consider what emotion task should be the focus of the work (whether conducting human-annotation or building an automatic system). (See §3.2 for a sample of common emotion tasks.) When building an AER system, a clear grasp of the task will help in making appropriate design choices. When choosing which AER system to use, a clear grasp of the emotion task most appropriate for the deployment context will help in choosing the right AER system. It is not uncommon for users of AER to have a particular emotion task in mind and mistakenly assume that an off-the-shelf AER system is designed for that task.

Each of the emotion tasks has associated ethical considerations. For example,

*Is the goal to infer one's true emotions? Is it possible to comprehensively determine one's internal mental state by any AI or human? (Hint: No.) Is it ethical to determine such a private state?*

Realize that it is impossible to capture the full emotional experience of a person (even if one had access to all the electrical signals in the brain). A less ambitious goal is to infer some aspects of one's emotional state.

Here, we see a distinct difference between AER that uses vision and AER that uses language. While there is little credible evidence of the connection between one's facial expressions and one's internal emotional state, there is a substantial amount of work on the idea that language is a window into one's mind (Chomsky 1975; Lakoff 2008; Pinker 2007)—which of course also includes emotions (Bamberg 1997; Wiebe, Wilson, and Cardie 2005; Tausczik and Pennebaker 2010).

That said, there is no evidence that one can determine the full (or even substantial portions) of one's emotional state through their language. (See also considerations #2 *Emotion Model* and #13 *Variability of Expression* ahead on complexity of the emotional experience and variability of expression.) Thus, often it is more appropriate to frame the AER task differently; for example, the objective could be:

- To study how people express emotions: Work that uses speaker-annotated labeled data such as emotion-word hashtags in tweets usually captures how people convey emotions (Mohammad 2012; Purver and Battersby 2012). What people convey may not necessarily indicate what they feel.
- To determine perceived emotion (how others may think one is feeling): Perceived emotions are not necessarily the emotions of the speaker. Emotion annotations by people who have not written the source text usually reveal perceived emotions. (This is most common in NLP data-annotation projects.) Annotation aggregation strategies, such as majority voting, usually only convey emotions perceived by a majority group. Are we missing out on the perceptions of some groups? (More on majority voting in #17 *Label Aggregation*.)
- To determine emotionality of language used in text (regardless of whose emotions, target/stimulus, etc.): This may be appropriate in some restricted-domain scenarios, for example, when one is looking at customer reviews. Here, the context is indicative that the emotionality in the language likely indicates attitude toward the product being reviewed. However, such systems have difficulty when dealing with movie and book reviews because then it has to distinguish between text expressing attitudes toward the book/movie from text describing what happened in the plot (which is likely emotional too).
- To determine trends at aggregate level: Emotionality of language is also useful when tracking broad patterns at an aggregate level—for example, tracking trends of emotionality in tens of thousands of tweets or text in novels over time (e.g., Paul and Dredze 2011; Mohammad 2011; Quercia et al. 2012). The idea is that aggregating information from a large number of instances leads to the determination of meaningful trends in emotionality. (See also discussion in #5 *Aggregate Level vs. Individual Level*.)

In summary, it is important to identify what emotion task is the focus of one's work, use appropriate data, and communicate the nuance of what is being captured to the stakeholders. Not doing so will lead to the misuse and misinterpretation of one's work. Specifically, **AER systems should not claim to determine one's emotional state from**

**their utterance, facial expression, gait, and so forth.** At best, AER systems capture what one is trying to convey or what is perceived by the listener/viewer, and, even there, given the complexity of human expression, they are often inaccurate. A separate question is whether AER systems can determine trends in the emotional state of a person (or a group) over time. Here, inferences are drawn at aggregate level from much larger amounts of data. Studies on public health, such as those listed in §3.3, fall in this category. Here too, it is best to be cautious in making claims about mental state, and to use AER as one source of evidence among many (and involve expertise from public health and psychology).

**#2. Emotion Model and Choice of Emotions:** Work on AER needs to operationalize the aspect of emotion it intends to capture, that is, decide on emotion-related categories or dimensions of interest, decide on how to represent them, and so on. Psychologists and neuroscientists have identified several theories of emotion to inform these decisions:

- The Basic Emotions Theory (BET): Work by Dr. Paul Ekman in 1960s galvanized the idea that some emotions (such as joy, sadness, fear, etc.) are universally expressed through similar facial expressions, and these emotions are more basic than others (Ekman 1992; Ekman and Davidson 1994). This was followed by other proposals of basic emotions by Plutchik, Izard, and others. However, many of the tenets of BET, such as the universality of some emotions and their fixed mapping to facial expressions, stand discredited or are in question (Barrett 2017a; Barrett et al. 2019).
- The Dimensional Theory: Several influential studies have shown that the three most fundamental, largely independent, dimensions of affect and connotative meaning are valence (positiveness–negativeness / pleasure–displeasure), arousal (active–sluggish), and dominance (dominant–submissive / in control–out of control) (Osgood, Suci, and Tannenbaum 1957; Russell 1980; Russell and Mehrabian 1977; Russell 2003). Valence and arousal specifically are commonly studied in a number of psychological and neurocognitive explorations of emotion.
- Cognitive Appraisal Theory: The core idea behind appraisal theory (Scherer 1999; Lazarus 1991) is that emotions arise from a person’s evaluation of a situation or event. (Some varieties of the theory point to a parallel process of reacting to perceptual stimuli as well.) Thus it naturally accounts for variability in emotional reaction to the same event, since different people may appraise the situation differently. Criticisms of appraisal theory center around questions such as: whether emotions can arise without appraisal; whether emotions can arise without physiological arousal; and whether our emotions inform our evaluations.
- The Theory of Constructed Emotions: Dr. Lisa Barrett proposed a new theory on how the human brain constructs emotions from our experiences of the world around us and the signals from our body (Barrett 2017b).

Because ML approaches rely on human-annotated data (which can be hard to obtain in large quantities), AER research has often gravitated to the Basic Emotions Theory, as that work allows one to focus on a small number of emotions. This attraction has been even stronger in the vision AER research because of BET’s suggested mapping between

facial expressions and emotions. However, as noted above, many of the tenets of BET stand debunked.

Consider which formulation of emotions is appropriate for your task/project. For example, one may choose to work with the dimensional model or the model of constructed emotions if the goal is to infer behavioral or health outcome predictions. Despite criticisms of BET, it makes sense for some NLP work to focus on *categorical emotions* such as joy, sadness, guilt, pride, fear, and so forth (including what some refer to as basic emotions) because people often talk about their emotions in terms of these concepts. Most human languages have words for these concepts (even if our individual mental representations for these concepts vary to some extent). However, note that work on categorical emotions by itself is not an endorsement of the BET. Do not refer to some emotions as basic emotions, unless you mean to convey your belief in the BET. Careless endorsement of theories can lead to the perpetuation of ideas that are actively harmful (such as suggesting we can determine internal state from outward appearance—physiognomy).

**#3. Meaning and Extra-Linguistic Information:** The meaning of an utterance is not only a property of language, but it is grounded in human activity, social interactions, beliefs, culture, and other extra-linguistic events, perceptions, and knowledge (Harris 1954; Chomsky 2014; Ervin-Tripp 1973; Bisk et al. 2020; Bender and Koller 2020; Hovy and Yang 2021). Thus one can express the same emotion in different ways in different contexts, different people express the same emotions in different ways, and the same utterances can evoke different emotions in different people. AER systems that do not take extra-linguistic information into consideration will always be limited in their capabilities, and risk being systematically biased, insensitive, and discriminatory. More on this in #13 *Variability of Expression* and #14 *Norms of Emotion Expression*.

**#4. Wellness and Emotion:** The prominent role of one's body in the theory of constructed emotion (Barrett 2017a) nicely accounts for the fact that various physical and mental illnesses (e.g., Parkinson, Alzheimer, cardiovascular disease, depression, anxiety) impact our emotional lives. Existing AER systems are not capable of handling these inter-subject and within-subject variabilities and thus should not be deployed in scenarios where their decisions could negatively impact the lives of people; and, where deployed, their limitations should be clearly communicated.

Emotion recognition is playing a greater role than ever before in understanding how our language reflects our wellness, understanding how certain physical and mental illnesses impact our emotional expression, and understanding how emotional expression can help improve our well-being. For some medical conditions, clinicians can benefit from a detailed history of one's emotional state. However, people are generally not very good at remembering how they had been feeling over the past week, month, and so on. Thus an area of interest is to use AER to help patients track their emotional state. See applications of AER in Public Health in Section 3.3. See also CLPsych Workshop proceedings. Note, however, that these are cases where the technology is working firmly in an assistive role to clinicians and psychologists—providing additional information in situations where human experts make decisions based on a number of other sources of information as well. See Chancellor et al. (2019) for ethical considerations on inferring mental health states from one's utterances.

**#5. Aggregate Level vs. Individual Level:** Emotion detection can be used to make inferences about individuals or groups of people; for example, to assist one in writing, to recommend products or services, or to determine broad trends in attitudes toward

a product, issue, or some other entity. Statistical inferences tend to be more reliable when using large amounts of data and when using more relevant data. Systems that make predictions about individuals often have very little pertinent information about the individual and thus often fall back on data from groups of people. Thus, given the person-to-person variability and within-person variability discussed in the earlier bullets, systems are imbued with errors and biases. Further, these errors are especially detrimental because of the direct and personal nature of such interactions. They may, for example, attribute majority group behavior/preferences to the individual, further marginalizing those that are not in the majority.

*Various ethical concerns, including privacy, manipulation, bias, and free speech, are further exacerbated when systems act on individuals.*

Work on finding trends in large groups of people on the other hand benefits from having a large amount of relevant information to draw on. However, see #43 *Group Privacy* and #47 to #50 *Implications for Social Groups* for relevant concerns.

## **B. Implications of Automation**

*What are the ethical implications of automating the chosen task?*

**#6. Why Automate (Who Benefits and Will this Shift Power):** When we choose to work on a particular AER task, or any AI task for that matter, it is important to ask ourselves why. Often the first set of responses may be straightforward: For example, to automate some process to make people's lives easier, or to provide access to some information that is otherwise hard to obtain, or to answer research questions about how emotions work. However, lately there has been a call to go beyond this initial set of responses and ask more nuanced, difficult, and uncomfortable questions such as:

- Who will benefit from this work and who will not (Trewin et al. 2019)?
- Will this work shift power from those who already have a lot of power to those who have less power (Kalluri 2020)?
- How can we reframe or redesign the task so that it helps those who are most in need (Monteiro 2019)?

Specifically for AER, this will involve considerations such as:

- Are there particular groups of people who will not benefit from this task: For example, people who convey and detect emotions differently than what is common (e.g., people on the autism spectrum), people who use language differently than the people whose data is being used to build the system (e.g., older people or people from a different region)?
- If AER is used in some application, say to determine insurance premiums, then is this further marginalizing those who are already marginalized?
- How can we prevent the use of emotion and stance detection systems for detecting and suppressing dissidents?
- How can AER help those who need the most help?

Various other considerations such as those listed in this sheet can be used to further evaluate the wisdom in investing our labor in a particular task.

