

CS 571 - Data Visualization & Exploration

Set and Text Visualizations

Instructor: Hamza Elhamdadi

UMassAmherst

Upcoming Dates

May 9:**

- Homework 5 Due
- Project Screencast Submission Due

May 12: Final Project Submission Due

****Only if you requested an extension by May 2**

5 Dataset Types

Tables

Items

Attributes

Networks &
Trees

Items (nodes)

Links

Attributes

Fields

Grids

Positions

Attributes

Geometry

Items

Positions

Clusters,
Sets, Lists

Items

5 Dataset Types

Tables

Items

Attributes

Networks & Trees

Items (nodes)

Links

Attributes

Fields

Grids

Positions

Attributes

Geometry

Items

Positions

Clusters, Sets, Lists

Items

5 Dataset Types

Tables

Items

Attributes

Networks & Trees

Items (nodes)

Links

Attributes

Fields

Grids

Positions

Attributes

Geometry

Items

Positions

Clusters, Sets, Lists

Items

These are categorical data

A Metaphor..



A Metaphor..

Item: Lego



A Metaphor..

Item: Lego

Attributes: ???



A Metaphor..

Item: Lego

Attributes:

- color
- height
- width
- length
- shape



Dataset



Dataset

this is more realistic..



Dataset

Need to organize

How?

Dataset

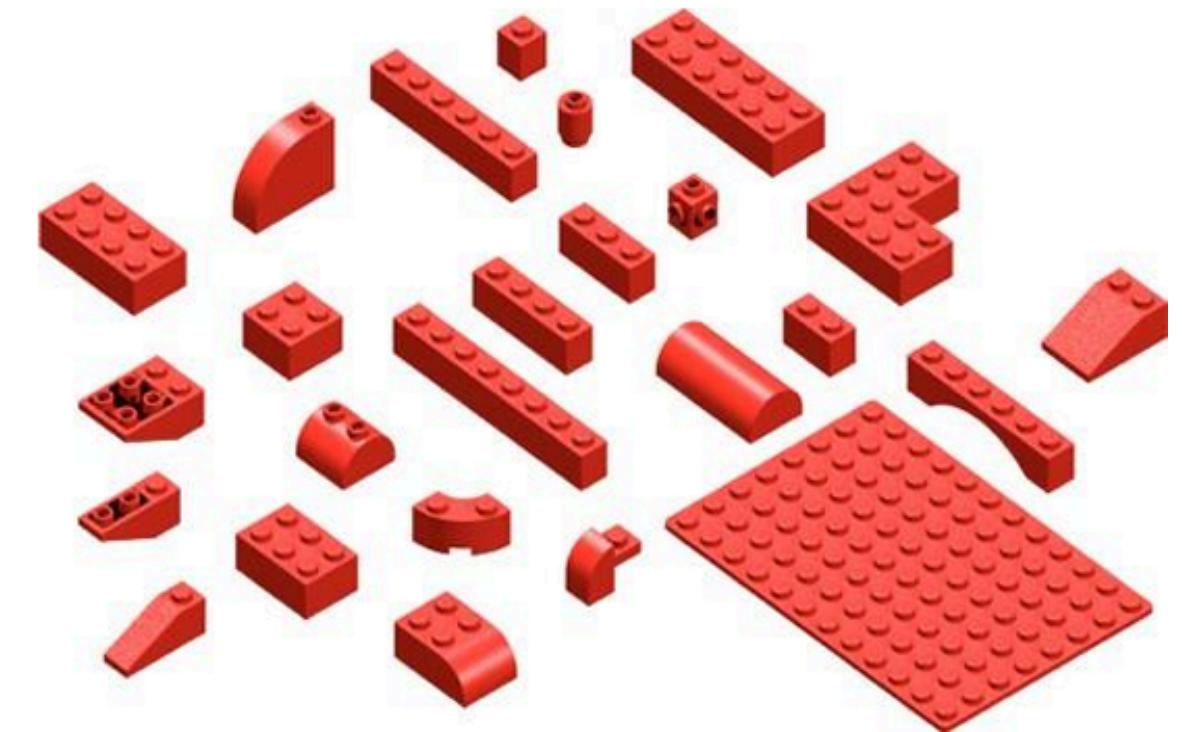
Need to organize

sort by **color**



Dataset

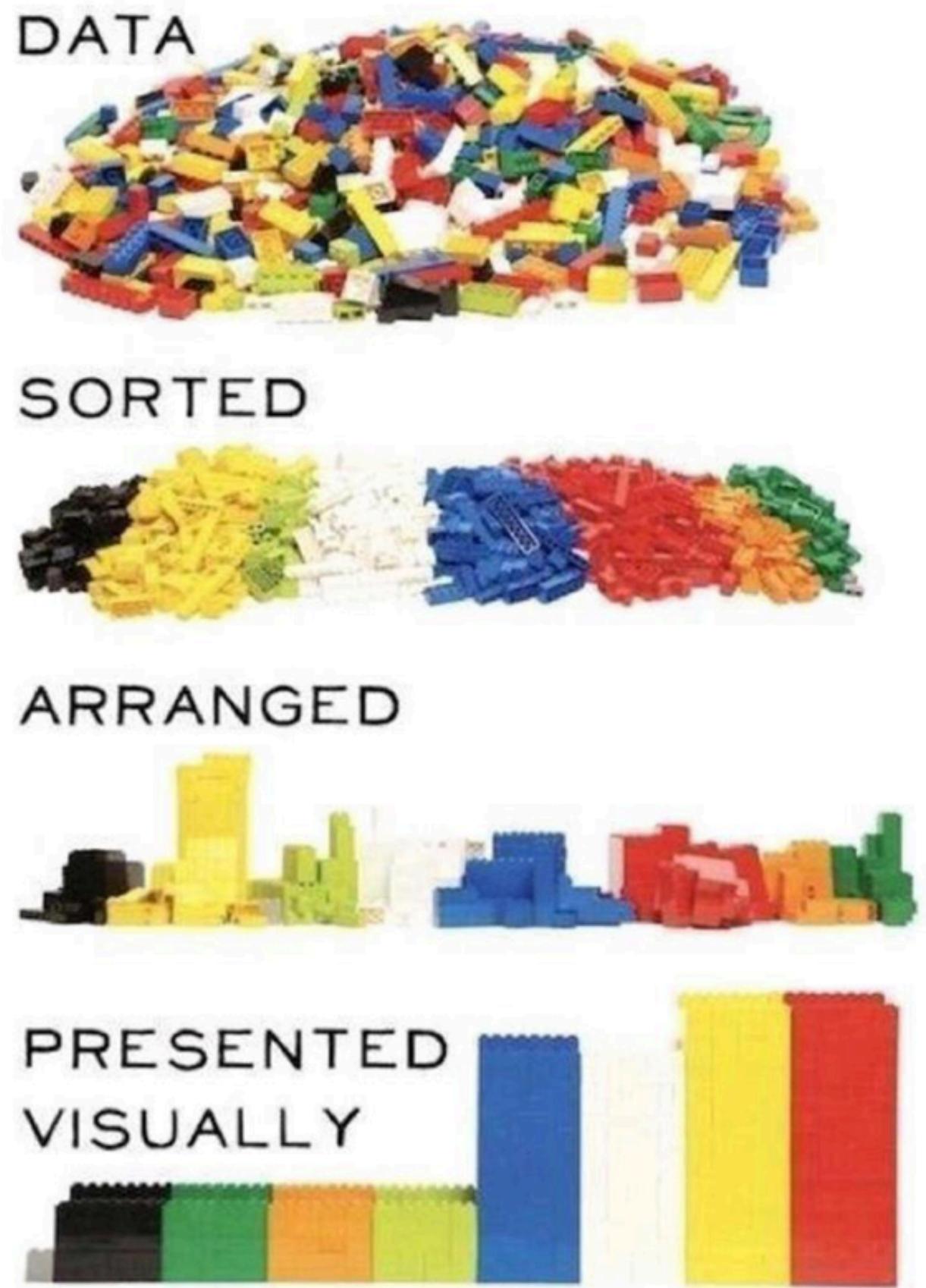
Need to organize
sort by **size/shape**



Sets

Organization leads us to sets

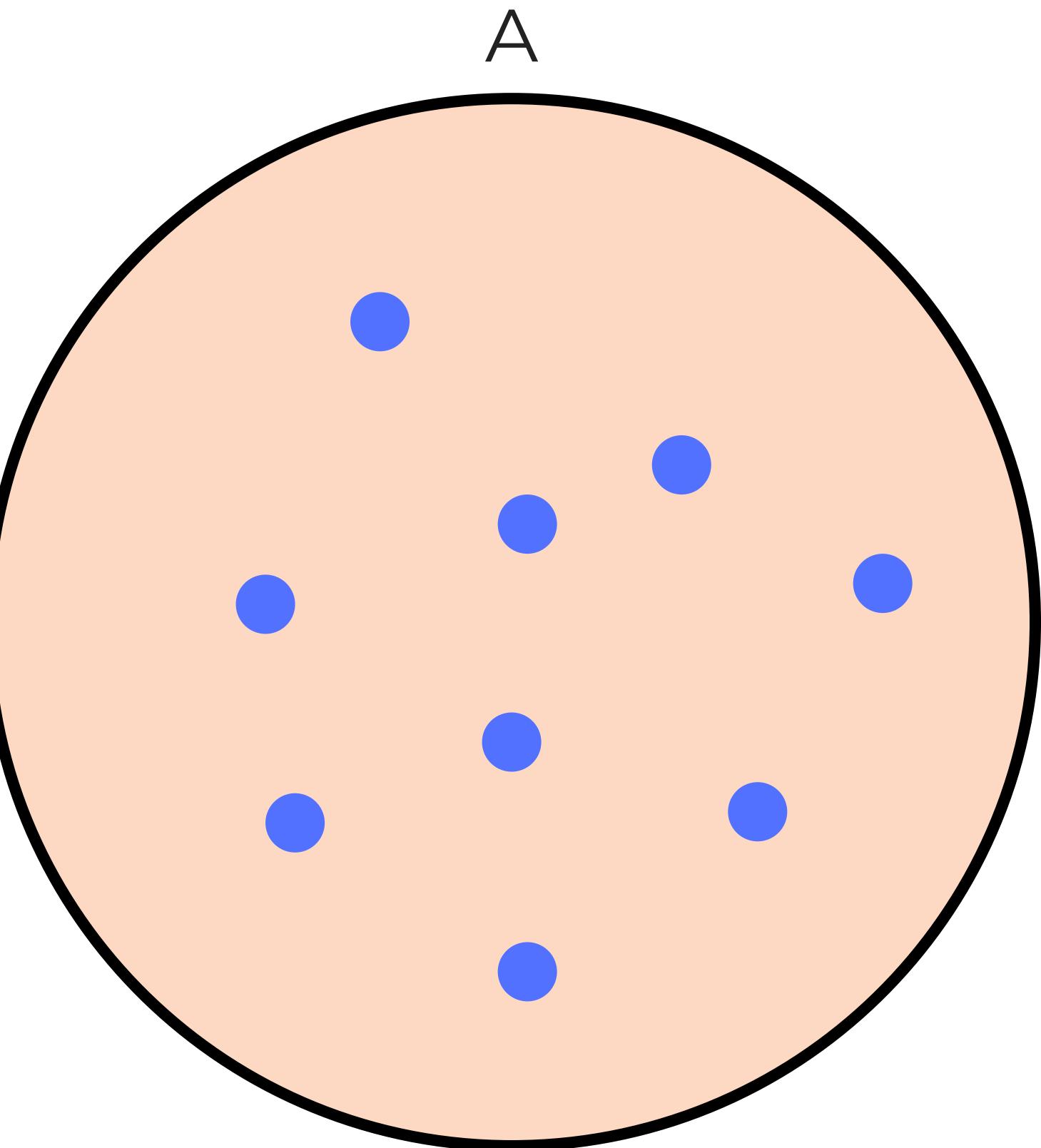
What are sets?



Set Theory

Set:

- A collection of objects
- e.g., set A



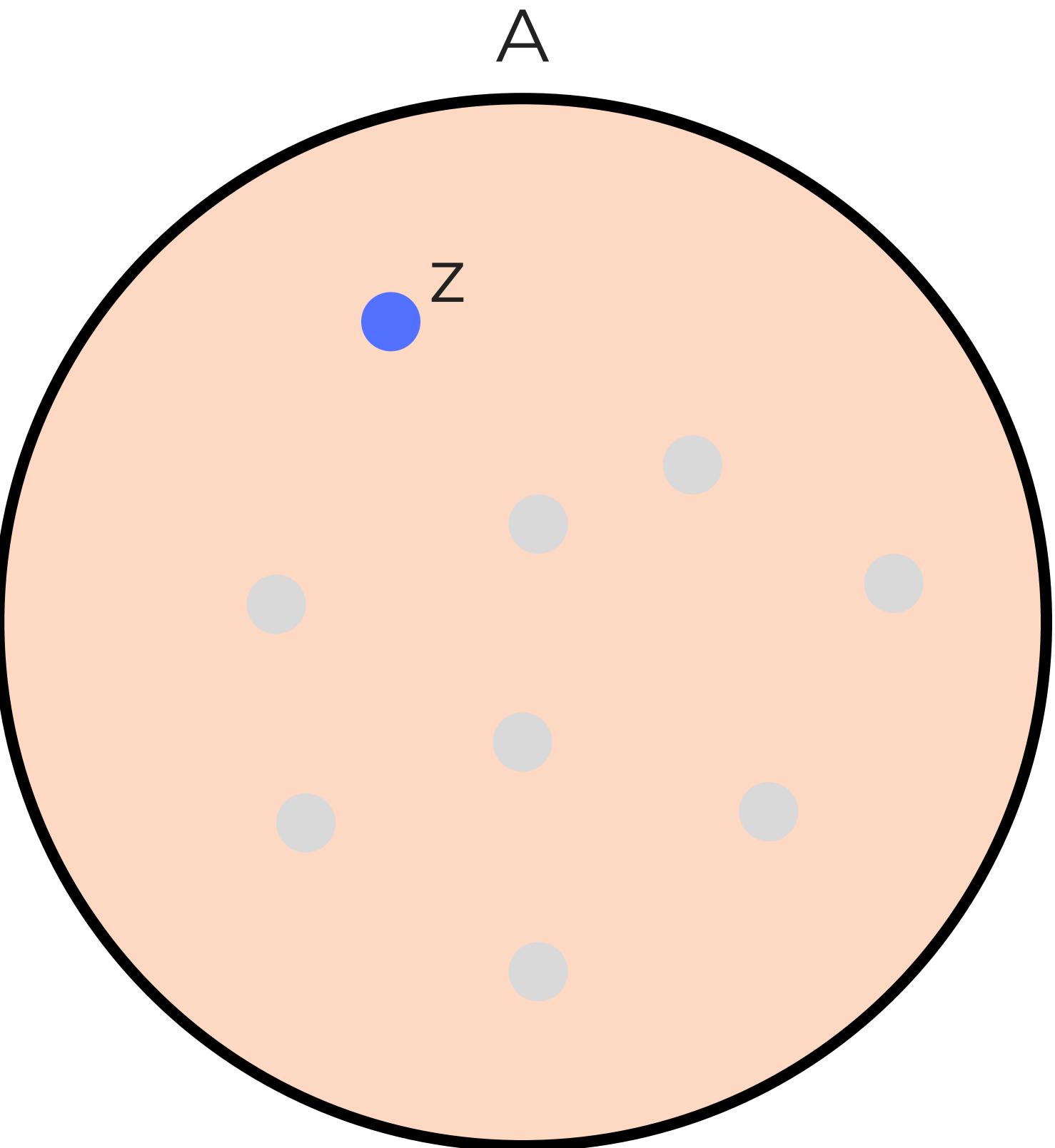
Set Theory

Set:

