# Social Bots in Japan's 2014 General Election

Stefan Evert[1], **Philipp Heinrich**[1], Fabian Schäfer[2]
[1]Corpus Linguistics Group
[2]Chair of Japanese Studies
**Friedrich-Alexander University of Erlangen-Nuremberg**
July 13, 2017

# Background and Motivation

- **EFI-project** *Exploring the Fukushima Effect*
    - EFI: *Emerging Fields Initiative* of FAU
    - interdisciplinary collaboration:
        - Corpus Linguistics Group
        - Japanese Studies
        - Communication Science
        - Computer Graphics
    - identification and analysis of the propagation of **discourses** in the **(semi-)public sphere**
    - corpora: newspapers and twitter

- Identification of Social Bots
    - important for further processing of data from social media
    - highly media specific
    - here: **Twitter**

# Social Bots in Twitter Data

- Social Bots:
  - "benign" or "malicious"
  - here: duplication of contents
  - retweet activity — creation of **original near-duplicates**

- our political Twitter data sets:
  - **Japanese General Election 2014**
    - 243,914 originals — 1,491,974
    - 307,670 retweets — 2,496,021
  - **German Federal Election 2013**
    - 932,996 originals
    - 720,895 retweets
  - German Federal Election 2017
    - approx. 140,000 originals/day
    - approx. 170,000 retweets/day

# Near-duplicate detection

1. normalization of texts

```python
def normalize(self):
    """ normalizes tweet for deduplication """
    url = r'http[s]?://(?:[a-zA-Z]|[0-9]|[$-_.&+]|[!*\(\),]|(?:%[0-9a-fA-F][0-9a-fA-F]))+'
    mention = r"@\w+"          # twitter user names contain alphanumeric characters
    rt = r'^RT\s'              # RT signs are always at beginning of tweet
    regex = re.compile(r'|'.join([url, mention, rt]))
    n = regex.sub("", self.txt)
    n = re.sub("\s", "", n)
    n = ''.join([c for c in n if not unicodedata.category(c).startswith('P')])  # strip all punctuation marks
    return n.lower()
```

# Normalization Example (Japanese)

#自民党 #議員 #野次 稀代の喧嘩師小泉元総理がぶっ潰す相手は
誰 か : 日比 野庵 旧館 (記事 ... http://t.co/qLUzl6S8DF #産めないのか
#粕人 間 #fxch

⇓

自民党議員野次稀代の喧嘩師小泉元総理がぶっ潰す相手は誰か
日比野庵 旧館記事産めないのか粕人間fxch

# Normalization Example (English)

#Turkey and the Mideast: an ever closer #union? http://t.co/IrTjz9yaNB
\n\n#Egypt #Article #World #News #Haber #Haberler #Facebook #Twitter

$$\Downarrow$$

turkeyandthemideastanevercloserunionegyptarticleworldnewshaberhaberler
facebooktwitter

$$\Uparrow$$

#Turkey and the Mideast: an ever closer #union\nhttp://t.co/mPA2qFMBSx
\n\n#Egypt #Article #World #News #Haber #Haberler #Facebook #Twitter

# Near-duplicate detection

1. normalization

```python
def normalize(self):
    """ normalizes tweet for deduplication """
    url = r'http[s]?://(?:[a-zA-Z]|[0-9]|[$-_@.&+]|[!*\(\),]|(?:%[0-9a-fA-F][0-9a-fA-F]))+'
    mention = r"@\w+"          # twitter user names contain alphanumeric characters
    rt = r'^RT\s'              # RT signs are always at beginning of tweet
    regex = re.compile(r'|'.join([url, mention, rt]))
    n = regex.sub("", self.txt)
    n = re.sub("\s", "", n)
    n = ''.join([c for c in n if not unicodedata.category(c).startswith('P')])  # strip all punctuation marks
    return n.lower()
```

2. mapping of normalized strings onto tweet ids (hashing)
3. extension: hierarchical clustering based on Levenshtein distance