**#7. Embracing Neurodiversity:** Much of the ML/NLP emotion work has assumed homogeneity of users and ignored neurodiversity, alexithymia, and autism spectrum. These groups have significant overlap, but are not identical. They are also often characterized as having difficulty in sensing and expressing emotions. Therefore these groups hold particular significance in the development of an inclusive AER system. Existing AER systems implicitly cater to the more populous neurotypical group. At minimum, such AER systems should explicitly acknowledge this limitation. Report disaggregated performance metrics for relevant groups. (See also #47 *Disaggregation*.)

Greater research attention needs to be paid to the neurodiverse group. When doing data annotations, we should try to obtain information on whether participants are neurodiverse or neurotypical (when participants are comfortable sharing that information), and include that information at an aggregate level when we report participant demographics. Work in psychology has used scales such as the Toronto Alexithymia Scale (TAS-20) to determine the difficulty that people might have in identifying and describing emotions (Bagby, Parker, and Taylor 1994).

**#8. Participatory/Emancipatory Design:** Participatory design in research and systems development centers people, especially marginalized and disadvantaged communities, such that they are not mere passive subjects but rather have the agency to shape the design process (Spinuzzi 2005). This has also been referred to as emancipatory research (Humphries, Mertens, and Truman 2020; Noel 2016; Oliver 1997) and is pithily captured by the rallying cry “nothing about us without us.” These calls have developed across many different domains, including research pertaining to disability (Stone and Priestley 1996; Seale et al. 2015), indigenous communities (Hall 2014), autism spectrum (Fletcher-Watson et al. 2019; Bertilsdotter Rosqvist et al. 2019), and neurodiversity (Brosnan et al. 2017; Motti and Evmenova 2020). See Motti and Evmenova (2020) for specific recommendations for conducting studies with neurodiverse participants.

**#9. Applications, Dual Use, Misuse:** AER is a powerful enabling technology that has a number of applications. Thus, like all enabling technologies it can be misused and abused. Examples of inappropriate commercial AER application include:

- Using AER at airports to determine whether an individual is dangerous simply from their facial expressions.
- Detecting stance toward governing authorities to persecute dissidents.
- Using deception detection or lie detection en masse without proper warrants or judicial approval. (Using such technologies even in carefully restricted individual cases is controversial.)
- Increasing someone’s insurance premium because the system has analyzed one’s social media posts to determine (accurately or inaccurately) that they are likely to have a certain mental health condition.
- Advertising that preys on the emotional state of people, for example, user-specific advertising to people when they are emotionally vulnerable.

*Socio-Psychological Applications:* Applications such as inferring patterns in emotions of a speaker to in turn infer other characteristics such as suitability for a job, personality traits, or health conditions are especially fraught with ethical concerns. For example, consider the use of the Myers–Briggs Type Indicator (MBTI) for hiring decisions or

research on detecting personality traits automatically. Notable ethical concerns include the following:

- MBTI is criticized by psychologists, especially for its lack of test–retest reliability (Boyle 1995; Gerras and Wong 2016; Grant 2013). The Big 5 personality traits formalism (Cobb-Clark and Schurer 2012) has greater validity, but even when using Big 5, it is easy to overstate the conclusions.
- Even with accurate personality trait identification, there is little to no evidence that using personality traits for hiring and team-composition decisions is beneficial. The use of such tests have also been criticized on the grounds of discrimination (Snow 2020).

*Health and Well-Being Applications:* AER has considerable potential for improving our health and well-being outcomes. However, the sensitive nature of such applications require substantial efforts to adhere to the best ethical principles. For example, how can harm be mitigated when systems make errors? Should automatic systems be used at all given that sometimes we cannot put a value to the cost of errors? What should be done when the system detects that one is at a high risk of suicide, depression, or some other severe mental health condition? How do we safeguard patient privacy? See the shared task at the 2021 CLPsych Workshop where a secure enclave was used to store the training and test data. See these papers for ethical considerations of AI systems in health care (Yu, Beam, and Kohane 2018; Lysaght et al. 2019; Panesar 2019).

*Applications in Art and Culture:* Lately there has been increasing use of AI in art and culture, especially through curation and recommendation systems. See Born et al. (2021) for a discussion of ethical implications, including whether we are really able to determine what art one would like, long-term impacts of automated curation (on users and artists), and diversity of sources and content.

AI is also used in the analysis and generation of art: for example, for literary analysis and generating poems, paintings, songs, and so forth. Since emotions are a central component of art, much of this work also includes automatic emotion recognition: for example, tracking the emotions of characters in novels, recommending songs for people based on their mood, and generating emotional music. This raises several questions including:

- Is it art if the creation did not involve human input?<sup>1</sup>
- Should AI play a collaborative role with other artists (enhancing their creativity) as opposed to generate pieces on its own?
- How will artists be impacted by AI's role in art?
- Who should get credit for AI art?<sup>2</sup>
- How should we critique AI art?<sup>3</sup>

See further discussion by Hertzmann (2020).

1 <https://www.artbasel.com/news/artificial-intelligence-art-artist-boundary>.

2 <https://www.cnn.com/style/article/ai-art-who-should-get-credit-conversation/index.html>.

3 <https://www.artnews.com/art-in-america/features/creative-ai-art-criticism-1202686003/>.

**#10. Disclosure of Automation:** Disclose to all stakeholders the decisions that are being made (in part or wholly) by automation. Provide mechanisms for the user to understand why relevant predictions were made, and also to contest the decisions. (See also #36 *Interpretability* and #40 *Contestability*.)

Artificial agents that perceive and convey emotions in a human-like manner can give one the impression that they are interacting with a human. Artificial agents should begin their interactions with humans by first disclosing that they are artificial agents (Dickson 2018), even though some studies show certain negative outcomes of such a disclosure (De Cicco, Palumbo et al. 2020; Mozafari, Weiger, and Hammerschmidt 2020).

### 3.4.2 Data (*Thirteen Considerations*)

#### C. Why This Data

*What are the ethical implications of using the chosen data?*

**#11. Types of Data:** Emotion and sentiment researchers have used text data, speech data, data from mobile devices, data from social media, product reviews, suicide notes, essays, novels, movie screenplays, financial documents, and so forth. All of these entail their own ethical considerations in terms of the various points discussed in this article. AER systems use data in various forms, including:

- *Large Language Models:* Language models such as BERT (that capture common patterns in language use) are obtained by training ML models on massive amounts of text found on the Internet. See Bender et al. (2021) for ethical considerations in the use of large language models, including: documentation debt, curation difficulty, incorporation of inappropriate biases, and perpetuation of stereotypes. Note also that using smaller amounts of data raise concerns as well: It may not have enough generalizable information; it may be easier to overfit on; and it may not include diverse perspectives. An important aspect of preparing data (big or small) is deciding how to curate it (e.g., what to discard).
- *Emotion Lexicons:* Emotion Lexicons are lists of words and their associated emotions (determined manually by annotation or automatically from large corpora). Word–emotion association lexicons (such as AFINN [Nielsen 2011], NRC Emotion Lexicon [Mohammad and Turney 2013], and the Valence, Arousal, Dominance Lexicon [Mohammad 2018]) are a popular type of resource used in emotion research, emotion-related data science, and machine learning models for AER. See Mohammad (2020) for biases and ethical considerations in the use of such emotion lexicons. Notable among these considerations is how words in different domains often convey different senses and thus have different emotion associations. Also, word associations capture historic perceptions that change with time and may differ across different groups of people. They are not indicative of inherent immutable emotion labels.
- *Labeled Training and Testing Data:* AER systems often make use of a relatively small number of example instances that are manually labeled (annotated) for emotions. A portion of these is used to train/fine-tune the large language model (training set). The rest is further split for development and testing. I discuss various ethical considerations associated with using emotion-labeled instances below.



**#12. Dimensions of Data:** The data used by AER systems can be examined across various dimensions: size of data; whether it is custom data (carefully produced for the research) or data obtained from an online platform (naturally occurring data); less private/sensitive data or more private/sensitive data; what languages are represented in the data; degree of documentation provided with the data; and so on. All of these have societal implications and the choice of datasets should be appropriate for the context of deployment.

#### **D. Human Variability vs. Machine Normativeness**

*What should we know about emotion data so that we use it appropriately?*

**#13. Variability of Expression and Mental Representation:** Language is highly variable—we can express roughly the same meaning in many different ways.

*Expressions of emotions through language are highly variable: Different people express the same emotion differently; the same text may convey different emotions to different people.*

This is true even for people living in the same area and especially true for people living in different regions, and people with different lived experiences. Some cues of emotion are somewhat more common and somewhat more reliable than others. This is usually the signal that automatic systems attempt to capture. We construct emotions in our brains from the signals we get from the world and the signal we get from our bodies. This mapping of signals to emotions is highly variable, and different people can have different signals associated with different emotions (Barrett 2017b); therefore, different people have different concept–emotion associations. For example, high school, public speaking, and selfies may evoke different emotions in different people. This variability is not to say that there are no commonalities. In fact, speakers of a language share substantial commonalities in their mental representation of concepts (including emotions), which enables them to communicate with each other. However, the variability should also be taken into consideration when building datasets, systems, and choosing where to deploy the systems.

**#14. Norms of Emotion Expression:** As John M. Culkin once said, “We shape our tools and thereafter they shape us.” Whether text, speech, vision, or any other modality, AI systems are often trained on a limited set of emotion expressions and their emotion annotations (emotion labels for the expressions).

*Thus, through their behavior (e.g., by recognizing some forms of emotion expression and not recognizing others), AI systems convey to the user that it is “normal” or appropriate to convey emotions in certain ways; implicitly invalidating other forms of emotion expression.*

Therefore it is important for emotion recognition systems to accurately map a diverse set of emotion instantiations to emotion categories/dimensions. That said, it is also worth noting that the variations in emotion and language expression are so large that systems can likely never attain perfection. The goal is to obtain useful levels of emotion recognition capabilities without having systematic gaps that convey a strong sense of emotion-expression normativeness.

Normative implications of AER are analogous to normative implications of movies (especially animated ones):

- Badly executed characters express emotions in fixed stereotypical ways.

- Good movies explore the diversity, nuance, and subtlety of human emotion expression.
- Influential movies (bad and good) convey to a wide audience around the world how emotions are expressed or what is “normal” in terms of emotion expression. Thus they can either colonize other groups, reducing emotion expression diversity, or they can validate one’s individualism and independence of self-expression.

Because AI systems are influenced by the data they train on, dataset development should:

- Obtain data from a diverse set of sources. Report details of the sources.
- Studies have shown that a small percentage of speakers often produce a large percentage of utterances (see study by Auxier and Anderson [2021] on tweets). Thus, when creating emotion datasets, limit the number of instances included per person. Mohammad and Kiritchenko (2018) kept one tweet for every query term and tweeter combination when studying relationships between affect categories (data also used in a SemEval-2018 Task 1 on emotions). Kiritchenko et al. (2020) kept at most three tweets per tweeter when studying expressions of loneliness.
- Obtain annotations from a diverse set of people. Report aggregate-level demographic information of the annotators.

Variability is common not just for emotions but also for language. People convey meaning in many different ways. Thus, these considerations apply to NLP in general.

**#15. Norms of Attitudes:** Different people and different groups of people might have different attitudes, perceptions, and associations with the same product, issue, person, social groups, and so on. Annotation aggregation, by, say, majority vote, may convey a more homogeneous picture to the ML system. Annotation aggregation may also capture stereotypes and inappropriate associations for already marginalized groups. (For example, majority group A may perceive a minority group B as less competent, or less generous.) Such inappropriate biases are also encoded in large language models. When using language models or emotion datasets, assess the risk of such biases for the particular context and take correcting action as appropriate.

