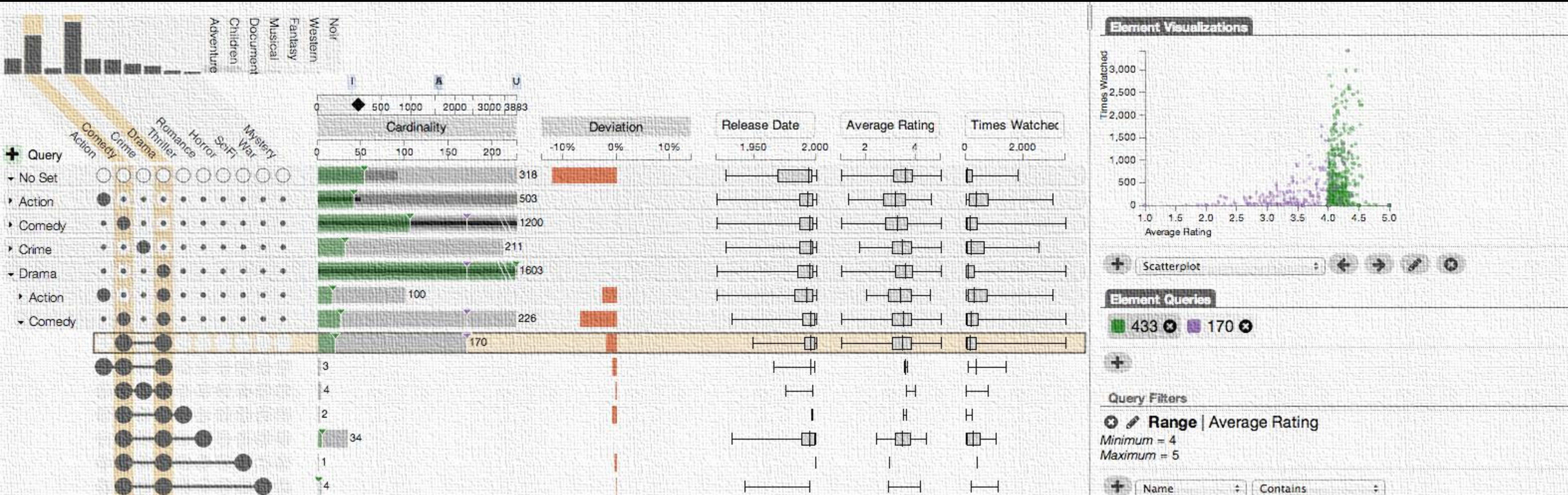




Visual Data Analysis for Biology & Pharmacology



visualization

pictures

The purpose of computing is insight, not numbers.

- Richard Wesley Hamming, 1962

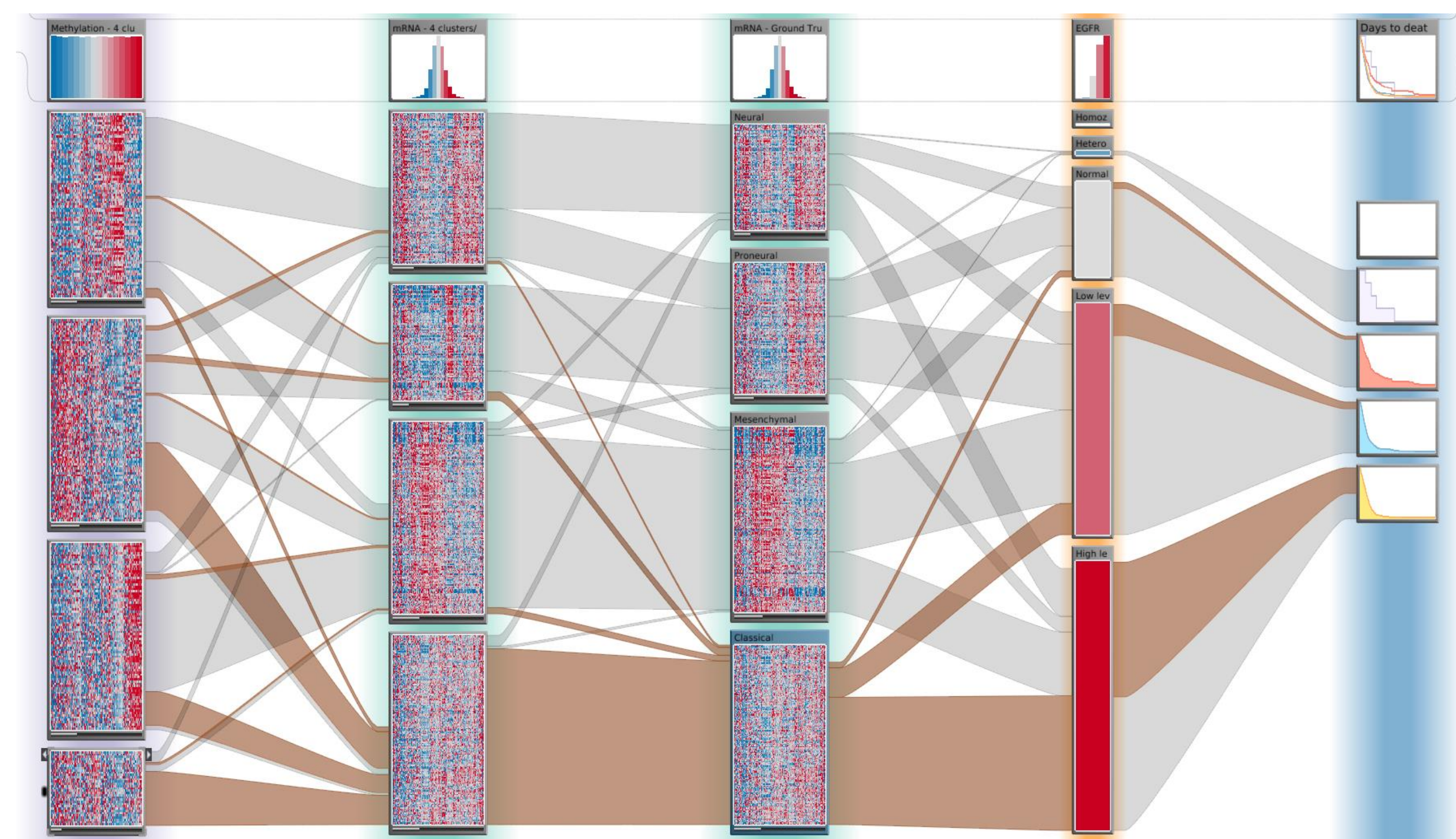


[D'Hont et al., Nature, 2012]

Good Data Visualization

- ... makes data **accessible**
- ... combines strengths of
humans and computers
- ... enables **insight**
- ... **communicates**

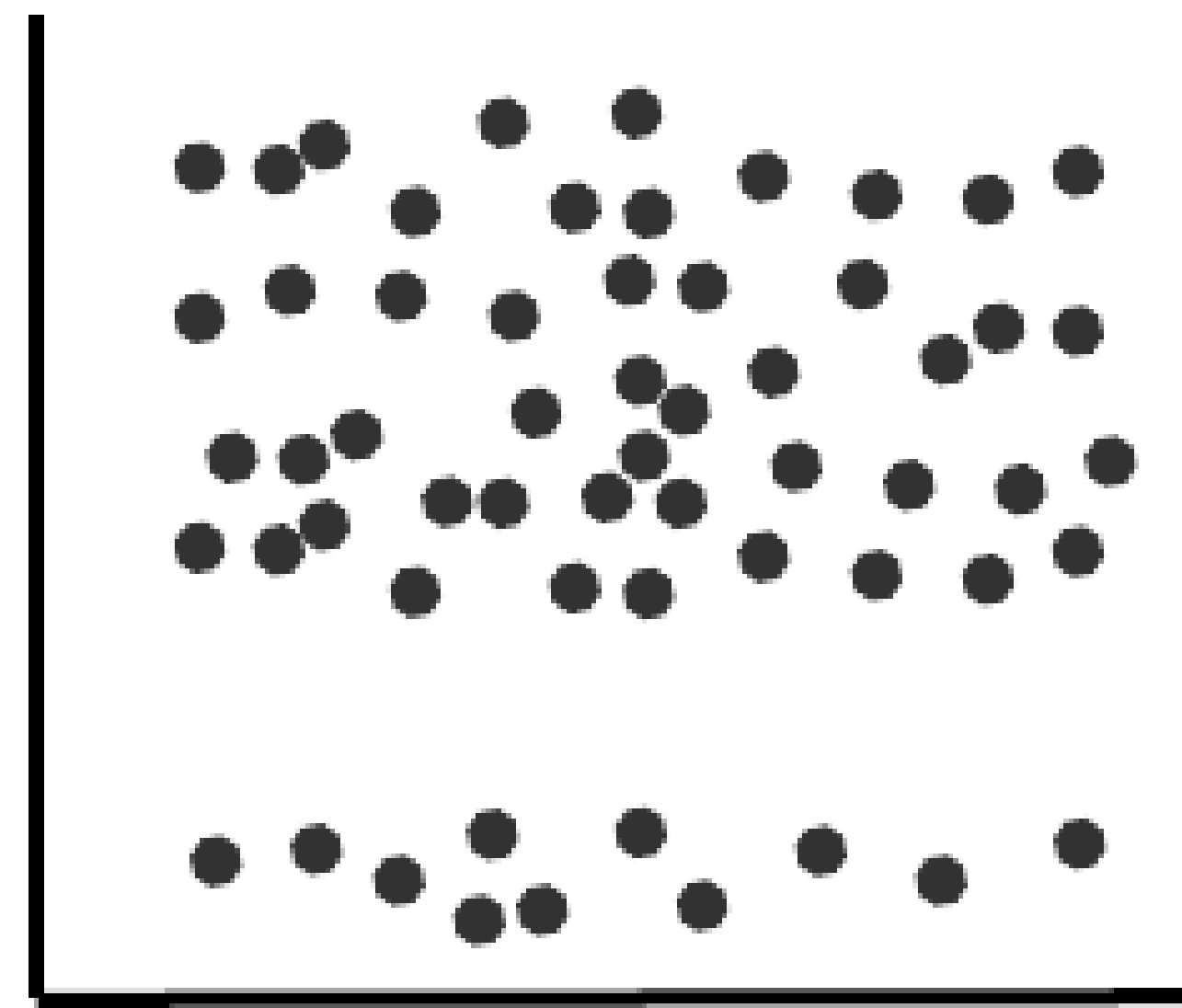
Purpose of Visualization



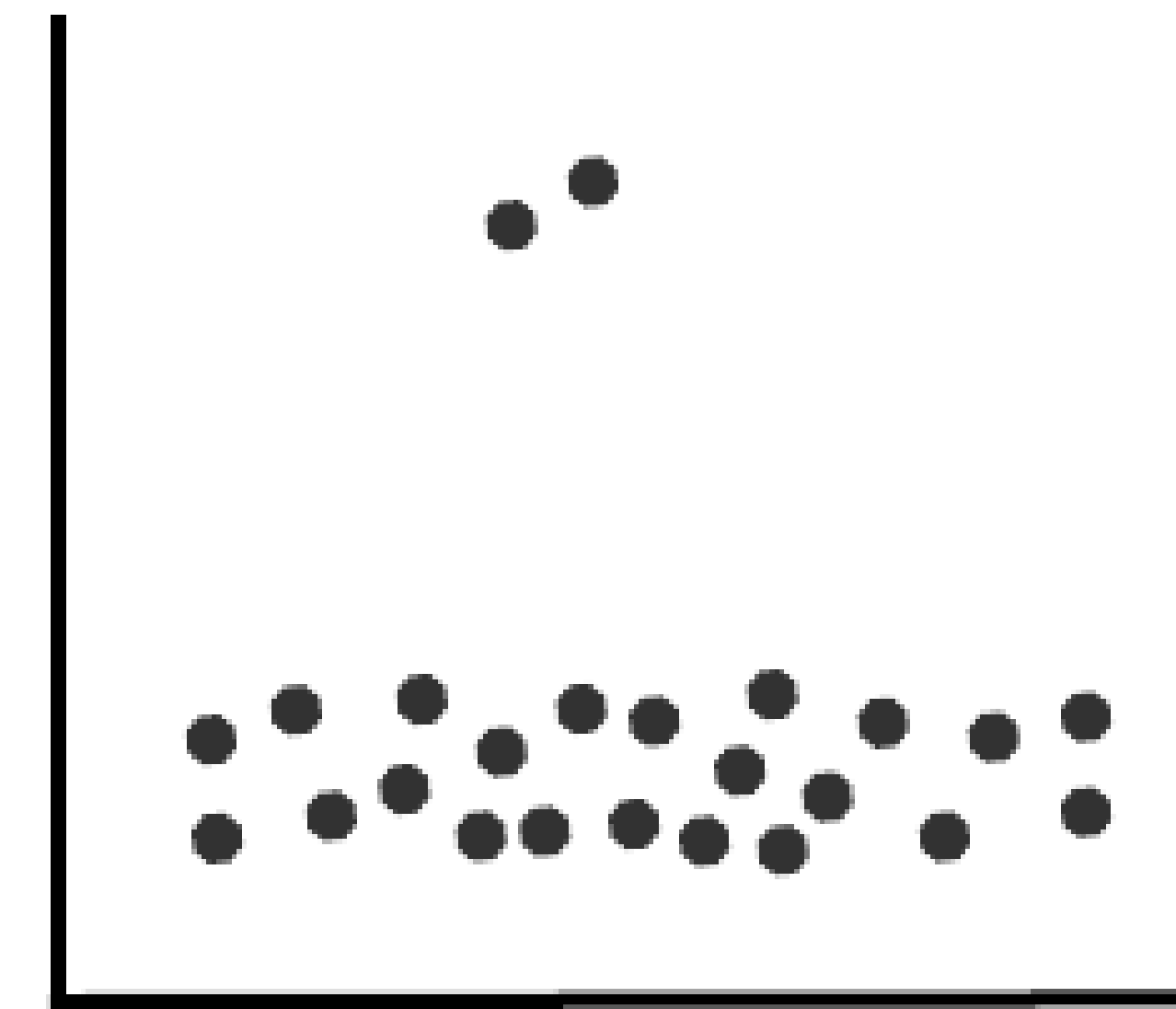
Visualization for Pattern Discovery



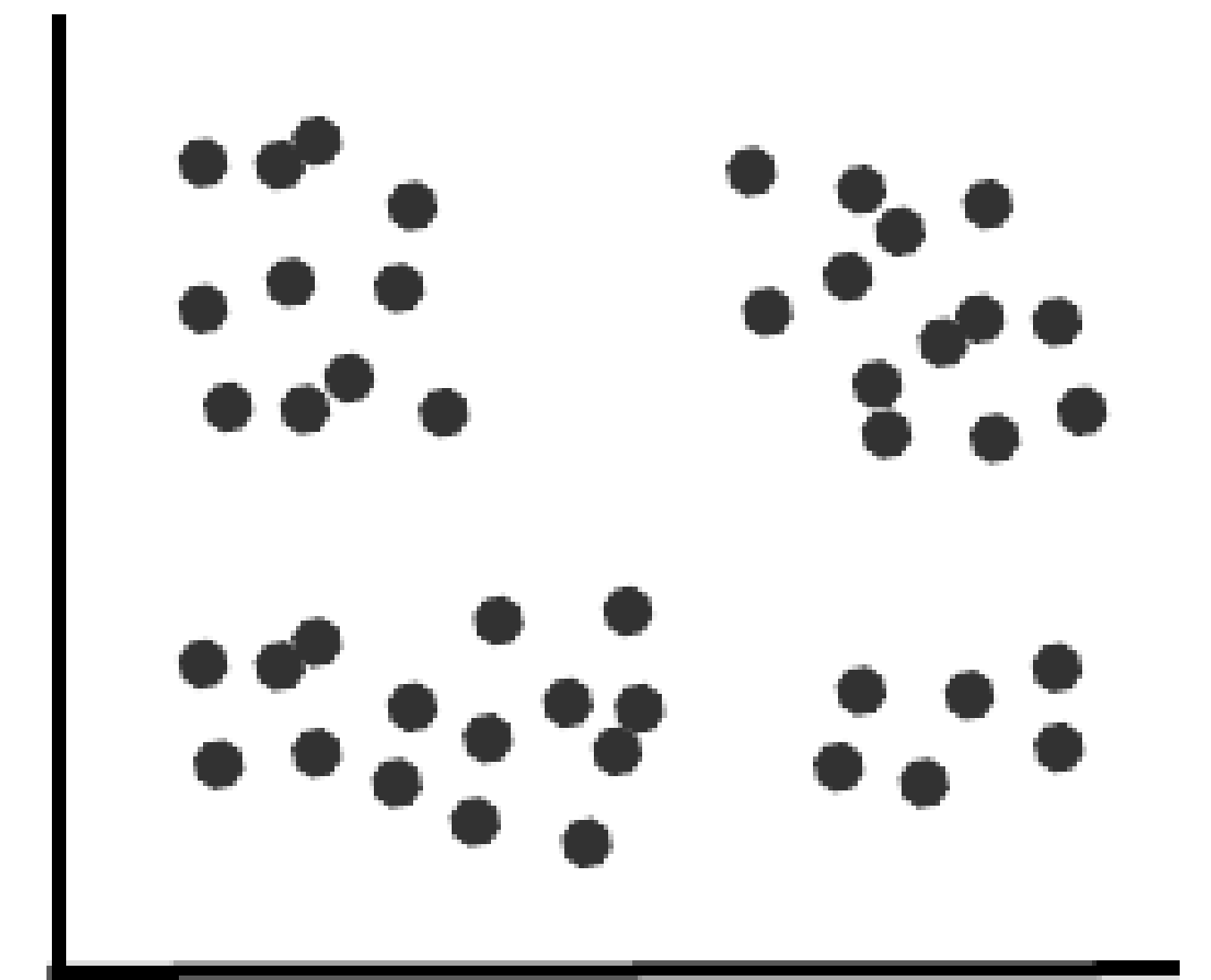
trends



gaps



outliers



clusters

Can We Trust Statistics?

I		II		III		IV	
x	y	x	y	x	y	x	y
10	8.04	10	9.14	10	7.46	8	6.58
8	6.95	8	8.14	8	6.77	8	5.76
13	7.58	13	8.74	13	12.74	8	7.71
9	8.81	9	8.77	9	7.11	8	8.84
11	8.33	11	9.26	11	7.81	8	8.47
14	9.96	14	8.1	14	8.84	8	7.04
6	7.24	6	6.13	6	6.08	8	5.25
4	4.26	4	3.1	4	5.39	19	12.5
12	10.84	12	9.13	12	8.15	8	5.56
7	4.82	7	7.26	7	6.42	8	7.91
5	5.6					8	6.89

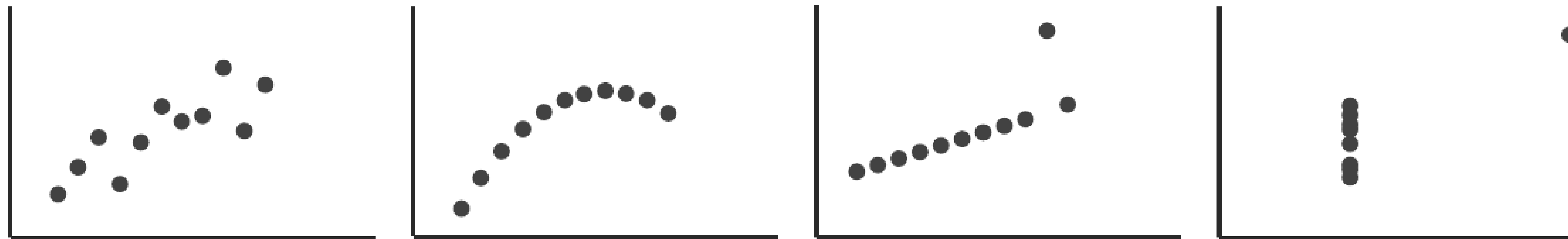
Mean x: 9 y: 7.50

Variance x: 11 y: 4.122

Correlation x – y: 0.816

Linear regression: $y = 3.00 + 0.500x$

Anscombe's Quartett



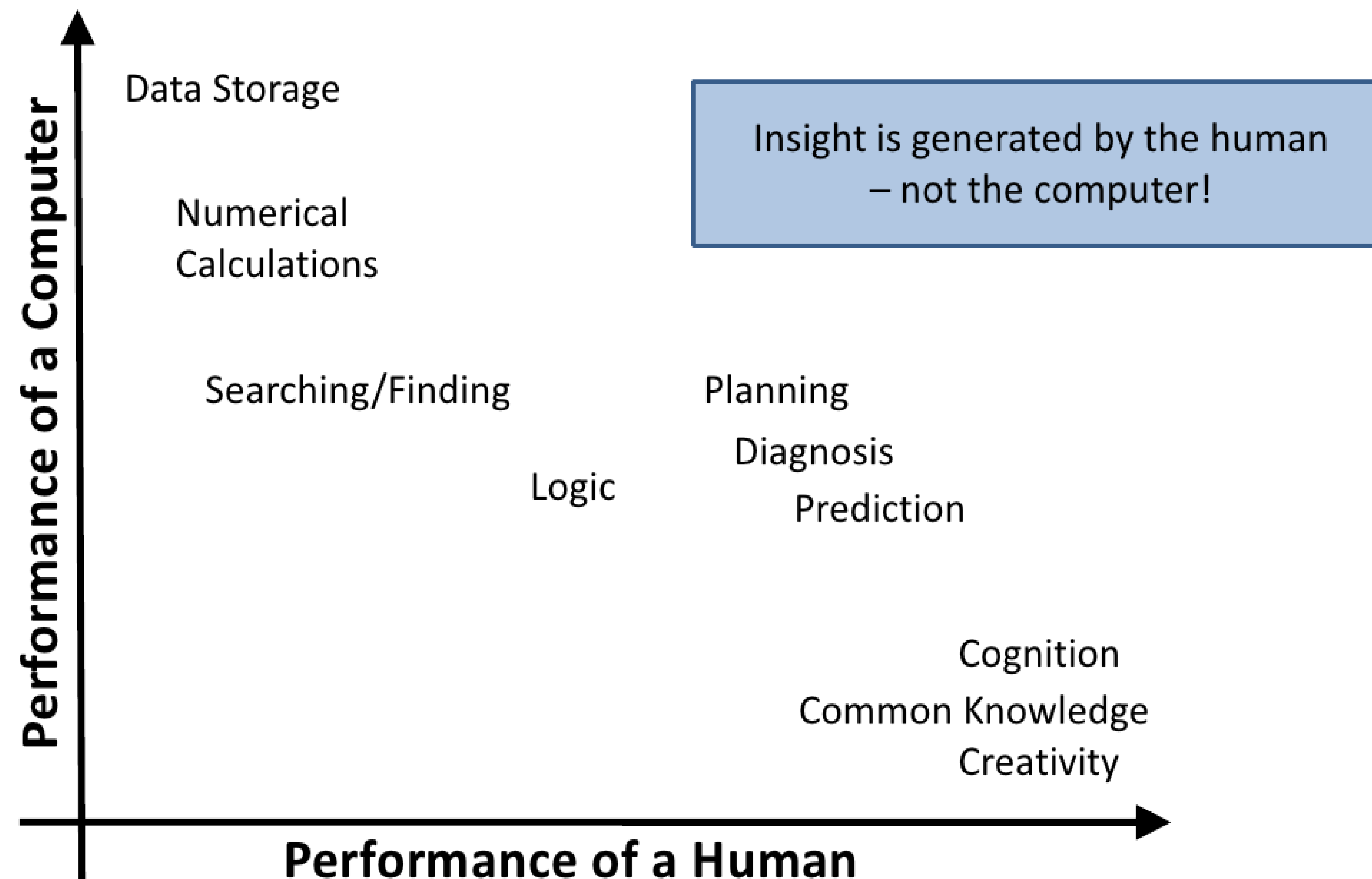
Mean x: 9 y: 7.50

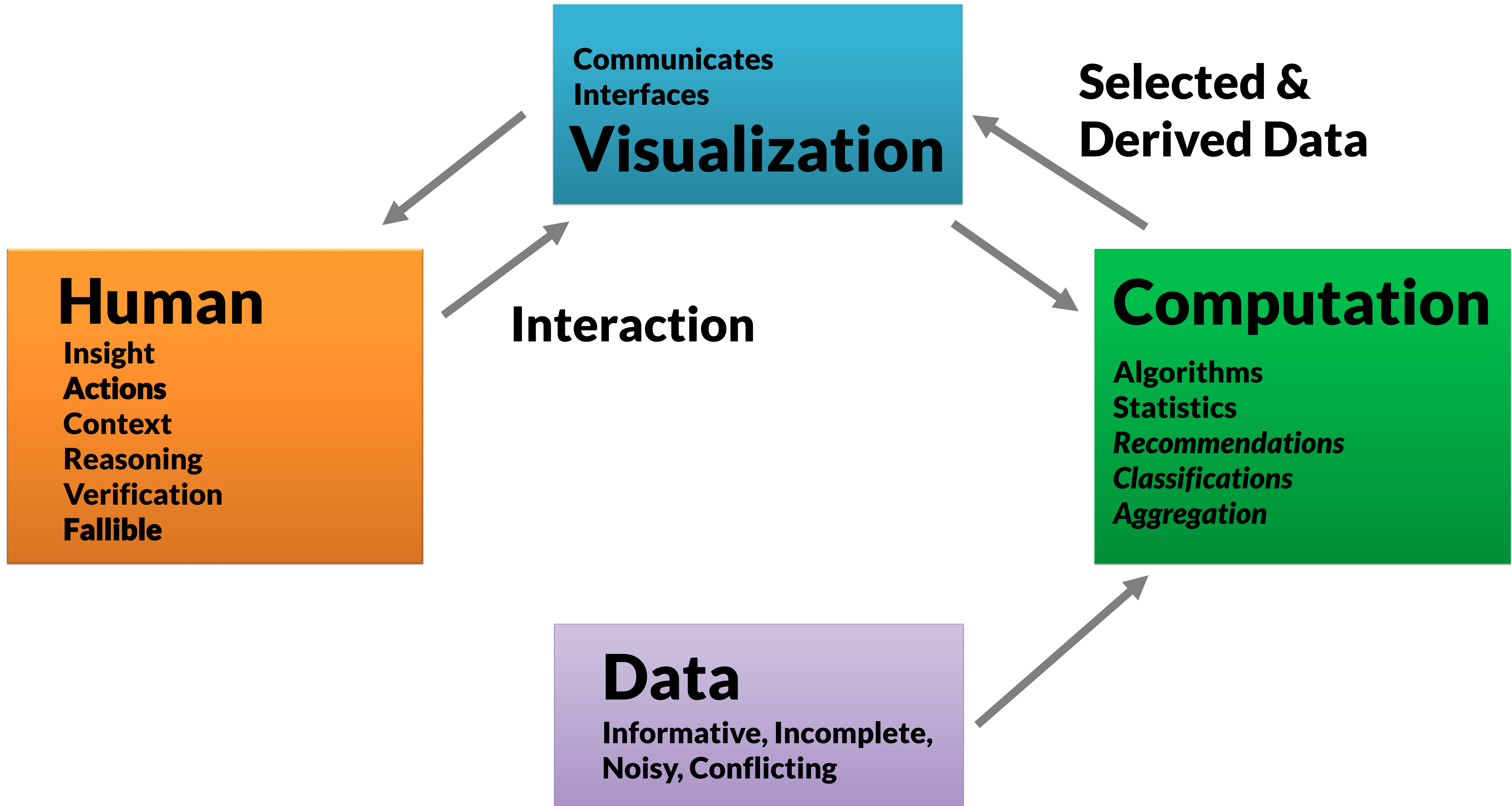
Variance x: 11 y: 4.122

Correlation x - y: 0.816

Linear regression: $y = 3.00 + 0.500x$

The Ability Matrix





Who am I?

alexander-lex.com

🐦 @alexander_lex

PostDoc & Lecturer @ Harvard,
Visual Computing Group,
PI: Hanspeter Pfister

Associated PostDoc @ Novartis
Institute of BioMedical Research,
Scientific Data Analysis,
PI: Mark Borowsky

How it looks like at home ↘



Credits

Marc Streit

Nils Gehlenborg

Christian Partl

Samuel Gratzl

Hendrik Strobelt

Dieter Schmalstieg

Mark Borowsky

Hanspeter Pfister

Johannes Kepler University Linz, AT

Harvard Medical School, Boston, USA

Graz University of Technology, AT

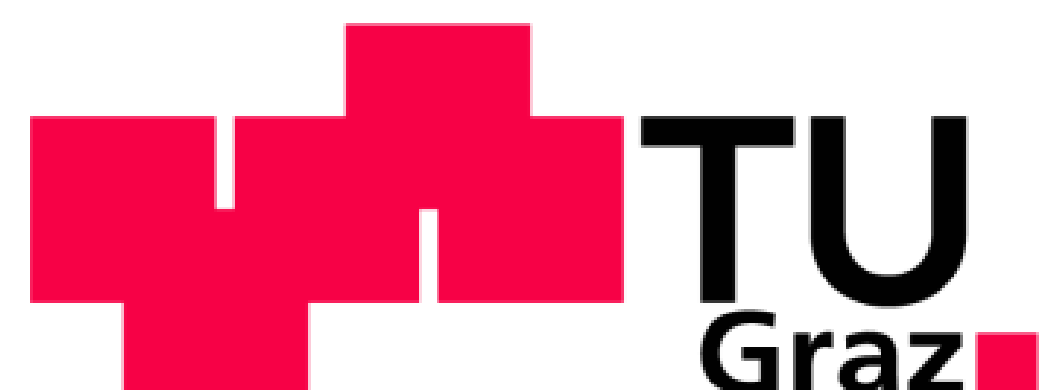
Johannes Kepler University Linz, AT

Harvard University, Cambridge, USA

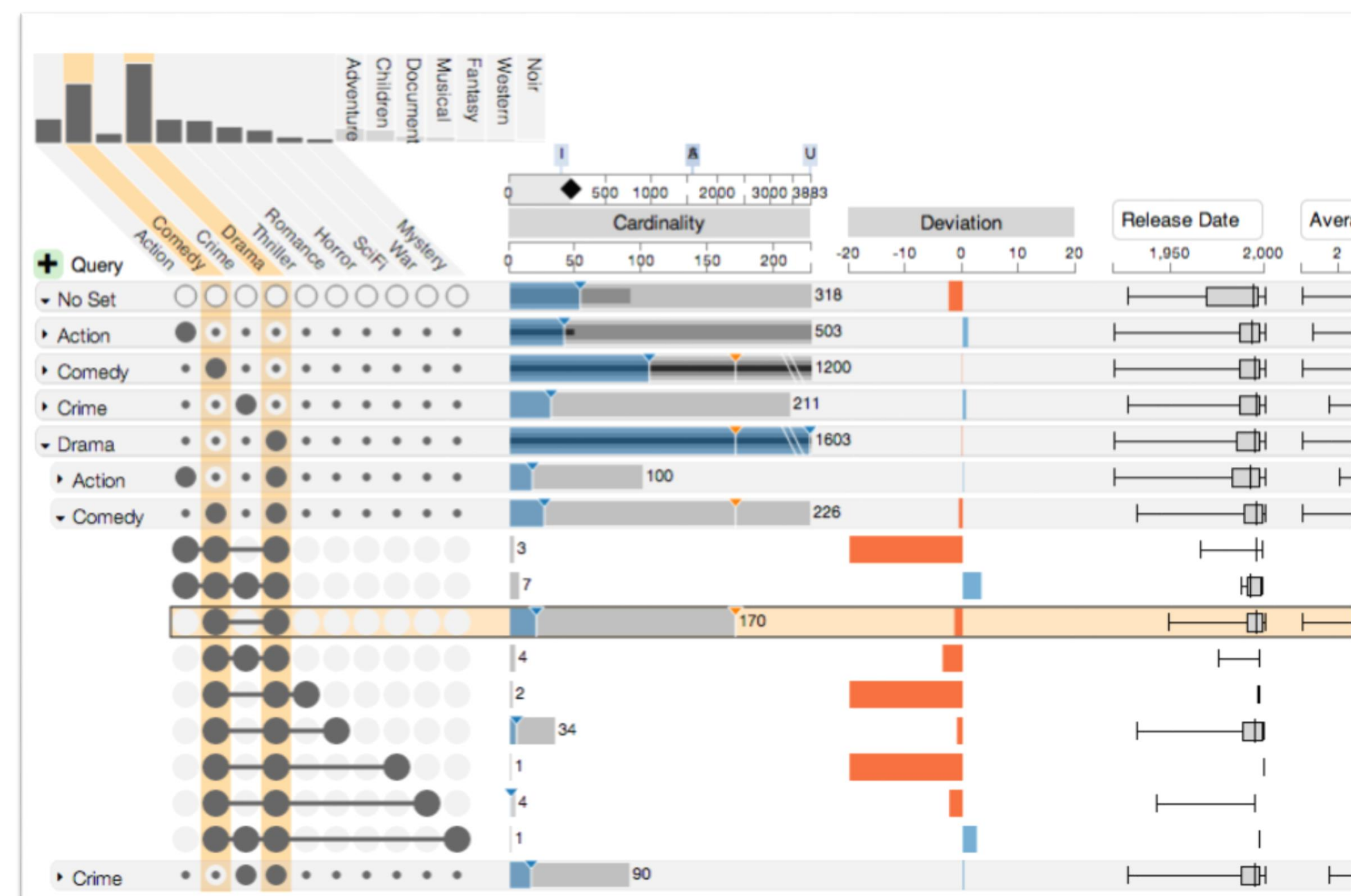
Graz University of Technology, AT

Novartis (NIBR), Cambridge, USA

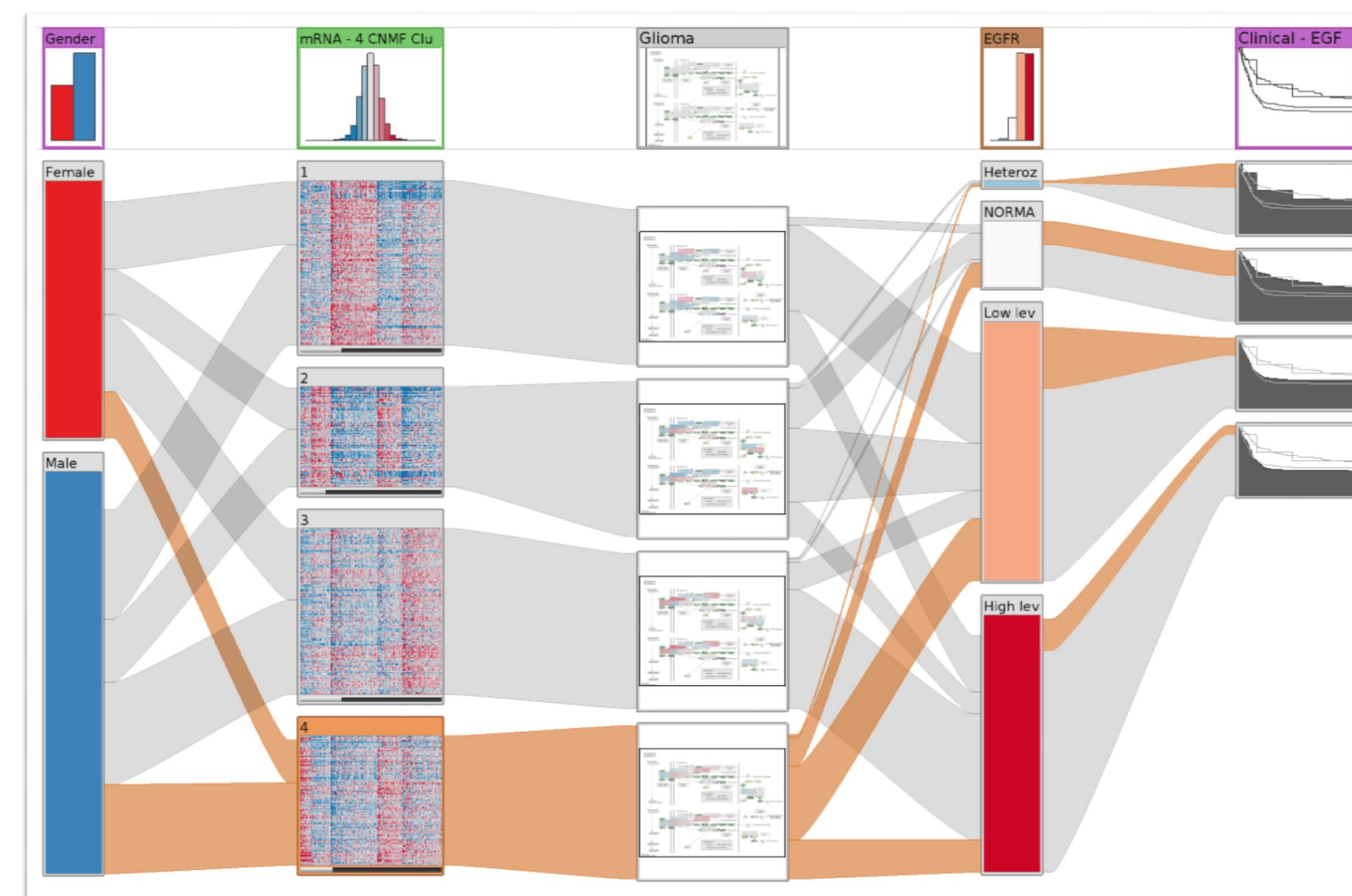
Harvard University, Cambridge, USA



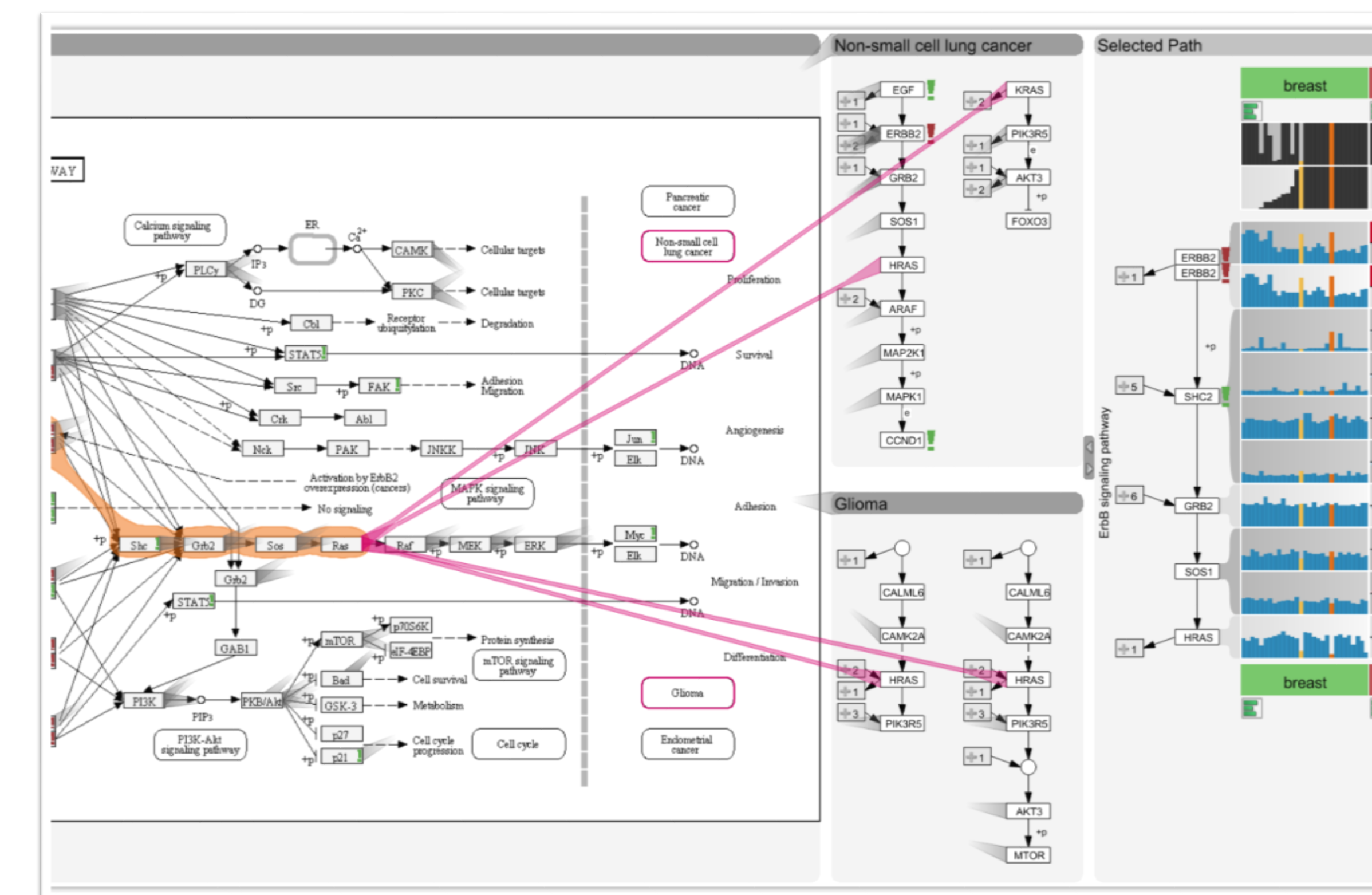
Content



Set Visualization



Heterogeneous Datasets
Cancer Subtype Analysis



Multivariate Networks
Pathway Analysis

Subset Visualization

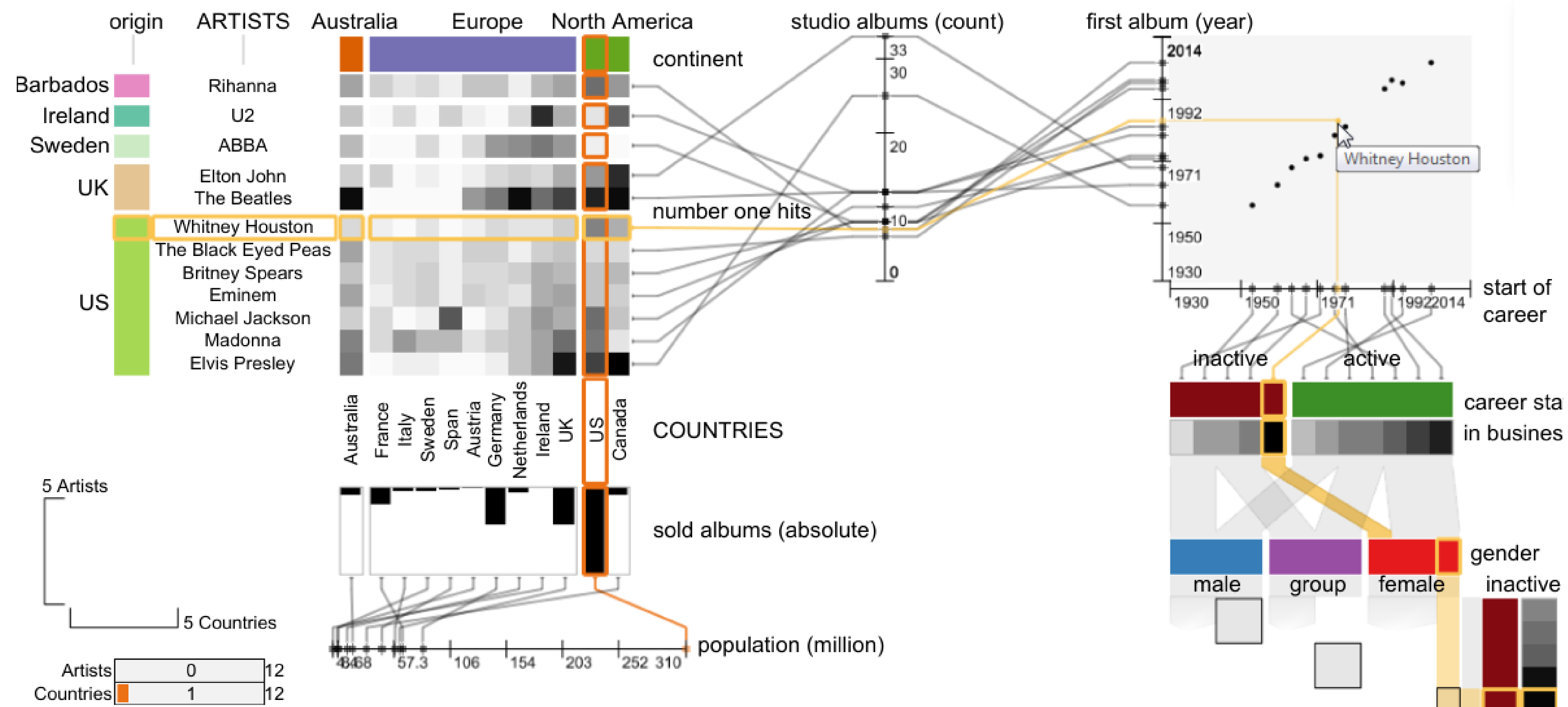
**Domino: Extracting, Comparing,
and Manipulating Subsets
across Multiple Tabular Datasets**

<http://domino.caleydo.org>

**Samuel Gratzl, Nils Gehlenborg,
Alexander Lex, Hanspeter Pfister
and Marc Streit**

IEEE TVCG (InfoVis '14)

HONORABLE MENTION AWARD



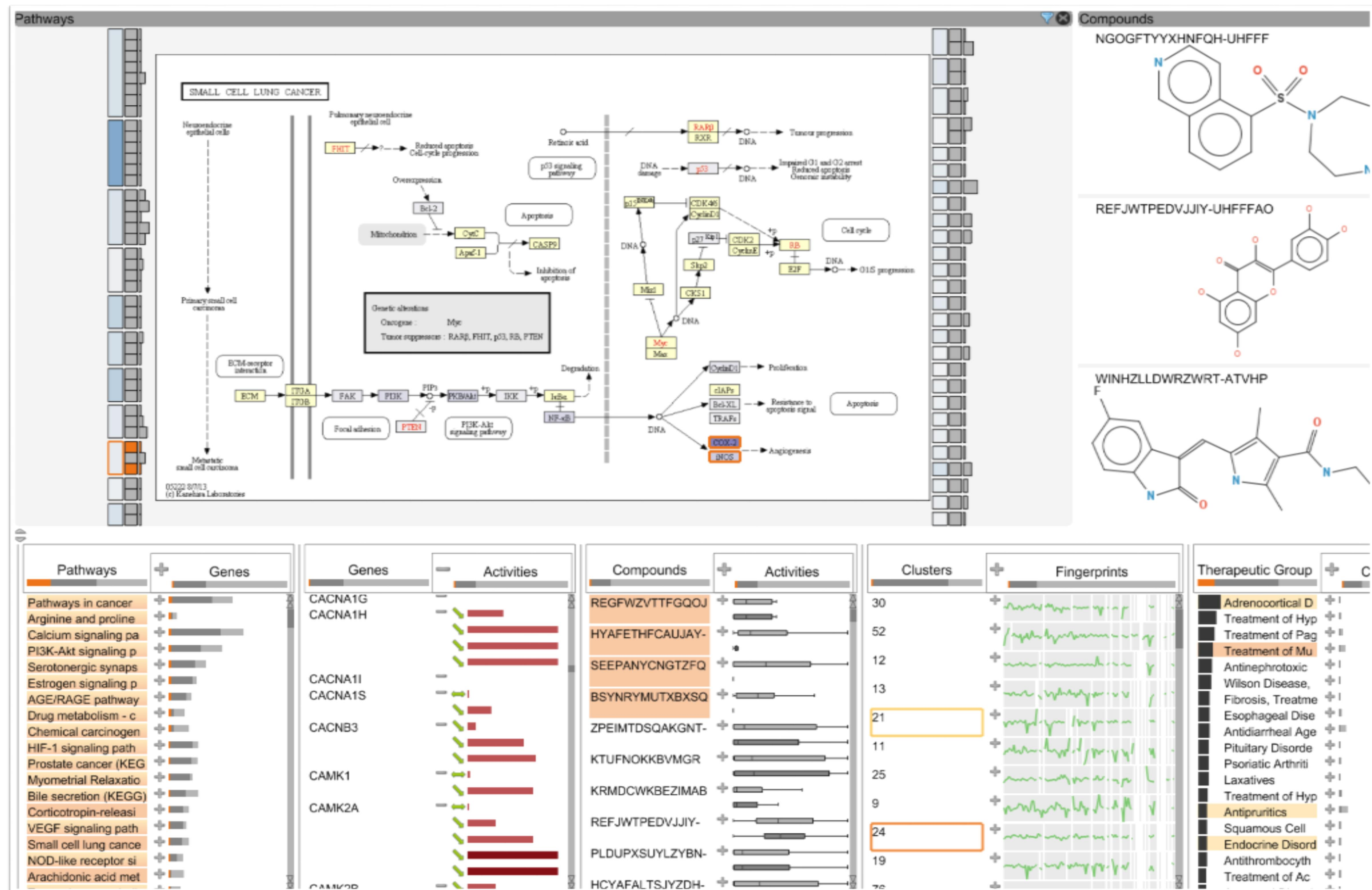
Multi-Relational Data Drug Discovery

ConTour: Data-Driven Exploration of Multi-Relational Datasets for Drug Discovery

<http://contour.caleydo.org>

Christian Partl, Alexander Lex,
Marc Streit, Hendrik Strobel,
Anne-Mai Wasserman,
H. Pfister, and Dieter Schmalstieg

IEEE TVCG (VAST '14)



Sequence Visualization

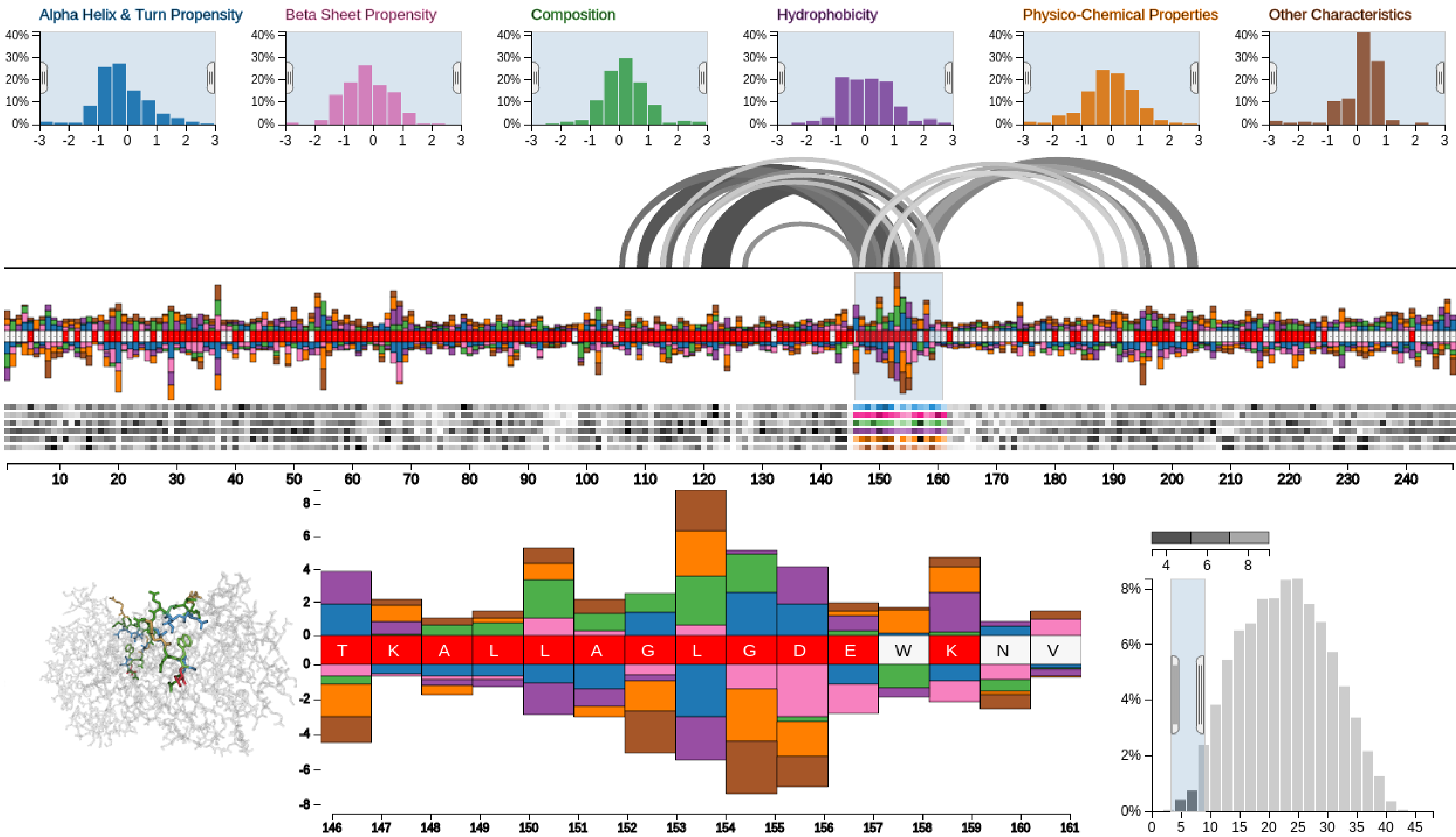
Protein Mutations

**Mu-8: Visualizing Differences
between Proteins and their Families**

<http://mu-8.com>

**Johnathan D Mercer, Balaji Pandian,
Alexander Lex*, Nicolas Bonneel,
and Hanspeter Pfister**

BMC Proceedings 2014



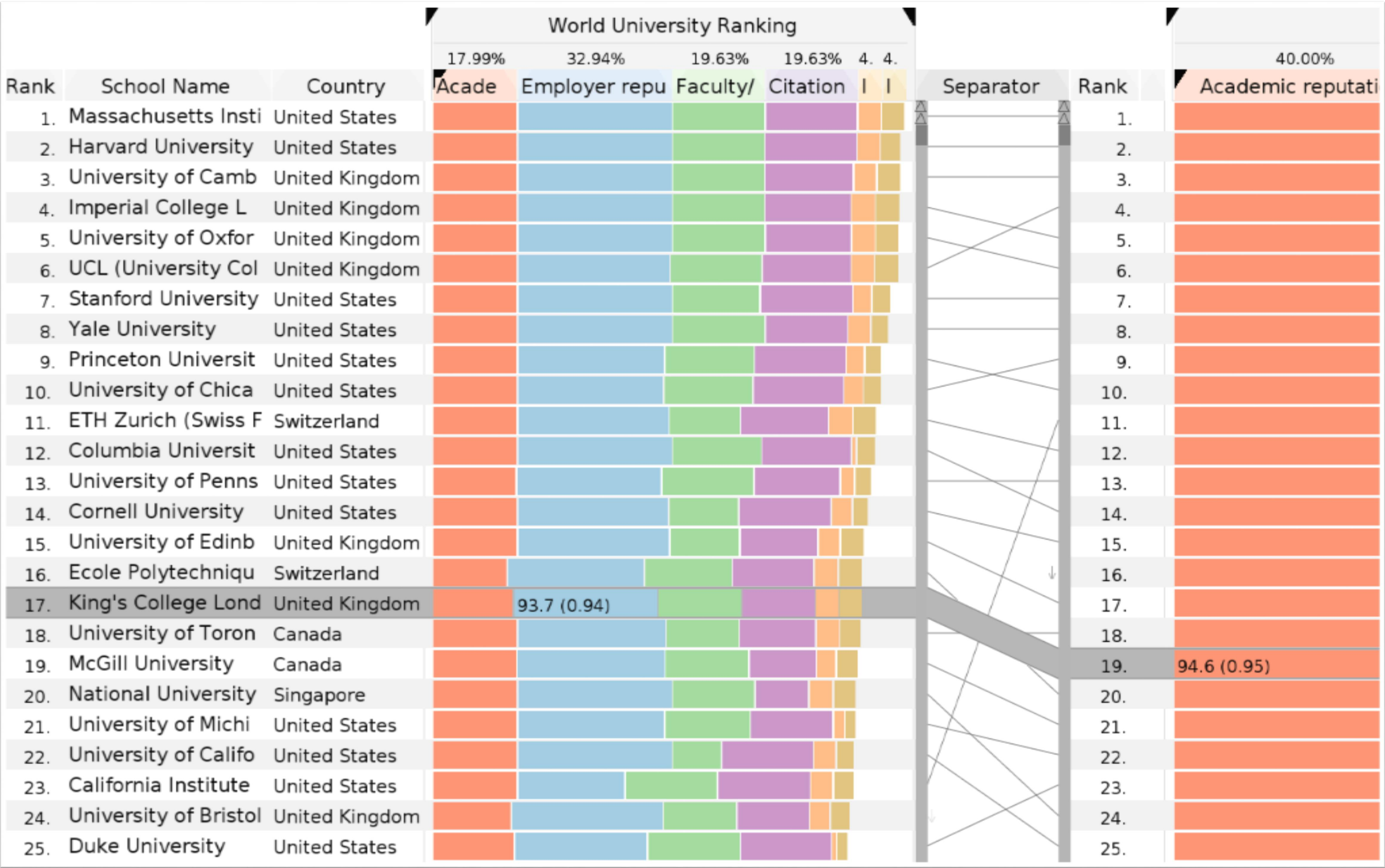
Multi-Attribute Rankings

LineUp: Visual Analysis of Multi-Attribute Rankings

<http://lineup.caleydo.org>

Samuel Gratzl, Alexander Lex,
Nils Gehlenborg, Hanspeter Pfister,
and Marc Streit

