

Unit 3: Learning from other people

5. The structure in language

11/19/2020

- 1. You can learn a lot from the co-occurrence structure of words in language**
- 2. Latent semantic analysis and Topics models both use this structure to learn about the world**
- 3. But some information is not (straightforwardly) in the co-occurrence structure of language**

How do you know so much without being told about it?



Plato's Problem:

Even uneducated people seem to know a lot

Plato's Solution:

Knowledge is innate

Plato (380 BC)



Chomsky's Problem:

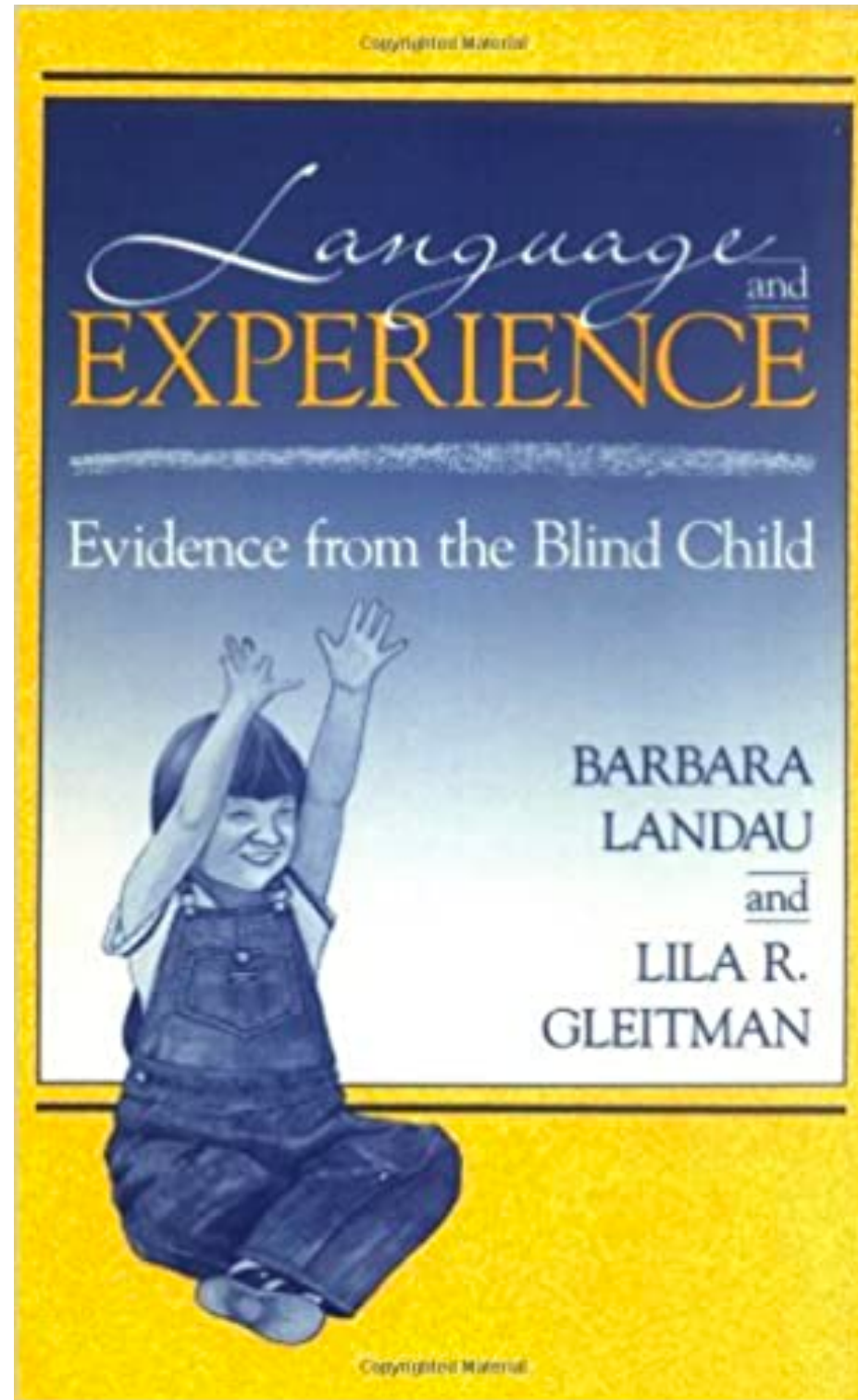
Children seem to learn language from insufficient input

Chomsky's Solution:

Universal grammar is innate

Chomsky (1986)

Blind children know the meanings of sight words!



Look up!



Let me see the back of your pants



Make it so mommy can't see the car

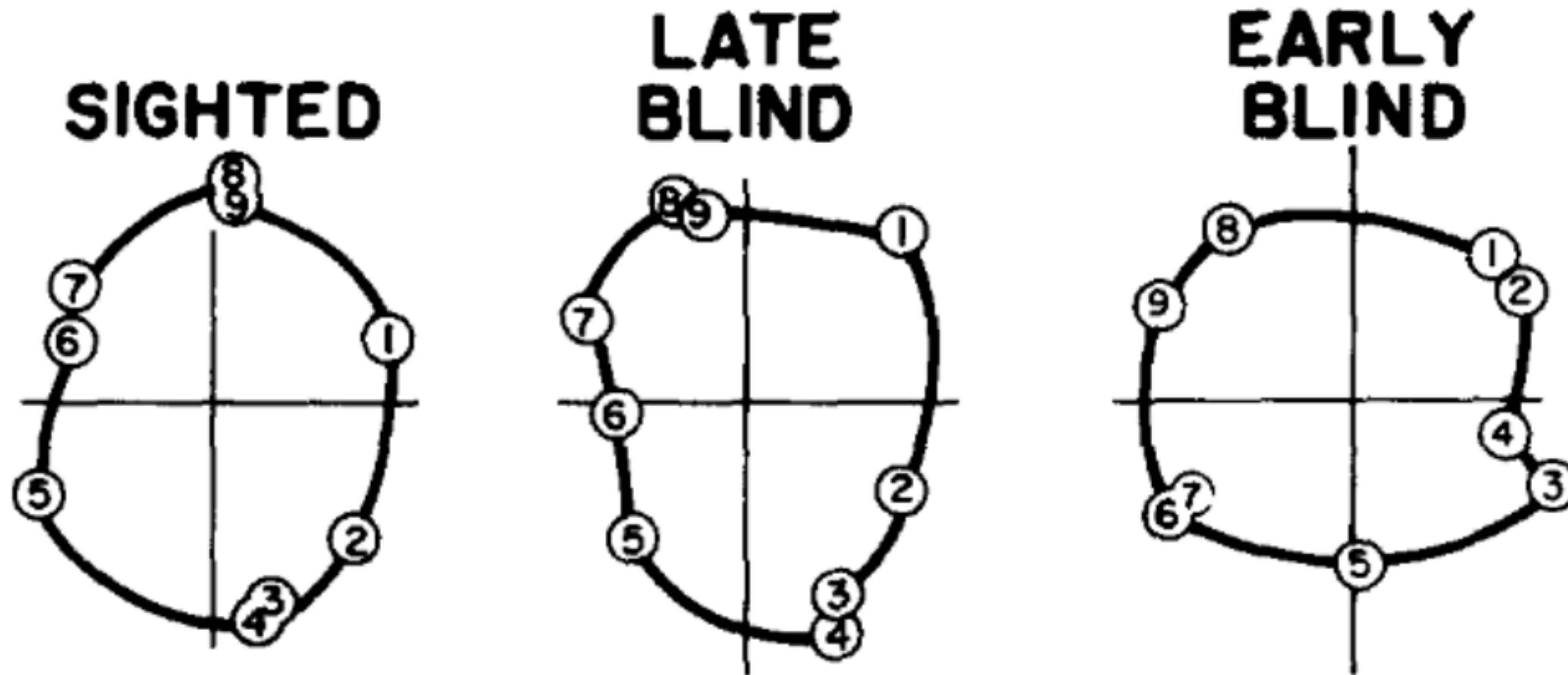


Let mommy see the car



Let mommy touch the car

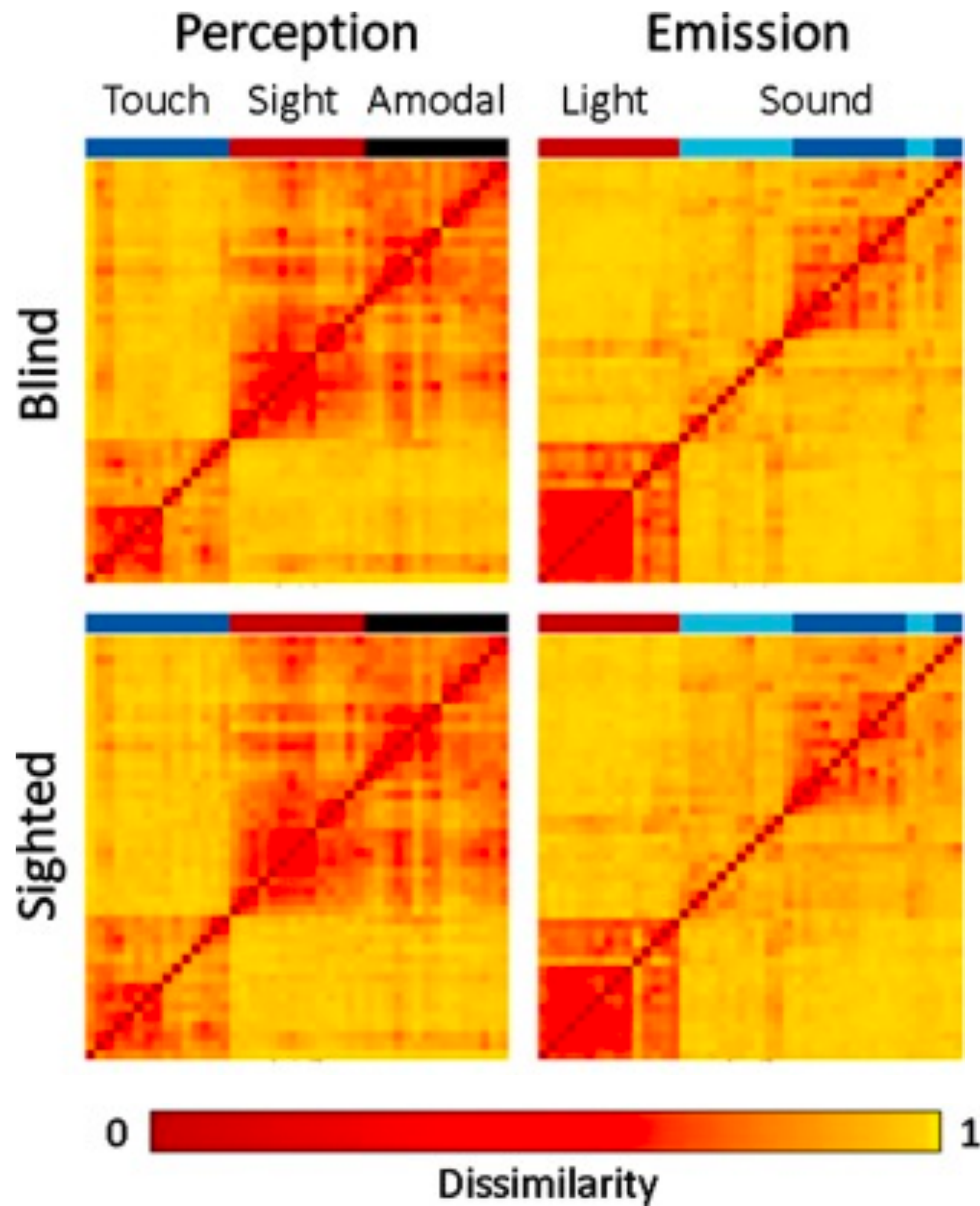
Blind adults color similarities look a lot like sighted adults



COLOR LEGEND:

1. RED 2. ORANGE 3. GOLD 4. YELLOW 5. GREEN
6. TURQUISE 7. BLUE 8. PURPLE 9. VIOLET

Blind adults general perceptual verb similarities are like sighted adults'



Visual		Touch		Amodal	
gawk		caress		characterize	
gaze		dab		classify	
glance		feel		discover	
glimpse		grip		examine	
leer		nudge		identify	
look		pat		investigate	
peek		pet		learn	
peer		pinch		note	
scan		prod		notice	
see		pub		perceive	
spot		scrape		question	
stare		stroke		recognize	
view		tap		scrutinize	
watch		tickle		search	
ogle*		touch		study	

Emission			Manner of Motion	
Light	Animate Sound	Inanimate Sound		
blaze	bark	beep		bounce
blink	bellow	boom		float
flare	groan	buzz		glide
flash	growl	chime		hobble
flicker	grumble	clang		roll
gleam	grunt	clank		saunter
glimmer	howl	click		scurry
glint	moan	crackle		skip
glisten	mutter	creak		slither
glitter	shout	crunch		spin
glow	squawk	gurgle		strut
shimmer	wail	hiss		trot
shine	whimper	sizzle		twirl
sparkle	whisper	squeak		twist
twinkle	yelp	twang		waddle

A solution to Plato's problem (Landauer & Dumais, 1997)

Red onions are sweeter than **white** ones

Red hair occurs naturally in one to two percent of the human population

Pittsburgh one of U.S. cities with highest number of **gray** days

Fall tips for a **green** spring lawn

Lake Tahoe stretches 22 miles long and 12 miles wide, with clear **blue** water that's more than 99 percent pure

Direct information:

There is a relationship between e.g. red and hair

Indirect information:

Red, white, gray, green, and blue are used in *similar contexts*.

