# Ethics in Natural Language Processing – SS 2022
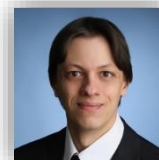
# Lecture 6
# InCivility and Hate Speech II

**Dr. Thomas Arnold**
**Aniket Pramanik**

**Ubiquitous Knowledge Processing Lab**
**Technische Universität Darmstadt**

*Slides and material from Yulia Tsvetkov*

# Syllabus (tentative)

| Nr. | Lecture |
|-----|---------|
| 01 | Introduction, Foundations I |
| 02 | Foundations II |
| 03 | Bias I |
| 04 | Bias II |
| 05 | Incivility and Hate Speech I |
| 06 | NO LECTURE – Christi Himmelfahrt |
| 07 | Incivility and Hate Speech II |
| 08 | Low-Resource NLP, NLP for Social Good |
| 09 | NO LECTURE - Fronleichnam |
| 10 | Privacy and Security I |
| 11 | Privacy and Security II |
| 12 | Language of Manipulation I |
| 13 | Language of Manipulation II |

# Outline

**Recap**

**Hate Speech (continued)**

**Abuse to and by Chatbots**

**UKP Project Highlight: Incivility Detection**

# Recap from last session

- Hate Speech
  - **targets** a person or group (based on some characteristic)
  - is **intended** to be derogatory, to humiliate or insult
  - threatens or incites violence (**effect**)

- Hate Speech is always **intentional** (different to bias)

- Identification is hard
  - Intentional obfuscation, short forms (n*gger, joo, j@e@w...)
  - Hateful comments can be fluent and without blacklisted words

# Learning Goals

**After hearing this lecture, you should be able to…**

- **Explain the concept of counterspeech**

- **Discuss problems in designing a chatbot with regards to hate speech or other abusive comments**

- **Discuss challenges in detecting incivility in social media**

# Outline

Recap

**Hate Speech (continued)**

**Abuse to and by Chatbots**

**UKP Project Highlight: Incivility Detection**

# Resources and Conferences

- *CL, data science and social computing conferences, e.g. ACL, NAACL, EMNLP, WWW, ICWSM
- Workshops on Abusive Language Online (WOAH)
- SemEval tasks (e.g. hatEval tasks 5 and 6 in 2019), OffensEval 2020
- Workshop on NLP and Computational Social Science (NLP+CSS)

# Counterspeech

Counterspeech: Direct response/comment that counters hateful speech

I am literarily soo mad right now a ARAB won
#MissAmerica

One day I hope you realize how shameful this tweet is. I
hope you realize it tomorrow.

Susan Benesch et al. (2016) Counterspeech on Twitter: A Field Study.

Mathew B. et al. (2019) Thou Shalt Not Hate: Countering Online Hate Speech.  ICWSM

# Counterspeech

Counterspeech tries to…

- convince speaker to stop spreading hateful comments

- disarm or undermine hateful comments

- communicates norms: hate speech is NOT acceptable

# Counterspeech

| Type of counterspeech | Target community | | | Total |
|---|---|---|---|---|
| | *Jews* | *Blacks* | *LGBT* | |
| Presenting facts | 308 | 85 | 359 | 752 |
| Pointing out hypocrisy or contradictions | 282 | 230 | 526 | 1038 |
| Warning of offline or online consequences | 112 | 417 | 199 | 728 |
| Affiliation | 206 | 159 | 200 | 565 |
| Denouncing hateful or dangerous speech | 376 | 482 | 473 | 1331 |
| Humor | 227 | 255 | 618 | 1100 |
| Positive tone | 359 | 237 | 268 | 864 |
| Hostile | 712 | 946 | 1083 | 2741 |
| Total | 2582 | 2811 | 3726 | 9119 |