IEEE TVCG (InfoVis '13)
BEST PAPER AWARD

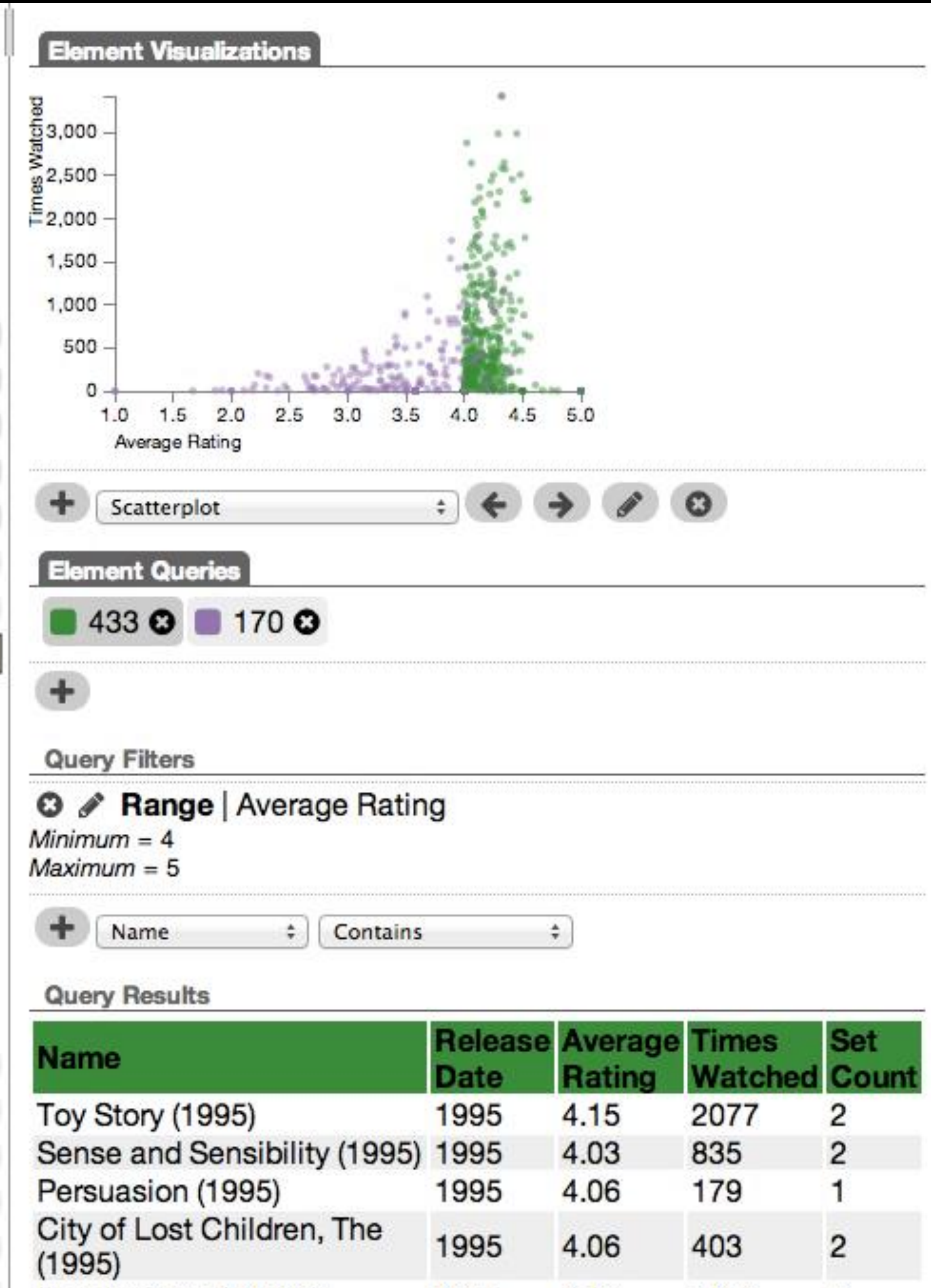
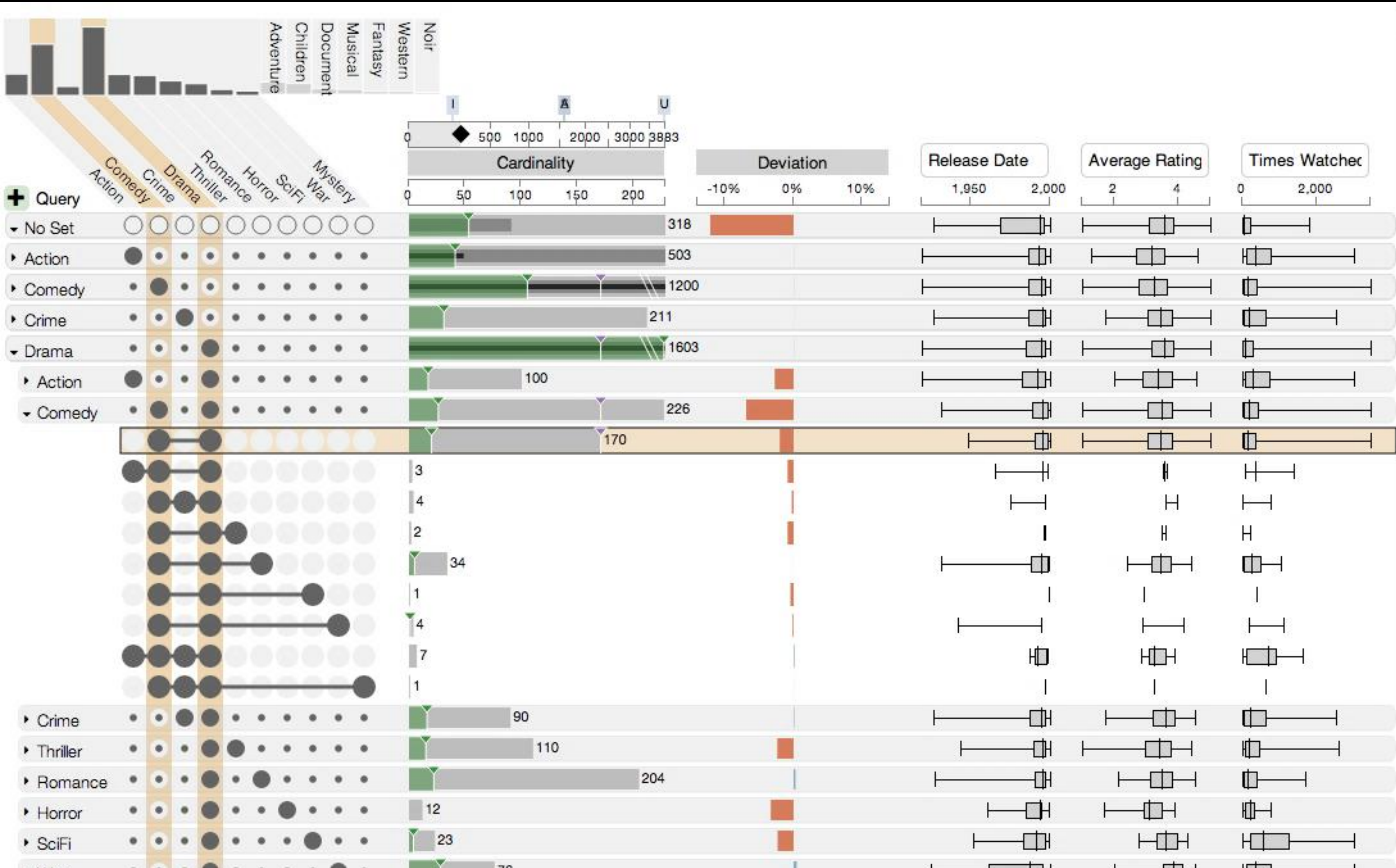


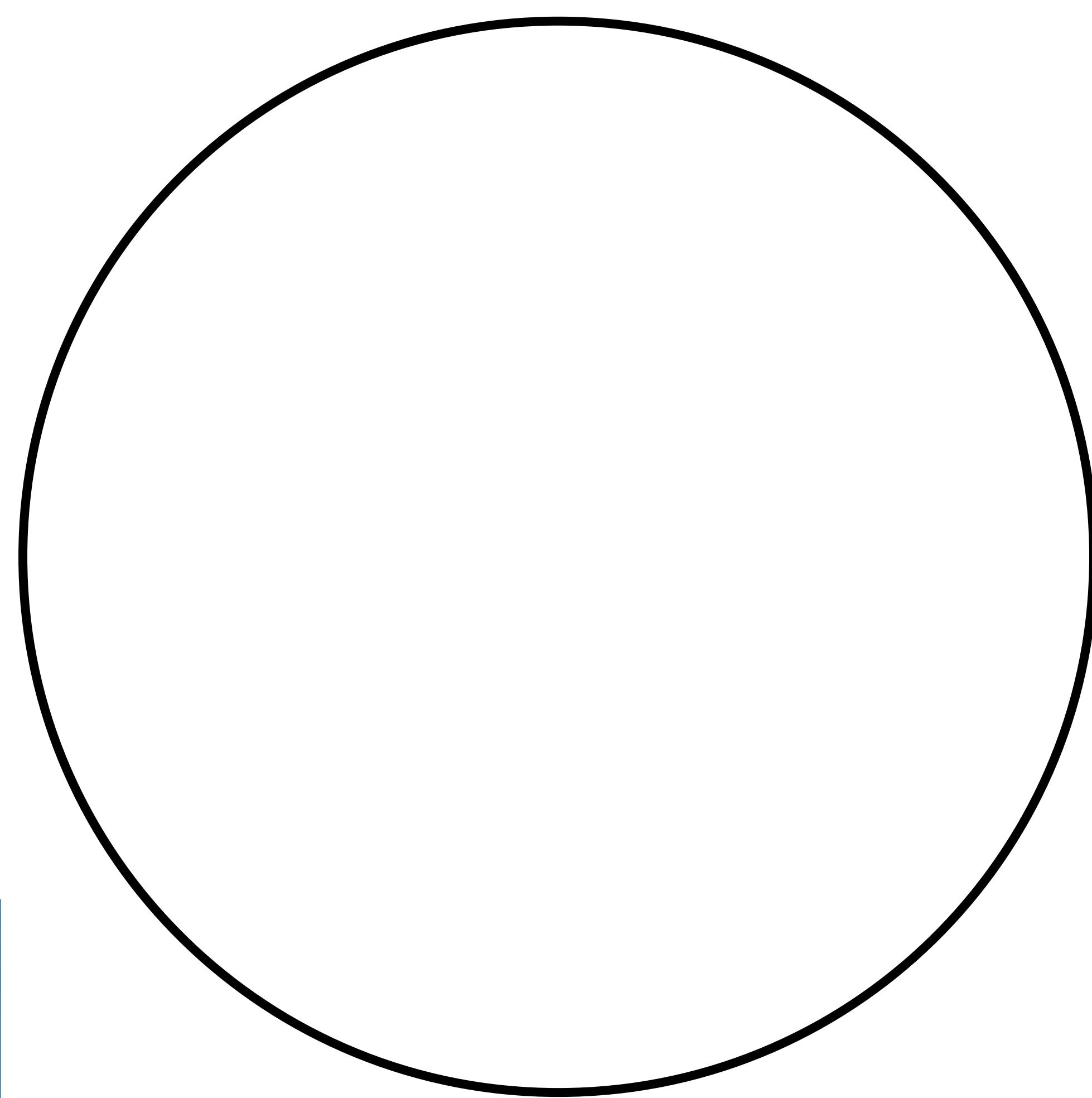
Alexander Lex, Nils Gehlenborg, Hendrik Strobel,
Romain Vuillemot, and Hanspeter Pfister

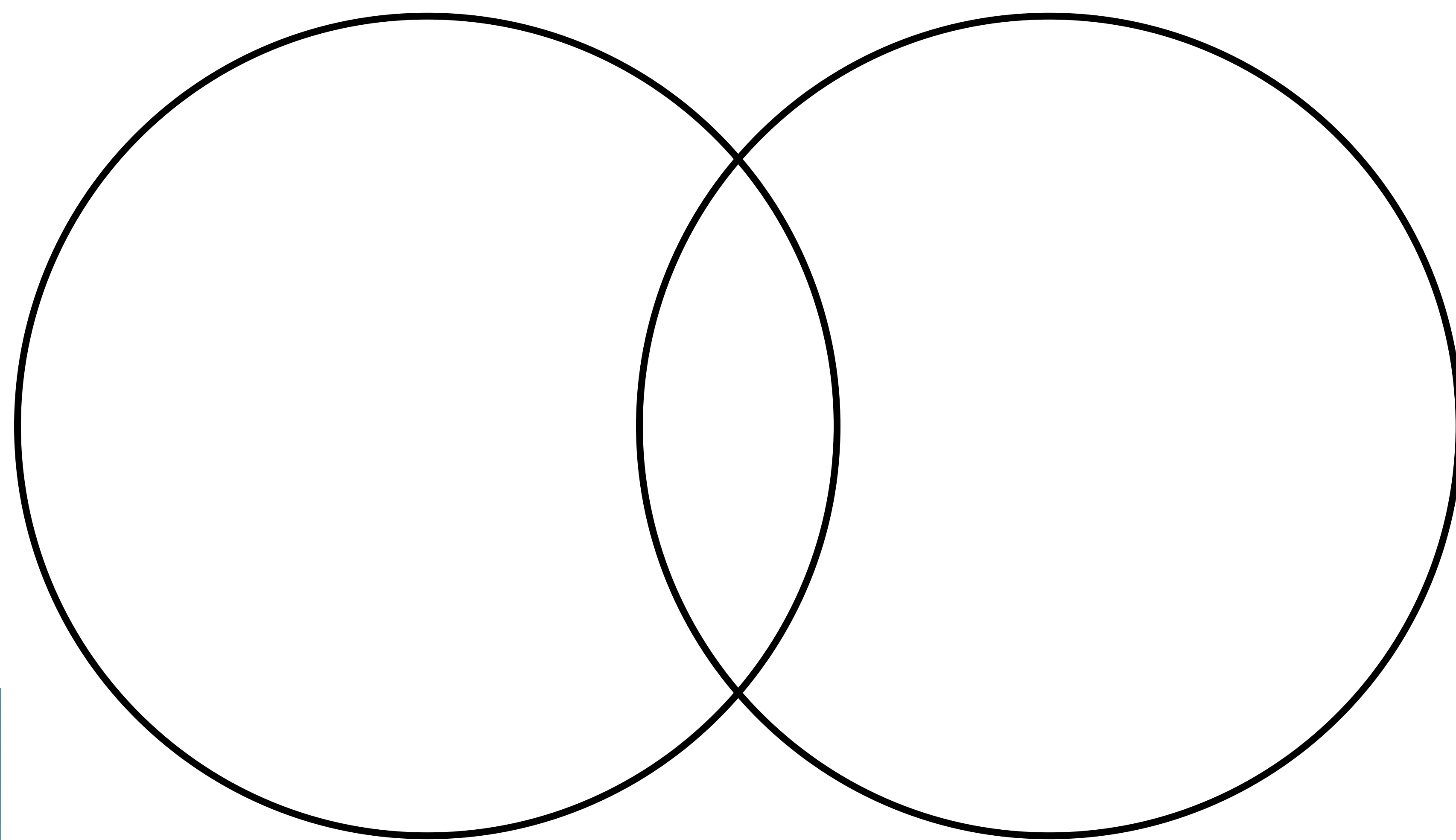
UpSet

InfoVis'14

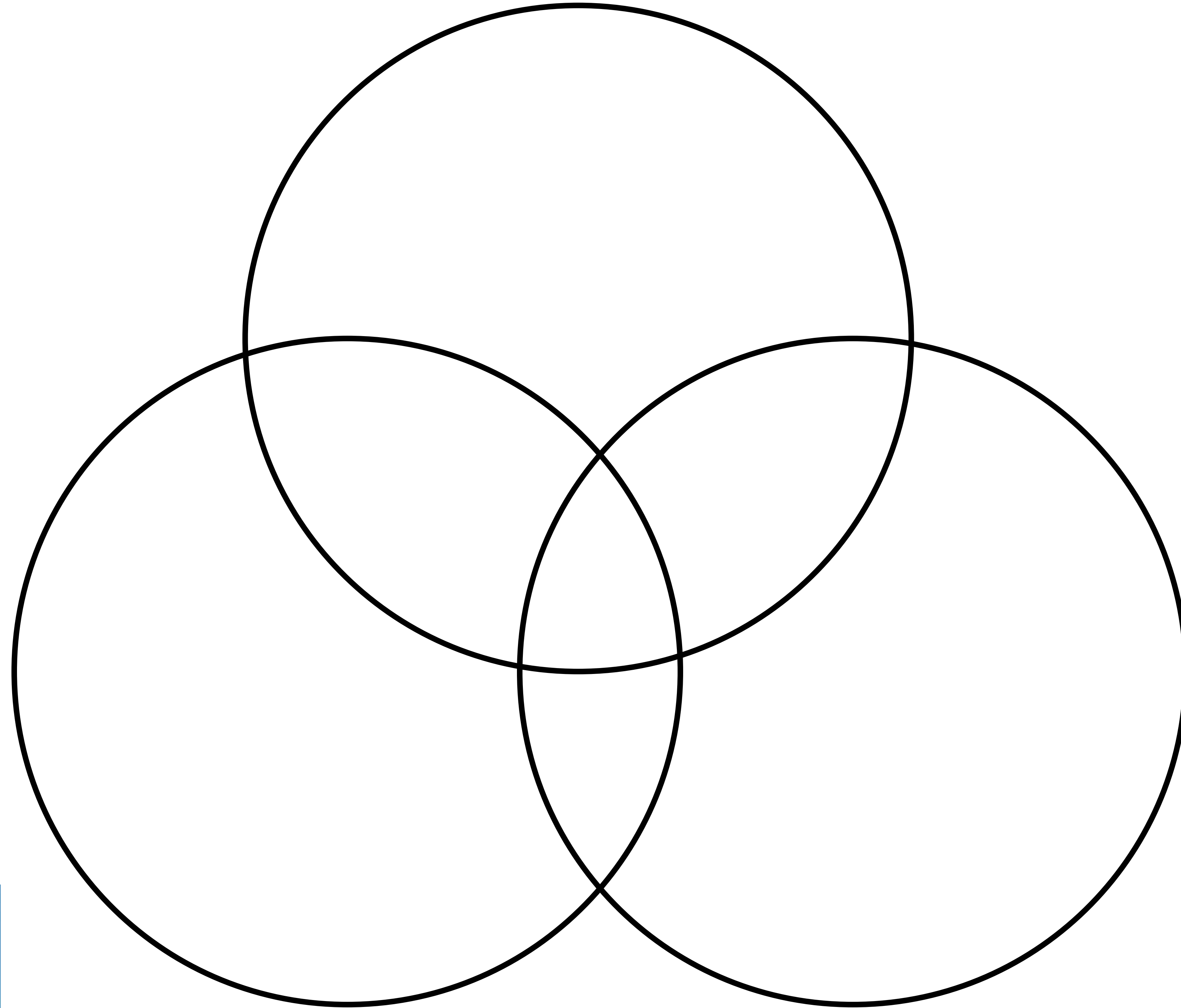
Visualizing Intersecting Sets





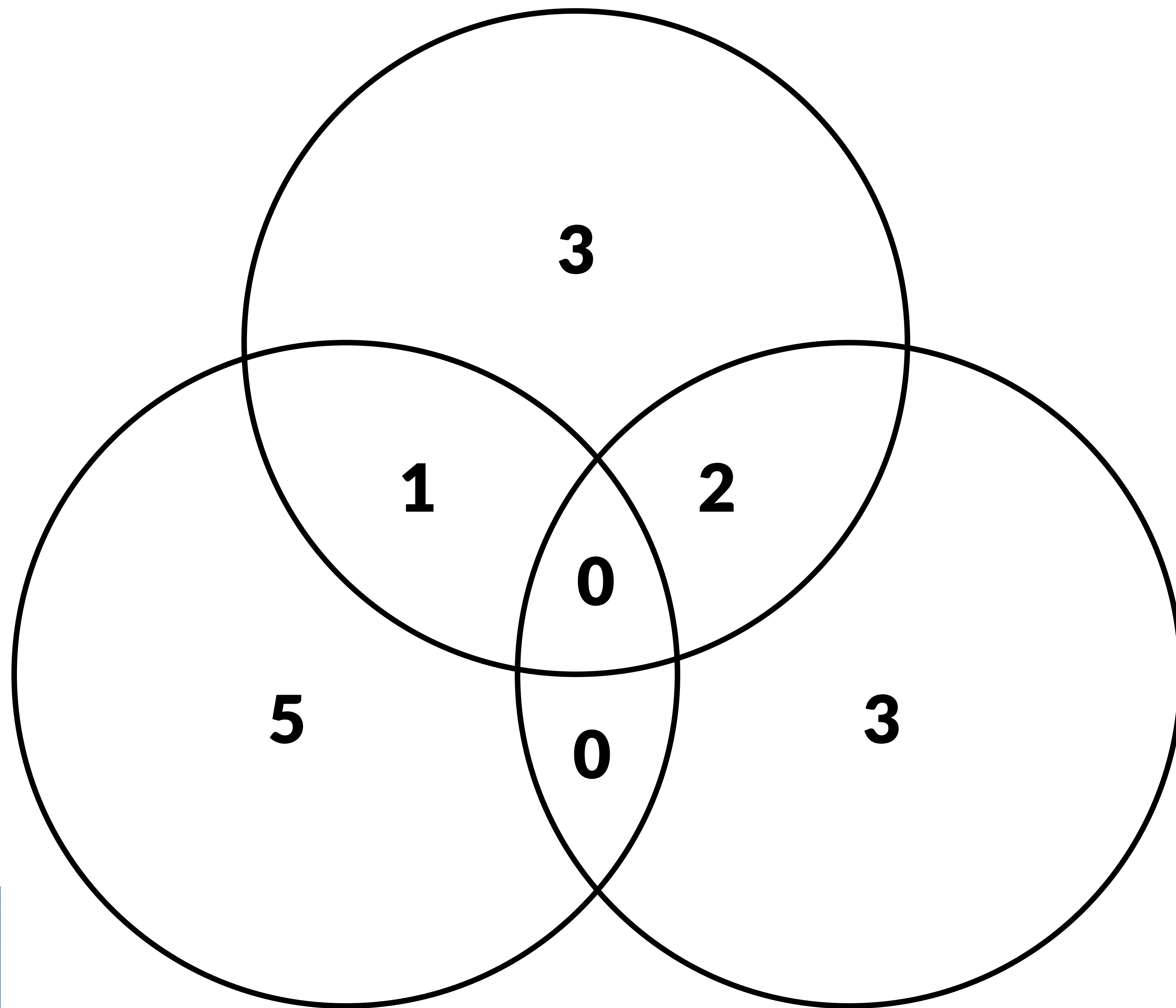


Evil





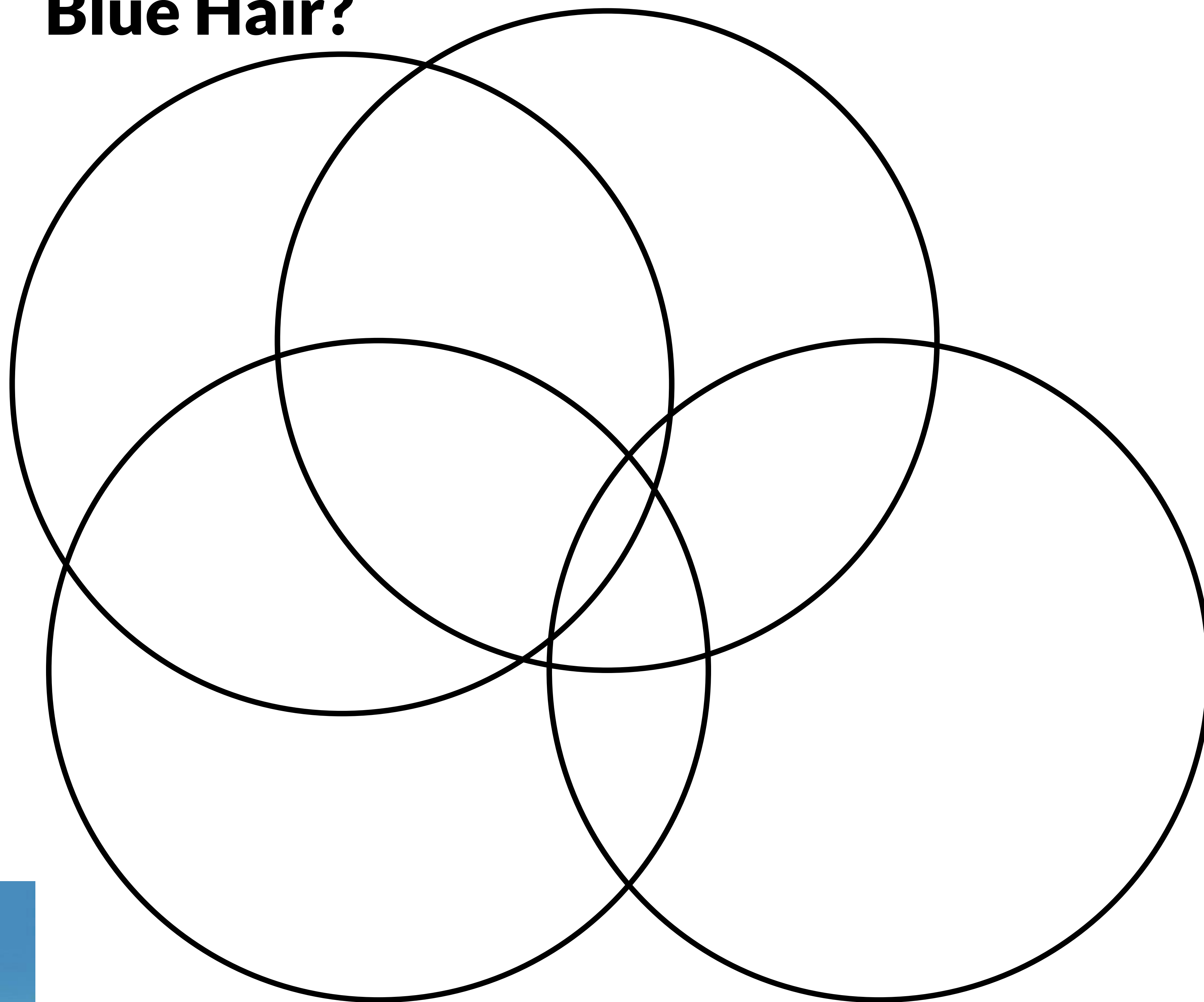
Evil





Blue Hair?

Evil

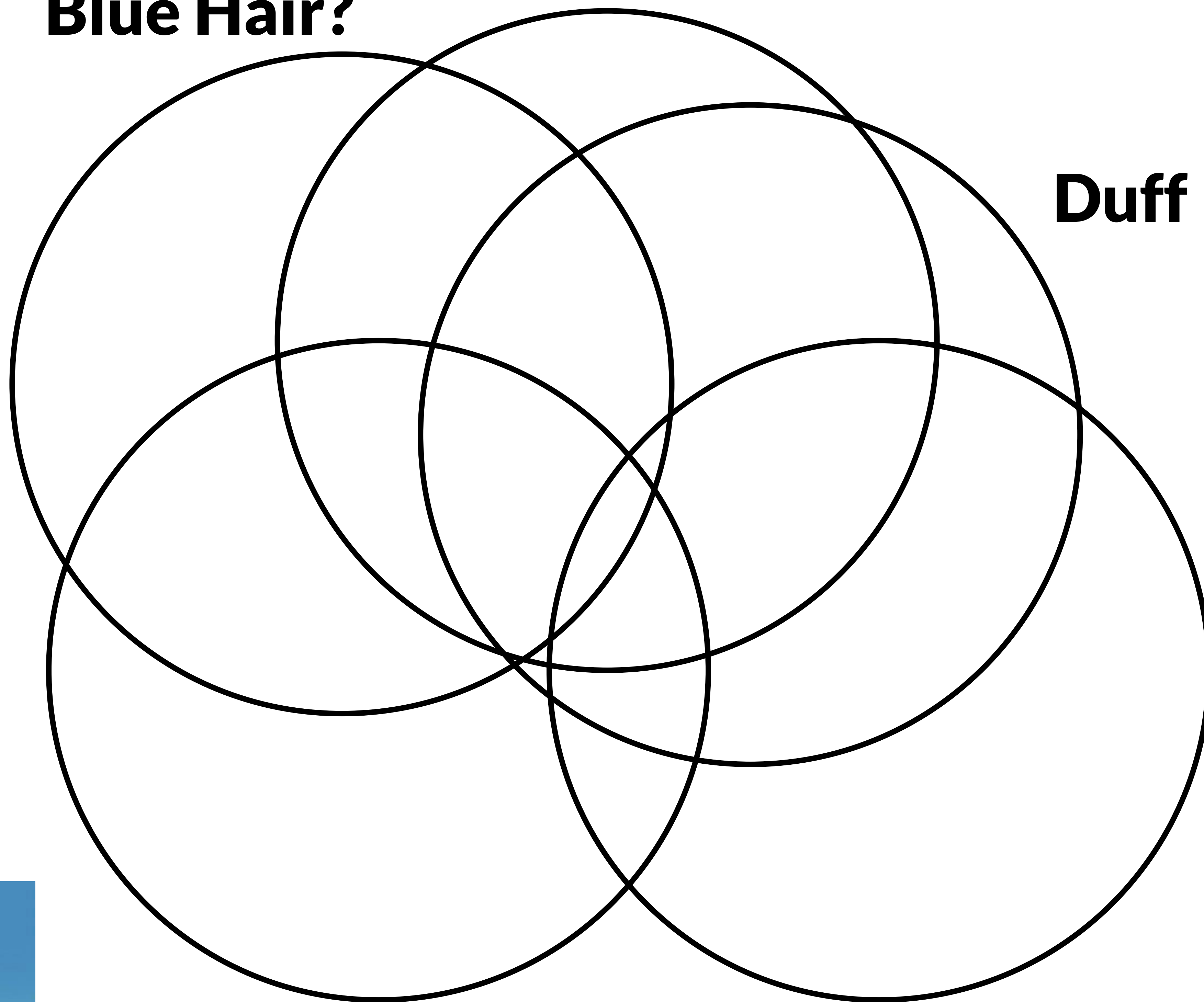


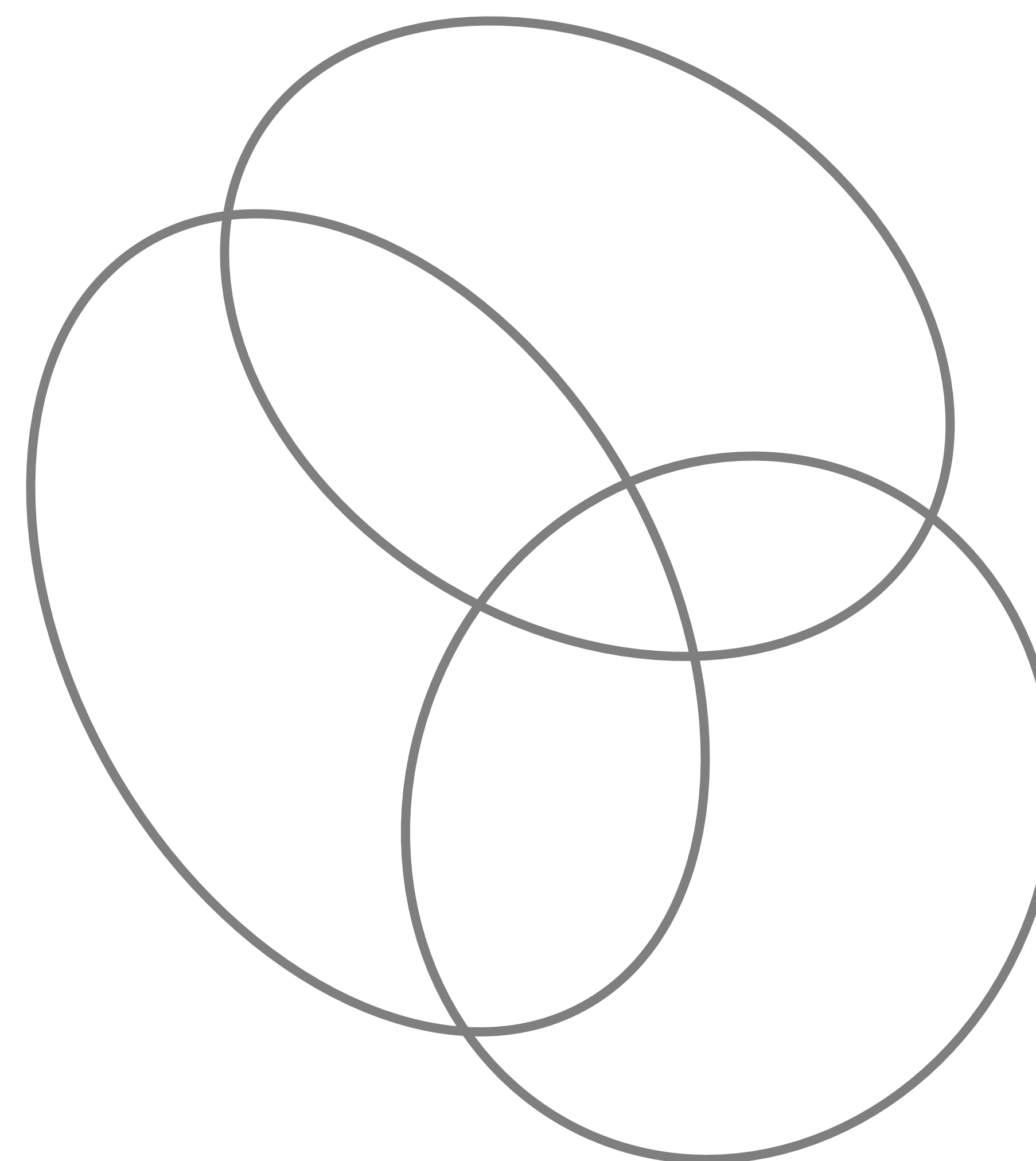


Blue Hair?

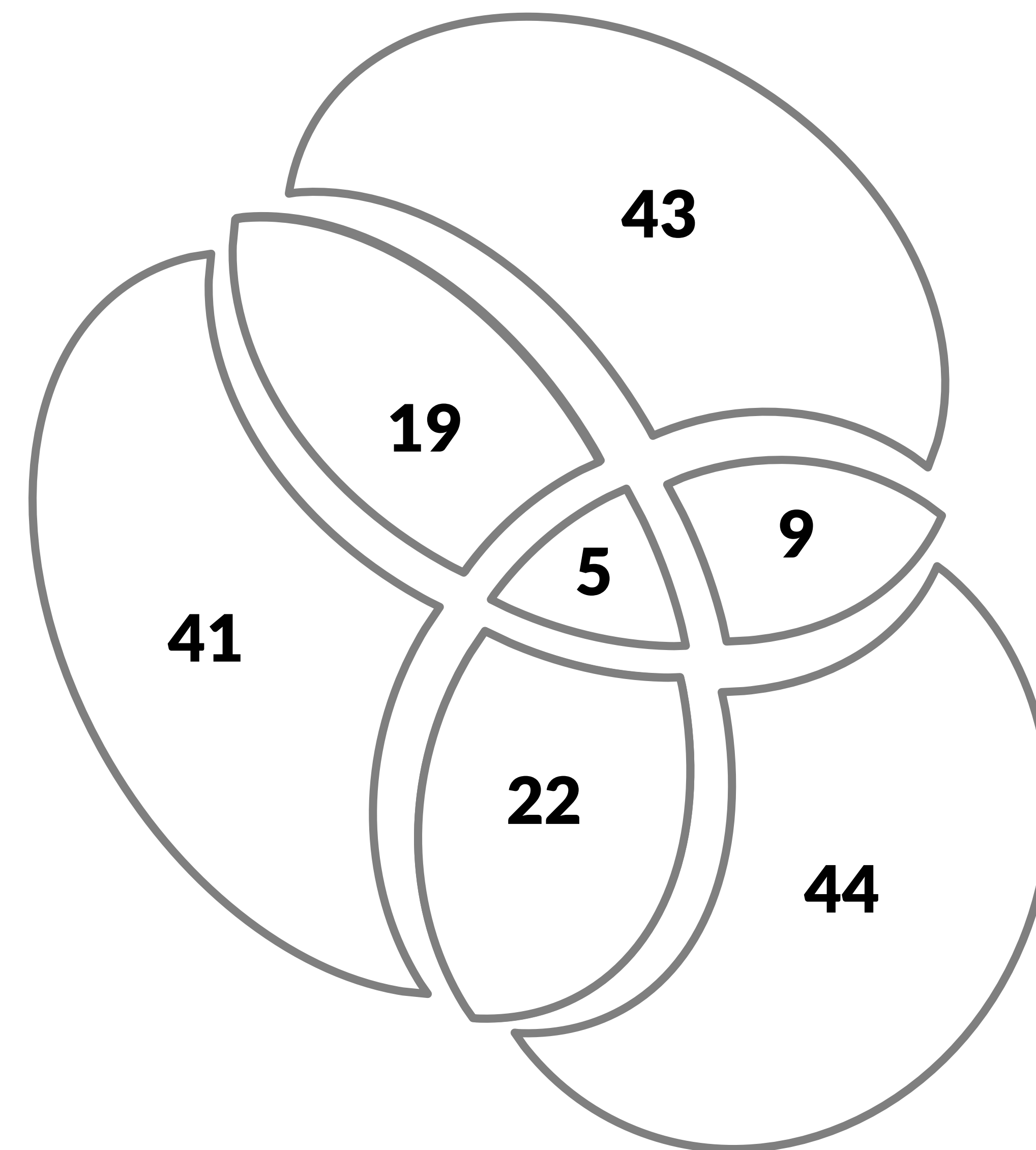
Evil

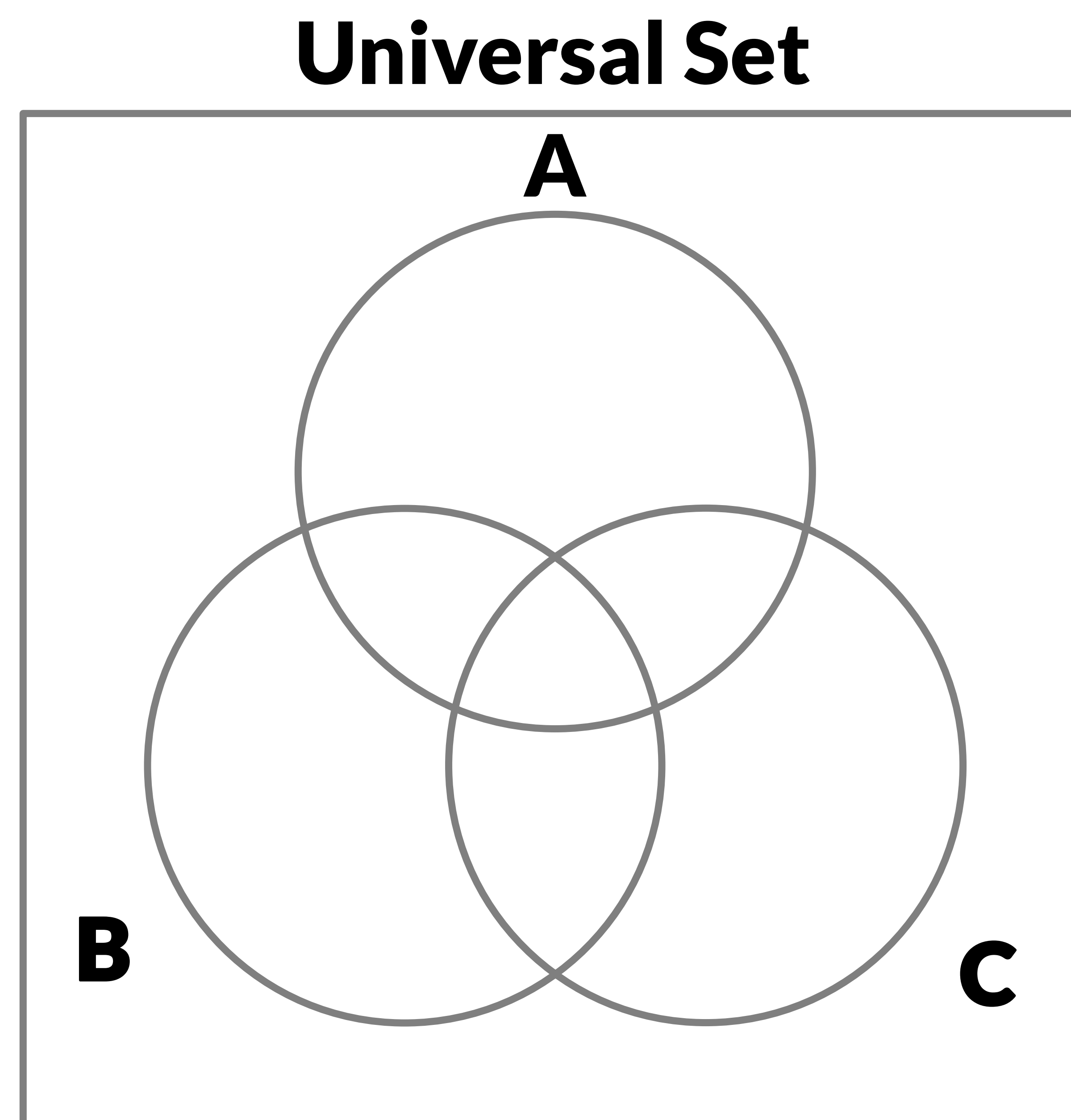
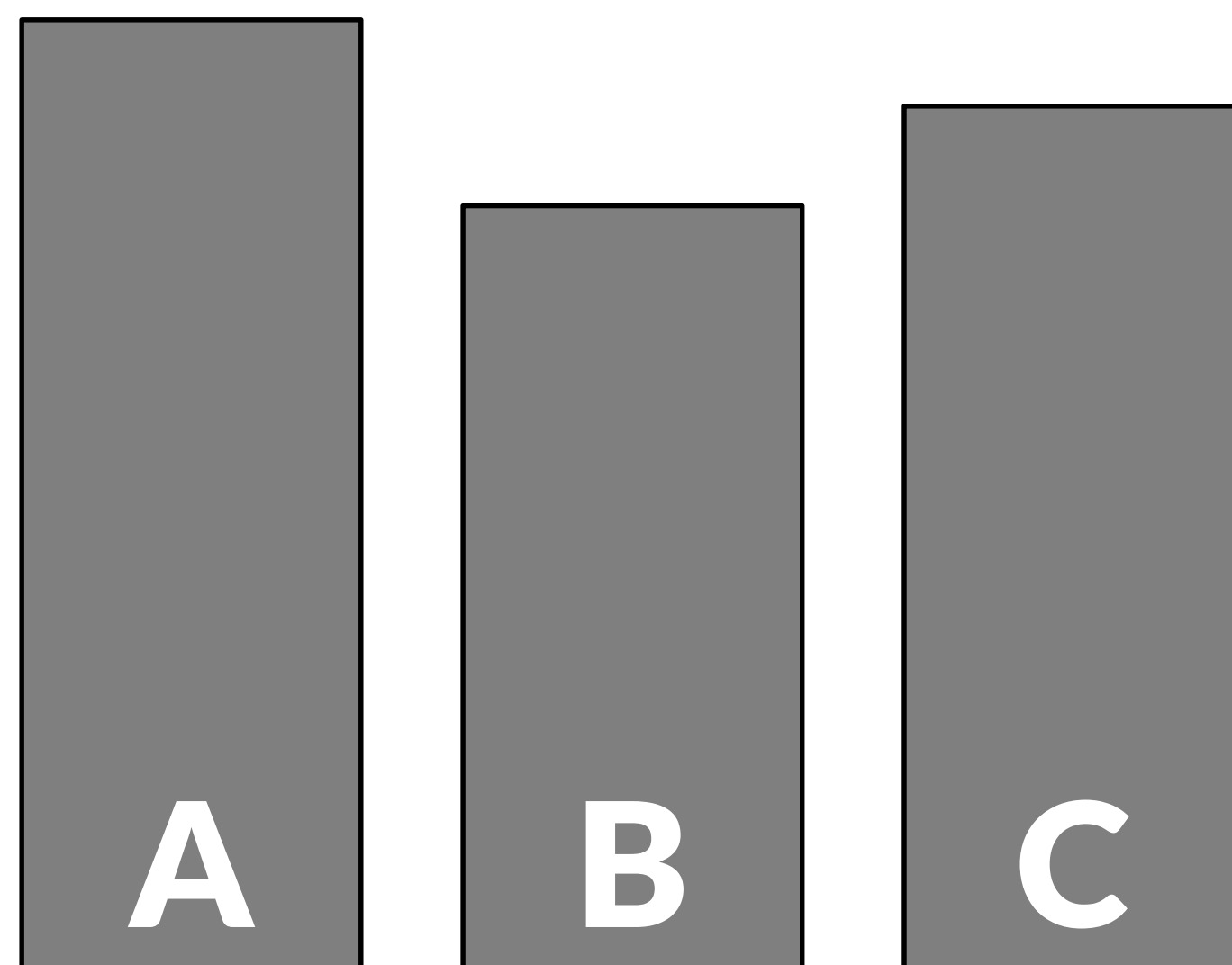
Duff Beer?

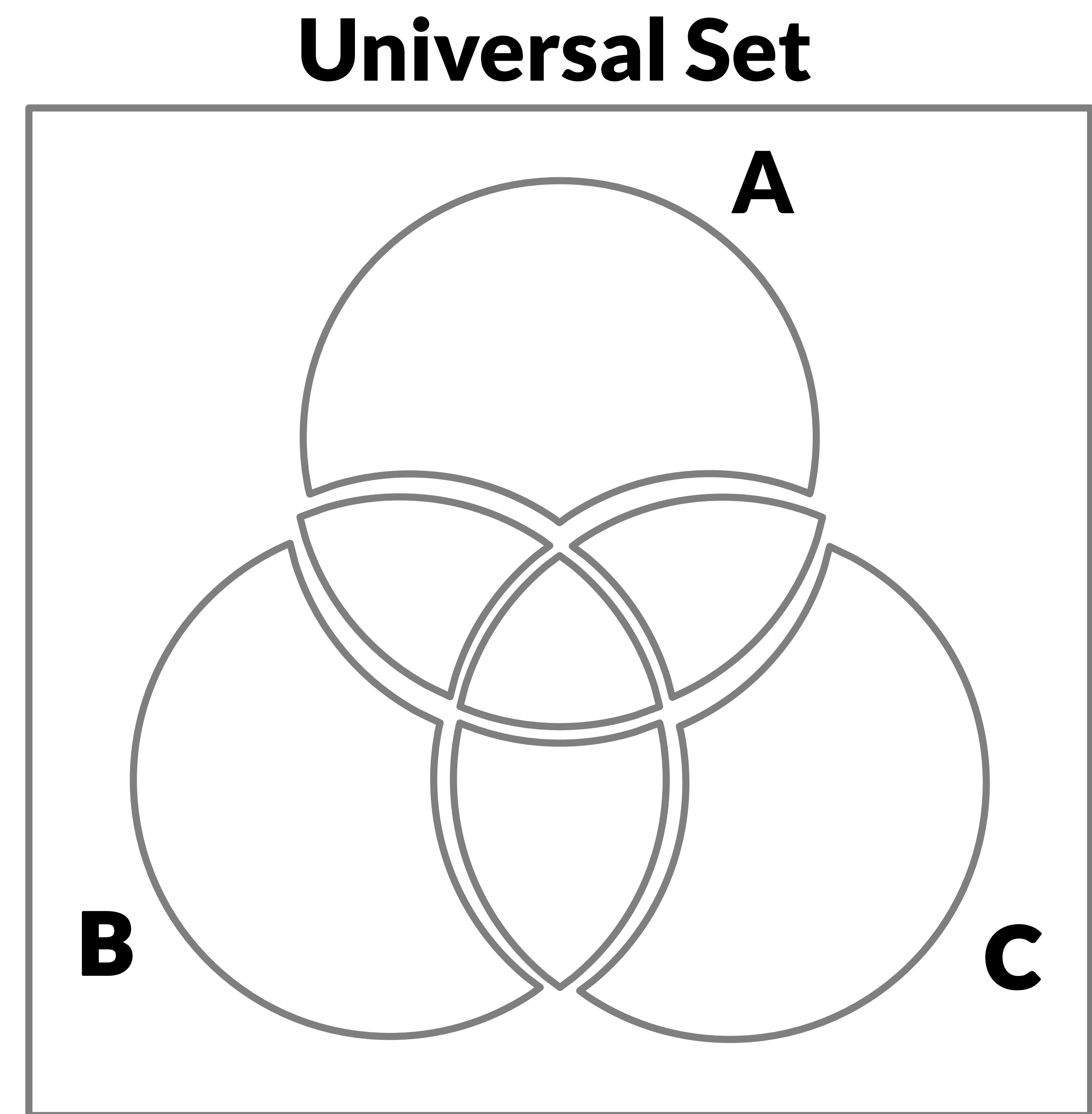
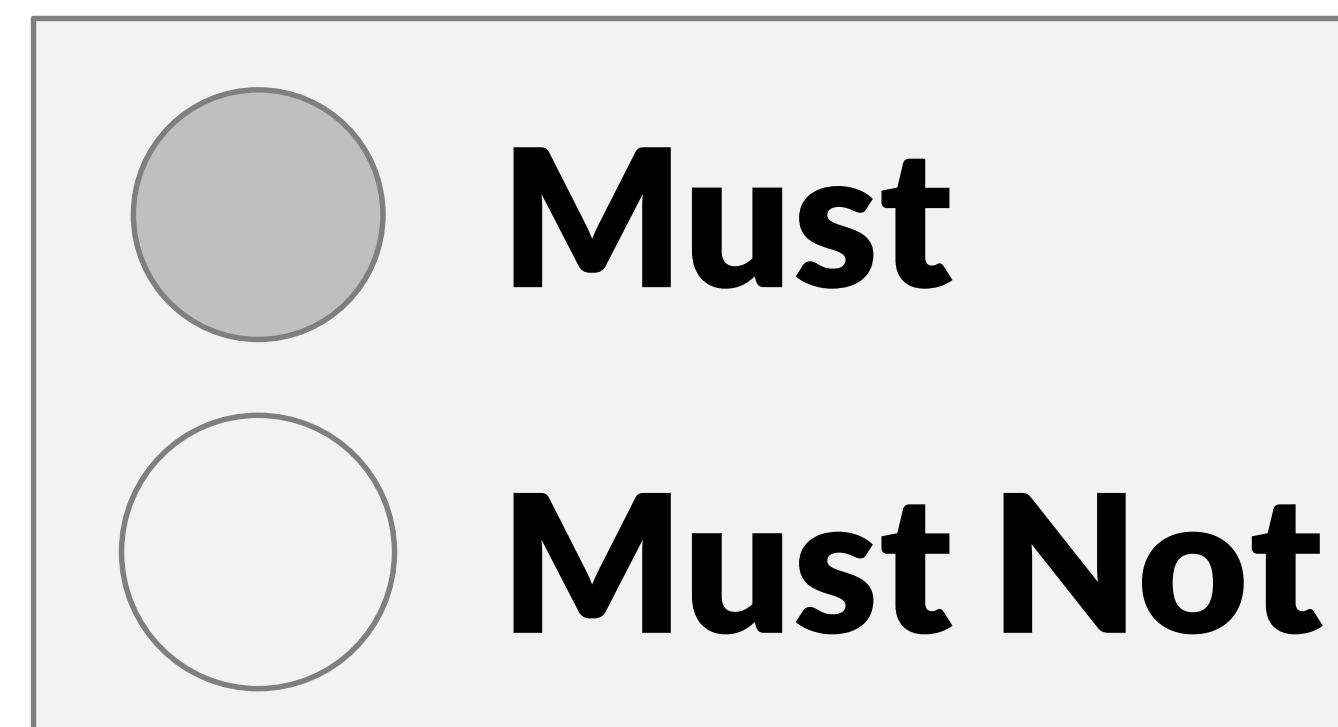
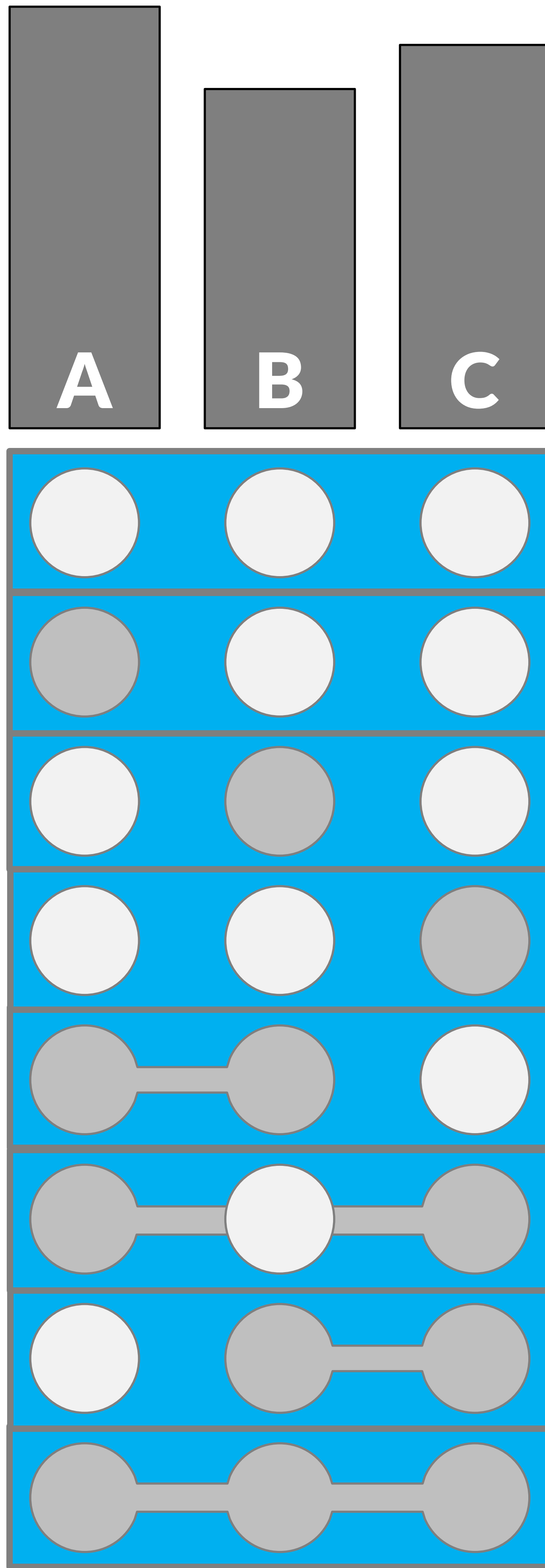


