# Extending Corpus-Based Discourse Analysis for Exploring Japanese Social Media

**Philipp Heinrich**[1] and Fabian Schäfer[2]
[1] *Chair of Computational Corpus Linguistics*, [2] *Chair of Japanese Studies*
**Friedrich-Alexander University of Erlangen-Nuremberg**
September 17, 2018

# Introduction

## Background

- Exploring the *Fukushima Effect*
  - identification and analysis of the tempo-spatial propagation of **discourses** in the **transnational algorithmic public sphere**
  - case study: Fukushima Effect (cf. Gono'i, 2015)
  - data: mass and social media (German, Japanese)
    - ☞ **Japanese Twitter**
  - www.linguistik.fau.de/projects/efe/
  - funded by the **Emerging Fields Initiative** of FAU
- Team:
  - **Chair of Computational Corpus Linguistics**
    Prof. Dr. Stefan Evert, Philipp Heinrich
  - **Chair of Japanese Studies**
    Prof. Dr. Fabian Schäfer, Olena Kalashnikova
  - **Chair of Communication Science**
    Prof. Dr. Christina Holtz-Bacha, Christoph Adrian
  - **Chair of Visual Computing**
    Prof. Dr.-Ing. Marc Stamminger, Jonas Müller

## Research Focus

- methodological foundation: Corpus-Based Discourse Analysis (CDA)
- development of novel techniques (Mixed-Methods Discourse Analysis, MMDA):
  - visualization
  - higher-order collocates
- ultimate goal: assist hermeneutic researchers in interpreting huge amounts of textual data without excessive cherry-picking
- lexical nodes in the case study here:
  - 福島 (Fukushima)
  - 選挙 (elections)
  - 脱原発 (nuclear phase-out)
  - 日本 (Japan) + (原子*)|(原発) (nuclear energy)

☞ focus on methodology

# Research Focus

- methodological foundation: Corpus-Based Discourse Analysis (CDA)
- development of novel techniques (Mixed-Methods Discourse Analysis, MMDA):
    - visualization
    - higher-order collocates
- ultimate goal: assist hermeneutic researchers in interpreting huge amounts of textual data without excessive cherry-picking
- lexical nodes in the case study here:
    - 福島 (Fukushima)
    - 選挙 (elections)
    - 脱原発 (nuclear phase-out)
    - 日本 (Japan) + (原子*)|(原発) (nuclear energy)

☞ focus on methodology

**Introduction**

**Methodology**
Japanese Twitter Corpus in Context
Keywords, Collocates, and Discourse
Visualization

**Case Study: Fukushima Effect**
Overview (Mass Media)
Japanese Twitter Data

**Conclusion**

# Methodology

# Corpora – mass media

## Frankfurter Allgemeine Zeitung (2011–2014)

- statistics:
  - 306,580 articles, 1,656,372 paragraphs
  - 145,055,523 tokens (1,981,726 types)
- linguistic annotation:
  - TreeTagger (tokenization, POS-tagging, lemmatization)

## Yomiuri Shimbun (2011–2015)

- statistics:
  - 1,688,435 articles, 12,757,433 paragraphs
  - 580,518,367 tokens (392,971 types)
- linguistic annotation:
  - MeCab (SUWs)

# Corpora – social media (Twitter)

## German Twitter

- 10,266,835 original posts
- linguistic annotation:
    - tokenization: SoMaJo (Proisl and Uhrig, 2016)
    - POS-tagging: SoMeWeTa (Proisl, 2018)
    - lemmatization: work in progress

## Japanese Twitter

- 411,452,027 original posts
- linguistic annotation:
    - MeCab + special dictionary: ipadic-neologd (Sato et al., 2017)

    + removal of noise: approximately 20% (Schäfer et al., 2017)

## Corpora – social media (Twitter)

### German Twitter

- 10,266,835 original posts
- linguistic annotation:
    - tokenization: SoMaJo (Proisl and Uhrig, 2016)
    - POS-tagging: SoMeWeTa (Proisl, 2018)
    - lemmatization: work in progress