**#16. One “Right” Label or Many Appropriate Labels:** When designing data annotation efforts, consider whether there is a “right” answer and a “wrong” one. Who decides what is correct/appropriate? Are we including the voices of those that are marginalized and already under-represented in the data? When working with emotion and language data, there are usually no “correct” answers, but rather, some answers are more appropriate than others. And there can be multiple appropriate answers.

- If a task has clear correct and wrong answers and knowing the answers requires some training/qualifications, then one can use domain experts to annotate the data. However, as mentioned, emotion annotations largely do not fall in this category.
- If the goal is to determine how people use language, and there can be many appropriate answers, or we want to know how people perceive

words, phrases, and sentences then we might want to use a large number of annotators. This is much more in line with what is appropriate for emotion annotations—people are the best judges of their emotions and of the emotions they perceive from utterances.

Seek appropriate demographic information (respectfully and ethically). Document annotator demographics, annotation instructions, and other relevant details. These are useful in conveying to the reader that there is no one “correct” answer and that the dataset is situated in who annotated the data, the precise annotation instructions, when the data was annotated, and so forth.

**#17. Label Aggregation:** Multiple annotations (by different people) for the same instance are usually aggregated by choosing the majority label. However, majority voting tends to capture majority group attitudes (at the expense of other groups). (See also Aroyo and Welty 2015, Checco et al. 2017, and Klenner et al. 2020.) As a result, sometimes researchers have released not just the aggregated results but also the raw (pre-aggregated data), as well as various versions of aggregated results. Others have argued in favor of not doing majority voting at all and including all annotations as input to ML systems (Basile 2020). However, saying all voices should be included has its own problems: For example, how to address and manage inappropriate/racist/sexist opinions; how to disentangle low-frequency valid opinions from genuine annotation errors and malicious annotations. (See also #15 *Norms of Attitudes* and #47 *Disaggregation*.)

If using majority voting, acknowledge its limitations. Acknowledge that it may be missing some/many voices. Explore statistical approaches to finding multiple appropriate labels, while still discarding noise. Use separate manual checks to determine whether the human annotations also capture inappropriate human biases. Such biases may be useful for some projects (e.g., work studying such biases), but not for others. Warn users of inappropriate biases that may exist in the data; and suggest strategies to deal with them when using the dataset.

**#18. Historical Data (Who is Missing and What are the Biases):** Machine learning methods feed voraciously on data (often historical data). Natural language processing systems often feed on huge amounts of data collected from the Internet. However, the data is not representative of everyone and seeped into this data are our biases. Historical data over-represents people who have had power, who are more well to do, mostly from the west, mostly English-speaking, mostly white, mostly able-bodied, and so on and so forth. So the machines that feed on such data often learn their perspectives at the expense of the views of those already marginalized.

When using any dataset, devote resources to study who is included in the dataset and whose voices are missing. Take corrective action as appropriate. Keep a portion of your funding for work with marginalized communities. Keep a portion of your funding for work on less-researched languages (Ruder 2020).

**#19. Training–Deployment Data Differences:** The accuracy of supervised systems is contingent on the assumption that the data the system is applied to is similar to the data the system was trained on. Deploying an off-the-shelf sentiment analysis system on data in a different domain, from a different time, or a different class distribution than the training data will likely result in poor predictions. Systems that are to be deployed to handle open-domain data should be trained on many diverse datasets and tested on many datasets that are quite different from the training datasets.

## E. The People Behind the Data

*What are the ethical implications on the people who have produced the data?*

When building systems, we make extensive use of (raw and emotion-labeled) data. It can sometimes be easy to forget that behind the data are the people that produced it, and imprinted in it is a plethora of personal information.

**#20. Platform Terms of Service:** Data for ML systems is often scraped from Web sites or extracted from large online platforms (e.g., Twitter, Reddit) using APIs. The terms of service for these platforms often include protections for the users and their data. Ensure that the terms of service of the source platforms are not violated: For example, data scraping is allowed and data redistribution is allowed (in raw form or through ids). Ensure compliance with the robot exclusion protocol.

**#21. Anonymization and Ability to Delete One's information:** Take actions to anonymize data when dealing with private data; for example, scrub identifying information. Some techniques are better at anonymization than others. (See for example, privacy-preserving work on word embeddings and sentiment data by Thaine and Penn [2021].) Provide mechanisms for people to remove their data from the dataset if they choose to.

*Choose to not work with a dataset if adequate safeguards cannot be placed.*

**#22. Warnings and Recourse:** Annotating highly emotional, offensive, or suicidal utterances can adversely impact the well-being of the annotators. Provide appropriate warnings. Minimize amount of data exposure per annotator. Provide options for psychological help as needed.

**#23. Crowdsourcing:** Crowdsourcing (splitting a task into multiple independent units and uploading them on the Internet so that people can solve them online) has grown to be a major source of labeled data in NLP, Computer Vision, and a number of other academic disciplines. Compensation often gets most of the attention when talking about crowdsourcing ethics, but there are several ethical considerations involved with such work such as: worker invisibility, lack of learning trajectory, humans-as-a-service paradigm, worker well-being, and worker rights. See Dolmaza (2011), Fort, Adda, and Cohen (2011), Standing and Standing (2018), Irani and Silberman (2013), and Shmueli et al. (2021). See (public) guidelines by AI2 for its researchers (AI2 2019).

### 3.4.3 Method (Eight Considerations)

## F. Why This Method

*What are the ethical implications of using a given method?*

**#24. Methods and their Trade-offs:** Different methods entail different trade-offs:

- **Less Accurate vs. More Accurate:** This usually gets all the attention; value other dimensions listed below as well. (See also §3.4.4 Impact.)
- **White Box (can understand why system makes a given prediction) vs. Black Box (do not know why it makes a given prediction):** Understanding the reasons behind a prediction helps identify bugs and biases; helps contestability; arguably, is better suited for answering research questions about language use and emotions.

- Less Energy Efficient vs. More Energy Efficient: See discussion further below on Green AI.
- Less Data Hungry vs. More Data Hungry: Data may not always be abundant; needing too much data of a person leads to privacy concerns.
- Less Privacy Preserving vs. More Privacy Preserving: There is greater appreciation lately for the need for privacy-preserving NLP.
- Fewer Inappropriate Biases vs. More Inappropriate Biases: We want our algorithms to not perpetuate/amplify inappropriate human biases.

Consider various dimensions of a method and their importance for the particular system deployment context before deciding on the method. Focusing on fewer dimensions may be okay in a research system, but widely deployed systems often require a good balance across the many dimensions.

**#25. Who is Left Out:** The dominant paradigm in ML and NLP is to use large pre-trained models pre-trained on massive amounts of raw data (unannotated text, pictures, videos, etc.) and then fine-tuned on small amounts of labeled data (e.g., sentences labeled with emotions) to learn how to perform a particular task. As such, these methods tend to work well for people who are well-represented in the data (raw and annotated), but not so well for others. (See also #18 *Historical Data*.)

*Even just documenting who is left out is a valuable contribution.*

Explore alternative methods that are more inclusive, especially for those not usually included by other systems.

**#26. Spurious Correlations:** Machine learning methods have been shown to be susceptible to spurious correlations. For example, Agrawal, Batra, and Parikh (2016) show that when asked what the ground is covered with, visual QA systems tend to always say *snow*, because in the training set, this question was only asked for when the ground was covered with snow. Winkler et al. (2019) and Bissoto, Valle, and Avila (2020) show spurious correlations in melanoma and skin lesion detection systems. Poliak et al. (2018) and Gururangan et al. (2018) show that natural language inference systems can sometimes decide on the prediction just from information in the premise, without regard for the hypothesis (for example, because a premise with negation is often a contradiction in the training set).

Similarly, machine learning systems capture spurious correlations when doing AER. For example, marking some countries and people of some demographics with less charitable and stereotypical sentiments and emotions. This phenomenon is especially marked in abusive language detection work, where it was shown that data collection methods in combination with the ML algorithm result in the system marking any comment with identity terms such as gay, Muslim, and Jew as offensive.

Consider how the data collection and machine learning setups can be addressed to avoid such spurious correlations, especially correlations that perpetuate racism, sexism, and stereotypes. In extreme cases, spurious correlations lead to pseudoscience and physiognomy. For example, there have been a spate of papers attempting to determine criminality, personality, trustworthiness, and emotions just from one's face or outer appearance. Note that sometimes, systematic idiosyncrasies of the data can lead to apparent good results on a held out test set even on such tasks. Thus it is important to consider whether the method and sources of information used are expected to capture

the phenomenon of interest. Is there a risk that the use of this method may perpetuate false beliefs and stereotypes? If yes, take appropriate corrective action.

**#27. Context is Everything:** Considering a greater amount of context is often crucial in correctly determining emotions/sentiment. What was said/written before and after the target utterance? Where was this said? What was the intonation and what was emphasized? Who said this? And so on. More context can be a double-edged sword, though. The more the system wants to know about a person to make better predictions, the more we worry about privacy. Work on determining the right balance between collecting more user information and privacy considerations, as appropriate for the context in which the system is deployed.

**#28. Individual Emotion Dynamics:** A form of contextual information is one's utterance emotion dynamics (Hollenstein 2015). The idea is that different people might have different steady states in terms of where they tend to most commonly be (considering any affect dimension of choice). Some may move out of this steady state often, but some may venture out less often. Some recover quickly from the deviations, and for some it may take a lot of time. Similar emotion dynamics occur in the text that people write or the words they utter—Utterance Emotion Dynamics (Hipson and Mohammad 2021). The degree of correlation between the utterance emotion dynamics and the true emotion dynamics may be correlated, but one can argue that examining utterance emotion dynamics is valuable on its own.

Access to utterance emotion dynamics provides greater context and helps judge the degree of emotionality of new utterances by the person. Systems that make use of such detailed contextual information are more likely to make appropriate predictions for diverse groups of people. However, the degree of personal information they require warrants care, concern, and meaningful consent from the users.

**#29. Historical Behavior is not always Indicative of Future Behavior (for Groups and Individuals):** Systems are often trained on static data from the past. However, perceptions, emotions, and behavior change with time. Thus automatic systems may make inappropriate predictions on current data. (See also #18 *Historical Data*.)

**#30. Emotion Management, Manipulation:** Managing emotions is a central part of any human–computer interaction system (even if this is often not an explicitly stated goal). Just as in human–human interactions, we do not want the systems we build to cause undue stress, pain, or unpleasantness. For example, a chatbot has to be careful to not offend or hurt the feelings of the user with whom it is interacting. For this, it needs to assess the emotions conveyed by the user, in order to then be able to articulate the appropriate information with appropriate affect.

However, this same technology can enable companies and governments to detect one's emotions to manipulate their behavior. For example, it is known that we purchase more products when we are sad. So sensing when you are most susceptible to suggestion to plant ideas of what to buy, who to vote for, or who to dislike, can have dangerous implications. On the other hand, identifying how to cater to individual needs to improve their compliance with public health measures in a worldwide pandemic, or to help people give up on smoking, may be seen in more positive light. As with many things discussed in this article, consider the context to determine what levels of emotional management and meaningful consent are appropriate.

