

# **Competition, selection and communicative need in language change**

**Võistlus, valik ja vajadus keele muutumises**

**Andres Karjus**

ERA Chair for Cultural Data Analytics, Tallinn University  
[andreskarjus.github.io](https://andreskarjus.github.io)

In collaboration with: Kenny Smith, Richard A. Blythe, Simon Kirby  
(University of Edinburgh)

TÜLING @ Tartu University | 27.04.2021

Started postdoc in 2020 in



<http://cudan.tlu.ee>

PhD (2020) from



MSc from



BA & MA from



# All living languages keep changing

- All the time
- Eventually diverge into different languages
- This is weird
- This research: focus on lexical change and competition therein
- What happens when new words are introduced into language?
- Why some semantic domains are more complex than others?
- Massive centuries-spanning corpora open up an unprecedented avenue of possible investigations into language dynamics.
- Variant usage frequencies but also meaning (and change) using distributional semantics methods
- Individual-level processes can be probed using human experiments

# In this talk

- Semantic similarity, colexification, and communicative need
- Experiment: communicative need modulates colexification
- Corpus studies:
  - Topical advection as a baseline model of frequency change and a proxy to communicative need
  - Communicative need modulates lexical competition
- Future directions: competition, complexity and informativeness

# Some concepts

a semantic space

words

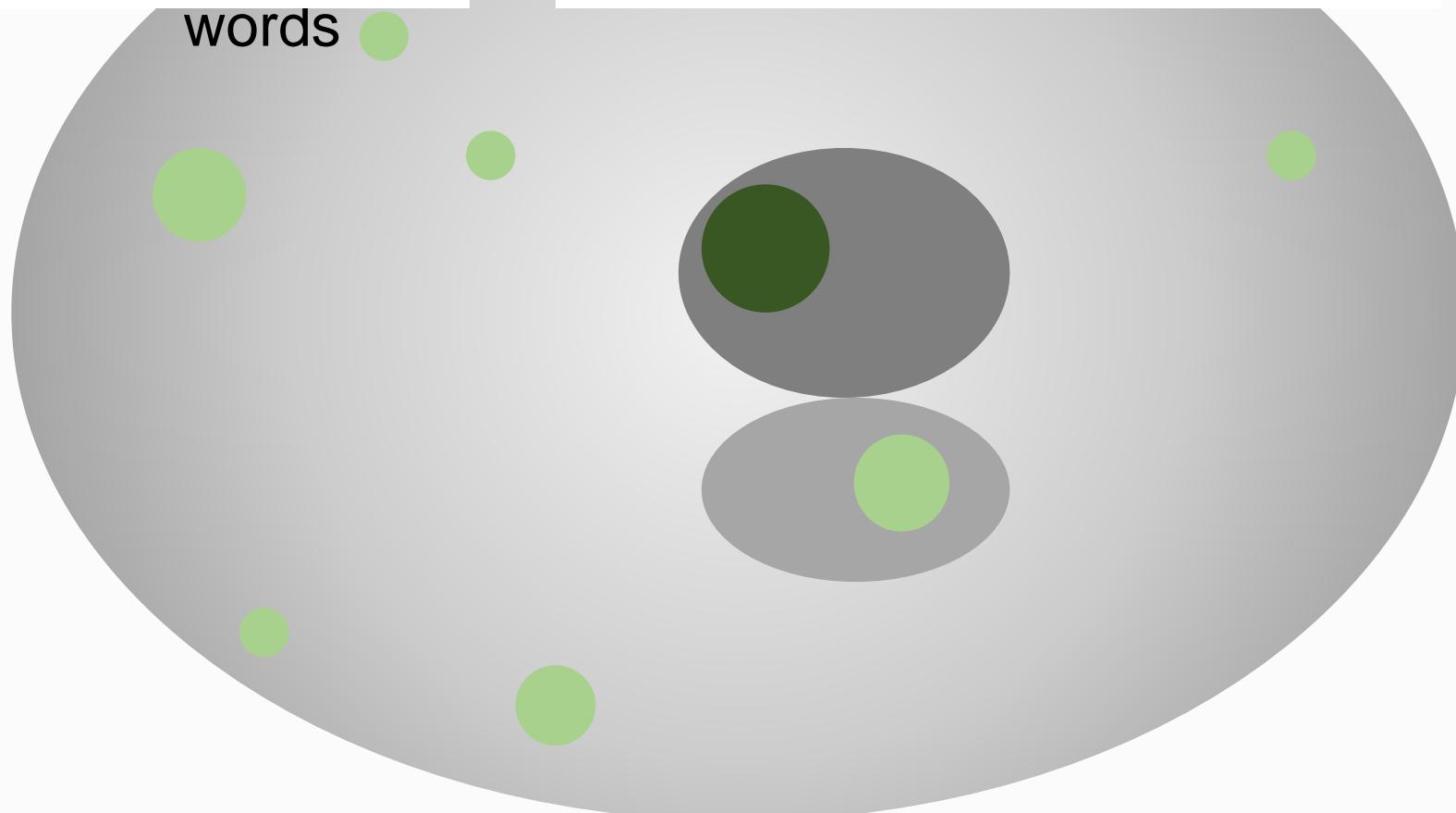
lexifies  
a meaning  
another meaning  
another word

“colexification”

# Complexity and informativeness

inverse of simplicity  
relates to learning  
cognitive cost

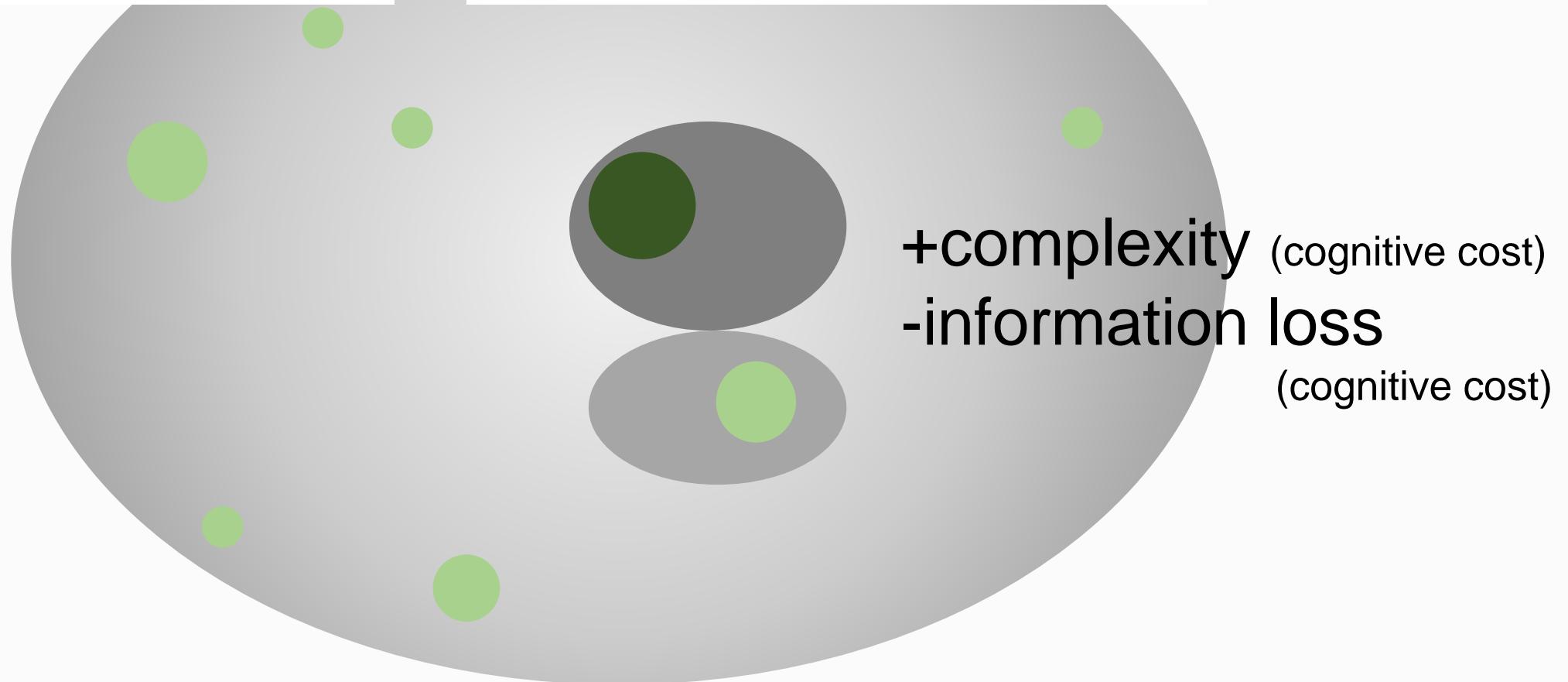
inverse of information loss  
accuracy, expressivity  
communicative cost



# Complexity and informativeness

inverse of simplicity;  
relates to learning  
cognitive cost

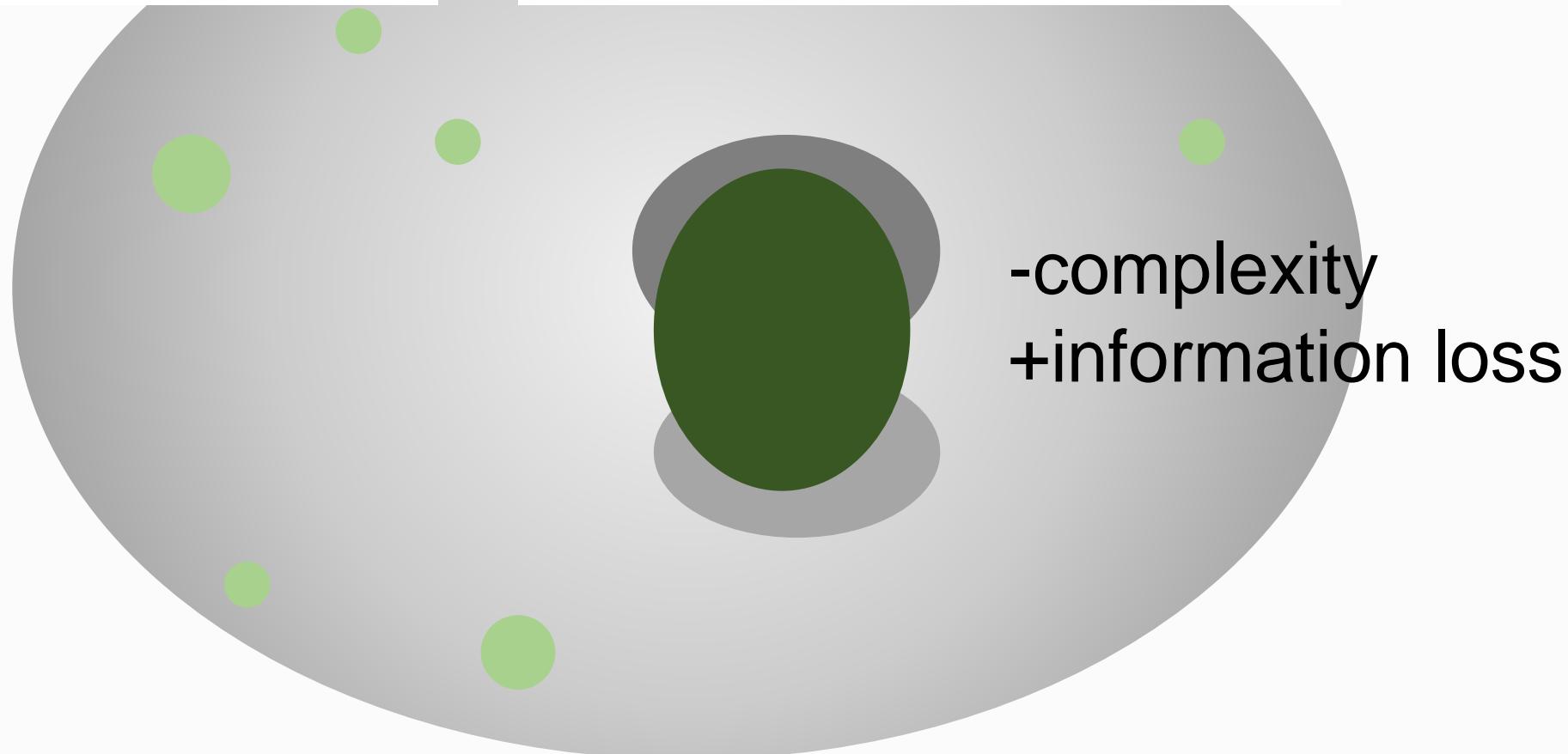
inverse of information loss;  
accuracy, expressivity  
communicative cost



# Complexity and informativeness

inverse of simplicity  
relates to learning  
cognitive cost

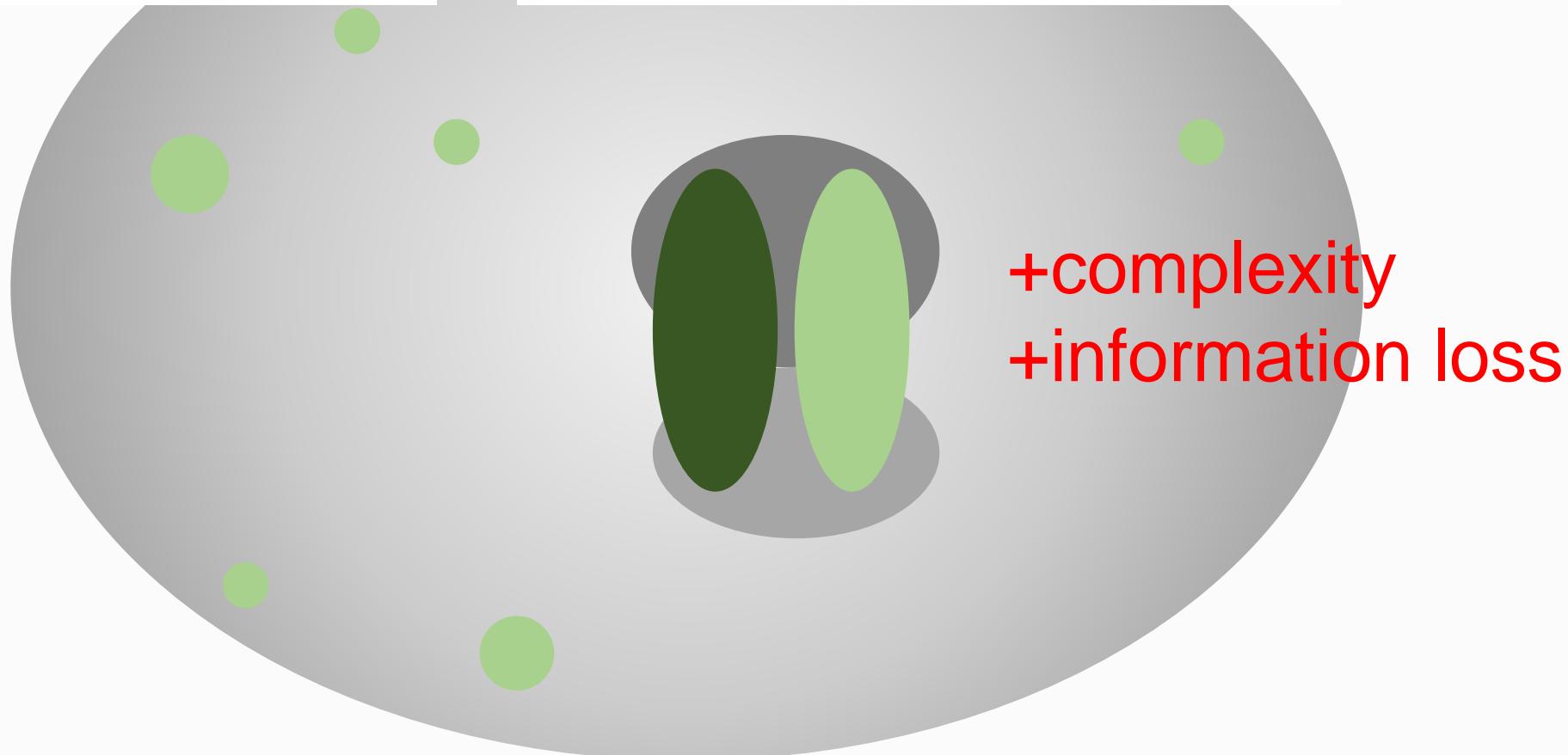
inverse of information loss  
accuracy, expressivity  
communicative cost



