

# Identifying and Mapping the Spread of Emerging Words

Jack Grieve

[j.grieve1@aston.ac.uk](mailto:j.grieve1@aston.ac.uk)

[@JWGrieve](#)

Centre for Forensic Linguistics  
Aston University

2 December 2016

BAULT 2016

University of Helsinki

# Emerging Words

Very large corpora of social media provide us with new opportunities to track newly emerging words as they enter into the general lexicon for the first time.

This study is based on a  $\sim 9$  billion word corpus of geocoded American Tweets ( $\sim 1$  billion Tweets) collected between October 2013 and November 2014.

Emerging words are identified and their properties are considered and their usage mapped to identify common sources of lexical innovation in American English.

# Rising Words

To find newly emerging words, the 97,246 words types that occur at least 500 times in the corpus were extracted from the corpus.

The word **relative frequency per day** was then compared to **day of the year** (across 399 days) using a **Spearman's rank correlation** coefficient to assess the degree to which the usage of each word in the corpus had been rising over the 13 month period.

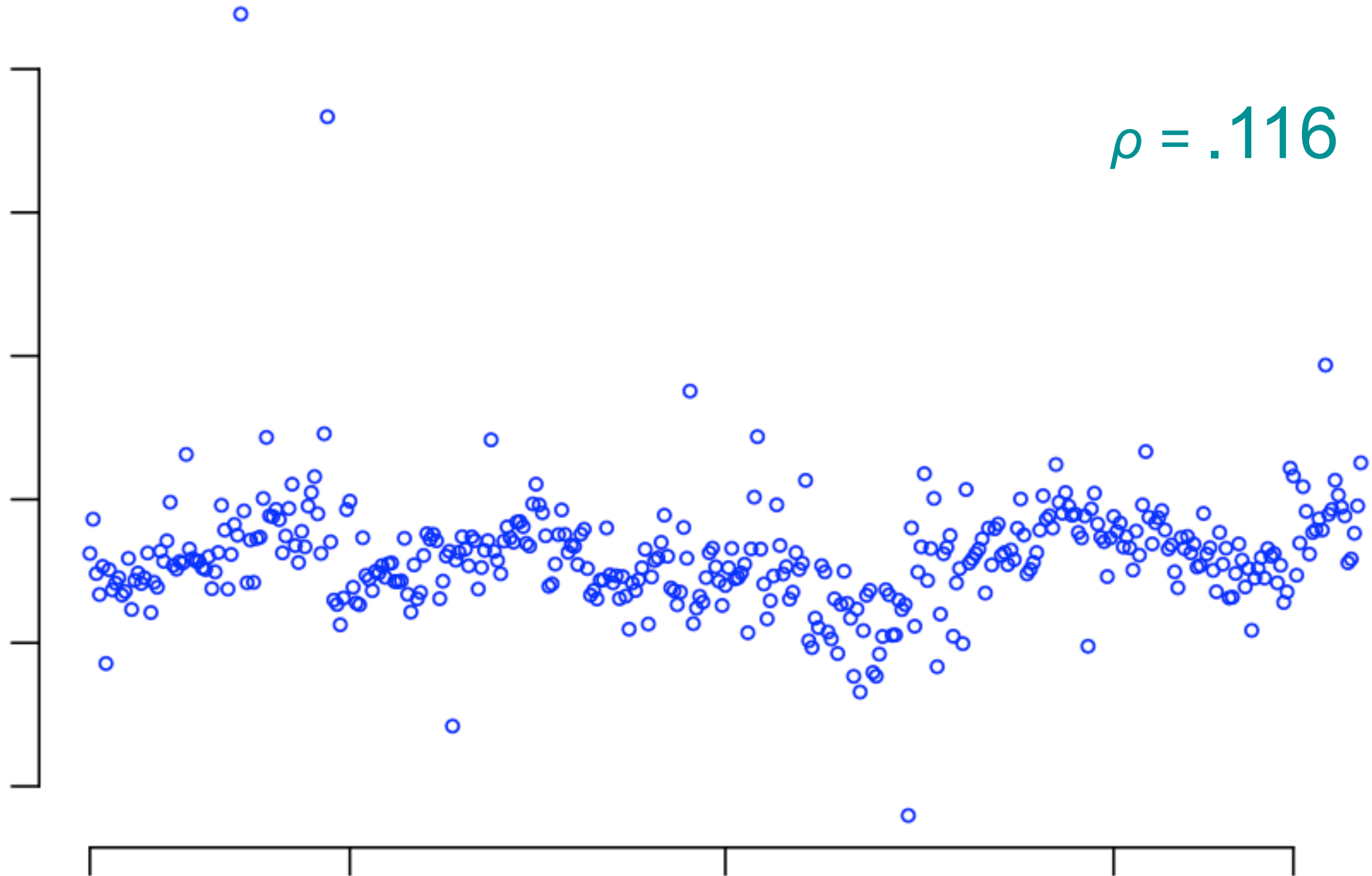
**AND**

Frequency per Billion Words

11500000 12500000 13500000

$\rho = .116$

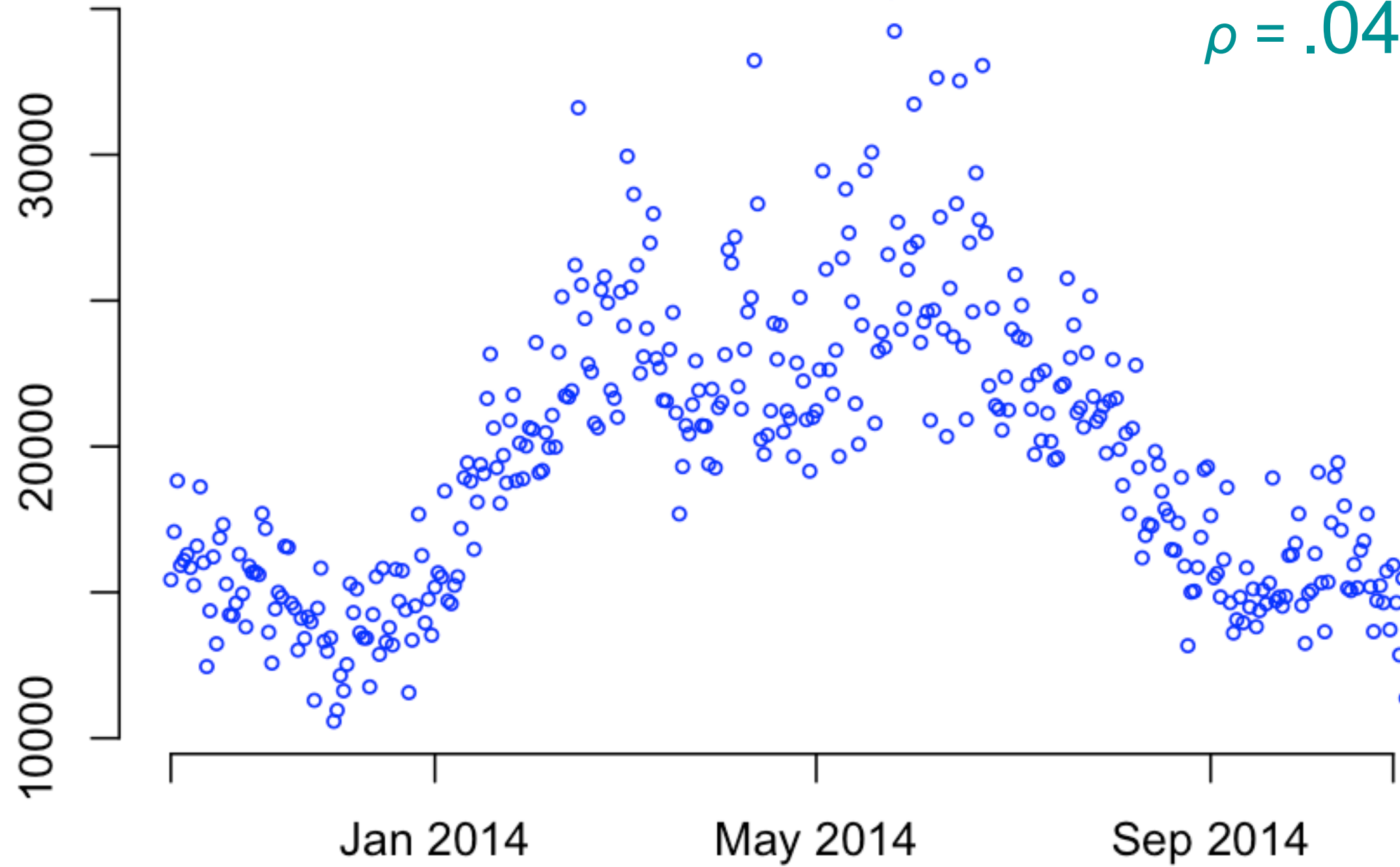
Jan 2014 May 2014 Sep 2014





# STRAWBERRY

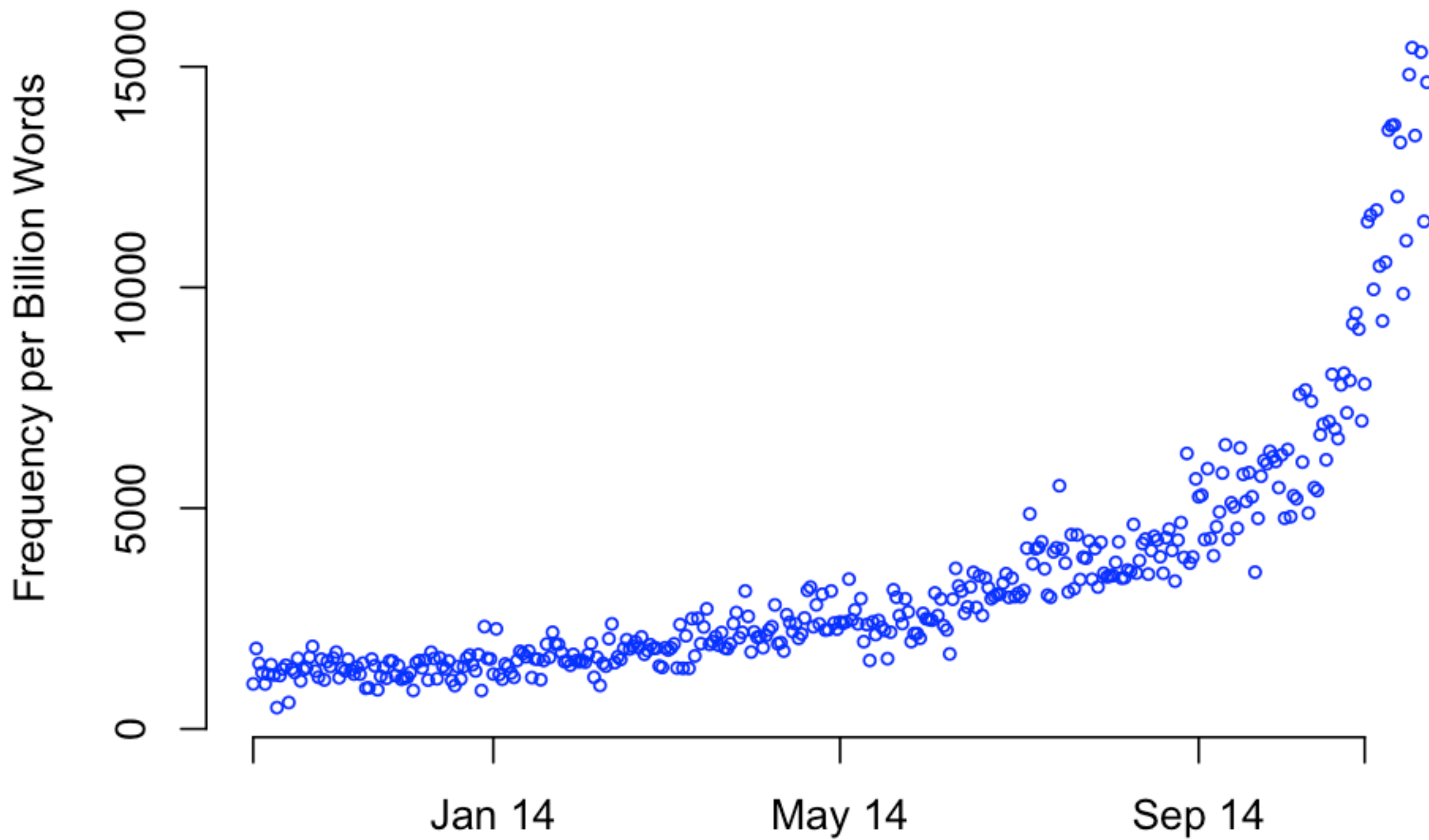
Frequency per Billion Words



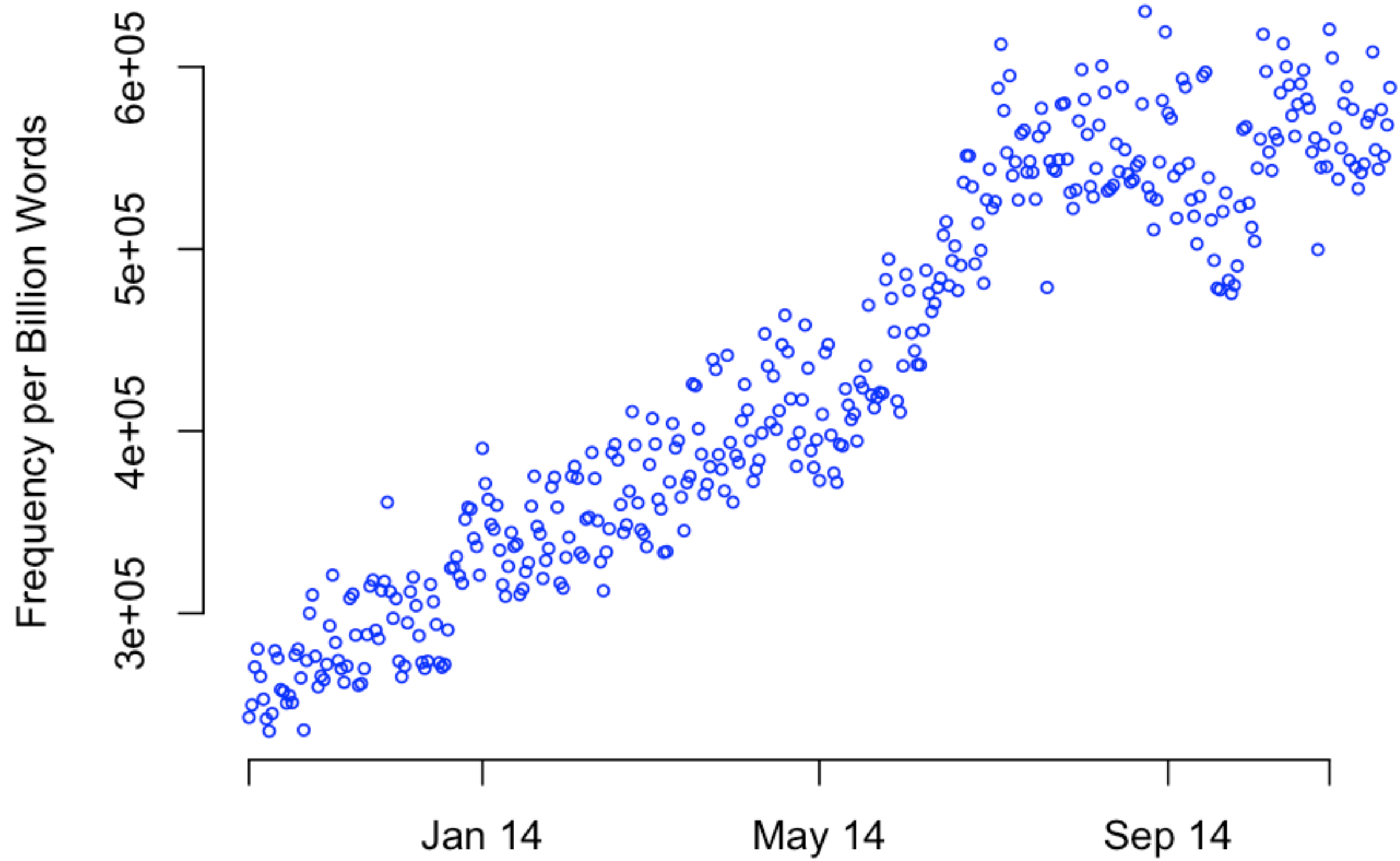
# Top 10 Rising Words on Twitter 2014

| Word              | $\rho$ | Definition                              |
|-------------------|--------|---|
| <i>fuckboy</i>    | 0.947  | Asshole, Jerk, Poser, Tool, etc.        |
| <i>rn</i>         | 0.938  | Right Now ( <b>Top Riser 2013 too</b> ) |
| <i>hbd</i>        | 0.928  | Happy Birthday                          |
| <i>fw</i>         | 0.927  | Fuck with                               |
| <i>unbothered</i> | 0.926  | Unconcerned & Disengaged                |
| <i>ft</i>         | 0.925  | Face time                               |
| <i>gmfu</i>       | 0.924  | Get me fucked up                        |
| <i>sm</i>         | 0.919  | So Much                                 |
| <i>squad</i>      | 0.919  | Group of friends                        |
| <i>asf</i>        | 0.918  | As fuck                                 |