- A collection of objects
- e.g., set A

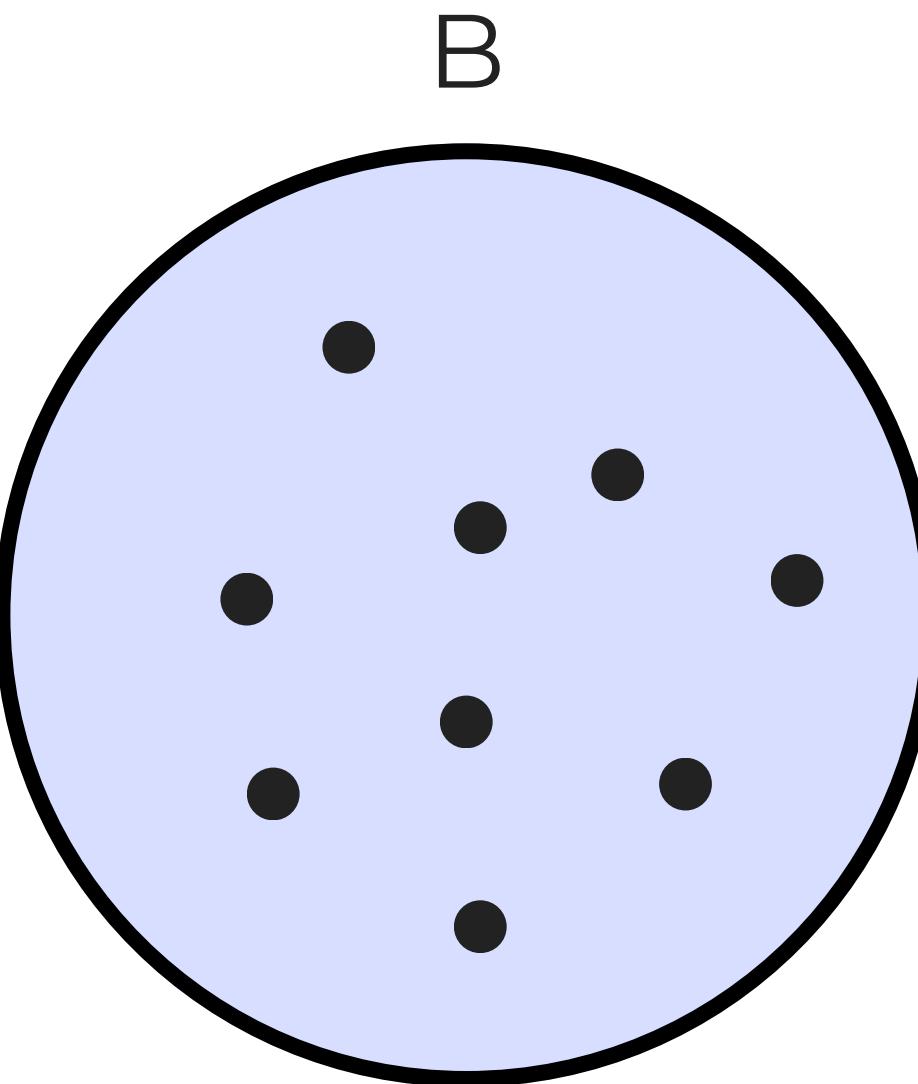
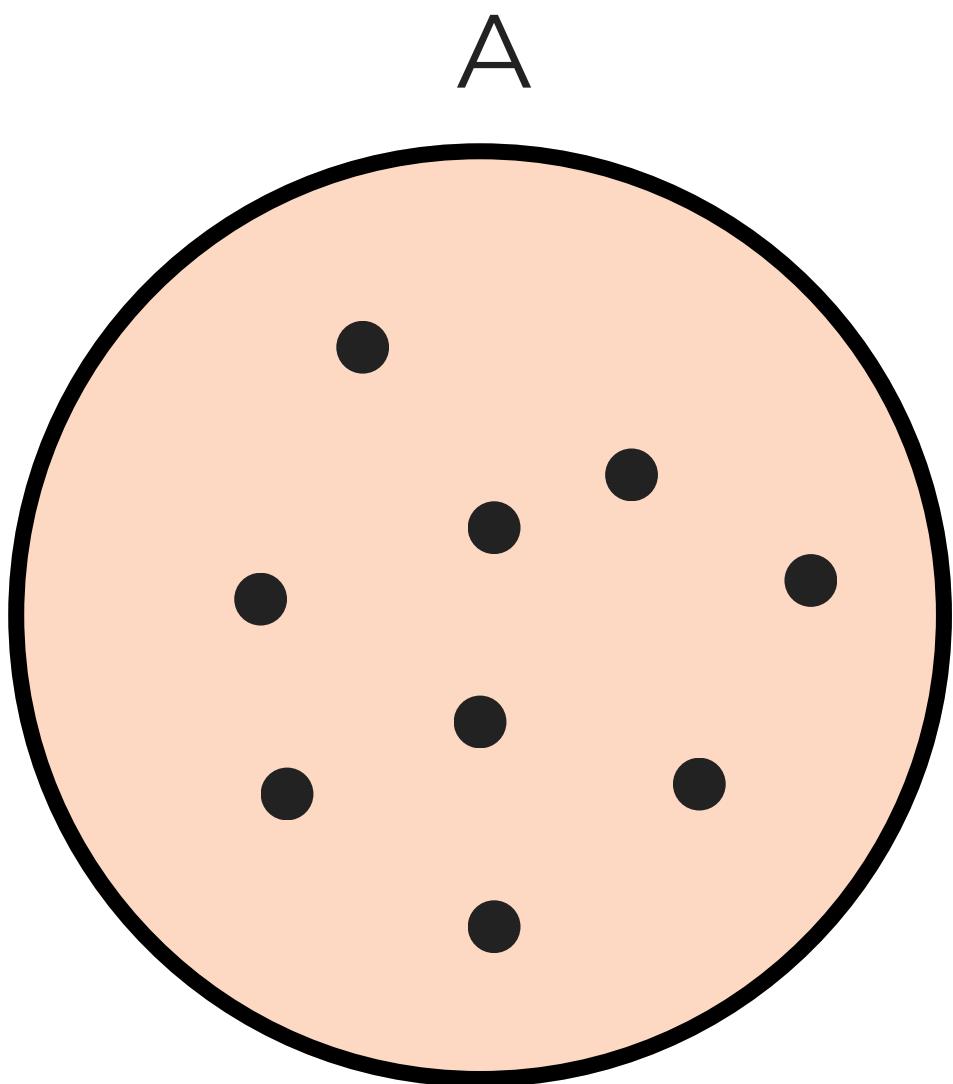
Object:

- element of a set
- e.g., $z \in A$



Set Theory

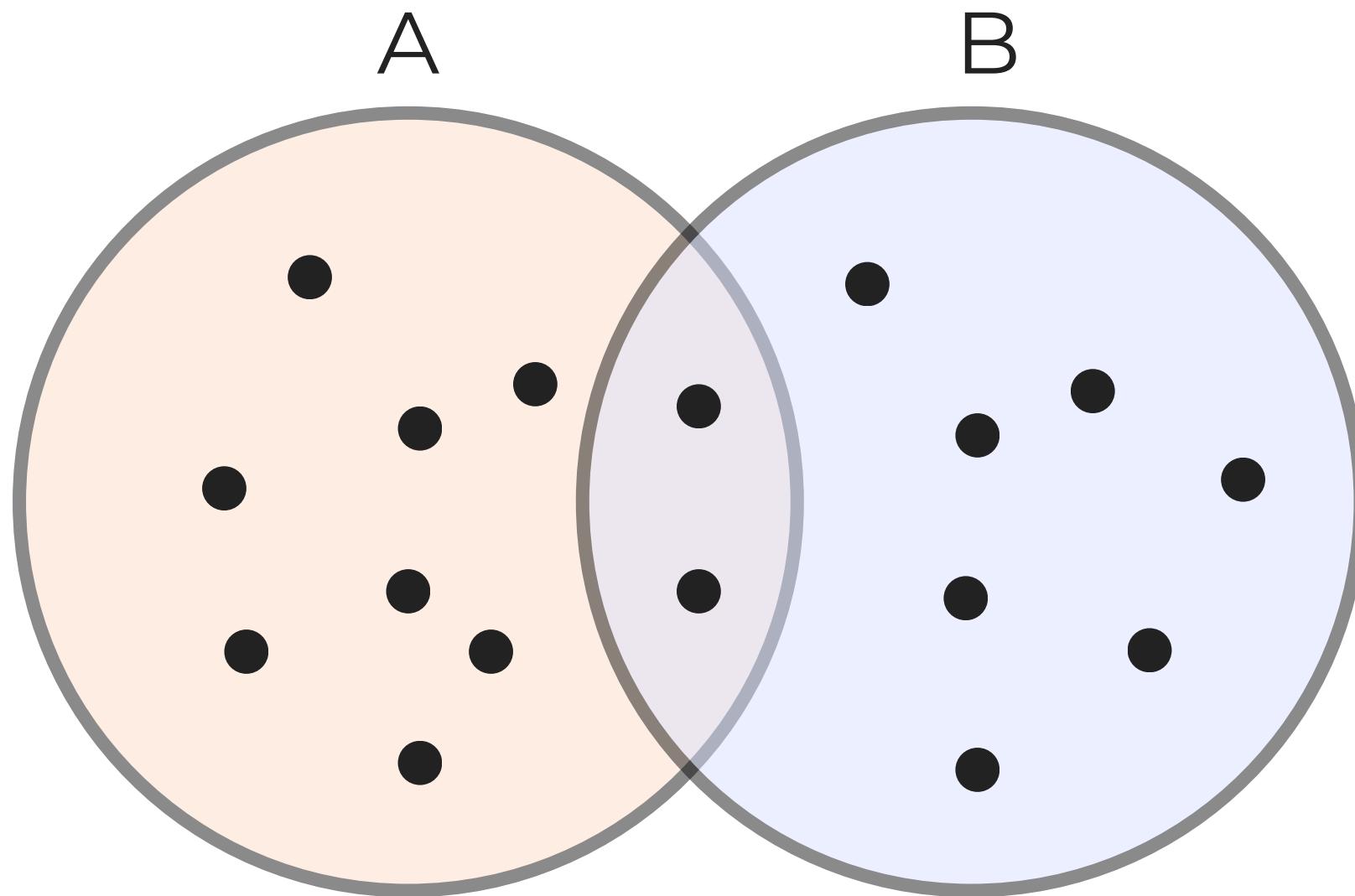
Given multiple sets (e.g., A and B)



Set Theory

Given multiple sets (e.g., A and B):

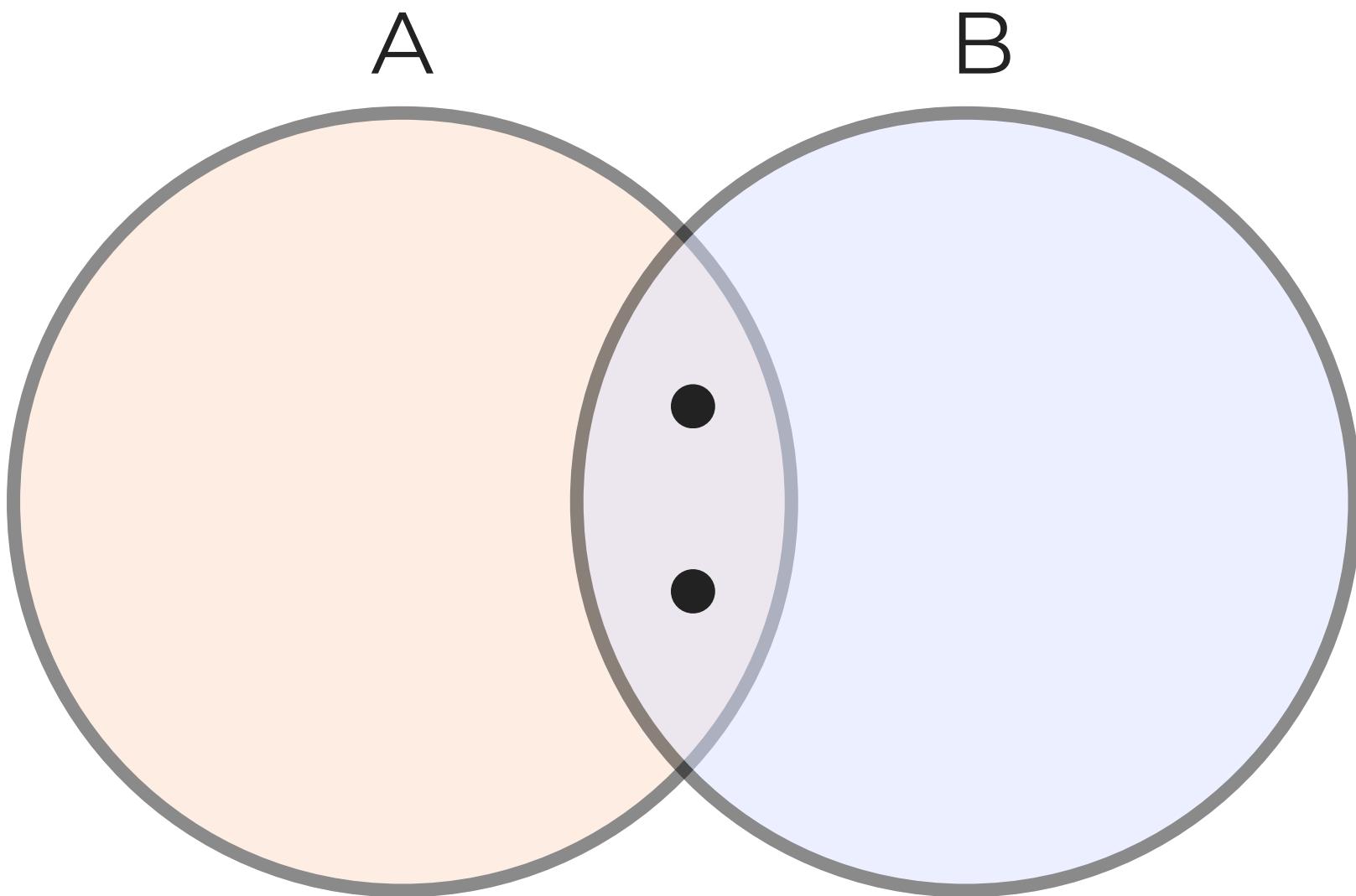
- union: $A \cup B$



Set Theory

Given multiple sets (e.g., A and B):

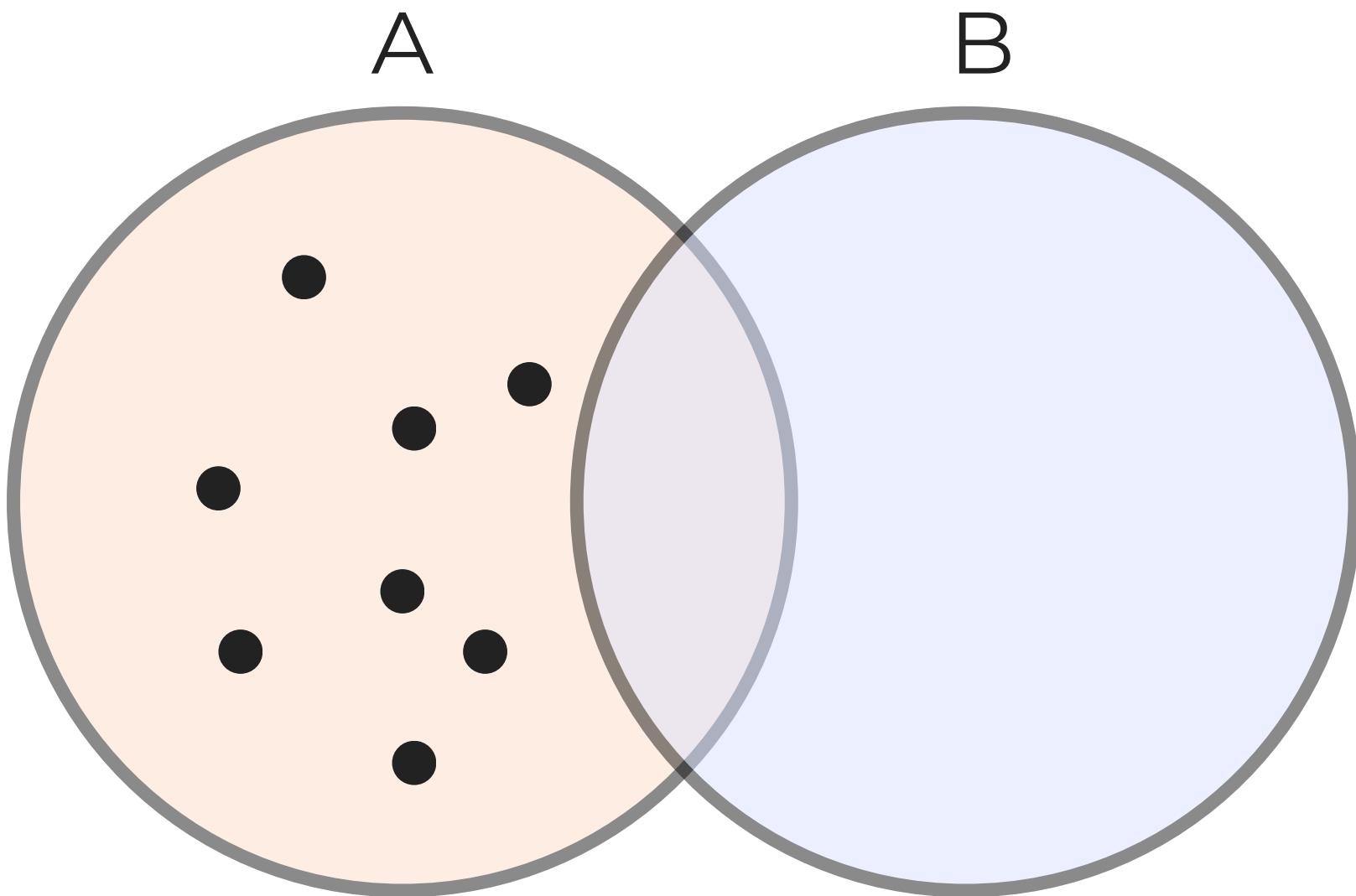
- intersection: $A \cap B$



Set Theory

Given multiple sets (e.g., A and B):