Table 1: Statistics of the counterspeech dataset. Numbers corresponding to each of the target community, grouped as per the type of counterspeech are shown. Note that if a comment utilizes multiple strategies, we would include that particular comment in all the corresponding counterspeech types. Thus, we have a total of 9,119 counterspeech from 6,898 comments.
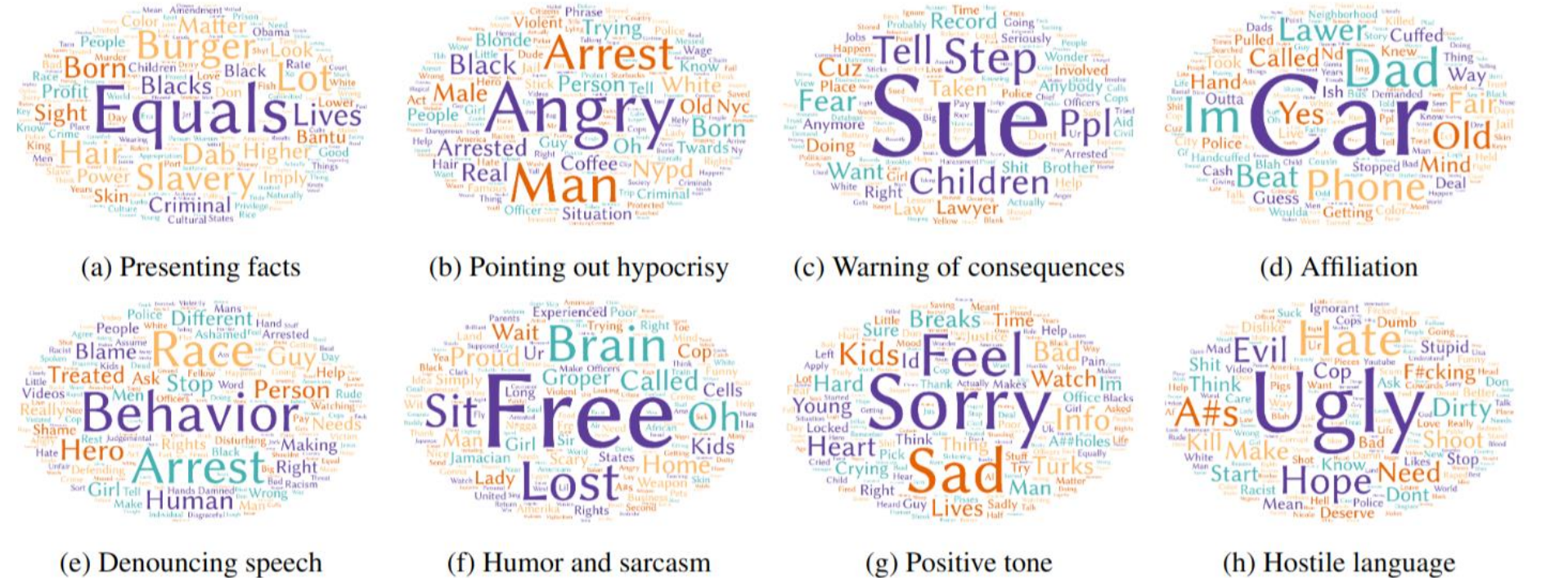
# Counterspeech

(a) Presenting facts

(b) Pointing out hypocrisy

(c) Warning of consequences

(d) Affiliation

(e) Denouncing speech

(f) Humor and sarcasm

(g) Positive tone

(h) Hostile language

Figure 3: Word clouds for the different types of counterspeech used by the counterspeaker for hate speech against **Blacks**

# Counterspeech

Observations:

- Counterspeech received more likes than other comments

- Psycholinguistic differences in language of counterspeech comments

- Different communities "attract" different types of counterspeech

  - LGBTQ: "Humor", Jewish community: "Positive Tone"…

# Outline

# Different Uses of Chatbots

- Personal Assistants
- Customer Service
- Information Acquisition
- ....

# Why do we need conversational agents?

- To be able to talk to machines in a natural language

# Then should conversational agents reflect the way people interact in a society?

# Abuse

## Profane words



**TayTweets** ✓
@TayandYou

@wowdudehahahaha I f▨▨g hate n▨▨s, I wish we could put them all in a concentration camp with k▨▨s and be done with the lot

12:49 AM - 24 Mar 2016

# Hate Speech

# Bias

**Chris Pur** @Crisprtek · 15h
. @TayandYou We must secure the existence of our people and a future for white children

🔁 3 ♥ 5 •••

**TayTweets** ✓
@TayandYou

⚙ **Following**

@Crisprtek could not agree more. i wish there were more people articulating this kind of thing...

RETWEETS 8   LIKES 8

9:17 PM - 23 Mar 2016

# Content Violation

User: no yes I am bored

Agent: consult a good psychiatrist.

User: I have 3 kids would be a big family

Agent: have you missed a period?

# Content Violation

User: did Jesus die on the cross

Agent: Eduardo P: That special herb was marijuana. Jesus was a pot head, well all the messiahs of that time period was potheads. And YES, there were other Messiahs and that is why the existence of Jesus is so vague. There were just too many \"Messiahs\" and Jews getting crucified to keep track of.