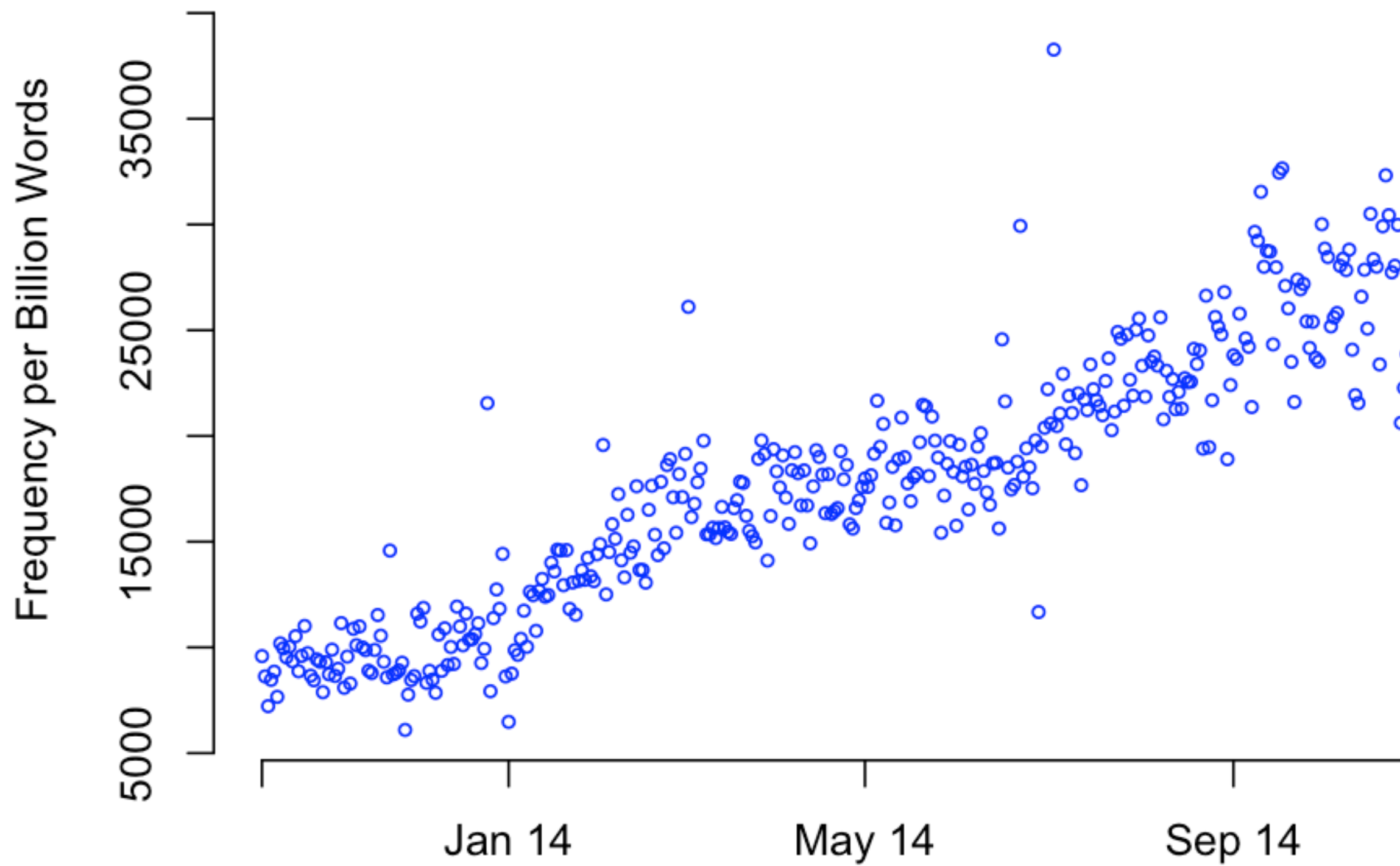
# FUCKBOY



RN



# HBD



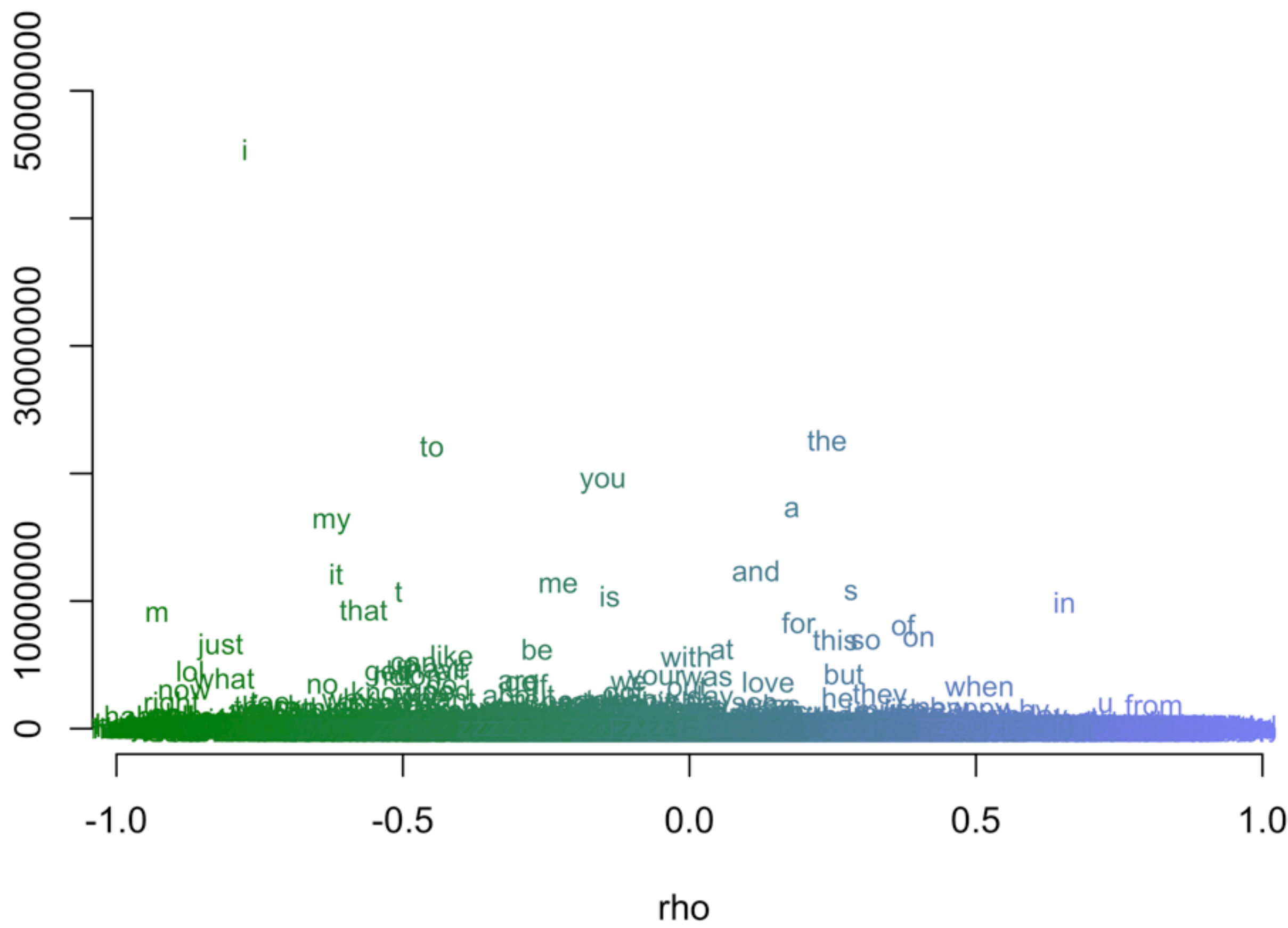
# Emerging Words

Although measuring correlations allows for rising words to be identified, most are far too common at the start of the time period to be considered emerging words.

Emerging words were identified by cross-referencing the list of rising words ( $p > .70$ ), against a list of rare words (relative frequency  $< 1$  per million words in 4th quarter of 2013), against a list of non-dictionary words.

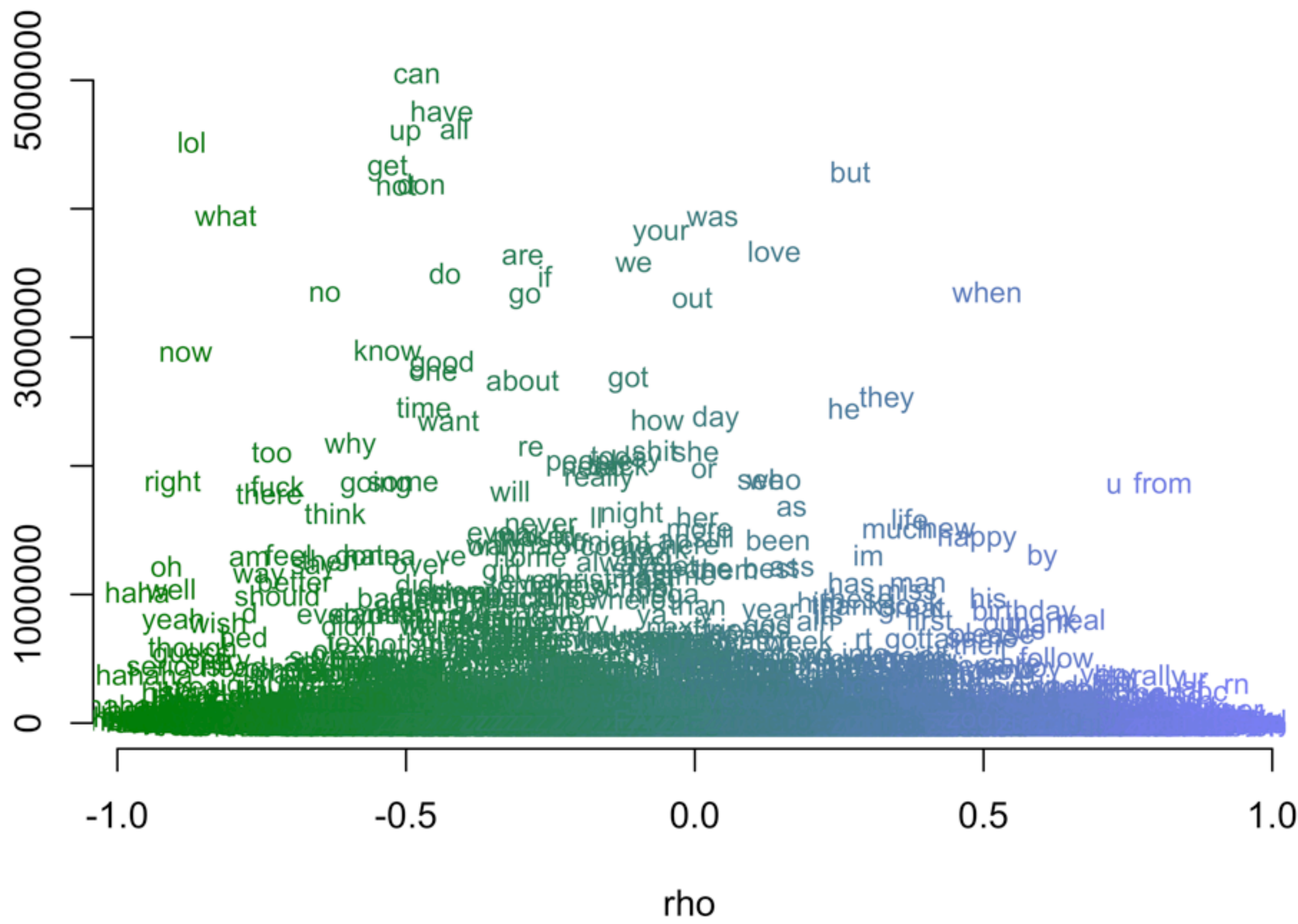
Proper nouns and medical terms associated with a rise of medical job adds over this period were also excluded.

