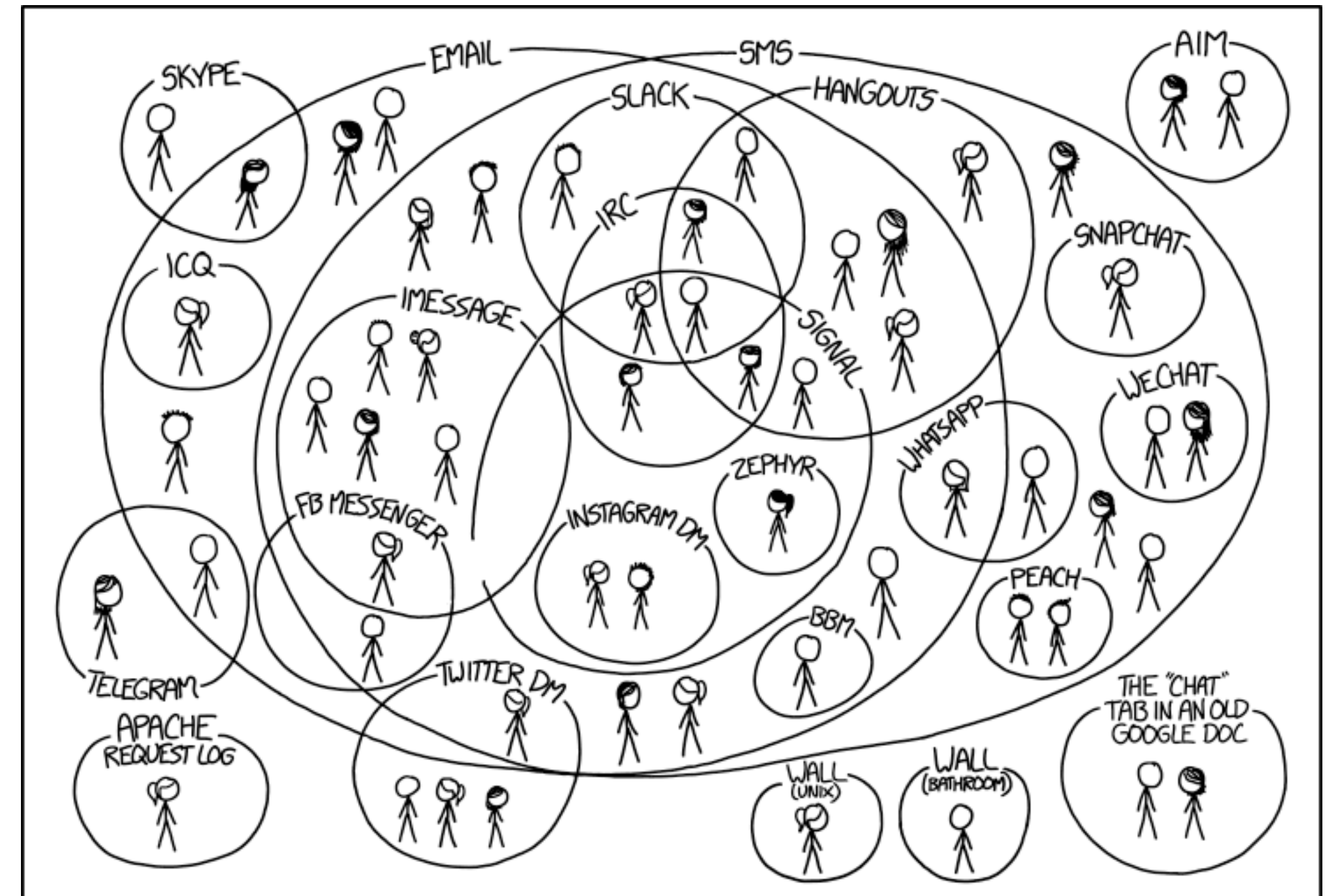


CS-5630 / CS-6630 Visualization for Data Science Set Visualization

Alexander Lex
alex@sci.utah.edu



I HAVE A HARD TIME KEEPING TRACK OF WHICH CONTACTS USE WHICH CHAT SYSTEMS.

The banana (*Musa acuminata*) genome and the evolution of monocotyledonous plants

Angélique D'Hont^{1*}, France Denoeud^{2,3,4*}, Jean-Marc Aury², Franc-Christophe Baurens¹, Françoise Carreel^{1,5}, Olivier Garsmeur¹, Benjamin Noel², Stéphanie Bocs¹, Gaëtan Droc¹, Mathieu Rouard⁶, Corinne Da Silva², Kamel Jabbari^{2,3,4}, Céline Cardi¹, Julie Poulain², Marlène Souquet¹, Karine Labadie², Cyril Jourda¹, Juliette Lengellé¹, Marguerite Rodier-Goud¹, Adriana Alberti², Maria Bernard², Margot Correa², Saravanaraj Ayyampalayam⁷, Michael R. Mckain⁷, Jim Leebens-Mack⁷, Diane Burgess⁸, Mike Freeling⁸, Didier Mbéguié-A-Mbéguié⁹, Matthieu Chabannes⁵, Thomas Wicker¹⁰, Olivier Panaud¹¹, Jose Barbosa¹¹, Eva Hribova¹², Pat Heslop-Harrison¹³, Rémy Habas⁵, Ronan Rivallan¹, Philippe Francois¹, Claire Poirion¹, Andrzej Kilian¹⁴, Dheema Burthia¹, Christophe Jenny¹, Frédéric Bakry¹, Spencer Brown¹⁵, Valentin Guignon^{1,6}, Gert Kema¹⁶, Miguel Dita¹⁹, Cees Waalwijk¹⁶, Steeve Joseph¹, Anne Dievart¹, Olivier Jaillon^{2,3,4}, Julie Leclercq¹, Xavier Argout¹, Eric Lyons¹⁷, Ana Almeida⁸, Mouna Jeridi¹, Jaroslav Dolezel¹², Nicolas Roux⁶, Ange-Marie Risterucci¹, Jean Weissenbach^{2,3,4}, Manuel Ruiz¹, Jean-Christophe Glaszmann¹, Francis Quétier¹⁸, Nabila Yahiaoui¹ & Patrick Wincker^{2,3,4}

Bananas (*Musa* spp.), including dessert and cooking types, are giant perennial monocotyledonous herbs of the order Zingiberales, a sister group to the well-studied Poales, which include cereals. Bananas are vital for food security in many tropical and subtropical countries and the most popular fruit in industrialized countries¹. The *Musa* domestication process started some 7,000 years ago in Southeast Asia. It involved hybridizations between diverse species and subspecies, fostered by human migrations², and selection of diploid and triploid seedless, parthenocarpic hybrids thereafter widely dispersed by vegetative propagation. Half of the current production relies on somaclones derived from a single triploid genotype (Cavendish)¹. Pests and diseases have gradually become adapted, representing an imminent danger for global banana production^{3,4}. Here we describe the draft sequence of the 523-megabase genome of a *Musa acuminata* doubled-haploid genotype, providing a crucial stepping-stone for genetic improvement of banana. We detected three rounds of whole-genome duplications in the *Musa* lineage, independently of those previously described in the Poales lineage and the one we detected in the Arecales lineage. This first monocotyledon high-continuity whole-genome sequence reported outside Poales represents an essential bridge for comparative genome analysis in plants. As such, it clarifies commelinid-

sequence errors. The assembly consisted of 24,425 contigs and 7,513 scaffolds with a total length of 472.2 Mb, which represented 90% of the estimated DH-Pahang genome size. Ninety per cent of the assembly was in 647 scaffolds, and the N50 (the scaffold size above which 50% of the total length of the sequence assembly can be found) was 1.3 Mb (Supplementary Text and Supplementary Tables 1–3). We anchored 70% of the assembly (332 Mb) along the 11 *Musa* linkage groups of the Pahang genetic map. This corresponded to 258 scaffolds and included 98.0% of the scaffolds larger than 1 Mb and 92% of the annotated genes (Supplementary Text, Supplementary Table 4 and Supplementary Fig. 1).

We identified 36,542 protein-coding gene models in the *Musa* genome (Supplementary Tables 1 and 5). A total of 235 microRNAs from 37 families were identified, including only one of the eight microRNA gene (*MIR*) families found so far solely in Poaceae⁸ (Supplementary Tables 6 and 7).

Viral sequences related to the banana streak virus (BSV) dsDNA plant pararetrovirus were found to be integrated in the Pahang genome, with 24 loci spanning 10 chromosomes (Supplementary Text and Supplementary Fig. 2). They belonged to a badnavirus phylogenetic group that differed from the endogenous BSV species (eBSV) found in *M. balbisiana*⁹ and most of them formed a new

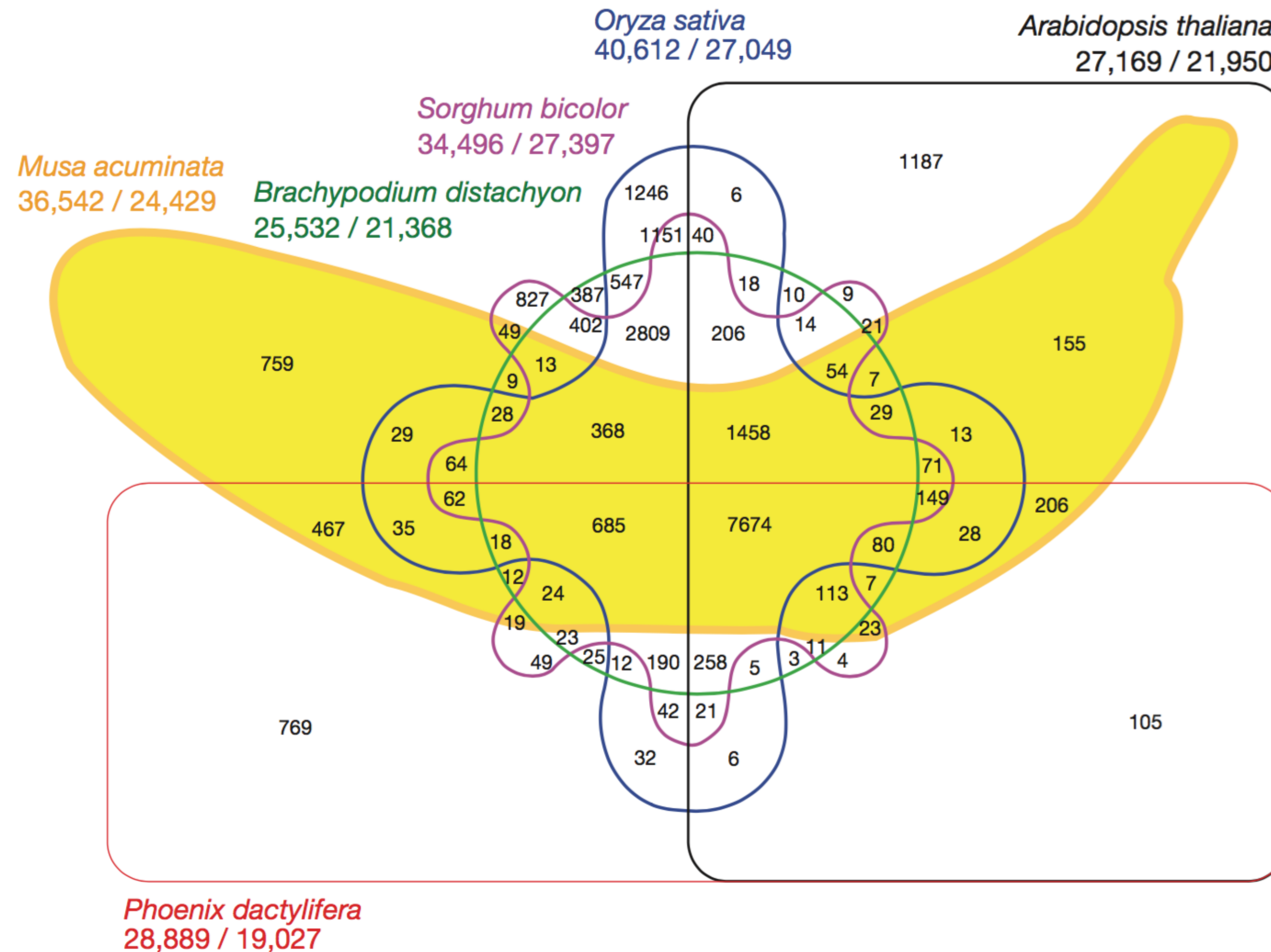
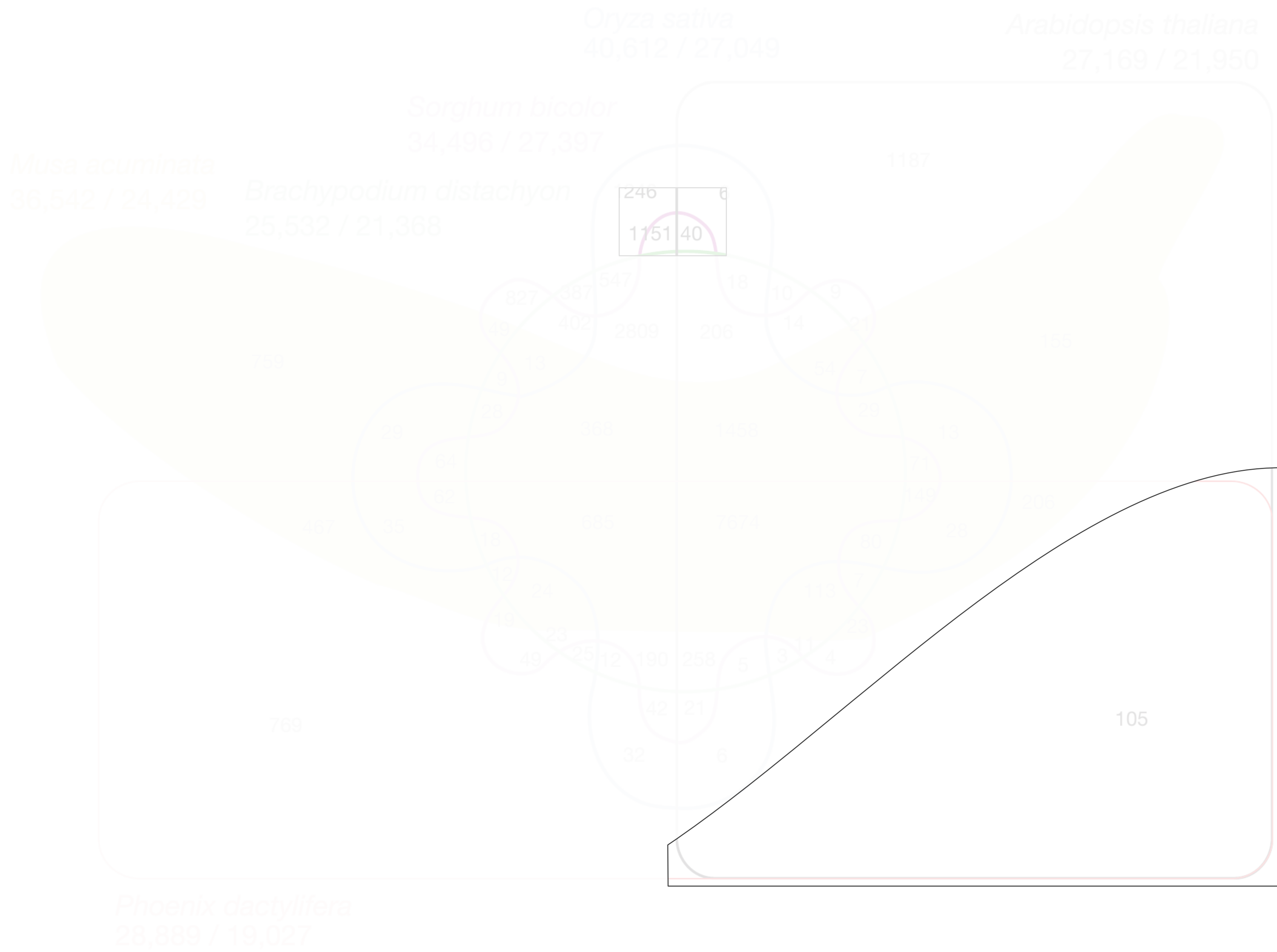


Figure 4 | Six-way Venn diagram showing the distribution of shared gene families (sequence clusters) among *M. acuminata*, *P. dactylifera*, *Arabidopsis thaliana*, *Oryza sativa*, *Sorghum bicolor* and *Brachypodium distachyon* genomes. Numbers of clusters are provided in the intersections. The total number of sequences for each species is provided under the species name (total number of sequences/total number of clustered sequences).



A

Dicots

Arabidopsis thaliana: 26304 / 24766
Glycine max: 36271 / 35969
Populus trichocarpa: 35516 / 33358
Ricinus communis: 30314 / 24039
Theobroma cacao: 28222 / 27154
Vitis vinifera: 24479 / 21795

Basal

Amborella trichopoda: 24611 / 21191

Early land plants

Selaginella moellendorffii:
 16832 / 15909
Physcomitrella patens:
 25938 / 19359

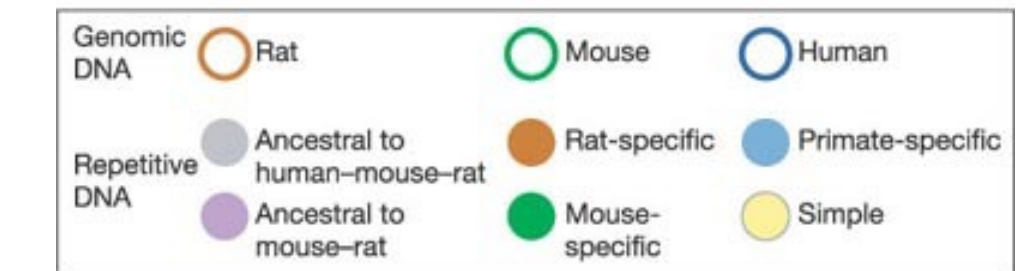
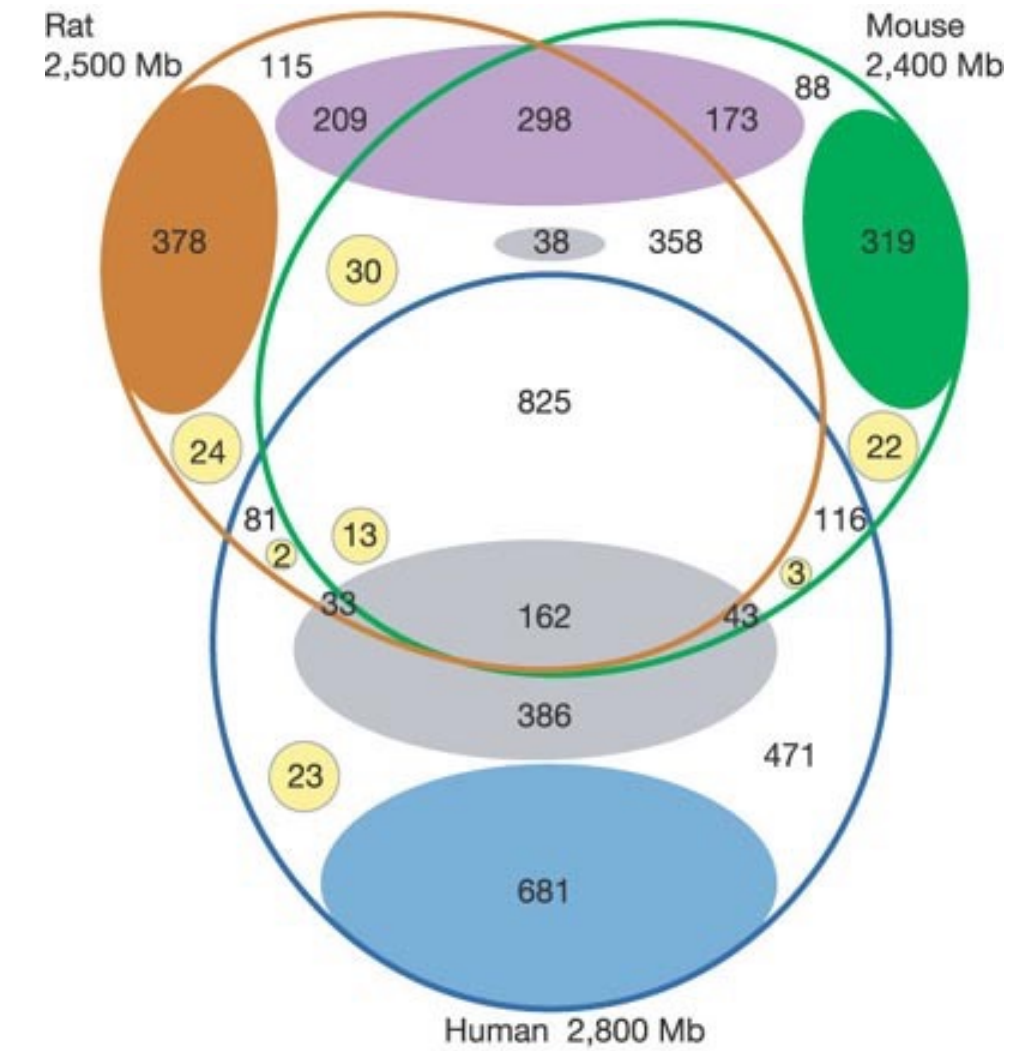
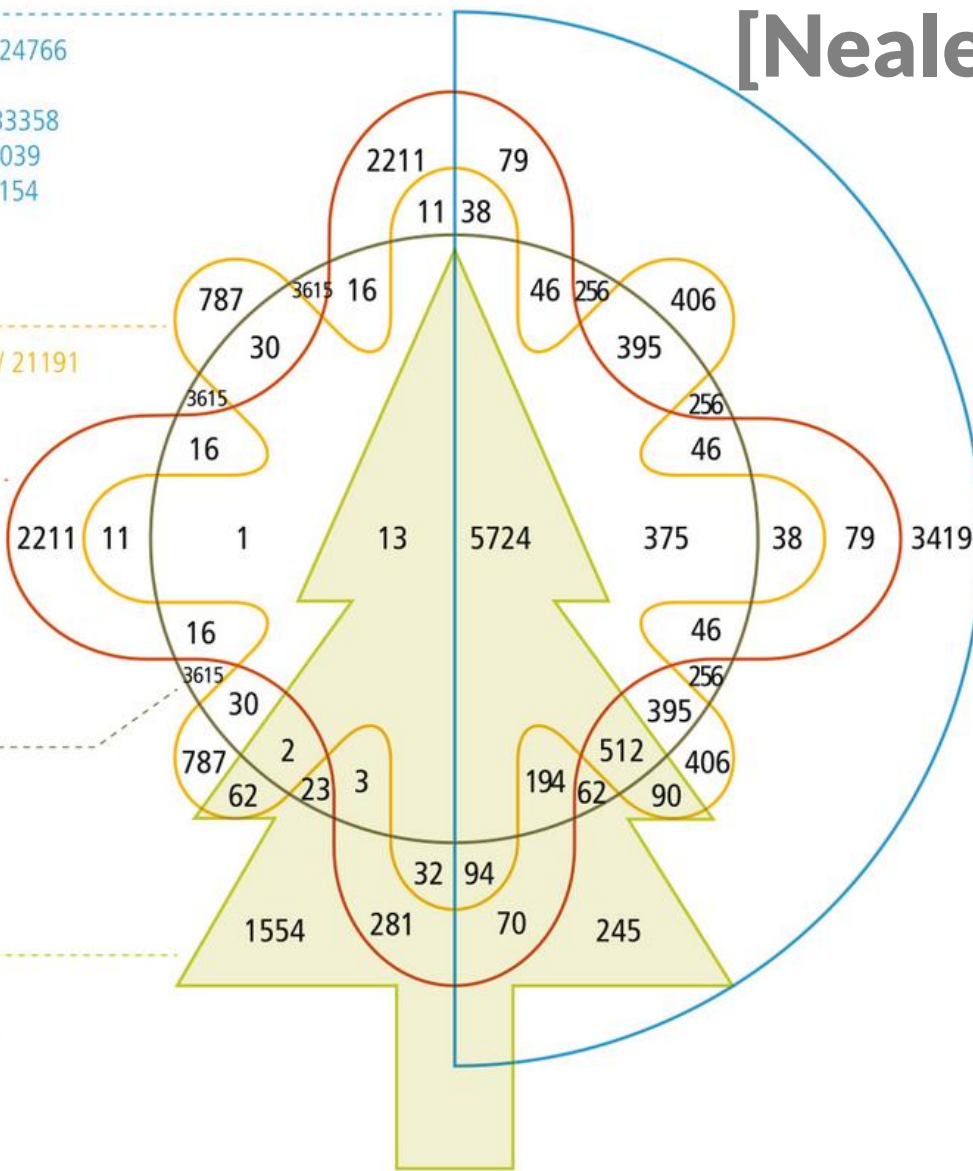
Monocots

Oryza sativa: 39459 / 32660
Zea mays: 34586 / 30799

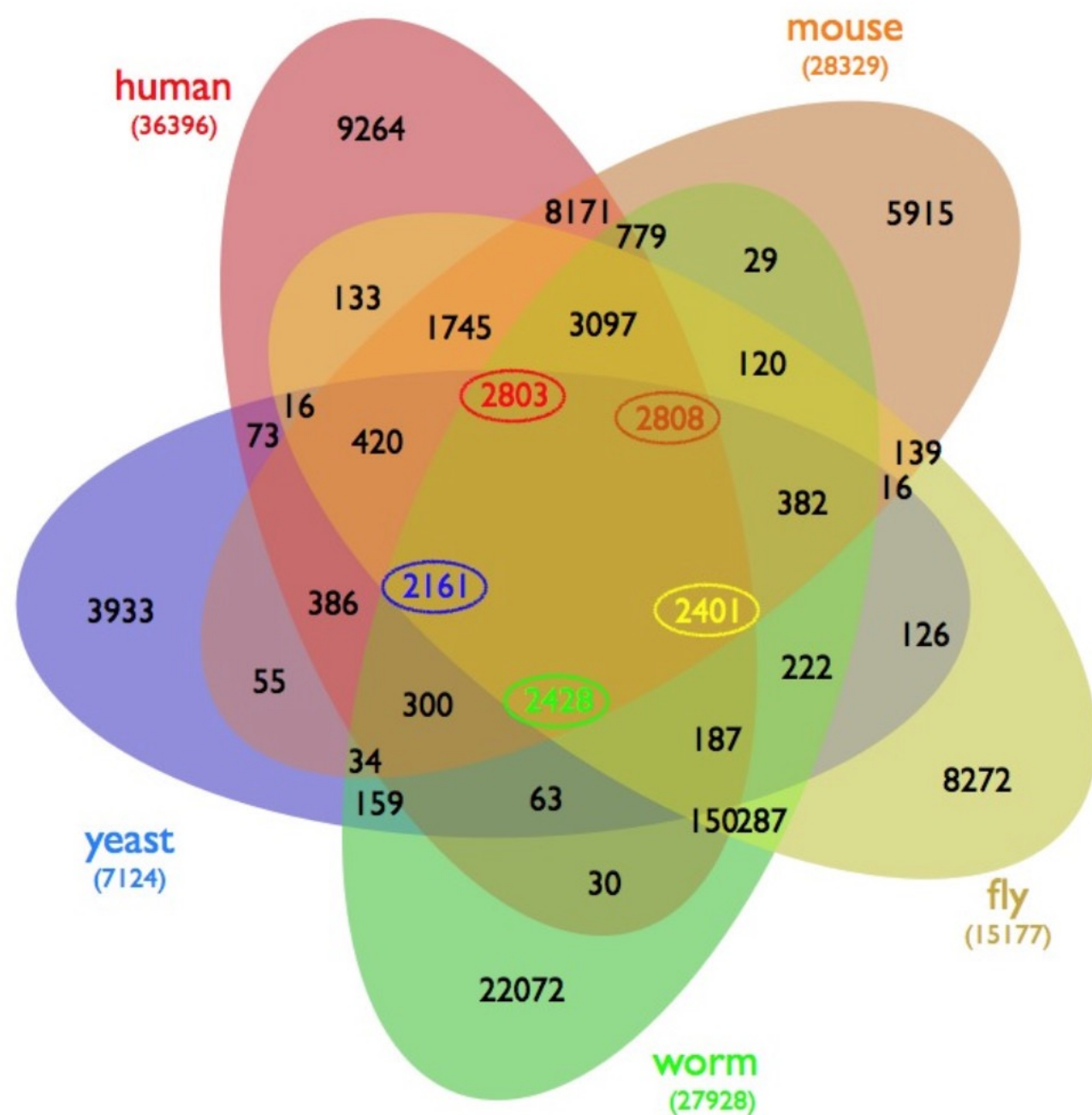
Conifers

Picea abies: 20861 / 19934
Picea sitchensis: 8758 / 7780
Pinus taeda: 47207 / 46720

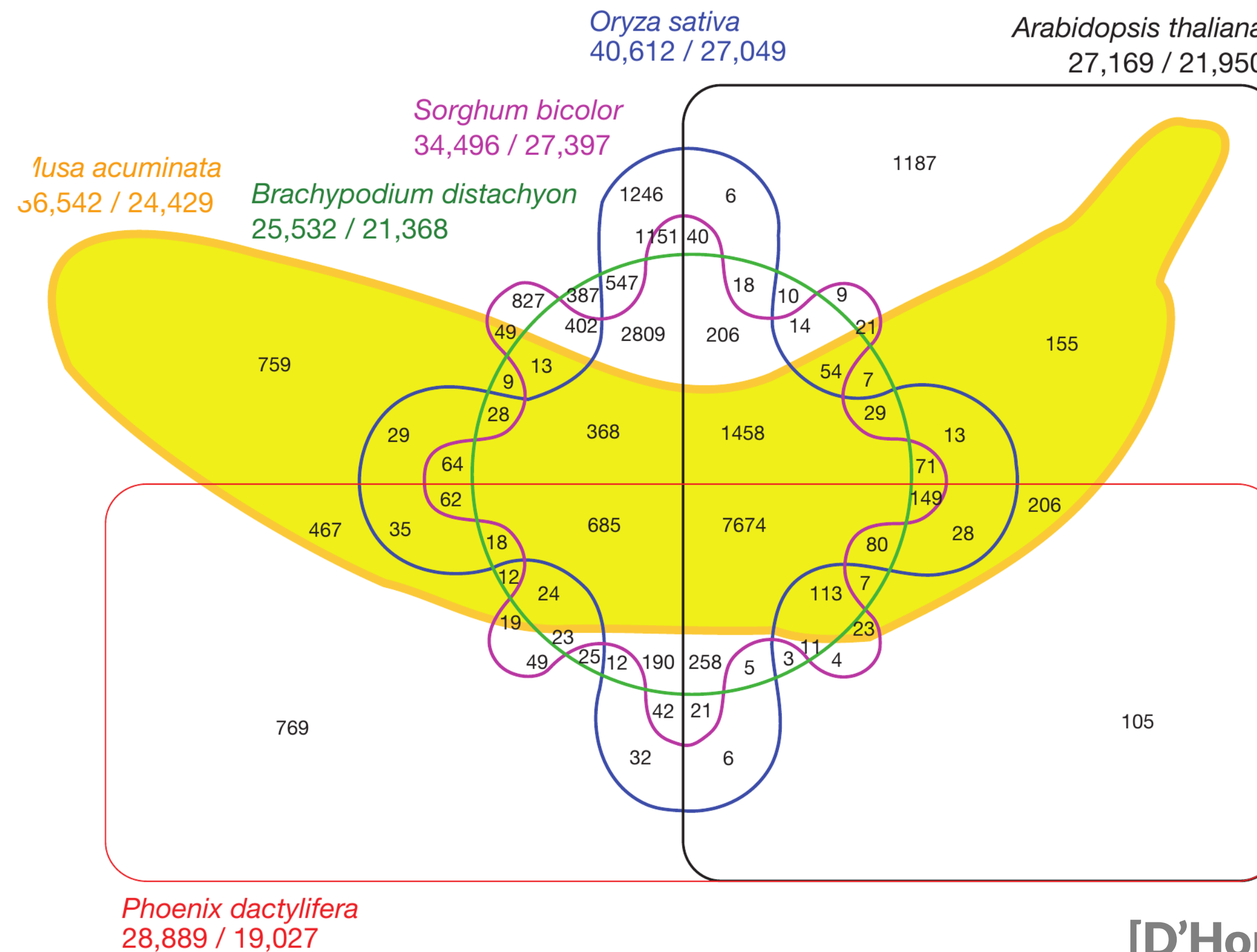
[Neale et al., BMC Genome Biology, 2014]



[Gibbs et al., Nature, 2004]



[Wiles et al., BMC Systems Biology]



[D'Hont et al., Nature, 2012]

Element ID	Sets	Attribute(s)
Name	Characteristics	Age
Lisa	School, Female	8
Bart	School, Male	10
Homer	Power Plant, Male	40
Mr. Burns	Evil, Power Plant, Male	90

What are some questions we'd like to ask?