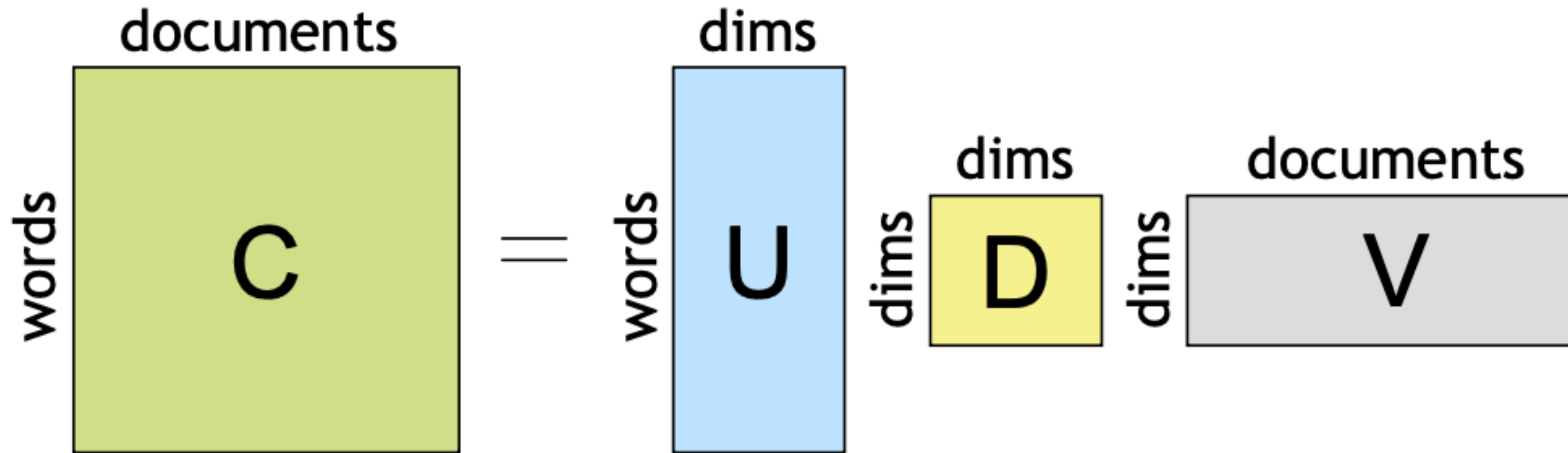
Contexts for e.g. blue and green are more similar than blue and red

Start with a term x document matrix

		Document					
Term		d1	d2	d3	d4	d5	d6
	rock	2	1	0	2	0	1
	granite	1	0	1	0	0	0
	marble	1	2	0	0	0	0
	music	0	0	0	1	2	0
	song	0	0	0	1	0	2
	band	0	0	0	0	1	0

Words that occur in similar documents are probably similar

Latent semantic analysis: Finding the hidden structure in documents



Instead of representing words by co-occurrence in documents, we want to represent their co-occurrence in semantic dimensions

Semantics are a low-dimensional compression of documents

Technical Memo Titles

c1: *Human machine interface for ABC computer applications*

c2: *A survey of user opinion of computer system response time*

c3: *The EPS user interface management system*

c4: *System and human system engineering testing of EPS*

c5: *Relation of user perceived response time to error measurement*

m1: *The generation of random, binary, ordered trees*

m2: *The intersection graph of paths in trees*

m3: *Graph minors IV: Widths of trees and well-quasi-ordering*

m4: *Graph minors: A survey*

An example of LSA (Landauer, Foltz, & Laham, 1998)

	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	1	0	0	1	0	0	0	0	0
interface	1	0	1	0	0	0	0	0	0
computer	1	1	0	0	0	0	0	0	0
user	0	1	1	0	1	0	0	0	0
system	0	1	1	2	0	0	0	0	0
response	0	1	0	0	1	0	0	0	0
time	0	1	0	0	1	0	0	0	0
EPS	0	0	1	1	0	0	0	0	0
survey	0	1	0	0	0	0	0	0	1
trees	0	0	0	0	0	1	1	1	0
graph	0	0	0	0	0	0	1	1	1
minors	0	0	0	0	0	0	0	1	1

$$r(\text{human}, \text{user}) = -.38$$

$$r(\text{human}, \text{minors}) = -.29$$

Decomposing the matrix

C

	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	1	0	0	1	0	0	0	0	0
interface	1	0	1	0	0	0	0	0	0
computer	1	1	0	0	0	0	0	0	0
user	0	1	1	0	1	0	0	0	0
system	0	1	1	2	0	0	0	0	0
response	0	1	0	0	1	0	0	0	0
time	0	1	0	0	1	0	0	0	0
EPS	0	0	1	1	0	0	0	0	0
survey	0	1	0	0	0	0	0	0	1
trees	0	0	0	0	0	1	1	1	0
graph	0	0	0	0	0	0	1	1	1
minors	0	0	0	0	0	0	0	1	1

U

D

V

0.22	-0.11	0.29	-0.41	-0.11	-0.34	0.52	-0.06	-0.41
0.20	-0.07	0.14	-0.55	0.28	0.50	-0.07	-0.01	-0.11
0.24	0.04	-0.16	-0.59	-0.11	-0.25	-0.30	0.06	0.49
0.40	0.06	-0.34	0.10	0.33	0.38	0.00	0.00	0.01
0.64	-0.17	0.36	0.33	-0.16	-0.21	-0.17	0.03	0.27
0.27	0.11	-0.43	0.07	0.08	-0.17	0.28	-0.02	-0.05
0.27	0.11	-0.43	0.07	0.08	-0.17	0.28	-0.02	-0.05
0.30	-0.14	0.33	0.19	0.11	0.27	0.03	-0.02	-0.17
0.21	0.27	-0.18	-0.03	-0.54	0.08	-0.47	-0.04	-0.58
0.01	0.49	0.23	0.03	0.59	-0.39	-0.29	0.25	-0.23
0.04	0.62	0.22	0.00	-0.07	0.11	0.16	-0.68	0.23
0.03	0.45	0.14	-0.01	-0.30	0.28	0.34	0.68	0.18

3.34

2.54

2.35

1.64

1.50

1.31

0.85

0.56

0.36

0.20	0.61	0.46	0.54	0.28	0.00	0.01	0.02	0.08
-0.06	0.17	-0.13	-0.23	0.11	0.19	0.44	0.62	0.53
0.11	-0.50	0.21	0.57	-0.51	0.10	0.19	0.25	0.08
-0.95	-0.03	0.04	0.27	0.15	0.02	0.02	0.01	-0.03
0.05	-0.21	0.38	-0.21	0.33	0.39	0.35	0.15	-0.60
-0.08	-0.26	0.72	-0.37	0.03	-0.30	-0.21	0.00	0.36
0.18	-0.43	-0.24	0.26	0.67	-0.34	-0.15	0.25	0.04
-0.01	0.05	0.01	-0.02	-0.06	0.45	-0.76	0.45	-0.07
-0.06	0.24	0.02	-0.08	-0.26	-0.62	0.02	0.52	-0.45

A low-dimensional reconstruction using the first 2 dimensions

C

	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	1	0	0	1	0	0	0	0	0
interface	1	0	1	0	0	0	0	0	0
computer	1	1	0	0	0	0	0	0	0
user	0	1	1	0	1	0	0	0	0
system	0	1	1	2	0	0	0	0	0
response	0	1	0	0	1	0	0	0	0
time	0	1	0	0	1	0	0	0	0
EPS	0	0	1	1	0	0	0	0	0
survey	0	1	0	0	0	0	0	0	1
trees	0	0	0	0	0	1	1	1	0
graph	0	0	0	0	0	0	1	1	1
minors	0	0	0	0	0	0	0	1	1

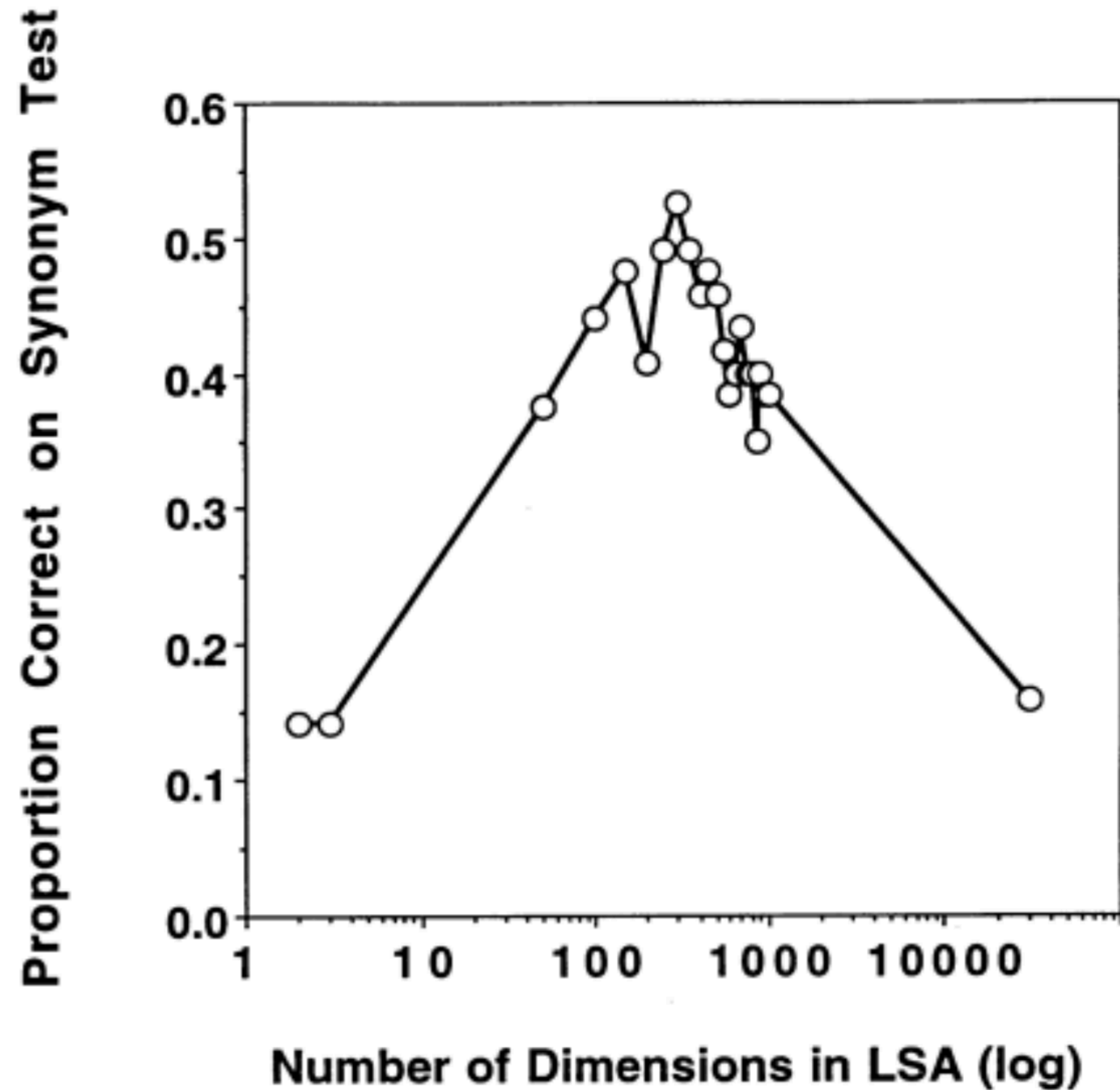
$$r(\text{human, user}) = .94$$

$$r(\text{human, minors}) = -.83$$

C'

	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	0.16	0.40	0.38	0.47	0.18	-0.05	-0.12	-0.16	-0.09
interface	0.14	0.37	0.33	0.40	0.16	-0.03	-0.07	-0.10	-0.04
computer	0.15	0.51	0.36	0.41	0.24	0.02	0.06	0.09	0.12
user	0.26	0.84	0.61	0.70	0.39	0.03	0.08	0.12	0.19
system	0.45	1.23	1.05	1.27	0.56	-0.07	-0.15	-0.21	-0.05
response	0.16	0.58	0.38	0.42	0.28	0.06	0.13	0.19	0.22
time	0.16	0.58	0.38	0.42	0.28	0.06	0.13	0.19	0.22
EPS	0.22	0.55	0.51	0.63	0.24	-0.07	-0.14	-0.20	-0.11
survey	0.10	0.53	0.23	0.21	0.27	0.14	0.31	0.44	0.42
trees	-0.06	0.23	-0.14	-0.27	0.14	0.24	0.55	0.77	0.66
graph	-0.06	0.34	-0.15	-0.30	0.20	0.31	0.69	0.98	0.85
minors	-0.04	0.25	-0.10	-0.21	0.15	0.22	0.50	0.71	0.62

Latent semantic analysis predicts human similarity judgments

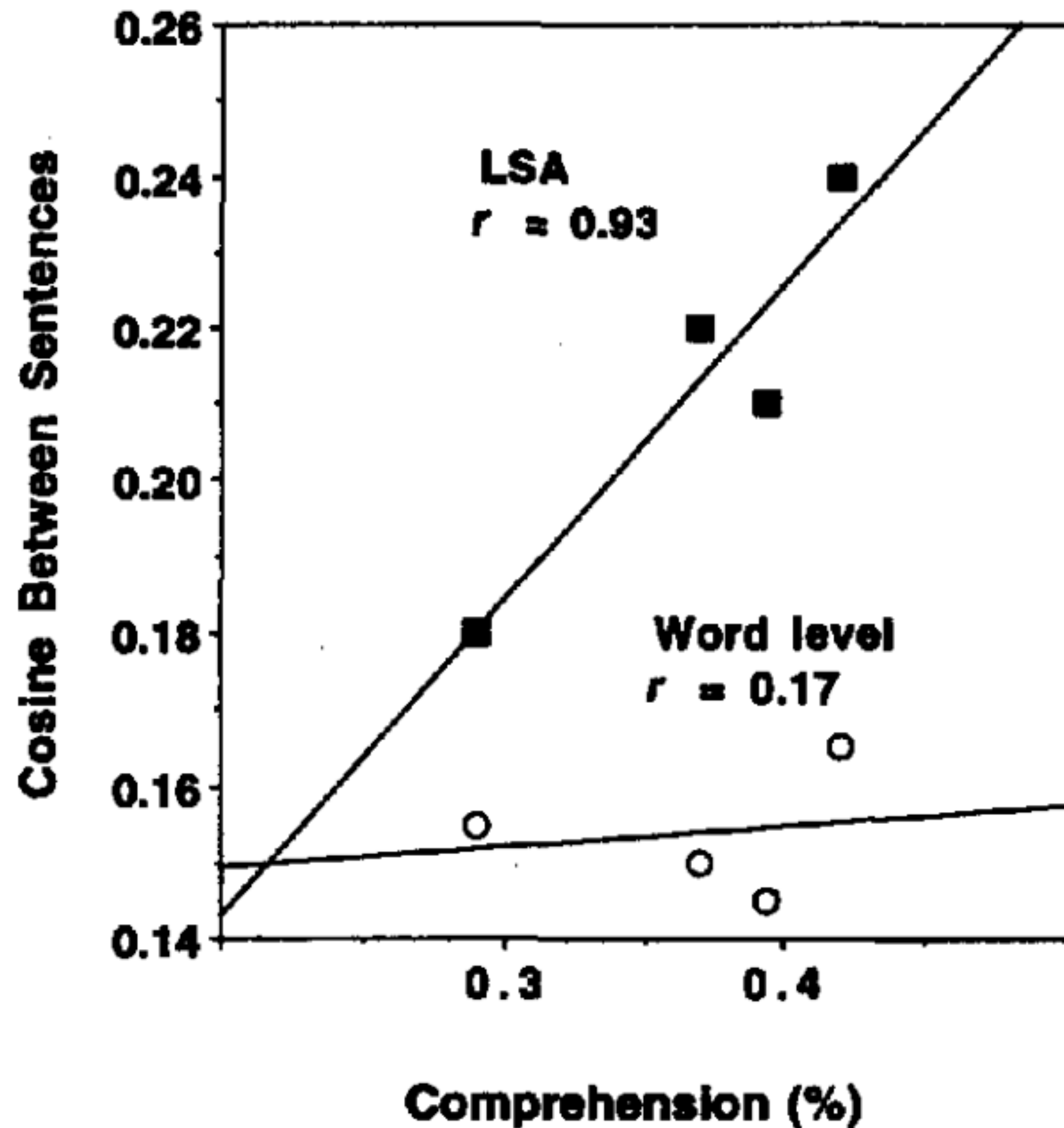


LSA can pass the synonym part of the TOEFL!