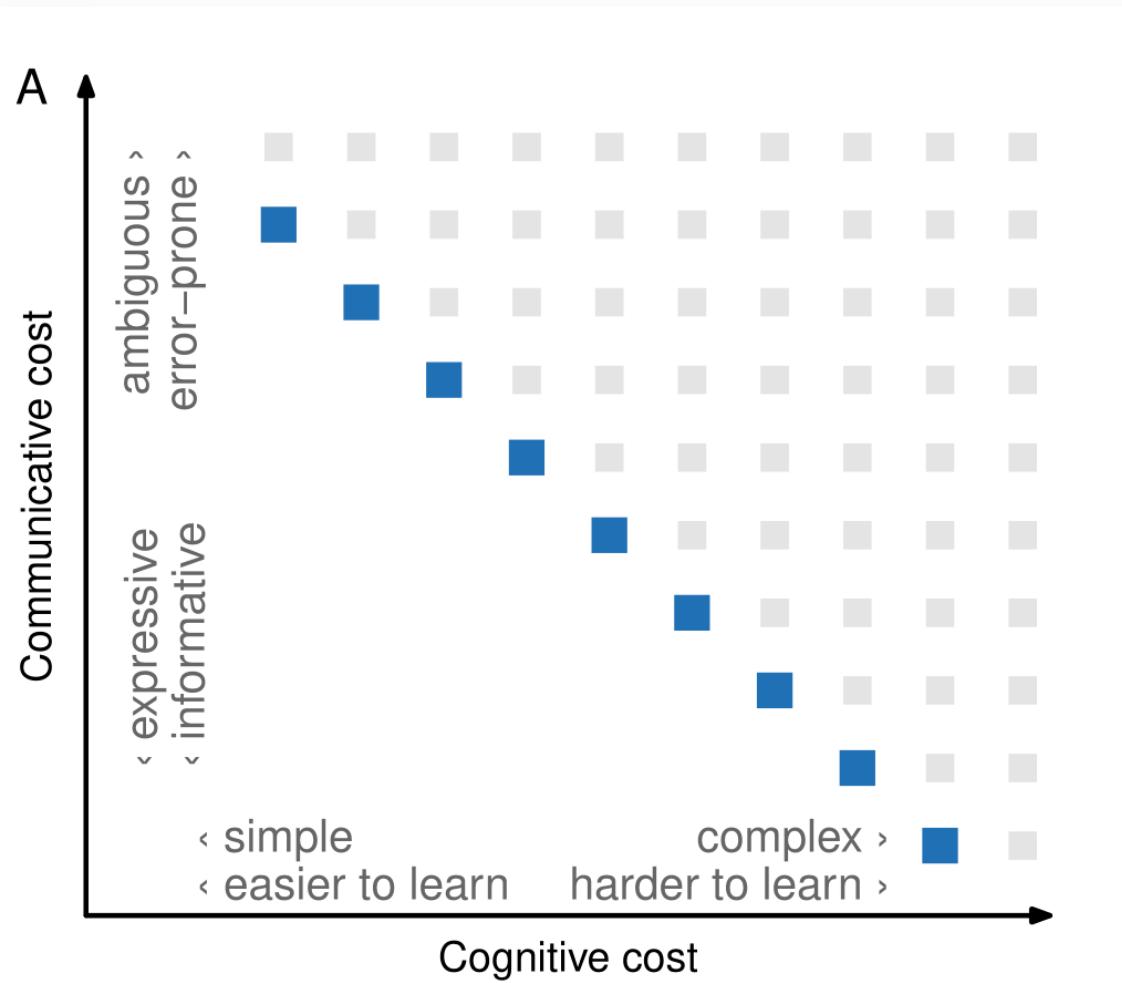
# Complexity and informativeness

inverse of simplicity  
relates to learning  
cognitive cost

inverse of information loss  
accuracy, expressivity  
communicative cost

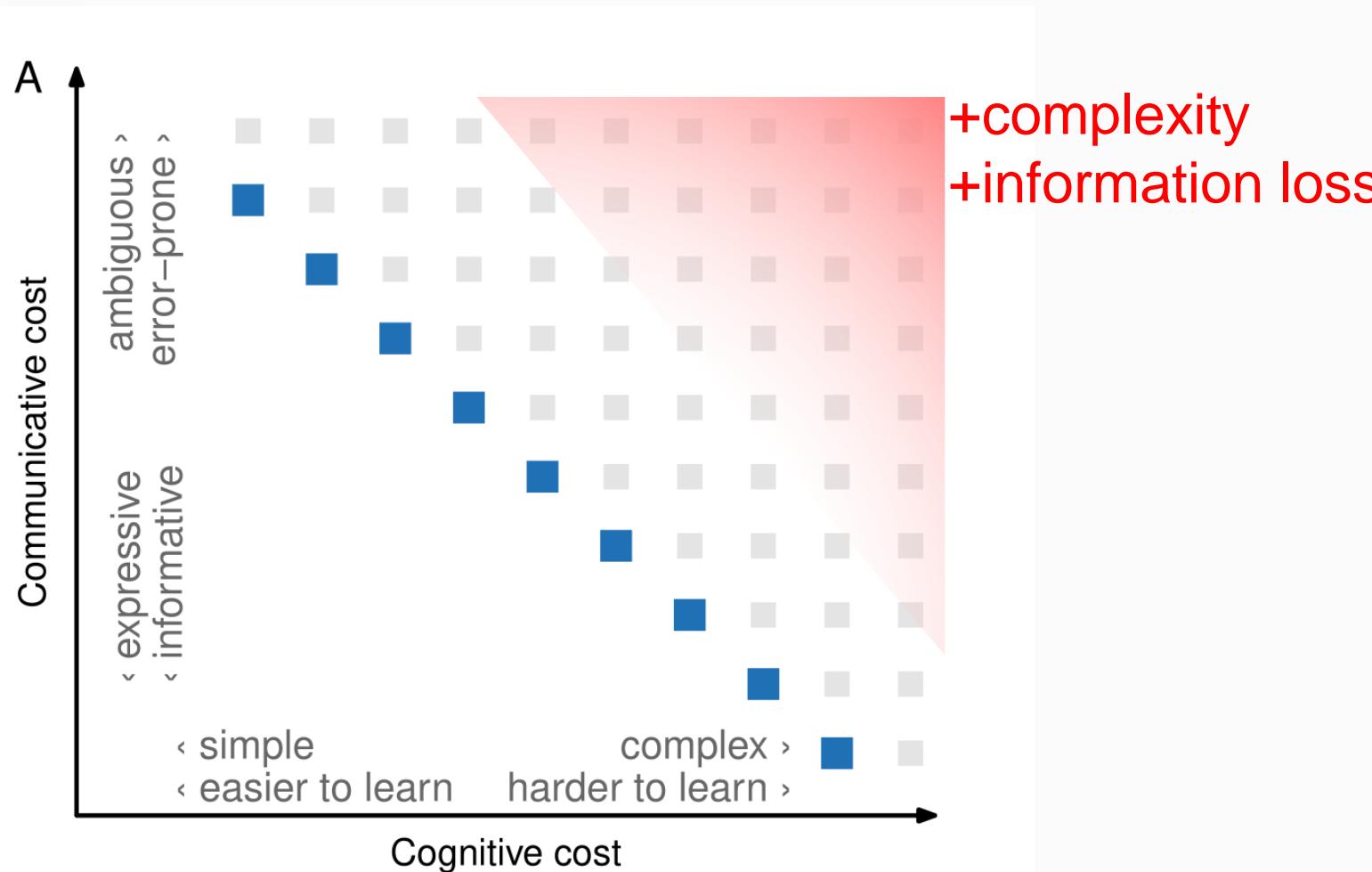


# The complexity-informativeness tradeoff and the optimal front



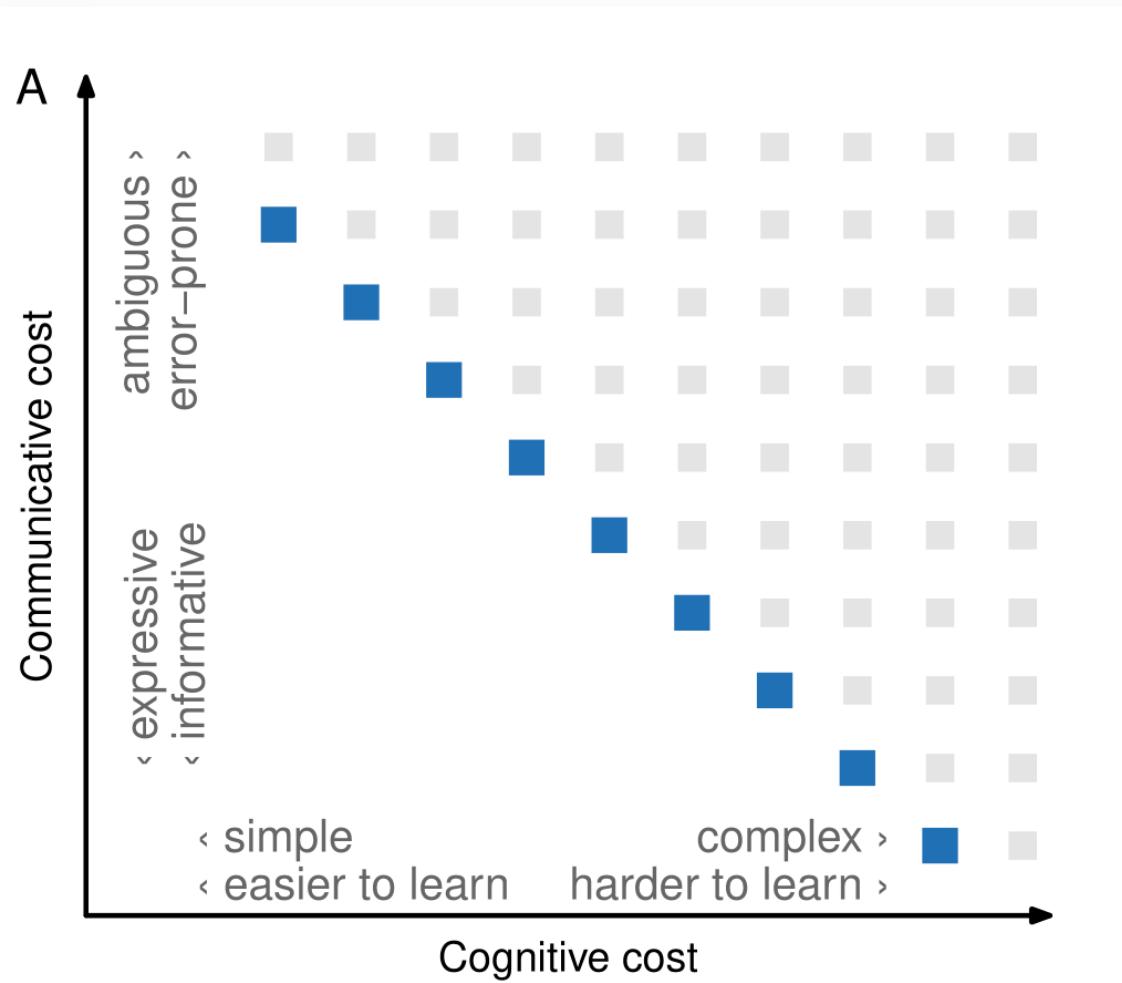
cf. Kemp et al 2012, Kemp et al 2018, Carr et al 2020

# The complexity-informativeness tradeoff and the optimal front



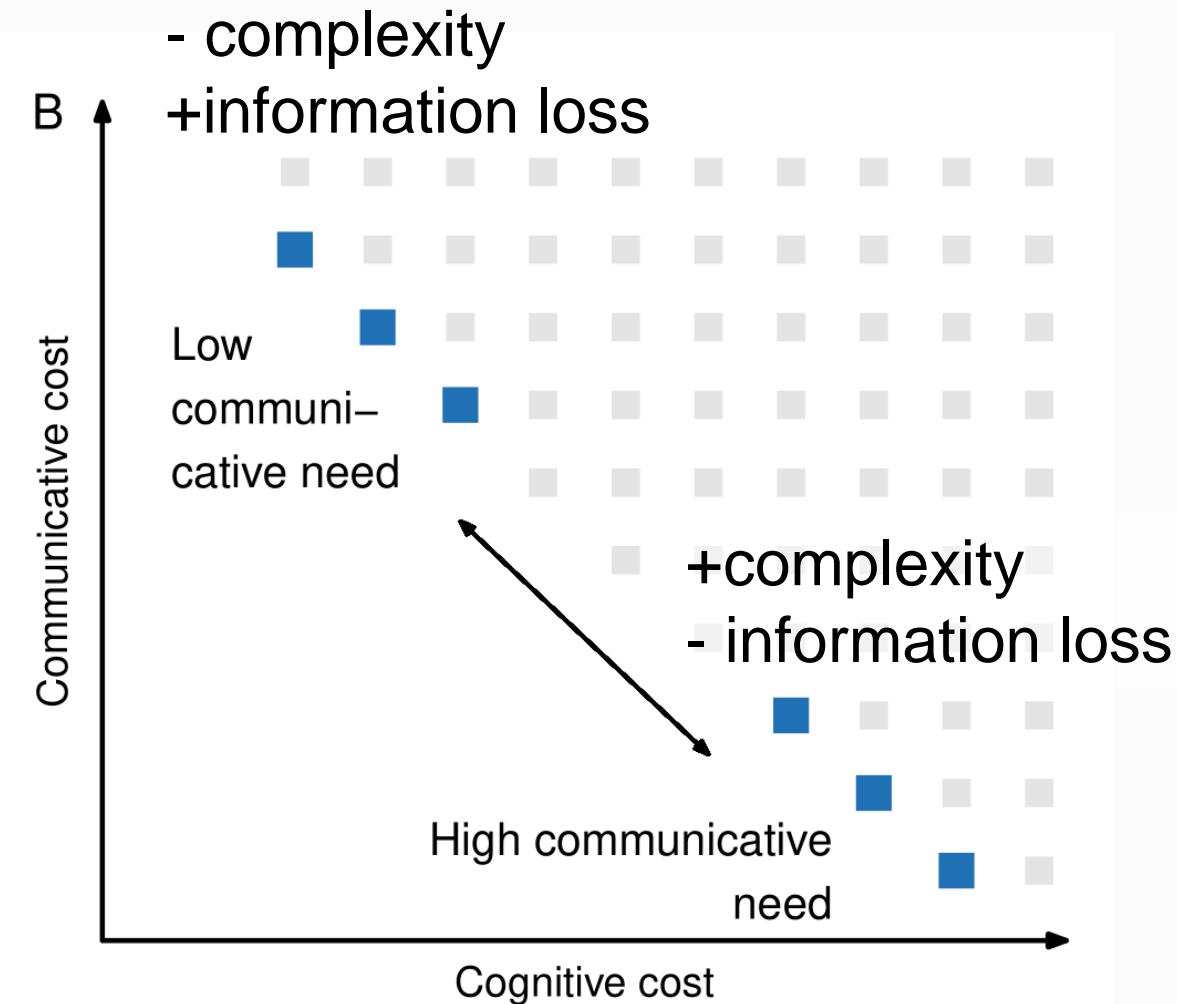
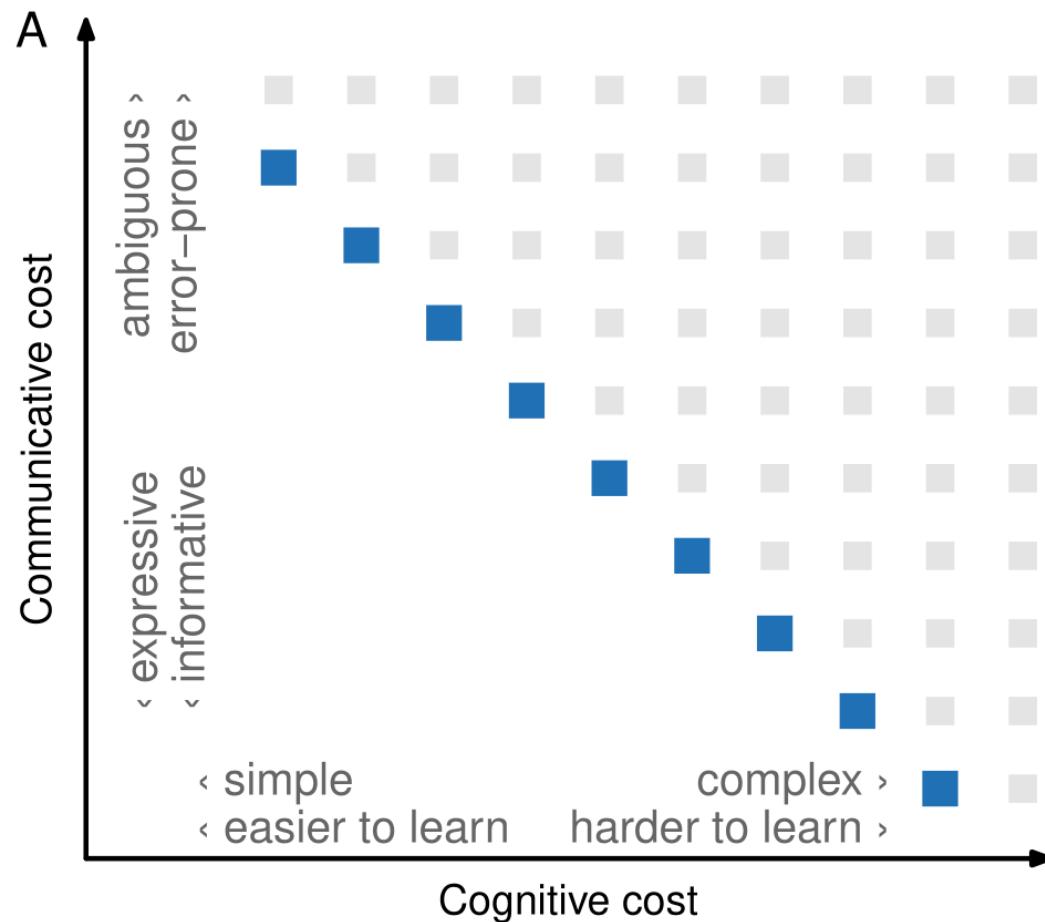
cf. Kemp et al 2012, Kemp et al 2018, Carr et al 2020

# The complexity-informativeness tradeoff and the optimal front



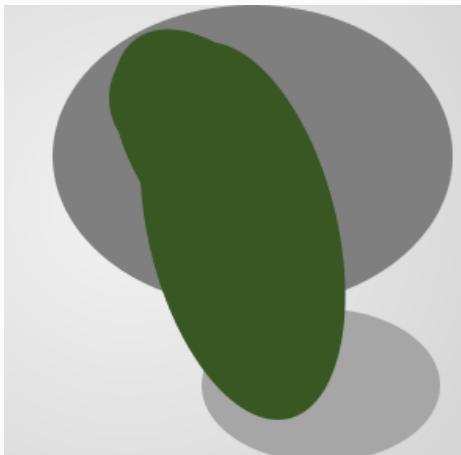
Describes lexicons of kinship terms, colour, numeral systems, negation; similar optimization effects in artificial language experiments

# The complexity-informativeness tradeoff and the optimal front



# Conceptual similarity and communicative need shape colexification: an experimental study

- Karjus, Blythe, Kirby, Wang, Smith, 2021 <https://arxiv.org/abs/2103.11024>
- Xu et al 2020, “Conceptual relations predict colexification across languages”, using 200+ languages
- Similar and associated senses (e.g. FIRE and FLAME) are more frequently **colexified** in world’s languages than unrelated or weakly associated meanings (like FIRE and SALT)



# Conceptual similarity and communicative need shape colexification: an experimental study

- ...but culture specific **communicative needs** should affect likelihood of colexification – e.g. if it is necessary for efficient communication to distinguish some similar meanings
- E.g. ICE and SNOW: less likely to be colexified in cold climates (Regier et al 2016)

# Conceptual similarity and communicative need shape colexification: an experimental study

- What is the cognitive mechanism though that leads to this cross-linguistic tendency?
- Maybe we can test these two claims experimentally?
- 4 experiments: initial one with student sample, replication on Mechanical Turk, two more follow-up experiments
- Dyadic communication game setup, 2 players, take turn sending and guessing messages (cf. Kirby et al 2008, Winters et al 2015)
- 135 rounds each (data from the first 1/3 of the game excluded)

- 10 meanings total
- 4 distractor meanings
- from Simlex999
  
- 6 target meanings
- 3 pairs
  
- Baseline: pairs co-occur uniformly
- Target condition: **similar ones occur together more often!**

WARRIOR	neme
THEFT	quto
STATE	nopo
RHYTHM	fita
TASK	mefa
JOB	mumi
PAIR	honи
COUPLE	
SHORE	
COAST	
	7 signals

# The game

Player 1

Players connected: 2. Score: 0/2

area fashion

Communicate *area* using...

piti

wuli

liha

naru

mano

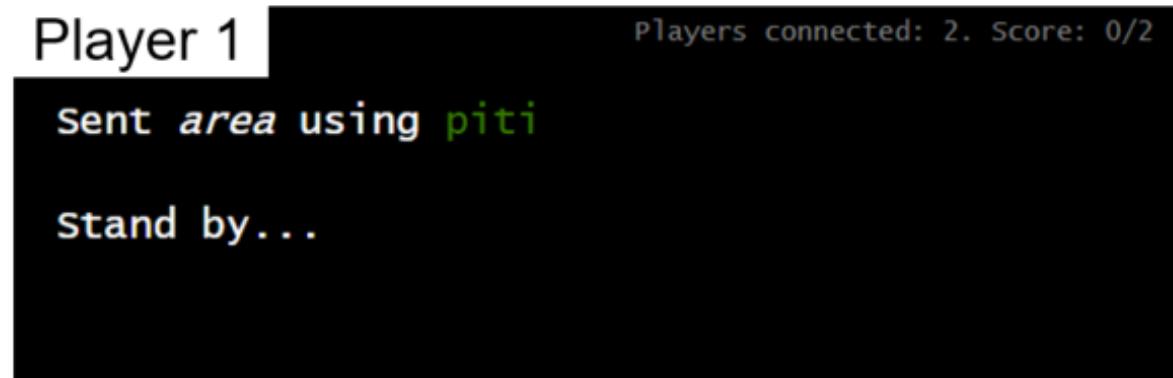
himu

qata

# The game



# The game



# The game

Player 1

Players connected: 2. Score: 0/2

area fashion

Communicate *area* using...

piti

wuli

liha

naru

mano

himu

qata

Player 2

Players conn

area fashion

Waiting for message...

Player 1

Players connected: 2. Score: 0/2

Sent *area* using piti

Stand by...

Player 2

Players connected: 2. Score: 0/2

area fashion

Message: piti

This means:

area

fashion

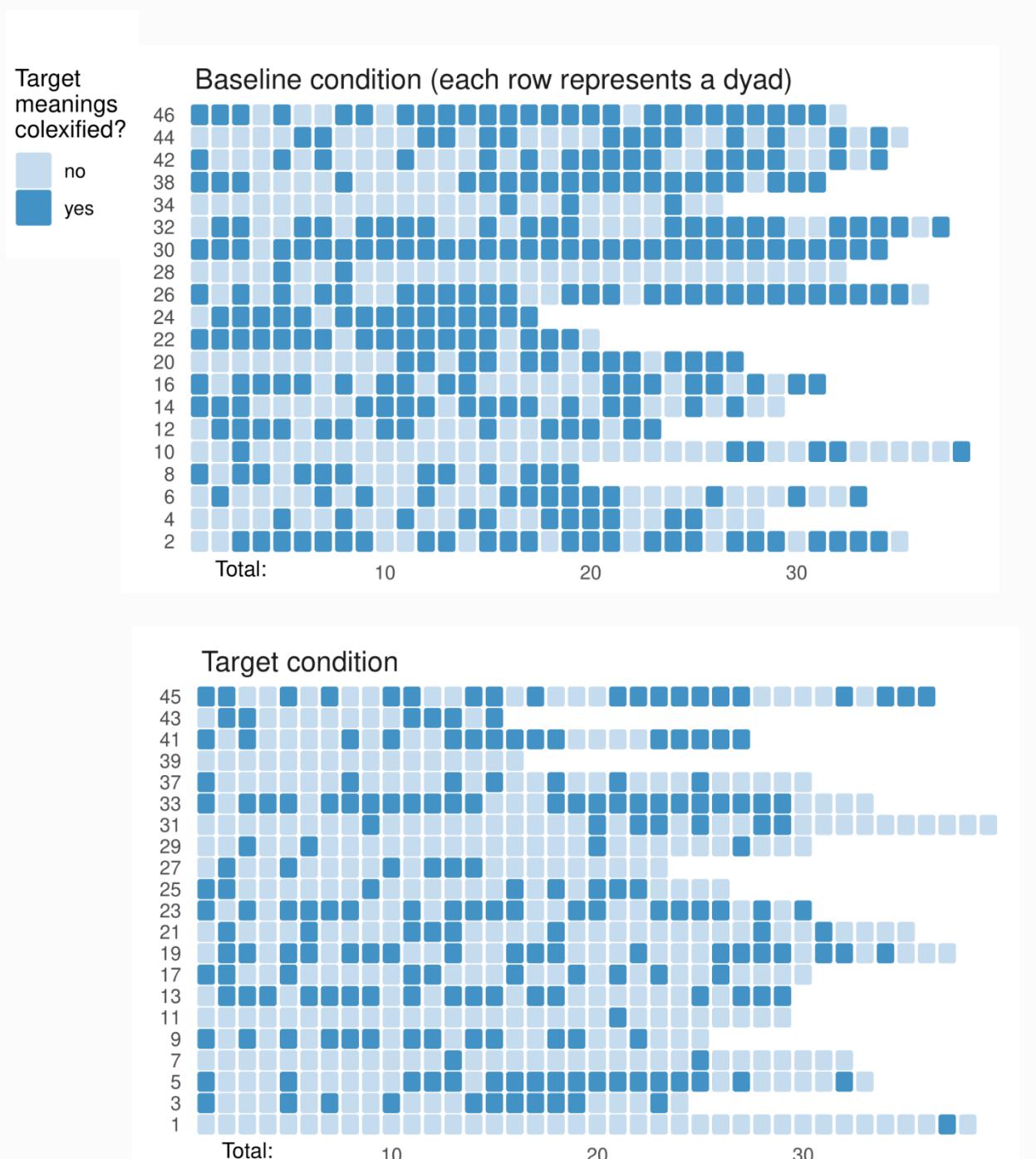
Expm no. 38, baseline condition, 96%, counts

WARRIOR	2						7
THEFT				9			
STATE					9		
RHYTHM						9	
TASK		2	4	2			1
JOB			9				
PAIR	8						
COUPLE	10						
SHORE		7				1	
COAST		10					
	neme	quto	nopo	fita	mefa	mumi	hon

