

NLP of Microblogs, Day 2: Working with Metadata

ESLLI 2016

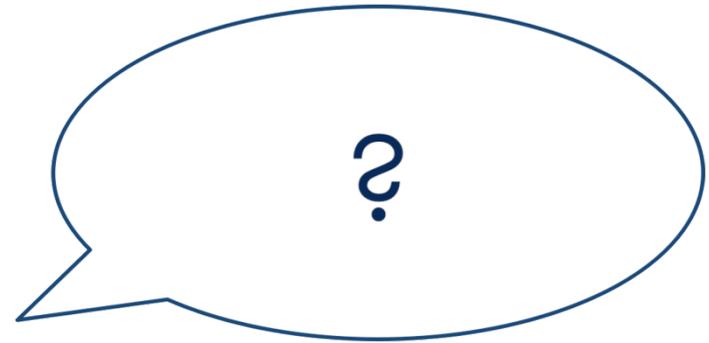
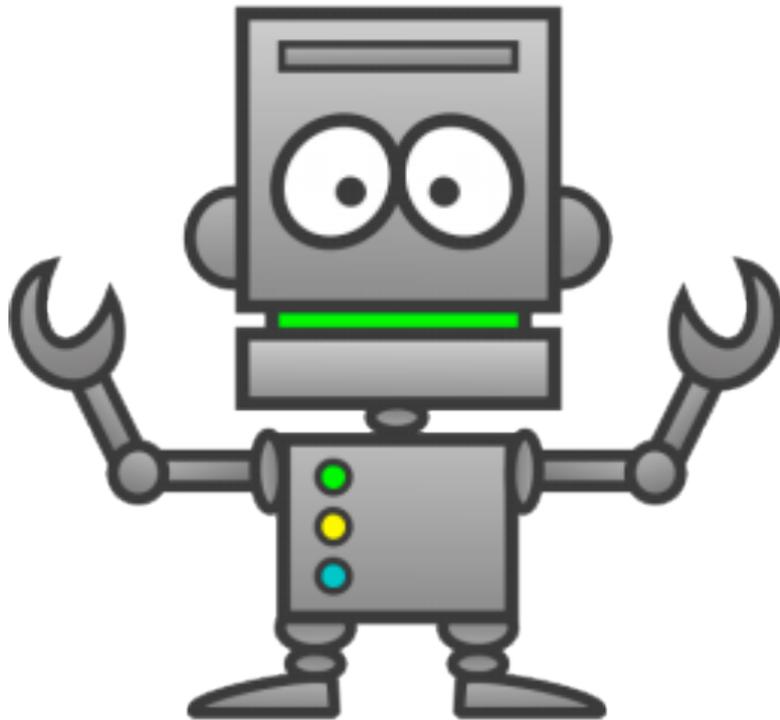
Bolzano/Bozen, Italy

Tatjana Scheffler & Manfred Stede

`tatjana.scheffler@uni-potsdam.de / stede@uni-potsdam.de`

August 23, 2016

Build Twitter bots?



Twitter Metadata

time

geo_coordinates

user profile info

id_str	from_user	text	created_at	time	geo_coord	user_lang	in_reply_to_u	in_reply_to_sc	from_user_id_str	in_reply_to	source	profile_image_un	user_followers_count	user_friends_count	user
7672580	Janelle_Meinung	obwohl das Buch "Wunderbare... https://t.co/jnpd8IdzT	Sun Aug 21 07:17:38 +0000 2016	21/08/2016 08:12:29		de			127902112			http://pbs.twimg.com/profile_images/2712712712/...	27	28	Über...
7672575	strosch9	...	Sun Aug 21 07:10:32 +0000 2016	21/08/2016 08:10:32		de			4011535756			http://pbs.twimg.com/profile_images/432912912/...	43	29	
7672573	strosch9	...	Sun Aug 21 07:09:50 +0000 2016	21/08/2016 08:09:50		de			4011535756			http://pbs.twimg.com/profile_images/432912912/...	43	29	
7672565	Flign2806	@c_m_barnhart Gute nacht :)	Sun Aug 21 07:06:26 +0000 2016	21/08/2016 08:06:26		en	7165361763	c_m_barnhart	491850402	767256094		http://pbs.twimg.com/profile_images/493460493/...	493	460	Pod...
7672563	TunedGlorette	Good Night! All animals friend! 🐾🐾🐾 https://t.co/889Qzr3Dw	Sun Aug 21 07:05:52 +0000 2016	21/08/2016 08:05:52		en-gb			701440480769404928			http://pbs.twimg.com/profile_images/383783783/...	383	344	
7672560	roadmendes	@lovelambuae gute nacht	Sun Aug 21 07:04:39 +0000 2016	21/08/2016 08:04:39			2872615833	lovelambuae	2839455557	767230588		http://pbs.twimg.com/profile_images/782782782/...	382	106	
7672558	hMAMmeyer2	...	Sun Aug 21 07:03:54 +0000 2016	21/08/2016 08:03:54		en			4920008313			http://pbs.twimg.com/profile_images/187187187/...	187	484	belj...
7672558	AlexMeins	Ich lag mich z pannen, Gute Nacht xd	Sun Aug 21 07:03:51 +0000 2016	21/08/2016 08:03:51		de			1901319685			http://pbs.twimg.com/profile_images/748748748/...	748	58	Bel...
7672555	strosch9	...	Sun Aug 21 07:02:26 +0000 2016	21/08/2016 08:02:26		de			4777377435			http://pbs.twimg.com/profile_images/192192192/...	192	182	
7672544	rianaTJ_	RT @DeanisTJ_ Gute Nacht :)	Sun Aug 21 06:58:10 +0000 2016	21/08/2016 07:58:10		de			2861956357			http://pbs.twimg.com/profile_images/362362362/...	362	91	#Ma...
7672541	Stilleforse	@Josenweifer @henris @winterschuppe gute Nacht	Sun Aug 21 06:57:08 +0000 2016	21/08/2016 07:57:08		de	2293432255	Dosenwerfer	3292060903	767254101		http://pbs.twimg.com/profile_images/791791791/...	791	265	
7672535	strosch9	@strosch9 ich hatte gute Nacht gesagt! .. Du solltest doch schlafen :)	Sun Aug 21 06:54:38 +0000 2016	21/08/2016 07:54:38		de	1235051204	katastrawfairy	3355544741	767207712		http://pbs.twimg.com/profile_images/161616161/...	16	59	#de...
7672527	Meine_TJ_	RT @DeanisTJ_ Gute Nacht :)	Sun Aug 21 06:51:17 +0000 2016	21/08/2016 07:51:17		de			3028395099			http://pbs.twimg.com/profile_images/147147147/...	147	16	#de...
7672526	germanykooal	@smla_bach kann übermorgen besser du mal die Schicht Gute nacht!	Sun Aug 21 06:50:59 +0000 2016	21/08/2016 07:50:59		en	2803486586	carola_bach	269804632	767254222		http://pbs.twimg.com/profile_images/448448448/...	448	289	Prag...
7672522	Sir_Natans	...	Sun Aug 21 06:49:27 +0000 2016	21/08/2016 07:49:27		de			3590643748			http://pbs.twimg.com/profile_images/202020202/...	20	16	
7672511	strosch9	RT @DeanisTJ_ Gute Nacht :)	Sun Aug 21 06:45:05 +0000 2016	21/08/2016 07:45:05		de			3238305772			http://pbs.twimg.com/profile_images/797979797/...	79	113	
7672506	strosch9	RT @DeanisTJ_ Gute Nacht :)	Sun Aug 21 06:42:53 +0000 2016	21/08/2016 07:42:53		de			2971698659			http://pbs.twimg.com/profile_images/126126126/...	126	90	bab...
7672499	KreativesName	RT @philinked Gute Nacht	Sun Aug 21 06:40:10 +0000 2016	21/08/2016 07:40:10		de			2520654340			http://pbs.twimg.com/profile_images/557557557/...	557	45	Ham...
7672498	KreativesName	RT @philinked Gute Nacht	Sun Aug 21 06:40:07 +0000 2016	21/08/2016 07:40:07		de			2520654340			http://pbs.twimg.com/profile_images/557557557/...	557	45	Ham...
7672498	KreativesName	RT @philinked Gute Nacht	Sun Aug 21 06:40:04 +0000 2016	21/08/2016 07:40:04		de			2520654340			http://pbs.twimg.com/profile_images/557557557/...	557	45	Ham...
7672495	carstenvic	Gute Nacht!	Sun Aug 21 06:38:46 +0000 2016	21/08/2016 07:38:46		de			2742242917			http://pbs.twimg.com/profile_images/141141141/...	141	230	köln
7672487	Cook_Gourmet	guter tips https://t.co/3fnduh	Sun Aug 21 06:35:39 +0000 2016	21/08/2016 07:35:39		tr			94107162			http://pbs.twimg.com/profile_images/132132132/...	132	143	
7672486	Ivylita	Gute Nacht!	Sun Aug 21 06:34:58 +0000 2016	21/08/2016 07:34:58		de			4817308534			http://pbs.twimg.com/profile_images/132132132/...	132	204	
7672483	Poweranz500	RT @Cmarfan15 Gute Nacht #Ostsee https://t.co/cVcVOMSPv	Sun Aug 21 06:33:54 +0000 2016	21/08/2016 07:33:54		de			2191795613			http://pbs.twimg.com/profile_images/27812781278/...	2781	2211	Deut...
7672478	libber_1	...	Sun Aug 21 06:31:51 +0000 2016	21/08/2016 07:31:51		de			2800511612			http://pbs.twimg.com/profile_images/969696969/...	96	41	GRU...
7672477	1900969Heiduk	da bleibt nur ein Kerzenlichtabendrot ☺☺☺	Sun Aug 21 06:31:22 +0000 2016	21/08/2016 07:31:22		de			4192867587			http://pbs.twimg.com/profile_images/101101101/...	101	120	Halle...
7672465	xlenasaaxi	RT @DeanisTJ_ Gute Nacht :)	Sun Aug 21 06:26:48 +0000 2016	21/08/2016 07:26:48		de			3214710001			http://pbs.twimg.com/profile_images/323232323/...	32	24	Nord...
7672450	heilmadame	Gute Nacht!	Sun Aug 21 06:20:59 +0000 2016	21/08/2016 07:20:59		de			102331010			http://pbs.twimg.com/profile_images/115115115/...	115	64	L I E...
7672450	heilmadame	Gute Nacht!	Sun Aug 21 06:20:59 +0000 2016	21/08/2016 07:20:59		de			102331010			http://pbs.twimg.com/profile_images/115115115/...	115	64	L I E...
7672448	smileoo	RT @DeanisTJ_ Gute Nacht :)	Sun Aug 21 06:20:11 +0000 2016	21/08/2016 07:20:11		de			3435112155			http://pbs.twimg.com/profile_images/152152152/...	152	133	E + S...
7672447	jilane79_	RT @DeanisTJ_ Gute Nacht :)	Sun Aug 21 06:19:44 +0000 2016	21/08/2016 07:19:44		de			754398992948232192			http://pbs.twimg.com/profile_images/101010101/...	10	12	Berli...
7672445	tsukku_gilbert	@sm_13_kun Gute Nacht. ああああ 今日もお疲れさま また明日な!	Sun Aug 21 06:18:57 +0000 2016	21/08/2016 07:18:57		ja	2598978290	kim_13_kun	1181550674	767244256		http://pbs.twimg.com/profile_images/575575575/...	575	471	国...
7672445	vanessarager_	RT @DeanisTJ_ Gute Nacht :)	Sun Aug 21 06:18:45 +0000 2016	21/08/2016 07:18:45		de			3569475135			http://pbs.twimg.com/profile_images/231231231/...	231	100	Colo...
7672435	tsuki_yonehana	RT @haublu42 Gute Nacht und danke für https://t.co/FFagUjARDR	Sun Aug 21 06:15:00 +0000 2016	21/08/2016 07:15:00		ja			2456330712			http://pbs.twimg.com/profile_images/204420442/...	2044	2954	
7672433	smileoo3333	https://t.co/3c3xT2Gc	Sun Aug 21 06:14:13 +0000 2016	21/08/2016 07:14:13		en			2196122099			http://pbs.twimg.com/profile_images/428042804/...	4280	4237	
7672425	AndreasAndy3131	Liebe @Hagenellin	Sun Aug 21 06:10:55 +0000 2016	21/08/2016 07:10:55		de			2344913904			http://pbs.twimg.com/profile_images/102810281/...	1028	674	4816...

