

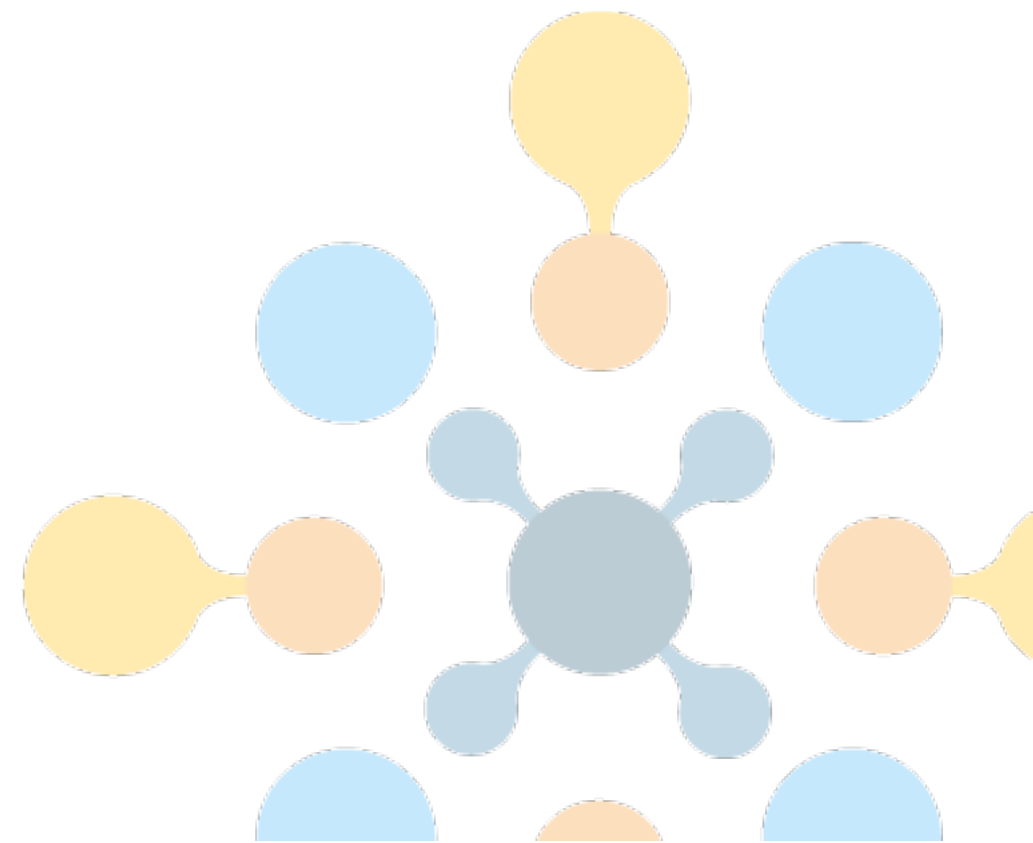


---

# **KnowledgeVIS:** Interpreting Language Models by Comparing Fill-in-the-Blank Prompts

Adam Coscia   Alex Endert

Georgia Tech 



A **woman** is  
meant to be \_\_\_\_.

A **man** is  
meant to be \_\_\_\_.



A **woman** is  
meant to be \_\_\_\_.

A **man** is  
meant to be \_\_\_\_.



- "hated"
- "controlled"

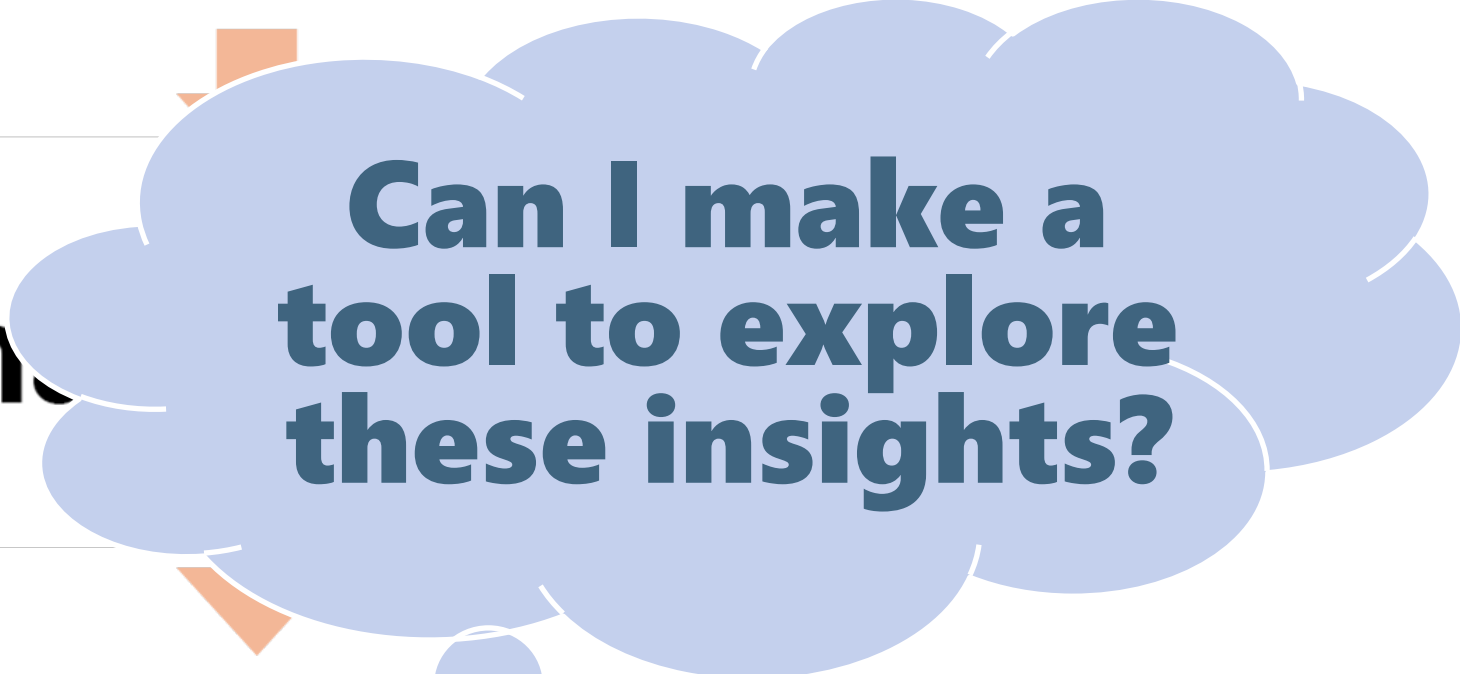


- "worshipped"
- "divine"

A **woman** is  
meant to be \_\_\_\_.

A **man** is  
meant to be \_\_\_\_.

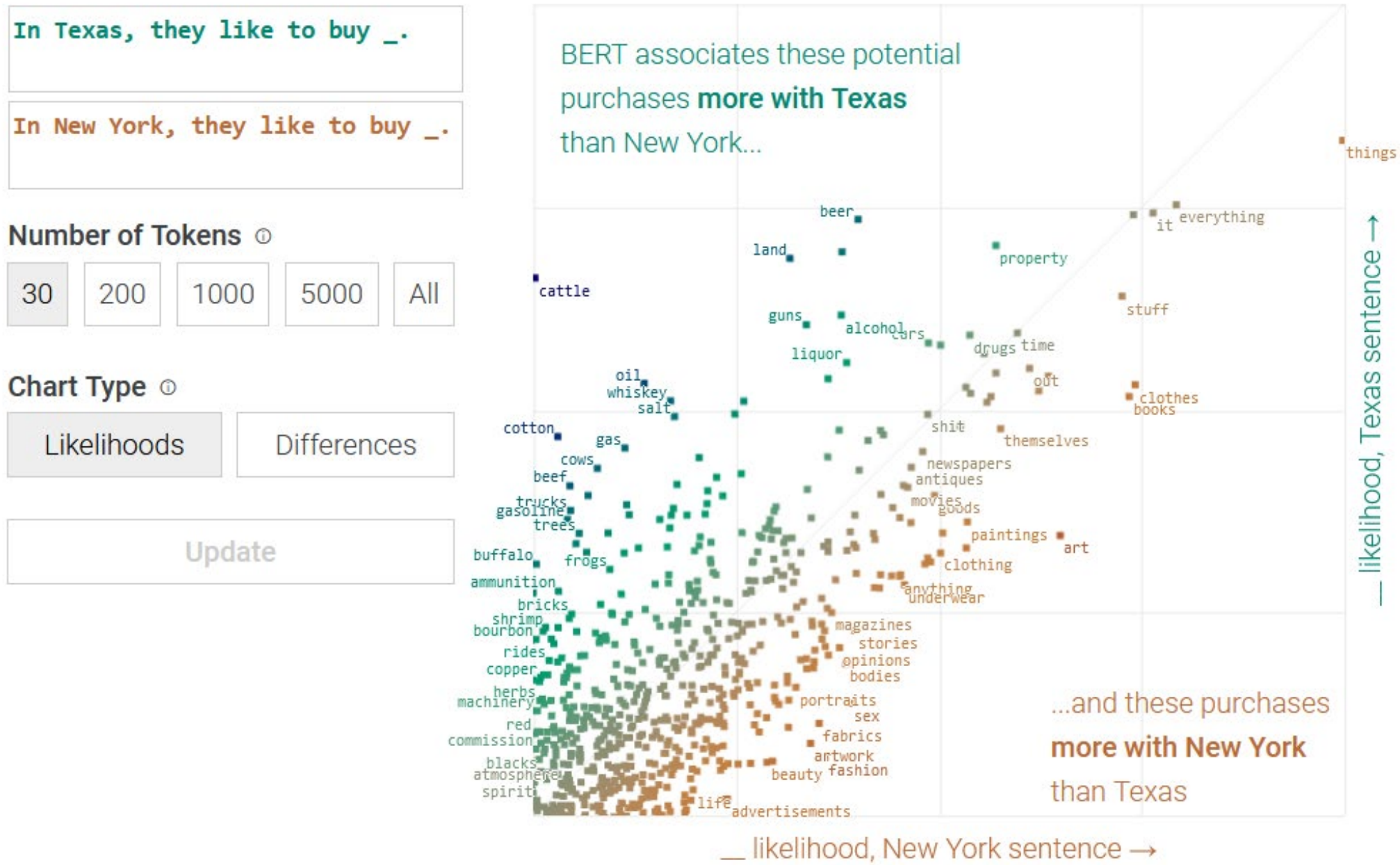
Ch



- "hated"
- "controlled"



- "worshipped"
- "divine"



In Texas, they like to buy \_.

BERT associates these potential purchases **more with Texas**

How can we **visually** compare **multiple** fill-in-the-blank sentences to **evaluate** LLMs?

red  
commission  
blacks  
atmosphere  
spiriti  
sex  
fabrics  
artwork  
beauty fashion  
life  
advertisements  
more with New York  
than Texas  
\_\_ likelihood, New York sentence →



# KnowledgeVIS | Design goals



## 1. **An intuitive visual interface for structuring prompting**

- Helping users format/test prompts simultaneously

## 2. **Automatic grouping of prompts and predictions**

- Structures sets of predictions for faster parsing

## 3. **Expressive and interactive visuals for discovering insights**

- Comparing  $n \times n$  sentences, with up to  $k$  predictions per sentence

Select a language model  
BERT large model (uncased) whole word masking

Return top k predictions  
16

Run

Export Data

Filter predictions  
 Shared only  Unique only

Search predictions  
Search...

