

Scatterchron: Visualizing Language Change in Diachronic Corpora

Code for the interactive visualizations:

<https://github.com/JasonKessler/KeynessToolsTalk/blob/main/BBC%20News%20One%20Year.ipynb>

Repo:

<https://github.com/JasonKessler/KeynessToolsTalk>

Introduction

- Why would we like to visualize how language changes over time?
- **Temporal.** News alerts are ordered by date.
 - Create a timeline of events over a year.
- **Sequential.** Chapters or pages in a novel.
 - Visually represent the plot and presence of characters.
- **Product review data.** Reviews are ordered by star rating.
 - See what features or issues are in different rating regions.

Outline

- Related work in visualizing language change over time
- Interactive text visualizations with **Scattertext**
- Binary difference: term frequencies before and after time-change
- Most associated terms in each time step.
 - Clustering time steps to make them more manageable
- Detecting terms that clump together at different time steps or are more dispersed
- Average time position
- Putting it all together: **Scatterchron**

Related Work

Parallel Tag Clouds

Key terms are vertically arranged at each time step

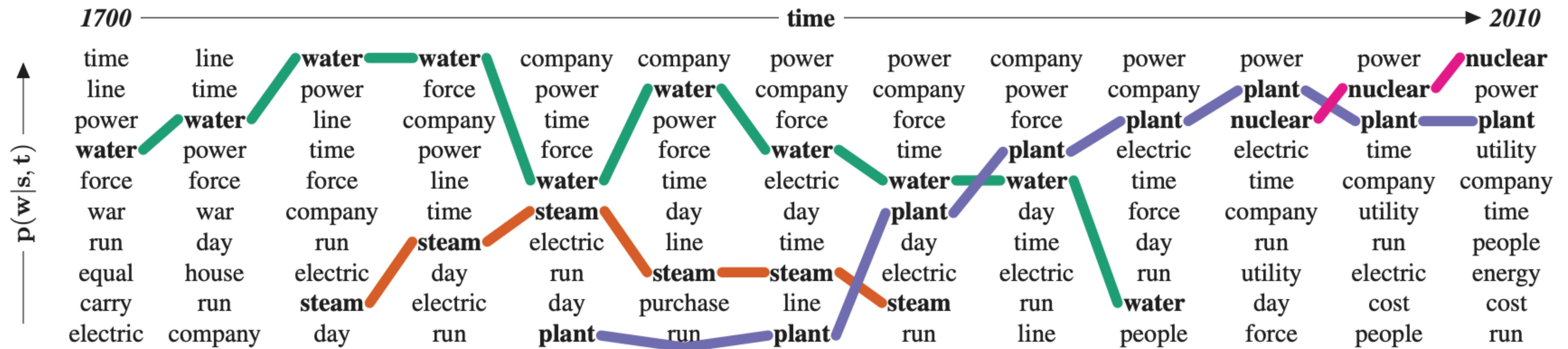
Size indicates time-step association (G^2)

Order is alphabetical

Connectors to the same words appearing in different time steps



Related Work

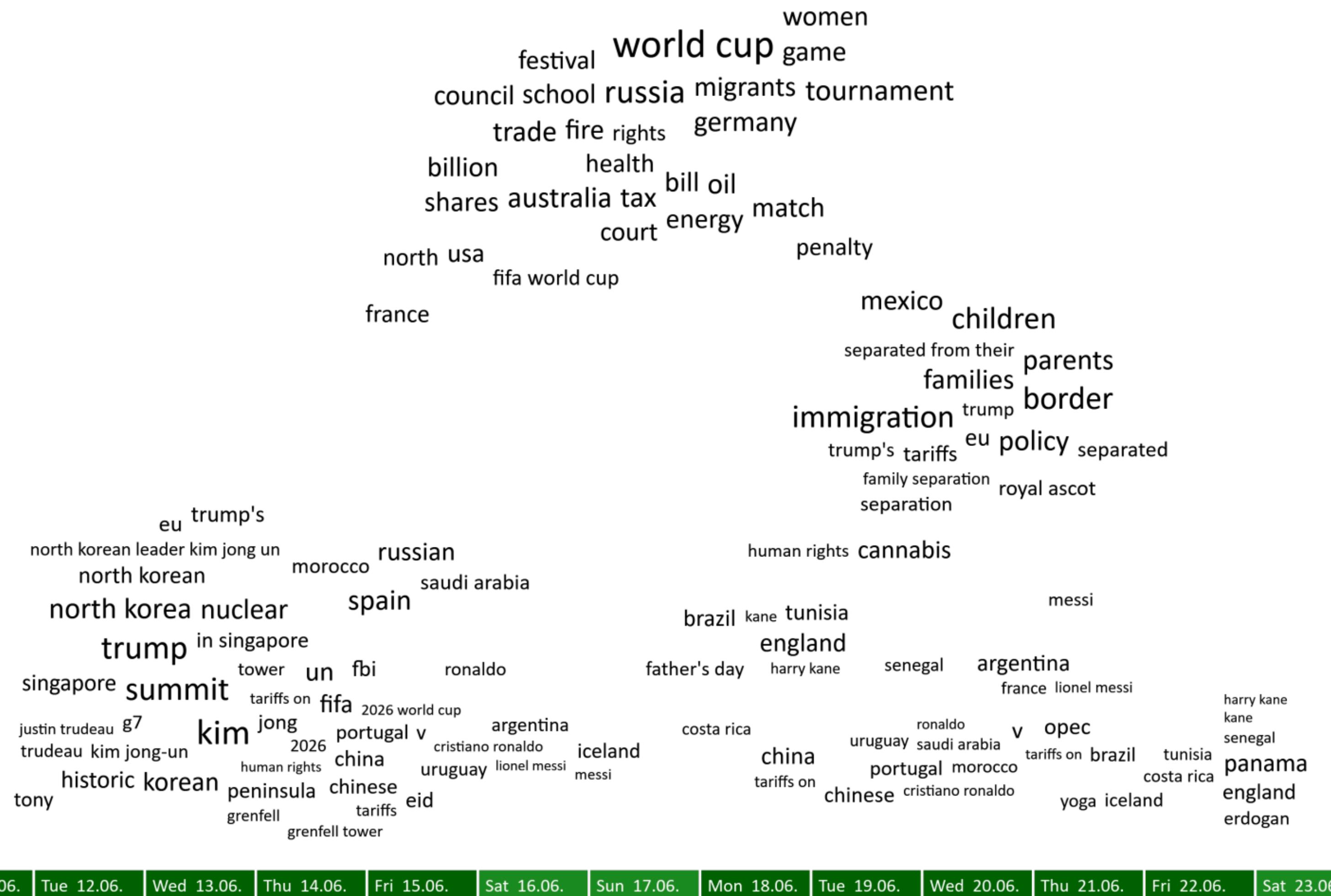


- “Words most associated with the energy-sense of the word ‘power’”
- Influenced by Parallel Tag Clouds
- More readable layout
- Ordered by association, not alphabetically