in_reply_to

user network

Metadata

1. geographic coordinates
2. time stamps
3. reply info (Thursday!)
4. user network (not covered in this course)
 - ▣ followers
 - ▣ friends
 - ▣ retweets
 - ▣ mentions
5. user profile information

1. geographic coordinates

Geolocation information in Twitter

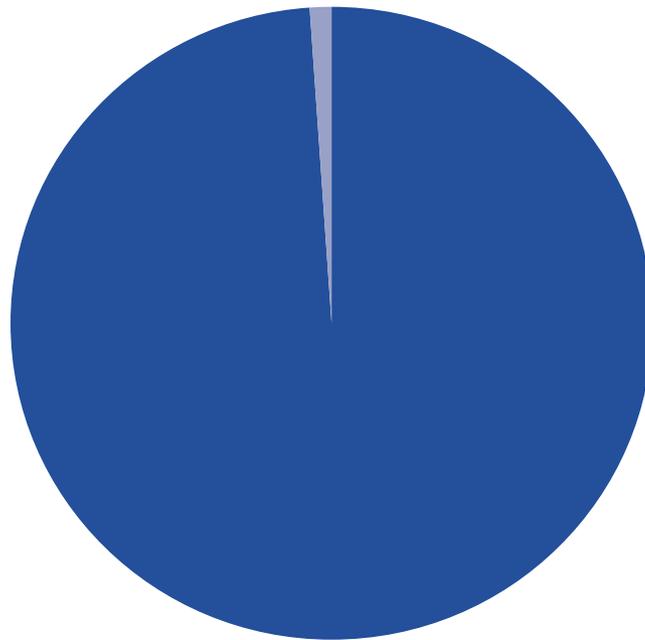
```
place (  
| country = "Germany"  
| place_type = "city"  
| country_code = "DE"  
| name = "Stuttgart"  
| full_name = "Stuttgart, Stuttgart"  
| url = "http://api.twitter.com/1/geo/id/  
e385d4d639c6a423.json"  
| id = "e385d4d639c6a423"  
| bounding_box (  
| | coordinates => Array (1) (  
| | | [...] )  
| | type = "Polygon" )  
| attributes ( ) )
```

Where do Germans tweet from?



images by Norman Rosner

Geo-tagged tweets



German tweets
April 2013

- non geo-tagged:
23,915,825
- geo-tagged:
263,364



Simone Biles @Simone_Biles · 19h

bye Rio De Janeiro 

good things must end, Thank you RIO for unforgettable memories  Being the flag bearer was a cherry on top! Congrats to all the other USA athletes as well. We killed it out there! GO TEAM USA  



  1.3K  8.4K

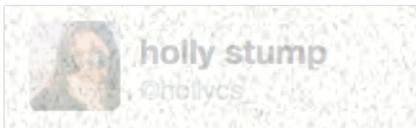


 Follow

Il giardino ad Alexanderplatz, in pieno centro, parla chiaro. Ovunque è così...



 Follow



Last night uberd from bar to wawa, bought the driver a hoagie, uberd back to the bar and my wawa at the bar



Simone Biles @Simone_Biles · Aug 20

beach day 

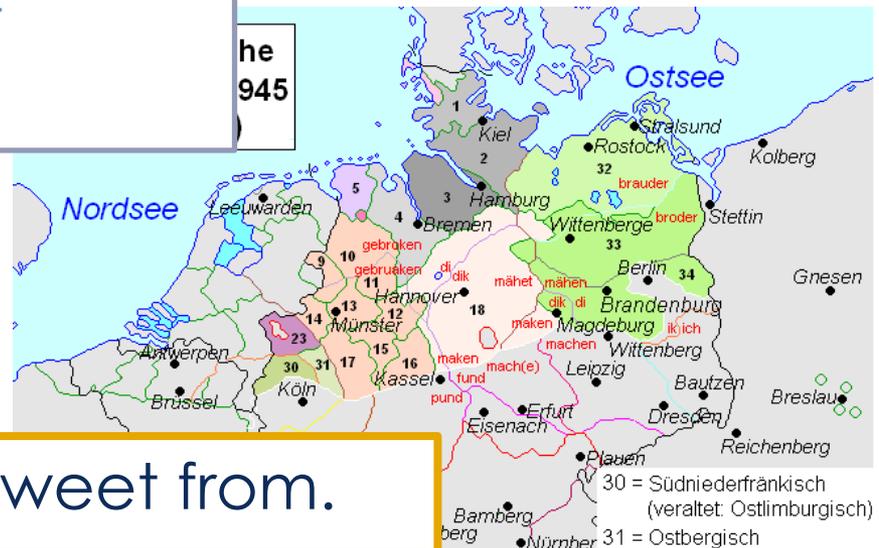
 Alexandra Raisman and Madison Kocian

Goal


Nina
@OddNina
⚙️
+ Follow

Wenn der alte Toaster von Mama, aus dem du schon als lütte Deern den Toast mampftes, langsam aufgibt :(

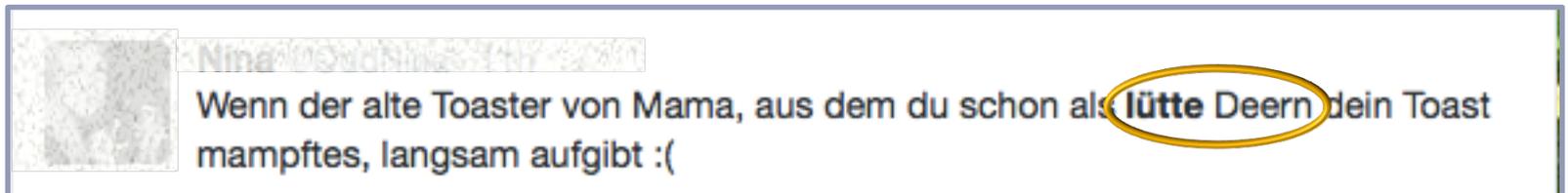
↩ Reply
↻ Retweet
★ Favorite
⋮ More



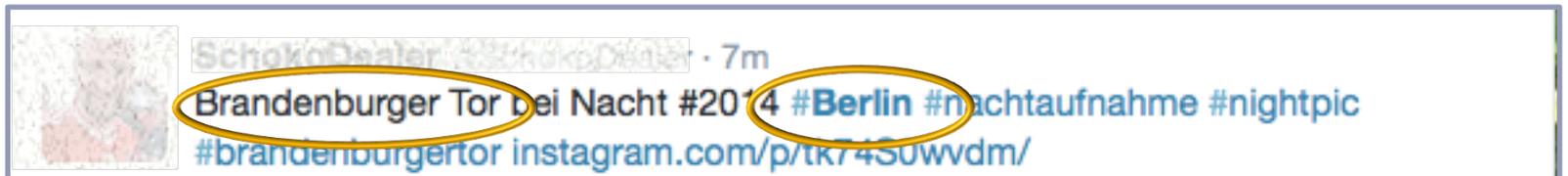
Find out where users tweet from.

Regional influences on tweets

1. Dialect origin



2. Current location





Simone Biles @Simone_Biles · Aug 20

beach day 

 Alexandra Raisman and Madison Kocian

Assumptions for identifying tweet location

- ▣ dialectal region and current location are independent
- ▣ But:
 - ▣ regionally diverging tweets should be relatively rare
 - ▣ probabilistic model of regional salience allows for assignment of tweets to several regions equally
- ▣ use only a tweet's text as features
 - ▣ ignore user profile and other metadata (for now)
- ▣ majority of tweets are “non-regional”

Related work

- Sociolinguistic Twitter studies on dialects (Eisenstein et al., 2012; Grieve, 2014)
- Localizing Twitter users by aggregating all of their tweets (Cheng et al., 2010; Hecht et al., 2011)
- Deriving location specific words from geotagged tweets (Leetaru et al., 2013)
- Thesaurus based geolocation, using known dialectal words (Scheffler et al., 2014)
- Classifying tweets into dialect regions (Scheffler et al., 2014)

Geographic Mapping of Tweets

Joint work with Johannes Gontrum

Approach

- Reconstructing the point of origin for German tweets.
 - Text based
 - Language independent
 - Using simple statistic properties
 - Exact coordinates