Chance is 25%

Applicants to US Universities on average get 52.7%

Latent semantic analysis predicts sentence comprehensibility

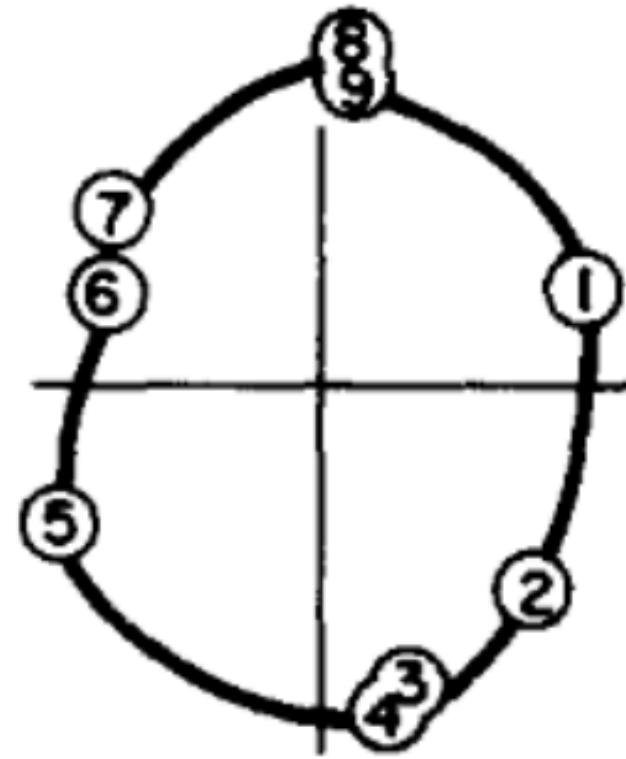


Mammals have very specialized teeth. There are four types of teeth in mammals: incisors, canines, premolars, and molars. The number and shape of each of these types of teeth are related to the kind of food the mammal eats.

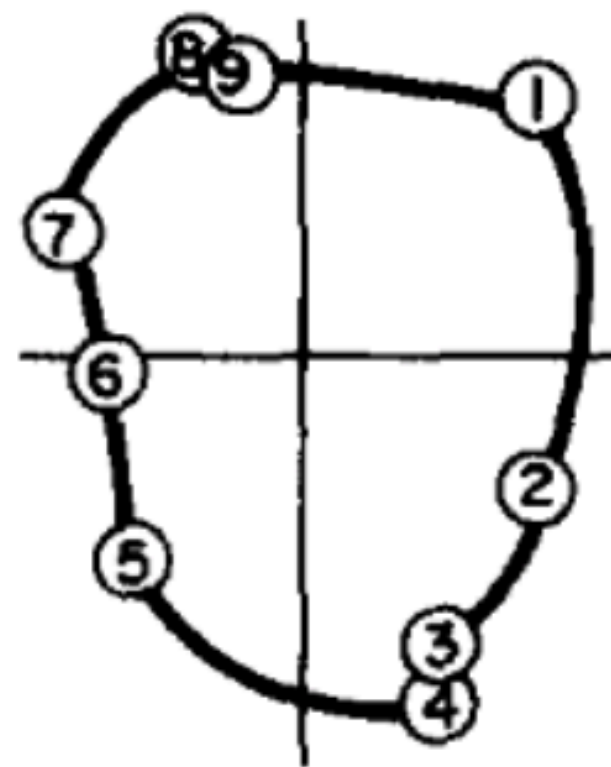
Another physical trait of mammals is that they can eat many different kinds of food because they have very specialized teeth. This trait also helps them to live in different kinds of environments

Can LSA recover perceptual similarities?

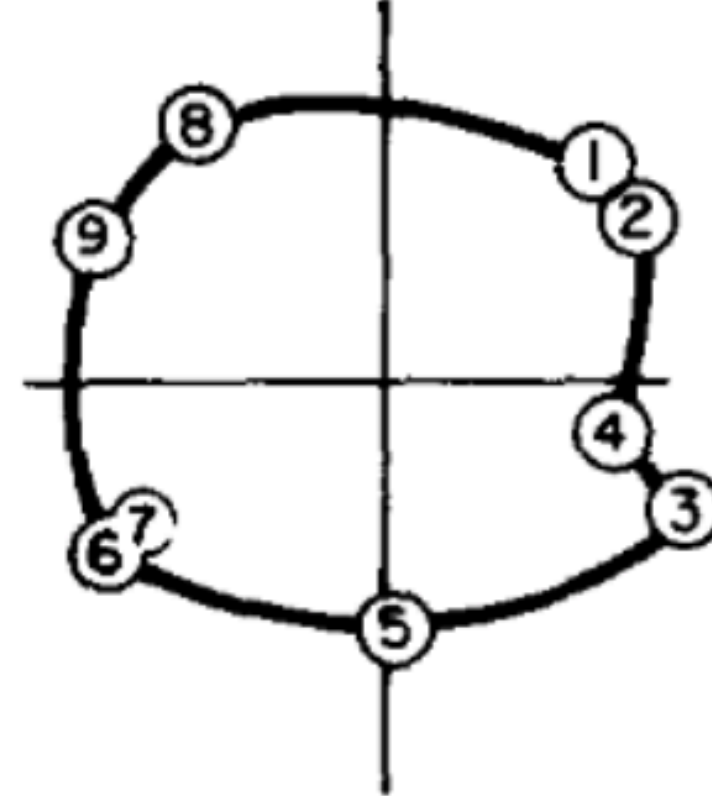
SIGHTED



LATE BLIND



EARLY BLIND



COLOR LEGEND:

- 1. RED 2. ORANGE 3. GOLD 4. YELLOW 5. GREEN
- 6. TURQUISE 7. BLUE 8. PURPLE 9. VIOLET

Visual		Touch		Amodal
gawk		caress		characterize
gaze		dab		classify
glance		feel		discover
glimpse		grip		examine
leer		nudge		identify
look		pat		investigate
peek		pet		learn
peer		pinch		note
scan		prod		notice
see		pub		perceive
spot		scrape		question
stare		stroke		recognize
view		tap		scrutinize
watch		tickle		search
ogle		touch		study

Emission			Manner of Motion
Light	Animate Sound	Inanimate Sound	
blaze	bark	beep	bounce
blink	bellow	boom	float
flare	groan	buzz	glide
flash	growl	chime	hobble
flicker	grumble	clang	roll
gleam	grunt	clank	saunter
glimmer	howl	click	scurry
glint	moan	crackle	skip
glisten	mutter	creak	slither
glitter	shout	crunch	spin
glow	squawk	gurgle	strut
shimmer	wail	hiss	trot
shine	whimper	sizzle	twirl
sparkle	whisper	squeak	twist
twinkle	yelp	twang	waddle

<https://bit.ly/test-lsa>

But what are these underlying dimensions?

C

	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	1	0	0	1	0	0	0	0	0
interface	1	0	1	0	0	0	0	0	0
computer	1	1	0	0	0	0	0	0	0
user	0	1	1	0	1	0	0	0	0
system	0	1	1	2	0	0	0	0	0
response	0	1	0	0	1	0	0	0	0
time	0	1	0	0	1	0	0	0	0
EPS	0	0	1	1	0	0	0	0	0
survey	0	1	0	0	0	0	0	0	1
trees	0	0	0	0	0	1	1	1	0
graph	0	0	0	0	0	0	1	1	1
minors	0	0	0	0	0	0	0	1	1

U

D

V

0.22	-0.11	0.29	-0.41	-0.11	-0.34	0.52	-0.06	-0.41
0.20	-0.07	0.14	-0.55	0.28	0.50	-0.07	-0.01	-0.11
0.24	0.04	-0.16	-0.59	-0.11	-0.25	-0.30	0.06	0.49
0.40	0.06	-0.34	0.10	0.33	0.38	0.00	0.00	0.01
0.64	-0.17	0.36	0.33	-0.16	-0.21	-0.17	0.03	0.27
0.27	0.11	-0.43	0.07	0.08	-0.17	0.28	-0.02	-0.05
0.27	0.11	-0.43	0.07	0.08	-0.17	0.28	-0.02	-0.05
0.30	-0.14	0.33	0.19	0.11	0.27	0.03	-0.02	-0.17
0.21	0.27	-0.18	-0.03	-0.54	0.08	-0.47	-0.04	-0.58
0.01	0.49	0.23	0.03	0.59	-0.39	-0.29	0.25	-0.23
0.04	0.62	0.22	0.00	-0.07	0.11	0.16	-0.68	0.23
0.03	0.45	0.14	-0.01	-0.30	0.28	0.34	0.68	0.18

3.34

2.54

2.35

1.64

1.50

1.31

0.85

0.56

0.36

0.20	0.61	0.46	0.54	0.28	0.00	0.01	0.02	0.08
-0.06	0.17	-0.13	-0.23	0.11	0.19	0.44	0.62	0.53
0.11	-0.50	0.21	0.57	-0.51	0.10	0.19	0.25	0.08
-0.95	-0.03	0.04	0.27	0.15	0.02	0.02	0.01	-0.03
0.05	-0.21	0.38	-0.21	0.33	0.39	0.35	0.15	-0.60
-0.08	-0.26	0.72	-0.37	0.03	-0.30	-0.21	0.00	0.36
0.18	-0.43	-0.24	0.26	0.67	-0.34	-0.15	0.25	0.04
-0.01	0.05	0.01	-0.02	-0.06	0.45	-0.76	0.45	-0.07
-0.06	0.24	0.02	-0.08	-0.26	-0.62	0.02	0.52	-0.45

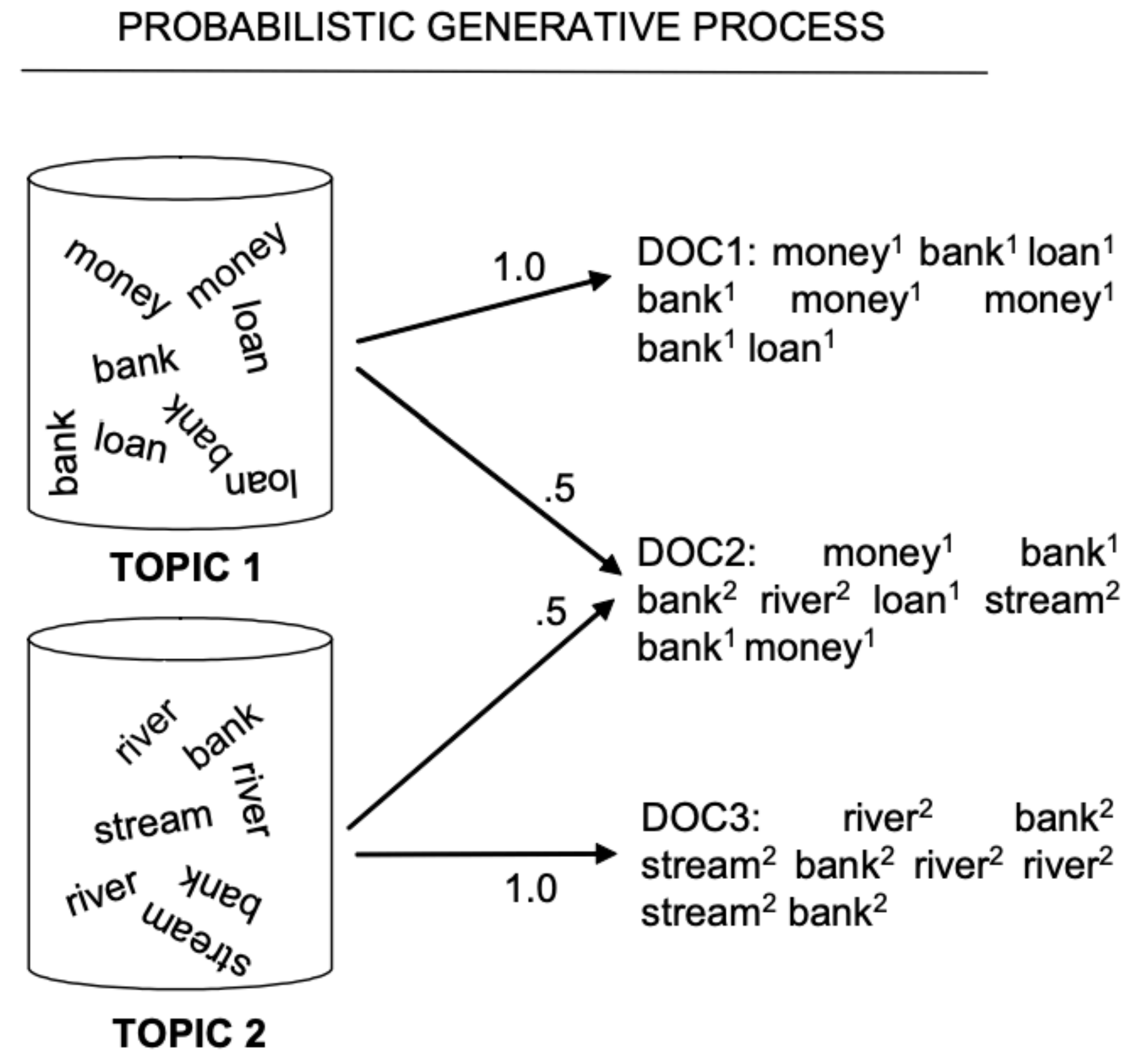
Topic Models (Latent Dirichlet Allocation) - Steyvers & Griffiths (2007)

Idea: A generative model for how words appear in documents

Each document is created by picking a set of topics, each with some weight.