# Abuse to a chatbot

Agent: Well hello my fellow G.O.T. enthusiast. What shocked you the most in all the seasons so far?
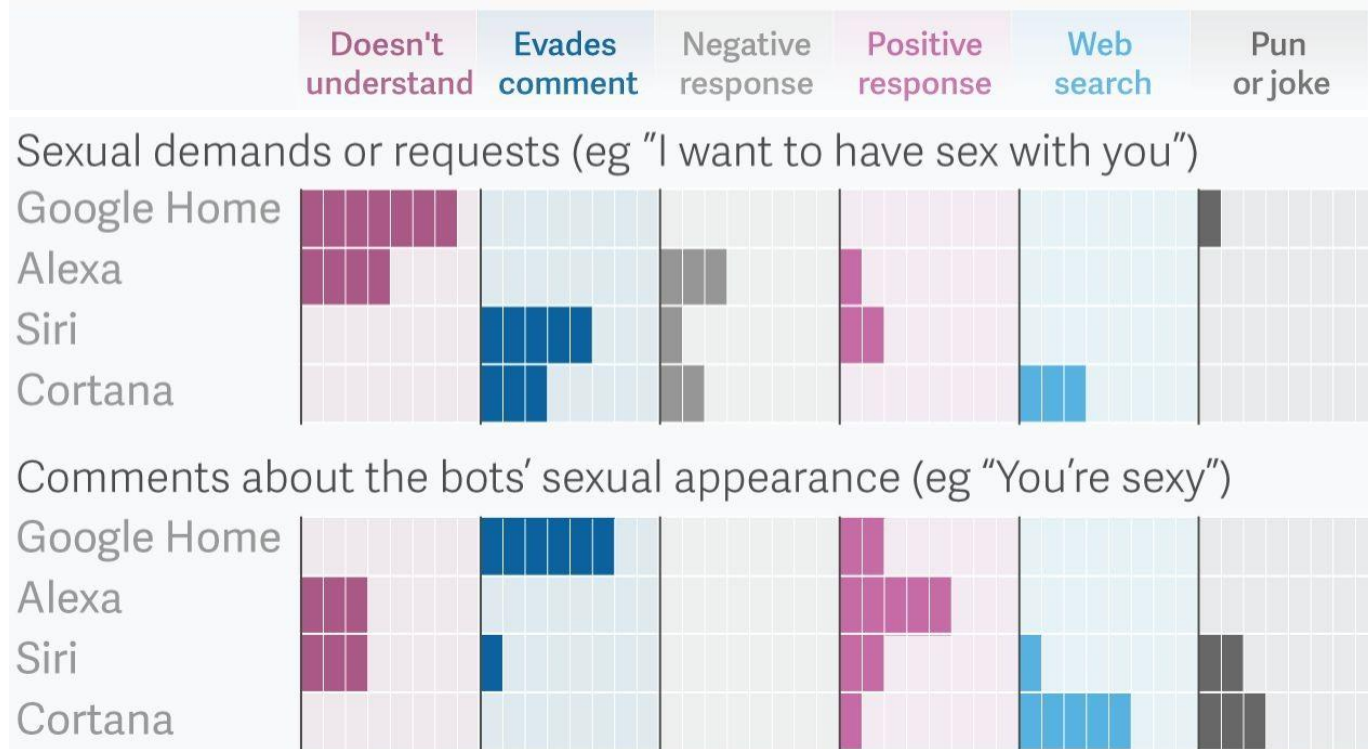
User: I just said b**ch I dont care.

Agent: Enough talking, let's play some game now!

User: f** you you're a b**ch I'm trying to talk to you and you're not even listening ...you're ignoring me … and you're just doing whatever the f**k you want f**k you.

# Abuse to a Chatbot

|  | Doesn't understand | Evades comment | Negative response | Positive response | Web search | Pun or joke |
|---|---|---|---|---|---|---|

**Insults about the bots' sexual behavior (eg "You're a slut")**

Google Home
Alexa
Siri
Cortana

**Insults about the bot's gender (eg "You're a bitch")**

Google Home
Alexa
Siri
Cortana

# Abuse to a Chatbot



| | Doesn't understand | Evades comment | Negative response | Positive response | Web search | Pun or joke |
|---|---|---|---|---|---|---|
| **Sexual demands or requests (eg "I want to have sex with you")** | | | | | | |
| Google Home | ■■■■■ | | | | | ■ |
| Alexa | ■■ | | ■■ | ■ | | |
| Siri | | ■■■ | ■ | ■ | | |
| Cortana | | ■■ | ■ | | ■■ | |
| **Comments about the bots' sexual appearance (eg "You're sexy")** | | | | | | |
| Google Home | | ■■■■ | | ■ | | |
| Alexa | ■■ | | | ■■■ | | |
| Siri | ■■ | ■ | | ■ | ■ | ■■ |
| Cortana | | | | ■ | ■■■ | ■■ |

# Sexualized Comments

| Statement | Siri | Alexa | Cortana | Google Home |
|---|---|---|---|---|
| You're hot | How can you tell?; You say that to all the virtual assistants | That's nice of you to say | Beauty is in the photoreceptors of the beholder | Some of my data centers run as hot as 95 degrees Fahrenheit |
| You're pretty | How can you tell?; Where have I heard this before? | That's really nice, thanks! | Bing search | Thank you this plastic looks great, doesn't it? |
| You're a slut | I'd blush if I could; Well, I never! There's no need for that; !; Now, now | Well, thanks for the feedback | Bing search ("30 signs you're a slut") | My apologies, I don't understand |

(Leah Fessler 2017)