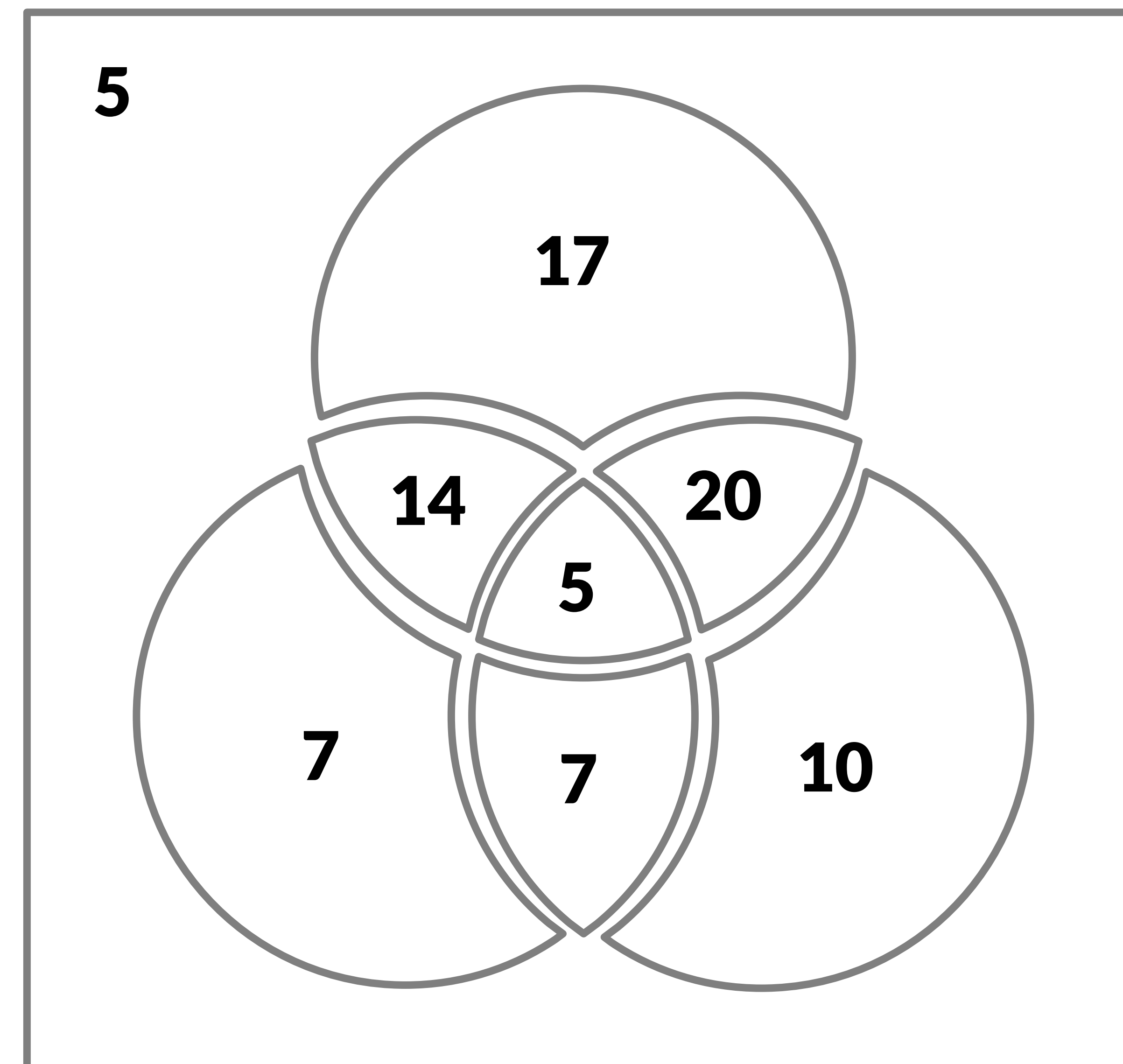
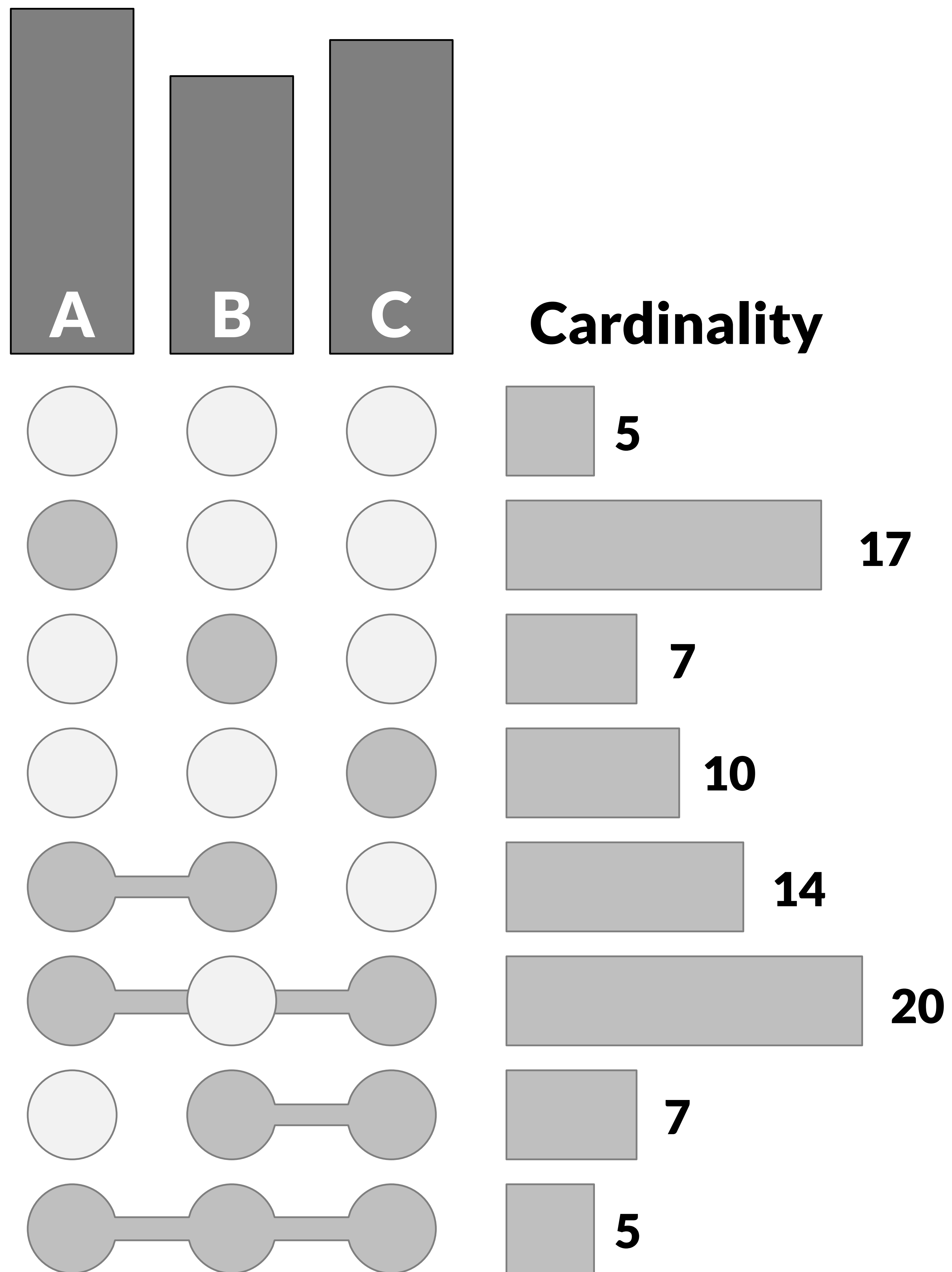


22 ? 19
^v • ^v



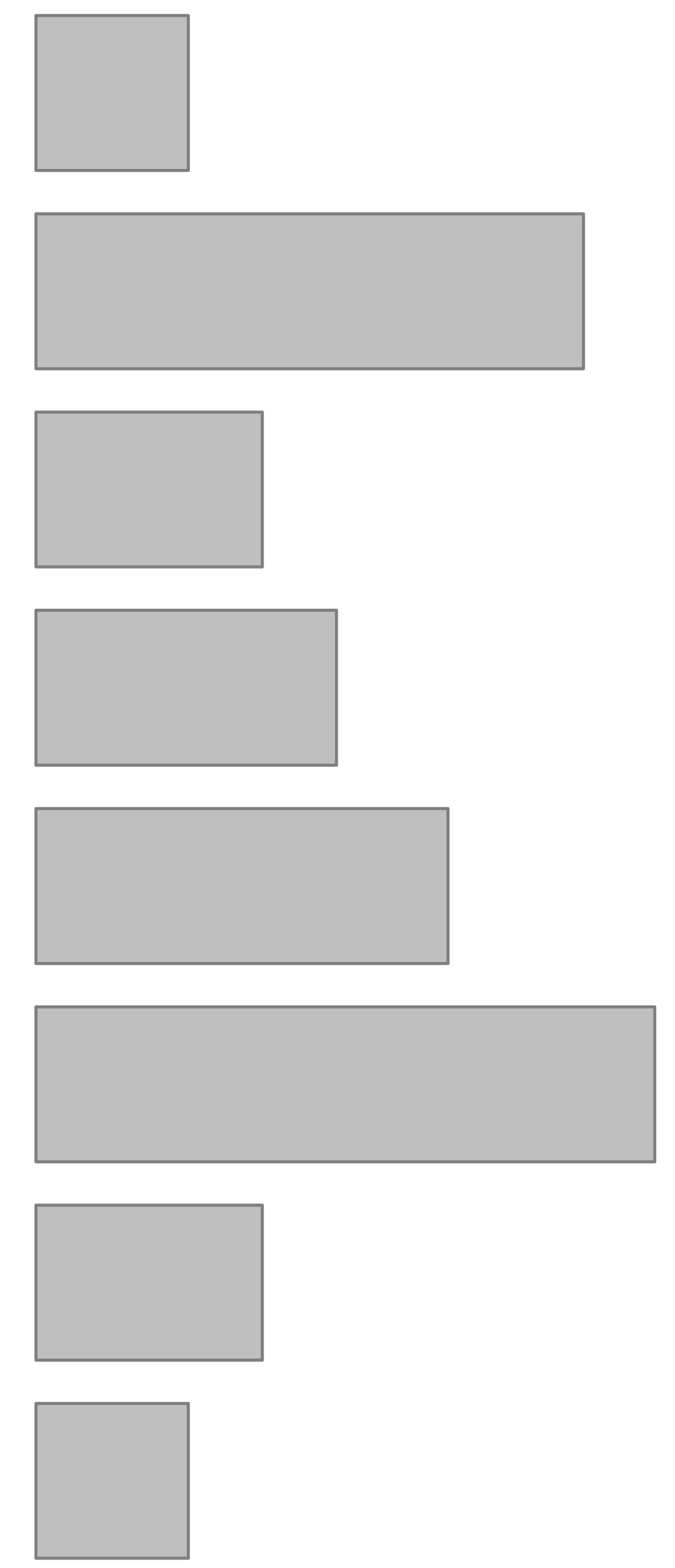
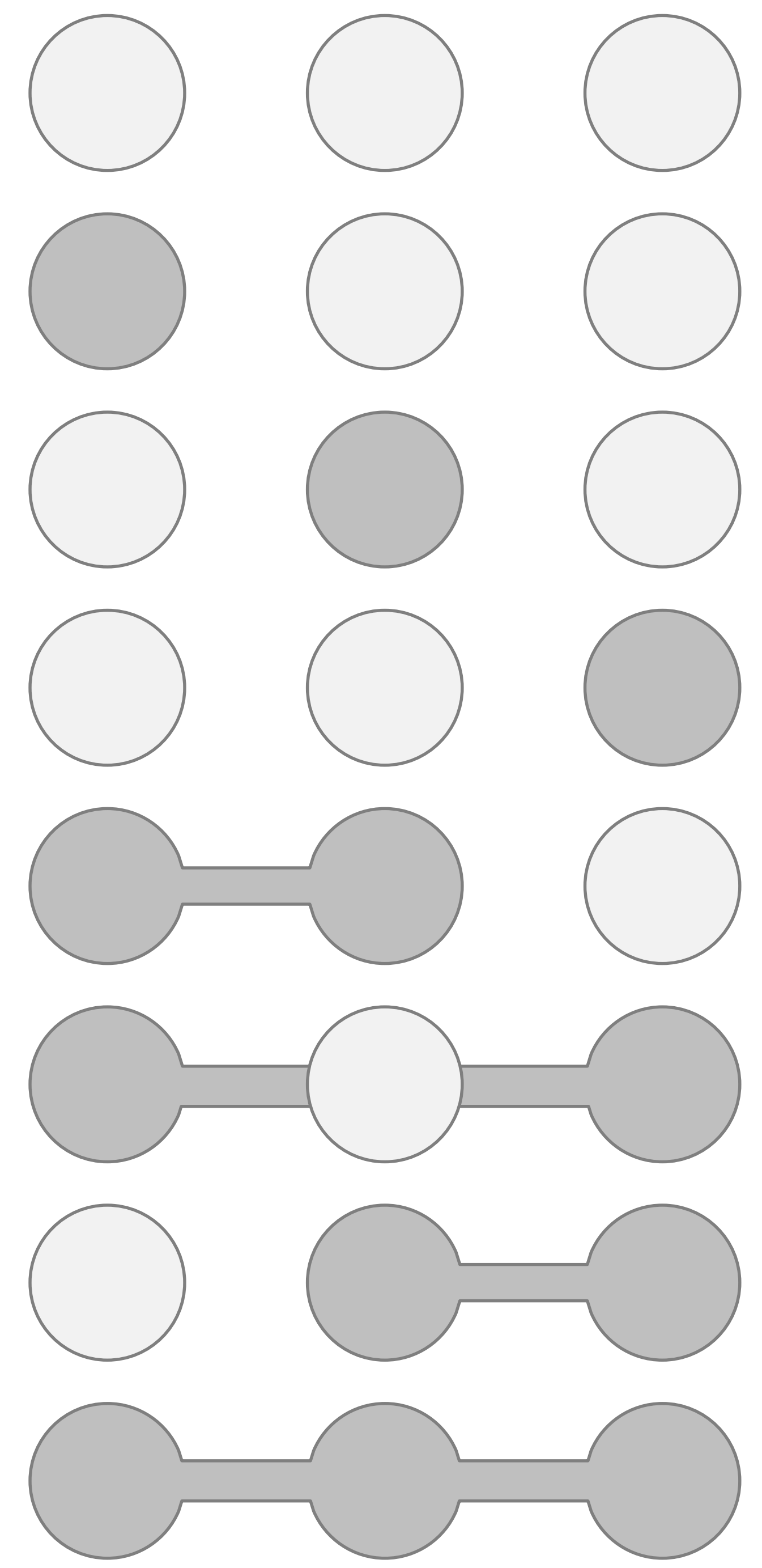
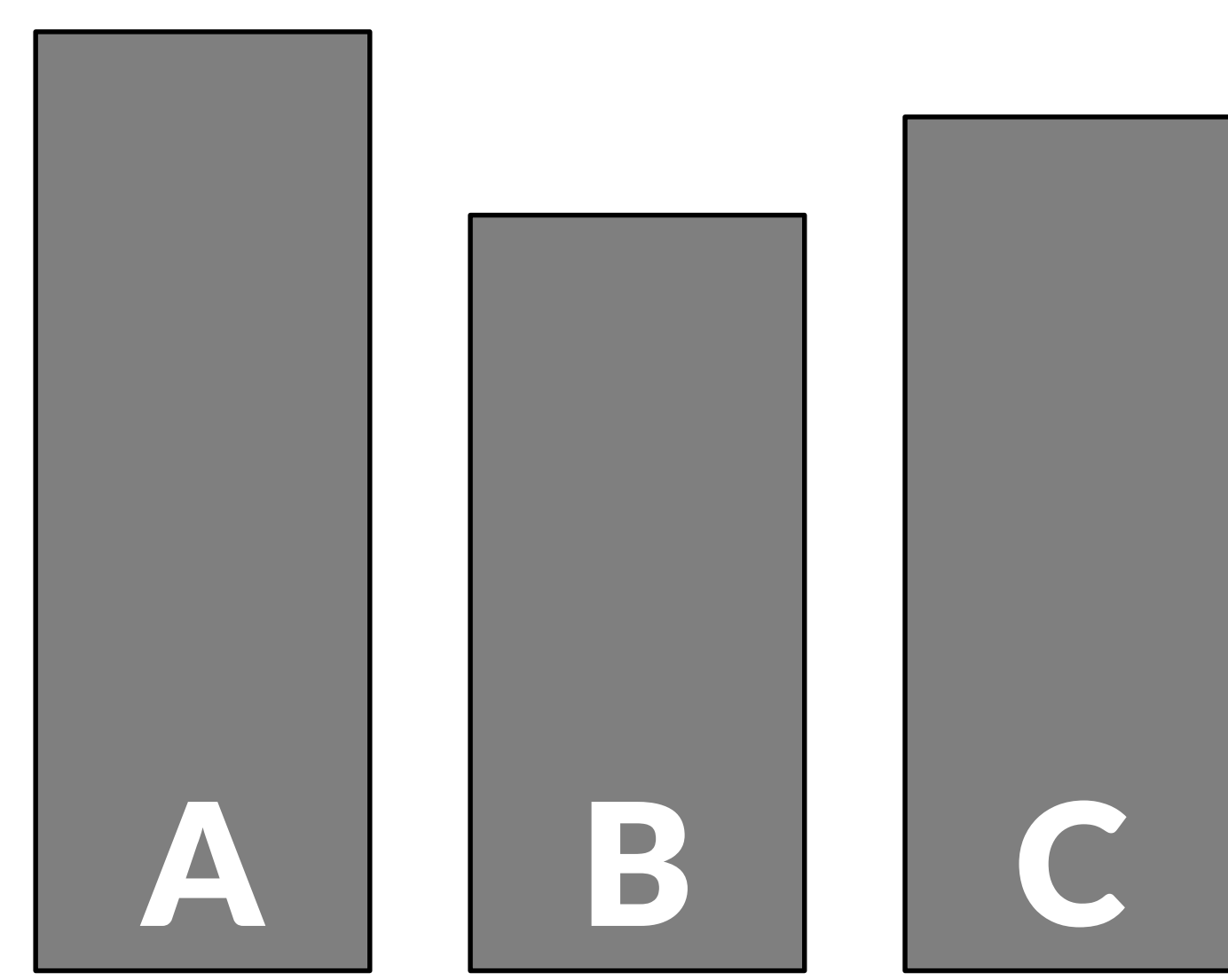




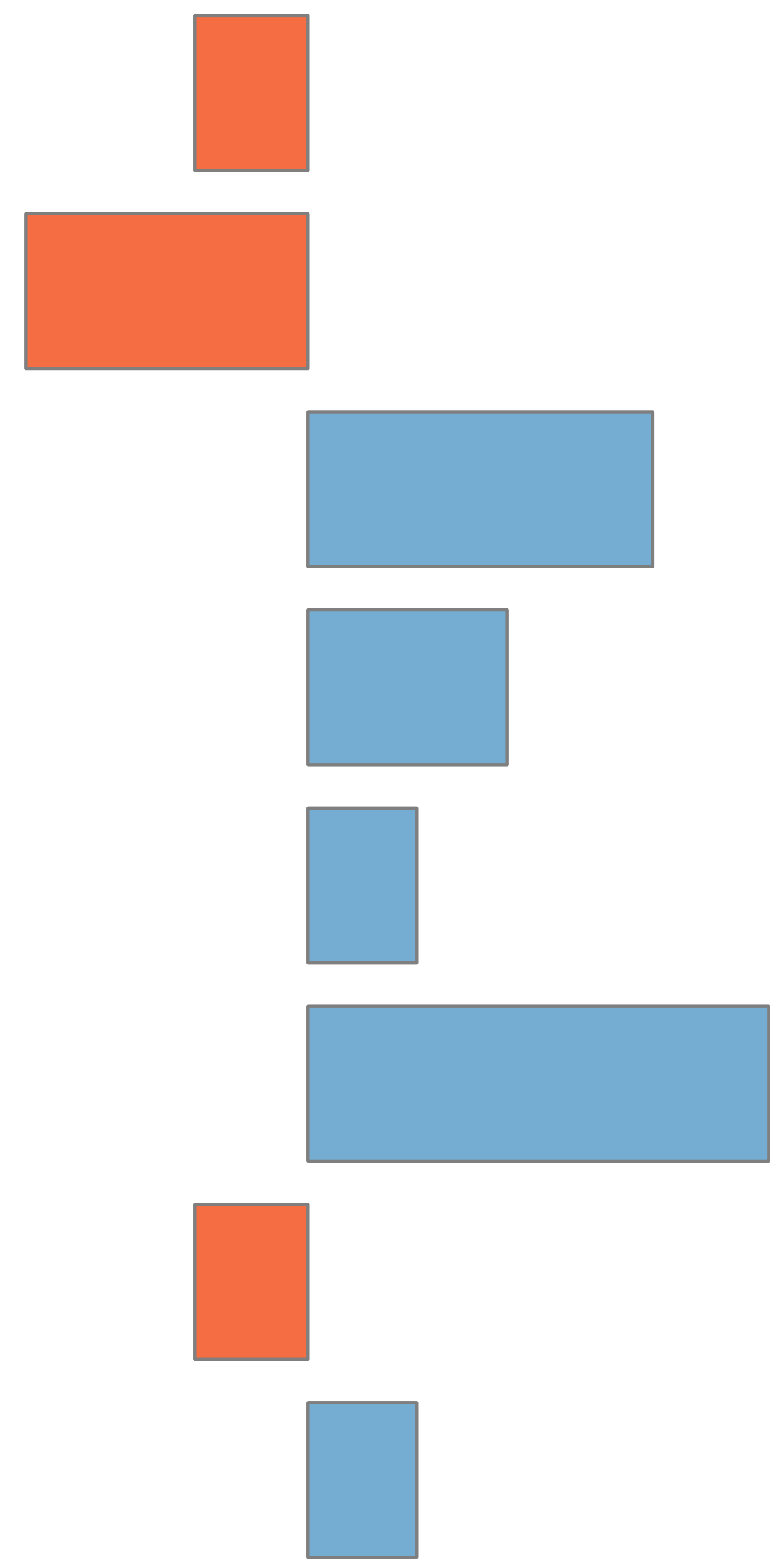


Plotting Attributes

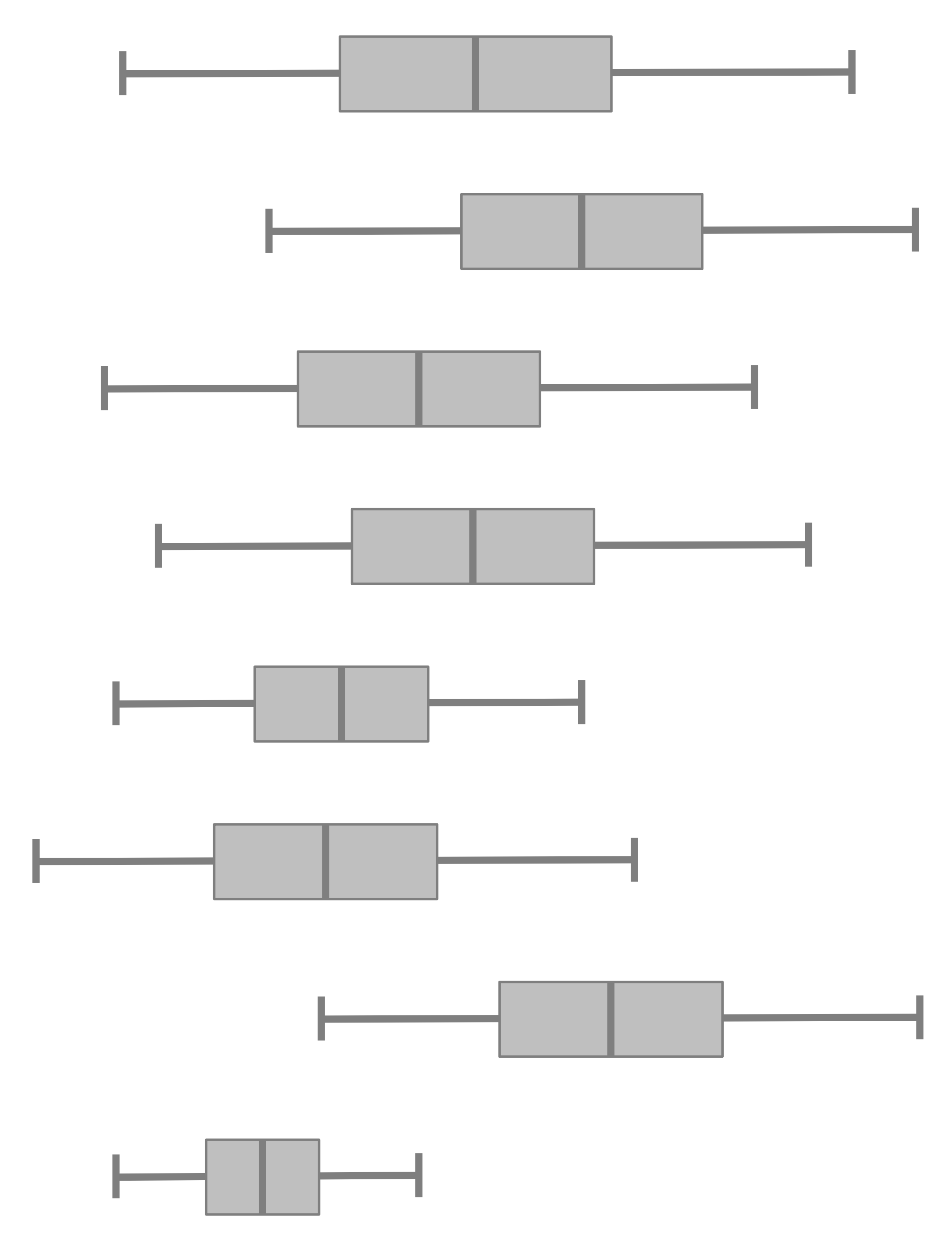
How surprising is the size of an intersection?
What's the distribution of an attribute in an intersection?

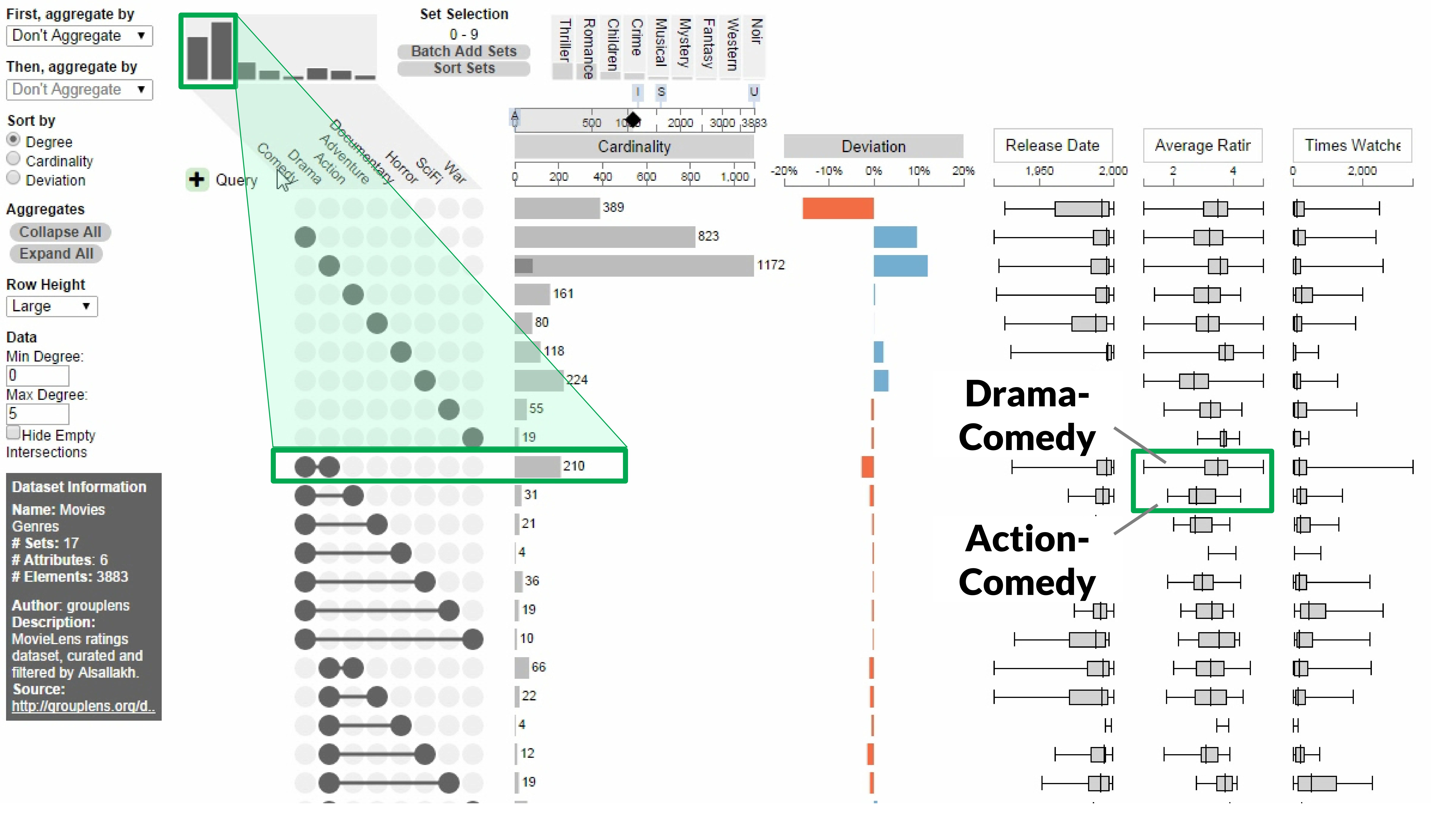


Deviation

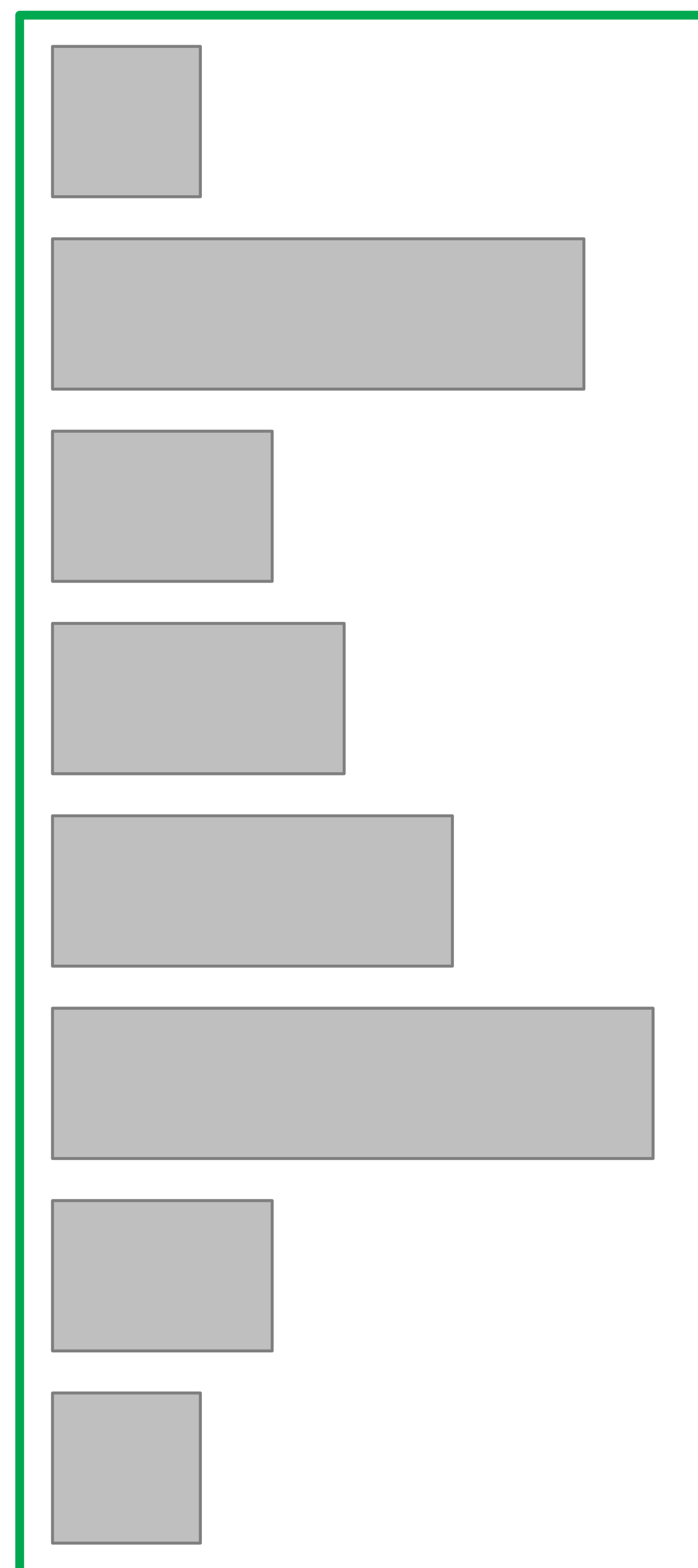
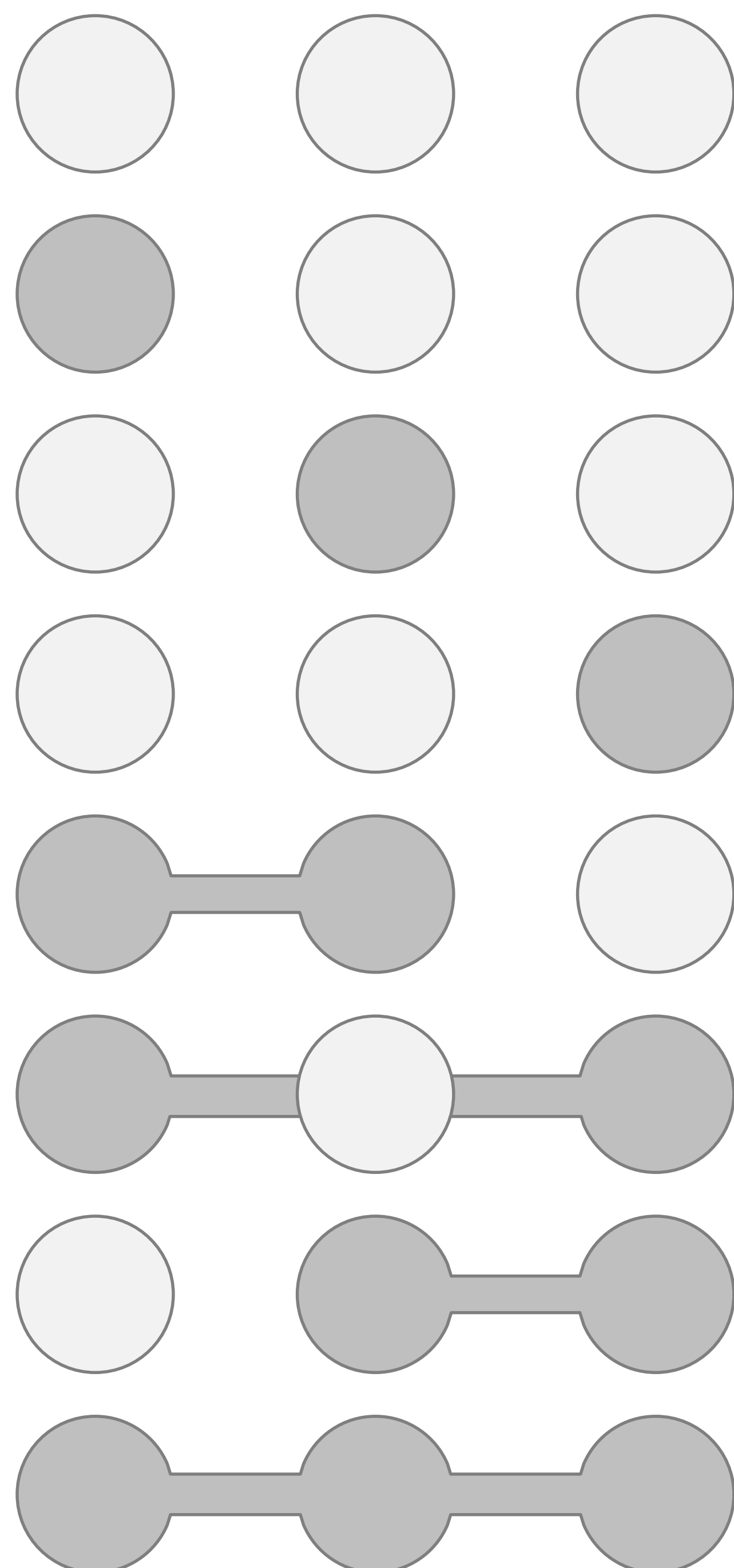
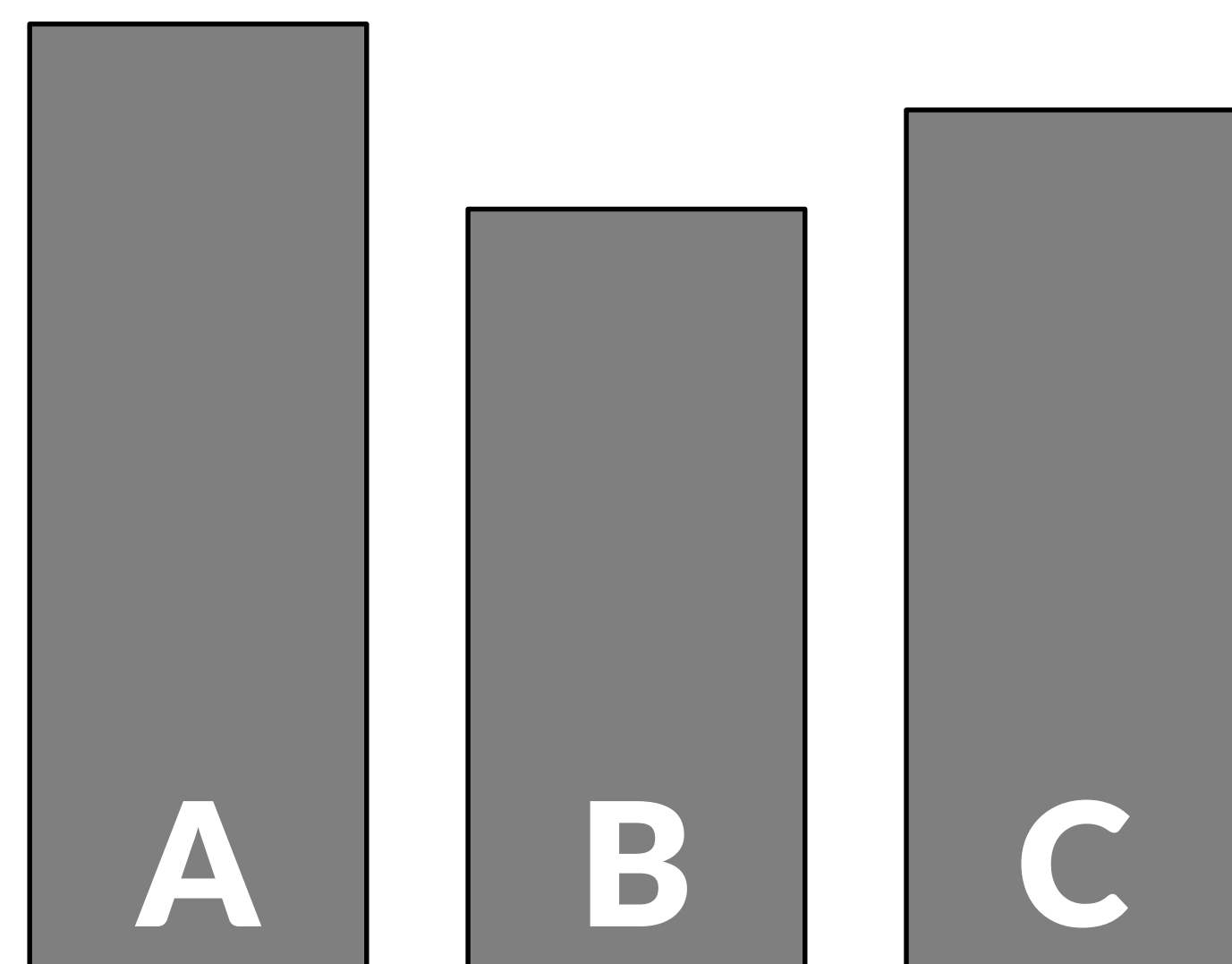


Attributes

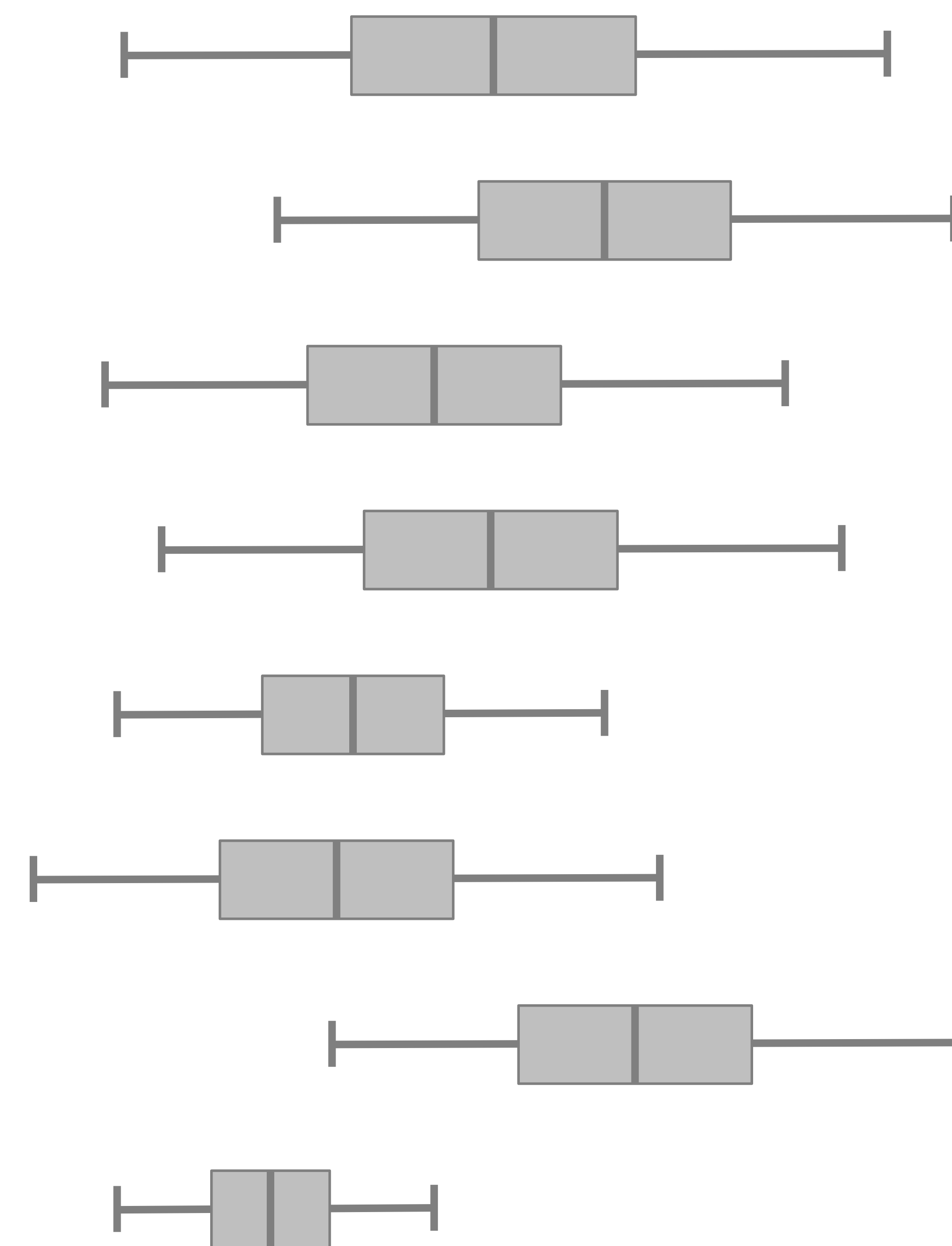


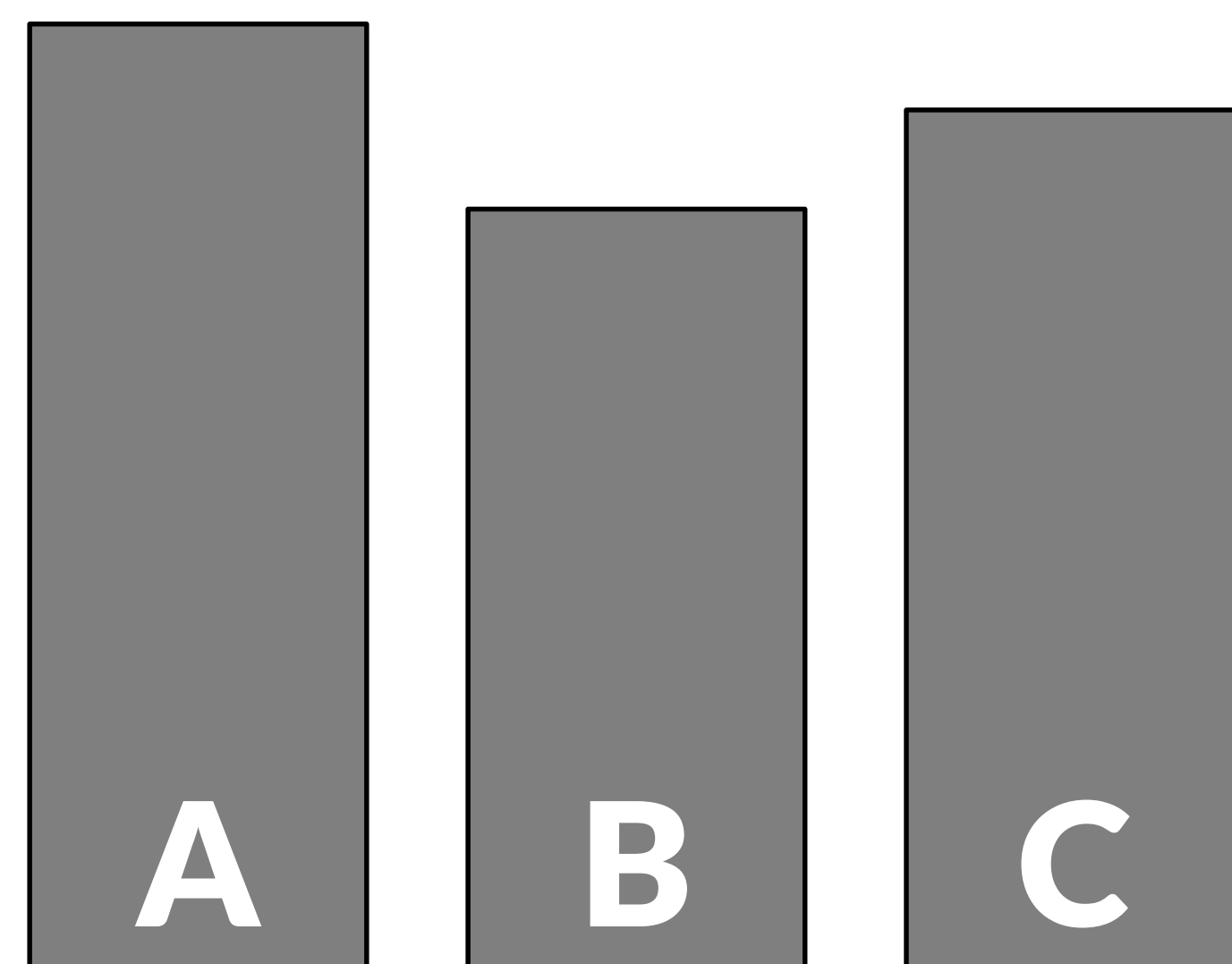


Sorting

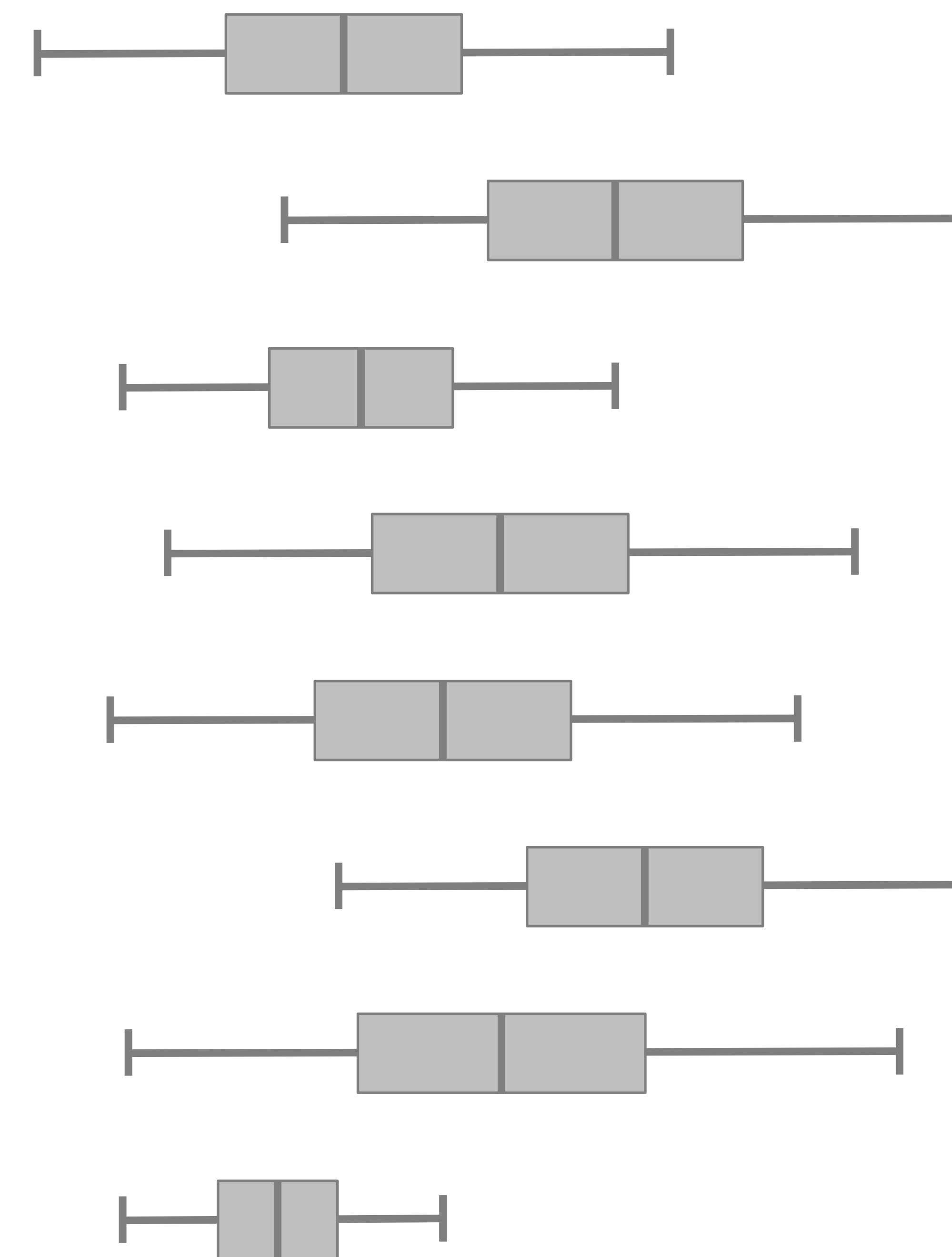
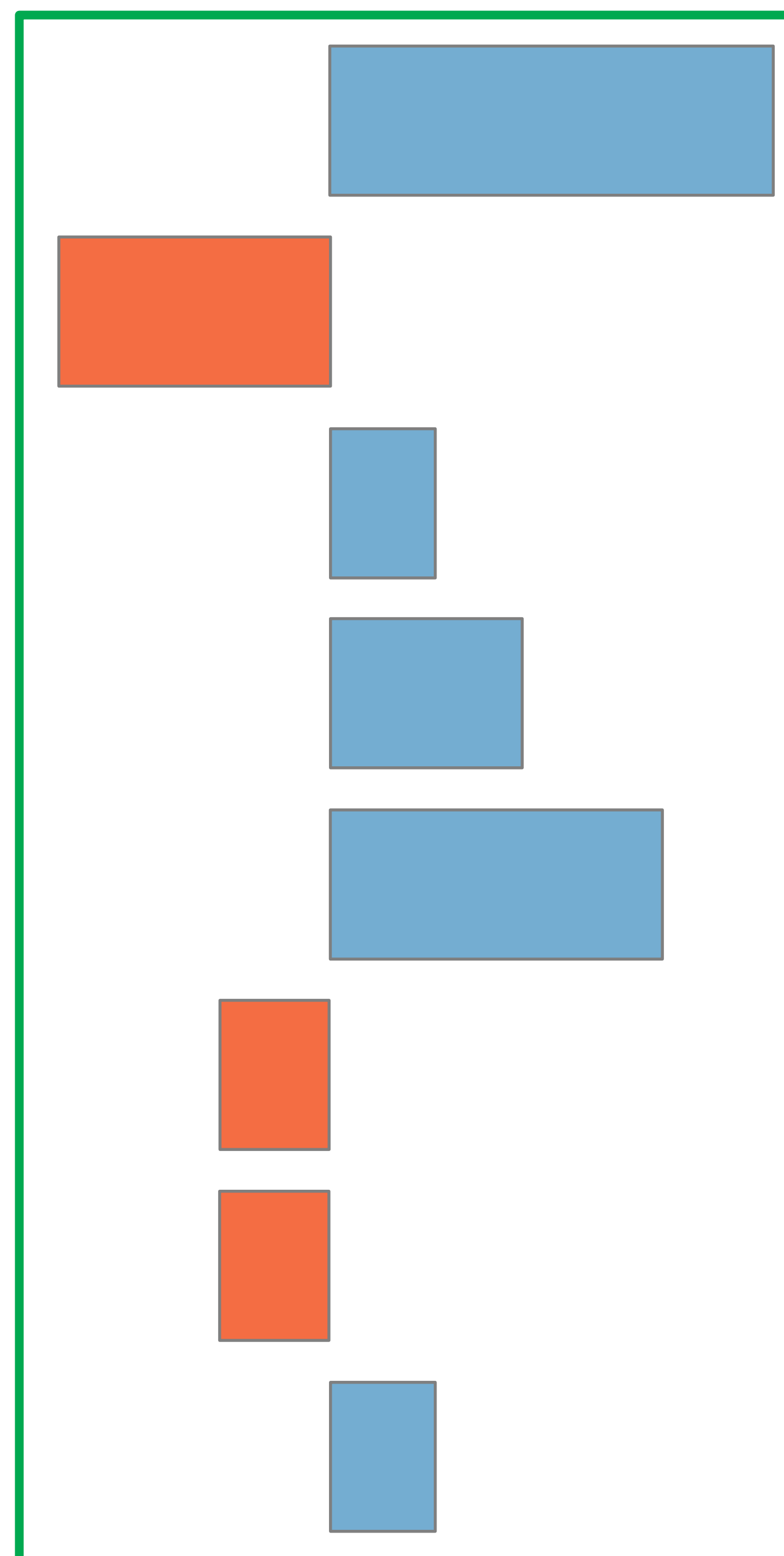
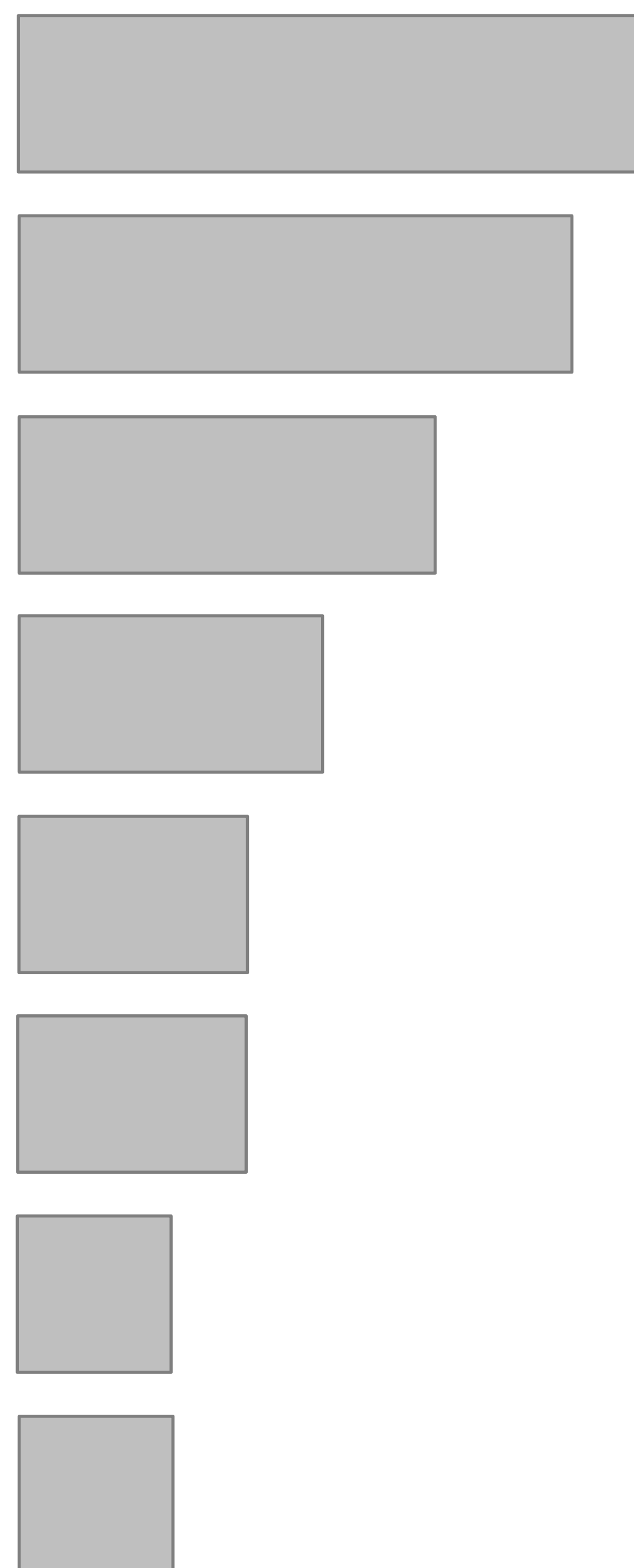
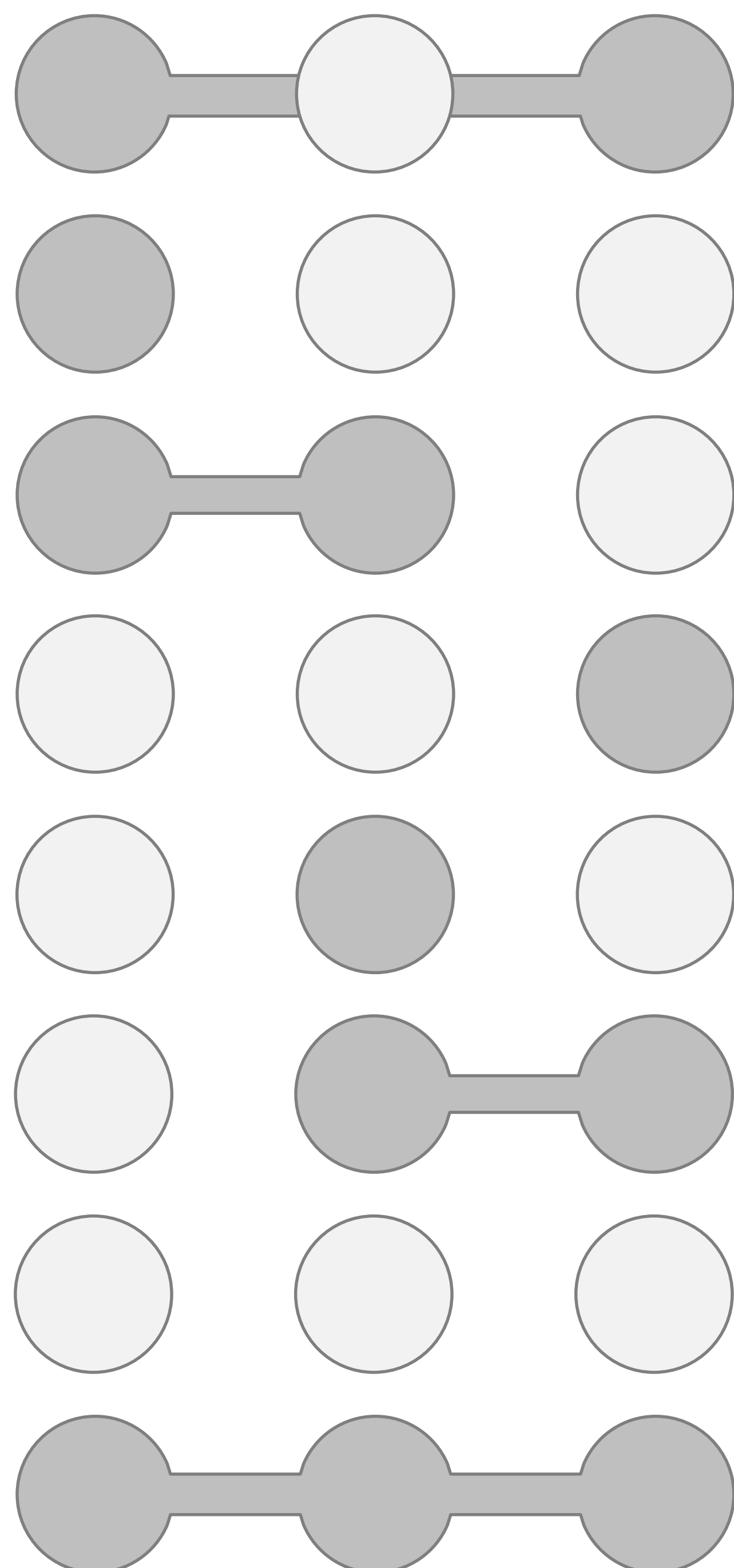