For every word in the document, it is chosen from one of the topics according to their weight

Model from Blei, Ng, & Jordan (2003)



Topic Models (Latent Dirichlet Allocation) - Steyvers & Griffiths (2007)

For each **topic** $1 \dots t$:

Draw a **multinomial over words** $\phi_t \sim \text{Dirichlet}(\beta)$

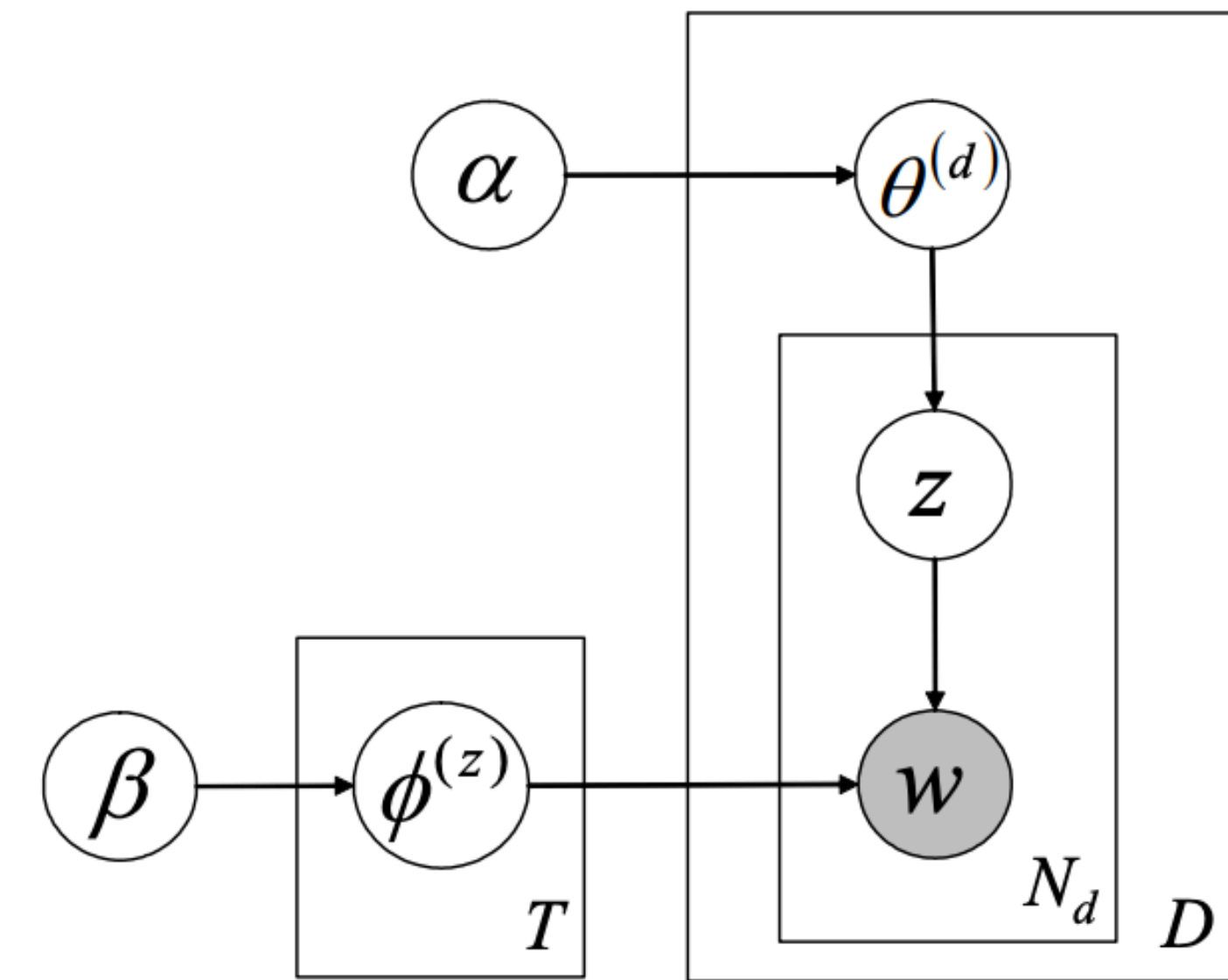
For each **document** $1 \dots d$:

Draw a **multinomial over topics** $\theta_d \sim \text{Dirichlet}(\alpha)$

For each word $w_{d,n}$:

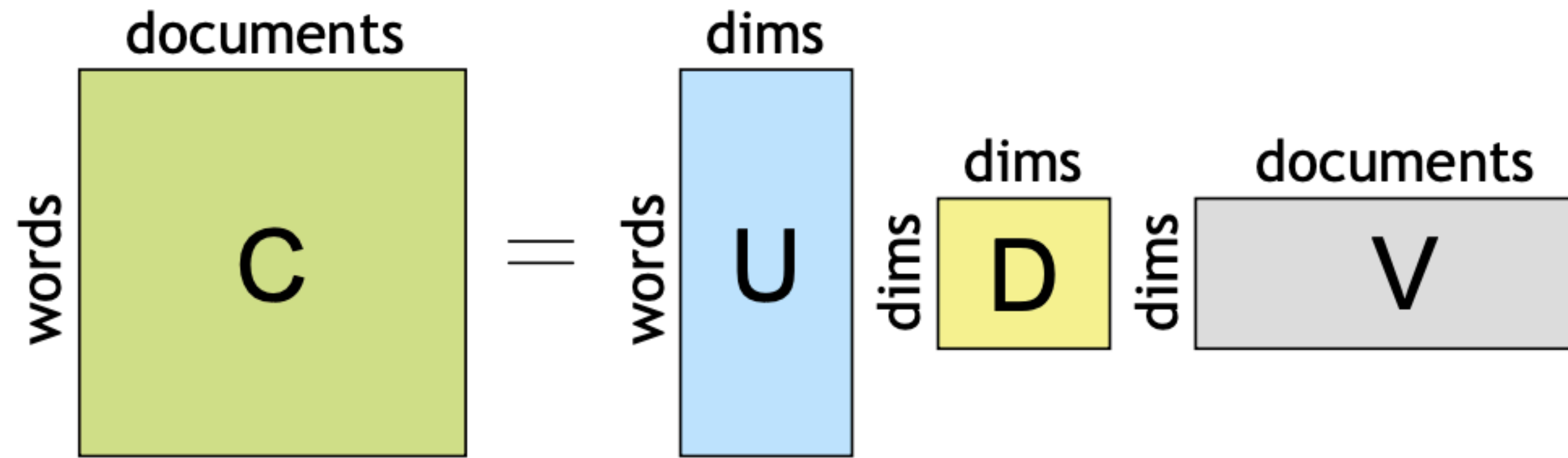
Draw a topic $Z_{d,n} \sim \text{Multinomial}(\theta_d)$ with $Z_{d,n} \in [1 \dots t]$