Design Workshop

work in groups

get to know the data (5 mins)

create two (rapid!) prototypes (2x5 mins)

Write up your two favorites (5 mins)

Upload to “Bonus” Canvas Dropbox by EOD

Element ID	Sets	Attribute(s)
Name	Characteristics	Age
Lisa	School, Female	8
Bart	School, Male	10
Homer	Power Plant, Male	40
Mr. Burns	Evil, Power Plant, Male	90

1. What is the biggest intersection?
2. Which sets make up an intersection?
3. How big is an intersection?
4. Does it work for more than four sets?
5. Does attribute value correlate with intersection

Tip: Don't always try to show all individuals

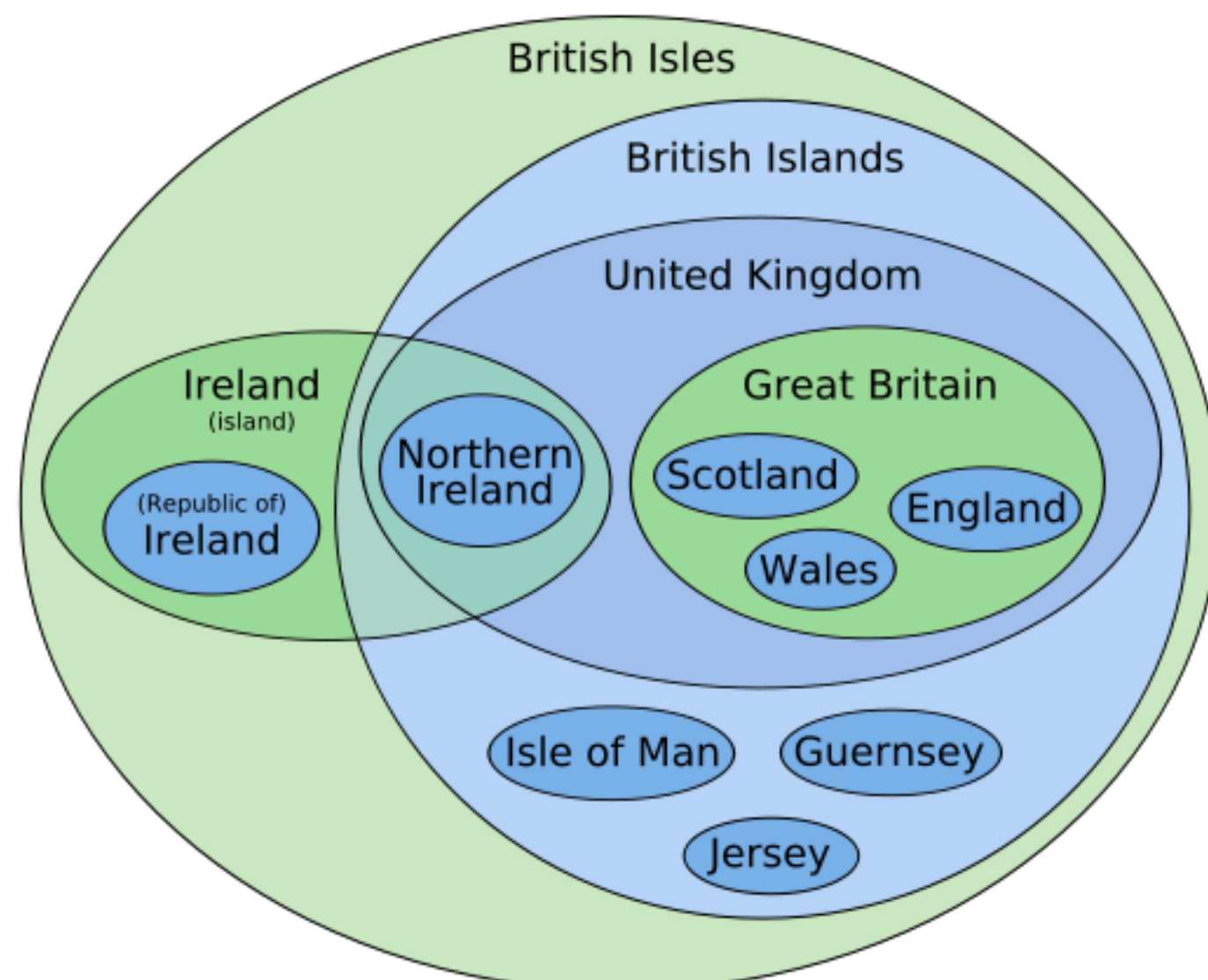
Venn and Euler Diagrams

Venn vs Euler

Euler Diagram

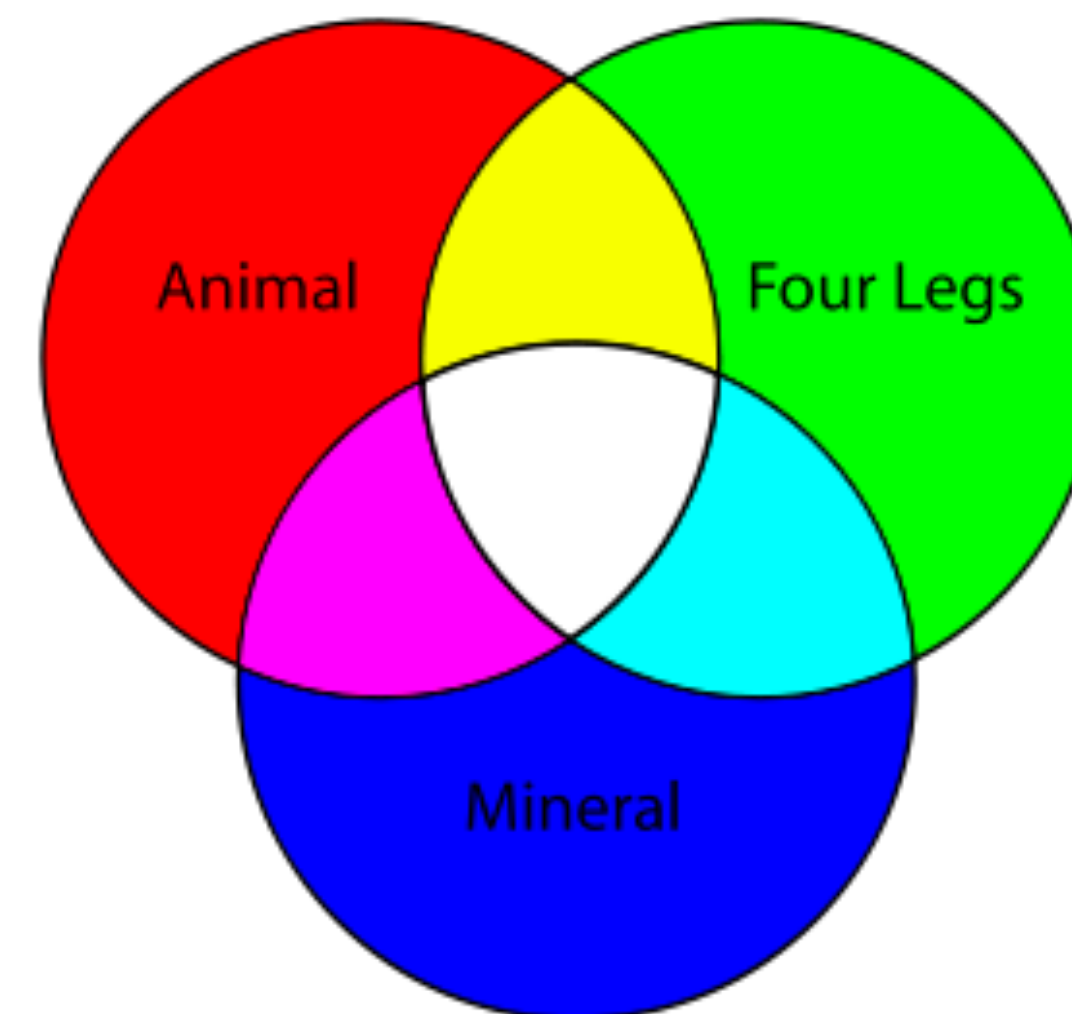
Shows logical relations

May omit empty intersections



Venn Diagram

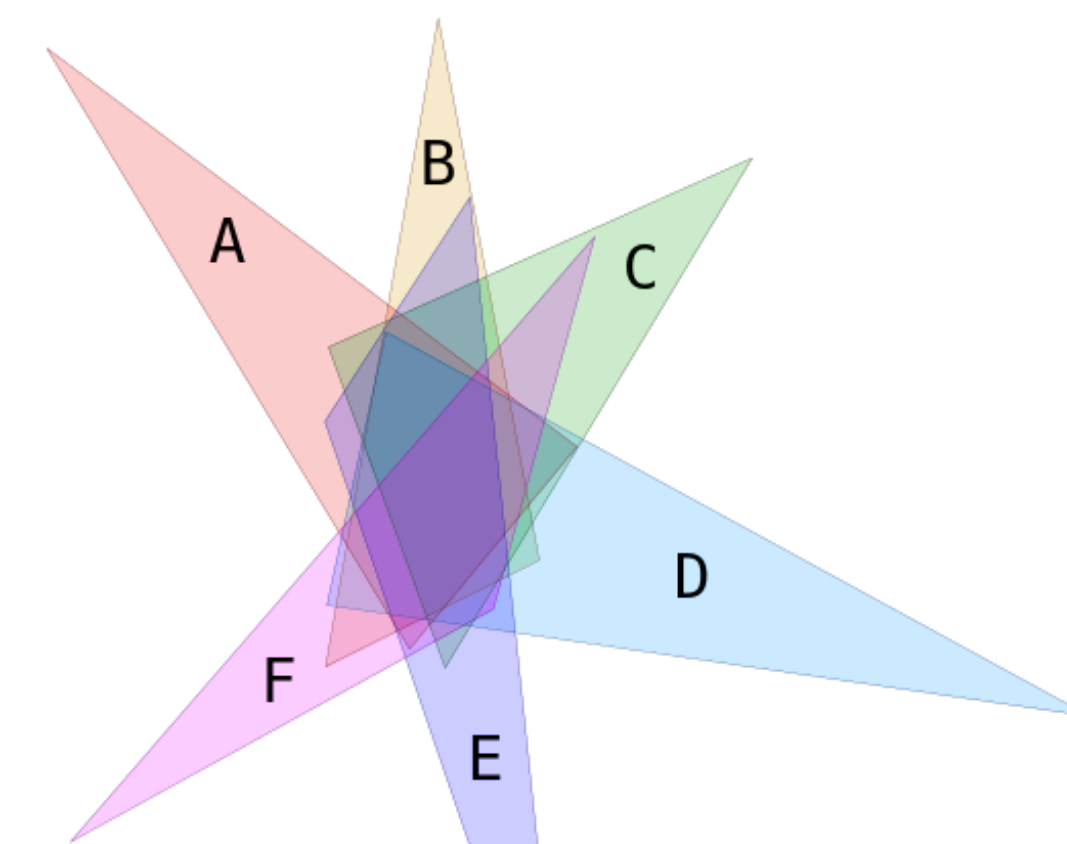
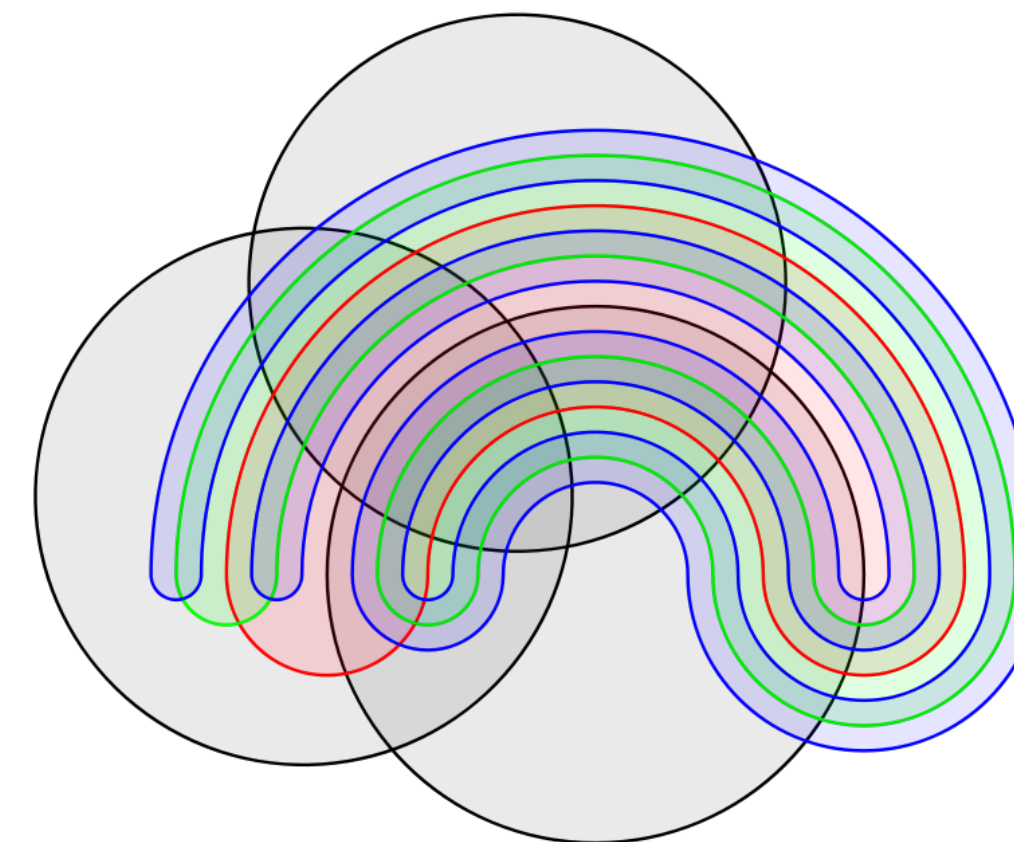
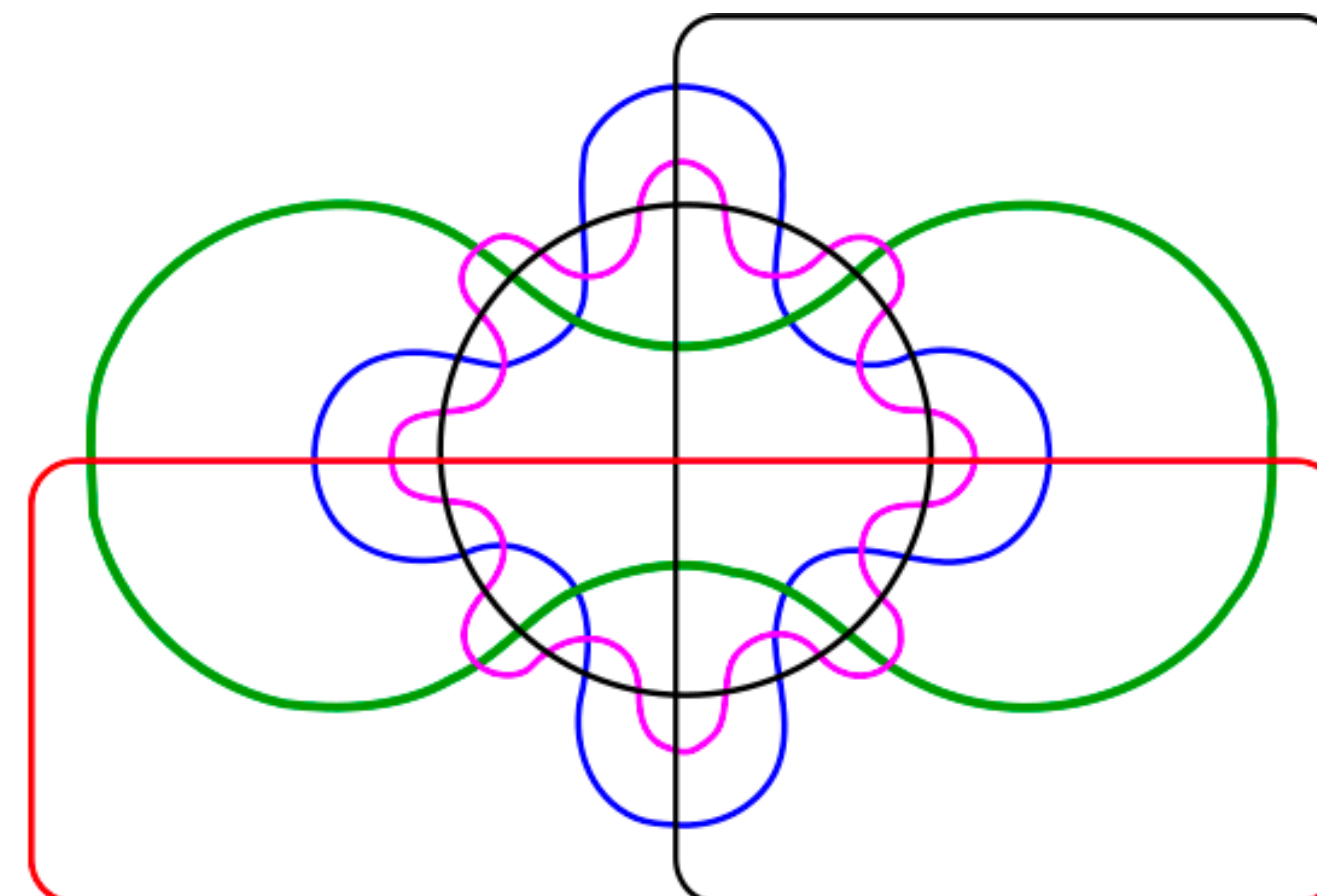
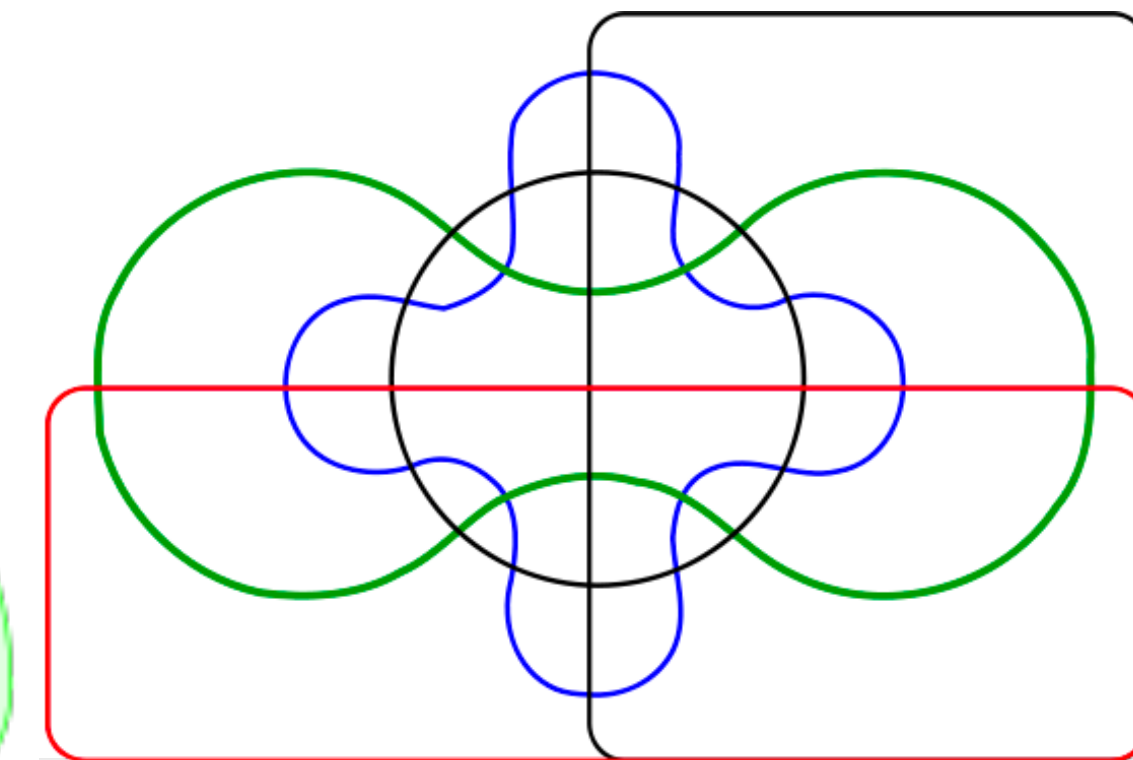
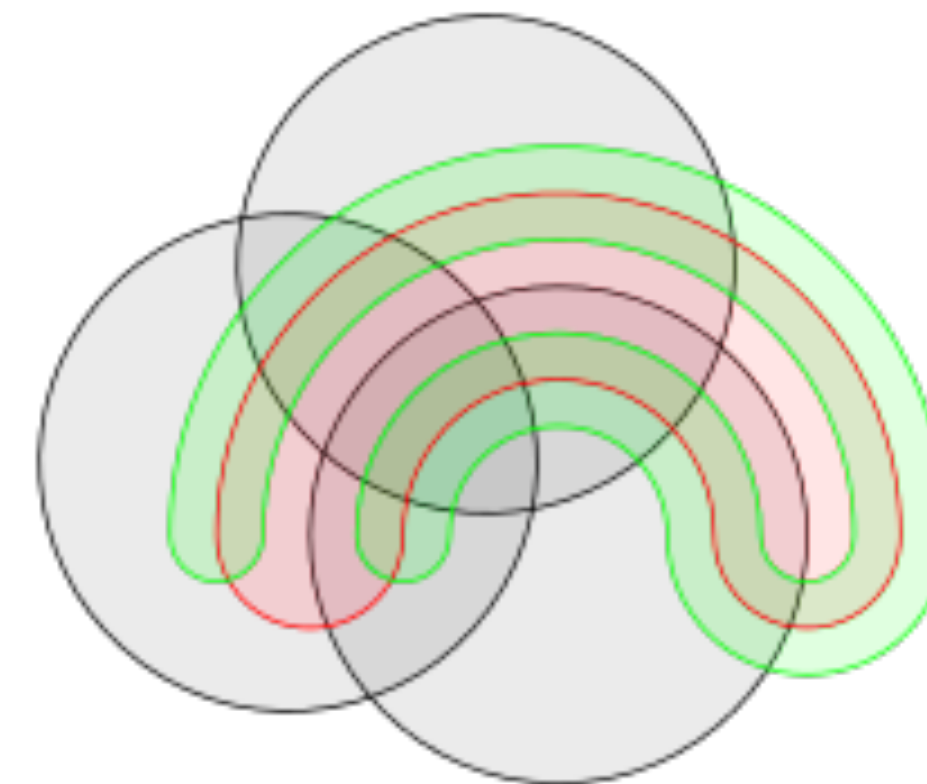
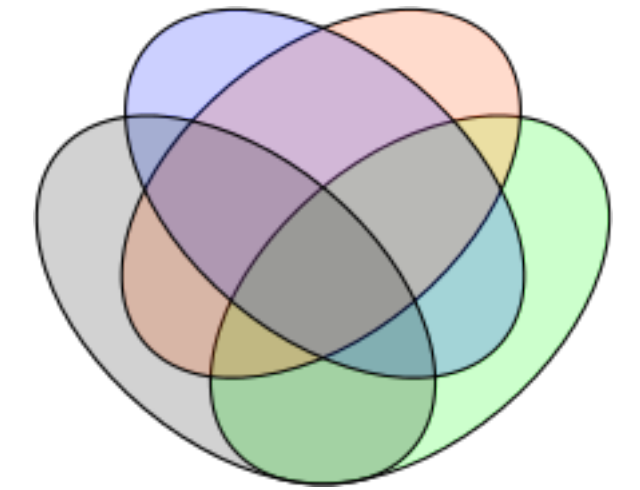
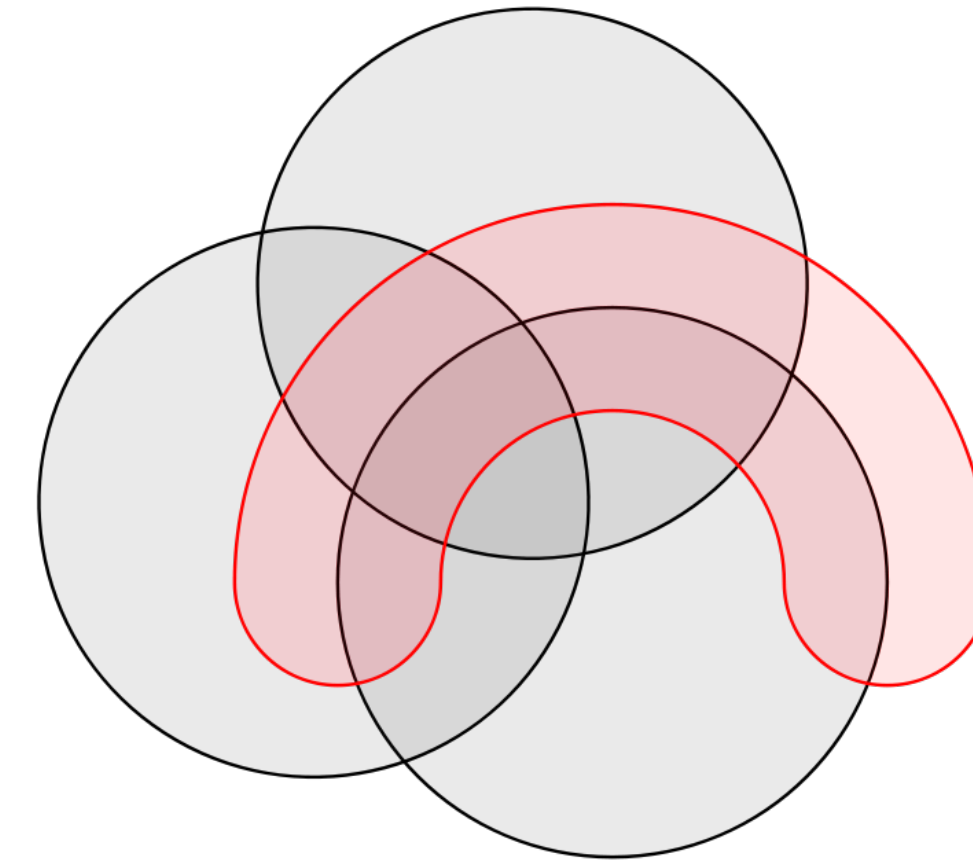
Shows all possible logical relations between sets (even if empty)



Venn Diagrams

Venn diagrams for many sets are hard

of intersections is 2^n



Area-Proportional Euler Diagrams

Problem with Venn: size doesn't correspond to the data.

Creating area-proportional Euler diagrams is hard.

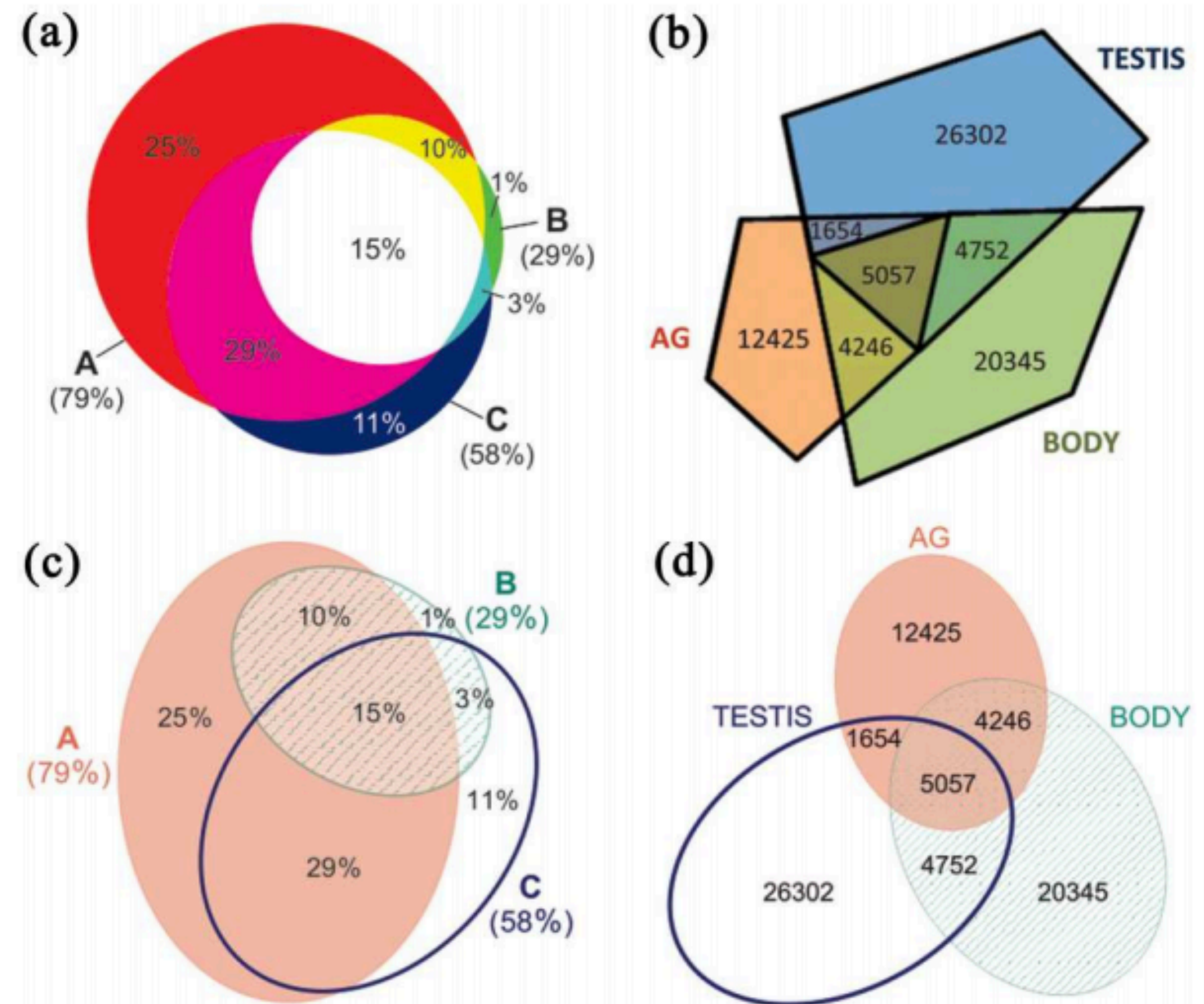
Layout criteria:

- area proportional

- simple curves (circles are best)

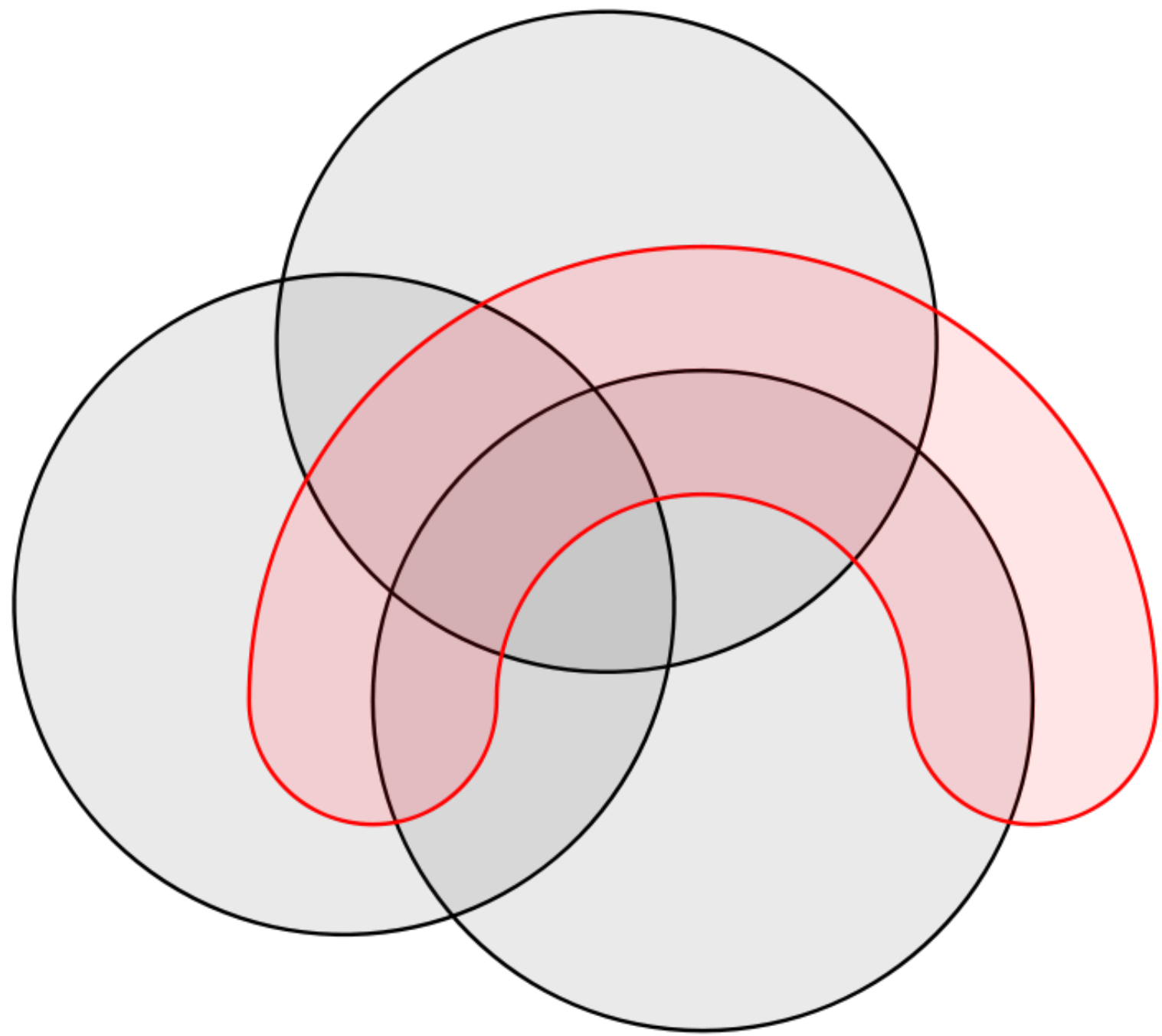
- makes it easy to identify which sets are participating in intersection