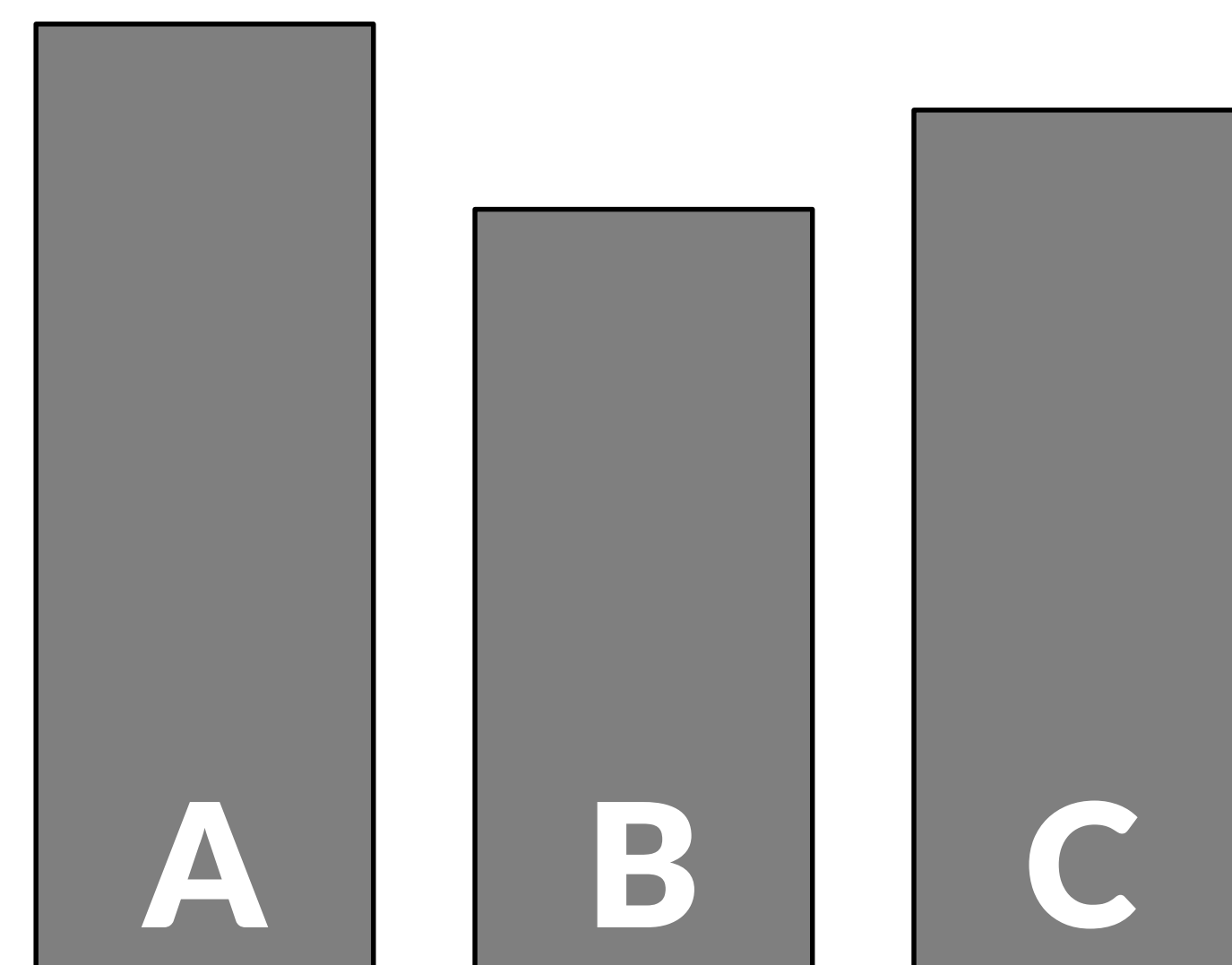
Which is the biggest intersection?
Sort By: Cardinality



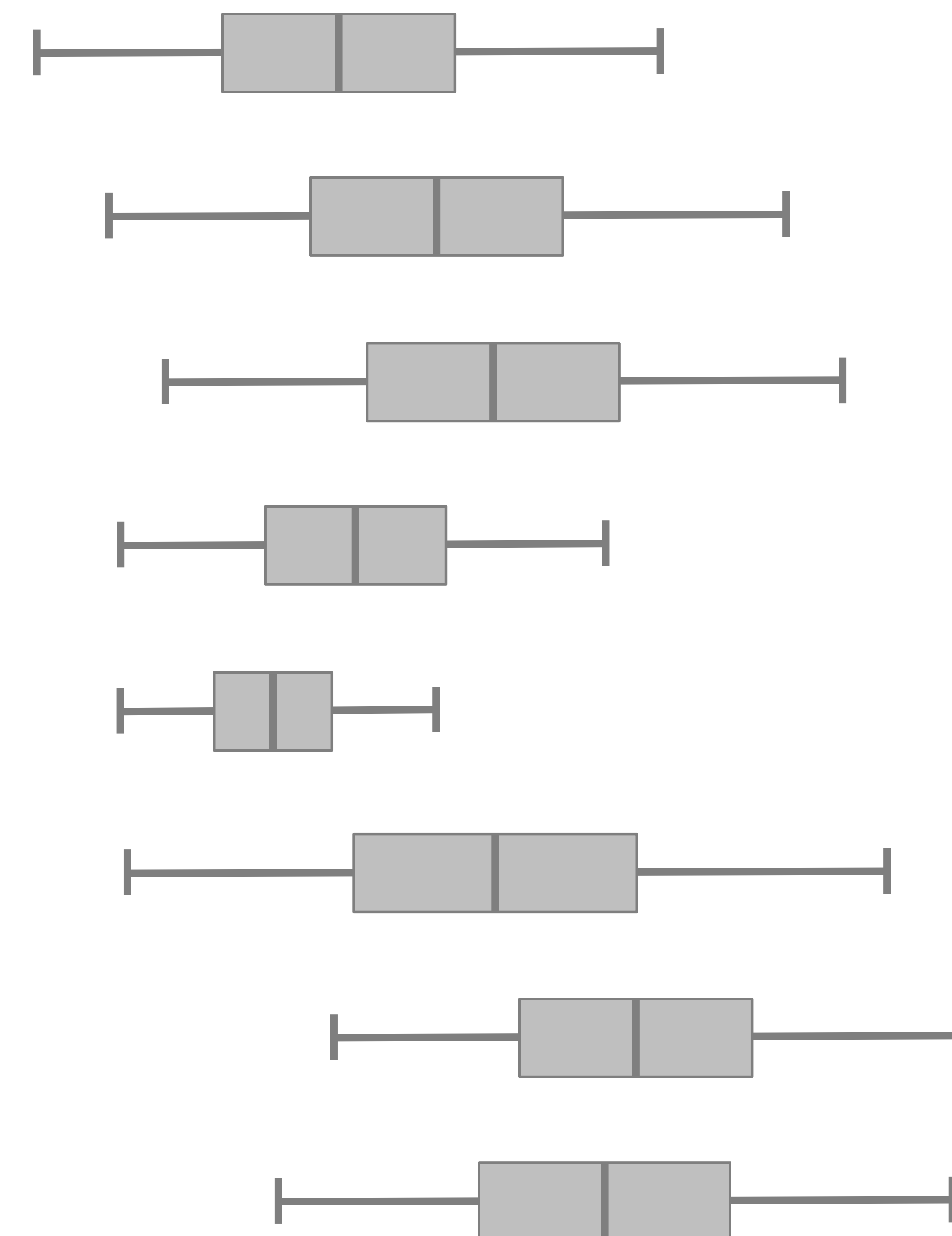
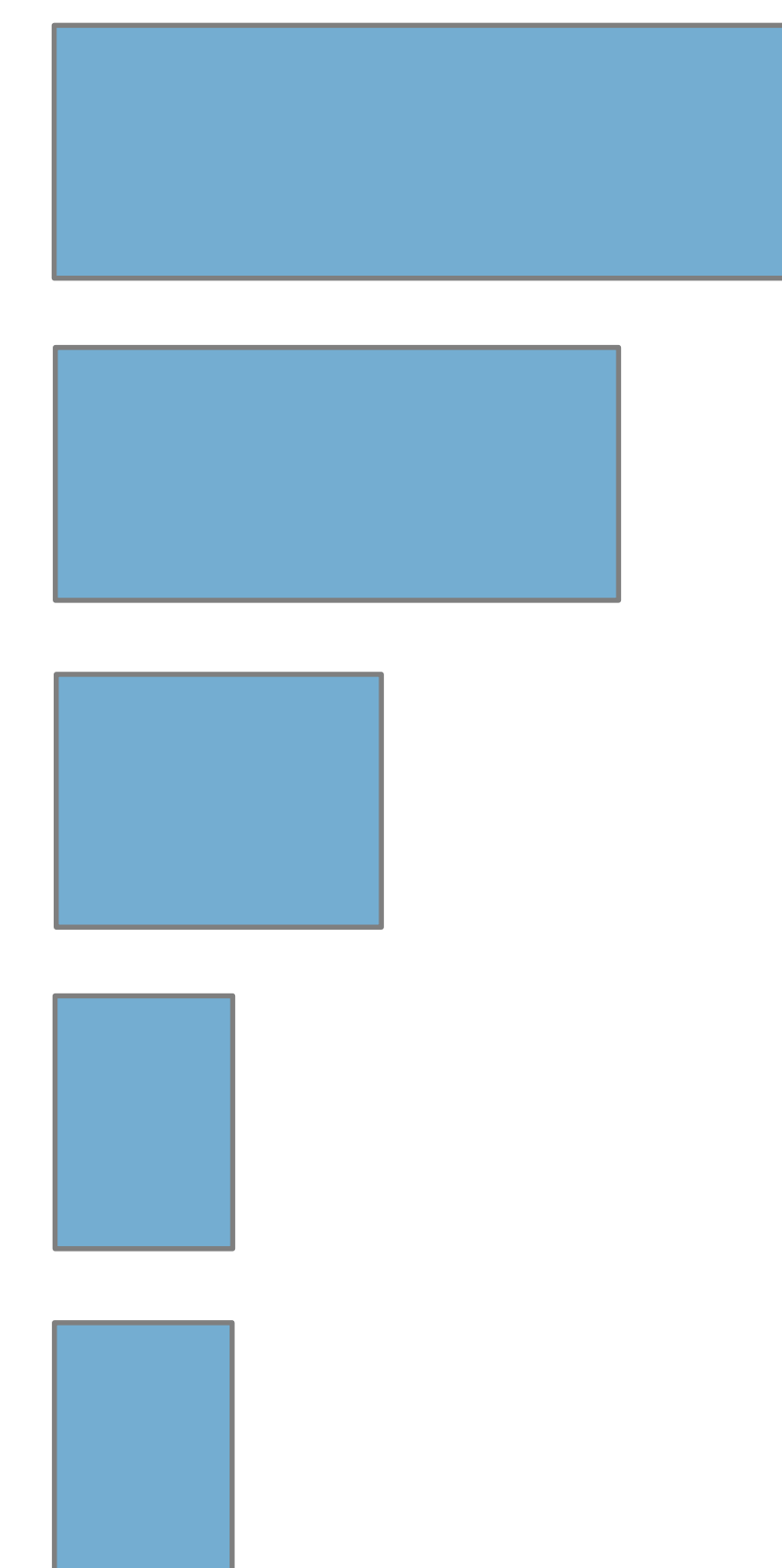
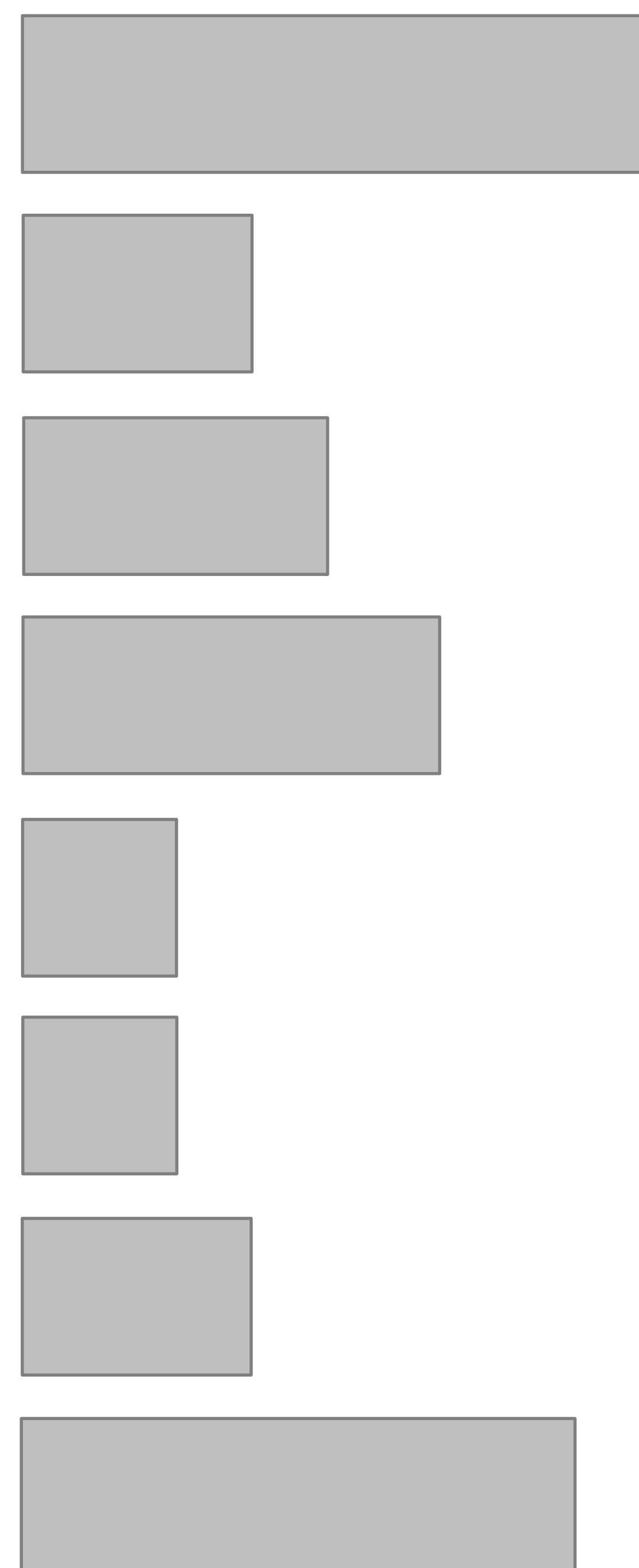
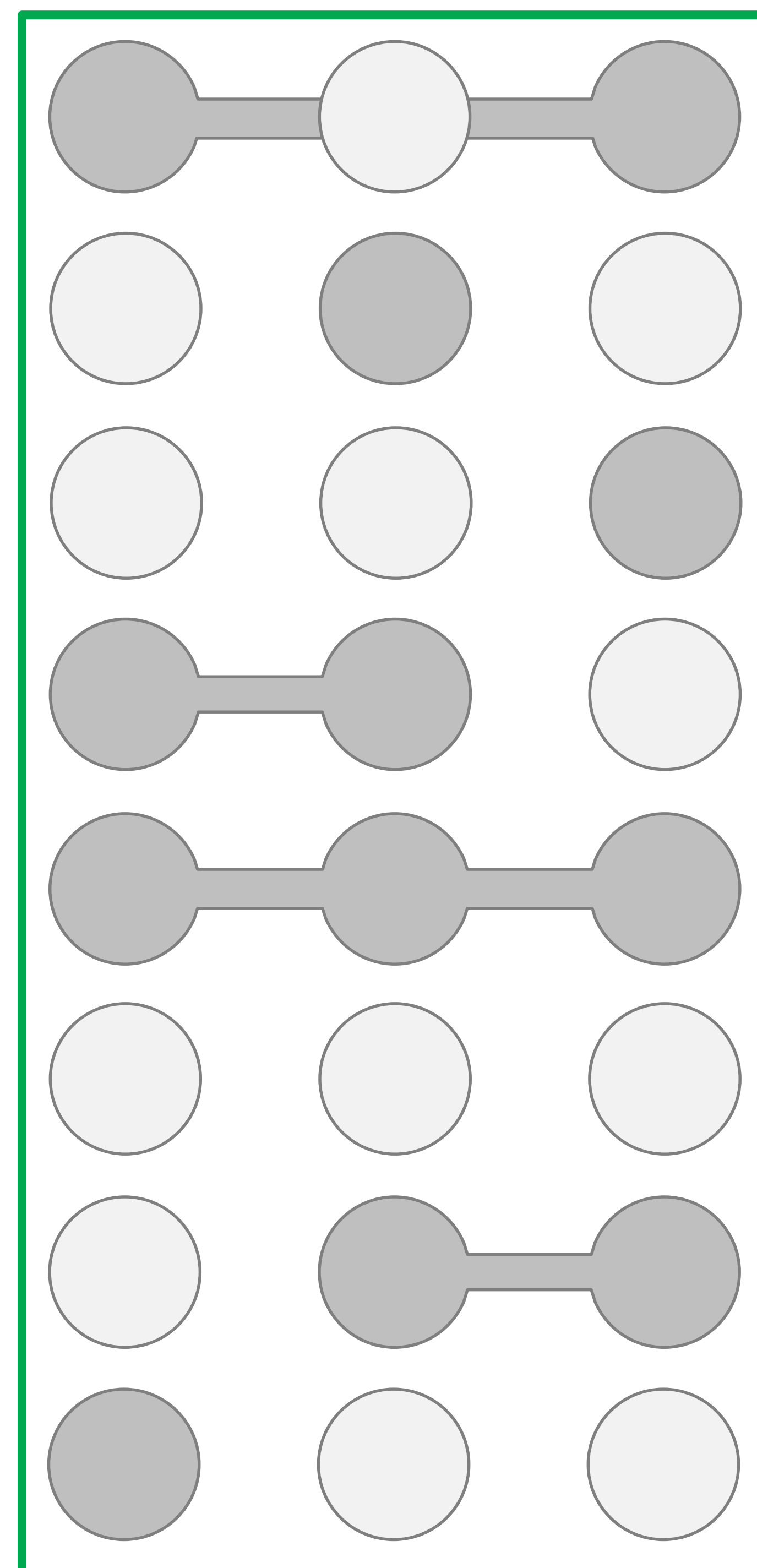


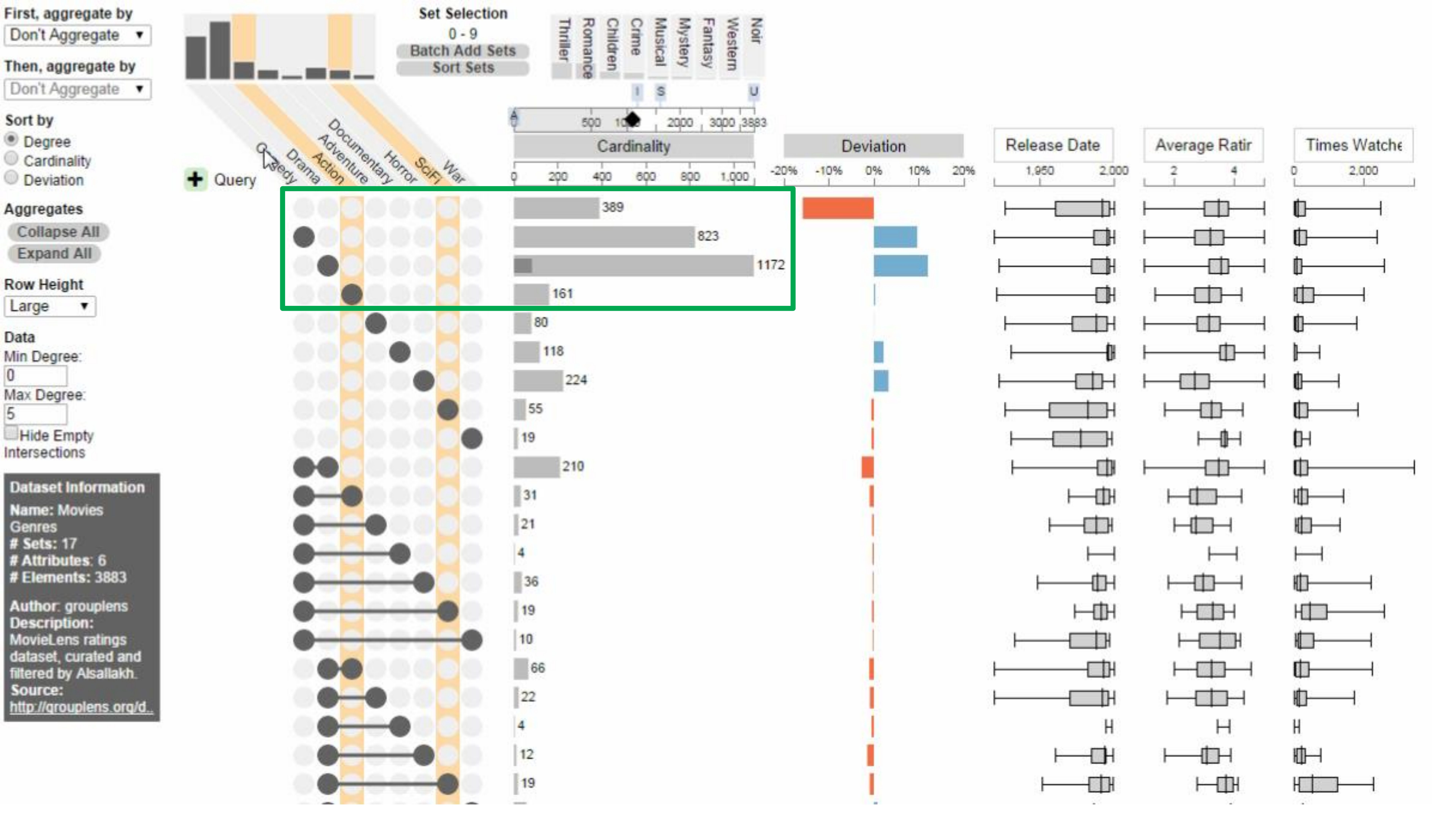
Which is the most 'surprising' intersection?
Sort By: Deviation



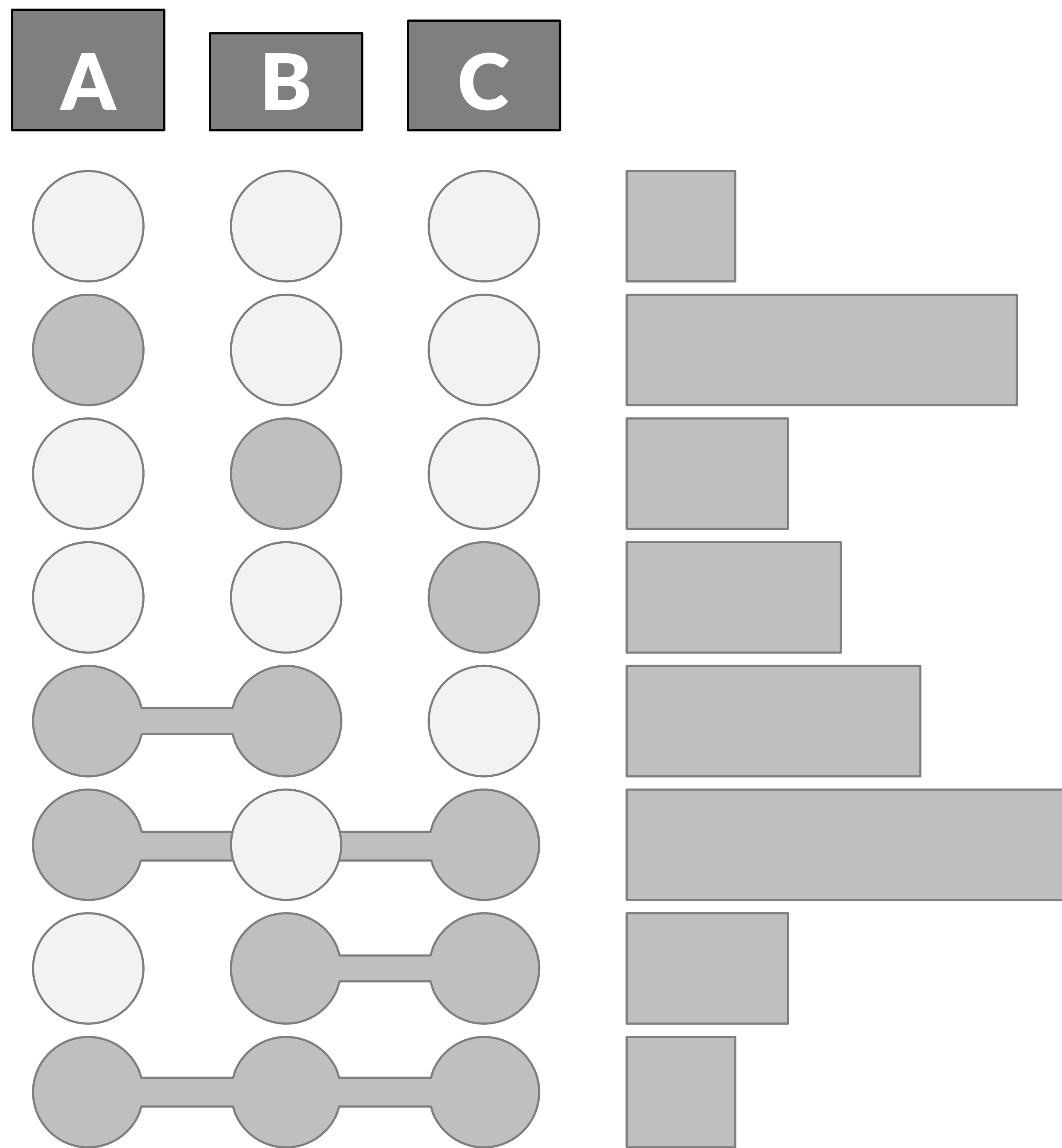


What are the properties of the intersections involving 'A'?
Sort By: Set

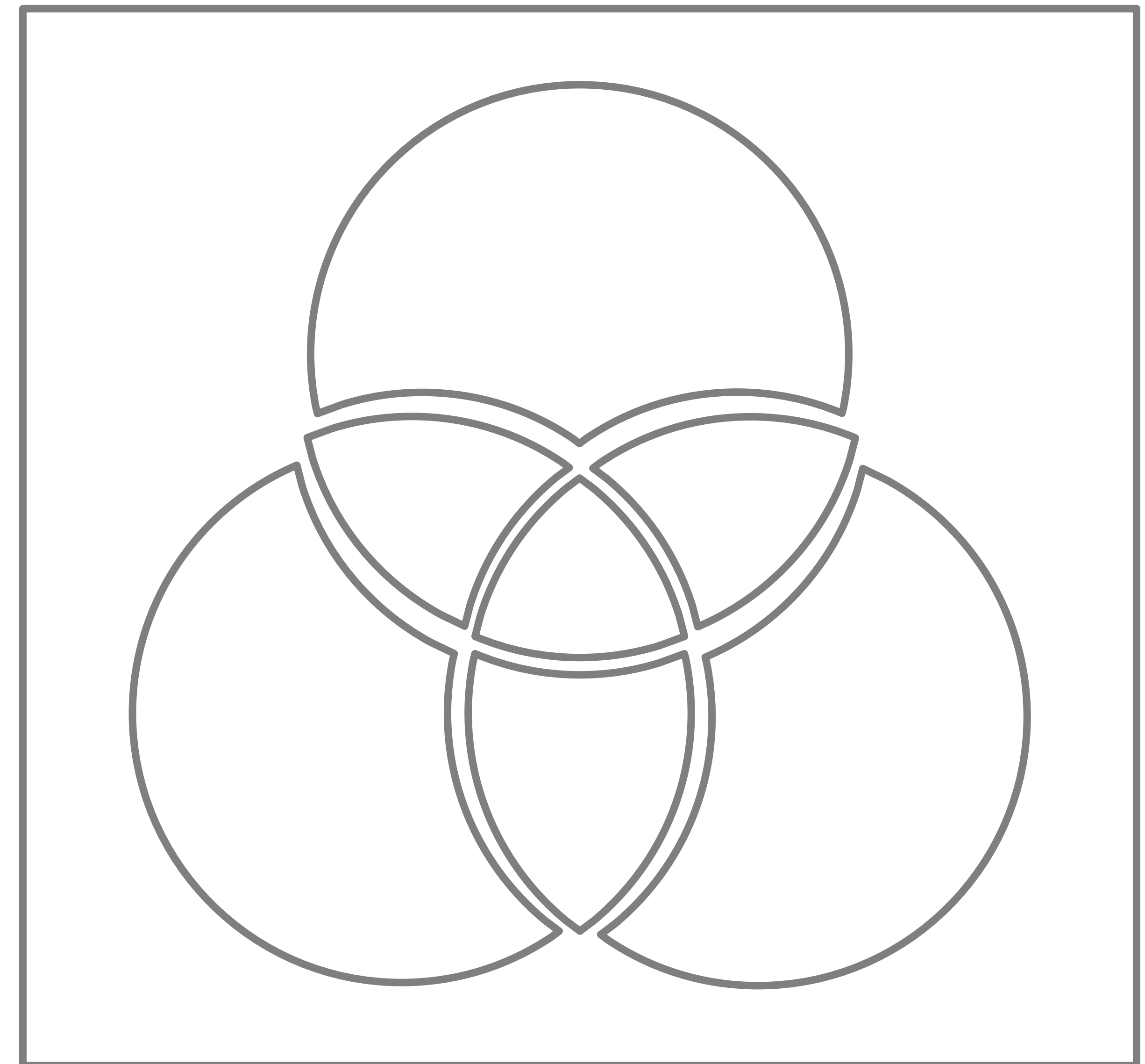


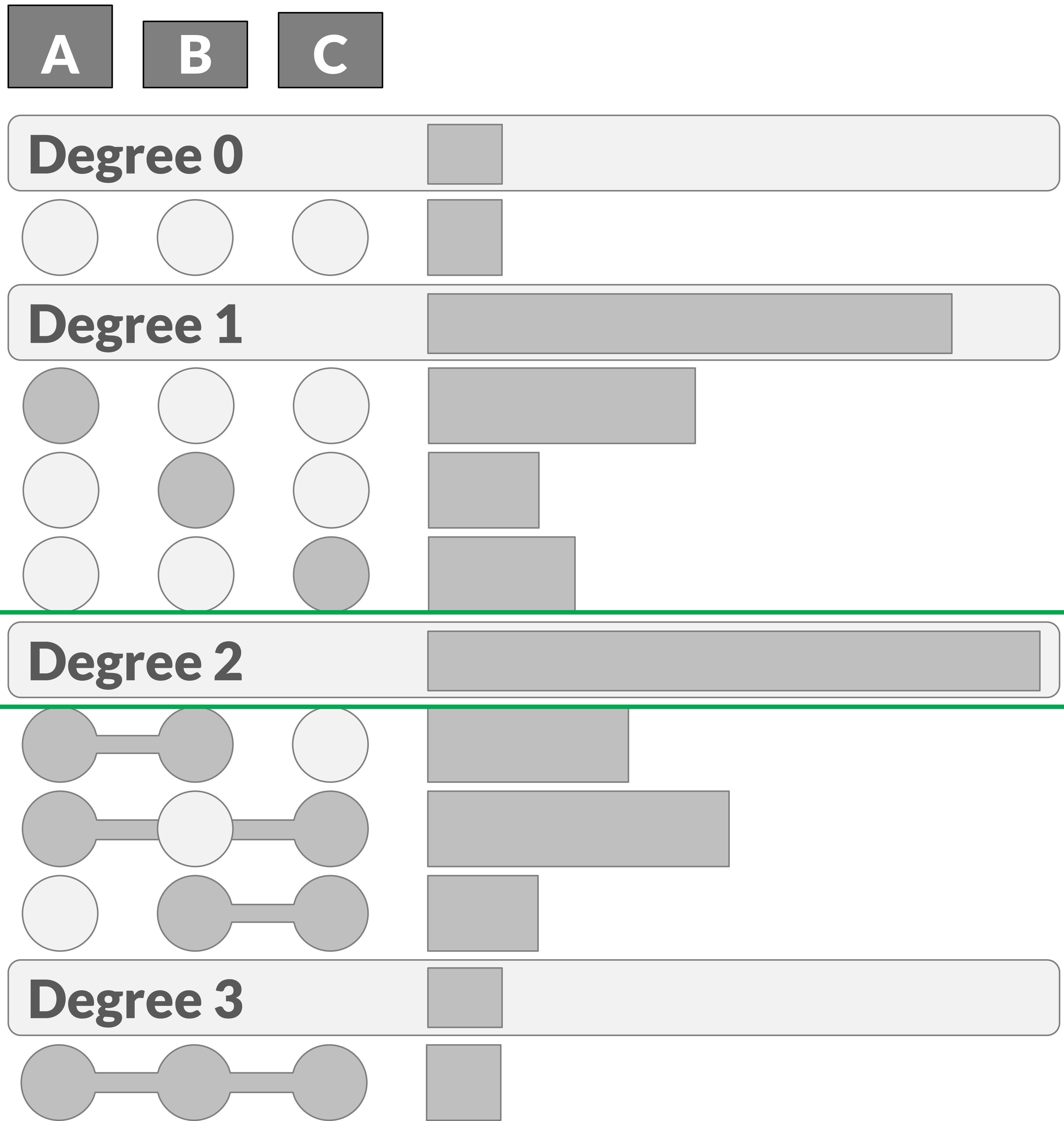


Aggregation

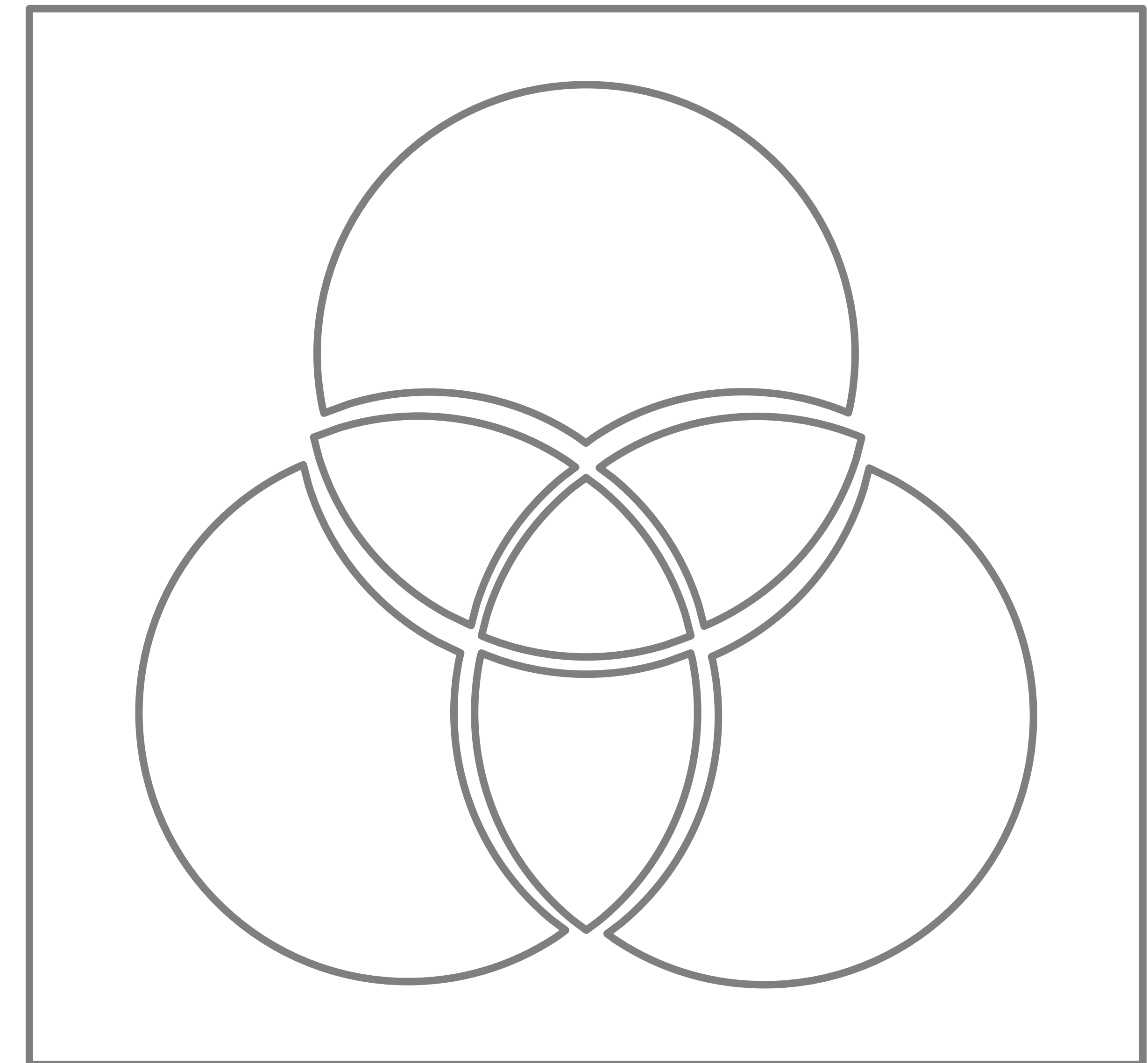


Are most elements in only one set?
Aggregate By: Degree





Are many items shared between two sets?
Aggregate By: Degree



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

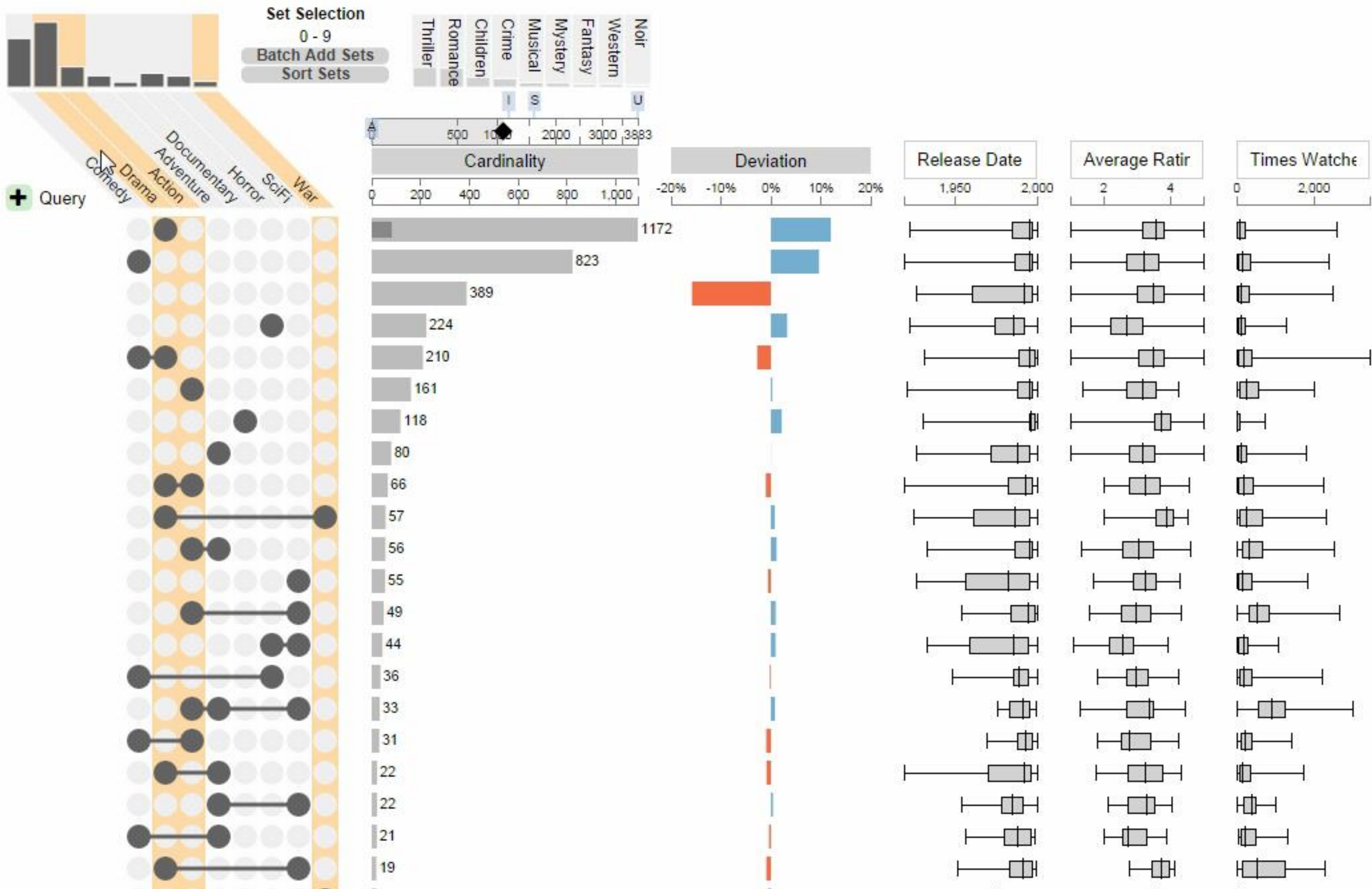
Author: grouplens

Description:

MovieLens ratings dataset, curated and filtered by Alsallakh.

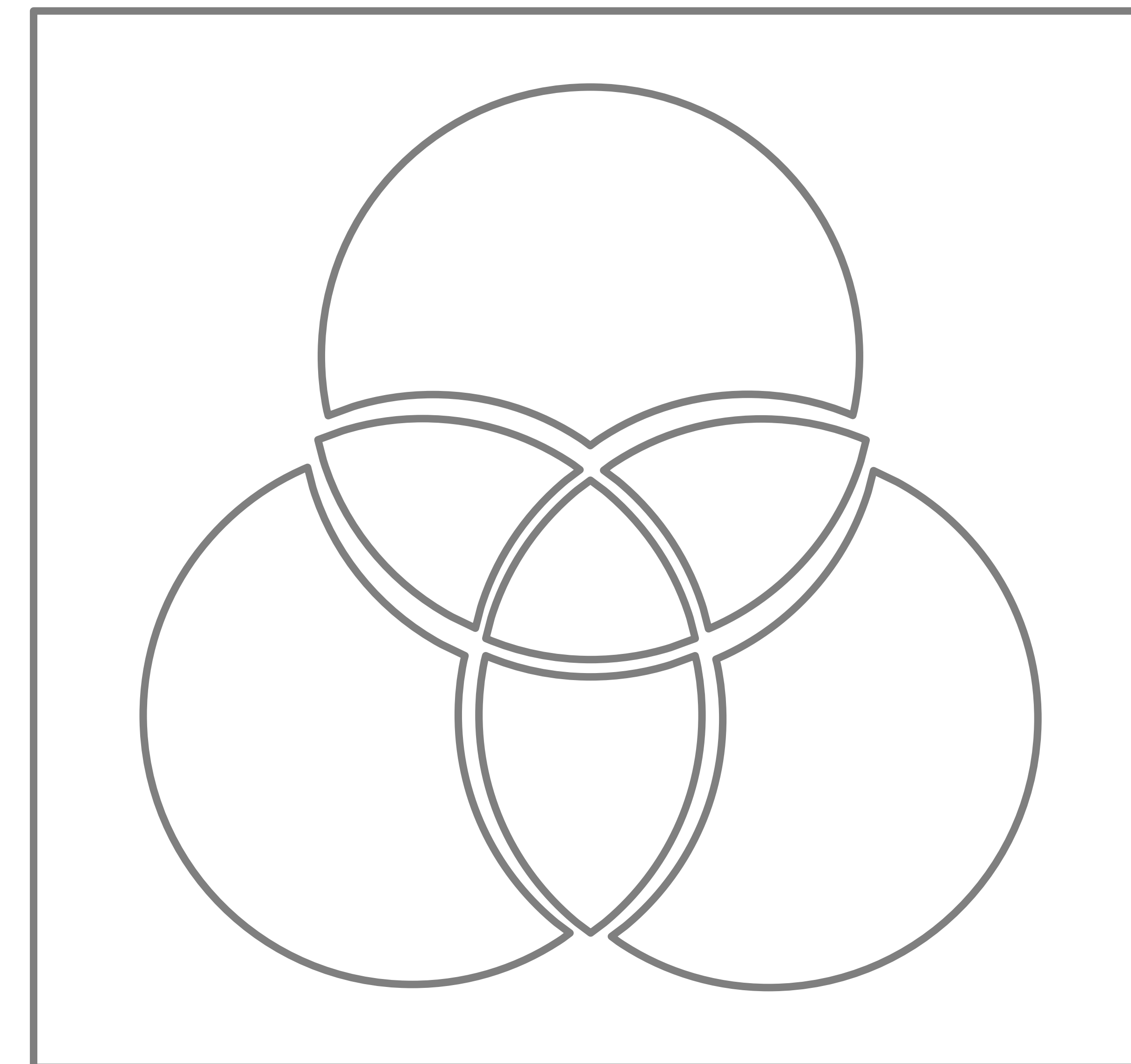
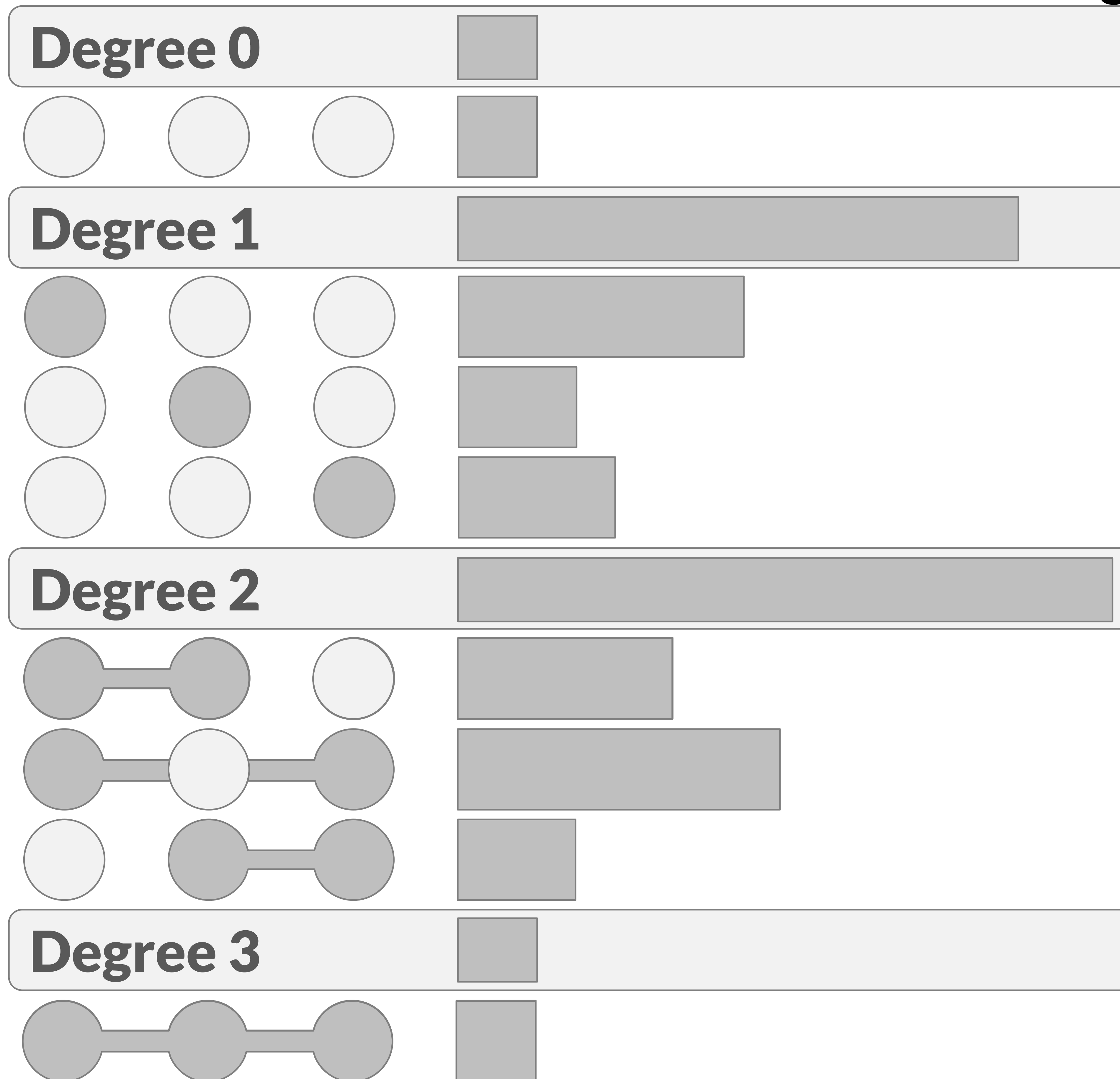
Source:

<http://grouplens.org/d...>

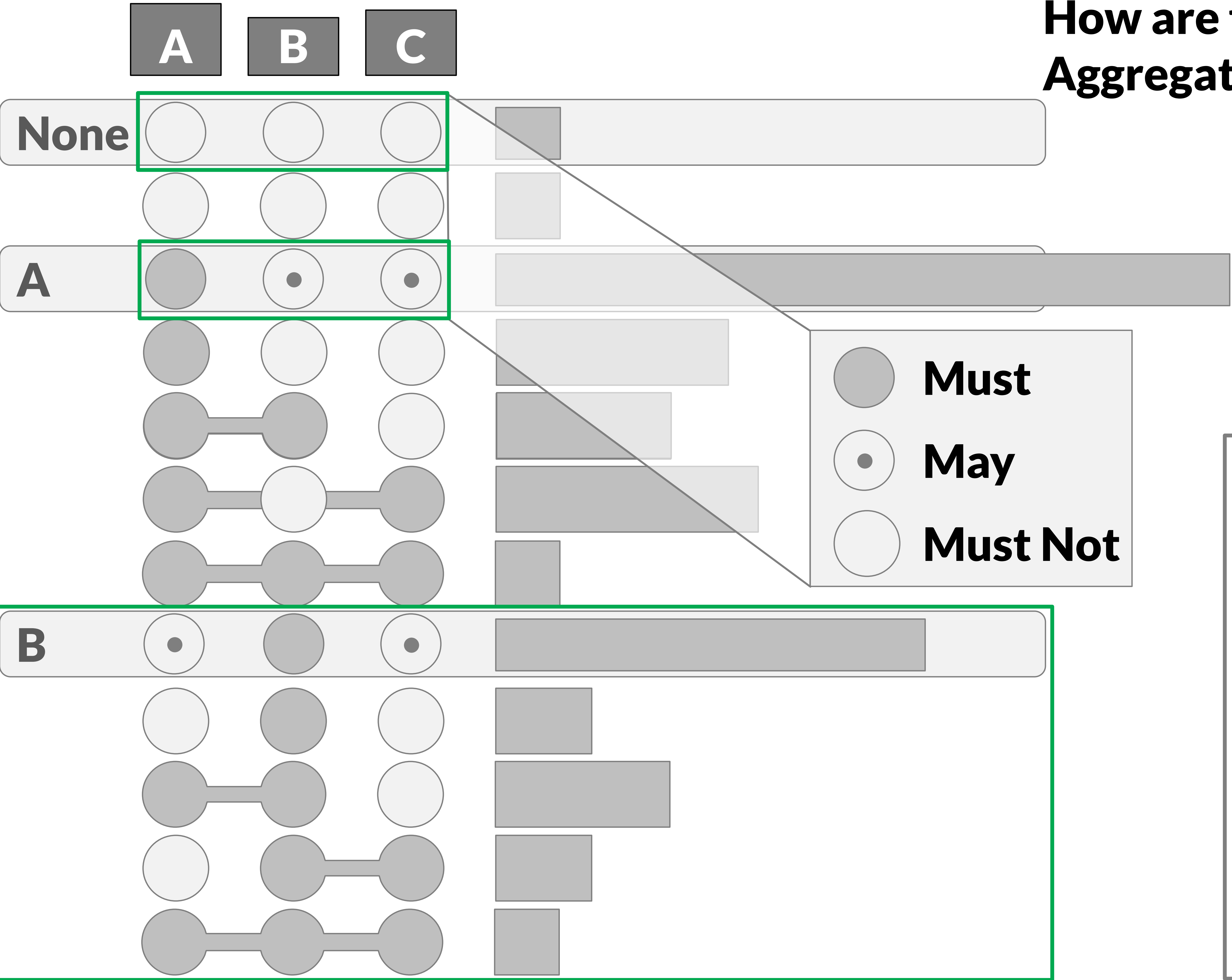




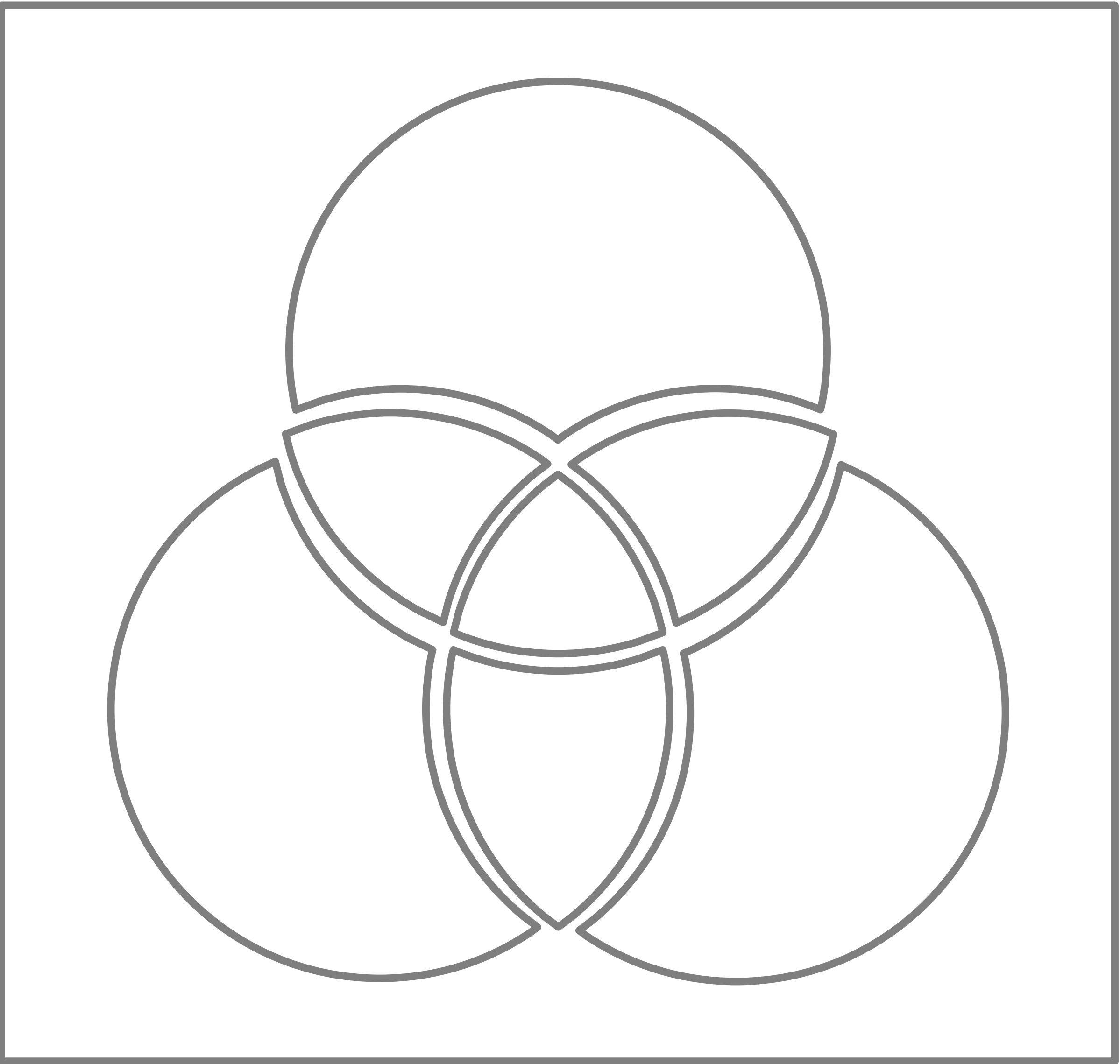
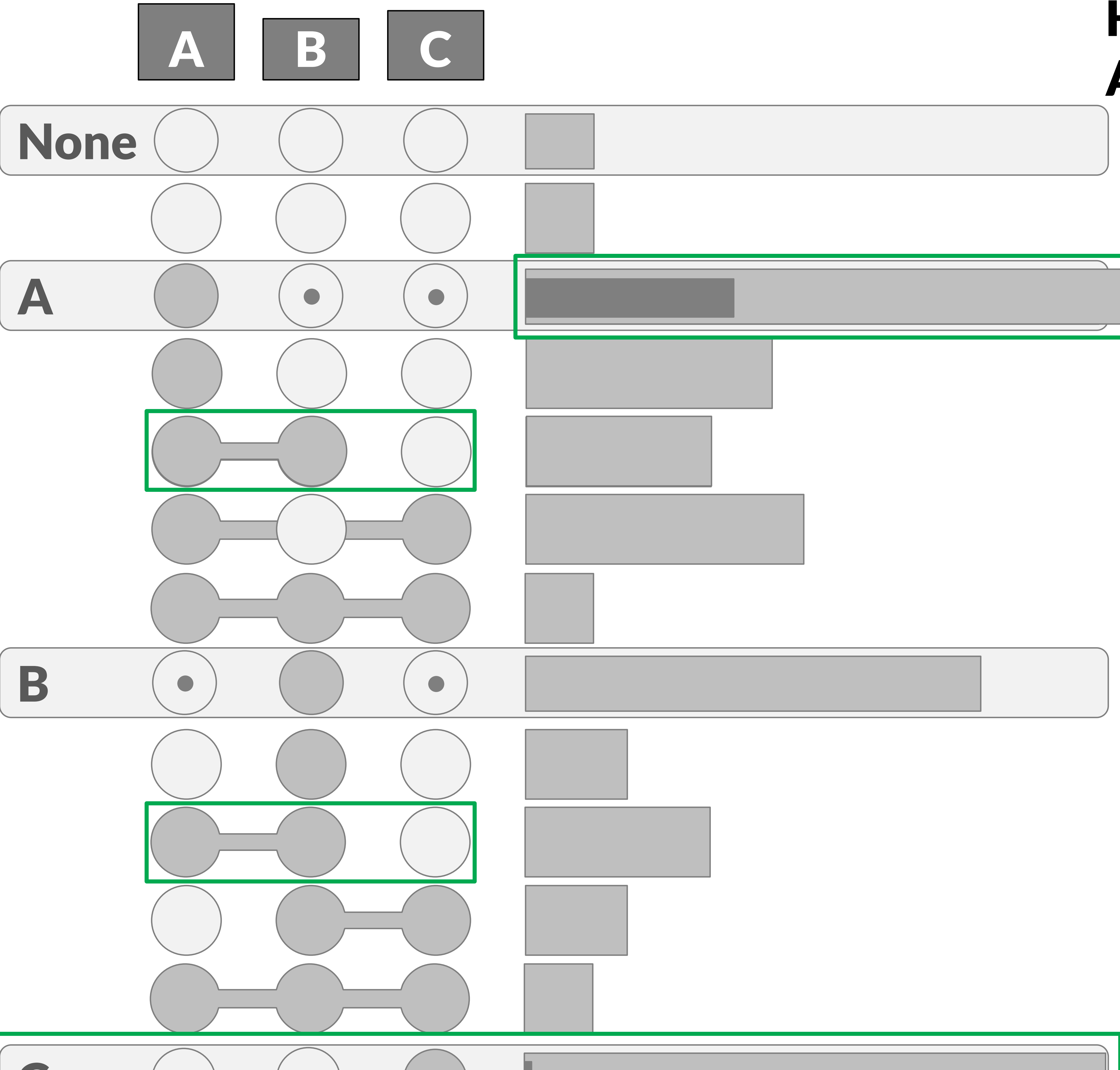
How are the elements of 'B' distributed?
Aggregate By: Set