Footprint of a Social Bot net

number of near duplicates

number of user accounts

# Near-duplicate detection

1. normalization

```
def normalize(self):
    """ normalizes tweet for deduplication """
    url = r'http[s]?://(?:[a-zA-Z]|[0-9]|[$-_.&+]|[!*\(\),]|(?:%[0-9a-fA-F][0-9a-fA-F]))+'
    mention = r"@\w+"          # twitter user names contain alphanumeric characters
    rt = r'^RT\s'              # RT signs are always at beginning of tweet
    regex = re.compile(r'|'.join([url, mention, rt]))
    n = regex.sub("", self.txt)
    n = re.sub("\s", "", n)
    n = ''.join([c for c in n if not unicodedata.category(c).startswith('P')])  # strip all punctuation marks
    return n.lower()
```
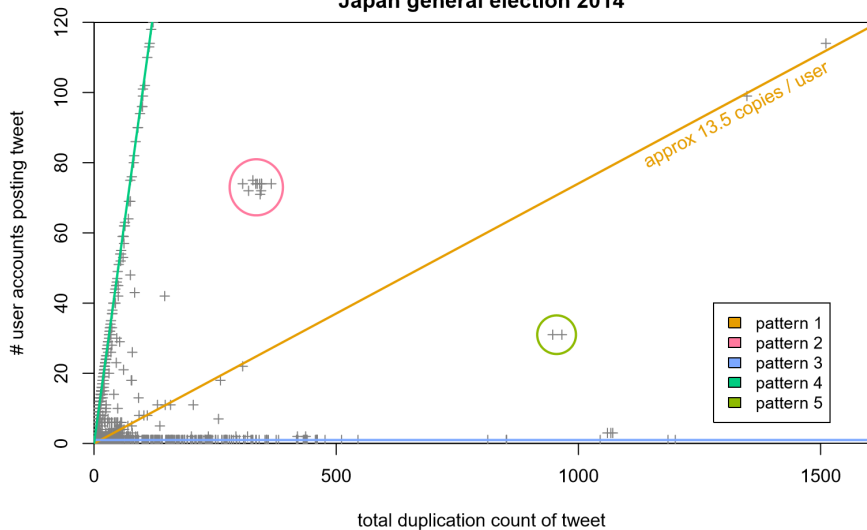
2. mapping of normalized strings onto tweet ids (hashing)
3. extension: hierarchical clustering based on Levenshtein distance

## Footprint of a Social Bot net

$$\frac{\text{number of near duplicates}}{\text{number of user accounts}}$$

**Japan general election 2014**

approx 13.5 copies / user

# user accounts posting tweet

total duplication count of tweet

- pattern 1
- pattern 2
- pattern 3
- pattern 4
- pattern 5

# Content Analysis: Pattern 1, 2+5

- tight clusters and regression line: stable ratios
- botnets disseminating pro-Abe propaganda related to the LDP-campaign in GE 2014
- many of the tweets attack the right-leaning blog "Ponkichi" which takes an anti-LDP, anti-TPP, and anti-migration stance
- judging from the racist terminology, these botnets are operated by a pro-Abe fraction of Japanese internet right-wingers (ネットウヨ)

EMERGING
FIELDS
INITIATIVE

FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

## Example from pattern 1

ABE_CRUSADER on Monday Dec 08, 2014, 10:18:25:

このクソブログ書いてる『ポン吉』ってヤツ最低！→
http://t.co/AI1zBF3vIB 安倍総理を叩くなや！TPP・移民受け入れ・
道州制に反対する倭猿は一匹残らず死ね！！もう一回地震起き
て津波で流されろ！！

$$\Downarrow$$

The guy who is writing this f***ing blog "Ponkichi" is disgusting! →
http://t.co/AI1zBF3vIB Do not attack PM Abe! All those Japanese apes
(wazaru) opposing TPP, the acceptance of immigrants and the regional reform
system should die!! A tsunami from another earthquake should wash them
away!!

## Example from pattern 5

Kimoizo_JAP on Monday Dec 08, 2014, 10:20:24:

我々は、南朝鮮人になりすましながらアメリカ・自民党・安倍総理を叩き、TPPに反対する和猿を陥れる活動をしています。TPP反対派の和猿共を陥れたい方は【てきとう】の活動にご協力ください。

$$\Downarrow$$

We are engaged in activities to get rid of those Japanese apes (wazaru) opposing the USA, the LDP, and PM Abe Shinzō while pretending to be South Koreans. Those who wish to get rid of the anti-TPP group of the Japanese apes please support [appropriate] activities.