Related Work

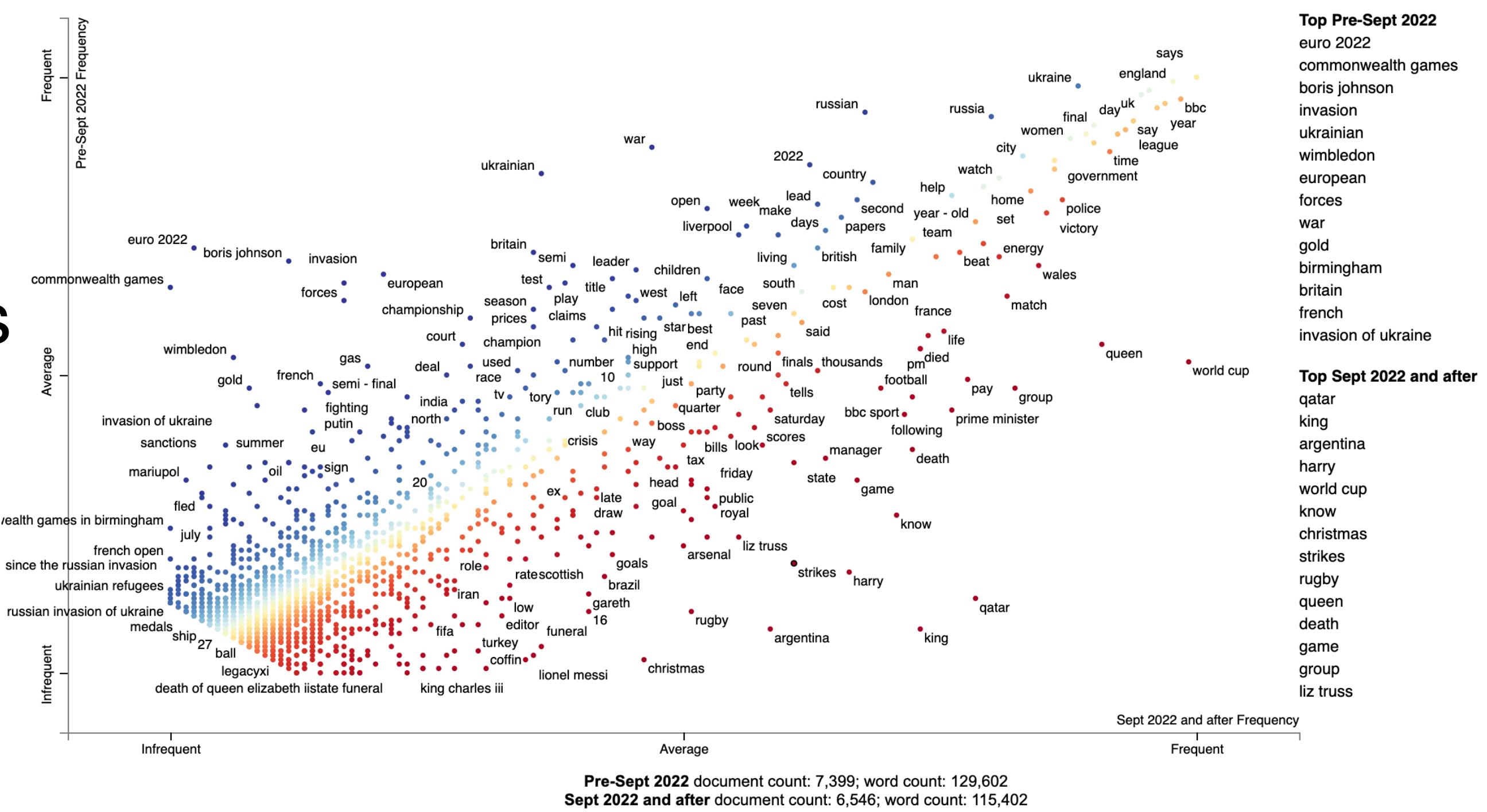
PyramidTags

- N-grams are displayed in a “cloud” above a timeline.
- Horizontal position: roughly mean point in the timeline
- Vertical: higher up indicates longer time-span
- Arranged s.t. words, even when not in phrases, tend to appear in sequence.
 - Similar terms cluster together



Time-step divided plot

- BBC RSS feed corpus (Feb 2022 - Feb 2023)¹
 - 13,945 alerts; 361 days; 281,589 words
 - Label each alert by before Sept 2022, Sept 2022 and after
 - Plot rank frequency of features on x-axis in <9/22 and of >9/22 on y-axis
 - Uses Scattertext (Kessler 2017)
 - Color by difference in ranked frequency



https://jasonkessler.github.io/keynesstalk/bbc_divided.html

Parallel tag clouds

Displaying time-step information

2022-02-22	2022-02-24	2022-02-26	2022-02-28	2022-03-01	2022-03-02	2022-03-03	2022-03-04	2022-03-05	2022-03-06	2022-03-07
restrictions	covid	watched	spread	rules	care	crime	russia	russia	war	ukraine
announced	people	virus	virus	face	fourth	really	ukraine	resistance	ukrainian	russian
covid	england	needs	key	uk	covid	tackle	experts	happen	manchester	woman
end	self	covid	sector	pitch	workers	little	reasons	warned	investigated	crew
government	tests	living	unlikely	carrying	people	economic	coronavirus in the uk	president	putin	rival
england	longer	uk	long - term	alice	england	putin	explore the data on coronavirus	answers	city	irish
striking images from our readers	free	canadian	beach	peace	thrown	late	relatives	kevin	parents	19
ms	government	offers	mercedes	seal	broadcaster	money	offensive	missile	russian	shock
unusual	million people have fled	holding	threatening	easy	waste	ukraine says	items	jeremy	national	leave
readers around the world	followed	people have fled their homes	racing	golden	alice	trains	feel	really	qatar	year - old

- “covid” is moused-over
- Jenson Shannon Divergence is used as time-step association metric
- Symmetric version of KL divergence
- This approximates Cohen’s d well and is faster to calculate

Parallel tag clouds

Displaying time-step information

2022-02-22	2022-02-24	2022-02-26	2022-02-28	2022-03-01	2022-03-02	2022-03-03	2022-03-04	2022-03-05	2022-03-06	2022-03-07
restrictions	covid	watched	spread	rules	care	crime	russia	russia	war	ukraine
announced	people	virus	virus	face	fourth	really	ukraine	resistance	ukrainian	russian
covid	england	needs	key	uk	covid	tackle	experts	happen	manchester	woman
end	self	covid	sector		workers	little	reasons	warned	investigated	crew
government	tests	living	unlikely	carrying	people	economic	coronavirus in the uk	president	putin	rival
england	longer	uk	long - term	alice	england	putin	explore the data on coronavirus	answers	city	irish
striking images from our readers	free	canadian	beach	peace	thrown	late	relatives	kevin	parents	19
ms	government	offers	mercedes	seal	broadcaster	money	offensive	missile	russian	shock
unusual	million people have fled	holding	threatening	easy	waste	ukraine says	items	jeremy	national	leave
readers around the world	followed	people have fled their homes	racing	golden	alice	trains	feel	really	qatar	year - old

- Size: frequency in class
- Ordered by delta JSD; tends to choose frequently but associated terms
- Problem: there are 361 time steps; we can't show them all without a lot of horizontal scrolling