How are the items of 'B' distributed?
Aggregate By: Set



How are the items of 'B' distributed?
Aggregate By: Set



First, aggregate by
Degree

Then, aggregate by
Don't Aggregate

Sort by
☒ Degree
☐ Cardinality
☐ Deviation

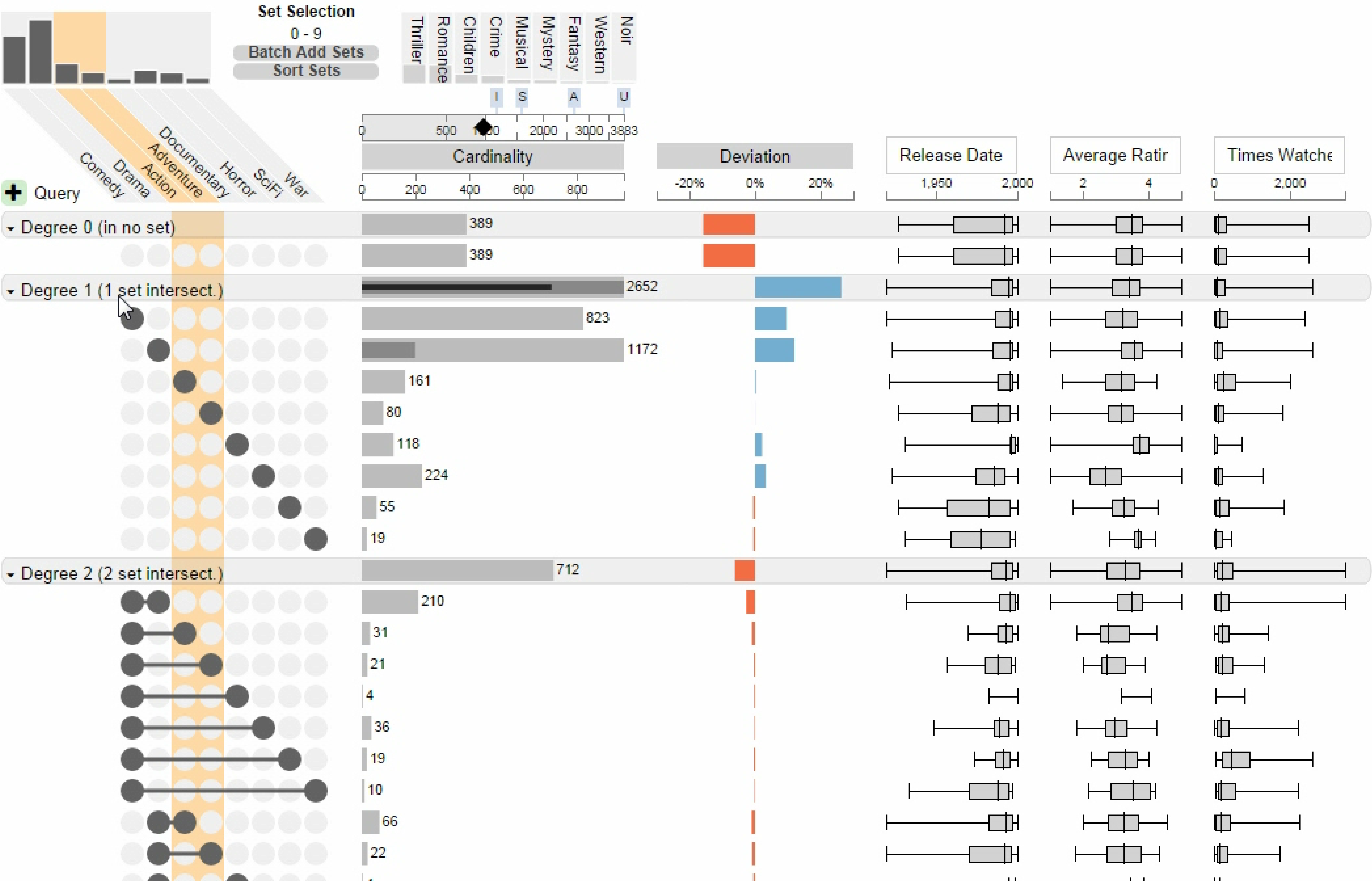
Aggregates

Row Height
Large

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

Dataset Information
Name: Movies
Genres
Sets: 17
Attributes: 6
Elements: 3883

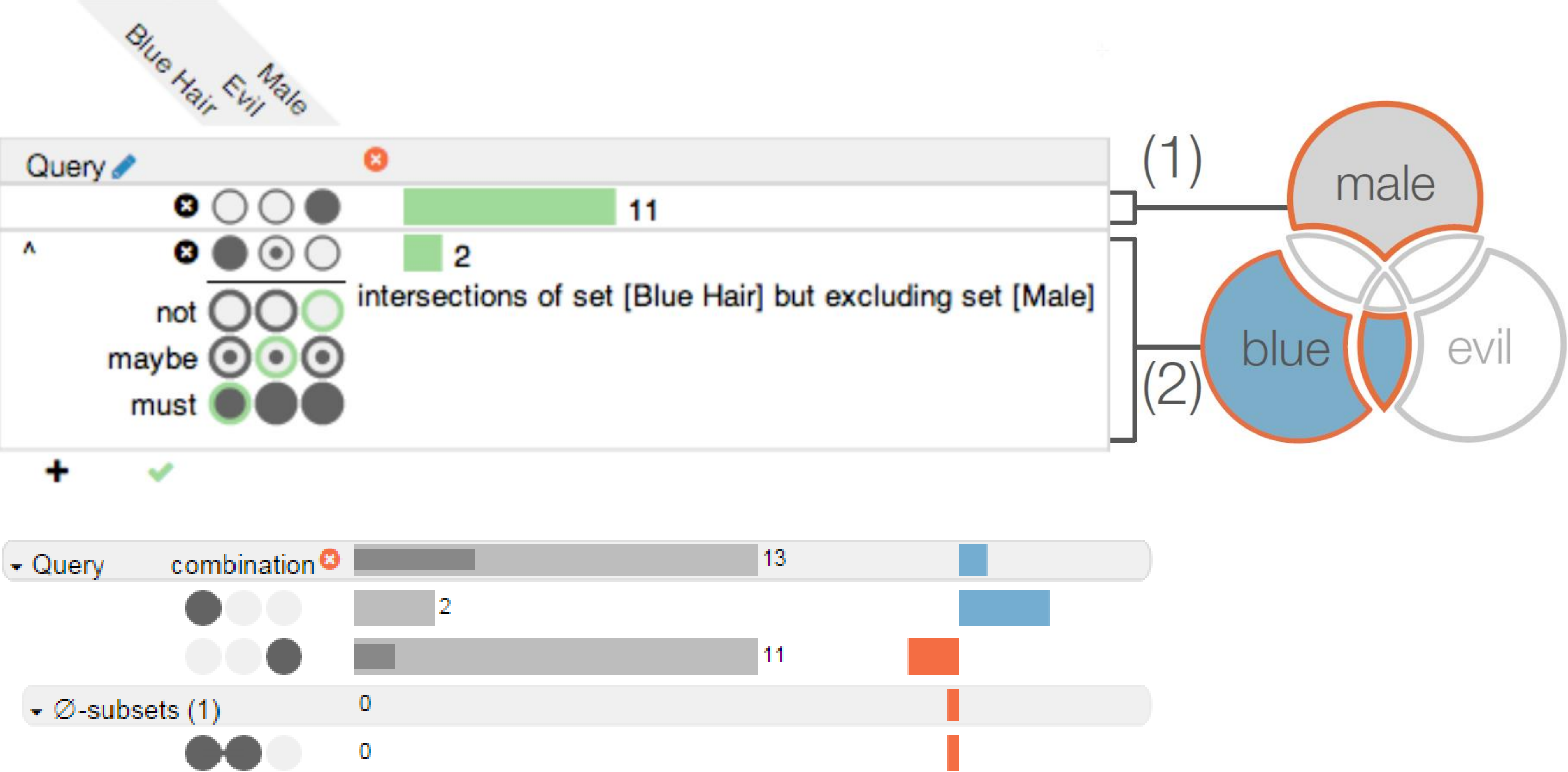
Author: grouplens
Description:
MovieLens ratings
dataset, curated and
filtered by Alsallakh.
Source:
<http://grouplens.org/d..>

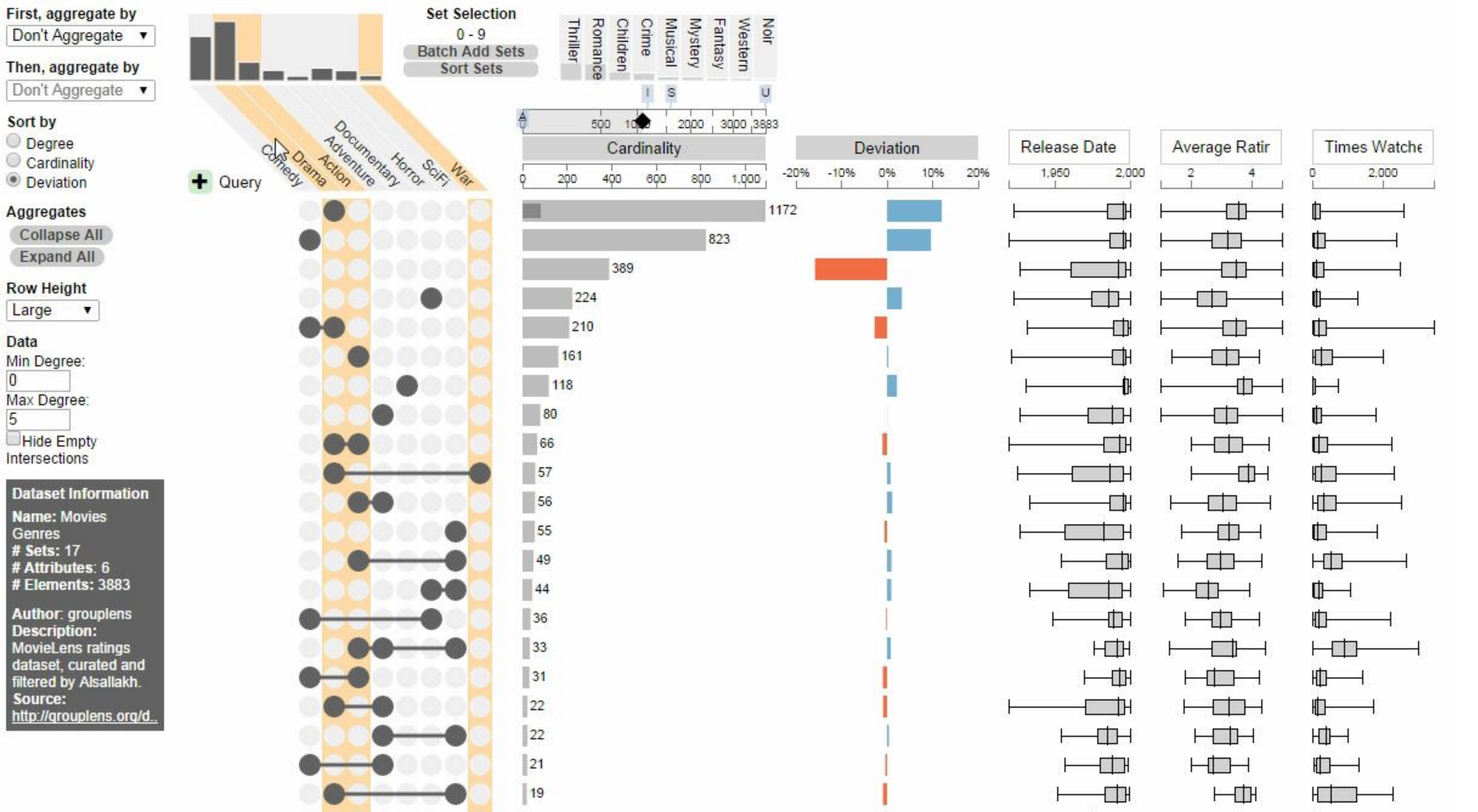


Queries

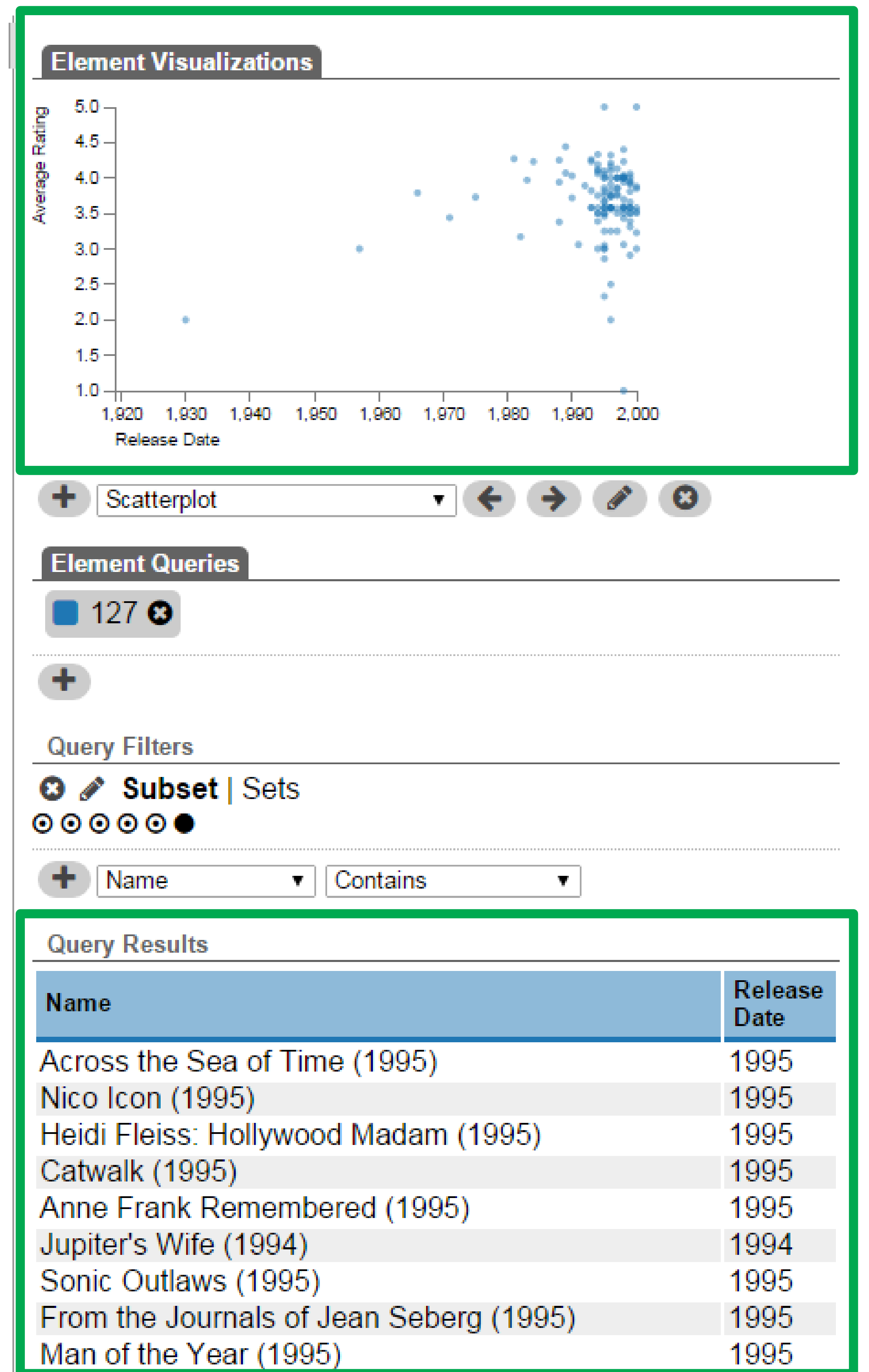
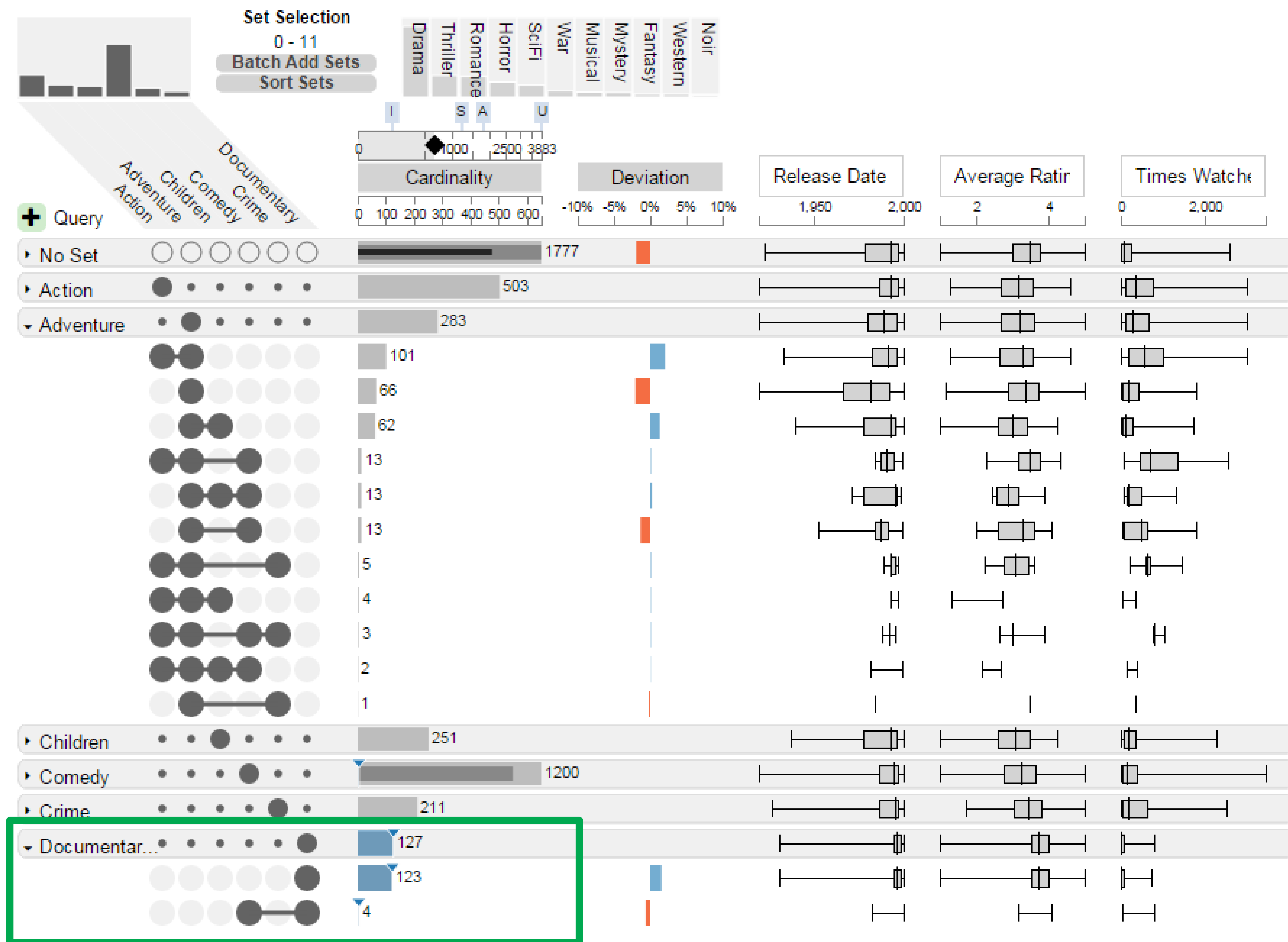
Which Simpsons characters are either male but not evil and don't have blue hair, or are female and have blue hair?

Answer with Queries

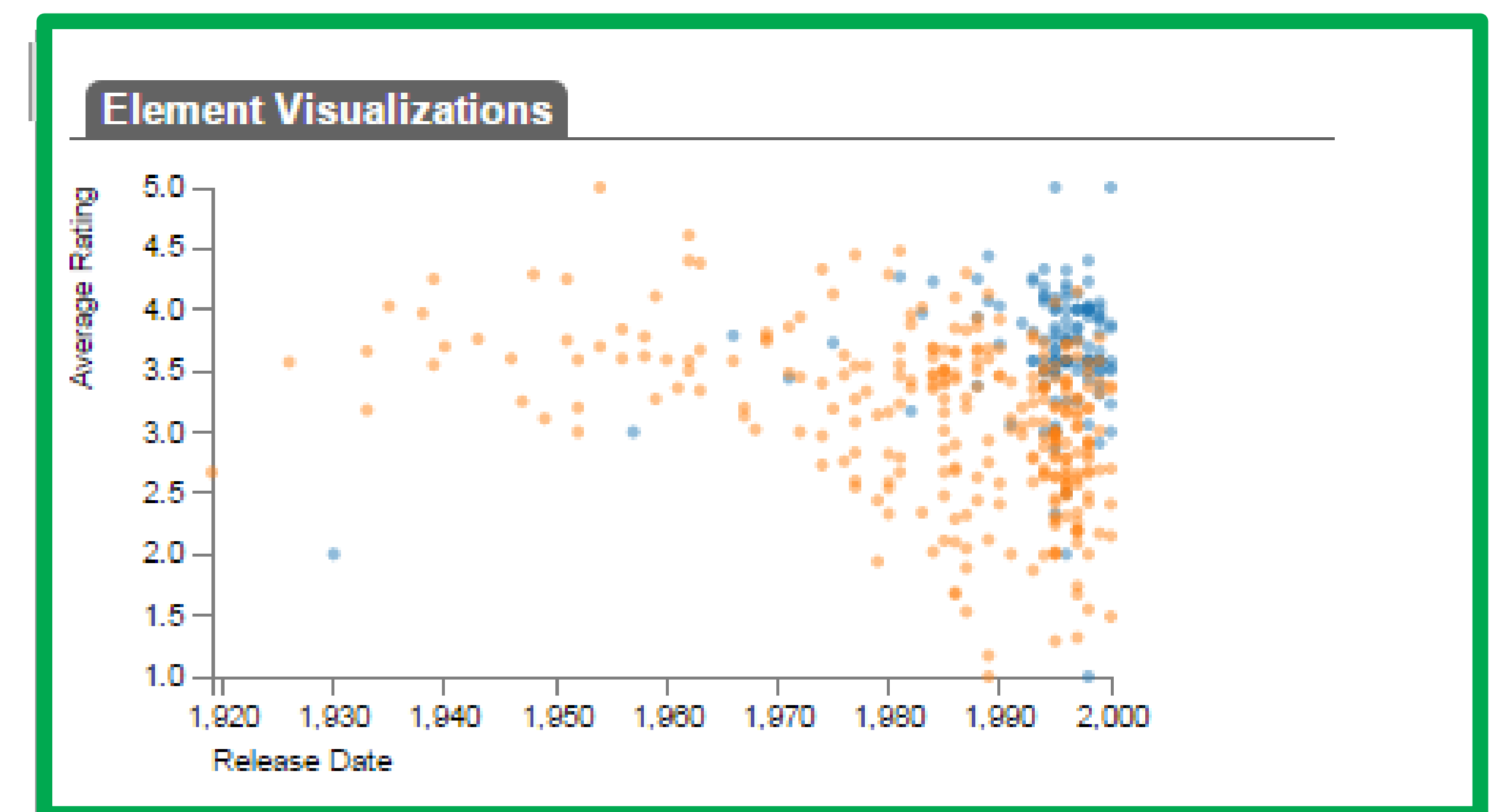
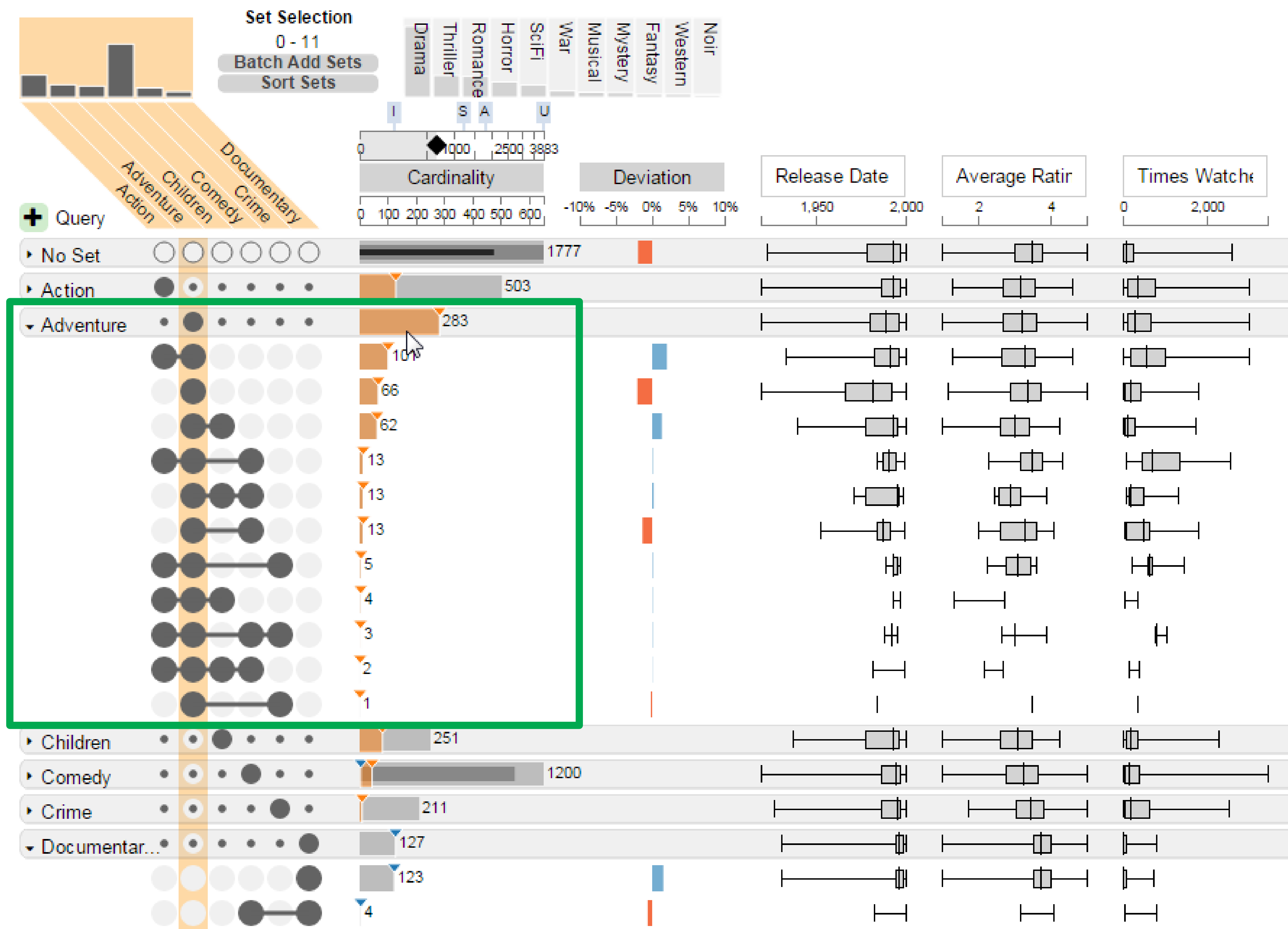




Exploring Attributes



How do the attribute of two intersections compare?



Scatterplot

Element Queries

127 283

+

Query Filters

Subset | Sets

+

Name Contains

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 the attribute of two intersections compare?

Applications

Molecular Biology

Social Network Analysis

Economics

Drug Discovery

....

<http://vcg.github.io/upset>

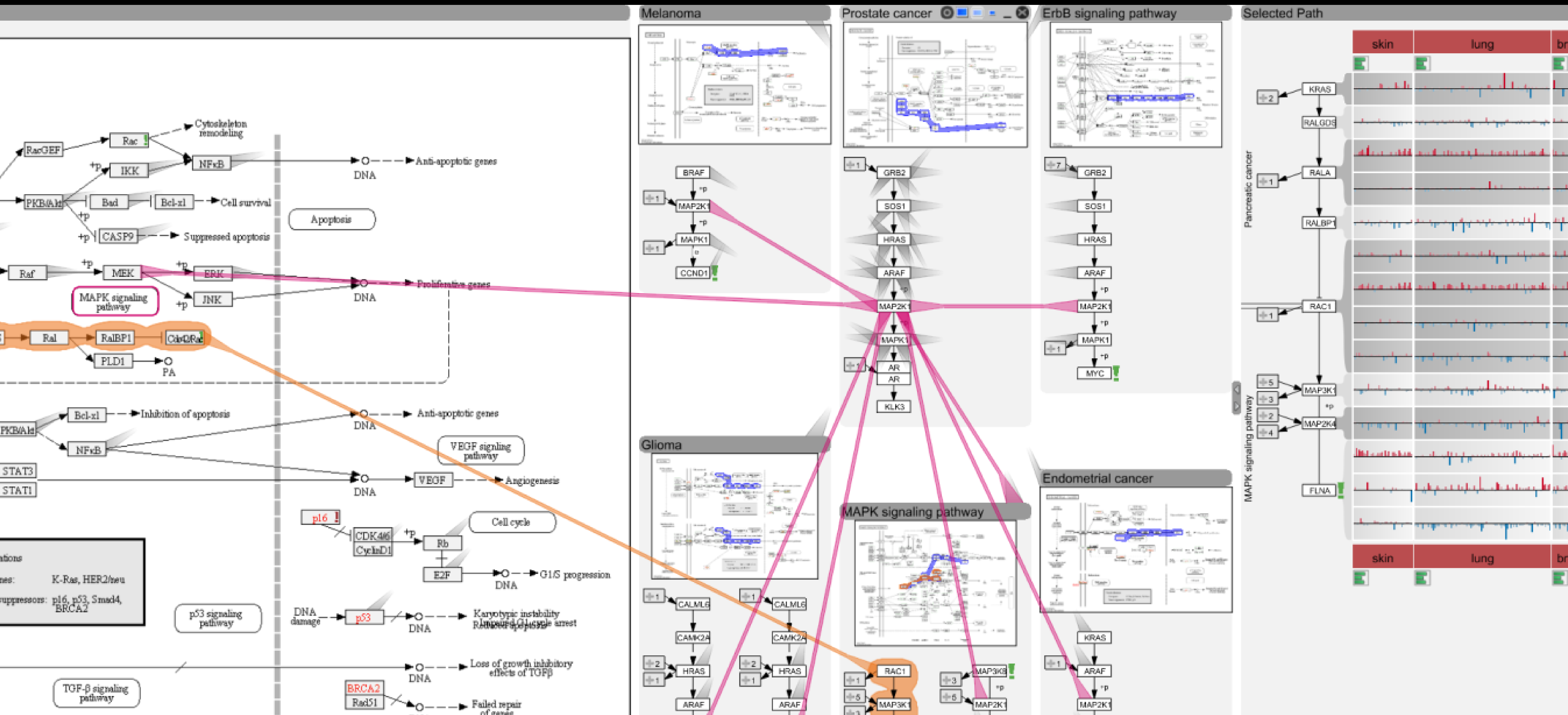
[PartI, BioVis '12]

Best Paper Award

[PartI, BMC Bioinf. '13]

[Lex, InfoVis '13]

Pathways – Entourage

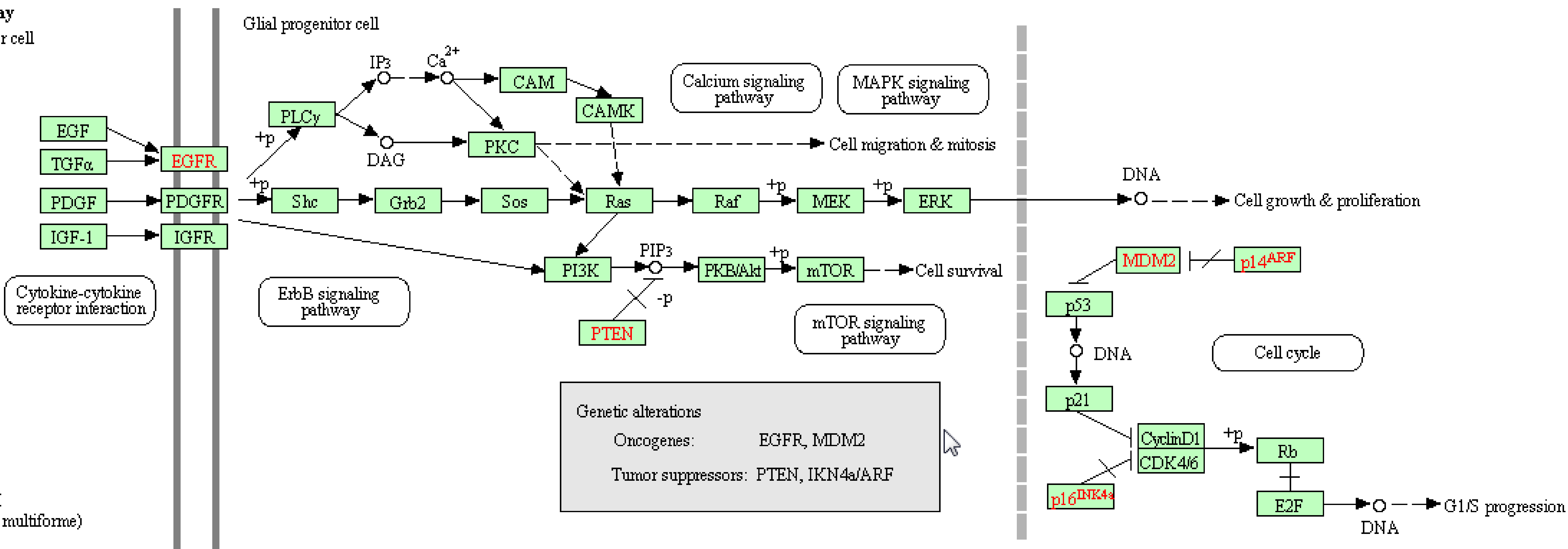


A Pathway

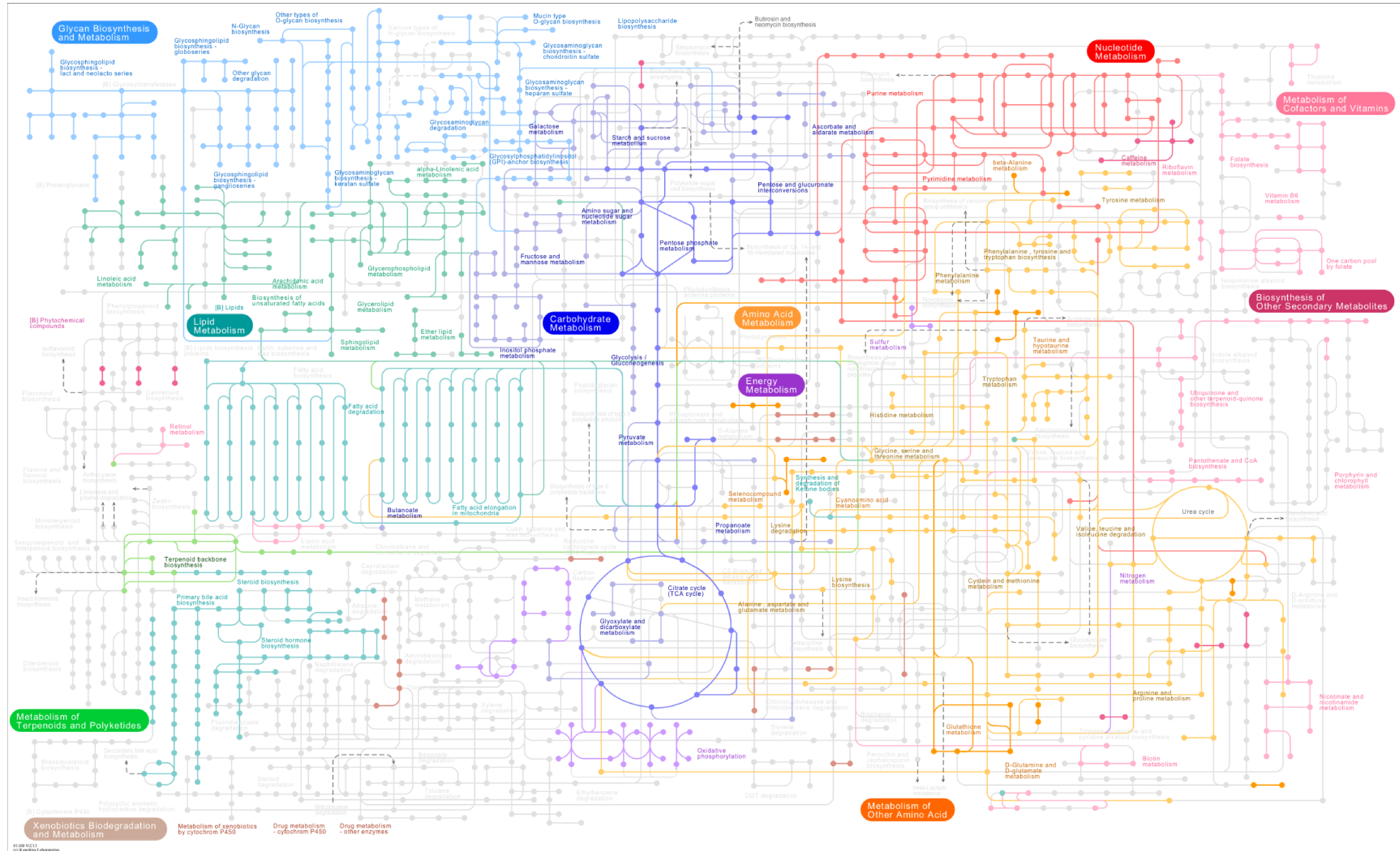
GLIOMA

De Novo pathway

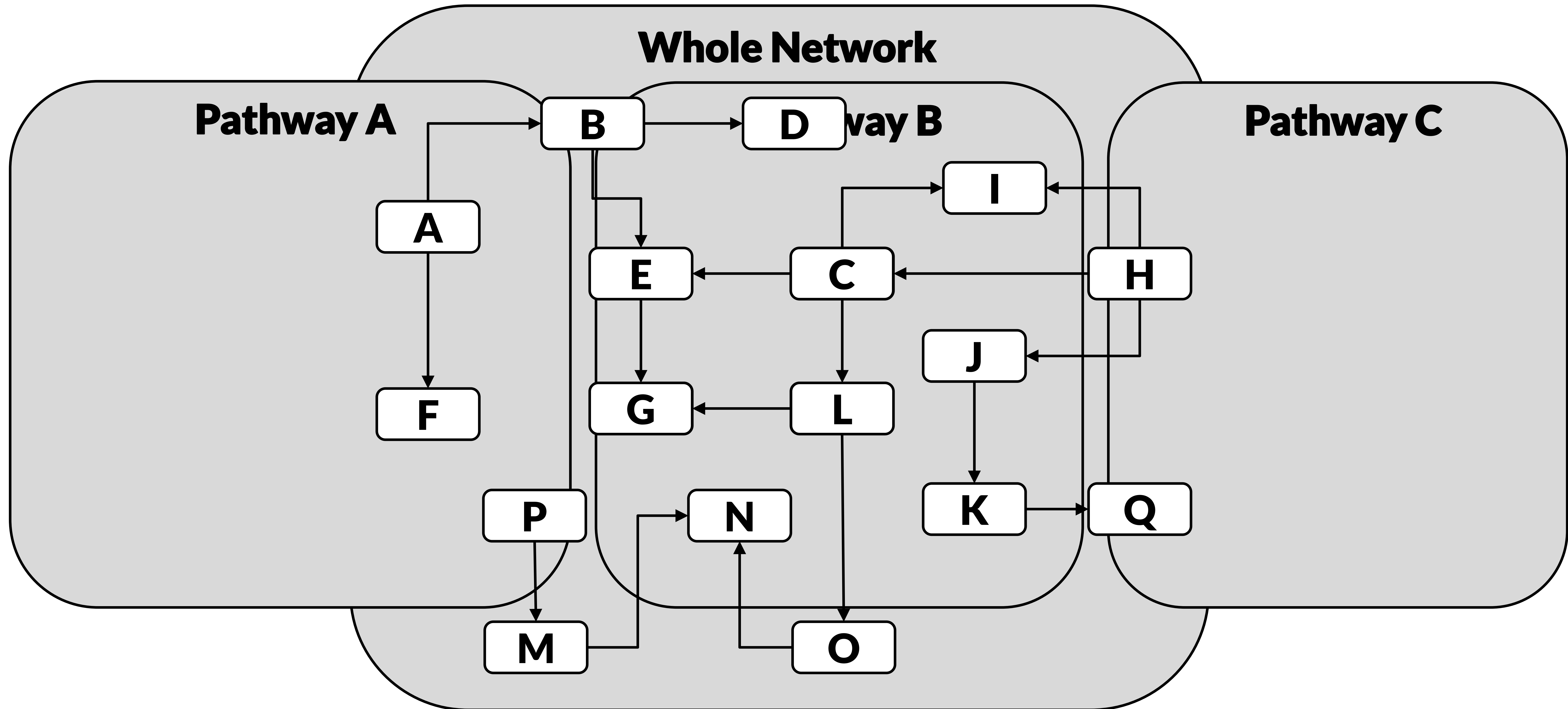
Glial progenitor cell



The bigger picture

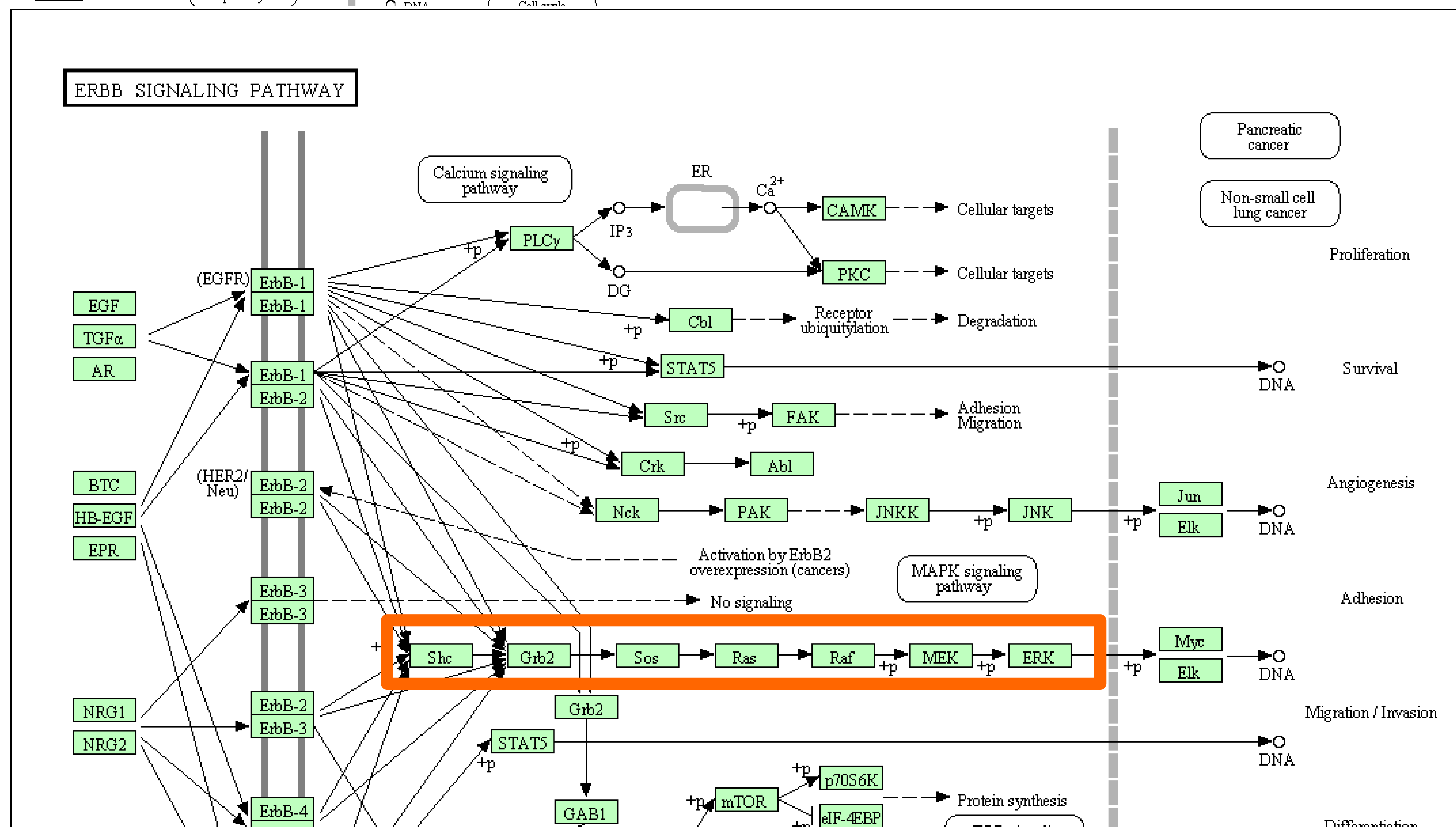
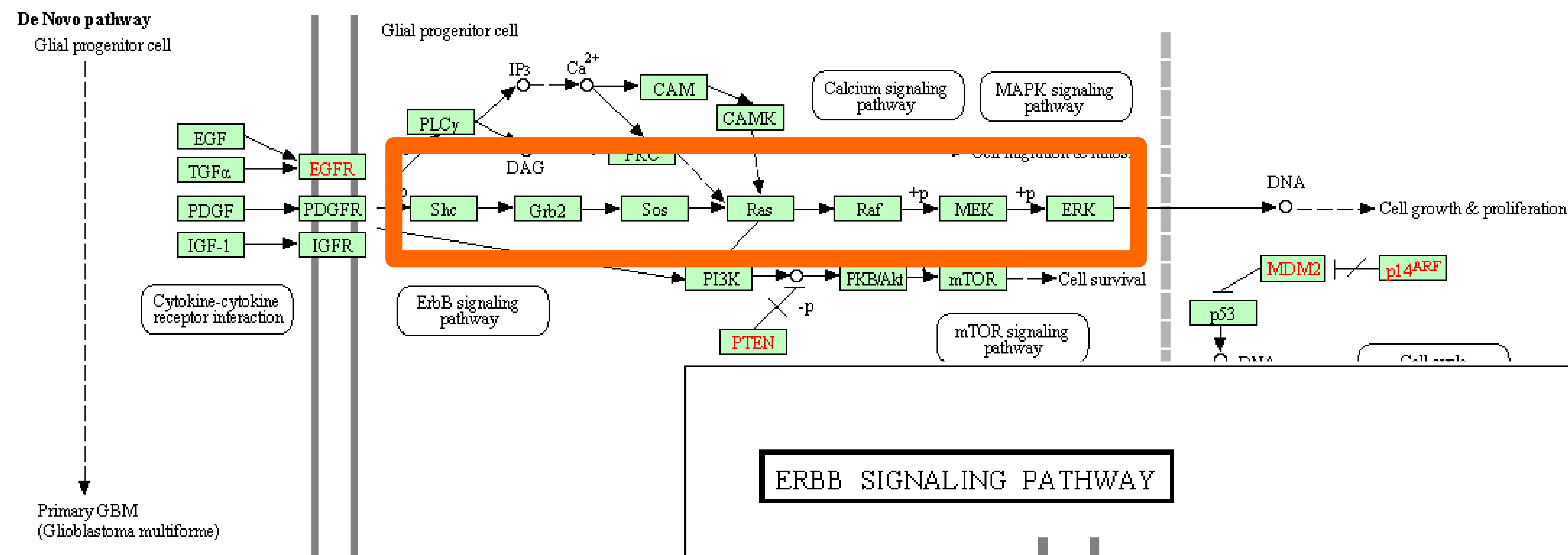


Background

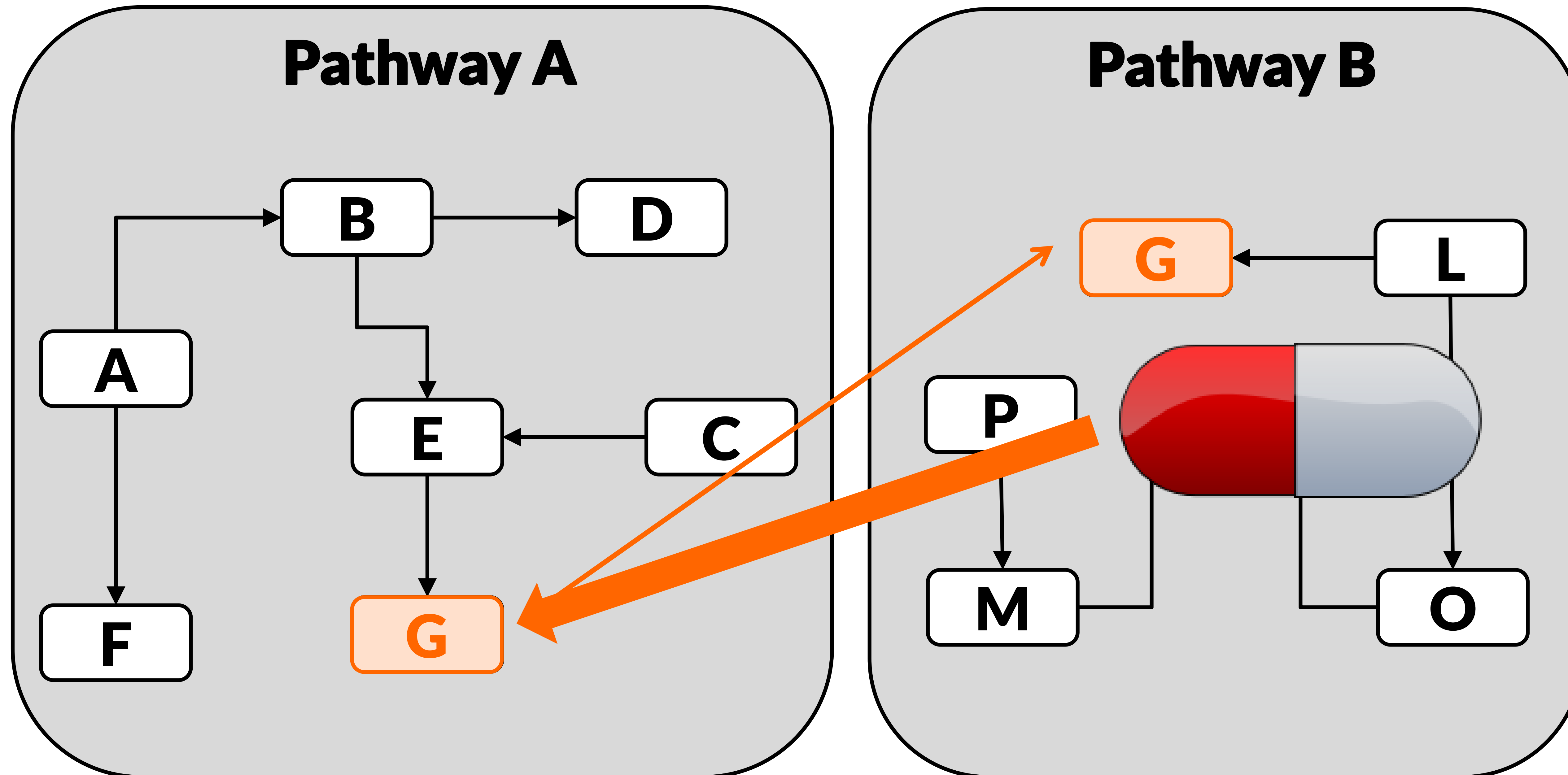


GLIOMA

A Pathway

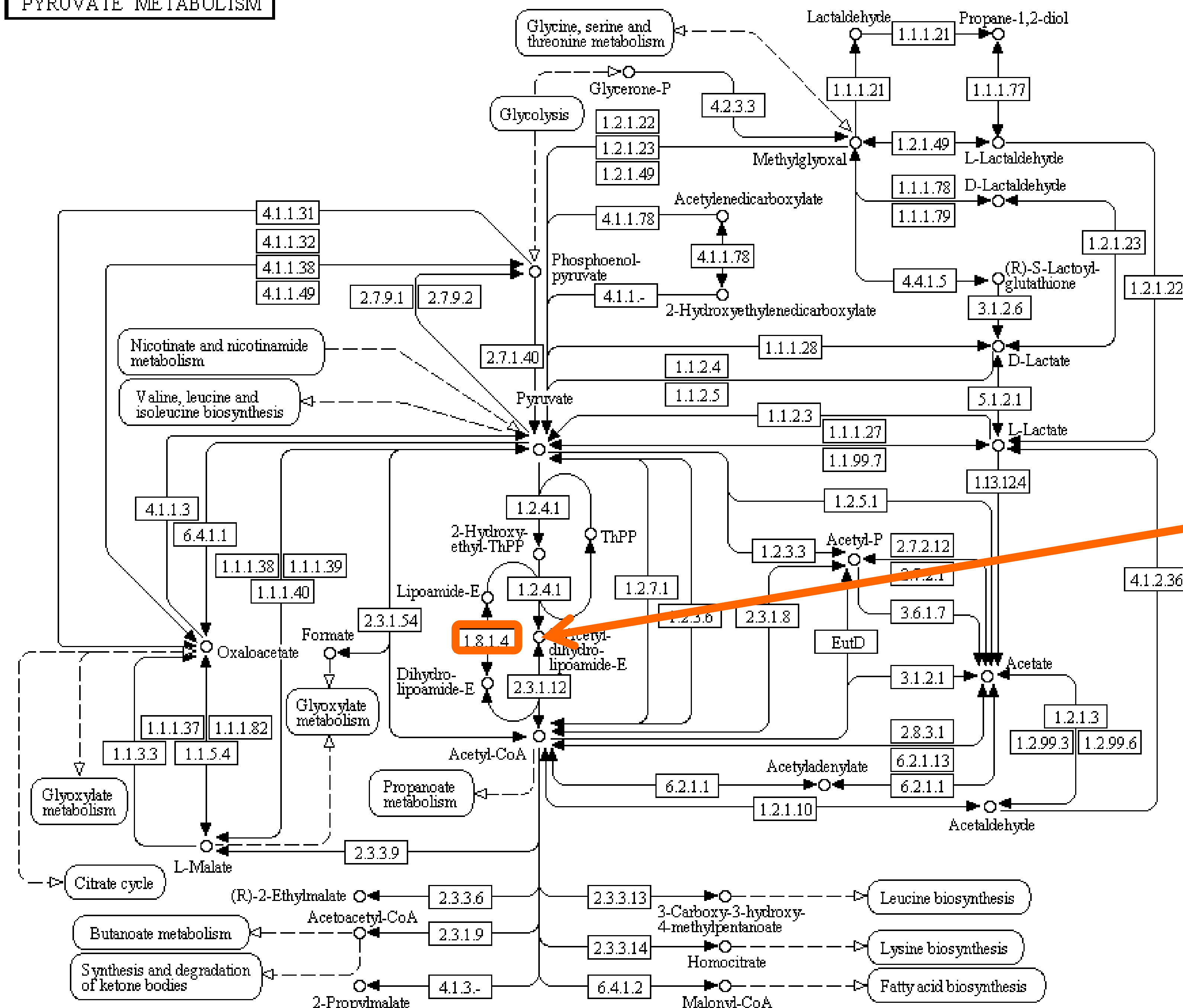


Challenges



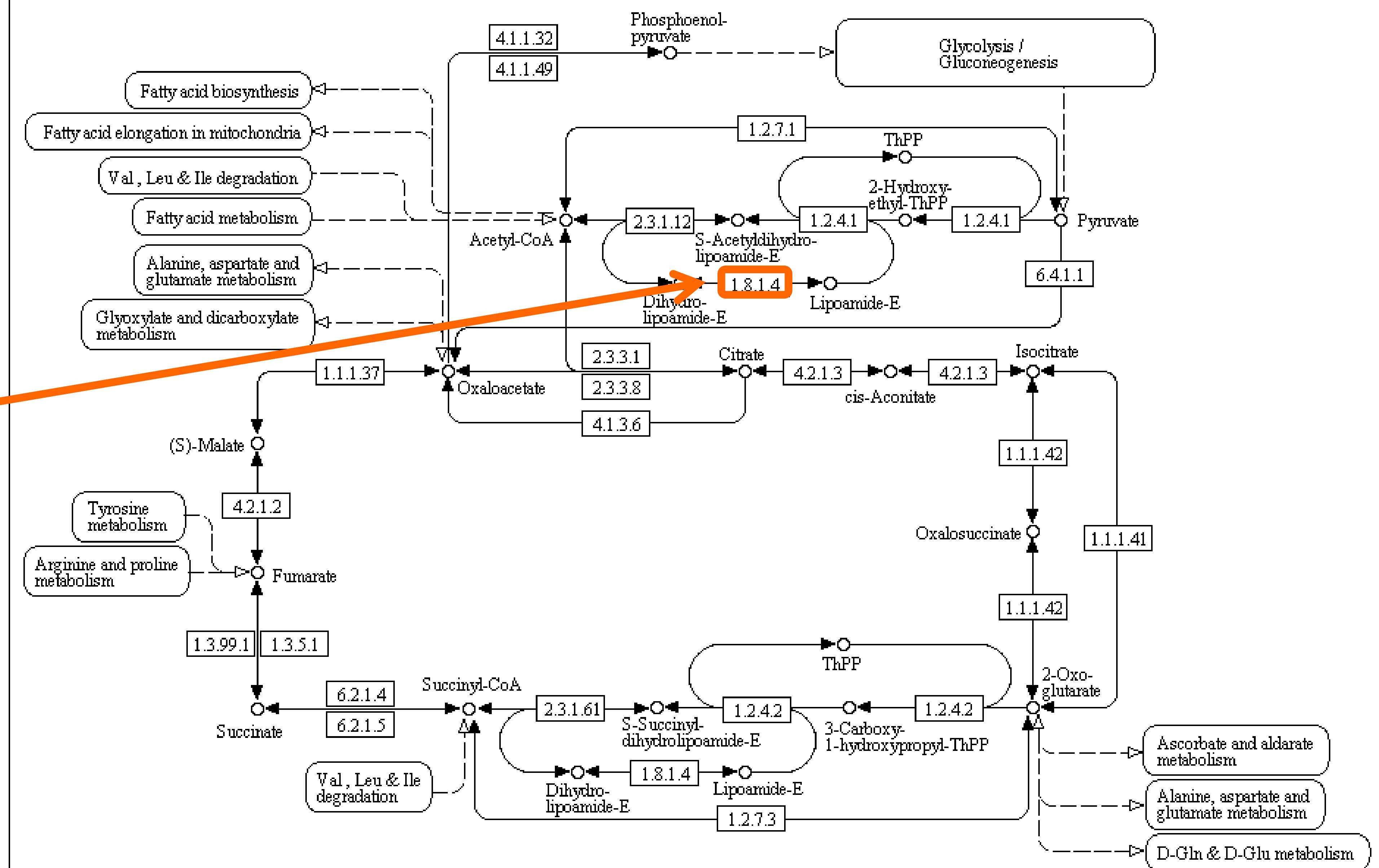
Challenges

PYRUVATE METABOLISM



00620 6/6/12
(c) Kanehisa Laboratories

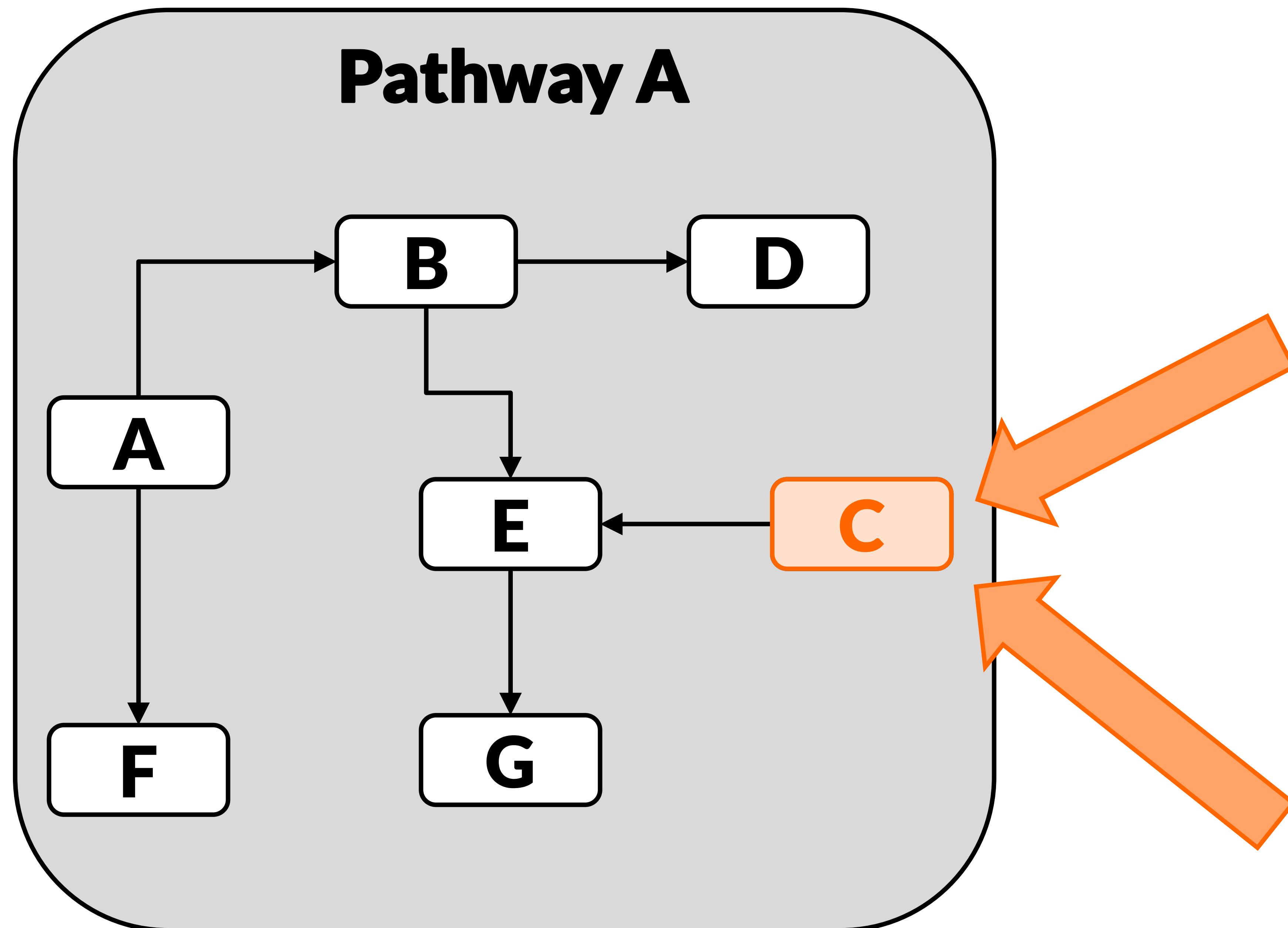
CITRATE CYCLE (TCA CYCLE)



00020 5/31/12
(c) Kanehisa Laboratories

How to visualize pathway relationships?