7 signals

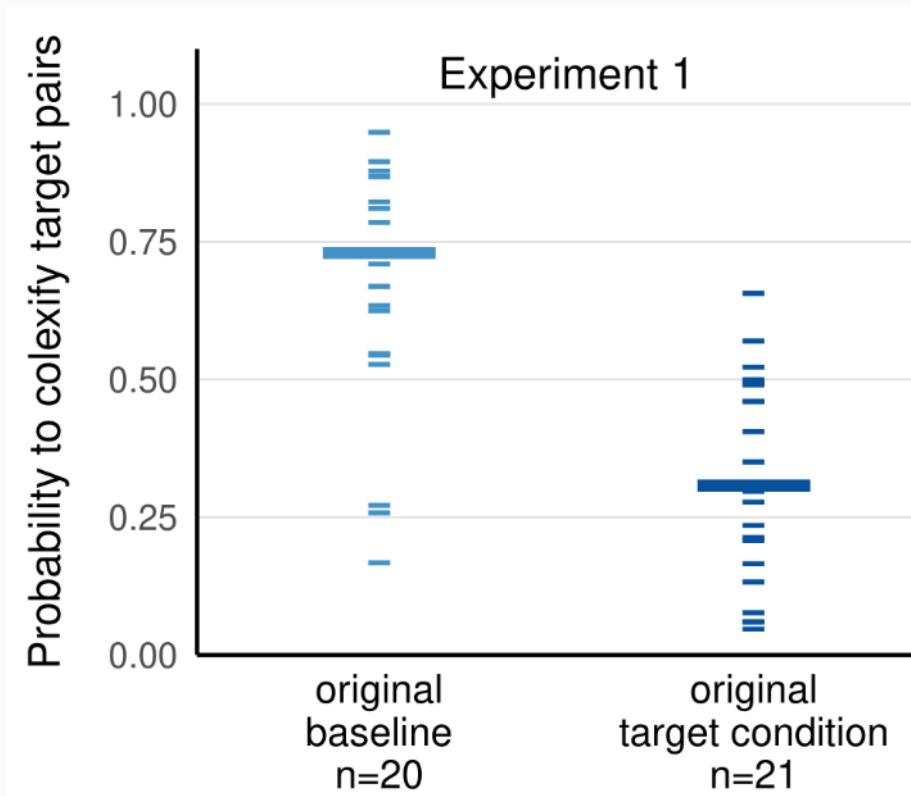
# Analysis

- Exclude low-accuracy dyads (41 left)
- Iterate through each experiment, record each instance of colexification (same signal, different meaning) involving a target meaning; n=1218.
- Logistic mixed effects regression (control for speaker/dyad, meaning pair)
- Colexification ~ condition\*round (dyads may change preferences over the course of the game)
- **Are similar meanings less likely to be colexified in the target condition?**



# Results

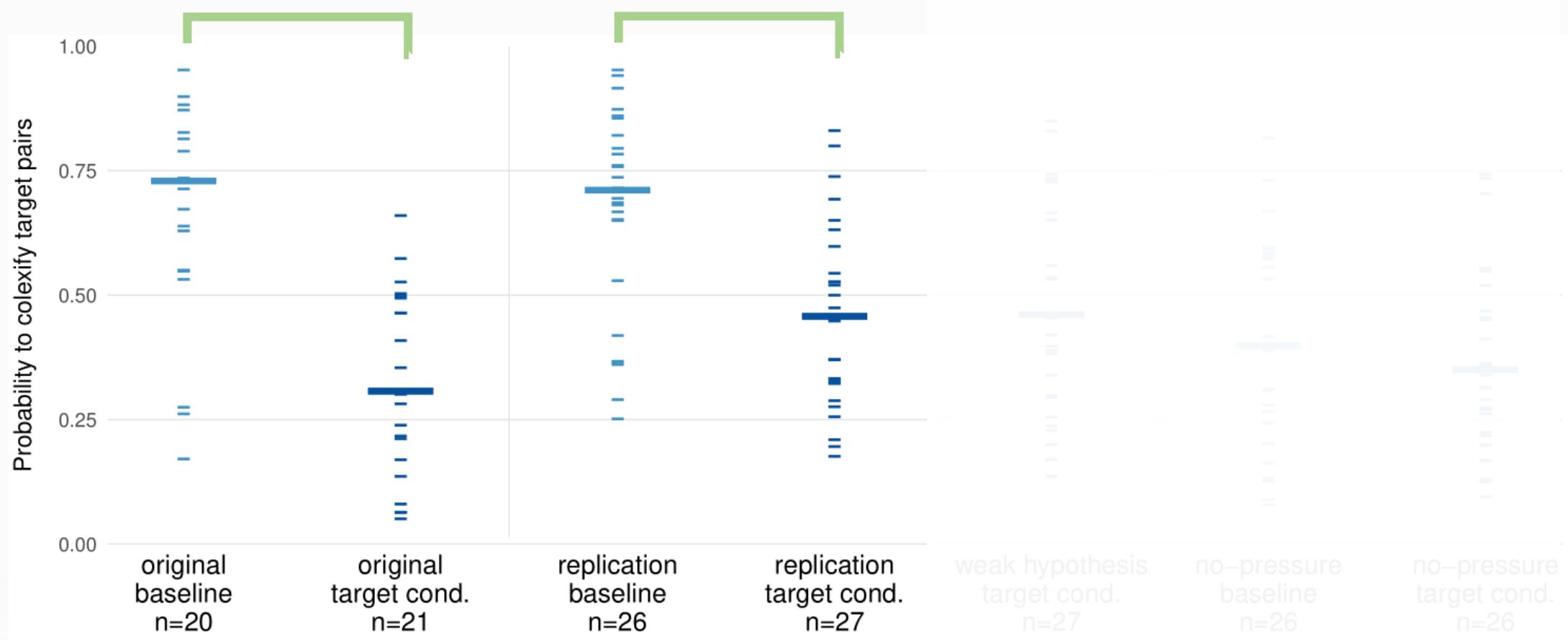
- Yes ( $p=0.001$ ).
- When need arises to distinguish similar meanings (target condition), speakers less likely to colexify them
- Confirms hypothesis that communicative needs may block colexification of related concepts.
- When no pressure to distinguish particular meanings (baseline), prefer to colexify similar meanings (confirms main finding of Xu et al 2020)



colexification ~	Estimate	SE	<i>z</i>	<i>p</i>
intercept (baseline condition, mid-game)	-0.22	0.37	-0.59	0.56
+ condition (target)	-0.52	0.51	-1.03	0.3
+ round	1.02	0.27	3.84	<0.01
+ condition (target) × round	-1.17	0.37	-3.17	<0.01

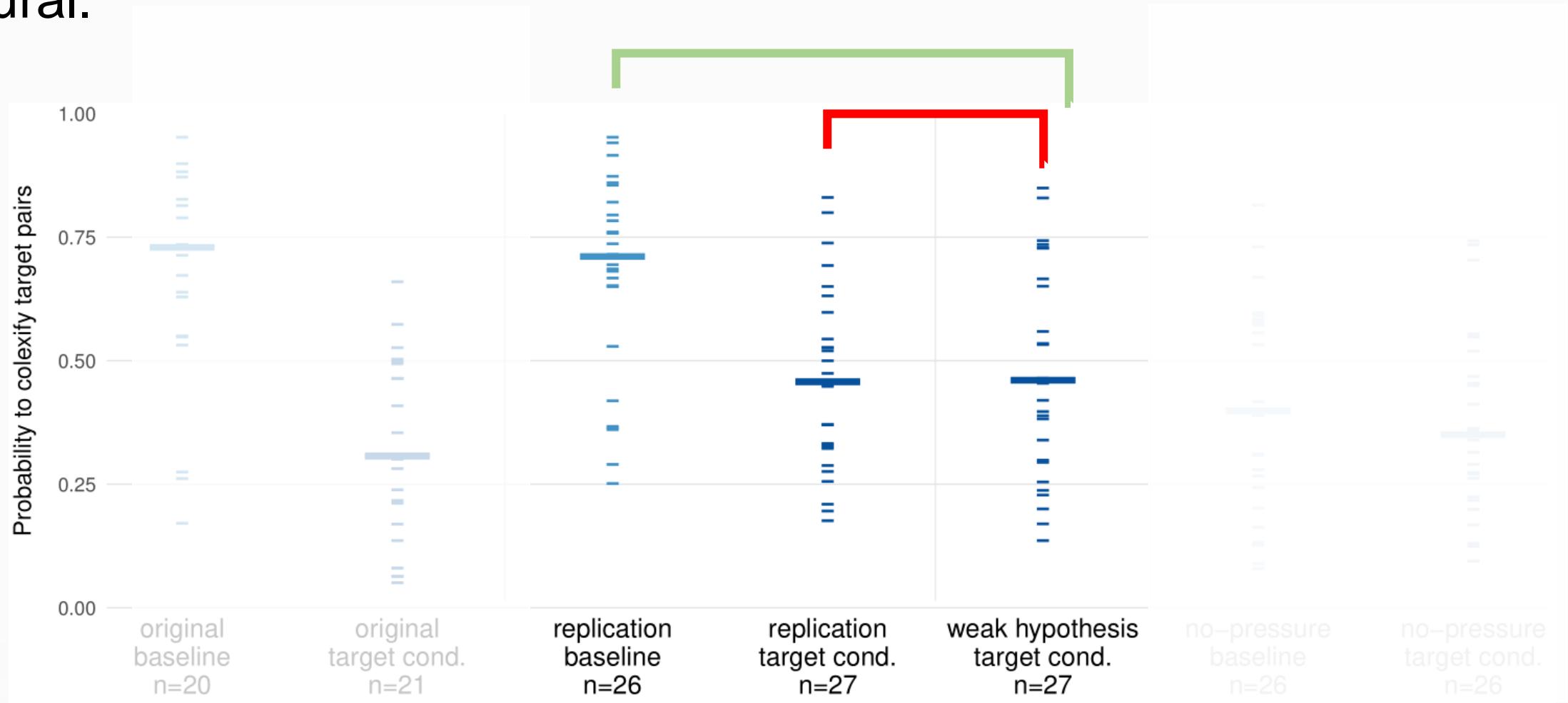
# Follow-up experiments

- Switch to Mechanical Turk (more flexible recruitment)
- Experiment 2, replication on MTurk, same setup
- Lower accuracy: 79 dyads, could use data only from 53.



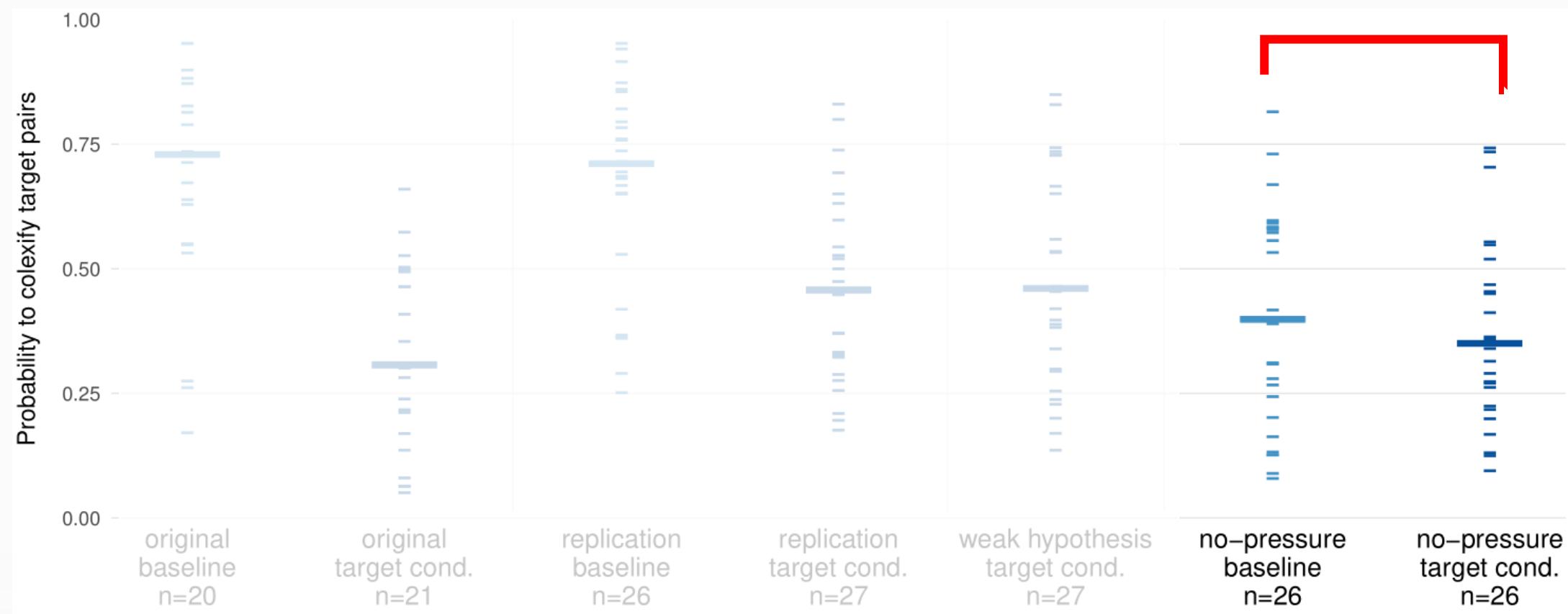
# Follow-up experiments

- Experiment 3, target condition only: introduce similar-meaning pairs into the distractor set to make colexifying them more natural.

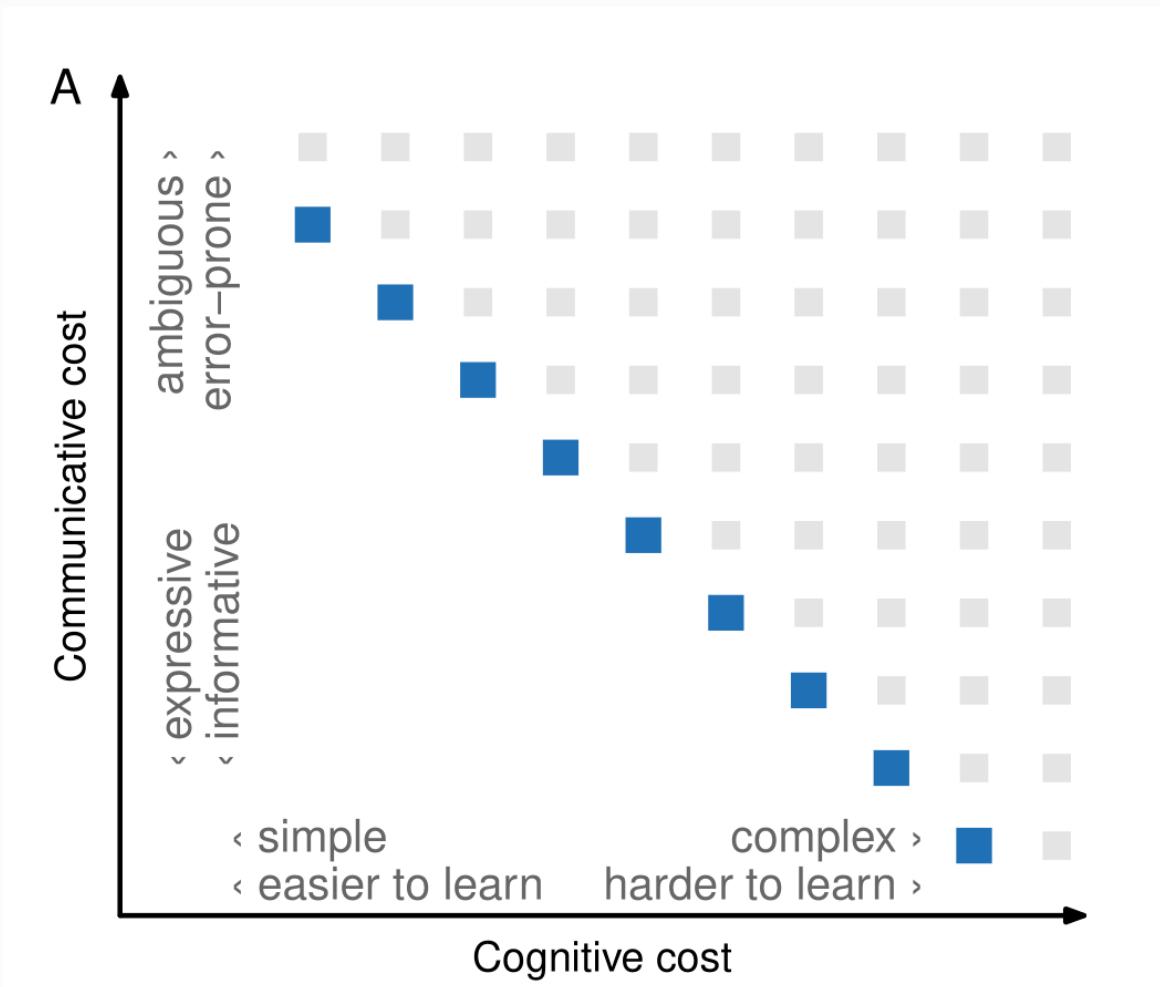


# Follow-up experiments

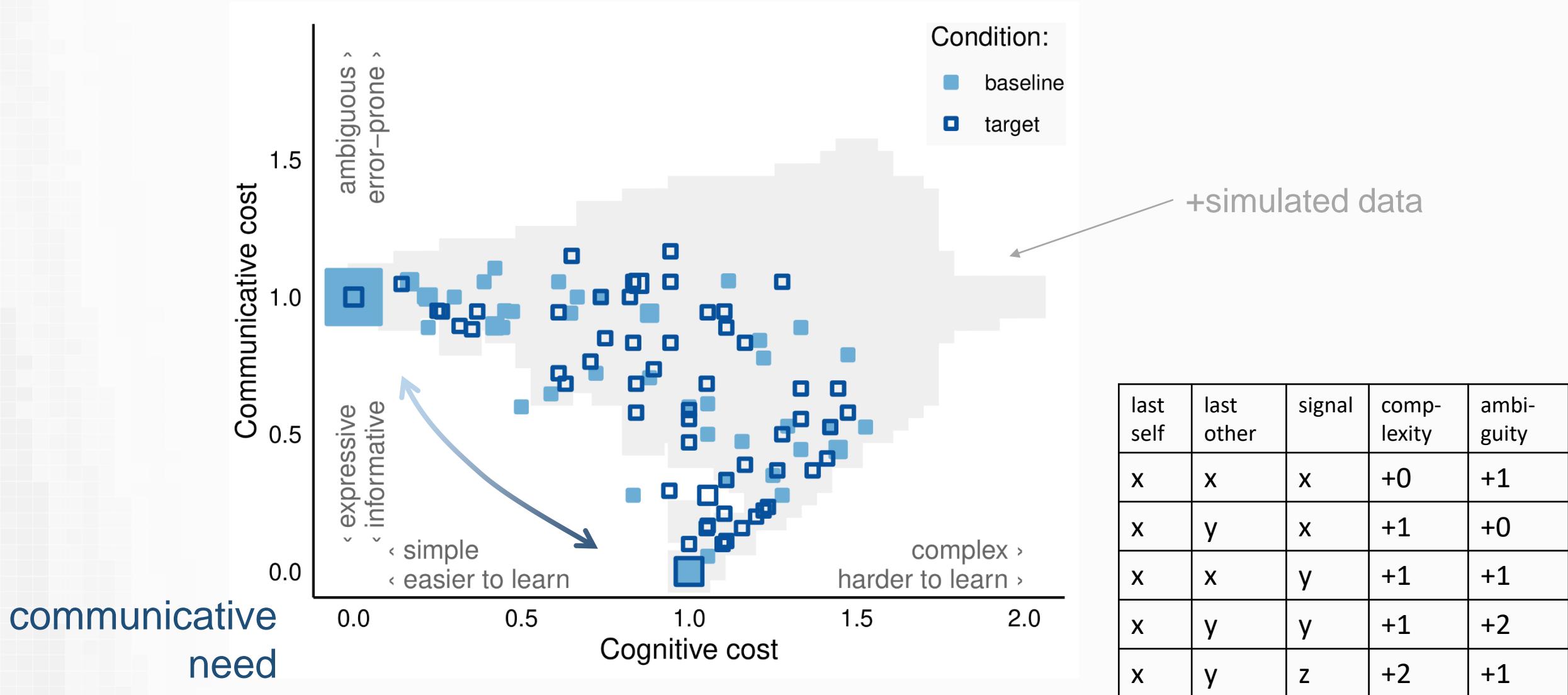
- Experiment 4: no pressure to colexify (10 signals for 10 meanings). No effect, and participants make significantly more use of the bigger signal space. But: natural language does have pressure to simplify (can't have infinite lexicons).