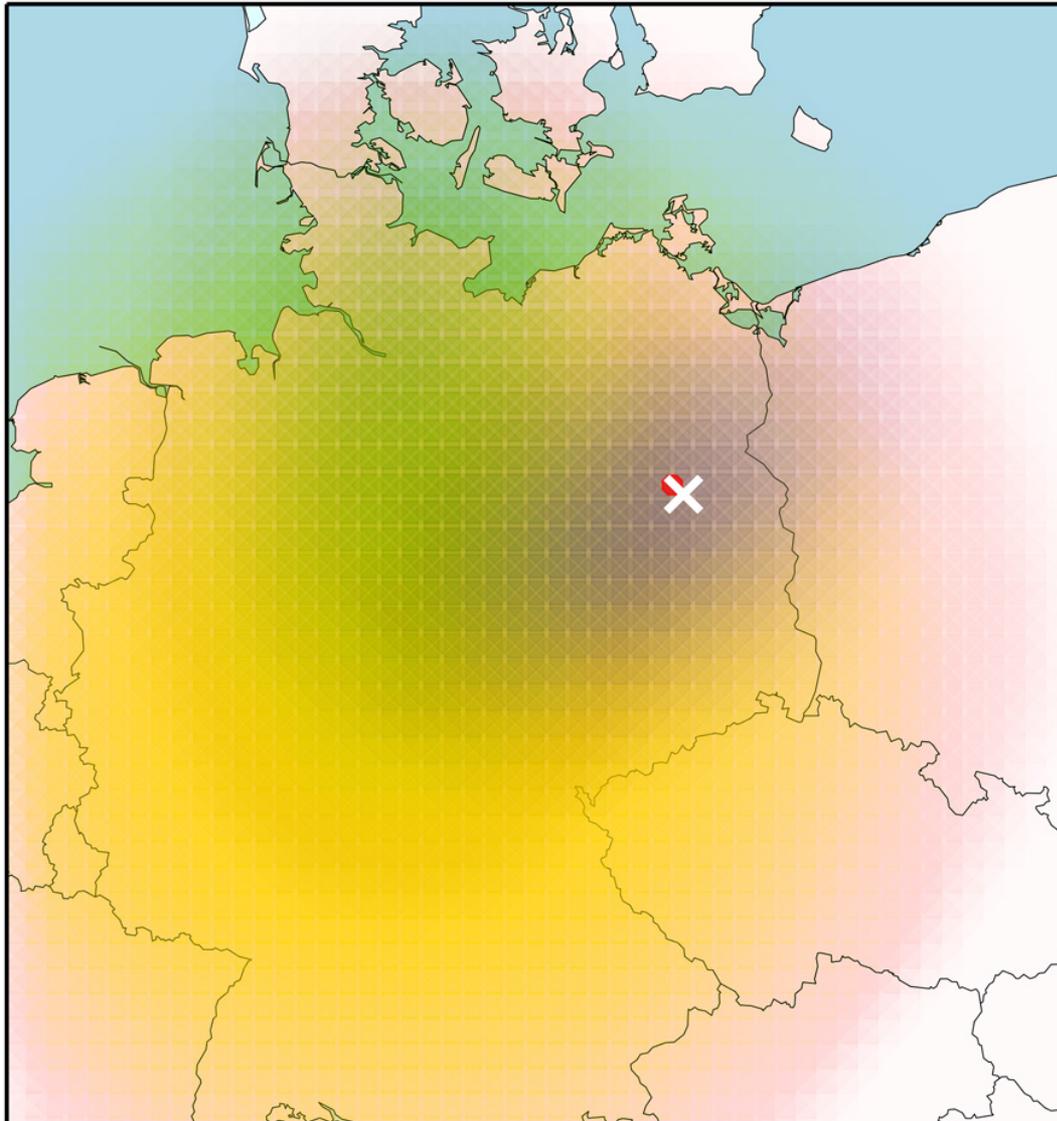
- Using the best 40% of all tokens, tweets can be reliably classified.
 - Median: 7km
 - Mean: 93km

Data

- 65 mio. German tweets
- Filtered by common German words
- Time span: February - May 2015
- Filtering corpus:
 - Removing tweets without geo tag
 - Removing tweets not from Germany, Austria, Switzerland
 - Removing generated tweets from bots
 - Removing numbers, URLs, mentions etc.
- Remaining: 360k tweets (0.55%)

Idea

- Some words are used at certain locations more often than others
- Derive a probability distribution for each word
- High variance vs. low variance
- Common words vs. highly informative words



green: hhwahl
blue: berlin
red: nordbahnhof
yellow: rest

x = true location

“balken gucken und so **hhwahl** pa **nordbahnhof** in **berlin**”

Reconstructing the location

- The original position of a tweet is dependent on the variance and median location of each token.
 1. Calculate variance and median location for each token.
 2. Remove common/non-regional words by variance threshold
 3. Compute the weighted midpoint for all other words in the tweet.
 4. Weight: Inverse variance.

Location of a tweet

t = tweet with tokens $t_0 \dots t_n$

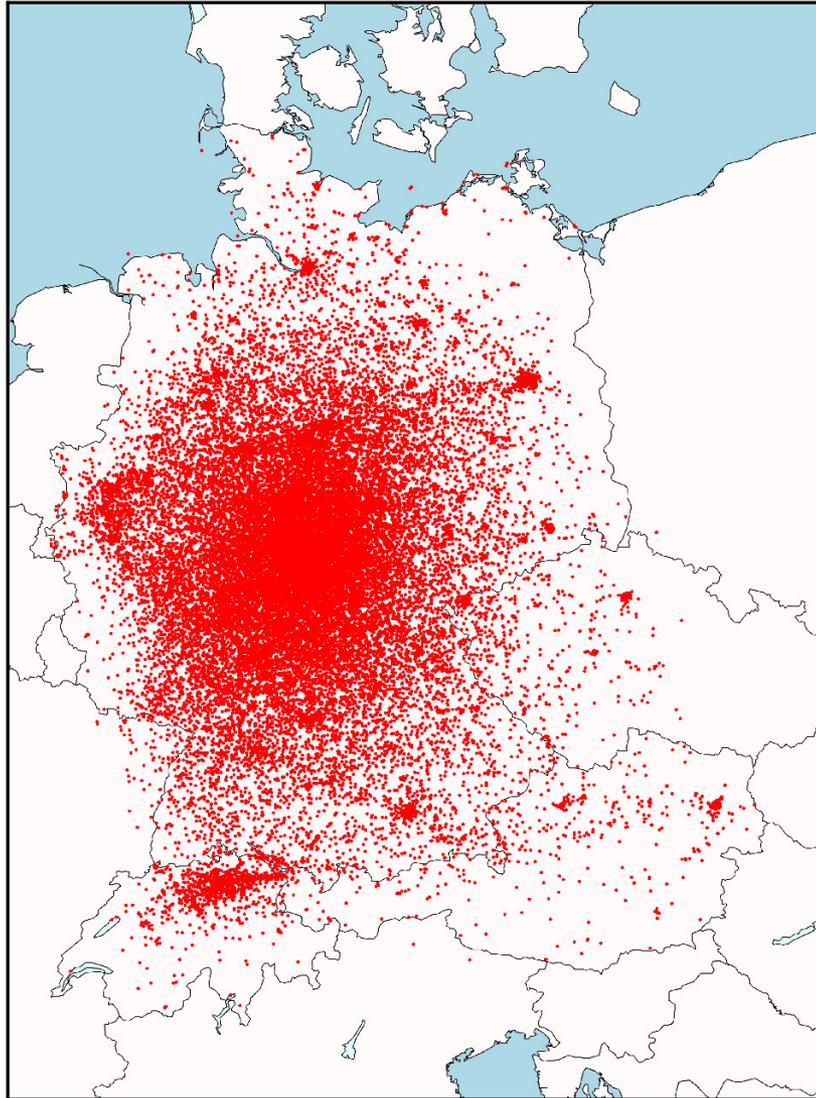
variance $\sigma_0 \dots \sigma_n$

median location $m_0 \dots m_n$

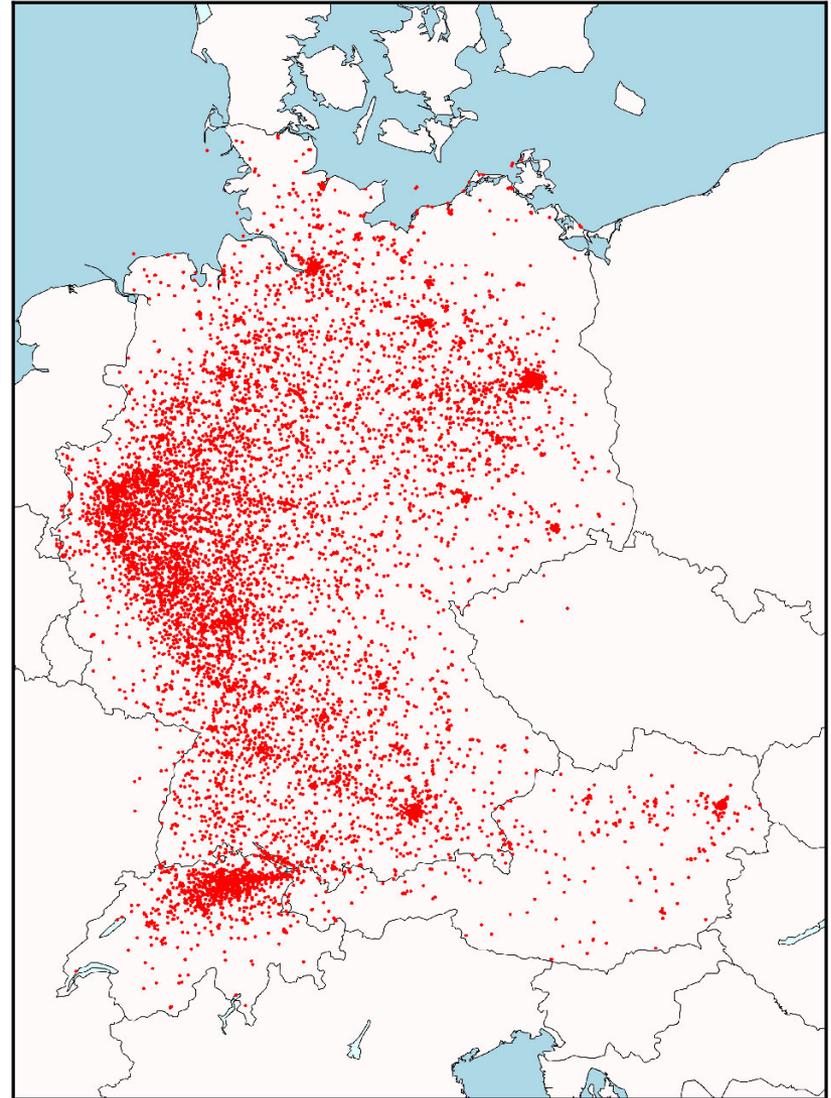
$$\text{Loc}(t) = \frac{\sum_{i=0}^n \vec{\sigma}_i^{-1} * \vec{m}_i}{\sum_{i=0}^n \vec{\sigma}_i^{-1}}$$

“balken gucken und so **hhwahl** pa **nordbahnhof** in **berlin**”

Midpoints of tokens (training corpus)



all tokens (100%)