2013 Q4 Rel Freq Per Billion Words



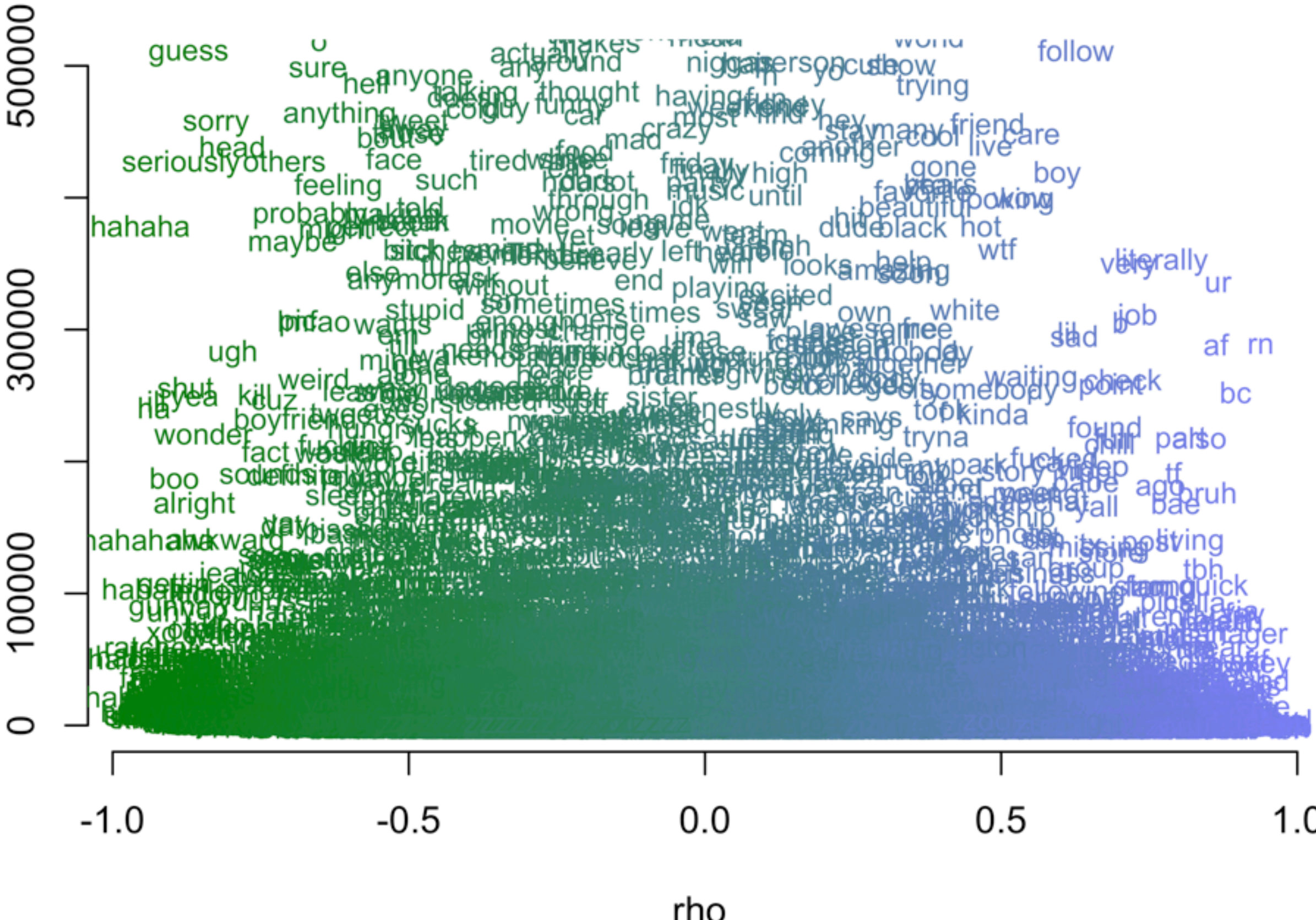


2013 Q4 Rel Freq Per Billion Words





2013 Q4 Rel Freq Per Billion Words





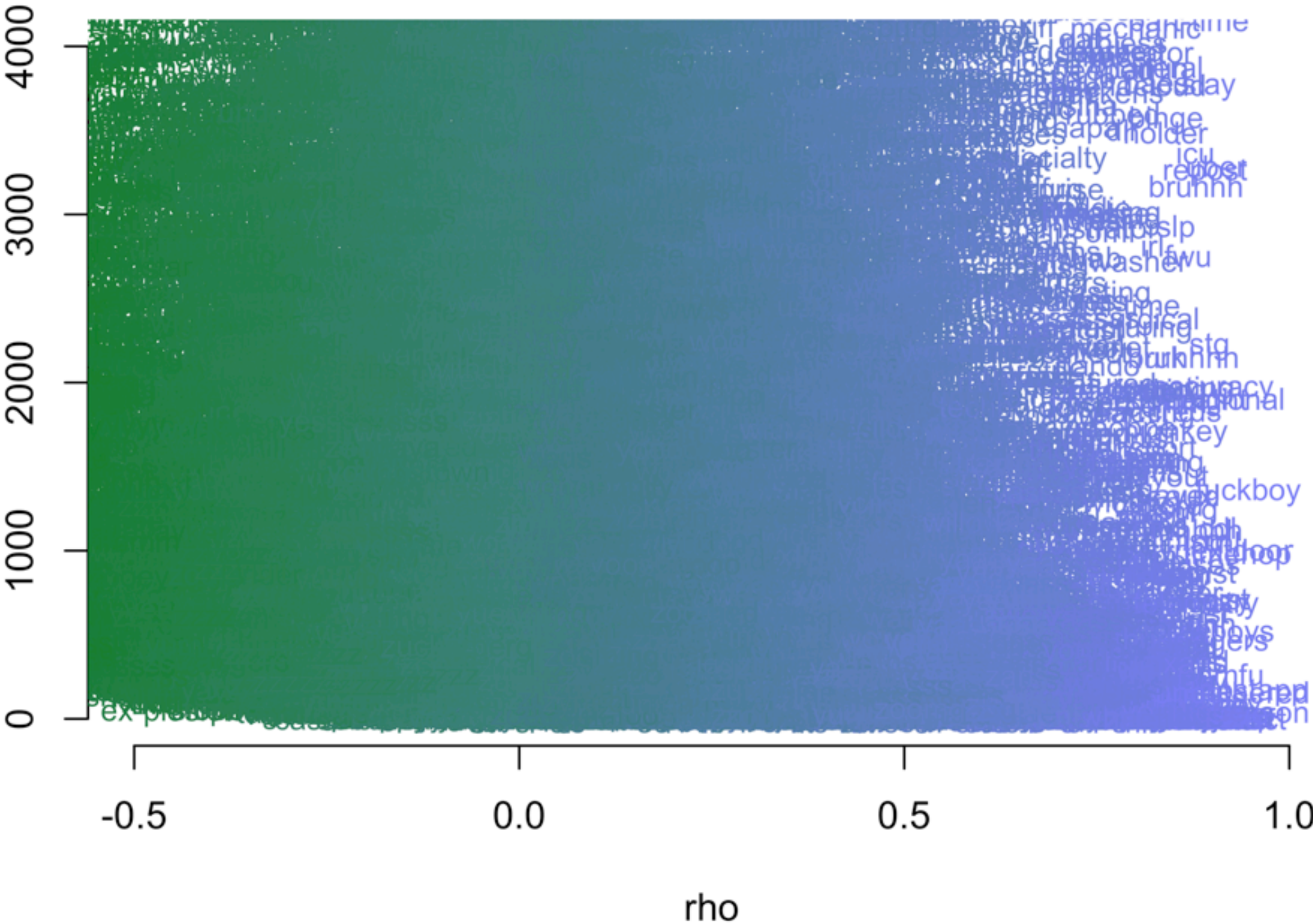




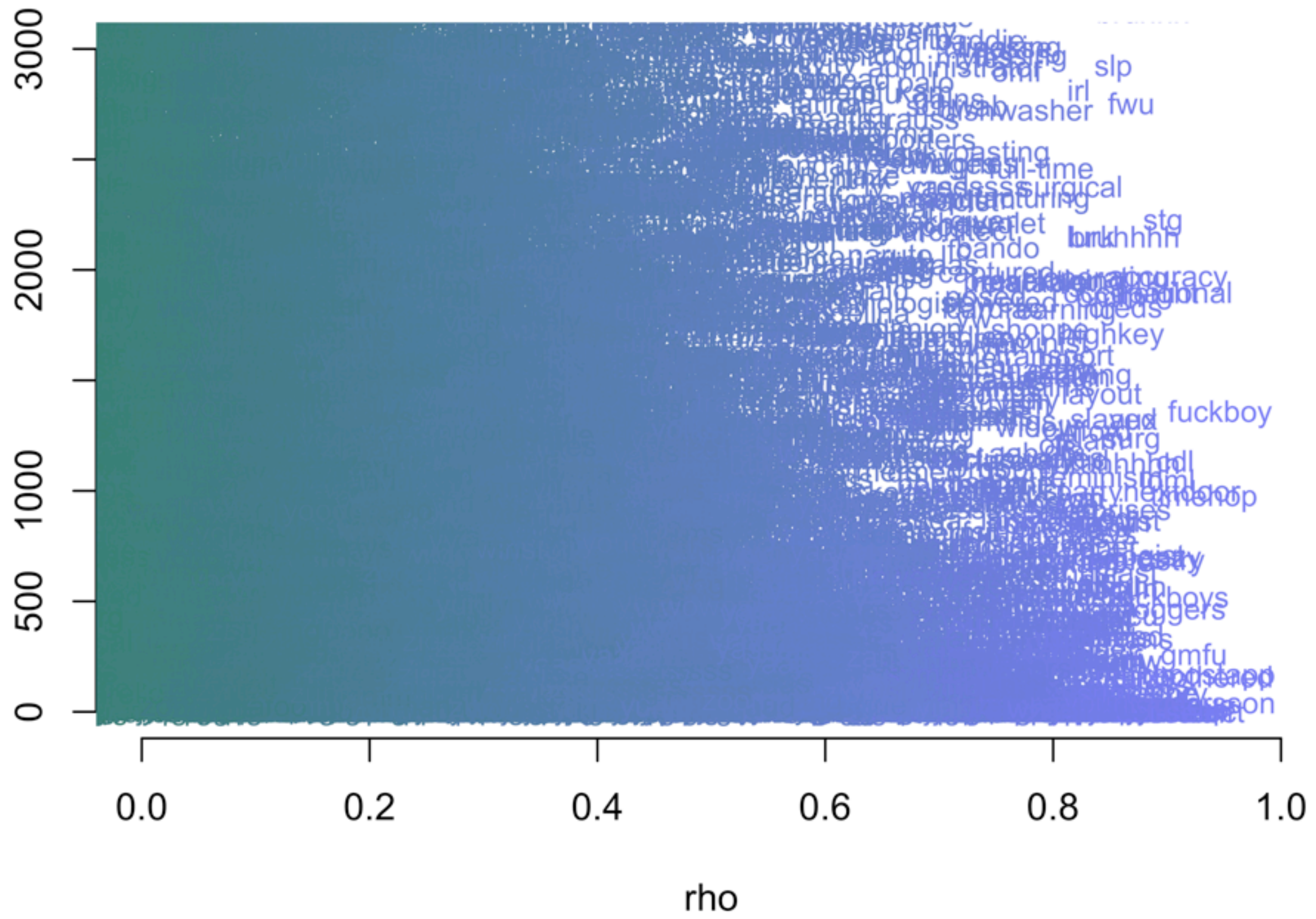




2013 Q4 Rel Freq Per Billion Words

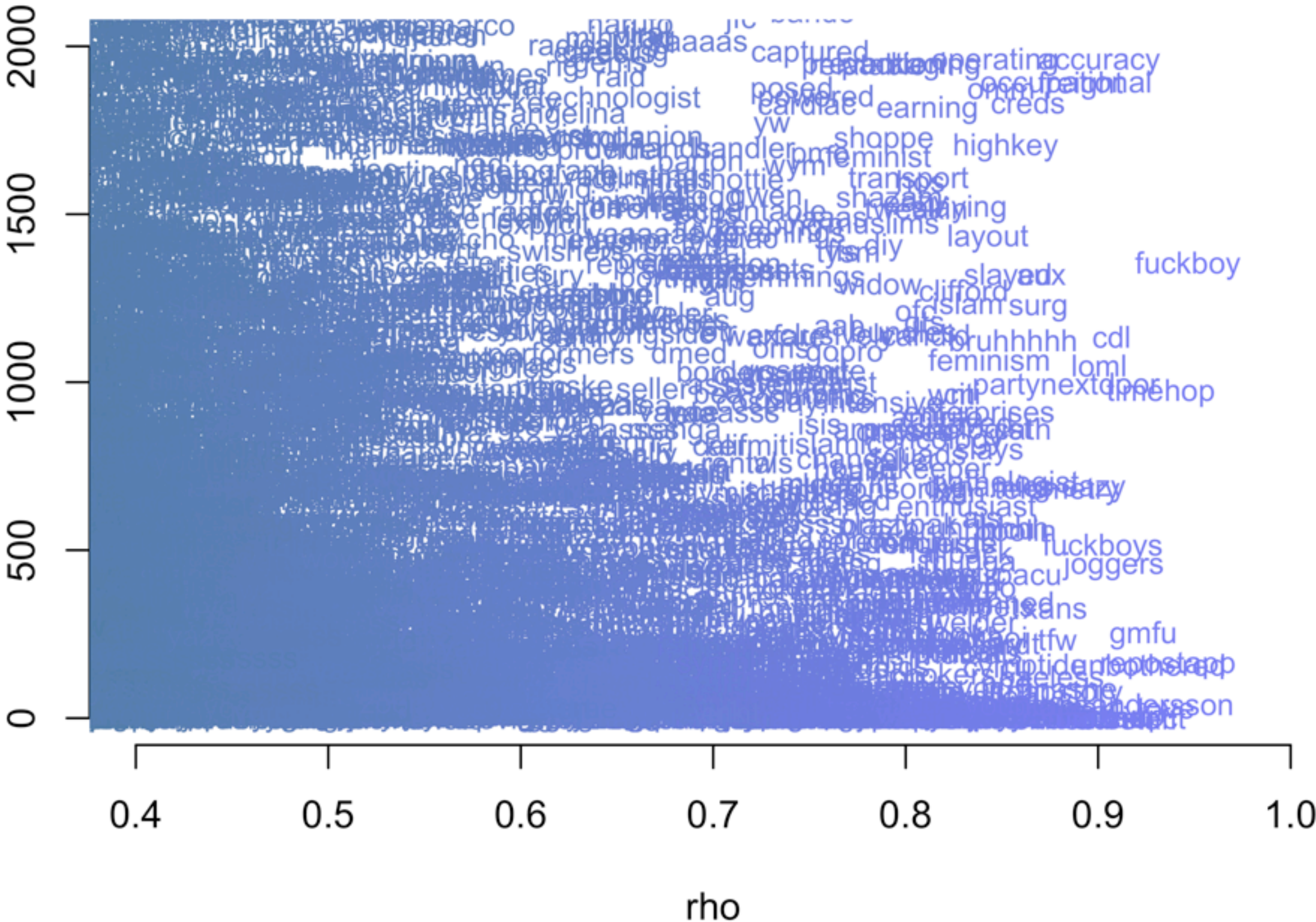


2013 Q4 Rel Freq Per Billion Words



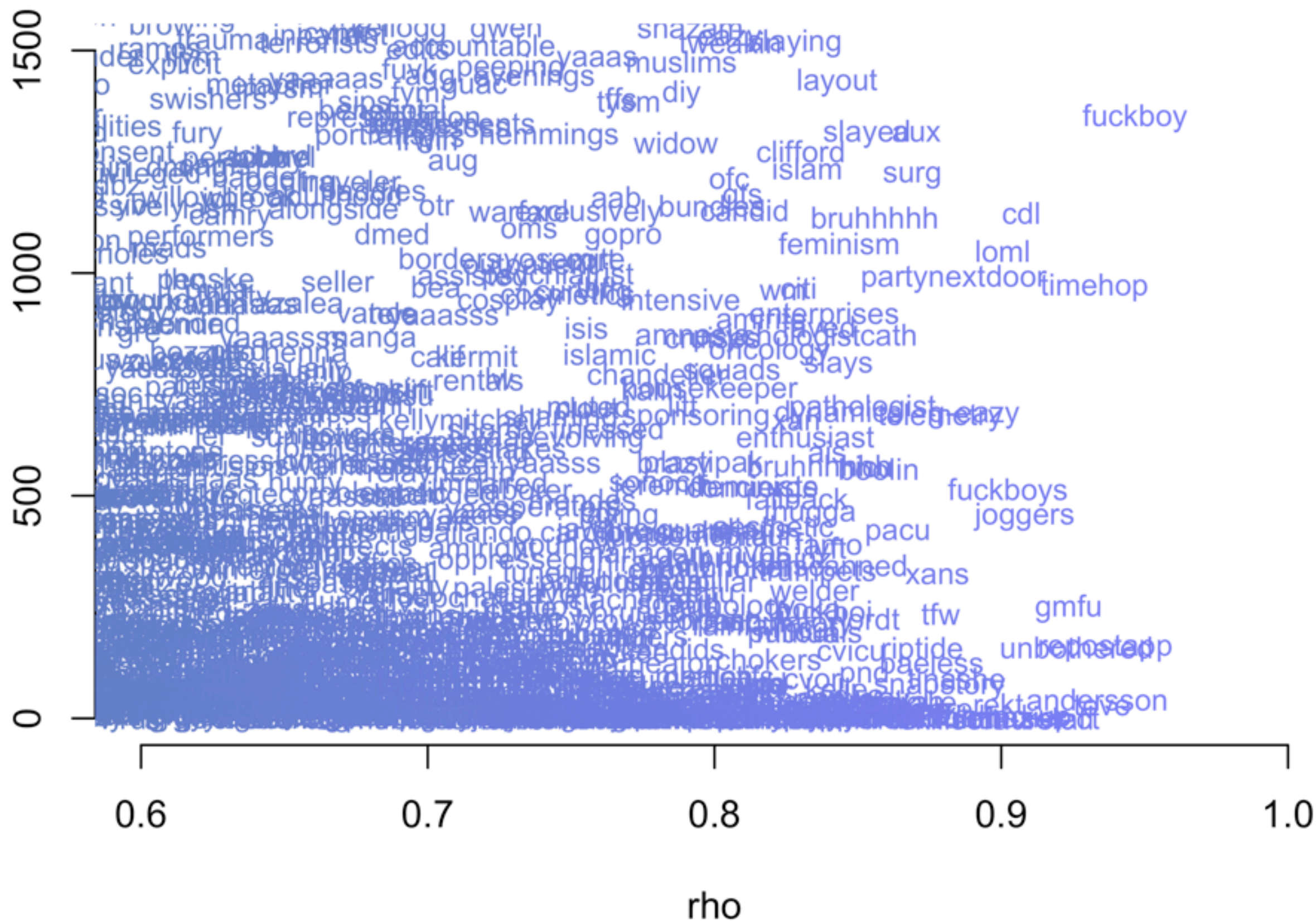


2013 Q4 Rel Freq Per Billion Words





2013 Q4 Rel Freq Per Billion Words



# Emerging Words

Through this process 81 emerging words were identified, including 27 spelling variants or inflected forms, leaving 54 unique emerging words for further analysis.

In general these words can be characterized as everyday slang; most are not restricted to Twitter or to CMC more generally (including acronyms).

Common topical domains include family and friends, sex and relationships, intoxication, technology, and Japanese culture.



# Emerging Words on Twitter 2014

| <b>Words</b>           | <b>Definition</b>                            |
|------------------------|--|
| <i>amirite</i>         | Am I right?                                  |
| <i>baeless</i>         | Single                                       |
| <i>baeritto</i>        | Bae (i.e. significant other)                 |
| <i>balayage</i>        | Hair style                                   |
| <i>boolin</i>          | Hanging out                                  |
| <i>brazy</i>           | Crazy  |
| <i>bruuh</i>           | Bro  |
| <i>candids</i>         | Candid public picture                        |
| <i>celfie</i>          | Selfie                                       |
| <i>cosplay</i>         | Costume role playing                         |
| <i>dwk</i>             | Driving While Kissing                        |
| <i>fallback (game)</i> | Skillful at talking one's way out of trouble |
| <i>famo</i>            | Family and friend                            |
| <i>faved</i>           | Favorited                                    |

# Emerging Words on Twitter 2014

| Words               | Definition                   |
|---------------------|------------------------------|
| <i>fhritp</i>       | Fuck Her Right In The Pussy  |
| <i>figgity</i>      | Intoxicated; Very            |
| <i>(on) fleek</i>   | (On) point                   |
| <i>fuckboys</i>     | Assholes                     |
| <i>gainz</i>        | Earnings                     |
| <i>gmfu</i>         | Get Me Fucked Up             |
| <i>goalz</i>        | Goals (i.e. life goals)      |
| <i>idgt</i>         | I Don't Get Tired            |
| <i>lfie</i>         | Life                         |
| <i>lifestyleeee</i> | Lifestyle                    |
| <i>litt</i>         | Lit (i.e. intoxicated, good) |
| <i>litty</i>        | Lit (i.e. intoxicated, good) |
| <i>lituation</i>    | A lit situation              |
| <i>lordt</i>        | Lord (esp. as exclamation)   |

# Emerging Words on Twitter 2014

| Words                 | Definition                                  |
|-----------------------|---|
| <i>lw</i>             | Light Weight                                |
| <i>mce</i>            | Man Crush Everyday                          |
| <i>mmmmmmmuah</i>     | Laughter                                    |
| <i>mutuals</i>        | Mutual friends                              |
| <i>nahfr</i>          | Nah For Real                                |
| <i>notifs</i>         | Notifications (esp. online)                 |
| <i>pcd</i>            | Post Concert Depression                     |
| <i>pullout (game)</i> | Skillful at <i>coitus interruptus</i>       |
| <i>rekt</i>           | Wrecked (i.e. intoxicated; defeated esp. in |
| <i>rq</i>             | Real Quick                                  |
| <i>scute</i>          | Cute  |
| <i>senpai</i>         | Father figure                               |
| <i>shordy</i>         | Shorty (i.e. a young woman)                 |
| <i>slayin</i>         | Slaying                                     |

# Emerging Words on Twitter 2014

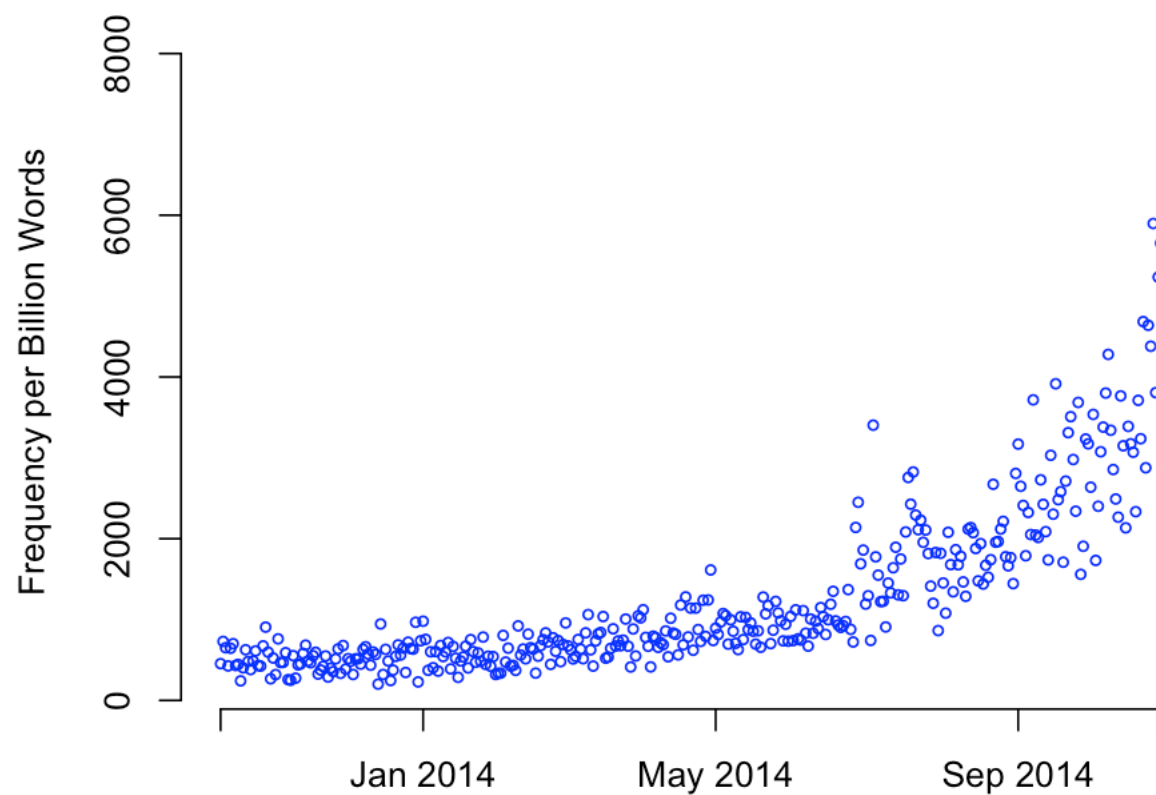
| Words            | Definition                          |
|------------------|-------------------------------------|
| <i>squad</i>     | Squad (i.e. a crew)                 |
| <i>tbfh</i>      | To Be Fucking Honest                |
| <i>tfw</i>       | That Feel When                      |
| <i>thotful</i>   | Thoughtful                          |
| <i>thottin</i>   | Looking for thots (i.e. promiscuous |
| <i>tookah</i>    | Marijuana                           |
| <i>traphouse</i> | Drug house                          |
| <i>unbae</i>     | Break up with                       |
| <i>waifu</i>     | Wife                                |
| <i>wce</i>       | Woman Crush Everyday                |
| <i>xans</i>      | Benzodiazapane pills                |
| <i>yaas</i>      | Yes                                 |

# General Findings

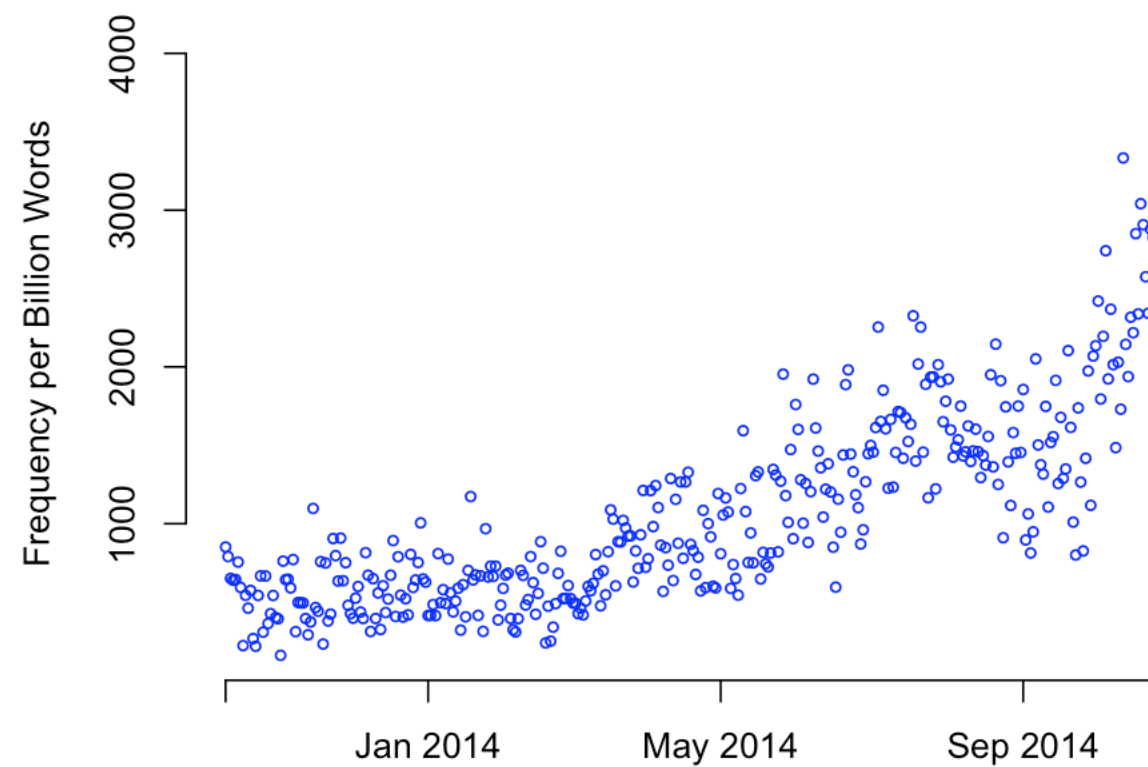
By providing a relatively large and unbiased set of emerging words, generalizations can be made about the process of lexical emergence:

1. Most of the emerging words were formed using standard **word formation processes**, with truncation, compounding and acronymization being most common.
2. The relative frequencies of most newly emerging words are characterized by **s-shaped curve** time series (or partial s-shaped curves).

