

# The Social Impact of Natural Language Processing

**Dirk Hovy**

Center for Language Technology  
University of Copenhagen  
Copenhagen, Denmark  
dirk.hovy@hum.ku.dk

**Shannon L. Spruit**

Ethics & Philosophy of Technology  
Delft University of Technology  
Delft, The Netherlands  
s.l.spruit@tudelft.nl

## Abstract

Medical sciences have long since established an ethics code for experiments, to minimize the risk of harm to subjects. Natural language processing (NLP) used to involve mostly anonymous corpora, with the goal of enriching linguistic analysis, and was therefore unlikely to raise ethical concerns. As NLP becomes increasingly wide-spread and uses more data from social media, however, the situation has changed: the outcome of NLP experiments and applications *can* now have a direct effect on individual users' lives. Until now, the discourse on this topic in the field has not followed the technological development, while public discourse was often focused on exaggerated dangers. This position paper tries to take back the initiative and start a discussion. We identify a number of social implications of NLP and discuss their ethical significance, as well as ways to address them.

## 1 Introduction

After the Nuremberg trials revealed the atrocities conducted in medical research by the Nazis, medical sciences established a set of rules to determine whether an experiment is ethical. This involved incorporating the principles of biomedical ethics as a *lingua franca* of medical ethics (Beauchamp and Childress, 2001).

These guidelines were designed to balance the potential value of conducting an experiment while preventing the exploitation of human subjects. Today, any responsible research institution uses these—or comparable—criteria to approve or reject experiments before any research can be conducted. The administrative body governing these decisions is the Institutional Review Board (IRB).

IRBs mostly pertain to experiments that directly involve human subjects, though, and so NLP and other data sciences have not employed such guidelines. Work on existing corpora is unlikely to raise any flags that would require an IRB approval.<sup>1</sup>

Data sciences have therefore traditionally been less engaged in ethical debates of their subject, even though this seems to be shifting, see for instance Wallach (2014), Galaz et al. (2015), or O'Neil (2016). The public outcry over the “emotional contagion” experiment on Facebook (Kramer et al., 2014) further suggests that data sciences now affect human subjects in real time, and that we might have to reconsider the application of ethical considerations to our research (Puschmann and Bozdag, 2014). NLP research not only involves similar data sets, but also works with their content, so it is time to start a discussion of the ethical issues specific to our field.

Much of the ethical discussion in data sciences to date, however, has centered around privacy concerns (Tse et al., 2015). We do not deny the reality and importance of those concerns, but they involve aspects of digital rights management/access control, policy making, and security, which are not specific to NLP, but need to be addressed in the data sciences community as a whole. Steps towards this have been taken by Russell et al. (2015).

Instead, we want to move beyond privacy in our ethical analysis and look at the wider **social impact** NLP may have. In particular, we want to explore the impact of NLP on social justice, i.e., equal opportunities for individuals and groups (such as minorities) within society to access resources, get their voice heard, and be represented in society.

---

<sup>1</sup>With few exceptions, such as dialogue research (Joel Tetreault, pers. comm.)

**Our contributions** We believe ethical discussions are more constructive if led by practitioners, since the public discussion of ethical aspects of IT and data sciences is often loaded with fear of the unknown and unrealistic expectations. For example, in the public discourse about AI (Hsu, 2012; Eadicicco, 2015; Khatchadourian, 2015), people either dismiss the entire approach, or exaggerate the potential dangers (see Etzioni (2014) for a practitioner’s view point). This paper is an attempt to take back the initiative for NLP.

At the same time, we believe that the field of ethics can contribute a more general framework, and so this paper is an interdisciplinary collaboration between NLP and ethics researchers.

To facilitate the discussion, we also provide some of the relevant **terminology** from the literature on ethics of technology, namely the concepts of *exclusion*, *overgeneralization*, *bias confirmation*, *topic under- and overexposure*, and *dual use*.