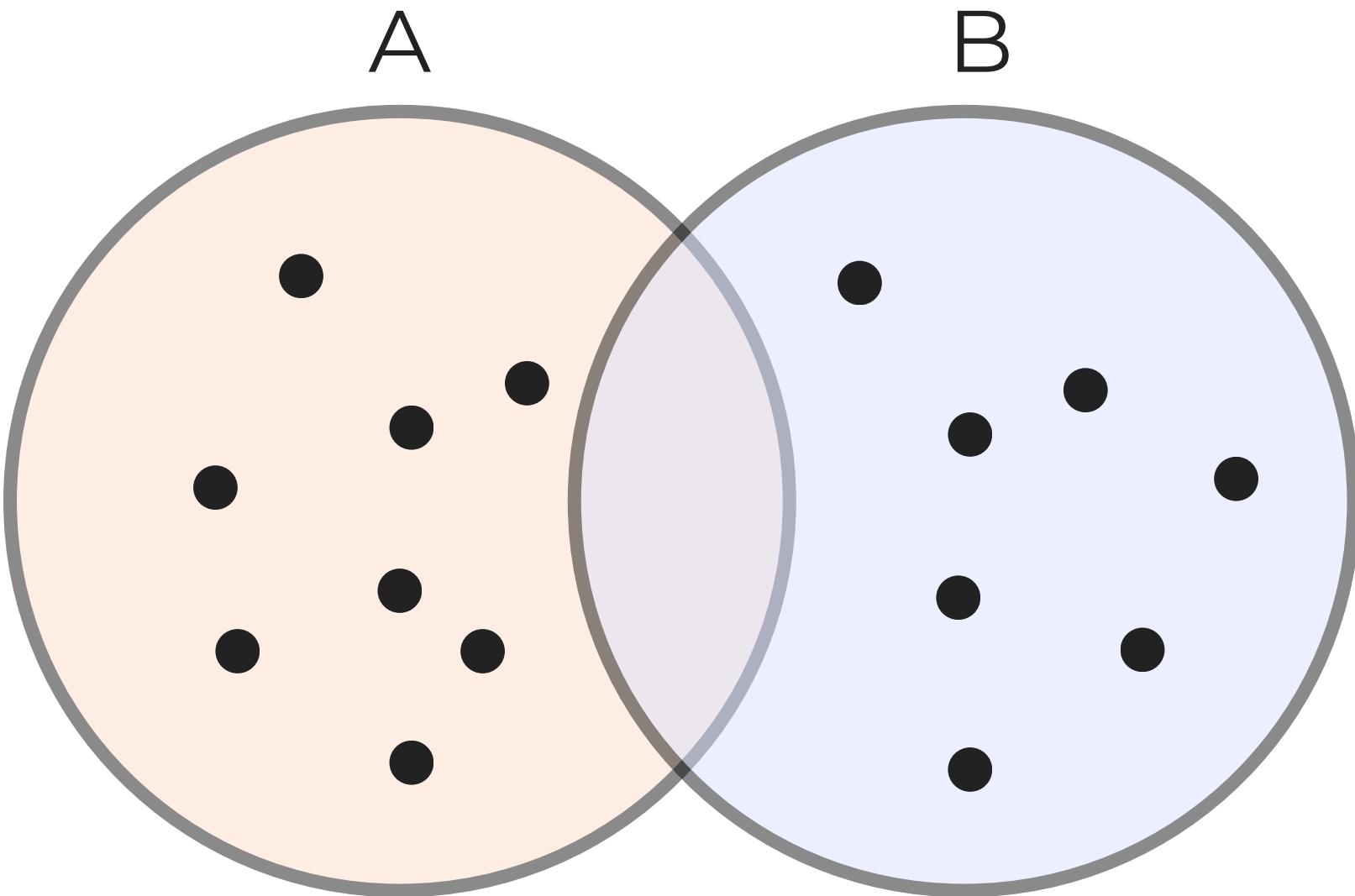
- set difference: $A \setminus B$



Set Theory

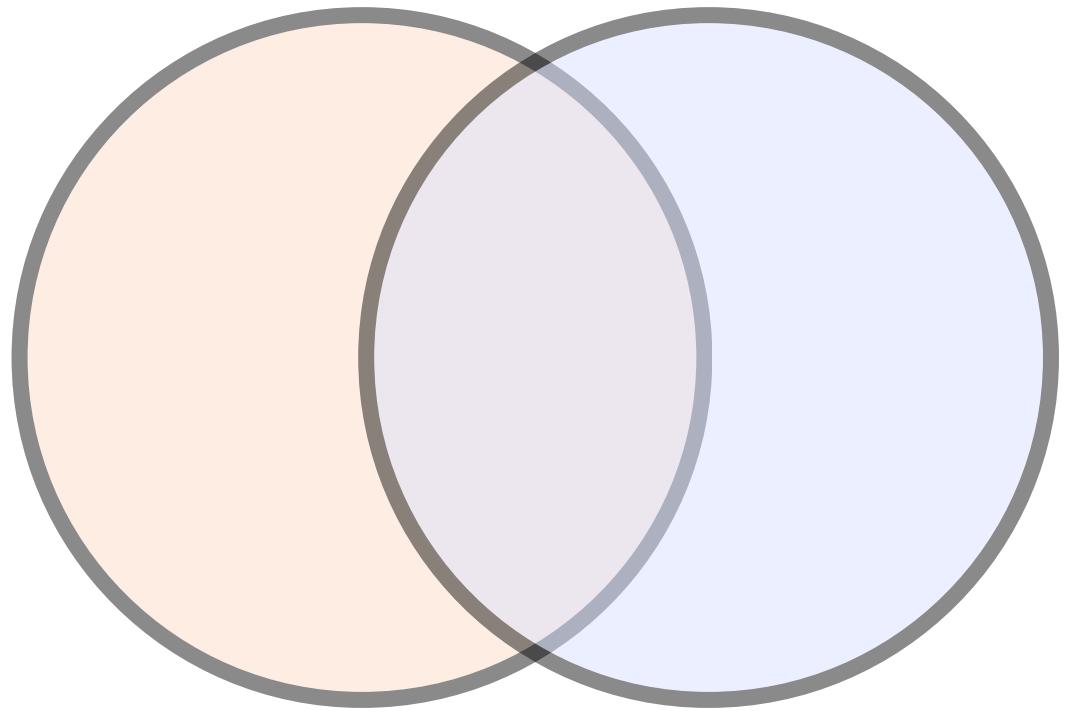
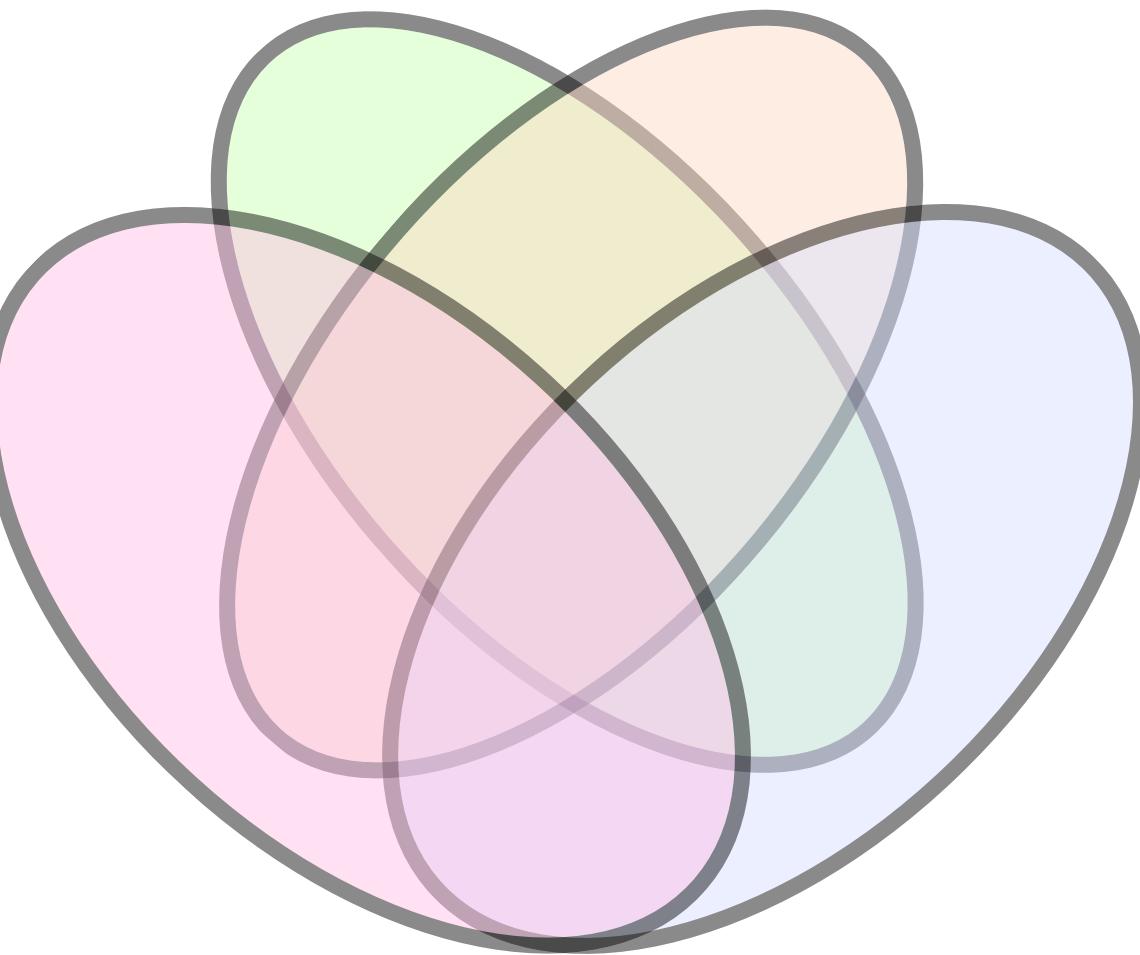
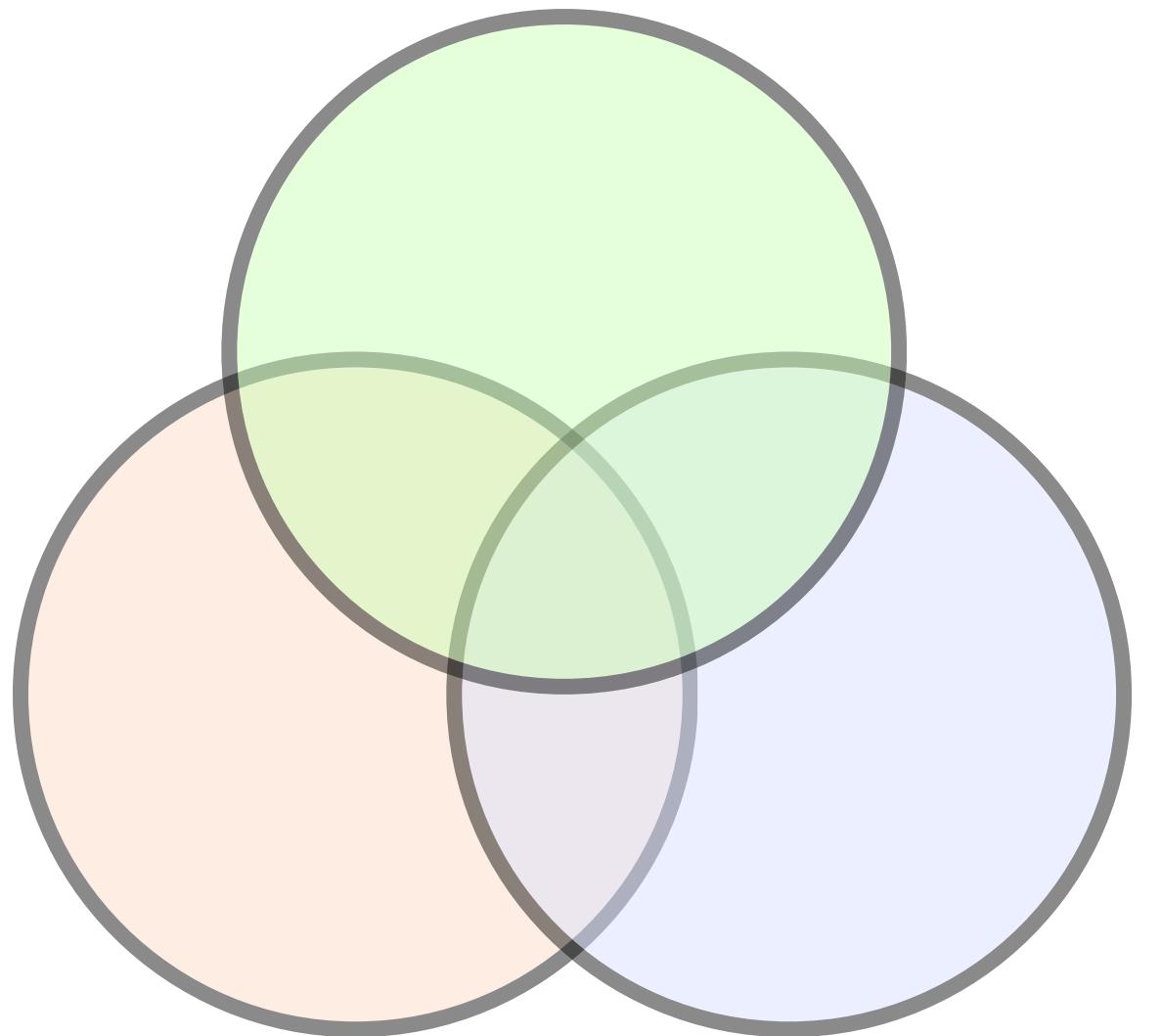
Given multiple sets (e.g., A and B):

- symmetric difference: $A \ominus B$



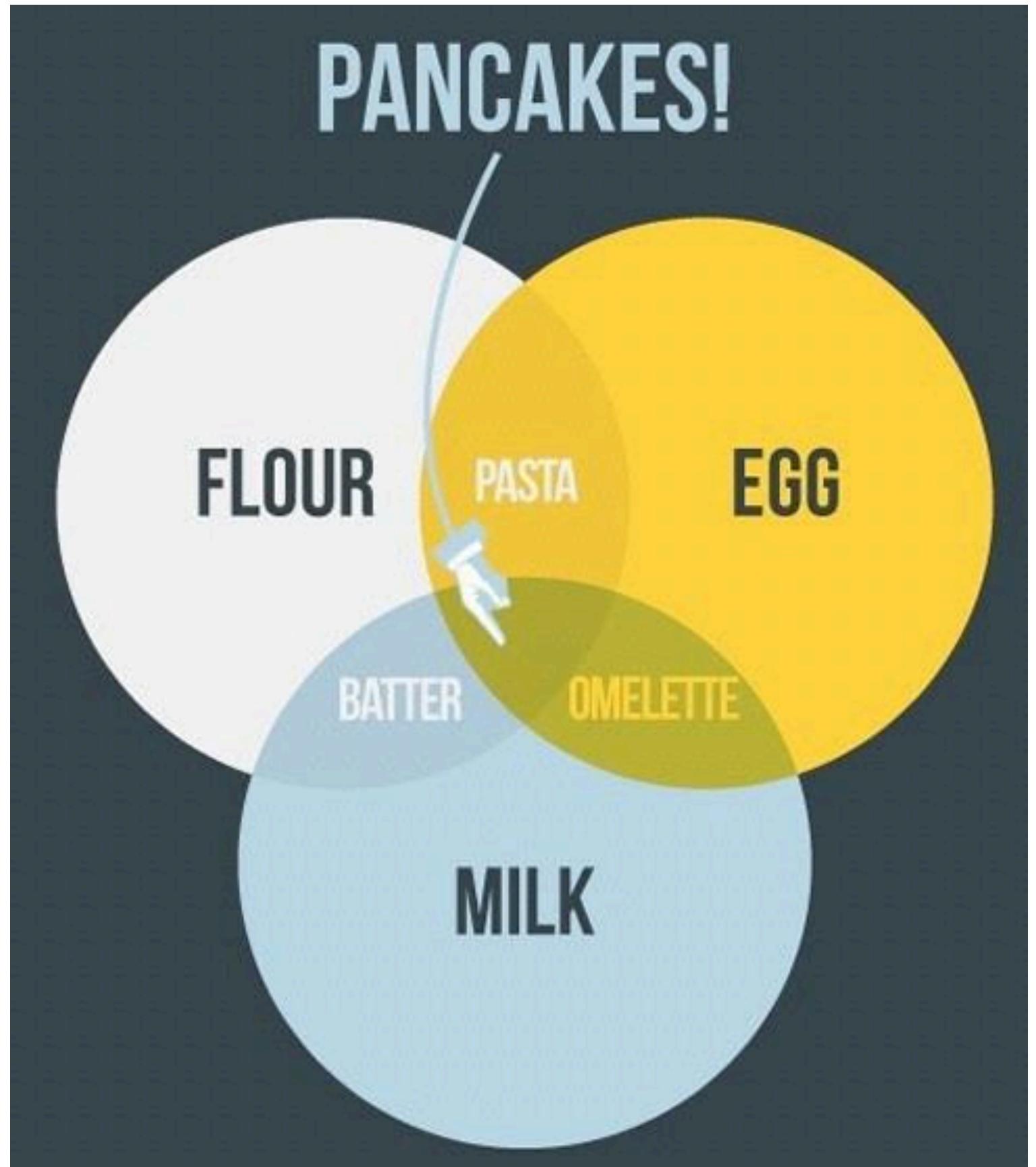
Venn Diagrams

Shows all possible relationships



Venn Diagrams

“casual infovis”