- Gestalt-principle: good continuation

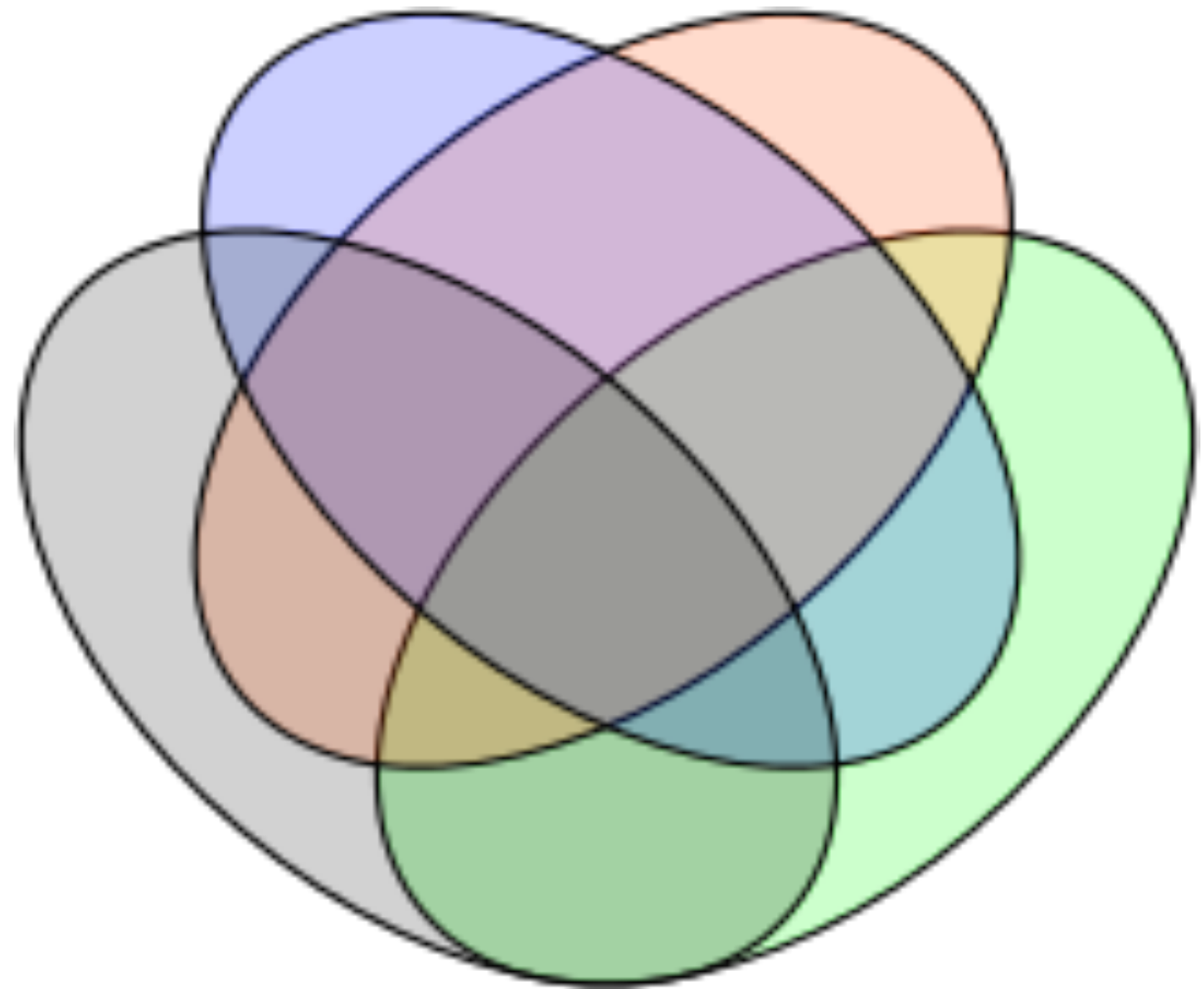


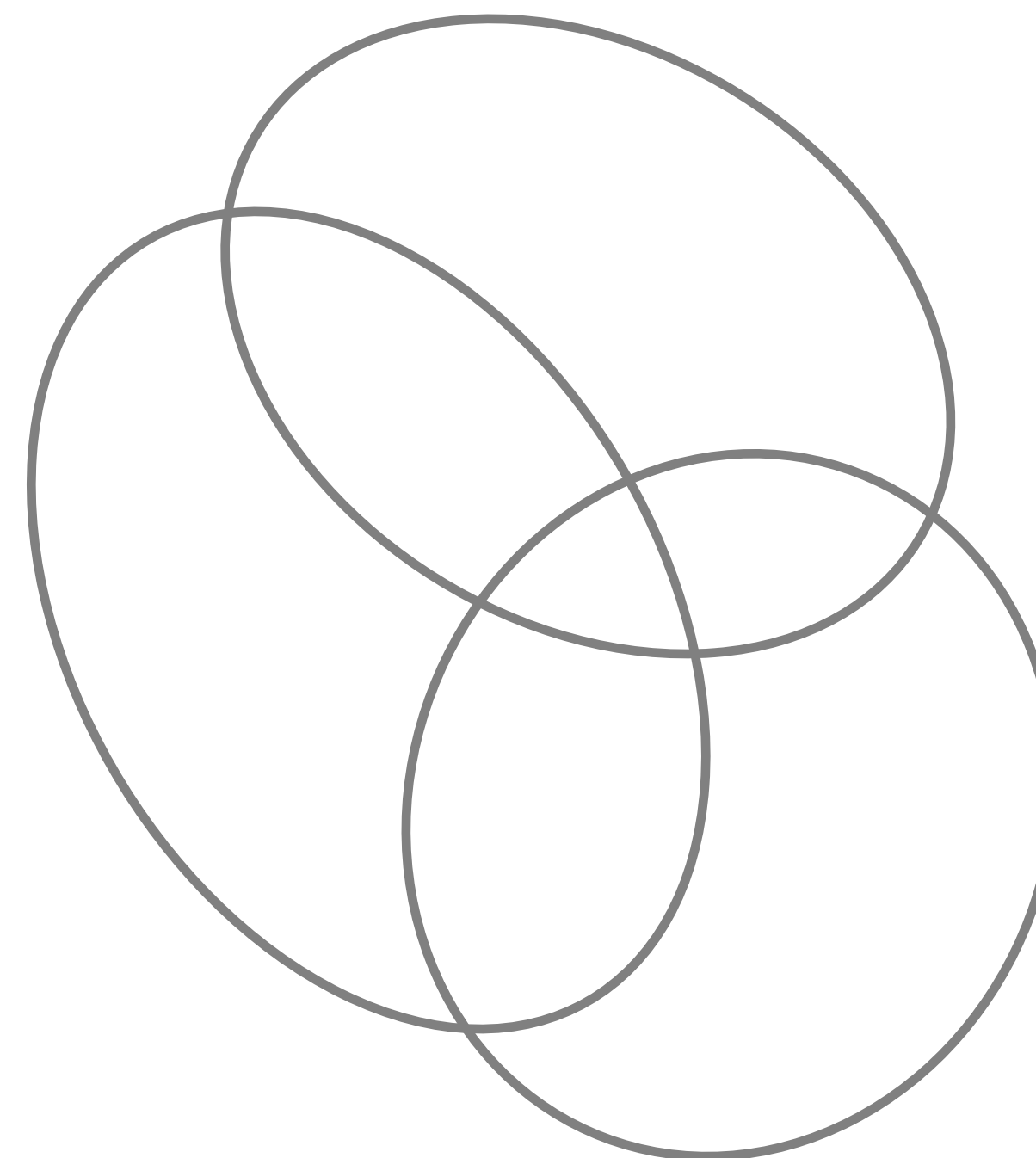
Compare Simple vs Complex Shape

Complex



Simple

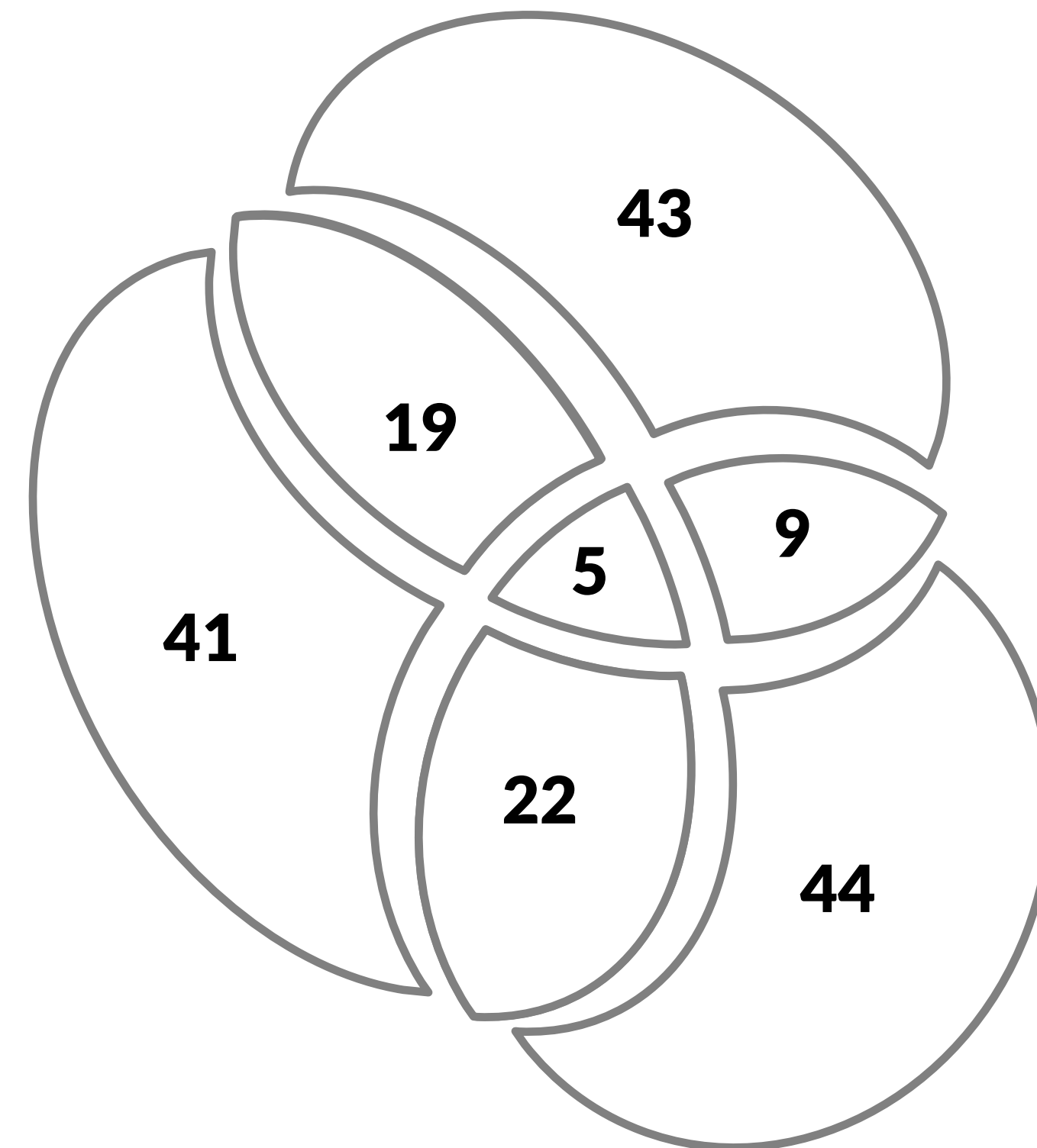




22

?

19



Venn-Euler Pros/Cons

Pros

Familiar

Intuitive

Work well for 2-4 sets

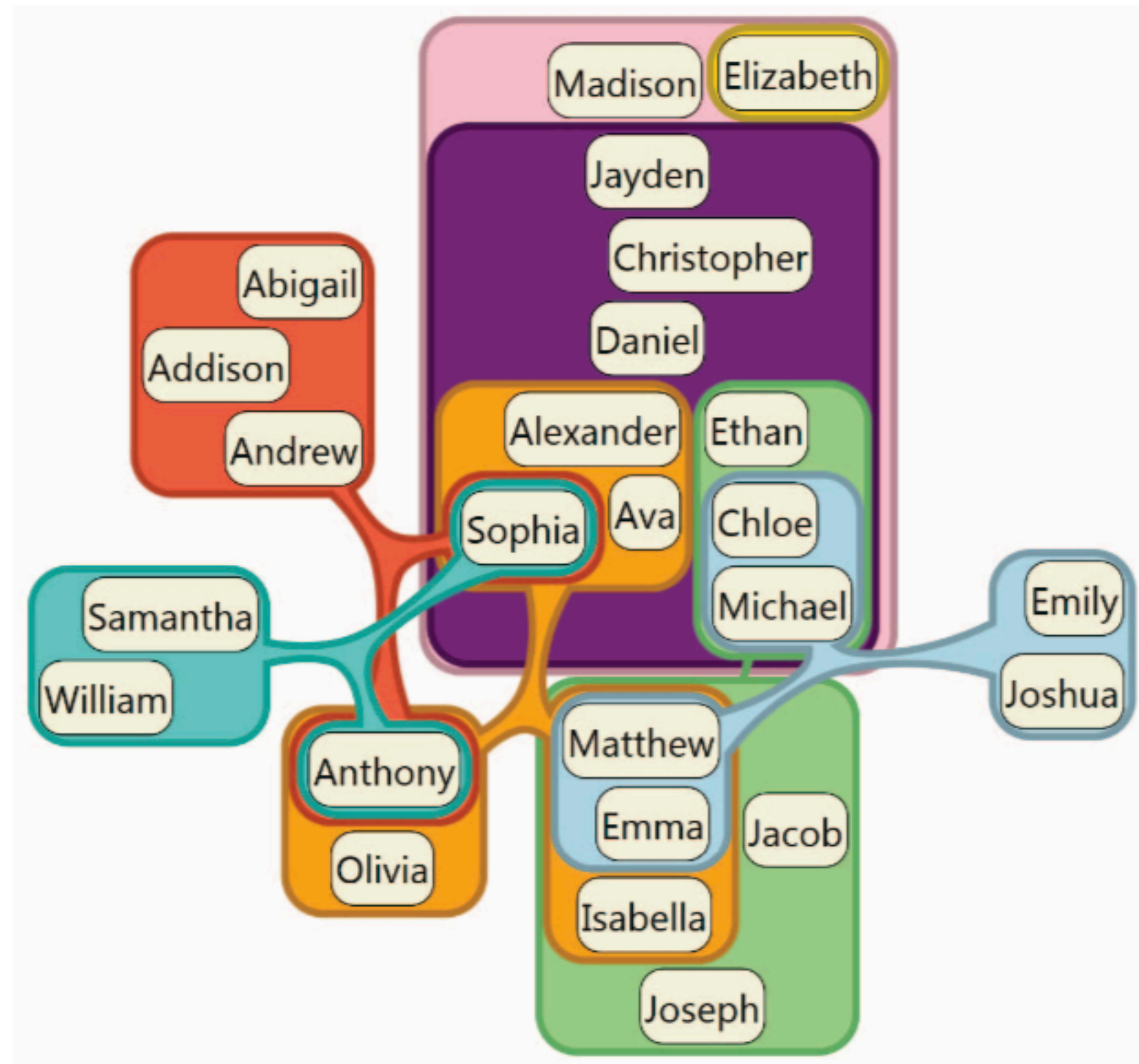
Cons

Doesn't work well for more than 4 sets

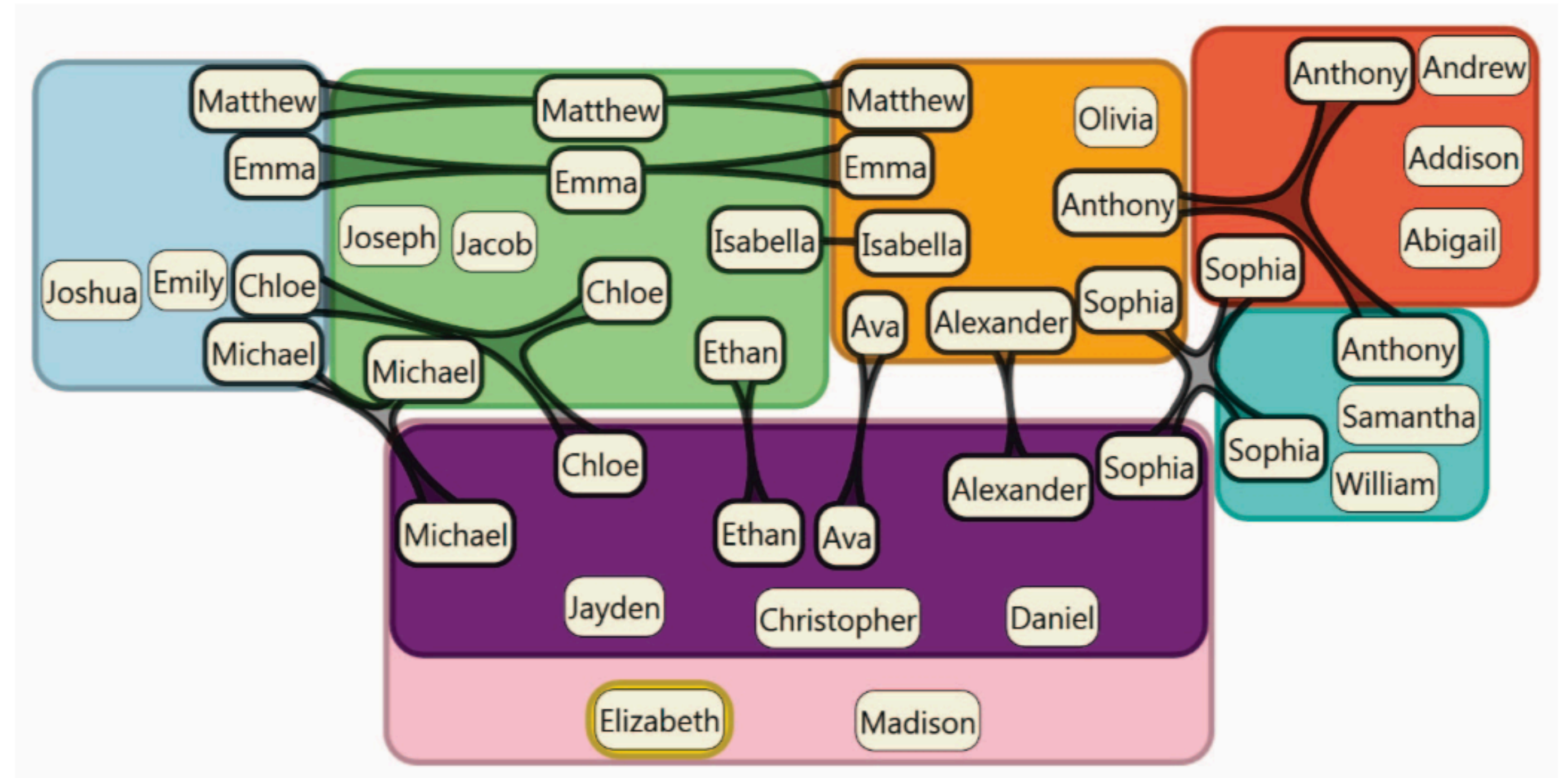
Area proportionality hard to do

Not well suited to show attributes

Relationships for specific Items



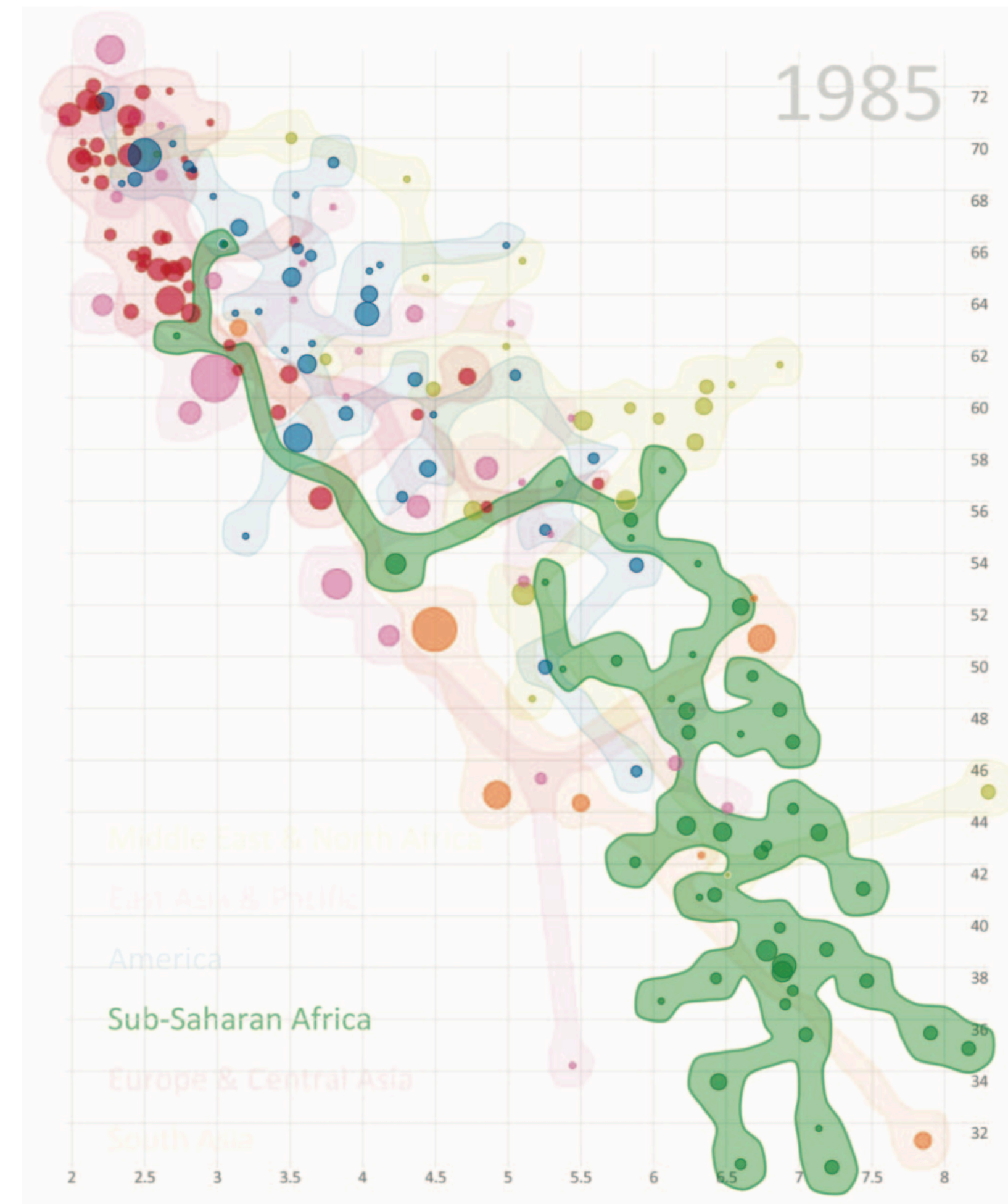
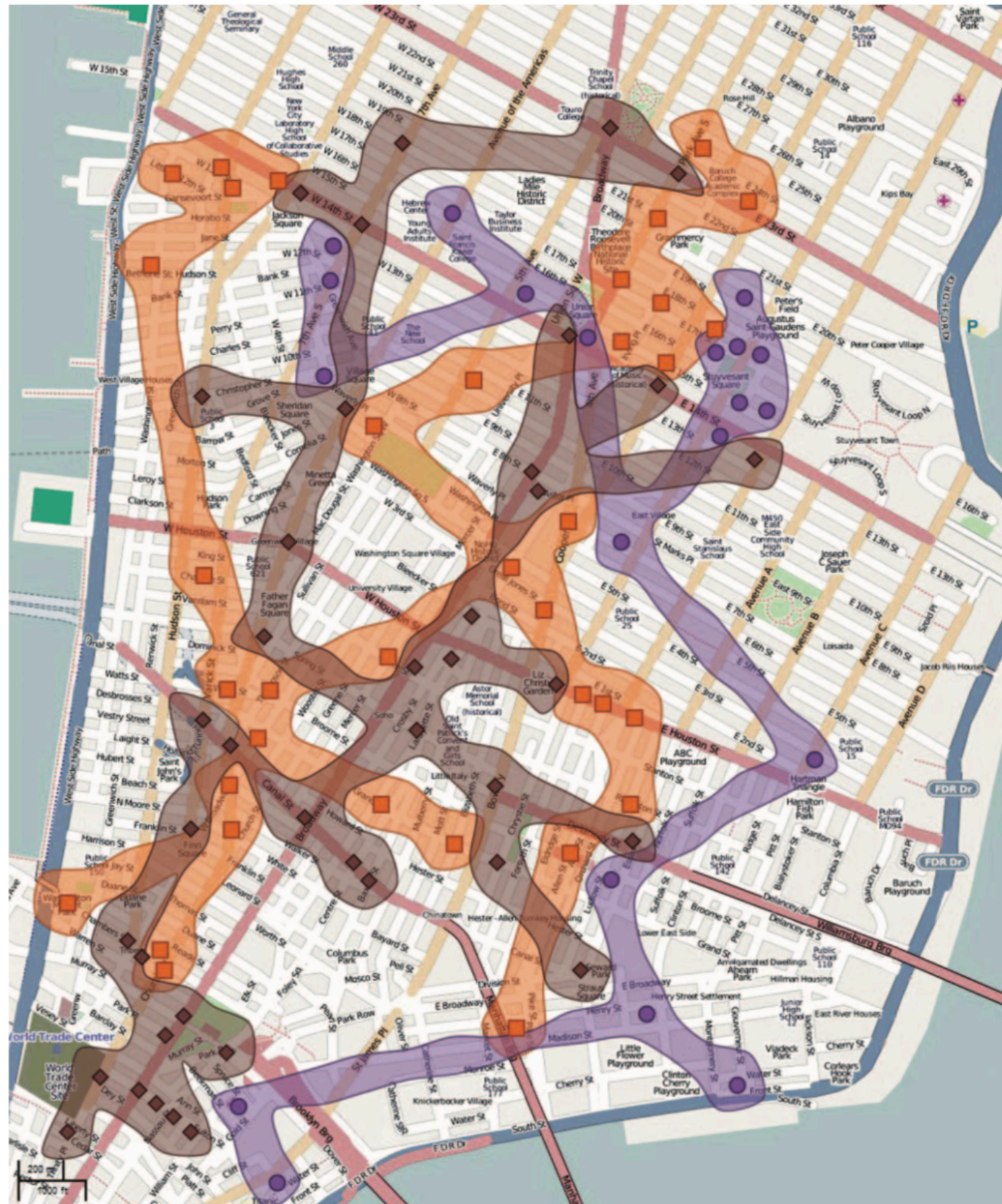
No Duplicate Nodes
Complex Shapes
Notice the Nesting



Duplicate Nodes
Simple Shapes

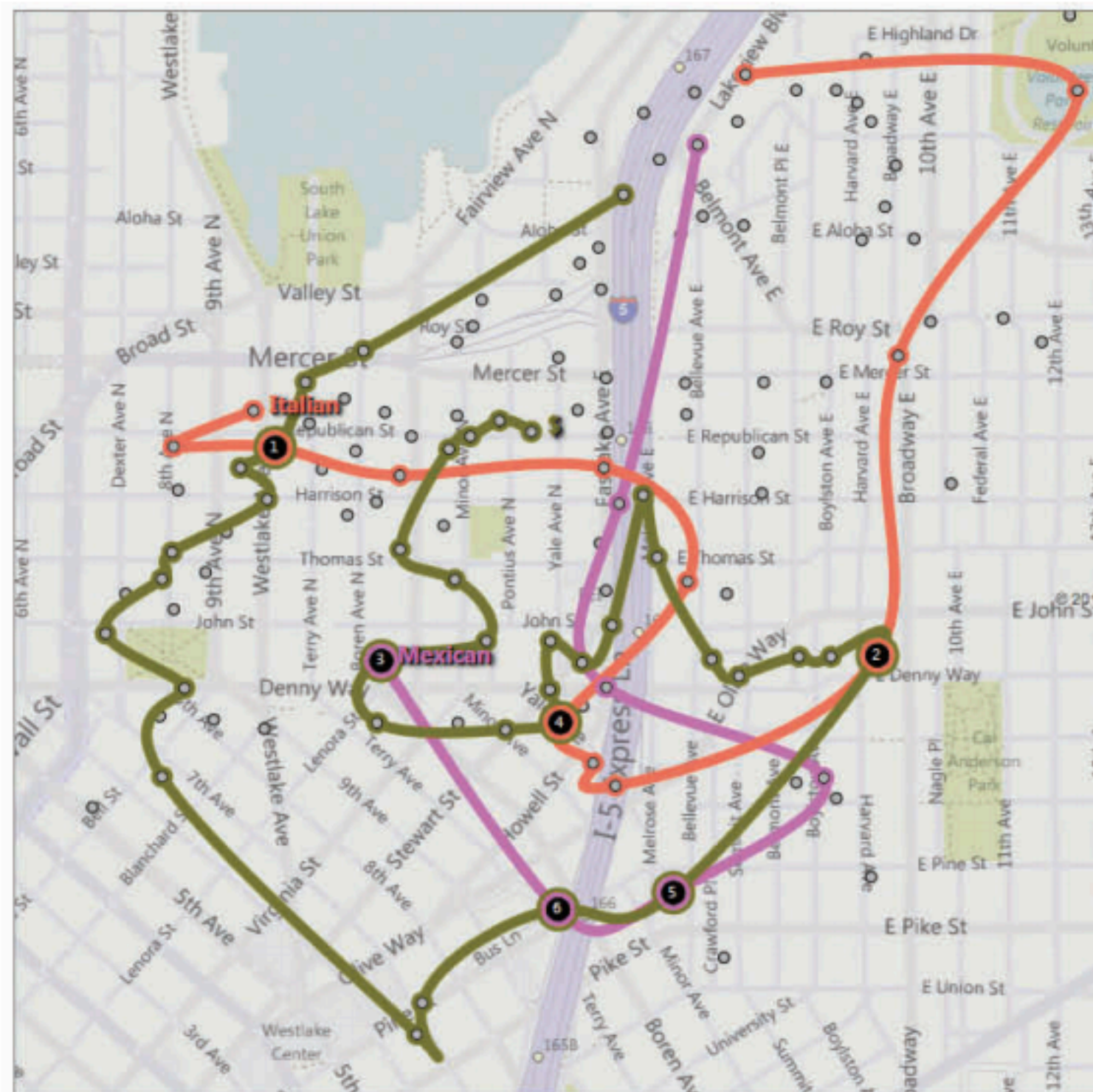
[Riche 2010]

Sets on top of a fixed layout



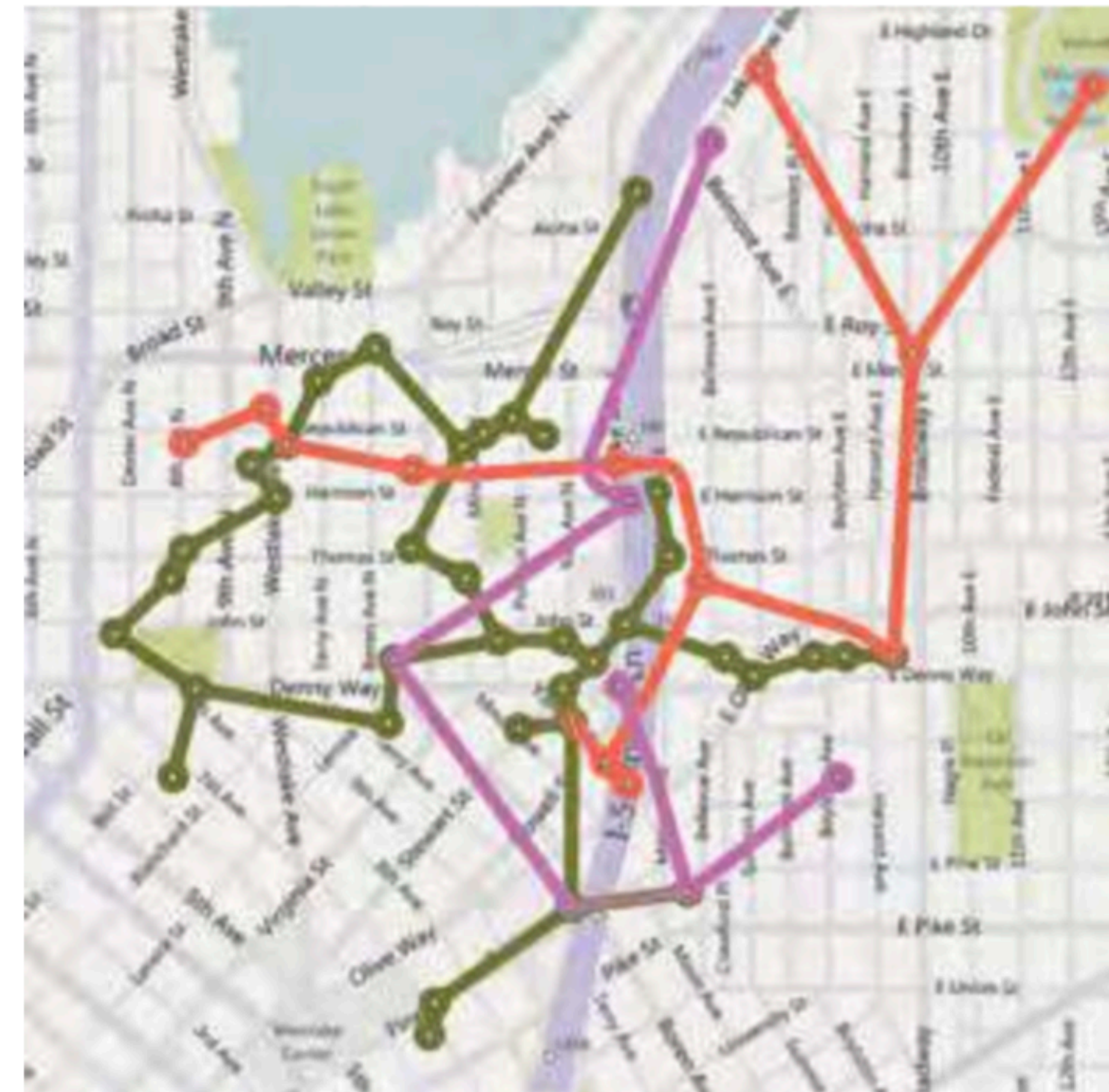
Sets on top of a fixed layout

LineSets



[Alper 2011]