## 2 Does NLP need an ethics discussion?

As discussed above, the makeup of most NLP experiments so far has not obviated a need for ethical considerations, and so, while we are aware of individual discussions (Strube, 2015), there is little discourse in the community yet. A search for “*ethic\**” in the ACL anthology only yields three results. One of the papers (McEnery, 2002) turns out to be a panel discussion, another is a book review, leaving only Couillault et al. (2014), who devote most of the discussion to legal and quality issues of data sets. We know social implications have been addressed in some NLP curricula,<sup>2</sup> but until now, no discipline-wide discussion seems to take place.

The most likely reason is that NLP research has not directly involved human subjects.<sup>3</sup> Historically, most NLP applications focused on further enriching existing text which was not strongly linked to any particular author (newswire), was usually published publicly, and often with some temporal distance (novels). All these factors created a distance between text and author, which prevented the research from directly affecting the authors’ situation.

<sup>2</sup>Héctor Martínez Alonso, personal communication

<sup>3</sup>Except for annotation: there are a number of papers on the status of crowdsource workers (Fort et al., 2011; Pavlick et al., 2014). Couillault et al. (2014) also briefly discuss annotators, but mainly in the context of quality control.

This situation has changed lately due to the increased use of social media data, where authors are current individuals, who can be directly affected by the results of NLP applications. Couillault et al. (2014) touch upon these issues under “traceability” (i.e., whether individuals can be identified): this is undesirable for experimental subjects, but might be useful in the case of annotators.

Most importantly, though: the subject of NLP—language—is a *proxy for human behavior*, and a strong signal of individual characteristics. People use this signal consciously, to portray themselves in a certain way, but can also be identified as members of specific groups by their use of subconscious traits (Silverstein, 2003; Agha, 2005; Johannsen et al., 2015; Hovy and Johannsen, 2016).

Language is always *situated* (Bamman et al., 2014), i.e., it is uttered in a specific situation at a particular place and time, and by an individual speaker with all the characteristics outlined above. All of these factors can therefore leave an imprint on the utterance, i.e., the texts we use in NLP carry *latent* information about the author and situation, albeit to varying degrees.

This information can be used to predict author characteristics from text (Rosenthal and McKeown, 2011; Nguyen et al., 2011; Alowibdi et al., 2013; Ciot et al., 2013; Liu and Ruths, 2013; Volkova et al., 2014; Volkova et al., 2015; Plank and Hovy, 2015; Preotiuc-Pietro et al., 2015a; Preotiuc-Pietro et al., 2015b), and the characteristics in turn can be detected by and influence the performance of our models (Mandel et al., 2012; Volkova et al., 2013; Hovy, 2015).

As more and more language-based technologies are becoming available, the ethical implications of NLP research become more important. What research is carried out, and its quality, directly affect the functionality and impact of those technologies.

The following is meant to start a discussion addressing ethical issues that can emerge in (and from) NLP research.

## 3 The social impact of NLP research

We have outlined the relation between language and individual traits above. Language is also a political instrument, though, and an instrument of power. This influence stretches into politics and everyday competition, for example for turn-taking (Laskowski, 2010; Bracewell and Tomlinson, 2012; Prabhakaran and Rambow, 2013; Prab-

hakaran et al., 2014; Tsur et al., 2015; Khouzami et al., 2015, inter alia), .

The mutual relationships between language, society, and the individual are also the source for the societal impact factors of NLP: failing to recognize group membership (Section 3.1), implying the wrong group membership (see Section 3.2), and overexposure (Section 3.3). In the following, we discuss sources of these problems in the data, modeling, and research design, and suggest possible solutions to address them.

### 3.1 Exclusion

As a result of the situatedness of language, any data set carries a **demographic bias**, i.e., latent information about the demographics in it. Overfitting to these factors can have severe effects on the applicability of findings. In psychology, where most studies are based on western, educated, industrialized, rich, and democratic research participants (so-called WEIRD, Henrich et al. (2010)), the tacit assumption that human nature is so universal that findings on this group would translate to other demographics has led to a heavily biased corpus of psychological data. In NLP, overfitting to the demographic bias in the training data is due to the *i.i.d.* assumption. I.e., models implicitly assume all language to be identical to the training sample. They therefore perform worse or even fail on data from other demographics.