# The complexity-informativeness tradeoff and the optimal front

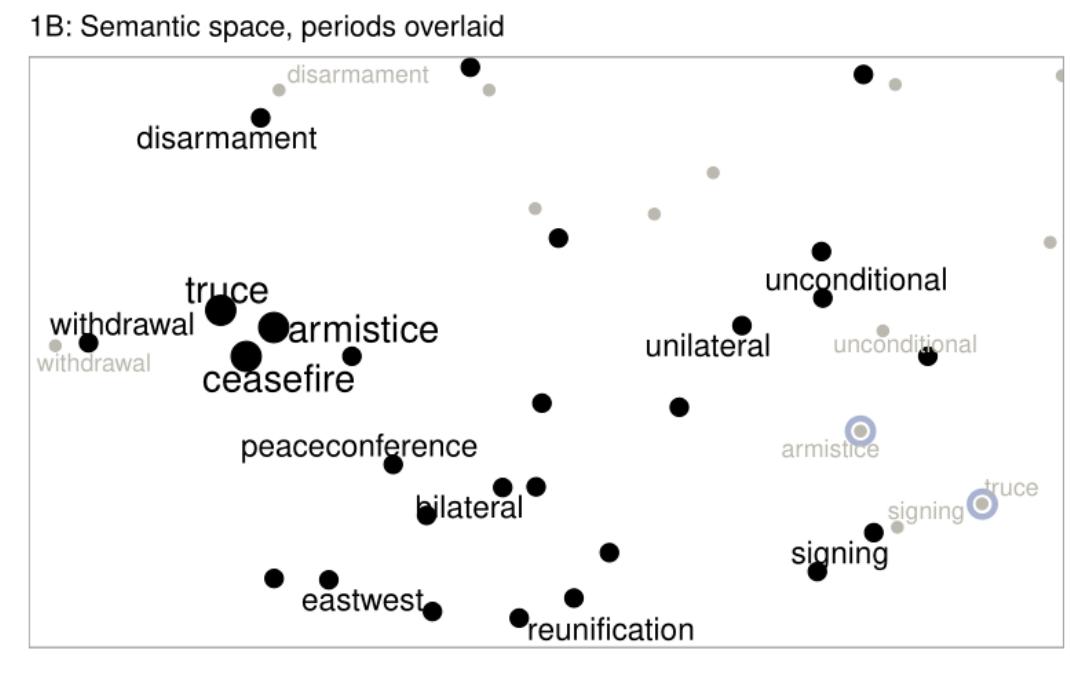
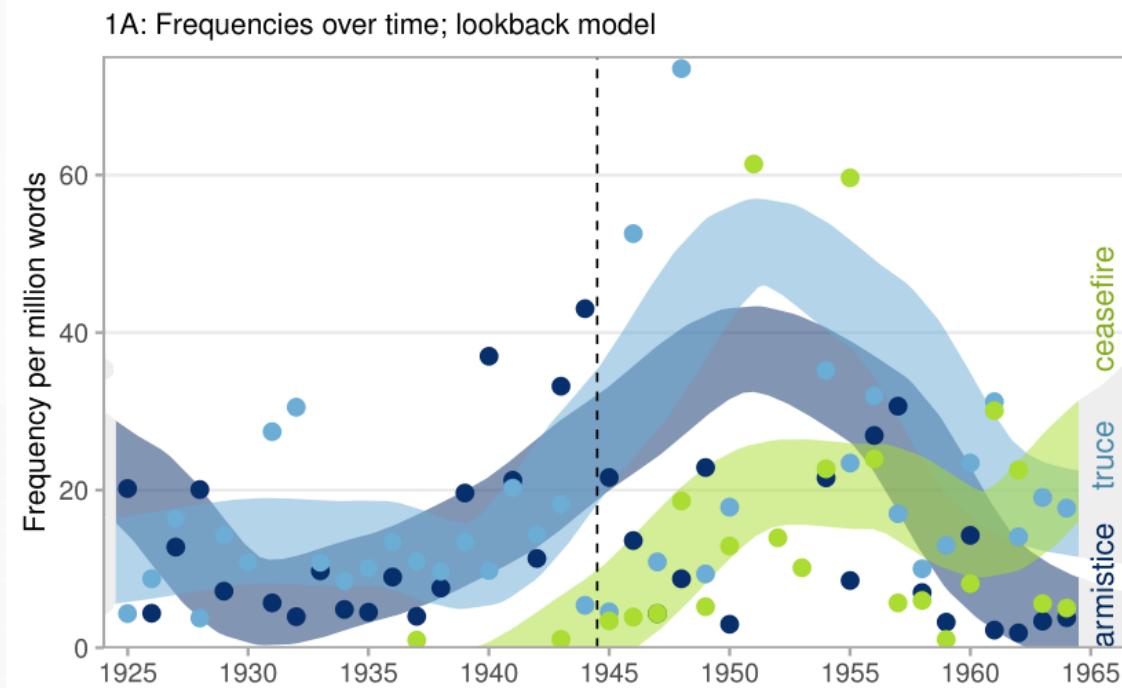


# The complexity-informativeness tradeoff and the optimal front



# Discussion

- Experimental results describe an individual-level lexical choice mechanism which produces results in line with typological colexification tendencies (Xu et al 2020) as well as the communicative need hypothesis
- Work in process: a model of lexical density (~extent of colexification) applied to word embeddings trained on diachronic corpora





# Communicative need modulates competition in language change

- Preprint: Karjus, Blythe, Kirby, Smith 2020  
<https://arxiv.org/abs/2006.09277>
- As new words, e.g. neologisms & borrowings are selected for, what happens to their older synonyms? Does direct competition always follow local frequency changes?
- Hypothesis:
  - frequency increase in a word will lead to direct competition with (and possibly replacement of) near-synonym(s)
  - unless the lexical subspace experiences high communicative need

# Communicative need modulates competition in language change

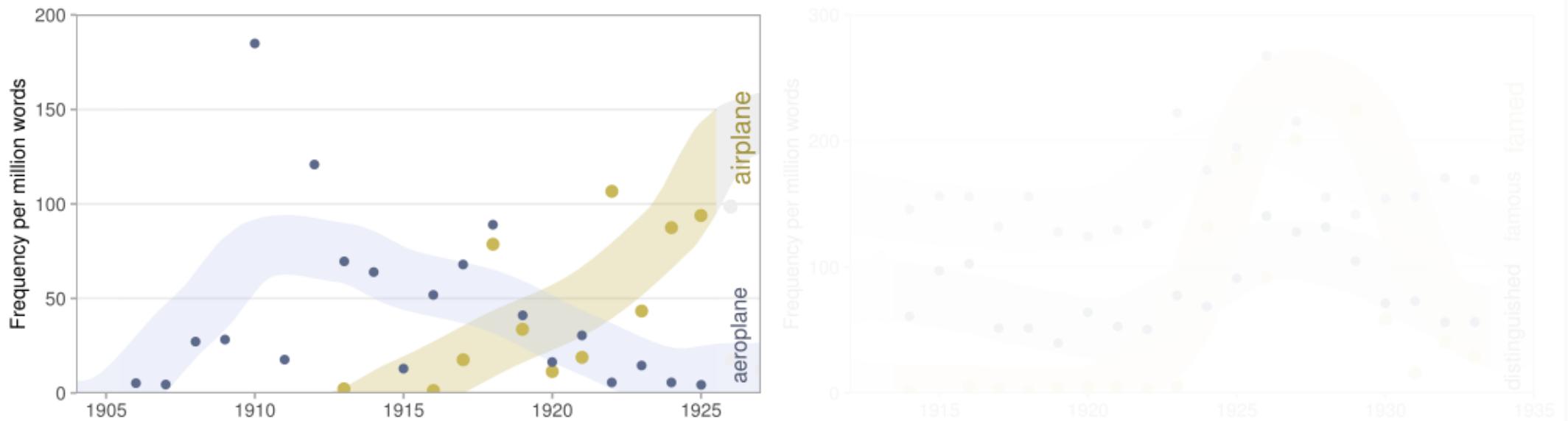


Figure 1: Example time series from the Corpus of Historical American English (COHA). Two decades after the invention of heavier-than-air powered aircraft, *airplane* replaced the initial term *aeroplane* (left side panel; the points are normalized yearly frequencies, with the lines representing smoothed averages for visual aid). Around the same time, *famed* appears to be increasing in usage. Yet it does not replace any semantically close words — *famed*, *famous*

# Communicative need modulates competition in language change

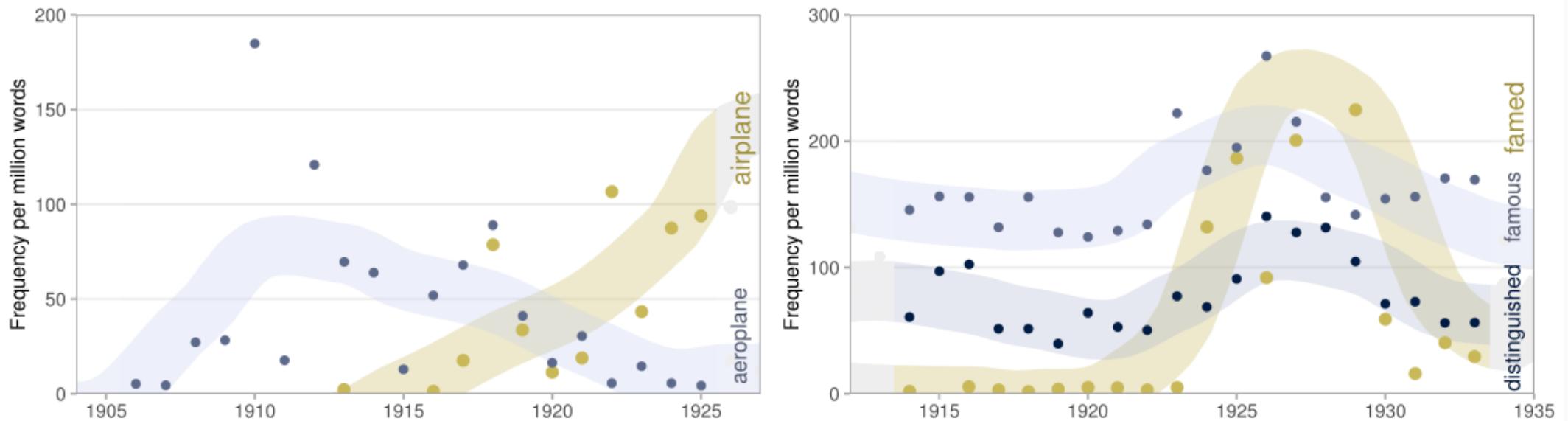
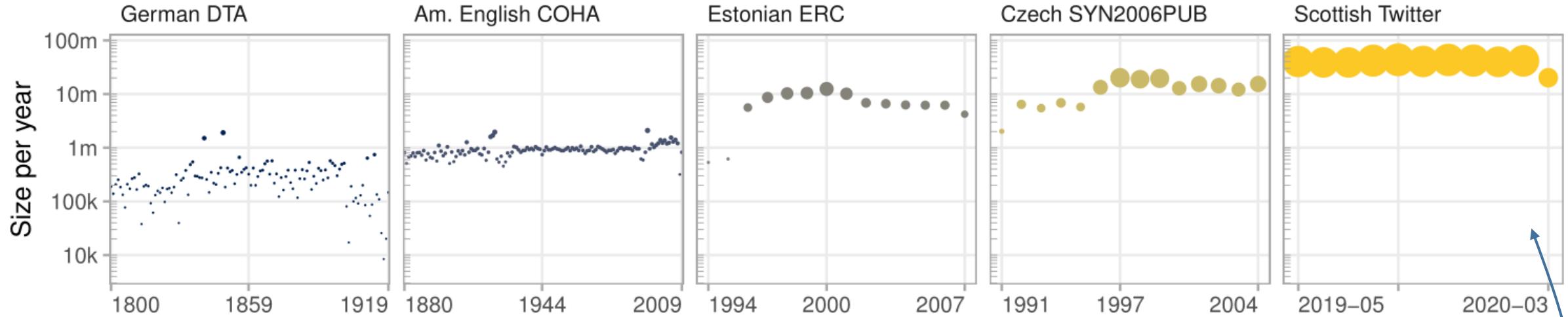


Figure 1: Example time series from the Corpus of Historical American English (COHA). Two decades after the invention of heavier-than-air powered aircraft, *airplane* replaced the initial term *aeroplane* (left side panel; the points are normalized yearly frequencies, with the lines representing smoothed averages for visual aid). Around the same time, *famed* appears to be increasing in usage. Yet it does not replace any semantically close words — *famed*, *famous*

# The corpora



- COHA&DTA: 10-year bins (5 for ERC, Czech, month for Twitter)
- Targets: min +2 log change, occurs min 100x and in enough years



Paul Anderson  
@acereject



Gaul Plancy  
@paul\_glancy

Follow

Ryanair are fly bastards they lure you in with lit 90 quid flights but aw ye want a case? 45 beans. Sit next to yer pal? Tenner mate. Yer grans got legs? Extra score.  
12:58 PM - 10 Jun 2019

If you leave a child in your car during this hot Glasgow weather please ensure a window is open so they can at least have a fag

- Need:
- A model of competition
- A model of communicative need

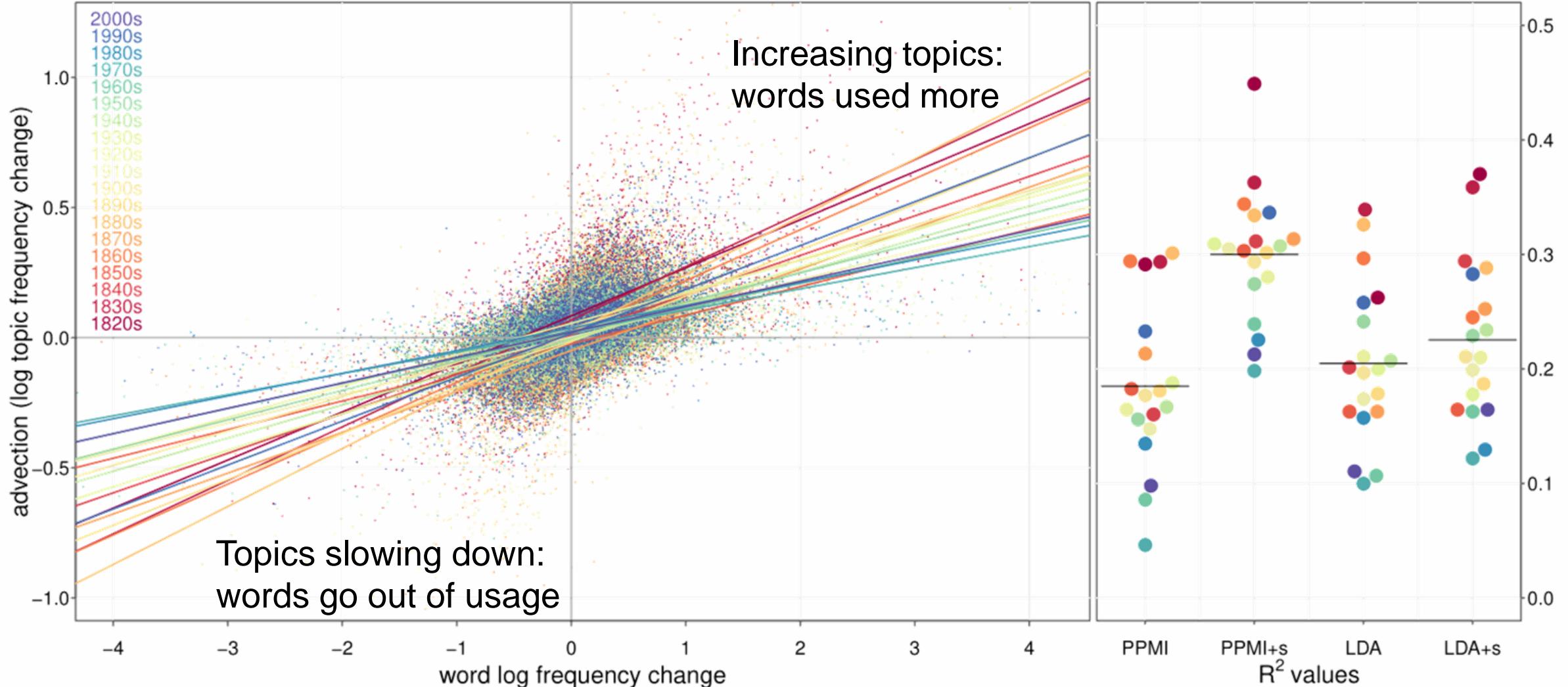
# A model of communicative need

- Karjus, Blythe, Kirby, Smith 2020, Quantifying the dynamics of topical fluctuations in language. *Language Dynamics and Change* <https://doi.org/10.1163/22105832-01001200>
- Idea: see how much the topic of a target word changes (weighted mean of the log frequency changes of the relevant topic (context) words of the target)
- Discourse topic prevalence ~ how much something needs to be talked about ~ communicative need
- Topics as the latent flow of language, dragging words along
- *advection* - the transfer of matter (or heat) by the flow of a fluid

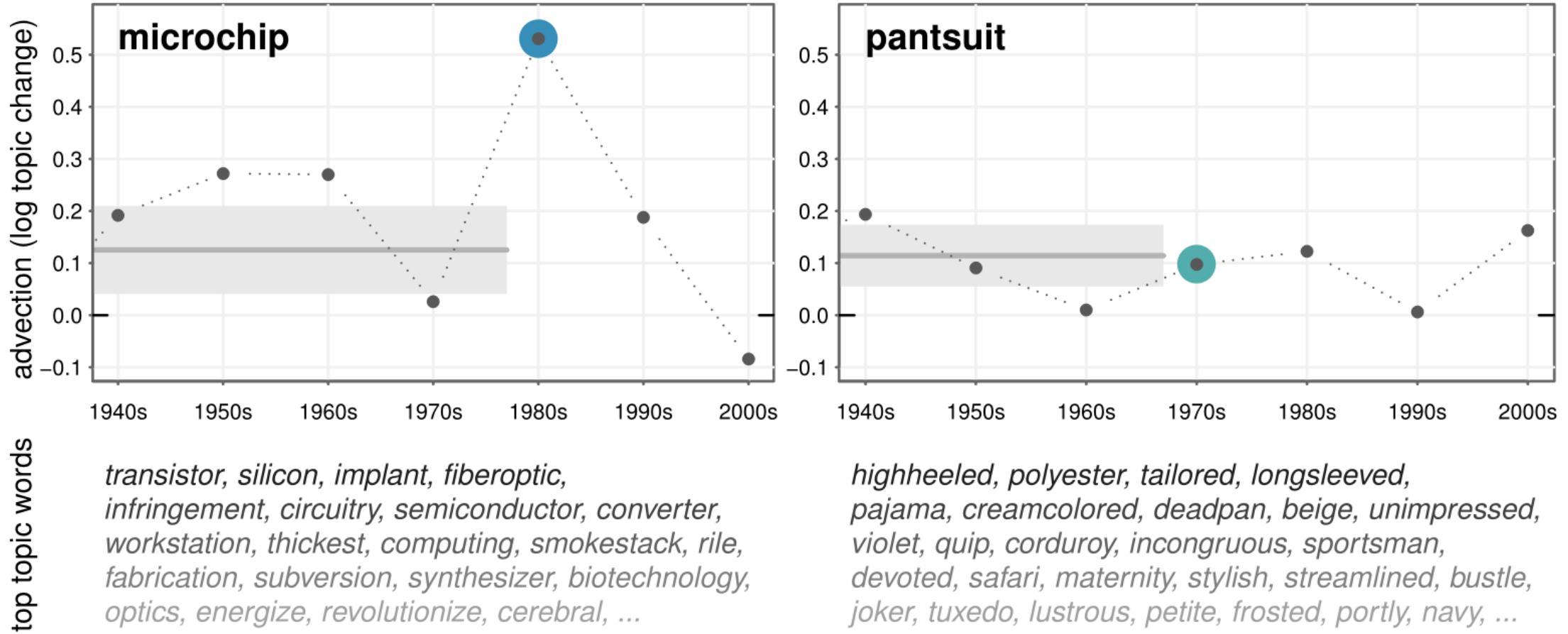
# A model of communicative need

- *advection* - the transfer of matter (or heat) by the flow of a fluid

# Quantifying the dynamics of topical fluctuations in language



# Advection a proxy to communicative need

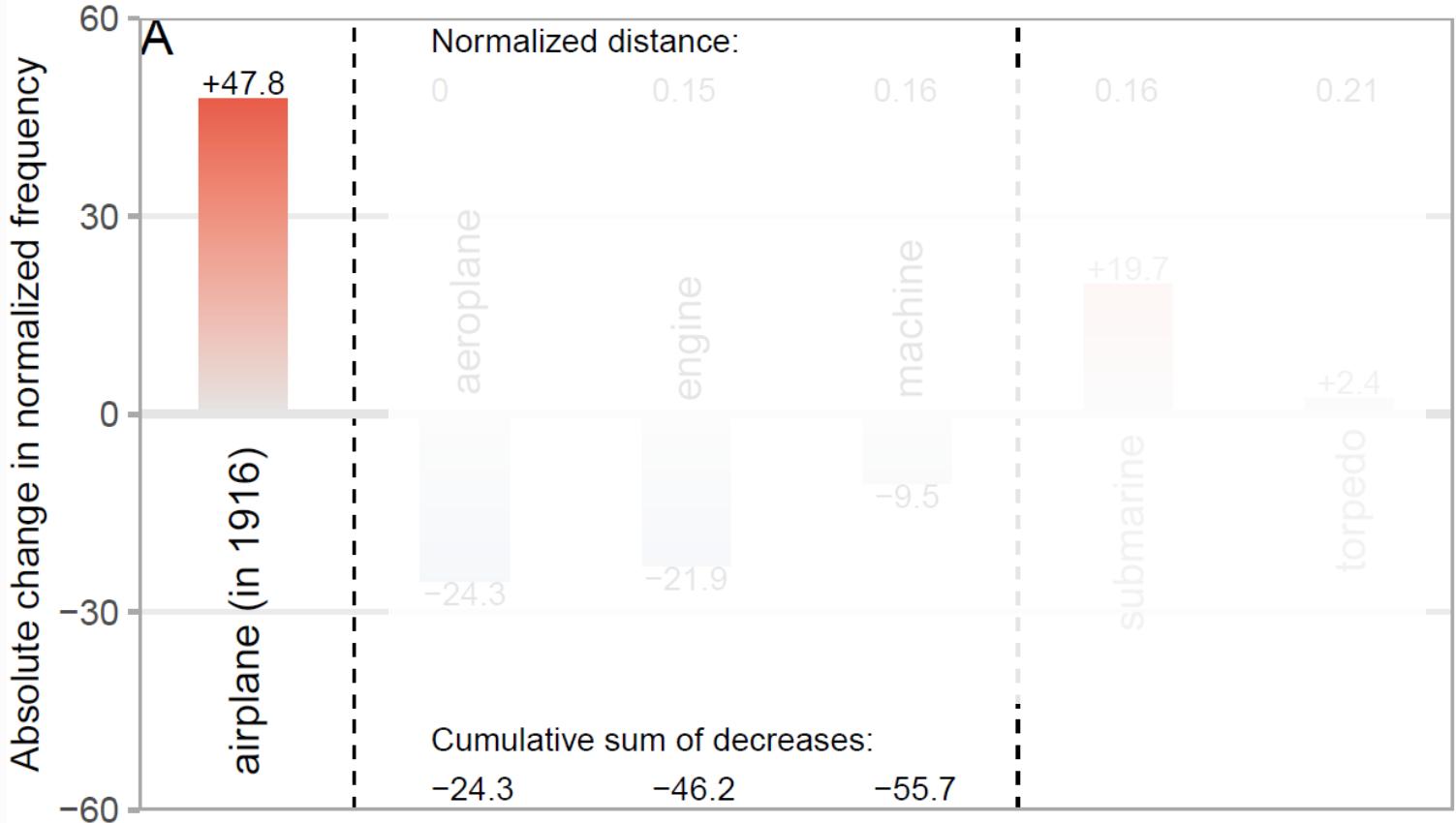


# A model of linguistic competition

- Meaning from word embeddings; **equalization range**: norm. cosine distance from target where the sum of (normalized) frequency decreases match the increase of the target
- Normalized corpus frequencies sum to 1
- Increase somewhere => decrease somewhere else
- A realistic model of language? Yes: time is finite and learning pressure biases for simpler lexicons. Can't have infinitely many words.
- Semantics: inferred from LSA, trained for each target word based on (PPMI-weighted) co-occurrence matrix of the preceding time bin, fit target vector into this model; this yields neighbours of the *position* where the new word will appear in