Challenges



Node	Sample 1	Sample 2	Sample 3	...
A	0.55	0.95	0.83	...
B	0.12	0.42	0.16	...
C	0.33	0.65	0.38	...
...

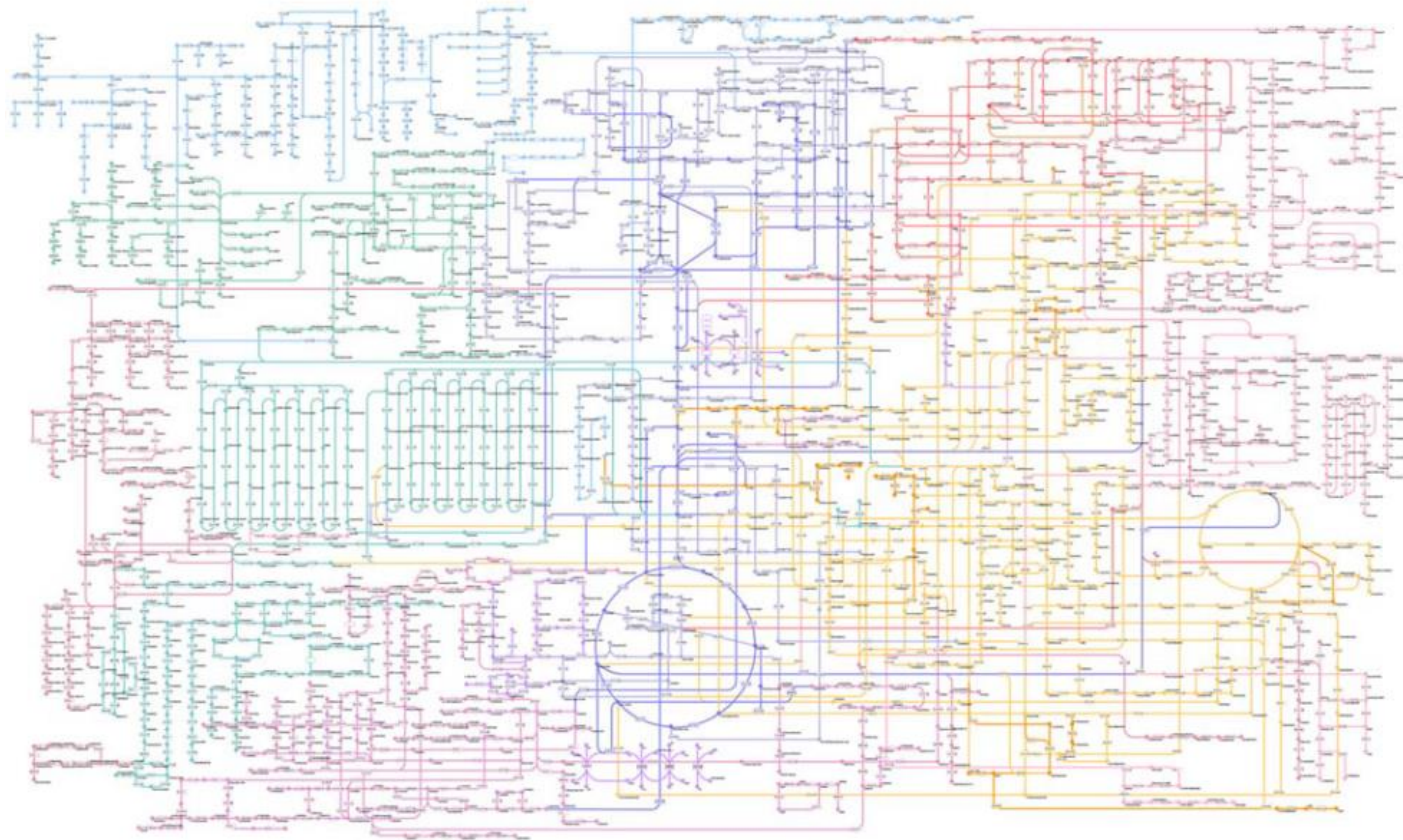
Node	Sample 1	Sample 2	Sample 3	...
A	low	low	very high	...
B	normal	low	high	...
C	high	very low	normal	...
...

How to visualize experimental data on pathways?

**How to visualize
pathway relationships?**

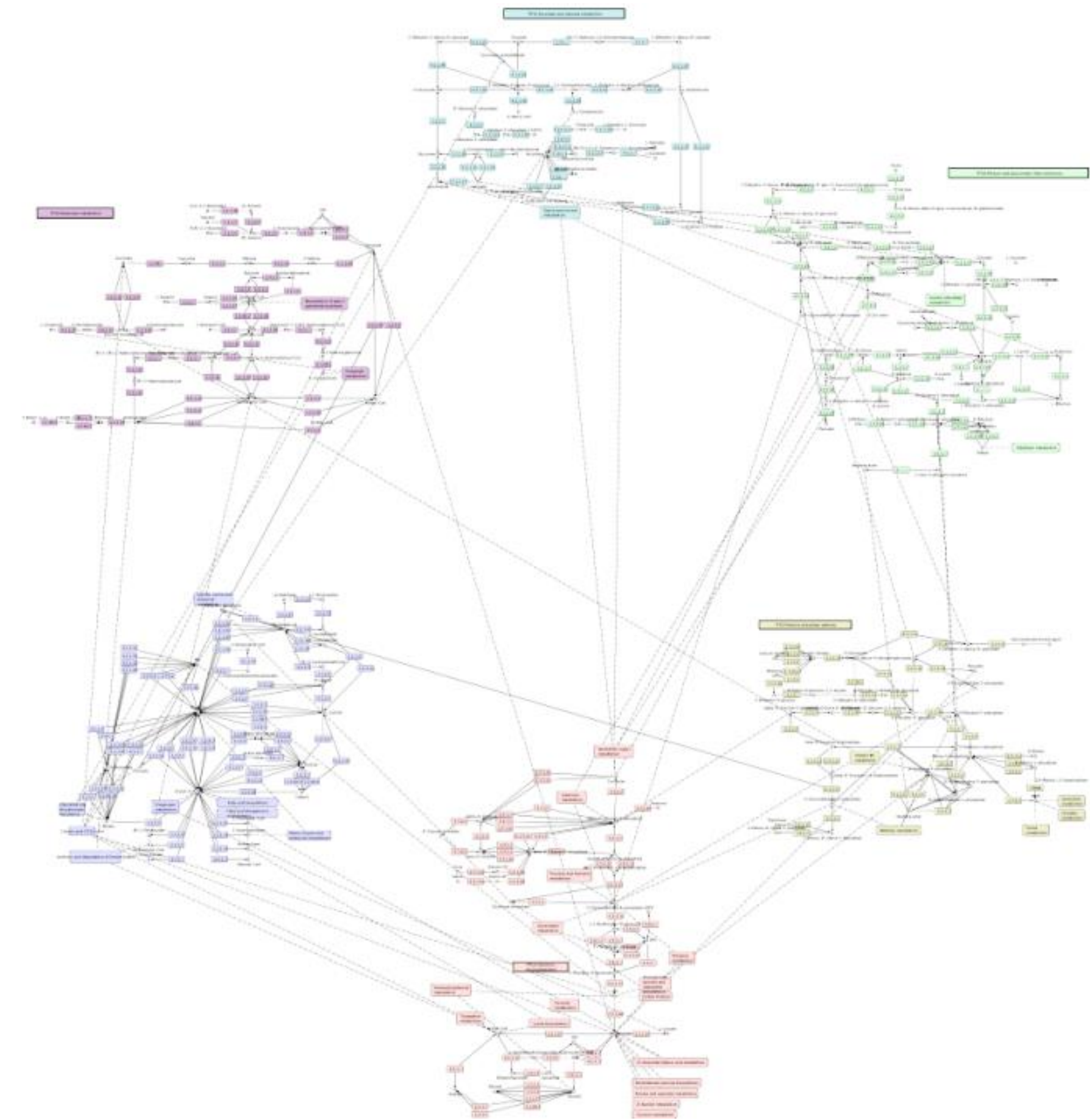
Approaches

Whole Network



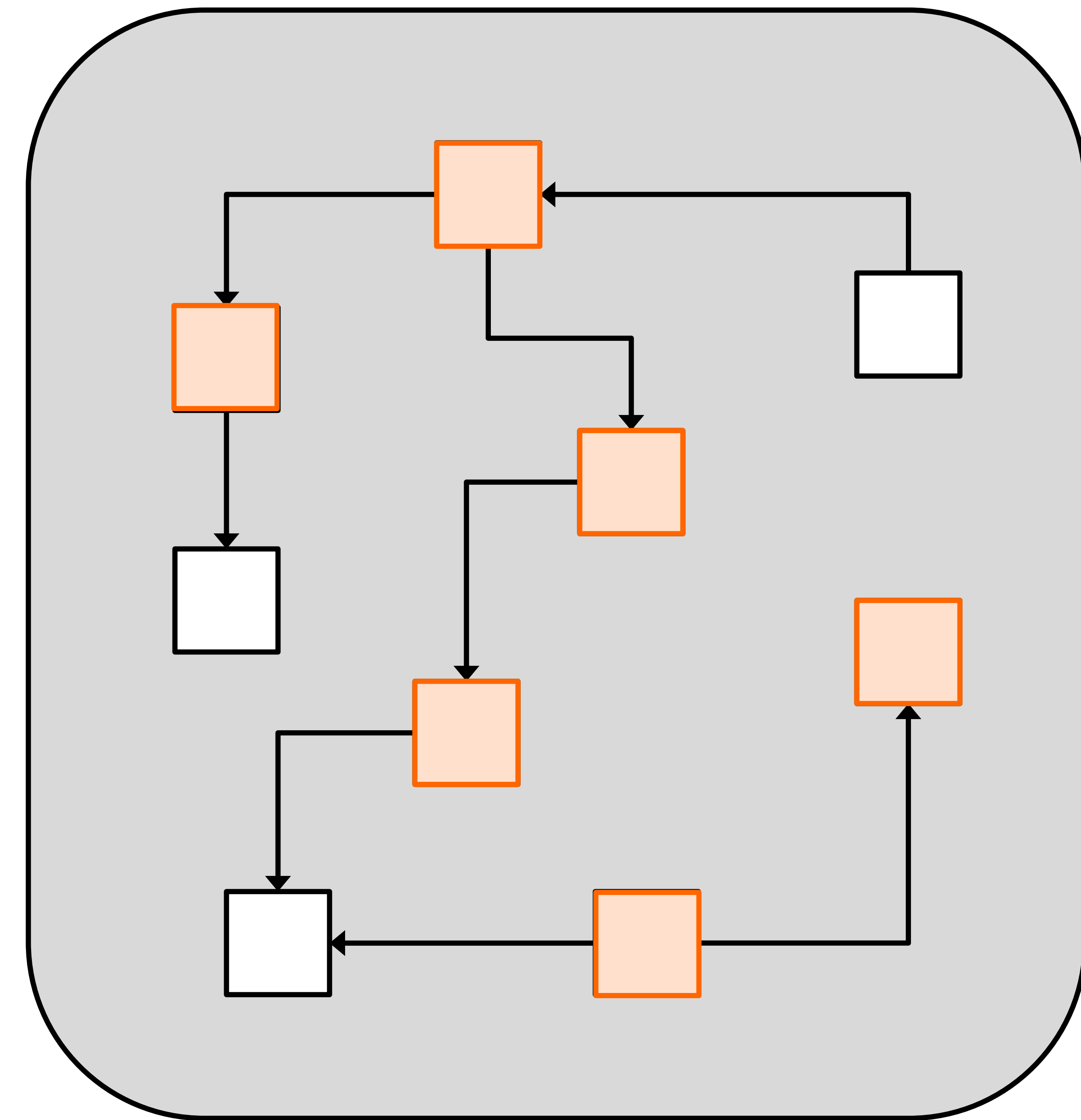
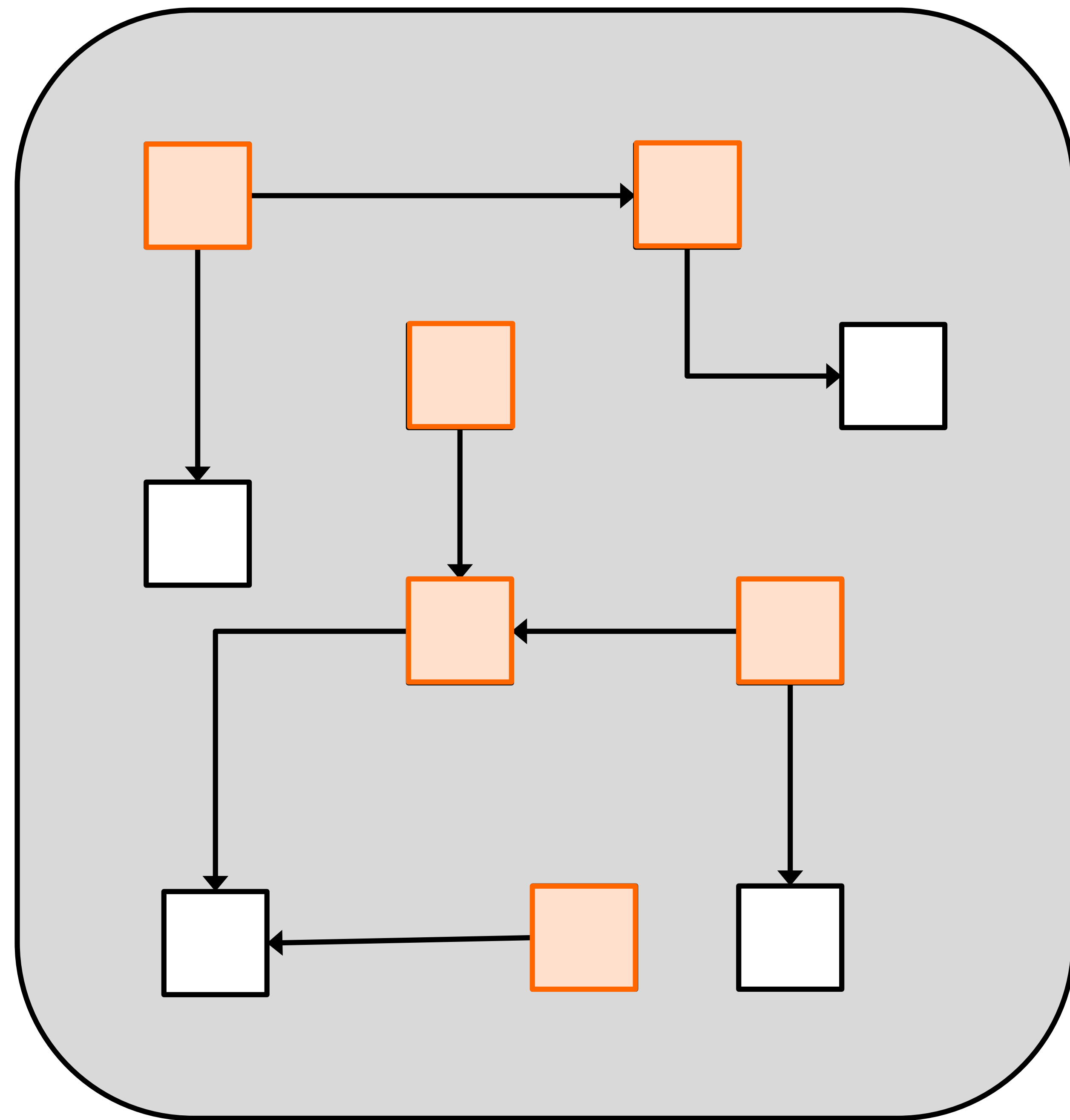
[Kono2009]

Connected Pathways

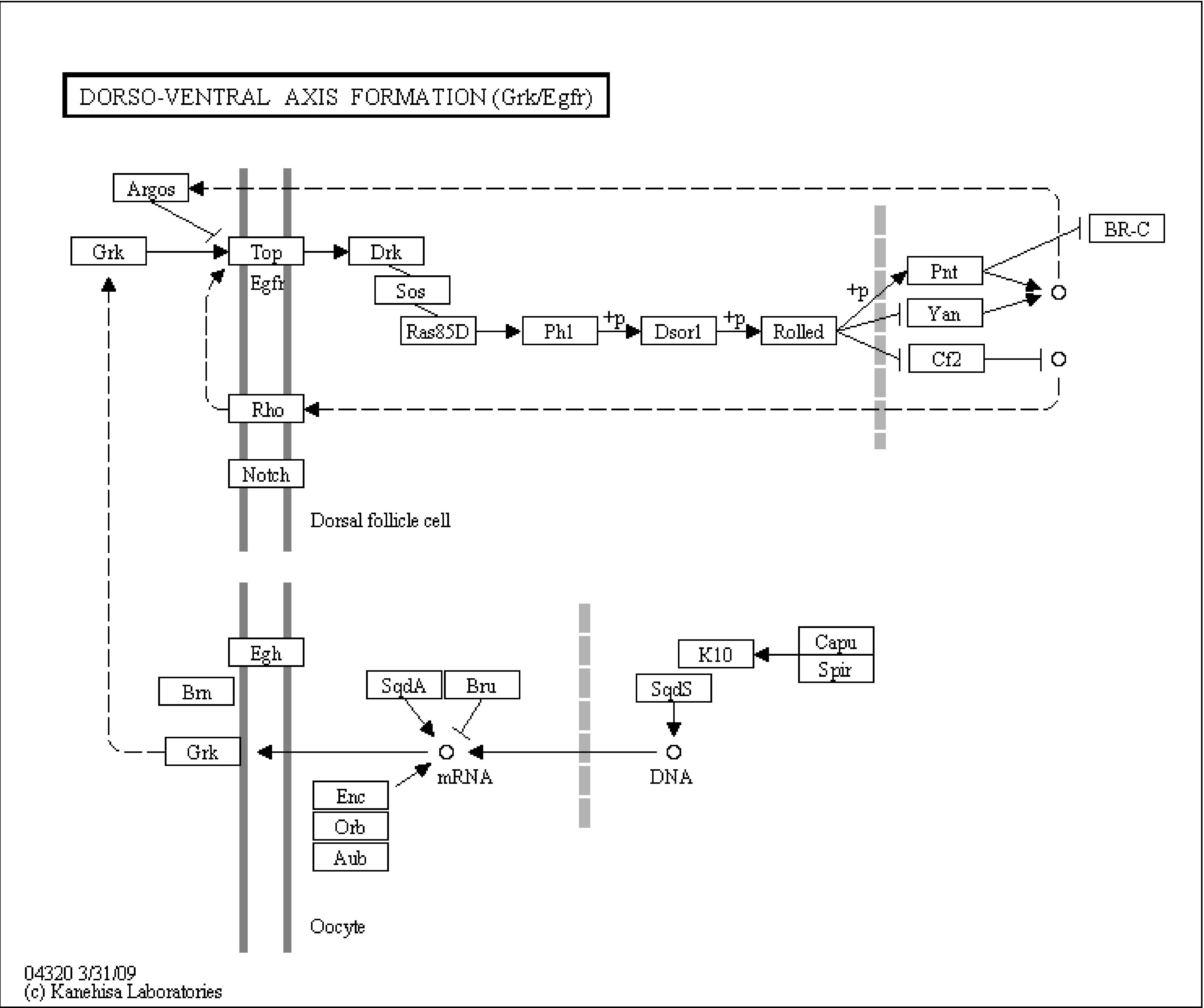


[Klukas2006]

Finding Related Pathways

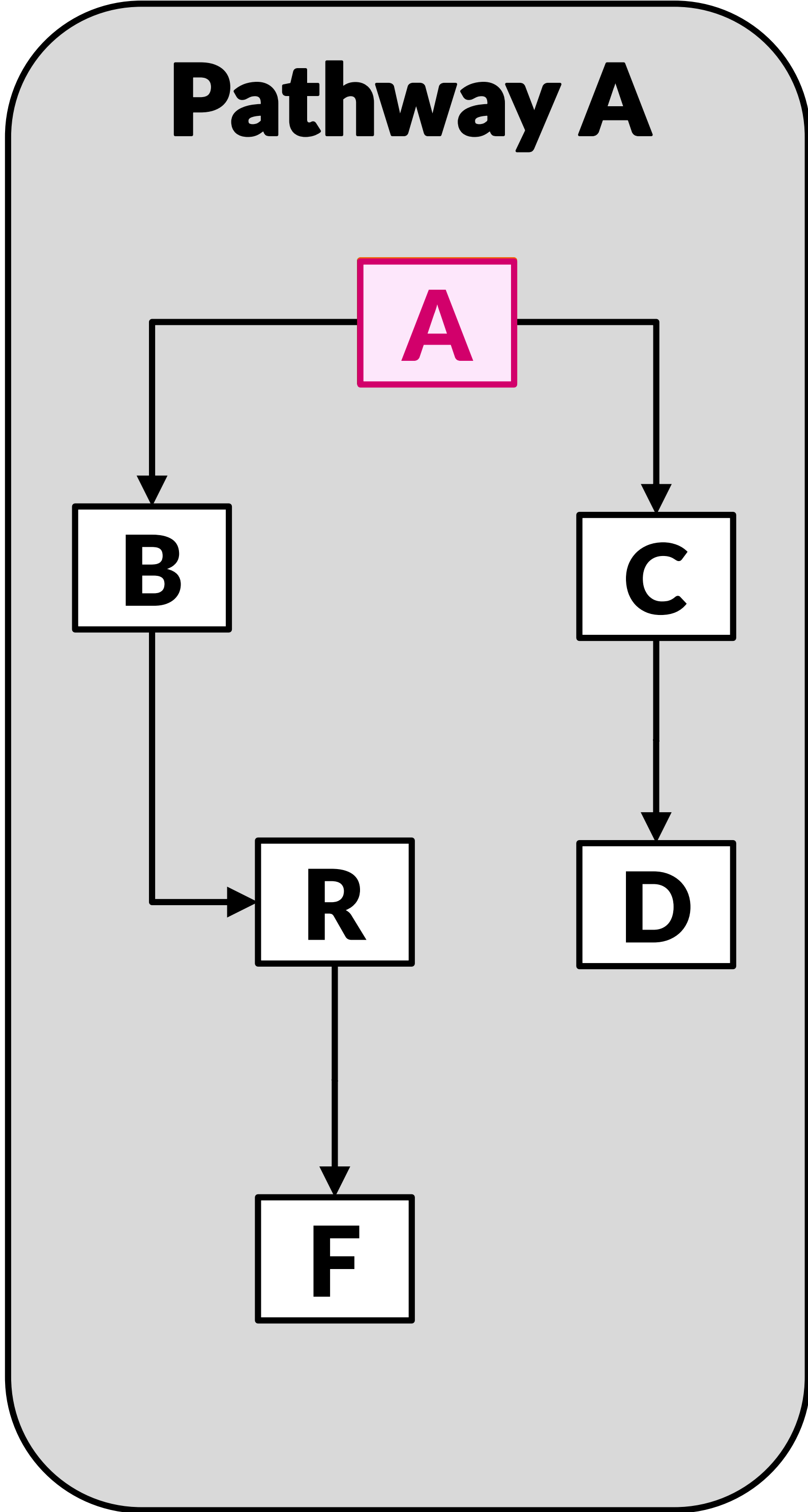


Finding Related Pathways

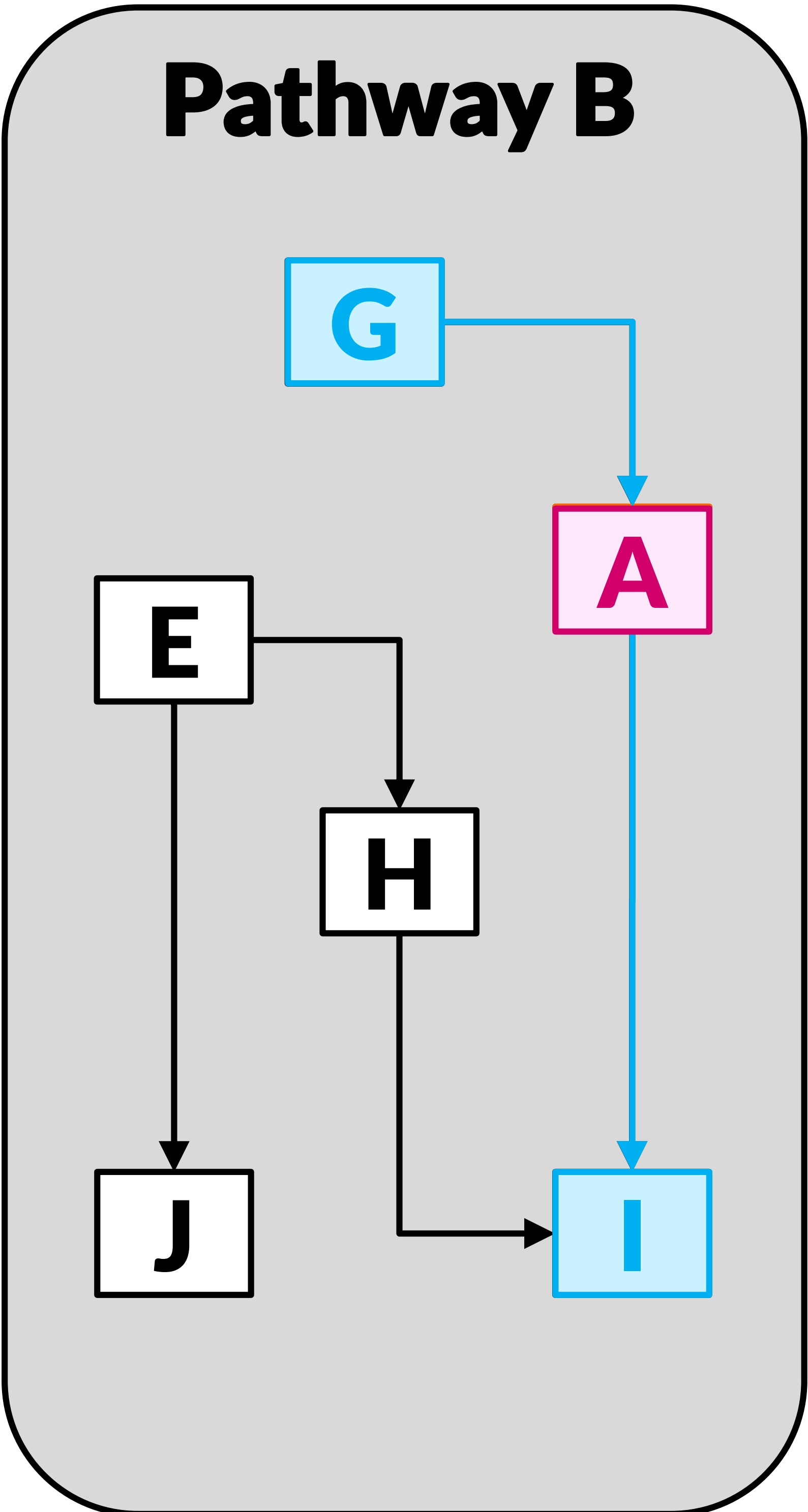


Pathways	
Pathway	Commo
Nifedipine Activity	
Endometrial cancer	
Thyroid cancer	
Melanoma	
Acute myeloid leukem	
Signaling of Hepatocy	
Non-small cell lung c	
Renal cell carcinoma	
Bladder cancer	
Gap junction	
GnRH signaling path	
Glioma	
IL-9 signaling pathwa	
Progesterone-mediate	
Estrogen signaling pa	
EPO Receptor Signali	
Prostate cancer	
IL-5 signaling pathwa	
Chronic myeloid leuk	
ErbB signaling pathw	
Fc epsilon RI signali	
Osteopontin Signaling	
Serotonin Receptor 2	
ErbB signaling pathw	
IL-3 Signaling Pathwa	
B cell receptor signal	
Prion diseases	
Prolactin signaling pa	
Focal adhesion	

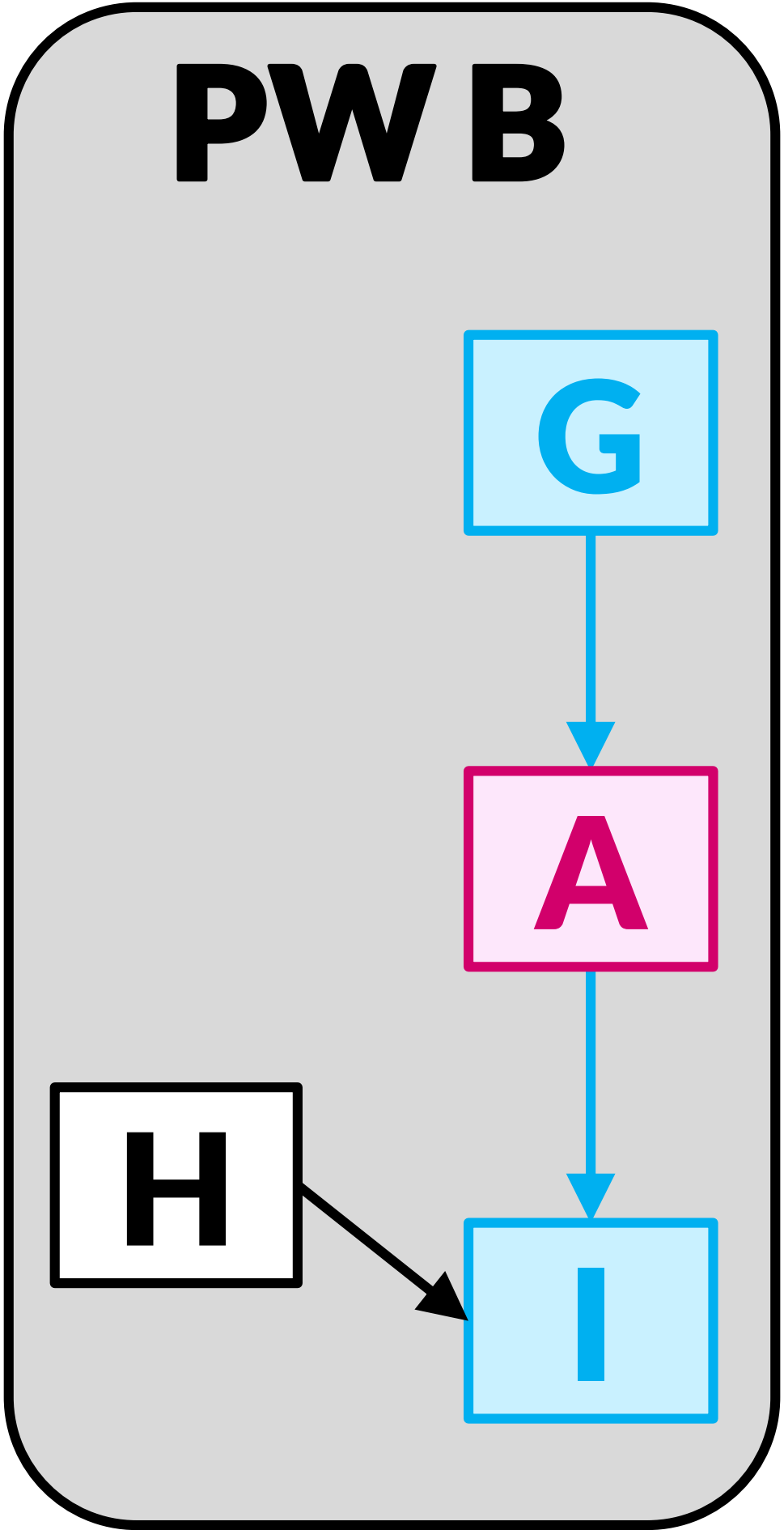
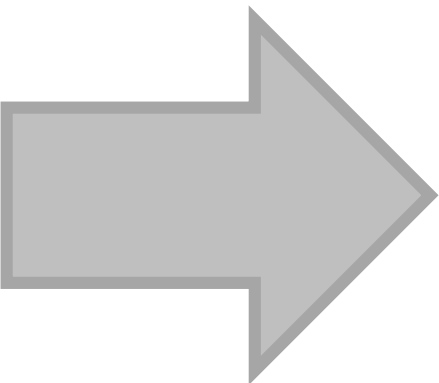
Contextual Subsets



Focus Pathway

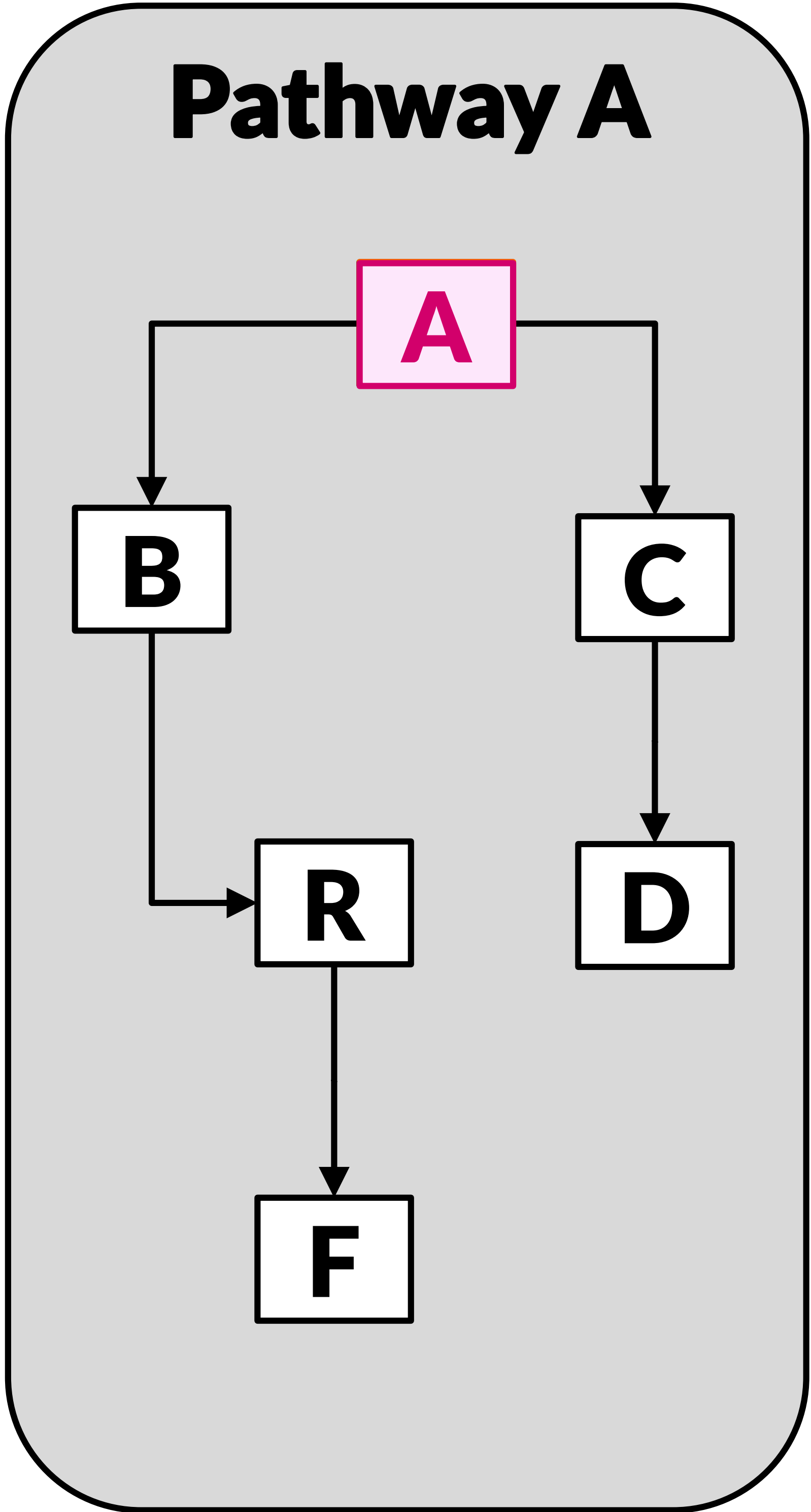


Context Pathway

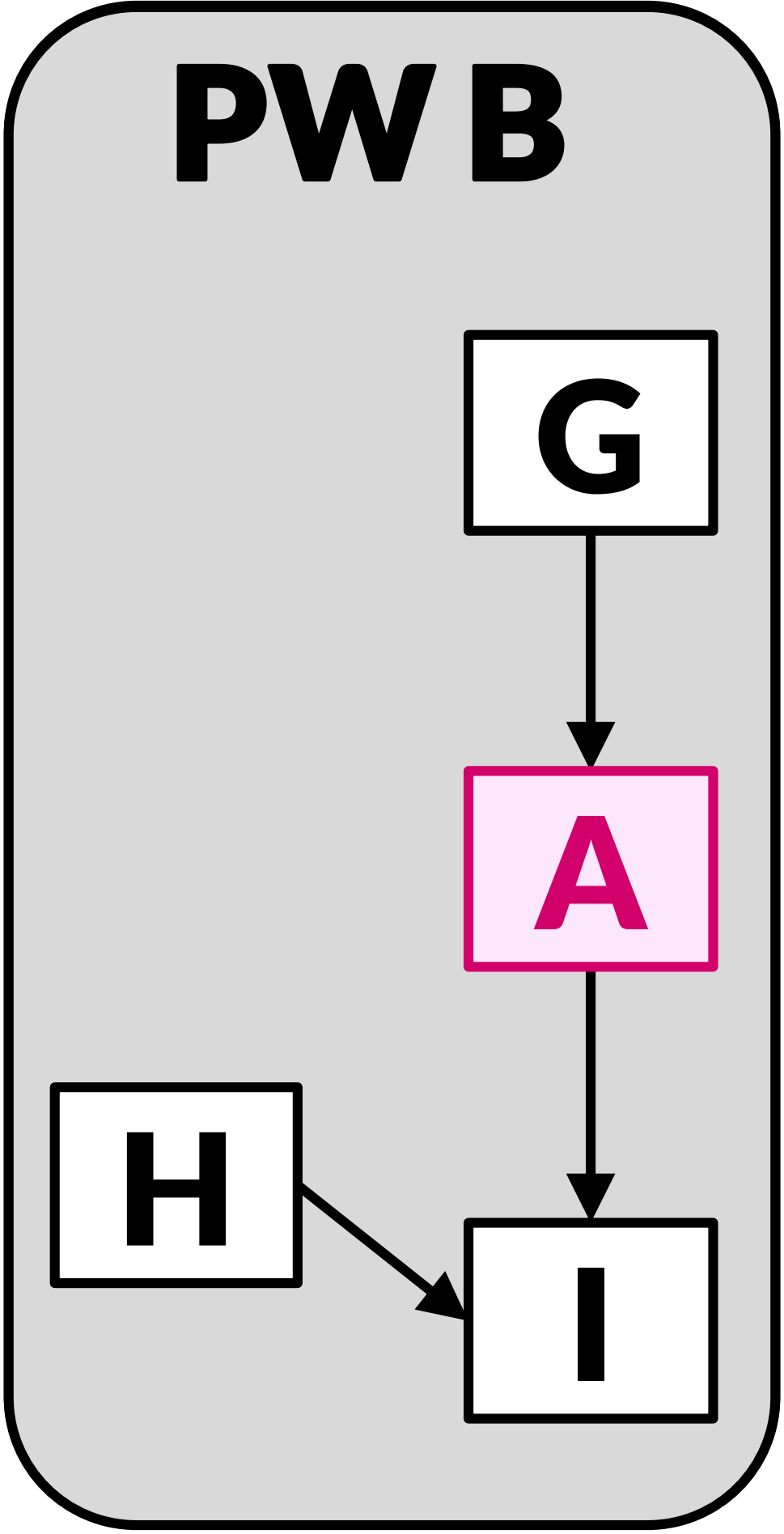
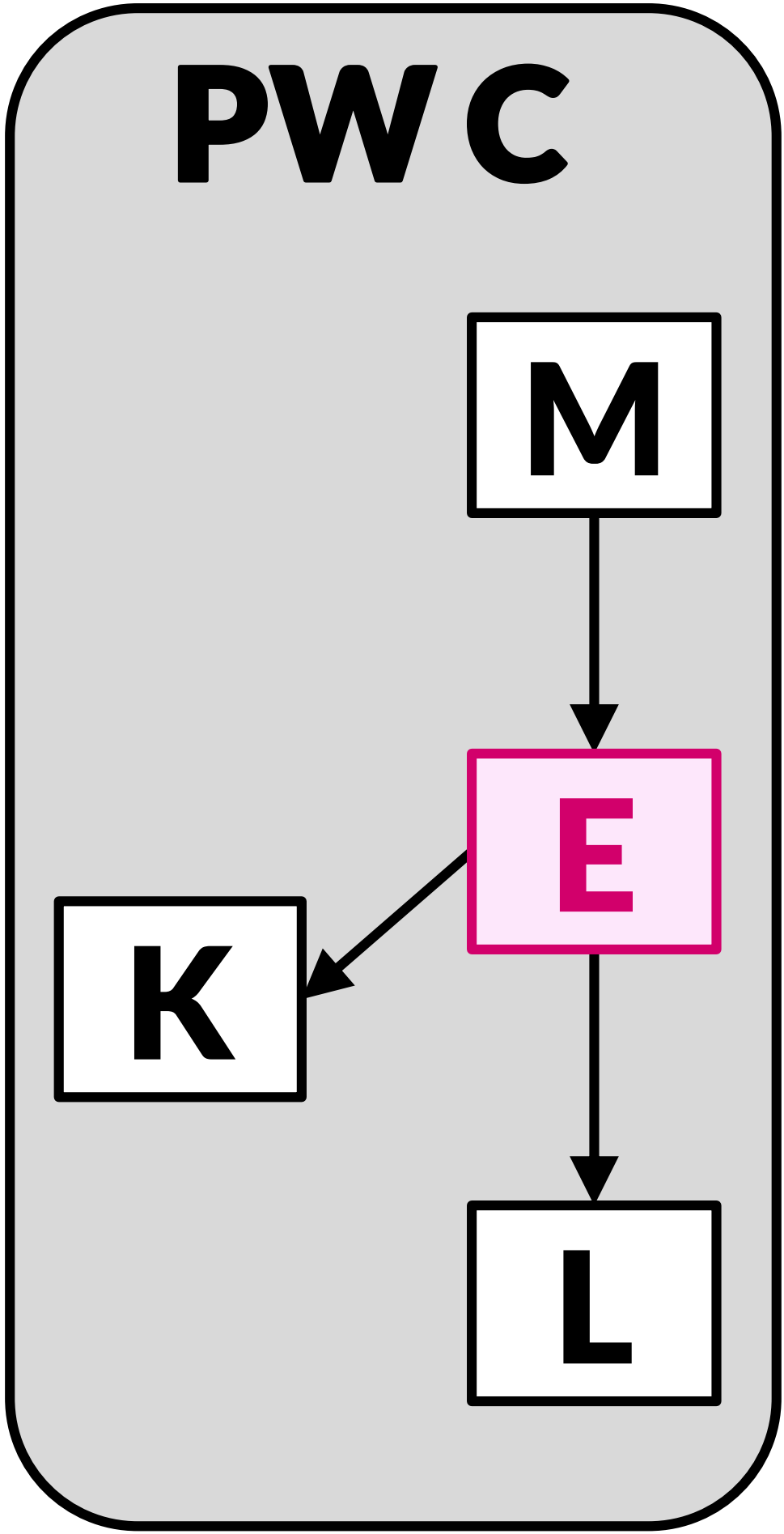


Contextual Subset

Contextual Subsets



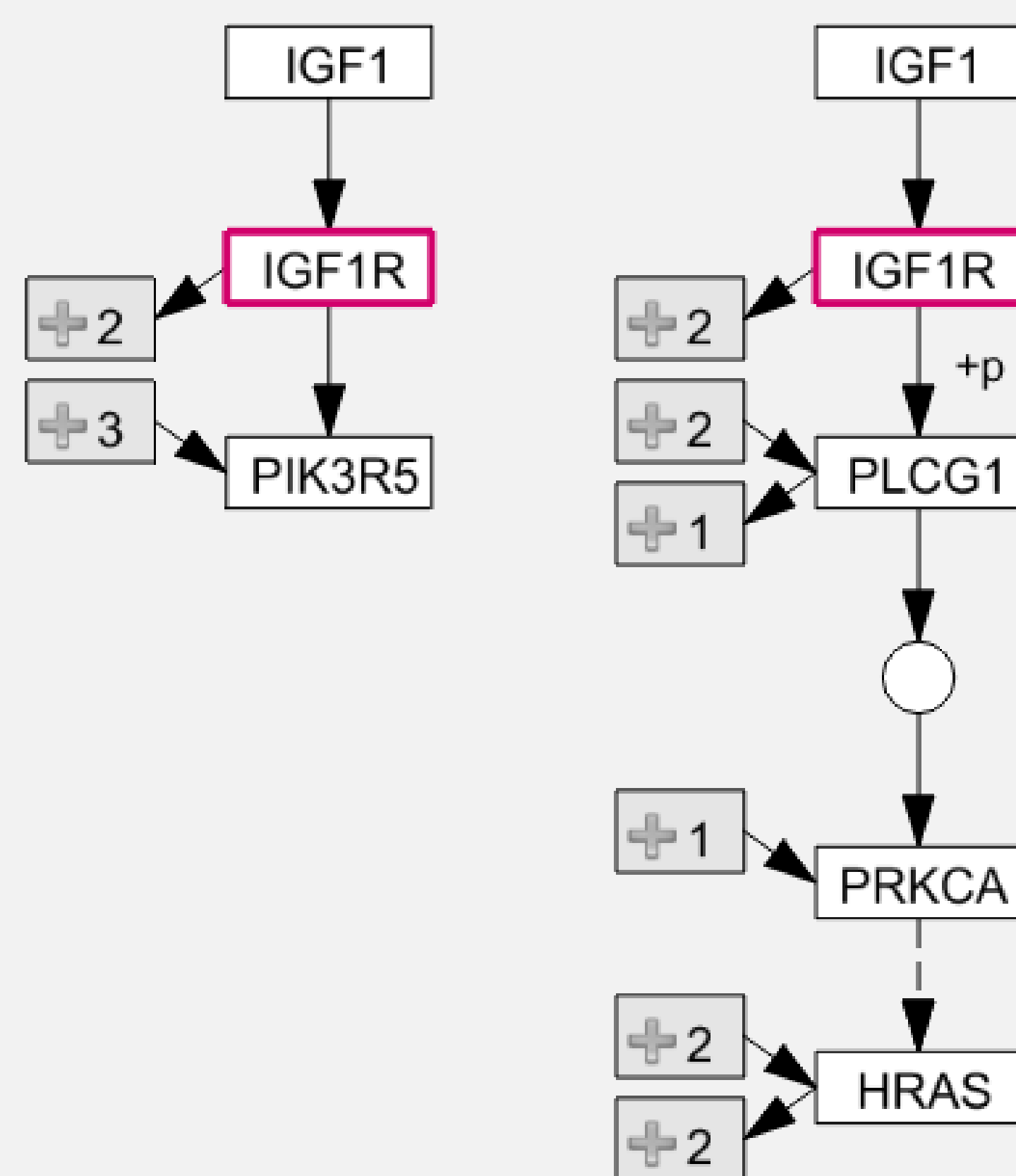
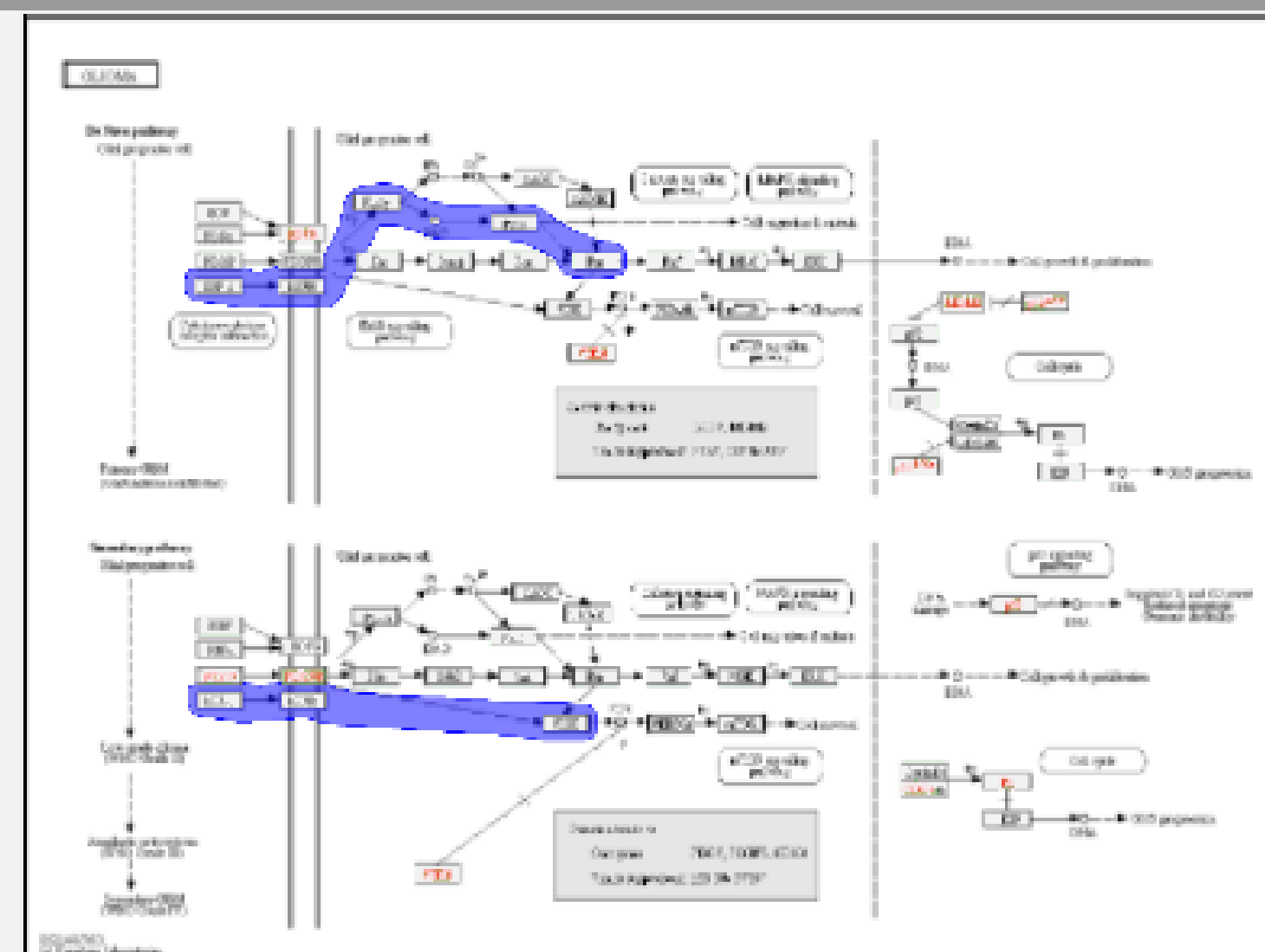
Focus Pathway