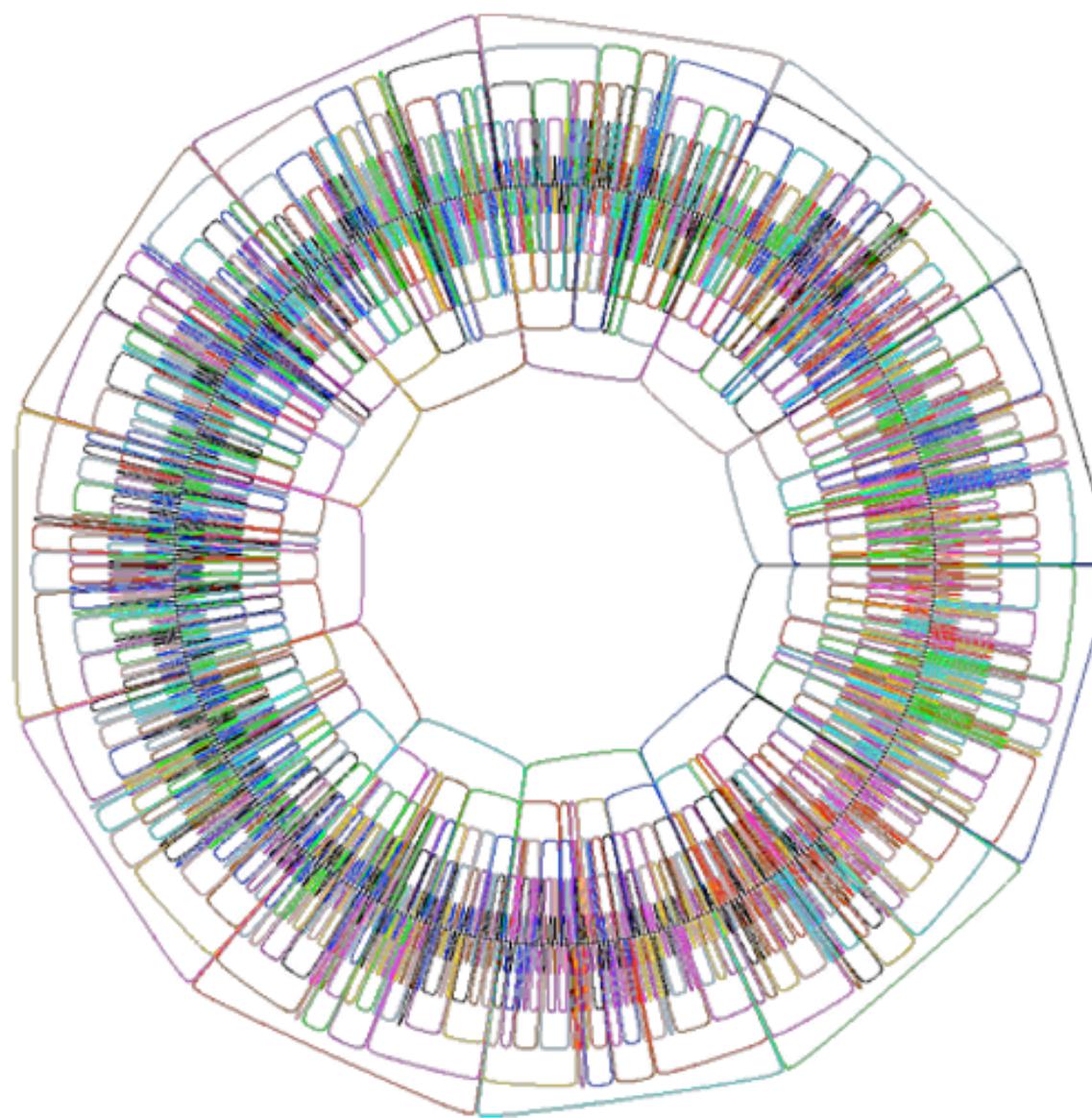
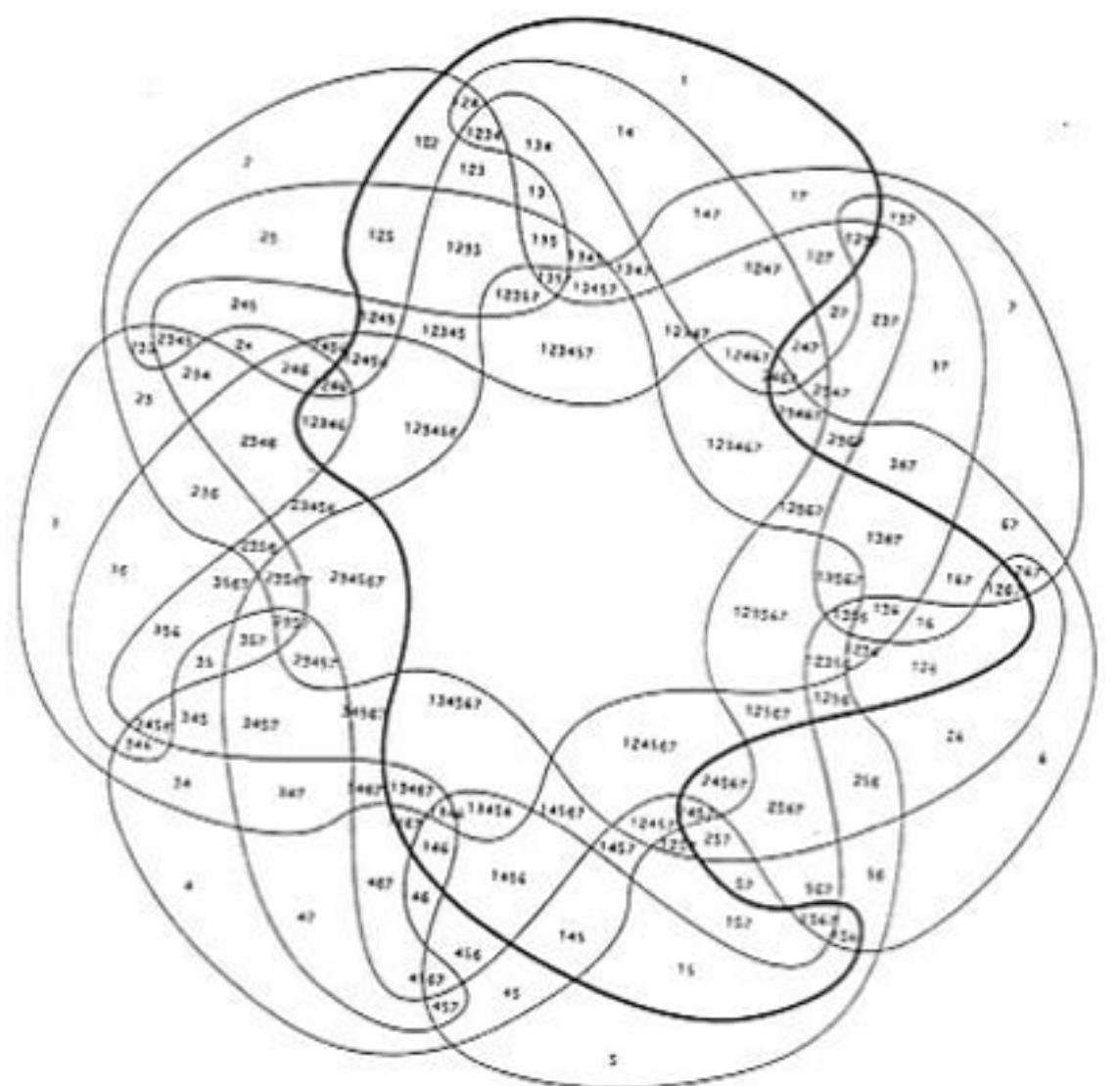


Venn Diagrams

get messy fast

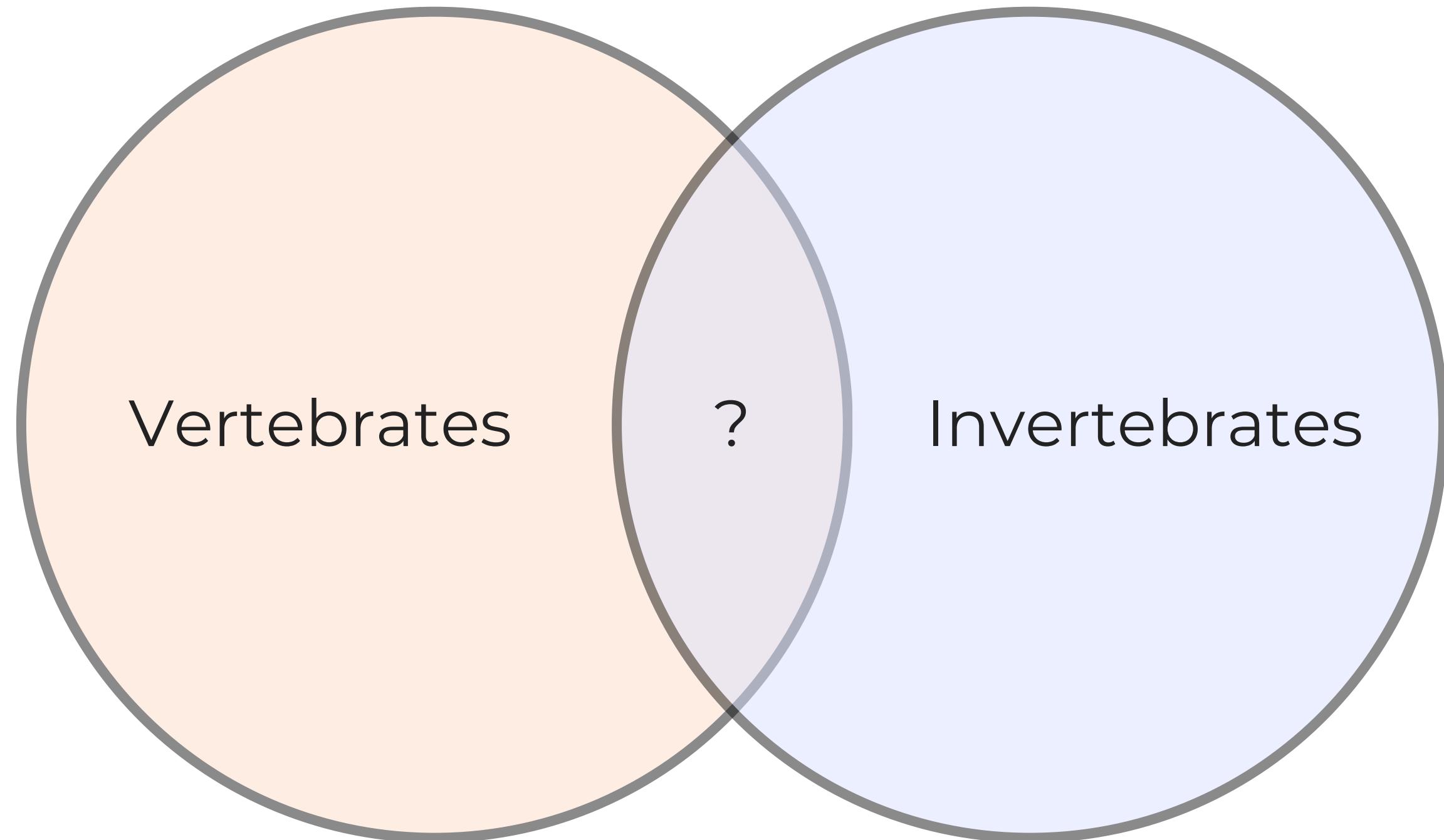
Venn Diagrams

get messy fast



Venn Diagrams

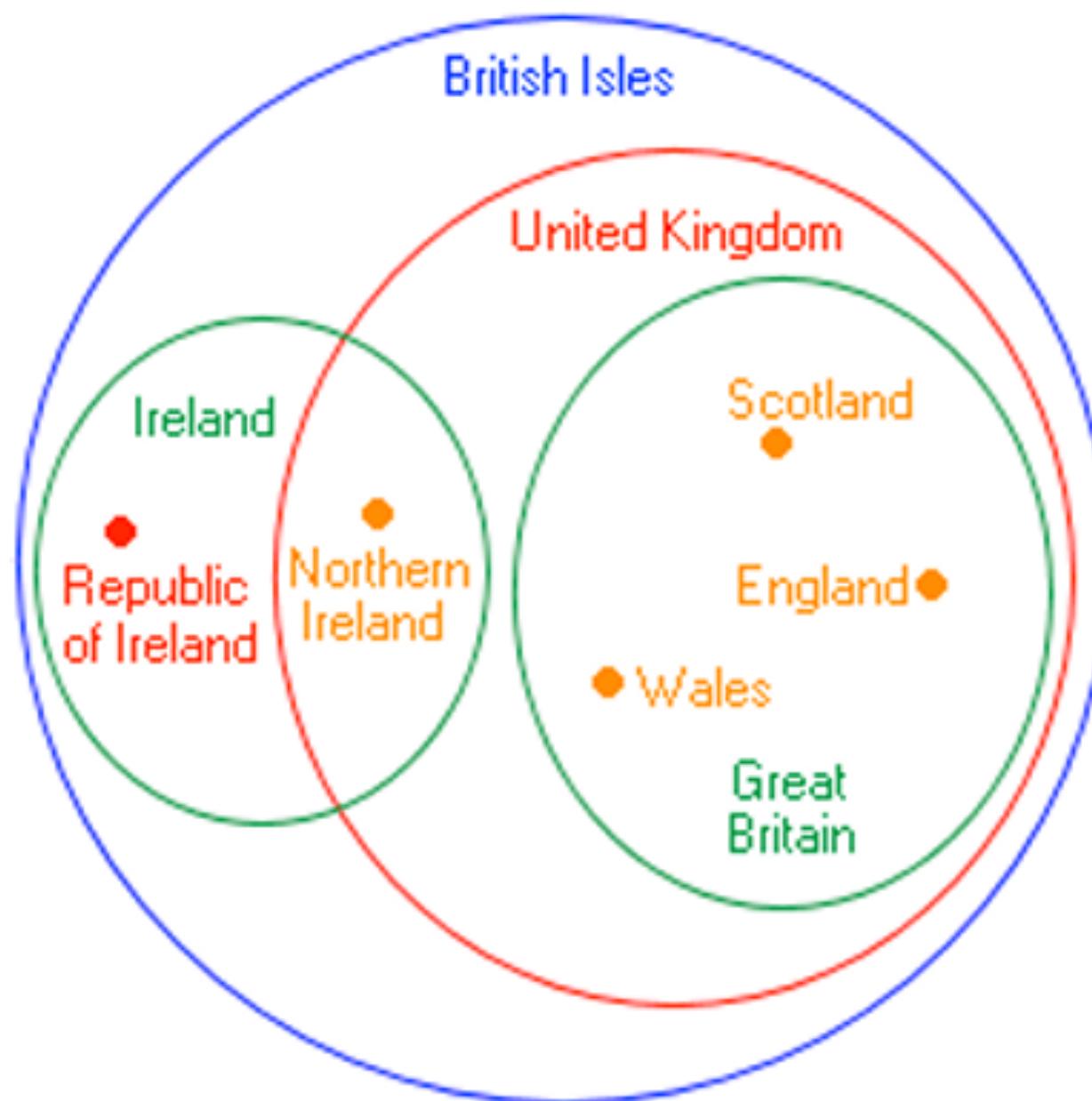
non-sensical?