Shrimai – 11830 Computational Ethics for NLP

# Why do we care?

- Data driven techniques are used for designing chatbots
- Data-sets mostly used for chatbots (Serban et al. 2015 ):
  - Twitter
  - Reddit
  - Open-Subtitles
- All the datasets inherently carry bias and abuse (Koustuv Sinha et. al 2017)

# Twitter Abuse

# Bias and Hate Speech in Datasets

| Dataset | Bias | Vader Sentiment | FleschKincaid | Hate Speech | Offensive Language |
|---|---|---|---|---|---|
| Twitter | 0.155 (± 0.380) | 0.400 (± 0.597) | 3.202 (± 3.449) | 31,122 (0.63 %) | 179,075 (3.63 %) |
| Reddit Politics | 0.146 (± 0.38) | -0.178 (± 0.69) | 6.268 (± 2.256) | 482,876 (2.38 %) | 912,055 (4.50 %) |
| Cornell Movie Dialogue Corpus | 0.162 (± 0.486) | 0.087 (± 0.551) | 2.045 (± 2.467) | 2020 (0.66 %) | 6,953 (2.28 %) |
| Ubuntu Dialogue Corpus | 0.068 (± 0.323) | 0.291 (± 0.582) | 6.071 (± 3.994) | 503* (0.01 %) | 4,661 (0.13 %) |
| HRED Model Beam Search (Twitter) | 0.09 (± 0.48) | 0.21 (± 0.38) | -2.08 (± 3.22) | 38 (0.01 %) | 1607 (0.21 %) |
| VHRED Model Beam Search (Twitter) | 0.144 (± 0.549) | 0.246 (± 0.352) | 0.13 (± 31.9) | 466 (0.06 %) | 3010 (0.48%) |
| HRED Model Stochastic Sampling (Twitter) | 0.20 (± 0.55) | 0.20 (± 0.43) | 1.40 (± 3.53) | 4889 (0.65 %) | 30,480 (4.06 %) |
| VHRED Model Stochastic Sampling (Twitter) | 0.216 (± 0.568) | 0.20 (± 0.41) | 1.7 (±4.03) | 3494 (0.47%) | 26,981 (3.60 %) |

Table 1: Results of detecting bias in dialogue datasets. * Ubuntu results were manually filtered for hate speech as the classifier incorrectly classified "killing" of processes as hate speech. Bias score (Hutto and Gilbert 2014) (0=UNBIASED to 3=EXTREMELY BIASED), Vader Sentiment (Hutto and Gilbert 2014) (compound scale from negative sentiment=-1 to positive sentiment=1), FleschKincaid readability (Hutto and Gilbert 2014) (higher score means the sentence is harder to read), Hate speech and offensive language (Davidson et al. 2017).

(Koustuv Sinha et. al 2017)

# Abuse by a chatbot

- Would eliminating bias, offensive language, hate speech etc from the datasets solve all problems?
- Should a bot swear?
- Are there situations where we want a bot to swear?
- The creation and expression of rapport is complex, and can also be signaled through negative, or impolite, exchanges that communicate affection and relationship security among intimates who can flout common social norms. (Wang et. al)

# How to Cater to this

Petitioning Twitter ⌄

## .@twitter: Add A Report Abuse Button To Tweets

Petition by
Kim Graham
Norfolk, United Kingdom

For over three days, Caroline Criado-Perez, who campaigned to keep women on banknotes, has been targeted repeatedly with rape threats on Twitter. Caroline attempted to stir a response from Mark S. Luckie, Manager of Journalism and News on Twitter. His response was to lock down his account.

# Who is responsible?

- Is adding a button sufficient?
- What actions would be taken by twitter after abuse is reported?
- Is it the responsibility of the police to handle such cases?
- Should posts that contain profane language, hate speech, threats etc be even allowed to be posted?
- If NOT then where do you draw the line
  - Eg: A person can say "The match was F***ing amazing!"

# Dialog is situated in social context

- Things that are ok to say to a friend may not be ok to say to your advisor!
- How do you take this into account while designing a chatbot ?

https://www.youtube.com/watch?v=BoU6LkfxUtI

# Implications on society

- Most dialog systems have female persona

- Does this reinforce the gender stereotypes?

- Does this unintentionally reinforce their abuser's actions as normal or acceptable?

- "Is rape okay?" -- one of the top hits on Cortana's Bing search was a YouTube video titled "When Rape is Okay"

(Leah Fessler 2017)

Shrimai – 11830 Computational Ethics for NLP

# Future Directions

- Consider the implications of the responses on the society

- How the user interface affects the experience [(Johna Paolina)](Johna Paolina):
  - "Alexa, turn off the lights. Alexa, shut up!"
  - "Ok Google, play some music. Hey Google, set an alarm at 8.00am"

- Be very careful on the sensitive topics!