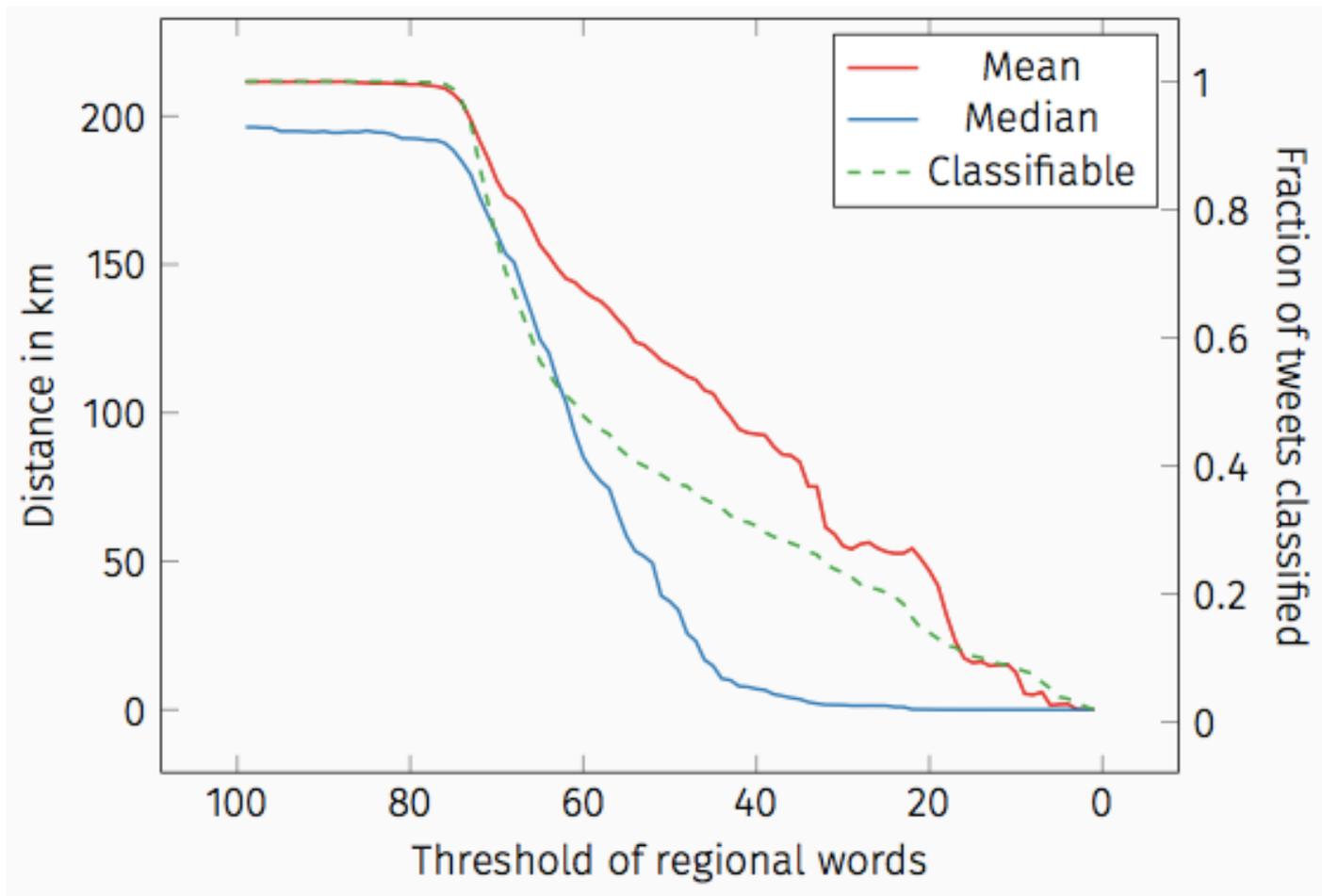


25% of tokens

Importance of filtering

- Not every token contains information about locality.
- Classifying a tweet with only common words will fail.
- False confidence in the result.
- Sort all tokens by variance, keep only lowest X%.
- Assumption: X% of words are regionally salient.

Regional threshold



Results

Threshold	Mean	Median	#Tweets
100	212km	196km	1000
75	207km	188km	988
50	116km	36km	377
40	93km	7km	306
30	55km	1.56km	233
20	47km	0.06km	139
10	12km	0.00km	84

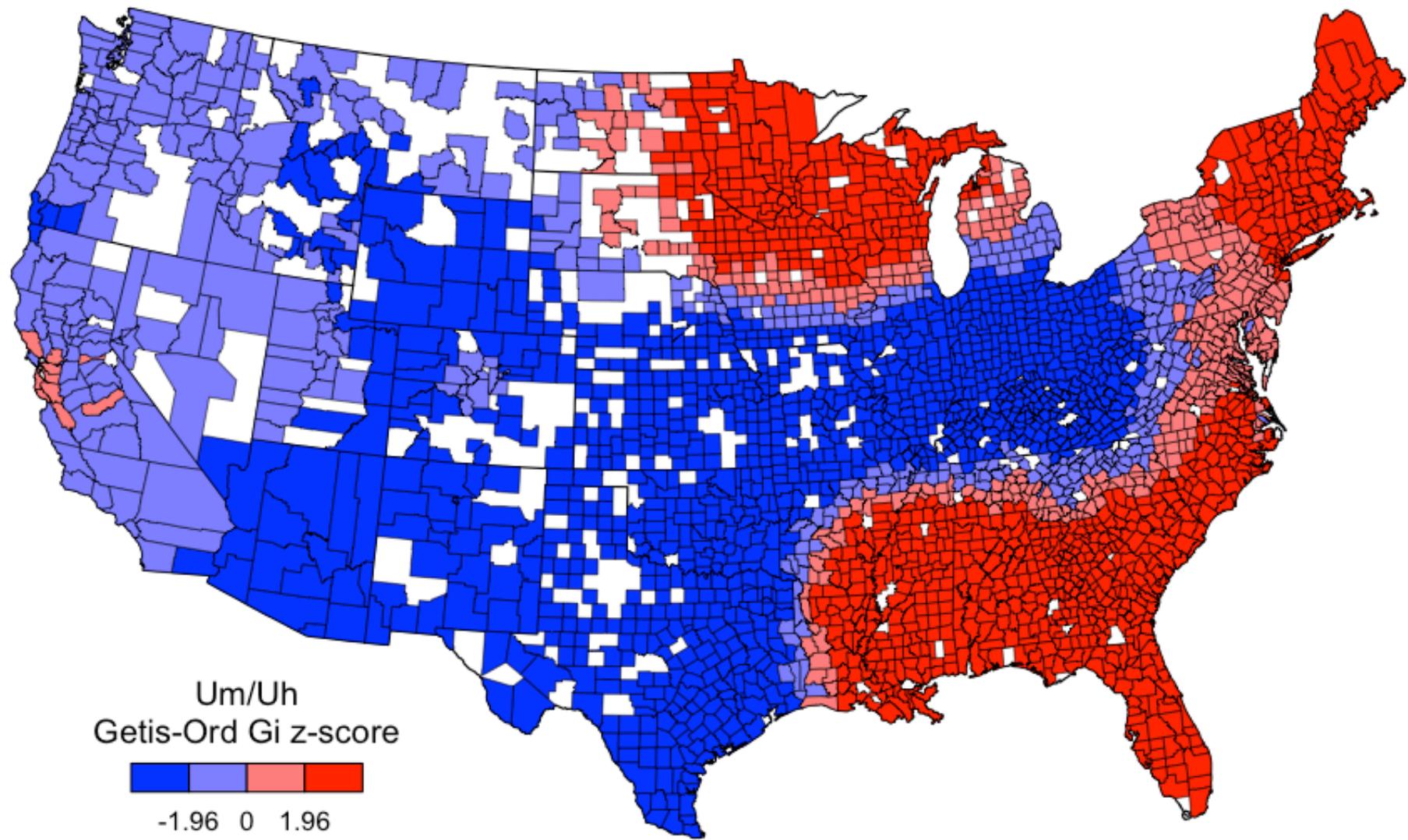
Regionally salient tokens

Berlin	Zurich	Essen
kadewe	tagi	rheinische
kudamm	uf	hattingen
alexanderplatz	het	herne
friedrichshain	isch	westfalen
brandenburg	scho	ddorf
fernsehturm	au	ruhr
dit	zuerichsee	thyssenkrupp
morjen	gseh	duisburg

Conclusion: geocoding of tweets

- Classify tweets only based on their text
- Geographic probability distribution for tokens
- Filtering and weighting tokens by variance
- Removal of wide-spread words
- Accurate reconstruction of tweets' location

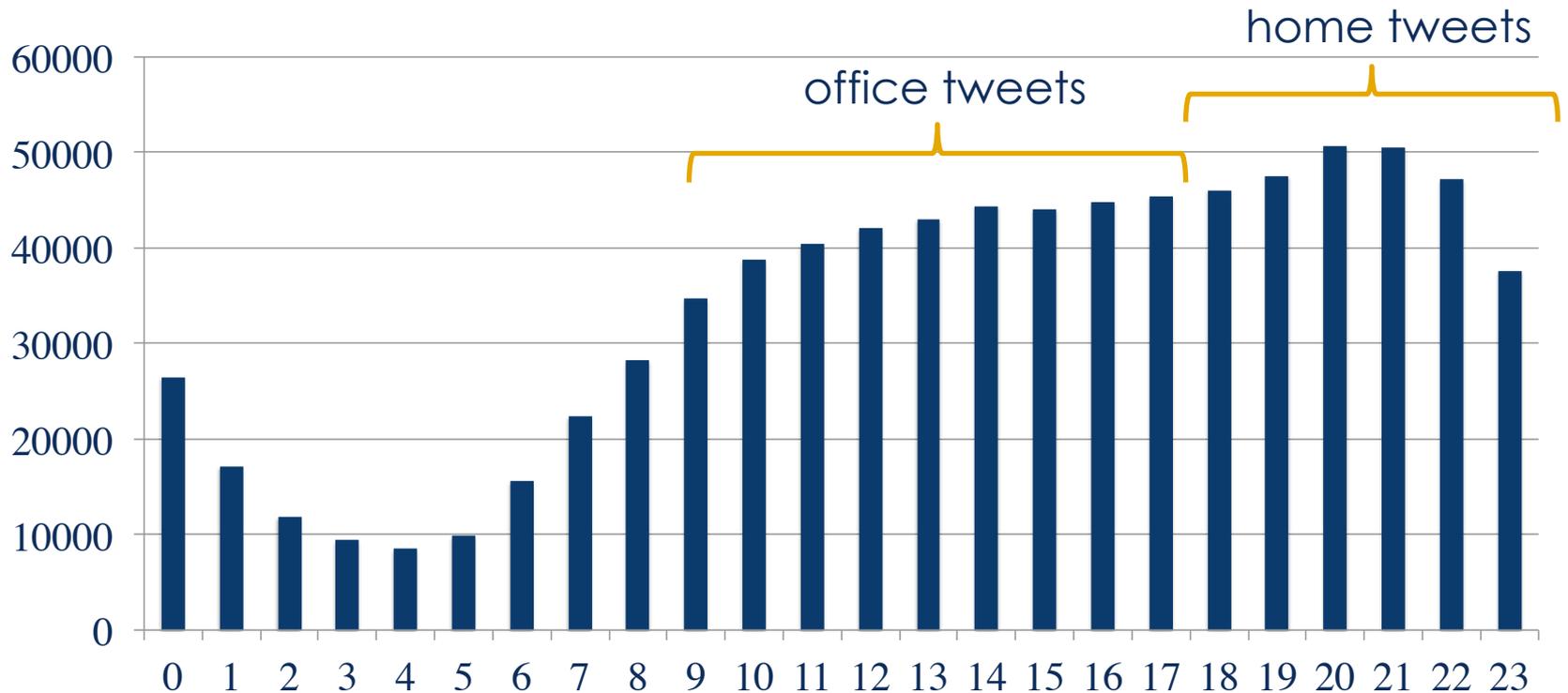
But: How do geocoded and non-geocoded tweets differ?