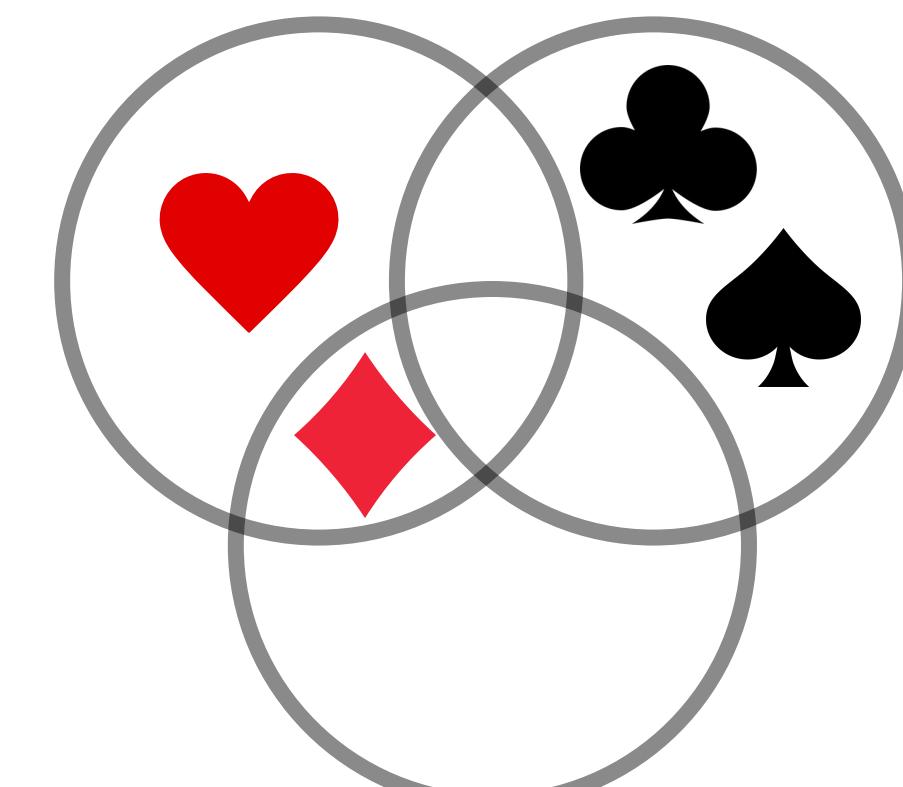
VENN

Euler Diagrams

only visualize existing relationships



Red Suits



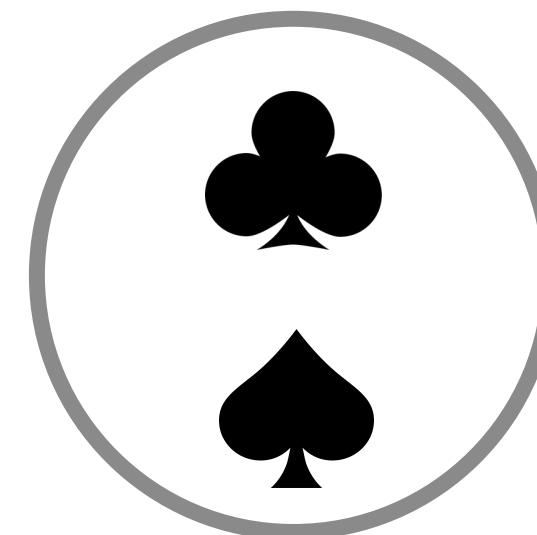
Diamonds

Red Suits



Diamonds

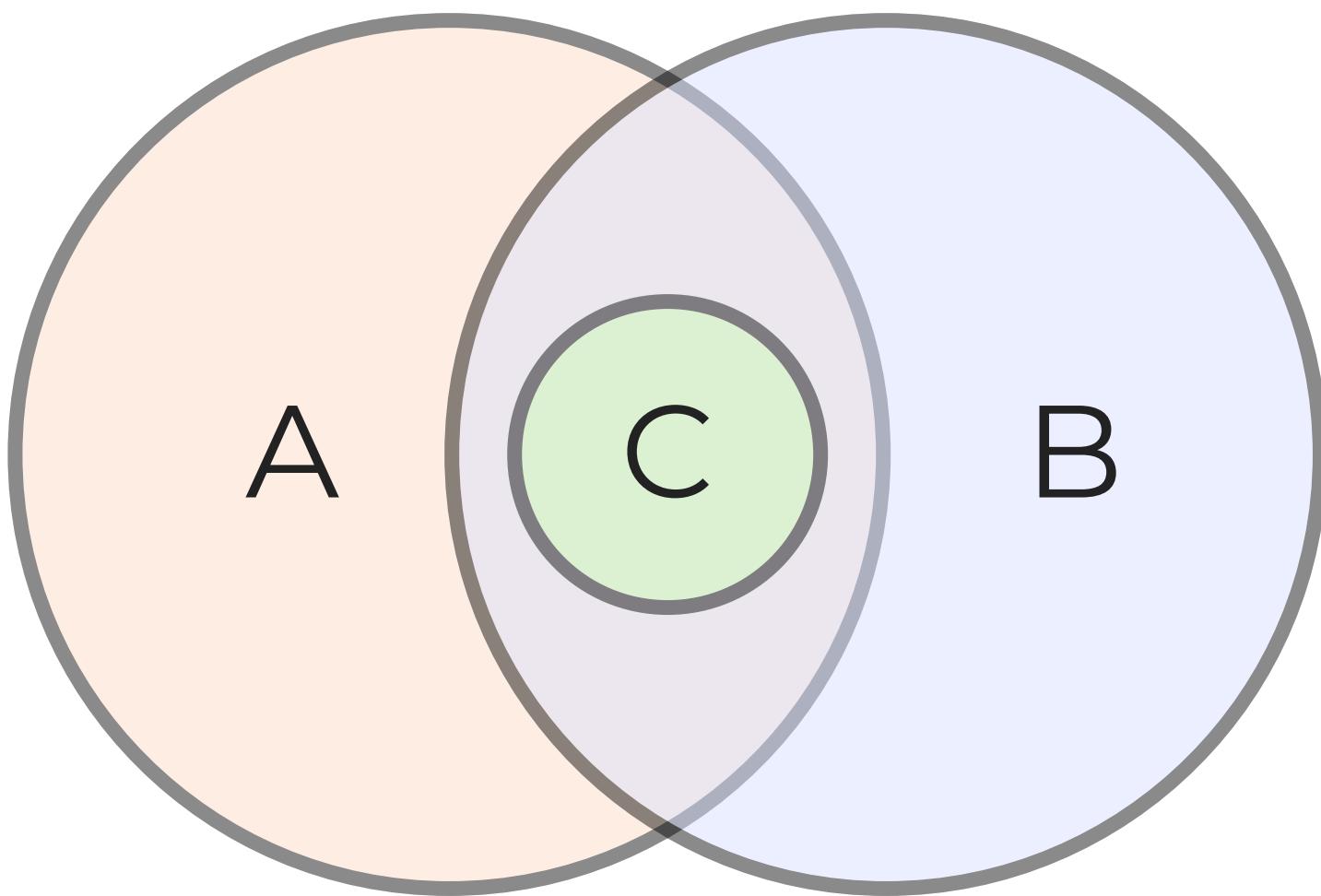
Black Suits



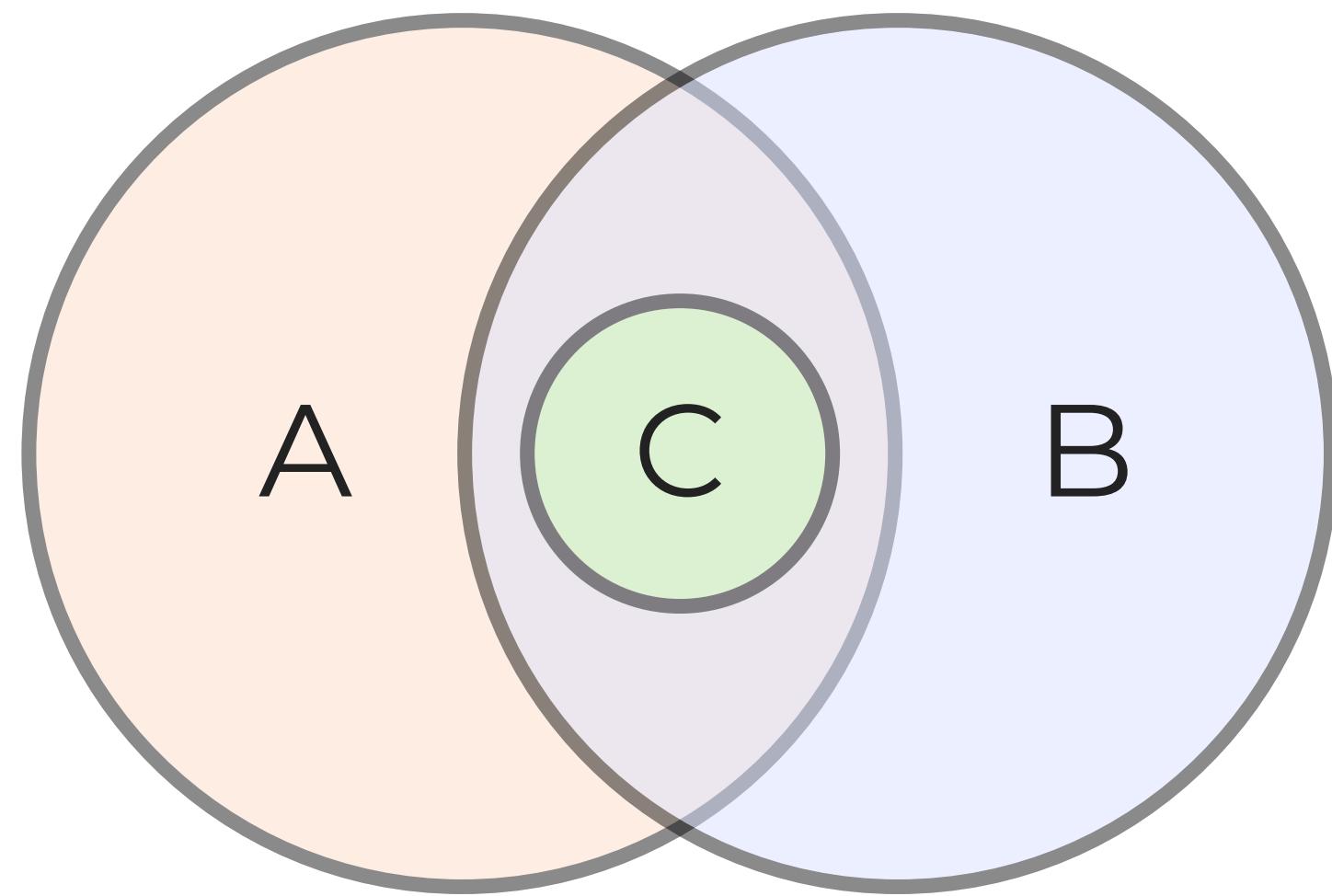
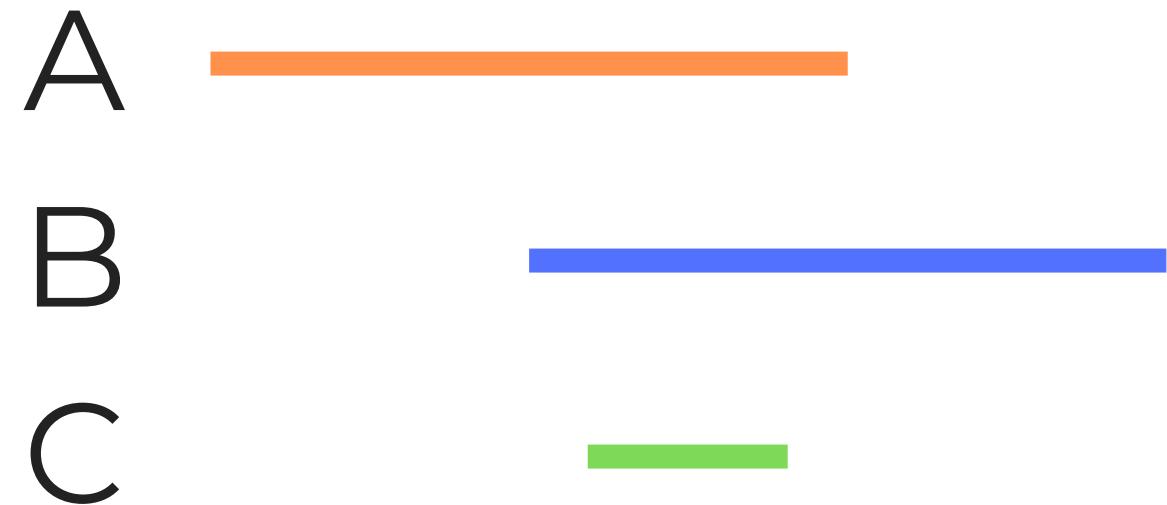
EULER

Euler Diagrams

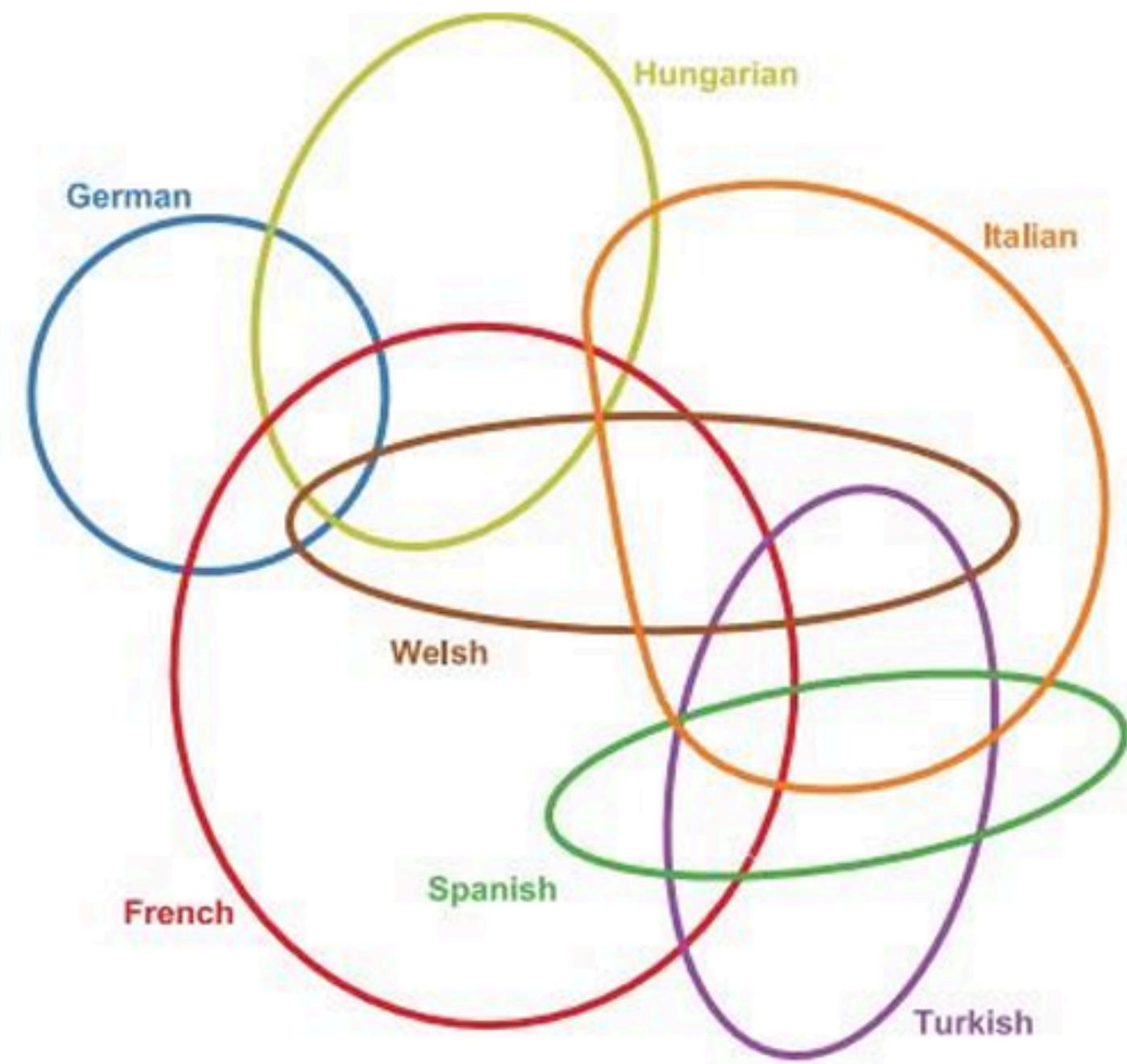
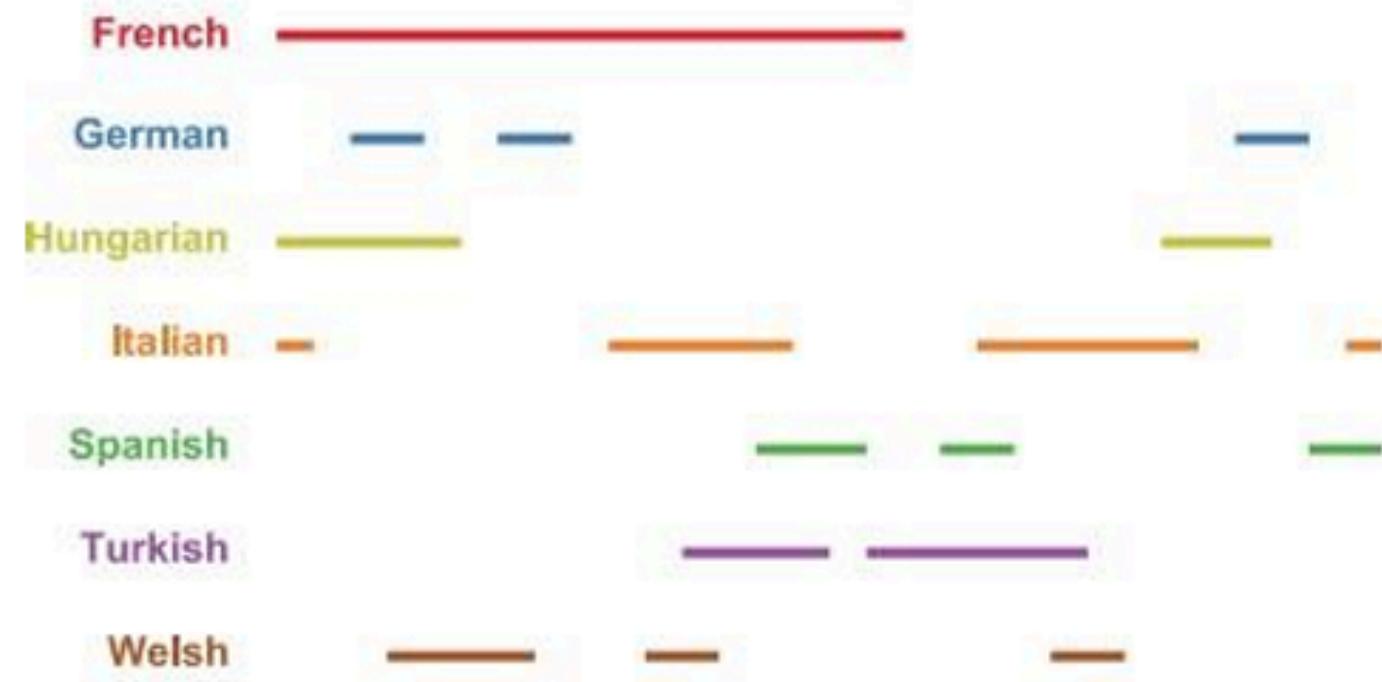
- A: People who **know what a Venn Diagram is**
- B: People who **know what an Euler Diagram is**
- C: People who **know the difference**