Kelp Diagrams

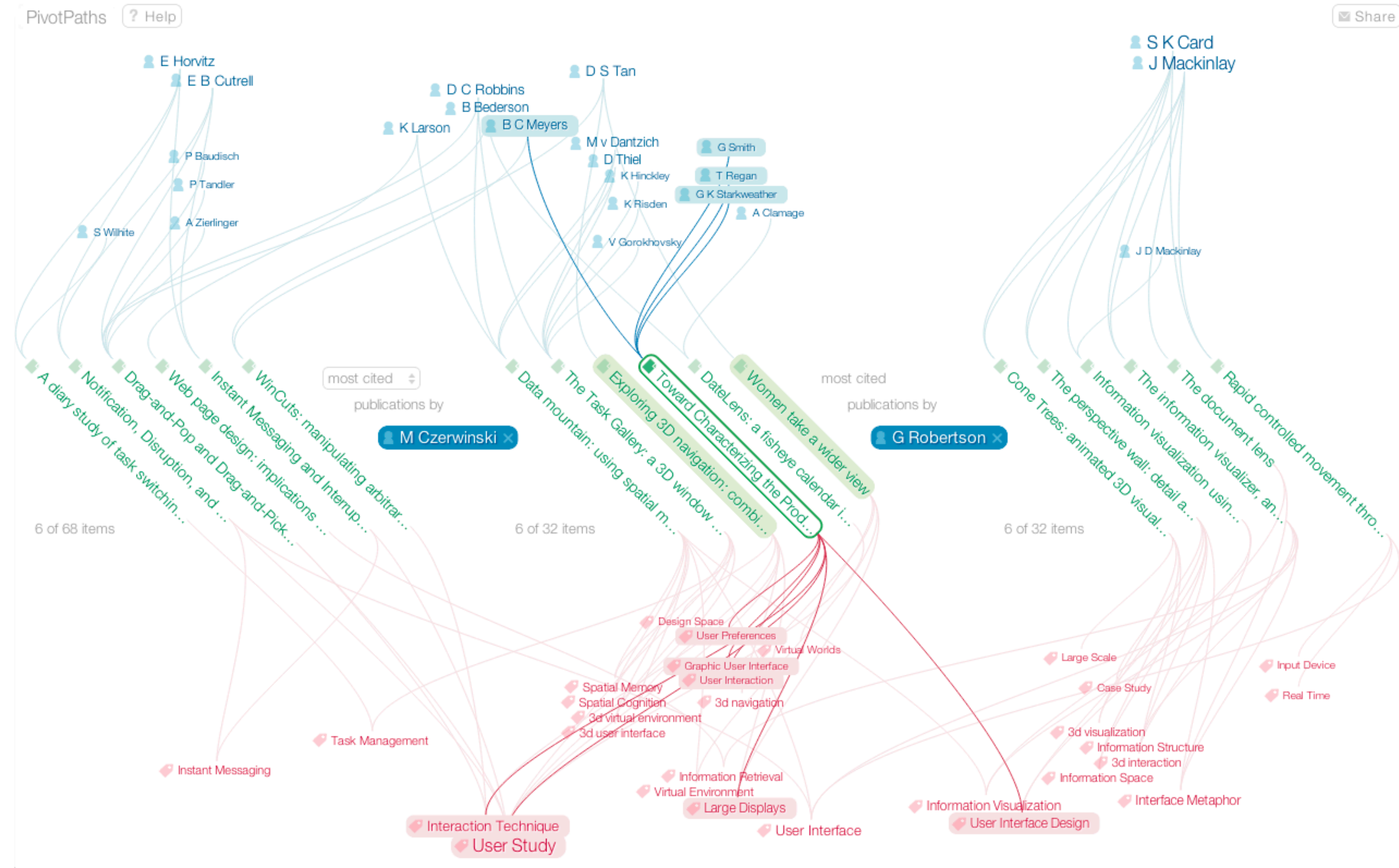


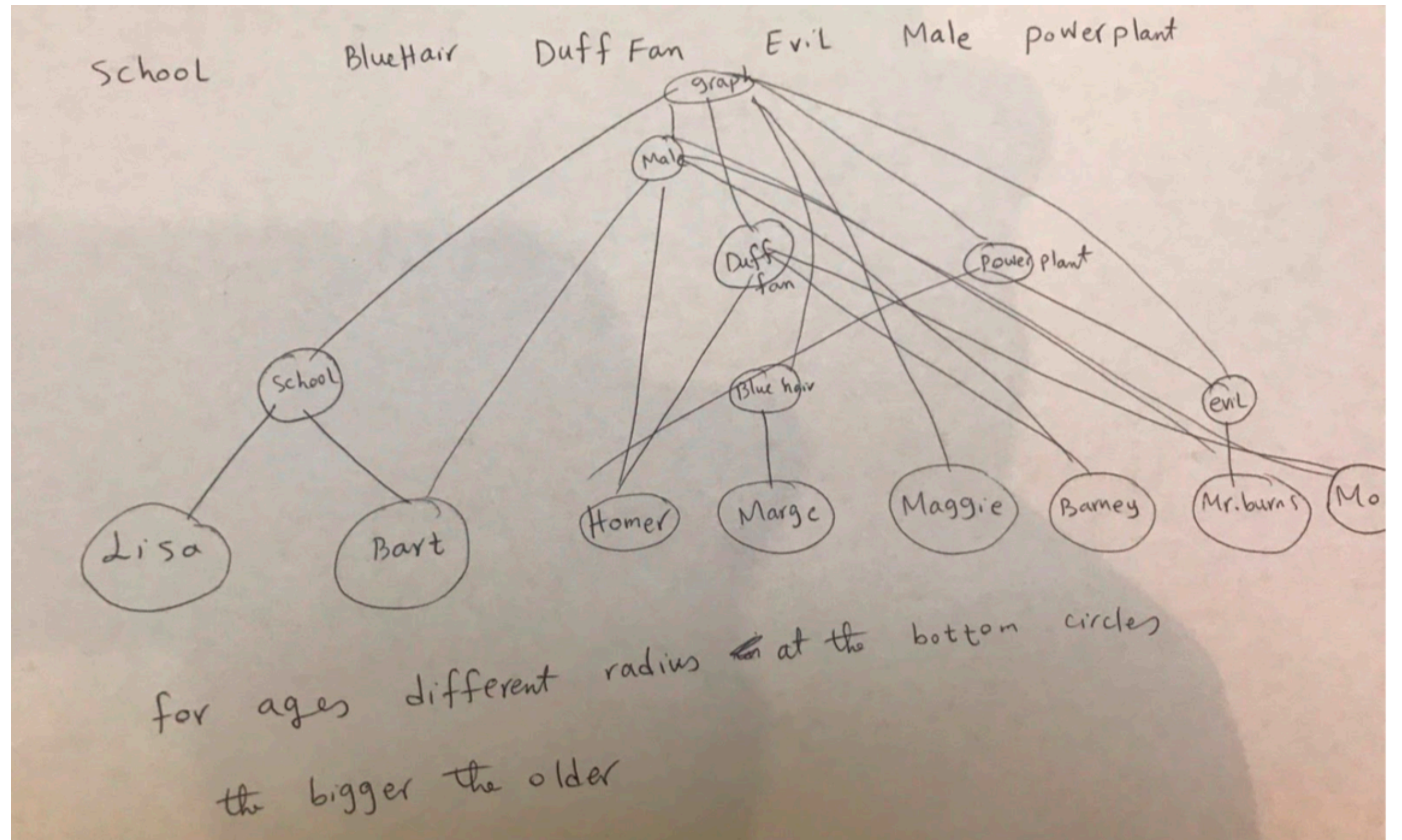
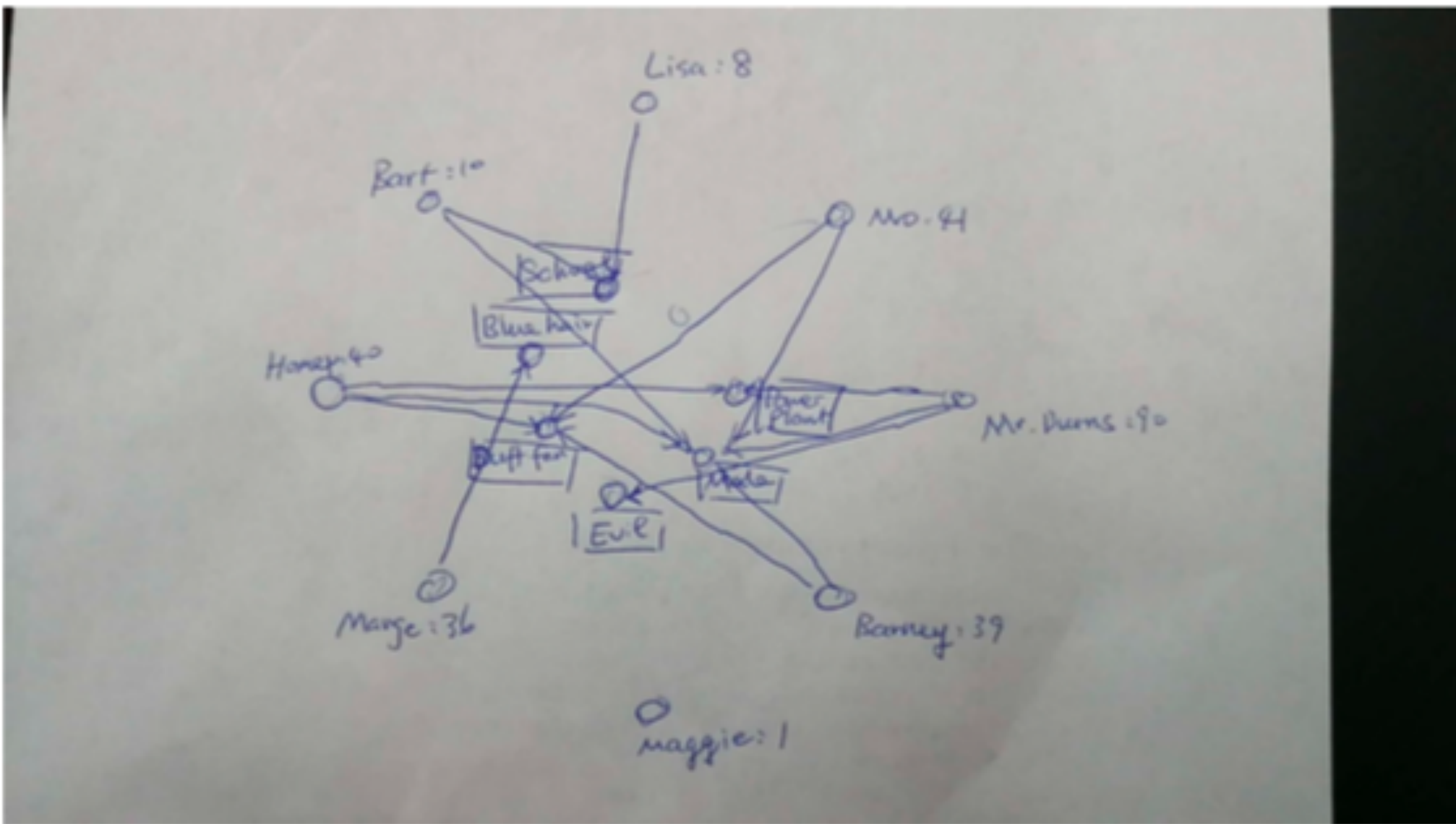
[Dinkla 2012]

Node-Link Techniques

Treat sets as
nodes

Connect to
elements that
are in set





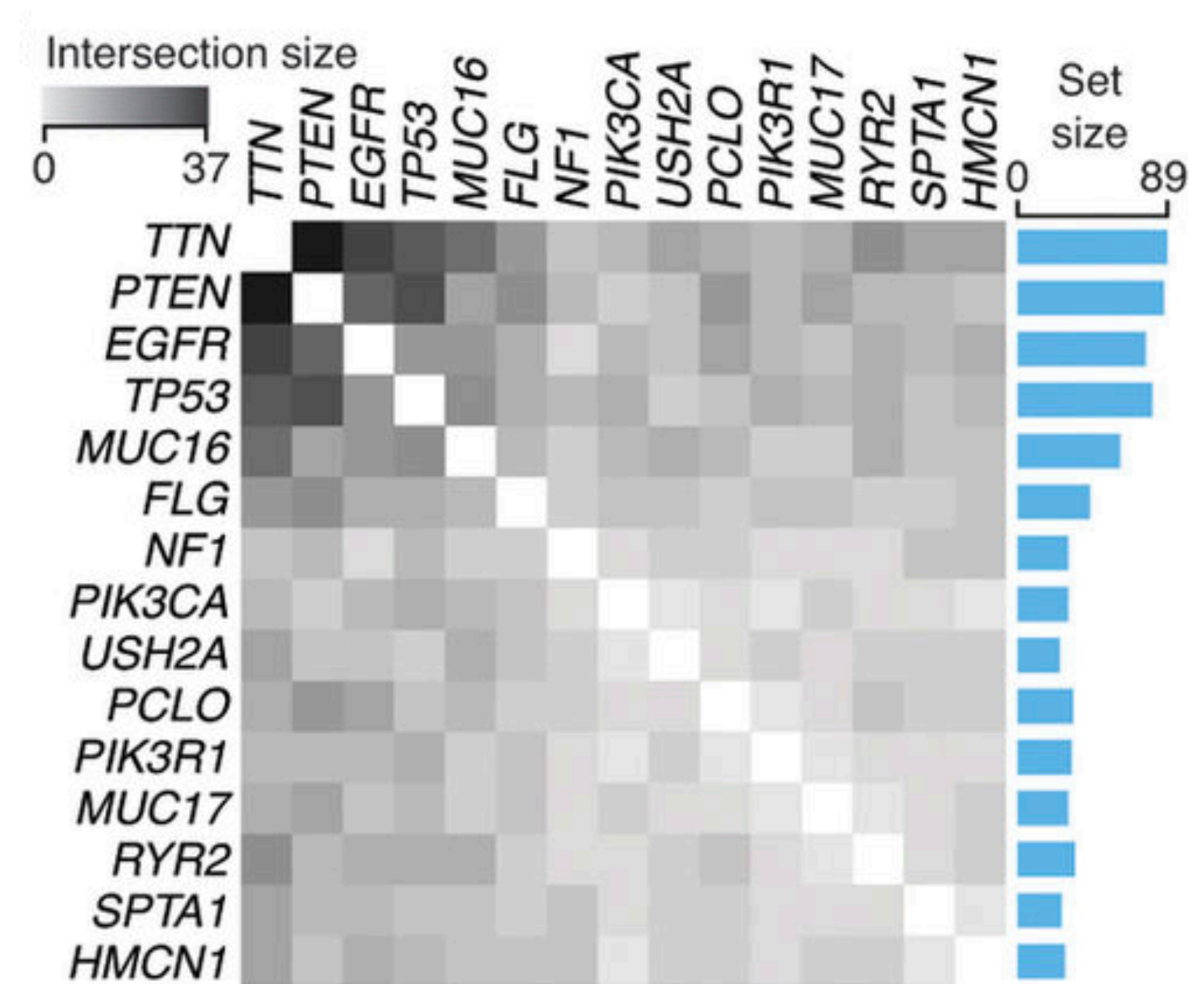
Showing Pairwise Overlap

Doesn't show higher-order overlaps

Very scalable

Can't show attributes

Co-Mutations of genes

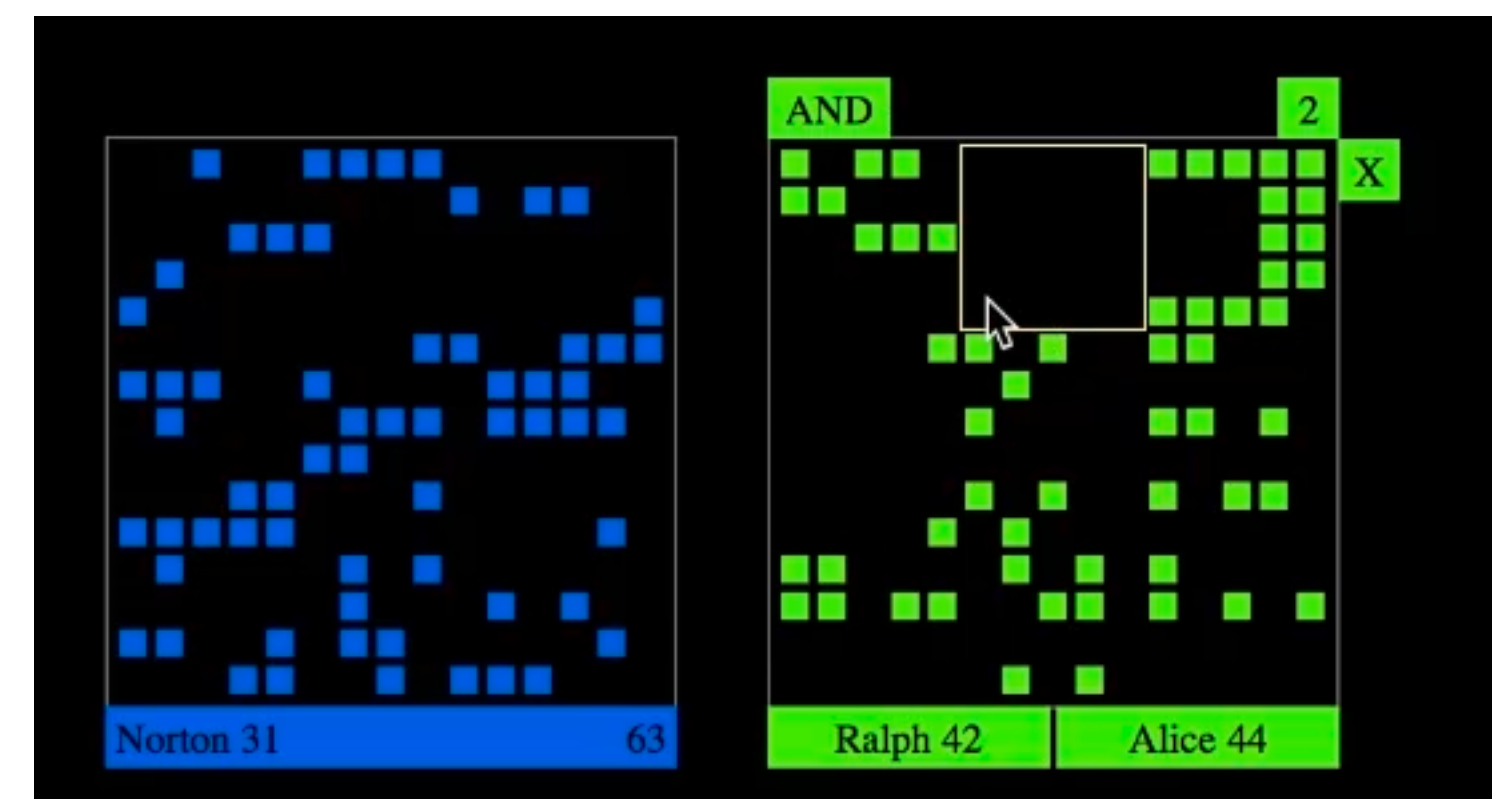
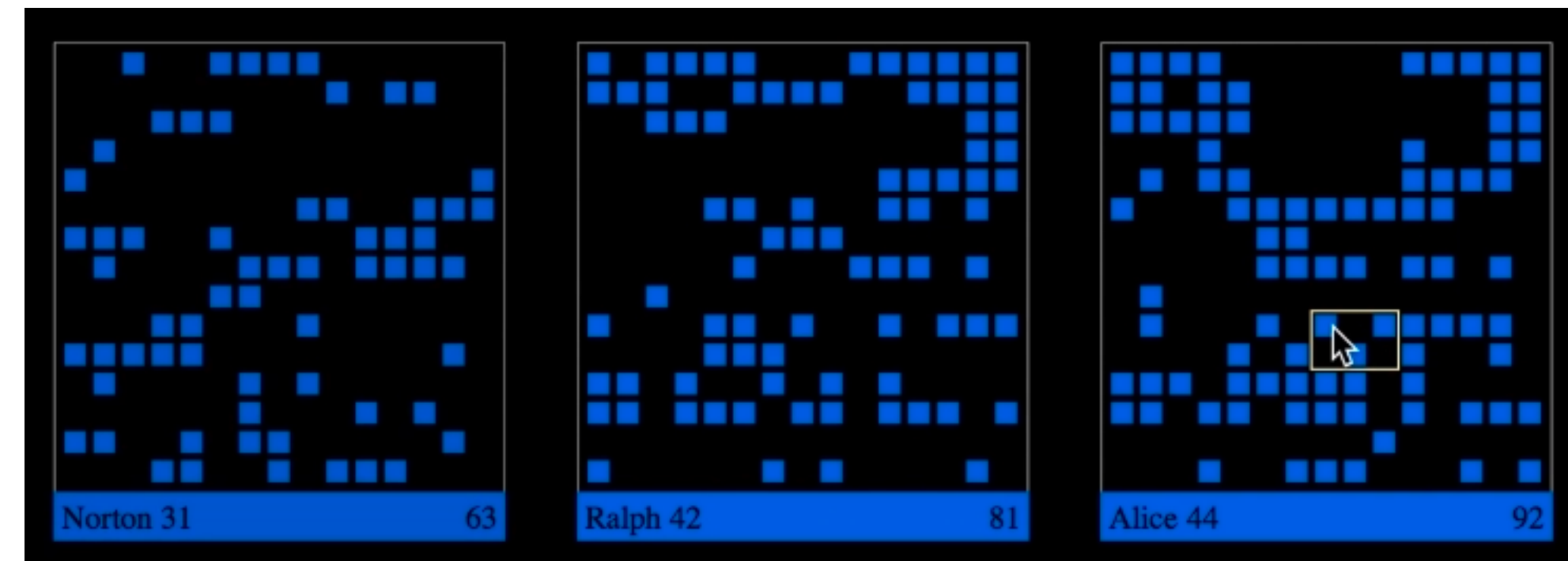


Set Matrices: OnSet

Set membership for each item shown in matrix

Comparisons can be made using AND or OR operations

Good for many sets and few items



Linear Diagrams



Fig. 1. Visualizing sets: linear diagrams.

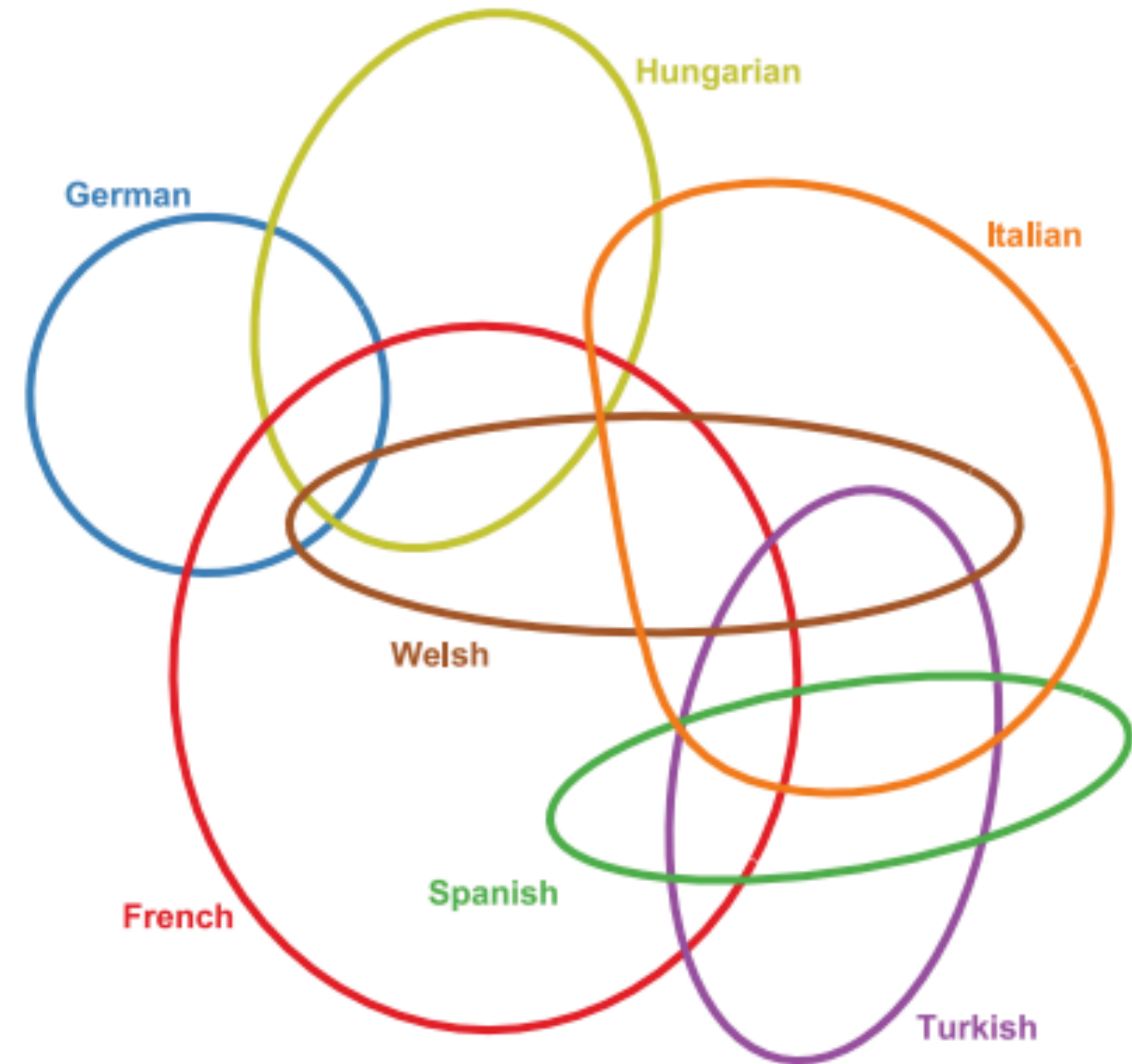
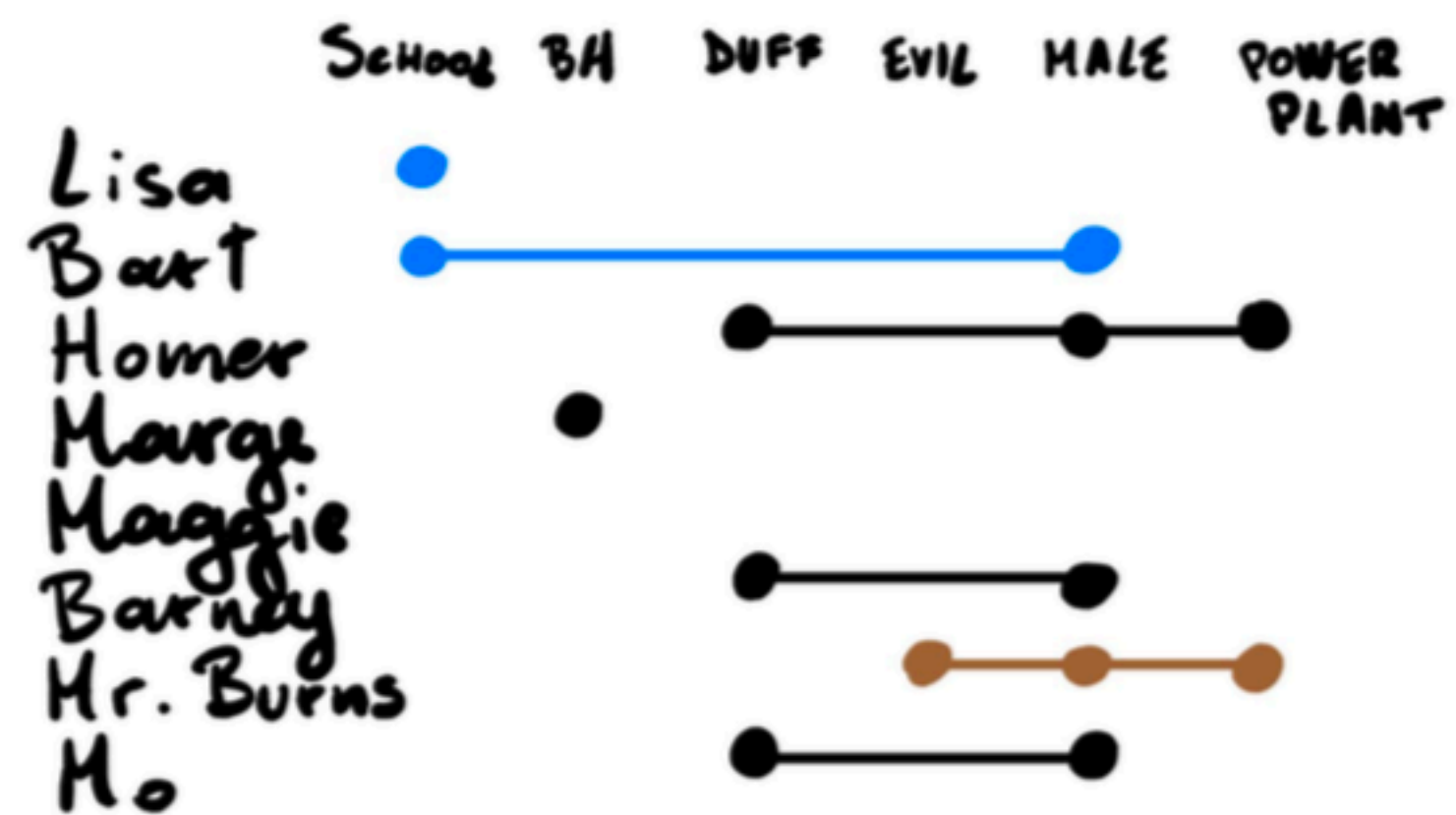


Fig. 2. Visualizing sets: Euler diagrams.



Marg
Burt
Homer

	School	Hair	Work	Beer
Marg				
Burt				
Homer				

Reasonable solutions will contain 2-3 characters

	School	B/H Hair	Duff fan	Evil	male	Power plant	Age
Lisa							
Bart							
Homer							
Marge							
Margee							
Maggie							
Barney							
Mr Burns							
Mo							