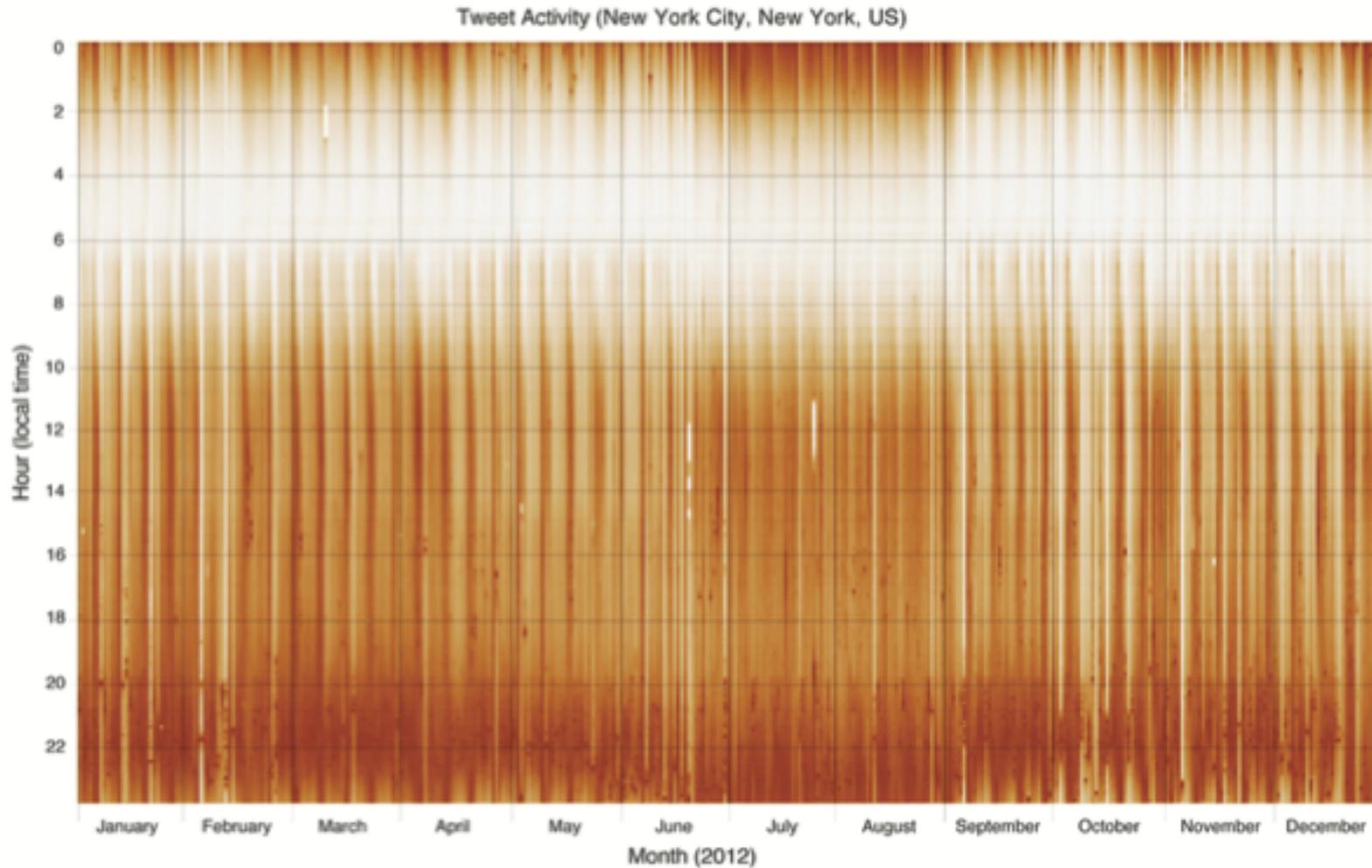


um / uh – Jack Grieve, <https://sites.google.com/site/jackgrieveaston/treesandtweets> , August 18, 2014

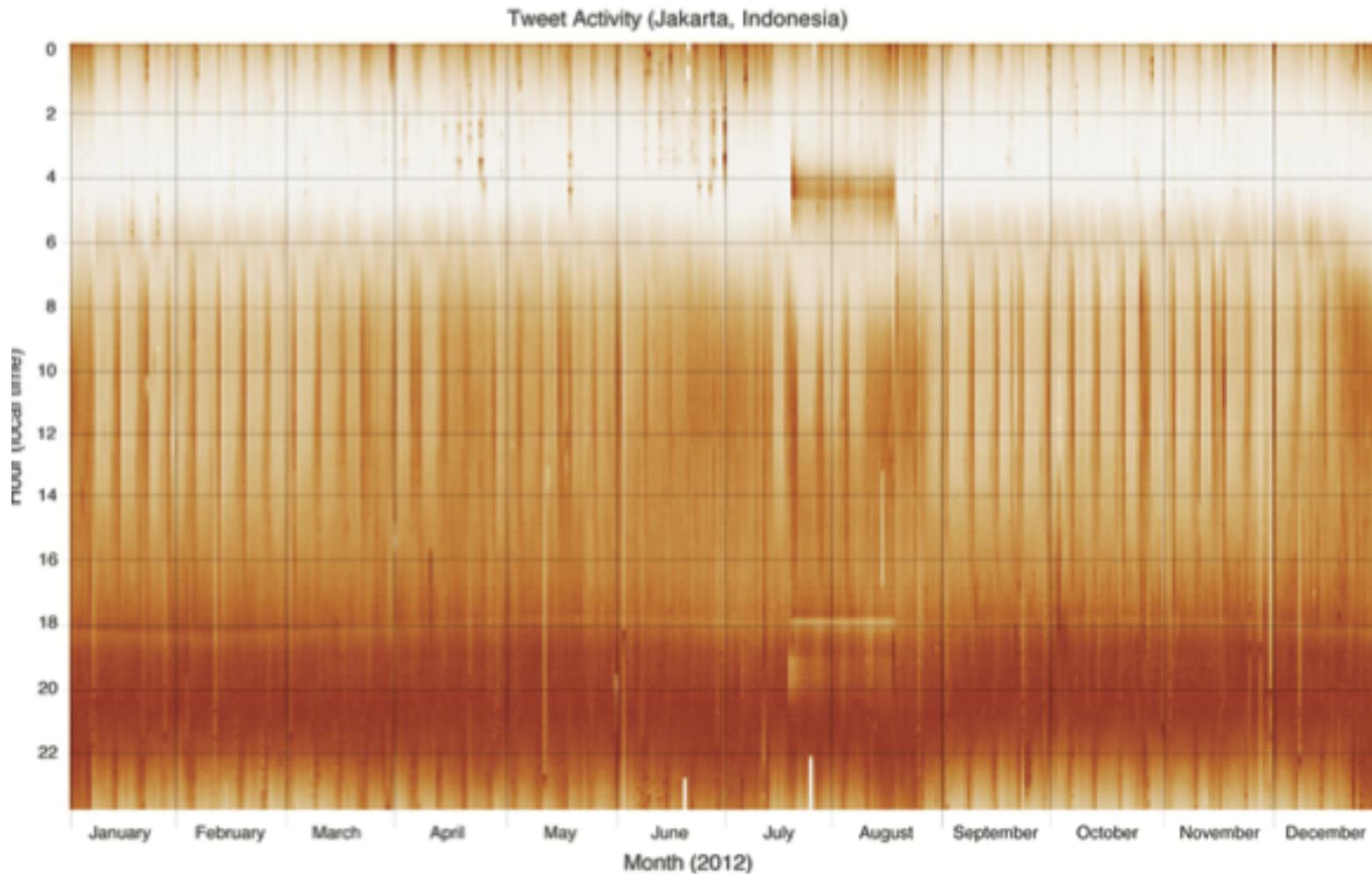
2. time stamps

When do users tweet?





Rios/Lin, 2013: Visualizing the “Pulse” of World Cities on Twitter. In: Proceedings of the Seventh International AAAI Conference on Weblogs and Social Media.

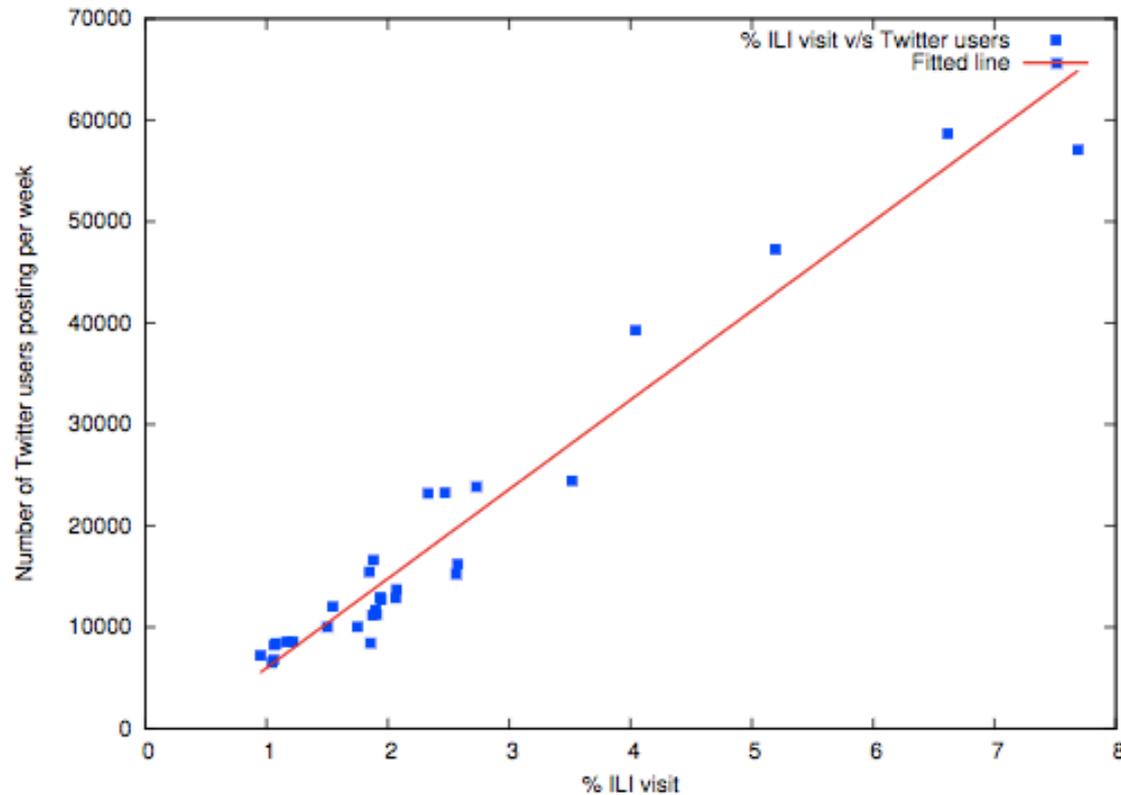


Rios/Lin, 2013: Visualizing the “Pulse” of World Cities on Twitter. In: Proceedings of the Seventh International AAAI Conference on Weblogs and Social Media.