Clustered parallel tag clouds

Displaying time-step information compactly

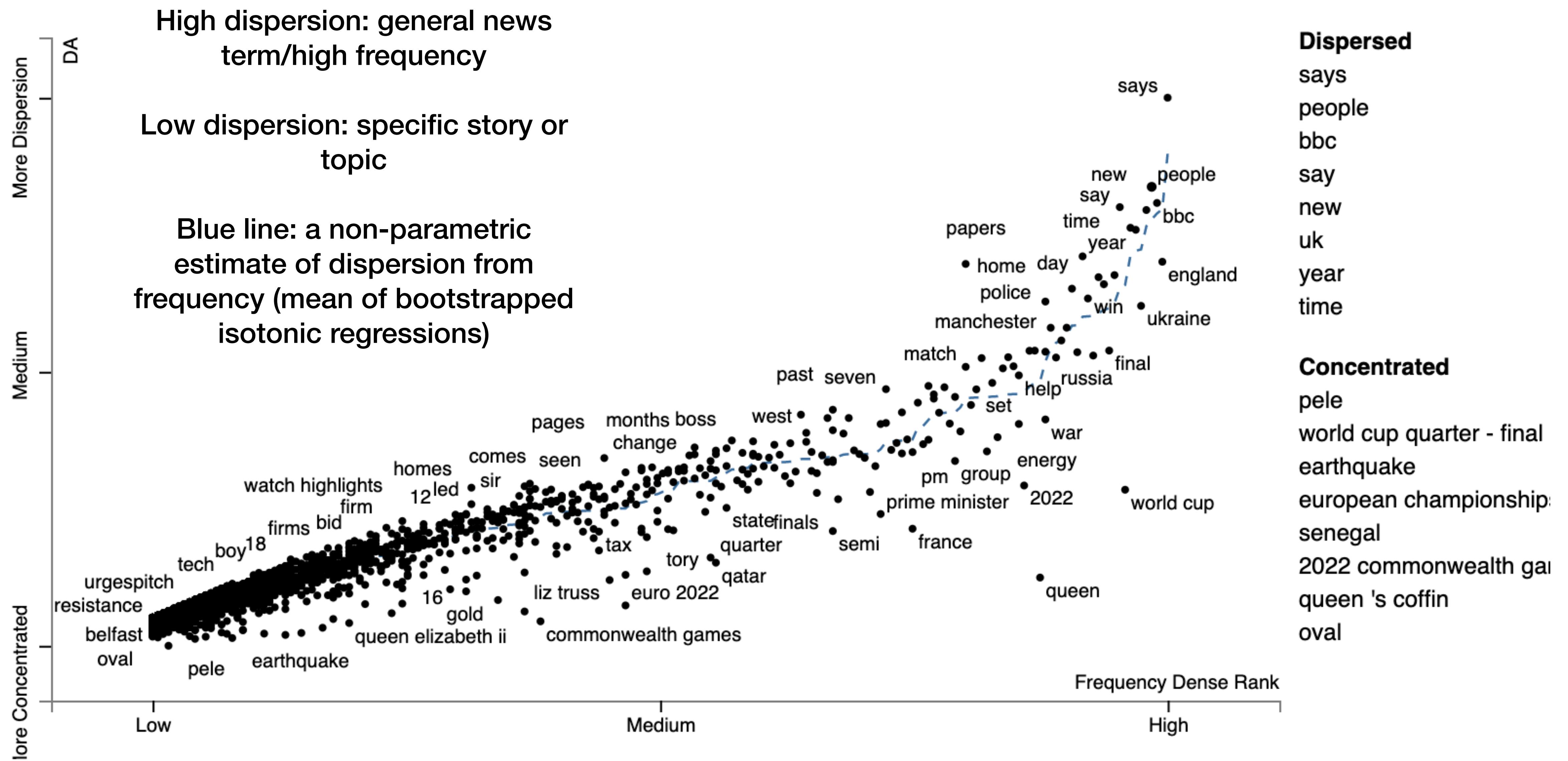
2022-02-22 to 2022-07-30	2022-07-31 to 2022-08-16	2022-08-17 to 2022-08-25	2022-08-26 to 2022-08-27	2022-08-28 to 2022-09-15	2022-09-16 to 2023-02-20
ukraine	commonwealth games	european	energy	queen	world cup
russian	gold	championships	ben	king	bbc
russia	birmingham	students	bills	coffin	qatar
war	england	britain	ben stokes	charles	argentina
ukrainian	commonwealth games in birmingham	euopean championships	days	queen elizabeth ii	christmas
boris johnson	wins	results	way	liz truss	group
euro 2022	2022	south	title	death	strikes
invasion	women	manchester	people	king charles iii	harry
city	win	great britain	rivals	state	rugby
wimbledon	watch	test	man	queen 's coffin	game

- Solution: create clusters of time steps
- Among terms displayed, compute Spearman's correlation coefficient on score ranks for adjacent time steps
- Select the K (above K=5) lowest correlations for breakpoints

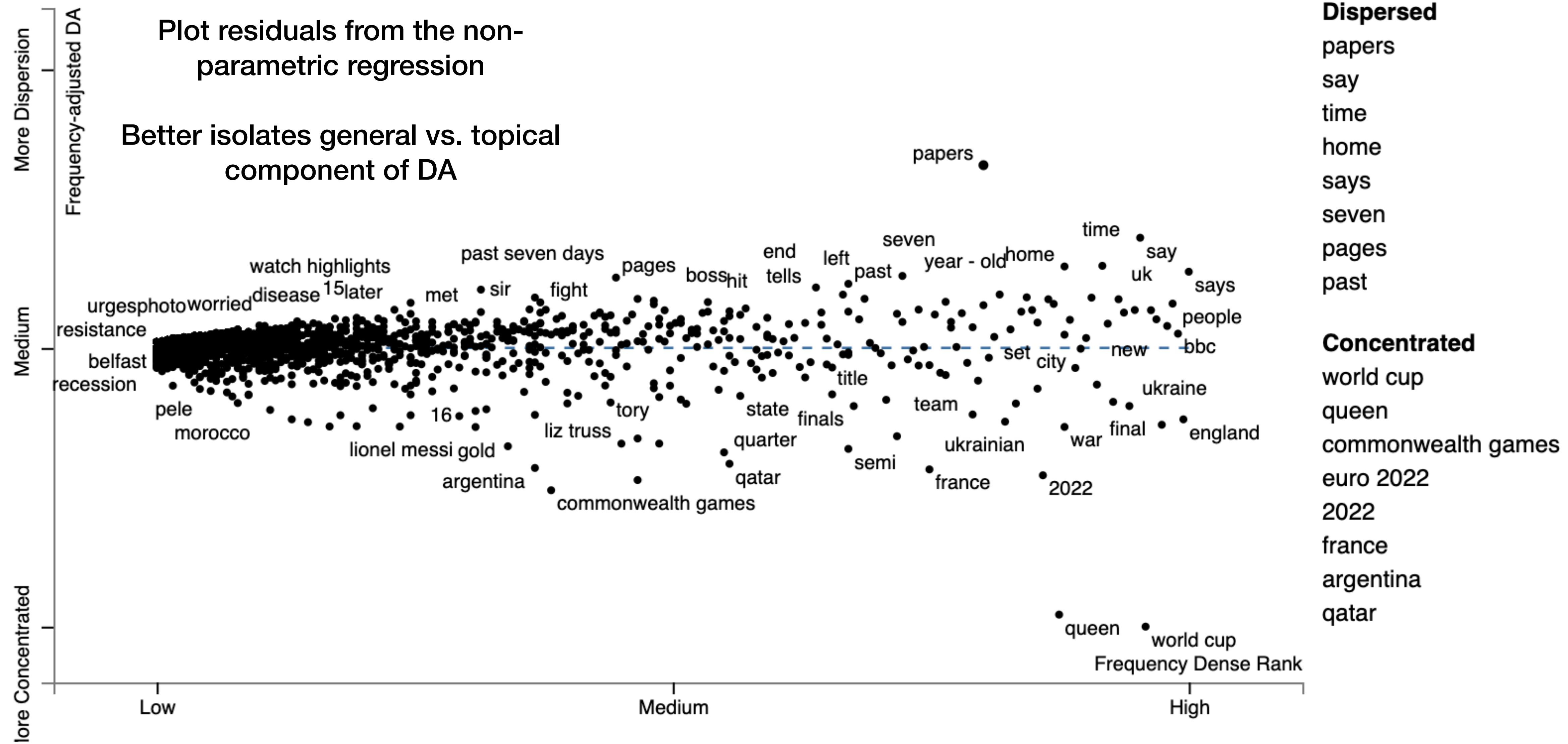
Plotting change over time: dispersion

- Dispersion: metric which indicates if terms appear at the same frequency distributed across corpus parts (time steps) or if they occur with disproportionate frequency in some rather than others
- Dispersion metric: DA (Burch et al., 2017)
- Mean absolute pairwise difference between part frequencies
- Ranges from 0 (a term only appears in one part) to 1 (same frequency in each part)
 - 0.5: the average pairwise difference is the mean frequency
- It takes $O(\text{time-steps}^2)$ to run. Gries' DP, etc. are faster but lower quality.