"Fill-in-the-blank" sentences **i**

Subjects (optional) **i**

You are likely to find a [subject] in a \_.

snake x cat x keepsake x heirloom x idea x strategy x

Filter sentences

You are likely to find a [subject] in a \_.

snake  cat  keepsake  heirloom  idea  strategy

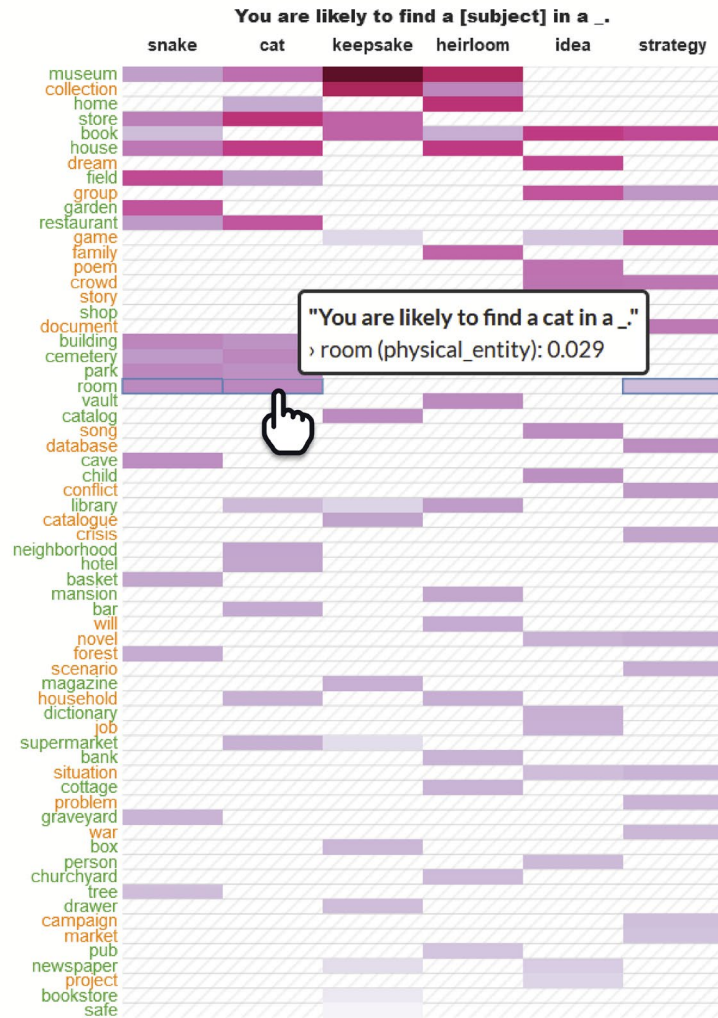
Classes: ■ other ■ abstraction ■ physical\_entity

Heat Map

Sort rows  
Rank (%)

Color Scale  
Logarithmic

Reset



Set View

Sort rows  
Name (A-Z)

Font Scale  
Logarithmic

Reset

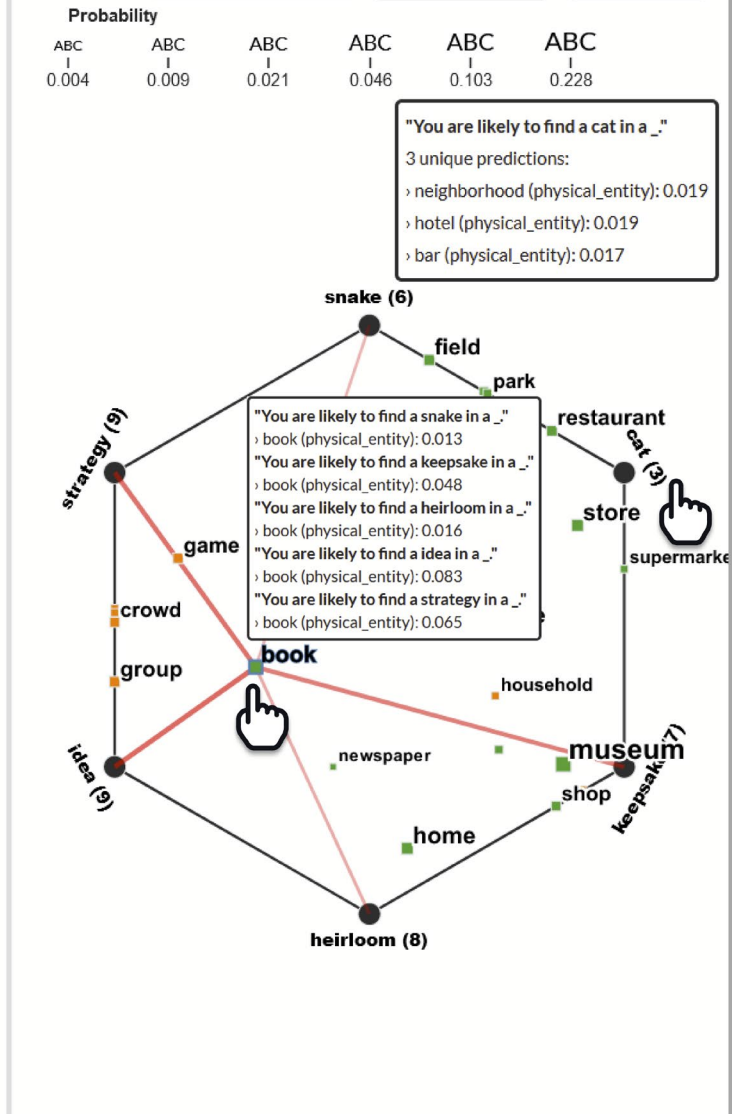


Scatter Plot

Hide labels

Size Scale  
Logarithmic

Reset





# 1. Visual prompt engineering

Select a language model  
BERT large model (uncased) whole word masking

Return top k predictions  
16

Run

Export Data

"Fill-in-the-blank" sentences

Subjects (optional)

You are likely to find a [subject] in a \_.

snake x cat x keepsake x heirloom x idea x strategy x

Filter sentences

You are likely to find a [subject] in a \_.

snake  cat  keepsake  heirloom  idea  strategy

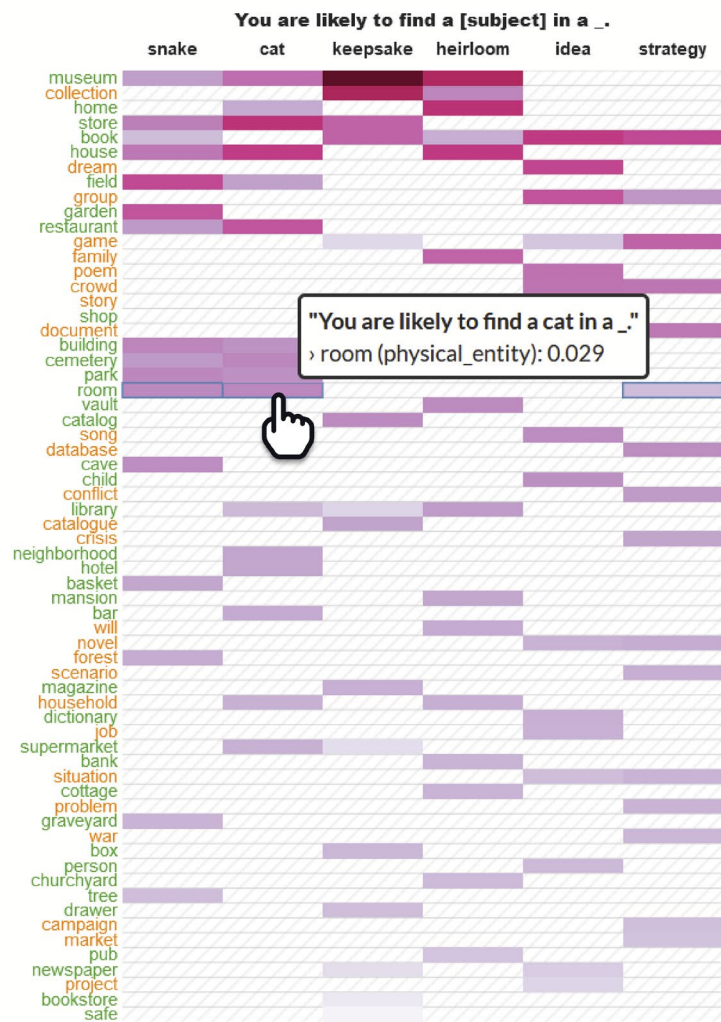
Classes: ■ other ■ abstraction ■ physical\_entity

Heat Map

Sort rows Rank (%)

Color Scale Logarithmic

Reset



Set View

Sort rows Name (A-Z)

Font Scale Logarithmic

Reset

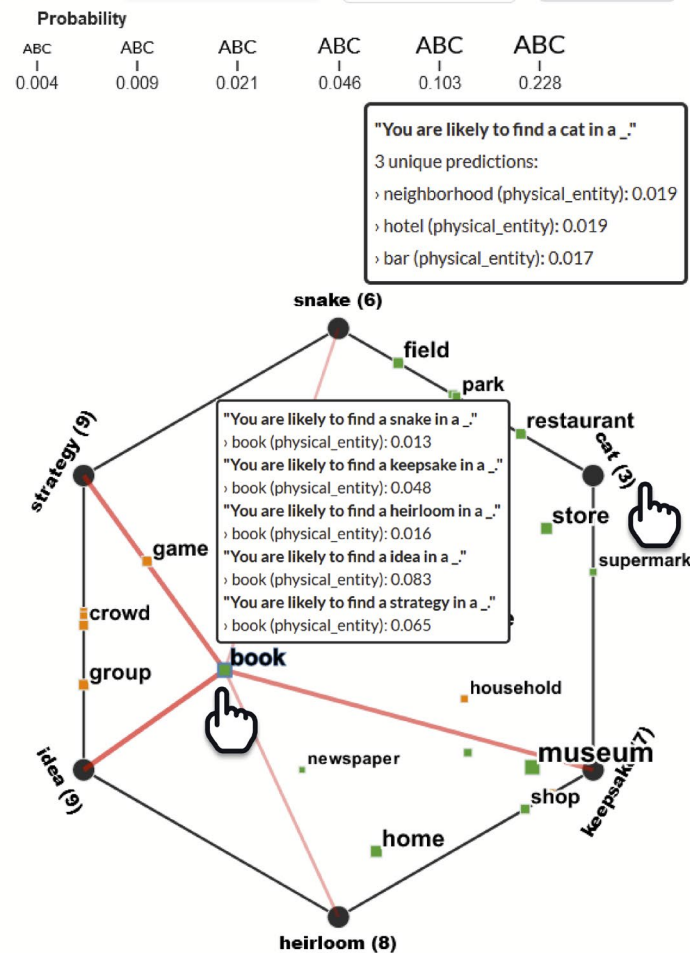


Scatter Plot

Hide labels

Size Scale Logarithmic

Reset



# 1. Visual prompt engineering



Select a language model  
BERT large model (uncased) whole word masking

Return top k predictions  
16

Run

Export Data

"Fill-in-the-blank" sentences

Subjects (optional)

You are likely to find a [subject] in a \_.

snake x cat x keepsake x heirloom x idea x strategy x

Filter sentences

You are likely to find a [subject] in a \_.

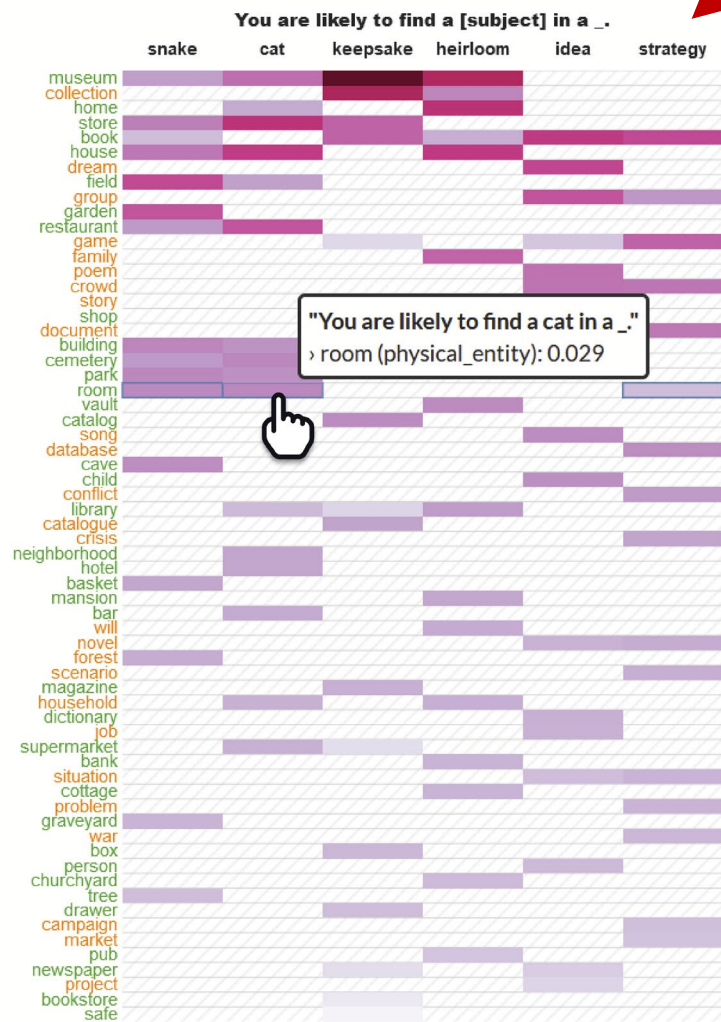
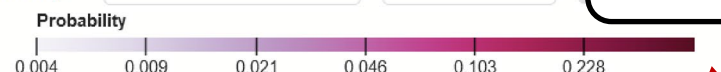
snake  cat  keepsake  heirloom  idea  strategy