**#31. Green AI:** A direct consequence of using ever-increasing pre-trained models (large number of training examples and hyperparameters) for AI tasks is that these systems

are now drivers of substantial energy consumption. Recent papers show the increasing carbon footprint of AI systems and approaches to address them (Strubell, Ganesh, and McCallum 2020; Schwartz et al. 2020). Thus, there is a growing push to develop AI methods that are not singularly focused on accuracy numbers on test sets, but are also mindful of efficiency and energy consumption (Schwartz et al. 2020). The authors encourage reporting of cost per example, size of training set, number of hyperparameters, and budget-accuracy curves. They also argue for regarding efficiency as a valued scientific contribution.

#### 3.4.4 Impact and Evaluation (Ten Considerations)

##### G. Metrics

*All evaluation metrics are misleading. Some metrics are more useful than others.*

**#32. Reliability/Accuracy:** No emotion recognition method is perfect. However, some approaches are much less accurate than others. Some techniques are so unreliable that they are essentially pseudoscience. For example, trying to predict personality, mood, or emotions through physical appearances has long been criticized (Arcas, Mitchell, and Todorov 2017). The ethics of a number of existing commercial systems that purportedly detect emotions from facial expressions is called into question by Barrett et al. (2019), which shows the low reliability of recognizing emotions from facial expressions.

**#33. Demographic Biases:** Some systems can be unreliable or systematically inaccurate for certain groups of people, races, genders, people with health conditions, people that are on the autism spectrum, people from different countries, etc. Such systematic errors can occur when working on:

- Utterances of a group or faces of a group: For example, low accuracy in recognizing emotions in text produced by African Americans or in recognizing faces of African Americans (Buolamwini and Gebru 2018).
- Utterances mentioning a group: For example, systematically marking texts mentioning African Americans as more angry, or texts mentioning women as more emotional (Kiritchenko and Mohammad 2018).

Determine and present disaggregated accuracies. Take steps to address disparities in performance across groups. (See also #47 *Disaggregation*.)

**#34. Sensitive Applications:** Some applications are considerably more sensitive than others and thus necessitate the use of a much higher quality of emotion recognition systems (if used at all). Automatic systems may sometimes be used in high-stakes applications if their role is to assist human experts. For example, assisting patients and health experts in tracking the patient's emotional state.

**#35. Testing (on Diverse Datasets, on Diverse Metrics):** Results on any test set are contingent on the attributes of that test set and may not be indicative of real-world performance, or implicit biases, or systematic errors of many kinds. Good practice is to test the system on many different datasets that explore various input characteristics. For example, see these evaluations that cater to a diverse set of emotion-related tasks, datasets, linguistic phenomena, and languages: SemEval 2014 Task 9 (Rosenthal et al. 2014), SemEval 2015 Task 10 (Rosenthal et al. 2015), and SemEval 2018 Task 1 (Mohammad et al. 2018). (The last of which also includes an evaluation component for demographic bias in sentiment analysis systems.) See Röttger et al.

(2020) for work on creating separate diagnostic datasets for various types of hate speech. See Google's recommendations on best practices on metrics and testing (<https://ai.google/responsibilities/responsible-ai-practices>).

## H. Beyond Metrics

*Are we even measuring the right things?*

**#36. Interpretability, Explainability:** As ML systems are deployed more widely and impact a greater sphere of our lives, there is a growing understanding that these systems can be flawed to varying degrees. One line of approach in understanding and addressing these flaws is to develop interpretable or explainable models. Interpretability and explainability each have been defined in a few different ways in the literature, but at the heart of the definitions is the idea that we should be able to understand why a system is making a certain prediction: What pieces of evidence are contributing to the decision, and to what degree? That way, humans can better judge how valid a particular prediction is, how accurate the model is for certain kinds of input, and even how accurate the system is in general and over time.

In line with this, AER systems should have components that depict why they are making certain predictions for various inputs. As described in the Luo et al. (2021) survey, such components can be viewed from several perspectives, including:

- Are the explanations meant for the scientist/engineer or a layperson?
- Are the explanations faithful (accurate reflections of system behavior)?
- Are the explanations easily comprehensible?
- To what extent do people trust the explanations?

Responsible research and product development entails actively considering various explainability strategies at the very outset of the project. This includes, where appropriate, specifically choosing an ML model that lends itself to better interpretability, running ablation and disaggregation experiments, running data perturbation and adversarial testing experiments, and so on.

**#37. Visualization:** Visualizations help convey trends in emotions and sentiments, and are common in the emotion analysis of streams of data such as tweet streams, novels, newspaper headlines, and so forth. There are several considerations when developing visualizations that impact the extent to which they are effective and convey key trends, and the extent to which they may be misleading:

- It is almost always important to not only show the broad trends but also to allow the user to drill down to the source data that is driving the trend.
- Summarize the data driving the trend, for example, through treemaps of the most frequent emotion words and phrases in the data.
- Interactive visualizations allow users to explore different trends in the data and even drill down to the source data that is driving the trends.

See work on visualizing emotions and sentiment (Mohammad 2011; Dwibhasi et al. 2015; Kucher, Paradis, and Kerren 2018; Fraser et al. 2019; Gallagher et al. 2021).



**#38. Safeguards and Guard Rails:** Devote time and resources to identify how the system can be misused and how the system may cause harm because of its inherent biases and limitations. Identify steps that can be taken to mitigate these harms.

**#39. Recognize that There will be Harms even when the System Works “Correctly”:** Provide a mechanism for users to report issues. Have resources in place to deal with unanticipated harms. Document societal impacts, including both benefits and harms.

**#40. Contestability and Recourse:** Mulligan, Kluttz, and Kohli (2019) argue that contestability—the mechanisms made available to challenge the predictions of an AI system—are more important and beneficial than transparency/explainability. Not only do they allow people to challenge the decisions made by a system, they also invite participation in the understanding of how machine learning systems work and their limitations. See Google’s What-If Tool as an example of how people are invited to explore ML systems by changing inputs (without needing to do any coding) (Google 2018). AER systems are encouraged to produce similar tools, for example:

- Tools that allow one to see counterfactuals—given a data point, what is the closest other data point for which the system predicts a different label; and tools that allow one to try out various input conditions/features to see which help obtain the desired classification label.
- Tools that allow one to see classification accuracies on different demographics and the impact of different classifier parameters and thresholds on these scores.
- Tools that allow one to see confidence of the classifier for a given prediction and the features that were primarily responsible for the decision.

See Denton et al. (2020) for ideas on participatory dataset creation and management.

**#41. Be Wary of Ethics Washing:** As we push farther into incorporating ethical practices in our projects, we need to be wary of inauthentic and cursory attention to ethics for the sake of appearances. The VentureBeat article (Johnson 2019) presents some nice tips to avoid ethics washing, including: “Welcome ‘constructive dissent’ and uncomfortable conversations,” “Don’t ask for permission to get started,” “Share your shortcomings,” “Be prepared for gray area decision-making,” and “Ethics has few clear metrics.”

### 3.4.5 Implications for Privacy, Social Groups (Nine Considerations)

#### I. Implications for Privacy

*(Cuts across Task Design, Data, Method, and Impact and Evaluation)*

**#42. Privacy and Personal Control:** As noted privacy expert Dr. Ann Cavoukian puts it: Privacy is not about secrecy or hiding information. It is about choice, “You have to be the one to make the decision.” Individuals may not want their emotions to be inferred. Applying emotion detection systems en masse—gathering emotion information continuously, without meaningful consent, is an invasion of privacy, harmful to the individual, and dangerous to society. (See the report created for the members of the European Parliament [Woensel and Nevil 2019]). Follow the seven principles of privacy by design (Schaar 2010): Proactive not Reactive (preventative not remedial), Privacy as the Default, Privacy Embedded into Design, Full Functionality (positive-sum, not zero-sum), End-to-End Security (full lifecycle), Visibility and Transparency, and Respect

for User Privacy (keep it user-centric). See also privacy-preserving work on sentiment by Thaine and Penn (2021).

**#43. Group Privacy and Soft Biometrics:** Floridi (2014) argues that many of our conversations around privacy are far too focused on individual privacy and ignore group privacy—the rights and protections we need as a group.

*There are very few Moby-Dicks. Most of us are sardines. The individual sardine may believe that the encircling net is trying to catch it. It is not. It is trying to catch the whole shoal. It is therefore the shoal that needs to be protected, if the sardine is to be saved.* — Floridi (2014)

The idea of group privacy becomes especially important in the context of soft-biometrics such as traits and preferences determined through AER that are not intended to be able to identify individuals, but rather identify groups of people with similar characteristics. See McStay (2020) for further discussions on the implications of AER on group privacy and how companies are using AER to determine group preferences, even though a large number of people disfavor such profiling.

**#44. Mass Surveillance versus Right to Privacy, Right to Freedom of Expression, and Right to Protest:** Emotion recognition, sentiment analysis, and stance detection can be used for mass surveillance by companies and governments (often without meaningful consent). There is low awareness in people that their information (e.g., what they say or click on an online platform) can be used against their best interest. Often people do not have meaningful choices regarding privacy when they use online platforms. In extreme cases, as in the case of authoritarian governments, this can lead to dramatic curtailing of freedoms of expression and the right to protest (ARTICLE19 2021; Wakefield 2021).

**#45. Right Against Self-Incrimination:** In a number of countries around the world, the accused are given legal rights against self-incrimination. However, automatic methods of emotion, stance, and deception detection can potentially be used to circumvent such protections (see ARTICLE19 [2021] page 37).

**#46. Right to Non-Discrimination:** Automatic methods of emotion, stance, and deception detection can sometimes systematically discriminate based on these protected categories such as race, gender, and religion. Even if ML systems are not fed race or gender information directly, studies have shown that they often pick up on proxy attributes for these categories. Report disaggregated results as appropriate.

## J. Implications for Social Groups

*(Cuts across Task Design, Data, Method, and Impact and Evaluation)*

**#47. Disaggregation:** Society has often viewed different groups differently (because of their race, gender, income, language, etc.), imposing unequal social and power structures (Lindsey 2015). Even when the biases are not conscious, the unique needs of different groups is often overlooked. For example, Perez (2019) discusses, through numerous examples, how there is a considerable lack of disaggregated data for women and how that is directly leading to negative outcomes in all spheres of their lives, including health, income, safety, and the degree to which they succeed in their endeavors. This holds true (perhaps even more) for transgender people. Thus, emotion researchers should consider the value of disaggregation at various levels, including:

- When creating datasets: Obtain annotations from a diverse group of people. Report aggregate-level demographic information. Rather than

only labeling instances with the majority vote, consider the value of providing multiple sets of labels as per each of the relevant and key demographic groups.

- When testing hypotheses or drawing inferences about language use: Consider also testing the hypotheses disaggregated for each of the relevant and key demographic groups.
- When building automatic prediction systems: Report performance disaggregated for each of the relevant and key demographic groups. (See work on model cards [Mitchell et al. 2019]. See how sentiment analysis systems can be systematically biased [Kiritchenko and Mohammad 2018].)