Plotting change over time: dispersion



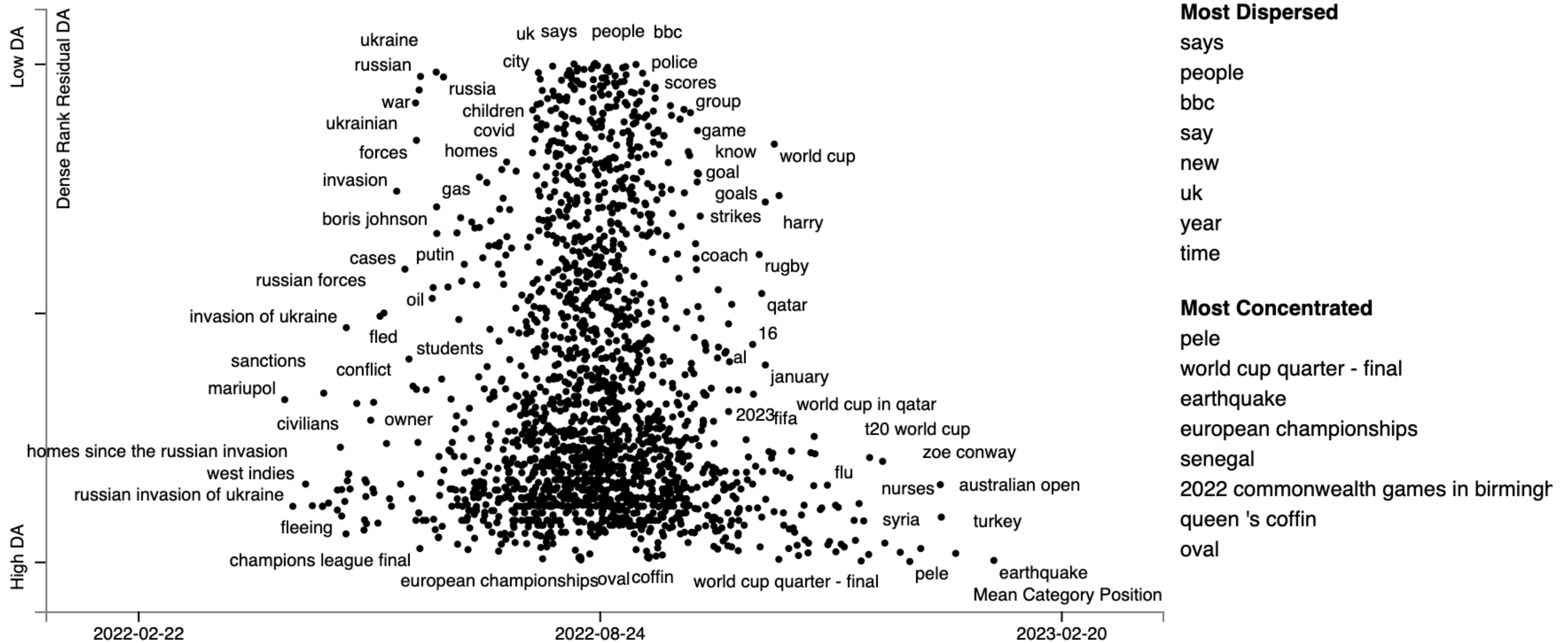
Residual DA: frequency normalization



Incorporating Average Position

- Residual dispersion: what's event-related or bursty vs. what's more general
- Missing: where on the timeline did the event occur
- Solution: center terms on their average time steps, similar to PyramidTags

Incorporating Average Position

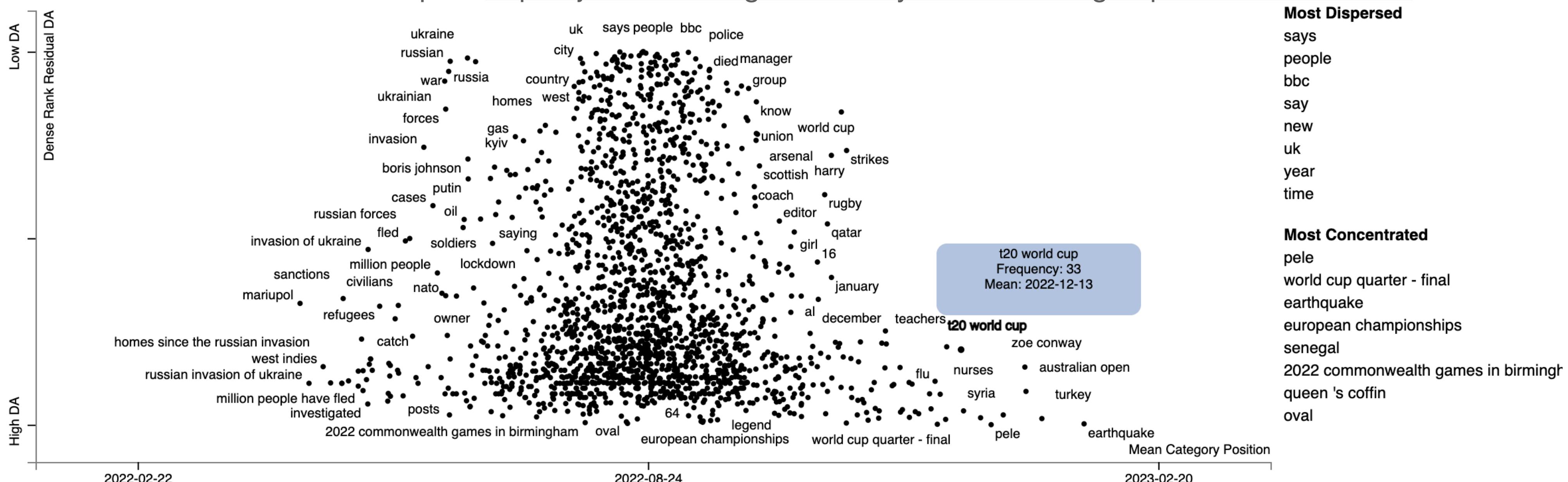


Scatterchron: stacking the visualizations

2022-02-22 to 2022-07-30 2022-07-31 to 2022-08-16 2022-08-17 to 2022-08-25 2022-08-26 to 2022-08-27 2022-08-28 to 2022-09-15 2022-09-16 to 2023-02-20

ukraine	commonwealth games	european	energy	queen	world cup
russian	gold	championships	ben	king	bbc
russia	birmingham	students	bills	coffin	qatar
war	england	britain	ben stokes	charles	argentina
ukrainian	commonwealth games in birmingham	euopean championships	days	queen elizabeth ii	christmas
boris johnson	wins	results	way	liz truss	group
euro 2022	2022	south	title	death	strikes
invasion	women	manchester	people	king charles iii	harry
city	win	great britain	rivals	state	rugby
wimbledon	watch	test	man	queen 's coffin	game

Interactive plot: https://jasonkessler.github.io/keynesstalk/bbc_grouped_time_table.html