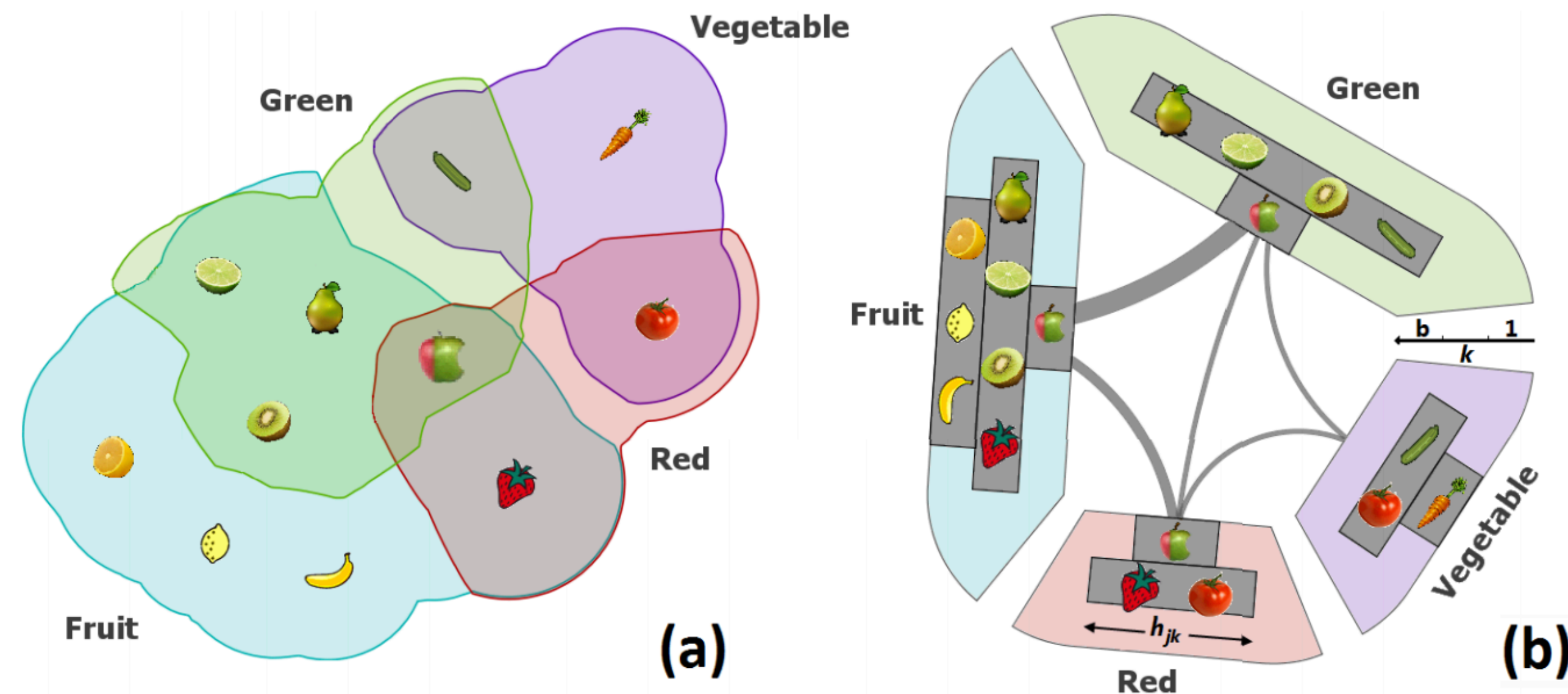
Radial Sets

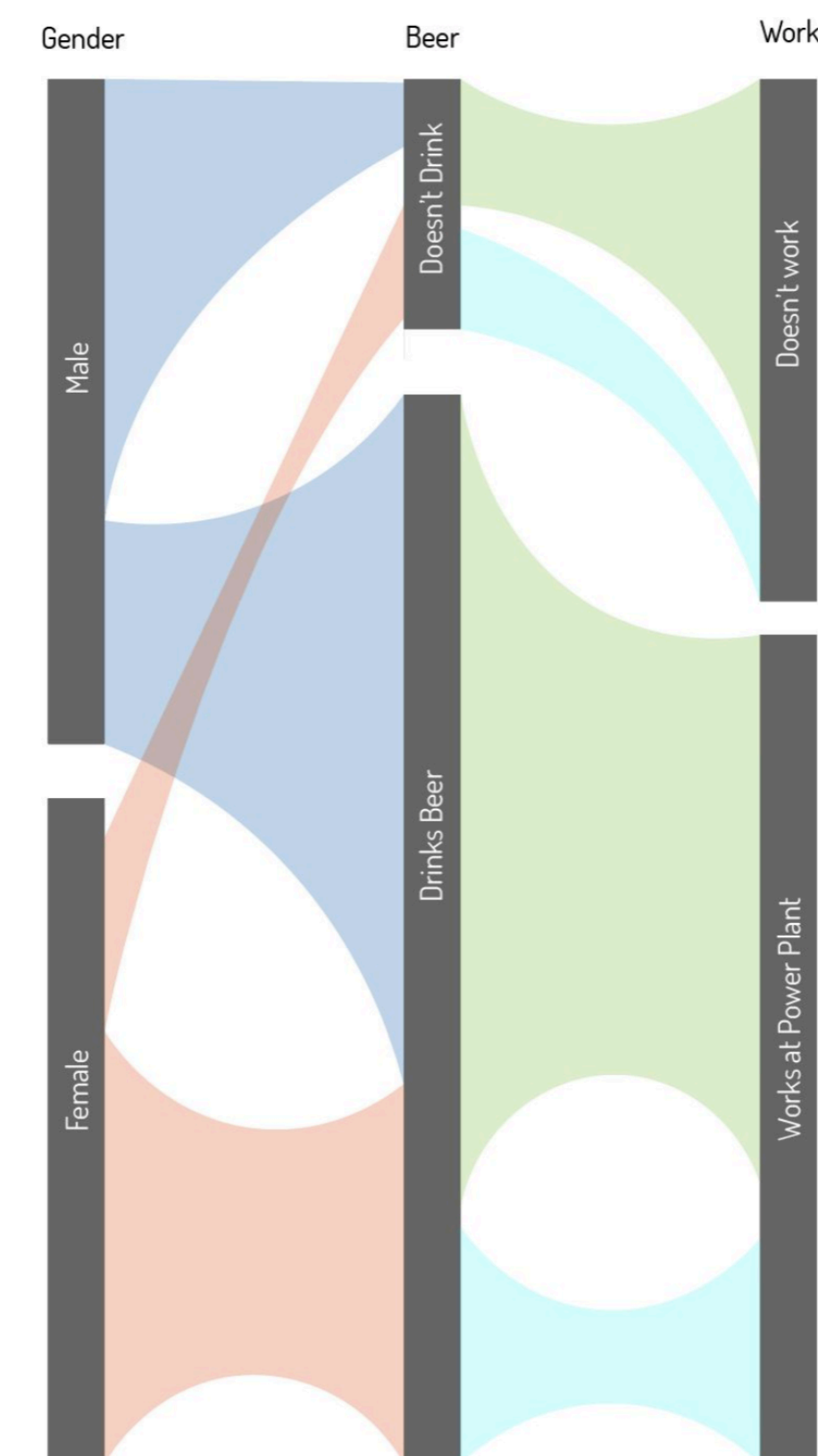
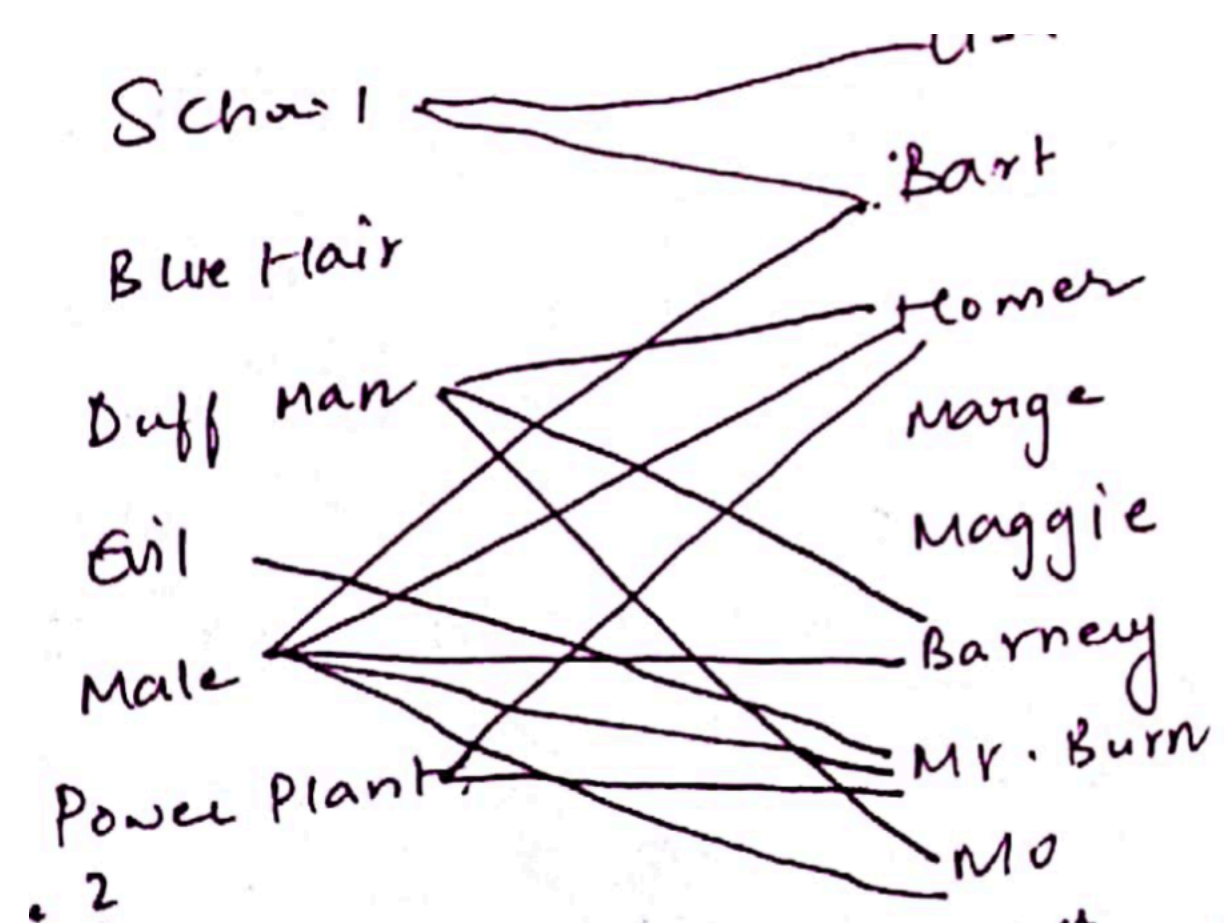
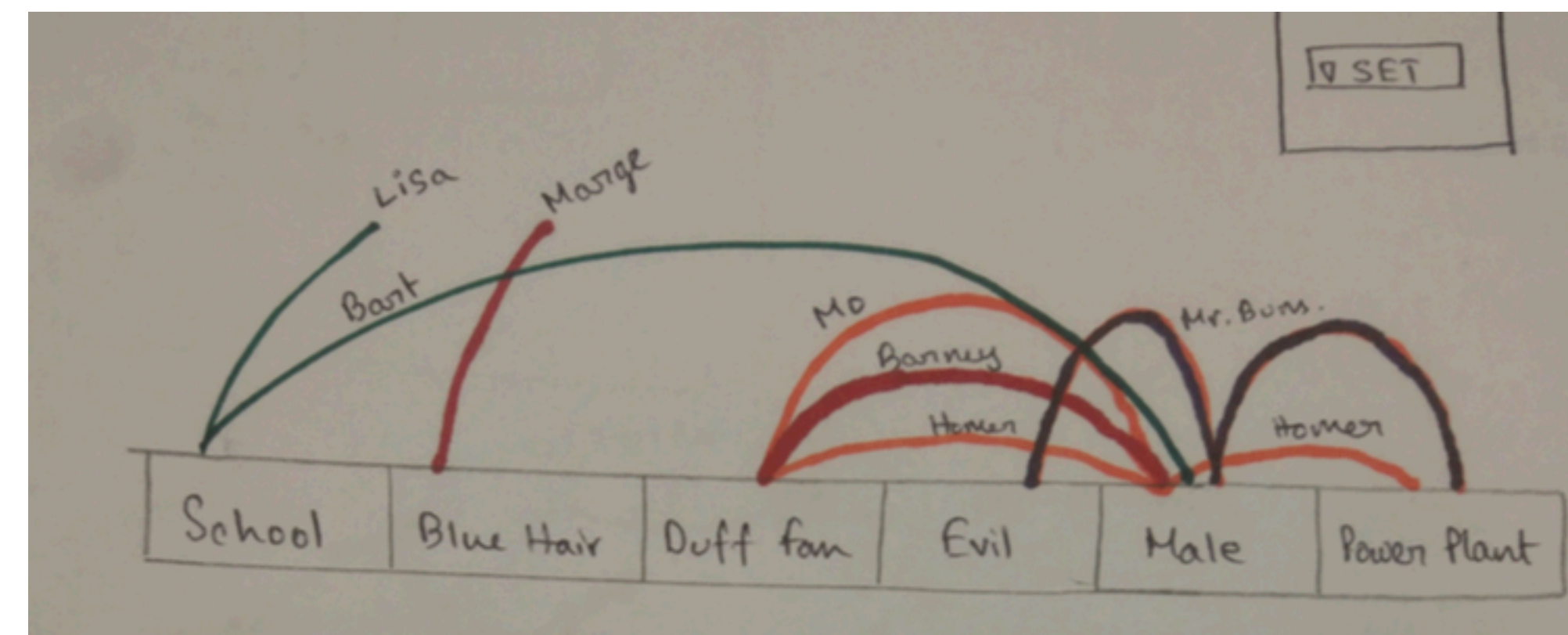
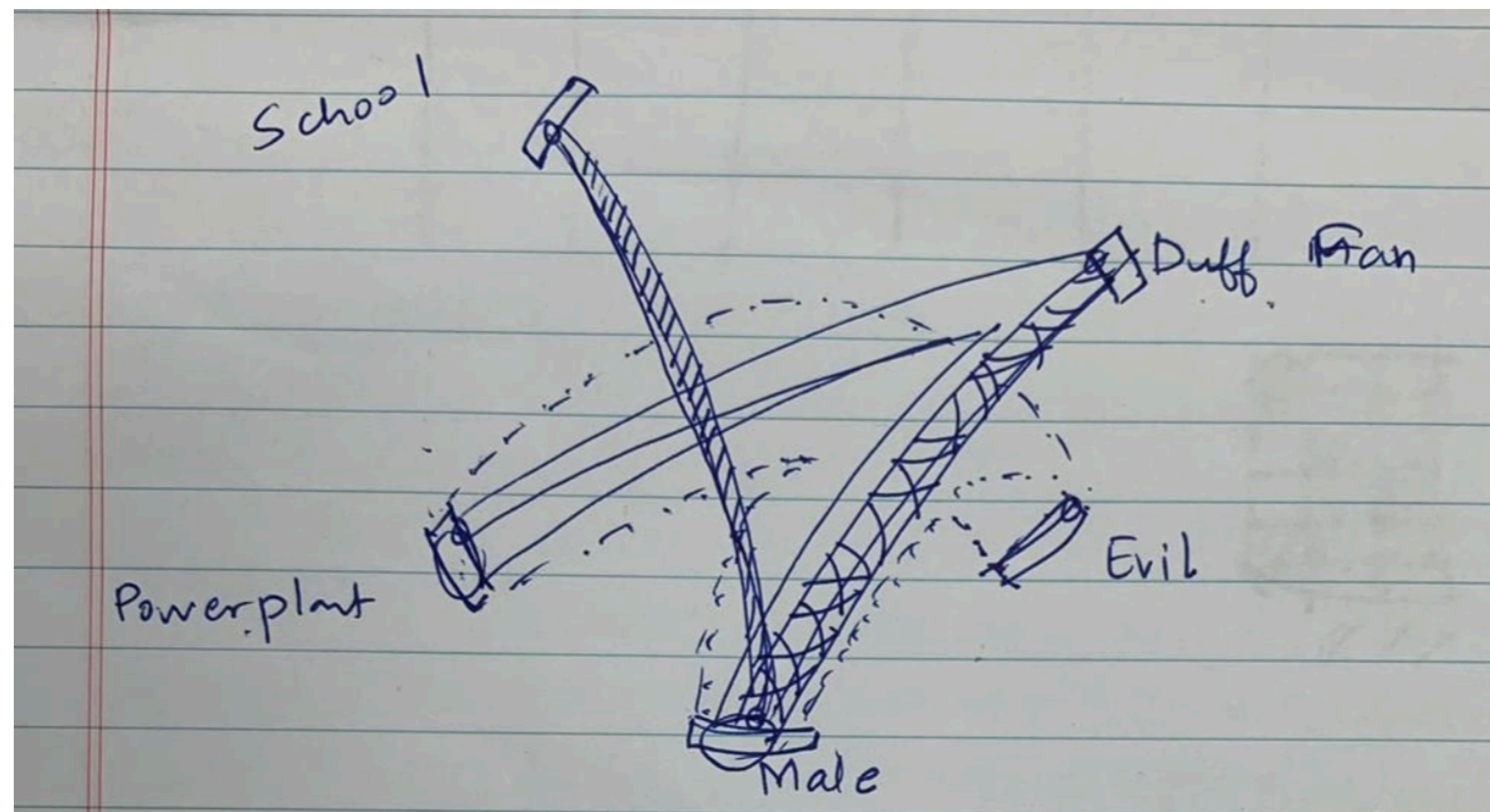
Sets are segments on a “circle”

Relationships are encoded as ribbons

Size of segments encodes size of sets

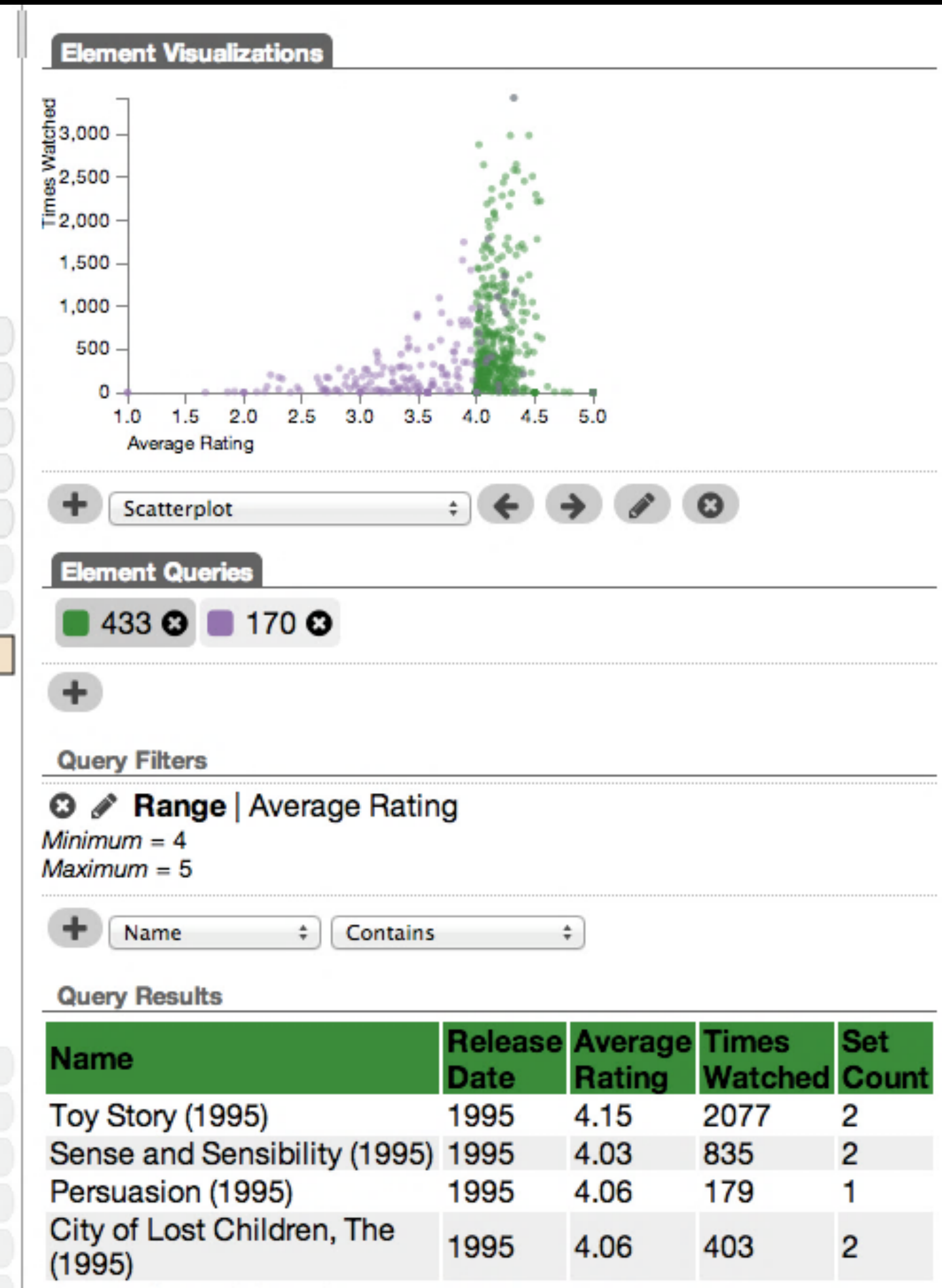
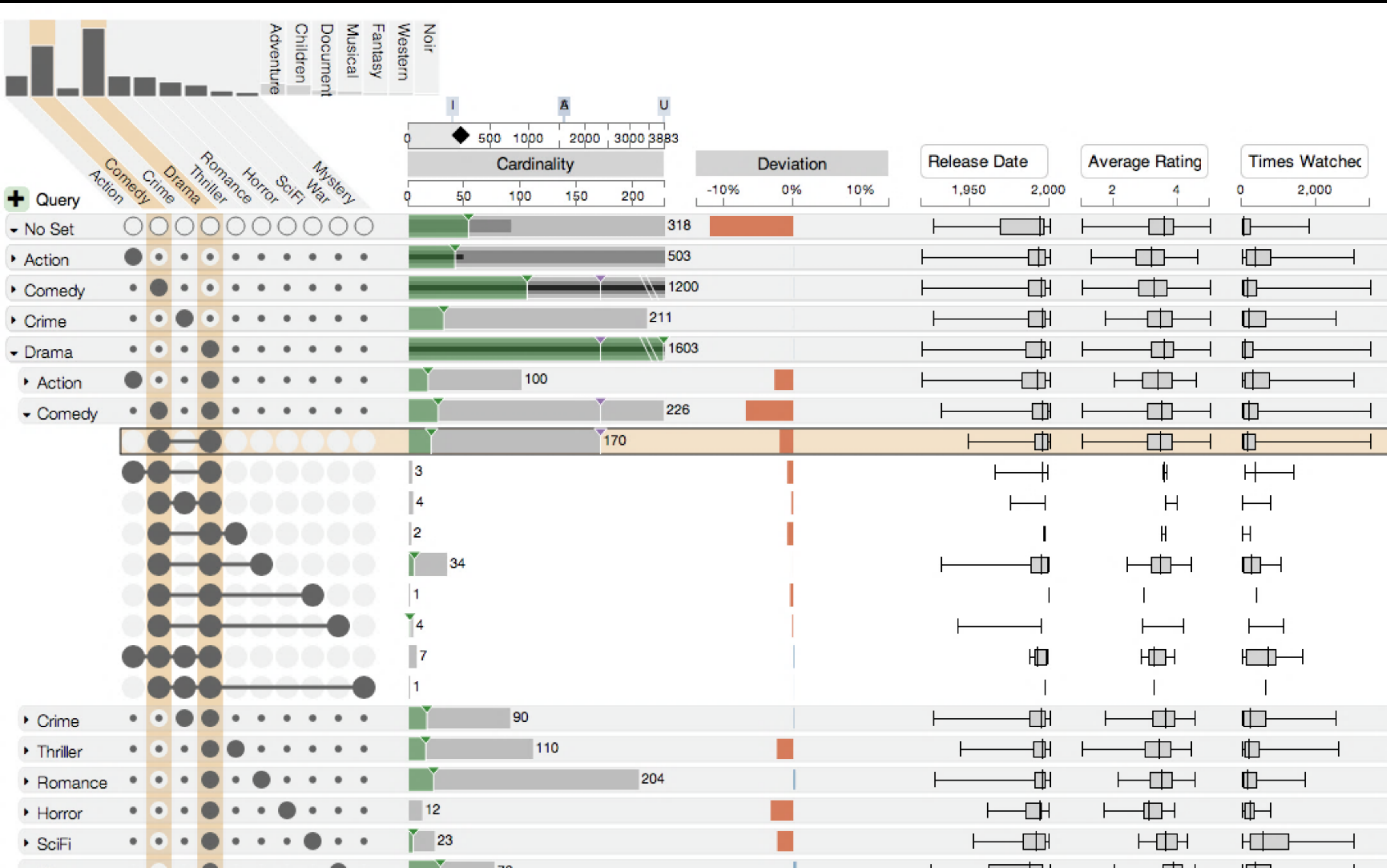
Histograms in segments show degrees





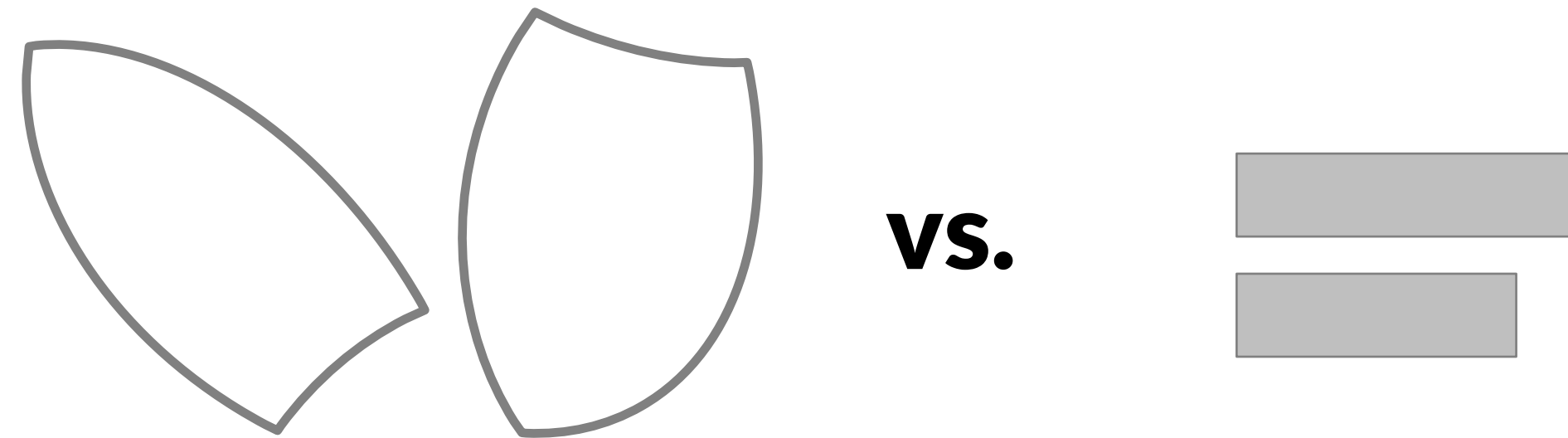
[InfoVis'14]

UpSet Visualizing Intersecting Sets

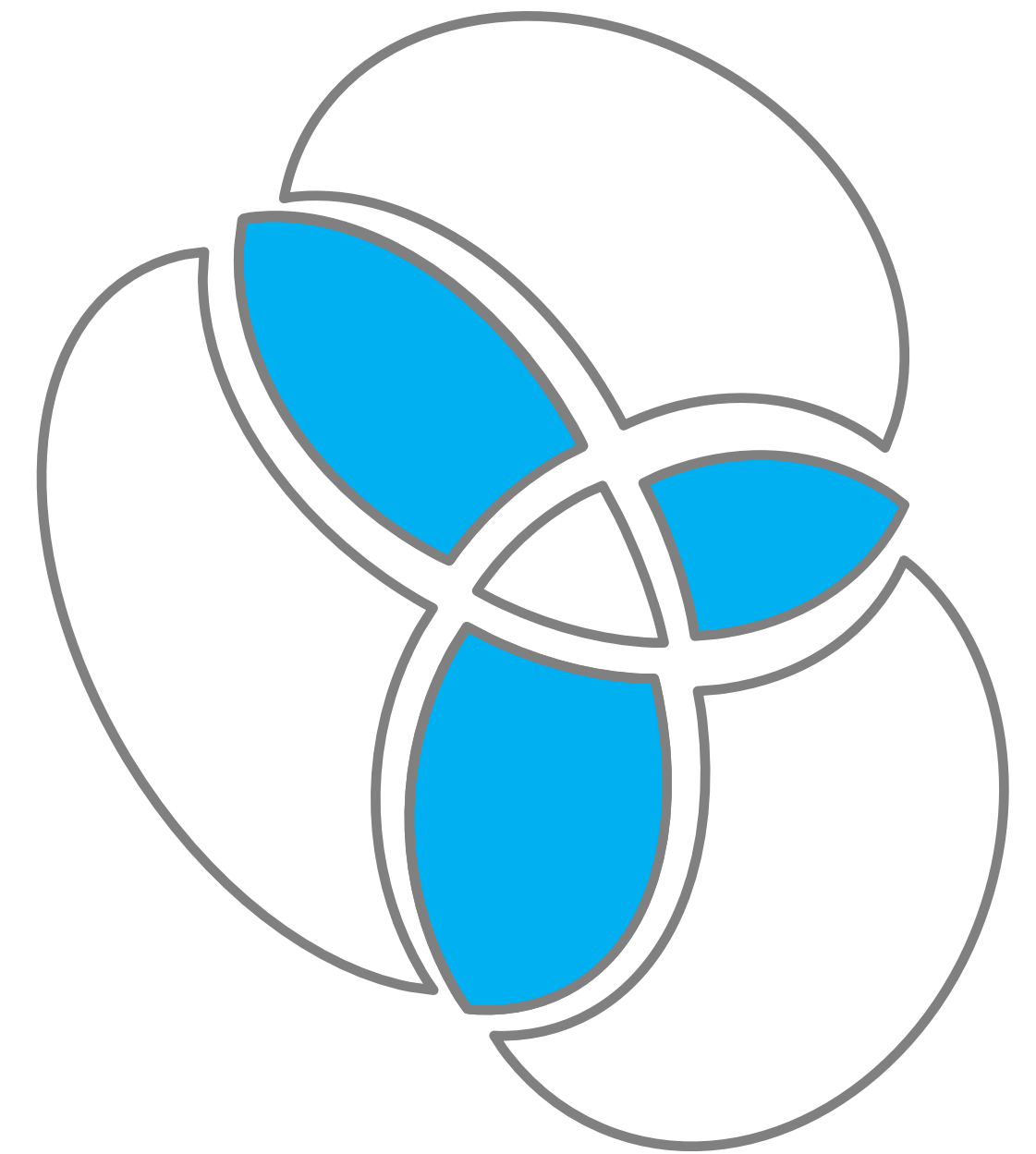


Set Vis Goals

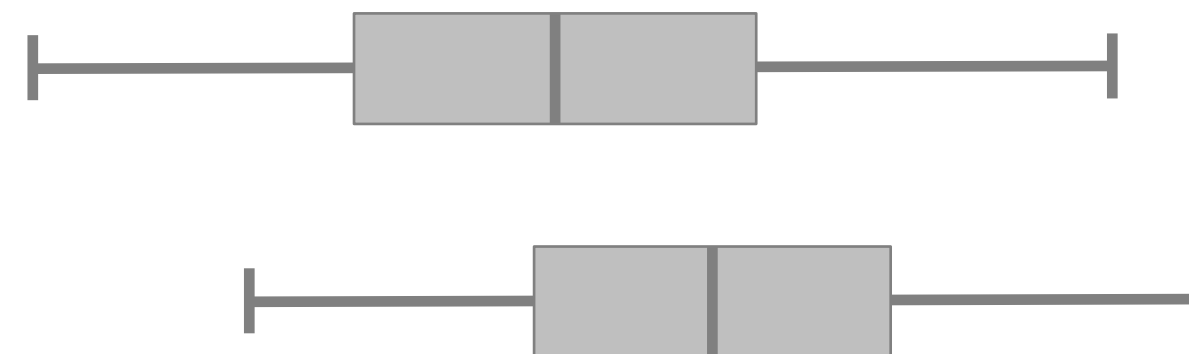
1. Efficient visual encoding

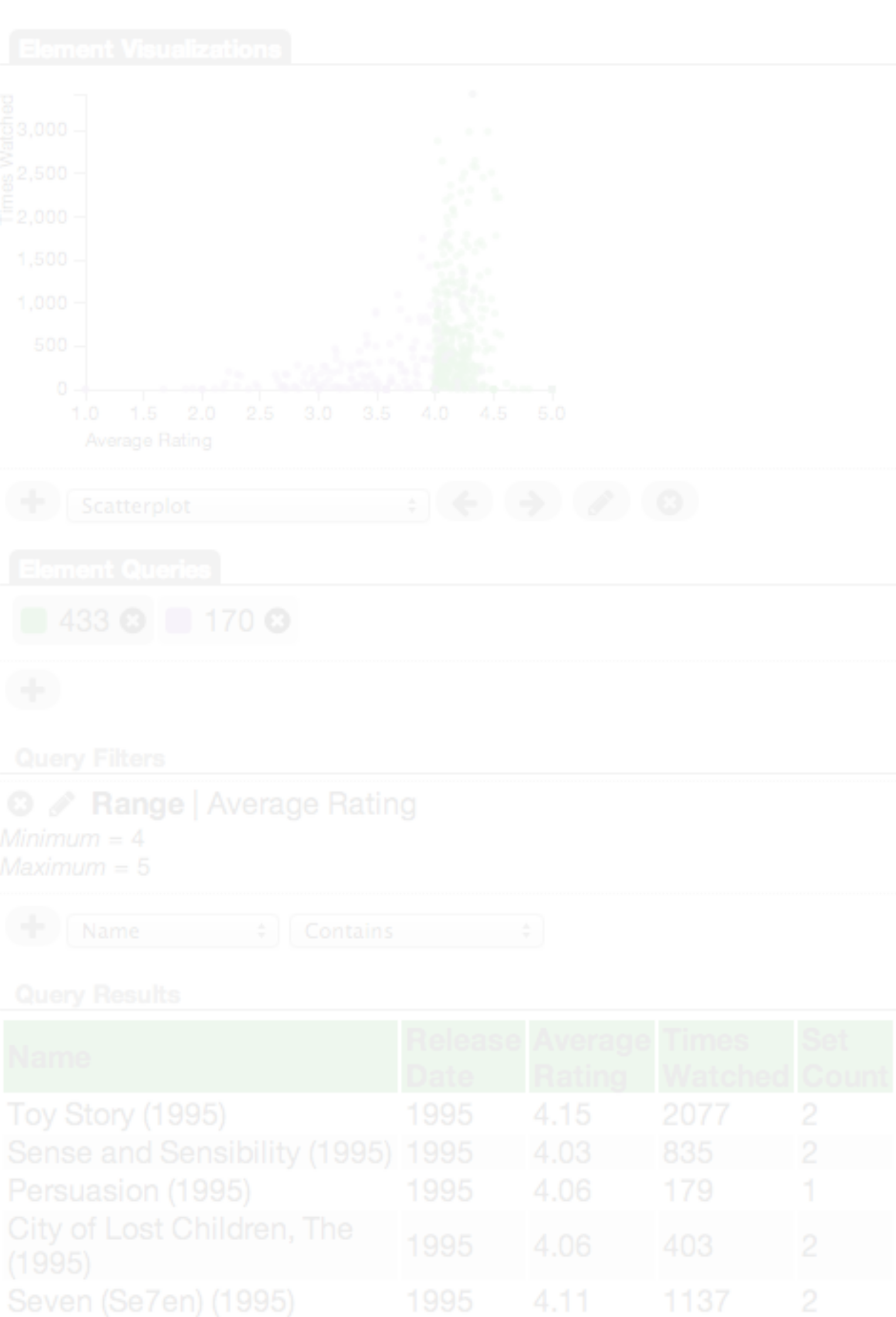


2. Creating complex slices of a dataset



3. Visualize attributes



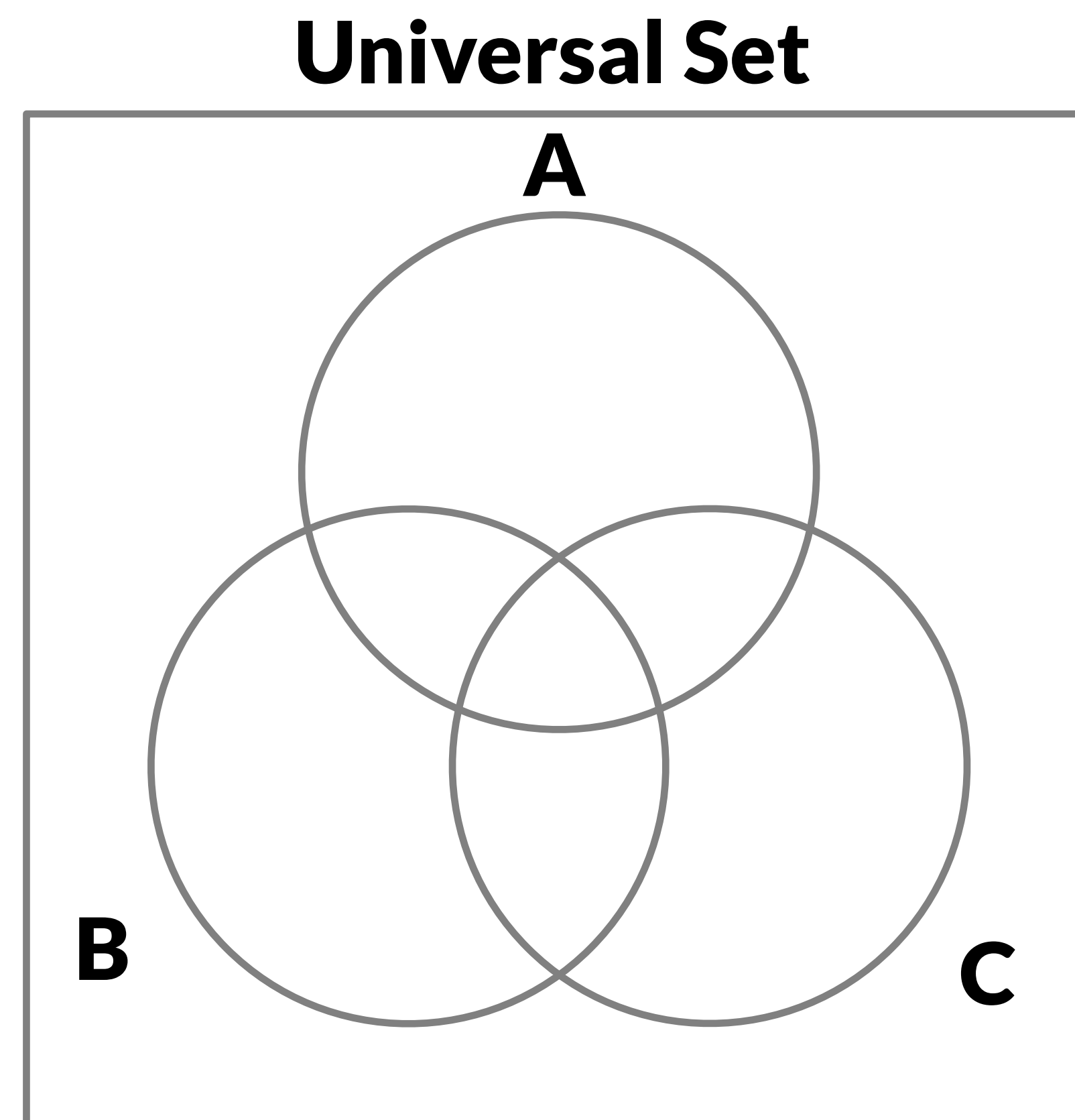
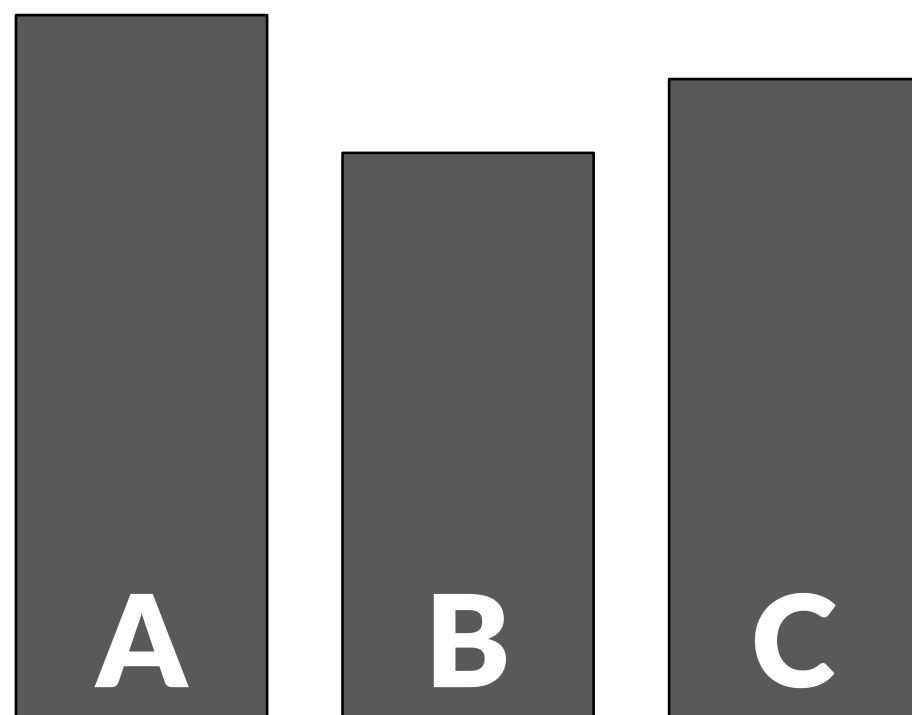