Twitter data as a social sensor

- earthquake observations (Sakaki et al., 2010)
- flu trend prediction (Achrekar et al., 2011)
- effect of weather on mood (Hannak et al., 2012)
- event detection (Ritter et al., 2012)

Flu trends (Achrekar et al., 2011)



Twitter + Circadian Rhythm

Joint work with Christopher Kyba, GFZ Potsdam

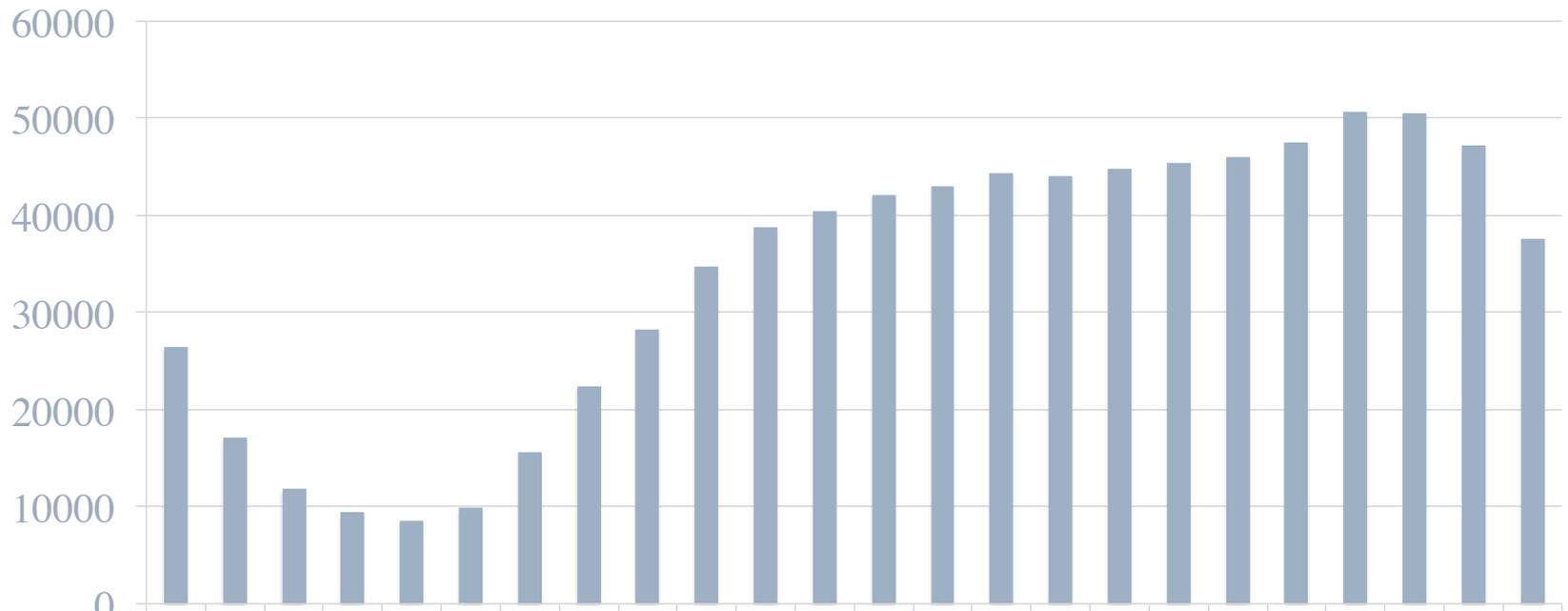


Image: Fotolia

Circadian rhythm

- Human schedules are determined by biological and social constraints
- Kantermann et al. (2007) showed that wake timing on free days tracks sunrise during standard time, but not during daylight savings time
- Sleep research is done in lab-based sleep studies or through surveys
- Roenneberg (2013) argues for much larger sample sizes of real-world data to improve our understanding of human sleep-wake patterns

Goal

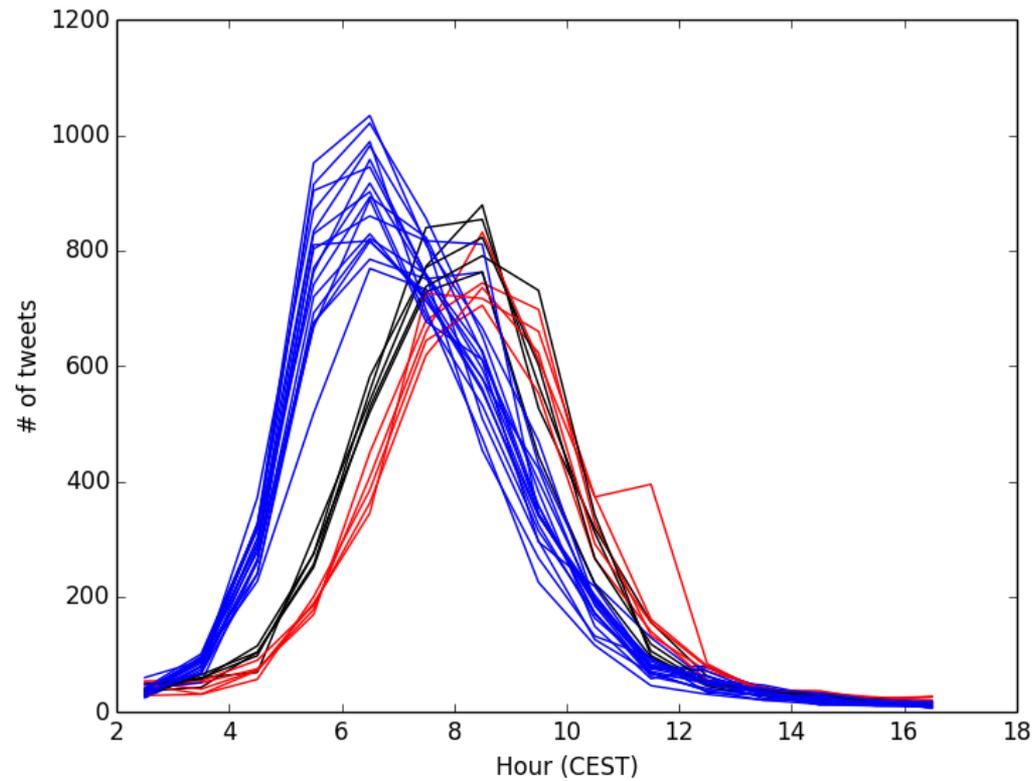


Use a Twitter time series to measure the sleep-wake pattern in humans and its interaction with DST.

Dataset

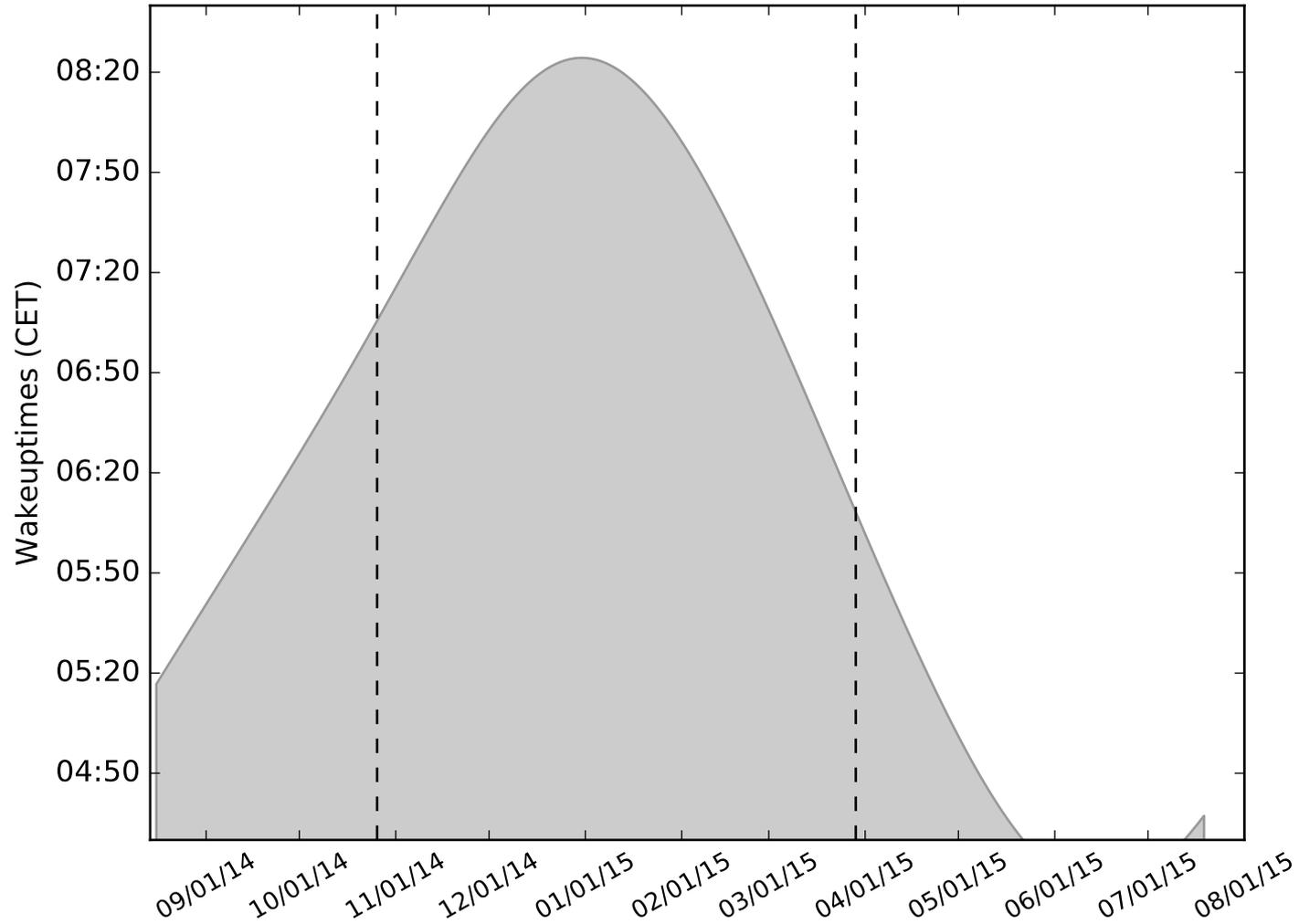
- one year of German tweets containing the phrase “guten morgen”
- August 15, 2014 to August 14, 2015
- dataset: 1,443,004 unique tweets from 206,633 individual users (retweets excluded)
- Tweets were binned by 15-minute windows