Future work

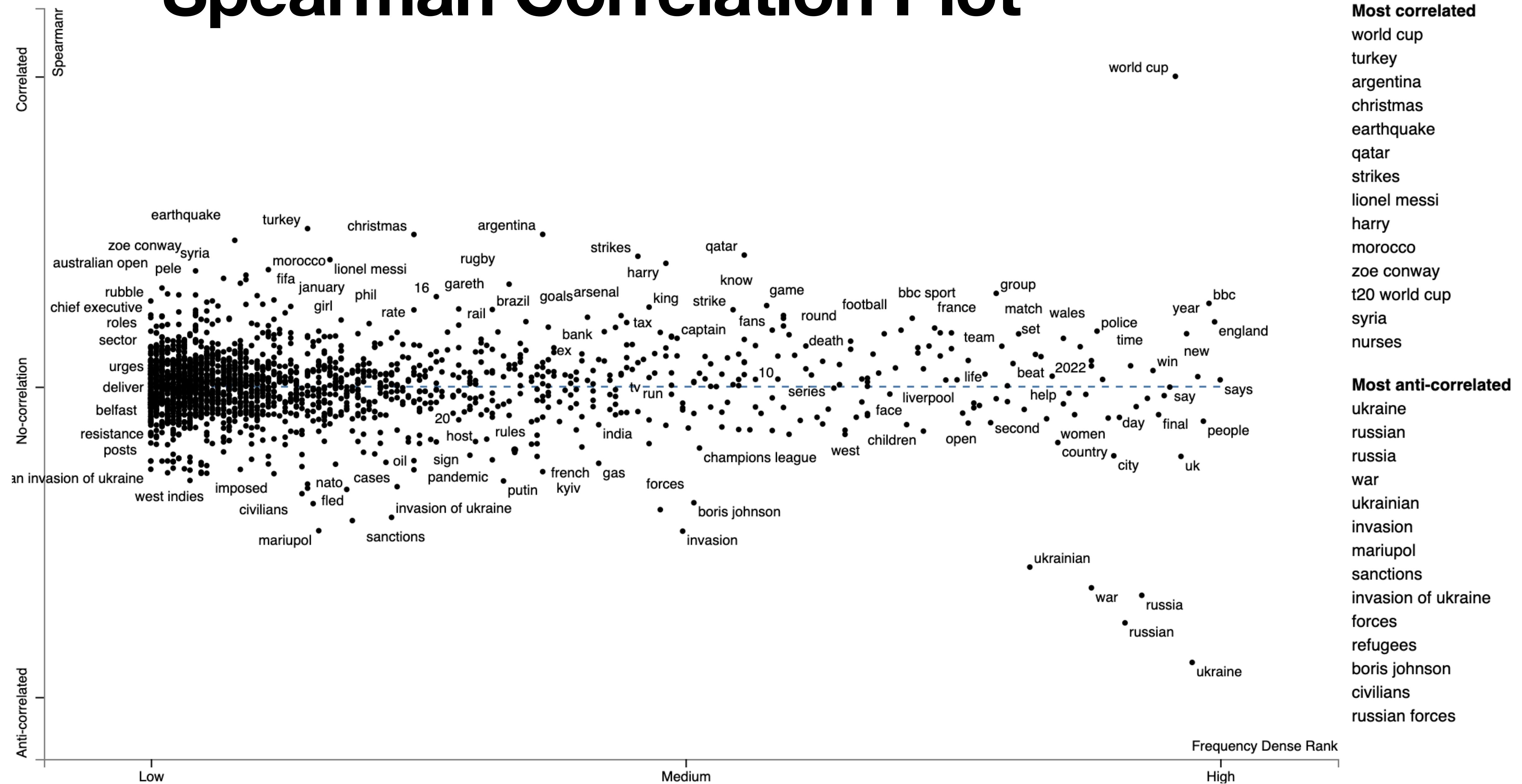
- Position-sensitive dispersion metrics
- Ways of factoring out or calling out periodicity
- Chart interactivity to show a time step-specific plot

Thank you! Any questions?

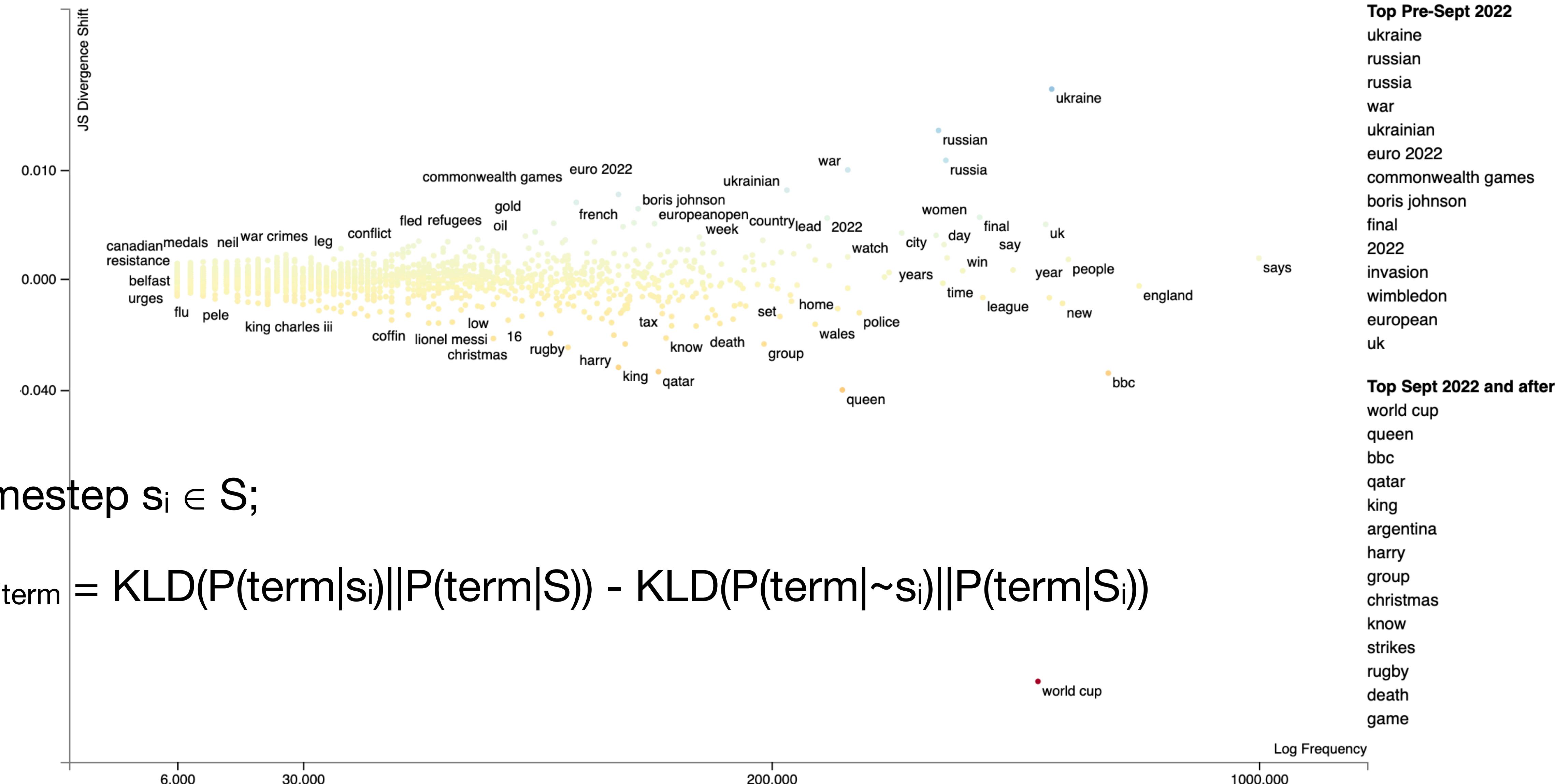
Phrase Selection

- Take 1-5 grams, keeping those which
 - start and end with a content word
 - appear at least three times
 - keep the top 2000 highest NPMI phrases per phrase length in tokens
 - does not makeup at least 60% of a longer (token-wise) selected phrase
- Select 2000 phrases for visualization by round-robin selection of the highest JSD association for each year
 - Jenson-Shanon Divergence tends to favor topical but high-frequency terms