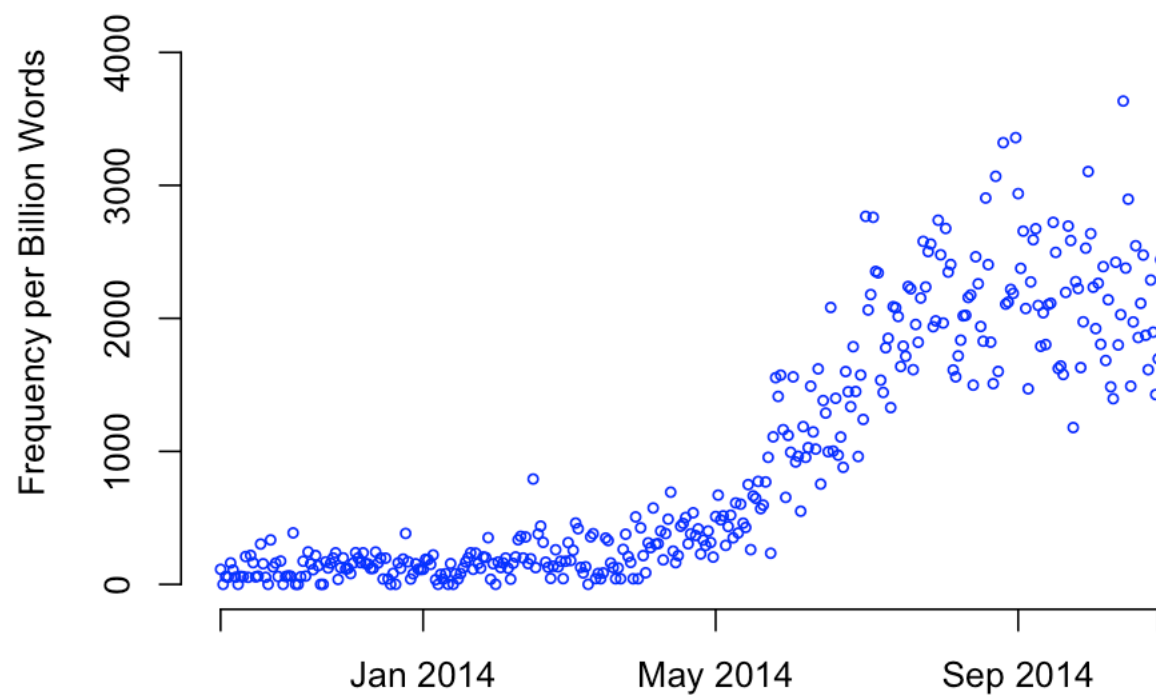
**FUCKBOYS**



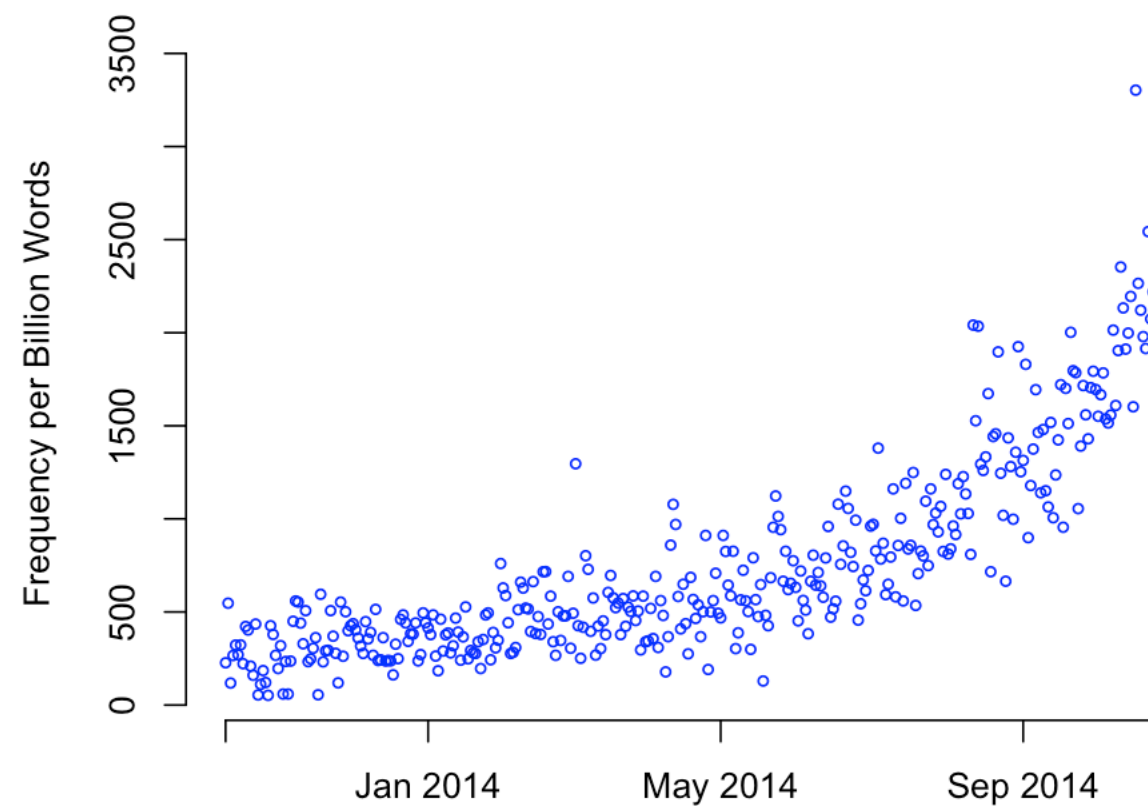
**BOOLIN**



**BAELESS**



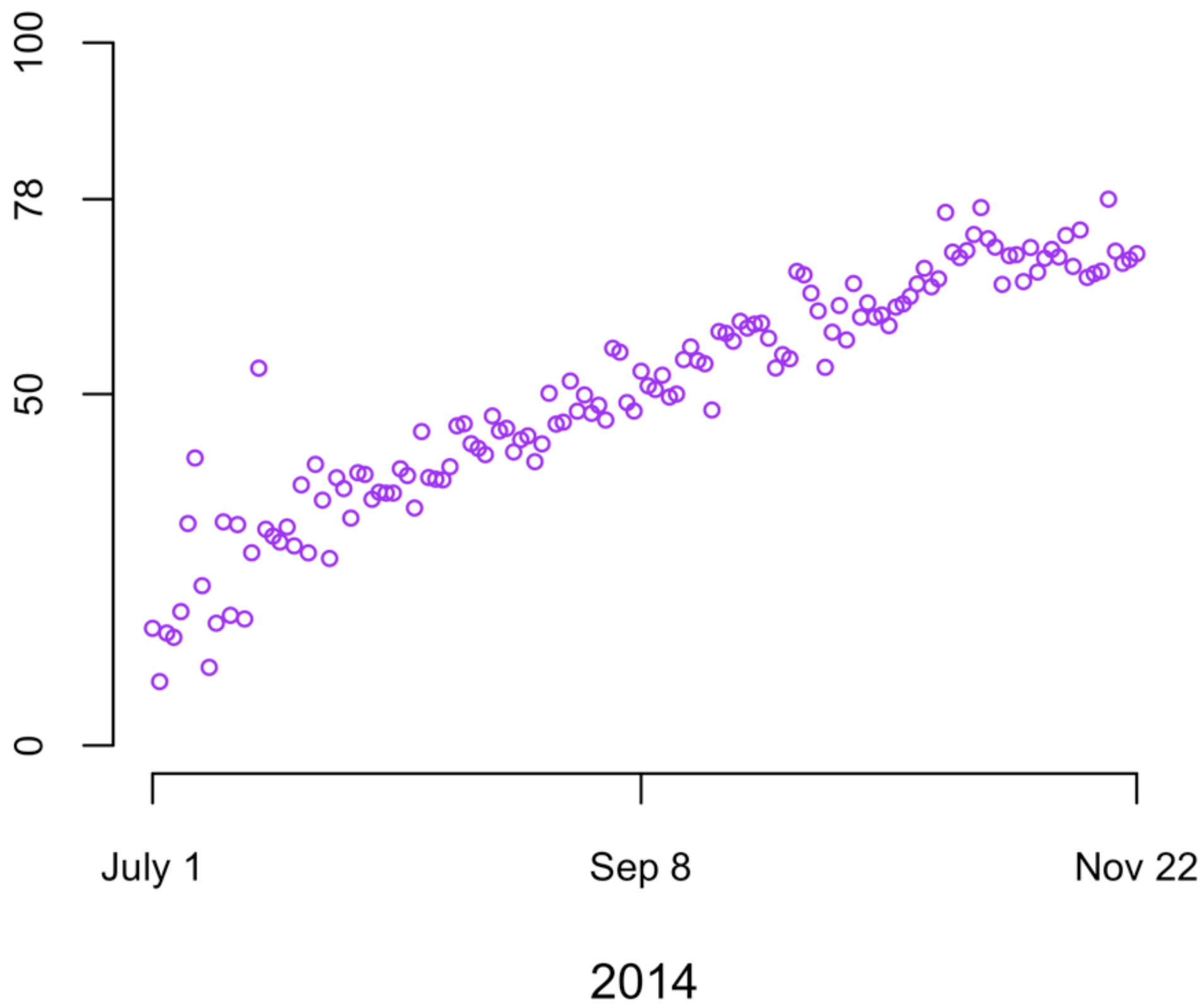
**XANS**



# General Findings

3. Numerous words emerge into the language each year; some persist, whereas others quickly die.
4. Most of the newly emerging words are not really new; rather, they have been laying dormant since the 2000s.
5. The meanings of a number of the newly emerging words appear to generalize over time.

Percentage of Tweets containing Fleek  
that do not contain (Eye)brow(s)





# Publication and Future Research

For more information on these results please see

Grieve, Nini & Guo. 2016. Analyzing lexical emergence in Modern American English online. *Forthcoming in English Language and Linguistics*. (Open access pre-print available on publisher's website).

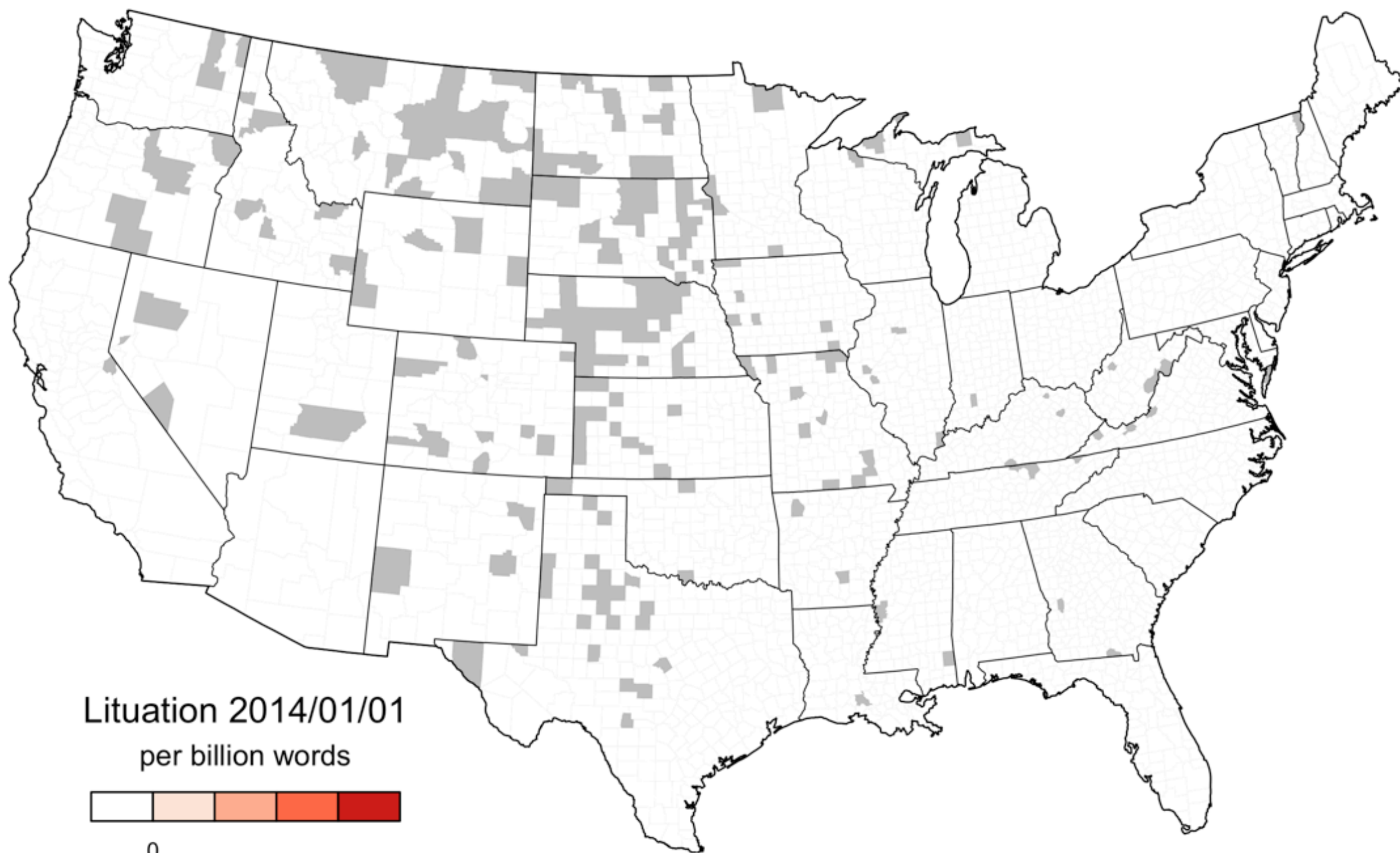
Future research: What triggers the emergence of a new word? What predicts the success or failure of new word? How can you model the rise in the frequency of a new word over time?

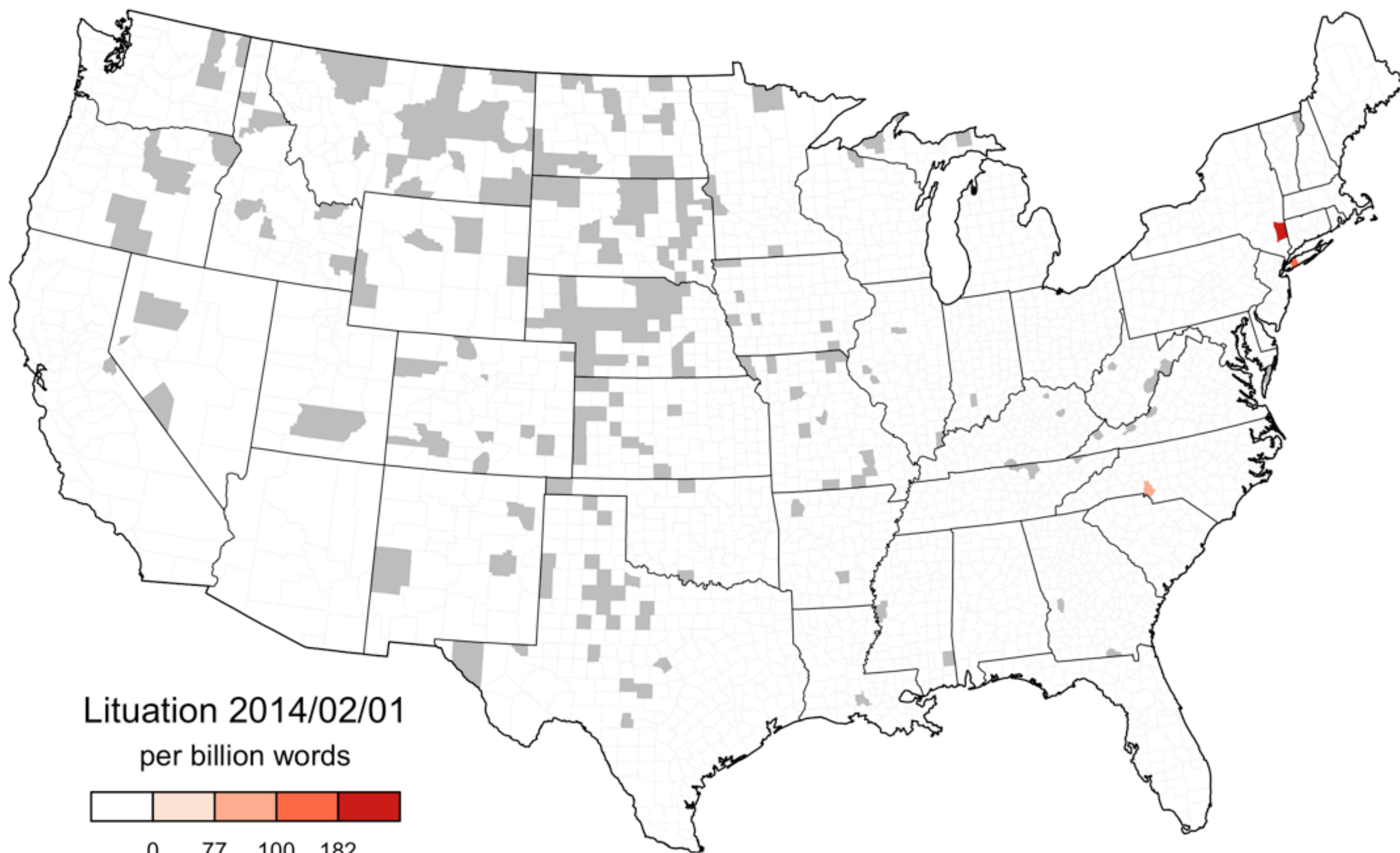
# Mapping Lexical Emergence

In addition to identifying and charting the spread of new words over time, because the Twitter corpus contains precise **geocoding** information, it is also possible to analyze lexical emergence from a regional perspective, including identifying regional **hubs of lexical innovation**.

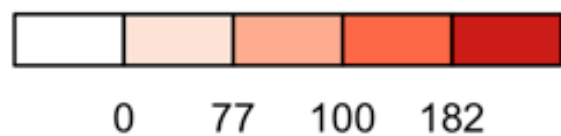
An initial problem, however, is how to map the spread of an emerging word.

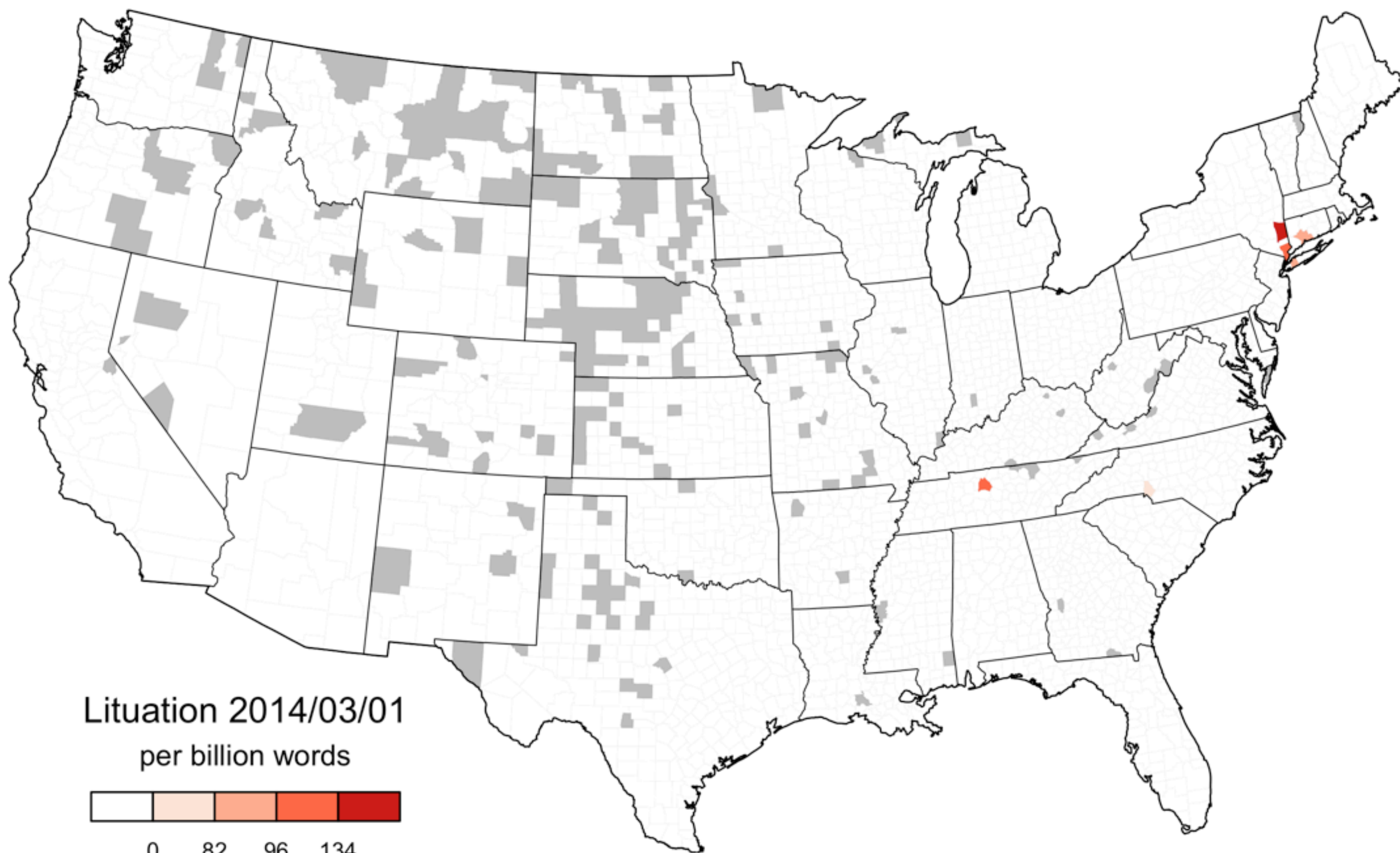
A basic approach is to map the **cumulative relative frequency** of the word across regionally defined sub-corpora (e.g. counties) over a series of points in time.