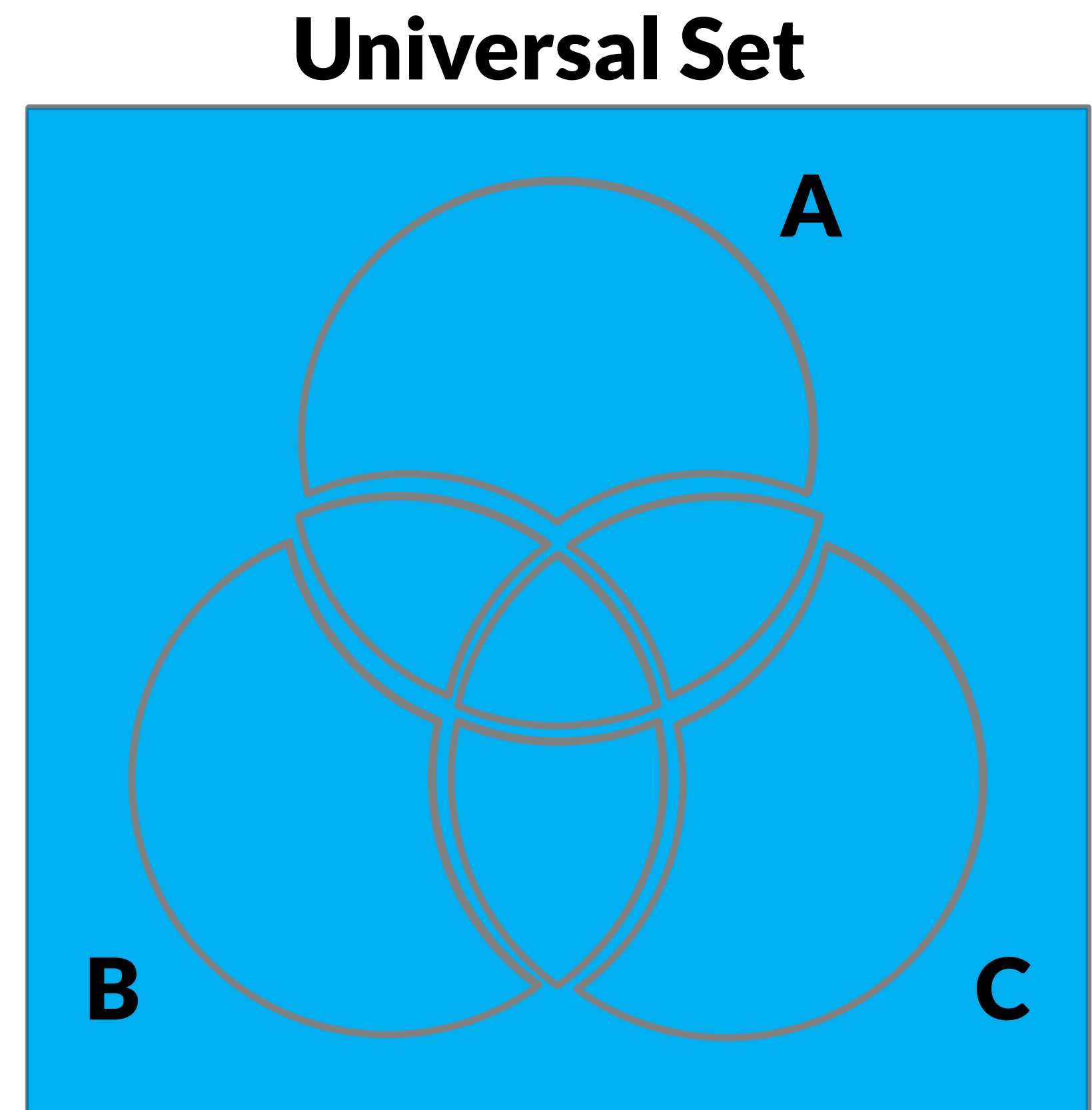
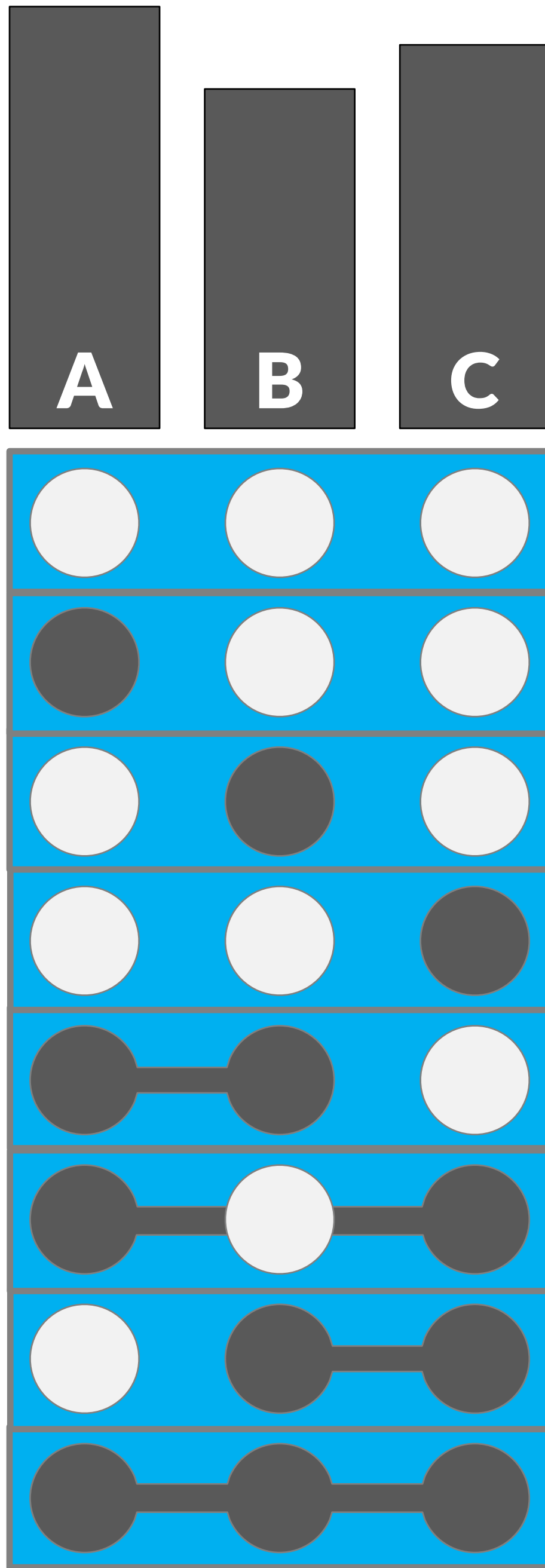
Visualizing Intersections

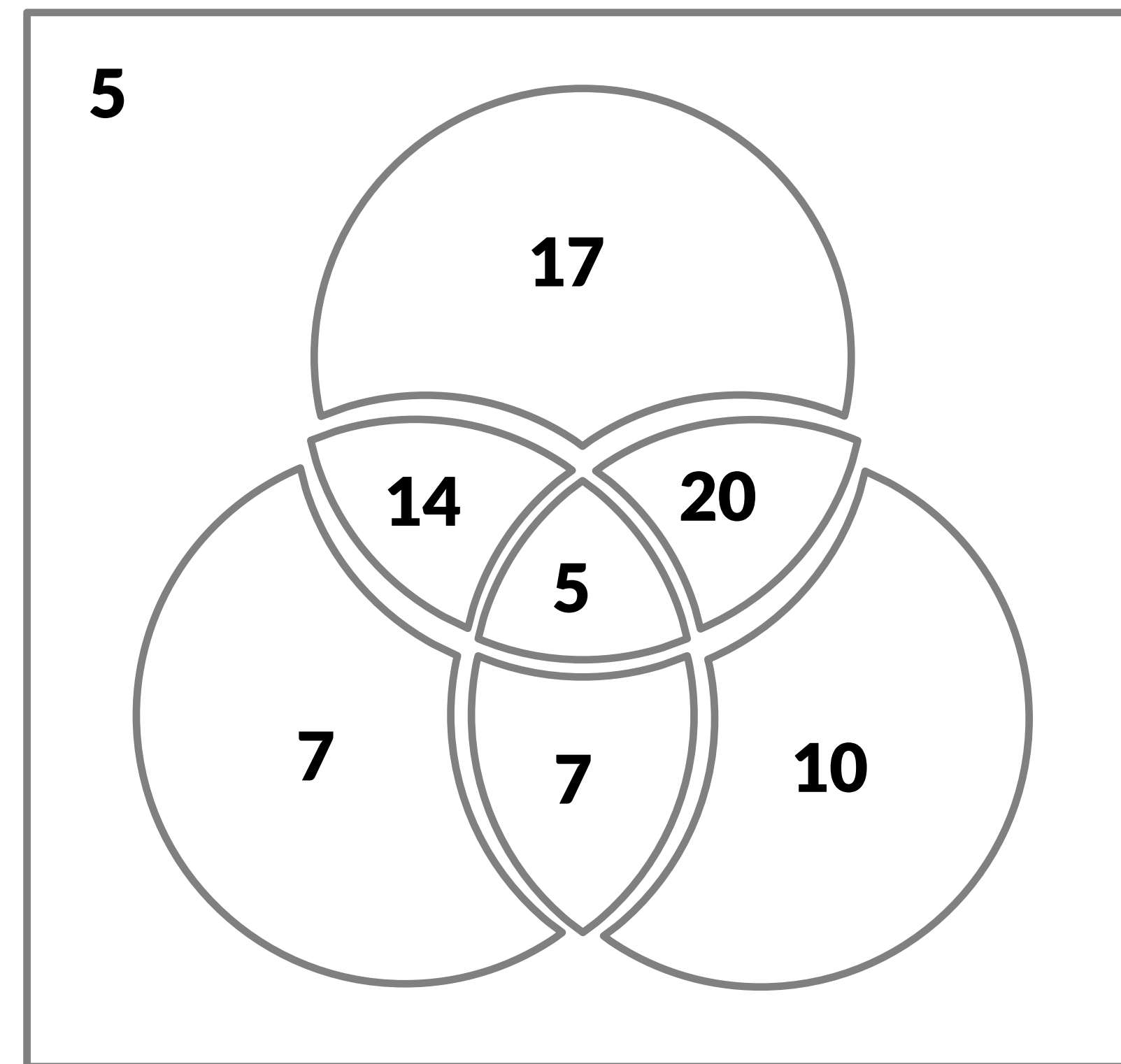
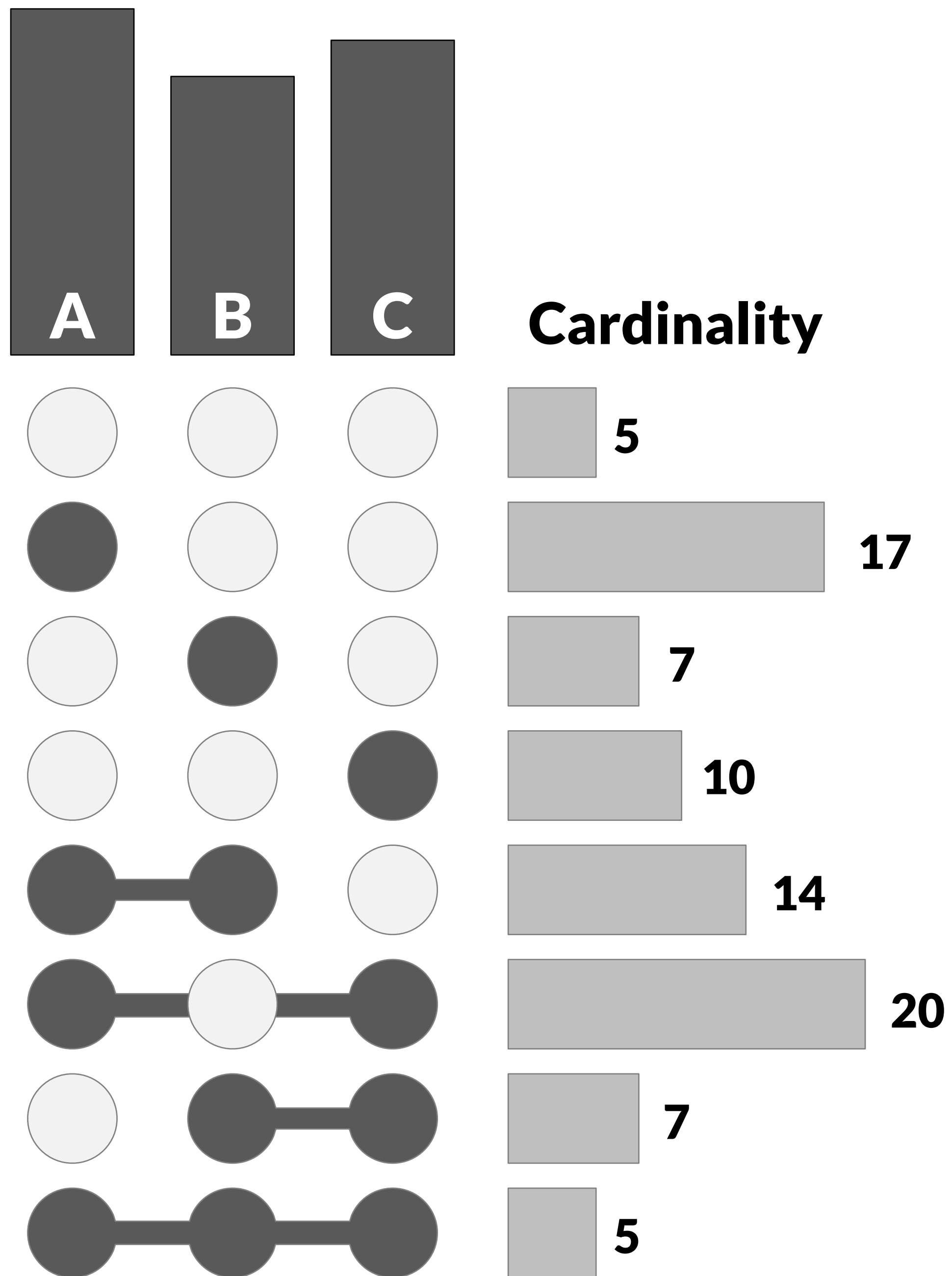
Visualizing Properties

Element List & Queries

Visualizing Intersections



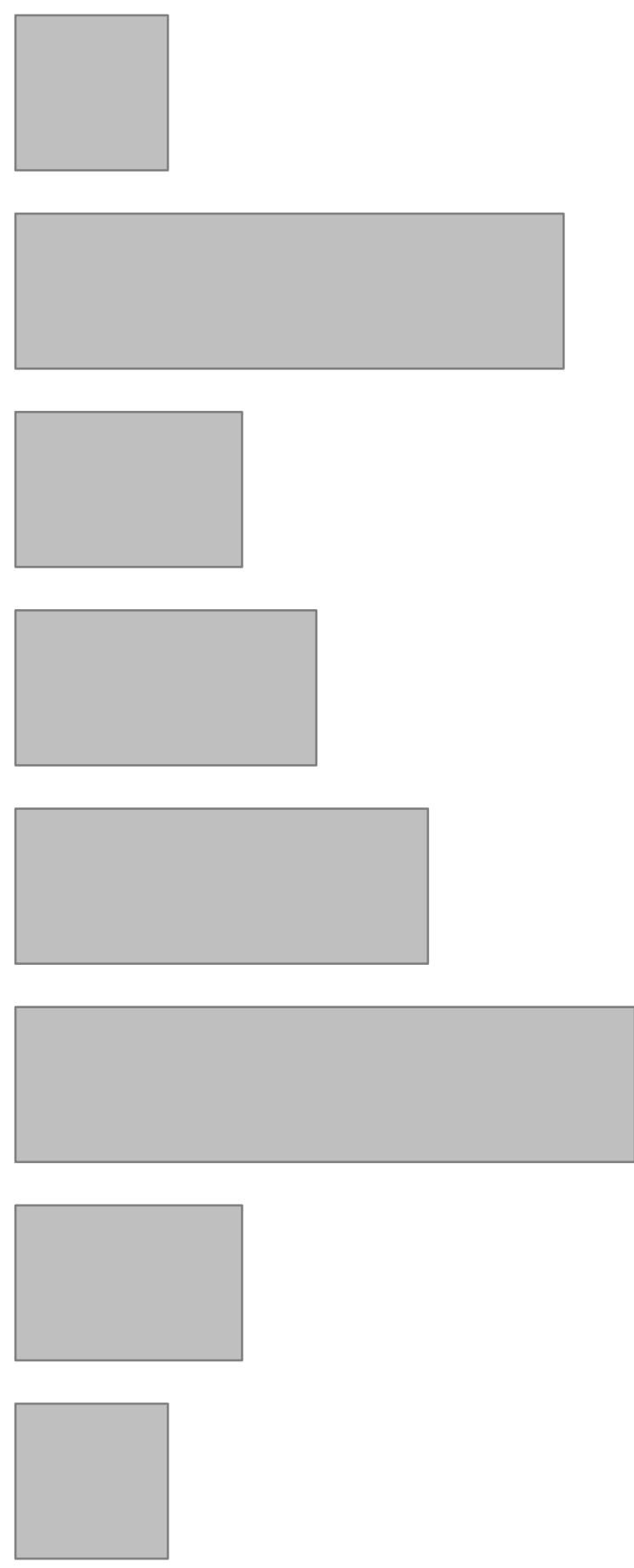
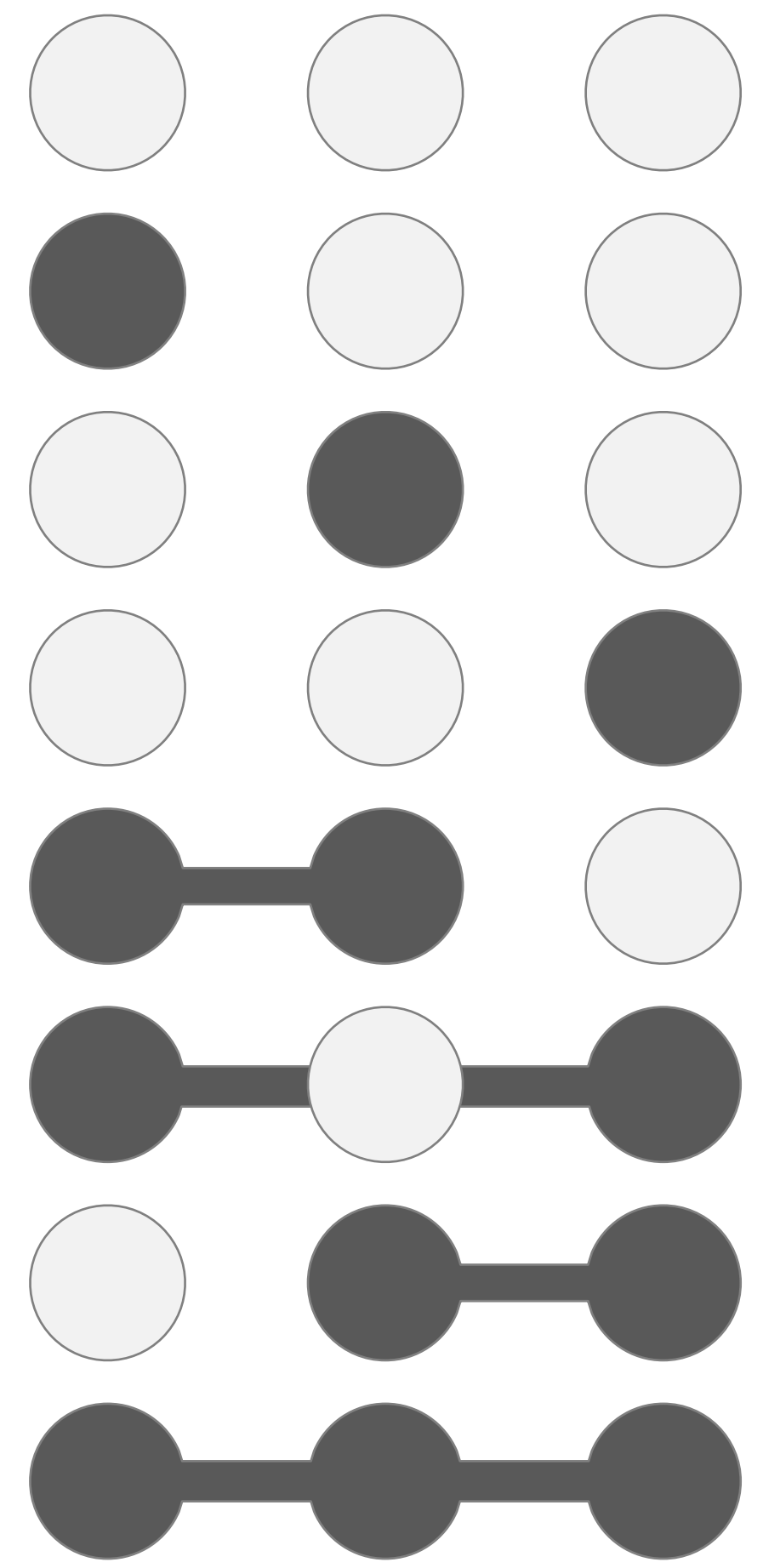
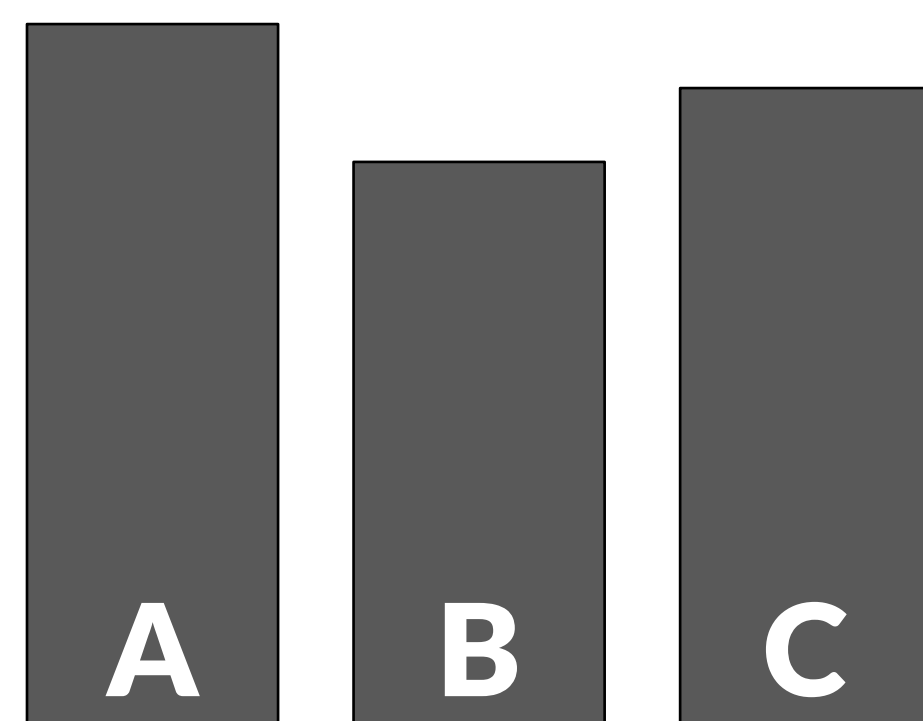




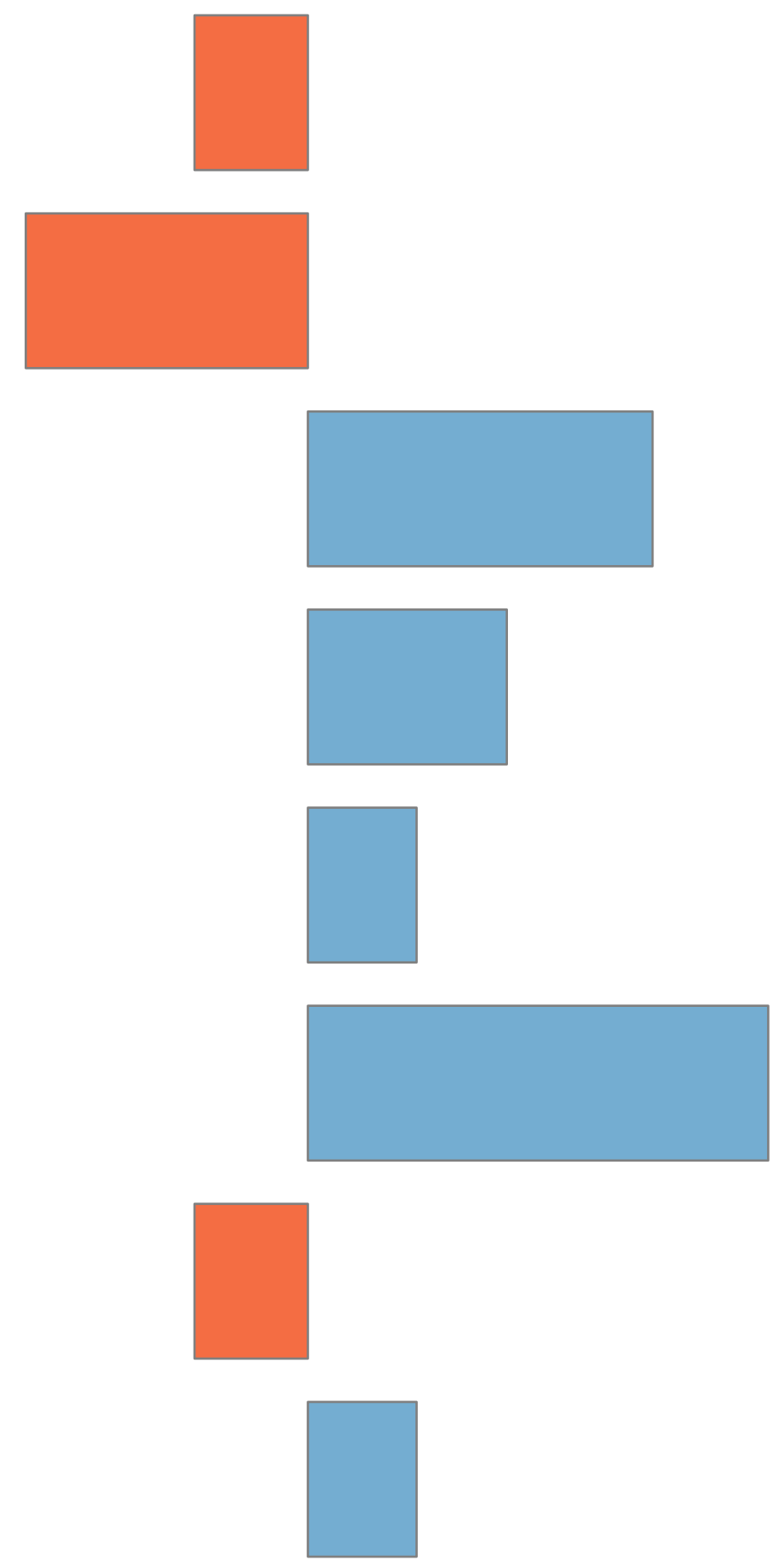
Plotting Attributes

What's the distribution of the size of an intersection?
 attribute in an intersection?

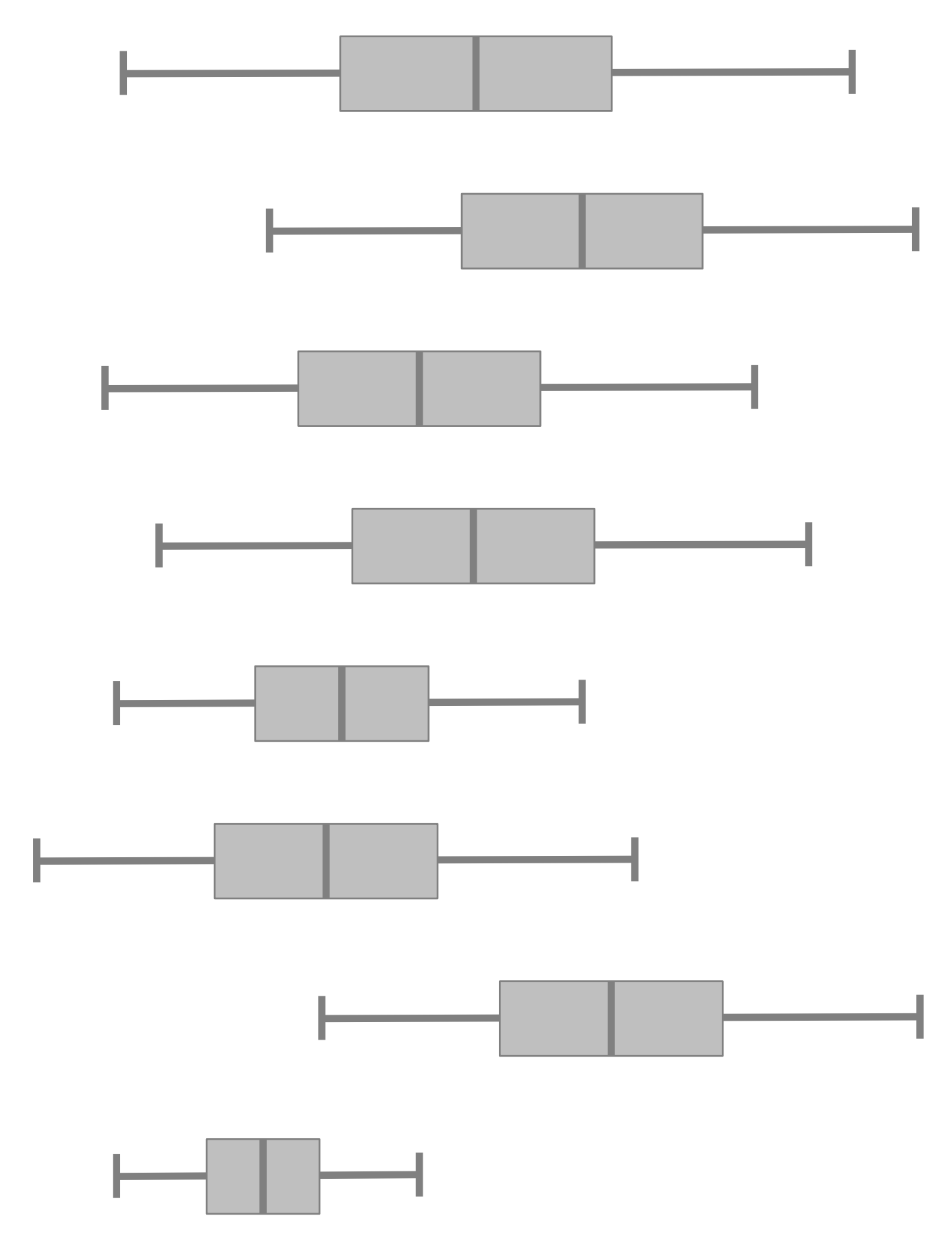
← Additional Plots →



Deviation



Attributes



First, aggregate by
Don't Aggregate

Then, aggregate by
Don't Aggregate

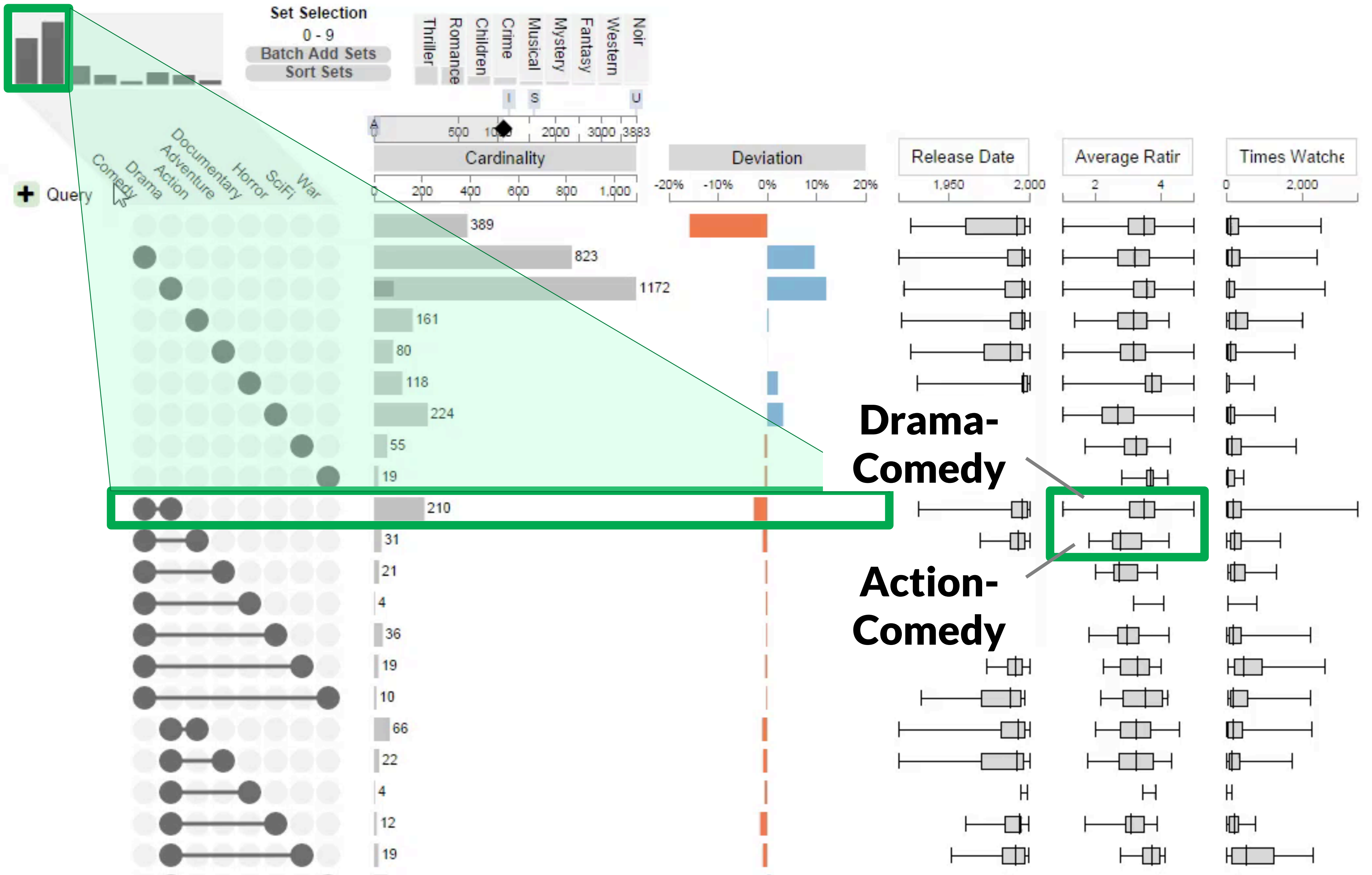
Sort by
☒ Degree
☐ Cardinality
☐ Deviation