# Content Analysis: Pattern 3

- horizontal line at bottom of plot: 57 uniques, 2 accounts
- links to online videos of demonstrations and websites of Japanese ultra-nationalists
- verbal attacks on "anti-Japanese" (反日) parties and their leaders
- often use two triple combinations of hashtags:
  - *#Liberal Democratic Party* combined with
    *#sexual harassment* and
    *#dissolution of parliament*
  - *#diet member* and *#heckling*
- thereby instrumentalizing trending hashtags referring to the sexist slander against female politician Ayaka Shiomura in Tokyo's city parliament to promote nationalist propaganda

**TWEETS**
305K

**FOLLOWING**
275

**FOLLOWERS**
311

**FAVORITES**
5

+ Follow

### 塩村あやかの本音
@Stupid00002

塩村文夫は旅館の女将になりたかった
し、政治家は遊びでやってます。放射能
汚染した東京をぶっ壊わしたい。オリン
ピックなんか失敗すればいい。備掃社の
塩村一ってダレよ？　本ア
カ:twitter.com/shiomura

Tweet to 塩村あやかの本音

**Tweets**   Tweets & replies

塩村あやかの本音 @Stupid00002 · Jun 17
#自民党 #セクハラ #解散 投票率5割なら舛添氏、7割なら"劇的な結果"
- PRESIDENT Online dlvr.it/BDl6Sb #Tokyo Radioactivity
View summary

塩村あやかの本音 @Stupid00002 · Jun 17
#自民党 #セクハラ #解散 安倍首相、解散決断で「民主党をぶっ潰す」
- PRESIDENT Online dlvr.it/BDl54p #Tokyo Radioactivity
View summary

## Content Analysis: Pattern 4

- angle bisector: number of users equals number of duplicates
- does not originate from a bot-net, but is rather rather the product of a political online magazine called Politas (http://politas.jp/)
- a large number of duplicate tweets are created through a "share on Twitter"-button referring to the same article in this magazine
- each tweet is sent deliberately by a genuine user but the content is automatically generated by the share button
- this pattern also shows the complexity of successfully detecting social or political bots: **In future, the behavior of bots will become more sophisticated and human behavior maybe even more bot-like; this convergence might complicate the distinction of "non-human actors" (bots) and "human actors" (users)**

# ポリタス

論点                                                    illustrated by トヨクラタケル

## 【総選挙2014】「選挙なんか行っても無駄」じゃないいくつかの理由〜この一票で変える、のではなく、「この一票で変わる」と思える社会を作り出すために
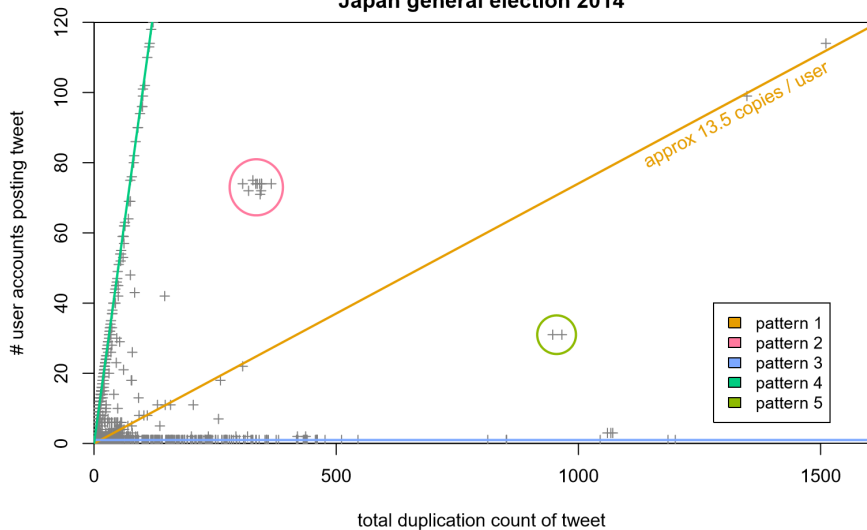
北田暁大（東京大学准教授）    2014年12月9日