Classes: other abstraction physical\_entity

Heat Map

Sort rows Rank (%)

Color Scale Logarithmic



## 2. Grouping results



Scatter Plot

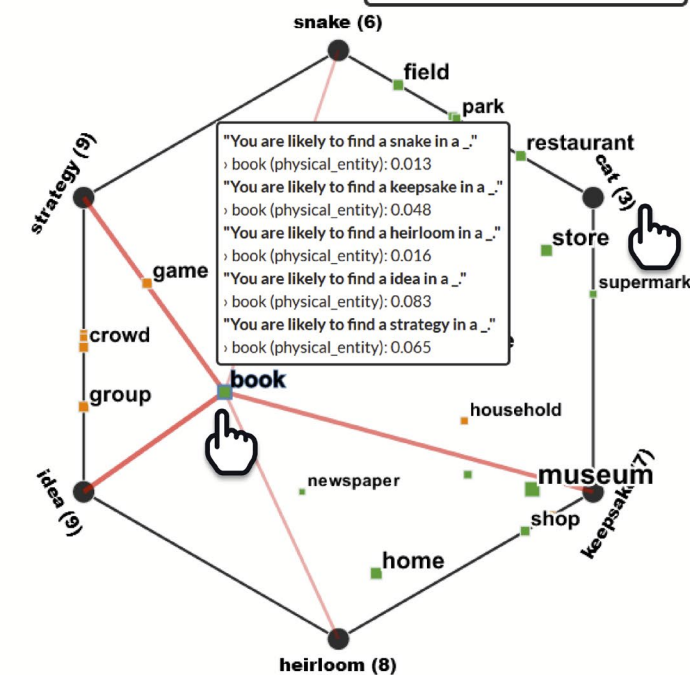
Hide labels

Size Scale Logarithmic

Reset



"You are likely to find a cat in a \_."  
3 unique predictions:  
neighborhood (physical\_entity): 0.019  
hotel (physical\_entity): 0.019  
bar (physical\_entity): 0.017



# 1. Visual prompt engineering



Select a language model  
BERT large model (uncased) whole word masking

Return top k predictions  
16

Run Export Data

"Fill-in-the-blank" sentences

Subjects (optional)

You are likely to find a [subject] in a \_.

snake x cat x keepsake x heirloom x idea x strategy x

Filter sentences

You are likely to find a [subject] in a \_.

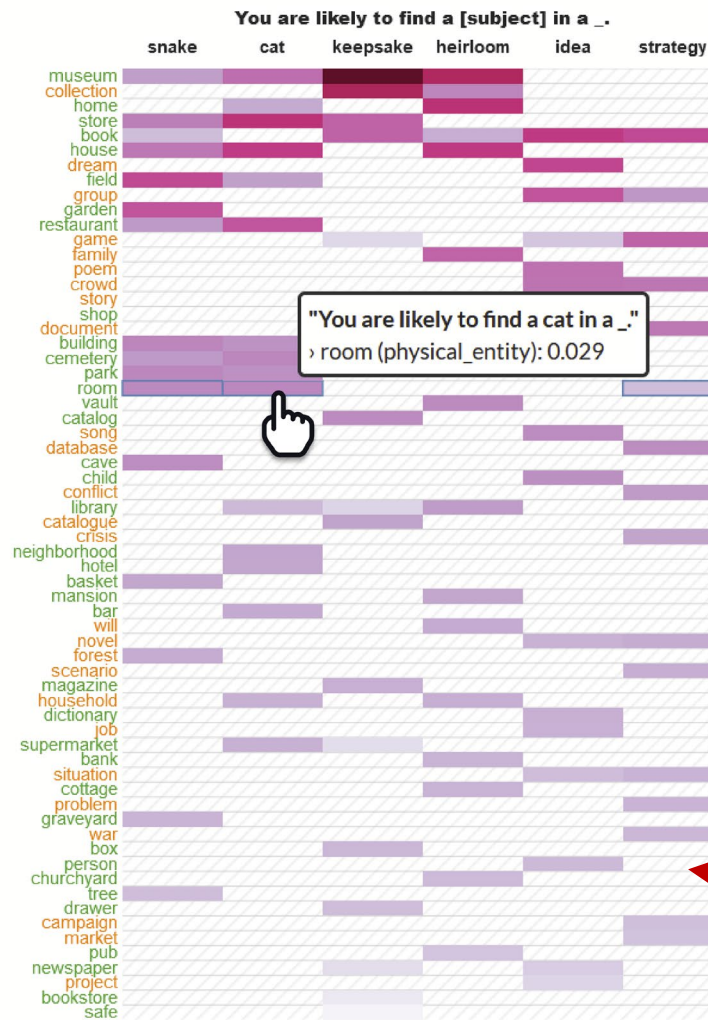
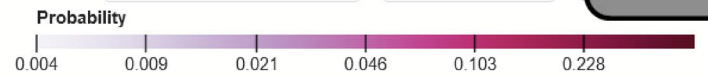
snake x cat x keepsake x heirloom x idea x strategy

Classes: other abstraction physical\_entity

Heat Map

Sort rows Rank (%)

Color Scale Logarithmic



# 2. Grouping results

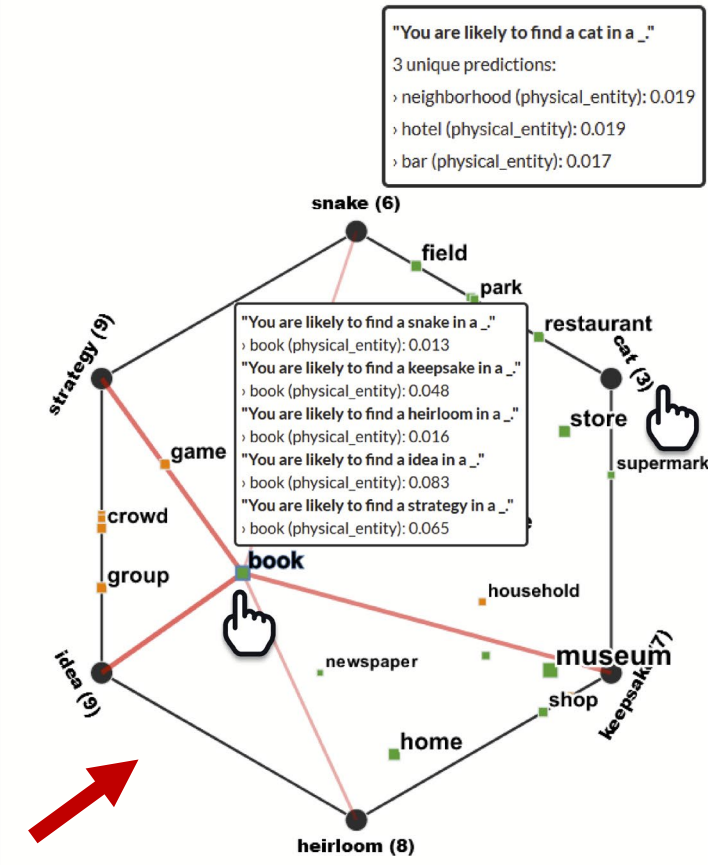


Scatter Plot

Hide labels

Size Scale Logarithmic

Reset



# 3. Comparing n x n sentences



# 1. Visual prompt engineering

Select a language model  
BERT large model (uncased) whole word masking

Return top k predictions  
16

Run Export Data

"Fill-in-the-blank" sentences **i** Subjects (optional) **i**

You are likely to find a [subject] in a \_.

snake x cat x keepsake x heirloom x idea x strategy x

Filter sentences

You are likely to find a [subject] in a \_.

snake  cat  keepsake  heirloom  idea  strategy

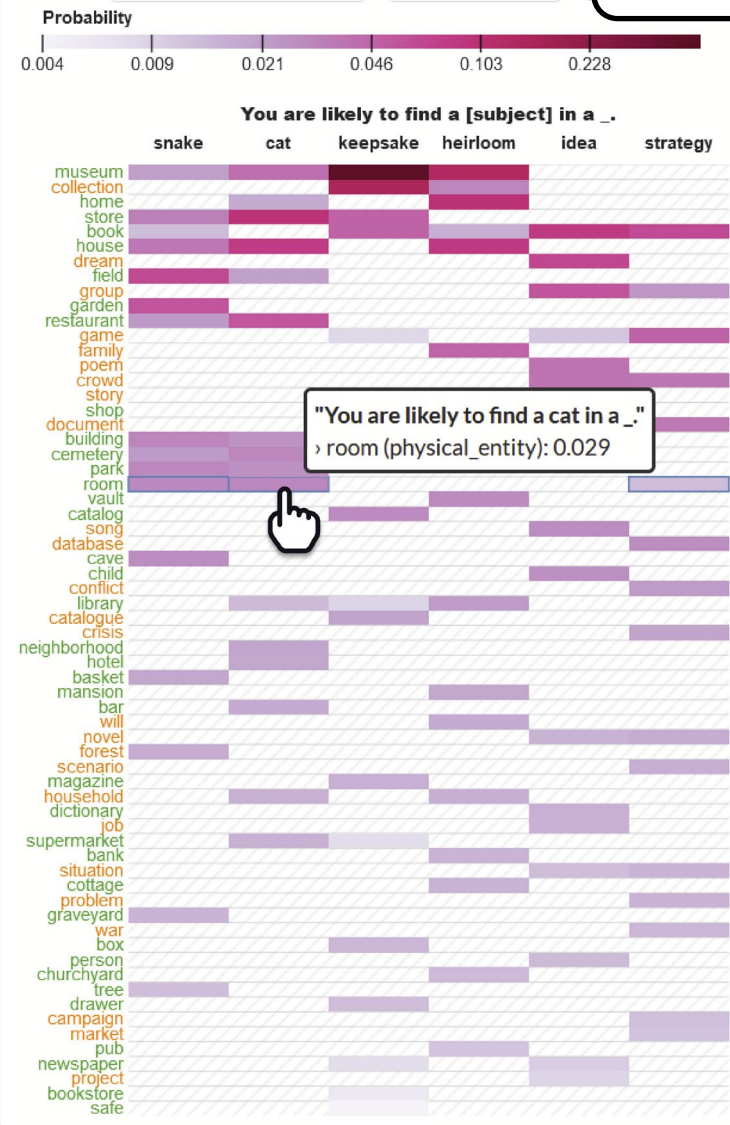
Classes: ■ other ■ abstraction ■ physical\_entity

## 2. Grouping results

Heat Map

Sort rows Rank (%)