Linear Diagrams



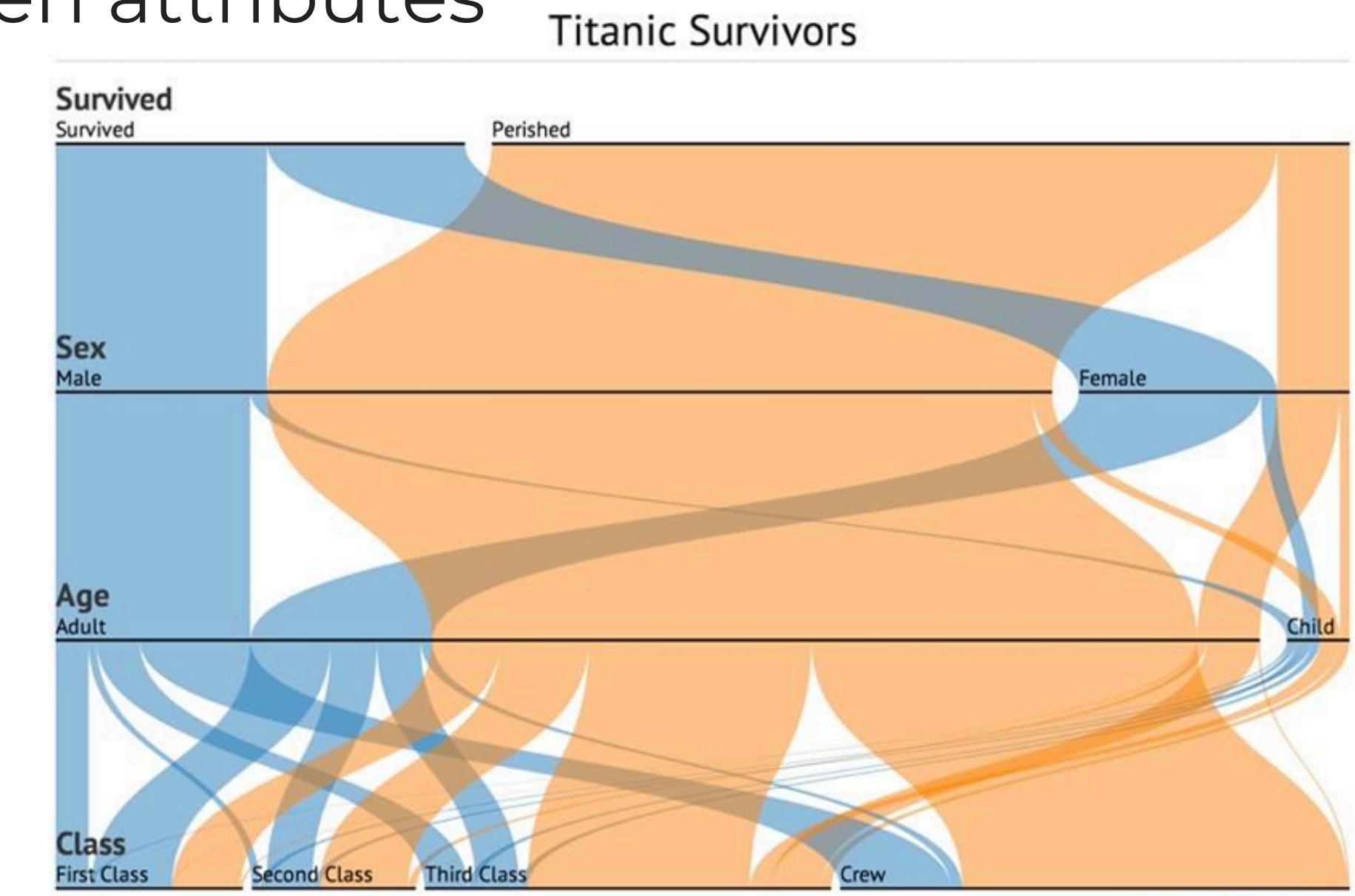
Linear Diagrams



Parallel Sets

Like parallel coordinates, but for categorical data

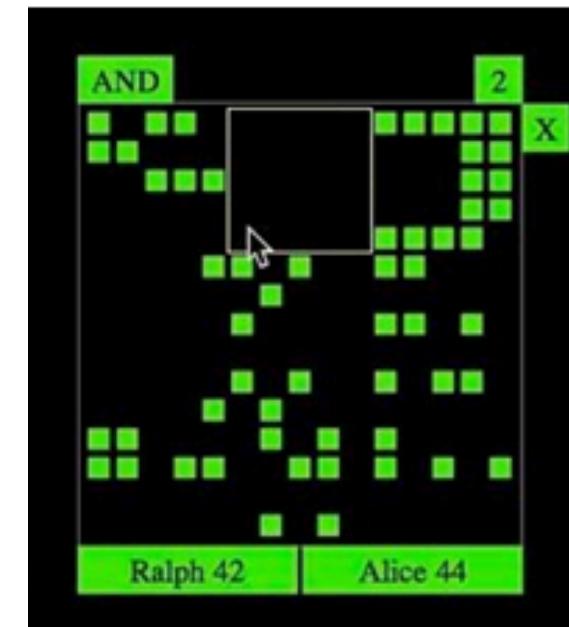
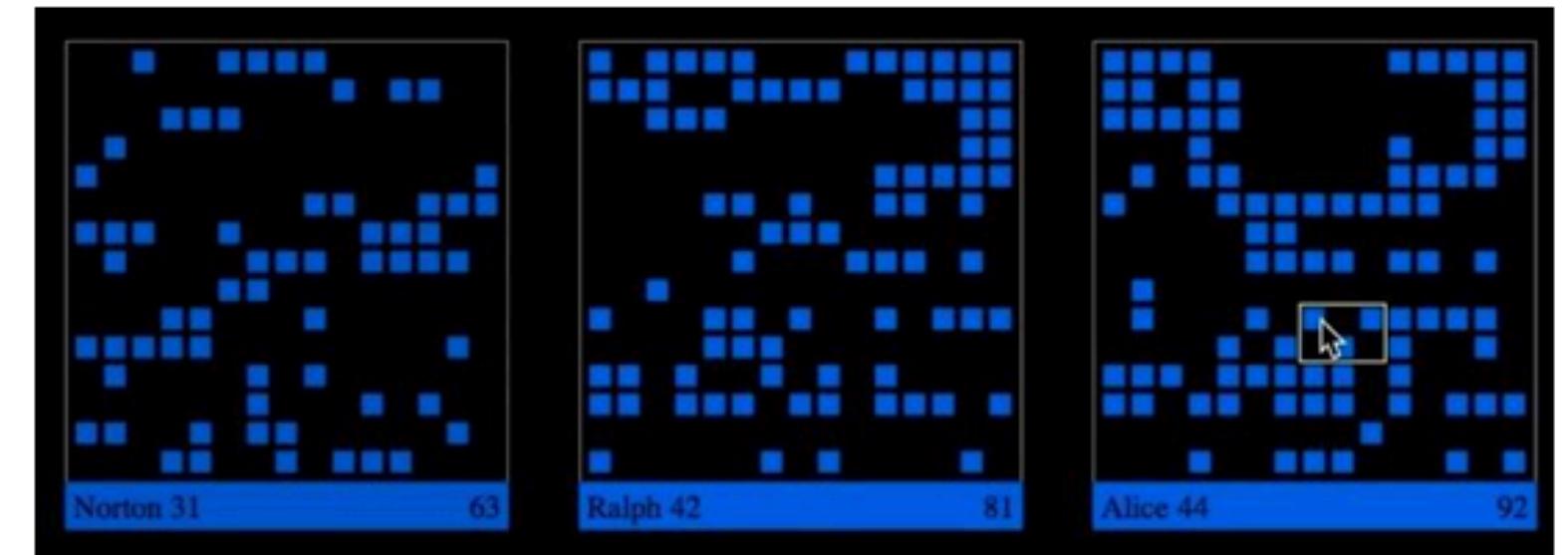
Task: Find relationship between attributes



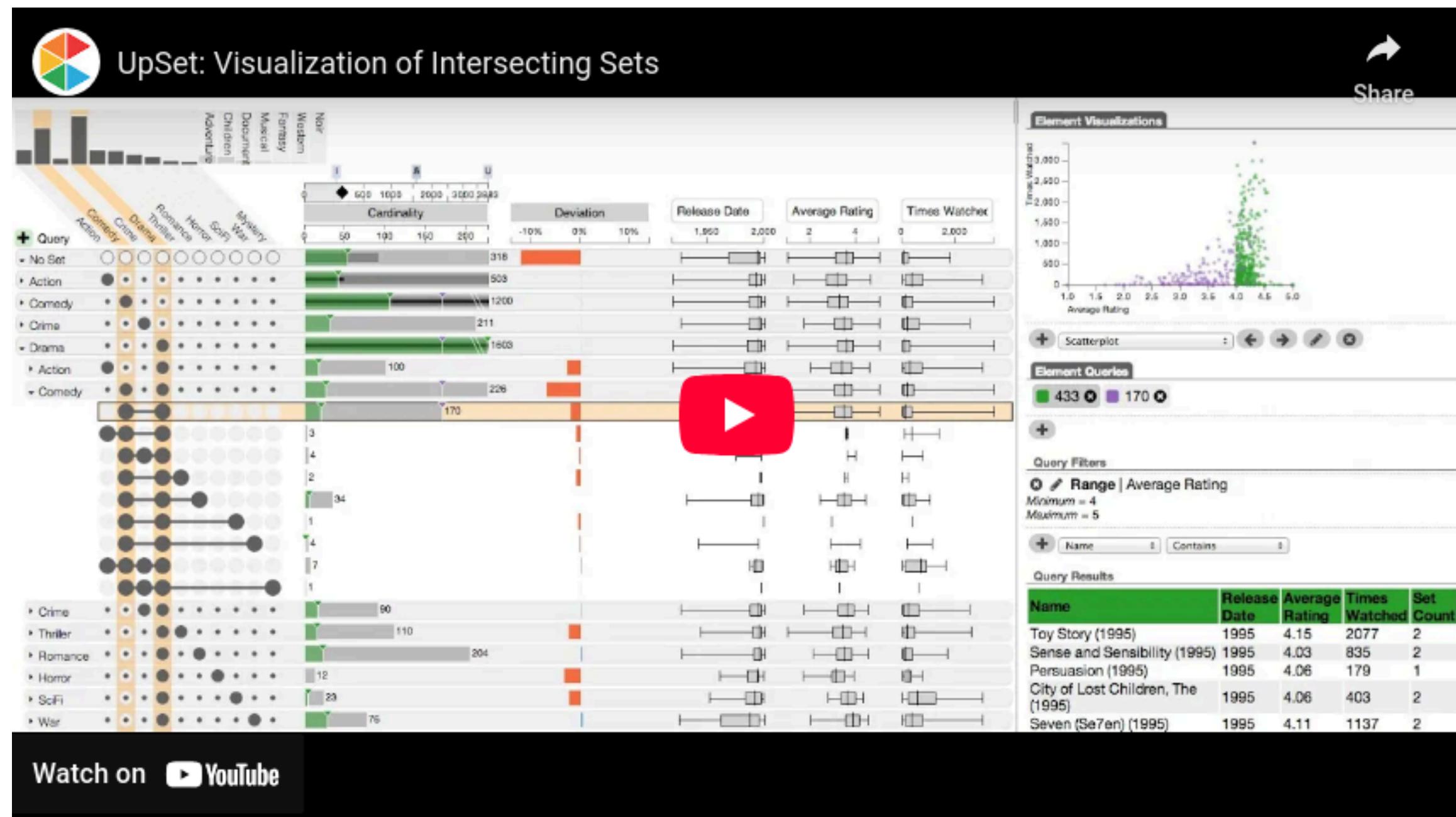
OnSet

Set membership for each item (formatted as a matrix)

supports visual logic operations



UpSet



<https://www.youtube.com/watch?v=-lfF2wGw7Qk>

Set Visualization Takeaways:

Applies to many datasets:

- Many categorical datasets can be viewed as sets

Limited number of sets = more tractable

5 Dataset Types

Tables

Items

Attributes

Networks &
Trees

Items (nodes)

Links

Attributes

Fields

Grids

Positions

Attributes

Geometry

Items

Positions

Clusters,
Sets, Lists

Items

5 Dataset Types (+1?)

Text is unique.

Tables

Items

Attributes

Networks &
Trees

Items (nodes)

Links

Attributes

Fields

Grids

Positions

Attributes

Geometry

Items

Positions

Clusters,
Sets, Lists

Items

Text

Text is Everywhere

**What kind of analysis questions might one ask
about text and documents?**

Which documents contain text about topic X?

Are there other documents that are similar to this one?

How are different words used in a document?

Which documents have an angry tone?

How does one set of documents differ from another set?

I want to understand the history of changes in a document.

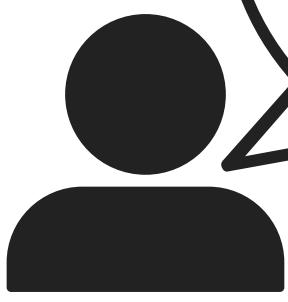
We can think of text as data

**Text is a different type of data than your typical
quantities and categories**

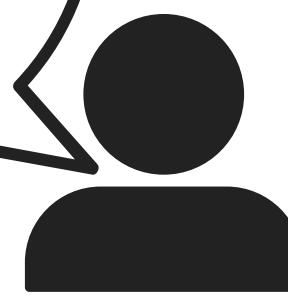
Text as Data

Visualizing text as data can include visualizing:

- individual documents
 - e.g., a news article, an email, a book, etc
- collections of documents
 - e.g., message exchanges, social media profiles, sets of academic publications, etc



Why should we visualize text?
Couldn't we just read the raw text...



We visualize text to show **patterns, trends, and themes**.
Then, we can **focus our time reading only the parts**
of the text that are relevant.

Benefits of Visualizing Text

Fast-track understanding: Visualizing text allows us to get the **“gist” of a document** in a less linear fashion.

Group: Visualizing text allows us to **cluster** the documents to get an **overview**.

Compare: Visualizing text allows us to **compare collections of documents** or **inspect the same collection over time**.

Correlate: We might want to see **how patterns in text relate to other data** (e.g., frequency of a company mentioned in the media with that company’s stock prices).

Example: We can compare the transcripts of speeches from Clinton and Obama to see how their stances on healthcare differed

September 10, 2009

TEXT

Obama's Health Care Speech to Congress

Following is the prepared text of President Obama's speech to Congress on the need to overhaul health care in the United States, as released by the White House.

Madame Speaker, Vice President Biden, Members of Congress, and the American people:

When I spoke here last winter, this nation was facing the worst economic crisis since the Great Depression. We were losing an average of 700,000 jobs per month. Credit was frozen. And our financial system was on the verge of collapse.

As any American who is still looking for work or a way to pay their bills will tell you, we are by no means out of the woods. A full and vibrant recovery is many months away. And I will not let up until those Americans who seek jobs can find them; until those businesses that seek capital and credit can thrive; until all responsible homeowners can stay in their homes. That is our ultimate goal. But thanks to the bold and decisive action we have taken since January, I can stand here with confidence and say that we have pulled this economy back from the brink.

I want to thank the members of this body for your efforts and your support in these last several months, and especially those who have taken the difficult votes that have put us on a path to recovery. I also want to thank the American people for their patience and resolve during this trying time for our nation.

But we did not come here just to clean up crises. We came to build a future. So tonight, I return to speak to all of you about an issue that is central to that future – and that is the issue of health care.

I am not the first President to take up this cause, but I am determined to be the last. It has now been nearly a century since Theodore Roosevelt first called for health care reform. And ever since, nearly every President and Congress, whether Democrat or Republican, has attempted to meet this challenge in some way. A bill for comprehensive health reform was first introduced by John Dingell Sr. in 1943. Sixty-five years later, his son continues to introduce that same bill at the beginning of each session.

Clinton's Health Plan: 'Security, Simplicity, Savings, Choice, Quality and Responsibility'

Transcript of President's Address to Congress on Health Care

Following is a transcript of President Clinton's address to a joint session of Congress last night on his plan to overhaul the nation's health care system, as recorded by The New York Times.

Mr. Speaker, Mr. President, members of Congress, distinguished guests, my fellow Americans: Before I begin my words tonight, I would like to ask that we all bow in a moment of silent prayer for the memory of those who were killed and those who have been injured in the tragic train accident in Alabama today.

My fellow Americans, tonight we come together to write a new chapter in the American story.

Our forebears enshrined the American dream — life, liberty, the pursuit of happiness. Every generation of Americans has worked to strengthen that legacy, to make our country a place of freedom and opportunity, a place where people who work hard can rise to their full potential and a place where their children can have a better future.

From the settling of the frontier to the landing on the moon, ours has been a continuous story of challenges defined, obstacles overcome, new horizons secured. That is what makes America what it is and Americans what we are.

Now we are in a time of profound change and opportunity. The end of the cold war, the information age, the global economy have brought us both opportunity and hope and strife and uncertainty.

Our purpose in this dynamic age must be to change, to make change our friend and not our enemy.

To achieve that goal, we must face all our challenges with confidence, with faith and with discipline, whether we're reducing the deficit, creating tomorrow's jobs and training our people to fill them, converting from a high-tech defense to a high-tech domestic economy, expanding trade, reinventing government, making our streets safer or rewarding work over idleness. All these challenges require us to change.

If Americans are to have the courage to change in a difficult time, we must first be secure in our most basic needs.

Fixing a Broken System

Tonight I want to talk to you about the most critical thing we can do to build that security.

This health care system of ours is badly broken and it is time to fix it.

Despite the dedication of literally millions of talented health care professionals, our health care is too uncertain and too expensive, too bureaucratic and too wasteful. It has too much fraud and too much greed.

At long last, after decades of false starts, we must make this our most urgent priority: giving every American health security



June R. Loper/The New York Times

President Clinton acknowledging applause last night. Behind him were Vice President Al Gore and Speaker Thomas S. Foley.

look up to. They are what is right with this health care system. But we also know that we can no longer afford to continue to ignore what is wrong.

Despite the dedication of literally millions of talented health care professionals, our health care is too uncertain and too expensive, too bureaucratic and too wasteful. It has too much fraud and too much greed.

At long last, after decades of false starts,

we must make this our most urgent priority: giving every American health security

Under our plan every American would receive a health care security card that will guarantee a comprehensive package of benefits over the course of an entire lifetime roughly comparable to the benefit package offered by most Fortune 500 companies.

This health care security card will offer this package of benefits in a way that can

than they would if they had regular treatment in a way that only adequate medicine can provide.

I also believe that over time we should phase in long-term care for the disabled and the elderly on a comprehensive basis.

Toward a Simpler System

and we can change it.

And doctors, nurses and consumers shouldn't have to worry about the fine print. If we have this one simple form there won't be any fine print. People will know what it means.

Toward Savings In the System

The third principle is savings.

Reform must produce savings in this health care system — it has to. We are spending over 14 percent of our income on health care; Canada's at 10; nobody else is over 9. We're competing with all these people for the future and the other major countries, they cover everybody. And they cover them with services as generous as the best company policies here in this country.

Rampant medical inflation is eating away at our wages, our savings, our investment capital, our ability to create new jobs in the private sector and this public treasury.

You know the budget we just adopted had steep cuts in defense, a five-year freeze on the discretionary spending so critical to reeducating America and investing in jobs and helping us to convert from a defense to a domestic economy. But we passed the budget, which has Medicaid increases of between 16 and 11 percent a year over the next five years and Medicare increases of between 11 and 9 percent in an environment where we assume inflation will be at 4 percent or less.

We cannot continue to do this. Our competitiveness, our whole economy, the integrity of the way the Government works, and ultimately our living standards, depend upon our ability to achieve savings without harming the quality of health care. Unless we do this, our workers will lose \$655 in income each year by the end of the decade. Small businesses will continue to face skyrocketing premiums, and a full third of small businesses now covering their employees say they will be forced to drop their insurance.

Large corporations will bear bigger disadvantages in global competition, and health care costs will devour more and more and more of our budget. Pretty soon all of you or the people who succeed you will be showing up here, and writing out checks for health care and interest on the debt and worrying about whether we've got enough defense, and that'll be it unless we have the courage to achieve the savings that are plainly there before us. Every state and local government will continue to cut back on everything from education to law enforcement to pay more and more for the same health care.

These rising costs are a special nightmare for our small businesses, the engine of our entrepreneurship and our job creation in America today. Health care premiums for small businesses are 35 percent higher than those of large corporations today, and they will keep rising at double-digit rates unless we act.

Word (Tag) Clouds

Visualize the frequency of words in a document



Word (Tag) Clouds

Visualize the frequency of words in a document



Obama's Speech (2009)

people will care every system us insurance one way

Portrait of the candidate as a pile of words

What's the most frequent word on McCain's blog? "Obama."

BY JOHN SCHWENKEL

BOTH OF THIS year's major party presidential candidates have made official campaign blogs a central part of their virtual headquarters. John McCain's campaign hired former Weekly Standard blogger Michael Gerson as its deputy communications director and made the "McCain Report" one-tile-newspaper above The Obama website, has taken a more collaborative approach, allowing supporters to create their own blogs and then moving some of that content, together with snippets of news coverage and messages from campaign officials, to the official "ObamaBlog" front page. The results are revealing. Online software at Wordle.net makes it easy to parse the blogs visually: with the three most frequently used words displayed largest. A snapshot from last week:

MCCAIN'S BLOG



obama Given the way the McGovern campaign has been complaining about unbalanced media attention, it's puzzling to see that it has been as fixated on Obama as everyone else. Plausible (evening media) coverage is one thing; successfully turning the election into a national referendum on the relatively inexperienced junior senator from Illinois just might yield McCain an improbable victory. It would also saddle, though, if voters tire of excessively negative politicking and end up looking kindly on the winner of the two races.

visit The ongoing attempt to stir controversy over the cancellation of Obama's scheduled visit to the Landstuhl Army Hospital has been a consistent part of the blogging strategy at the McGa Report, and in this case the growing backlash against these attacks from Republican bloggers have called them out. He has only managed to fuel that tendency, with each round of explanations being picked through, point by point.

drilling Following the lead of the national GOP, the McCain campaign has recently started hammering away at the need to drill for oil. This is a bit of an unusual move for McCain – it hasn't been long since his much-publicized break from Republican orthodoxy on exactly this issue – but it may be one he has to make as both candidates look to establish a compelling message on the economy.

Barack Obama's opponents offstage, a solidifying holiday has led the McCain team to "Ex. No. 1." Most analysts are unlikely to have any significant impact, but McCain's propensity for the widespread desire for a "step up and do something," while difficult to sustain, is a potent factor.

OBAMA'S BLOG



future Barack Obama's youth and relative newcomer to national politics make him the perfect candidate to carry his message on leaving the past behind and moving into a better future. This sort of progressive rhetoric has the effect of keeping people optimistic while still reminding them that the status quo just isn't good

change Obama's trademark theme of change is still figuring prominently in his campaign rhetoric, and with good reason: the national mood is deeply pessimistic, and everyone is looking for signs of something different. While this sort of campaigning is decidedly vague, it has worked well enough so far that there's no re-

donate	The Obama campaign's official blogging constantly emphasizes the importance of donating to supporters' candidacy. While Obama's decision to reverse his position on public financing drew a fair amount of criticism, the perception of a campaign funded by small donors rather than wealthy special interests
supporters	In addition to keeping an already-focussed-the candidate himself, Obama's campaign blog often recounts the personal stories of his individual supporters. This is a way of keeping that enthusiasm high, and making his grassroots supporters feel they are doing something important.