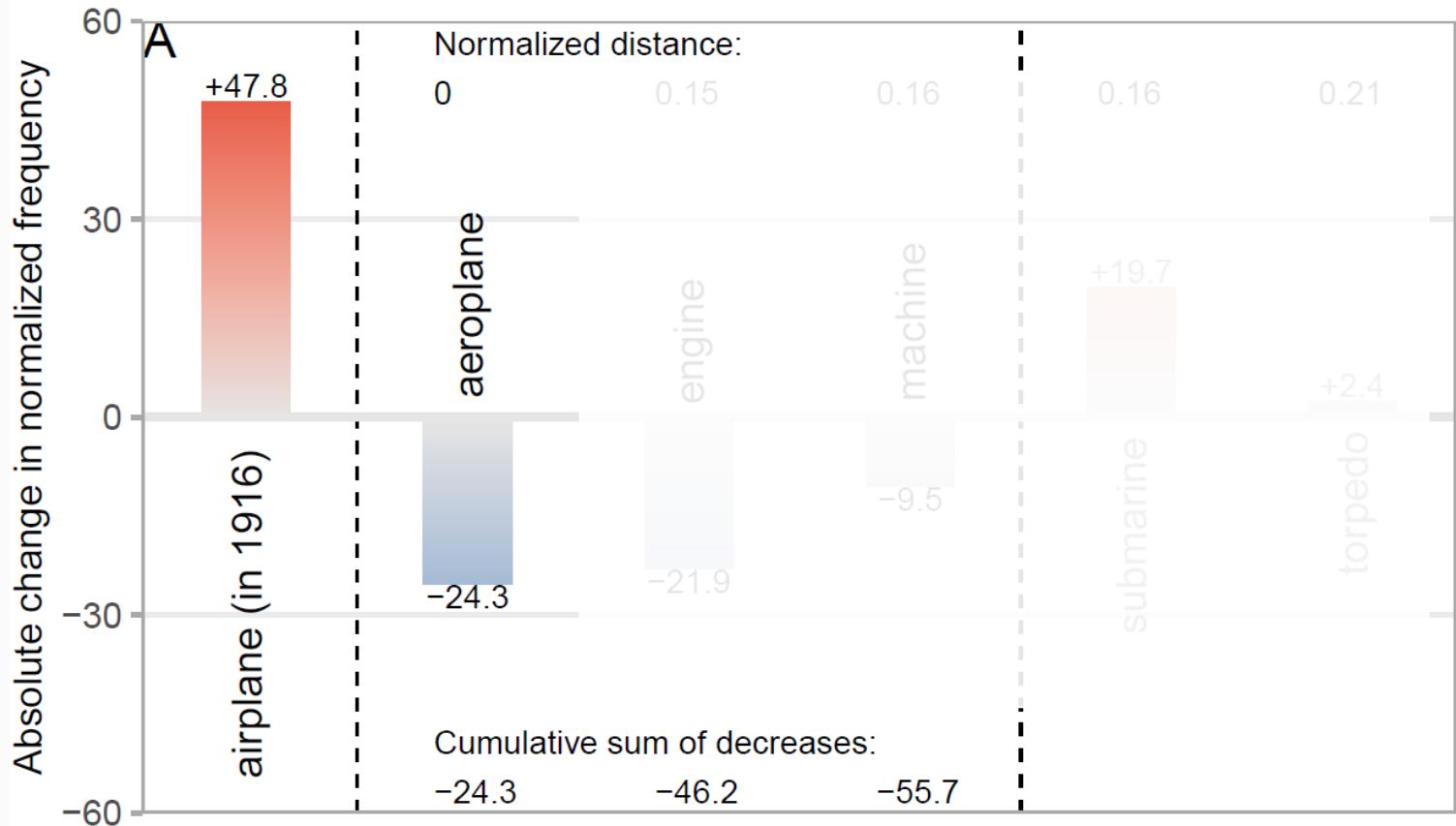
# A model of linguistic competition

- Meaning from word embeddings; **equalization range**: norm. cosine distance from target where the sum of (normalized) frequency decreases match the increase of the target



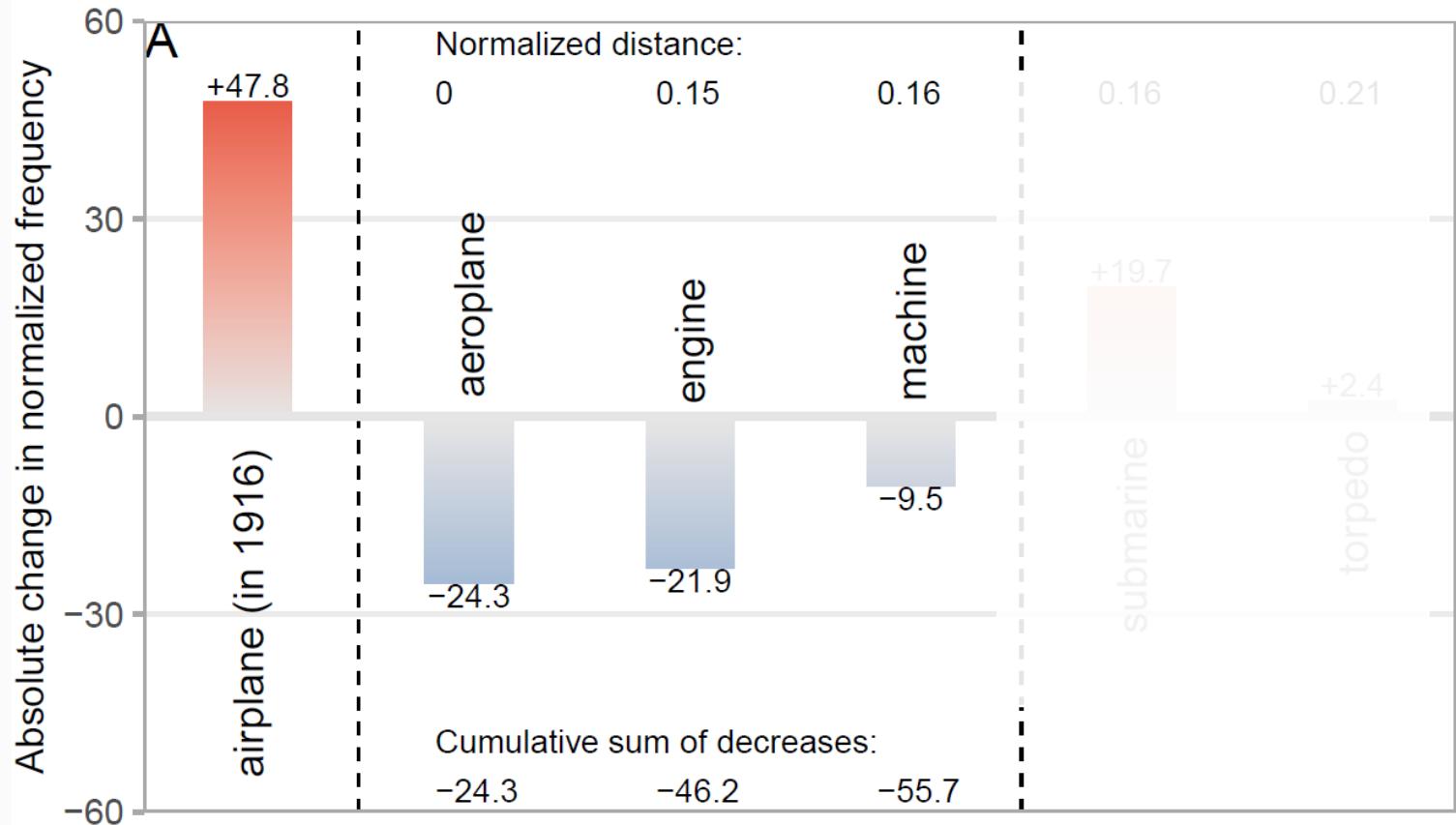
# A model of linguistic competition

- Meaning from word embeddings; **equalization range**: norm. cosine distance from target where the sum of (normalized) frequency decreases match the increase of the target



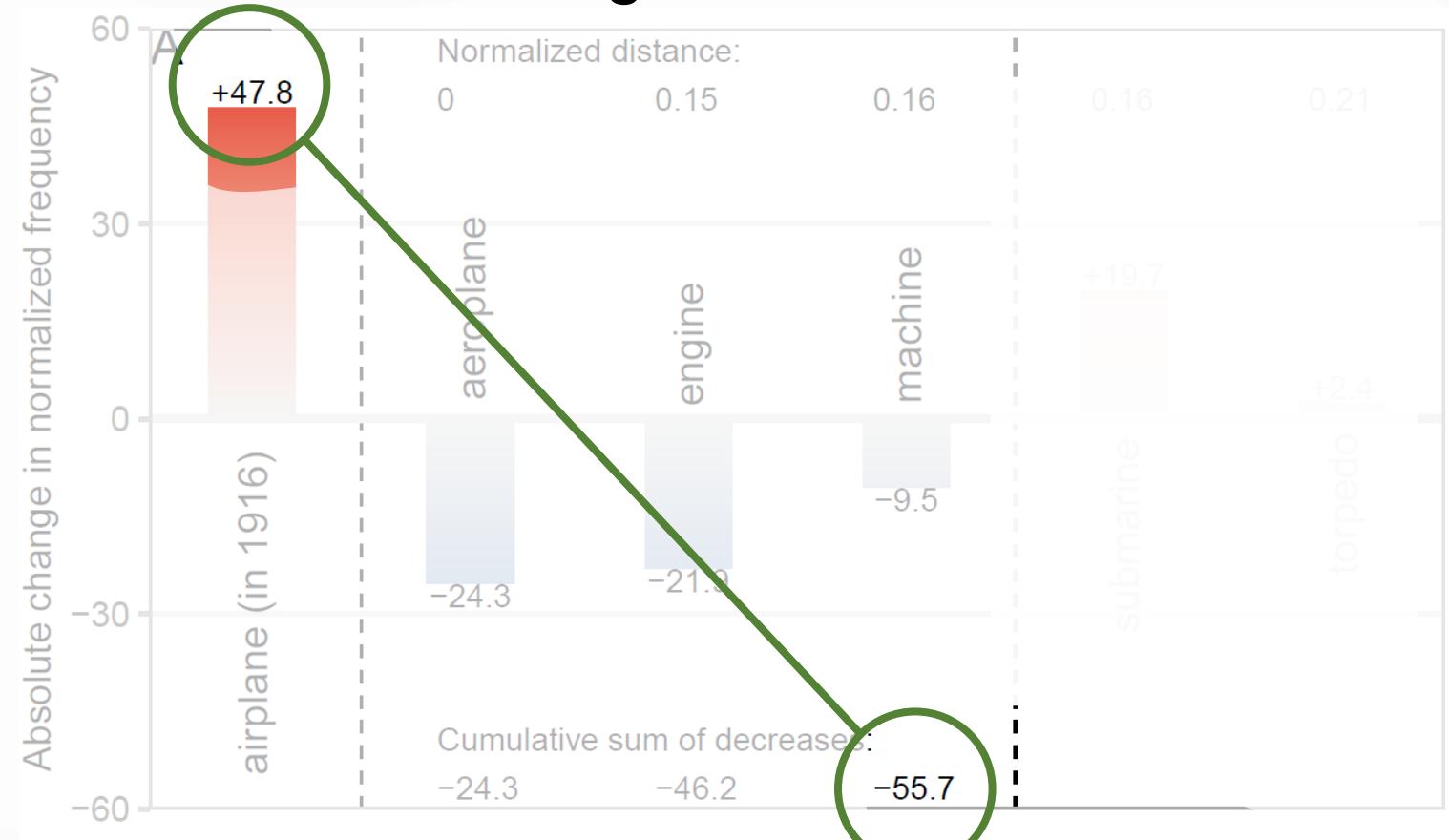
# A model of linguistic competition

- Meaning from word embeddings; **equalization range**: norm. cosine distance from target where the sum of (normalized) frequency decreases match the increase of the target



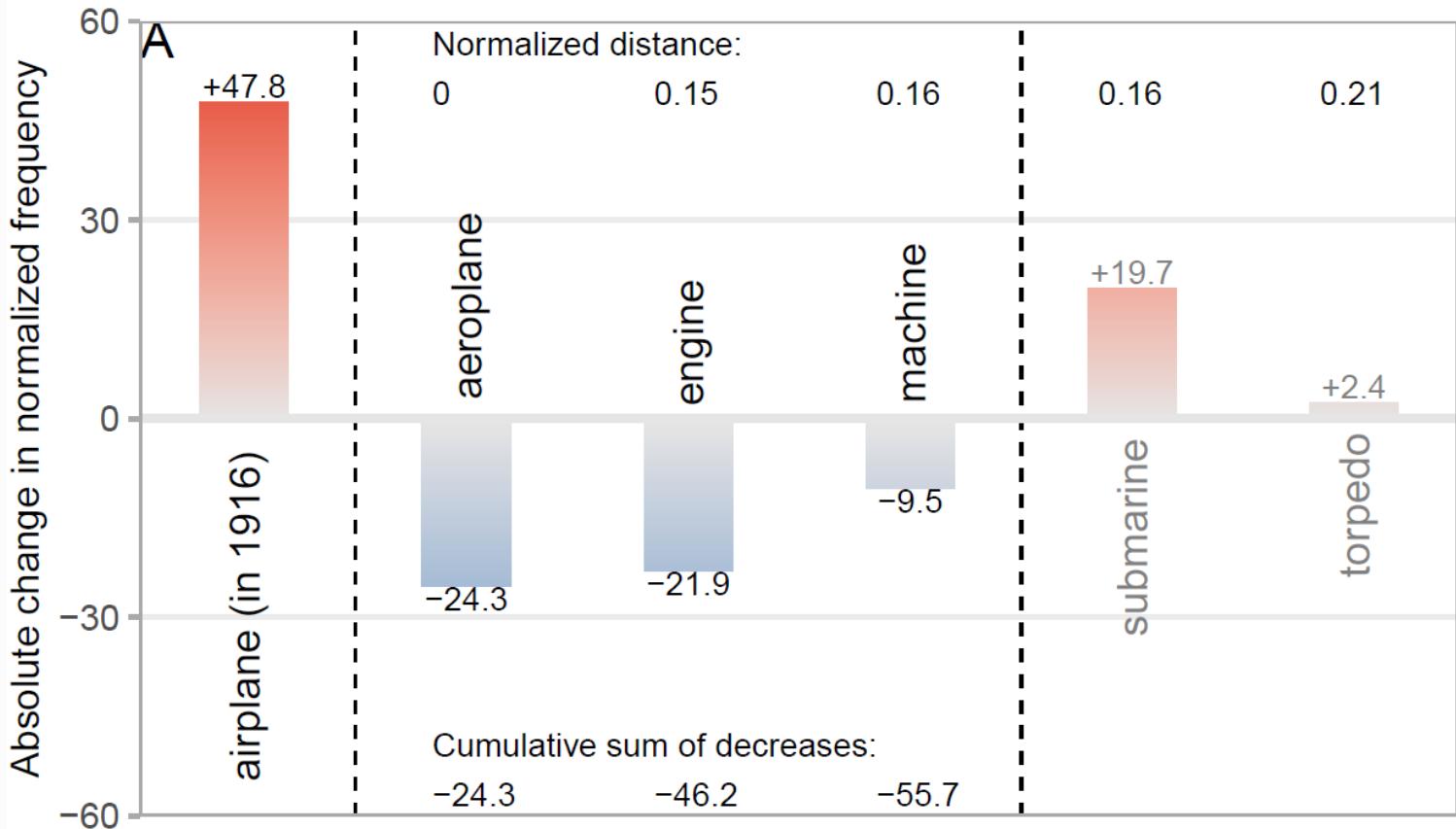
# A model of linguistic competition

- Meaning from word embeddings; **equalization range**: norm. cosine distance from target where the sum of (normalized) frequency decreases match the increase of the target

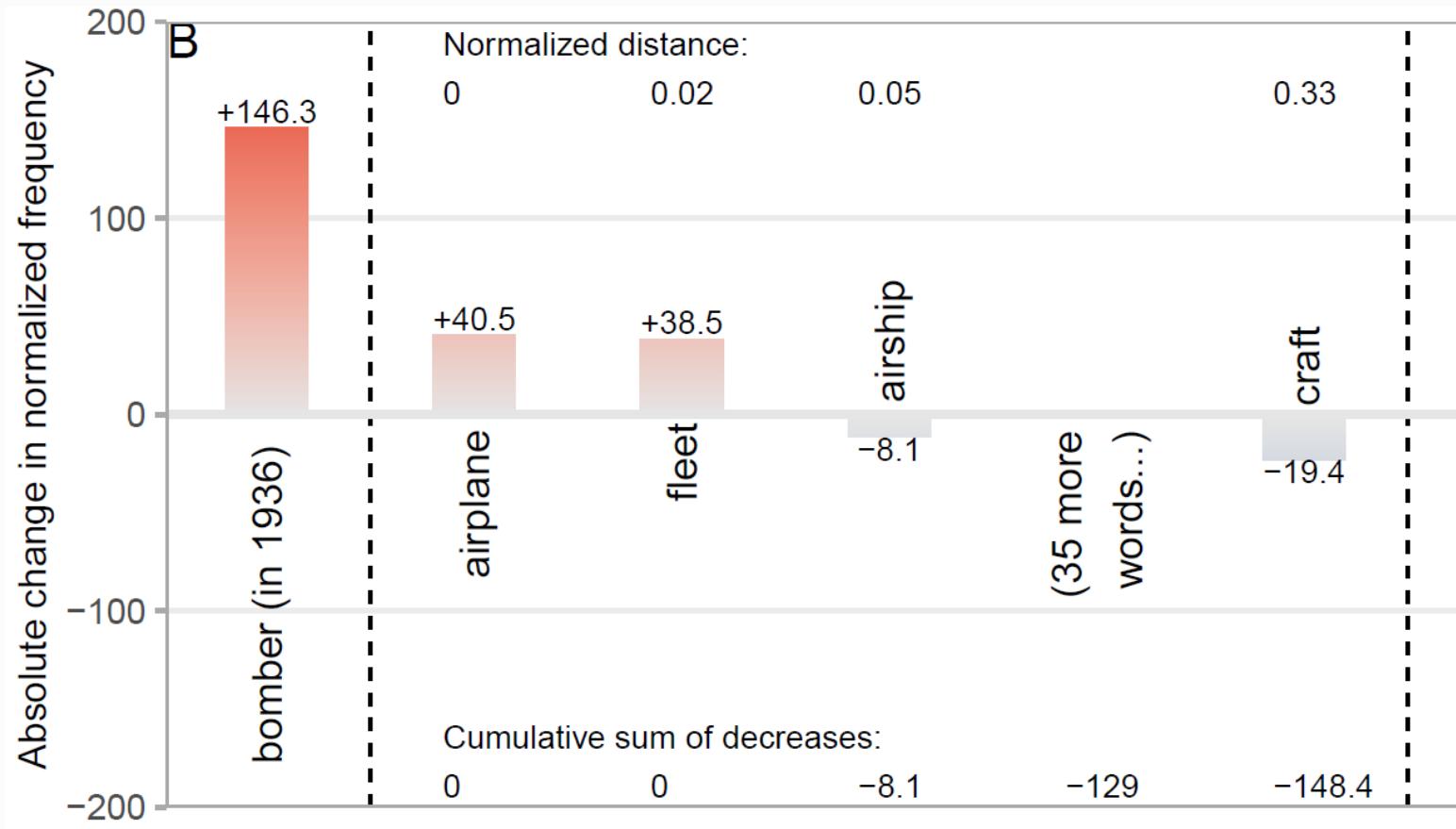


# A model of linguistic competition

- Meaning from word embeddings; **equalization range**: norm. cosine distance from target where the sum of (normalized) frequency decreases match the increase of the target