Aggregates
Collapse All
Expand All

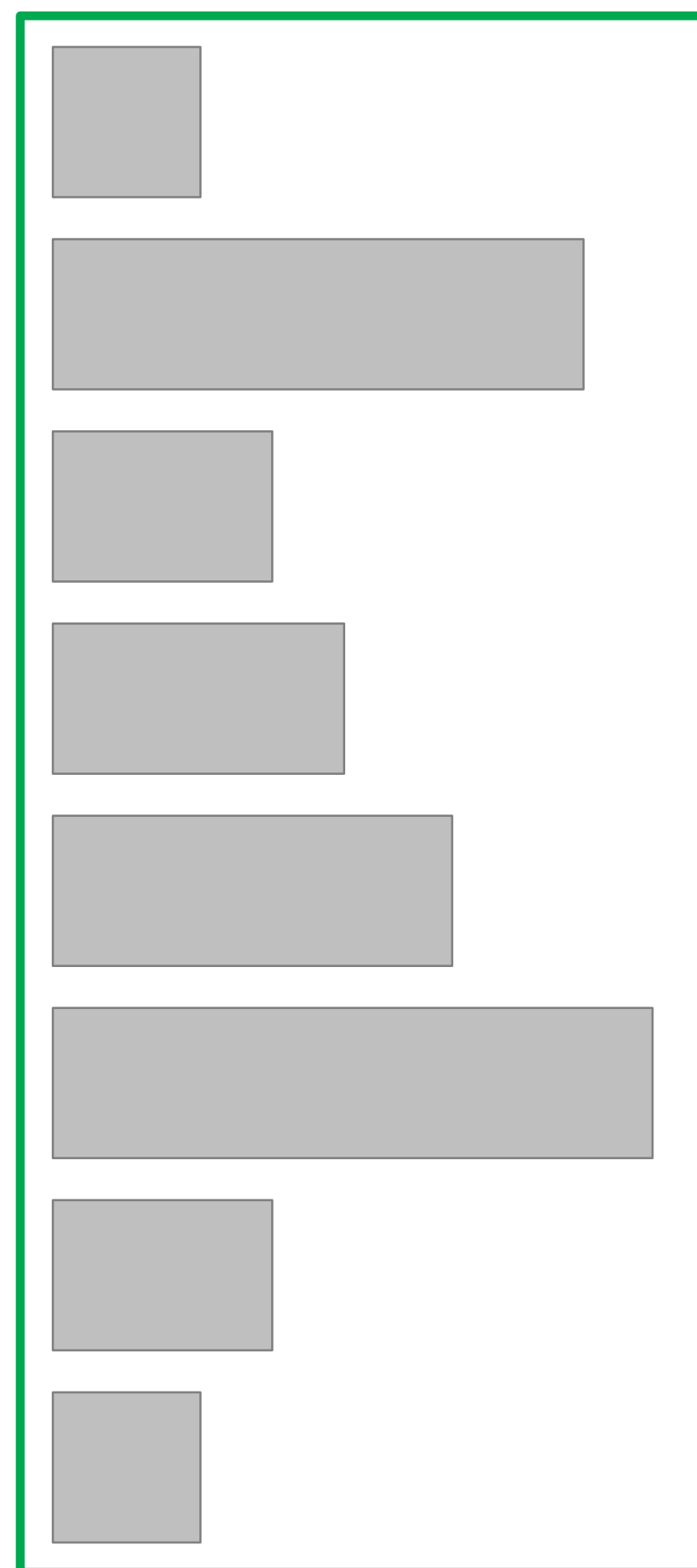
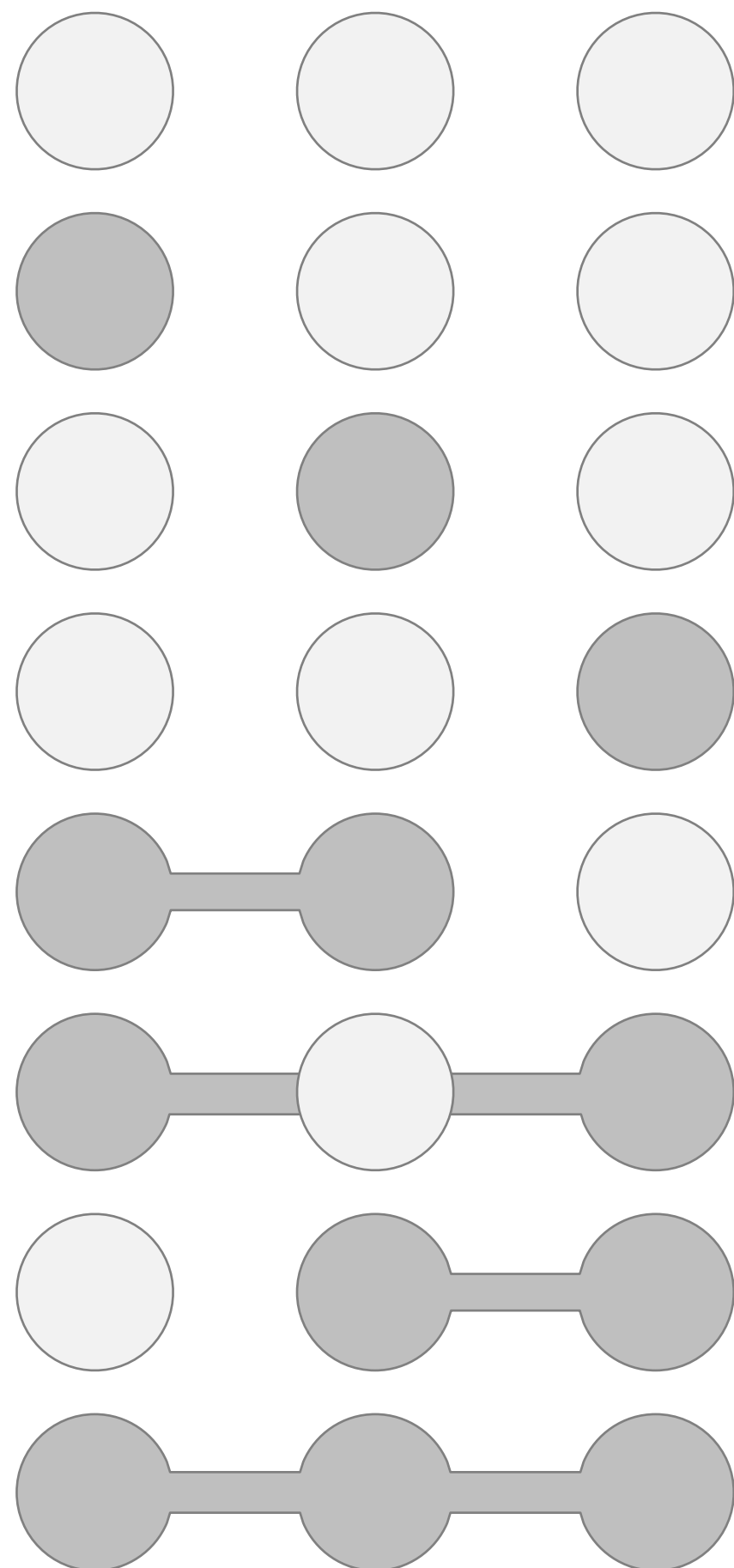
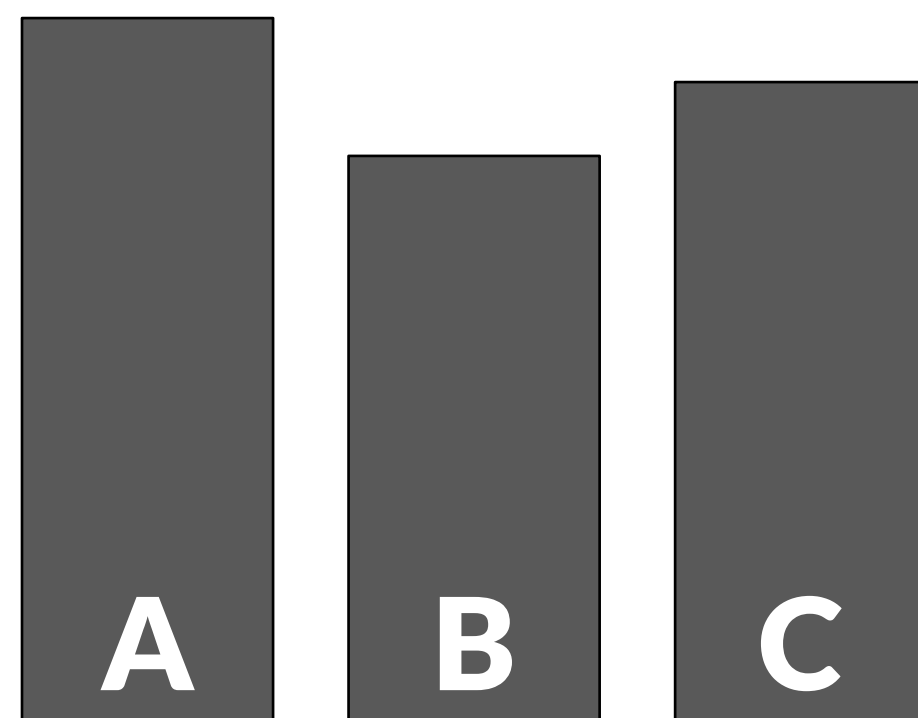
Row Height
Large

Data
Min Degree:
0
Max Degree:
5
☐ Hide Empty
Intersections

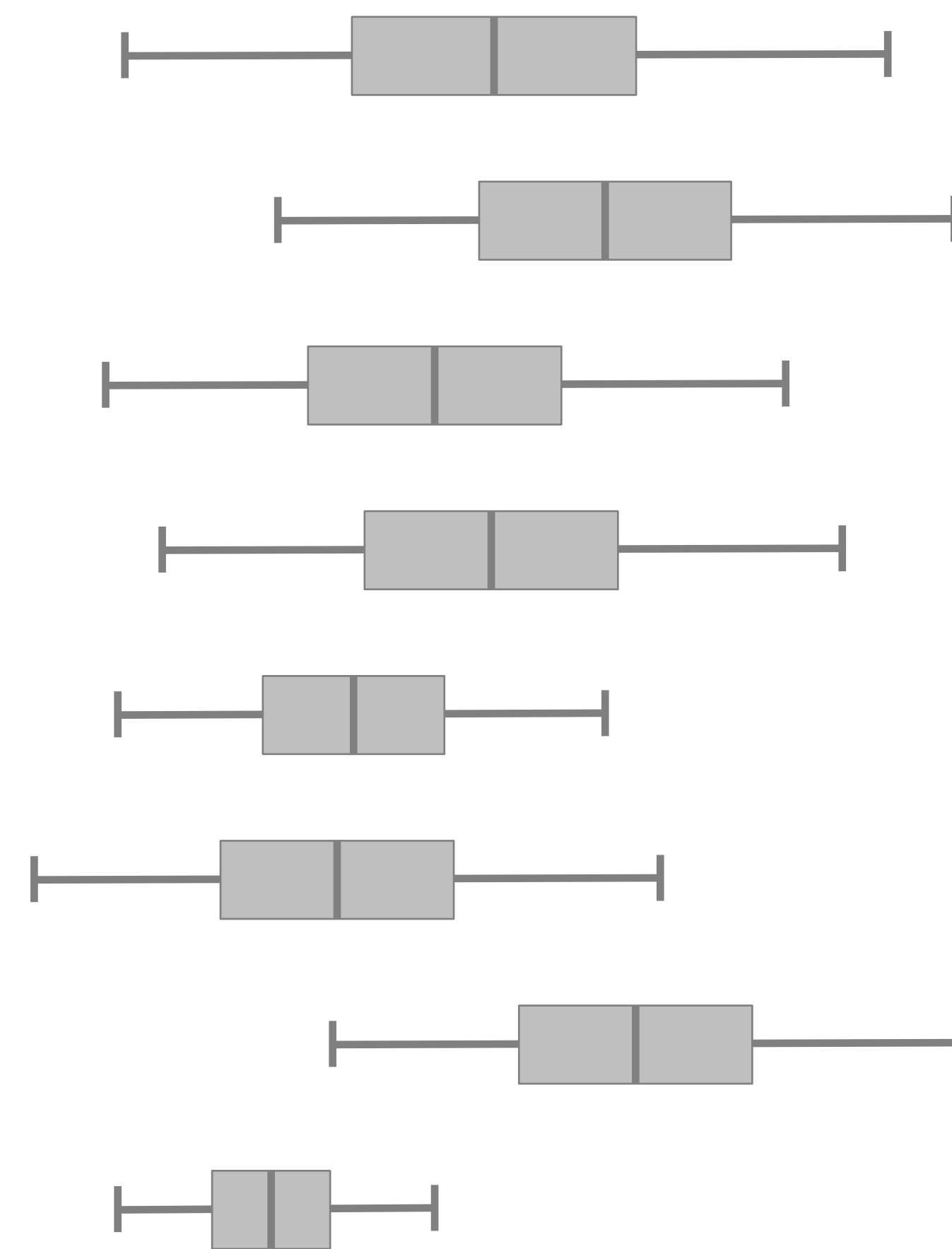
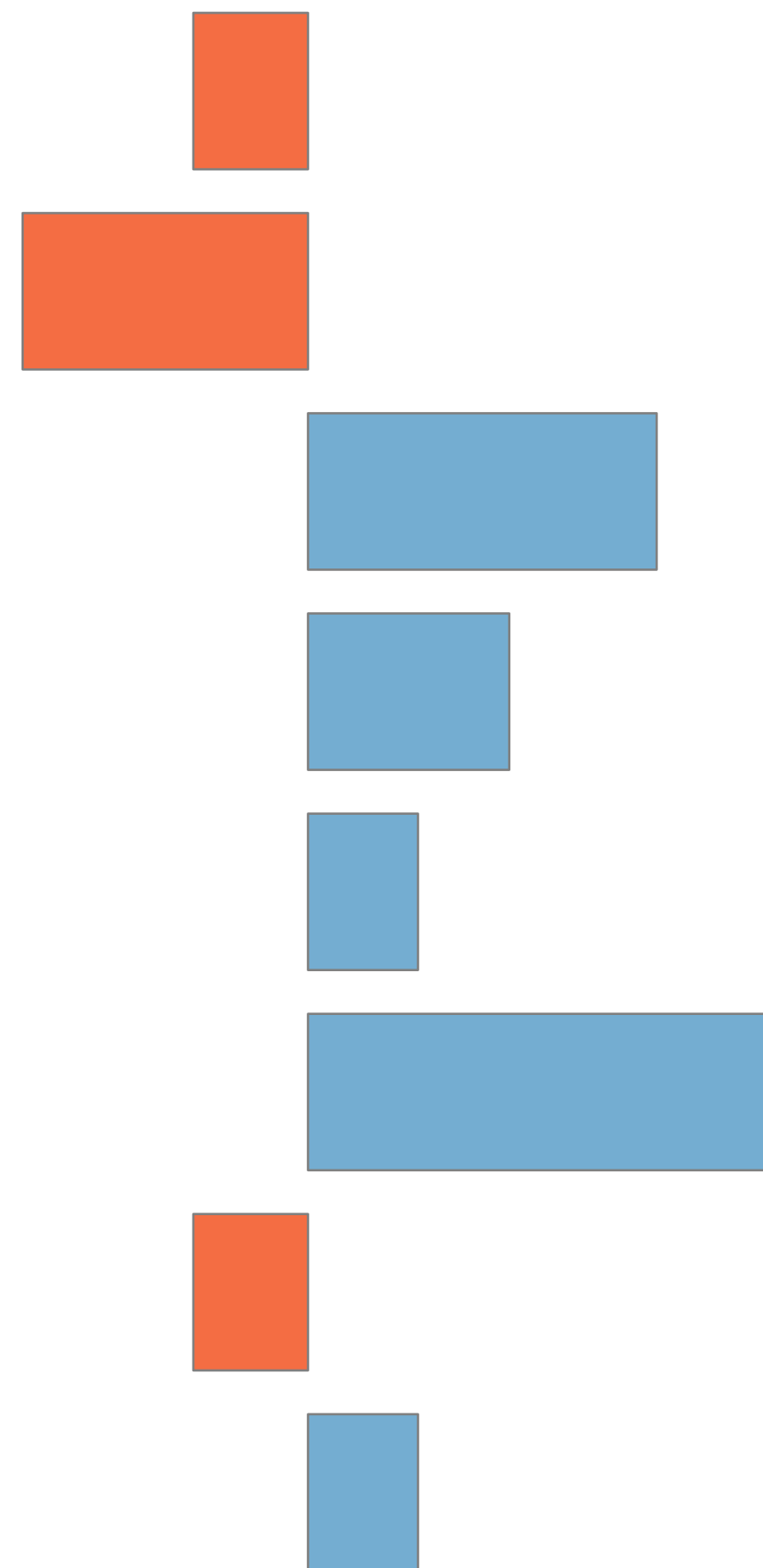
Dataset Information
Name: Movies
Genres
Sets: 17
Attributes: 6
Elements: 3883
Author: grouplens
Description:
MovieLens ratings
dataset, curated and
filtered by Alsallakh.
Source:
<http://grouplens.org/d..>

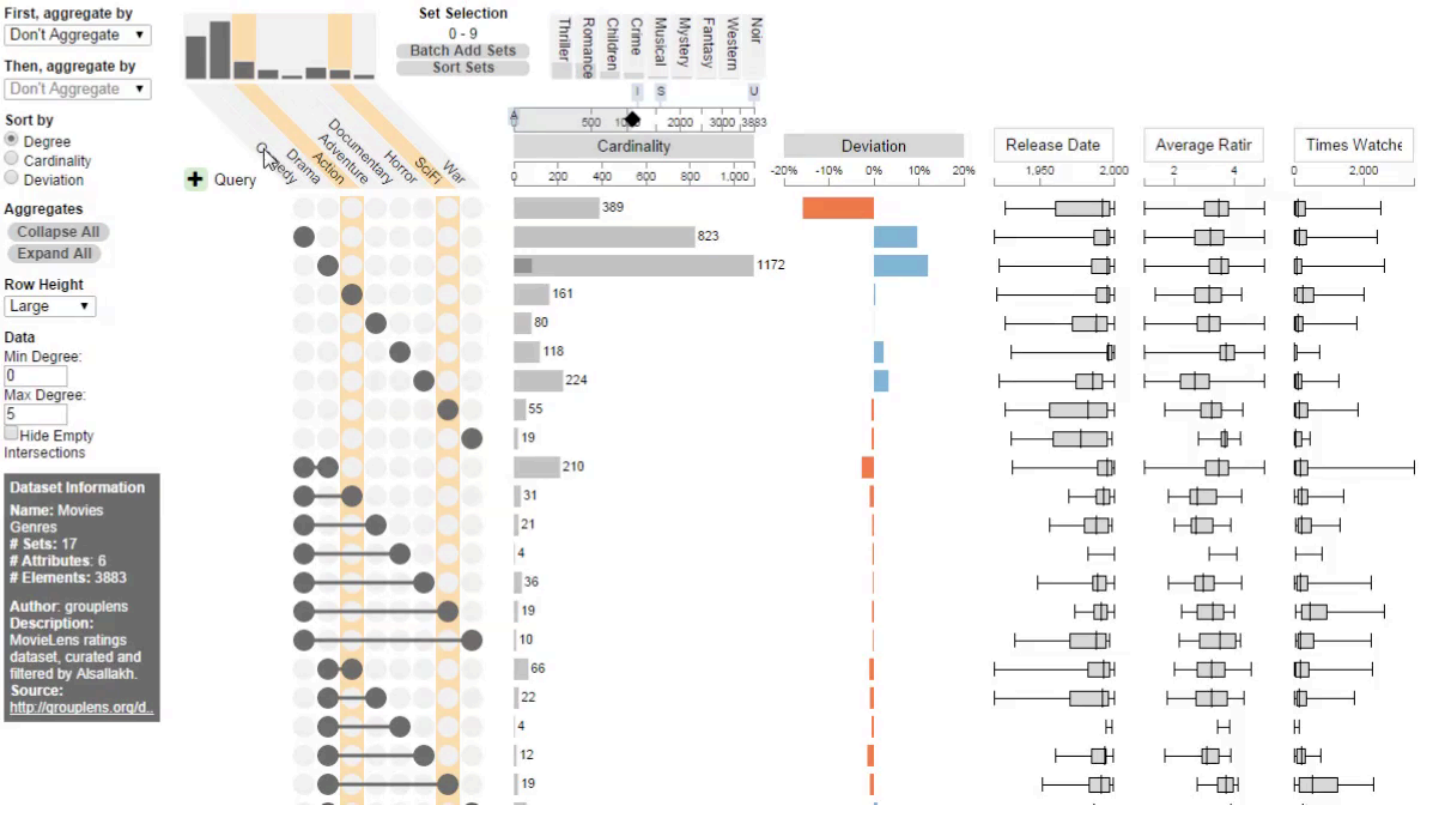


Sorting

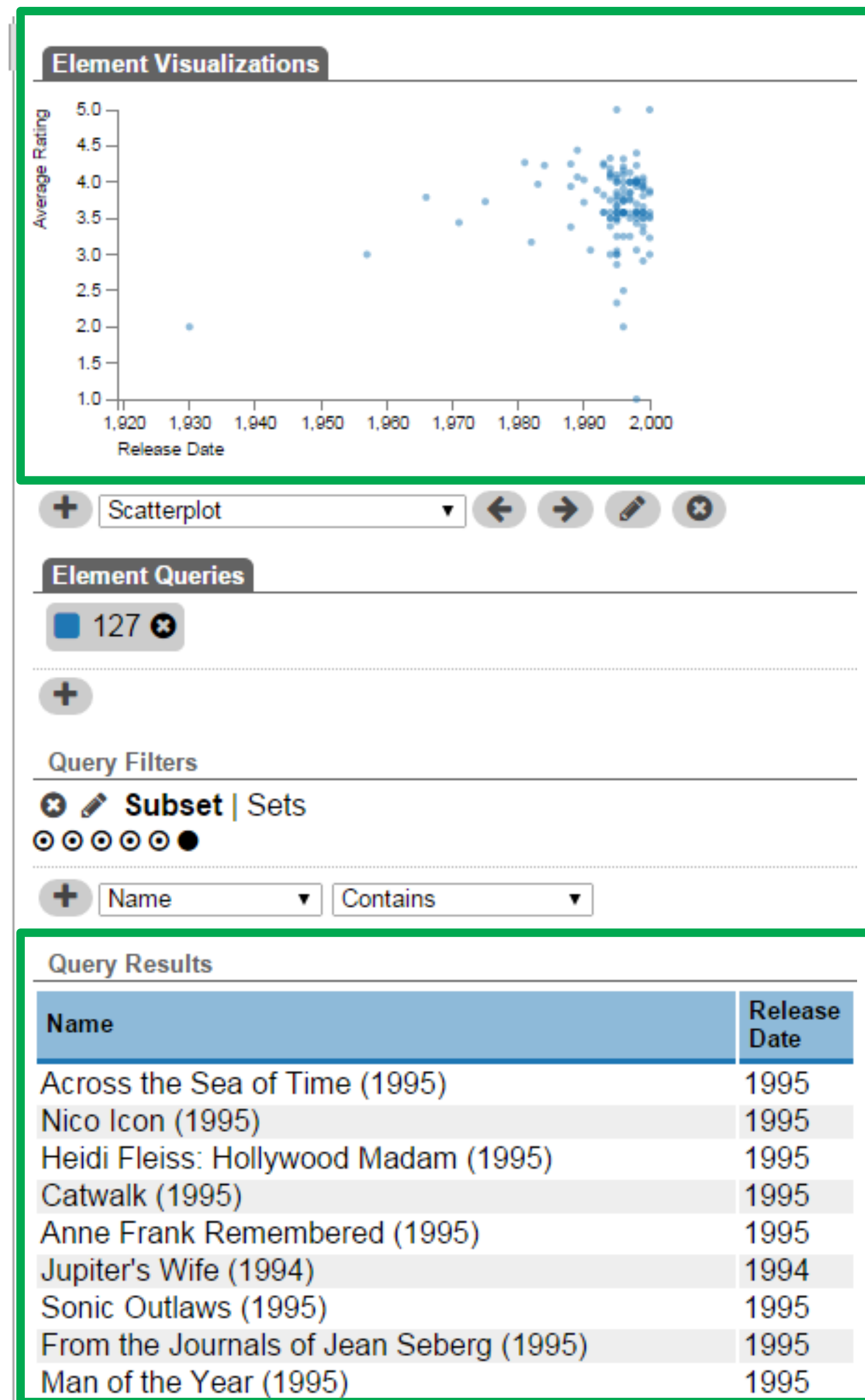
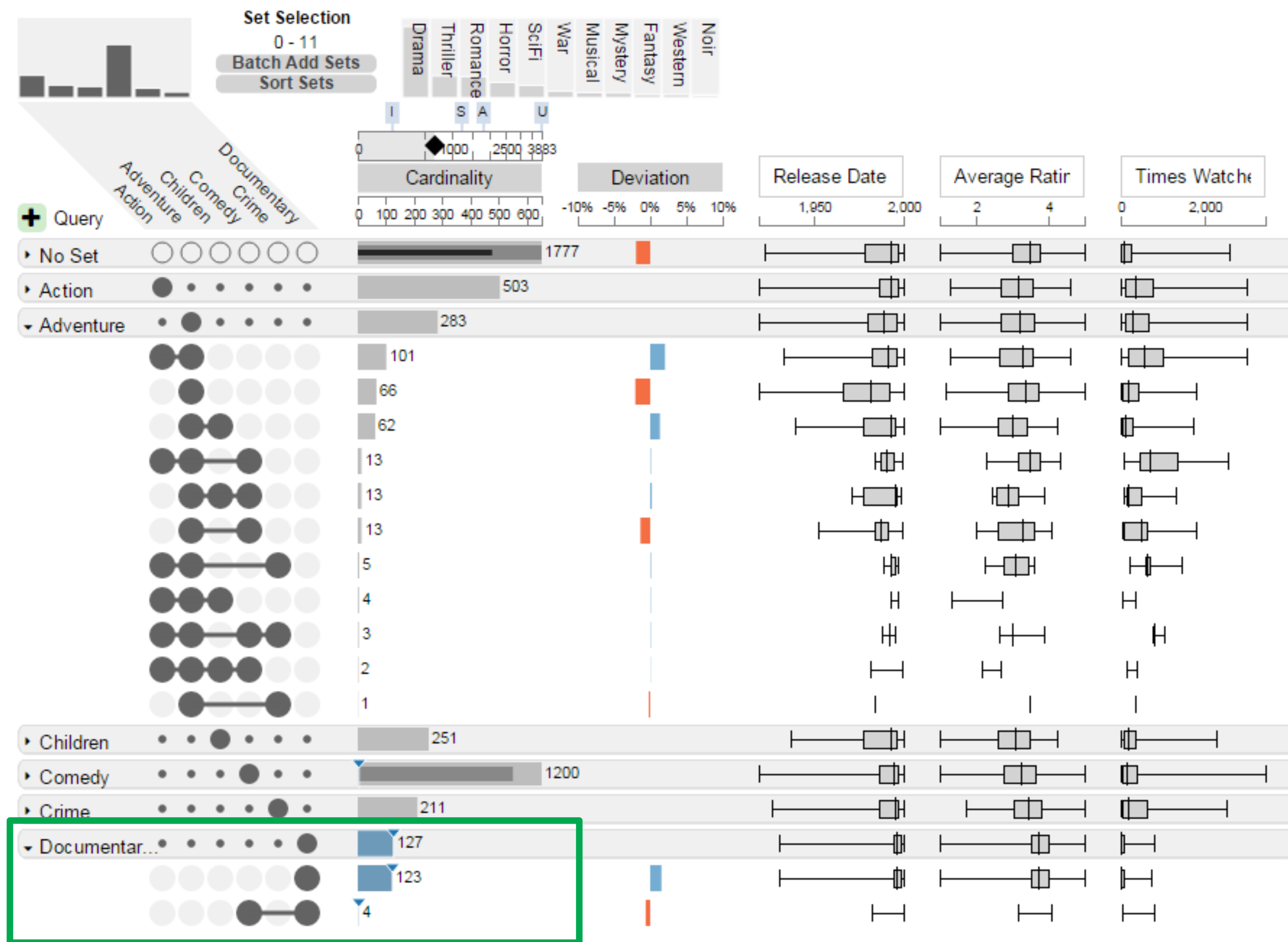


Which is the biggest intersection?
Sort By: Cardinality

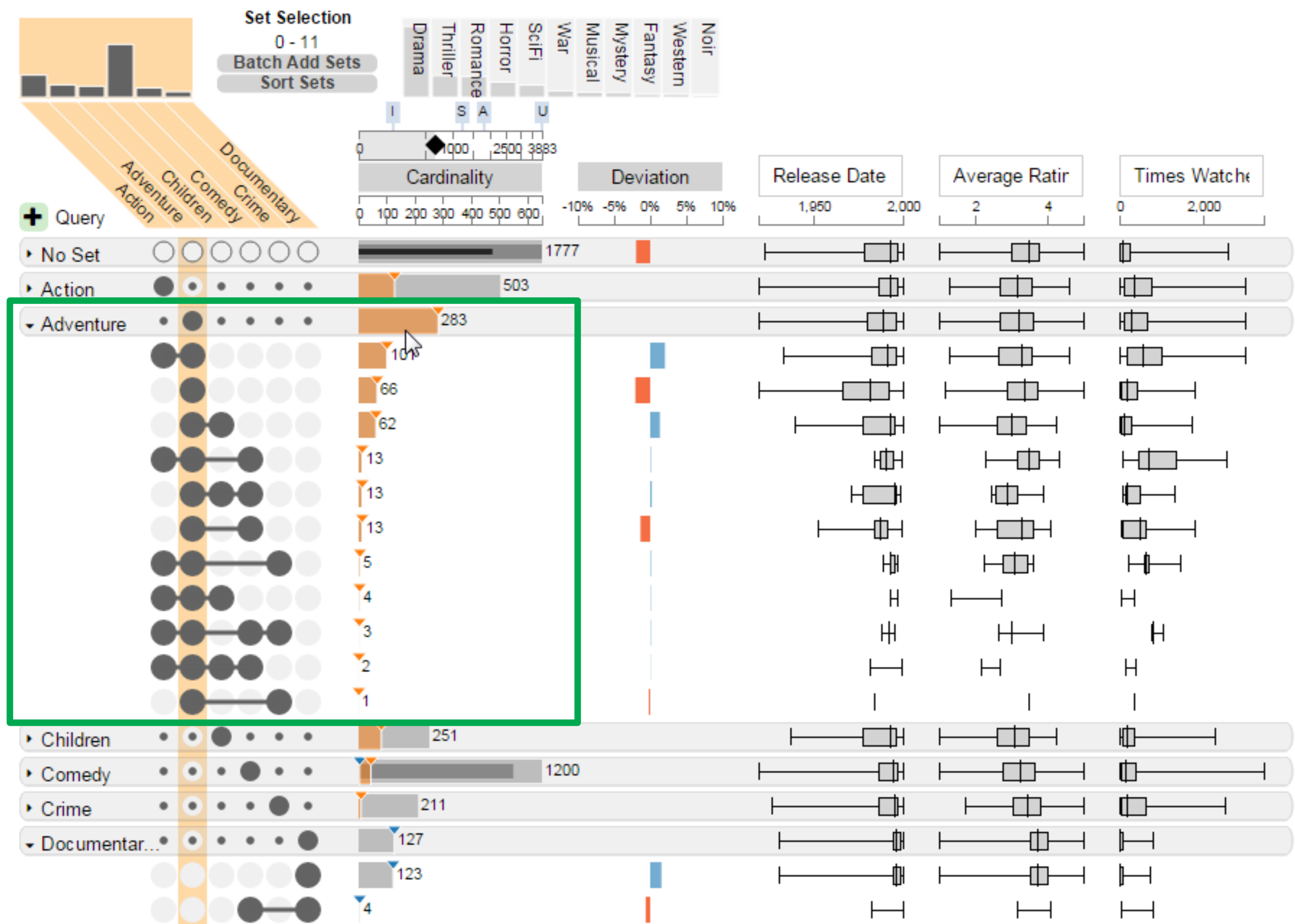




Elements & Attributes



How do documentaries compare to adventure movies?



Scatterplot

Element Queries
127 283

Query Filters
Subset | Sets

Query Results

Name	Release Date
Jumanji (1995)	1995
Tom and Huck (1995)	1995
GoldenEye (1995)	1995
Cutthroat Island (1995)	1995
City of Lost Children, The (1995)	1995
Wings of Courage (1995)	1995
Mortal Kombat (1995)	1995
Kids of the Round Table (1995)	1995
Indian in the Cupboard, The (1995)	1995
White Squall (1996)	1996
Muppet Treasure Island (1996)	1996

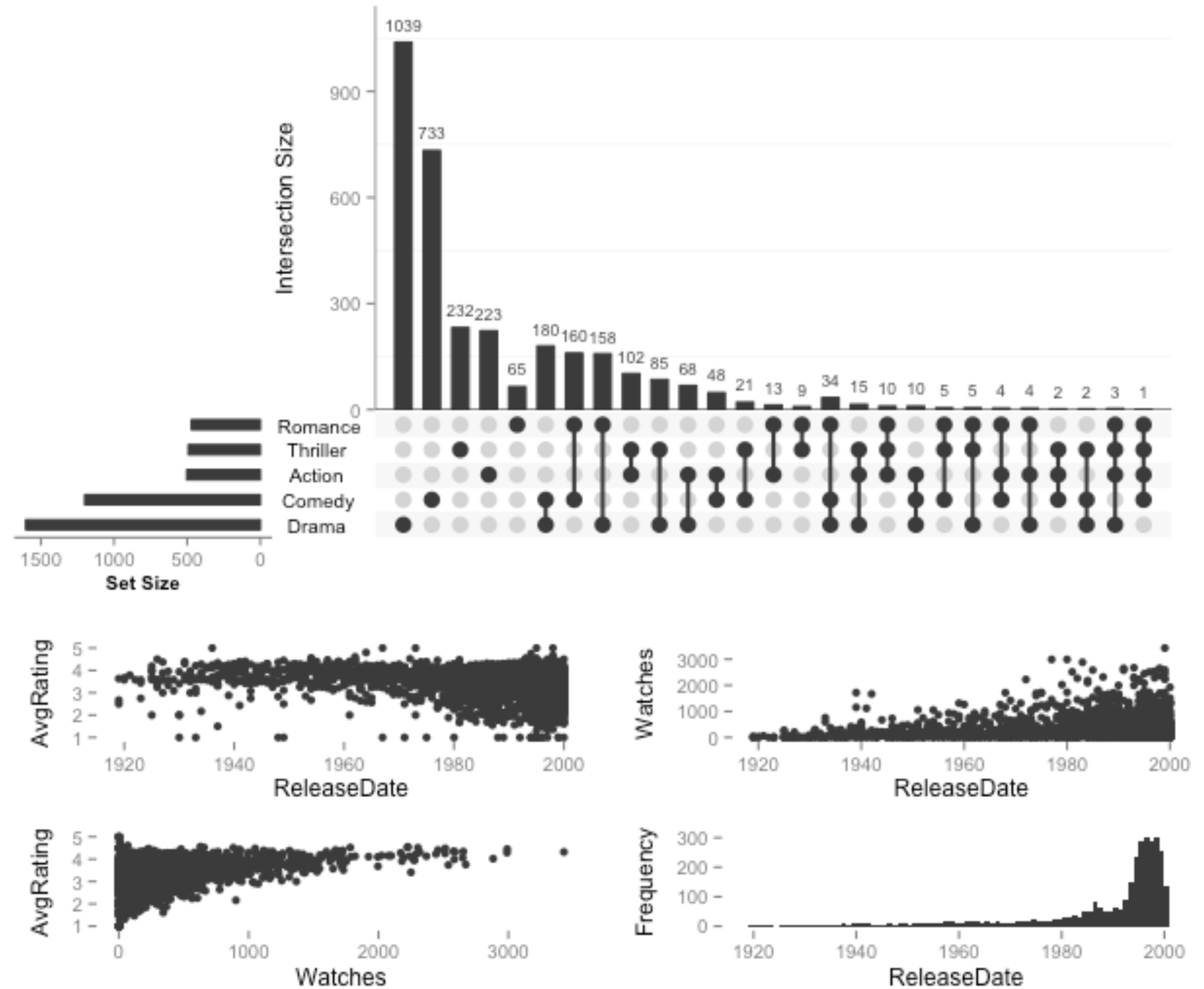
How do documentaries compare to adventure movies?