Spearman Correlation Plot

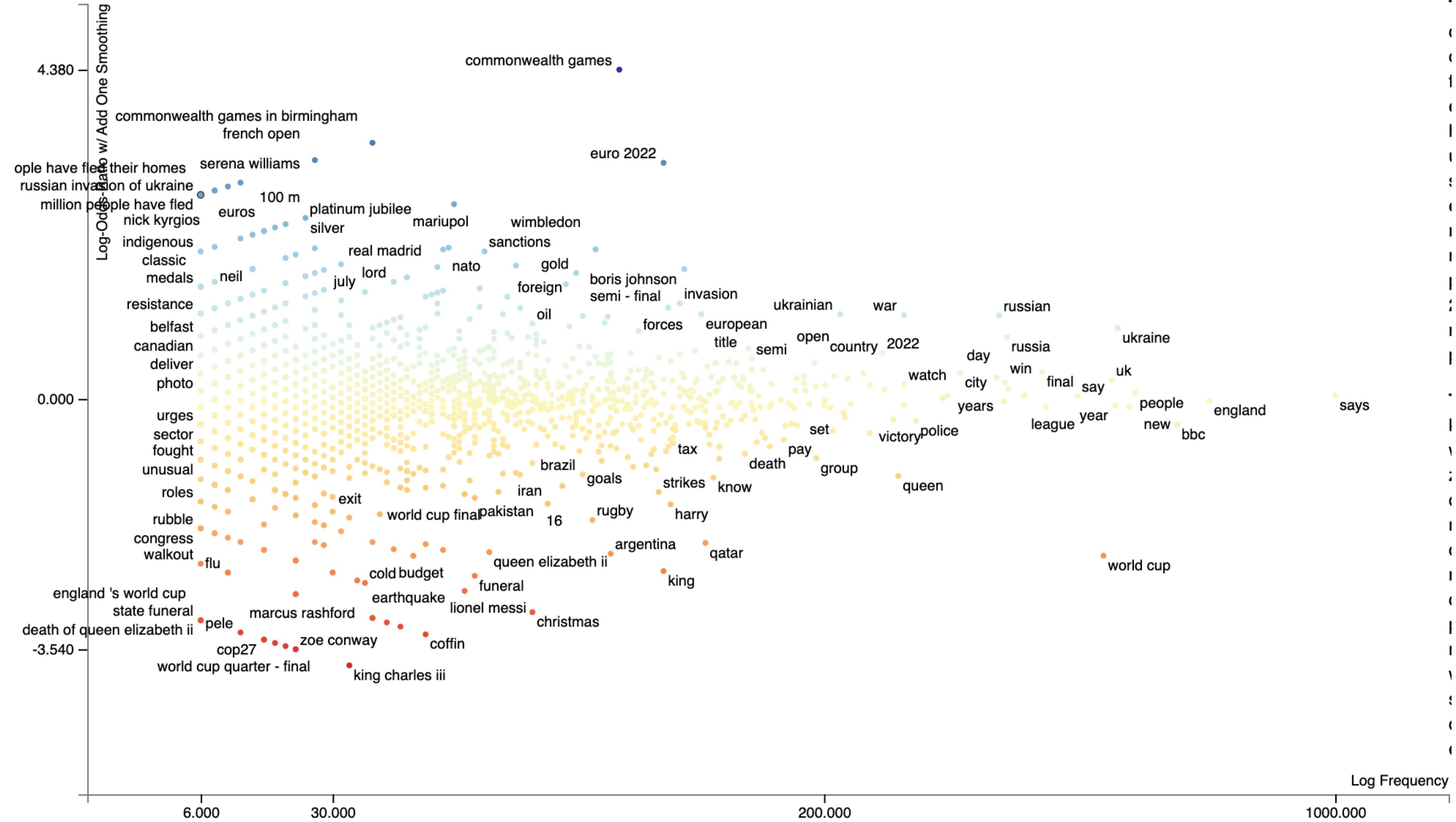


Delta JS Divergence



- For a timestep $s_i \in S$;
- $\Delta JSD_{term} = KLD(P(term|s_i) || P(term|S)) - KLD(P(term|\sim s_i) || P(term|S_i))$

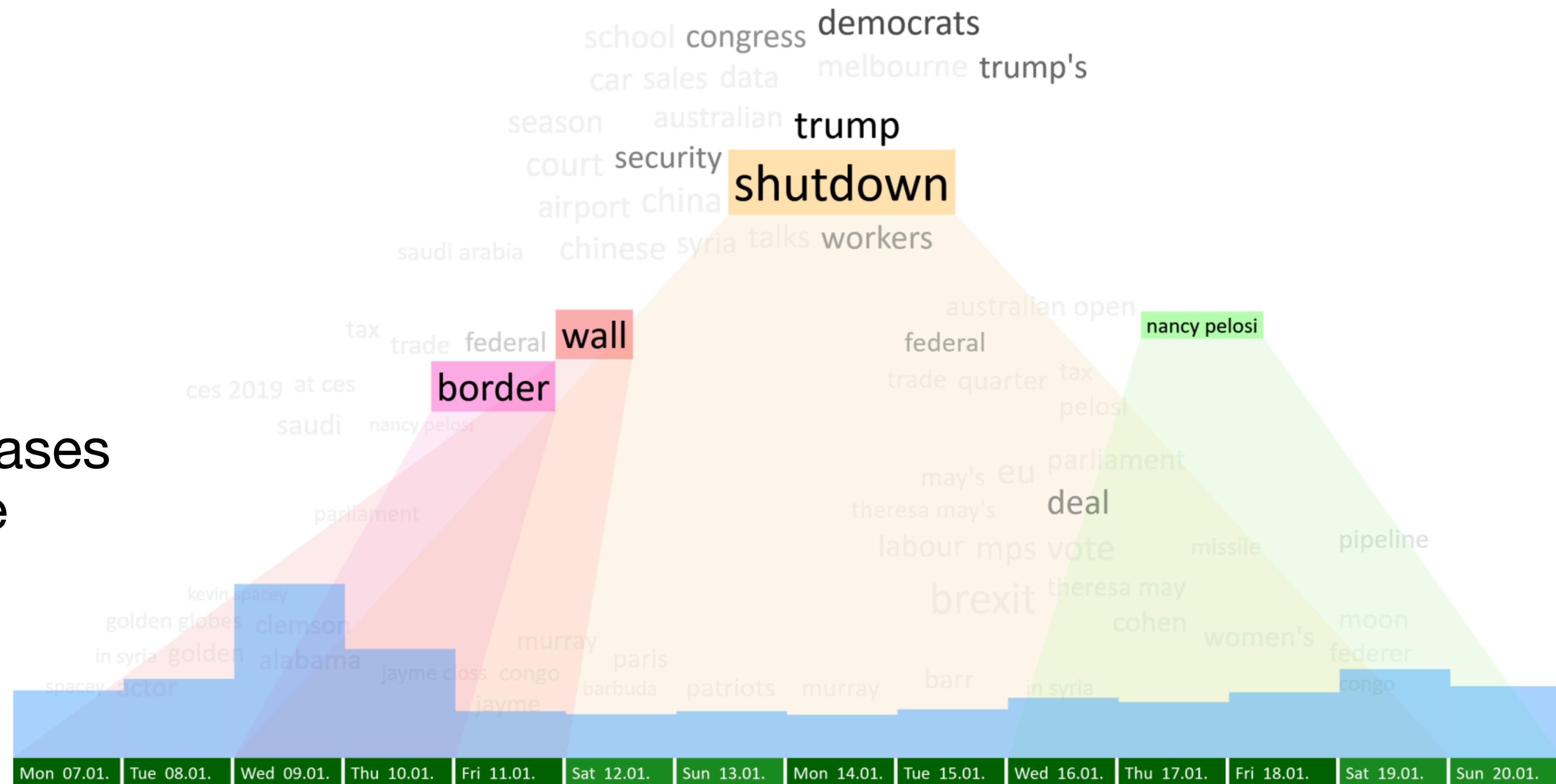
C/f log odds ratio



Related Work

PyramidTags

- Hovering over a phase shows its timespan
- Related phrases are highlighted proportional relatedness
 - Relatedness: two phrases tend to occur in same documents in close proximity