# A model of linguistic competition

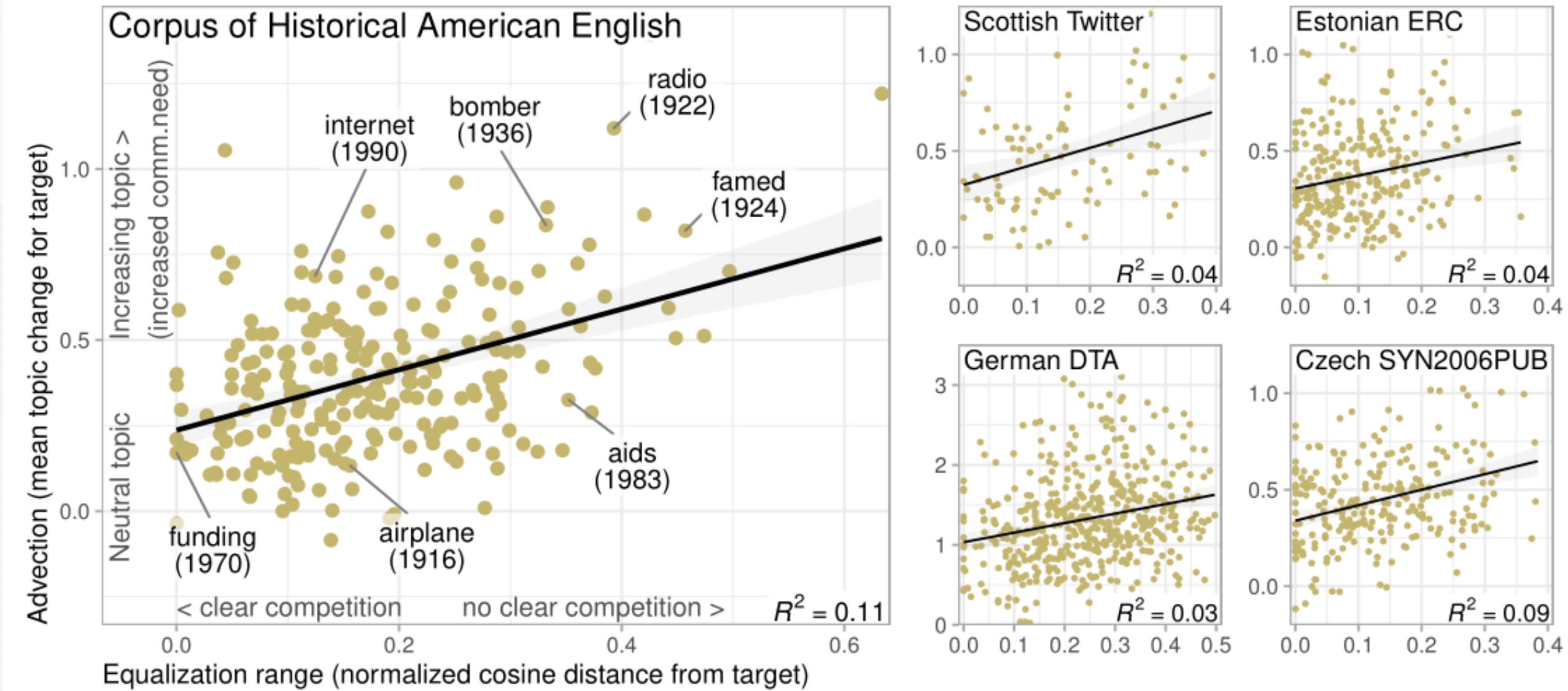


# Important

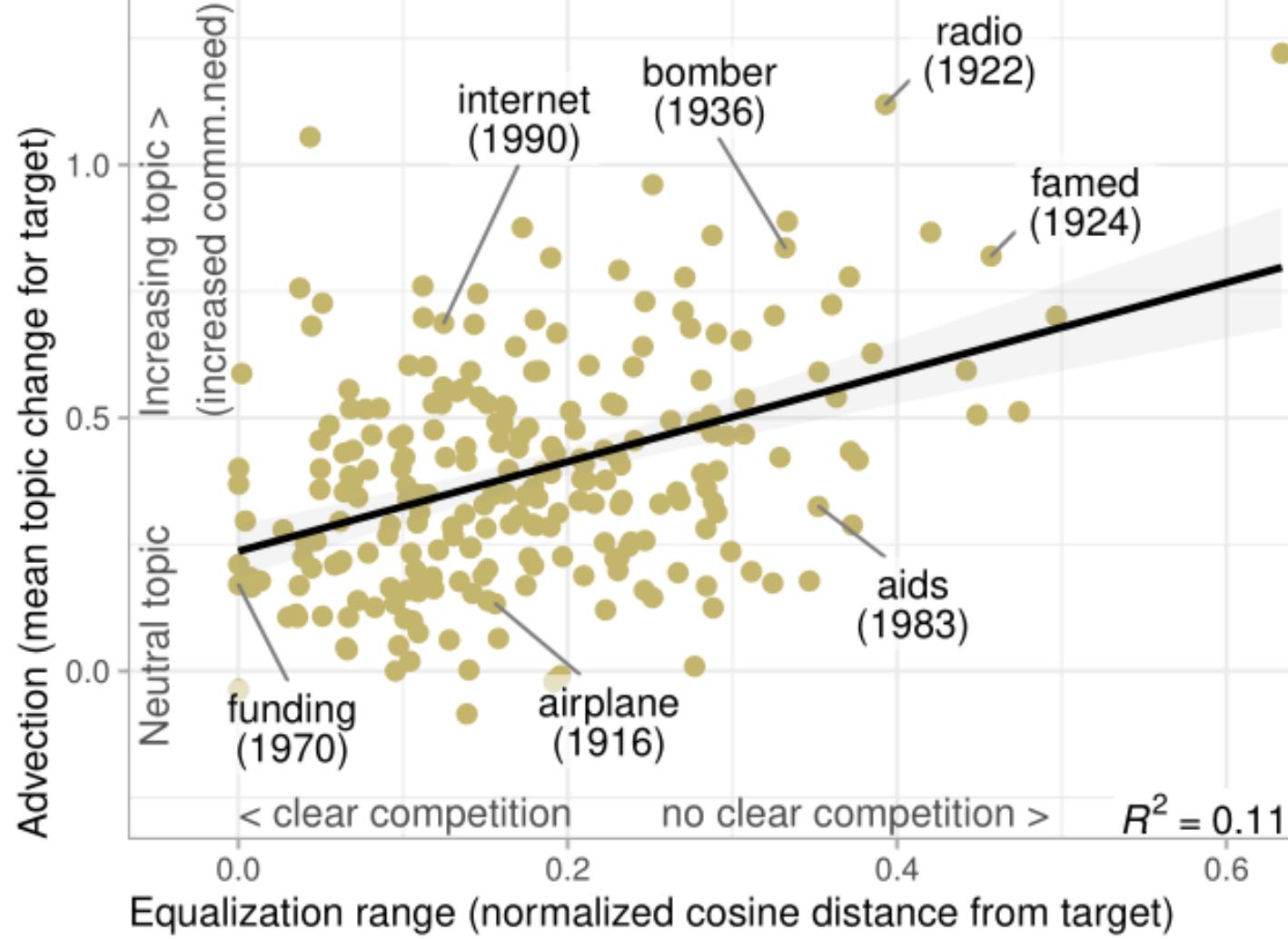
- Both models based on lists of words, but decorrelated:
  - advection: weighted list of associated, co-occurring words (1<sup>st</sup> order similarity)
  - competition: list of all words, ordered by embedding cosine similarity (2<sup>nd</sup> order similarity), minus any words in the advection list for a given target
- Necessary, but can weaken the competition model accuracy, if closest neighbours (~synonyms) also co-occur with target:
  - *airplane / aeroplane airship aerial propeller balloon engine machine submarine biplane wireless torpedo*

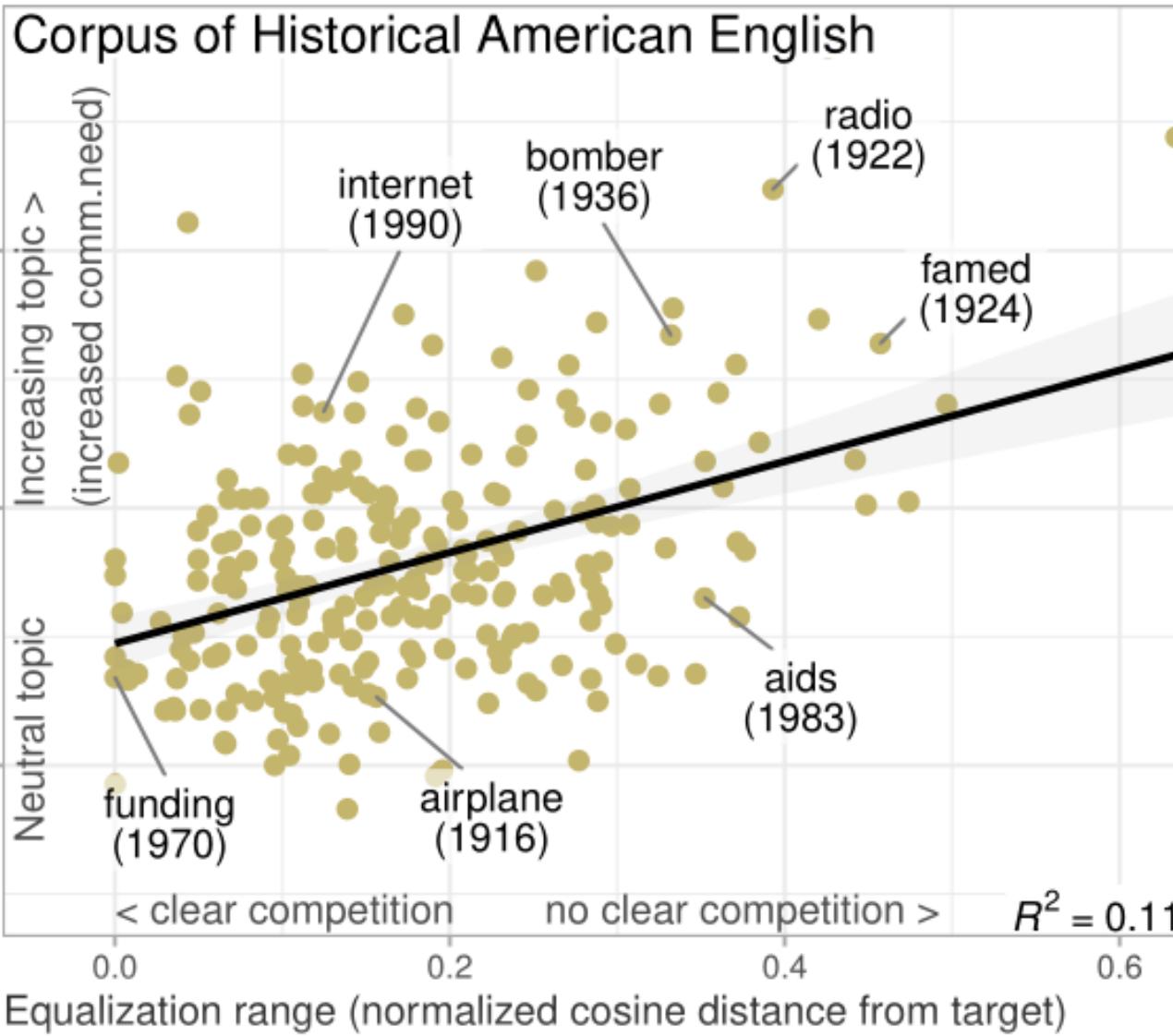
# Results

- Topical advection (proxy to communicative need) correlates with
- Equalization range (proxy to extent of competition)

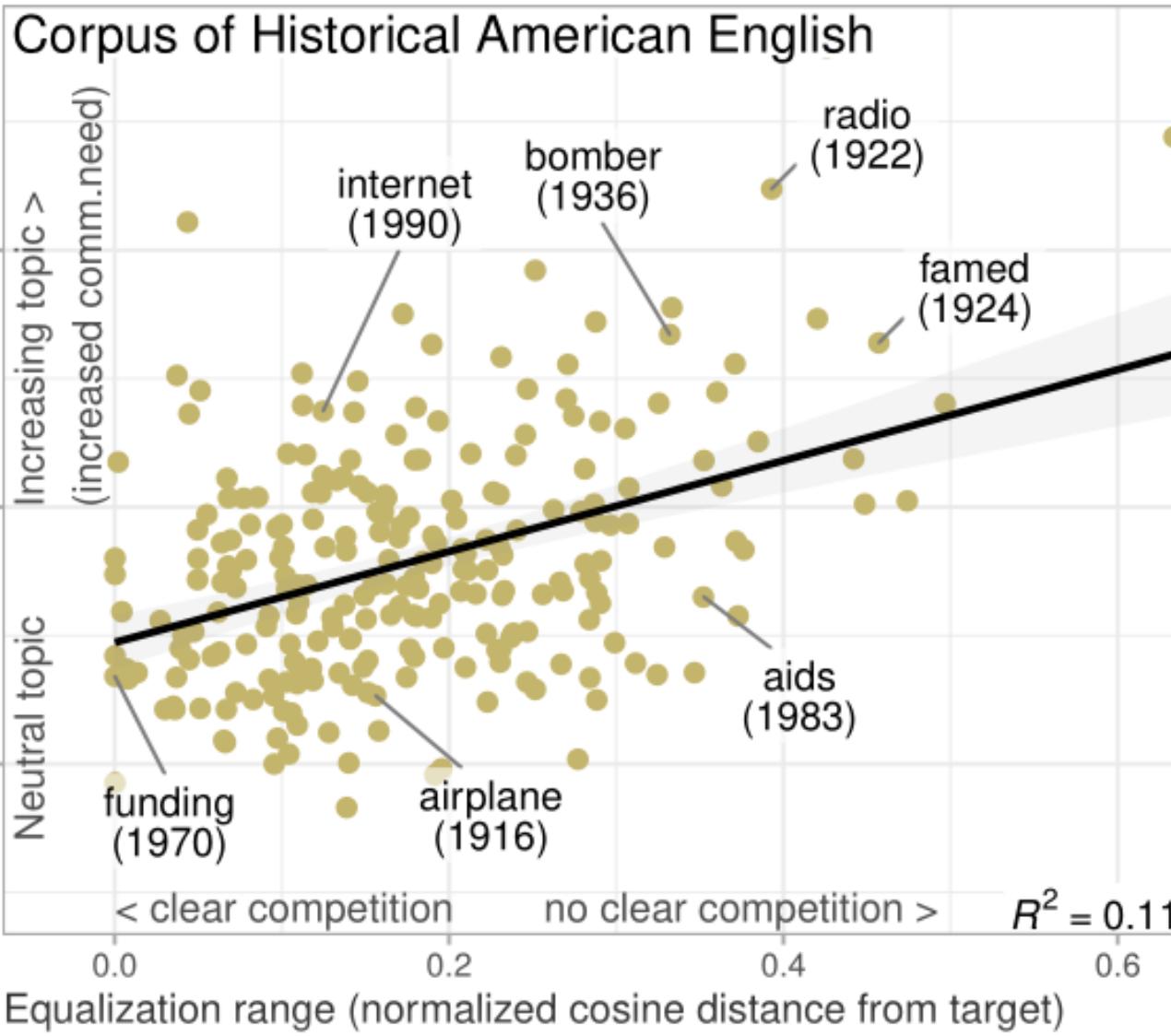


# Corpus of Historical American English





- Lower communicative need:  
competition more likely

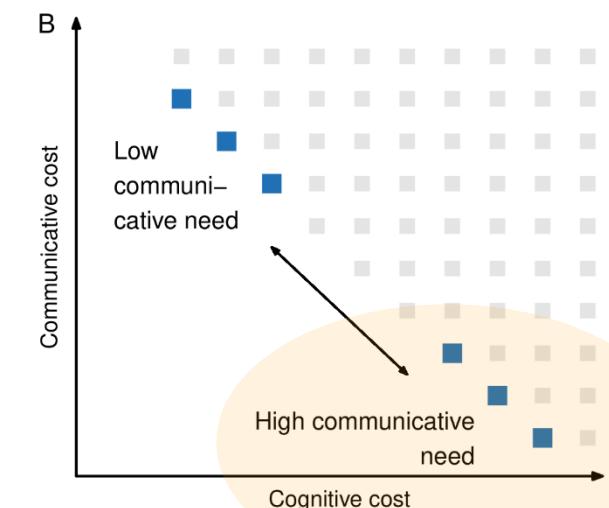
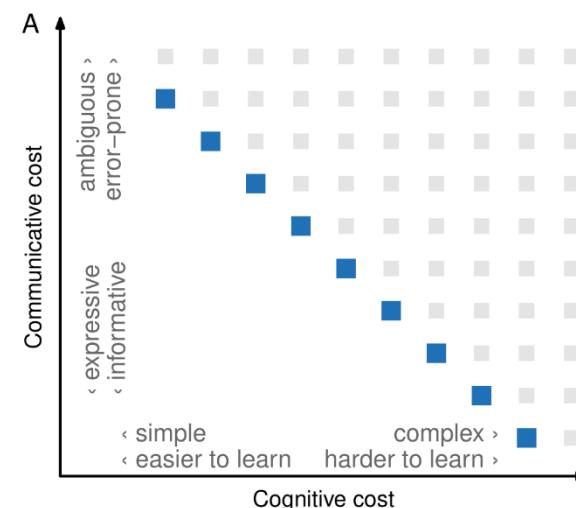


- High communicative need: similar words more likely to coexist

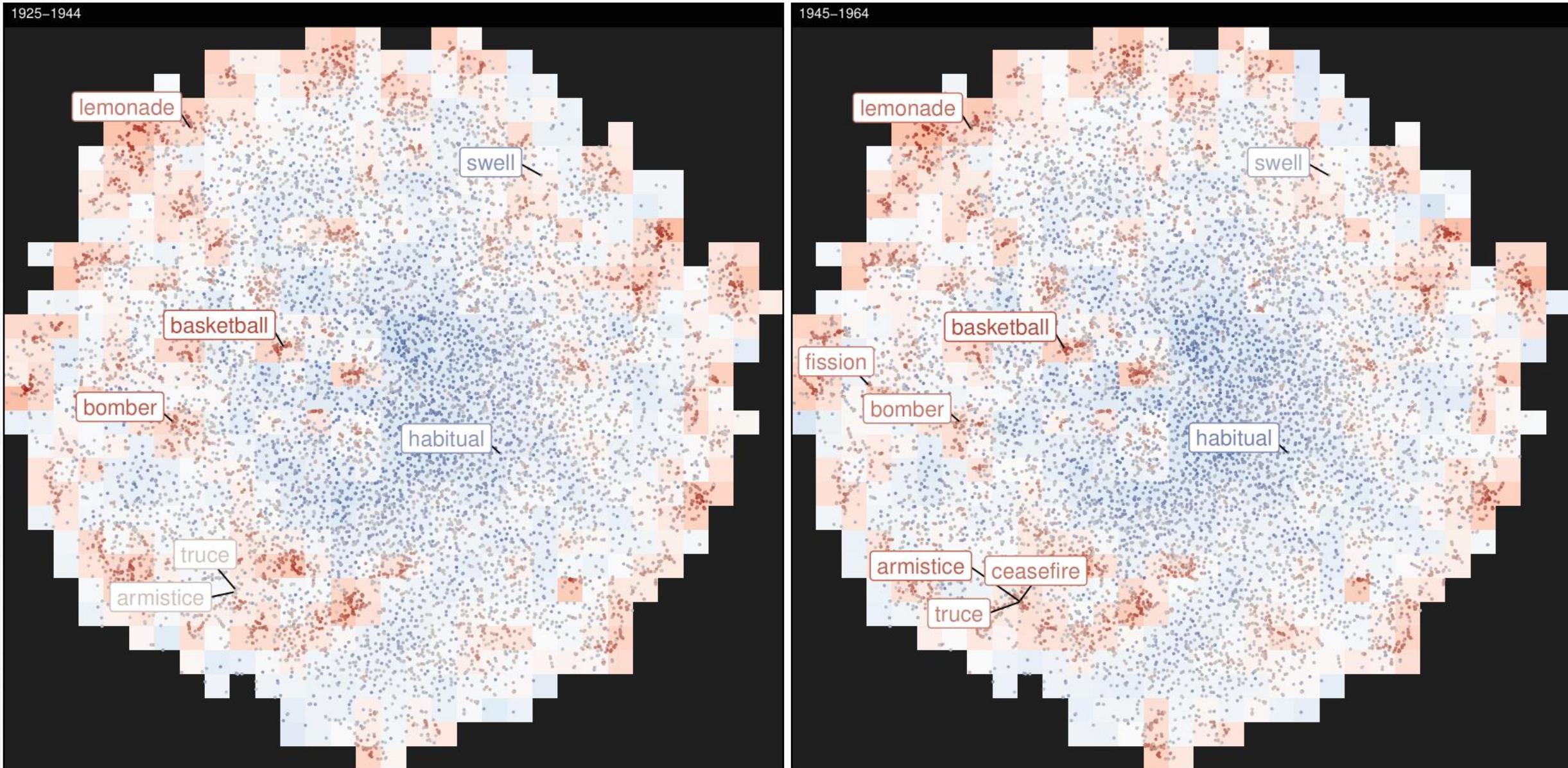
- Lower communicative need: competition more likely

# Discussion

- Communicative need, after controlling for a slew of other lexicostatistical variables, describes a small amount of variance in competitive interactions
- Small effect, but consistent across languages and genres
- High communicative need facilitates the co-existence of similar words (more complex lexical subspace)
- But: this method has some shortcomings...