Context Pathways

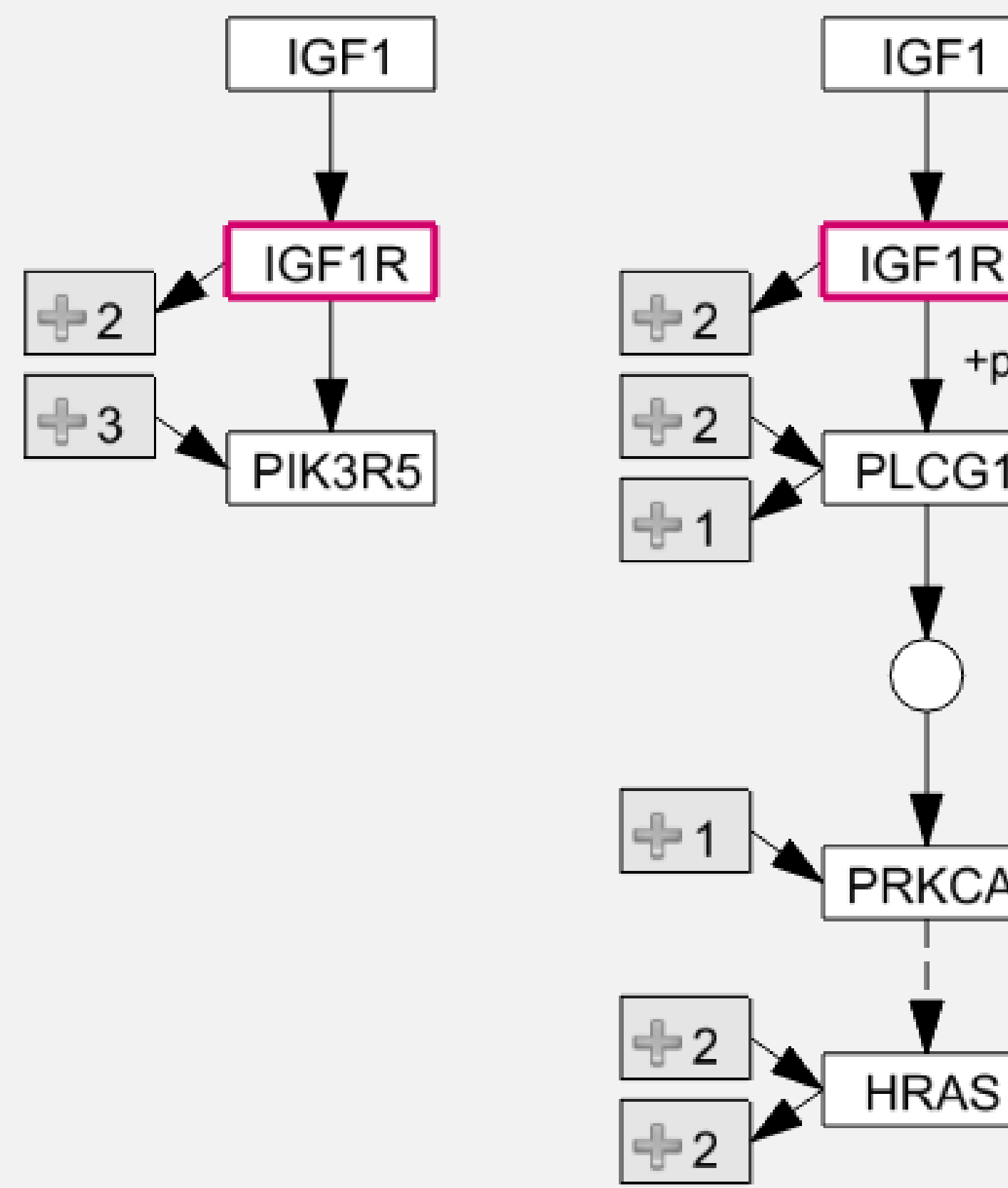
Levels of Detail

Glioma



High

Glioma



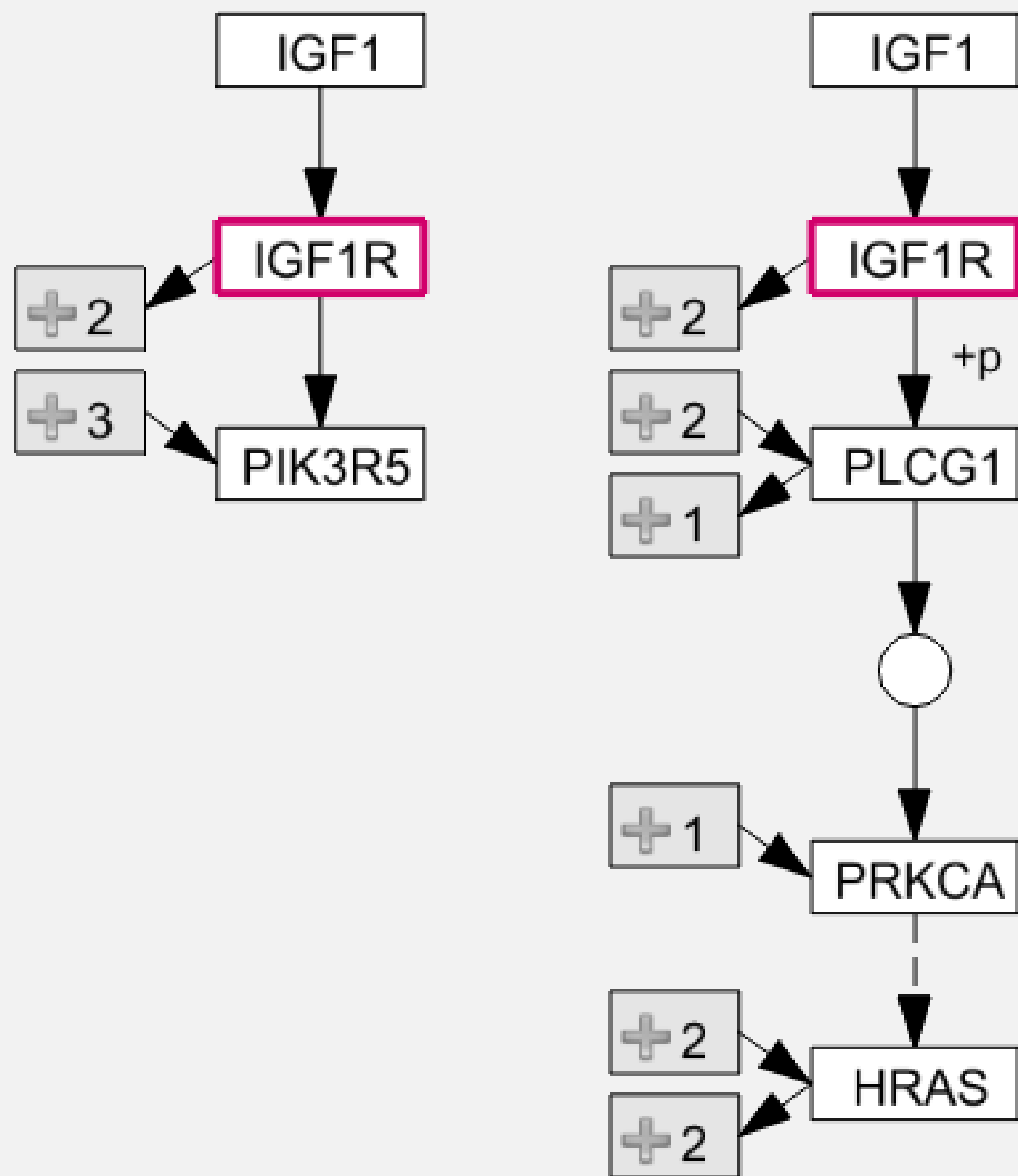
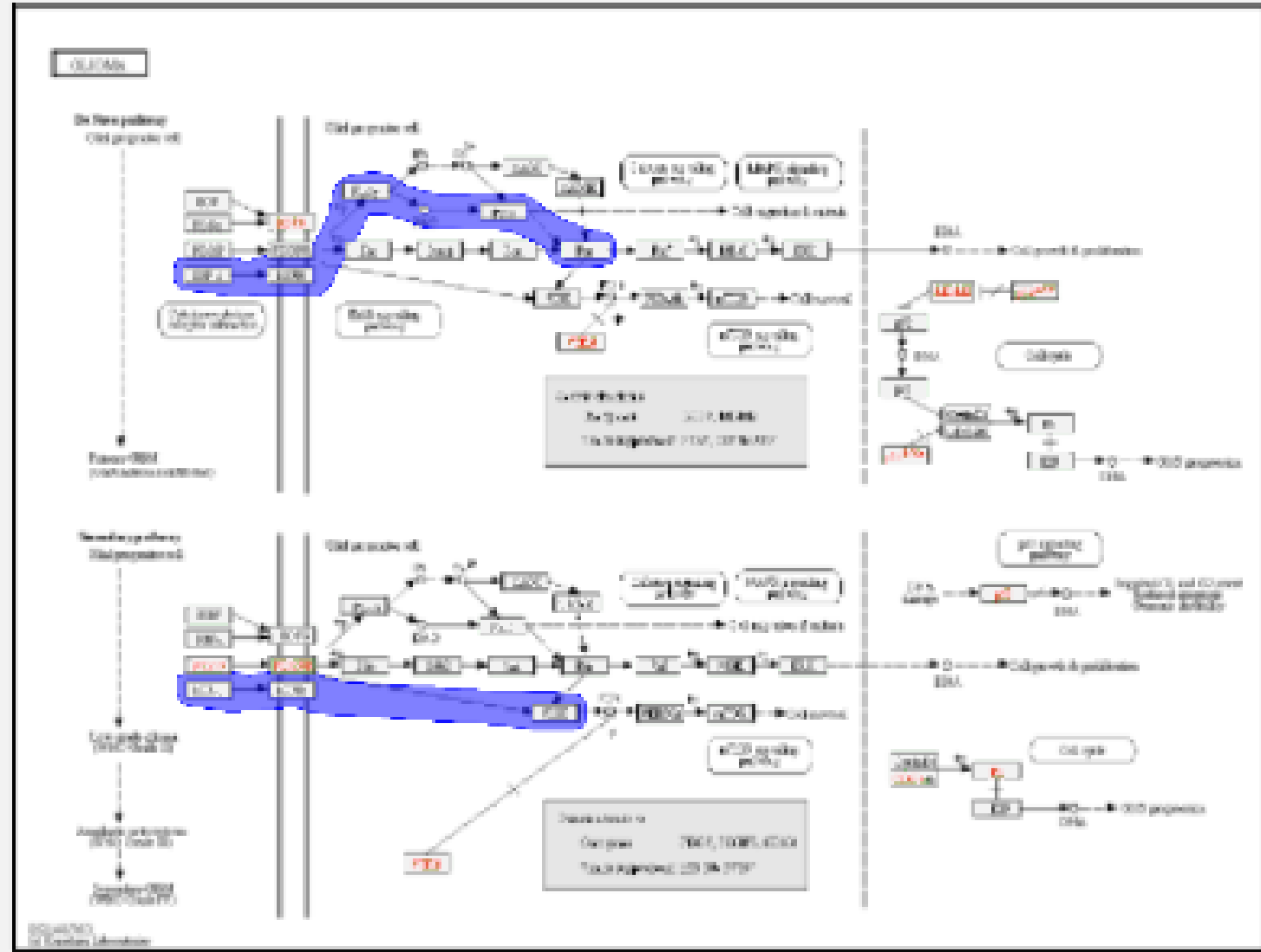
Medium

Glioma

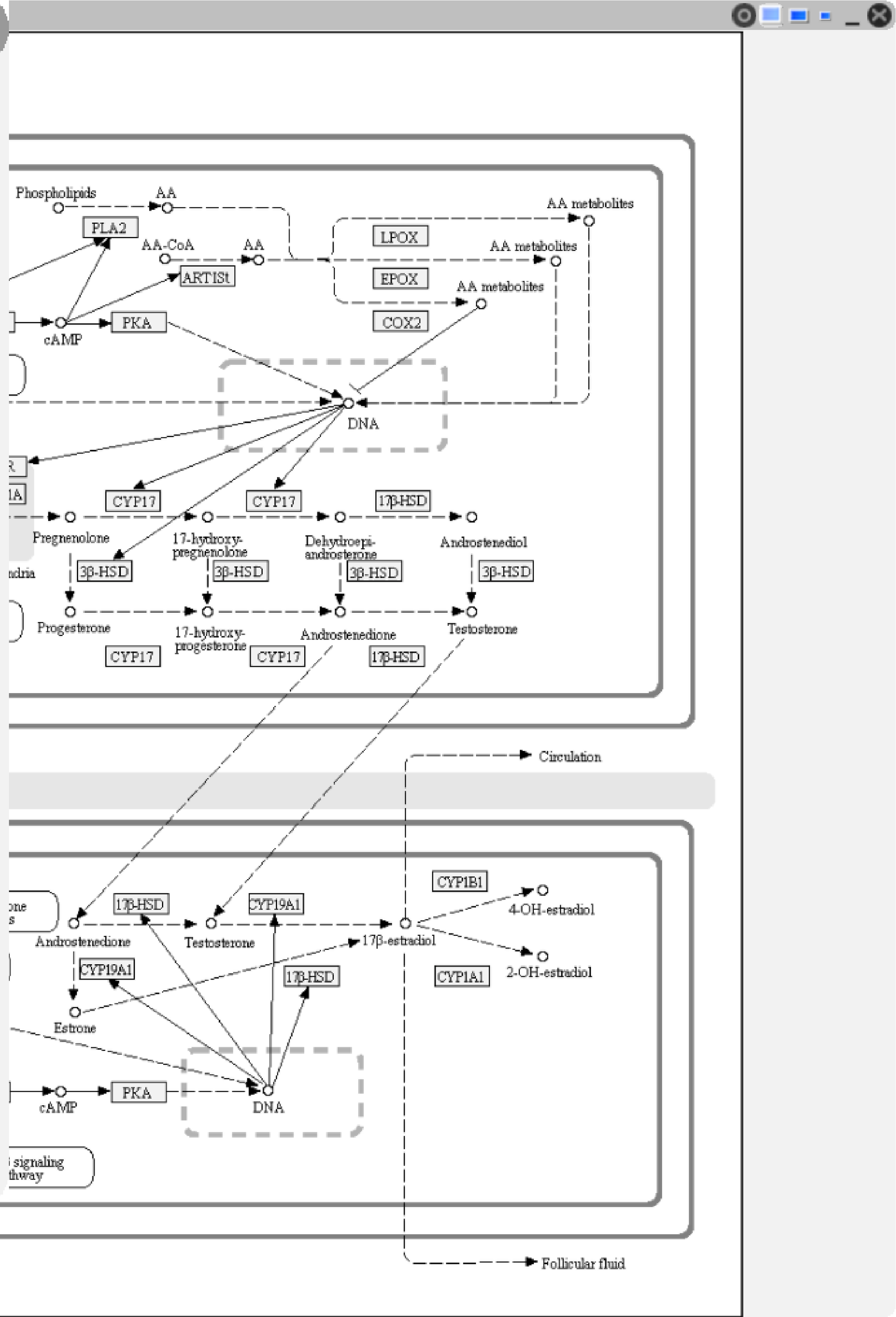
Low

Layout

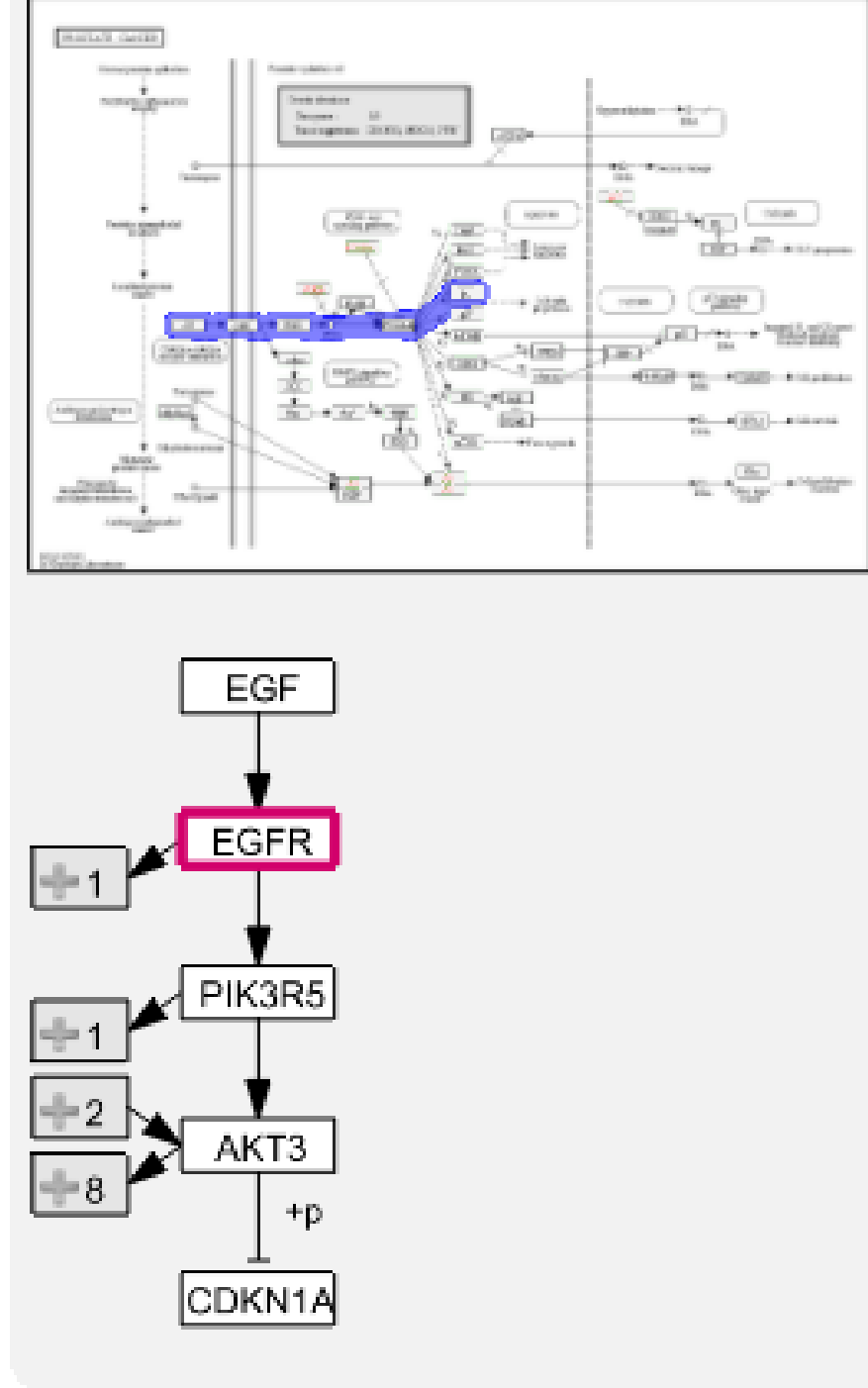
Glioma



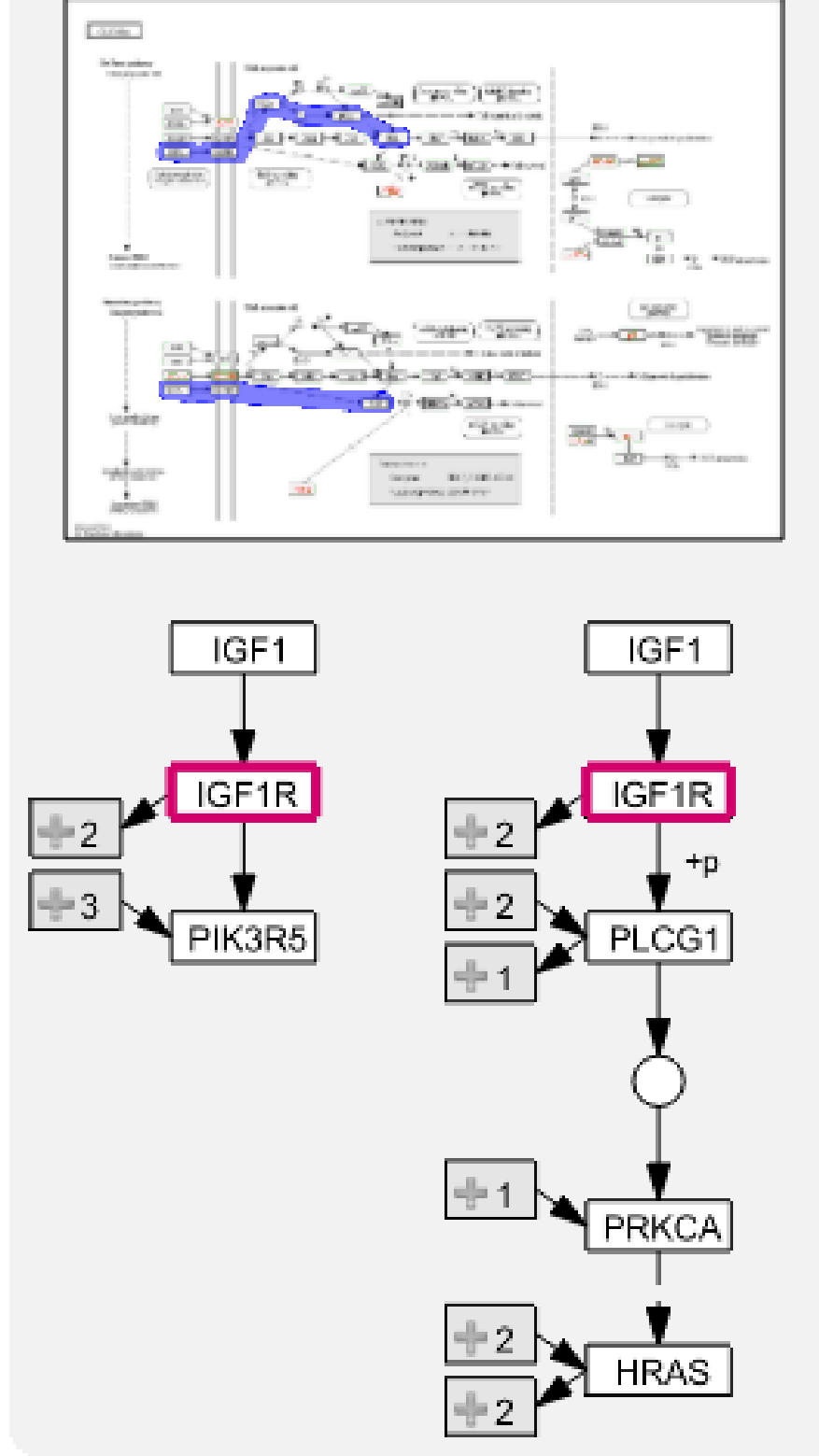
04913 4/22/13
(c) Kanehisa Laboratories



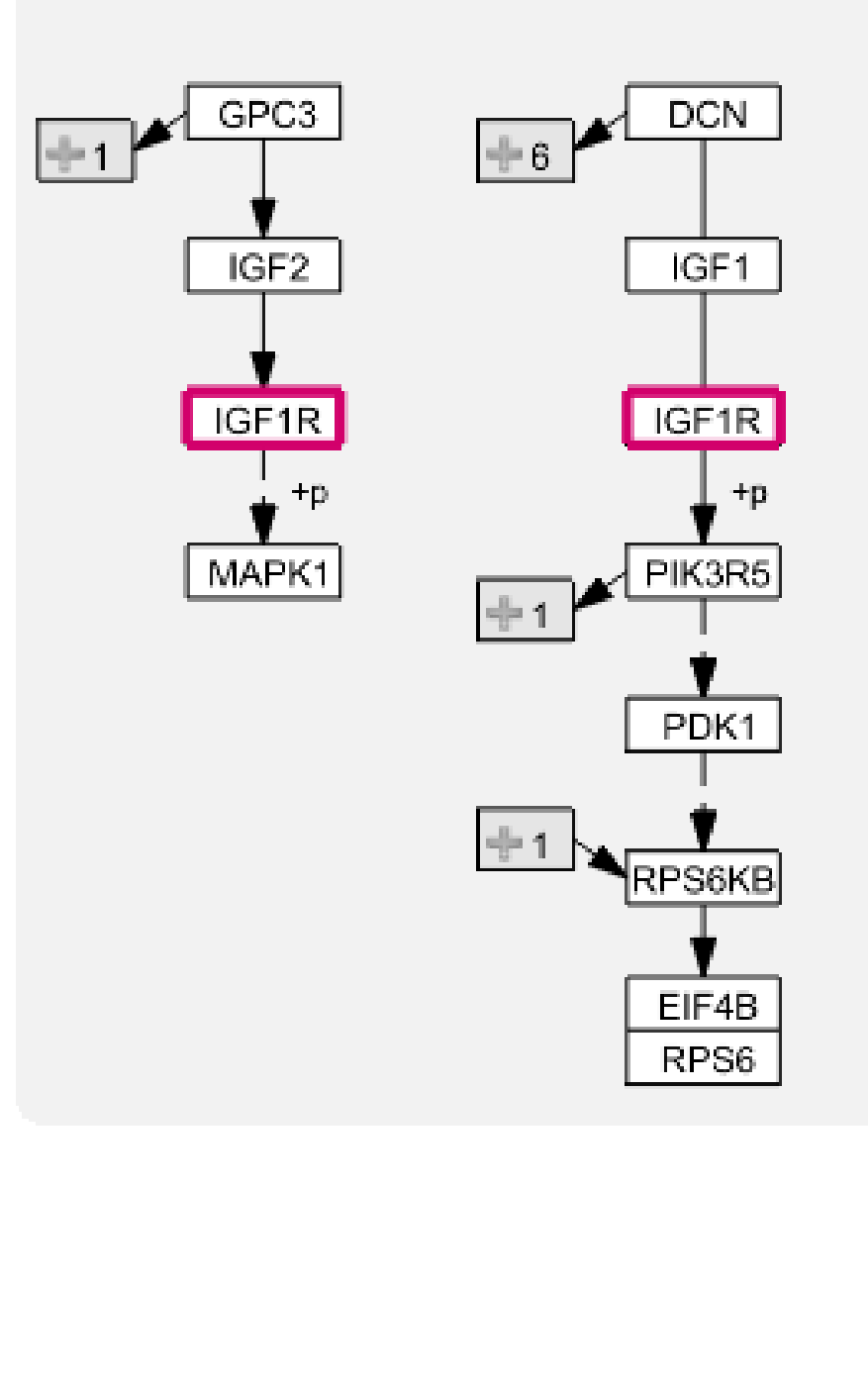
Prostate cancer



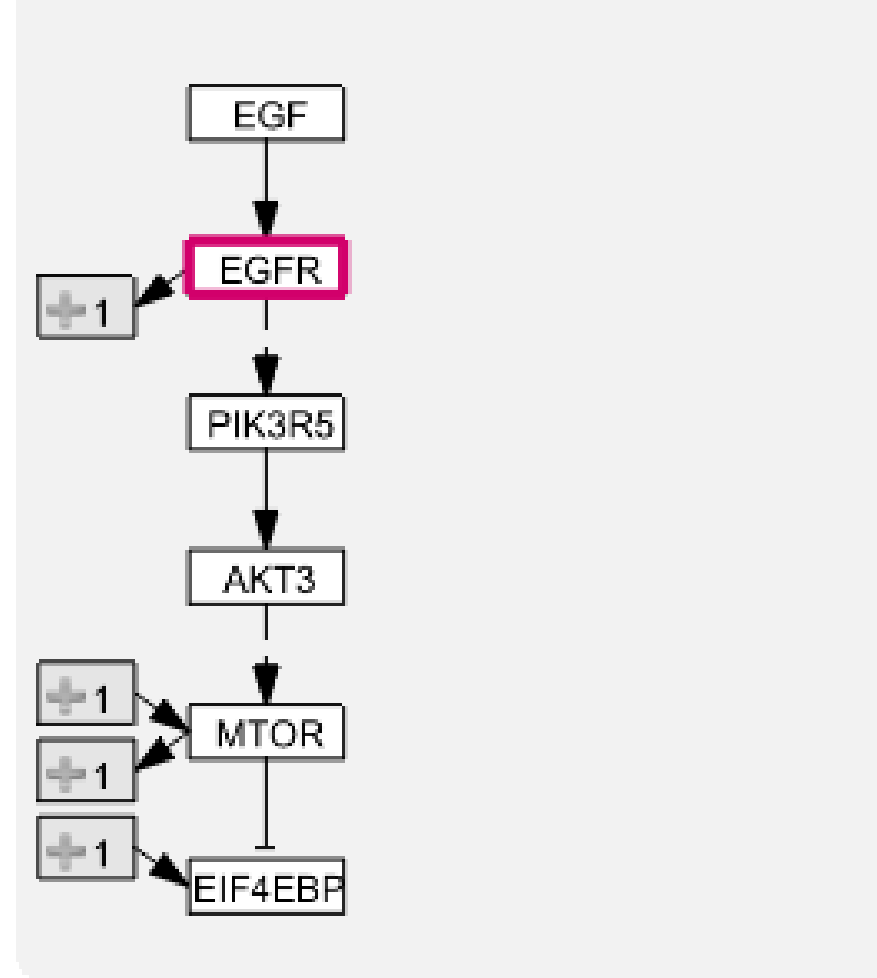
Glioma



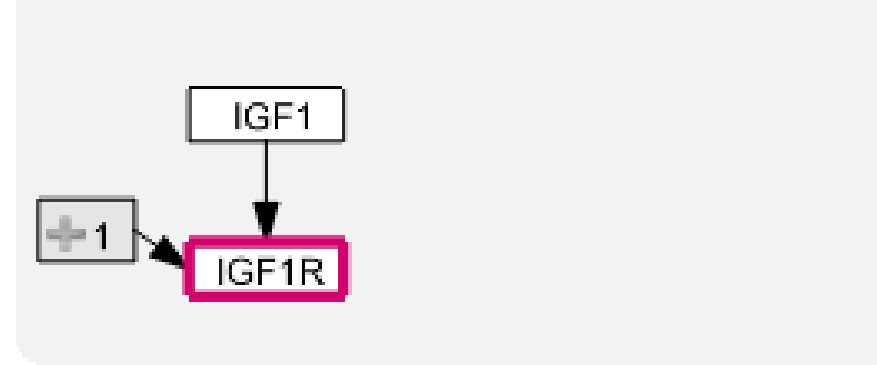
Proteoglycans in cancer



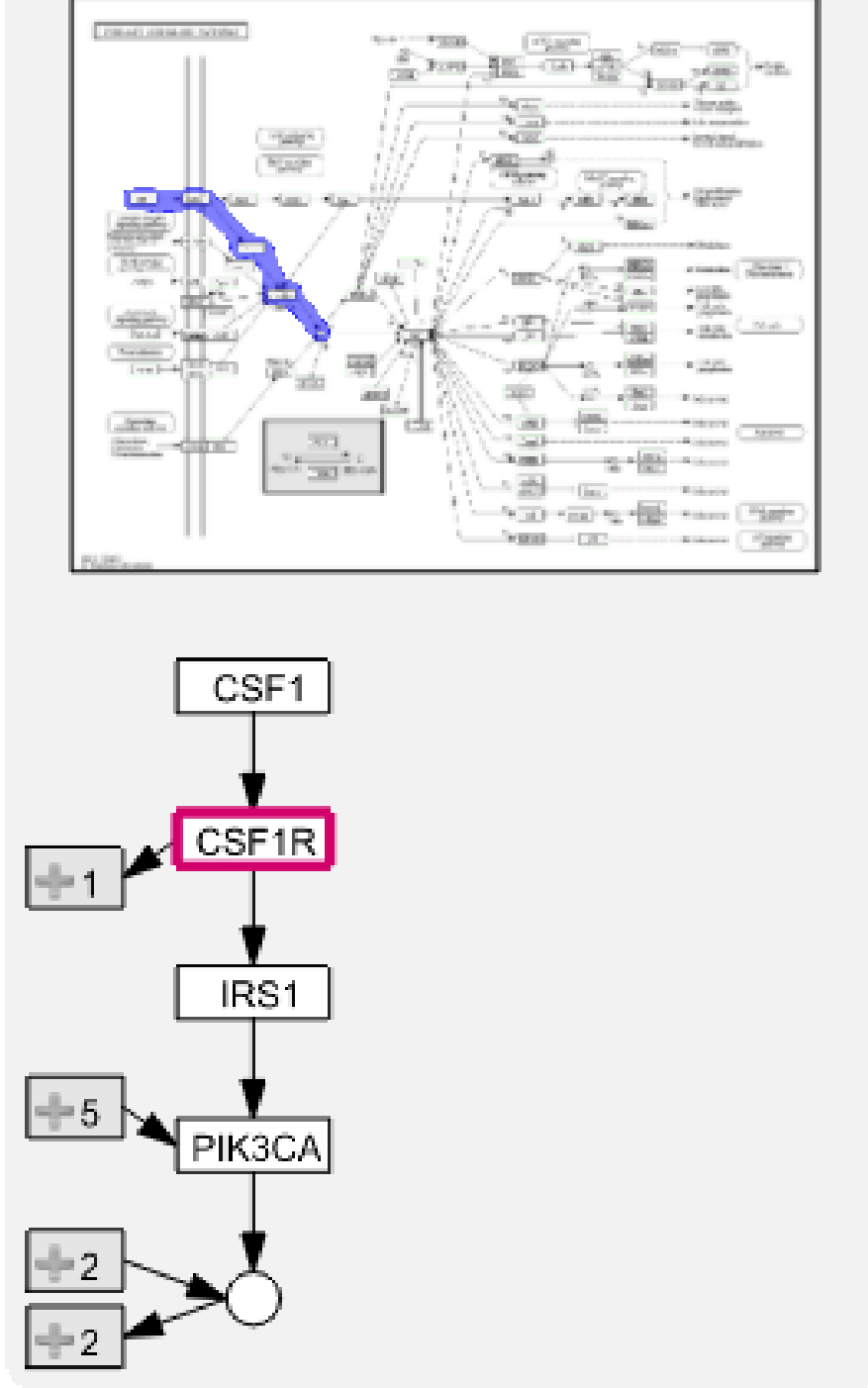
HIF-1 signaling pathway



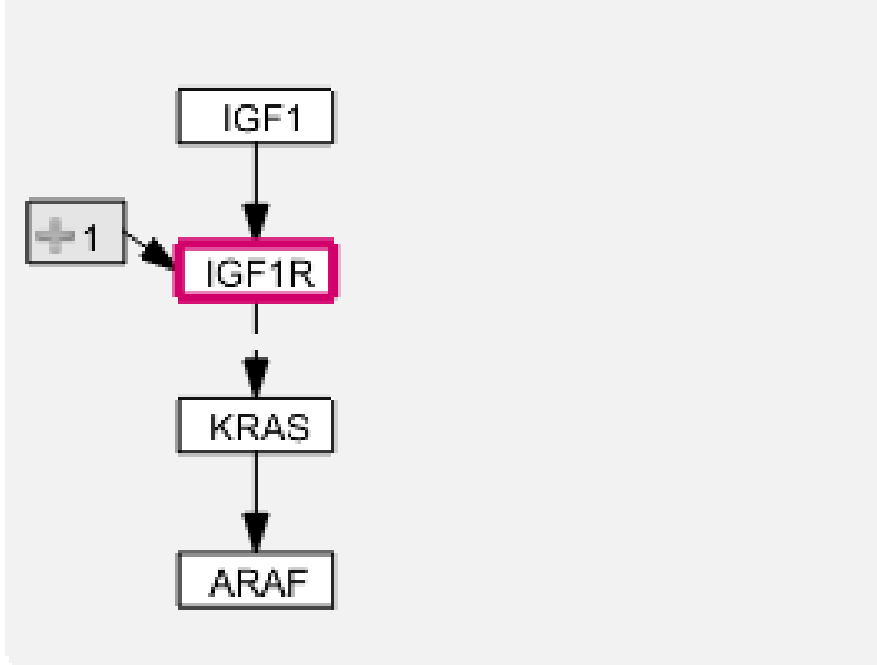
Oocyte meiosis



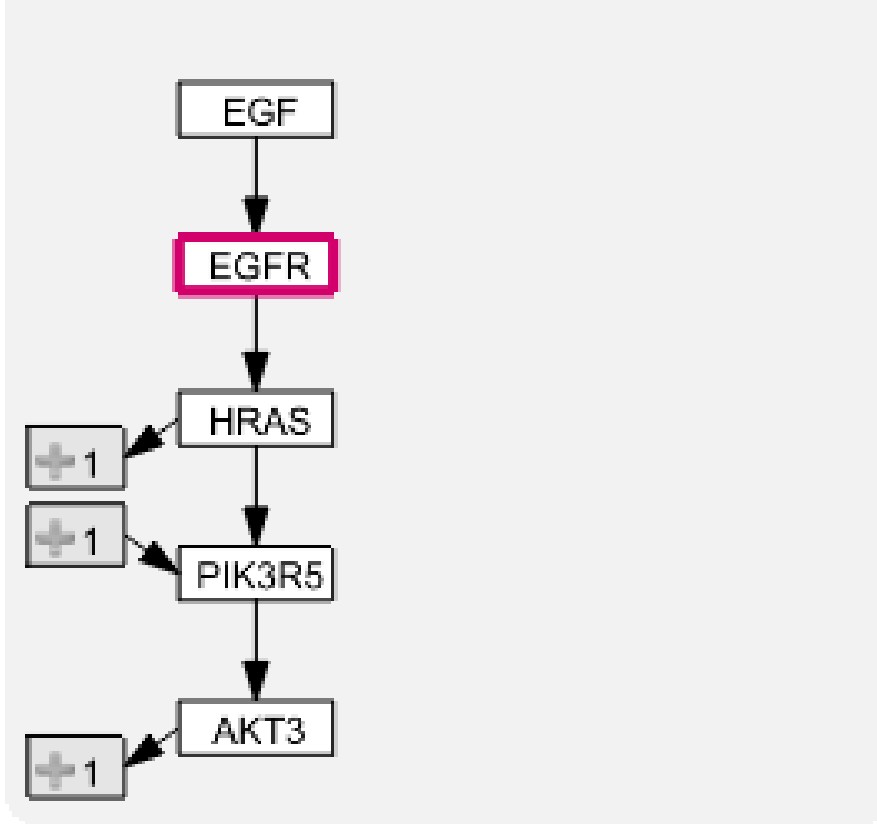
PI3K-Akt signaling pathway



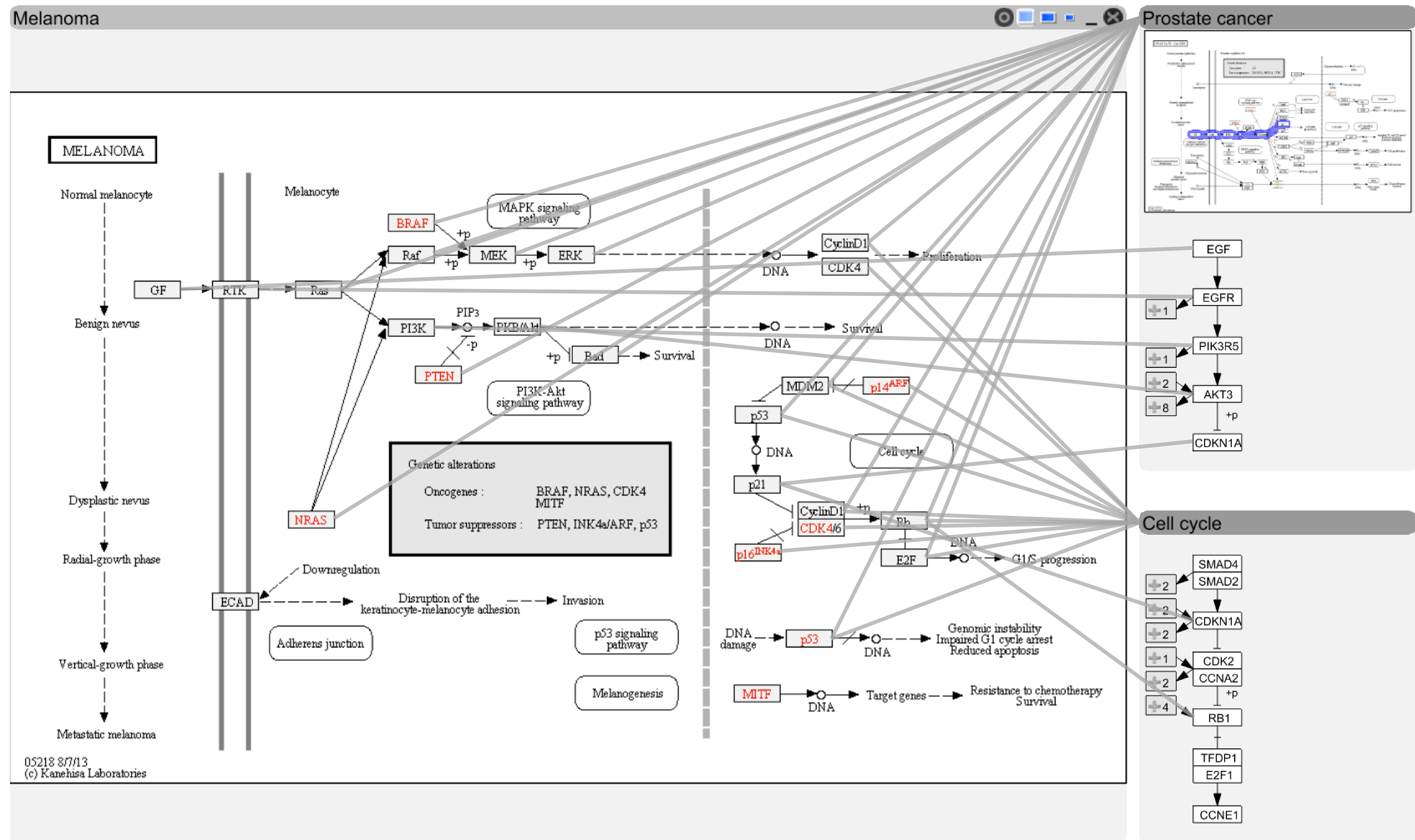
Progesterone-mediated oocyte



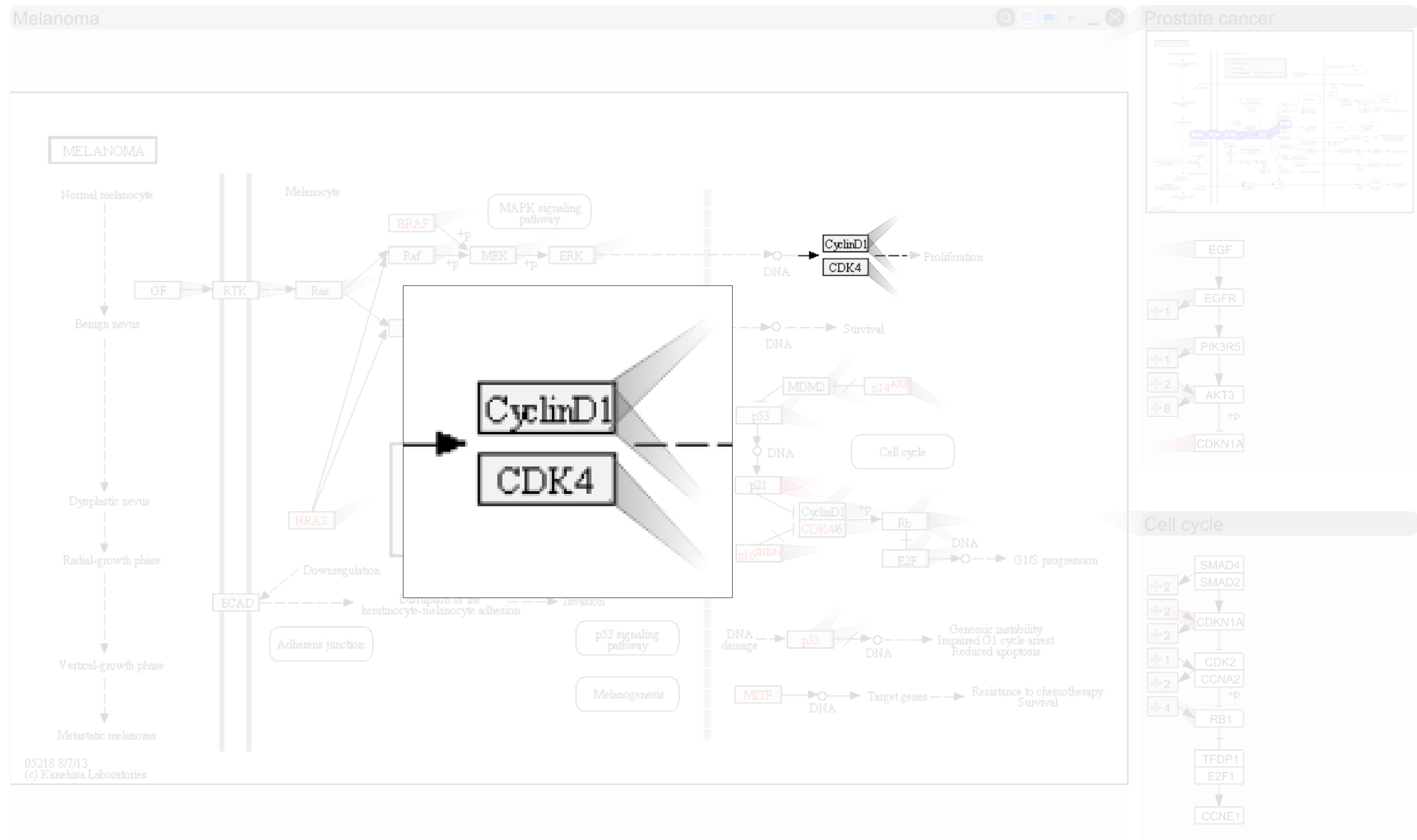
Melanoma



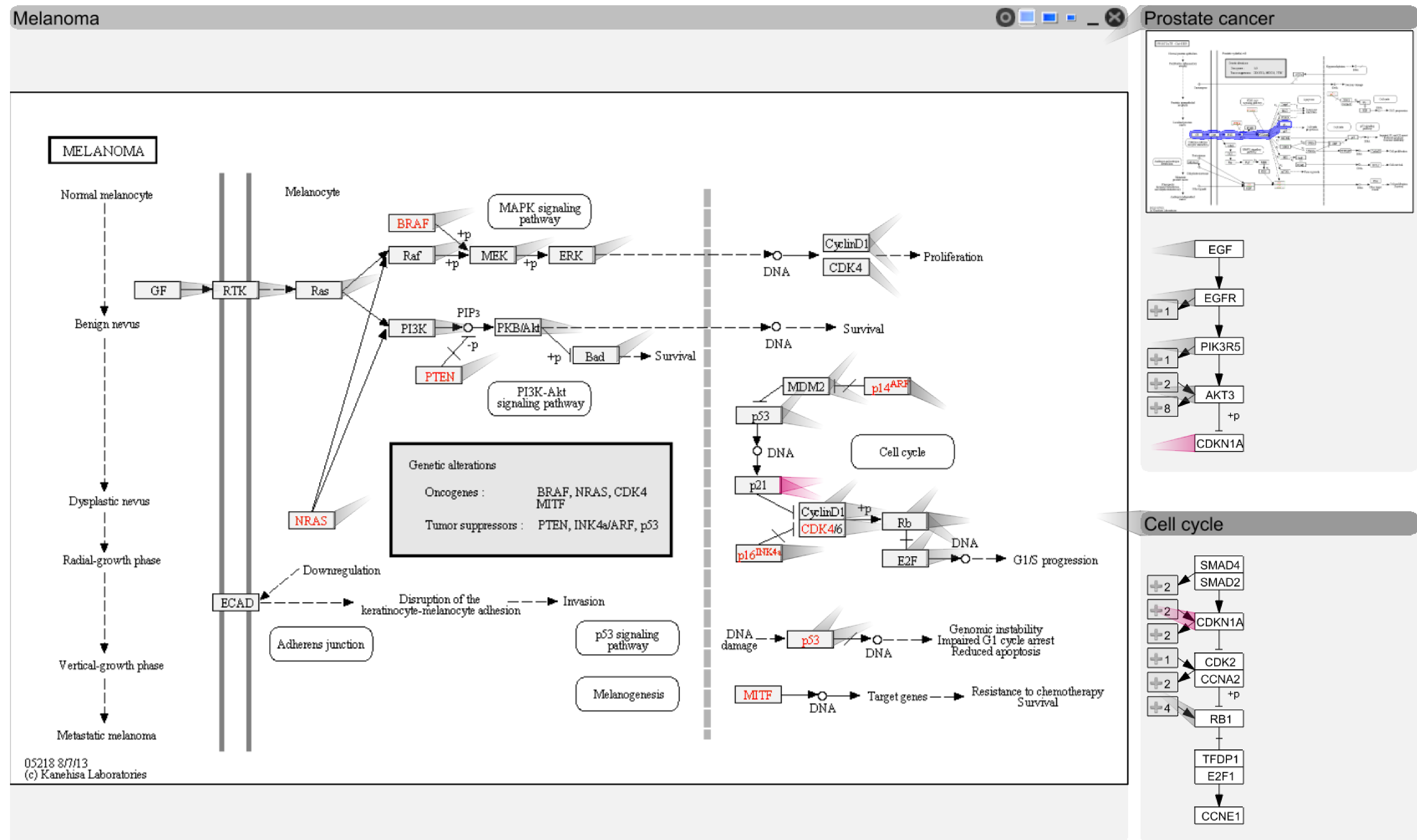
Visualizing Relationships



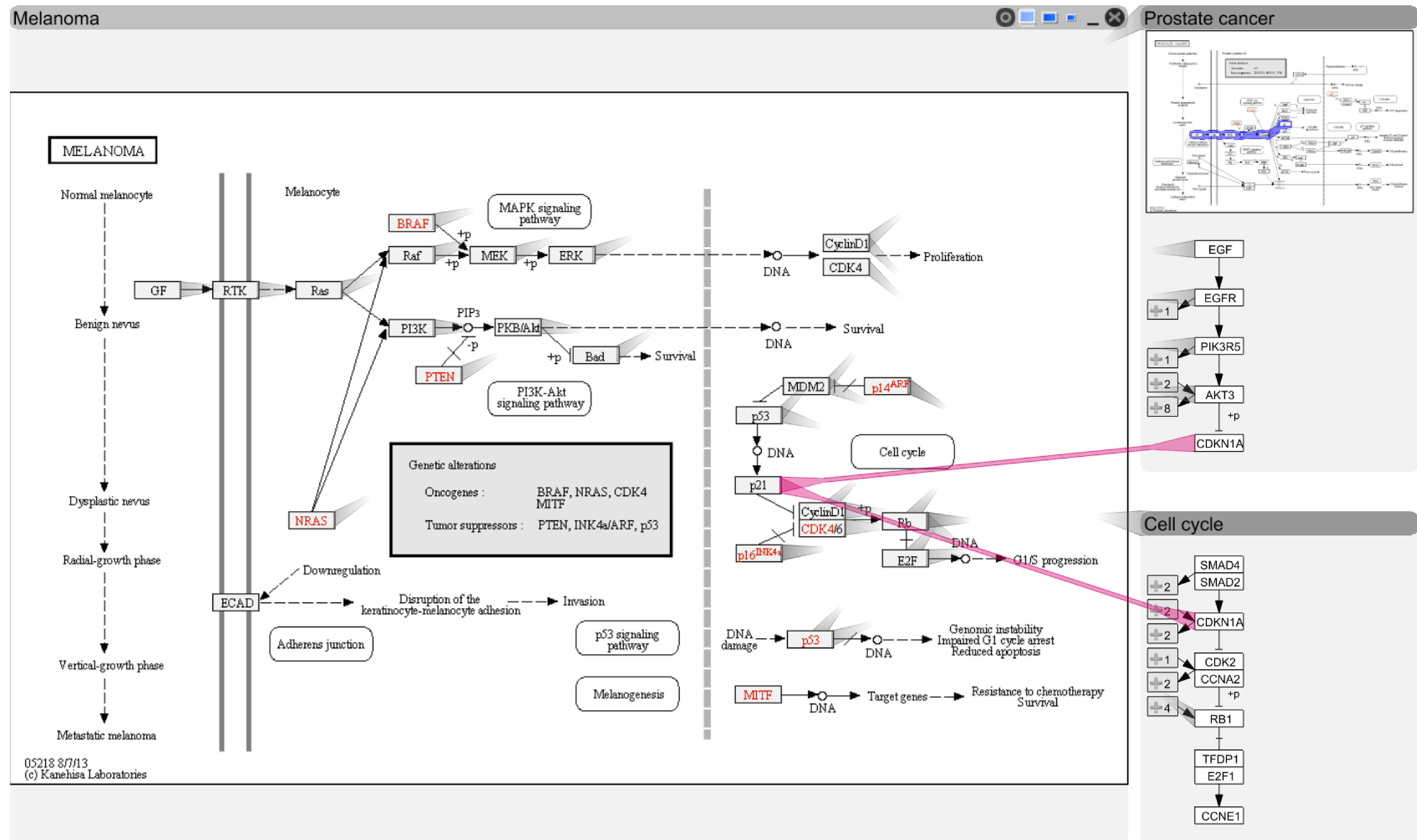
Visualizing Relationships



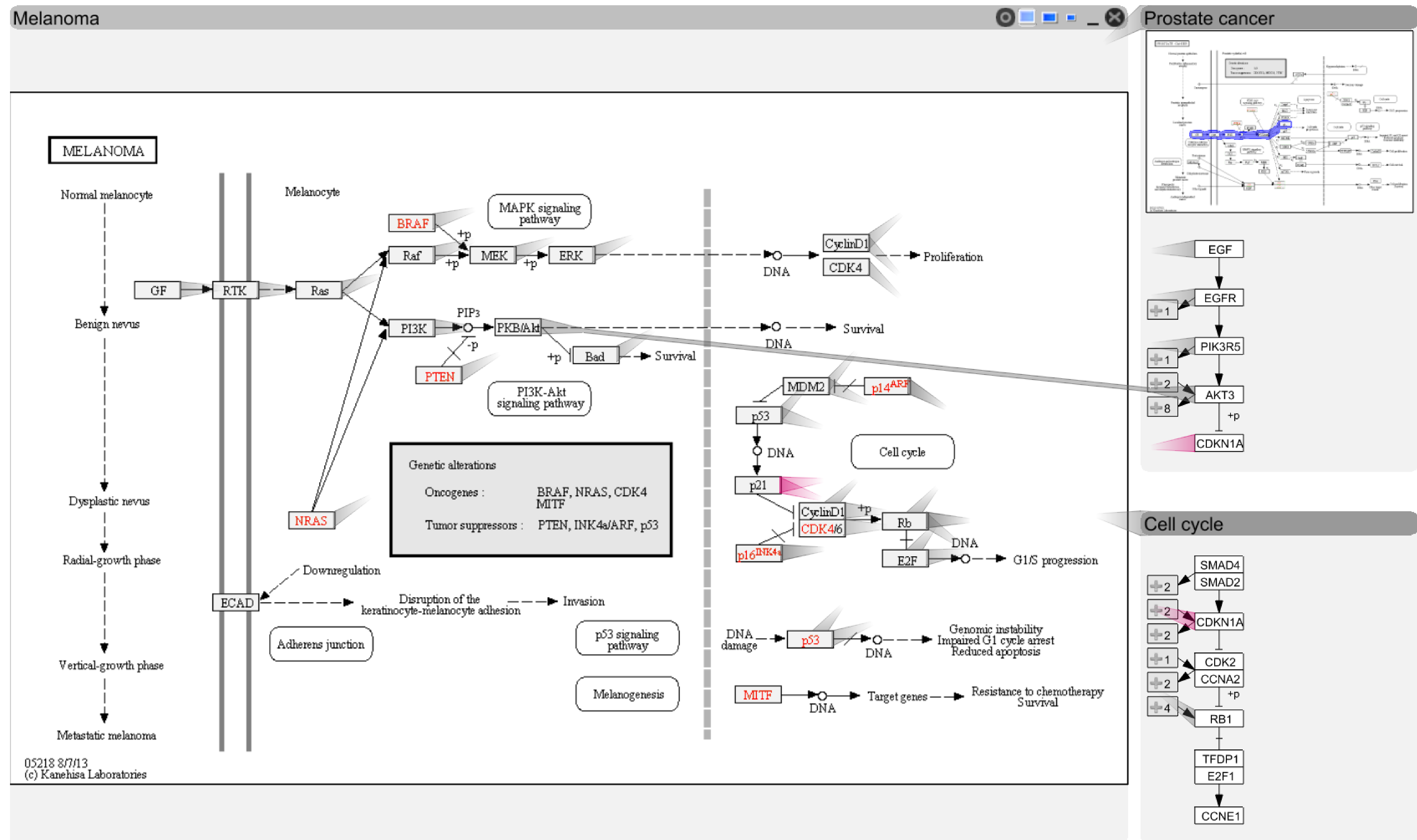
Visualizing Relationships



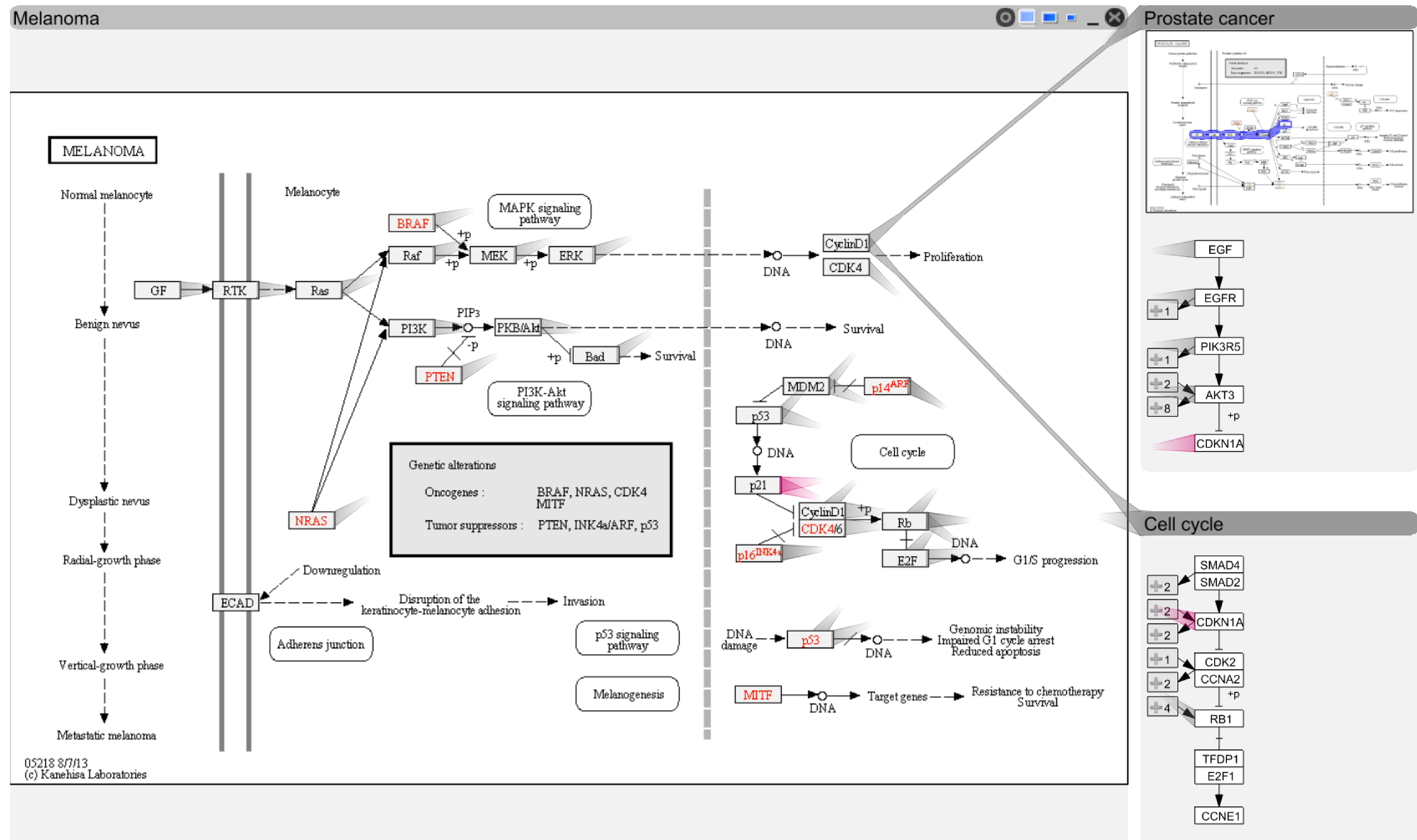
Visualizing Relationships



Visualizing Relationships



Visualizing Relationships



**How to visualize
experimental data on pathways?**

Experi- mental Data and Pathways

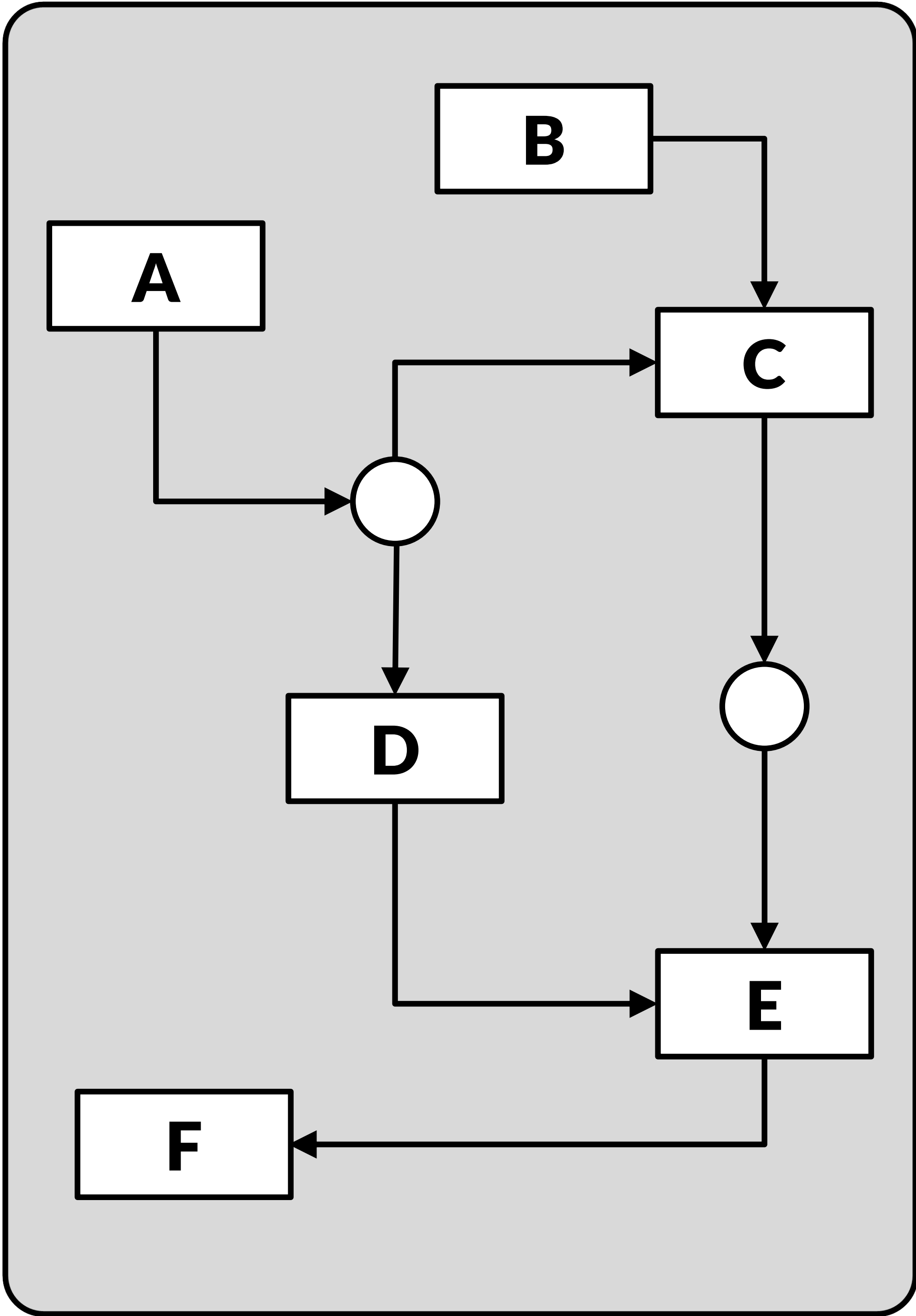
enRoute

[PartI, BioVis '12]

Cannot account for **variation** found in
real-world data

Branches can be **(in)activated** due to
mutation,
changed gene expression,
modulation due to drug treatment,
etc.