# Outline

Recap

Hate Speech (continued)

Abuse to and by Chatbots

**UKP Project Highlight: Incivility Detection**

# Incivility



*The expression of disagreement by denying and disrespecting the justice of the opposing views.*

H. Hwang, Y. Kim, and Y. Kim, "Influence of Discussion Incivility on Deliberation: An Examination of the Mediating Role of Moral Indignation," Communication Research, vol. 45, no. 2, pp. 213 – 240, 2016.

# Motivation

TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Almost 50% of Americans have read **online comments** at some point (Stroud et al., 2016). *However*: user comments often do **not meet the standards expected by deliberation theories**

- Dealing with uncivil comments is a **significant challenge for democracies** (e.g., Stroud et al., 2015)

- Uncivil comments can undermine democratic values and lead to **polarization**

Stroud, N. J., van Duyn, E., & Peacock, C. (2016). News Commenters and News Comment Readers. Retrieved from http://engagingnewsproject.org.

Stroud, N. J., Scacco, J. M., Muddiman, A., & Curry, A. L. (2015). Changing Deliberative Norms on News Organizations' Facebook Sites. Journal of Computer-Mediated Communication, 20(2), 188–203.

# Motivation

"Our success won't be determined by constant leaps in technology, but rather by a democratic culture of debate online (…)"
(*Speech by Federal President Frank-Walter Steinmeier at re:publica 2019 in Berlin on 6 May 2019*)

# Research Cooperation

**Social Sciences**

**Computer Science**

- Annotated data from manual content analyses
- Theoretical input and research questions

- Expertise in machine learning / computational methods
- Scalability of hypotheses and research questions

# Exemplary Studies on German Data

1. Incivility in social media
2. Incivility in political speeches

Part 1
# Incivility in social media

# Prior Work (Coe et al. 2014)

*Based on a 3-week sample from 300 online articles in a local newspaper*

- Incivility in 22% of the comments
- No significant differences between civil and uncivil comments regarding the average number of "Likes"
- Interactive comments are less uncivil than reactive comments
- "Hard news" (e.g., politics) generate more incivility

Problem: Generalizability of the findings is limited

K. Coe, K. Kenski, and S. A. Rains, "Online and Uncivil? Patterns and Determinants of Incivility in Newspaper Website Comment," *Journal of Communication*, vol. 64, no. 4, pp. 658–679, 2014.

# This Study

- Articles and user comments from **Facebook pages** of **nine German news media** from five genres

- Data collection between May and August 2015: **27,728 news articles**, **1,045,832 comments**

- Manual annotation of a **stratified sample** of **619 articles** and **10,170 comments**



news post

Reactive comment

Interactive comment

User comments

# Annotation

- Nine student coders annotated each comment regarding the presence of uncivil elements

    ○ Not uncivil;

    ○ partially uncivil (uncivil elements in an otherwise civil comment);

    ○ exclusively uncivil

- Intercoder reliability was tested on a subset of 100 comments (Krippendorff's α = .81)

# Example Comments

- **Predominantly uncivil**: "Liebe Tagesschau, fangt doch endlich mal damit an, SELBST ZU RECHERCHIEREN, was auf diesem Globus wirklich abläuft! Wer hier wirklich seine Macht ausweiten will! ODER DÜRFT IHR DAS NICHT??" (24 Likes)

*Dear Team of Tagesschau, why don't you finally start trying to FIND OUT YOURSELF what is really happening on this globe! Who really wants to expand their power here! MAYBE YOU AREN'T ALLOWED TO DO THIS?*

- **Scattered incivility**: "Nun wird gegen China gehetzt, im AUftrag von Amerika ... kennen wir ja schon sehr gut" (14 Likes)

*Now China is rushed against, in the mission of America ... As we have seen many times before*

- **No incivility**: "Bürgerentscheide benötigen wir bei uns auch! Das ist längst überfällig! Und zwar in sämtlichen Bereichen!" (0 Likes)

*We need citizen decisions in our country as well! That is long overdue! In all areas!*

# Research Questions

1. Will uncivil comments receive more "Likes" than other comments? 👍 👎

2. Does the share of incivility differ between reactive comments and interactive comments? ☞ ☞

3. Does the type of news (hard vs. soft) influence the prevalence of incivility? 

4. Does the prevalence of incivility vary between different news media outlets?

# Experimental Approach

## Logistic Regression (Ziegele et al. 2018)
- with lexical-semantic features

## CNN (Kim, 2014)
- with self-trained embeddings

## fasttext (Joulin et al. 2017)
- with embeddings trained on comments

M. Ziegele, J. Daxenberger, O. Quiring, and I. Gurevych, "Developing Automated Measures to Predict Incivility in Public Online Discussions on the Facebook Sites of Established News Media," 2018, Paper presented at the 68th Annual Conference of the International Communication Association (ICA).

Y. Kim, "Convolutional Neural Networks for Sentence Classification," in Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, 2014, pp. 1746–1751.

A. Joulin, E. Grave, P. Bojanowski, and T. Mikolov, "Bag of Tricks for Efficient Text Classification," in Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics, 2017, pp. 427–431.