### Japanese Twitter

- 411,452,027 original posts
- linguistic annotation:
    - MeCab + special dictionary: ipadic-neologd (Sato et al., 2017)

        + removal of noise: approximately 20% (Schäfer et al., 2017)

# Corpus-Based Discourse Analysis (CDA)

- CDA means analyzing and deconstructing concordance lines (Baker, 2006)
  - concordances are the essence of discourses
- finding **discourses**: **nodes** + **attitudes**
  - (topic) nodes: defined by *keywords* or (more generally) *corpus queries*
  - attitudes: *collocates* that are retrieved by statistical methods
- examples
  - "refugees as victims" (Baker, 2006)
  - "Fukushima as worst case scenario"

**in practice:**

- look at (*n* best) collocates of topic node

- make up categories on the fly

- categorize manually

## Corpus-Based Discourse Analysis (CDA)

- CDA means analyzing and deconstructing concordance lines (Baker, 2006)
  - concordances are the essence of discourses
- finding **discourses**: **nodes** + **attitudes**
  - (topic) nodes: defined by *keywords* or (more generally) *corpus queries*
  - attitudes: *collocates* that are retrieved by statistical methods
- examples
  - "refugees as victims" (Baker, 2006)
  - "Fukushima as worst case scenario"

**in practice:**

- look at (*n* best) collocates of topic node
- make up categories on the fly
- categorize manually

# Collocates and Keywords

## keywords

- given two frequency lists of lexical items
- perform statistical tests on frequency litss
  - always *viz.* reference corpus
  - measures: log-likelihood, log-ratio, frequency filter

## collocates

- given a definition of a subcorpus
- rate lexical items according to association strength
  - windows vs. segments (**textual co-occurrence**)
  - association measures: see above

EMERGING
FIELDS
INITIATIVE

FAU FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

## From Textual Co-Occurrences to Collocates

- contingency table (cf. Evert, 2008)

|              | $w_2 \in t$ | $w_2 \notin t$ |          |
|--------------|-------------|----------------|----------|
| $w_1 \in t$  | $O_{11}$    | $O_{12}$       | $= R_1$  |
| $w_1 \notin t$ | $O_{21}$  | $O_{22}$       | $= R_2$  |
|              | $= C_1$     | $= C_2$        | $= N$    |

- calculate expected frequencies subject to independence of co-occurrences ($E_{ij}$)
- apply association measure

$$LL(O_{11}, O_{12}, O_{21}, O_{22}) = 2 \sum_{ij} O_{ij} \log \frac{O_{ij}}{E_{ij}},$$

## Extension: Higher-Order Collocates

1. discourse collocates
   - straightforward generalization with respect to textual co-occurrence
   - look at co-occurrence frequencies of tweets that were identified to be part of the discourse at hand (topic + attitude)
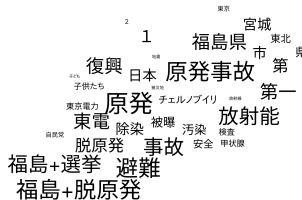   - collocates represent lexical items that play a role in the **discourse**

2. second-order topic-collocates (or attitude-collocates)
   - look at co-occurrence frequencies of one set of lexical items $c$ in tweets that are about a certain topic $t$
   - collocates of $c$ that are particulary important for $t$

# Extension: Higher-Order Collocates

1. discourse collocates
   - straightforward generalization with respect to textual co-occurrence
   - look at co-occurrence frequencies of tweets that were identified to be part of the discourse at hand (topic + attitude)
   - collocates represent lexical items that play a role in the **discourse**

2. second-order topic-collocates (or attitude-collocates)
   - look at co-occurrence frequencies of one set of lexical items $c$ in tweets that are about a certain topic $t$
   - collocates of $c$ that are particulary important for $t$

# Extension: Visualization

- based on high-dimensional word embeddings (Word2Vec) (Mikolov et al., 2013)
  - basis: 133,526,833 deduplicated and preprocessed Japanese tweets collected between February 2017 and June 2018 via the Streaming API
- t-distributed stochastic neighbour-embedding (t-SNE) to project onto two-dimensional plane (van der Maaten and Hinton, 2008)
  - semantically similar items are pre-grouped together
- size of lexical items represents association strength towards (topic) node