Apparently, both candidates mentioned Obama a lot

Word (Tag) Clouds

Visualize the frequency of words in a document

Word clouds are:

- social, friendly, approachable
- harmless when the goal is to have fun
- less harmless when used for communicating scientific or journalistic information

Word (Tag) Clouds: Issues

Word clouds:

- prioritize frequent words rather than important words
- prioritize words that have few alternatives
- ignores core words that are only explicitly mentioned a few times
- are perceptually difficult to parse
- confound the size channel with word length
- don't actively support comparison

Some alternatives to word clouds

Sorted Word Clouds



addition afford afghanistan ago agree ahead alive america american americans army auto back
benefits break bush business businesses Care cars century challenges chance change child children clean clear clinton college
companies country create cut daughters day days debate decades decent democrats deserve dignity dollars dreams done economic
economy education election end energy face failure families family finally find finish fix fundamental fundamentals future
generation george give giving good government grateful great hands hard health hear heard higher home hope idea ideas
invest iraq job jobs john judgment kennedy lead leave life lives long longer lost love made make makes making man
market mccain measure meet men michelle middle-class military million moment moments money moral nation new
night nuclear obligation oil part party past pay people percent plan plans plan politics poverty power president programs
progress promise protect proud provide pursue put ready renewable republicans require respect responsibility restore
reward rise safe security senator sense sick sights small stand standards start states stood strength student talk talking tax taxes teachers
technology ten things thirty threats time today tonight tough troops turn understand united veterans walk washington
watch watching whiners woman women work worked workers working world years young

Sorted Word Clouds

WORDCOUNT

◀ PREVIOUS WORD

NEXT WORD ►

the of and to a in that it is was i for on you he be with as by have given this no but with in they from y which now there is none equal to it

CURRENT WORD

FIND WORD:

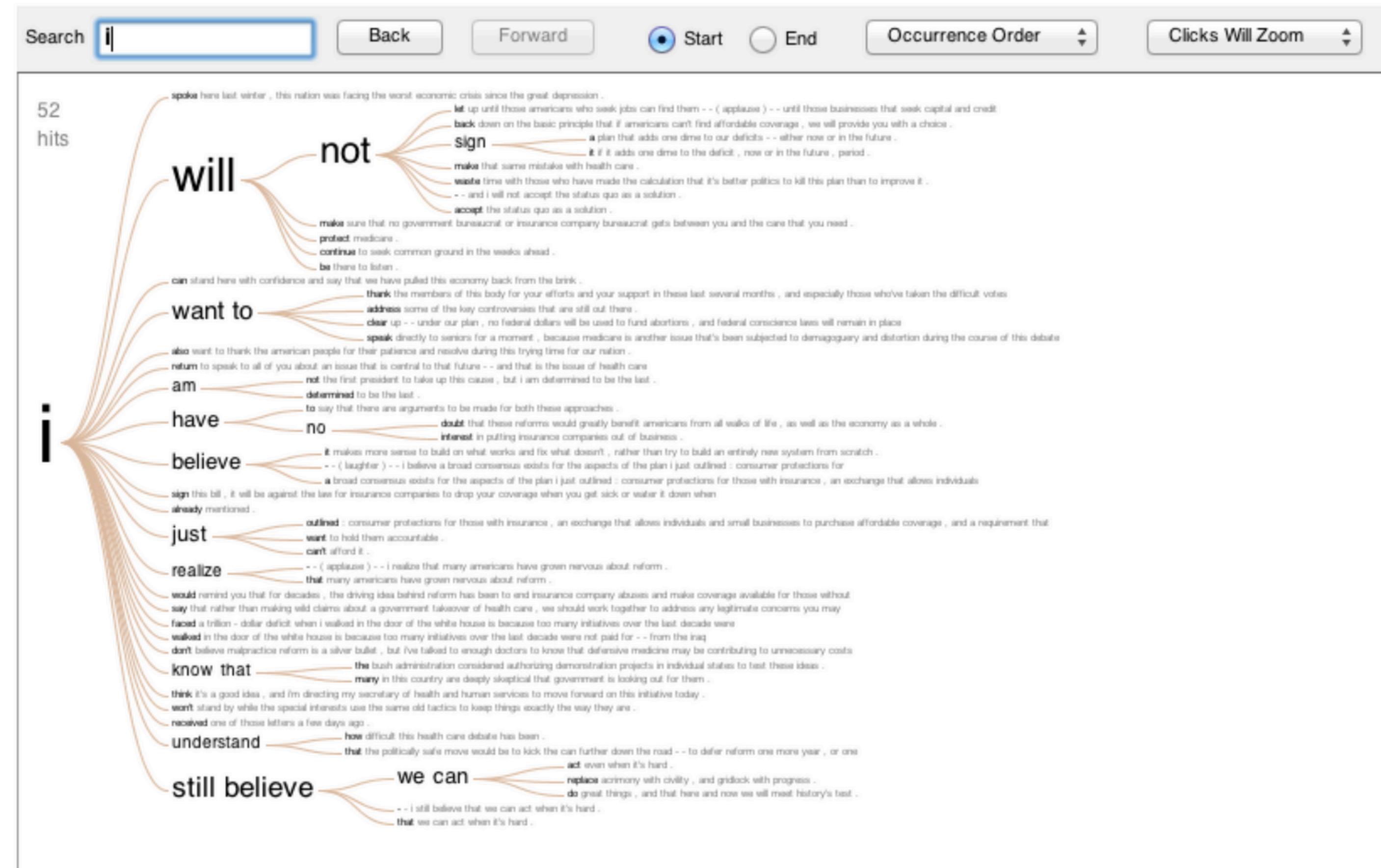
BY RANK:

REQUESTED WORD: THE

86800 WORDS IN ARCHIVE

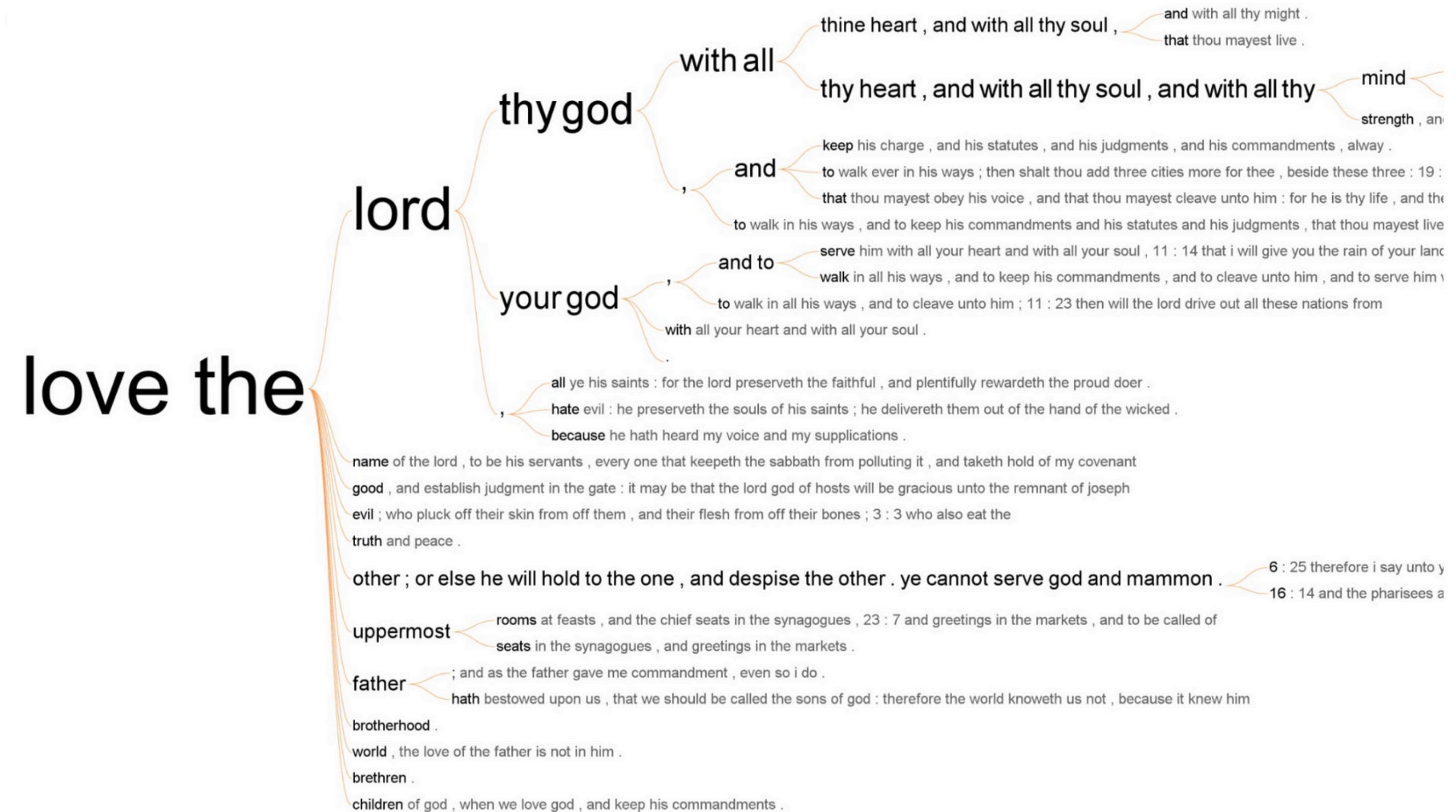
RANK: 1

WordTree



Obama's Speech (2009)

WordTree



The Bible

WordTree

if love be rough with you , be rough with love .

if love be blind , love cannot hit the mark .

if love be blind , it best agrees with night .

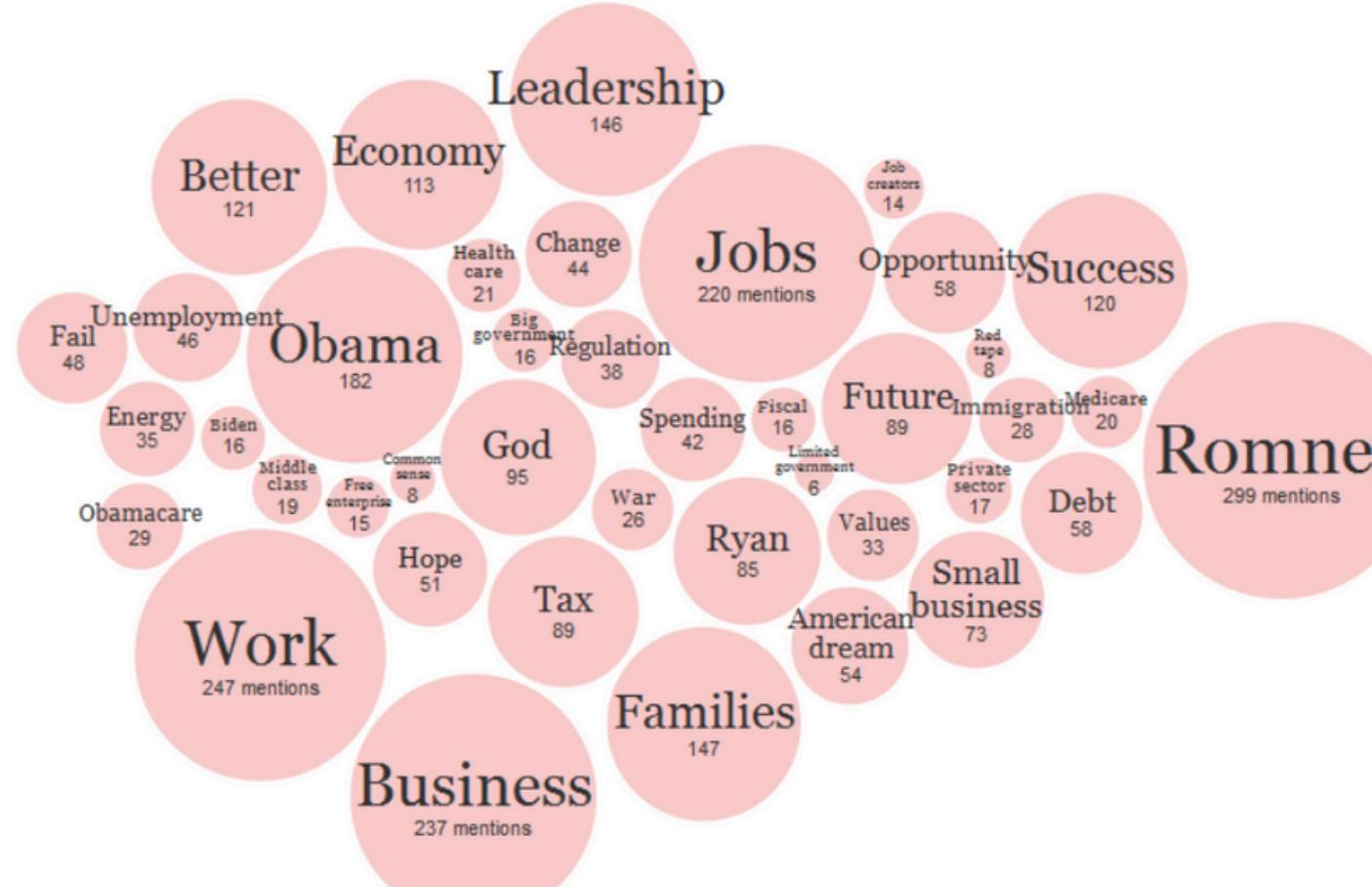
if love be rough with you , be rough with love .
if love be blind , love cannot hit the mark .
it best agrees with night .

Bubble Charts

At the Republican Convention, the Words Being Used

A look at how often speakers at the Republican National Convention have used certain words and phrases so far, based on an analysis of transcripts from the Federal News Service.