Applications

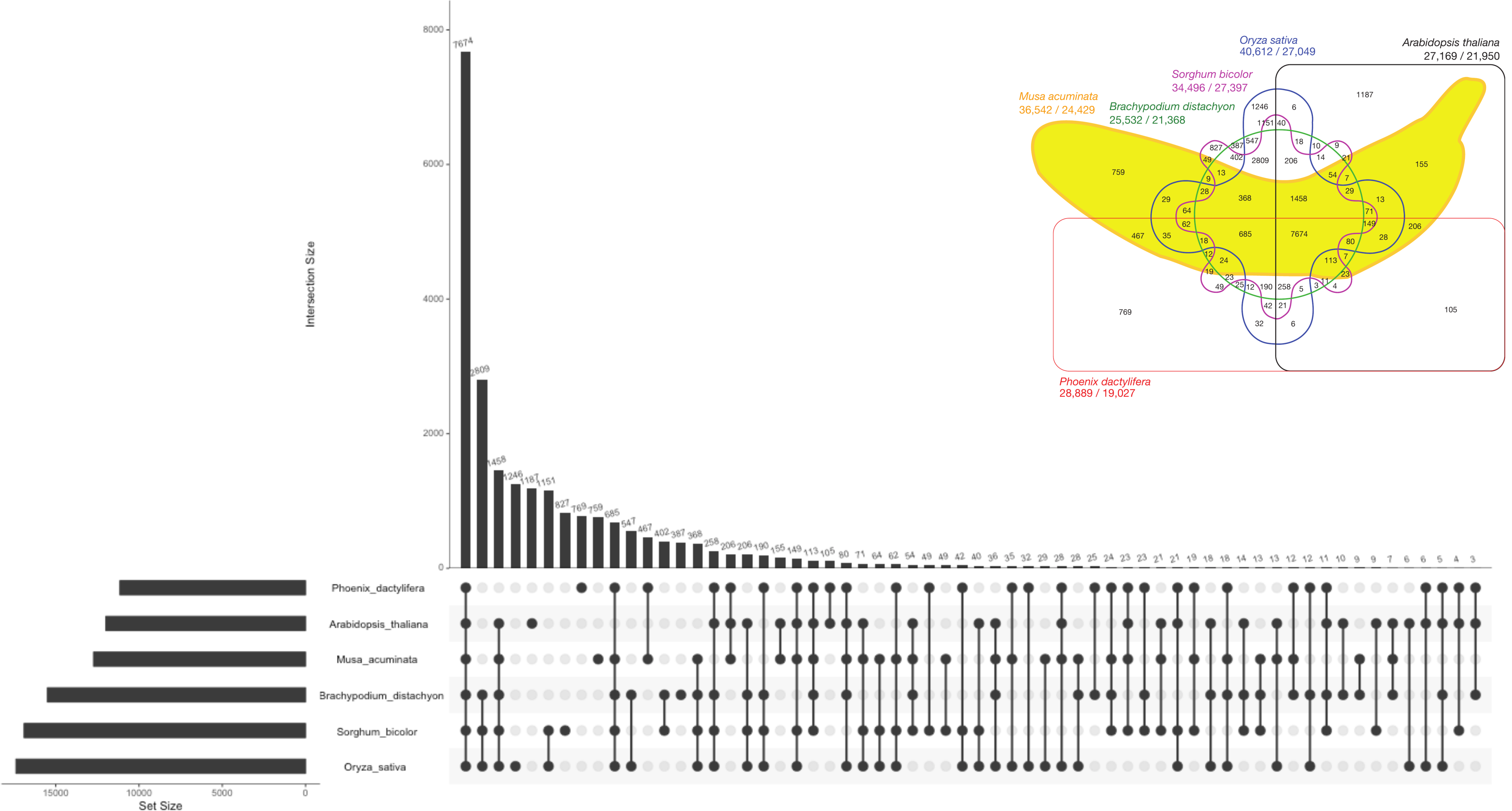
R-Version: UpSetR

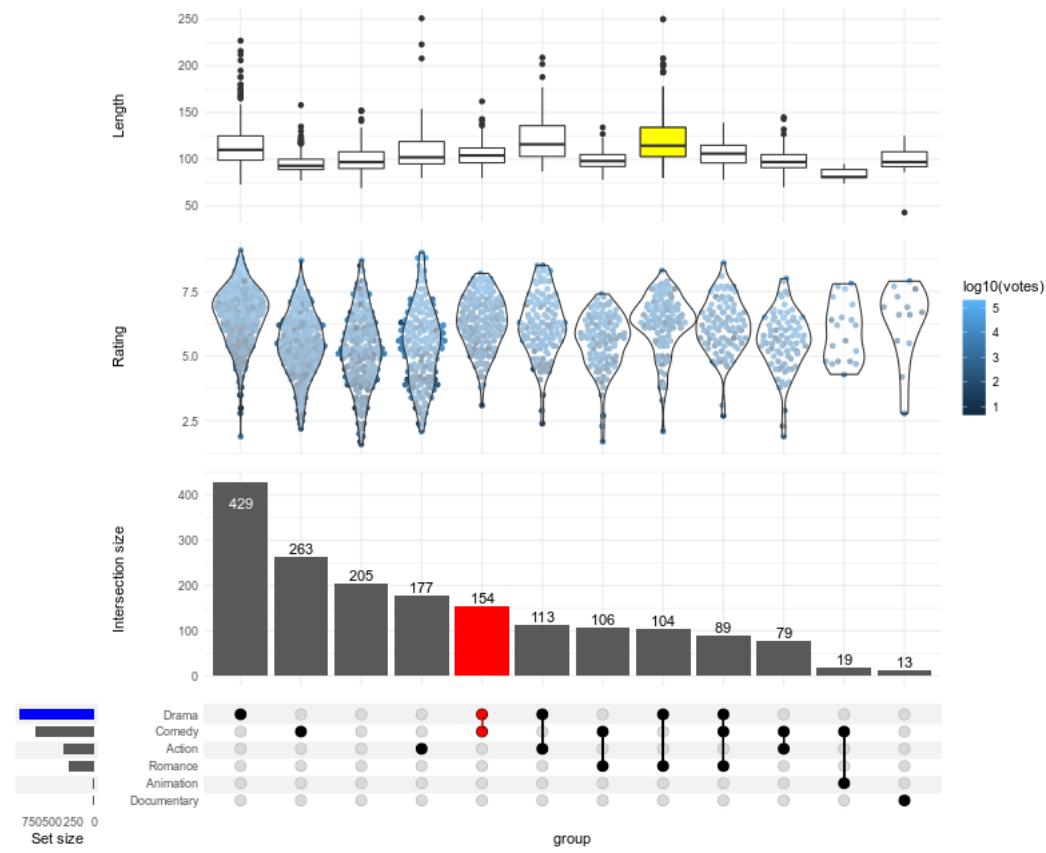
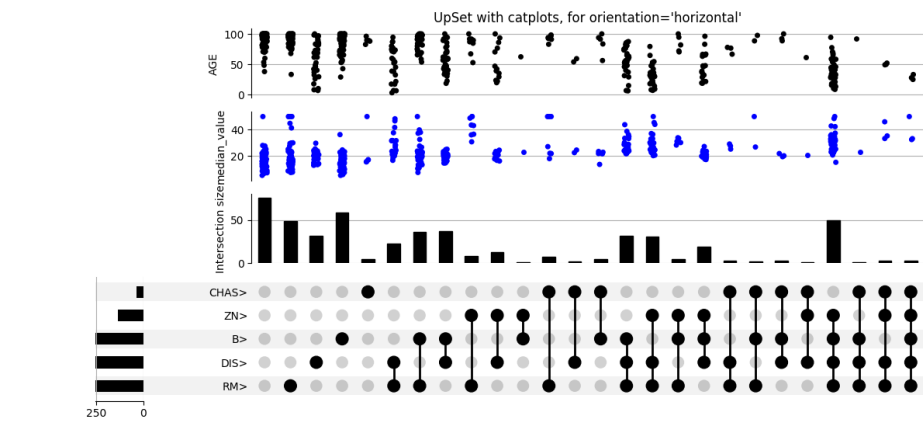
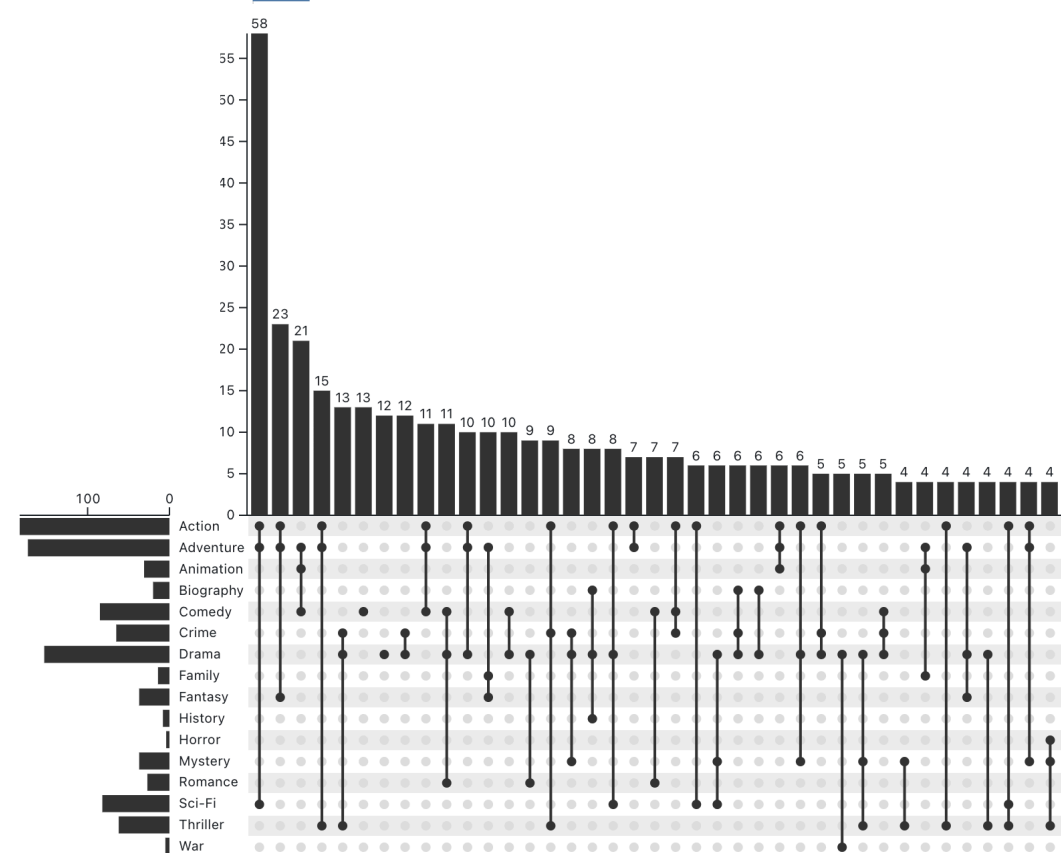
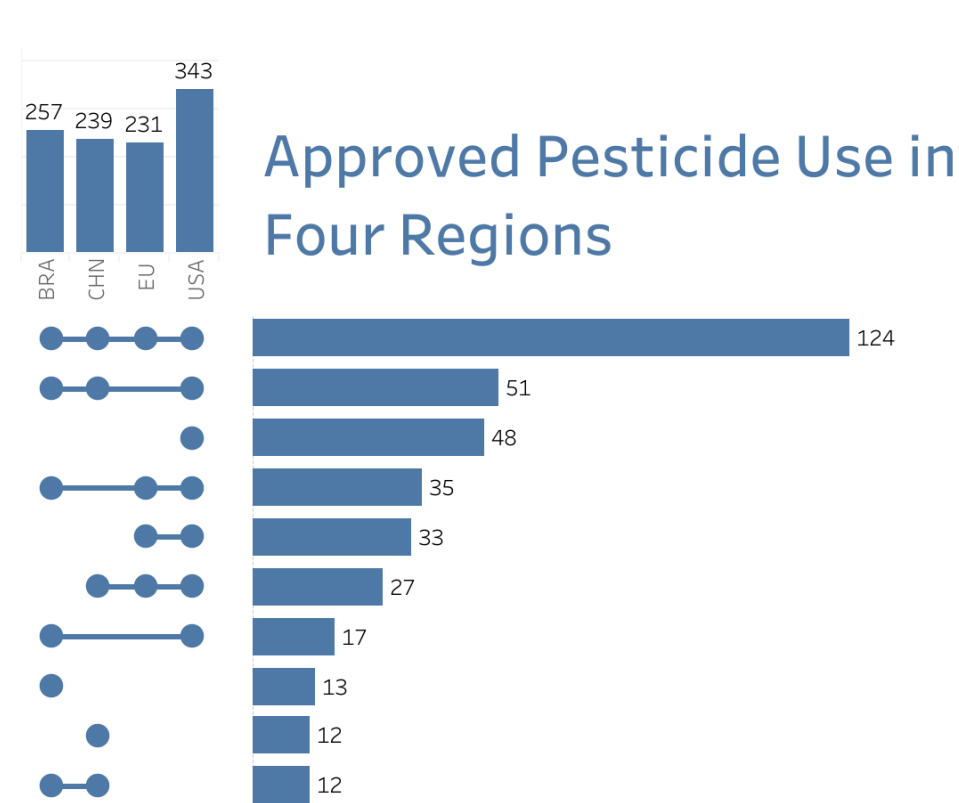
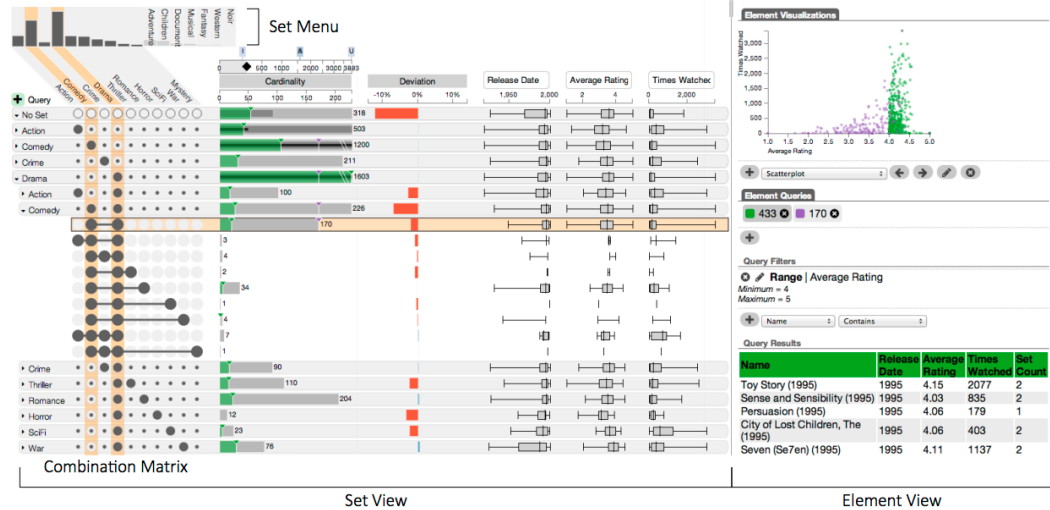
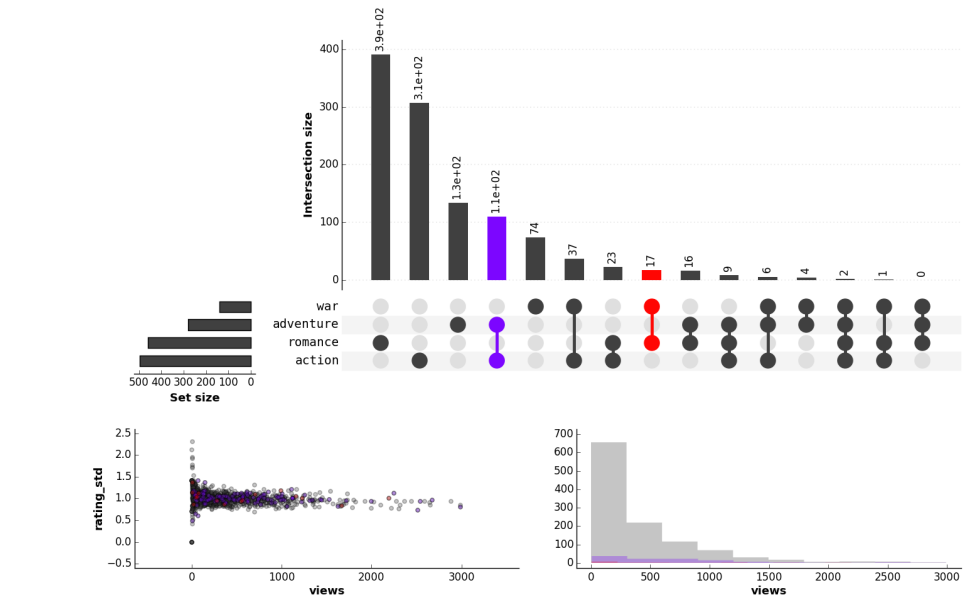
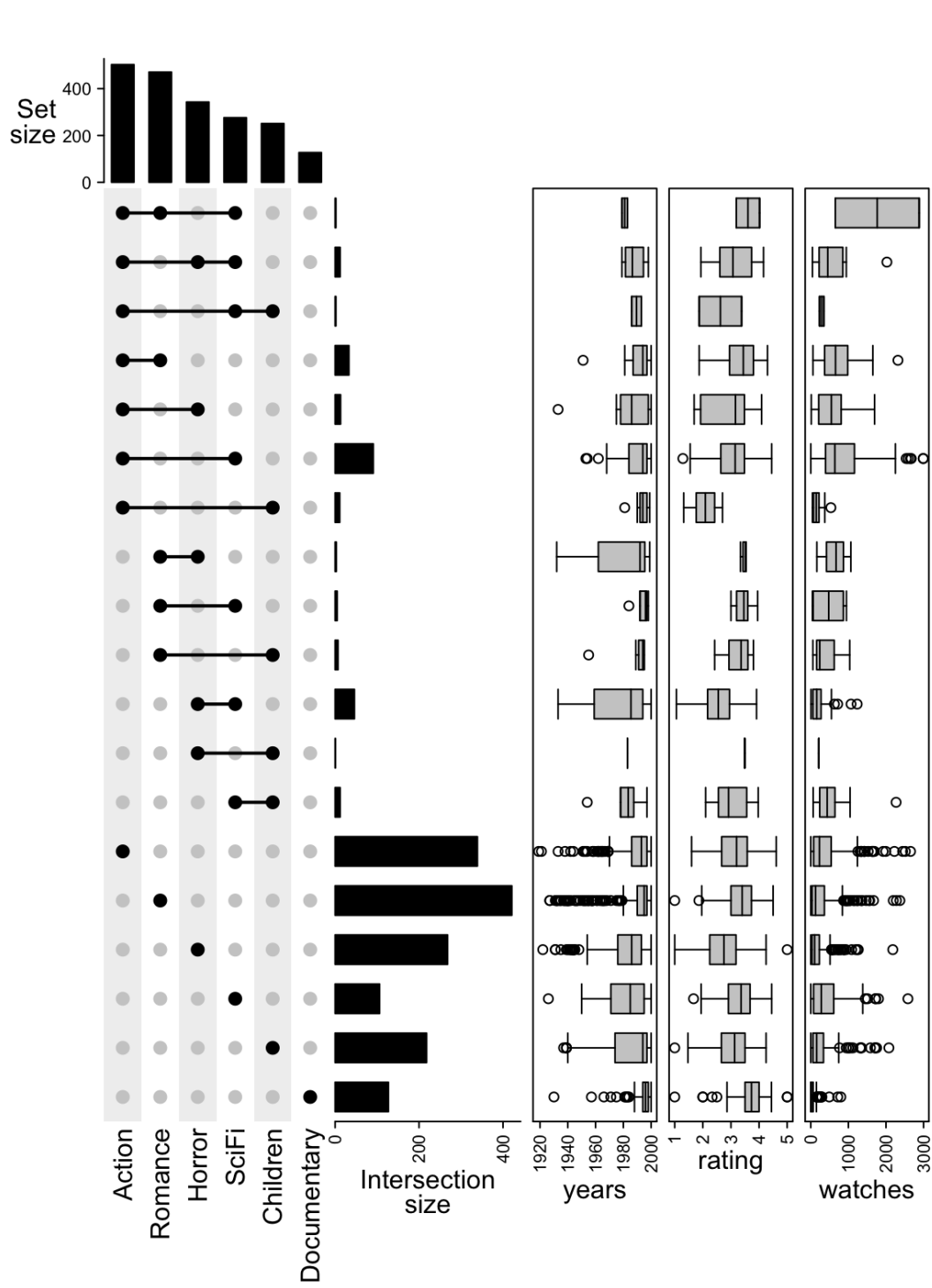
Developed at HMS

Some design adaptations



The Banana Chart Redesigned

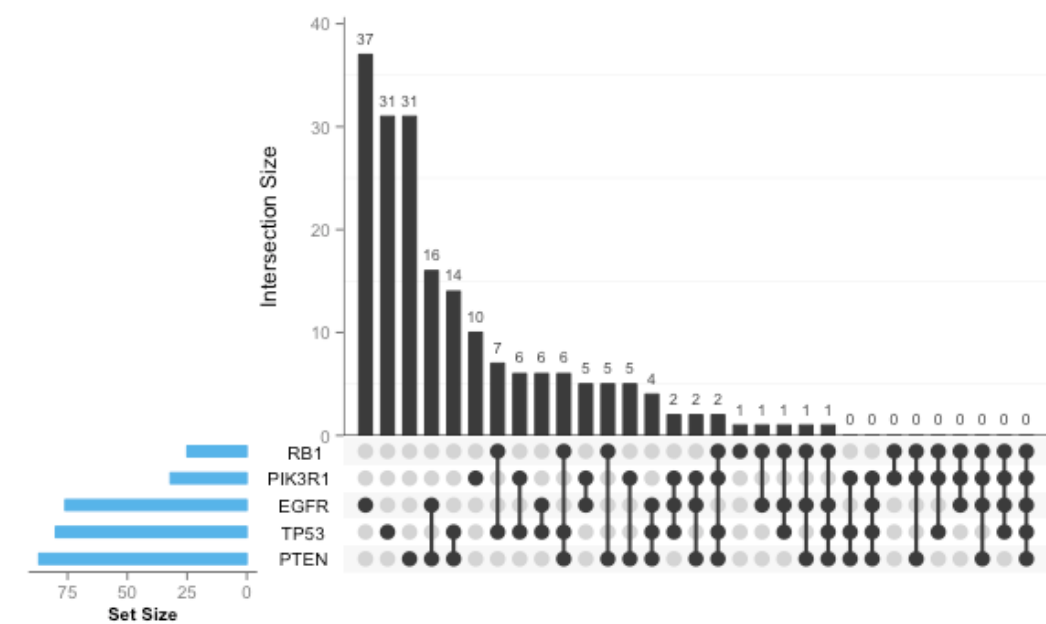
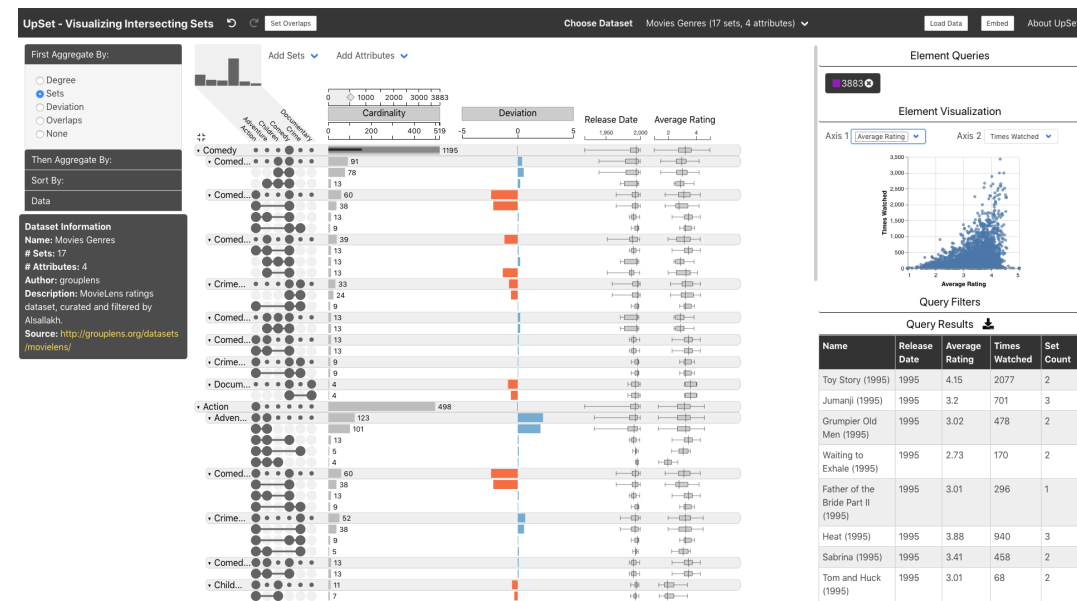
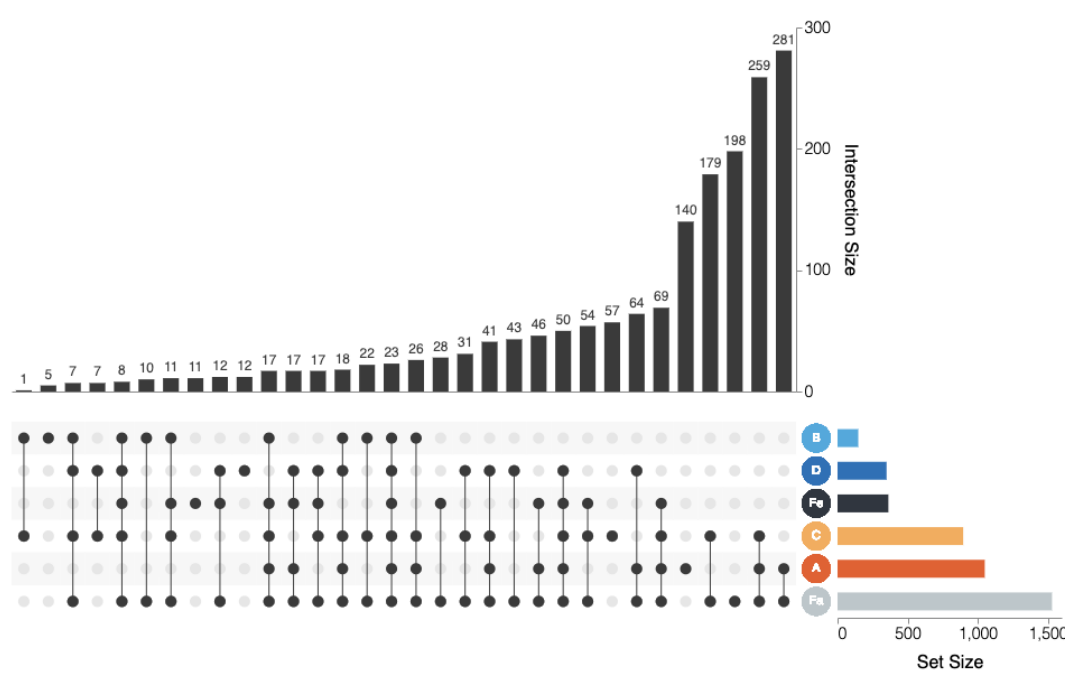
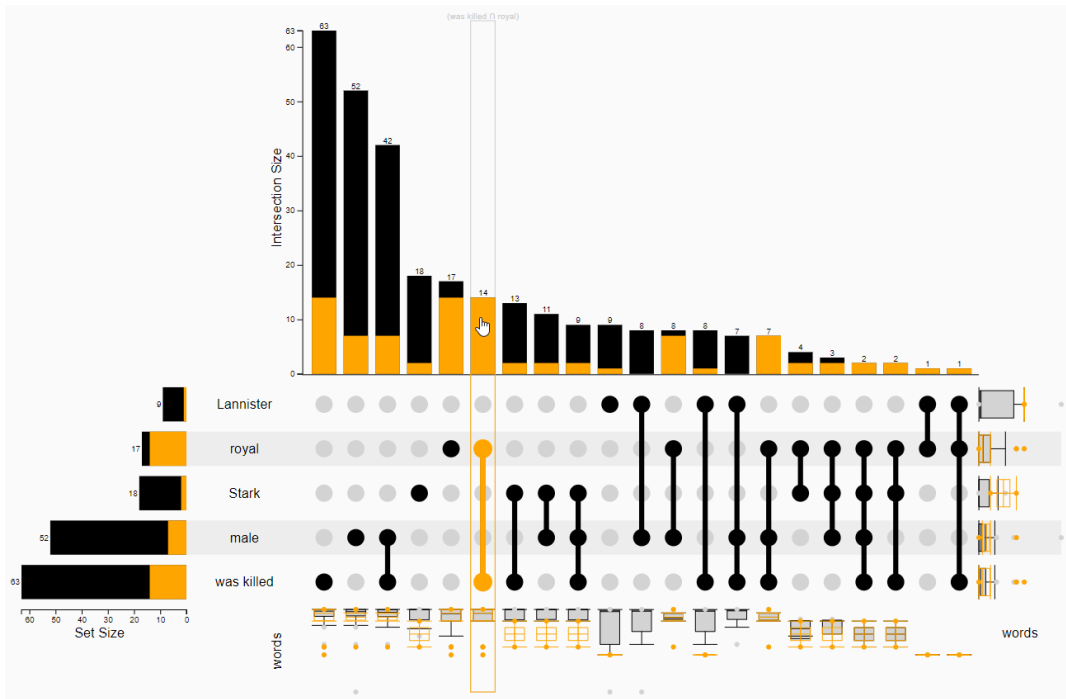




Symptoms Reported by Users of the COVID Symptom Tracker App

Story & Data: <https://www.nature.com/articles/s41586-020-00154-w>
Altair-based UpSet Plot: <https://github.com/hms-dbmi/upset-altair-notebook>

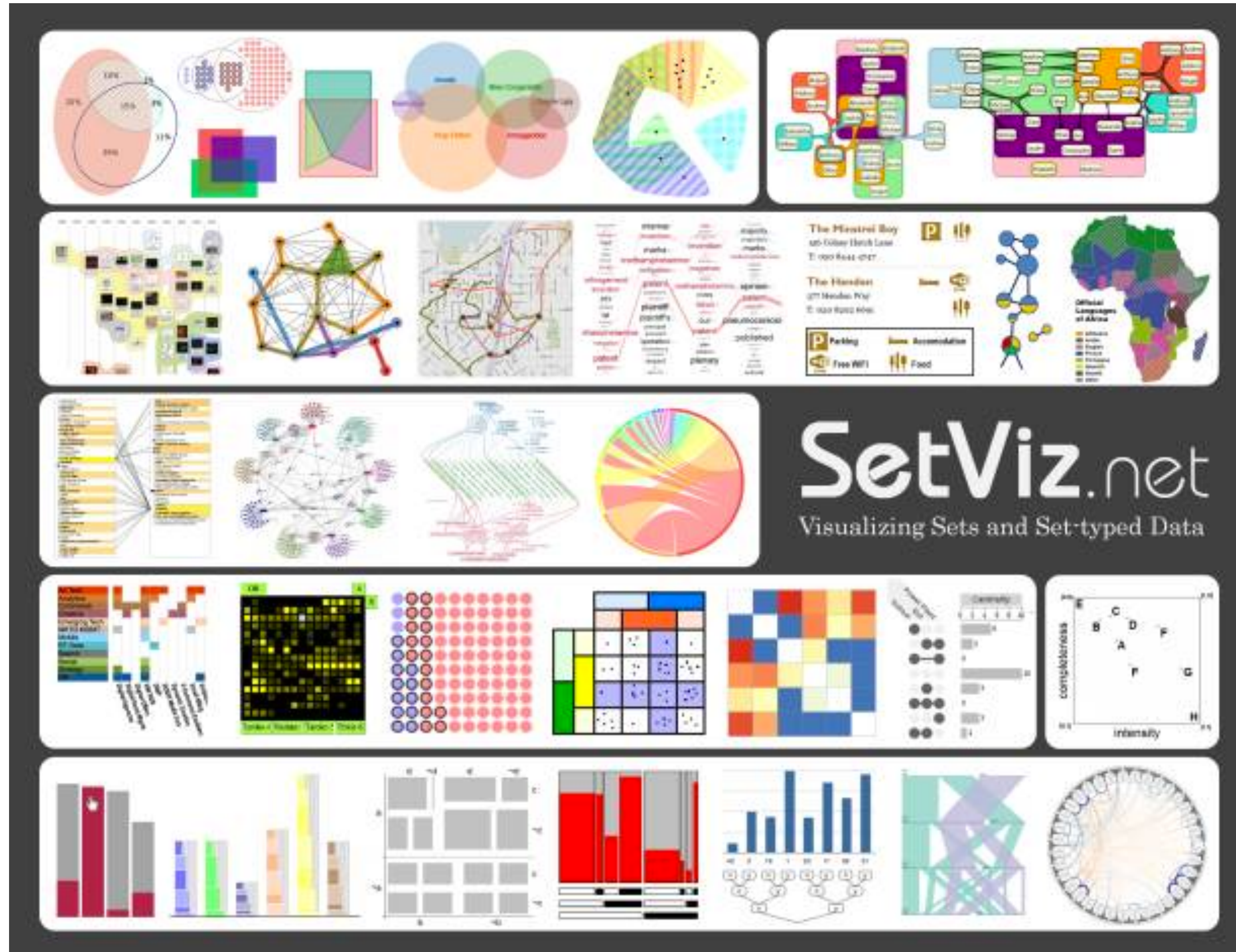
● Shortness of Breath ● Diarrhea ● Fever ● Cough ● Anosmia ● Fatigue



DESIGN 2

	Male	Duffan	Bluehair	Values
Matrix 1	●	●	0	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> 3
	●	0	0	<input type="checkbox"/> <input type="checkbox"/> 2
	0	0	0	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> 2
	0	0	●	<input type="checkbox"/> 1
	0	●	0	0
	0	●	●	0
	●	0	●	0
	●	●	●	0
Matrix 2	School	Efil	Power Plant	
	0	0	0	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> 4
	●	0	0	<input type="checkbox"/> <input type="checkbox"/> 2
	0	0	●	<input type="checkbox"/> 1
	0	●	●	<input type="checkbox"/> 1
	0	0	●	0
	●	●	0	0
	●	0	●	0
	●	●	●	0

Other Options



<http://setviz.net>