# Discussion

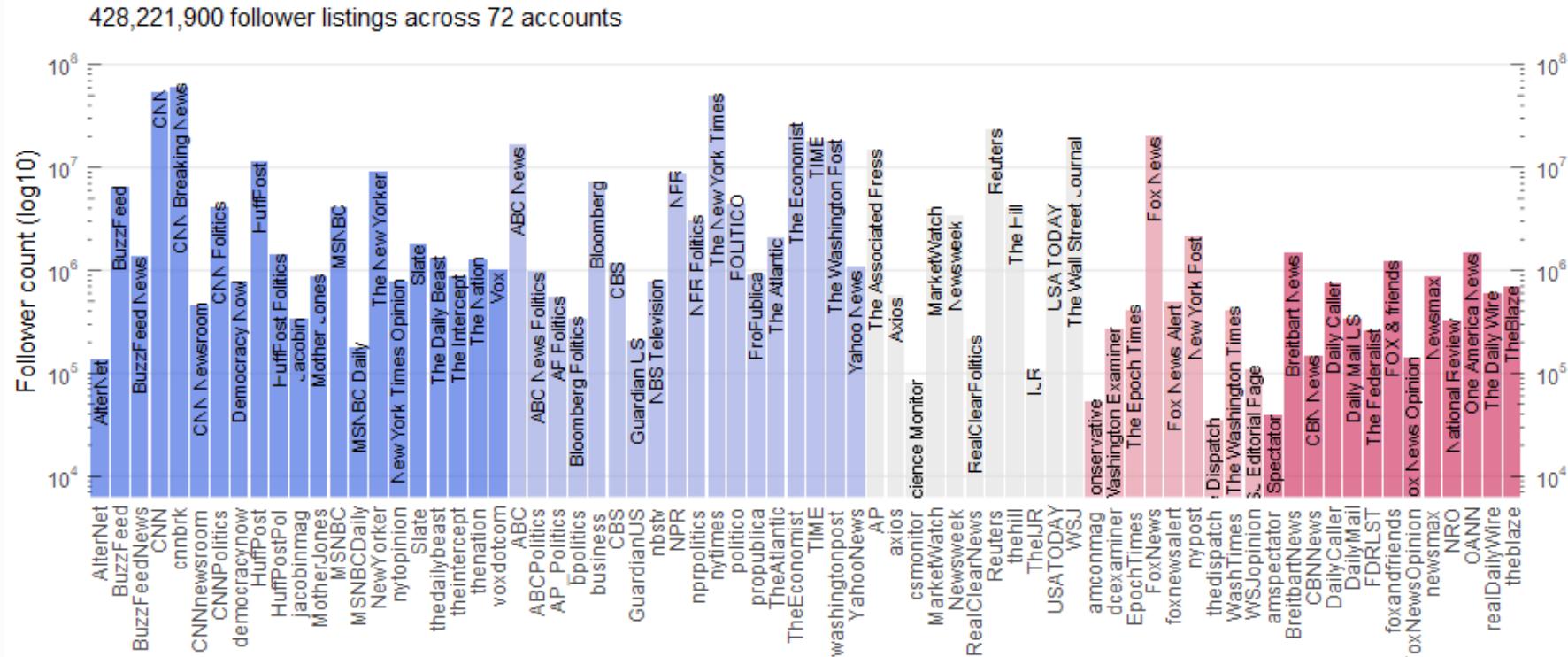


# Conclusions & future research

- Converging typological, experimental and corpus evidence supports the argument for the role of communicative need from earlier cross-linguistic research
- Many reasons why languages change; one is adaption to the changing needs of their speakers
- Future: Iron out the competition model (also apply to domains other than language)
- Future: semantic divergence on social media
- Future: application of complexity-informativeness other domains

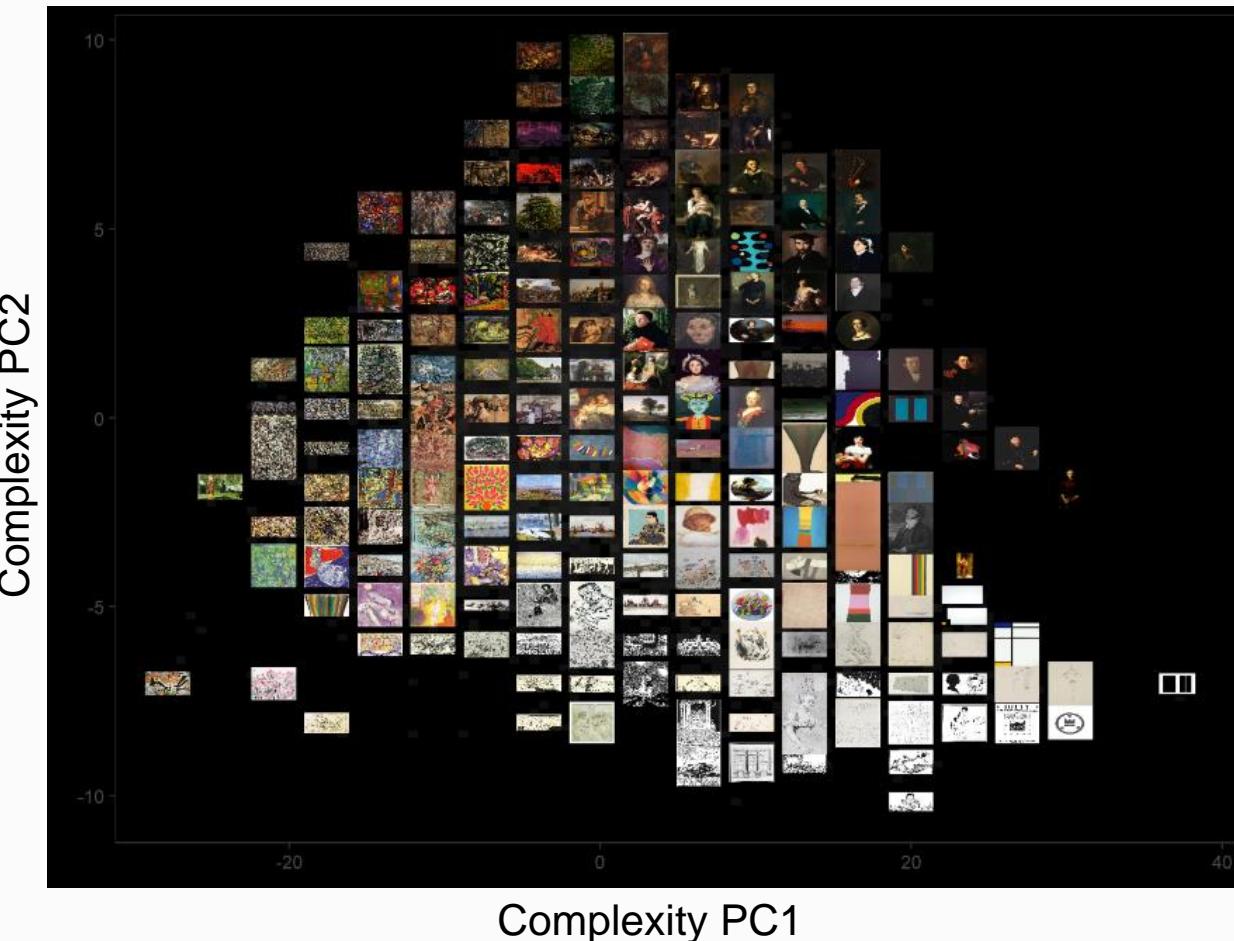
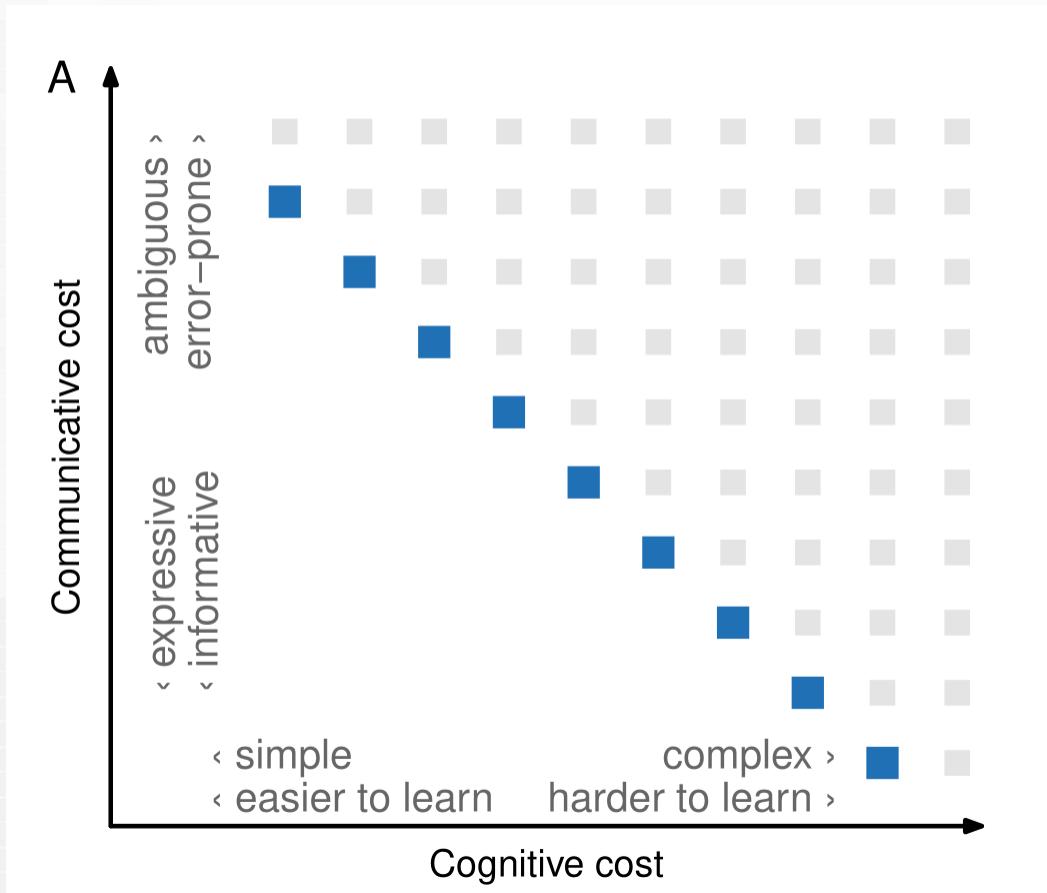
# Future: semantic divergences

- ...and semantics-driven misunderstanding in a polarizing world, using social media data (with Christine Cuskley, Newcastle)
- Infer political alignment of Twitter users from follower relations
- Compile a tweets corpus, grouped by alignment
- Model meaning using models from Lexical Semantic Change



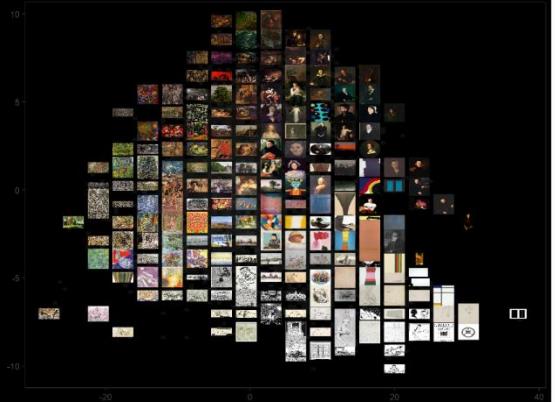
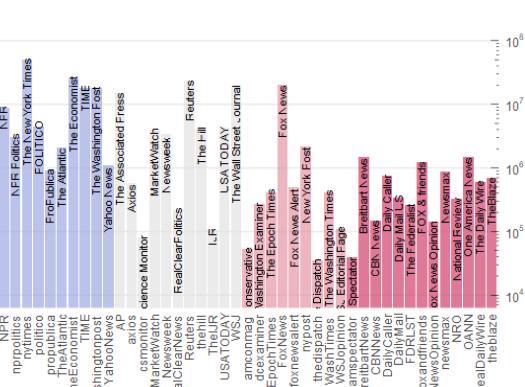
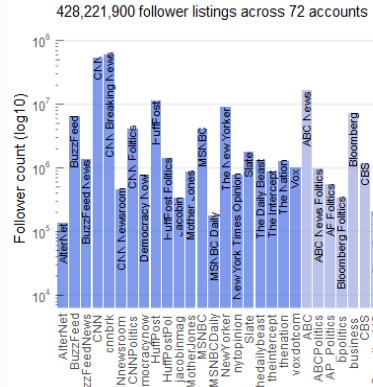
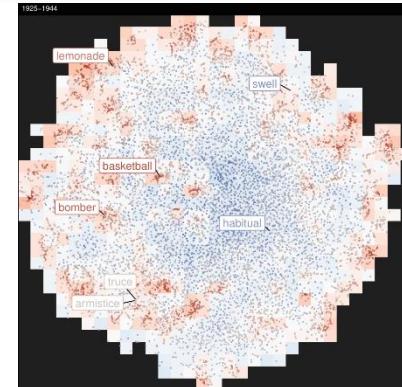
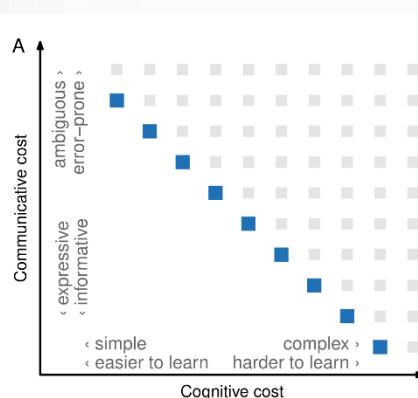
# Future: compression ensembles

- Would the complexity-informativeness approach be applicable/useful in domains of cumulative cultural evolution other than language? (with Maximilian Schich/Tallinn, Sebastian Ahnert/Cambridge/Turing)



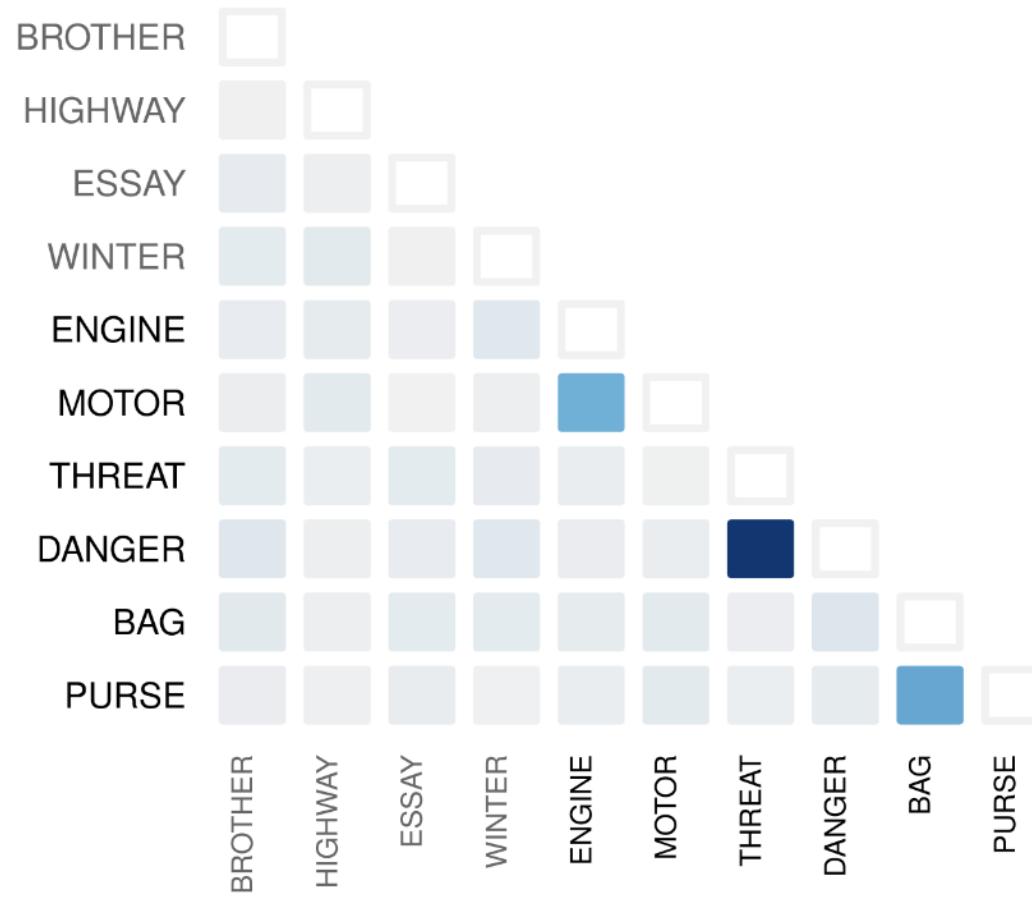
# Conclusions & future research

- Converging typological, experimental and corpus evidence supports the argument for the role of communicative need from earlier cross-linguistic research
- Many reasons why languages (cultures?) change; one is adaption to the changing needs of their speakers
- Future: Iron out the competition model
- Future: research into semantic divergence on social media
- Future: application of complexity-informativeness to art, etc.



# Appendix

A: Semantic similarity between stimuli



B: Formal similarity between stimuli

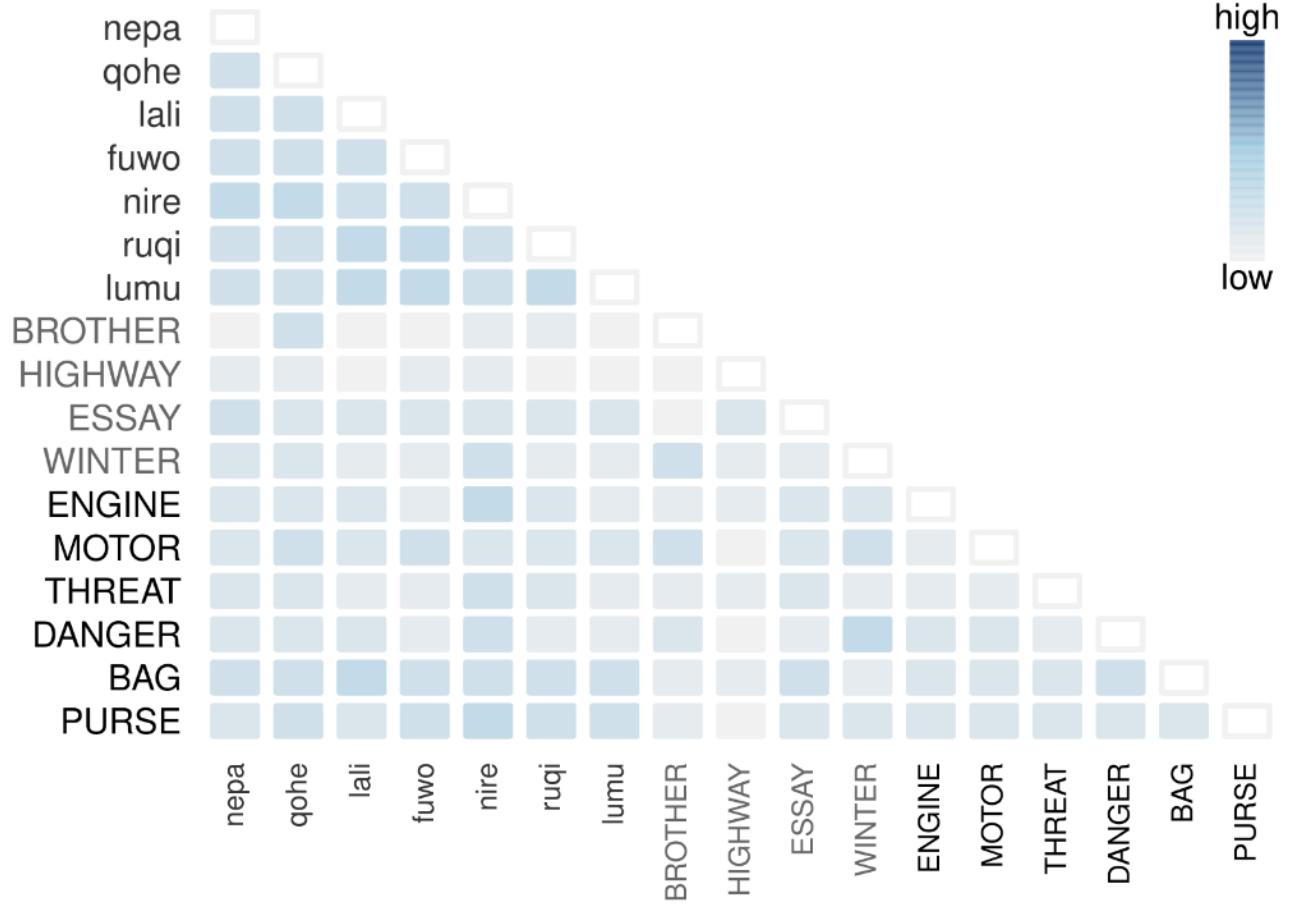
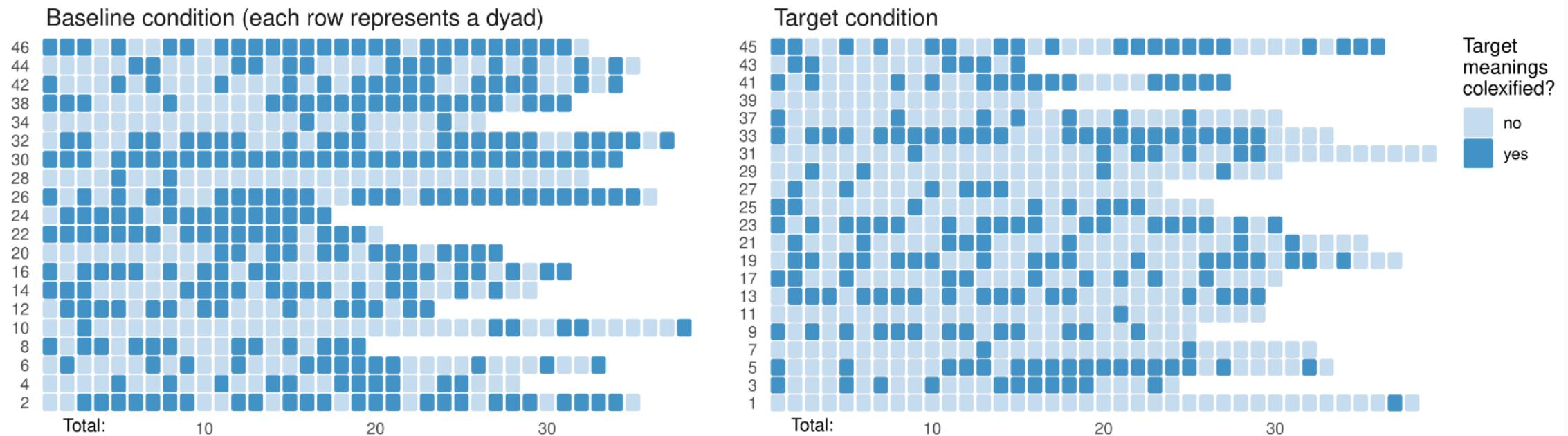


Figure 2: The meanings and signals used in one game (dyad no. 21 in Experiment 1). The left side A panel illustrates the meaning space: only the target pairs have high similarities (dark blue), with low similarities between all other meanings. The diagonal (self-similarity) is marked by white squares. The right side B panel shows the similarity of form (as inverse of edit distance), of both the meanings (in caps) and signals (lowercase). The stimuli in our experiments are generated in a way that ensures only target meanings are semantically similar to one another, and that form similarity remains low across the board.



	colexification ~	Estimate	SE	<i>z</i>	<i>p</i>
intercept (baseline condition, mid-game)	-0.22	0.37	-0.59	0.56	
+ condition (target)	-0.52	0.51	-1.03	0.3	
+ round	1.02	0.27	3.84	<0.01	
+ condition (target) × round	-1.17	0.37	-3.17	<0.01	

Table 2: The fixed effects from the mixed effects regression model applied to data from Experiment 1, predicting the value of the colexification variable (reference level: "no") by the interaction between condition (reference: baseline) and round number. The latter is included to account for progress over the course of the experiment. The results indicate a statistically significant difference in participant behavior between the two conditions, supporting the hypothesis that communicative need can drive lexification preferences above and beyond conceptual similarity.