Potential consequences are **exclusion** or demographic misrepresentation. This in itself already represents an ethical problem for research purposes, threatening the universality and objectivity of scientific knowledge (Merton, 1973). These problems exacerbate, though, once they are applied to products. For instance, standard language technology may be easier to use for white males from California (as these are taken into account while developing it) rather than women or citizens of Latino or Arabic descent. This will reinforce already existing demographic differences, and makes technology less user friendly for such groups, cf. authors like Bourdieu and Passeron (1990) have shown how restricted language, like class specific language or scientific jargon, can hinder the expression of outsiders' voices from certain practices. A lack of awareness or decreased attention for demographic differences in research stages can therefore lead to issues of exclusion of people along the way.

Concretely, the consequences of exclusion for NLP research have recently been pointed out by Hovy and Søgaard (2015) and Jørgensen et al. (2015): current state-of-the-art NLP models score a significantly lower accuracy for young people and ethnic minorities vis-à-vis the modeled demographics.

Better awareness of these mechanism in NLP research and development can help prevent problems further on. Potential counter-measures to demographic bias can be as simple as downsampling the over-represented group in the training data to even out the distribution. The work by Mohamady and Culotta (2014) shows another approach, by using existing demographic statistics as supervision. In general, measures to address overfitting or imbalanced data can be used to correct for demographic bias in data.

### 3.2 Overgeneralization

Exclusion is a side-effect of the data. **Overgeneralization** is a modeling side-effect.

As an example, we consider automatic inference of user attributes, a common and interesting NLP task, whose solution also holds promise for many useful applications, such as recommendation engines and fraud or deception detection (Badaskar et al., 2008; Fornaciari and Poesio, 2014; Ott et al., 2011; Banerjee et al., 2014).

The cost of false positives seems low: we might be puzzled or amused when receiving an email addressing us with the wrong gender, or congratulating us to our retirement on our 30th birthday.

In practice, though, relying on models that produce false positives may lead to bias confirmation and overgeneralization. Would we accept the same error rates if the system was used to predict sexual orientation or religious views, rather than age or gender? Given the right training data, this is just a matter of changing the target variable.

To address overgeneralization, the guiding question should be “would a false answer be worse than no answer?” We can use dummy variables, rather than take a *tertium non datur* approach to classification, and employ measures such as error weighting and model regularization, as well as confidence thresholds.

### 3.3 The problem of exposure

**Topic overexposure creates biases** Both exclusion and overgeneralization can be addressed algo-

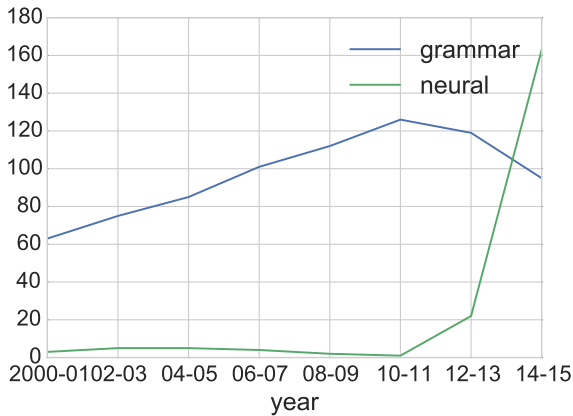


Figure 1: ACL title keywords over time

rhythmically, while **topic overexposure** originates from research design.

In research, we can observe this effect in waves of research topics that receive increased mainstream attention, often to fall out of fashion or become more specialized, cf. ACL papers with “grammars” vs. “neural” in the title (Figure 1).

Such topic overexposure may lead to a psychological effect called **availability heuristic** (Tversky and Kahneman, 1973): if people can recall a certain event, or have knowledge about specific things, they infer it must be more important. For instance, people estimate the size of cities they recognize to be larger than that of unknown cities (Goldstein and Gigerenzer, 2002).