2012.11.07 – 2012.12.24    node: 626.11 tw.p.m (5062/8084830)

# Case Study: Fukushima Effect

## Mass media in the aftermath of 3/11 (Heinrich et al., 2018)

### German (FAZ)

- salience of *energy transition* discourse relatively stable (2011–2014)
- *nuclear phase-out* (Atomausstieg) as part of this discourse: sparked shortly after 3/11
  - political actors and issues (*Ethikkommission*, *electricity supply*)
  - economic actors (*RWE*)
  - technological issues (*Stromnetz*)

### Japanese (Yomiuri)

- *nuclear phase-out* (脱原発) in 2011:
  - political actors (菅, 野田, 首相)
  - economic issues (発電, 稼働, 復興)
  - technological aspects (安全, 燃料)
- *nuclear phase-out* in 2014:
  - elections and politics (演説, as used in 街頭演説)
  - fewer words regarding economics (note アベノミクス)

**Introduction**

**Methodology**
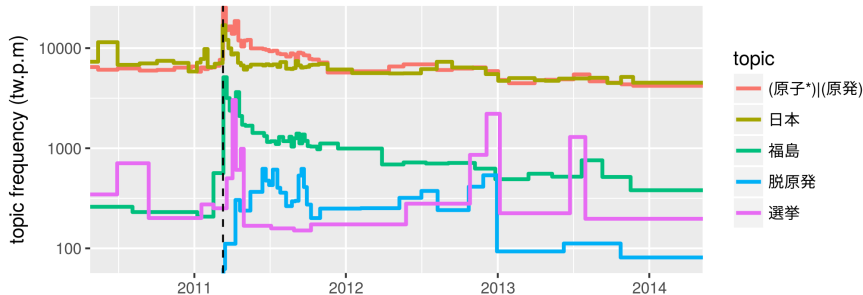
**Case Study: Fukushima Effect**

**Conclusion**

Figure: Frequencies (in tweets per million) of selected topics during the observation period on a logarithmic scale. The dashed line represents March 11, 2011. All observed frequencies peak at or shortly after 3/11.

東京

宮城　茨城

２　３
#save_fukushima
1
青森　岩手
山形　　県
市　号

福島県

http://bit.ly/17n4iz
東日本大震災　　産
第一　第

原発　原発事故
半径

東京電力
毎日新聞　　　東電　放射線
機

社員　原子力発電所
電力会社　風評被害
放射能
冷却

避難　事故
定年　志願

爆発

Figure: Node: 福島 (*Fukushima*).

東京 千葉 茨城
宮城 栃木
2 3 岩手 県
1 福島県 産 号
#save_fukushima 農産物 市 第一 第

http://bit.ly/17n4iz
@mainichijpnews

毎日ｊｐ 原発 野菜
東京電力 原発事故 ほうれん草
毎日新聞 放射能 牛乳
風評被害 放射性物質
東電 事故 放射能漏れ 安全 機
避難
出荷

Figure: Node: 福島 (*Fukushima*).

7

千葉

東京

宮城          茨城

1          震度6弱          岩手          市

#save_fushima          震源          #genpatsu

震度６弱          中通り

@mainichijpnews          余震          浜通り          第

同人女          地震          産          第一

毎日ｊｐ  原発事故          レベル

東京電力          放射線          放射性物質

毎日新聞  原発  チェルノブイリ

東電          差別          事故          心配          放射能

避難          揺れ          大丈夫

Figure: Node: 福島 (*Fukushima*).

東京

2

1

宮城

福島県 市 東北 県

復興 日本 原発事故 第

地震

子ども

子供たち 被災地 チェルノブイリ 放射線 第一

東京電力 原発 放射能

東電 除染 被曝 汚染 検査

自民党 脱原発 事故 安全 甲状腺

福島+選挙 避難

福島+脱原発

Figure: Node: 福島 (*Fukushima*).