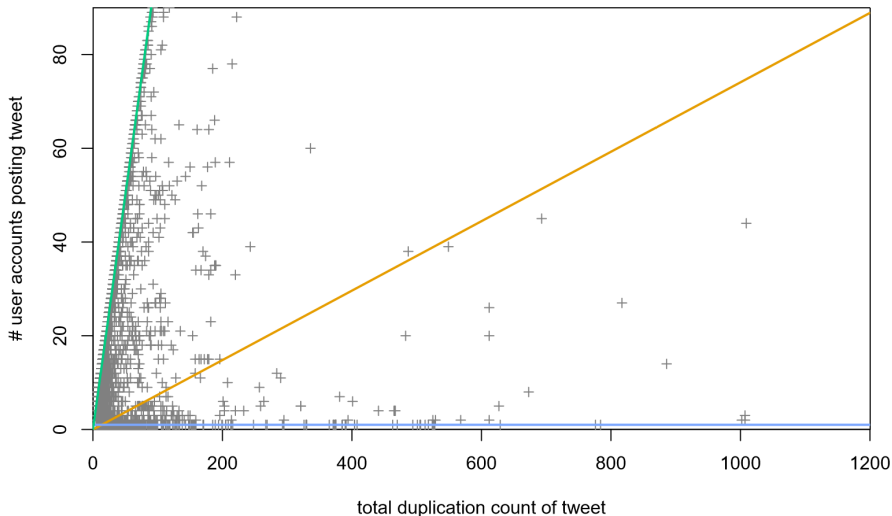Twitterでシェア    Facebookでシェア    HatenaBookmarkでシェア    Mailでシェア
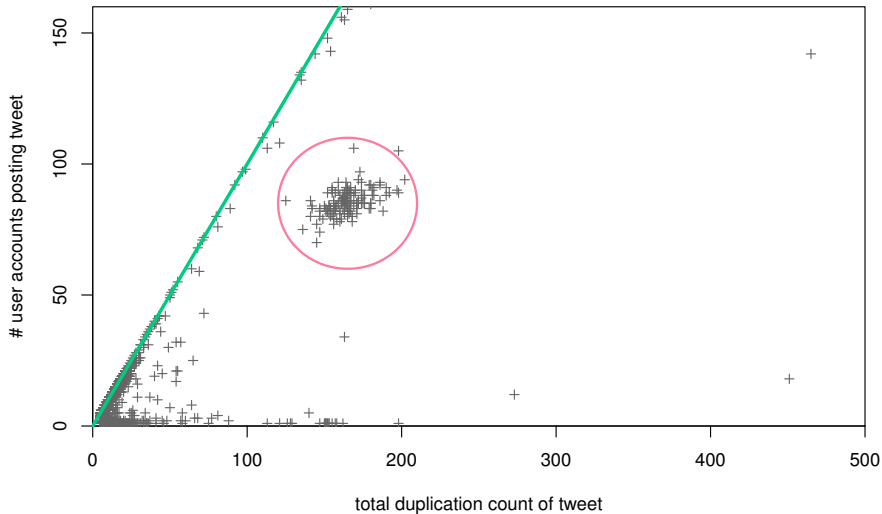
**Japan general election 2014**

approx 13.5 copies / user

# user accounts posting tweet

total duplication count of tweet

- pattern 1
- pattern 2
- pattern 3
- pattern 4
- pattern 5

Japan general election 2014 (Gnip data)

## Validation



Bundestagswahl 2013

# Validation on German data set – introducing Catbots
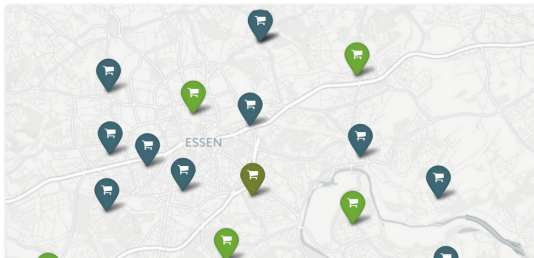
# OpenRuhr

Offene Daten für das Ruhrgebiet

open:ruhr

Hervorgehobener Beitrag

## OpenRuhr – offene Daten für das Ruhrgebiet

Willkommen! OpenRuhr ist eine offene Plattform für mehr OpenData und OpenGoverment im Ruhrgebiet. Mehr Infos, das Ratsinformationsysten-Projekt, Newsletter, Kontakt.

Dieser Beitrag wurde unter Allgemein abgelegt am 22. Juni 2012 von Ernesto Ruge.

Suche

### NEUESTE BEITRÄGE

Wo ist Markt in Essen?
OK Lab Ruhrgebiet: Es geht los! Wichtige Links und die nächsten Aktionen
Einladung: OpenDataDay im Unperfekthaus mit OK Lab Niederrhein
Einladung: OpenData-Hackathon in Bochum
Datenverarbeitung ohne Schnittstelle –

## Conclusion

- social bots are responsible for a large amount of tweets: **61.2% of originals in the Japanese data set are near-duplicates!**
- social bots create and / or use trending hashtags
- social bots are generally **not seen** by human twitterers (no interaction)
- **quantitative corpus analysis** (keyword analysis, co-occurrence analysis) **is prone to distortion by social bots**
- detection of near-duplicates is essential for further processing of data from social media
- **next steps** in this line of research:
    - monitor social bots in the time leading up to the German Federal Election 2017
    - analyze life cycle of social bots
    - combine the social bot footprint presented here with other available indicators

Thanks for listening.
**Questions?**