Dispersion

- Some terms are highly specific to particular type of movies. Tend to appear clumped in reviews for particular movies.
 - From low to high frequency: movie names, actors, and genres
- Some terms tend to appear at a roughly consistent rate in different sections of a corpus.
 - Are not specific to a particular type of movie.
 - The words “audience” or “quite” in movie reviews.
 - Function words (we’d think)

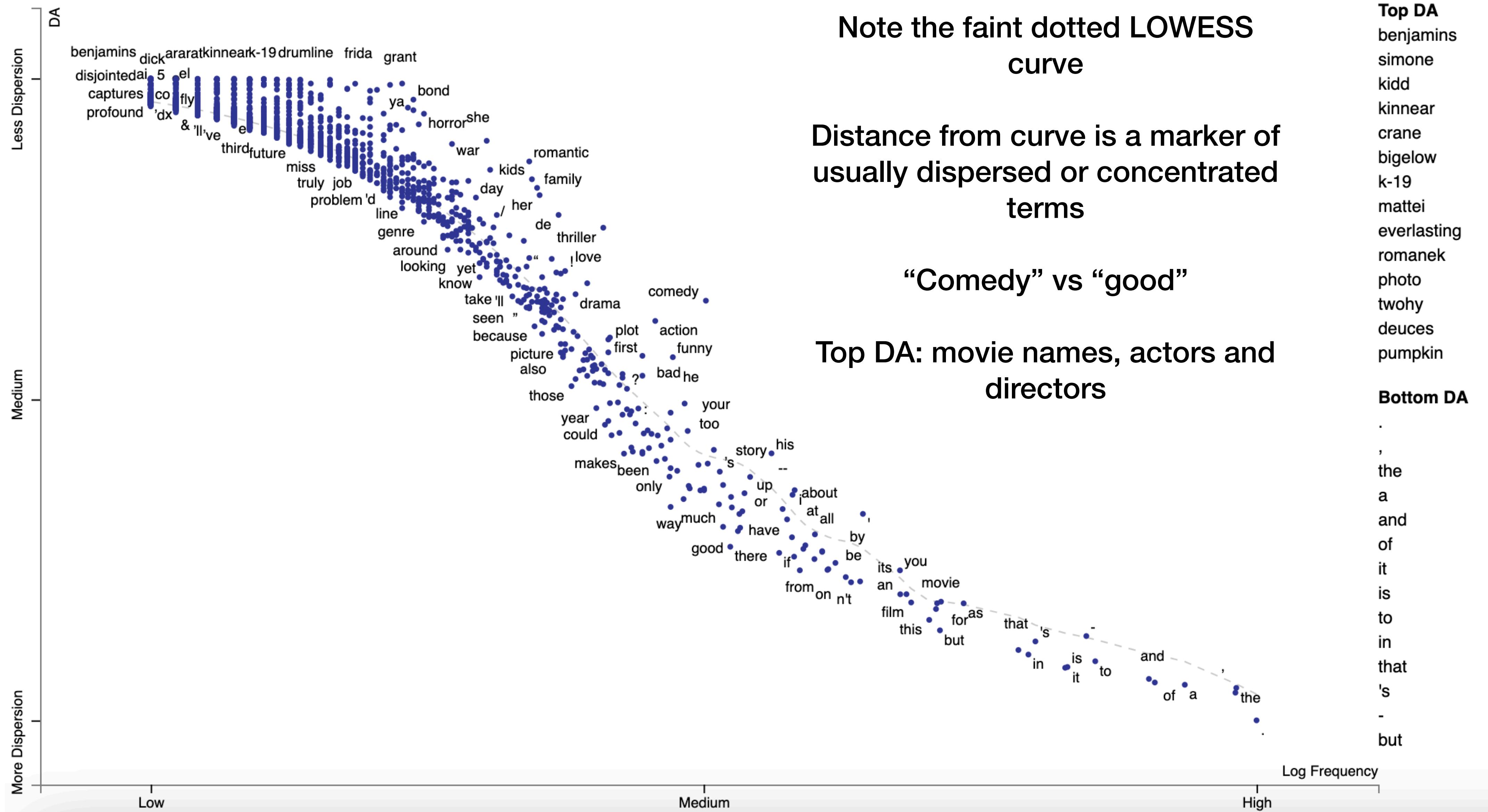
Dispersion

- Measures of dispersion
 - Lots of metrics
- Burch's DA (2017): pairwise difference of absolute deviations.
 - **Assumes roughly equal sized parts**
 - Range: [0, 1] (0: max dispersion)
 - DA = 0.5 when dispersion = the average term frequency
- Inversely correlated with term frequency
- $O(|\text{parts}|^2)$ required to compute

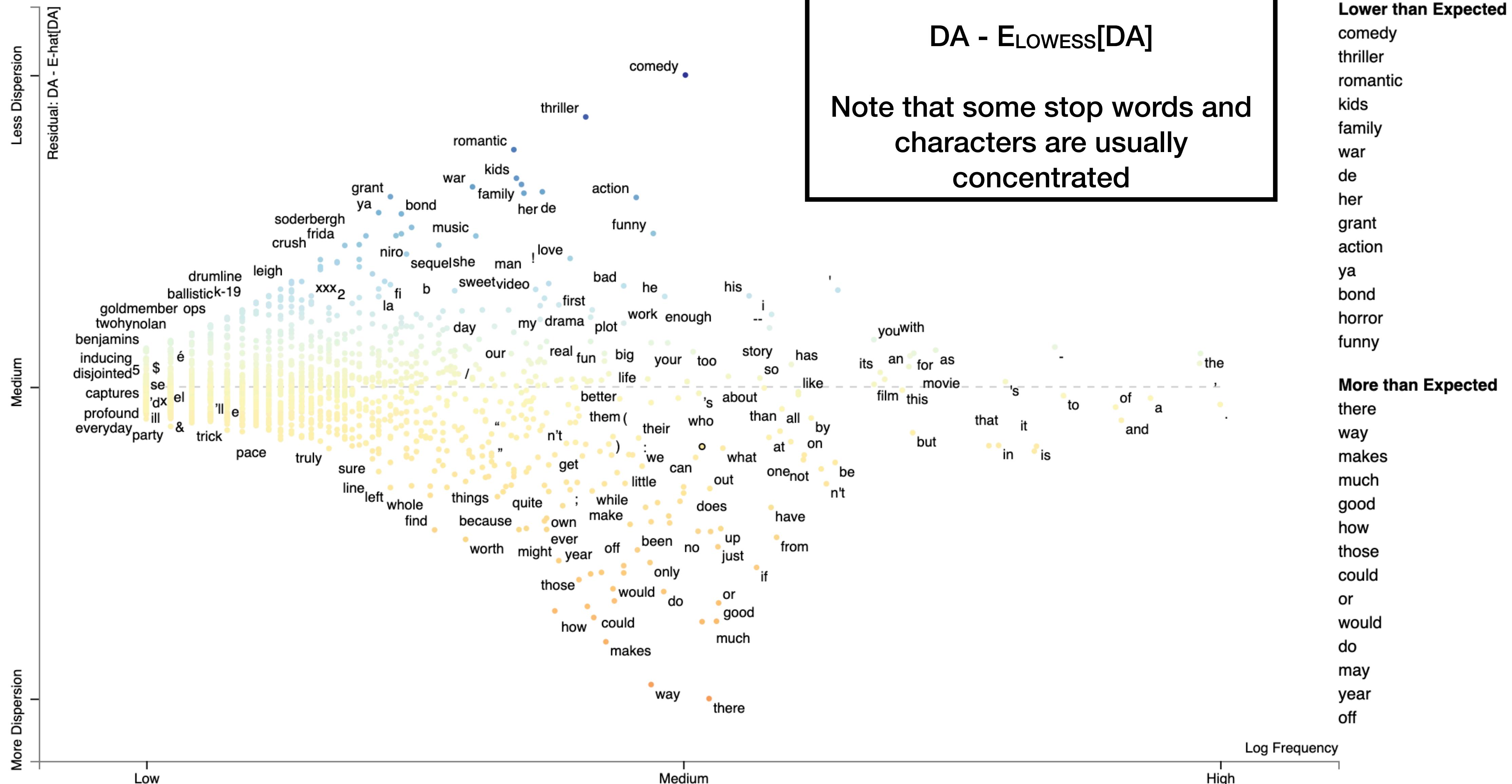
Y_i : term frequency in part i
k: # parts in corpus

$$D_A = 1 - \frac{\frac{1}{k(k-1)/2} \sum_{i=1}^{k-1} \sum_{j=i+1}^k |Y_i - Y_j|}{2\bar{Y}}$$