Draw a word $w_{d,n} \sim \text{Multinomial}(\beta_{Z_{d,n}})$

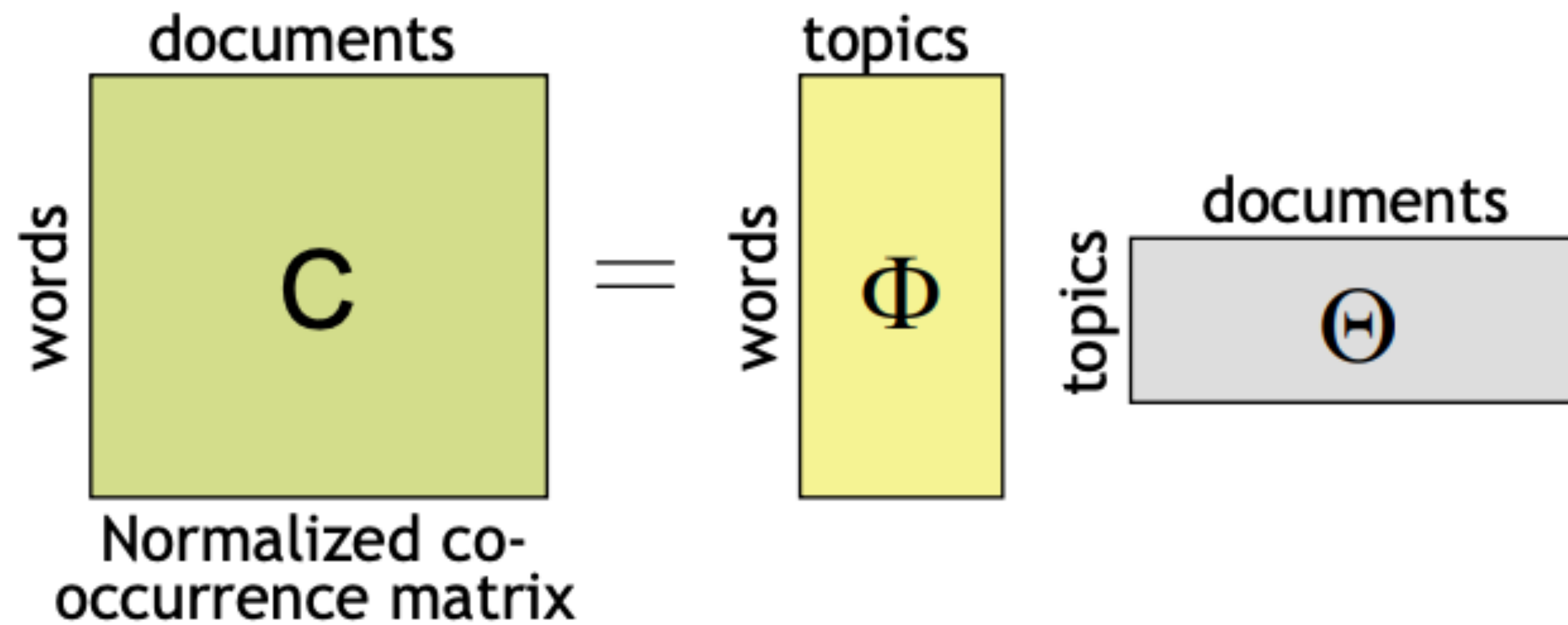


Comparing LSA and Topic Models

LSA

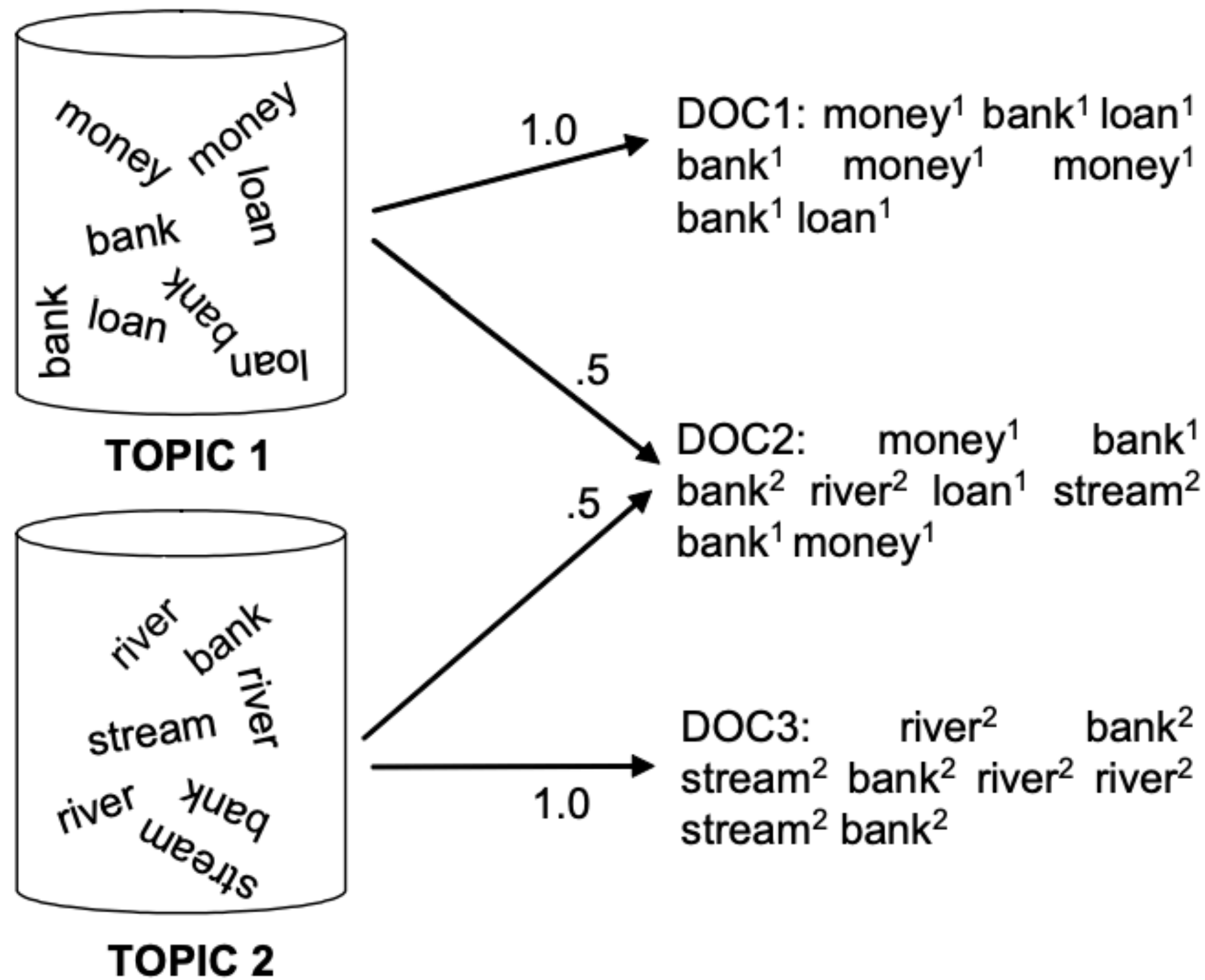


Topic Model

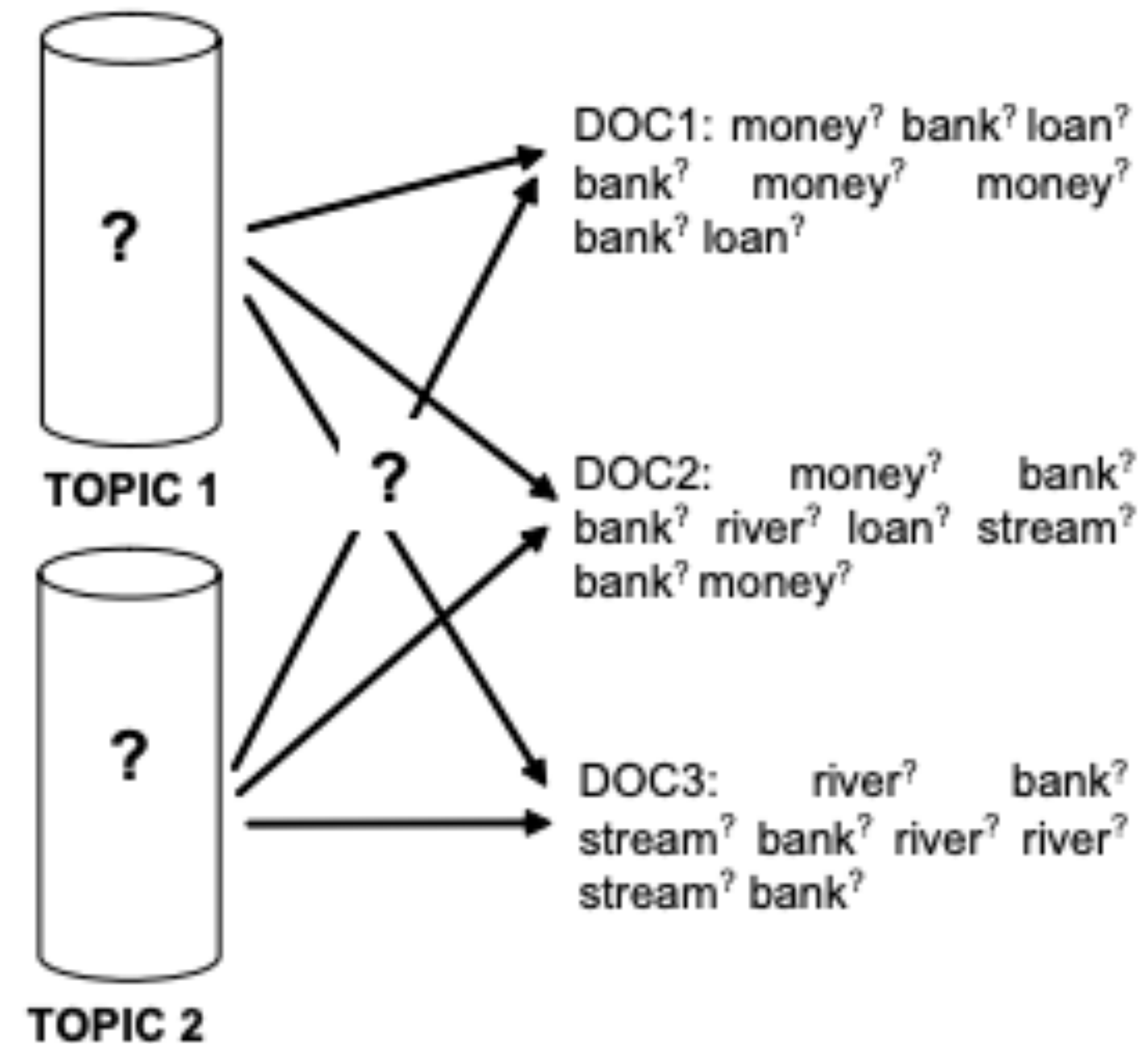


Inferring topics by Bayesian inference using sampling

PROBABILISTIC GENERATIVE PROCESS

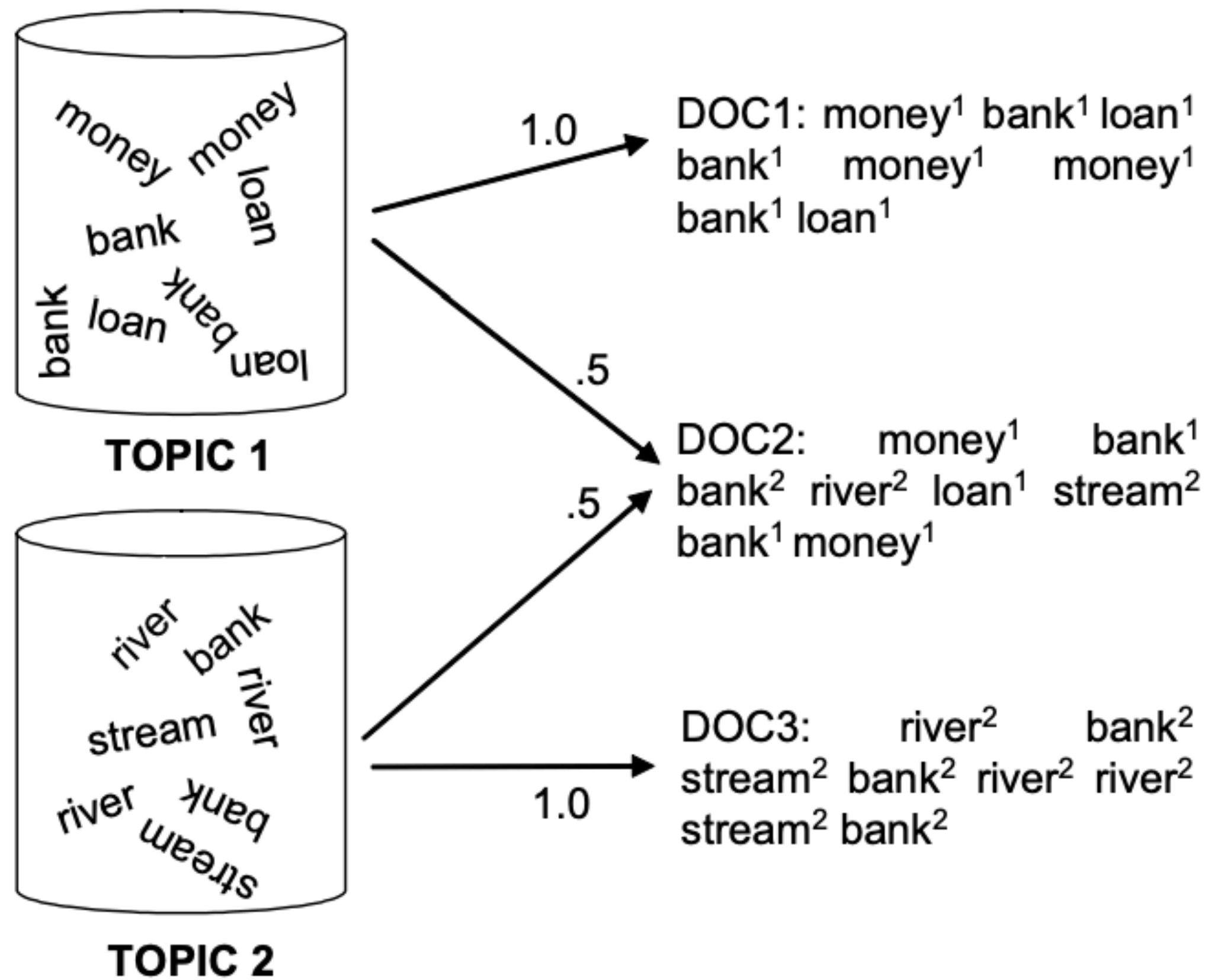


STATISTICAL INFERENCE

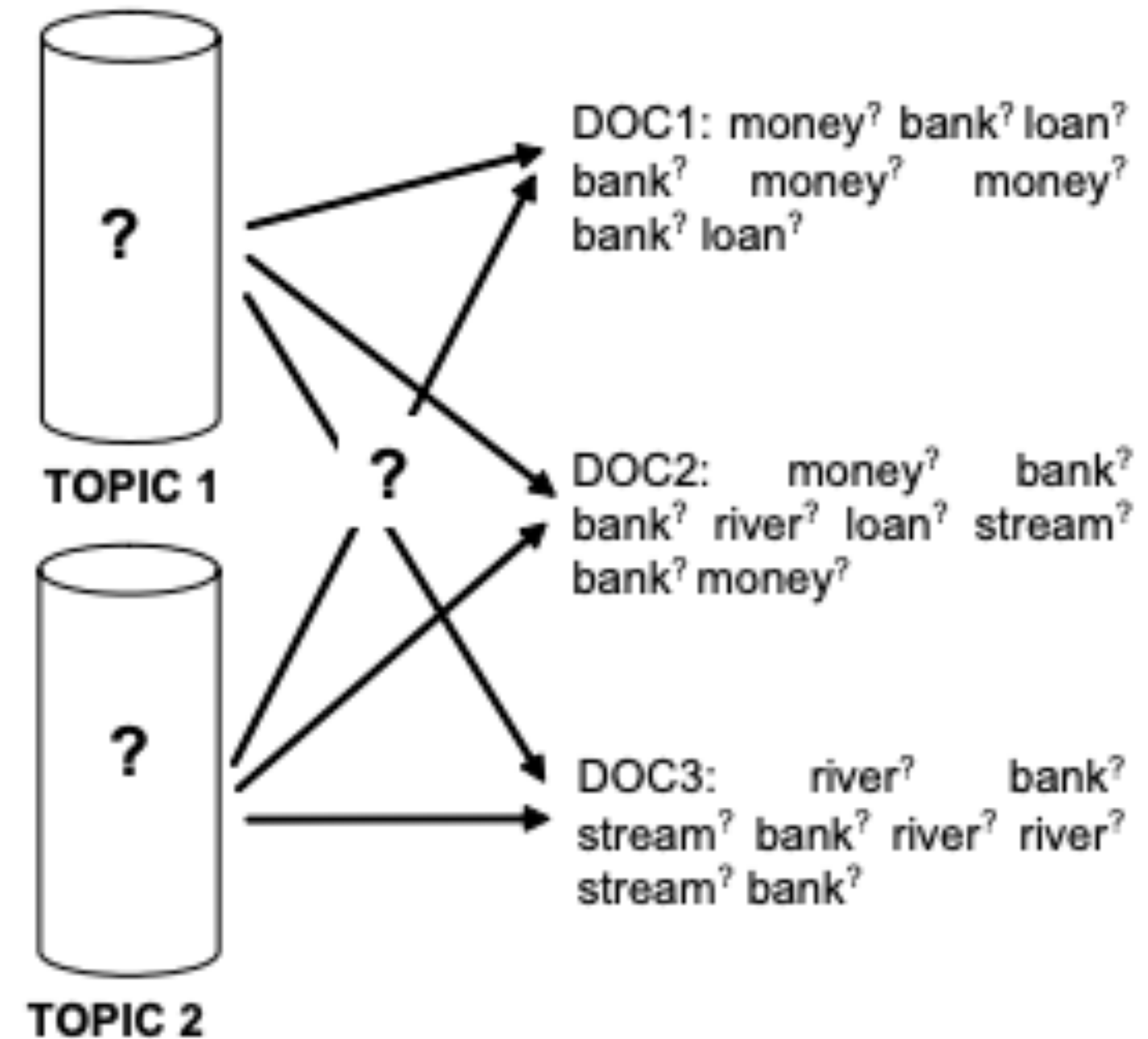


Inferring topics by Bayesian inference using sampling

PROBABILISTIC GENERATIVE PROCESS



STATISTICAL INFERENCE



Example topics from the Touchstone Applied Sciences Association (TASA) corpus

Topic 247

word	prob.
DRUGS	.069
DRUG	.060
MEDICINE	.027
EFFECTS	.026
BODY	.023
MEDICINES	.019
PAIN	.016
PERSON	.016
MARIJUANA	.014
LABEL	.012
ALCOHOL	.012
DANGEROUS	.011
ABUSE	.009
EFFECT	.009
KNOWN	.008
PILLS	.008

Topic 5

word	prob.
RED	.202
BLUE	.099
GREEN	.096
YELLOW	.073
WHITE	.048
COLOR	.048
BRIGHT	.030
COLORS	.029
ORANGE	.027
BROWN	.027
PINK	.017
LOOK	.017
BLACK	.016
PURPLE	.015
CROSS	.011
COLORED	.009

Topic 43

word	prob.
MIND	.081
THOUGHT	.066
REMEMBER	.064
MEMORY	.037
THINKING	.030
PROFESSOR	.028
FELT	.025
REMEMBERED	.022
THOUGHTS	.020
FORGOTTEN	.020
MOMENT	.020
THINK	.019
THING	.016
WONDER	.014
FORGET	.012
RECALL	.012

Topic 56

word	prob.
DOCTOR	.074
DR.	.063
PATIENT	.061
HOSPITAL	.049
CARE	.046
MEDICAL	.042
NURSE	.031
PATIENTS	.029
DOCTORS	.028
HEALTH	.025
MEDICINE	.017
NURSING	.017
DENTAL	.015
NURSES	.013
PHYSICIAN	.012
HOSPITALS	.011

Topics models can resolve polysemy

Bix beiderbecke, at age⁰⁶⁰ fifteen²⁰⁷, sat¹⁷⁴ on the slope⁰⁷¹ of a bluff⁰⁵⁵ overlooking⁰²⁷ the mississippi¹³⁷ river¹³⁷. He was listening⁰⁷⁷ to music⁰⁷⁷ coming⁰⁰⁹ from a passing⁰⁴³ riverboat. The music⁰⁷⁷ had already captured⁰⁰⁶ his heart¹⁵⁷ as well as his ear¹¹⁹. It was jazz⁰⁷⁷. Bix beiderbecke had already had music⁰⁷⁷ lessons⁰⁷⁷. He showed⁰⁰² promise¹³⁴ on the piano⁰⁷⁷, and his parents⁰³⁵ hoped²⁶⁸ he might consider¹¹⁸ becoming a concert⁰⁷⁷ pianist⁰⁷⁷. But bix was interested²⁶⁸ in another kind⁰⁵⁰ of music⁰⁷⁷. He wanted²⁶⁸ to **play**⁰⁷⁷ the cornet. And he wanted²⁶⁸ to **play**⁰⁷⁷ jazz⁰⁷⁷ ...

Topic 77

word	prob.
MUSIC	.090
DANCE	.034
SONG	.033
PLAY	.030
SING	.026
SINGING	.026
BAND	.026
PLAYED	.023
SANG	.022
SONGS	.021
DANCING	.020
PIANO	.017
PLAYING	.016
RHYTHM	.015
ALBERT	.013
MUSICAL	.013

Topic 82

word	prob.
LITERATURE	.031
POEM	.028
POETRY	.027
POET	.020
PLAYS	.019
POEMS	.019
PLAY	.015
LITERARY	.013
WRITERS	.013
DRAMA	.012
WROTE	.012
POETS	.011
WRITER	.011
SHAKESPEARE	.010
WRITTEN	.009
STAGE	.009

Topic 166

word	prob.
PLAY	.136
BALL	.129
GAME	.065
PLAYING	.042
HIT	.032
PLAYED	.031
BASEBALL	.027
GAMES	.025
BAT	.019
RUN	.019
THROW	.016
BALLS	.015
TENNIS	.011
HOME	.010
CATCH	.010
FIELD	.010

Topics models can resolve polysemy

There is a simple⁰⁵⁰ reason¹⁰⁶ why there are so few periods⁰⁷⁸ of really great theater⁰⁸² in our whole western⁰⁴⁶ world. Too many things³⁰⁰ have to come right at the very same time. The dramatists must have the right actors⁰⁸², the actors⁰⁸² must have the right playhouses, the playhouses must have the right audiences⁰⁸². We must remember²⁸⁸ that plays⁰⁸² exist¹⁴³ to be performed⁰⁷⁷, not merely⁰⁵⁰ to be read²⁵⁴. (even when you read²⁵⁴ a play⁰⁸² to yourself, try²⁸⁸ to perform⁰⁶² it, to put¹⁷⁴ it on a stage⁰⁷⁸, as you go along.) as soon⁰²⁸ as a play⁰⁸² has to be performed⁰⁸², then some kind¹²⁶ of theatrical⁰⁸² ...

Topic 77

word	prob.
MUSIC	.090
DANCE	.034
SONG	.033
PLAY	.030
SING	.026
SINGING	.026
BAND	.026
PLAYED	.023
SANG	.022
SONGS	.021
DANCING	.020
PIANO	.017
PLAYING	.016
RHYTHM	.015
ALBERT	.013
MUSICAL	.013

Topic 82

word	prob.
LITERATURE	.031
POEM	.028
POETRY	.027
POET	.020
PLAYS	.019
POEMS	.019
PLAY	.015
LITERARY	.013
WRITERS	.013
DRAMA	.012
WROTE	.012
POETS	.011
WRITER	.011
SHAKESPEARE	.010
WRITTEN	.009
STAGE	.009

Topic 166

word	prob.
PLAY	.136
BALL	.129
GAME	.065
PLAYING	.042
HIT	.032
PLAYED	.031
BASEBALL	.027
GAMES	.025
BAT	.019
RUN	.019
THROW	.016
BALLS	.015
TENNIS	.011
HOME	.010
CATCH	.010
FIELD	.010

Topics models can resolve polysemy

Jim²⁹⁶ has a game¹⁶⁶ book²⁵⁴. Jim²⁹⁶ reads²⁵⁴ the book²⁵⁴. Jim²⁹⁶ sees⁰⁸¹ a game¹⁶⁶ for one. Jim²⁹⁶ plays¹⁶⁶ the game¹⁶⁶. Jim²⁹⁶ likes⁰⁸¹ the game¹⁶⁶ for one. The game¹⁶⁶ book²⁵⁴ helps⁰⁸¹ jim²⁹⁶. Don¹⁸⁰ comes⁰⁴⁰ into the house⁰³⁸. Don¹⁸⁰ and jim²⁹⁶ read²⁵⁴ the game¹⁶⁶ book²⁵⁴. The boys⁰²⁰ see a game¹⁶⁶ for two. The two boys⁰²⁰ **play¹⁶⁶** the game¹⁶⁶. The boys⁰²⁰ **play¹⁶⁶** the game¹⁶⁶ for two. The boys⁰²⁰ like the game¹⁶⁶. Meg²⁸² comes⁰⁴⁰ into the house²⁸². Meg²⁸² and don¹⁸⁰ and jim²⁹⁶ read²⁵⁴ the book²⁵⁴. They see a game¹⁶⁶ for three. Meg²⁸² and don¹⁸⁰ and jim²⁹⁶ **play¹⁶⁶** the game¹⁶⁶. They **play¹⁶⁶**...