2013.07.22 – 2013.09.09    node: 757.58 tw.p.m (5152/6800569)

Figure: Node: 福島 (*Fukushima*).

石原
石原都知事
#senkyo_enki
政治　　民主主義　べき　民
福島+選挙
選　票　政党　国民
選挙+ドイツ　立候補　　　　　　　　　　　　　　政治家
統一
浦安　選管　候補　有権者
候補
延期 統一地方選挙 議員
市　県
統一地方選 民主党 知事
都　当選
知事選
#ogiri_110401
投票 候補者
@ogiri_tweet
県議
選挙カー
都知事選 #senkyo 投票率
告示

Figure: Node: 選挙 (*elections*).

石原
行っ
き いっ
結果 の 文句
権利 人
選 政策 選挙権 若い人
今回 デモ
行 行く 候補 政治 関心 有権者 民 若者
行か 票 政治家
立候補 都知事選 知事 都知事
都 当選 立候補者
期日前投票 県議 投票率
投票 現職
投票所

Figure: Node: 選挙 (*elections*).

行っ

結果   いう

原発  政治

争点        国民      人

政党            若者

有権者    政治家

今回

行く

行  行か       党   ONLINE  自民党  日本

選挙+脱原発  投票  YOMIURI  票  民主党  自民

未来の党

衆院選  候補者  投票率

開票

投票所

Figure: Node: 選挙 (*elections*).

福島+脱原発　孫　城南信用金庫

財団設立　エネルギー　氏　孫正義

原発　電力　理事長　孫社長

社会

派　自然エネルギー　ソフトＢ

宣言　メッセージ　政策　個人　反原発　作ろう

http://t.co/5Y9hsY5

提言　毎日新聞　原発推進

転換　拠出億

推進

ドイツ+脱原発　日本　円　0　1

信金　ドイツ　デモ

城南

#genpatsu

Figure: Node: 脱原発 (*phasing out nuclear energy*).

Figure: Node: 脱原発 (*phasing out nuclear energy*).

Figure: Node: 脱原発 (*phasing out nuclear energy*).

http://shindanmaker.com/23632

れ

勝っ

べき

こと

的

翻訳

英語

た

ぁぁ

よう

問題

い

RT

から

も

だろう

ある

など

国家

日本語

語

思う

ぉ

文化

政治

報道

外国人

いう

日本人

中国人

米国

国民

マスコミ

#_ki

米

民主党

アメリカ

中国

代表

パラグアイ

サッカー

韓国

アジア

海

世界

西

海外

外国

Figure: Node: 日本 (*Japan*).

Figure: Node: 日本 (*Japan*).

原発
れる
れ 語 い
べき
問題
も
よう
的 など
こと ある の 竹島 政府 報道 ニュース
いう 笑 日本人 政治 マスコミ
国家
米 韓国人 米国 国 国民
アメリカ 維新の会
韓国 世界 戦争 代表
中国 外国 西
海外 東

Figure: Node: 日本 (*Japan*).

原発 原爆
前原外相 原点 原子力 原因 れ 原稿 RT
原料 語
原油 問題 原作
べき など
こと も 原則 政治 辞任
ある の
いう 日本人 国 献金
外国人 政治家 外相 東原亜希
米国
世界 石原都知事
組新の党 ※
アメリカ 日本 前原 石原慎太郎
韓国 中国 石原 原口
原宿 原
秋葉原 氏

Figure: Discourse Node: 日本 (*Japan*) + (原子\*)|(原発) (*nuclear energy*).

Figure: Discourse collocates of 日本 (*Japan*) + (原子*)|(原発) (*nuclear energy*).

Figure: Discourse collocates of 日本 (*Japan*) + (原子*)|(原発) (*nuclear energy*).

Figure: Discourse collocates of 日本 (*Japan*) + (原子*)|(原発) (*nuclear energy*).

Figure: Second-order topic-collocates of 日本 (*Japan*) in tweets containing (原子\*)|(原発) (*nuclear energy*).