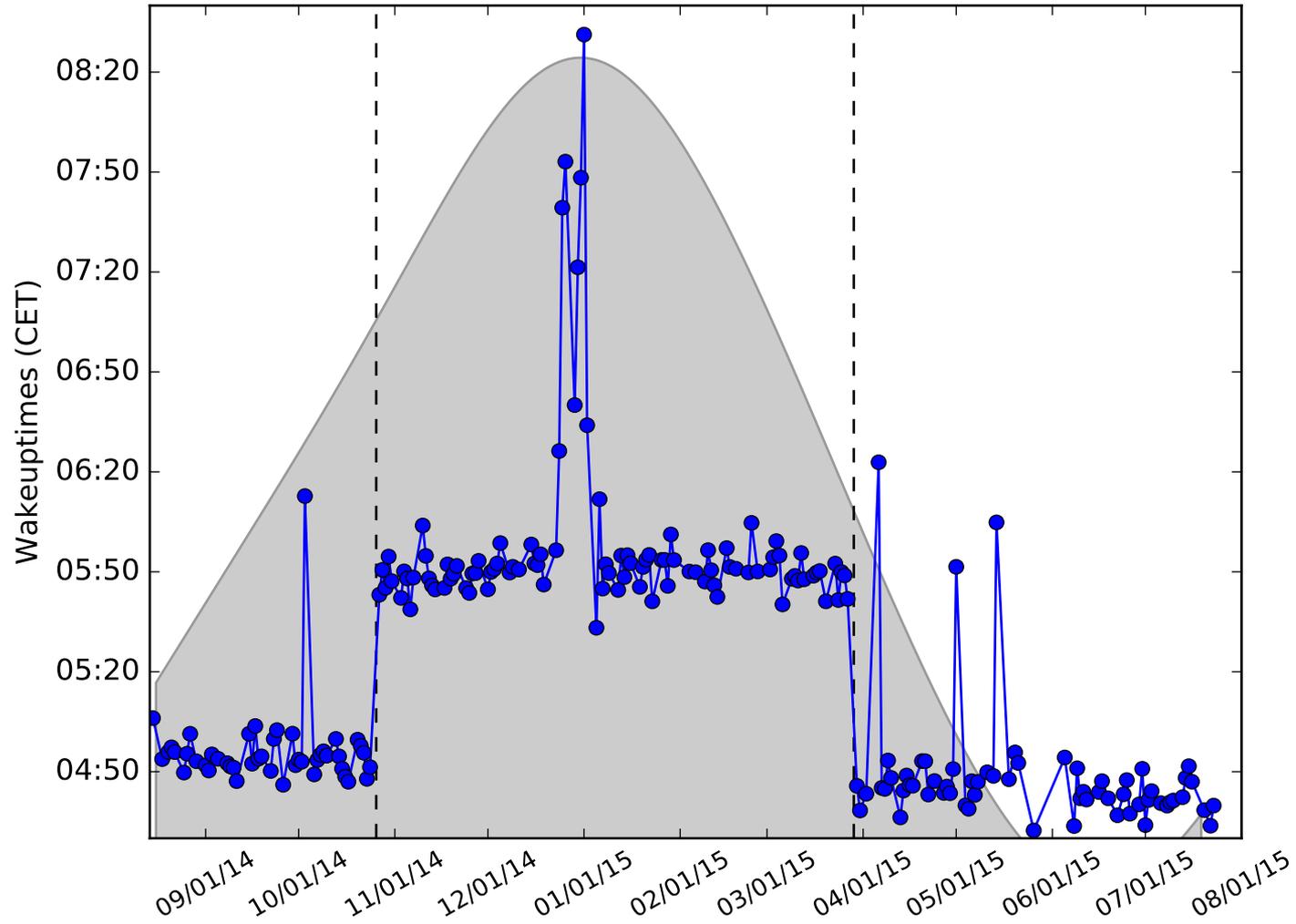
'Good morning'-tweets

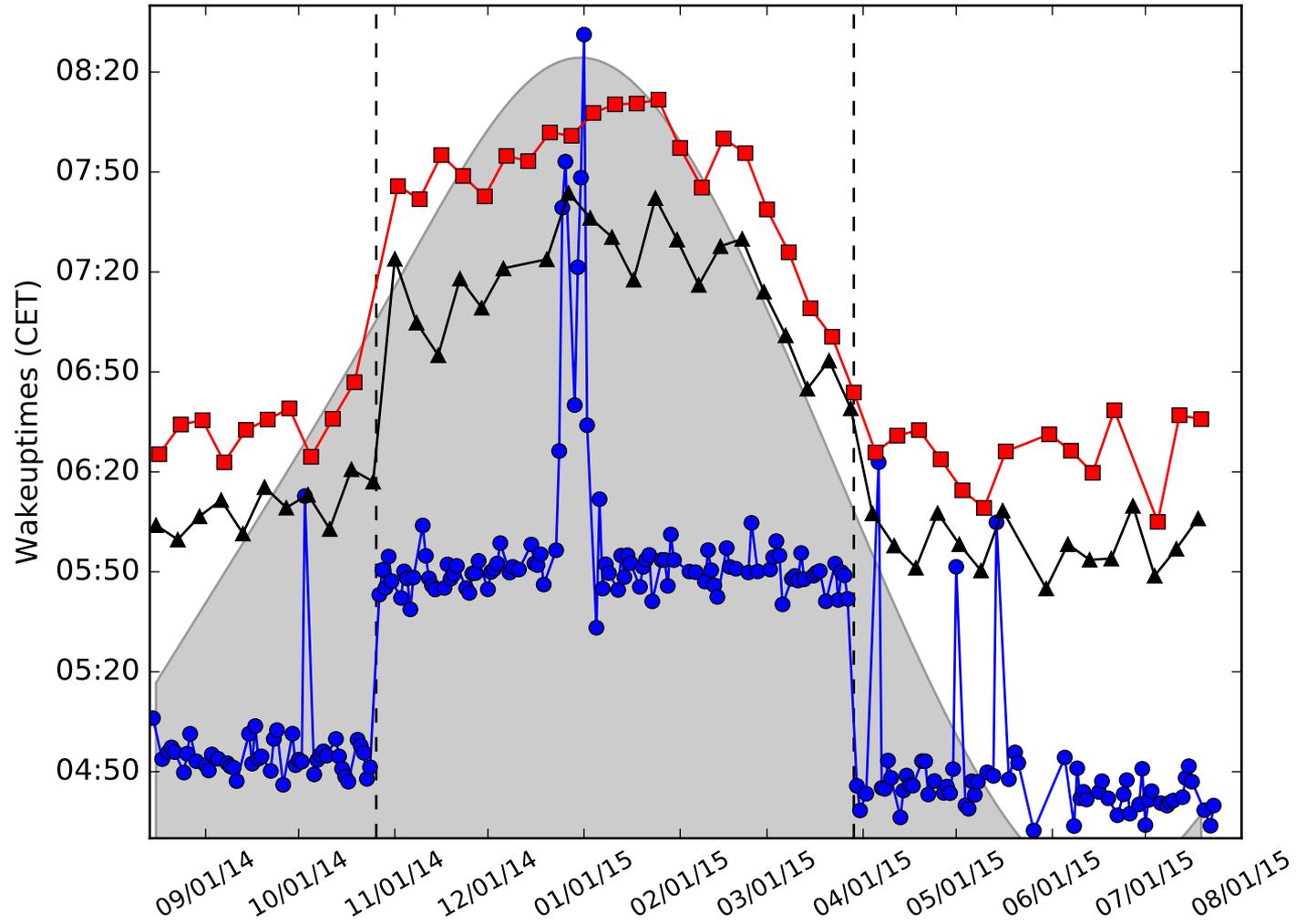


Analysis

- “Onset of Twitter activity” = time at which the rate of ‘good morning’-tweets reached half of the maximum
- The relation of this time to the sunrise time and social time was studied







General patterns

- morning greetings realistically reflect the onset of activity times (OTAT), between 4:30 and 8:30 local time on weekdays
- OTAT (determined geometrically) is related to the local sunrise and social time
- OTAT is much earlier on weekdays than on free days (public holidays resemble Saturdays)
- during the winter and standard time, OTAT on free days tracks dawn
- This relationship ends with the start of DST, and as a result the difference between wake times on free and work days grows considerably.

Interaction with social norms

We computed the difference in OTAT between work and free days for 42 weeks during the study period:

- Saturday: 79 ± 14 min
- Sunday: 103 ± 25 min
- weekday/weekend difference is largest in January (Saturday: ~ 99 mins, Sunday: ~ 140 mins)
- is smallest in Spring just up to the introduction of DST

Discussion

- Close tracking of dawn by OTAT on free days in Winter is disrupted by the introduction of DST.
- Large difference between Saturdays and Sundays
- Results are consistent with sleep survey data by Kantermann et al. (2007)
- Use of social media data in studying sleep patterns and other social norms?
- Implications for future policy change wrt. DST – could use social media to track effects?

3. Recovering user profile info

Demographic user data

- ▣ user profile:
 - ▣ language
 - ▣ location
 - ▣ name, description, profile image, ...

- ▣ demographic information:
 - ▣ gender
 - ▣ age
 - ▣ location
 - ▣ income/profession
 - ▣ personality, languages, interests, ...