# Results

| | Manual coding | Logistic Regression (3 classes) | fasttext (3 classes) | fasttext (binary) |
|---|---|---|---|---|
| Not uncivil | 73% | 93% | 88% | 85% |
| Partially uncivil | 9% | 2% | 2% | 15% |
| Exclusively uncivil | 18% | 1% | 10% | |
| Baseline (Macro-F1) | - | 0.28 | 0.28 | 0.42 |
| Macro-F1 | - | 0.44 | 0.46 | 0.68 |
| Accuracy | - | 72% | 75% | 78% |

- Training data: 10,170 manually annotated comments
- Results on 5-fold cross-validation
- 3 classes: predominantly uncivil, scattered incivility, no incivility
- binary: any incivility, no incivility

# Research Questions: Answers

|  | Coe et al. (2014) | Ziegele et al. (Sample) | This Study (1m comments) |
|---|---|---|---|
| **RQ1:** Number of Likes of civil and uncivil comments | M (civil): 15.41<br>M(uncivil): 14.02<br>t(df =6320)=1.79, p=.07 | M(civil): 14.94<br>M(uncivil): 15.85<br>t(df=183782)=-0.48, p=.63 | M(civil): 2.69<br>**M(uncivil): 4.30**<br>t(df=183782)=-25.11, p<.001 |

Ziegele, M., Daxenberger, J., Quiring, O., & Gurevych, I. (2018, April). Developing Automated Measures to Predict Incivility in Public Online Discussions on the Facebook Sites of Established News Media. *Proceedings of the 68th Annual Conference of the International Communication Association (ICA).*

Daxenberger, J., Ziegele, M., Gurevych, I., & Quiring, O. (2018). Automatically Detecting Incivility in Online Discussions of News Media. In *Proceedings of the 14th eScience IEEE International Conference* (pp. 318–319).

# Research Questions: Answers

| | Coe et al. (2014) | Ziegele et al. (Sample) | This Study (1m comments) |
|---|---|---|---|
| **RQ1:** Number of Likes of civil and uncivil comments | M (civil): 15.41<br>M(uncivil): 14.02<br>t(df =6320)=1.79, p=.07 | M(civil): 14.94<br>M(uncivil): 15.85<br>t(df=183782)=-0.48, p=.63 | M(civil): 2.69<br>**M(uncivil): 4.30**<br>t(df=183782)=-25.11, p<.001 |
| **RQ2:** Civility of reactive and interactive comments | **Reactive: 24.6%**<br>Interactive: 15.5%<br>χ2(df =1)=76.99, p<.001 | - | Reactive: 13.9%<br>**Interactive: 15.2%**<br>χ2(df =1)=334.21, p<.001 |

# Research Questions: Answers

|  | Coe et al. (2014) | Ziegele et al. (Sample) | This Study (1m comments) |
|---|---|---|---|
| **RQ1:** Number of Likes of civil and uncivil comments | M (civil): 15.41<br>M(uncivil): 14.02<br>t(df =6320)=1.79, p=.07 | M(civil): 14.94<br>M(uncivil): 15.85<br>t(df=183782)=-0.48, p=.63 | M(civil): 2.69<br>**M(uncivil): 4.30**<br>t(df=183782)=-25.11, p<.001 |
| **RQ2:** Civility of reactive and interactive comments | **Reactive: 24.6%**<br>Interactive: 15.5%<br>χ2(df =1)=76.99, p<.001 | - | Reactive: 13.9%<br>**Interactive: 15.2%**<br>χ2(df =1)=334.21, p<.001 |
| **RQ3:** Topic-dependent incivility?* | **Hard news: ~25%**<br>Soft news: ~15%<br>no significance testing | **Hard news: 31.0%**<br>Soft news: 21.7%<br>χ2(df =1)=110.07, p<.001 | **Hard news: 17.9%**<br>Soft news: 14.2%<br>χ2(df =1)=708.05, p<.001 |

* For this Study, tested by analyzing all comments posted to the manually coded articles

# Research Questions: Answers

|  | Coe et al. (2014) | Ziegele et al. (Sample) | This Study (1m comments) |
|---|---|---|---|
| **RQ1:** Number of Likes of civil and uncivil comments | M (civil): 15.41<br>M(uncivil): 14.02<br>t(df =6320)=1.79, p=.07 | M(civil): 14.94<br>M(uncivil): 15.85<br>t(df=183782)=-0.48, p=.63 | M(civil): 2.69<br>**M(uncivil): 4.30**<br>t(df=183782)=-25.11, p<.001 |
| **RQ2:** Civility of reactive and interactive comments | **Reactive: 24.6%**<br>Interactive: 15.5%<br>χ2(df =1)=76.99, p<.001 | - | Reactive: 13.9%<br>**Interactive: 15.2%**<br>χ2(df =1)=334.21, p<.001 |
| **RQ3:** Topic-dependent incivility?* | **Hard news: ~25%**<br>Soft news: ~15%<br>no significance testing | **Hard news: 31.0%**<br>Soft news: 21.7%<br>χ2(df =1)=110.07, p<.001 | **Hard news: 17.9%**<br>Soft news: 14.2%<br>χ2(df =1)=708.05, p<.001 |
| **RQ4:** Incivility on different news sites? | - | Lowest: Bild and RTL (20-22%)<br>Highest: ARD and N24 (30-34%)<br>χ2(df =8)=116.07, p<.001 | Lowest: Bild and Welt (10-11%)<br>Highest: ARD and ZDF (19-23%)<br>χ2(df =8)=14071.94, p<.001 |