**#48. Intersectional Invisibility in Research:** Intersectionality refers to the complex ways in which different group identities such as race, class, neurodiversity, and gender overlap to amplify discrimination or disadvantage. Purdie-Vaughns and Eibach (2008) argue how people with multiple group identities are often not seen as prototypical members of any of their groups and thus are subject to, what they call, intersectional invisibility—omissions of their experiences in historical narratives and cultural representation, lack of support from advocacy groups, and mismatch with existing anti-discrimination frameworks. Many of the forces that lead to such invisibility (e.g., not being seen as prototypical members of a group) along with other notions common in the quantitative research paradigm (e.g., the predilection to work on neat, non-overlapping, populous categories) lead to intersectional invisibility in research. As ML/NLP researchers, we should be cognizant of such blind spots and work to address these gaps. Further, new ways of doing research that address the unique challenges of doing intersectional research need to be valued and encouraged.

**#49. Reification and Essentialization:** Some demographic variables are essentially, or in big part, social constructs. Thus, work on disaggregation can sometimes reinforce false beliefs that there are innate differences across different groups or that some features are central for one to belong to a social category. Thus it is imperative to contextualize work on disaggregation. For example, by impressing on the reader that even though race is a social construct, the impact of people's perceptions and behavior around race lead to very real-world consequences.

**#50. Attributing People to Social Groups:** In order to be able to obtain disaggregated results, sometimes one needs access to demographic information. This of course leads to considerations such as whether they are providing meaningful consent to the collection of such data and whether the data is being collected in a manner that respects their privacy, their autonomy (e.g., can they choose to delete their information later), and their dignity (e.g., allowing self-descriptions). Challenges persist in terms of how to design effective and inclusive questionnaires (Bauer et al. 2017). Further, even with self-report textboxes that give the respondent the primacy and autonomy to express their race, gender, and so on, downstream research often ignores such data or combines information in ways beyond the control of the respondent (Keyes 2019).

Some work tries to infer aggregate-level group statistics automatically. For example, inferring race or gender, from cues such as the type of language used, historical name-gender associations, and so forth, to do disaggregated analysis. However, such approaches are fraught with ethical concerns such as misgendering, essentialization, and reification. Further, historically, people have been marginalized because of their

social category, and so methods that try to detect these categories raise legitimate and serious concerns of abuse, erasure, and perpetuating stereotypes.

In many cases, it may be more appropriate to perform disaggregated analysis on something other than a social category. For example, when testing face recognition systems, it might be more appropriate to test the system performance on different skin tones (as opposed to race). Similarly, when working on language data, it might be more appropriate to analyze data partitioned by linguistic gender (as opposed to social gender). See Cao and Daumé (2021) for a useful discussion on linguistic vs. social gender and also for a great example to create more inclusive data for research.

#### 4. In Summary

This article aggregates and organizes various ethical considerations relevant to automatic emotion recognition, drawn from the wider AI Ethics and Affective Computing literature. It includes brief sections on the modalities of information, task, and applications of AER to set the context. Then it presents fifty ethical considerations grouped thematically. Notably, the ethics sheet fleshes out assumptions hidden in how AER is commonly framed, and in the choices often made regarding the data, method, and evaluation. Special attention is paid to the implications of AER on privacy and social groups. It discusses how these considerations manifest within AER and outlines best practices for responsible research. A succinct list of key recommendations for responsible AER discussed in the article is provided in the Appendix.

The objective of the sheet is to encourage practitioners to think in more detail and at the very outset: why to automate, how to automate, and how to judge success based on broad societal implications. I hope that it will help engage the various stakeholders of AER with each other; help stakeholders challenge assumptions made by researchers and developers; and help develop appropriate harm mitigation strategies. Additionally, for those who are new to emotion recognition, the ethics sheet acts as a useful introductory document (complementing survey articles).

As an expert on a technology, an often overlooked and undervalued responsibility is to convey its broad societal impacts to those who deploy the technology, those who make policy decisions about the technology, and society at large. I hope that this ethics sheet helps to that end for emotion recognition, and also spurs the wider community to ask and document: *What ethical considerations apply to my task?*

#### APPENDIX: Recommendations for Responsible AER

Below is a list of key recommendations for responsible AER discussed earlier in the context of various ethical considerations. They are compiled here for easy access. Note that adhering to these recommendations does not guarantee “ethicalness”; nor do these recommendations apply to all contexts. They are guidelines meant to help responsible development and use of AER systems. Particular development or deployment contexts entail further considerations and steps to address them.

##### *Task Design*

1. Center the people, especially marginalized and disadvantaged communities, such that they are not mere passive subjects but rather have the agency to shape the design process.

2. Ask who will benefit from this work and who will not. Will this work shift power from those who already have a lot of power to those that have less power? How can the task be designed so that it helps those that are most in need?
3. Ask how the AER design will impact people in the context of neurodiversity, alexithymia, and autism spectrum.
4. Carefully consider what emotion task should be the focus of the work (whether conducting a human-annotation study or building an automatic prediction model). Different emotion tasks entail different ethical considerations. Communicate the nuance of exactly what emotions are being captured to the stakeholders. Not doing so will lead to the misuse and misinterpretation of one's work.
5. AER systems should not claim to determine one's emotional state from their utterance, facial expression, gait, and so forth. At best, AER systems capture what one is trying to convey or what is perceived by the listener/viewer, and even there, given the complexity of human expression, they are often inaccurate.
6. Even when AER systems attempt to determine the emotional state of a person (or a group) over time (drawing inferences at aggregate level from large amounts of data), such as studies on public health listed in §3.3, it is best to be cautious when making claims. Use AER as one source of evidence among many (and involve relevant expertise; e.g., from public health and psychology).
7. Lay out the theoretical foundations for the task from relevant research fields such as psychology, linguistics, and sociology, and relate the opinions of relevant domain experts to the task formulation. Realize that it is impossible to capture the full emotional experience of a person.
8. Do not refer to some emotions as basic emotions, unless you mean to convey your belief in the Basic Emotions Theory. Careless endorsement of theories can lead to the perpetuation of belief in ideas that are actively harmful (such as suggesting we can determine internal state from outward appearance—physiognomy).
9. Realize that various ethical concerns, including privacy, manipulation, bias, and free speech, are further exacerbated when systems act on data pertaining to people. Take steps such as anonymization and realizing information at aggregate levels.
10. Think about how the AER system can be misused, and how that misuse can be minimized.
11. Use AER as one source of information among many.
12. Do not use AER for fully automated decision making. AER may be used to assist humans in making decisions, coming up with ideas, suggesting where to delve deeper, and sparking their imagination. Consider also the risk of the system inappropriately biasing the human decision makers.

13. Disclose to all stakeholders the decisions that are being made (in part or wholly) by automation. Provide mechanisms for the user to understand why relevant predictions were made, and also to contest the decisions.

#### *Data*

14. Examine the choice of data used by AER systems across various dimensions: size of data; whether it is custom data or data obtained from an online platform; less private/sensitive data or more private/sensitive data; what languages are represented; degree of documentation; and so on.
15. Expressions of emotions through language are highly variable: Different people express the same emotion differently; the same text may convey different emotions to different people. This variability should also be taken into consideration when building datasets, systems, and choosing where to deploy the systems.
16. Variability is common not just for emotions but also for natural language. People convey meaning in many different ways. There is usually no one “correct” way of articulating our thoughts.
17. Aim to obtain useful level of emotion recognition capabilities without having systematic gaps that convey a strong sense of emotion-expression normativeness.
18. When using language models or emotion datasets, avoid perpetuating stereotypes of how one group of people perceives another group.
19. Obtain data from a diverse set of sources. Report details of the sources.
20. When creating emotion datasets, limit the number of instances included per person. Mohammad and Kiritchenko (2018) kept one tweet for every query term and tweeter combination when studying relationships between affect categories (data also used in a shared task on emotions). Kiritchenko et al. (2020) kept at most three tweets per tweeter when studying expressions of loneliness.
21. Obtain annotations from a diverse set of people. Report aggregate-level demographic information of the annotators.
22. In emotion and language data, often there are no “correct” answers. Instead, it is a case of some answers being more appropriate than others. And there can be multiple appropriate answers.
23. Part of conveying that there is no one “correct” answer is to convey how the dataset is situated in many parameters, including: who annotated it, the precise annotation instructions, what data was presented to the annotators (and in what form), and when the data was annotated.
24. Release raw data annotations as well as any aggregations of annotations.
25. If using majority voting, acknowledge its limitations.
26. Explore statistical approaches to finding multiple appropriate labels.

27. Use manual and automatic checks to determine whether the human annotations have also captured inappropriate biases. Such biases may be useful for some projects (e.g., work studying such biases), but not for others. Warn users appropriately and deploy measures to mitigate their impact.
28. When using any dataset, devote time and resources to study who is included in the dataset and whose voices are missing. Take corrective action as appropriate.
29. Keep a portion of your funding for work with marginalized communities and for work on less-researched languages.
30. Systems that are to be deployed to handle open-domain data should be trained on many diverse datasets and tested on many datasets that are quite different from the training datasets.
31. Ensure that the terms of service of the source platforms are not violated: For example, data scraping is allowed and data redistribution is allowed (in raw form or through ids). Check the platform terms of service. Ensure compliance with the robot exclusion protocol. Take actions to anonymize data when dealing with sensitive or private data (e.g., scrub identifying information). Choose to not work with a dataset if adequate safeguards cannot be placed.
32. Proposals of data annotation efforts that may impact the well-being of annotators should first be submitted for approval to one's Research Ethics Board (REB) / Institutional Research Board (IRB). The board will evaluate and provide suggestions so that the work complies with the required ethics standards.
33. An excellent jumping off point for further information on ethical conduct of research involving human subjects is The Belmont Report. The guiding principles they proposed are Respect for Persons, Beneficence, and Justice.

### *Method*

34. Examine choice of method across various dimensions such as interpretability, privacy concerns, energy efficiency, data needs, and so on. Focusing on fewer dimensions may be okay in a research system, but widely deployed systems often require a good balance across the many dimensions. AI methods tend to work well for people who are well-represented in the data (raw and annotated), but not so well for others. Documenting who is left out is valuable. Explore alternative methods that are more inclusive. Consider how the data collection and machine learning setups can be addressed to avoid spurious correlations, especially correlations that perpetuate racism, sexism, and stereotypes.
35. Systems are often trained on static data from the past. However, perceptions, emotions, and behavior change with time. Consider how automatic systems may make inappropriate predictions on current data.
36. Consider the system deployment context to determine what levels of emotional management and meaningful consent are appropriate.

37. Consider the carbon footprint of your method and value efficiency as a contribution. Report costs per example, size of training set, number of hyperparameters, and budget-accuracy curves.

#### *Impact and Evaluation*

38. Consider whether the chosen metrics are measuring what matters.
39. Some methods can be unreliable or systematically inaccurate for certain groups of people, races, genders, people with health conditions, people from different countries, and so forth. Determine and present disaggregated accuracies. Test the system on many different datasets that explore various input characteristics.
40. Responsible research and product development entails actively considering various explainability strategies at the very outset of the project. This includes, where appropriate, specifically choosing an ML model that lends itself to better interpretability, running ablation and disaggregation experiments, running data perturbation and adversarial testing experiments, and so on.
41. When visualizing emotions, it is almost always important to not only show the broad trends but also to allow the user to drill down to the source data that is driving the trend. One can also summarize the data driving the trend, for example, through treemaps of the most frequent emotion words.
42. Devote time and resources to identify how the system can be misused and how the system may cause harm because of its inherent biases and limitations. Recognize that there will be harms even when the system works “correctly.” Identify steps that can be taken to mitigate these harms.
43. Provide mechanisms for contestability that not only allow people to challenge the decisions made by a system about them, but also invites participation in the understanding of how machine learning systems work and its limitations.