What your language says about your occupation

- Hu et al., 2016
- Match data from different social media: Twitter + LinkedIn
- Twitter data used to represent styles, interests and personalities
- correlate linguistic info from Twitter and occupation info from LinkedIn

Data

Linguistic data:

- 3000 most recent tweets
- Twokenizer; remove too rare and too frequent n-grams

Occupation data:

- LDA topic model based on LinkedIn skill endorsements
- matrix of skill-job cluster relations
- matrix of person-job cluster relations

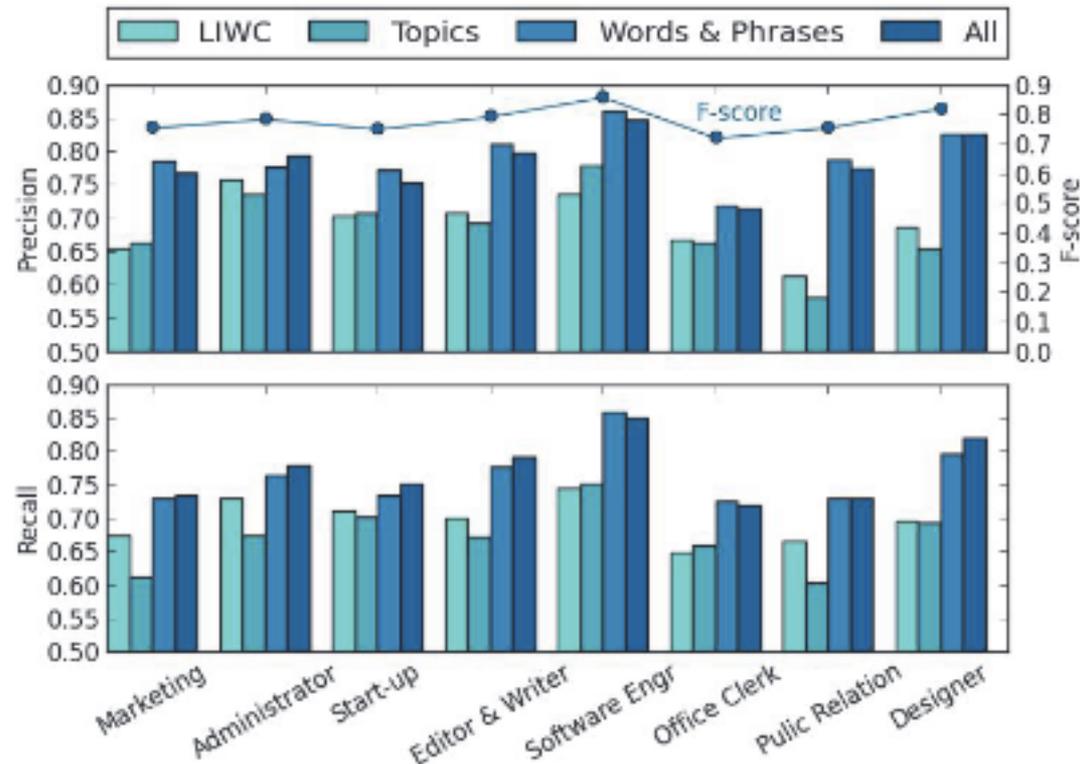
Jobs and skills

■ 8 largest job clusters

Job Name	First Skill	Second Skill	Third Skill	Fourth Skill	Fifth Skill
1. Marketing	Digital MKTG	Social Media MKTG	Online MKTG	Digital Strategy	Advertising
2. Administrator	Public Speaking	Leadership	Fundraising	Event Planning	Coaching
3. Start-up	Start-ups	Entrepreneurship	Strategy	Business DEV	Management
4. Editer&Writer	Blogging	Editing	Journalism	Copy Editing	Storytelling
5. Software Engr	MySQL	CSS	JavaScript	PHP	jQuery
6. Office Clerk	Microsoft Office	Microsoft Excel	PowerPoint	Microsoft Word	Customer Service
7. Public Relation	Public Relations	Media Relations	Press Releases	Strategic COMM	Corporate COMM
8. Designer	Graphic Design	Web Design	Photography	Illustrator	Photoshop

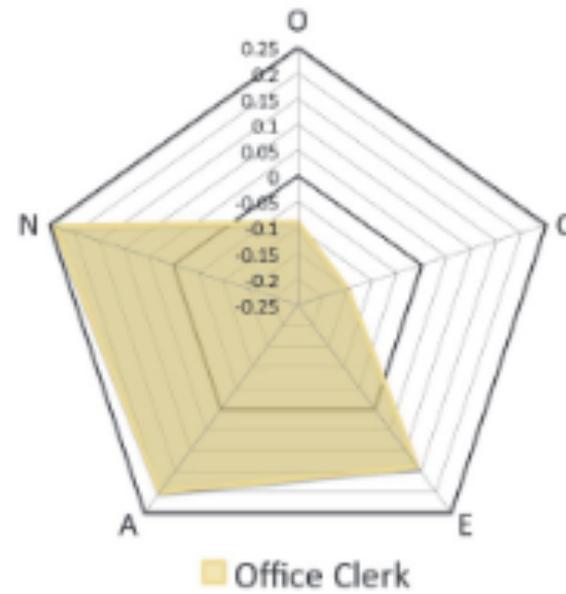
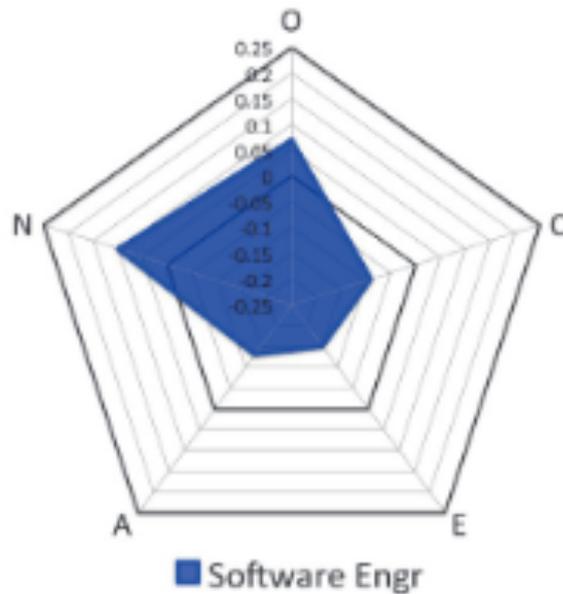
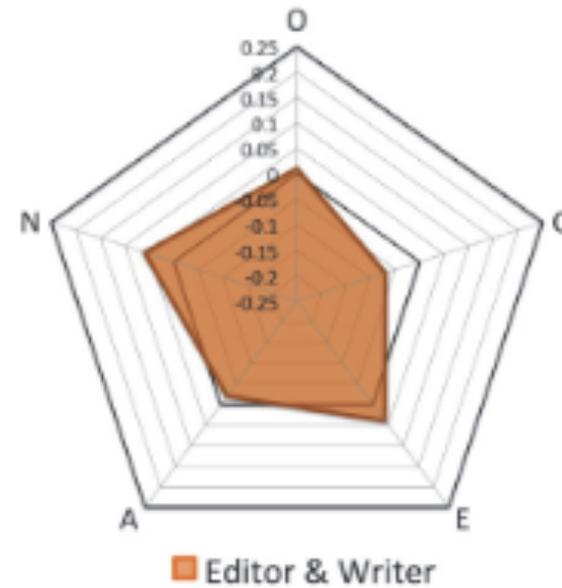
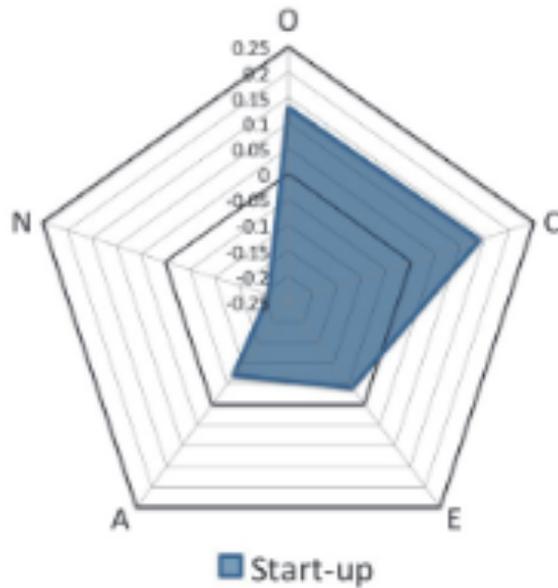
Predicting a job

- ▣ 5-fold cross validation
- ▣ assign each user to the job whose weight is > 0.8
- ▣ overall F-score of all 8 jobs is 0.78
- ▣ tweet n-grams are very predictive for these occupations



Personality

- BM Watson Personality Insights service API:
www.ibm.com/smarterplanet/us/en/ibmwatson/
- Personality traits based on a user's tweets
- Big Five:
 - Openness (to new things)
 - Conscientiousness (self-discipline)
 - Extraversion (social interactions)
 - Agreeableness (cooperativity, willing to compromise)
 - Neuroticism (instability of emotions)



Conclusions: Hu et al., 2016

- Combine data from different social media platforms
- Soft-clustering occupations to avoid uncertainty in self-reporting information
- Extract linguistic style information and Big Five personality traits from tweet text
- Pearson Correlation Coefficients used to explore differences in these properties across jobs

Summary Day 2

- Using metadata together with linguistic information
- Using linguistic information to recover metadata
- Social media data as a sensor for human behavior / real-world events

Coming up...

Wednesday: Sentiment (classification and clustering)

Thursday: Conversations and Discourse

Friday:

.... stay tuned!

Thank you.

tatjana.scheffler@uni-potsdam.de

Twitter Bots: to prepare...

1. Get a Twitter account!
2. Create a new Twitter application and receive consumer key and secret.
3. After the step above, you will be redirected to your app's page. Create an access token in the "Your access token" section.

If you want to get started:

Quick and easy Twitter bots: [Make your own @HydrateBot](#)

Python corpus based Twitter bots: [Creative Twitter bots](#)

Eisenstein J, B. O'Connor, N. Smith, and E.P. Xing. 2010. A latent variable model for geographic lexical variation. In Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, pages 1277–1287.

Gontrum J and Scheffler T. Text-based Geolocation of German Tweets. In: Proceedings of the NLP4CMC 2015 Workshop at GSCCL, Duisburg-Essen, Germany. 2015.

HU, T.; XIAO, H.; LUO, J.; NGUYEN, T.. What the Language You Tweet Says About Your Occupation. International AACL Conference on Web and Social Media, North America, mar. 2016. Available at: <<http://www.aclweb.org/ocs/index.php/ICWSM/ICWSM16/paper/view/13020/12738>>.

Kantermann T, Juda M, Merrow M, Roenneberg T (2007) The human circadian clock's seasonal adjustment is disrupted by daylight saving time. *Current Biology* 17: 1996–2000.

Pavalanathan U. and J. Eisenstein. 2015. Confounds and Consequences in Geotagged Twitter Data.

Rios, M and Lin, J (2013) Visualizing the “pulse” of world cities on Twitter. In Proceedings of ICWSM, Boston, USA: AACL.

Roenneberg, T (2013) Chronobiology: the human sleep project. *Nature* 498(7455):427–428.

Scheffler T (2014) A German Twitter snapshot, In Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC 2014), Reykjavik, Iceland: European Language Resources Association (ELRA).

Scheffler T, Kyba CCM (2016) Measuring Social Jetlag in Twitter Data, In Proceedings of the Tenth International AACL Conference on Web and Social Media (ICWSM 2016), Cologne, Germany: AACL. <http://www.aclweb.org/ocs/index.php/ICWSM/ICWSM16/paper/view/13080>

... and references therein