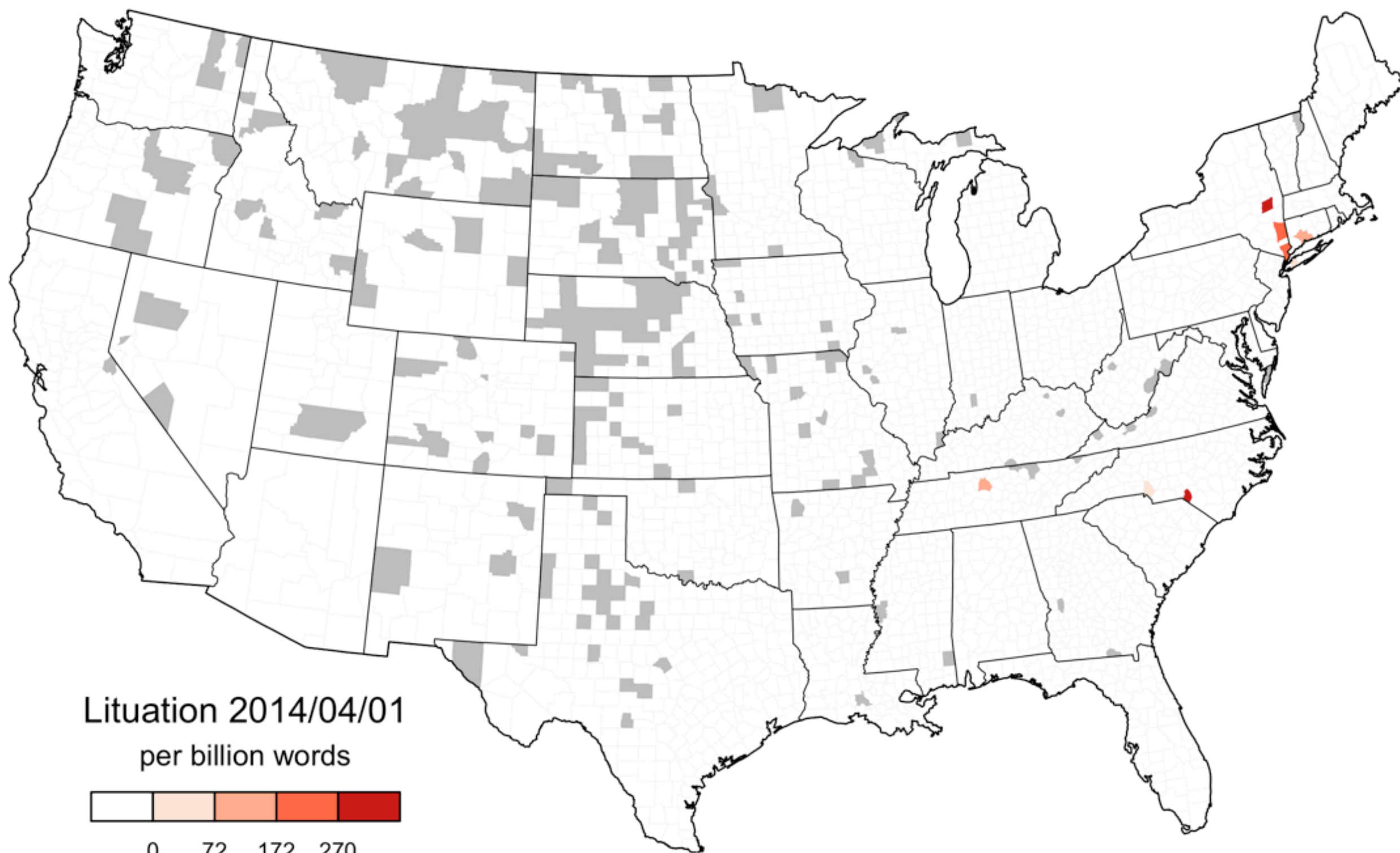


Litutation 2014/02/01  
per billion words

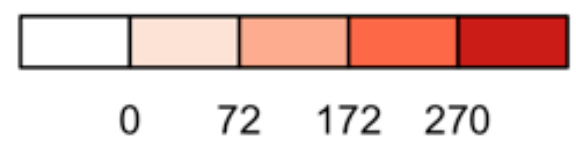


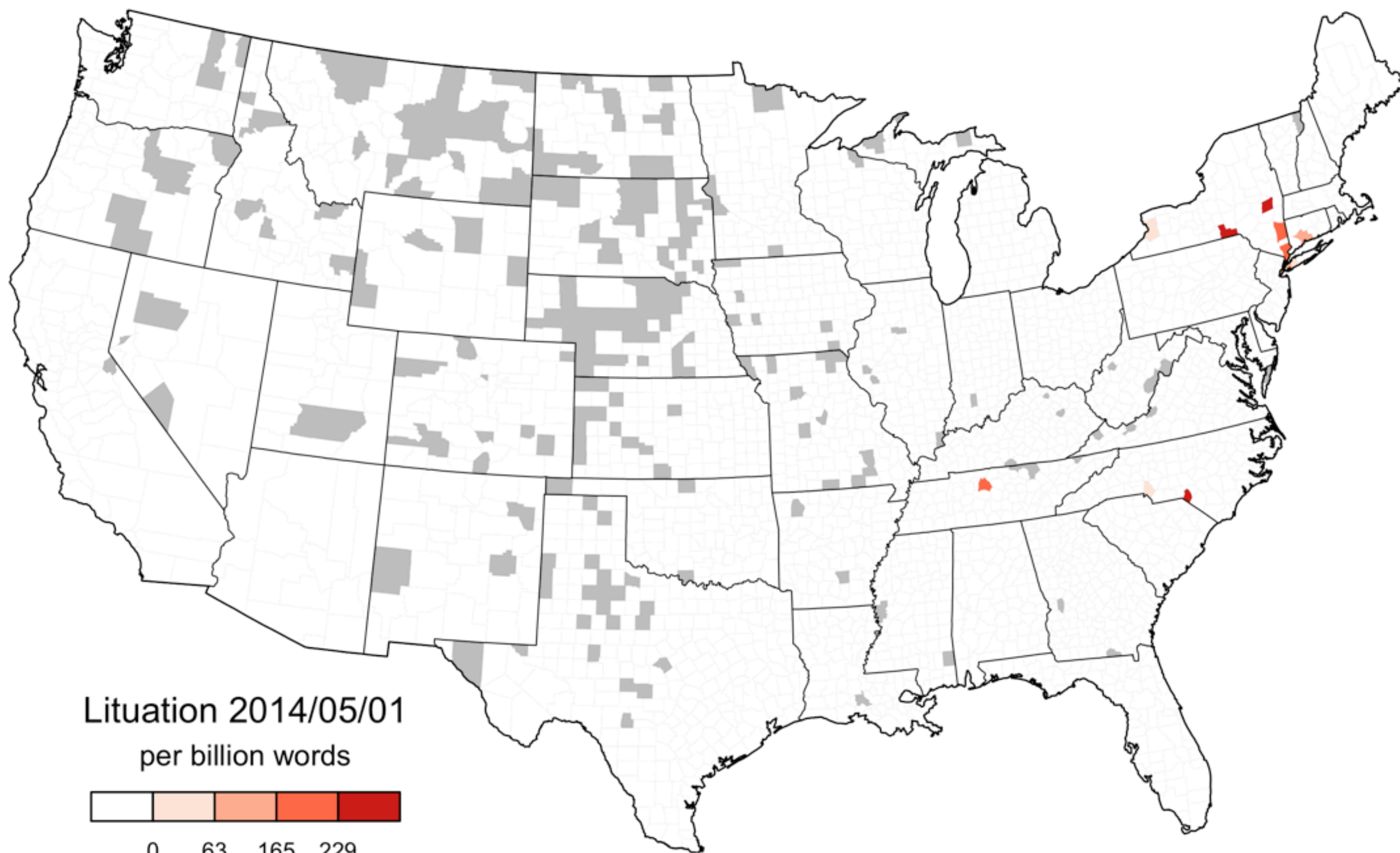




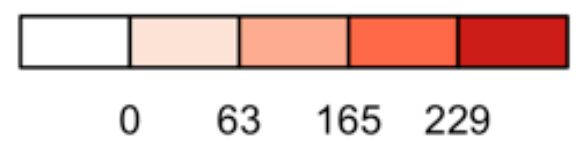


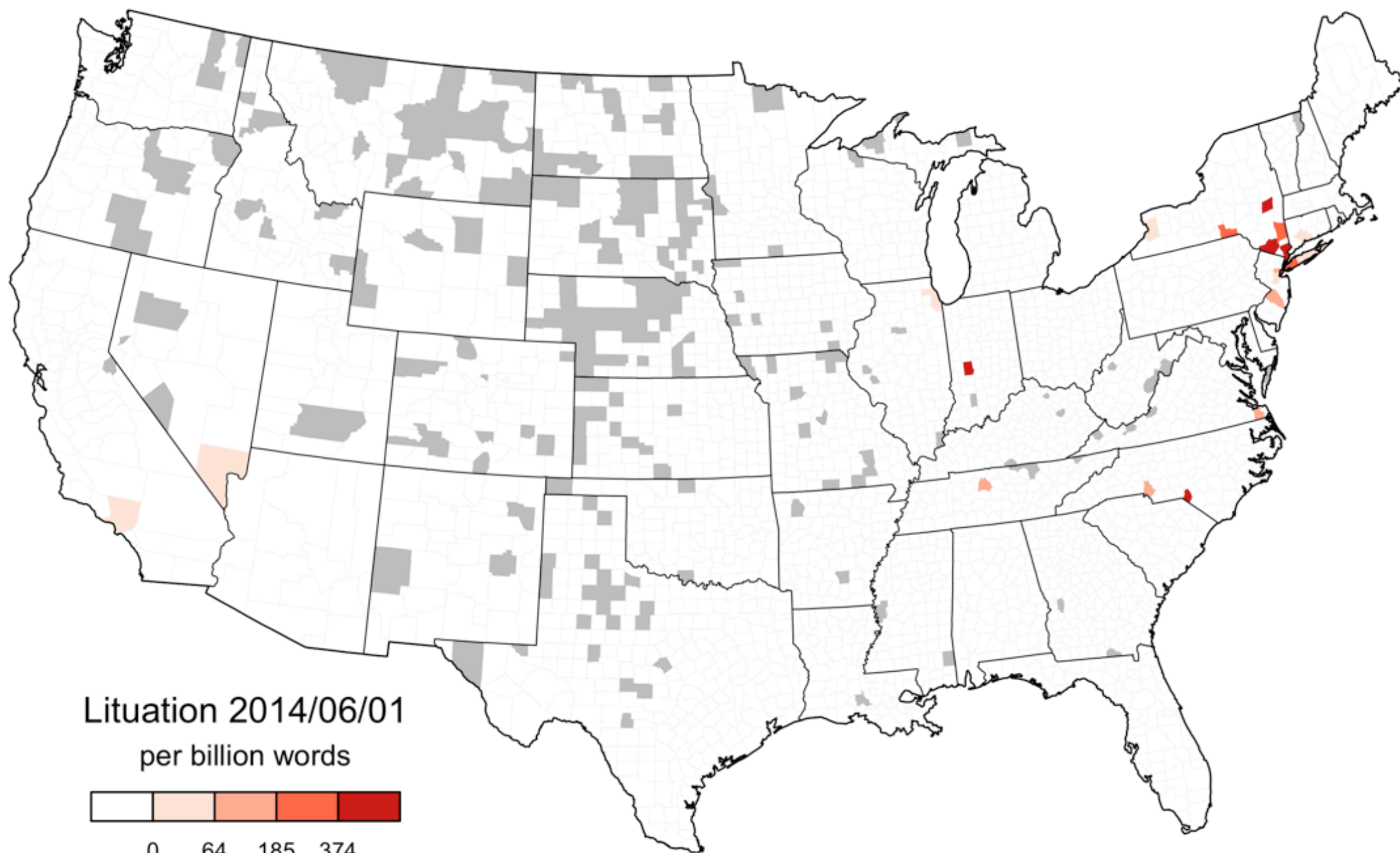
Lituation 2014/04/01  
per billion words



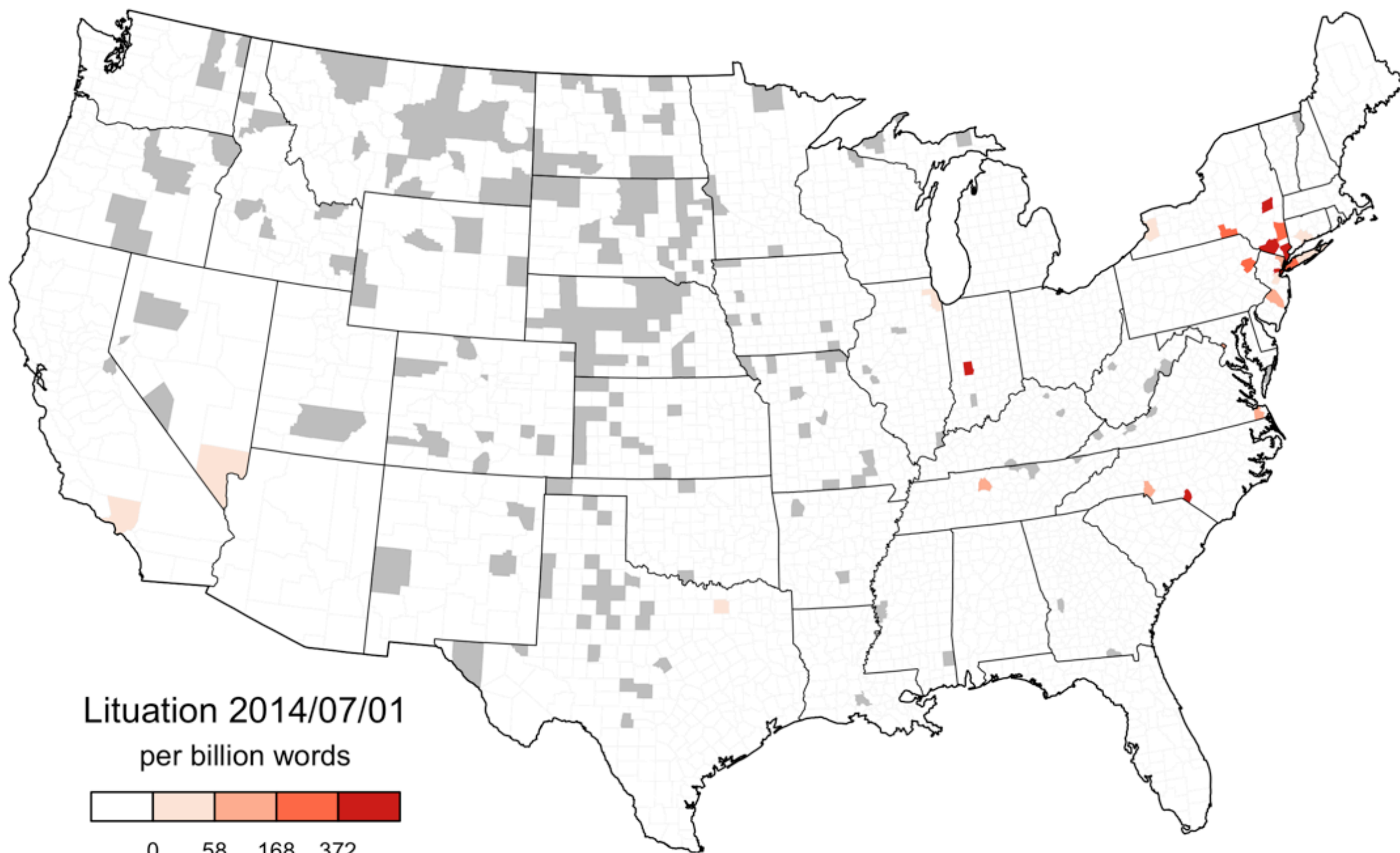


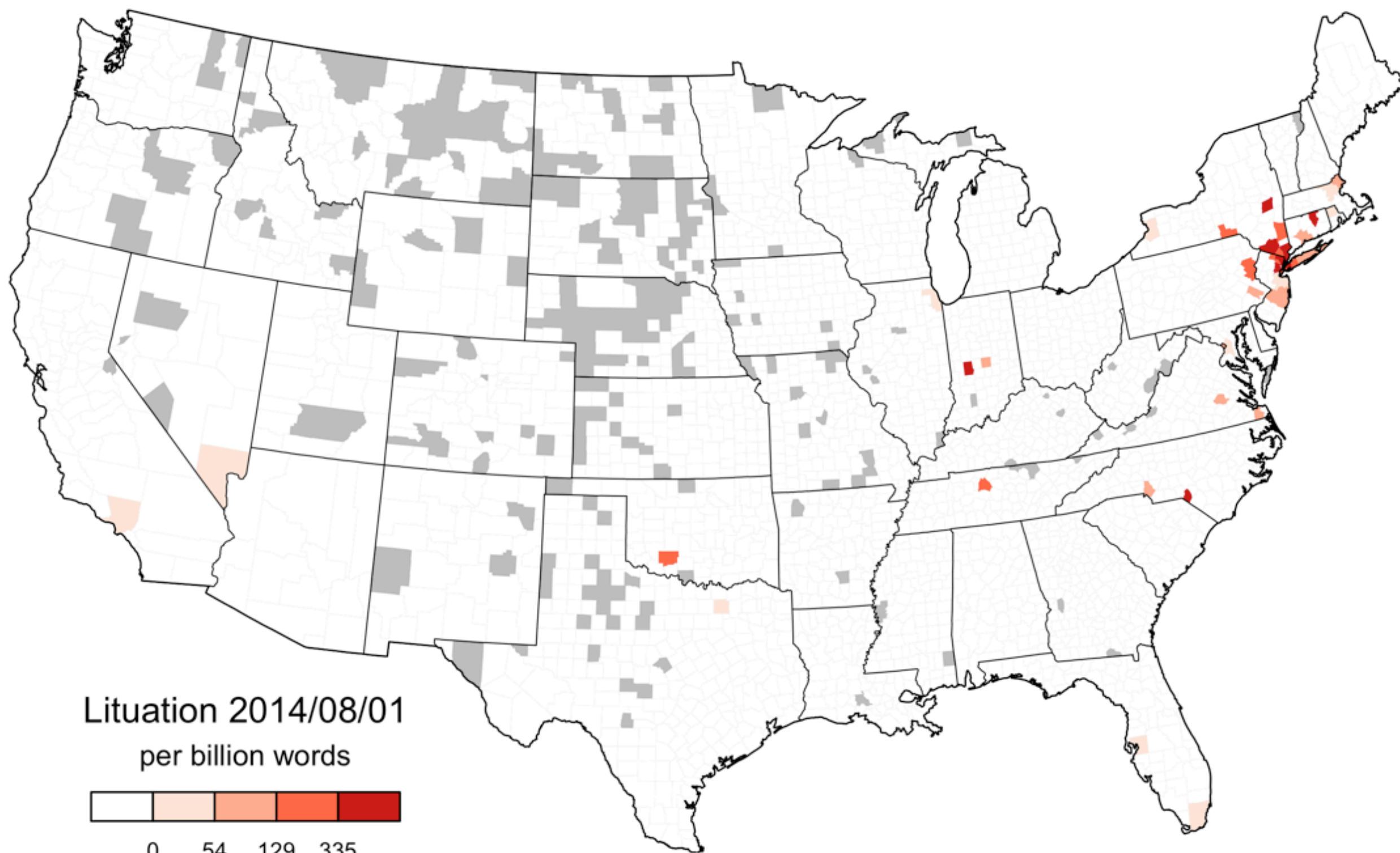
Lituation 2014/05/01  
per billion words