However, the same holds for individuals/groups/characteristics we research. The heuristics become ethically charged when characteristics such as violence or negative emotions are more strongly associated with certain groups or ethnicities (Slovic et al., 2007). If research repeatedly found that the language of a certain demographic group was harder to process, it could create a situation where this group was perceived to be difficult, or abnormal, especially in the presence of existing biases. The confirmation of biases through the gendered use of language, for example, has also been at the core of second and third wave feminism (Mills, 2012).

Overexposure thus creates biases which can lead to discrimination. To some extent, the frantic public discussion on the dangers of AI can be seen as a result of overexposure (Sunstein, 2004).

There are no easy solutions to this problem, which might only become apparent in hindsight. It can help to assess whether the research direction

of a project feeds into existing biases, or whether it overexposes certain groups.

**Underexposure can negatively impact evaluation.** Similar to the WEIRD-situation in psychology, NLP tends to focus on Indo-European data/text sources, rather than small languages from other language groups, for example in Asia or Africa. This focus creates an imbalance in the available amounts of labeled data. Most of the existing labeled data covers only a small set of languages. When analyzing a random sample of Twitter data from 2013, we found that there were no treebanks for 11 of the 31 most frequent languages, and even fewer semantically annotated resources (the ACE corpus covers only English, Arabic, Chinese, and Spanish).<sup>4</sup>

Even if there is a potential wealth of data available from other languages, most NLP tools are geared towards English (Schnoebelen, 2013; Munro, 2013). The prevalence of resources for English has created an **underexposure** to typological variety: both morphology and syntax of English are global outliers. Would we have focused on  $n$ -gram models to the same extent if English was as morphologically complex as, say, Finnish?

While there are many approaches to develop multi-lingual and cross-lingual NLP tools for linguistic outliers (Yarowsky and Ngai, 2001; Das and Petrov, 2011; Søgaard, 2011; Søgaard et al., 2015; Agić et al., 2015), there simply are more commercial incentives to overexpose English, rather than other languages. Even if other languages are equally (or more) interesting from a linguistic and cultural point of view, English is one of the most widely spoken language and therefore opens up the biggest market for NLP tools. This focus on English may be self-reinforcing: the existence of off-the-shelf tools for English makes it easy to try new ideas, while to start exploring other languages requires a higher startup cost in terms of basic models, so researchers are less likely to work on them.

#### 4 Dual-use problems

Even if we address all of the above concerns and do not intend any harm in our experiments, they can still have unintended consequences that negatively affect people’s lives (Jonas, 1984).

Advanced analysis techniques can vastly improve search and educational applications

<sup>4</sup>Thanks to Barbara Plank for the analysis!

(Tetreault et al., 2015), but can re-enforce prescriptive linguistic norms when degrading on non-standard language. Stylometric analysis can shed light on the provenance of historic texts (Mosteller and Wallace, 1963), but also endanger the anonymity of political dissenters. Text classification approaches help decode slang and hidden messages (Huang et al., 2013), but have the potential to be used for censorship. At the same time, NLP can also help uncovering such restrictions (Bamman et al., 2012). As recently shown by Hovy (2016), NLP techniques can be used to detect fake reviews, but also to generate them in the first place.

All these examples indicate that we should become more aware of the way other people appropriate NLP technology for their own purposes. The unprecedented scale and availability can make the consequences of NLP technologies hard to gauge.

The unintended consequences of research are also linked to the incentives associated with funding sources. The topic of government and military involvement in the field deserves special attention in this respect. On the one hand, Anderson et al. (2012) show how a series of DARPA-funded workshops have been formative for ACL as a field in the 1990s. On the other hand, there are scholars who refuse military-related funding for moral reasons.<sup>5</sup>

While this decision is up to the individual researcher, the examples show that moral considerations go beyond the immediate research projects. We may not directly be held responsible for the unintended consequences of our research, but we can acknowledge the ways in which NLP can enable morally questionable/sensitive practices, raise awareness, and lead the discourse on it in an informed manner. The role of the researcher in such ethical discussions has recently been pointed out by Rogaway (2015).