Dispersion in Movie Corpus



Residual Dispersion



ACL Paper Title Corpus

- 10,974 long ACL/EMNLP paper titles from 1979 to 2022
- Objective: see changes in the field of NLP
- Consider year published as category

Visualizing Over Time

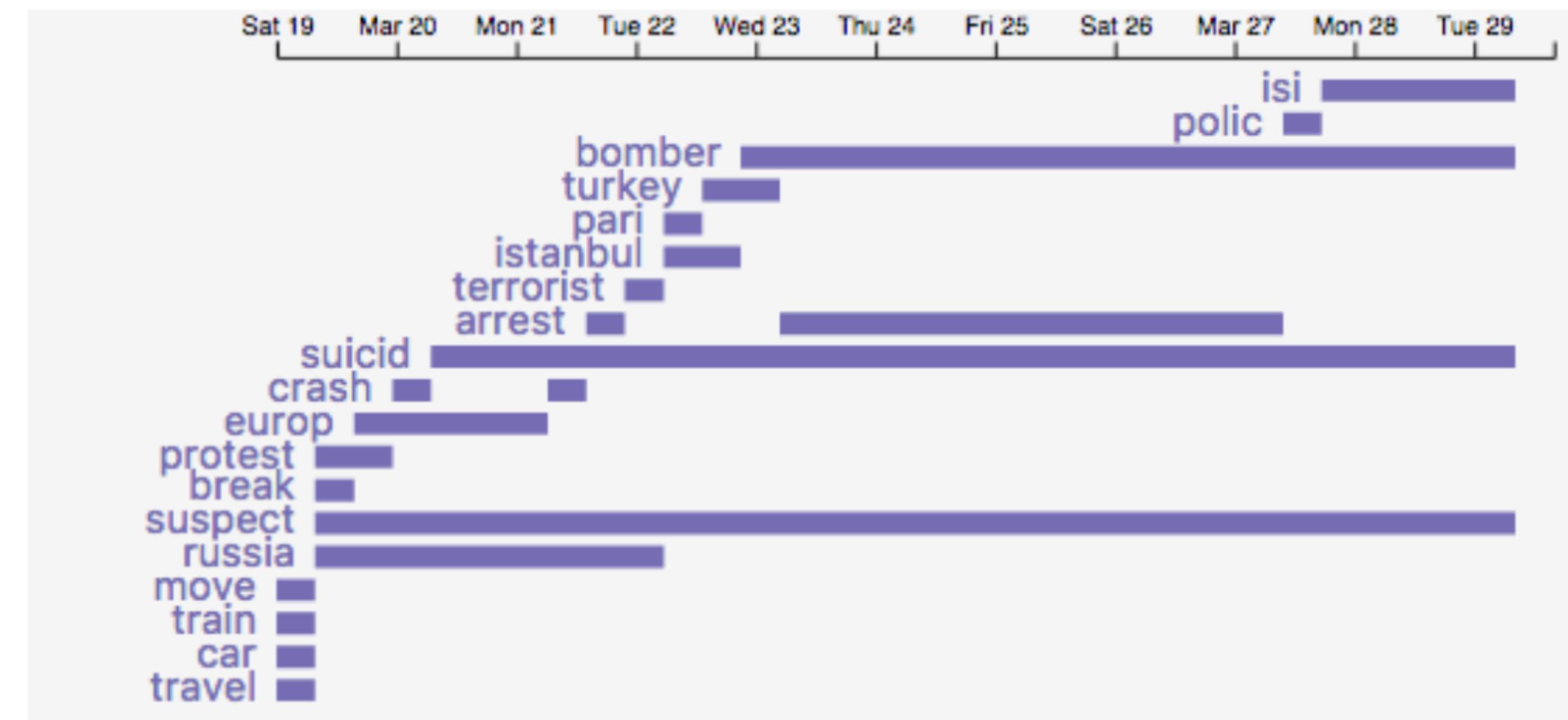
- We can make visualization similar to parallel tag clouds
- Ordered by keyness (log-odds-ratio)
- Size frequency in category (globally scaled)
- Too many years to view without scrolling

	1981	1982	1983	1984	1985
issue		access	properties	project	
natural language interface		database	formal	syntactic analysis	
production		technical	mathematical	simulation	
aspects		standard	tree adjoining grammars	design	
component		acl	definite	machine translation system	tree adjoining grammars
artificial		issue	treatment	computer	letter
interface		design	natural language interface	production	movement
strategies		experience	partially	conceptual	repair
experimental		solution	talk	linguistically	assignment
database		salience	technique	environment	binary

Related Work

ESTEEM

- Uses a Gantt chart to display words with similar embeddings to a query
- Embeddings are produced at each time step
- Simple, easy-to-follow visualization
- Few terms can be visualized this way



- Stemmed words with similar day-specific embeddings to “bomb” from German Tweets after the 2016 Brussels terrorist attack

Visualizing Over Time

- Solve scrolling problem: merge years together with high cosine similarity of keyness values of displayed terms

1979-1983	1984-1988	1989-2006	2007-2008	2009-2010	2011-2015	2016-2018	2019-2022
natural language interface project		unification - base	topic segmentation	web search	joint inference	attention model	bert
design	tree adjoining grammars	maximum entropy	phoneme conversion	expectation	distributional semantics	lstm	pre - train
issue	machine translation system	unification	semantic information	letter	segmentation use	neural semantic	pre - training
access	computer	automatic acquisition	online learning	phrase - base statistical machine	dual decomposition	deep reinforcement learning	transformers
production	sharing	lr	non - local	system combination	decipherment	sequence learning	few - shot
database	design	utterances	semantic relatedness	hierarchical phrase - base translation	acl	skip	multi - hop
computer	natural language interface	processor	context - sensitive	binarization	reorder	sequence models	pretrain
speech act	production	corpus - base	em	base statistical machine translation	recursive neural	deep neural networks	pretrained
interface	unification grammar	machine learning	improve statistical machine translation	phrase - base translation	unsupervised semantic	bandit	pre - train language model
experimental	categorial	constraint - base	subjectivity	semi - supervised learning	walk	character - level	contrastive learning

Use residual dispersion to see interesting terms which do not make the tag clouds

1979-1983	1984-1988	1989-2008	2009-2010	2011-2015	2016-2018	2019-2022
natural language interface project		machine learning web search		joint inference	attention model	bert
design	tree adjoining grammars	unification - base	expectation	distributional semantics	lstm	pre - train
issue		machine translation system	automatic acquisition	letter	segmentation use	neural semantic
access	computer	lr	phrase - base statistical machine	dual decomposition	deep reinforcement learning	transformers
production	sharing	maximum entropy	system combination	decipherment	sequence learning	few - shot
database	design	lexicalized	hierarchical phrase - base translation acl		skip	multi - hop
computer	natural language interface	corpus - base	binarization	reorder	sequence models	pretrain
speech act	production	large corpora	base statistical machine translation	recursive neural	deep neural networks	pretrained
interface	unification grammar	co - occurrence	phrase - base translation	unsupervised semantic	bandit	pre - train language model
experimental	categorial	unification	semi - supervised learning	walk	character - level	contrastive learning