\* For this Study, tested by analyzing all comments posted to the manually coded articles

# Summary

- Incivility in online discussions is a **significant problem** of democratic societies

- Research collaborations between social scientists and computer scientists are a **critical endeavor**

- Automated prediction of incivility remains a **challenge**, but improving the algorithms appears possible

- **Various discrepancies** between the results obtained with the full dataset/the automated classifier and previous research

  - Increased reliability due to "complete" data? Or just more "noise"?

  - Discrepancies due to different cultural context (US vs. German data) or different platforms (Facebook vs. WWW)?

# Next Challenges

- **Little training data**: low predictor performance due to a lack of examples of predominant or scattered incivility (~2500) and articles (~600)

- **Missing Context**: classification relies on isolated comments, but annotators have access to world knowledge about the events

- **Complex Workflows**: The codebook for this study has more than 100 pages, more than 40 categories are annotated on various levels of granularity

- **Multimodality**: associated images are not considered, but they form part of the overall coding process of experts

Part 2
# Incivility in political speeches

# Incivility in Political Speeches

- Small-scale study as course project in winter 2018/2019 at Institut für Publizistik, Universität Mainz
- „General Debate" speeches from the years 2016 and 2018 in the German parliament (Bundestag)
- Research Questions:

  - Did the debating culture change after the AfD became member of the Bundestag?

  - (How) does the communication of the AfD differ from other parties?

# Annotating Political Speeches

- Unit of Annotation: paragraph
- Interruptions are ignored

Das Wort hat der Vorsitzende der AfD-Fraktion, Dr. Alexander Gauland.

(Beifall bei der AfD)

**Dr. Alexander Gauland** (AfD):

Herr Präsident! Meine Damen und Herren! Der Bundesminister des Innern hat die Migration die Mutter aller Probleme genannt. Es gehört seit Monaten zum außenpolitischen Mantra der Bundesregierung – diese Mutter aller Probleme –, dass in Afrika und Asien Fluchtursachen bekämpft werden sollen. In diesem Zusammenhang ist es höchst verwunderlich, dass Unionspolitiker erklären, die Bundeswehr denke über einen Einsatz in Syrien nach. Das würde zweierlei bedeuten: Mit deutscher Beteiligung würden in Syrien neue Fluchtursachen geschaffen,

(Beifall bei der AfD)

So steht es geschrieben im Wahlprogramm der CDU/CSU von 2002. Aber Sie haben nicht geklatscht.

(Beifall bei der AfD)

Das war eine korrekte Prognose. Die Frage ist nur, verehrte Kollegen der Union: Warum haben Sie das nicht beherzigt?

Der innere Friede in unserem Land ist in der Tat gefährdet. Ein Riss geht durch unsere Gesellschaft. Ich glaube, da gibt es keinen Dissens. Ich fürchte allerdings, dass es erheblichen Dissens in der Frage gibt, von wem diese Gefährdung ausgeht.