## 5 Conclusion

In this position paper, we outlined the potential social impact of NLP, and discussed ways for the practitioner to address this. We also introduced *exclusion*, *overgeneralization*, *bias confirmation*, *topic overexposure*, and *dual use*. Countermeasures for exclusion include bias control techniques

<sup>5</sup>For a perspective in a related field see <https://web.eecs.umich.edu/~kuipers/opinions/no-military-funding.html>

like downsampling or priors; for overgeneralization: dummy labels, error weighting, or confidence thresholds. Exposure problems can only be addressed by careful research design, and dual-use problems seem hardly addressable on the level of the individual researcher, but require the concerted effort of our community.

We hope this paper can point out ethical considerations for collecting our data, designing the experimental setup, and assessing the potential application of our systems, and help start an open discussion in the field about the limitations and problems of our methodology.

## Acknowledgements

The authors would like to thank Joel Tetreault, Rachel Tatman, Joel C. Wallenberg, the members of the COASTAL group, and the anonymous reviewers for their detailed and invaluable feedback. The first author was funded under the ERC Starting Grant LOWLANDS No. 313695. The second author was funded by the Netherlands Organization for Scientific Research under grant number 016.114.625.

## References

- Asif Agha. 2005. Voice, footing, enregisterment. *Journal of linguistic anthropology*, pages 38–59.
- Željko Agić, Dirk Hovy, and Anders Søgaard. 2015. If all you have is a bit of the Bible: Learning POS taggers for truly low-resource languages. In *Proceedings of the 53rd annual meeting of the ACL*.
- Jalal S Alowibdi, Ugo A Buy, and Philip Yu. 2013. Empirical evaluation of profile characteristics for gender classification on twitter. In *Machine Learning and Applications (ICMLA), 2013 12th International Conference on*, volume 1, pages 365–369. IEEE.
- Ashton Anderson, Dan McFarland, and Dan Jurafsky. 2012. Towards a computational history of the ACL: 1980-2008. In *Proceedings of the ACL-2012 Special Workshop on Rediscovering 50 Years of Discoveries*, pages 13–21. Association for Computational Linguistics.
- Sameer Badaskar, Sachin Agarwal, and Shilpa Arora. 2008. Identifying real or fake articles: Towards better language modeling. In *Proceedings of the Third International Joint Conference on Natural Language Processing: Volume-II*.
- David Bamman, Brendan O’Connor, and Noah Smith. 2012. Censorship and deletion practices in Chinese social media. *First Monday*, 17(3).