Topic 77

word	prob.
MUSIC	.090
DANCE	.034
SONG	.033
PLAY	.030
SING	.026
SINGING	.026
BAND	.026
PLAYED	.023
SANG	.022
SONGS	.021
DANCING	.020
PIANO	.017
PLAYING	.016
RHYTHM	.015
ALBERT	.013
MUSICAL	.013

Topic 82

word	prob.
LITERATURE	.031
POEM	.028
POETRY	.027
POET	.020
PLAYS	.019
POEMS	.019
PLAY	.015
LITERARY	.013
WRITERS	.013
DRAMA	.012
WROTE	.012
POETS	.011
WRITER	.011
SHAKESPEARE	.010
WRITTEN	.009
STAGE	.009

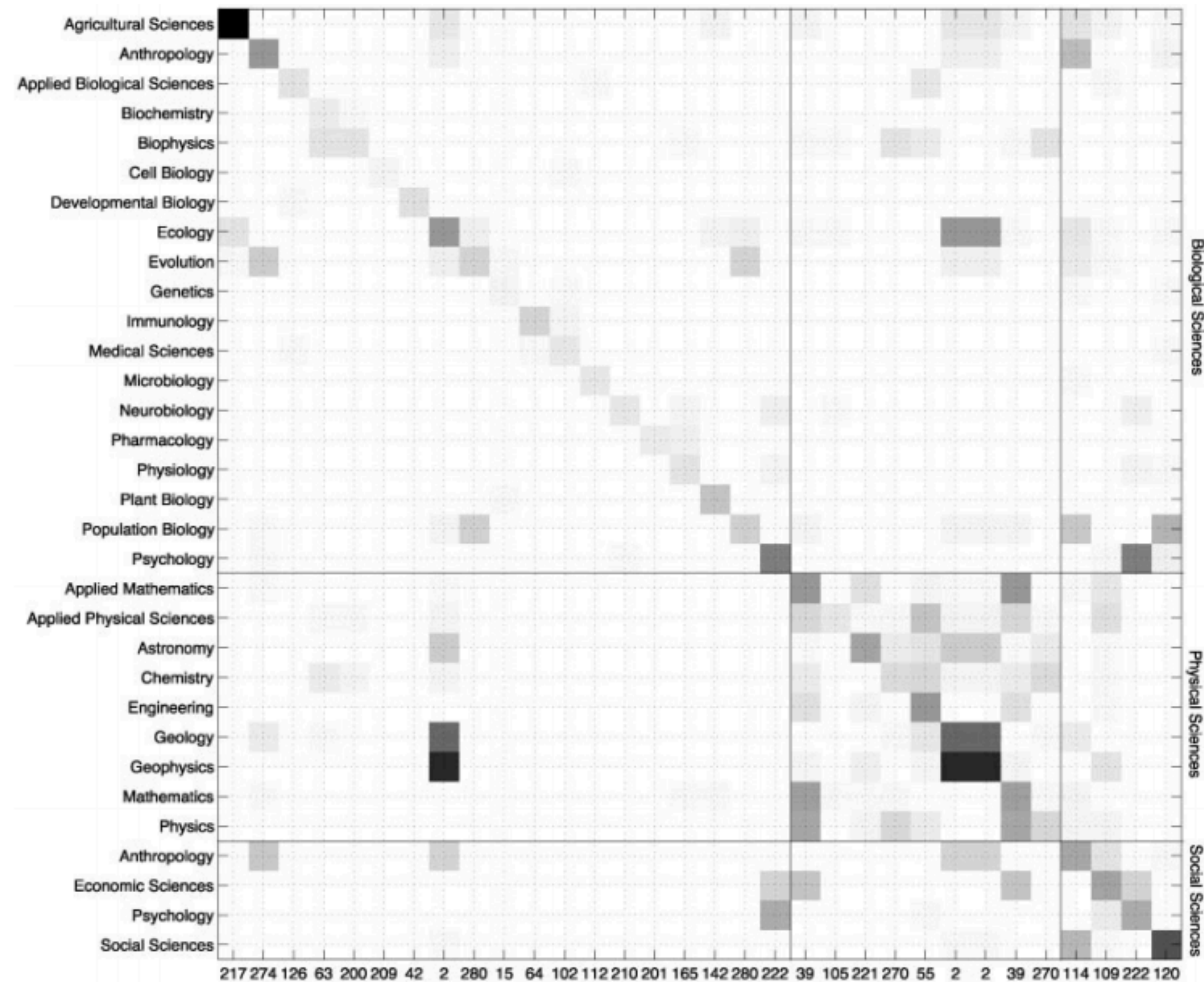
Topic 166

word	prob.
PLAY	.136
BALL	.129
GAME	.065
PLAYING	.042
HIT	.032
PLAYED	.031
BASEBALL	.027
GAMES	.025
BAT	.019
RUN	.019
THROW	.016
BALLS	.015
TENNIS	.011
HOME	.010
CATCH	.010
FIELD	.010

Training topic models on twitter

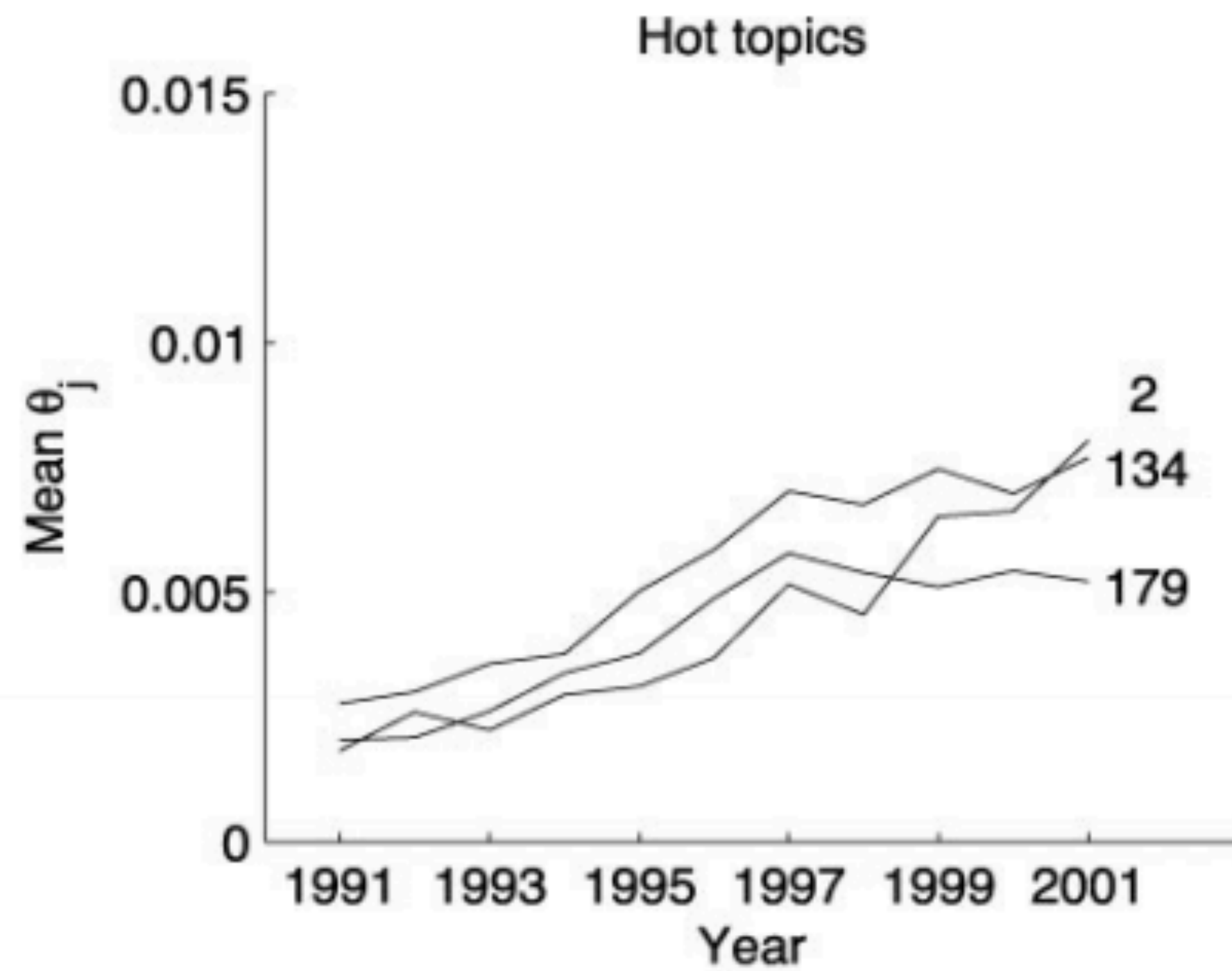
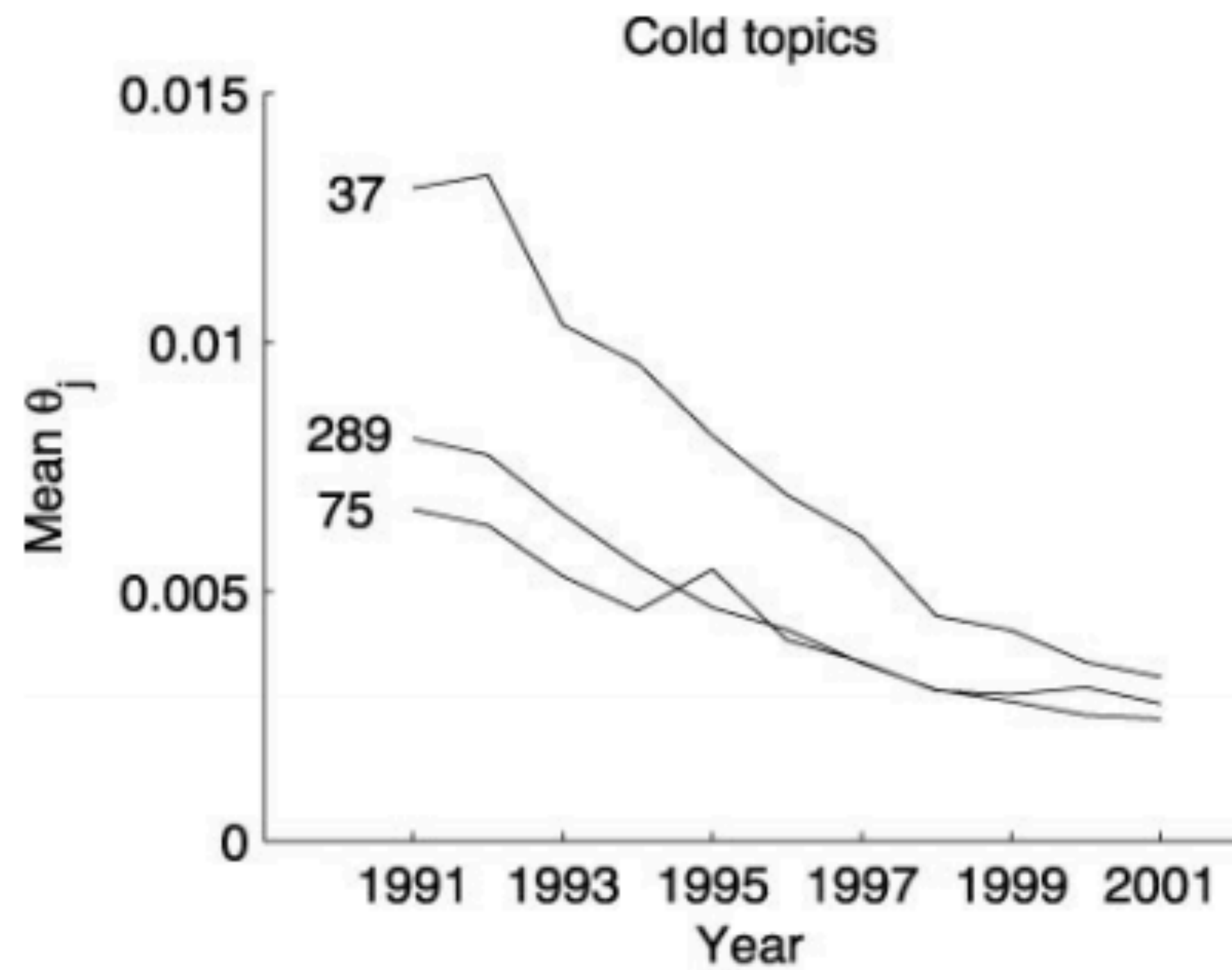
<https://ermooore.shinyapps.io/twittr/>

Finding scientific topics (Griffiths & Steyvers, 2004)



217 INSECT MYB PHEROMONE LENS LARVAE	274 SPECIES PHYLOGENETIC EVOLUTION EVOLUTIONARY SEQUENCES	126 GENE VECTOR VECTORS EXPRESSION TRANSFER	63 STRUCTURE ANGSTROM CRYSTAL RESIDUES STRUCTURES	200 FOLDING NATIVE PROTEIN STATE ENERGY	209 NUCLEAR NUCLEUS LOCALIZATION CYTOPLASM EXPORT
42 NEURAL DEVELOPMENT DORSAL EMBRYOS VENTRAL	2 SPECIES GLOBAL CLIMATE CO2 WATER	280 SPECIES SELECTION EVOLUTION GENETIC POPULATIONS	15 CHROMOSOME REGION CHROMOSOMES KB MAP	64 CELLS CELL ANTIGEN LYMPHOCYTES CD4	102 TUMOR CANCER TUMORS HUMAN CELLS
112 HOST BACTERIAL BACTERIA STRAINS SALMONELLA	210 SYNAPTIC NEURONS POSTSYNAPTIC HIPPOCAMPAL SYNAPSES	201 RESISTANCE RESISTANT DRUG DRUGS SENSITIVE	165 CHANNEL CHANNELS VOLTAGE CURRENT CURRENTS	142 PLANTS PLANT ARABIDOPSIS TOBACCO LEAVES	222 CORTEX BRAIN SUBJECTS TASK AREAS
39 THEORY TIME SPACE GIVEN PROBLEM	105 HAIR MECHANICAL MB SENSORY EAR	221 LARGE SCALE DENSITY OBSERVED OBSERVATIONS	270 TIME SPECTROSCOPY NMR SPECTRA TRANSFER	55 FORCE SURFACE MOLECULES SOLUTION SURFACES	114 POPULATION POPULATIONS GENETIC DIVERSITY ISOLATES
		109 RESEARCH NEW INFORMATION UNDERSTANDING PAPER	120 AGE OLD AGING LIFE YOUNG		

Finding scientific topics (Griffiths & Steyvers, 2004)