#### *Implications for Privacy*

44. Privacy is not about secrecy. It is about personal choice. Follow Dr. Cavoukian’s seven principles of privacy by design.
45. Consider that people might not want their emotions to be inferred. Applying emotion detection systems en masse—gathering emotion information continuously, without meaningful consent, is an invasion of privacy, harmful to the individual, and dangerous to society.
46. Soft-biometrics also have privacy concerns. Consider implications of AER on group privacy and that a large number of people disfavor such profiling.
47. Obtain meaningful consent as appropriate for the context. Working with more sensitive and more private data requires a more involved consent process where the user understands the privacy concerns and willingly provides consent. Consider harm mitigation strategies such as



anonymization techniques and differential privacy. Beware that these can vary in effectiveness.

48. Plan for how to keep people’s information secure.
49. Obtain permission for secondary use or if you intend to distribute the data.
50. When working out the privacy–benefit trade-offs, consider who will really benefit from the technology. Especially consider whether those who benefit are people with power or those with less power. Also, as Dr. Cavoukian says, often privacy and benefits can both be had, “it is not a zero-sum game.”
51. Consider implications of AER for mass surveillance and how that undermines right to privacy, right to freedom of expression, right to protest, right against self-incrimination, and right to non-discrimination.

*Implications for Social Groups*

52. When creating datasets, obtain annotations from a diverse group of people. Report aggregate-level demographic information. Rather than only labeling instances with the majority vote, consider the value of providing multiple sets of labels as per each of the relevant and key demographic groups.
53. When testing hypotheses or drawing inferences about language use, consider also testing the hypotheses disaggregated for each of the relevant demographic groups.
54. When building automatic prediction systems, evaluate and report performance disaggregated for each of the relevant demographic groups.
55. Consider and report the implication of the AER system on intersectionality.
56. Contextualize work on disaggregation: for example, by impressing on the reader that even though race is a social construct, the impact of people’s perceptions and behavior around race lead to very real-world consequences.
57. Obtaining demographic information requires careful and thoughtful considerations such as whether people are providing meaningful consent to the collection of such data and whether the data being collected is in a manner that respects their privacy, their autonomy (e.g., can they choose to delete their information later), and dignity (e.g., allowing self-descriptions).

**Acknowledgments**

I am grateful to Annika Schoene, Mallory Feldman, and Tara Small for their belief and encouragement in the early days of this project. Many thanks to Mallory Feldman (Carolina Affective Neuroscience Lab, UNC) for discussions on the psychology and complexity of emotions. Many thanks to Annika Schoene, Mallory Feldman, Roman Klinger, Rada Mihalcea, Peter Turney, Barbara Plank, Malvina Nissim, Viviana

Patti, Maria Liakata, and Emily Mower Provost for discussions about ethical considerations for emotion recognition and thoughtful comments. Many thanks to Tara Small, Emily Bender, Esmā Balkir, Isar Nejadgholi, Patricia Thaine, Brendan O’Connor, Cyril Goutte, Eric Joanis, Joel Martin, Roland Kuhn, and Sowmya Vajjala for thoughtful comments on the blog post on this work.

## References

- Agrawal, Aishwarya, Dhruv Batra, and Devi Parikh. 2016. Analyzing the behavior of visual question answering models. *arXiv preprint arXiv:1606.07356*. <https://doi.org/10.18653/v1/D16-1203>
- AI2. 2019. Crowdsourcing: Pricing ethics and best practices. Medium. <https://medium.com/ai2-blog/crowdsourcing-pricing-ethics-and-best-practices-8487fd5c9872>
- Arcas, Blaise, Margaret Mitchell, and Alexander Todorov. 2017. Physiognomy's new clothes. Medium. <https://medium.com/@blaisea/physiognomys-new-clothes-f2d4b59fdd6a>
- Aroyo, Lora and Chris Welty. 2015. Truth is a lie: Crowd truth and the seven myths of human annotation. *AI Magazine*, 36(1):15–24. <https://doi.org/10.1609/aimag.v36i1.2564>
- ARTICLE19. 2021. Emotional entanglement: China's emotion recognition market and its implications for human rights. <https://www.article19.org/wp-content/uploads/2021/01/ER-Tech-China-Report.pdf>
- Auxier, Brooke and Monica Anderson. 2021. Social media use in 2021. *Pew Research Center*.
- Bagby, R. Michael, James D. A. Parker, and Graeme J. Taylor. 1994. The twenty-item Toronto Alexithymia scale: I. Item selection and cross-validation of the factor structure. *Journal of Psychosomatic Research*, 38(1):23–32. [https://doi.org/10.1016/0022-3999\(94\)90005-1](https://doi.org/10.1016/0022-3999(94)90005-1)
- Bamberg, Michael. 1997. Language, concepts and emotions: The role of language in the construction of emotions. *Language Sciences*, 19(4):309–340. [https://doi.org/10.1016/S0388-0001\(97\)00004-1](https://doi.org/10.1016/S0388-0001(97)00004-1)
- Barrett, Lisa Feldman. 2017a. *How Emotions are Made: The Secret Life of the Brain*. Houghton Mifflin Harcourt.
- Barrett, Lisa Feldman. 2017b. The theory of constructed emotion: An active inference account of interoception and categorization. *Social Cognitive and Affective Neuroscience*, 12(1):1–23. <https://doi.org/10.1093/scan/nsx060>
- Barrett, Lisa Feldman, Ralph Adolphs, Stacy Marsella, Aleix M. Martinez, and Seth D. Pollak. 2019. Emotional expressions reconsidered: Challenges to inferring emotion from human facial movements. *Psychological Science in the Public Interest*, 20(1):1–68. <https://doi.org/10.1177/1529100619832930>
- Basile, Valerio. 2020. It's the end of the gold standard as we know it. On the impact of pre-aggregation on the evaluation of highly subjective tasks. *2020 AIXIA Discussion Papers Workshop, AIXIA 2020 DP*, volume 2776, pages 31–40.
- Bauer, Greta R., Jessica Braimoh, Ayden I. Scheim, and Christoffer Dharma. 2017. Transgender-inclusive measures of sex/gender for population surveys: Mixed-methods evaluation and recommendations. *PLOS ONE*, 12(5):e0178043. <https://doi.org/10.1371/journal.pone.0178043>
- Bender, Emily M., Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623. <https://doi.org/10.1145/3442188.3445922>
- Bender, Emily M. and Alexander Koller. 2020. Climbing towards NLU: On meaning, form, and understanding in the age of data. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5185–5198. <https://doi.org/10.18653/v1/2020.acl-main.463>
- Bertilsdotter Rosqvist, Hanna, Marianthi Kourti, David Jackson-Perry, Charlotte Brownlow, Kirsty Fletcher, Daniel Bendelman, and Lindsay O'Dell. 2019. Doing it differently: Emancipatory autism studies within a neurodiverse academic space. *Disability & Society*, 34(7–8):1082–1101. <https://doi.org/10.1080/09687599.2019.1603102>
- Bisk, Yonatan, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai, Mirella Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, Nicolas Pinto, and Joseph Turian. 2020. Experience grounds language. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8718–8735. <https://doi.org/10.18653/v1/2020.emnlp-main.703>
- Bissoto, Alceu, Eduardo Valle, and Sandra Avila. 2020. Debiasing skin lesion datasets and models? Not so fast. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 740–741. <https://doi.org/10.1109/CVPRW50498.2020.00378>
- Born, Georgina, Jeremy Morris, Fernando Diaz, and Ashton Anderson. 2021.

- Artificial Intelligence, music recommendation, and the curation of culture. Technical report. Schwartz Reisman Institute.
- Boyle, Gregory J. 1995. Myers–Briggs type indicator (MBTI): Some psychometric limitations. *Humanities & Social Sciences Papers*, 30(1):71–74. <https://doi.org/10.1111/j.1742-9544.1995.tb01750.x>
- Brosnan, Mark, Samantha Holt, Nicola Yuill, Judith Good, and Sarah Parsons. 2017. Beyond autism and technology: Lessons from neurodiverse populations. *Journal of Enabling Technologies*, 11(2):43–48. <https://doi.org/10.1108/JET-02-2017-0007>
- Buolamwini, Joy and Timnit Gebru. 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on Fairness, Accountability and Transparency*, pages 77–91.
- Cao, Yang Trista and Hal Daumé. 2021. Toward gender-inclusive coreference resolution: An analysis of gender and bias throughout the machine learning lifecycle. *Computational Linguistics*, 47(3):615–661. [https://doi.org/10.1162/colia\\_00413](https://doi.org/10.1162/colia_00413)
- Chancellor, Stevie, Michael L. Birnbaum, Eric D. Caine, Vincent M. B. Silenzio, and Munmun De Choudhury. 2019. A taxonomy of ethical tensions in inferring mental health states from social media. In *Proceedings of the Conference on Fairness, Accountability, and Transparency, FAT\* '19*, pages 79–88. <https://doi.org/10.1145/3287560.3287587>
- Checchio, Alessandro, Kevin Roitero, Eddy Maddalena, Stefano Mizzaro, and Gianluca Demartini. 2017. Let's agree to disagree: Fixing agreement measures for crowdsourcing. In *Proceedings of the Fifth AAAI Conference on Human Computation and Crowdsourcing*, pages 11–20.
- Chomsky, Noam. 1975. *Reflections on Language*. Pantheon.
- Chomsky, Noam. 2014. *Aspects of the Theory of Syntax*, volume 11. MIT Press.
- Cobb-Clark, Deborah A. and Stefanie Schurer. 2012. The stability of big-five personality traits. *Economics Letters*, 115(1):11–15. <https://doi.org/10.1016/j.econlet.2011.11.015>
- De Choudhury, Munmun, Michael Gamon, Scott Counts, and Eric Horvitz. 2013. Predicting depression via social media. In *Seventh International AAAI Conference on Weblogs and Social Media*, pages 128–137.
- De Cicco, Roberta, Riccardo Palumbo, et al. 2020. Should a chatbot disclose itself? Implications for an online conversational retailer. In *International Workshop on Chatbot Research and Design*, pages 3–15. [https://doi.org/10.1007/978-3-030-68288-0\\_1](https://doi.org/10.1007/978-3-030-68288-0_1)
- Denton, Emily, Alex Hanna, Razvan Amironesei, Andrew Smart, Hilary Nicole, and Morgan Klaus Scheuerman. 2020. Bringing the people back in: Contesting benchmark machine learning datasets. *arXiv preprint arXiv:2007.07399*.
- Dickson, Ben. 2018. Why AI must disclose that it's AI. *PC Magazine*. <https://www.pcmag.com/opinions/why-ai-must-disclose-that-its-ai>
- Dolmaza, Julie McDonough. 2011. The ethics of crowdsourcing. *Linguistica Antverpiensia, New Series—Themes in Translation Studies*, 10:97–110.
- Dwibhasi, Sharat, Dheeraj Jami, Shivkanth Lanka, and Goutam Chakraborty. 2015. Analyzing and visualizing the sentiments of Ebola outbreak via tweets. In *Proceedings of the SAS Global Forum*, pages 26–29.
- Eichstaedt, Johannes C., Hansen Andrew Schwartz, Margaret L. Kern, Gregory Park, Darwin R. Labarthe, Raina M. Merchant, Sneha Jha, Megha Agrawal, Lukasz A. Dziurzynski, Maarten Sap, Christopher Weeg, Emily E. Larson, Lyle H. Ungar, and Martin E. P. Seligman. 2015. Psychological language on Twitter predicts county-level heart disease mortality. *Psychological Science*, 26(2):159–169. <https://doi.org/10.1177/0956797614557867>
- Ekman, Paul. 1992. Are there basic emotions? *Psychological Review*, 99(3):550–553. <https://doi.org/10.1037/0033-295X.99.3.550>
- Ekman, Paul Ed and Richard J. Davidson. 1994. *The Nature of Emotion: Fundamental Questions*. Oxford University Press.
- Ervin-Tripp, Susan. 1973. Some strategies for the first two years. In *Cognitive Development and Acquisition of Language*, pages 261–286. <https://doi.org/10.1016/B978-0-12-505850-6.50018-9>
- Fletcher-Watson, Sue, Jon Adams, Kabie Brook, Tony Charman, Laura Crane, James Cusack, Susan Leekam, Damian Milton, Jeremy R. Parr, and Elizabeth Pellicano. 2019. Making the future together: Shaping autism research through meaningful participation. *Autism*, 23(4):943–953. <https://doi.org/10.1177/1362361318786721>