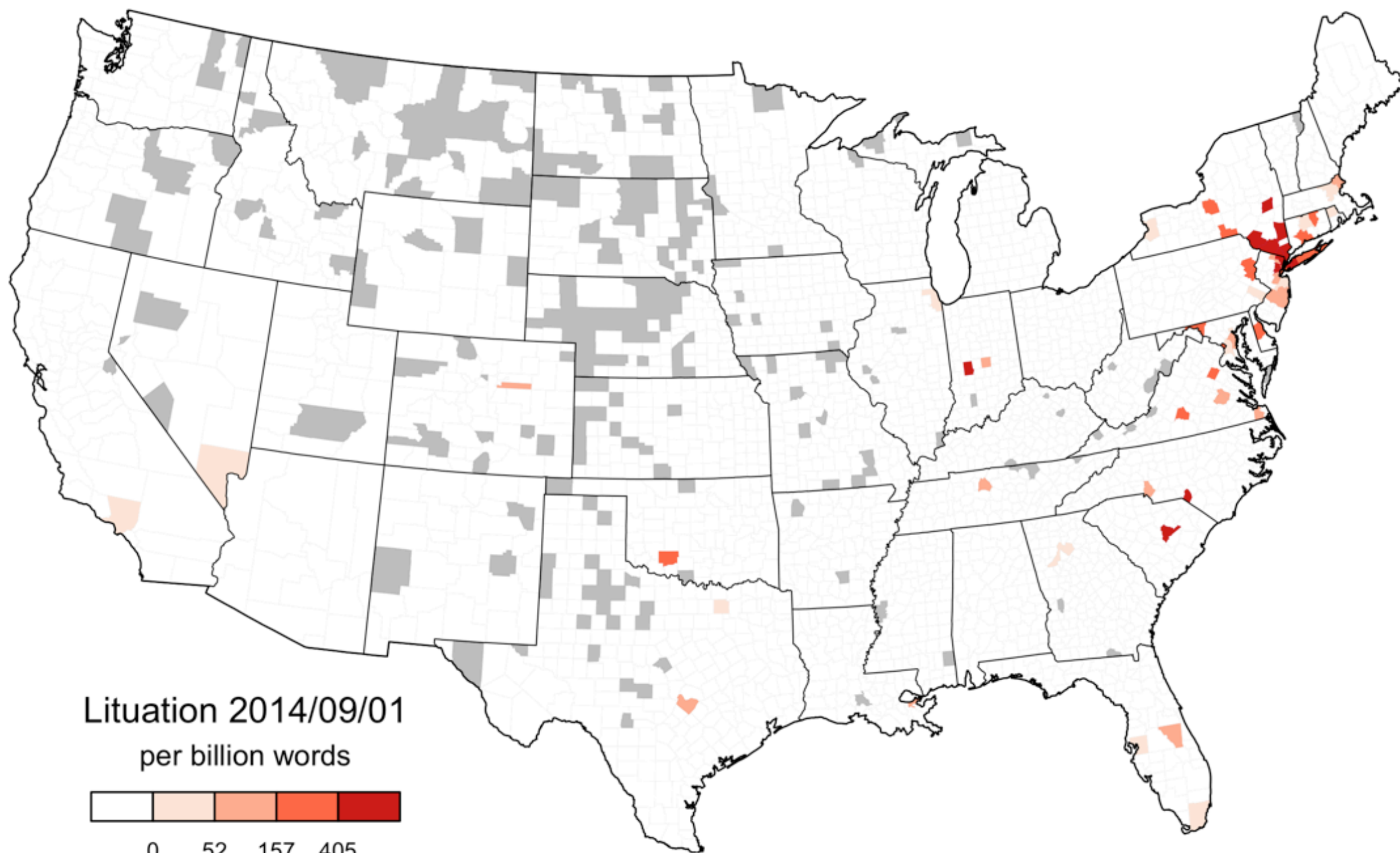




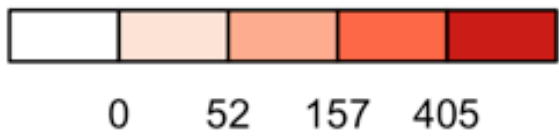




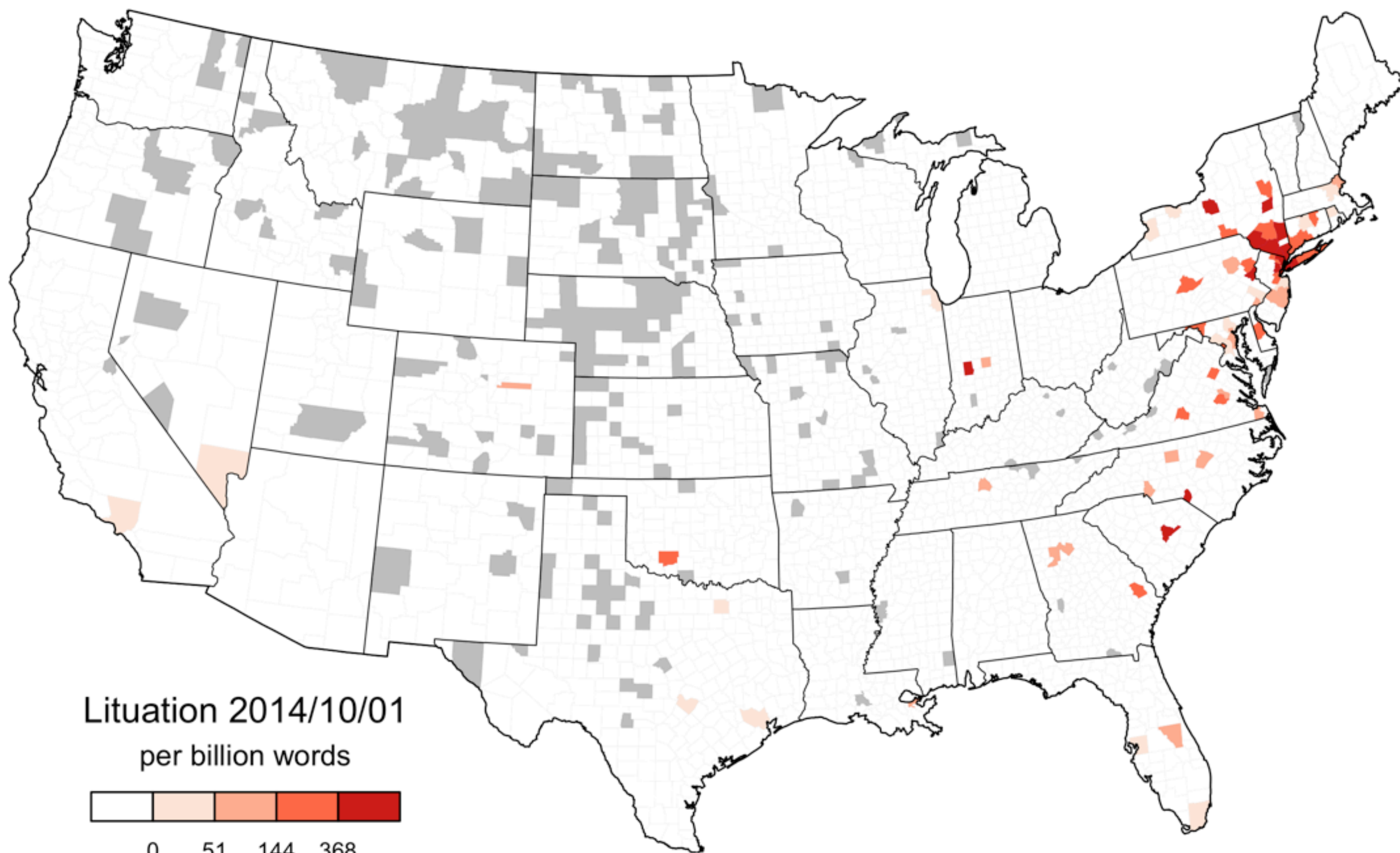


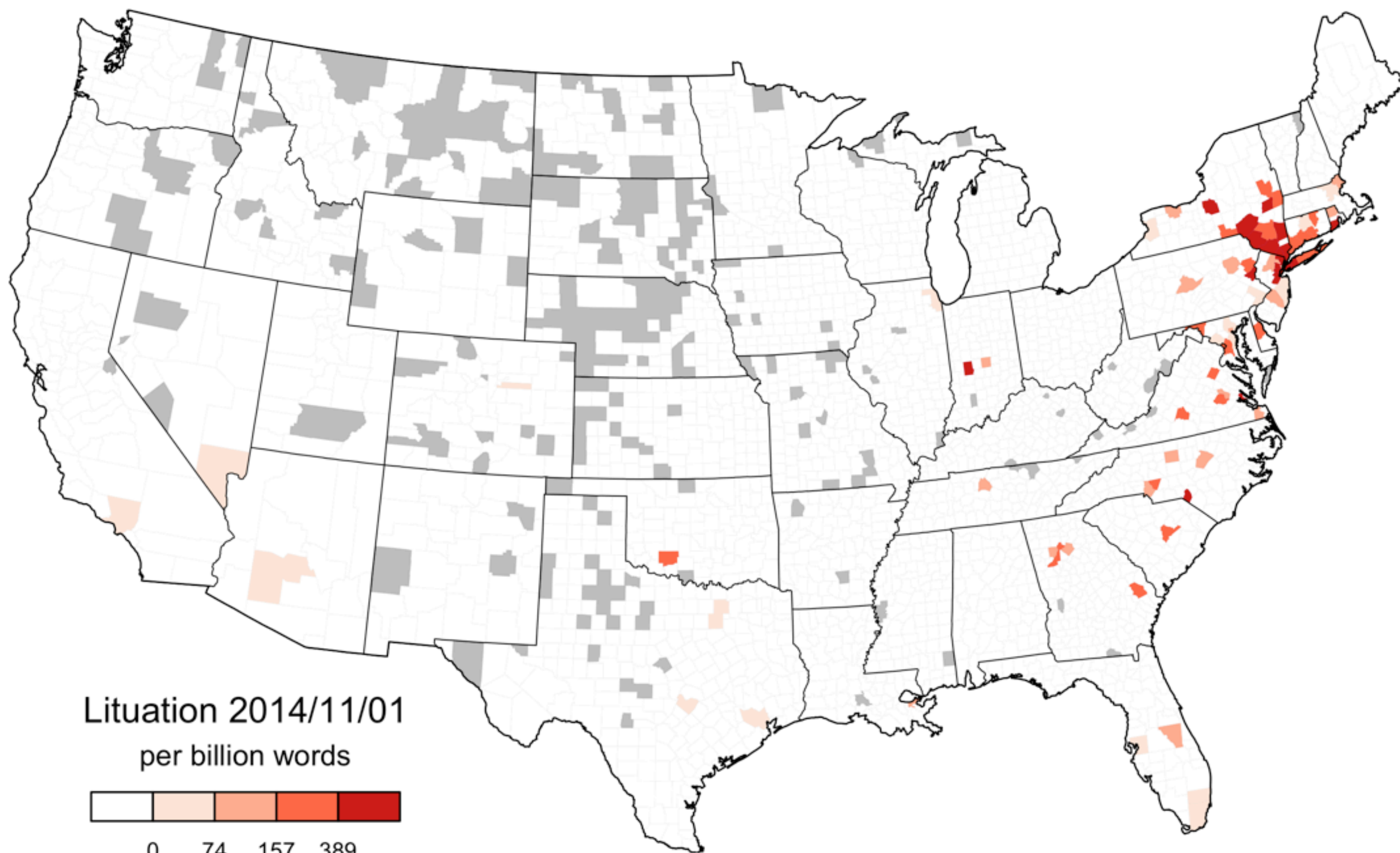


Lituation 2014/09/01  
per billion words







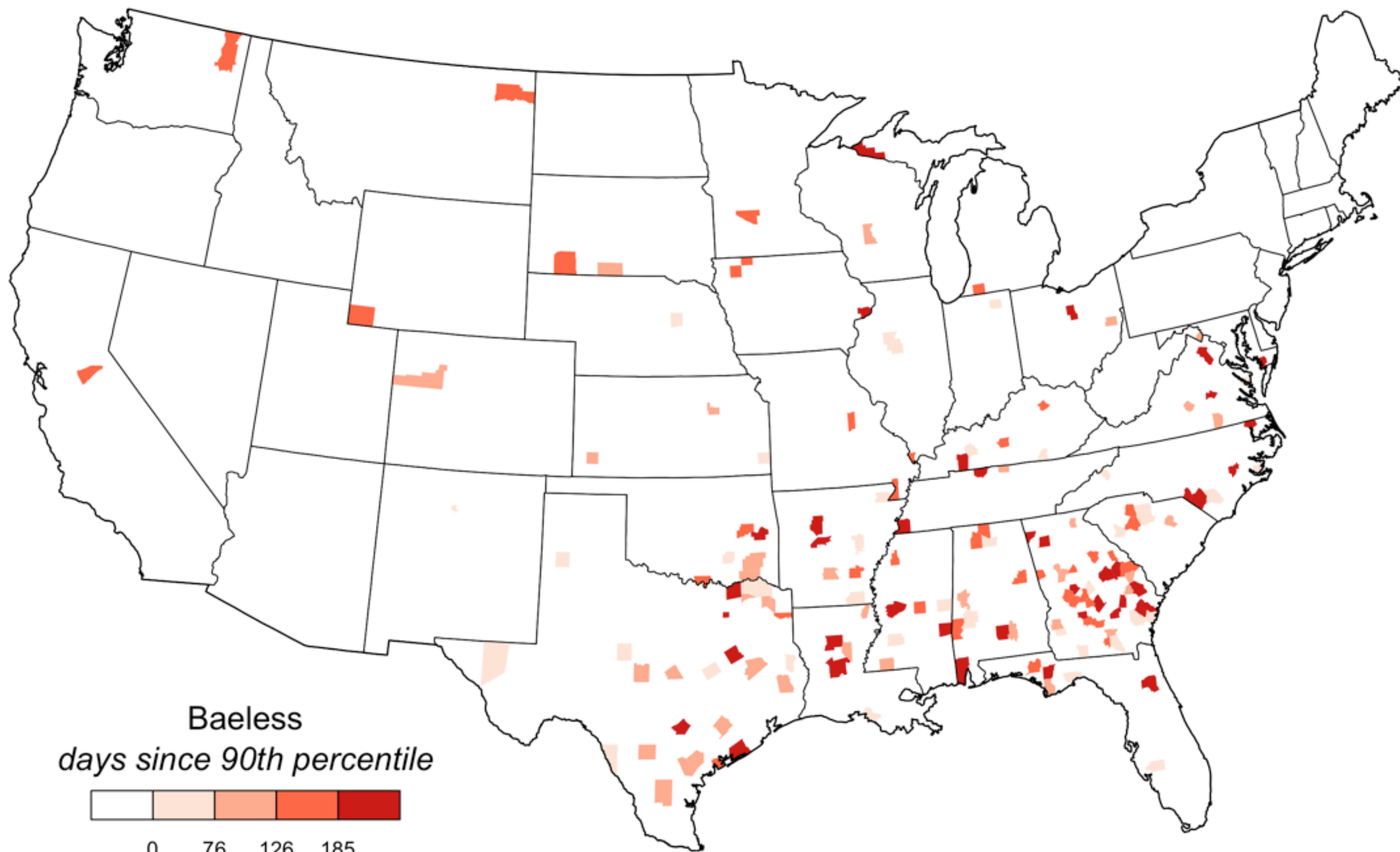


# Mapping Lexical Emergence

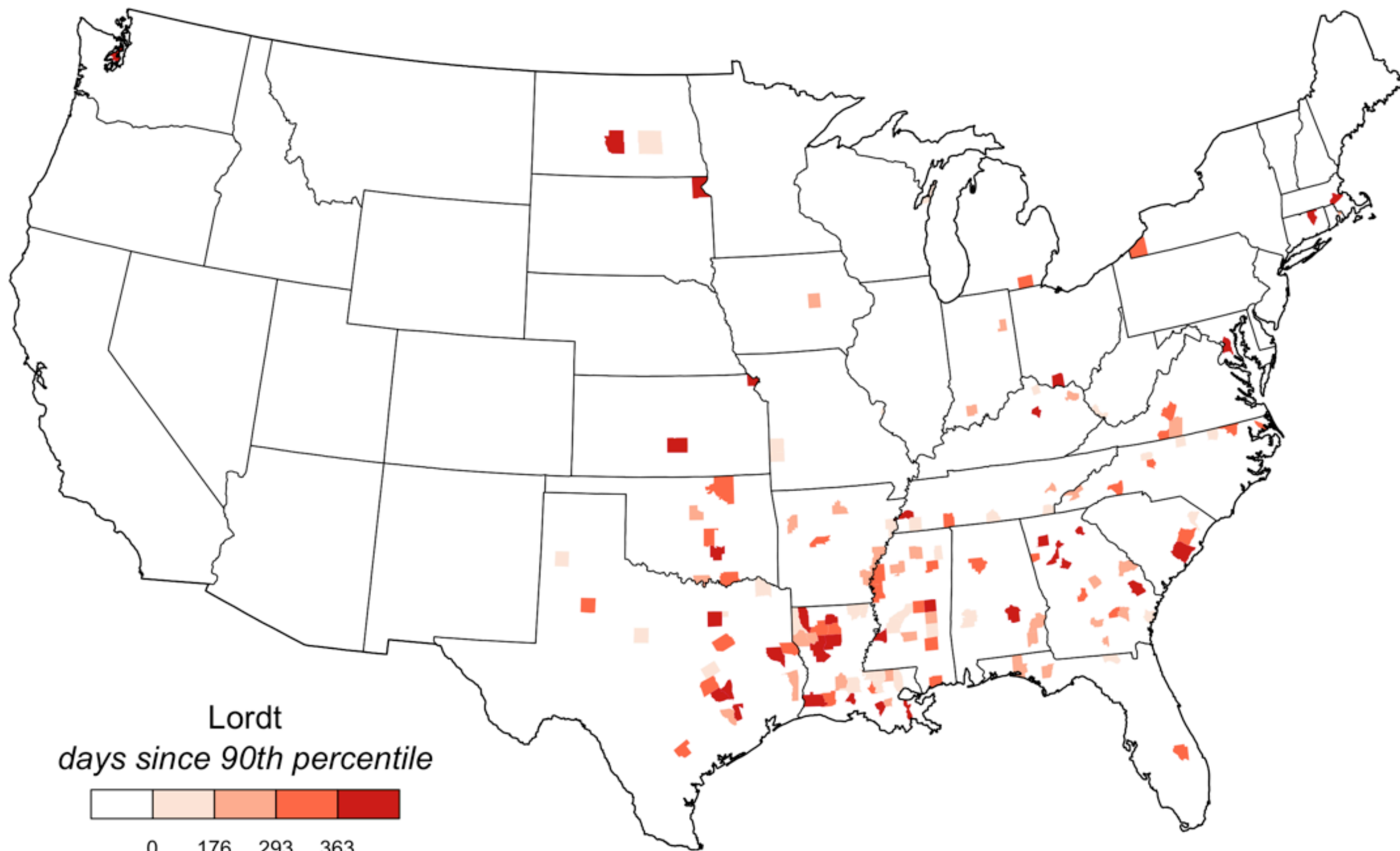
A problem with this approach is that it is unclear (at least to me) how to identify common regional patterns based on a series of maps for each word.

To solve this problem spatial time series for each word was reduced to a single map by measuring the **number of days since the word reached a specific relative frequency** (e.g. 90th percentile overall) by county.

This controls for variation in amount of data across counties (as opposed to mapping days since first occurrence, which follows population density).