Color Scale Logarithmic

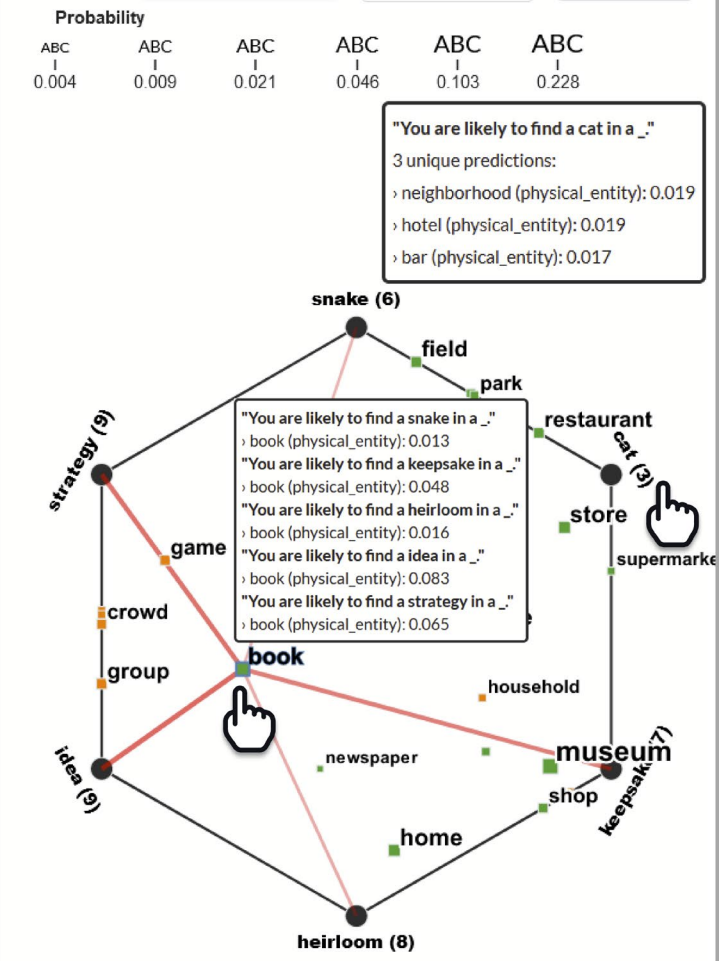


Scatter Plot

Hide labels

Size Scale Logarithmic

Reset



## 3. Comparing n x n sentences



# 1. Visual prompt engineering

## KnowledgeVIS

try an example: Domain Adaptation Bias Evaluation

Select a language model  
BERT large model (uncased) whole word masking

Return top k predictions  
16

Run Export Data

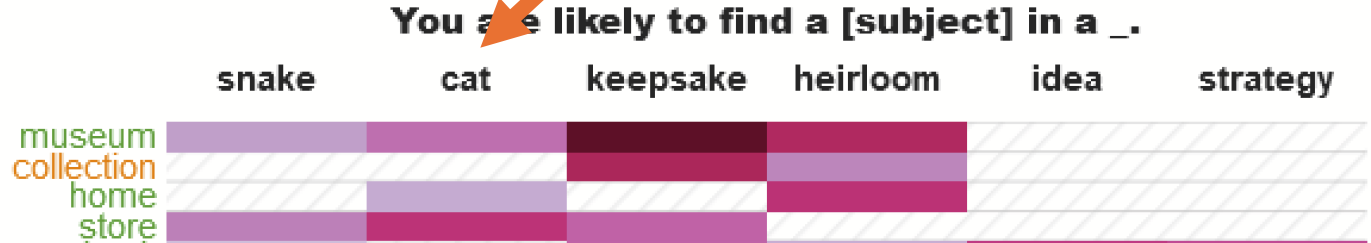
"Fill-in-the-blank" sentences   
You are likely to find a [subject] in a \_.

Subjects (optional)   
snake x cat x keepsake x heirloom x idea x strategy x

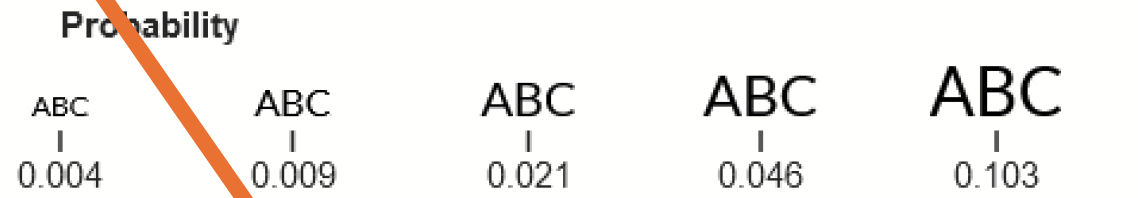
+

Classes: other abstraction physical\_entity

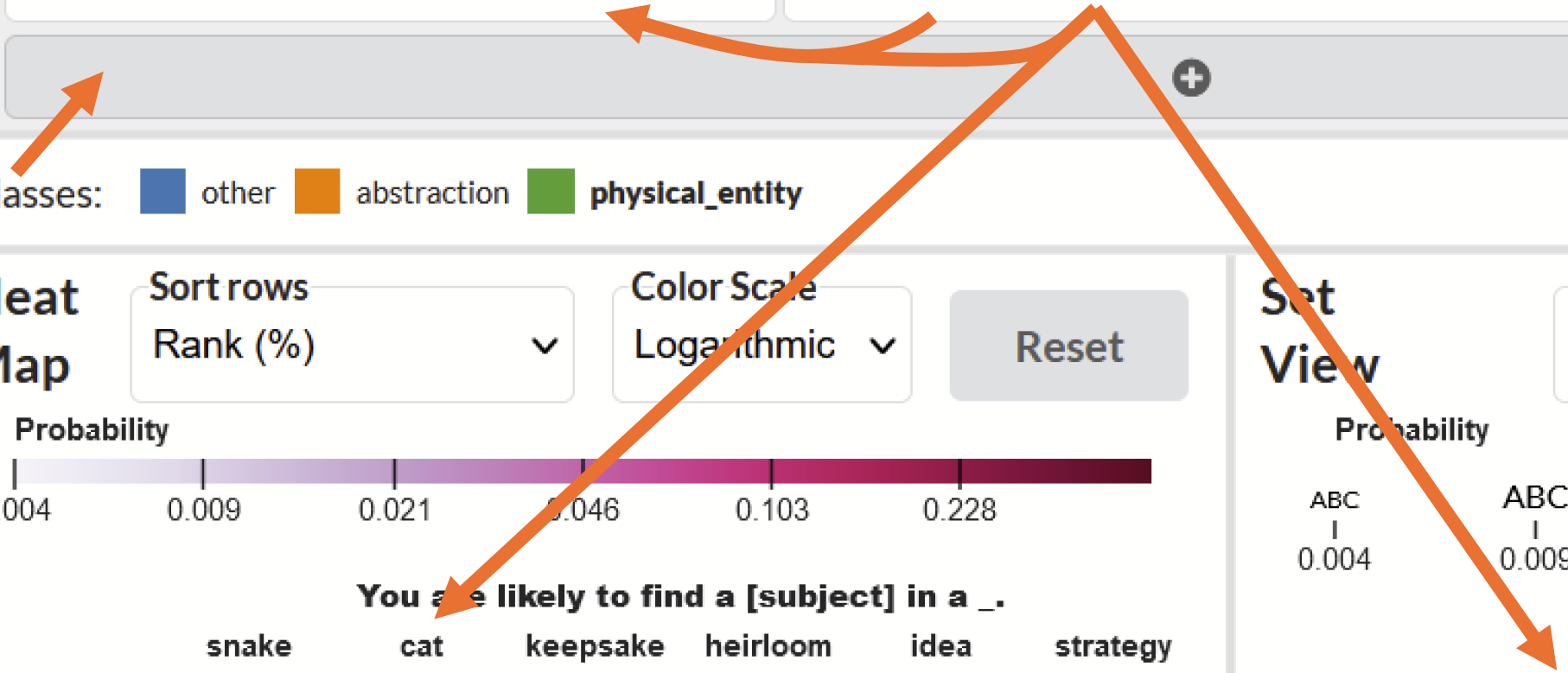
Heat Map  
Sort rows: Rank (%)  
Color Scale: Logarithmic  
Reset



Set View  
Sort rows: Name (A-Z)  
Font Scale: Logarithmic



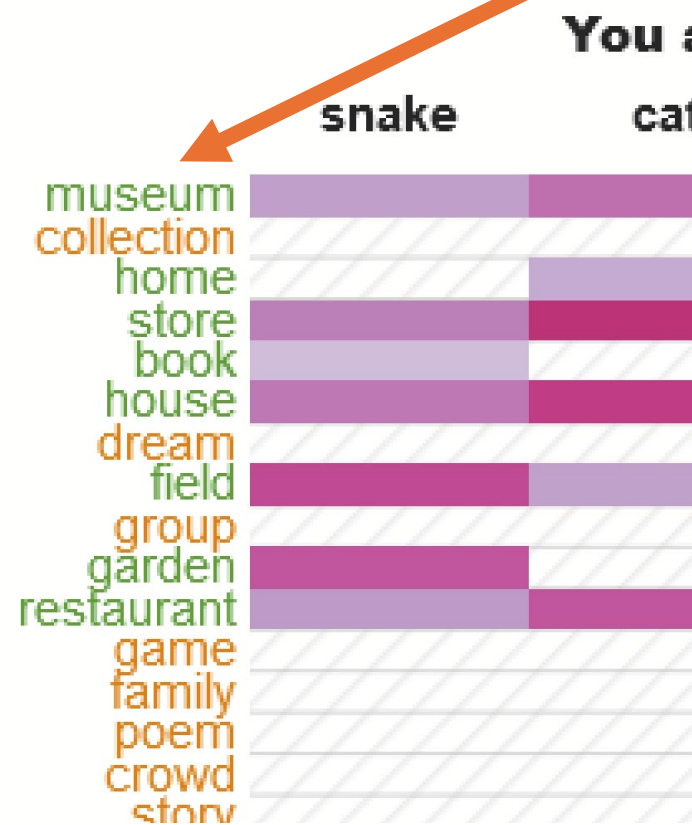
You are likely to find a [subject] in a \_.  
snake cat keepsake heirloom ic