- David Bamman, Chris Dyer, and Noah A. Smith. 2014. Distributed representations of geographically situated language. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, pages 828–834. Proceedings of ACL.
- Ritwik Banerjee, Song Feng, Seok Jun Kang, and Yejin Choi. 2014. Keystroke patterns as prosody in digital writings: A case study with deceptive reviews and essays. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1469–1473. Association for Computational Linguistics.
- Tom L. Beauchamp and James F. Childress. 2001. *Principles of biomedical ethics*. Oxford University Press, USA.
- Pierre Bourdieu and Jean-Claude Passeron. 1990. *Reproduction in education, society and culture*, volume 4. Sage.
- David Bracewell and Marc Tomlinson. 2012. The language of power and its cultural influence. In *Proceedings of COLING 2012: Posters*, pages 155–164. The COLING 2012 Organizing Committee.
- Morgane Ciot, Morgan Sonderegger, and Derek Ruths. 2013. Gender Inference of Twitter Users in Non-English Contexts. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, Seattle, Wash.*, pages 18–21.
- Alain Couillault, Karën Fort, Gilles Adda, and Hugues Mazancourt (de). 2014. Evaluating corpora documentation with regards to the Ethics and Big Data Charter. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC-2014)*. European Language Resources Association (ELRA).
- Dipanjan Das and Slav Petrov. 2011. Unsupervised part-of-speech tagging with bilingual graph-based projections. In *Proceedings of the 49th annual meeting of the ACL*.
- Lisa Eadicicco. 2015. Bill Gates: Elon Musk Is Right, We Should All Be Scared Of Artificial Intelligence Wiping Out Humanity. Business Insider, January 28. <http://www.businessinsider.com/bill-gates-artificial-intelligence-2015-1> Retrieved Feb 24, 2016.
- Oren Etzioni. 2014. Its Time to Intelligently Discuss Artificial Intelligence. Backchannel, December 9 <https://backchannel.com/ai-wont-exterminate-us-it-will-empower-us-5b7224735bf3#.eia6vtimy> Retrieved Feb 24, 2016.
- Tommaso Fornaciari and Massimo Poesio. 2014. Identifying fake amazon reviews as learning from crowds. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 279–287. Association for Computational Linguistics.
- Karën Fort, Gilles Adda, and K. Bretonnel Cohen. 2011. Amazon mechanical turk: Gold mine or coal mine? *Computational Linguistics*, 37(2):413–420.
- Victor Galaz, Fredrik Moberg, and Fernanda Torre. 2015. The Biosphere Code Manifesto. <http://thebiospherecode.com/index.php/manifesto> Retrieved Feb 24, 2016.
- Daniel G. Goldstein and Gerd Gigerenzer. 2002. Models of ecological rationality: the recognition heuristic. *Psychological review*, 109(1):75.
- Joseph Henrich, Steven J. Heine, and Ara Norenzayan. 2010. The weirdest people in the world? *Behavioral and brain sciences*, 33(2-3):61–83.
- Dirk Hovy and Anders Johannsen. 2016. Exploring Language Variation Across Europe - A Web-based Tool for Computational Sociolinguistics. In *Proceedings of LREC*.
- Dirk Hovy and Anders Sjøgaard. 2015. Tagging performance correlates with author age. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics*.
- Dirk Hovy. 2015. Demographic factors improve classification performance. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics*.
- Dirk Hovy. 2016. The Enemy in Your Own Camp: How Well Can We Detect Statistically-Generated Fake Reviews—An Adversarial Study. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics.
- Jeremy Hsu. 2012. Control dangerous AI before it controls us, one expert says. NBC News, March 1. [http://www.nbcnews.com/id/46590591/ns/technology\\_and\\_science-innovation](http://www.nbcnews.com/id/46590591/ns/technology_and_science-innovation) Retrieved Feb 24, 2016.
- Hongzhao Huang, Zhen Wen, Dian Yu, Heng Ji, Yizhou Sun, Jiawei Han, and He Li. 2013. Resolving entity morphs in censored data. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1083–1093. Association for Computational Linguistics.
- Anders Johannsen, Dirk Hovy, and Anders Sjøgaard. 2015. Cross-lingual syntactic variation over age and gender. In *Proceedings of CoNLL*.
- Hans Jonas. 1984. *The Imperative of Responsibility: Foundations of an Ethics for the Technological Age*. Original in German: Prinzip Verantwortung.) Chicago: University of Chicago Press.
- Anna Jørgensen, Dirk Hovy, and Anders Sjøgaard. 2015. Challenges of studying and processing dialects in social media. In *Workshop on Noisy User-generated Text (W-NUT)*.