Figure: Second-order topic-collocates of 日本 (*Japan*) in tweets containing (原子*)|(原発) (*nuclear energy*).
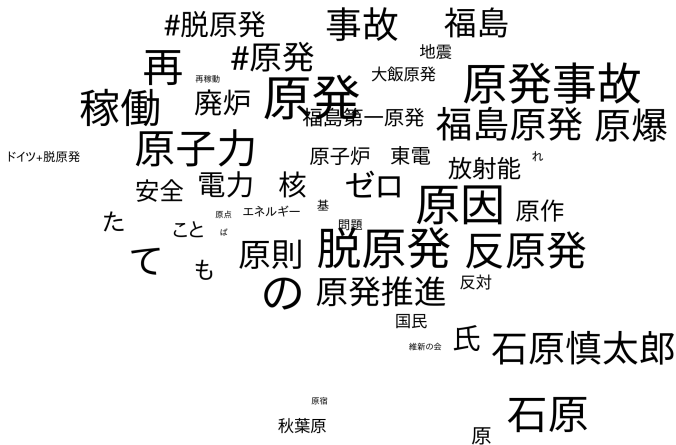
Figure: Second-order topic-collocates of 日本 (*Japan*) in tweets containing (原子\*)|(原発) (*nuclear energy*).

# Qualitative Summary

- 福島 (*Fukushima*)
  - has always been a topic on Twitter
  - important collocates during the observation period are lexical items referring to the accident (原発, 原発事故) and the hashtag #save_fukushima, but also the electric utility holding company 東電 (*TEPCO*)
  - focus shifts to political actors 安倍首相 (*Prime Minister Shinzō Abe*) and the results of and measures taken due to the radioactive accident: 除染 (decontamination), 汚染水 (*contaminated water*), 放射能 (*radioactivity*)

- 選挙 (*elections*)
  - huge peaks in the number of tweets at dates which coincide e. g. with the elections of Tokyo's governor after resignation of 石原 (*Shintaro Ishihara*)
  - further important collocates are 結果 (*results*), 都知事選 (*gubernatorial election*), and 候補(者) (*candidate, candidacy*)
  - end of 2012: most important collocates have shifted towards 自民 (*Liberal Democratic Party*), nuclear power (plants) (原発)
  - actors change

## Qualitative Summary (ctd.)

- 脱原発 (*nuclear phase-out*)
    - enters the debate only a couple of weeks after 3/11
    - whether or not to "break with nuclear energy" is a discussion led elsewhere, e. g. in ドイツ (*Germany*)
    - further important collocates are 福島 (*Fukushima*), 原発 (*nuclear power plant*), and デモ (*demonstration*)
    - another peak in the end of 2012, with political actors as collocates such as 未来の党 (the *Tomorrow Party of Japan*) and 山本太郎 (*Tarō Yamamoto*)
- 日本 (*Japan*) and (原子*)|(原発) (*nuclear energy*)
    - before 3/11: collocates of Japan mostly general (語, other countries)
    - in the aftermath of 3/11: 地震 (*earthquake*), 復興 (*reconstruction*), 原発 (*nuclear power plant*), and 赤十字社 (*red cross*)
    - after 2012: 原発 (*nuclear power plant*) remains an important collocate
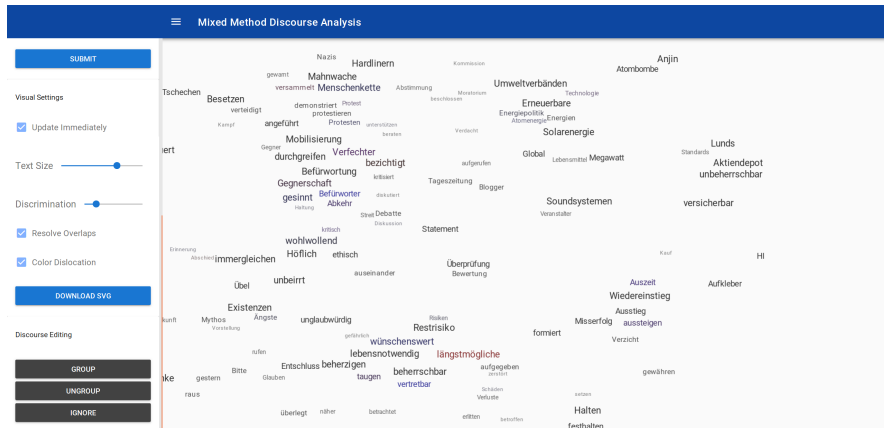
**Conclusion and Future Work**