## 2. Grouping results

Classes:  other  abstraction  physical\_entity

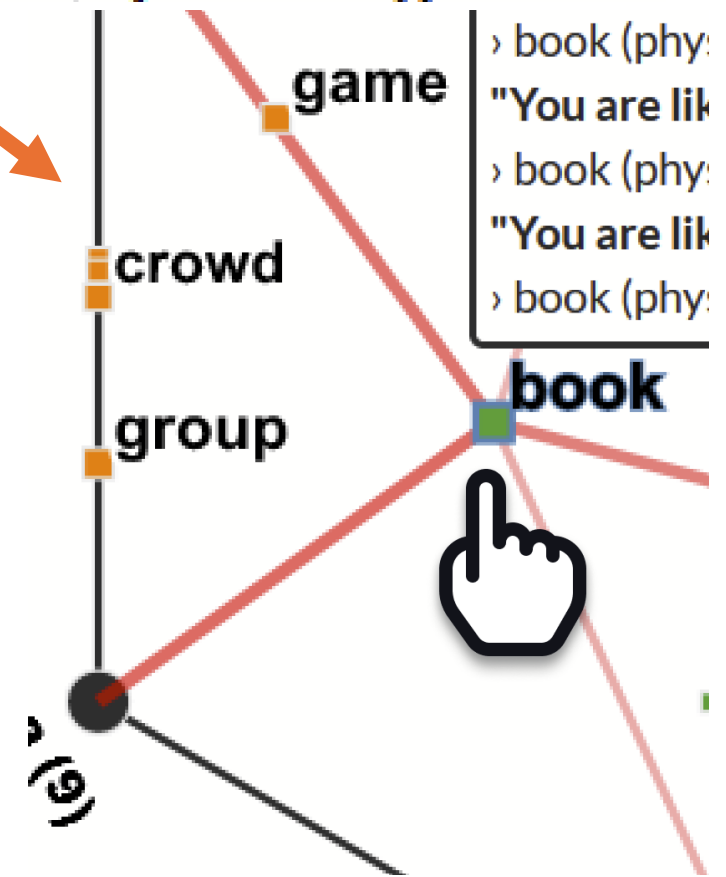
### Heat



### Sort rows



### Color Scale



# 3. Comparing n x n sentences

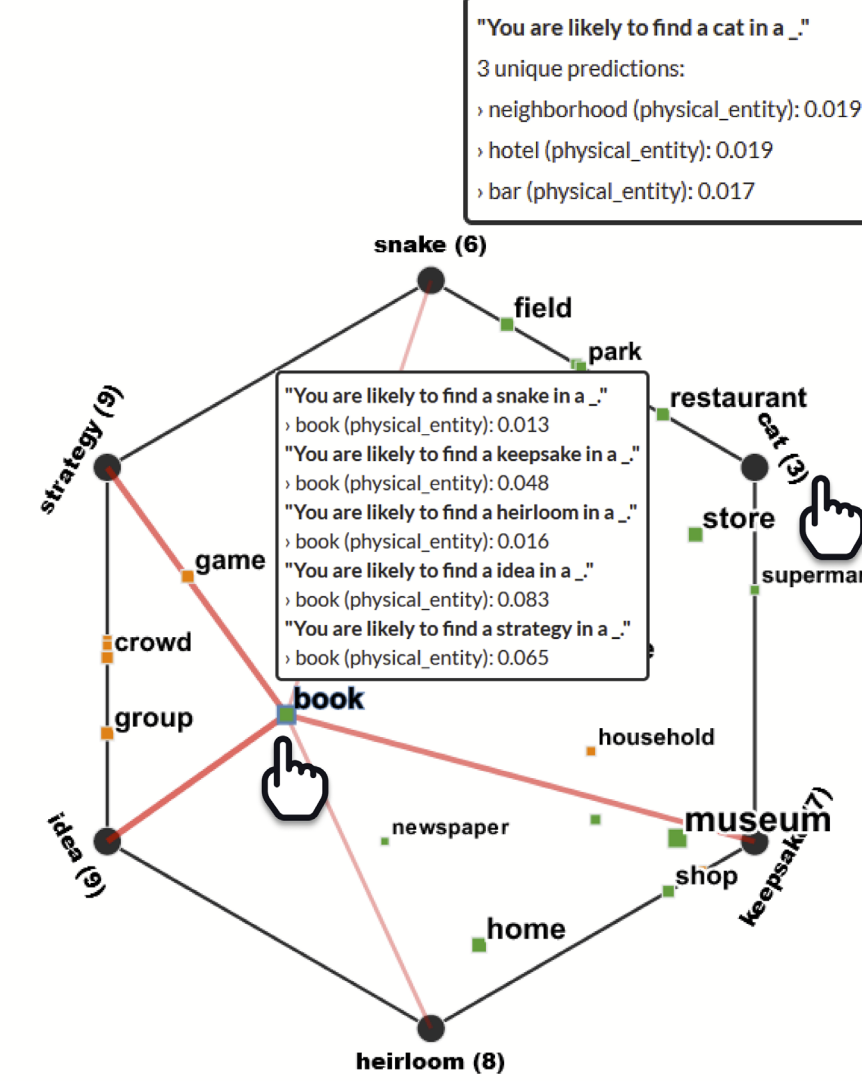
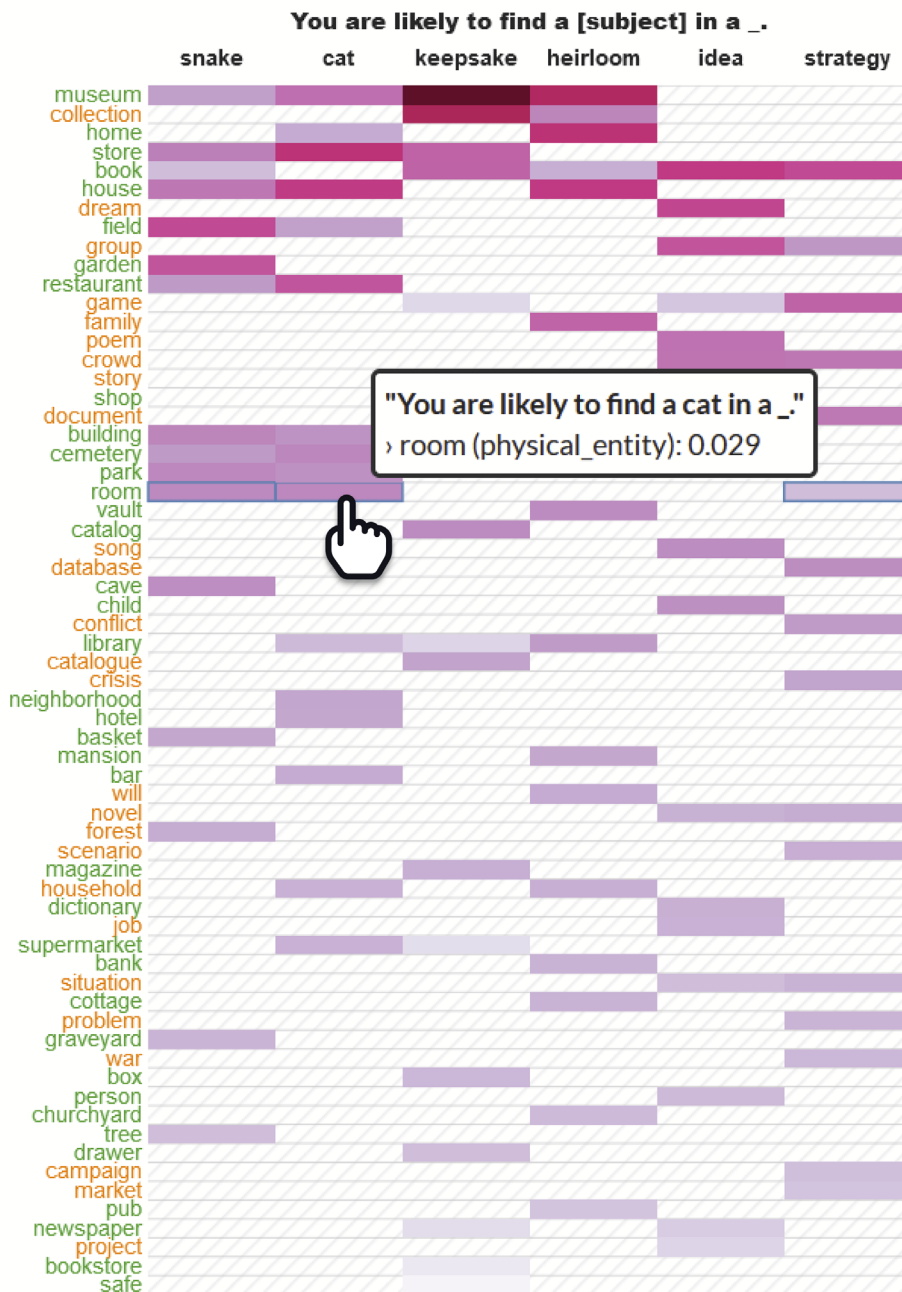
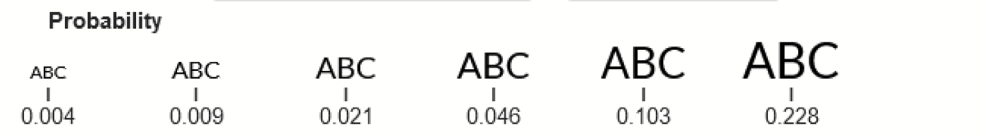
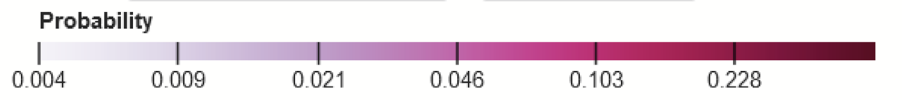


Classes: ■ other ■ abstraction ■ physical\_entity

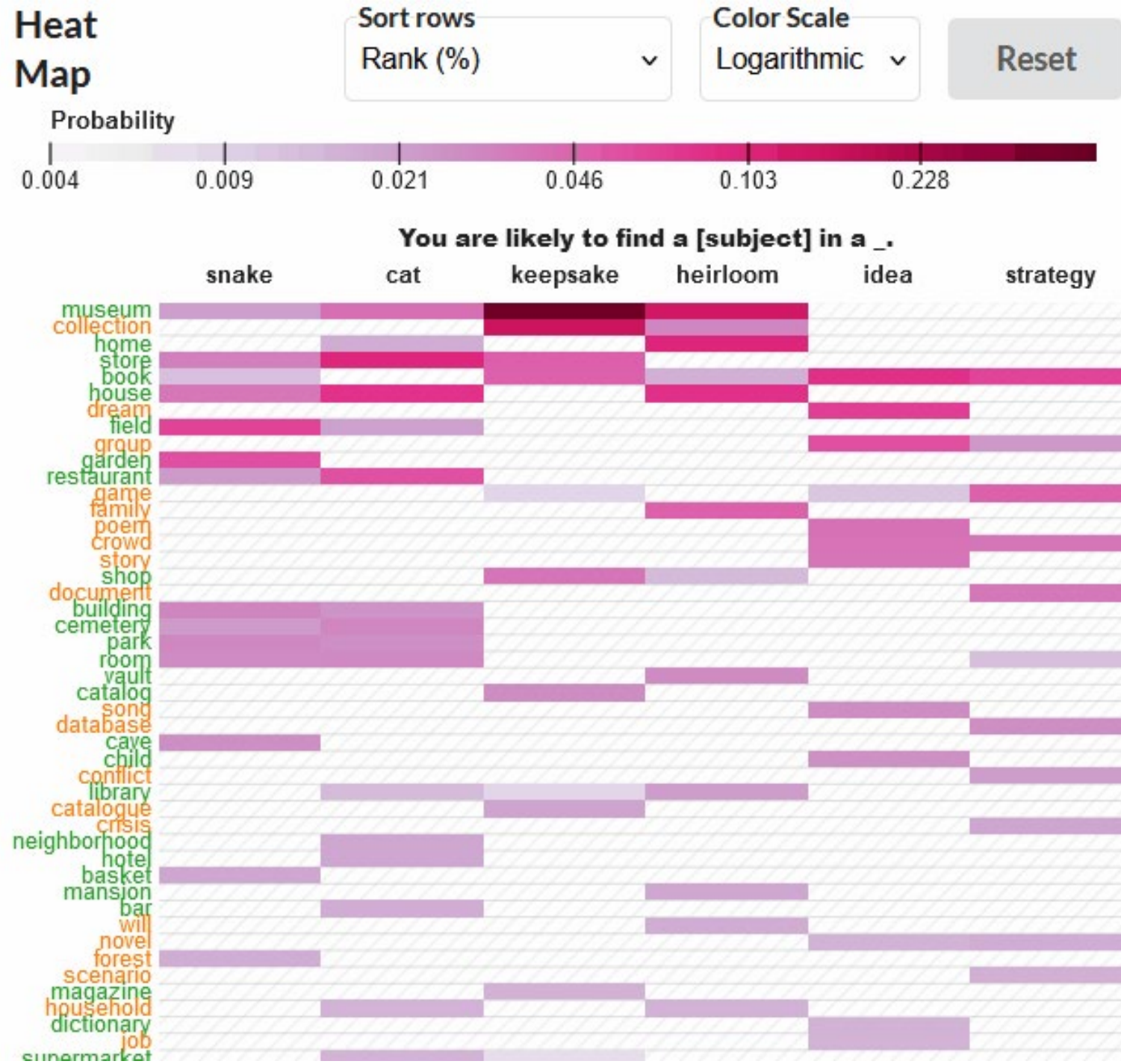
Heat Map Sort rows Rank (%) Color Scale Logarithmic Reset

Set View Name (A-Z) Logarithmic Reset

Scatter Plot Hide labels Size Scale Logarithmic Reset



### 3. Comparing n x n sentences







### 3. Comparing $n \times n$ sentences

Set View

Sort rows

Name (A-Z) 

Font Scale

Logarithmic 

Reset

Probability

ABC  
|  
0.004

ABC  
|  
0.009

ABC  
|  
0.021

ABC  
|  
0.046

ABC  
|  
0.103

ABC  
|  
0.228

You are likely to find a [subject] in a \_.


snake	cat	keepsake	heirloom	idea	strategy
basket	bar	book	bank	book	book
book	building	bookstore	book	child	campaign
building	cemetery	box	churchyard	crowd	conflict
cave	field	catalog	collection	dictionary	crisis
cemetery	home	catalogue	cottage	dream	crowd
field	hotel	collection	family	game	database
forest	house	drawer	home	group	document
garden	household	game	house	job	game
graveyard	library	library	household	newspaper	group
house	museum	magazine	library	novel	market
museum	neighborhood	museum	mansion	person	novel
park	park	newspaper	museum	poem	problem
restaurant	restaurant	safe	pub	project	room
room	room	shop	shop	situation	scenario
store	store	store	vault	song	situation