- Raffi Khatchadourian. 2015. The Doomsday Invention: Will artificial intelligence bring us utopia or destruction? *The New Yorker* (magazine), November 23. <http://www.newyorker.com/magazine/2015/11/23/doomsday-invention-artificial-intelligence-nick-bostrom> Retrieved Feb 24, 2016.
- Hatim Khouzami, Romain Laroche, and Fabrice Lefevre, 2015. *Proceedings of the 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, chapter Optimising Turn-Taking Strategies With Reinforcement Learning, pages 315–324. Association for Computational Linguistics.
- Adam DI Kramer, Jamie E Guillory, and Jeffrey T Hancock. 2014. Experimental evidence of massive-scale emotional contagion through social networks. *Proceedings of the National Academy of Sciences*, 111(24):8788–8790.
- Kornel Laskowski. 2010. Modeling norms of turn-taking in multi-party conversation. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 999–1008. Association for Computational Linguistics.
- Wendy Liu and Derek Ruths. 2013. What’s in a name? using first names as features for gender inference in twitter. In *Analyzing Microtext: 2013 AAAI Spring Symposium*.
- Benjamin Mandel, Aron Culotta, John Boulahanis, Danielle Stark, Bonnie Lewis, and Jeremy Rodrigue, 2012. *Proceedings of the Second Workshop on Language in Social Media*, chapter A Demographic Analysis of Online Sentiment during Hurricane Irene, pages 27–36. Association for Computational Linguistics.
- Tony McEnery. 2002. Ethical and legal issues in corpus construction. In *Proceedings of the Third International Conference on Language Resources and Evaluation (LREC’02)*. European Language Resources Association (ELRA).
- Robert K Merton. 1973. The normative structure of science. *The sociology of science: Theoretical and empirical investigations*, 267.
- Sara Mills. 2012. *Gender matters: Feminist linguistic analysis*. Equinox Pub.
- Ehsan Mohammady and Aron Culotta, 2014. *Proceedings of the Joint Workshop on Social Dynamics and Personal Attributes in Social Media*, chapter Using County Demographics to Infer Attributes of Twitter Users, pages 7–16. Association for Computational Linguistics.
- Frederick Mosteller and David L Wallace. 1963. Inference in an authorship problem: A comparative study of discrimination methods applied to the authorship of the disputed federalist papers. *Journal of the American Statistical Association*, 58(302):275–309.
- Robert Munro. 2013. NLP for all languages. *Idibon Blog*, May 22 <http://idibon.com/nlp-for-all> Retrieved May 17, 2016.
- Dong Nguyen, Noah A Smith, and Carolyn P Rosé. 2011. Author age prediction from text using linear regression. In *Proceedings of the 5th ACL-HLT Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities*, pages 115–123. Association for Computational Linguistics.
- Cathy O’Neil. 2016. The Ethical Data Scientist. *Slate*, February 4 [http://www.slate.com/articles/technology/future\\_tense/2016/02/how\\_to\\_bring\\_better\\_ethics\\_to\\_data\\_science.html](http://www.slate.com/articles/technology/future_tense/2016/02/how_to_bring_better_ethics_to_data_science.html) Retrieved Feb 24, 2016.
- Myle Ott, Yejin Choi, Claire Cardie, and T. Jeffrey Hancock. 2011. Finding deceptive opinion spam by any stretch of the imagination. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 309–319. Association for Computational Linguistics.
- Ellie Pavlick, Matt Post, Ann Irvine, Dmitry Kachaev, and Chris Callison-Burch. 2014. The language demographics of amazon mechanical turk. *Transactions of the Association of Computational Linguistics – Volume 2, Issue 1*, pages 79–92.
- Barbara Plank and Dirk Hovy. 2015. Personality traits on twitter or how to get 1,500 personality tests in a week. In *Proceedings of the 6th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 92–98.
- Vinodkumar Prabhakaran and Owen Rambow. 2013. Written dialog and social power: Manifestations of different types of power in dialog behavior. In *Proceedings of the Sixth International Joint Conference on Natural Language Processing*, pages 216–224. Asian Federation of Natural Language Processing.
- Vinodkumar Prabhakaran, Ashima Arora, and Owen Rambow. 2014. Staying on topic: An indicator of power in political debates. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1481–1486. Association for Computational Linguistics.
- Daniel Preotiuc-Pietro, Vasileios Lampos, and Nikolaos Aletras. 2015a. An analysis of the user occupational class through twitter content. In *ACL*.
- Daniel Preotiuc-Pietro, Svitlana Volkova, Vasileios Lampos, Yoram Bachrach, and Nikolaos Aletras. 2015b. Studying user income through language, behaviour and affect in social media. *PloS one*, 10(9):e0138717.