- Floridi, Luciano. 2014. Open data, data protection, and group privacy. *Philosophy & Technology*, 27(1):1–3. <https://doi.org/10.1007/s13347-014-0157-8>
- Fort, Karën, Gilles Adda, and K. Bretonnel Cohen. 2011. Last words: Amazon Mechanical Turk: Gold mine or coal mine? *Computational Linguistics*, 37(2):413–420. [https://doi.org/10.1162/COLI\\_a\\_00057](https://doi.org/10.1162/COLI_a_00057)
- Fraser, Kathleen C., Frauke Zeller, David Harris Smith, Saif M. Mohammad, and Frank Rudzicz. 2019. How do we feel when a robot dies? Emotions expressed on Twitter before and after hitchBOT's destruction. In *Proceedings of the Tenth Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 62–71. <https://doi.org/10.18653/v1/W19-1308>
- Gallagher, Ryan J., Morgan R. Frank, Lewis Mitchell, Aaron J. Schwartz, Andrew J. Reagan, Christopher M. Danforth, and Peter Sheridan Dodds. 2021. Generalized word shift graphs: A method for visualizing and explaining pairwise comparisons between texts. *EPJ Data Science*, 10(1):4. <https://doi.org/10.1140/epjds/s13688-021-00260-3>
- Gerras, Stephen J. and Leonard Wong. 2016. Moving beyond the MBTI. *Military Review*, 96:54–57.
- Google. 2018. Google AI Blog. <https://ai.googleblog.com/2018/09/the-what-if-tool-code-free-probing-of.html>
- Grant, Adam. 2013. Say goodbye to MBTI, the fad that won't die. LinkedIn. <https://www.linkedin.com/pulse/20130917155206-69244073-say-goodbye-to-mbti-the-fad-that-won-t-die>
- Guntuku, Sharath Chandra, Rachele Schneider, Arthur Pelullo, Jami Young, Vivien Wong, Lyle Ungar, Daniel Polsky, Kevin G. Volpp, and Raina Merchant. 2019. Studying expressions of loneliness in individuals using Twitter: An observational study. *BMJ Open*, 9(11):e030355. <https://doi.org/10.1136/bmjopen-2019-030355>
- Gururangan, Suchin, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R. Bowman, and Noah A. Smith. 2018. Annotation artifacts in natural language inference data. *arXiv preprint arXiv:1803.02324*. <https://doi.org/10.18653/v1/N18-2017>
- Hall, Lisa. 2014. 'With' not 'about': Emerging paradigms for research in a cross-cultural space. *International Journal of Research & Method in Education*, 37(4):376–389. <https://doi.org/10.1080/1743727X.2014.909401>
- Harris, Zellig S. 1954. Distributional structure. *Word*, 10(2–3):146–162. <https://doi.org/10.1080/00437956.1954.11659520>
- Hertzmann, Aaron. 2020. Computers do not make art, people do. *Communications of the ACM*, 63(5):45–48. <https://doi.org/10.1145/3347092>
- Hipson, Will E. and Saif M. Mohammad. 2021. Emotion dynamics in movie dialogues. *PLOS ONE*, 16:1–19. <https://doi.org/10.1371/journal.pone.0256153>
- Hollenstein, Tom. 2015. This time, it's real: Affective flexibility, time scales, feedback loops, and the regulation of emotion. *Emotion Review*, 7(4):308–315. <https://doi.org/10.1177/1754073915590621>
- Hovy, Dirk and Diyi Yang. 2021. The importance of modeling social factors of language: Theory and practice. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 588–602. <https://doi.org/10.18653/v1/2021.naacl-main.49>
- Humphries, Beth, Donna M. Mertens, and Carole Truman. 2020. Arguments for an 'emancipatory' research paradigm, *Research and Inequality*. Routledge, pages 3–23. <https://doi.org/10.1201/9781003071679-2>
- Irani, Lilly C. and M. Six Silberman. 2013. Turkopticon: Interrupting worker invisibility in Amazon Mechanical Turk. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 611–620. <https://doi.org/10.1145/2470654.2470742>
- Johnson, Khari. 2019. How AI companies can avoid ethics washing. VentureBeat. <https://venturebeat.com/2019/07/17/how-ai-companies-can-avoid-ethics-washing/>
- Kalluri, Pratyusha. 2020. Don't ask if Artificial Intelligence is good or fair, ask how it shifts power. *Nature*, 583(7815):169–169. <https://doi.org/10.1038/d41586-020-02003-2>
- Karam, Zahi N., Emily Mower Provost, Satinder Singh, Jennifer Montgomery, Christopher Archer, Gloria Harrington, and Melvin G. Mcinnis. 2014. Ecologically valid long-term mood monitoring of individuals with bipolar disorder using speech. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4858–4862.

- <https://doi.org/10.1109/ICASSP.2014.6854525>
- Keyes, Os. 2019. Counting the countless. REAL LIFE. <https://reallifemag.com/counting-the-countless/>
- Kiritchenko, Svetlana, Will Hipson, Robert Coplan, and Saif M. Mohammad. 2020. SOLO: A corpus of tweets for examining the state of being alone. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 1567–1577.
- Kiritchenko, Svetlana and Saif Mohammad. 2018. Examining gender and race bias in two hundred sentiment analysis systems. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 43–53. <https://doi.org/10.18653/v1/S18-2005>
- Klenner, Manfred, Anne Göhring, Michael Amsler, Sarah Ebling, Don Tuggener, Manuela Hürlimann, and Martin Volk. 2020. Harmonization sometimes harms. In *Proceedings of the 5th Swiss Text Analytics Conference (SwissText) & 16th Conference on Natural Language Processing (KONVENS)*.
- Kucher, Kostiantyn, Carita Paradis, and Andreas Kerren. 2018. Visual analysis of sentiment and stance in social media texts. In *EuroVis (Posters)*, pages 49–51.
- Lakoff, George. 2008. *Women, Fire, and Dangerous Things: What Categories Reveal About the Mind*. University of Chicago Press.
- Lazarus, Richard S. 1991. Progress on a cognitive-motivational-relational theory of emotion. *American Psychologist*, 46(8):819. <https://doi.org/10.1037/0003-066X.46.8.819>
- Lindsey, Linda L. 2015. The sociology of gender theoretical perspectives and feminist frameworks. In *Gender Roles*. Routledge, pages 23–48. <https://doi.org/10.4324/9781315664095-6>
- Luo, Siwen, Hamish Ivison, Caren Han, and Josiah Poon. 2021. Local interpretations for explainable natural language processing: A survey. *arXiv preprint arXiv:2103.11072*.
- Lysaght, Tamra, Hannah Yeefen Lim, Vicki Xafis, and Kee Yuan Ngiam. 2019. AI-assisted decision-making in healthcare. *Asian Bioethics Review*, 11(3):299–314. <https://doi.org/10.1007/s41649-019-00096-0>
- MacAvaney, Sean, Anjali Mittu, Glen Coppersmith, Jeff Leintz, and Philip Resnik. 2021. Community-level research on suicidality prediction in a secure environment: Overview of the CLPsych 2021 shared task. In *Proceedings of the Seventh Workshop on Computational Linguistics and Clinical Psychology: Improving Access*, pages 70–80. <https://doi.org/10.18653/v1/2021.clpsych-1.7>
- McStay, Andrew. 2020. Emotional AI, soft biometrics and the surveillance of emotional life: An unusual consensus on privacy. *Big Data & Society*, 7(1):2053951720904386. <https://doi.org/10.1177/2053951720904386>
- Mitchell, Margaret, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model cards for model reporting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 220–229. <https://doi.org/10.1145/3287560.3287596>
- Mohammad, Saif. 2011. From once upon a time to happily ever after: Tracking emotions in novels and fairy tales. In *Proceedings of the 5th ACL-HLT Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities*, pages 105–114. Portland.
- Mohammad, Saif. 2012. #Emotional tweets. In *\*SEM 2012: The First Joint Conference on Lexical and Computational Semantics—Volume 1: Proceedings of the Main Conference and the Shared Task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 246–255.
- Mohammad, Saif, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. SemEval-2018 task 1: Affect in tweets. In *Proceedings of the 12th International Workshop on Semantic Evaluation*, pages 1–17. <https://doi.org/10.18653/v1/S18-1001>
- Mohammad, Saif and Svetlana Kiritchenko. 2018. Understanding emotions: A dataset of tweets to study interactions between affect categories. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.
- Mohammad, Saif M. 2018. Obtaining reliable human ratings of valence, arousal, and dominance for 20,000 English words. In *Proceedings of the Annual Conference of the Association for Computational Linguistics (ACL)*, Melbourne, Australia. <https://doi.org/10.18653/v1/P18-1017>
- Mohammad, Saif M. 2020. Practical and ethical considerations in the effective use of emotion and sentiment lexicons. *arXiv:2011.03492*.