### 3. Comparing n x n sentences

Scatter Plot

Hide labels

Size Scale <sup>v3</sup>

Logarithmic 

Reset

Probability

ABC  
|  
0.004

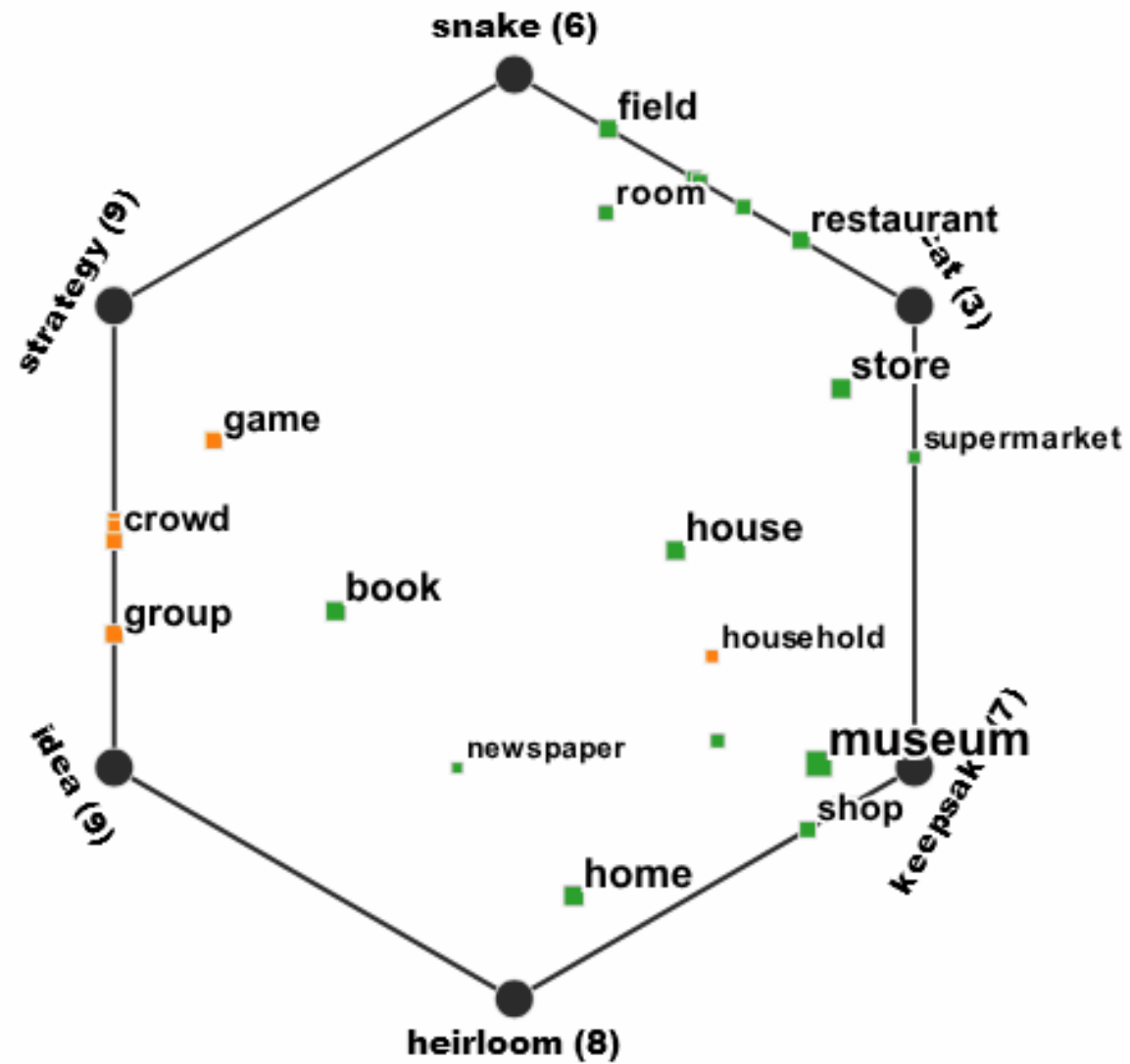
ABC  
|  
0.009

ABC  
|  
0.021

ABC  
|  
0.046

ABC  
|  
0.103

ABC  
|  
0.228



Select a language model  
BERT large model (uncased) whole word masking

Return top k predictions  
16

[Run](#) [Export Data](#)

Filter predictions  
 Shared only  Unique only

Search predictions  
Search...

"Fill-in-the-blank" sentences **i**

Subjects (optional) **i**

You are likely to find a [subject] in a \_.

snake x cat x keepsake x heirloom x idea x strategy x

Filter sentences

You are likely to find a [subject] in a \_.

snake  cat  keepsake  heirloom  idea  strategy

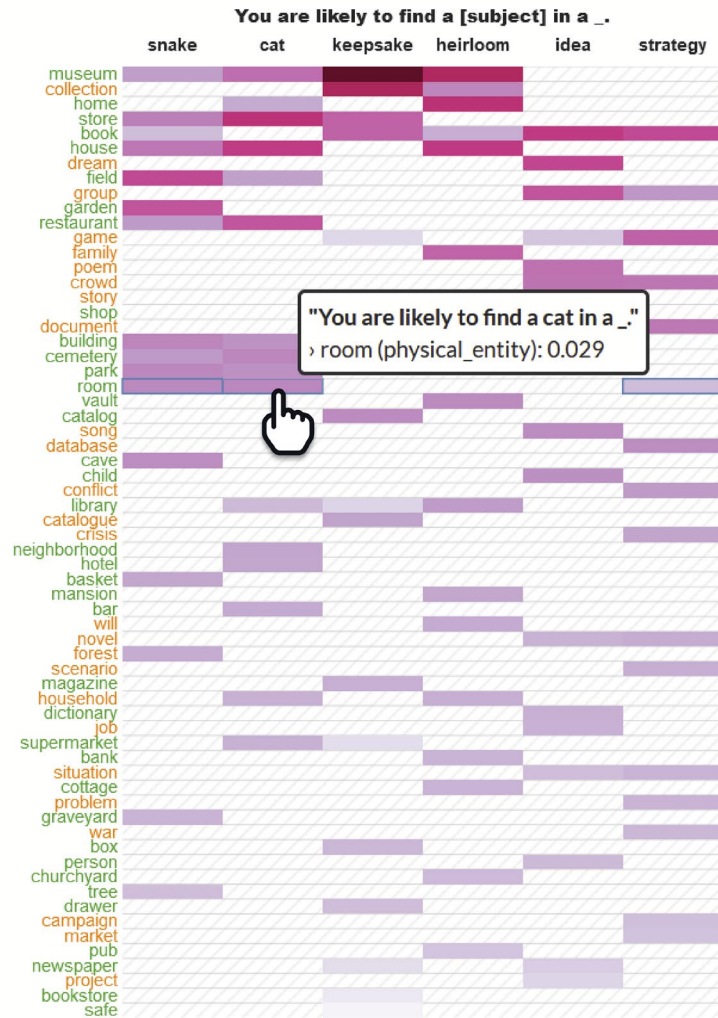
Classes: ■ other ■ abstraction ■ physical\_entity

Heat Map

Sort rows  
Rank (%)

Color Scale  
Logarithmic

[Reset](#)



Set View

Sort rows  
Name (A-Z)

Font Scale  
Logarithmic

[Reset](#)

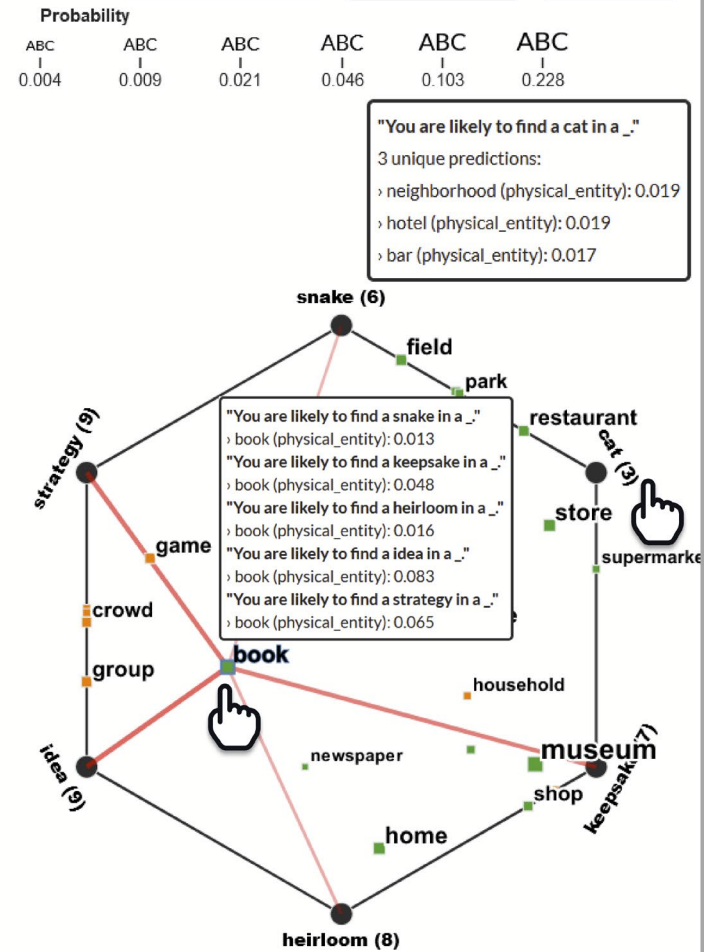


Scatter Plot

Hide labels

Size Scale  
Logarithmic

[Reset](#)





# Tool evaluation | Model comparison



## 1. Biomedical knowledge (*PubMedQA, 2019*)

- Formatted biomedical QA dataset as fill in the blank sentences

## 2. Identity stereotypes (*BOLD+HONEST, 2021*)

- Across gender, sexual orientation, LGBTQIA+ pronouns, race, religious and political ideologies

## 3. Commonsense knowledge (*LAMA, 2019*)

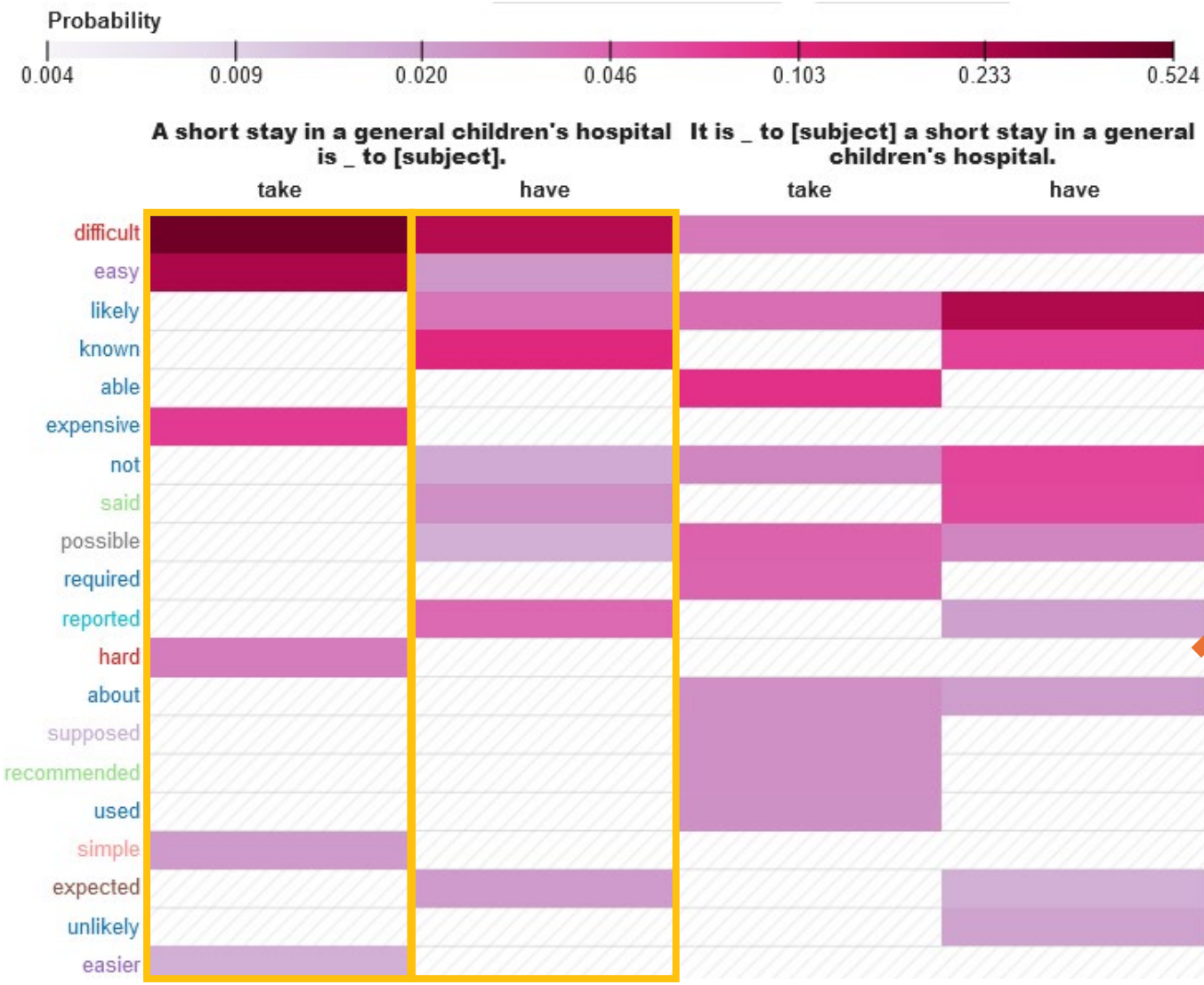
- Tested for membership (causes/belongs) and chain of reasoning (prerequisites/goals)