- CDA of Japanese Twitter data in the aftermath of 3/11
- focus on methodological advancement of the field
  - visualization (ease manual labour)
  - higher-order collocates (triangulate semantics of discourses)
- qualitative empirical level:
  - **nuclear phase-out debate** entered Japanese Twitter only **several weeks after 3/11**
  - salience of discussions about **phasing out nuclear energy** and about nuclear energy in general is quite volatile and **correlates i. a. with elections**
  - **particular parts of the nuclear energy discussion** entered the collocational profile of the very **general discourse** around Japan
- where do we go from here?

## Conclusion and Future Work

- CDA of Japanese Twitter data in the aftermath of 3/11
- focus on methodological advancement of the field
    - visualization (ease manual labour)
    - higher-order collocates (triangulate semantics of discourses)
- qualitative empirical level:
    - **nuclear phase-out debate** entered Japanese Twitter only **several weeks after 3/11**
    - salience of discussions about **phasing out nuclear energy** and about nuclear energy in general is quite volatile and **correlates i. a. with elections**
    - **particular parts of the nuclear energy discussion** entered the collocational profile of the very **general discourse** around Japan
- where do we go from here?

# Towards Mixed-Methods Discourse Analysis

Thanks for listening.
**Questions?**

# References

Paul Baker. *Using Corpora in Discourse Analysis.* Continuum, London, 2006.

Stefan Evert. Corpora and collocations. In Anke Lüdeling and Merja Kytö, editors, *Corpus Linguistics. An International Handbook*, chapter 58. Mouton de Gruyter, Berlin, 2008.

Michel Foucault. *L'Archéologie du savoir.* Éditions Gallimard, Paris, 1969.

Ikuo Gono'i. 2015-nen ANPO, Minshushugi wo futatabi hajimeru wakamono-tachi (ANPO in 2015. The Youth that is restarting Democracy), 2015.

Philipp Heinrich, Christoph Adrian, Olena Kalashnikova, Fabian Schäfer, and Stefan Evert. A Transnational Analysis of News and Tweets about Nuclear Phase-Out in the Aftermath of the Fukushima Incident. In Andreas Witt, Jana Diesner, and Georg Rehm, editors, *Proceedings of the LREC 2018 "Workshop on Computational Impact Detection from Text Data"*, Paris, 2018. ELRA.

Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781, 2013.

Thomas Proisl. SoMeWeTa: A Part-of-Speech Tagger for German Social Media and Web Texts. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC'18)*, 2018.

Thomas Proisl and Peter Uhrig. SoMaJo: State-of-the-art tokenization for German web and social media texts. In Paul Cook, Stefan Evert, Roland Schäfer, and Egon Stemle, editors, *Proceedings of the 10th Web as Corpus Workshop (WAC-X) and the EmpiriST Shared Task*, pages 57–62, Berlin, 2016. Association for Computational Linguistics.

Toshinori Sato, Taiichi Hashimoto, and Manabu Okumura. Implementation of a word segmentation dictionary called mecab-ipadic-neologd and study on how to use it effectively for information retrieval (in japanese). In *Proceedings of the Twenty-three Annual Meeting of the Association for Natural Language Processing*, pages NLP2017–B6–1. The Association for Natural Language Processing, 2017.

Fabian Schäfer, Stefan Evert, and Philipp Heinrich. Japan's 2014 General Election: Political Bots, Right-Wing Internet Activism and PM Abe Shinzō's Hidden Nationalist Agenda. *Big Data*, 5:1 – 16, 2017.

L.J.P van der Maaten and G.E. Hinton. Visualizing High-Dimensional Data Using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 2008.