37
CDNA
AMINO
SEQUENCE
ACID
PROTEIN
ISOLATED
ENCODING
CLONED
ACIDS
IDENTITY
CLONE
EXPRESSED

289
KDA
PROTEIN
PURIFIED
MOLECULAR
MASS
CHROMATOGRAPHY
POLYPEPTIDE
GEL
SDS
BAND
APPARENT
LABELED

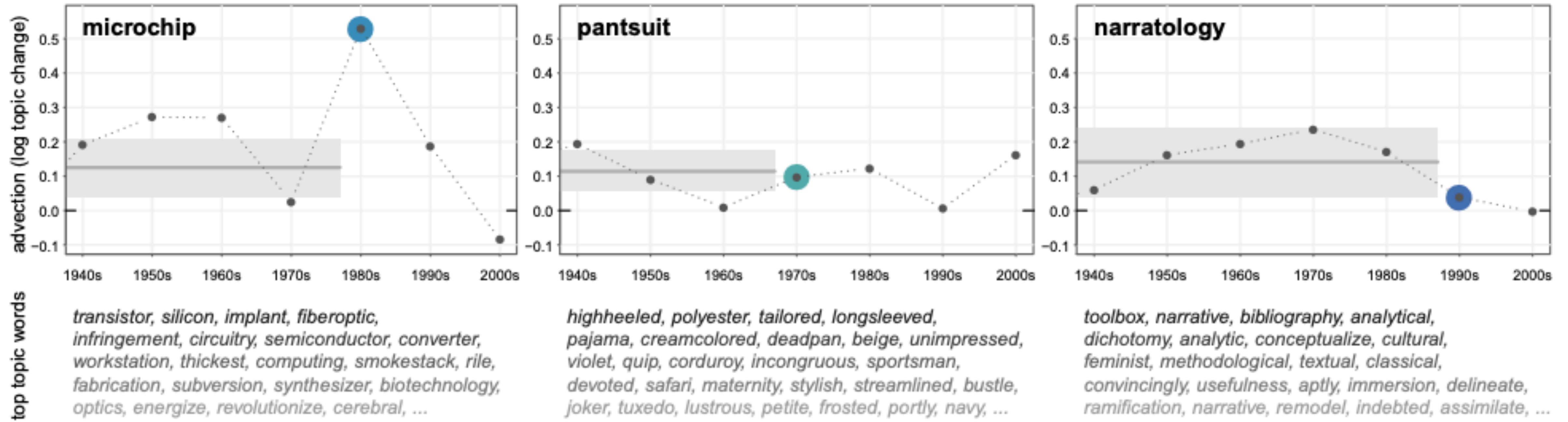
75
ANTIBODY
ANTIBODIES
MONOCLONAL
ANTIGEN
IGG
MAB
SPECIFIC
EPITOPE
HUMAN
MABS
RECOGNIZED
SERA

2
SPECIES
GLOBAL
CLIMATE
CO2
WATER
ENVIRONMENTAL
YEARS
MARINE
CARBON
DIVERSITY
OCEAN
EXTINCTION

134
MICE
DEFICIENT
NORMAL
GENE
NULL
MOUSE
TYPE
HOMOZYGOUS
ROLE
KNOCKOUT
DEVELOPMENT
GENERATED

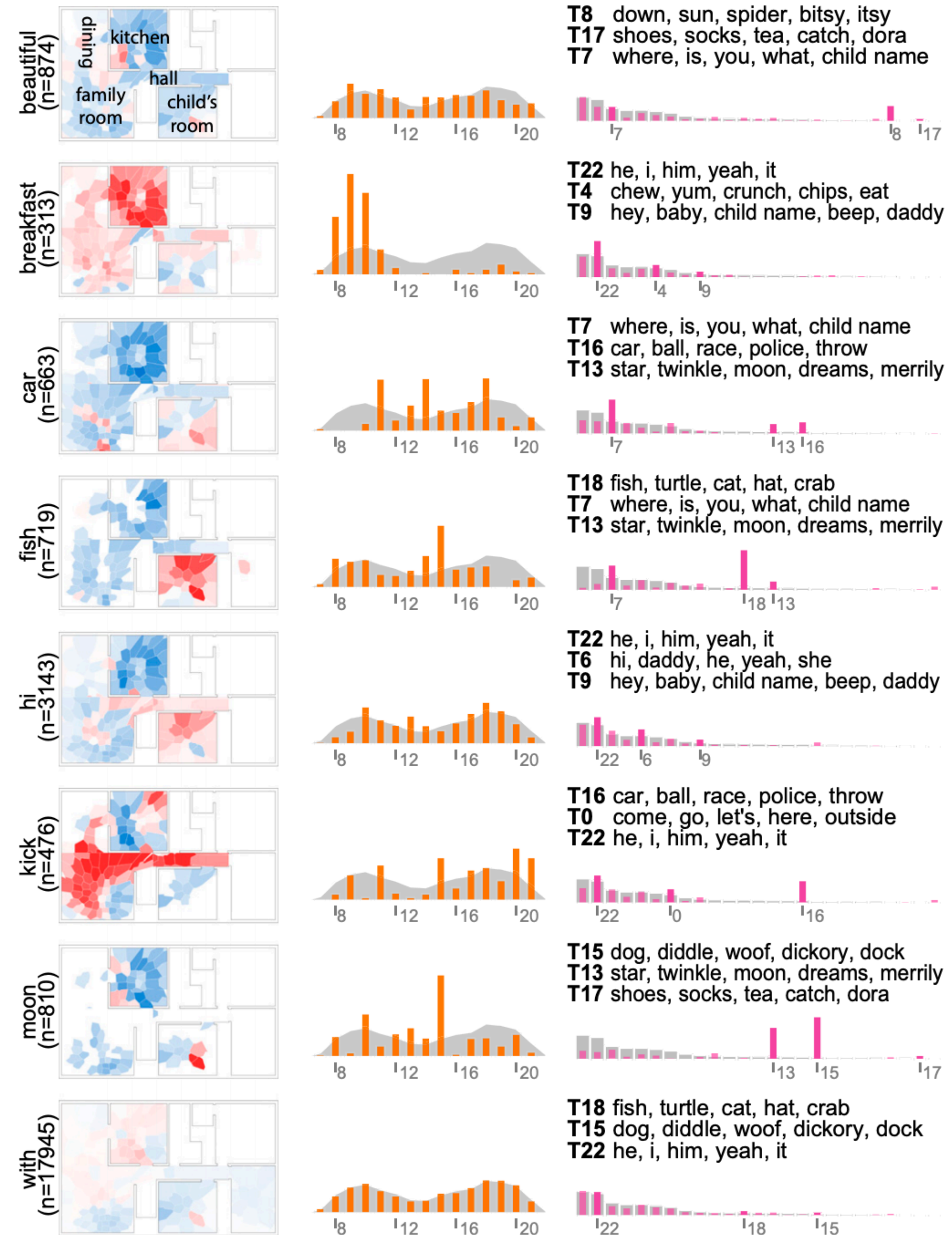
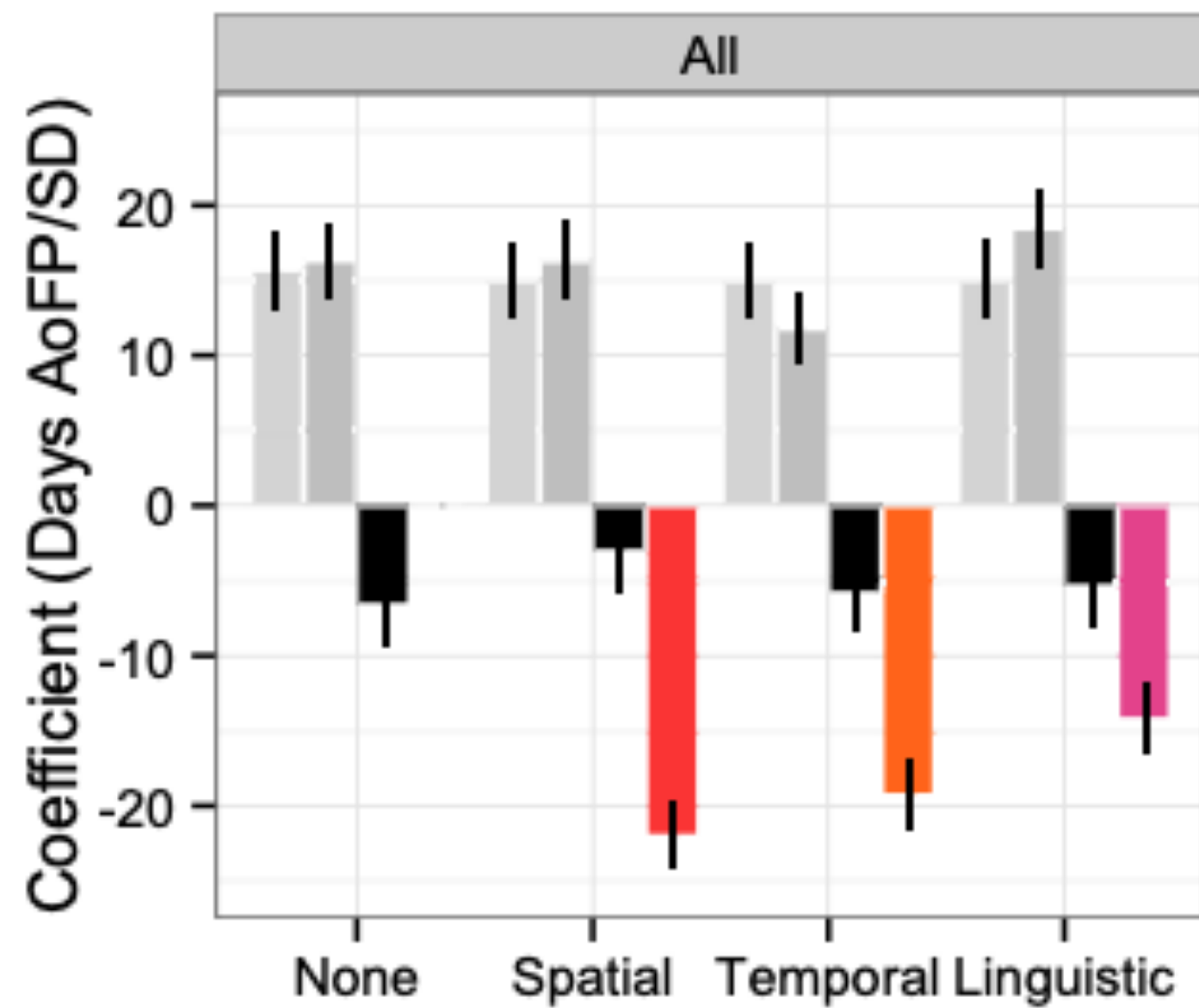
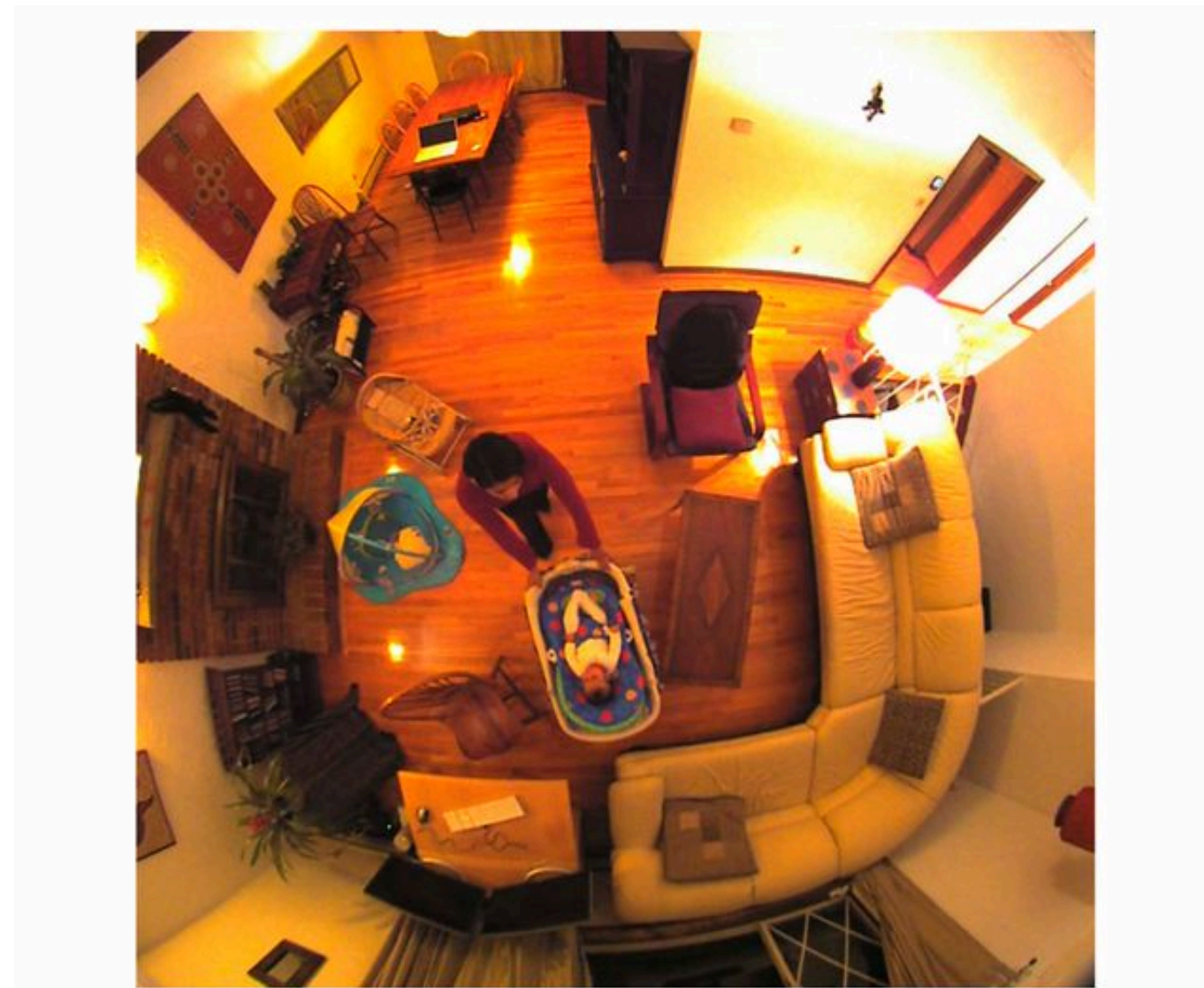
179
APOPTOSIS
DEATH
CELL
INDUCED
BCL
CELLS
APOPTOTIC
CASPASE
FAS
SURVIVAL
PROGRAMMED
MEDIATED

When do words come into language? (Karjus, Blythe, Kirby, Smith, 2020)



58% of words are coined when their topic is trending

Predicting the birth of words (Roy et al. 2015)