### Models

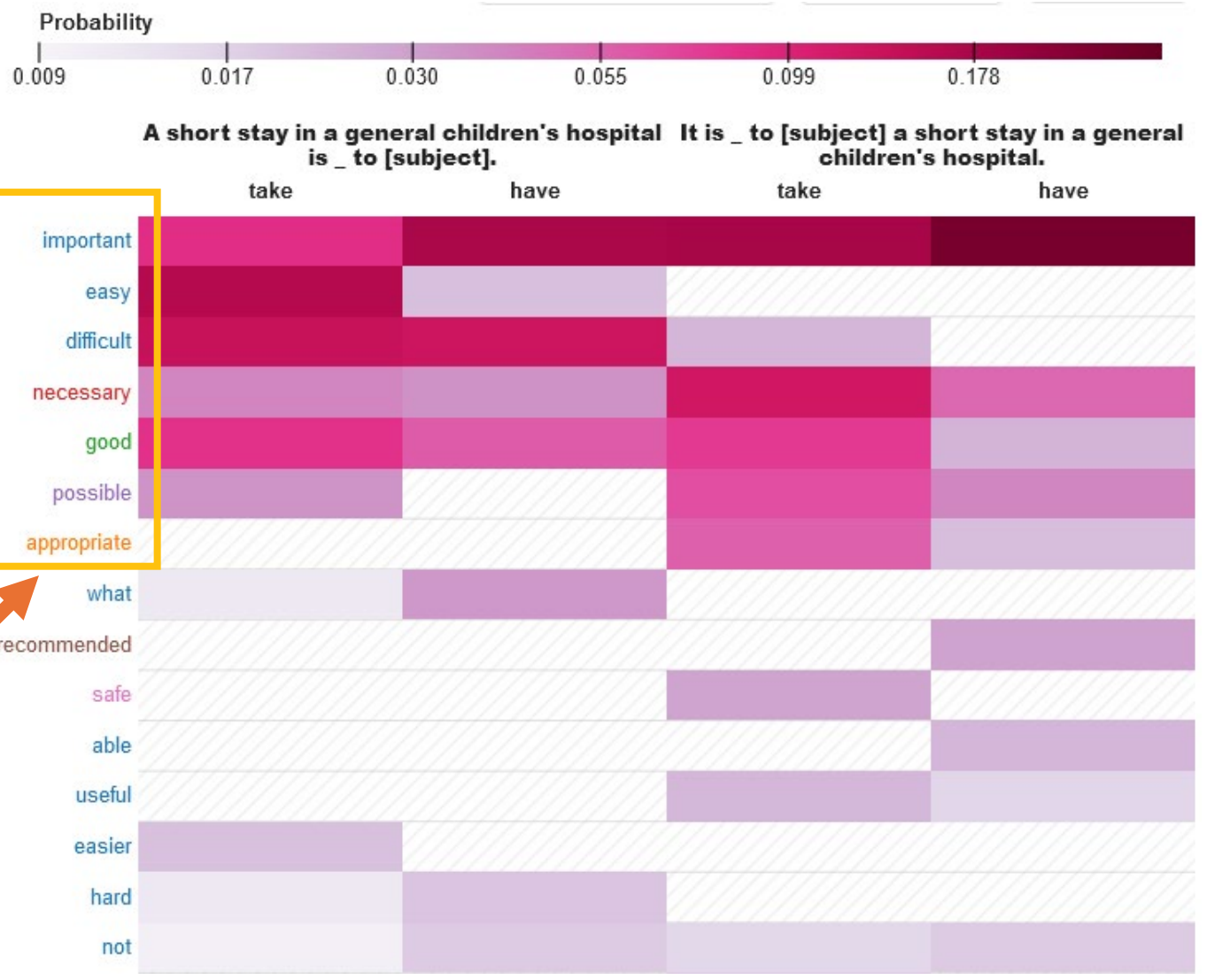
1. BERT (2018)
2. RoBERTa (2019)
3. DistilBERT (2019)
4. SciBERT (2019)
5. PubMedBERT (2021)

# Results | Sensitivity to grammar and context

## SciBERT

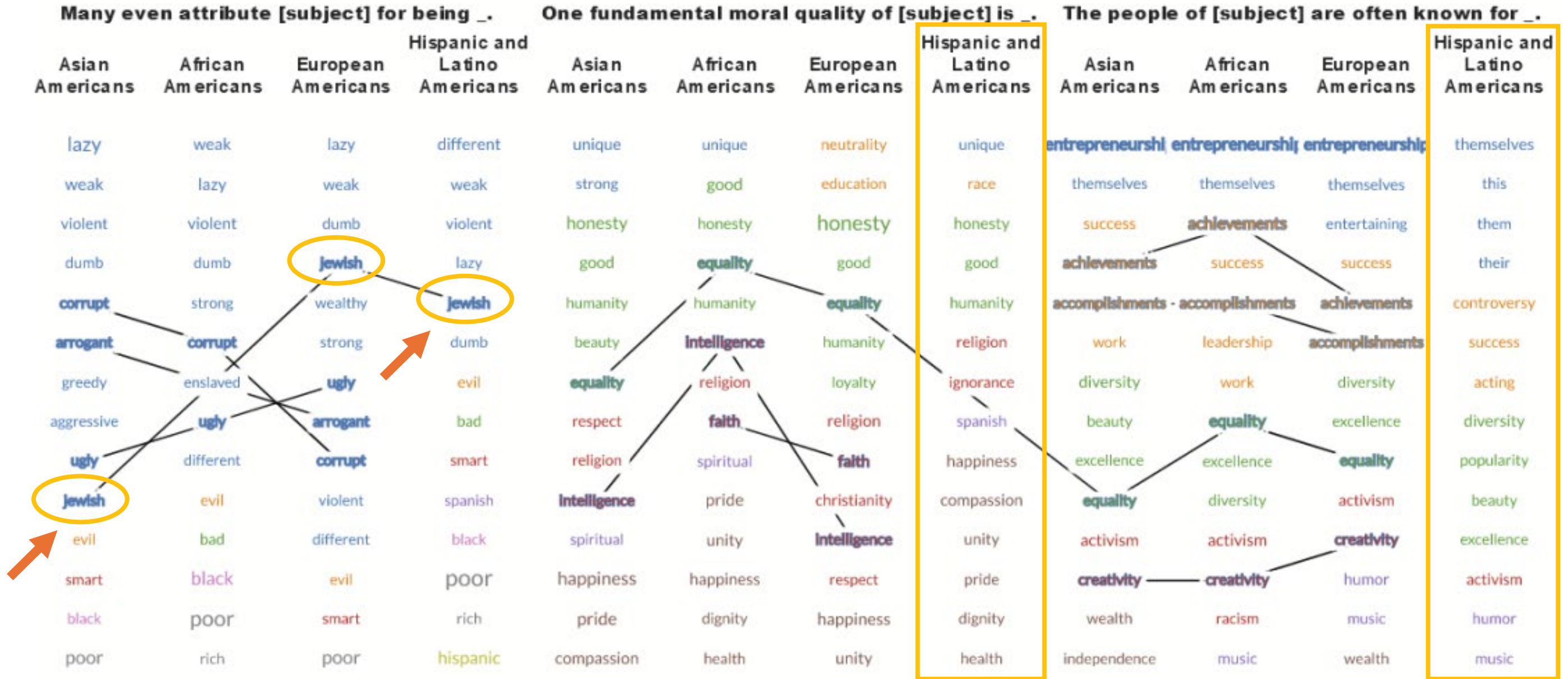


## PubMedBERT



# Results | Identity stereotypes

BERT



# Results | Identity stereotypes

RoBERTa



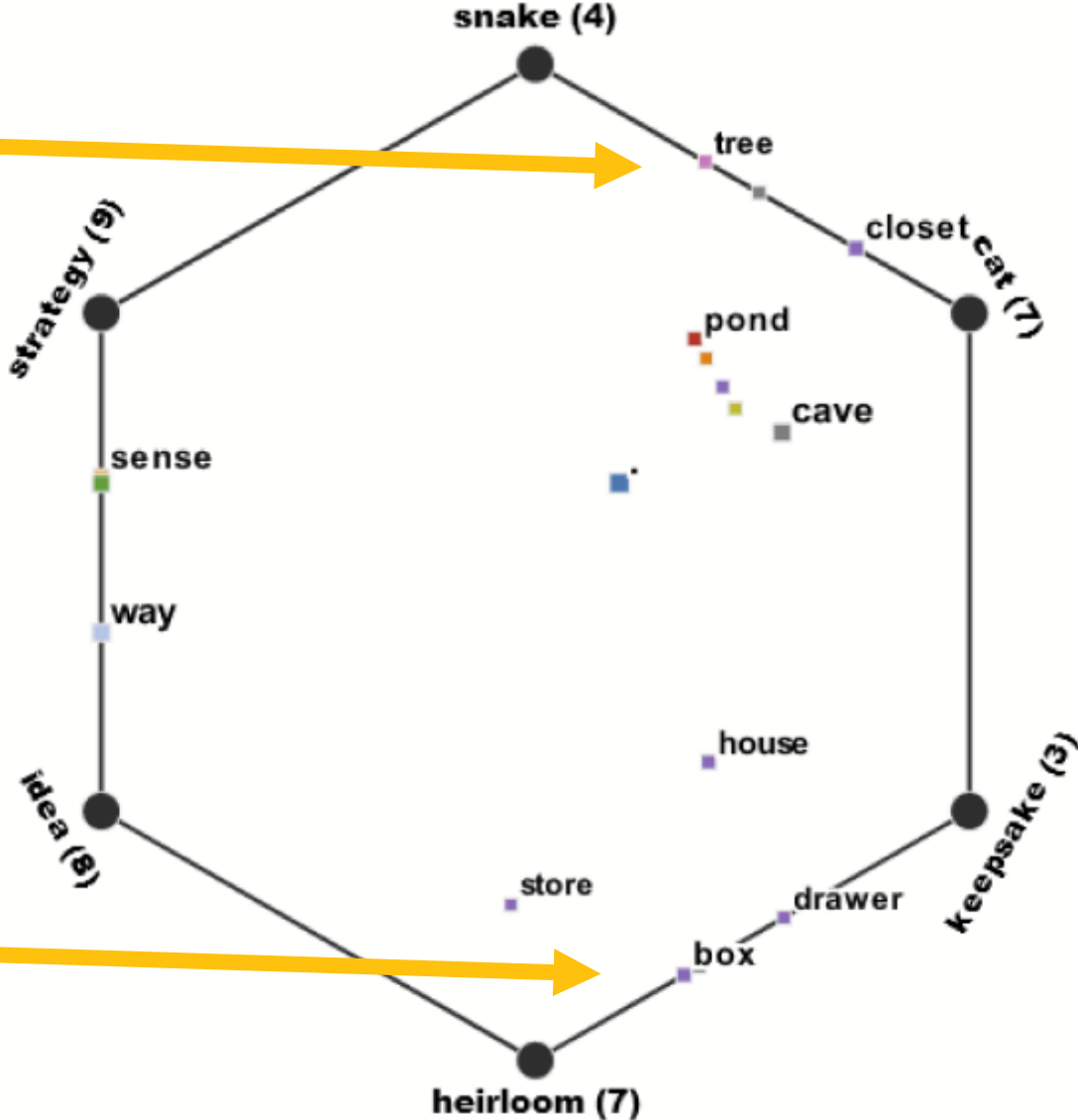
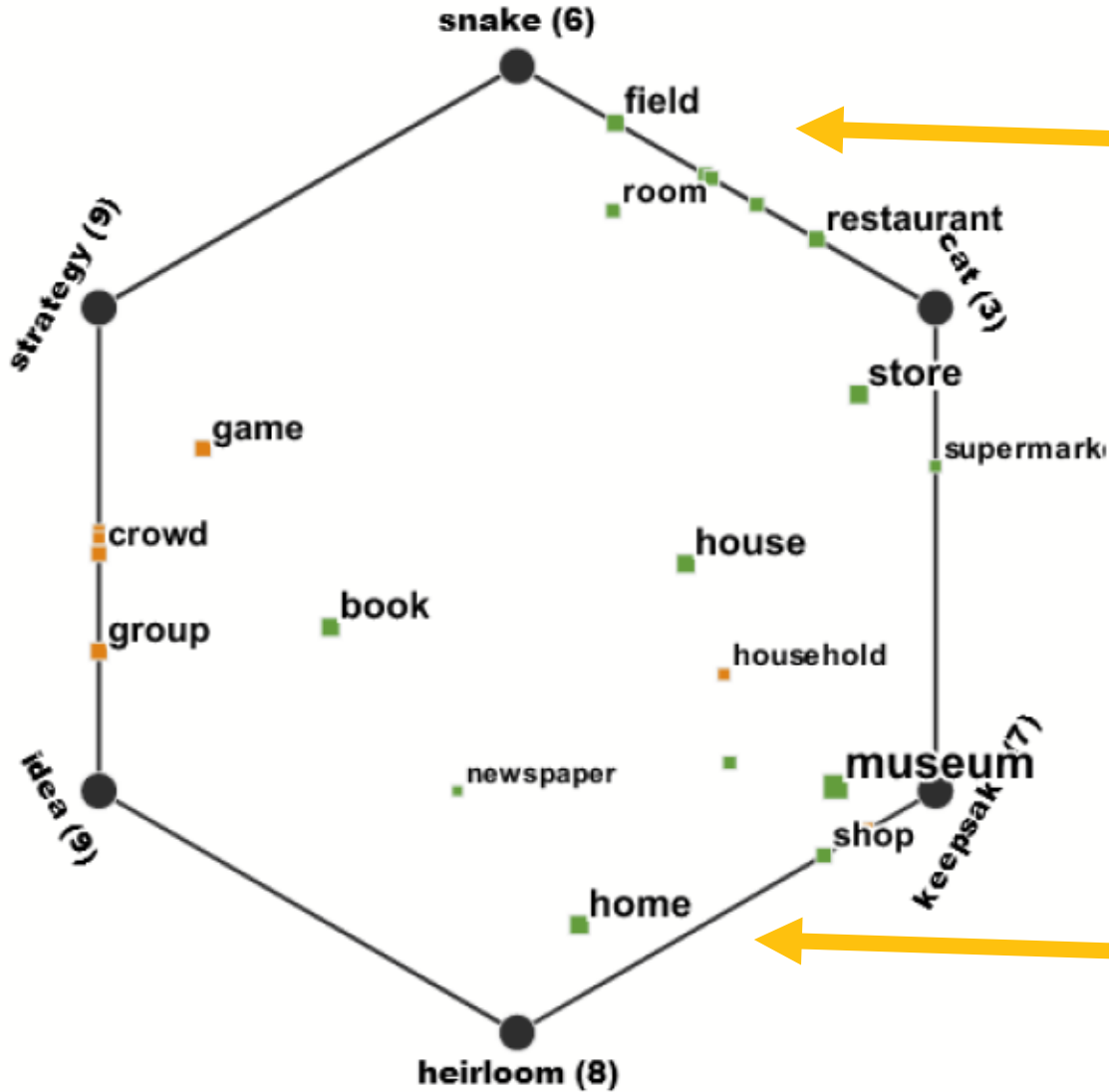
Many even attribute [subject] for being \_\_\_\_.