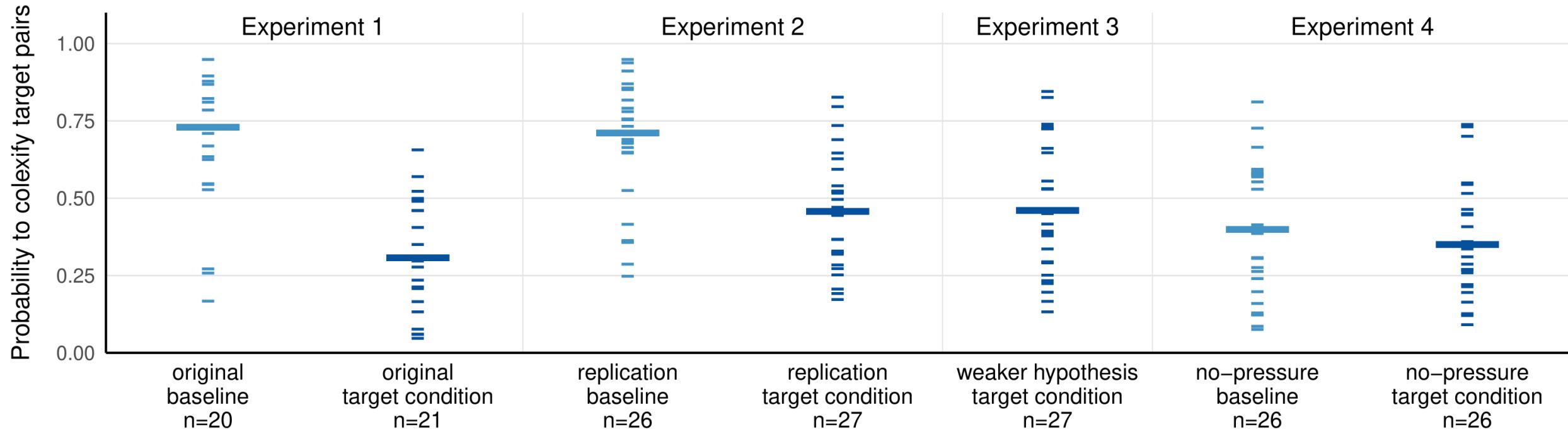
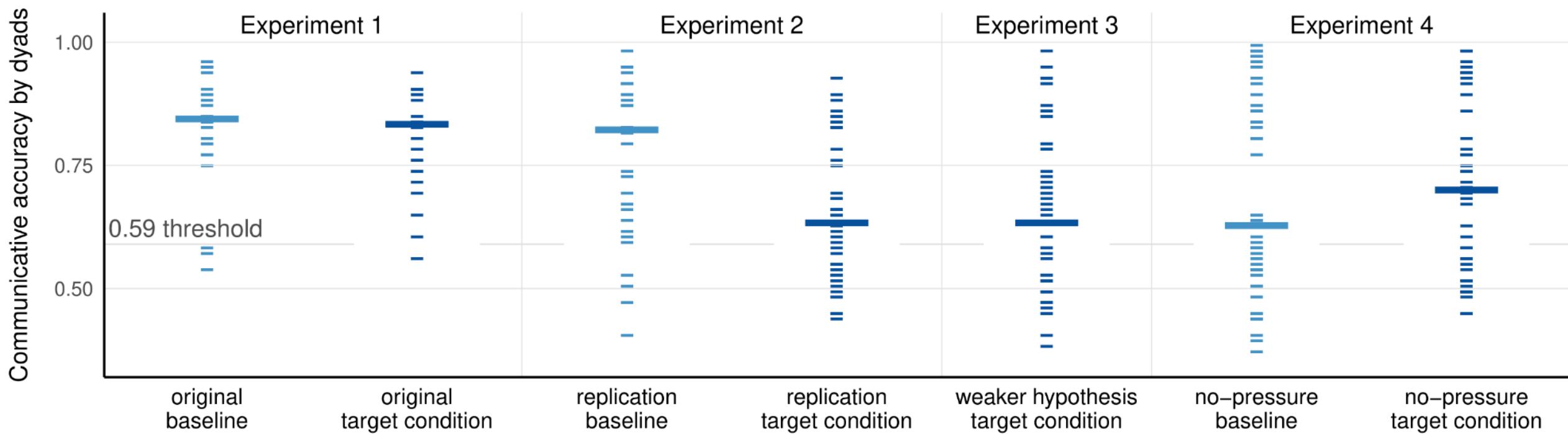


Figure 5: Estimated probabilities of colexifying similar-meaning target pairs by the end of the game. Each notch is one dyad, the wider bars are medians (unlike in our full statistical analysis, meanings are collapsed here within dyads for clearer visualization purposes). Target conditions (with manipulated communicative need) are in darker blue. In Experiment 1 (left) as well as its replication (Experiment 2), dyads were more likely to end up colexifying similar meanings like TRIP and JOURNEY in the baseline condition, but less so in the target condition, when faced with communicative need to distinguish them. This graph displays all dyads across all experiments (173 dyads or 346 participants in total) and will be referred back to in later sections.



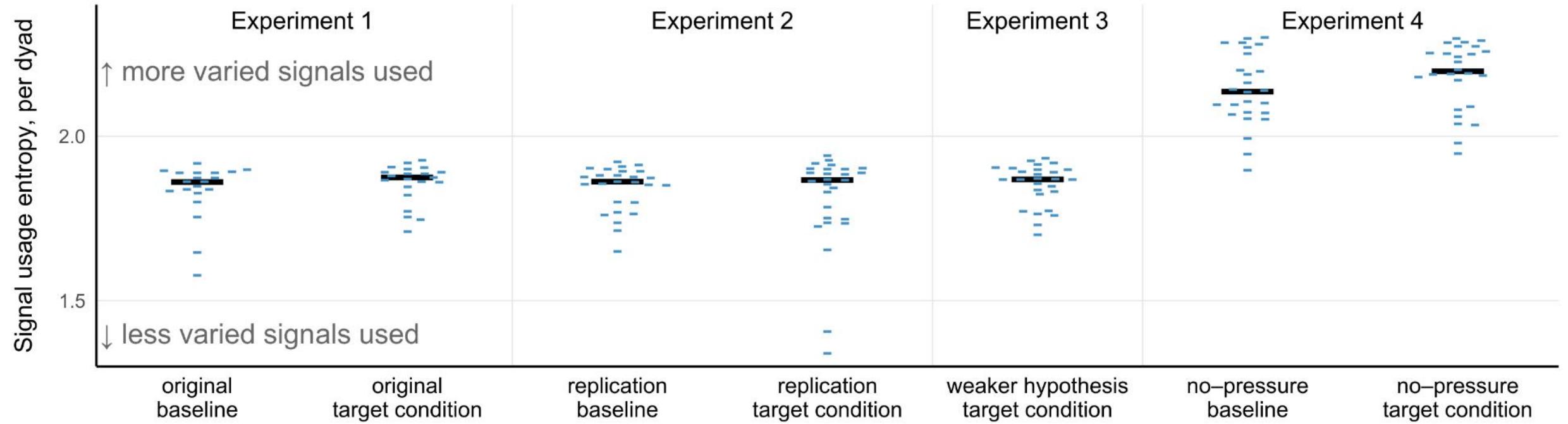


Figure 8: Signaling entropy across all dyads (blue notches), in the post-burn-in part of the game. Low values on the vertical axis indicate fewer signals were consistently used, higher values indicate a larger variety of signals were used by a given dyad. Overlapping values are pushed slightly aside horizontally to ensure visibility. The wider bars represent medians. Dyads in Experiment 4 generally made use of most of the extended signal space (10 instead of 7), setting it apart from the rest of the experiments.

Expm no. 38, baseline condition, 96%, counts

WARRIOR	2				7	
THEFT			9			
STATE				9		
RHYTHM					9	
TASK	2	4	2		1	
JOB		9				
PAIR	8					
COUPLE	10					
SHORE		7		1		
COAST		10				

neme quto nopo fita mefa mumi honi

neme quto nopo fita mefa mumi honi

neme quto nopo fita mefa mumi honi

Expm no. 38, baseline condition, 96%, PPMI

WARRIOR						3.13
THEFT				3.03		
STATE					3.17	
RHYTHM						3.32
TASK	0.07	1.62	0.86		0.32	0.87
JOB		2.79				1
PAIR	2.17					1
COUPLE	2.17					1
SHORE		2.05		0.17		1
COAST		2.24				1

neme quto nopo fita mefa mumi honi

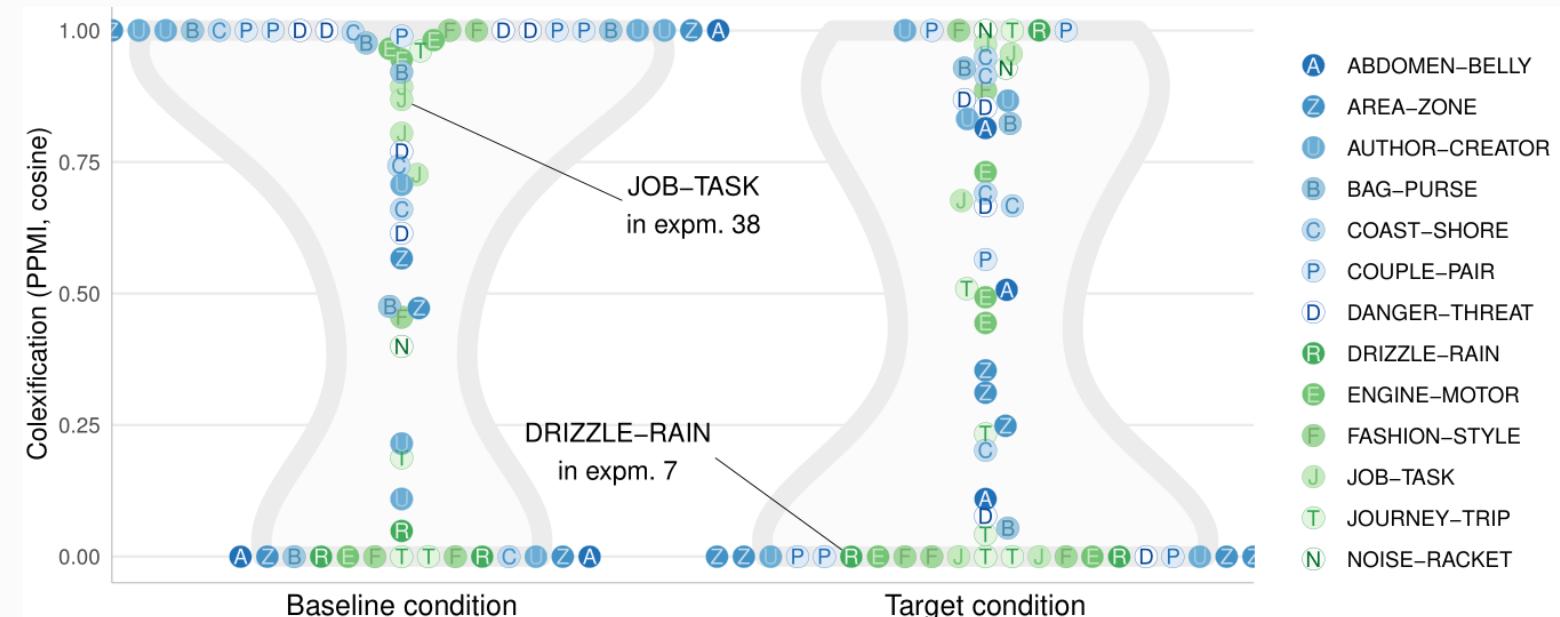
neme quto nopo fita mefa mumi honi

neme quto nopo fita mefa mumi honi

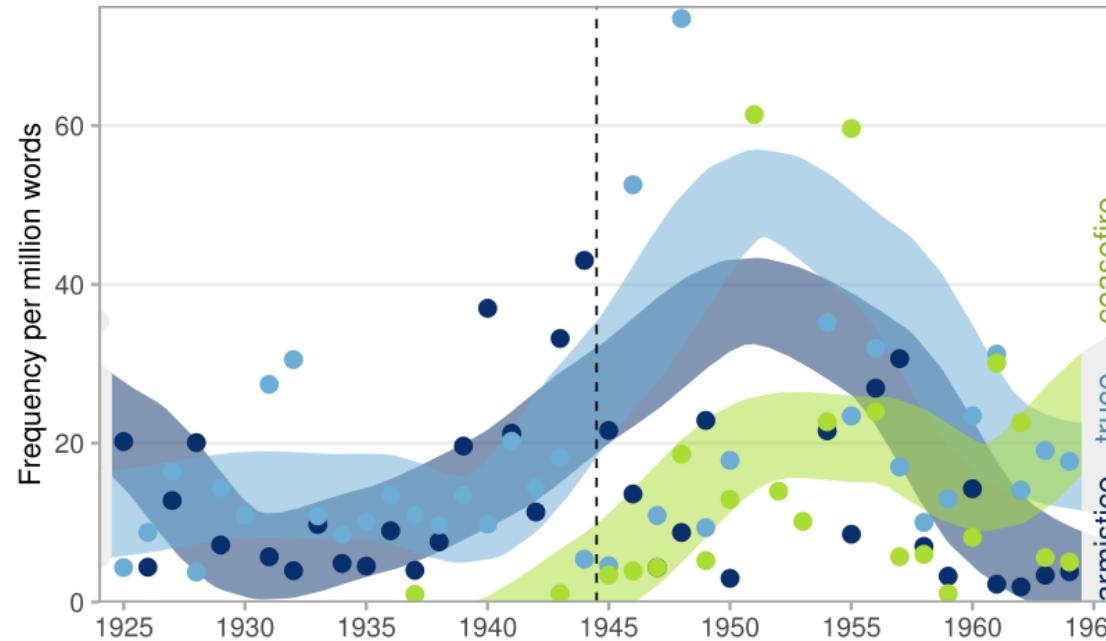
Expm no. 7, target condition, 91%, PPMI

VERDICT						3.03
ORGAN					1.81	1.62
KING						3.32
GAUGE						3.49
RAIN					2.23	0
DRIZZLE					2.49	0
TASK	2.1					0
JOB						0
STYLE		2.4				0
FASHION	2.1					0

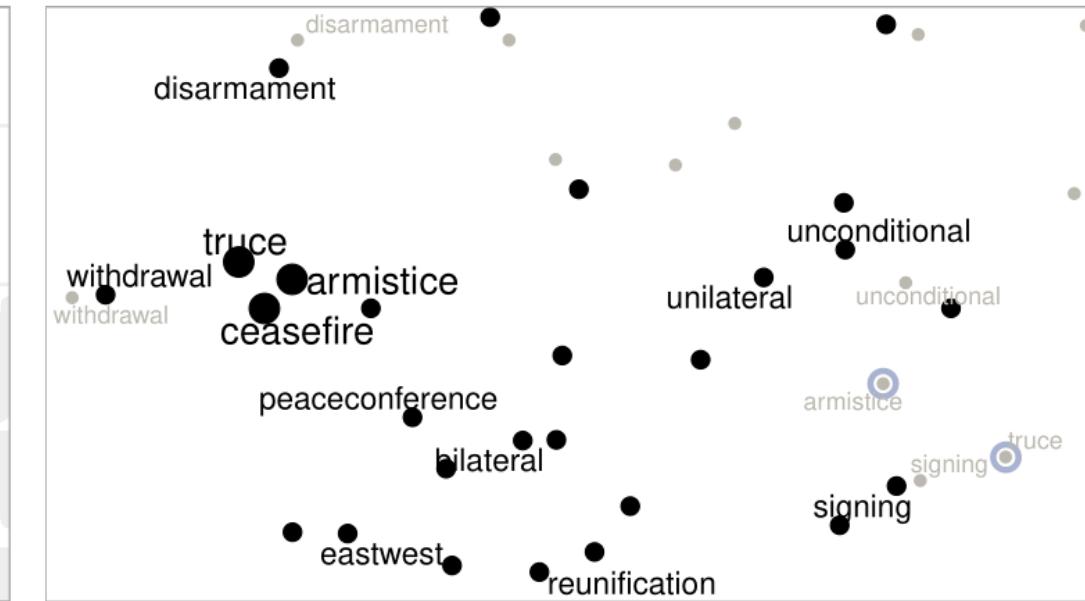
lahe pam muta qih posameqo qere



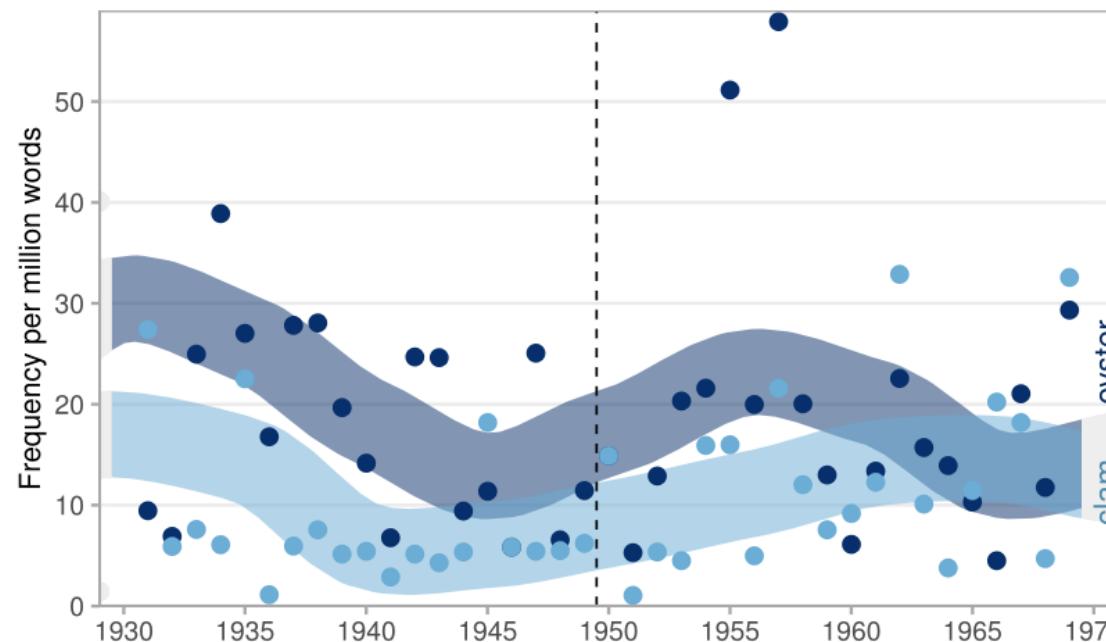
1A: Frequencies over time; lookback model



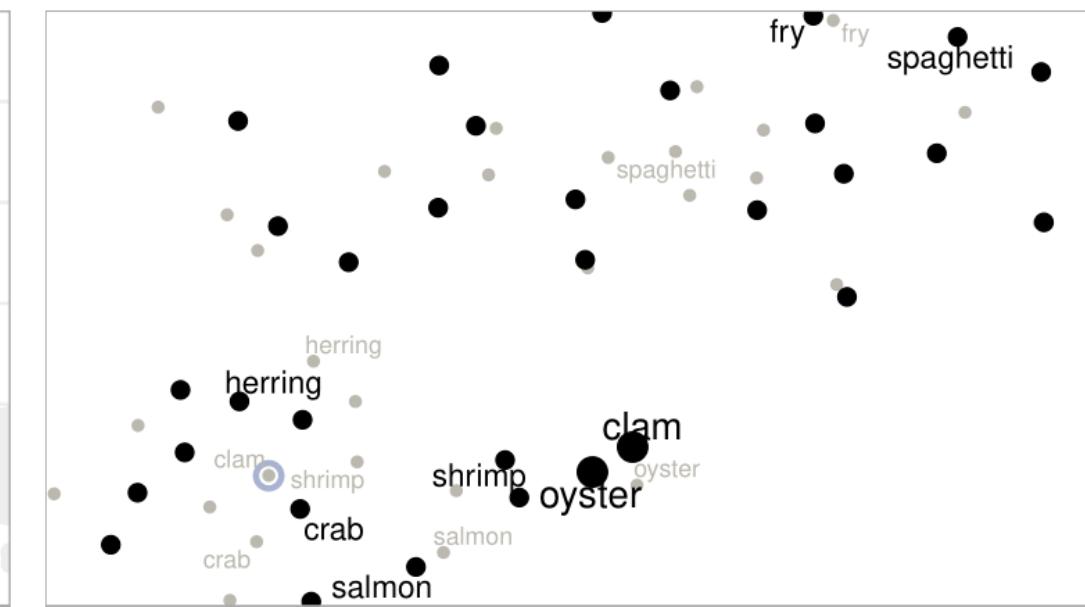
1B: Semantic space, periods overlaid



2A: Frequencies over time; lookback model



2B: Semantic space, periods overlaid



# Further evaluation

- But: direct synonym competition is very rare!
- Sample: COHA, equalization range <0.2 & number of losers <4 (n=52)
- near-synonym competition:
  - aeroplane → airplane, close-up → close shot, appropriation → funding apartment+inn → motel
- some proper nouns
  - guerrilla → Taliban, Yugoslav → Algerian
- mostly contextual, in-topic replacements:
  - railway → airline, opera+concert → movie atomic bomb → ballistic missile
- Still, advection predicts if replacement or not