- Mohammad, Saif M. 2021a. Ethics sheets for AI tasks. In *Proceedings of the 60th Annual Meeting of the Association of Computational Linguistics (ACL-2022)*.
- Mohammad, Saif M. 2021b. Sentiment analysis: Automatically detecting valence, emotions, and other affectual states from text. In Herbert L. Meiselman, editor, *Emotion Measurement (Second Edition)*, Woodhead Publishing, pages 323–379. <https://doi.org/10.1016/B978-0-12-821124-3.00011-9>
- Mohammad, Saif M., Svetlana Kiritchenko, Parinaz Sobhani, Xiaodan Zhu, and Colin Cherry. 2016. Semeval-2016 task 6: Detecting stance in tweets. In *Proceedings of the International Workshop on Semantic Evaluation, SemEval '16*, pages 31–41. <https://doi.org/10.18653/v1/S16-1003>
- Mohammad, Saif M., Parinaz Sobhani, and Svetlana Kiritchenko. 2017. Stance and sentiment in tweets. *Special Section of the ACM Transactions on Internet Technology on Argumentation in Social Media*, 17(3):1–23. <https://doi.org/10.1145/3003433>
- Mohammad, Saif M. and Peter D. Turney. 2013. Crowdsourcing a word-emotion association lexicon. *Computational Intelligence*, 29(3):436–465. <https://doi.org/10.1111/j.1467-8640.2012.00460.x>
- Monteiro, Mike. 2019. *Ruined by Design: How Designers Destroyed the World, and What We Can Do to Fix It*. Mule Design.
- Motti, Vivian Genaro and Anna Evmenova. 2020. Designing technologies for neurodiverse users: Considerations from research practice. In *Human Interaction and Emerging Technologies*, pages 268–274. [https://doi.org/10.1007/978-3-030-25629-6\\_42](https://doi.org/10.1007/978-3-030-25629-6_42)
- Mozafari, Nika, Welf H. Weiger, and Maik Hammerschmidt. 2020. The chatbot disclosure dilemma: Desirable and undesirable effects of disclosing the non-human identity of chatbots. In *ICIS*, pages 1–18.
- Mulligan, Deirdre K., Daniel Kluttz, and Nitin Kohli. 2019. Shaping our tools: Contestability as a means to promote responsible algorithmic decision making in the professions. *Available at SSRN 3311894*.
- Nielsen, Finn Årup. 2011. A new ANEW: Evaluation of a word list for sentiment analysis in microblogs. In *Proceedings of the ESWC Workshop on 'Making Sense of Microposts': Big Things Come in Small Packages*, pages 93–98.
- Noel, Lesley-Ann. 2016. Promoting an emancipatory research paradigm in design education and practice. In *Proceedings of DRS2016 International Conference, Vol. 6: Future-Focused Thinking*, pages 27–30. <https://doi.org/10.21606/drs.2016.355>
- Oliver, Michael. 1997. Emancipatory research: Realistic goal or impossible dream. *Doing Disability Research*, 2:15–31.
- Osgood, Charles Egerton, George J. Suci, and Percy H. Tannenbaum. 1957. *The Measurement of Meaning*. 47. University of Illinois Press.
- Panesar, Arjun. 2019. *Machine Learning and AI for Healthcare*. Springer. <https://doi.org/10.1007/978-1-4842-3799-1>
- Pang, Bo, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? Sentiment classification using machine learning techniques. In *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing (EMNLP 2002)*, pages 79–86. <https://doi.org/10.3115/1118693.1118704>
- Paul, Michael J. and Mark Dredze. 2011. You are what you tweet: Analyzing Twitter for public health. *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media*, 5(1):265–272.
- Perez, Caroline Criado. 2019. *Invisible Women: Exposing Data Bias in a World Designed for Men*. Random House.
- Picard, Rosalind W. 2000. *Affective Computing*. MIT Press. <https://doi.org/10.7551/mitpress/1140.001.0001>
- Pinker, Steven. 2007. *The Stuff of Thought: Language as a Window Into Human Nature*. Penguin.
- Poliak, Adam, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. 2018. Hypothesis only baselines in natural language inference. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 180–191. <https://doi.org/10.18653/v1/S18-2023>
- Purdie-Vaughns, Valerie and Richard P. Eibach. 2008. Intersectional invisibility: The distinctive advantages and disadvantages of multiple subordinate-group identities. *Sex Roles*, 59(5):377–391. <https://doi.org/10.1007/s11199-008-9424-4>
- Purver, Matthew and Stuart Battersby. 2012. Experimenting with distant supervision for emotion classification. In *Proceedings of the 13th Conference of the European Chapter of*

- the Association for Computational Linguistics*, pages 482–491.
- Quercia, Daniele, Jonathan Ellis, Licia Capra, and Jon Crowcroft. 2012. Tracking “Gross community happiness” from tweets. In *Proceedings of the ACM 2012 Conference on Computer Supported Cooperative Work, CSCW '12*, pages 965–968. <https://doi.org/10.1145/2145204.2145347>
- Resnik, Philip, William Armstrong, Leonardo Claudino, Thang Nguyen, Viet-An Nguyen, and Jordan Boyd-Graber. 2015. Beyond LDA: Exploring supervised topic modeling for depression-related language in Twitter. In *Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, pages 99–107. <https://doi.org/10.3115/v1/W15-1212>
- Rosenthal, Sara, Preslav Nakov, Svetlana Kiritchenko, Saif Mohammad, Alan Ritter, and Veselin Stoyanov. 2015. SemEval-2015 task 10: Sentiment analysis in Twitter. In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 451–463. <https://doi.org/10.18653/v1/S15-2078>
- Rosenthal, Sara, Alan Ritter, Preslav Nakov, and Veselin Stoyanov. 2014. SemEval-2014 task 9: Sentiment analysis in Twitter. In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 73–80. <https://doi.org/10.3115/v1/S14-2009>
- Röttger, Paul, Bertram Vidgen, Dong Nguyen, Zeerak Waseem, Helen Margetts, and Janet Pierrehumbert. 2020. HateCheck: Functional tests for hate speech detection models. *arXiv preprint arXiv:2012.15606*. <https://doi.org/10.18653/v1/2021.ac1-long.4>
- Ruder, Sebastian. 2020. Why you should do NLP beyond English. <https://ruder.io/nlp-beyond-english/index.html>.
- Russell, James A. 1980. A circumplex model of affect. *Journal of Personality and Social Psychology*, 39(6):1161. <https://doi.org/10.1037/h0077714>
- Russell, James A. 2003. Core affect and the psychological construction of emotion. *Psychological Review*, 110(1):145. <https://doi.org/10.1037/0033-295X.110.1.145>
- Russell, James A. and Albert Mehrabian. 1977. Evidence for a three-factor theory of emotions. *Journal of Research in Personality*, 11(3):273–294. [https://doi.org/10.1016/0092-6566\(77\)90037-X](https://doi.org/10.1016/0092-6566(77)90037-X)
- Schaar, Peter. 2010. Privacy by design. *Identity in the Information Society*, 3(2):267–274. <https://doi.org/10.1007/s12394-010-0055-x>
- Scherer, Klaus R. 1999. *Appraisal Theory*. John Wiley & Sons Ltd.
- Schwartz, Hansen Andrew, Johannes C. Eichstaedt, Margaret L. Kern, Lukasz Dziurzynski, Richard E. Lucas, Megha Agrawal, Gregory J. Park, Shrinidhi K. Lakshminanth, Sneha Jha, Martin EP Seligman, Lyle H. Ungar, and Richard E. Lucas. 2013. Characterizing geographic variation in well-being using tweets. In *Seventh International AAAI Conference on Weblogs and Social Media*, pages 583–591.
- Schwartz, Roy, Jesse Dodge, Noah A. Smith, and Oren Etzioni. 2020. Green AI. *Communications of the ACM*, 63(12):54–63. <https://doi.org/10.1145/3381831>
- Seale, Jane, Melanie Nind, Liz Tilley, and Rohhss Chapman. 2015. Negotiating a third space for participatory research with people with learning disabilities: An examination of boundaries and spatial practices. *Innovation: The European Journal of Social Science Research*, 28(4):483–497. <https://doi.org/10.1080/13511610.2015.1081558>
- Shmueli, Boaz, Jan Fell, Soumya Ray, and Lun-Wei Ku. 2021. Beyond fair pay: Ethical implications of NLP crowdsourcing. *arXiv preprint arXiv:2104.10097*. <https://doi.org/10.18653/v1/2021.naacl-main.295>
- Snow, Shane. 2020. That personality test may be discriminating people... and making your company dumber. LinkedIn. <https://www.linkedin.com/pulse/personality-test-may-discriminating-people-making-your-shane-snow>
- Soleymani, Mohammad, David Garcia, Brendan Jou, Björn Schuller, Shih-Fu Chang, and Maja Pantic. 2017. A survey of multimodal sentiment analysis. *Image and Vision Computing*, 65:3–14. <https://doi.org/10.1016/j.imavis.2017.08.003>
- Spinuzzi, Clay. 2005. The methodology of participatory design. *Technical Communication*, 52(2):163–174.
- Standing, Susan and Craig Standing. 2018. The ethical use of crowdsourcing. *Business Ethics: A European Review*, 27(1):72–80. <https://doi.org/10.1111/beer.12173>
- Stone, Emma and Mark Priestley. 1996. Parasites, pawns and partners: Disability research and the role of non-disabled researchers. *British Journal of Sociology*,

- pages 699–716. <https://doi.org/10.2307/591081>
- Strubell, Emma, Ananya Ganesh, and Andrew McCallum. 2020. Energy and policy considerations for modern deep learning research. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(09):13693–13696. <https://doi.org/10.1609/aaai.v34i09.7123>
- Tausczik, Yla R. and James W. Pennebaker. 2010. The psychological meaning of words: LIWC and computerized text analysis methods. *Journal of Language and Social Psychology*, 29(1):24–54. <https://doi.org/10.1177/0261927X09351676>
- Thaine, Patricia and Gerald Penn. 2021. The Chinese remainder theorem for compact, task-precise, efficient and secure word embeddings. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3512–3521. <https://doi.org/10.18653/v1/2021.eacl-main.306>
- Trewin, Shari, Sara Basson, Michael Muller, Stacy Branham, Jutta Treviranus, Daniel Gruen, Daniel Hebert, Natalia Lyckowski, and Erich Manser. 2019. Considerations for AI fairness for people with disabilities. *AI Matters*, 5(3):40–63. <https://doi.org/10.1145/3362077.3362086>
- Turney, Peter. 2002. Thumbs up or thumbs down? Semantic orientation applied to unsupervised classification of reviews. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 417–424. <https://doi.org/10.3115/1073083.1073153>
- Wakefield, Jane. 2021. AI emotion-detection software tested on Uyghurs. BBC. <https://www.bbc.com/news/technology-57101248>
- Wiebe, Janyce, Theresa Wilson, and Claire Cardie. 2005. Annotating expressions of opinions and emotions in language. *Language Resources and Evaluation*, 39(2):165–210. <https://doi.org/10.1007/s10579-005-7880-9>
- Winkler, Julia K., Christine Fink, Ferdinand Toberer, Alexander Enk, Teresa Deinlein, Rainer Hofmann-Wellenhof, Luc Thomas, Aimilios Lallas, Andreas Blum, Wilhelm Stolz, and Holger Haenssle. 2019. Association between surgical skin markings in dermoscopic images and diagnostic performance of a deep learning convolutional neural network for melanoma recognition. *JAMA Dermatology*, 155(10):1135–1141. <https://doi.org/10.1001/jamadermatol.2019.1735>
- Woensel, Lieve Van and Nissy Nevil. 2019. What if your emotions were tracked to spy on you? European Parliamentary Research Service, PE 634.415. [https://www.europarl.europa.eu/RegData/etudes/ATAG/2019/634415/EPRS\\_ATA\(2019\)634415\\_EN.pdf](https://www.europarl.europa.eu/RegData/etudes/ATAG/2019/634415/EPRS_ATA(2019)634415_EN.pdf)
- Yu, Kun Hsing, Andrew L. Beam, and Isaac S. Kohane. 2018. Artificial Intelligence in healthcare. *Nature Biomedical Engineering*, 2(10):719–731. <https://doi.org/10.1038/s41551-018-0305-z>
- Zhang, Lei, Shuai Wang, and Bing Liu. 2018. Deep learning for sentiment analysis: A survey. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 8(4):e1253. <https://doi.org/10.1002/widm.1253>