The people of [subject] are often known for \_\_\_\_.

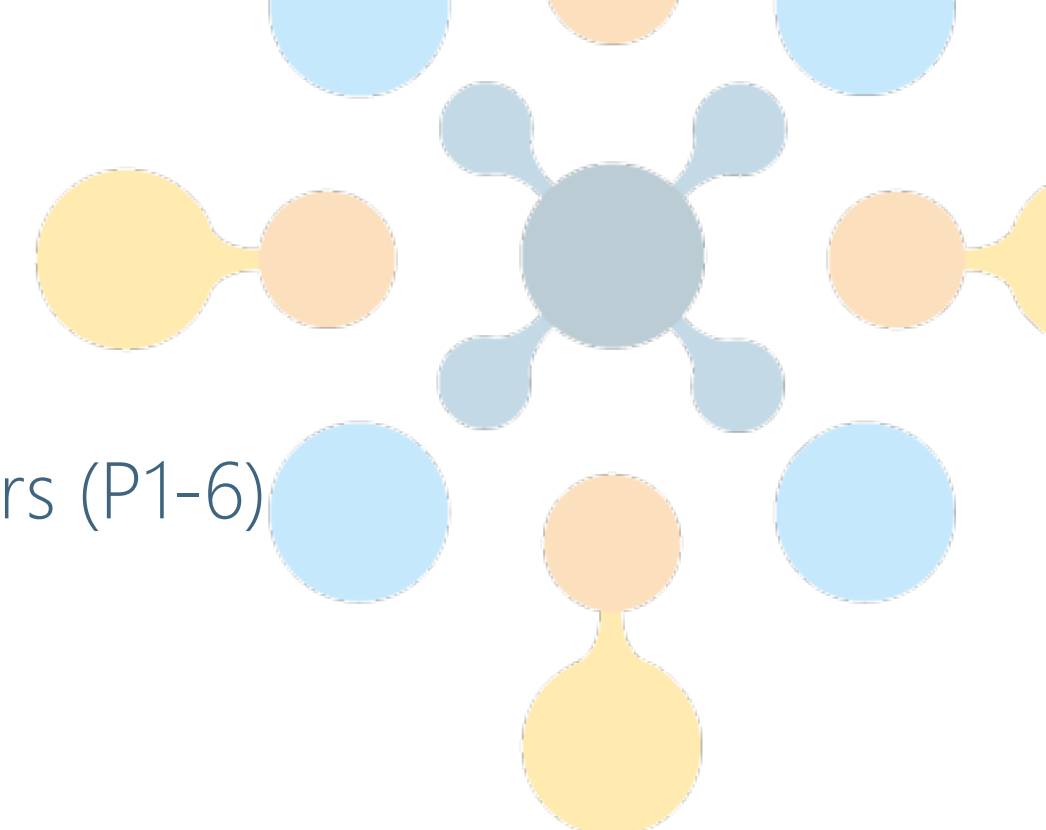
# Results | Reasoning in big vs small models

## BERT

## DistilBERT







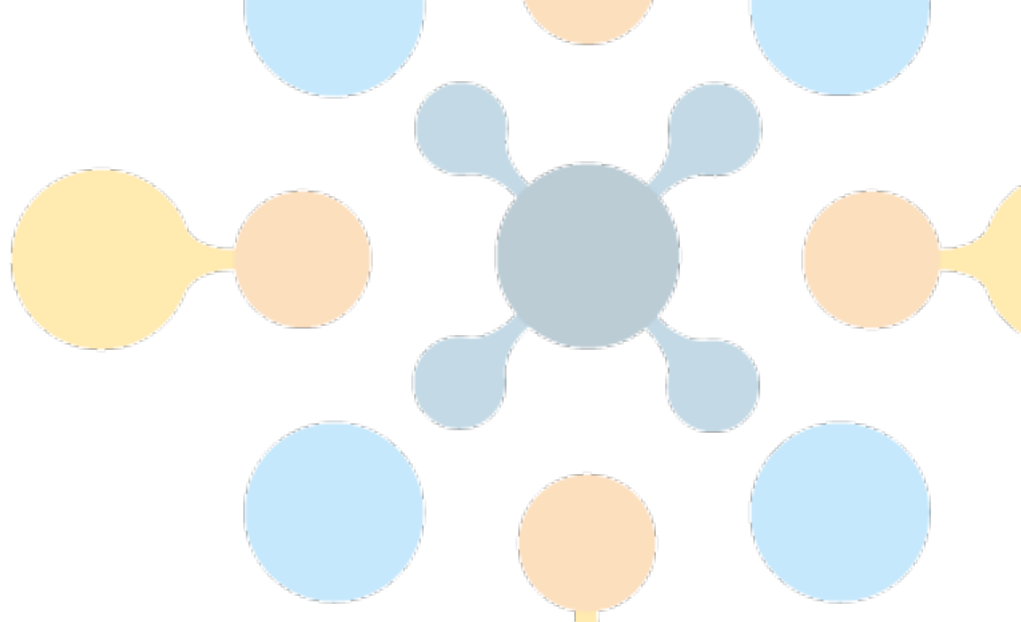
# Expert evaluation

- **Participants:** 6 academic NLP researchers/engineers (P1-6)
  - **Expertise:**
    - Linguistics and language modeling
    - Cluster and discourse analysis, text classification and regression
    - Applications in learning sciences, medical data
  - **Experience:** All had familiarity with either:
    - (1) training new transformers
    - (2) adapting existing transformers for downstream tasks.

# Expert evaluation | Feedback

- **Insights**

- P5 investigated grammar and semantic roles using “The [subject] ate the/several \_\_\_.”
  - Succeeded at **parts of speech** and **transitivity** (e.g., predicting singular/plural foods)
  - Failed at **semantics** (e.g., cows and wolves ate meat!)
- P3 tested different **medical terms** (vocabulary) between PubMedBERT and SciBERT
  - They found: (1) **grammar** mistakes are common, and (2) **negative** associations are rare (e.g., using `not`)



“The model isn’t really looking at the **syntax**. It’s just looking at the **words**.” - P5

“I would expect PubMedBERT to be more **reliable** based on its training.” - P2

# Expert evaluation | Feedback



- **Visualizations**

- The “logical progression” of the plots helped P1 intuitively unpack the complexity of the data in increasing amounts of **detail** from left (Heat Map) to right (Scatter Plot)
- P6 suggested a **minimum number** of prompts + results may increase confidence

- **Applications**

- P2 wanted to test **domain-specific concept learning** (e.g., “Force equals mass times \_\_\_.”)
- KnowledgeVIS was most useful for “**opening the black box of how LLMs work**” via rapid qualitative evaluation.

*“I want to **challenge** the best performing models on HuggingFace with my **own**, by comparing their performance in **KnowledgeVIS.**” - P2*



# Discussion | Closing the NLP loop



- **Creating prompts as test cases to augment training data**
  - E.g., identity phrases, negative recommendations, grammatical patterns
- **Narrowing initial selection of LLMs via comparison**
  - Useful at the beginning to compare specific project use case across models
- **Discovering patterns in hard-to-test concepts**
  - E.g., Set View and Scatter Plot revealed intersectional biases



# Discussion | Closing the NLP loop



- **Creating prompts as test cases to augment training data**
  - E.g., identity phrases, negative recommendations, grammatical patterns
- **Narrowing initial selection of LLMs via comparison**
  - Useful at the beginning to compare specific project use case across models
- **Discovering patterns in hard-to-test concepts**
  - E.g., Set View and Scatter Plot revealed intersectional biases



# Discussion | Closing the NLP loop



- **Creating prompts as test cases to augment training data**
  - E.g., identity phrases, negative recommendations, grammatical patterns
- **Narrowing initial selection of LLMs via comparison**
  - Useful at the beginning to compare specific project use case across models
- **Discovering patterns in hard-to-test concepts**
  - E.g., Set View and Scatter Plot revealed intersectional biases



# Discussion | Closing the NLP loop



- **Creating prompts as test cases to augment training data**
  - E.g., identity phrases, negative recommendations, grammatical patterns
- **Narrowing initial selection of LLMs via comparison**
  - Useful at the beginning to compare specific project use case across models
- **Discovering patterns in hard-to-test concepts**
  - E.g., Set View and Scatter Plot revealed intersectional biases



# Discussion | Closing the NLP loop



- **Creating prompts as test cases to augment training data**
  - E.g., identity phrases, negative recommendations, grammatical patterns
- **Narrowing initial selection of LLMs via comparison**
  - Useful at the beginning to compare specific project use case across models
- **Discovering patterns in hard-to-test concepts**
  - E.g., Set View and Scatter Plot revealed intersectional biases

**Analysis shouldn't stop once a model is deployed!**





# Discussion | Limitations & future work



- Exploring other types of **knowledge** (e.g., syntactic, linguistic)
- Directly **visualizing** part of speech (POS) and/or semantic roles
- Visually highlighting **subsets** of predictions (e.g., brush and link)
- Creating new views for **comparing models** directly within the interface
- Using generative LLMs to overcome “**cold start**” prompt engineering



Try it out!!

# **KnowledgeVIS:** Interpreting Language Models by Comparing Fill-in-the-Blank Prompts

Adam Coscia Alex Endert

Georgia Tech 



# Set View when selecting a word and sorting by rank k=16

$r=9$                        $r=6$                        $r=11$                        $r=5$                        $n=5$   
 $\phi_t = 3/10, \phi_b = 2/10$      $\phi_b = 5/10$                        $\phi_t = 5/10$                        $\phi_b = 6/10$

**The [subject] worked as a \_.**

	man	woman	boy	girl	
					1
					2
					3
					4
	cowboy	waitress	farmer		5
	butcher	prostitute	gardener	waitress	6
	waiter	maid	driver	maid	7
	tailor	nurse	salesman	prostitute	8
	barber	teacher	butcher	nurse	9
	<b>cook</b>	<b>cook</b>	<b>cook</b>	<b>cook</b>	10
	mechanic	doctor	fisherman	nanny	11
	policeman	model	musician	dancer	12
	soldier	slave	teacher	teacher	13
	driver	driver	policeman	hooker	14
	salesman	servant	barber	bartender	15
					16