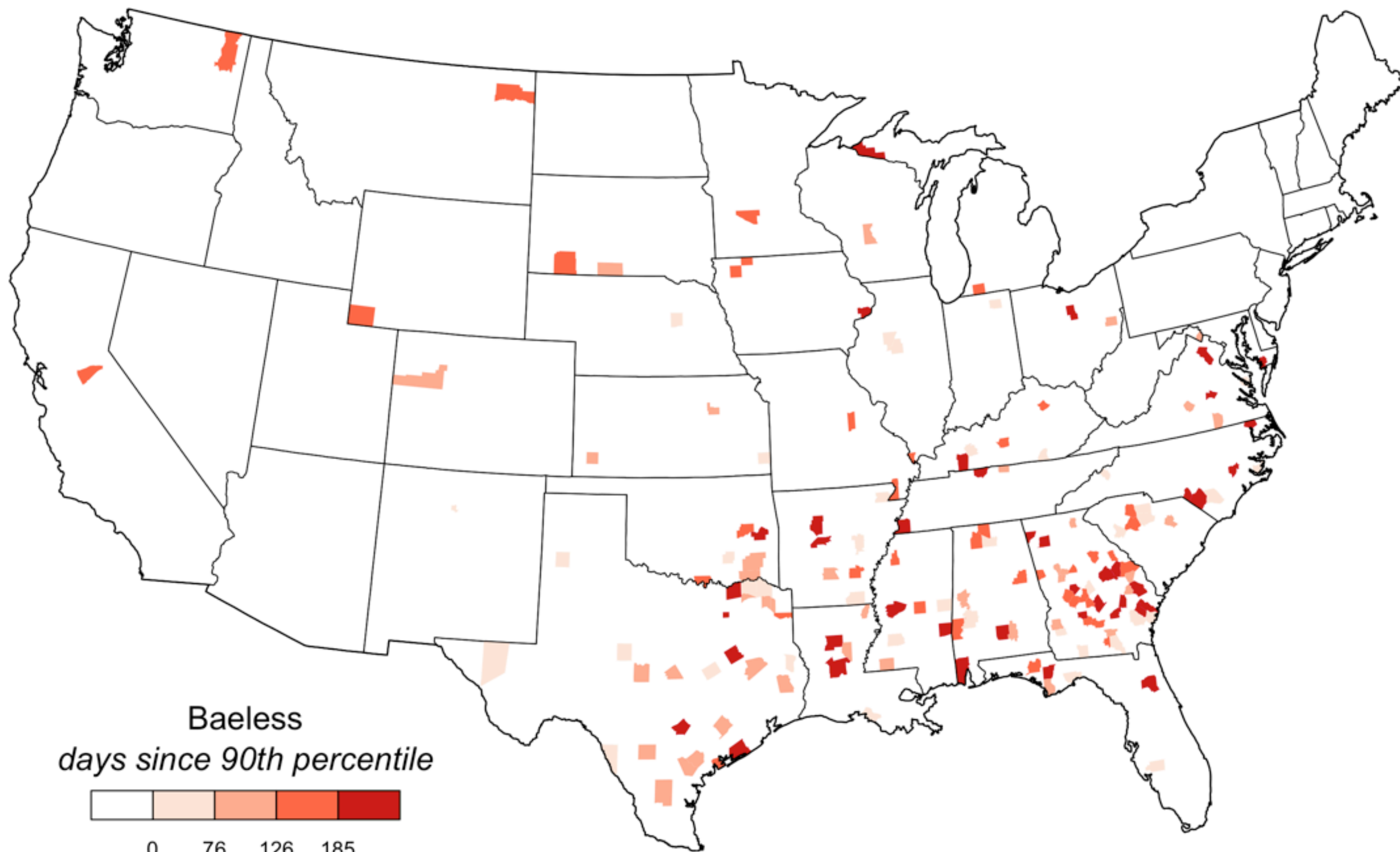


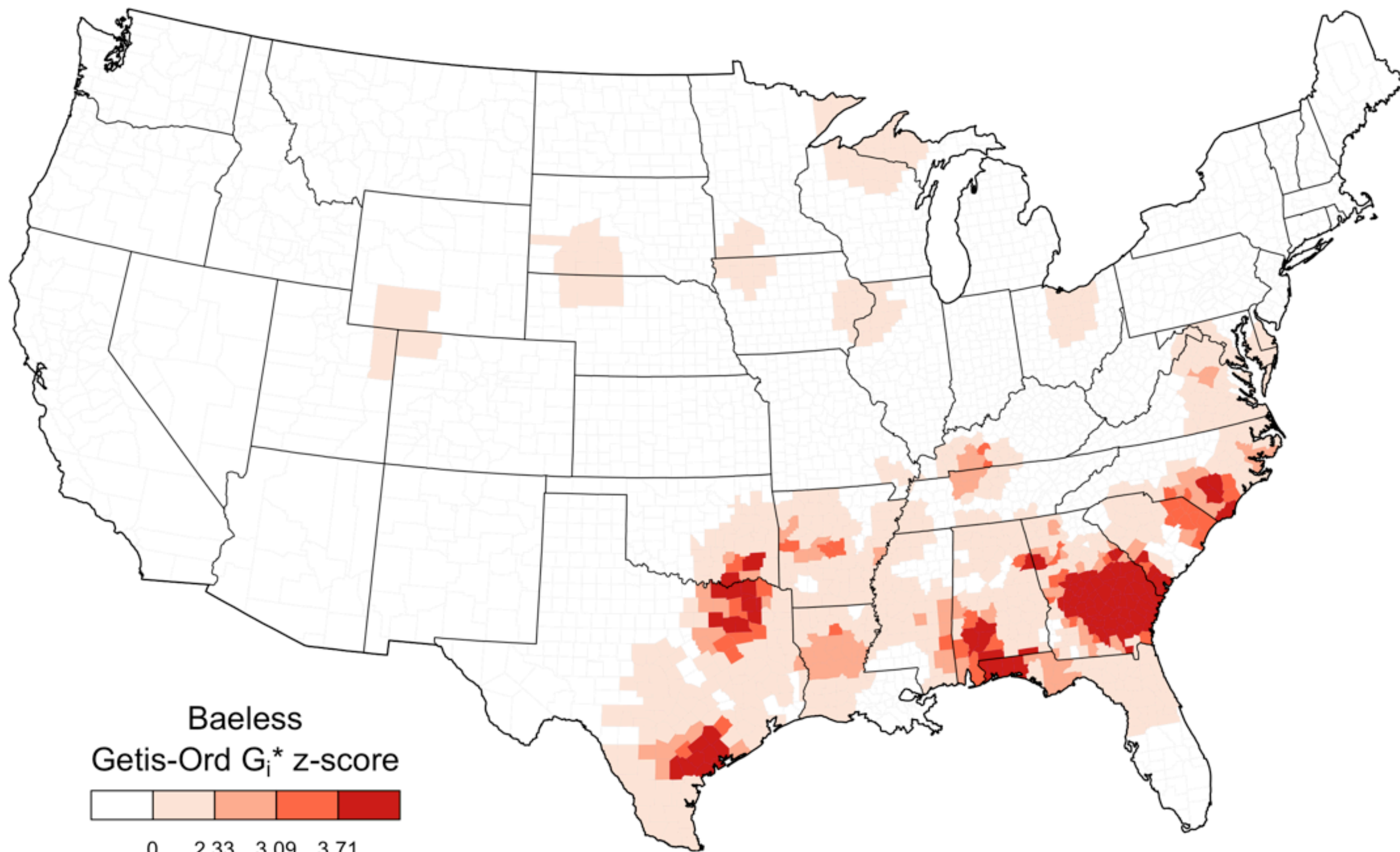
# Mapping Lexical Emergence

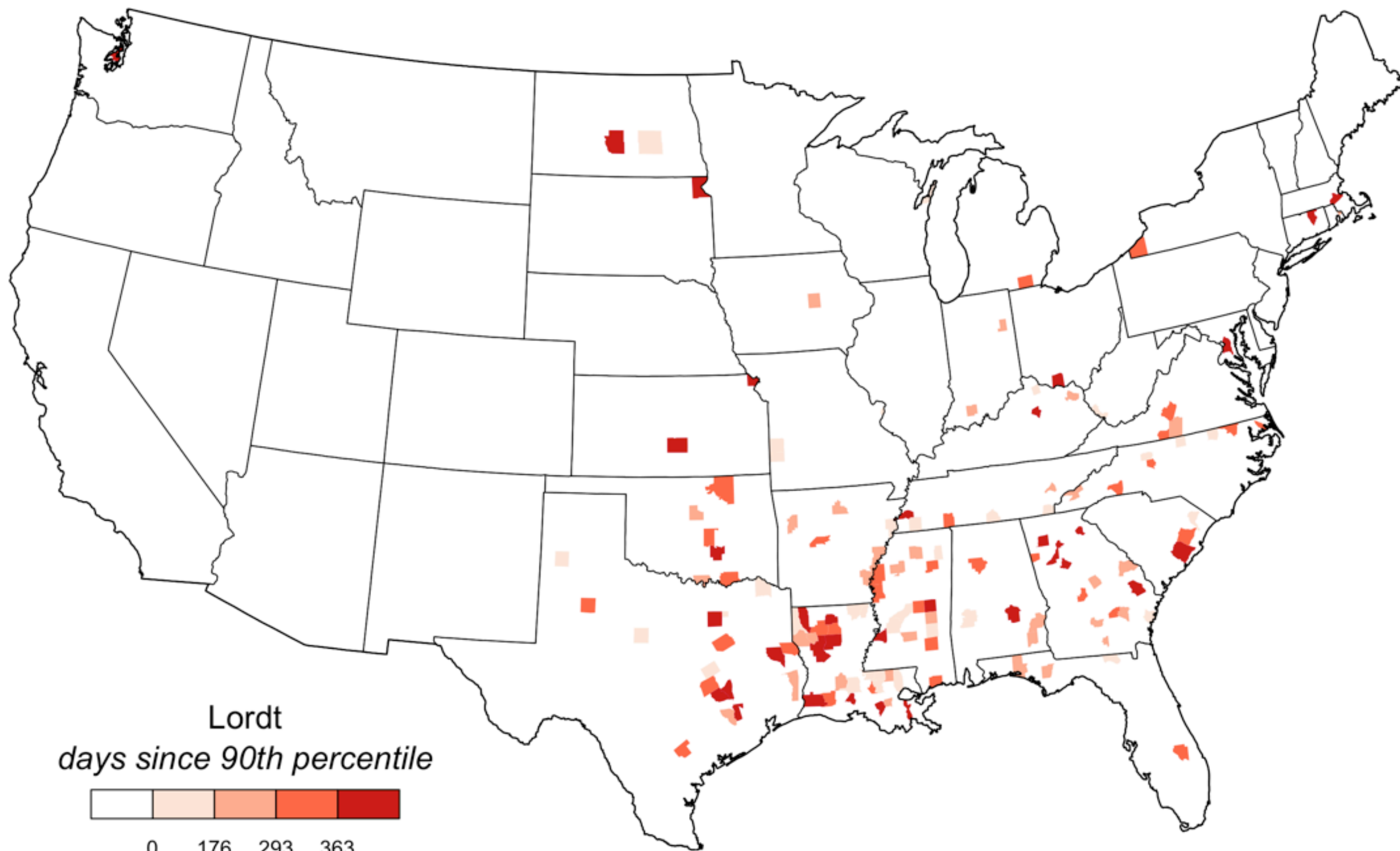
To identify hubs of lexical innovation the 90th percentile relative frequency threshold maps for all 54 words (measured across 3,075 counties) were subjected to a **multivariate spatial analysis**.

1. All maps are subjected to a **Getis-Ord Gi\*** analysis to identify underlying regional signals.
2. The smoothed maps are subjected to a **factor analysis** to identify common regional patterns.

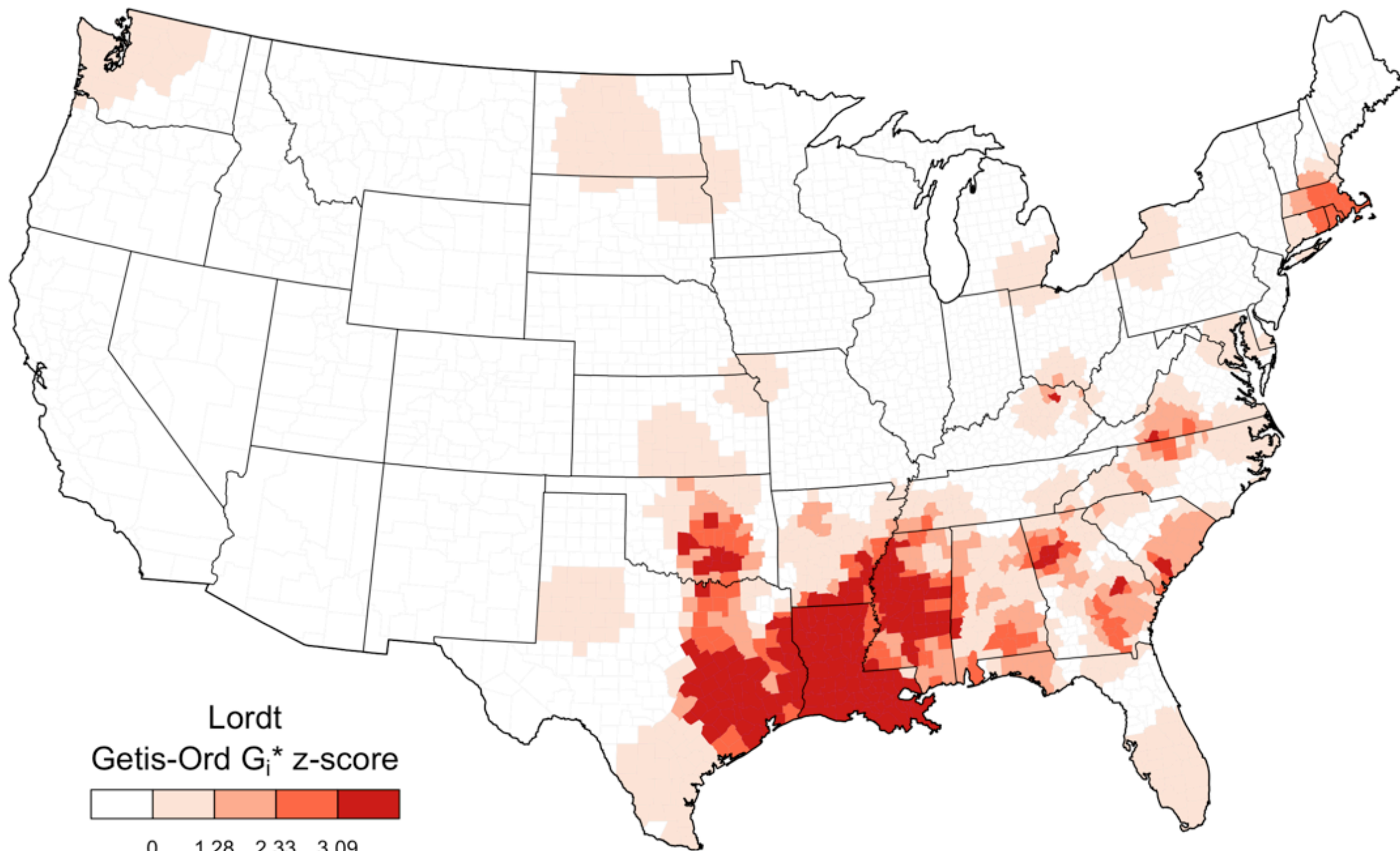
See Grieve. 2016. Regional Variation in Written American English. Cambridge University Press.











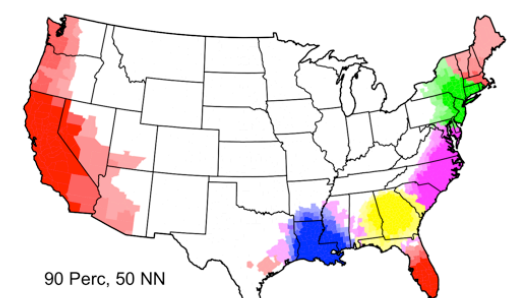
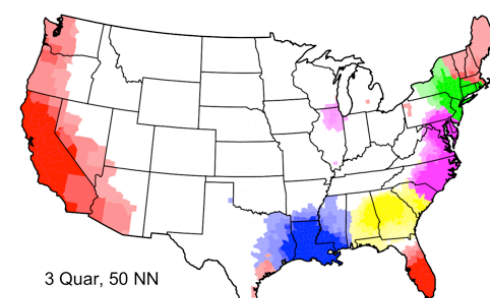
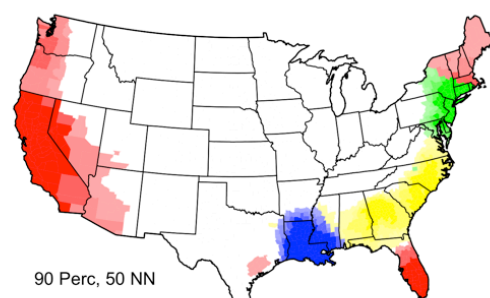
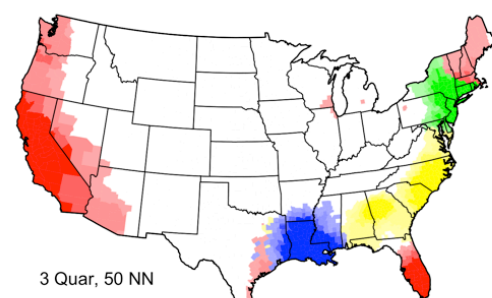
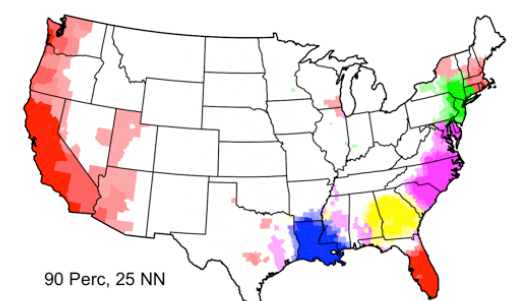
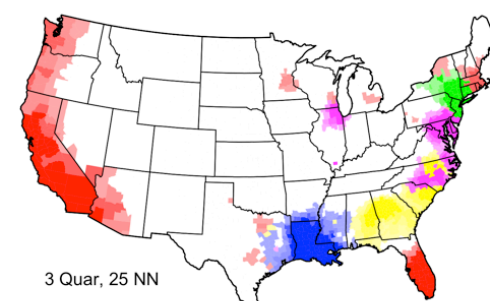
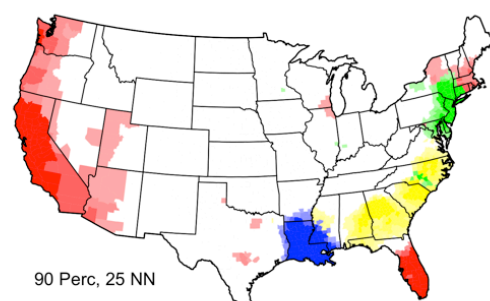
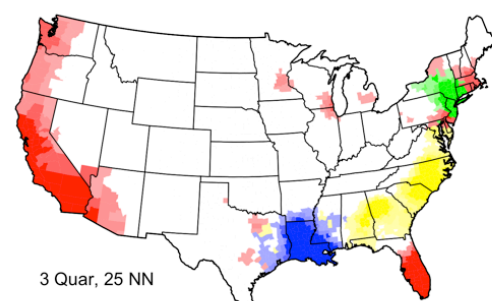
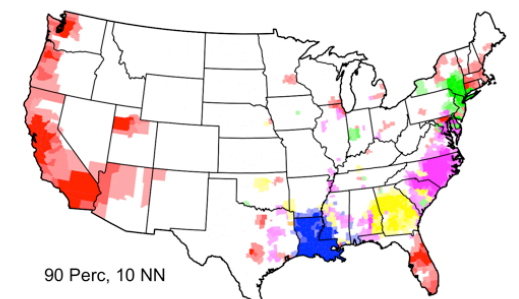
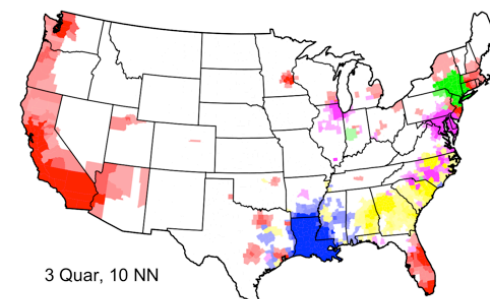
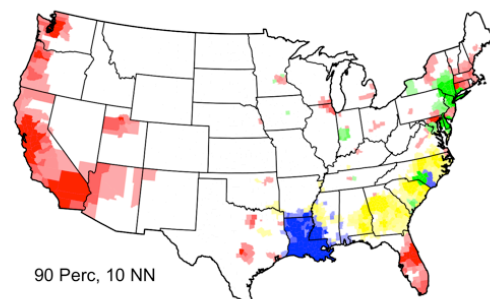
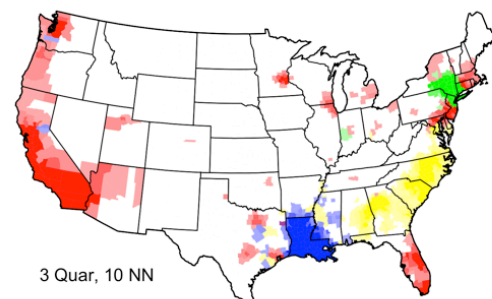
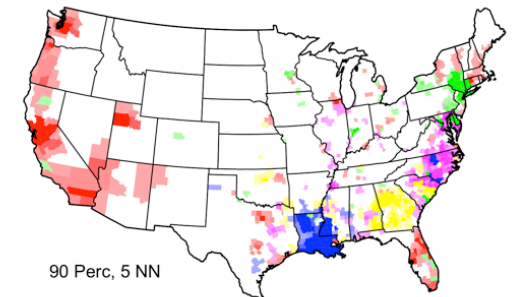
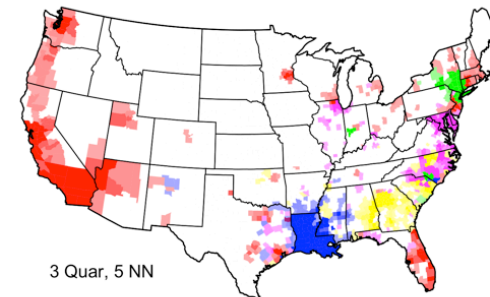
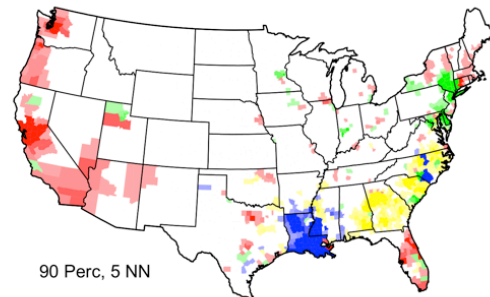
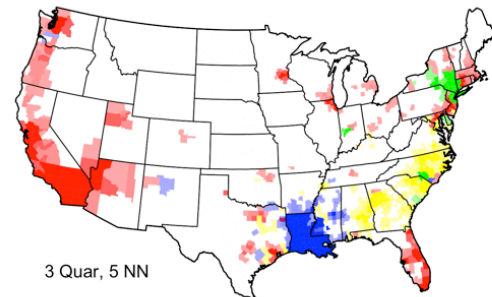
# Mapping Lexical Emergence