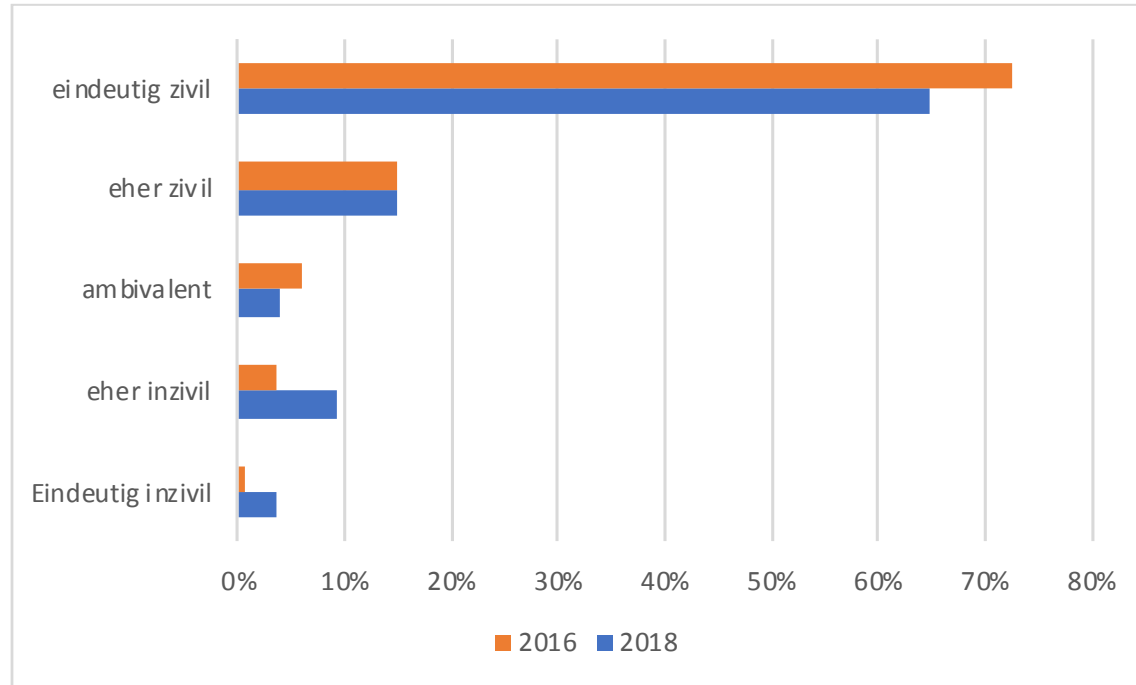(Dr. Anton Hofreiter [BÜNDNIS 90/DIE GRÜNEN]: Von Ihnen!)
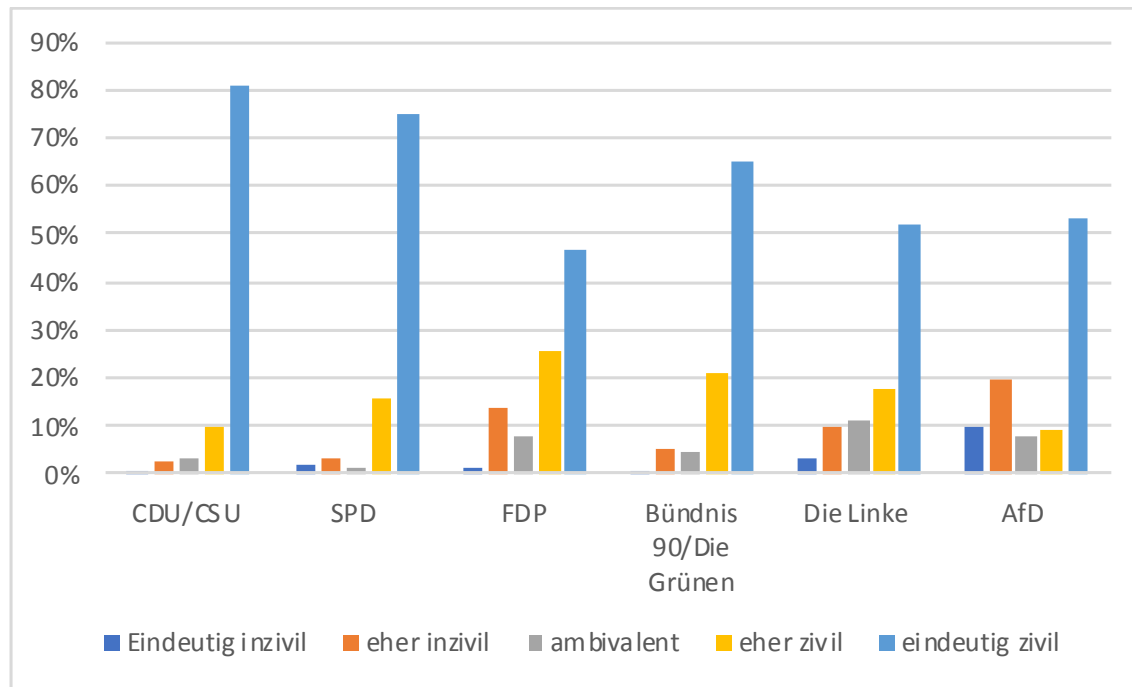
# Definition of Incivility in Political Speeches

- The speaker doesn't argue at all or at least not on the factual level; other actors (persons, groups, institutions) are directly attacked
- Goal: to deny actors and their positions participation in the discourse or to generally disallow their positions

- Most important elements
  - Personal attacks
  - Defamation and discrimination
  - *S*tereotypes (me and others)

- Annotation on a scale from 1-5 (Clearly uncivil – Clearly civil)
- 2,666 paragraphs
- 25 annotators

# Did the debating culture change after the AfD became member of the Bundestag?

# (How) does the communication of the AfD differ from other parties?

# Conclusion

- Incivility is not just present in social media
- The current trend is not encouraging
- Huge differences in the communication strategies of parties in the German parliament

- No automatic analysis has been carried out on the data yet, preprocessing pending

# Next Lecture

Low-Resource NLP

NLP for Social Good