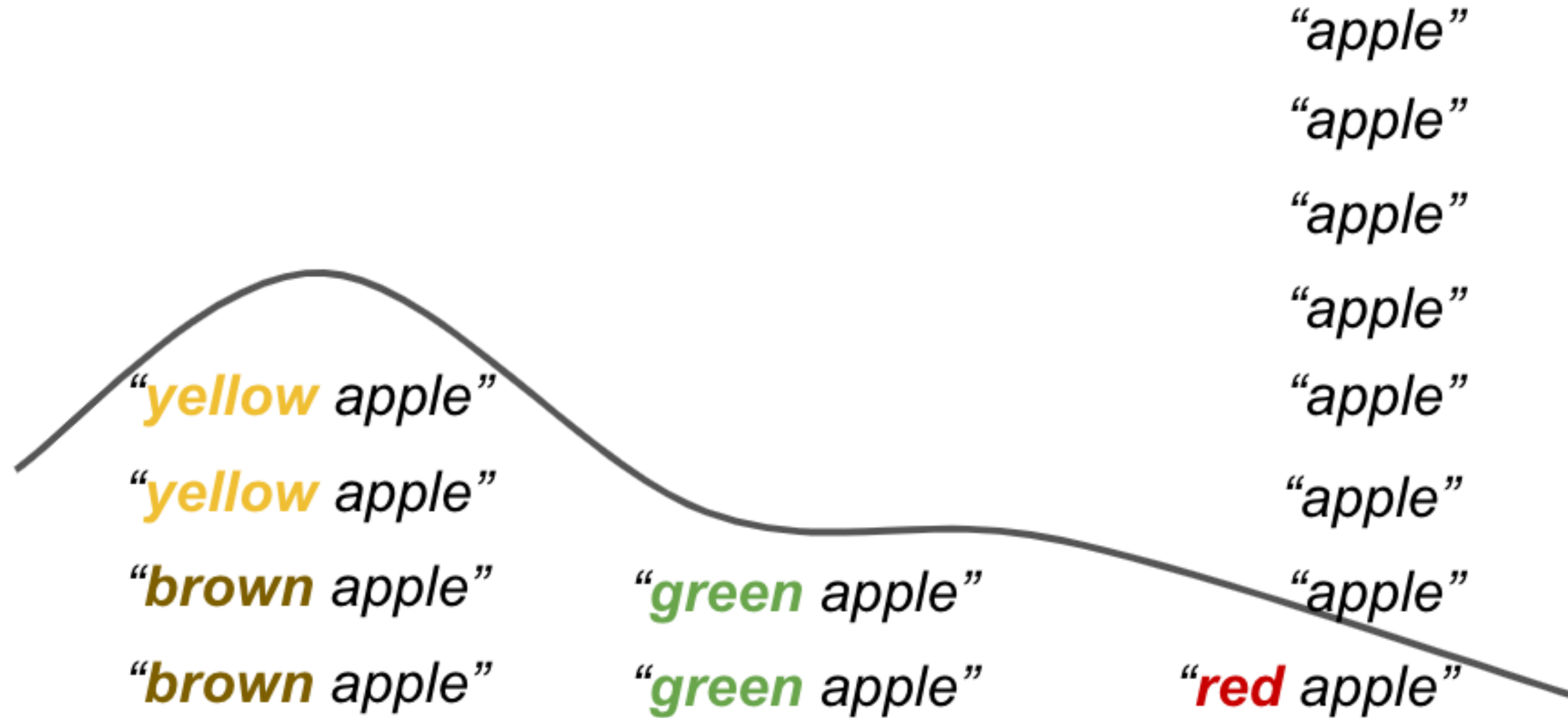
What's missing from language?



What's missing from language?



We don't narrate the world! We talk about atypical things.



What does speech to children look like? (Bergey, Morris, & Yurovsky, 2019)

don't put a **plastic cup** in there.

I like your **scaly skin**.

are **bananas yellow**?

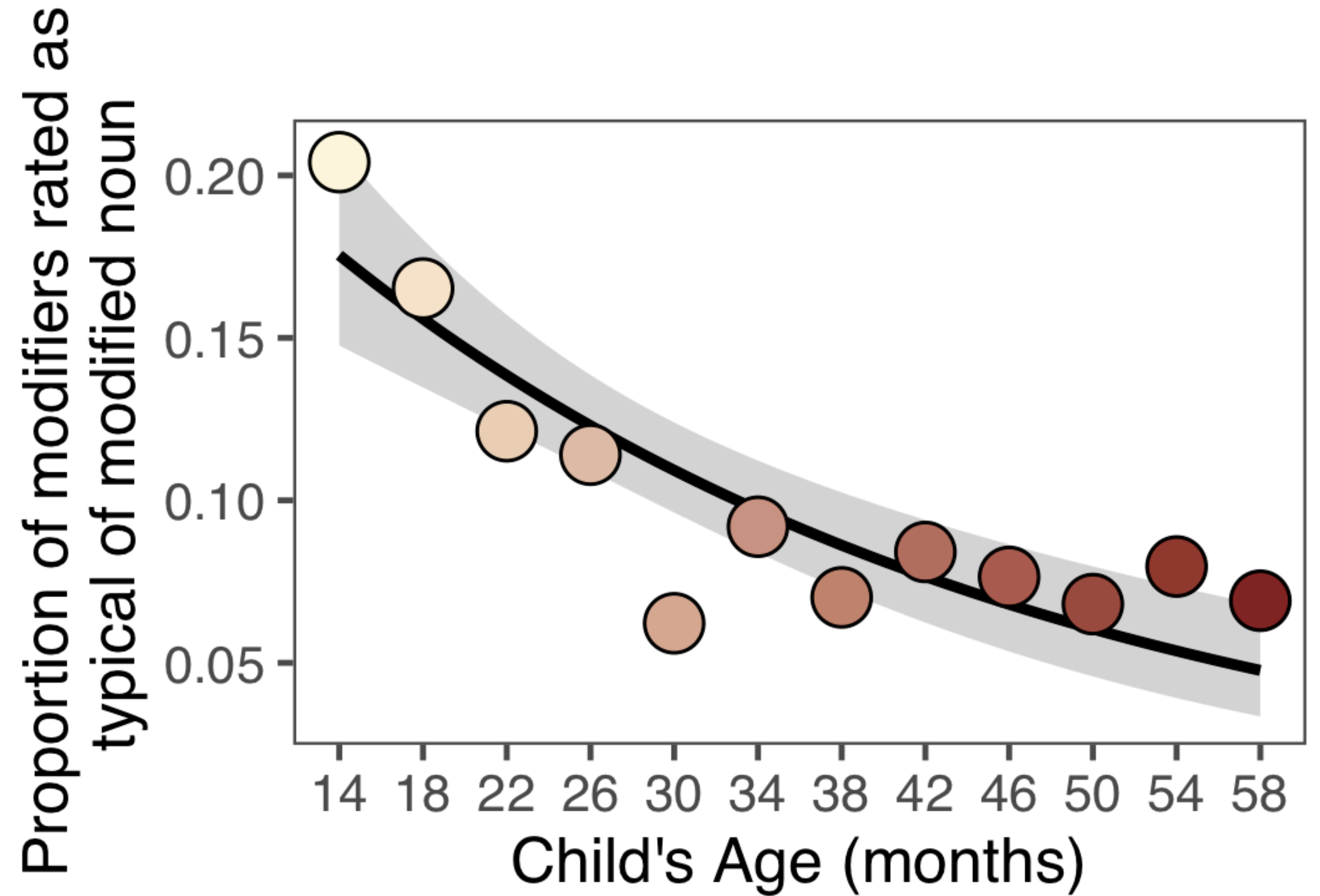
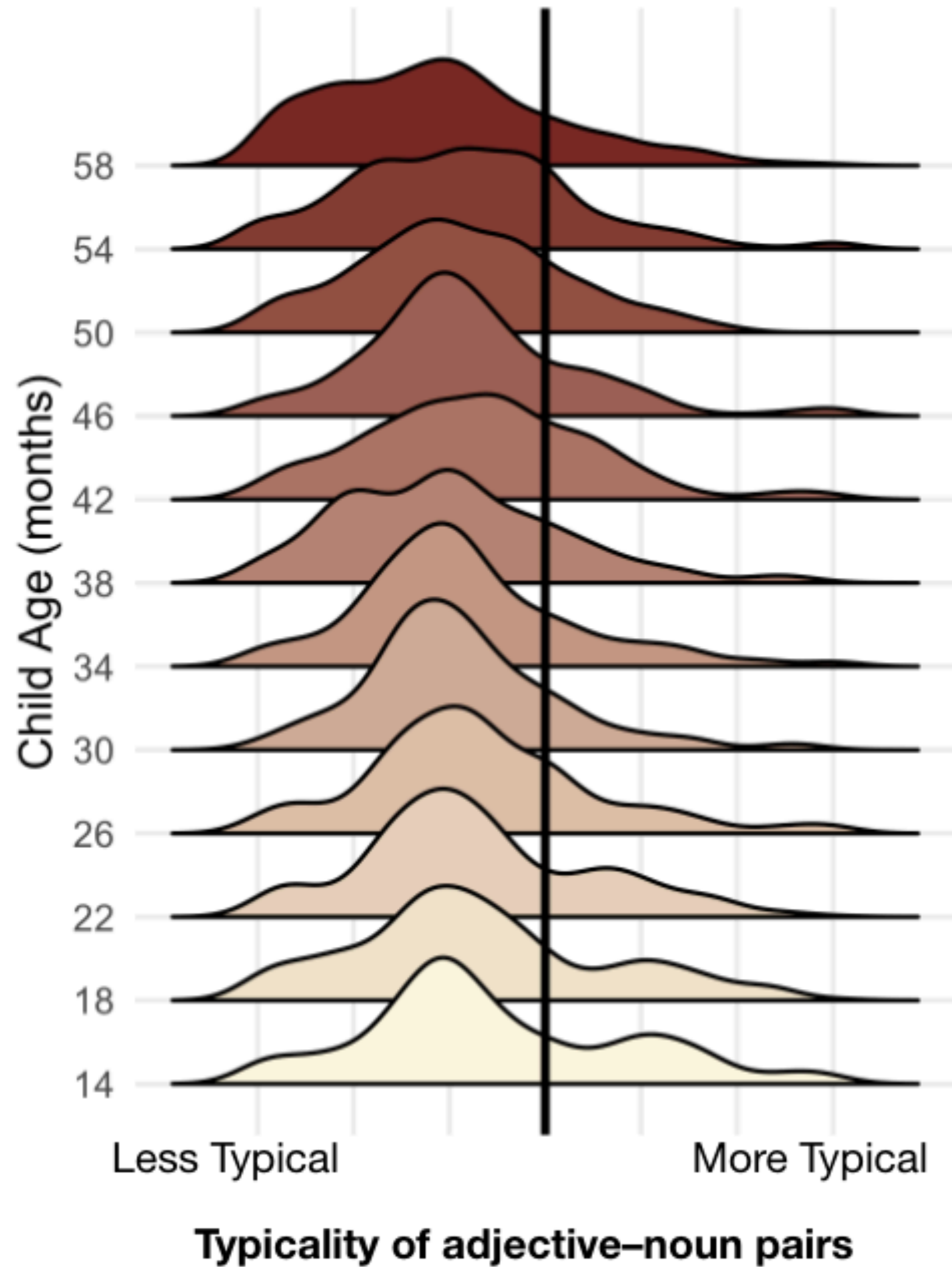
oh, look, are you giving him some nice **curly hair**?

here I'll put it on because your **hands** are **chalky**.

How common is it for an [apple] to be a
[yellow apple]?

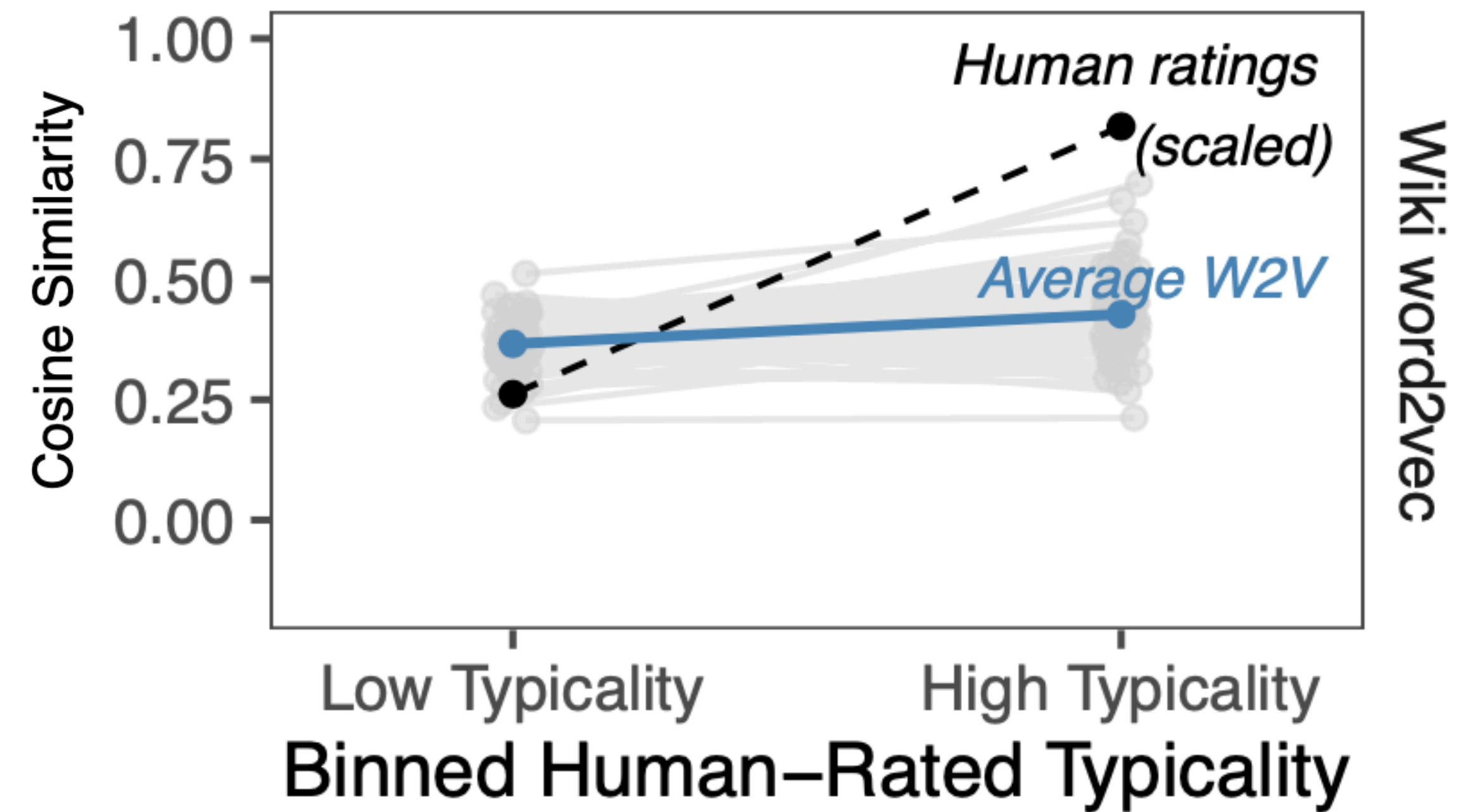
never	rarely	sometimes	about half the time	often	almost always	always
(1)	(2)	(3)	(4)	(5)	(6)	(7)

Children hear more about atypical colors



Semantic space models can't recover typical semantics

noun	typical adjective	atypical adjective
puzzle	flat	giant
apple	red	brown
bird	outside	purple
elephant	fat	pink
whale	wet	red
frog	green	purple



- 1. You can learn a lot from the co-occurrence structure of words in language**
- 2. Latent semantic analysis and Topics models both use this structure to learn about the world**
- 3. But some information is not (straightforwardly) in the co-occurrence structure of language**