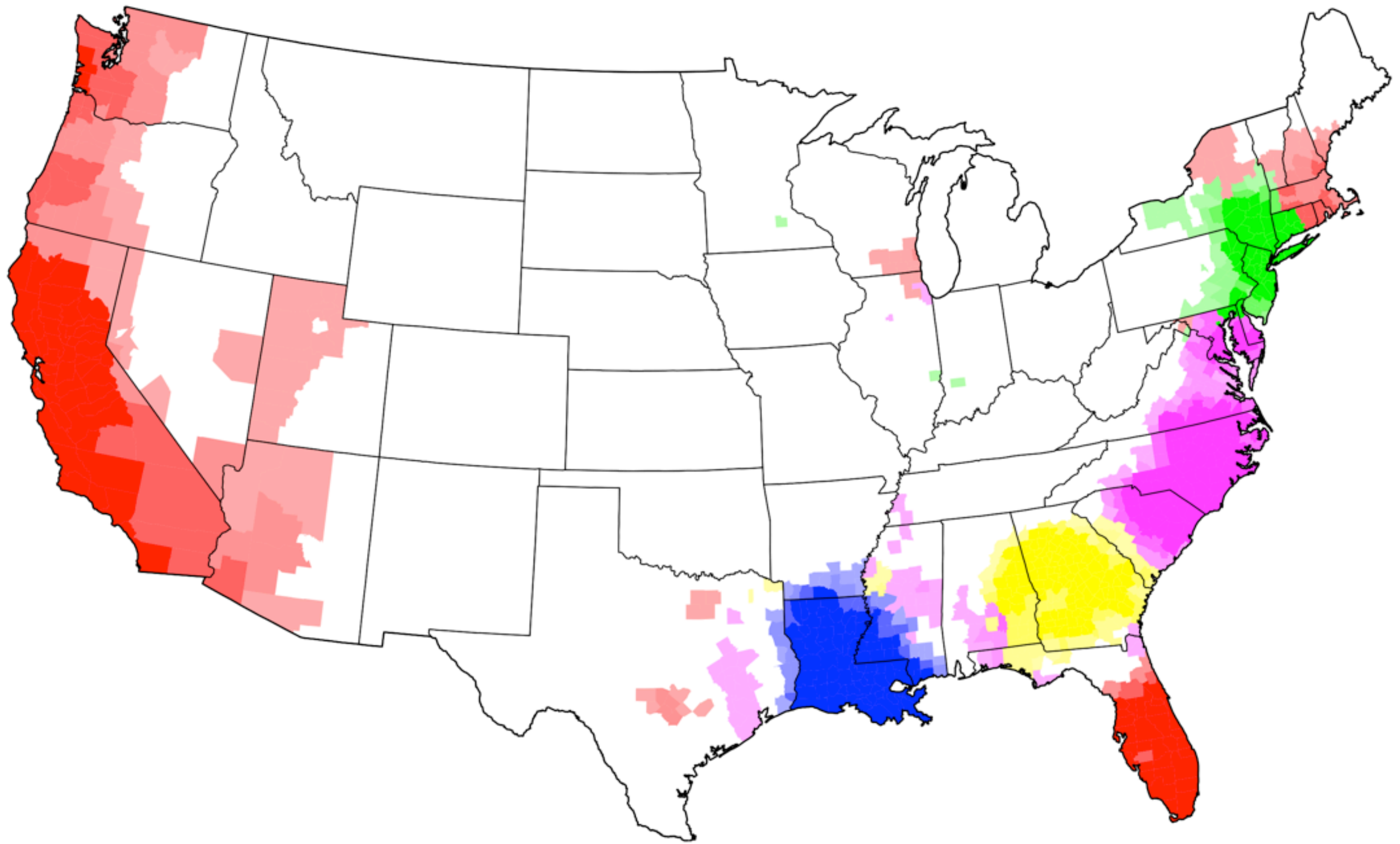
To identify hubs of lexical innovation the 90th percentile relative frequency threshold maps for all 54 words (measured across 3,075 counties) were subjected to a **multivariate spatial analysis**.

1. All maps are subjected to a **Getis-Ord Gi\*** analysis to identify underlying regional signals.
2. The smoothed maps are subjected to a **Factor Analysis** to identify common regional patterns.

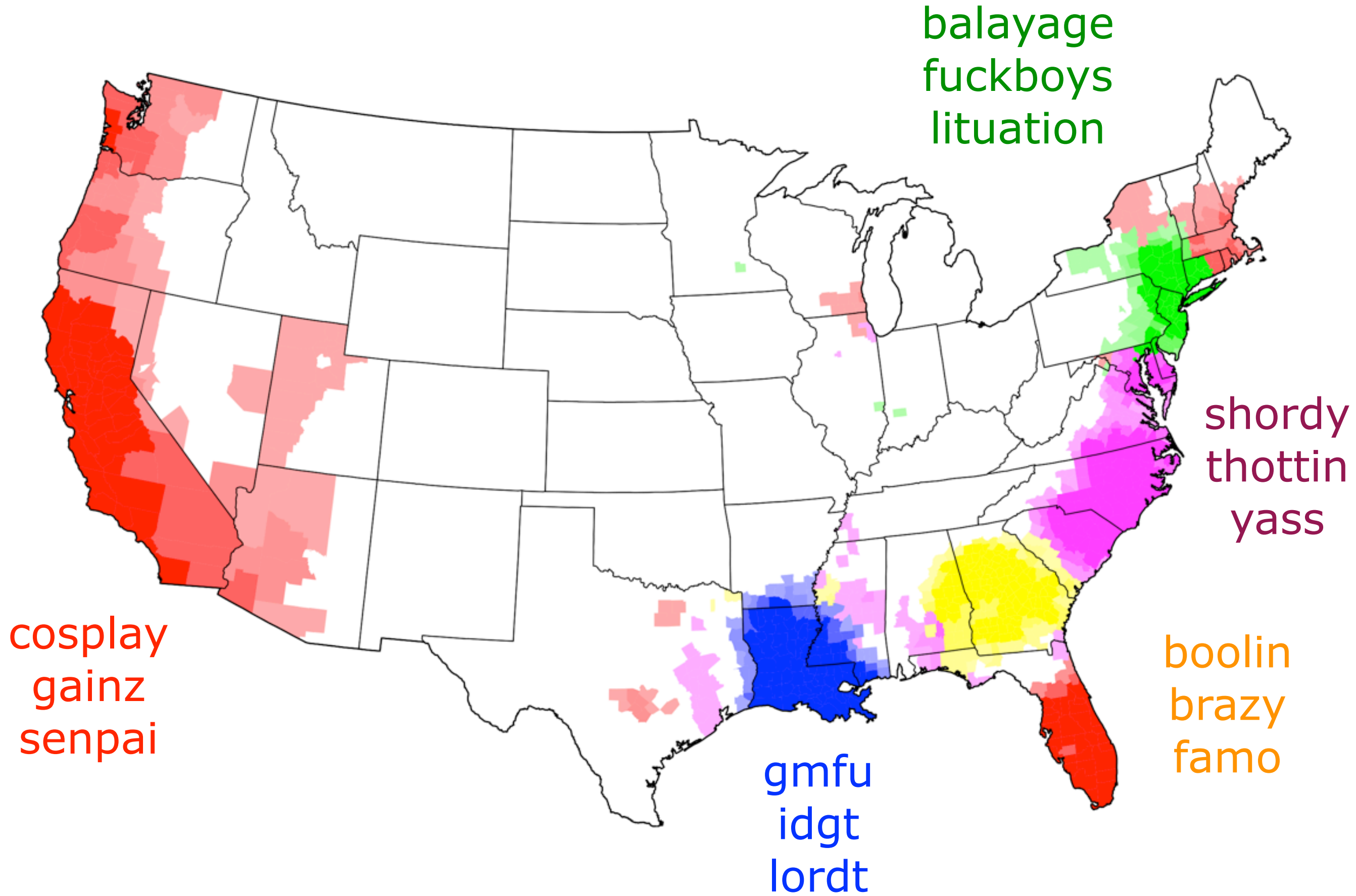
See Grieve. 2016. Regional Variation in Written American English. Cambridge University Press.

# Results Across Parameter Settings







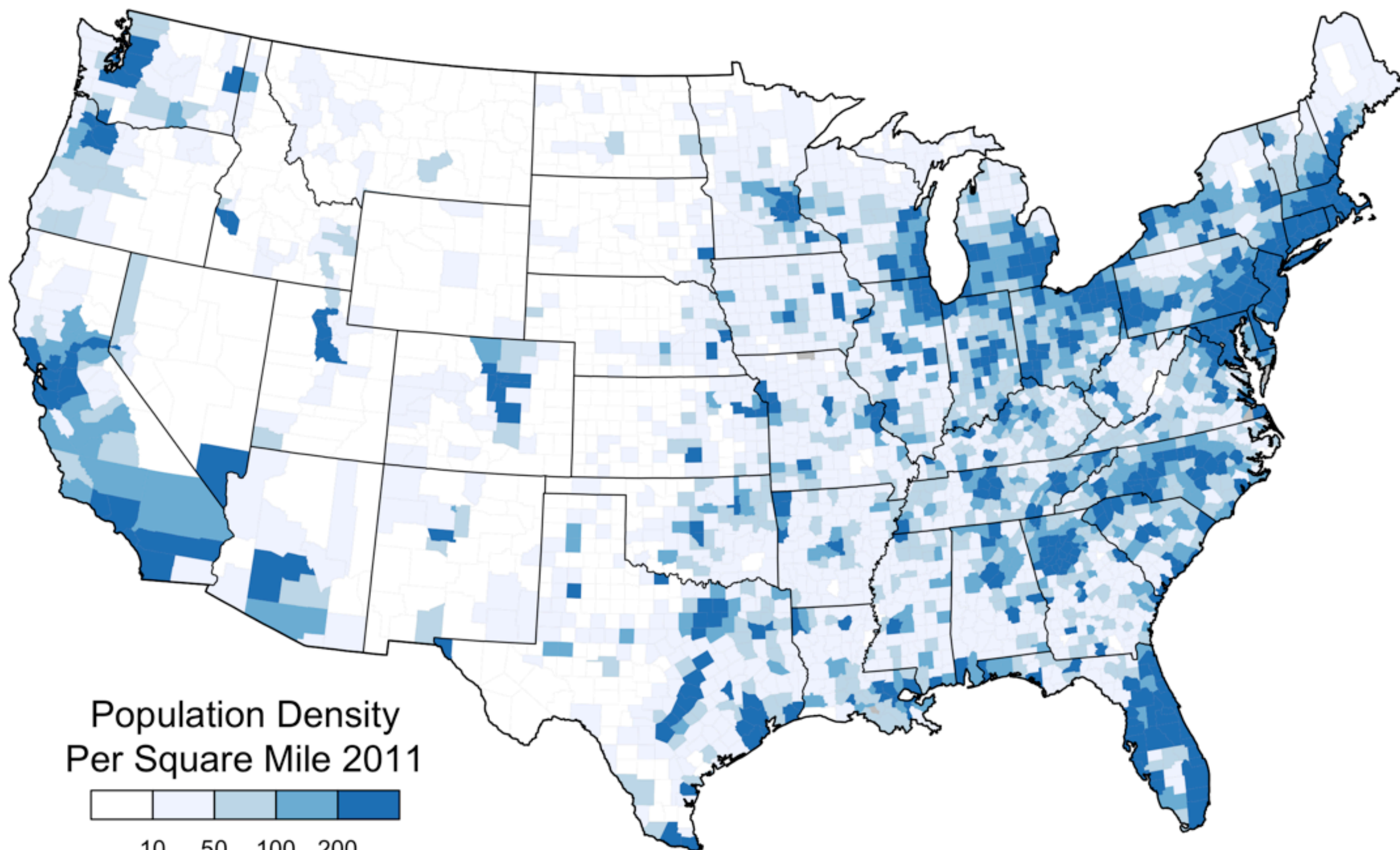


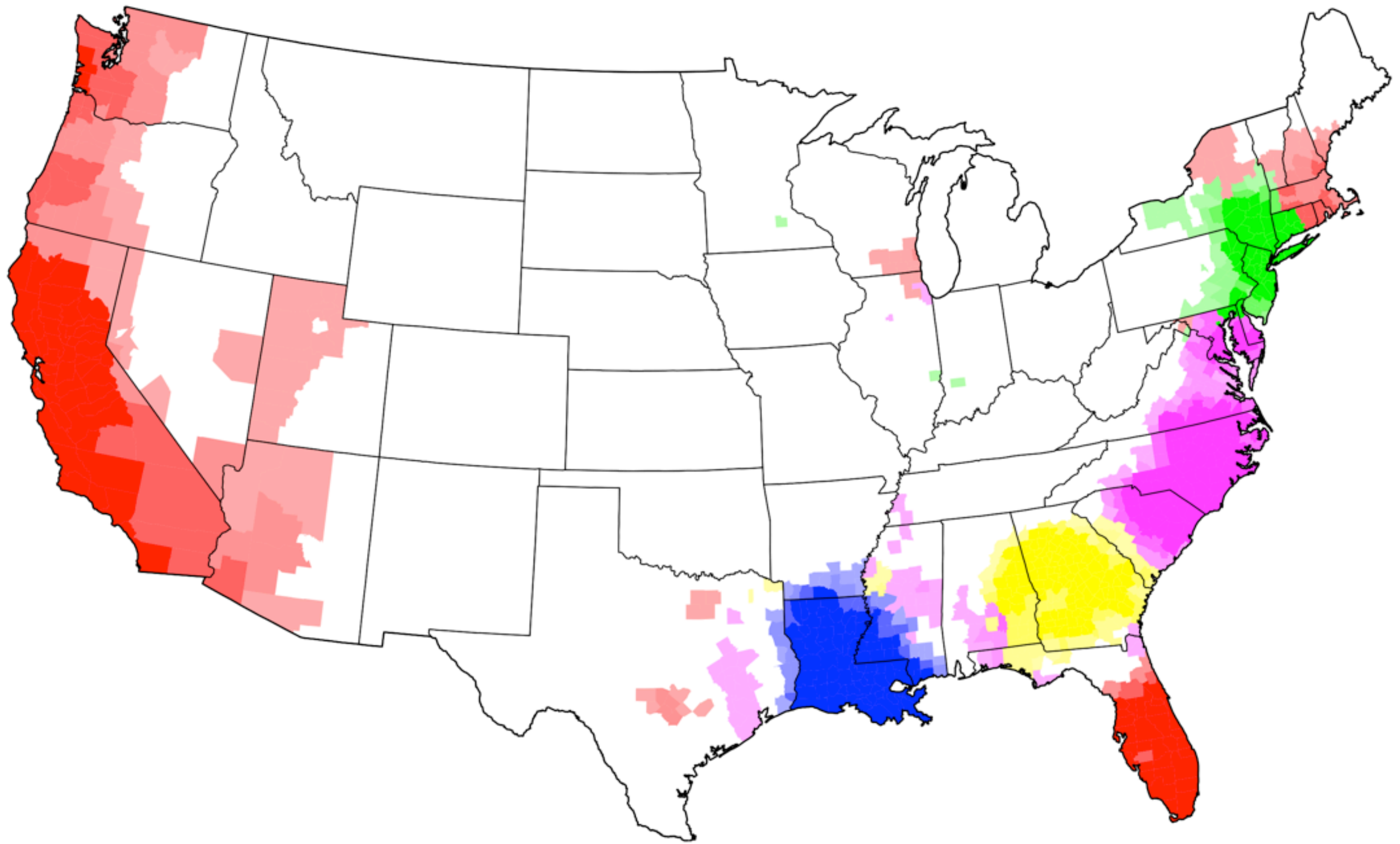
# Conclusions

Five main hubs of lexical innovation were identified on American Twitter from 2014.

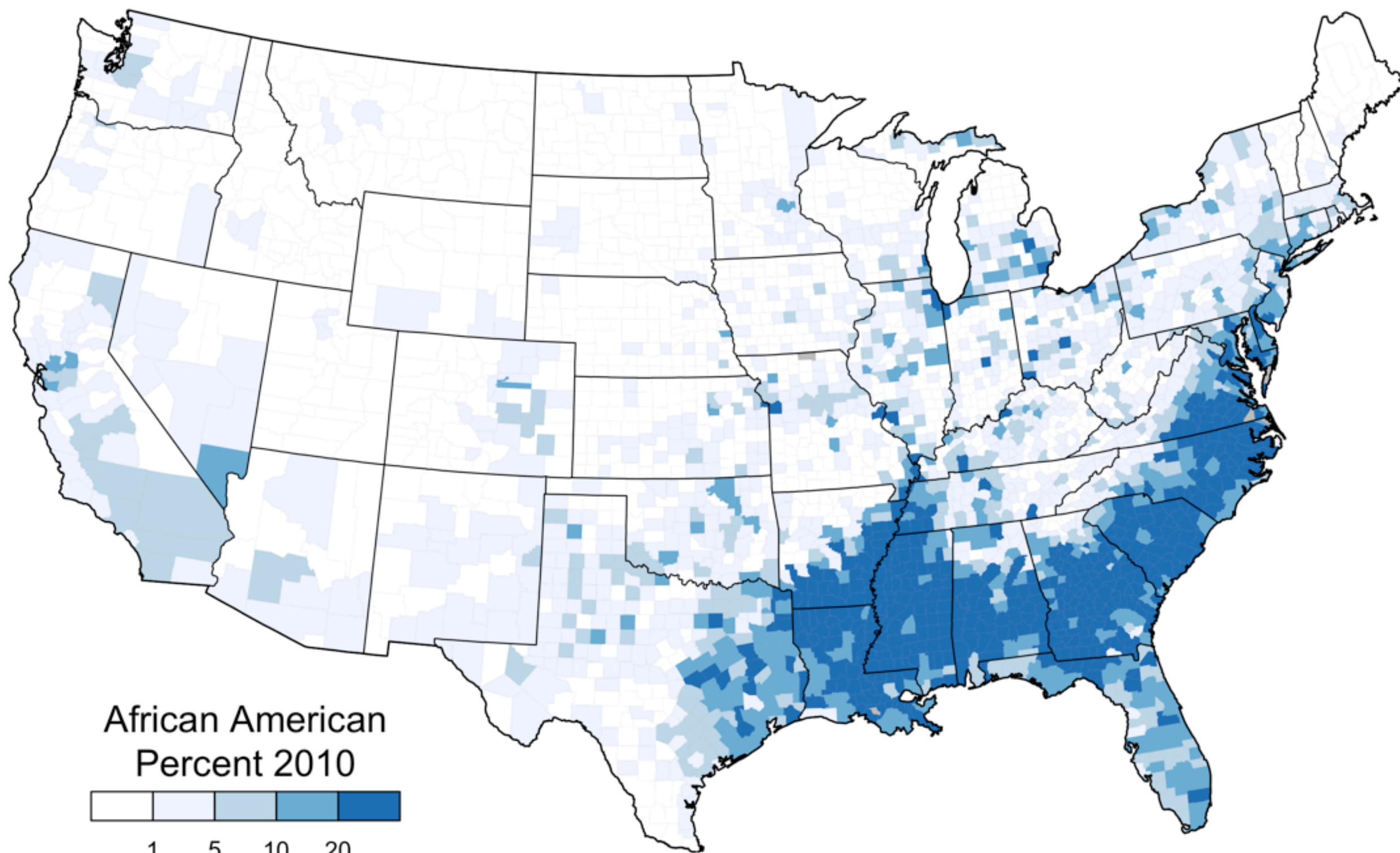
Newly emerging words do not originate from one location or at random, nor do they tend to simply follow patterns of population density.

In addition to California and New York clusters, three distinct hubs were surprisingly found in the Southeastern United States, attesting to the influence and diversity of African American English (and perhaps revealing AAE dialect regions).









# Future Research

How do these patterns change over time and over registers?

How do newly emerging words spread across space (e.g. wave vs. gravity models)?

See Grieve, Nini, Guo. 2016. Mapping lexical innovation on American social media. In review at *Journal of English Linguistics*.

# Identifying and Mapping the Spread of New Words

Jack Grieve

[j.grieve1@aston.ac.uk](mailto:j.grieve1@aston.ac.uk)

[@JWGrieve](#)

Centre for Forensic Linguistics  
Aston University

2 December 2016

BAULT 2016

University of Helsinki