- Cornelius Puschmann and Engin Bozdog. 2014. Staking out the unclear ethical terrain of online social experiments. *Internet Policy Review*, 3(4).
- Phillip Rogaway. 2015. The moral character of cryptographic work. Technical report, IACR-Cryptology ePrint Archive.
- Sara Rosenthal and Kathleen McKeown. 2011. Age prediction in blogs: A study of style, content, and online behavior in pre-and post-social media generations. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pages 763–772. Association for Computational Linguistics.
- Stuart Russell, Daniel Dewey, Max Tegmark, Janos Kramar, and Richard Mallah. 2015. Research priorities for robust and beneficial artificial intelligence. Technical report, Future of Life Institute.
- Tyler Schnoebelen. 2013. The weirdest languages. Idibon Blog, June 21 <http://idibon.com/the-weirdest-languages> Retrieved May 17, 2016.
- Michael Silverstein. 2003. Indexical order and the dialectics of sociolinguistic life. *Language & Communication*, 23(3):193–229.
- Paul Slovic, Melissa L. Finucane, Ellen Peters, and Donald G. MacGregor. 2007. The affect heuristic. *European Journal of Operational Research*, 177(3):1333 – 1352.
- Anders Søgaard, Željko Agić, Héctor Martínez Alonso, Barbara Plank, Bernd Bohnet, and Anders Johannsen. 2015. Inverted indexing for cross-lingual nlp. In *Proceedings of the 53rd annual meeting of the ACL*.
- Anders Søgaard. 2011. Data point selection for cross-language adaptation of dependency parsers. In *Proceedings of ACL*.
- Michael Strube. 2015. It is never as simple as it seems: The wide-ranging impacts of ethics violations. *Ethical Challenges in the Behavioral and Brain Sciences*, page 126.
- Cass R Sunstein. 2004. Precautions against what? the availability heuristic and cross-cultural risk perceptions. *U Chicago Law & Economics, Olin Working Paper*, (220):04–22.
- Joel Tetreault, Jill Burstein, and Claudia Leacock, 2015. *Proceedings of the Tenth Workshop on Innovative Use of NLP for Building Educational Applications*, chapter Proceedings of the Tenth Workshop on Innovative Use of NLP for Building Educational Applications. Association for Computational Linguistics.
- Jonathan Tse, Dawn E Schrader, Dipayan Ghosh, Tony Liao, and David Lundie. 2015. A bibliometric analysis of privacy and ethics in ieee security and privacy. *Ethics and Information Technology*, 17(2):153–163.
- Oren Tsur, Dan Calacci, and David Lazer. 2015. A frame of mind: Using statistical models for detection of framing and agenda setting campaigns. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1629–1638. Association for Computational Linguistics.
- Amos Tversky and Daniel Kahneman. 1973. Availability: A heuristic for judging frequency and probability. *Cognitive psychology*, 5(2):207–232.
- Svitlana Volkova, Theresa Wilson, and David Yarowsky. 2013. Exploring demographic language variations to improve multilingual sentiment analysis in social media. In *Proceedings of EMNLP*, pages 1815–1827.
- Svitlana Volkova, Glen Coppersmith, and Benjamin Van Durme. 2014. Inferring user political preferences from streaming communications. In *Proceedings of the 52nd annual meeting of the ACL*, pages 186–196.
- Svitlana Volkova, Yoram Bachrach, Michael Armstrong, and Vijay Sharma. 2015. Inferring latent user properties from texts published in social media (demo). In *Proceedings of the Twenty-Ninth Conference on Artificial Intelligence (AAAI)*, Austin, TX, January.
- Hanna Wallach. 2014. Big Data, Machine Learning, and the Social Sciences: Fairness, Accountability, and Transparency. <https://medium.com/@hannawallach/big-data-machine-learning-and-the-social-sciences-927a8e20460d#uhbcr4wa0> Retrieved Feb 24, 2016.
- David Yarowsky and Grace Ngai. 2001. Inducing multilingual POS taggers and NP bracketers via robust projection across aligned corpora. In *Proceedings of NAACL*.