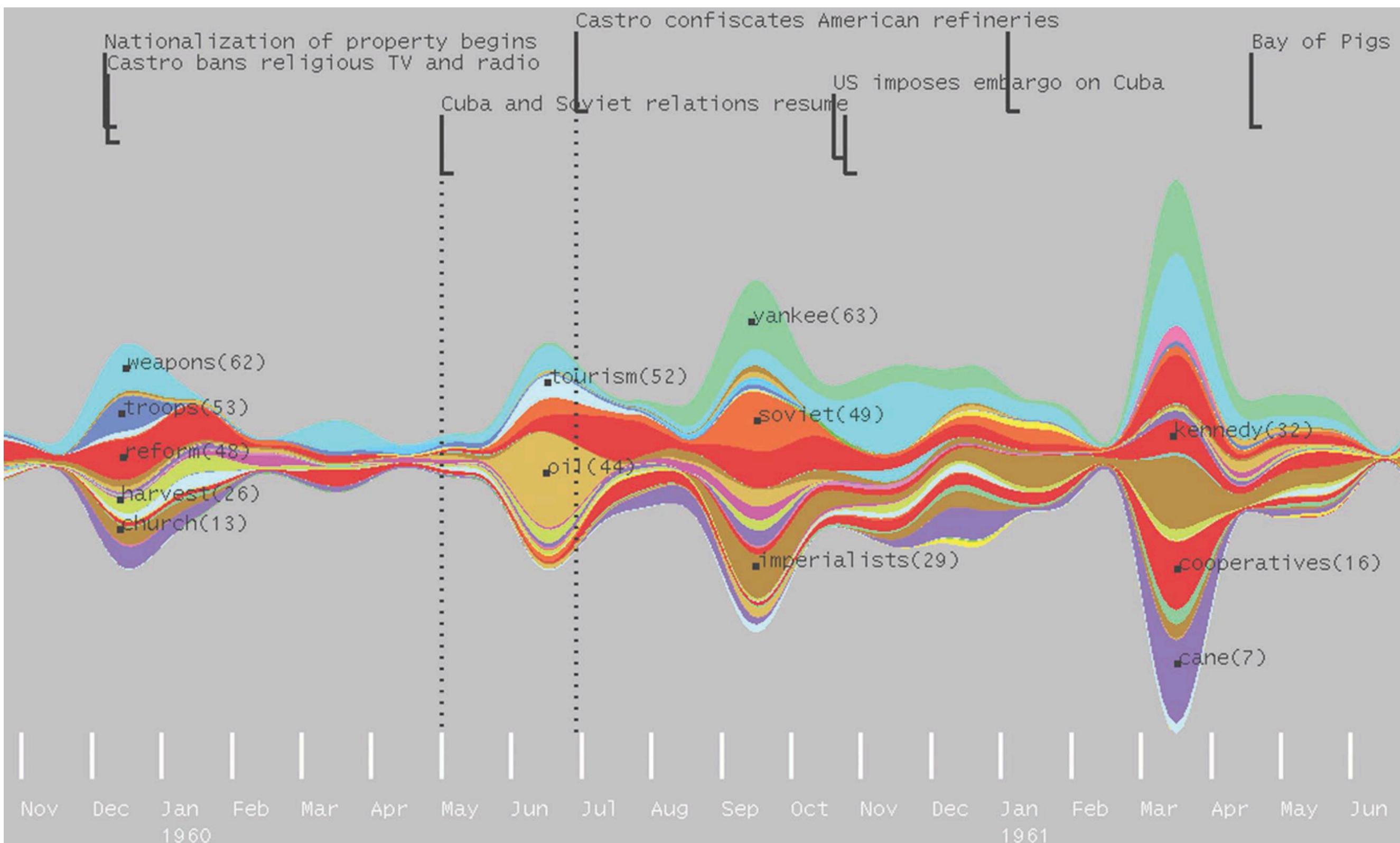
Add word or phrase



Parallel Tag Clouds



ThemeRiver

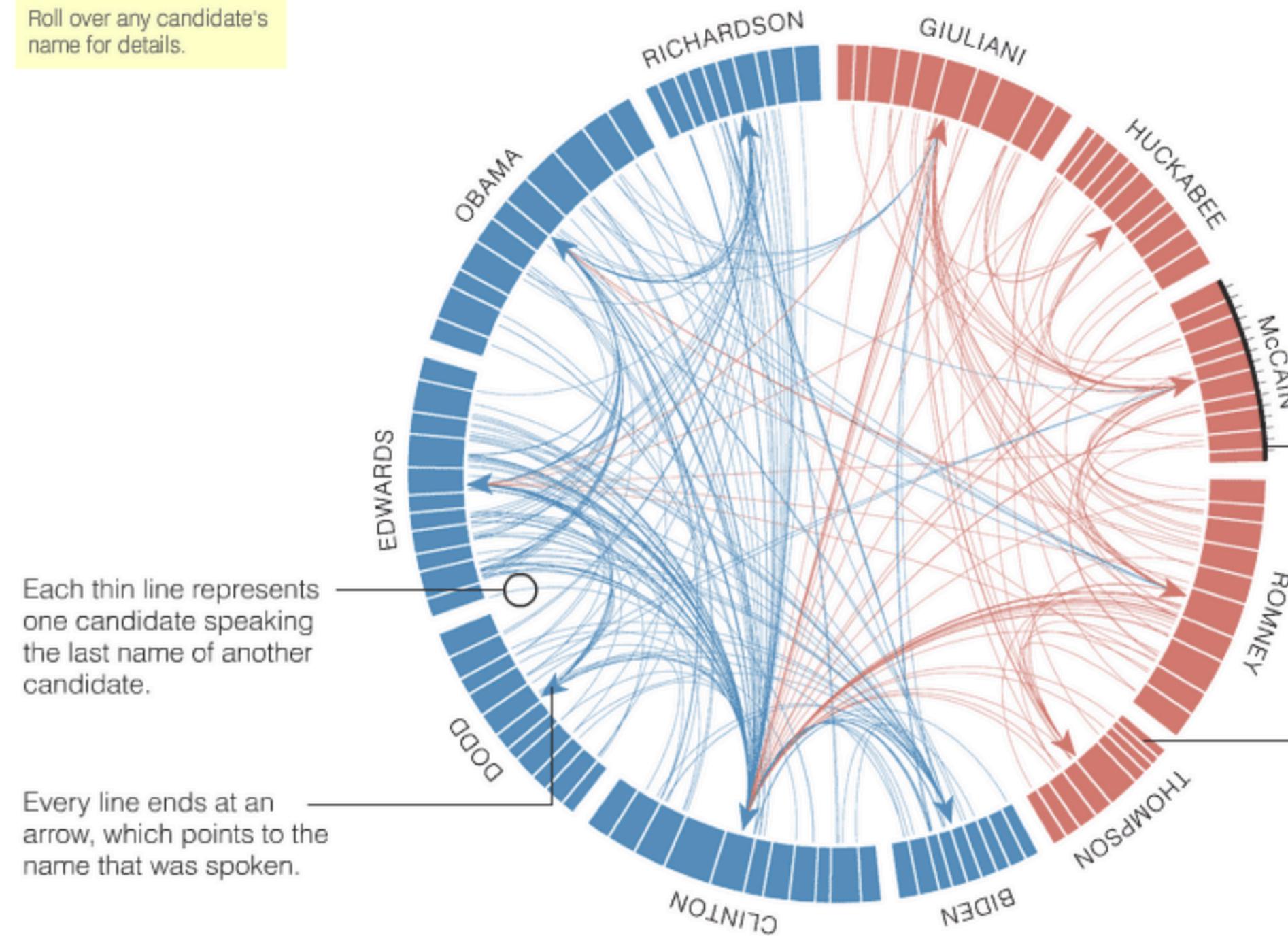


Chord/Arc Diagram

Naming Names

Names used by major presidential candidates in the series of Democratic and Republican debates leading up to the Iowa caucuses.

Roll over any candidate's name for details.



**This is called
a chord or arc
diagram.**

The length of each circle segment represents the total number of words spoken by the candidate during the debates. Each tick mark represents 1,000 words.

Each slice represents one debate, arranged clockwise from the first to the final debate.

Each thin line represents one candidate speaking the last name of another candidate.

Every line ends at an arrow, which points to the name that was spoken.

Challenges in Text Visualization

Challenges in Text Visualization

1. Text is mostly unstructured (categorical, but atypical)

1. Text is Mostly Unstructured

Text is categorical, but not your typical categorical

Words can be categorically different, but:

- semantically similar (e.g., tennis, swimming, running)
- synonymous (e.g., easy, effortless, facile)
- ordered (e.g., January, February, etc.)
- and so on..

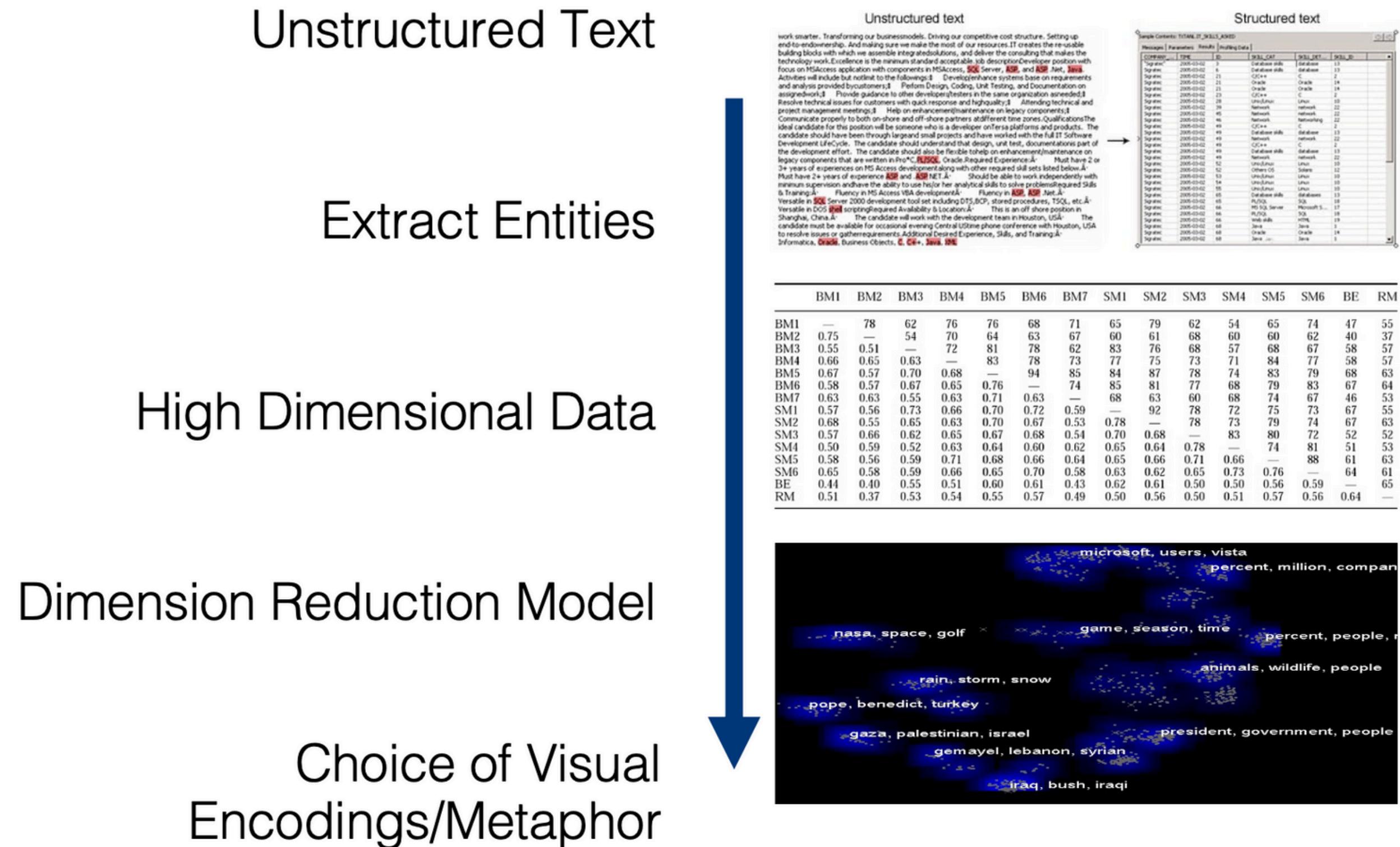
For example:

- Fall = Autumn (not the same word)
- Fall ?= Fall (not the same meaning)

Challenges in Text Visualization

1. Text is mostly unstructured (categorical, but atypical)
2. Textual Data can be very large and high-dimensional

2. Textual data is large and high-dimensional



The Text Processing Pipeline

Three Models:

- Bag of Words
- Keyword Weighting
- Term Commonness

Bag of Words

Text

Hogwarts School of Witchcraft and Wizardry, commonly shortened to Hogwarts, is a fictional British school of magic for students aged eleven to eighteen, and is the primary setting for the first six books in J. K. Rowling's Harry Potter series...

Albus Percival Wulfric Brian **Dumbledore** is a fictional character in J. K. Rowling's Harry Potter series. For most of the series, he is the headmaster of the wizarding school Hogwarts. As part of his backstory, it is revealed that he is the founder and leader of ...

Collinwood Mansion is a fictional house featured in the Gothic horror soap opera Dark Shadows (1966–1971). Built in 1795 by Joshua Collins, Collinwood has been home to the Collins family—and other sometimes unwelcome supernatural visitors...

Data Table

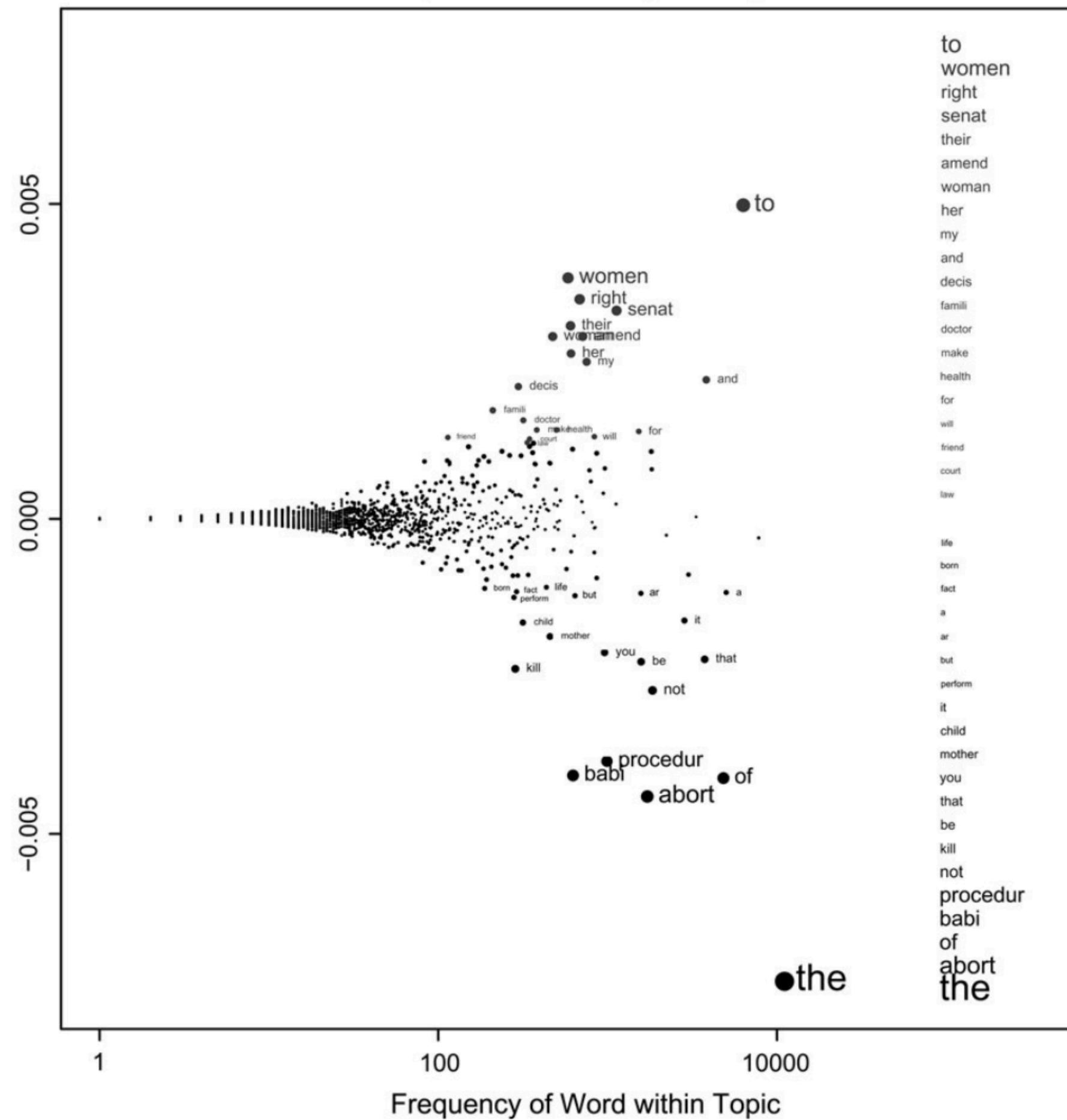
	document		
term	Hogwarts (Top)	Dumbledore (Middle)	Collinwood (Bottom)
a	1	1	1
of	1	2	
in	1	1	2
is	2	4	1
fictional	1	1	1
school	1		
rowling's	1	1	
harry	1	1	
potter	1	1	
series	1	1	
house			1

Keyword Weighting

Term frequency (tf_{td}) = number of times term t occurs in document d

$$\text{Normalized Term Frequency } (tf_{td_{norm}}) = \frac{tf_{td}}{\sum_t tf_{td}}$$

Partisan Words, 106th Congress, Abortion
(Difference of Proportions)



Democrats used more

Republicans used more

Keyword Weighting

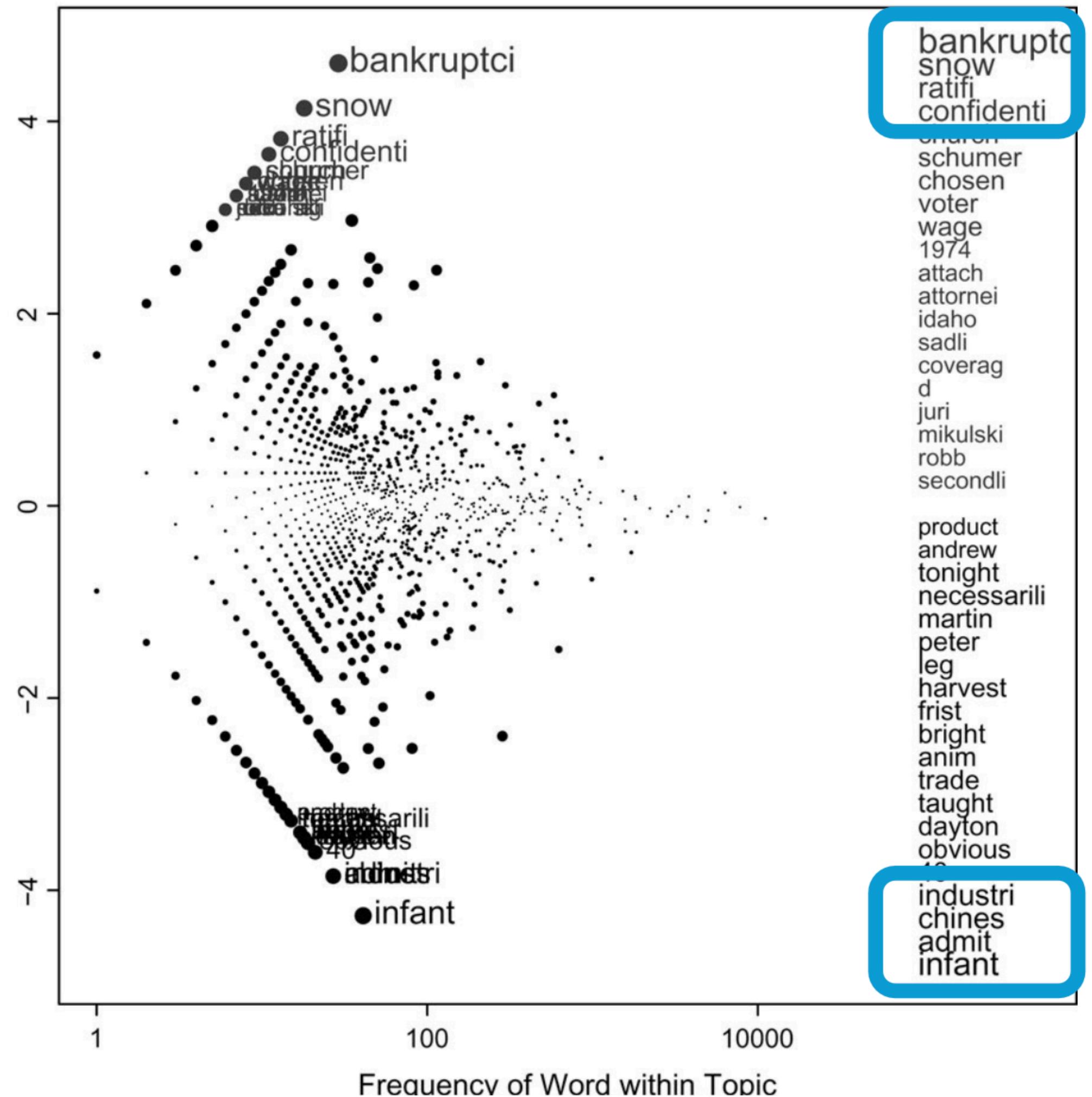
Normalized Term Frequency considering the whole corpus ($tf_{td.wc}$)

$$= \log(1 + tf_{td}) \times \log\left(\frac{N}{d_t}\right)$$

where d_t = number of documents containing term t

N = number of documents in the corpus

Partisan words, 106th Congress, Abortion
(Log-Odds-Ratio, Smoothed Log-Odds-Ratio)



Democrats used more

Republicans used more

Term Commonness

Normalized Term Frequency considering the entire english language ($tf_{td.e}$)

$$= \log (tf_{td}) \times \log (tf_{t.e})$$

The Text Processing Pipeline: Issues

All three models:

- rely only on term frequencies
- ignore contextual information like grammar, part of speech, negation, adjective-noun pairs, pronoun reference, etc.

Oblivion is a better game than Skyrim .

Is Oblivion a better game than Skyrim ?

Same frequencies, different meanings

Challenges in Text Visualization

- 1.Text is mostly unstructured (categorical, but atypical)
- 2.Textual Data can be very large and high-dimensional
- 3.Context is critical for understanding text
- 4.Summarizing, grouping, and parsing text is time-consuming

FIN

Upcoming Dates

May 9:**

- Homework 5 Due
- Project Screencast Submission Due

May 12: Final Project Submission Due

****Only if you requested an extension by May 2**