Good Old Color Coding



A	-3.4	4.2
B	2.8	1.8
C	3.1	-2.2
D	-3	-2.8
E	0.5	0.3
F	0.3	0.3

Challenge: Data Scale & Heterogeneity

Large **number of experiments**

Large datasets have more than 500 experiments

Multiple **groups/conditions**

Different **types** of data, require **different visualization techniques**

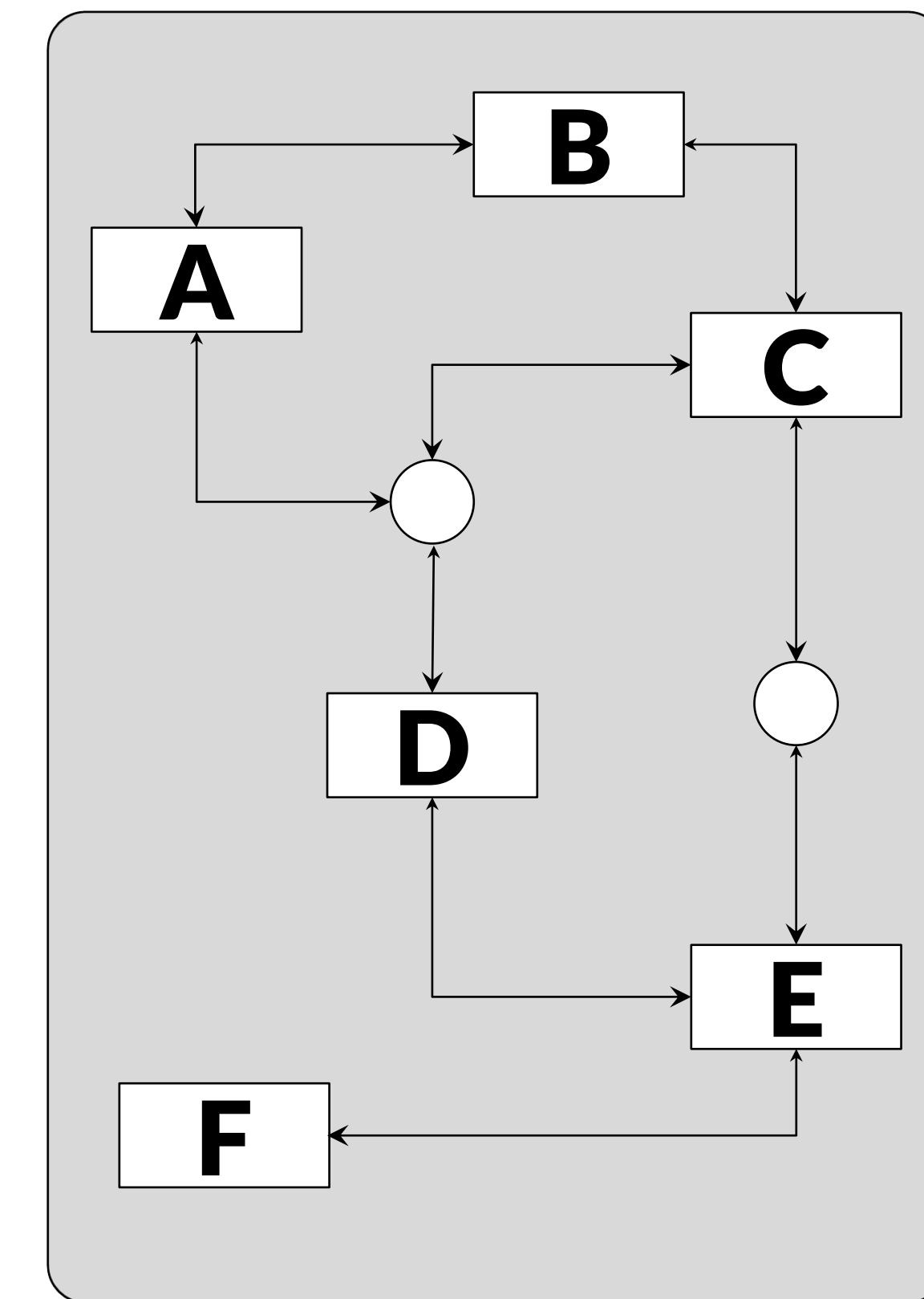
Challenge: Supporting Multiple Tasks

Two central tasks:

Explore **topology** of pathway

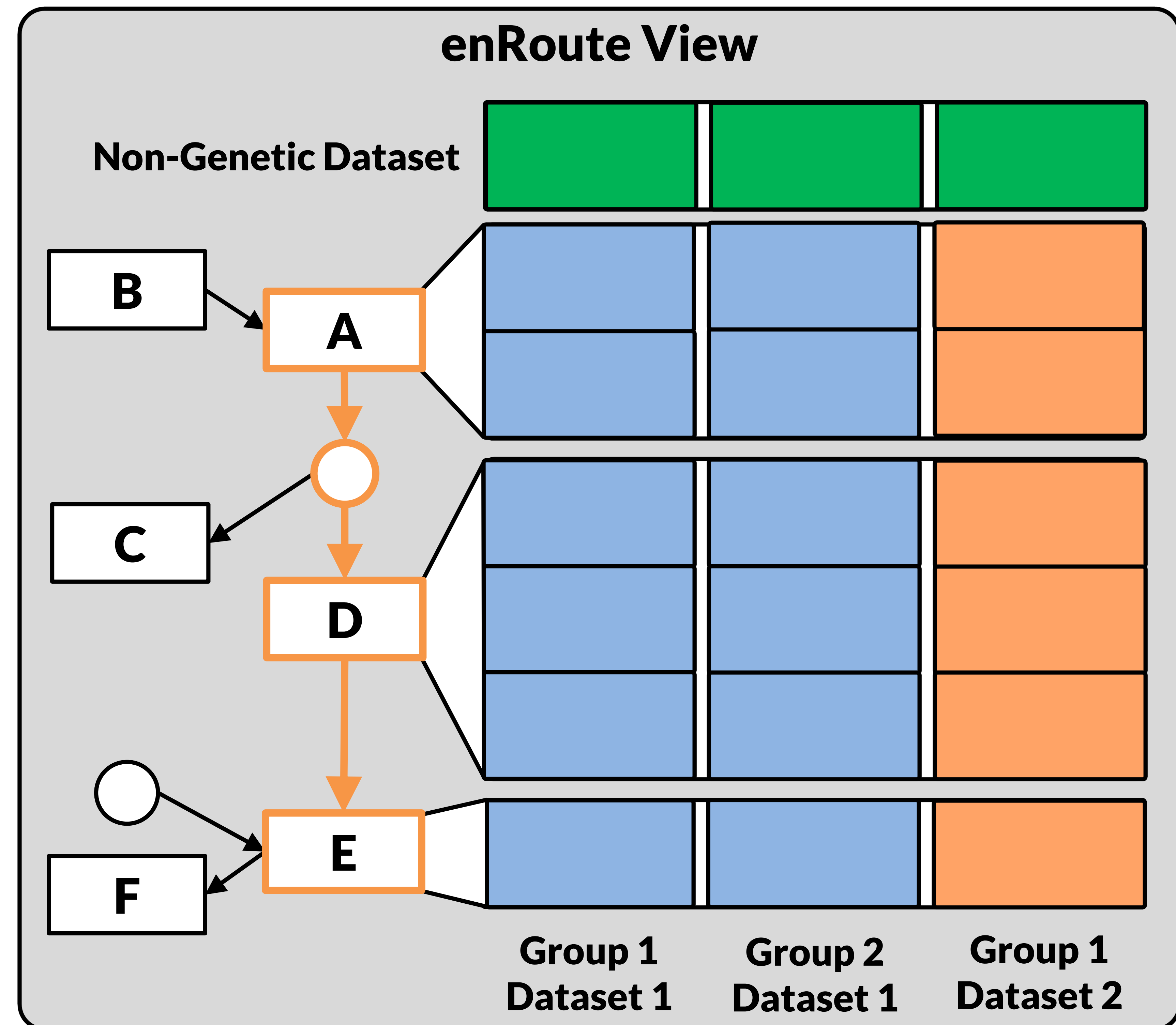
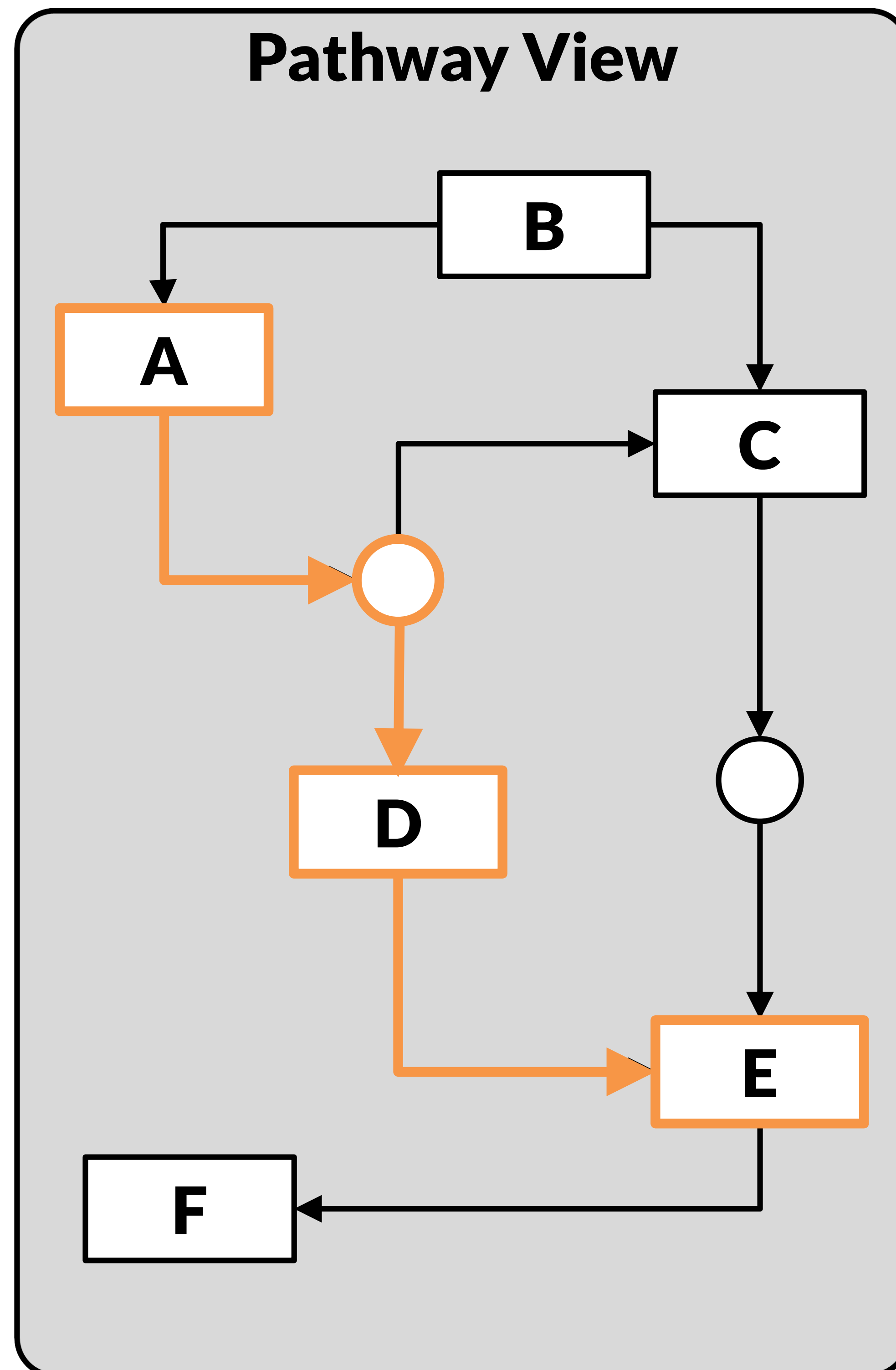
Explore the **attributes** of the nodes
(experimental data)

Need to support both!

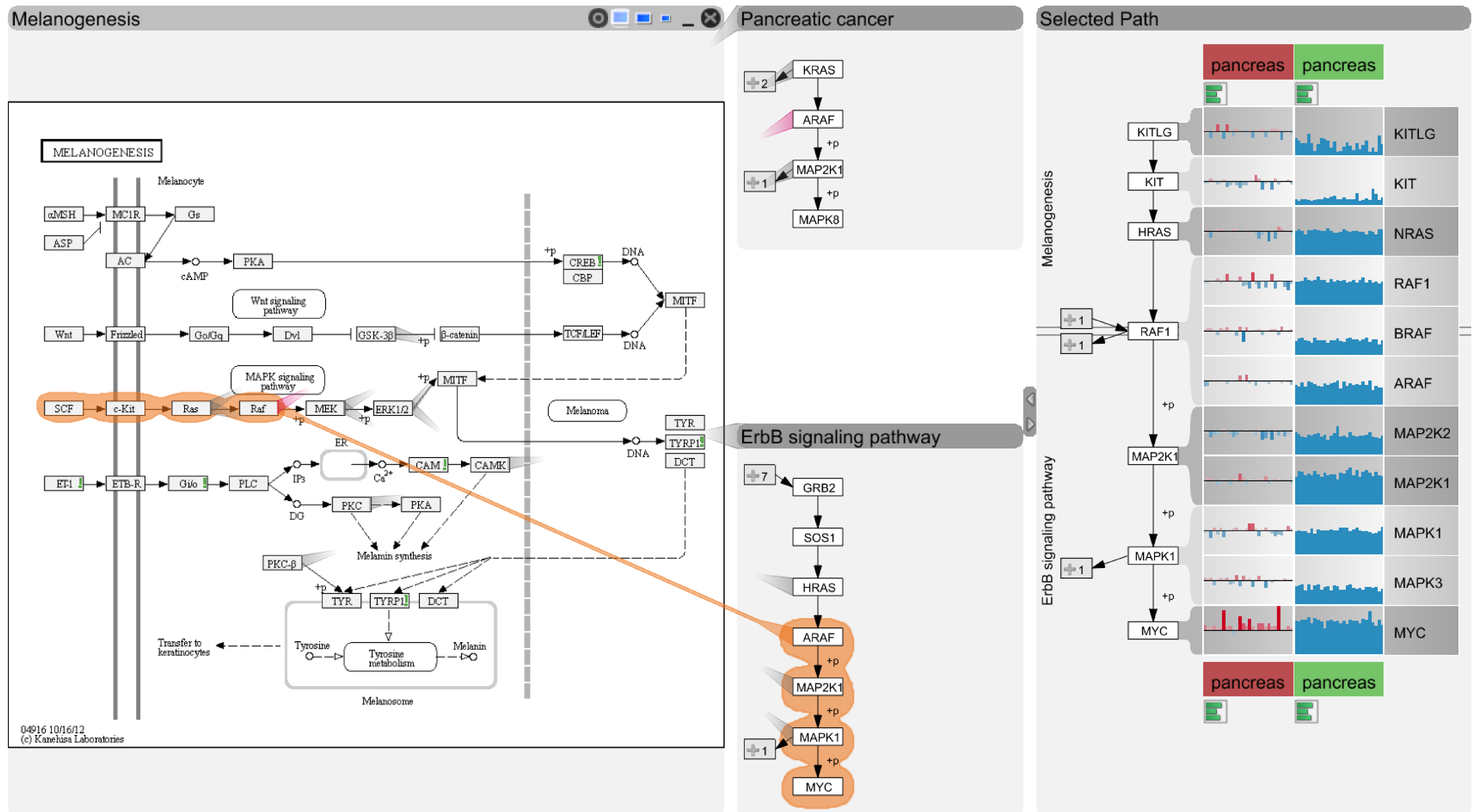


	Sample 1	Sample 2	Sample 3
Gene 1	1	1.1	0.4
Gene 2	2	0.5	1.2
Gene 3	1.4	0.2	0.5
Gene 4	0.3	0.5	0.7

Concept



enRoute



FileDataWindowViewHelp

Entourage

Pathways

Pathway

Filter:
<None>

1 C donor

2-Oxocarboxylic acid

ABC transporters

ABC-family proteins

ACE Inhibitor Pathwa

Acetylcholine Synthes

Acute myeloid leukem

Adherens junction

Adipocyte TarBase

Adipocytokine signali

Adipogenesis

Advanced glycosylatio

Aflatoxin B1 metaboli

African trypanosomias

AGE/RAGE pathway

AhR pathway

Alanine and aspartate

Alanine, aspartate an

Alcoholism

Aldosterone-regulated

Allograft rejection

Allograft rejection

Alpha 6 Beta 4 signal

alpha-Linolenic acid

Alzheimer's disease

Alzheimers Disease

amino acid conjugatio

amino acid conjugatio

Amino sugar and nucl

Aminoacyl-tRNA bios

Amoebiasis

Amphetamine addicti

AMPK signaling

Amyotrophic lateral sc

Androgen receptor si

Angiogenesis

Angiogenesis

angiogenesis overvie

Antigen processing an

APC/C-mediated degra

Apoptosis

Apoptosis

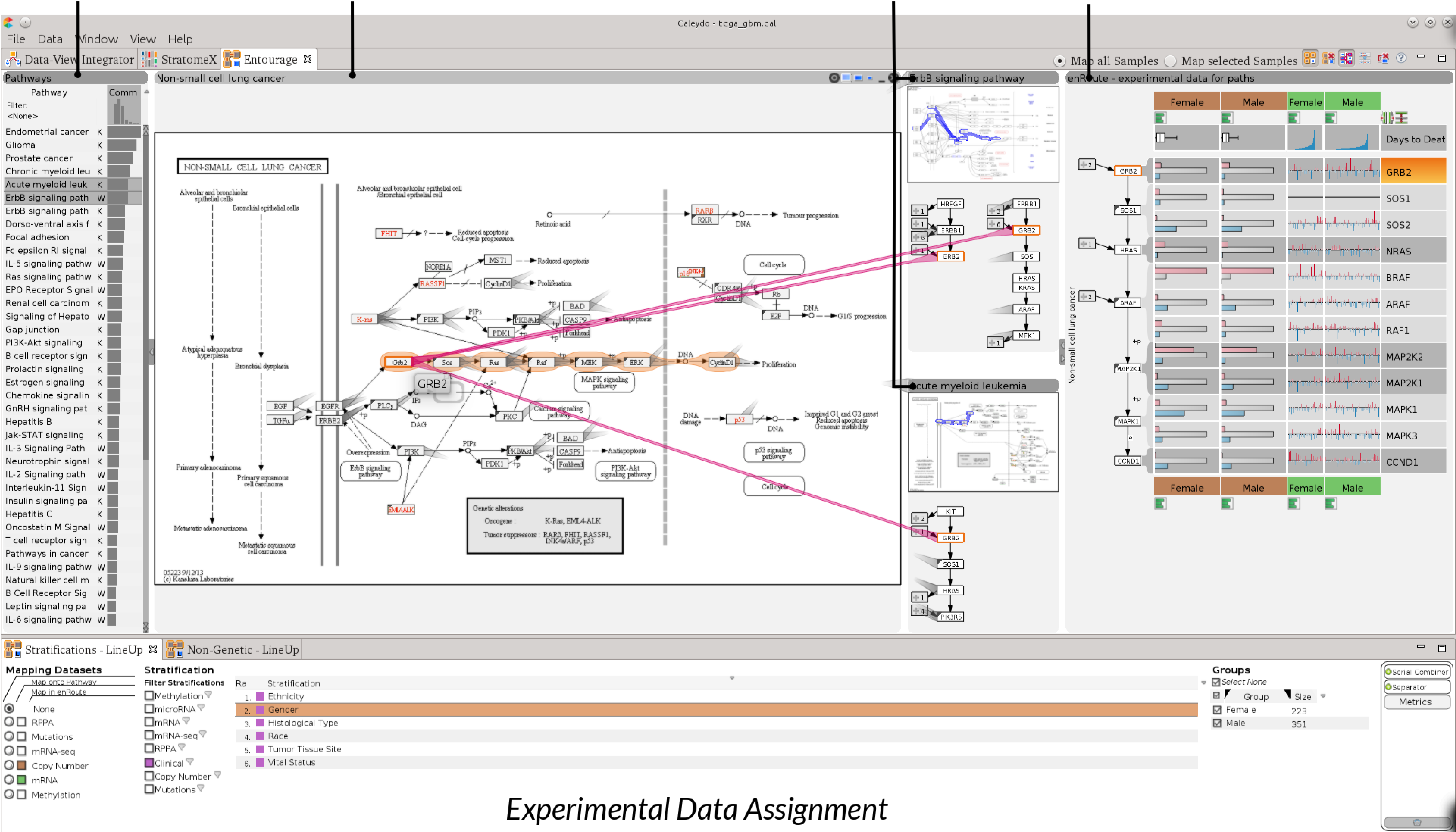
Apoptosis Meta Path

Apoptosis Modulation

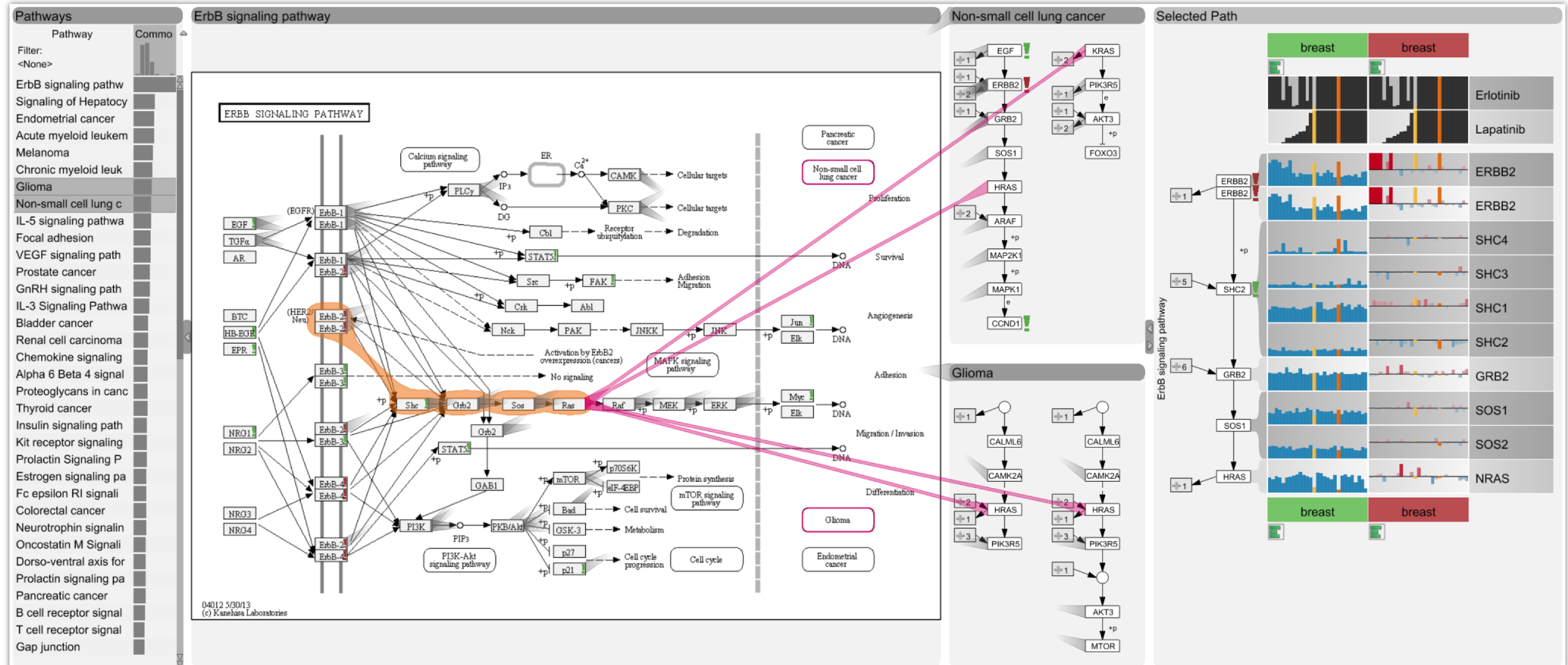
Apoptosis Modulation

Apoptosis, anoikis an

Selected Path



Case Study: CCLE Data

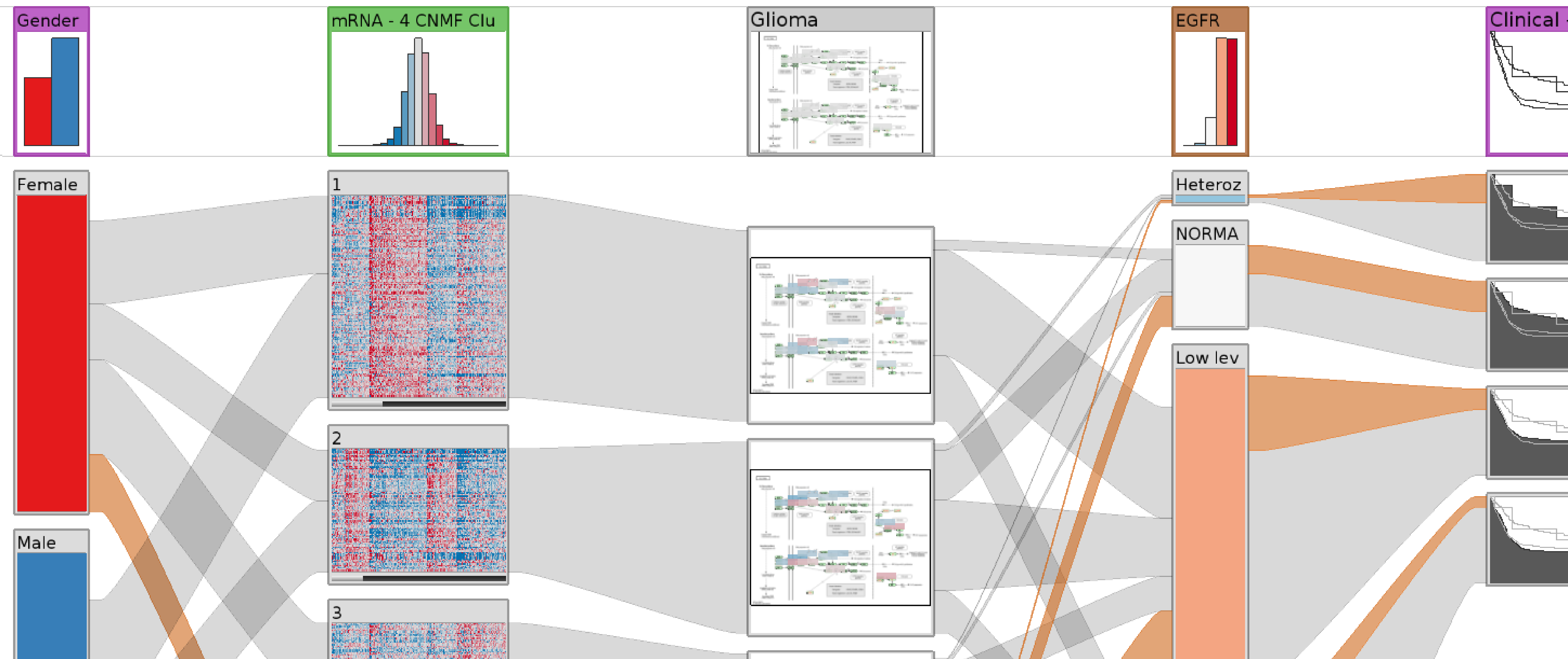


More Information:

<http://entourage.caleydo.org>

[Lex, EuroVis '12]
3rd Best Paper Award
[Streit, Nature Methods '14]
To appear

Heterogeneous Data Cancer Subtypes



Cancer Subtypes

Cancer is not homogeneous

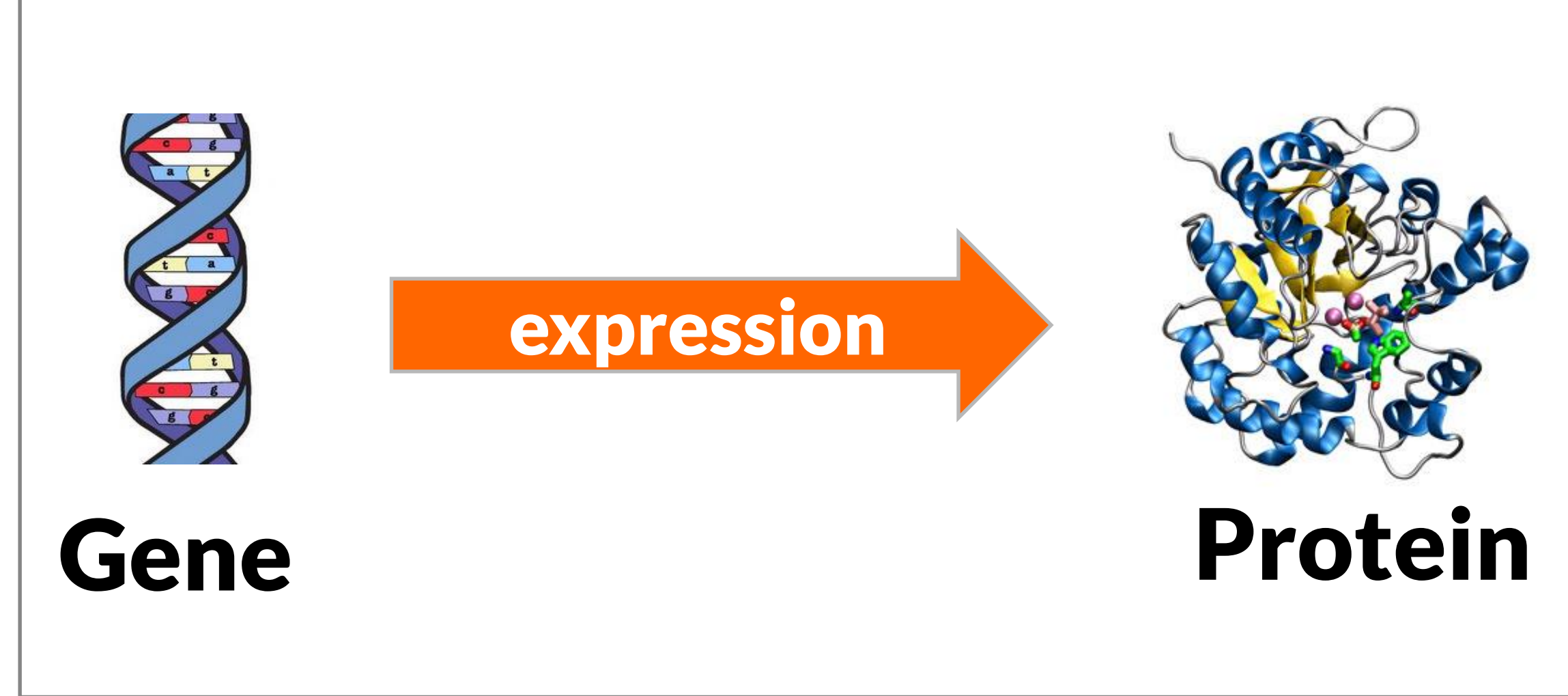
different histology

different **molecular alterations**

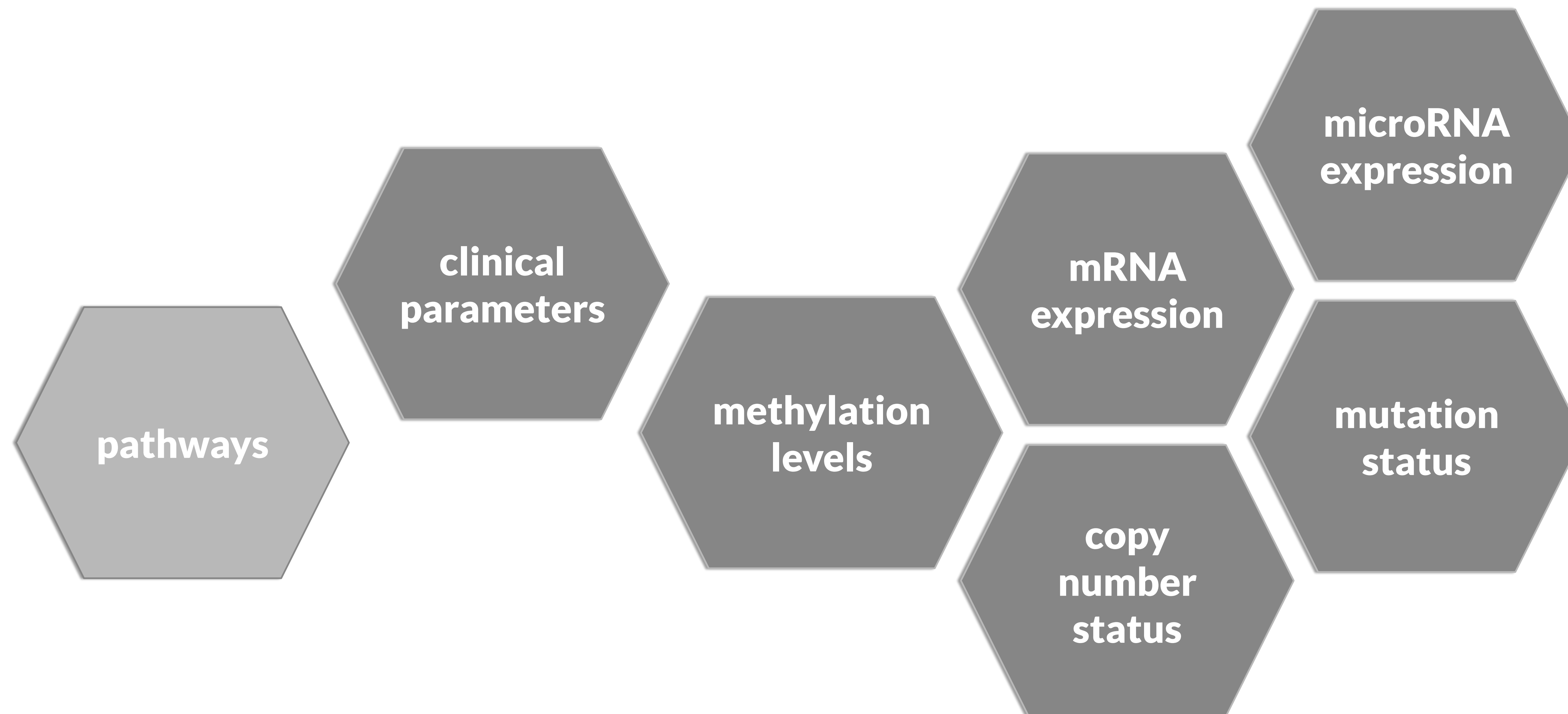
Subtypes have serious implications

different **treatment for subtypes**

****prognosis** varies between subtypes**



The Cancer Genome Atlas  *Understanding genomics
to improve cancer care*



StratomeX **visualizes...**

- ... the relationships between multiple heterogeneous datasets**
- ... the data within the datasets**
- ... alternative clusterings & groupings**
- ... the effect of groupings on clinical parameters & biological processes**

Stratifying Patients

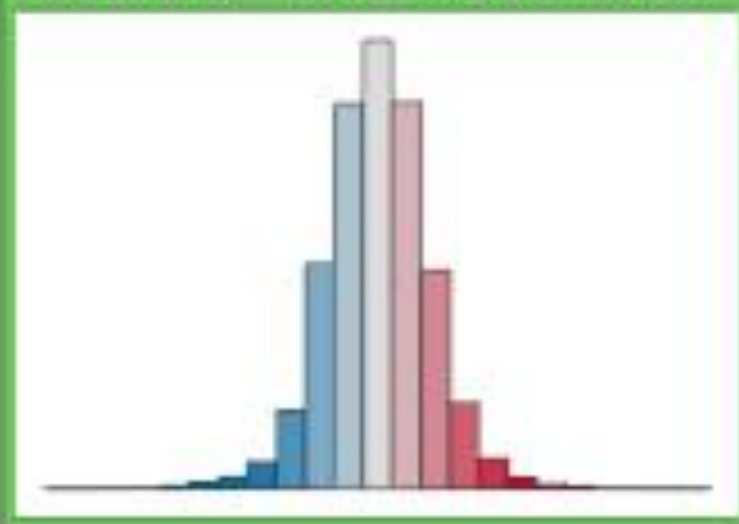
Cluster A1

Cluster A2

Cluster A3

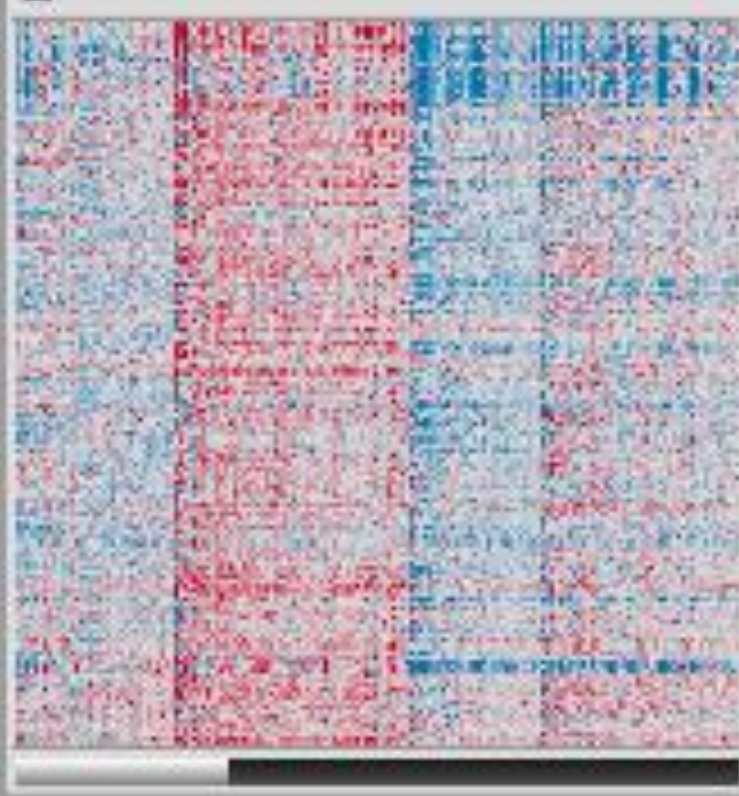
**Subtypes are identified
by stratifying datasets, e.g.,
based on an expression pattern
a mutation status
a copy number alteration
a combination of these**

mRNA - 4 CNMF Clu

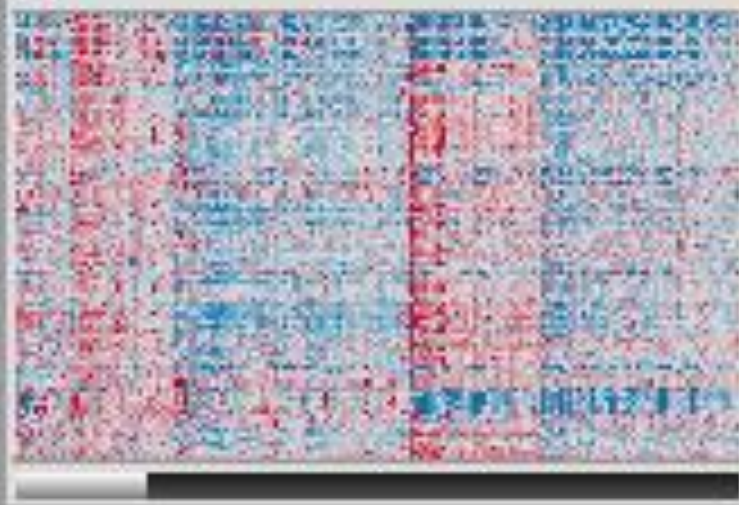


**Header /
Summary of
whole Stratification**

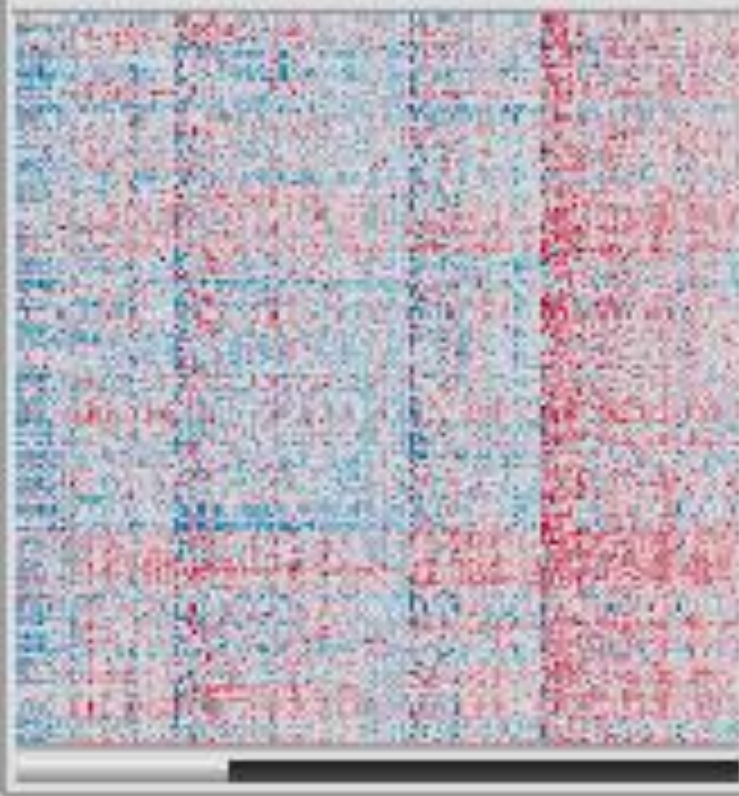
1



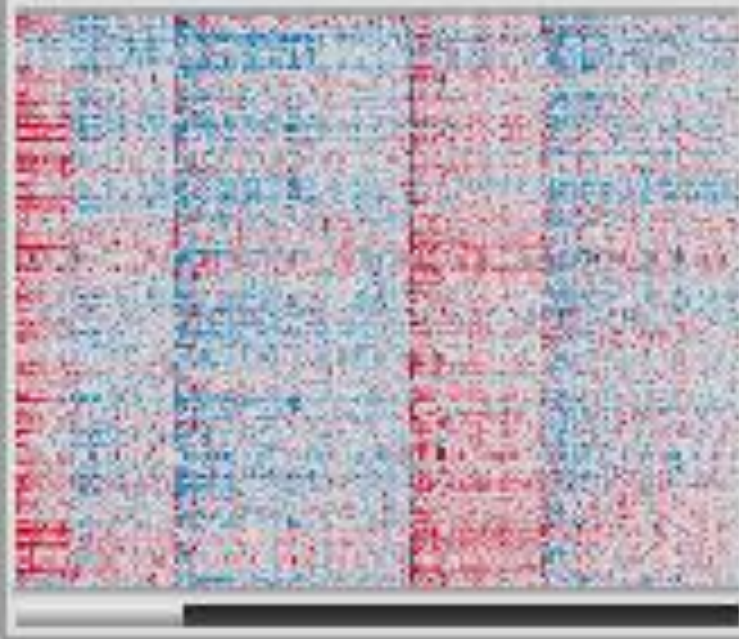
2



3



4



Patients

**Candidate Subtype /
Heat Map**

Genes

Comparing Stratifications

Cluster A1

Cluster A2

Cluster A3

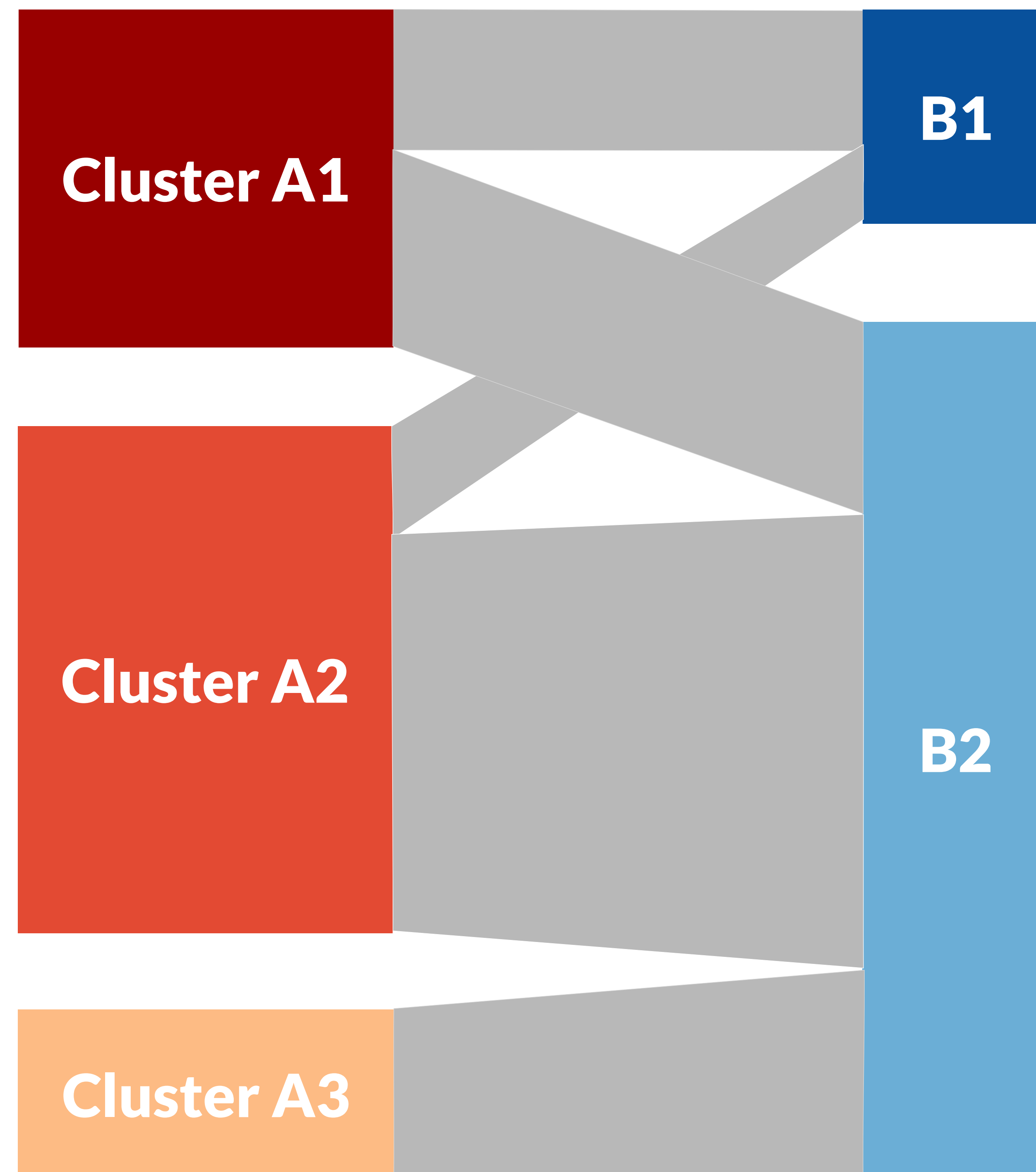
**Tabular
e.g., mRNA**

B1

B2

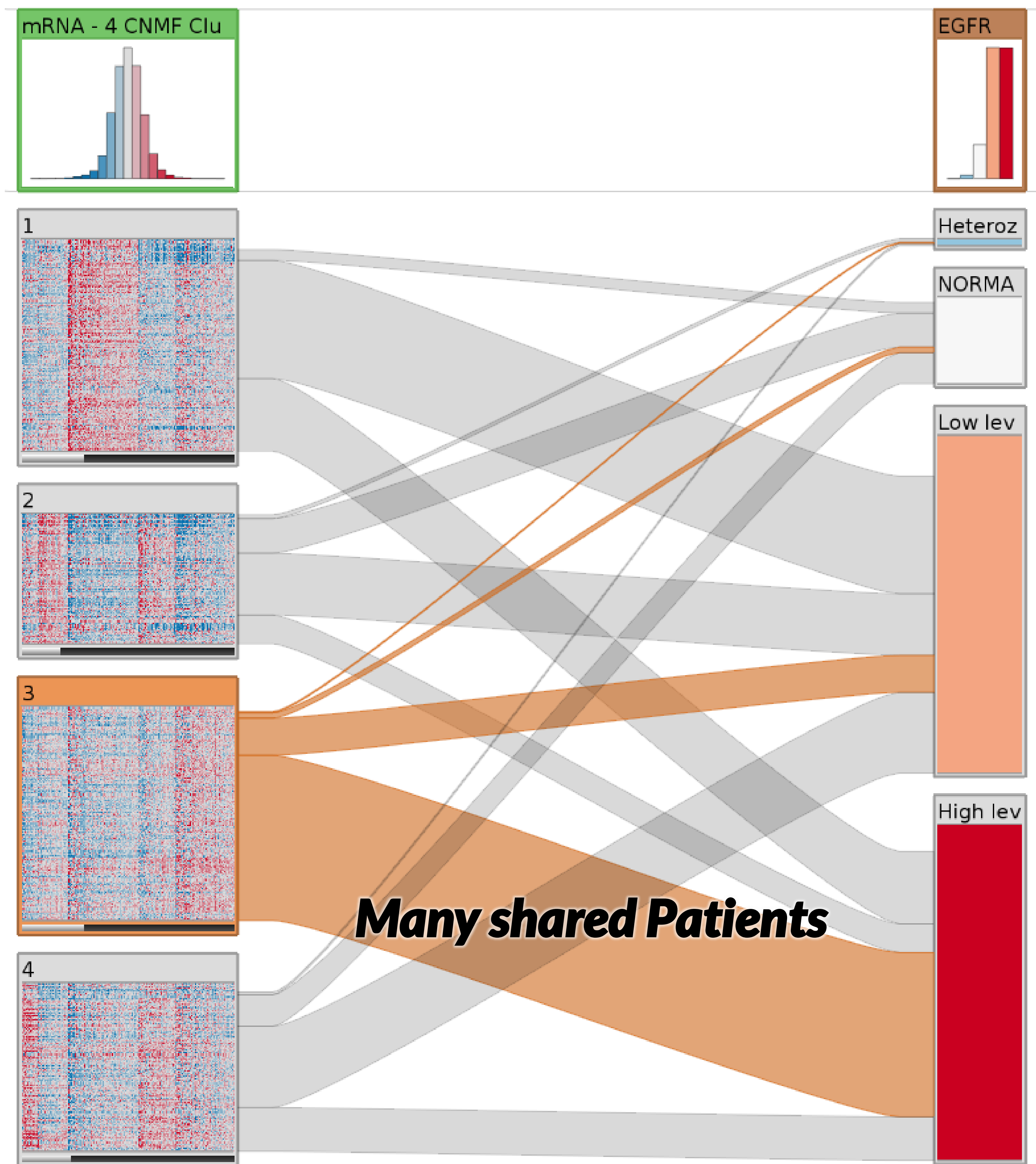
**Categorical,
e.g., mutation status**

Comparing Stratifications



Tabular
e.g., mRNA

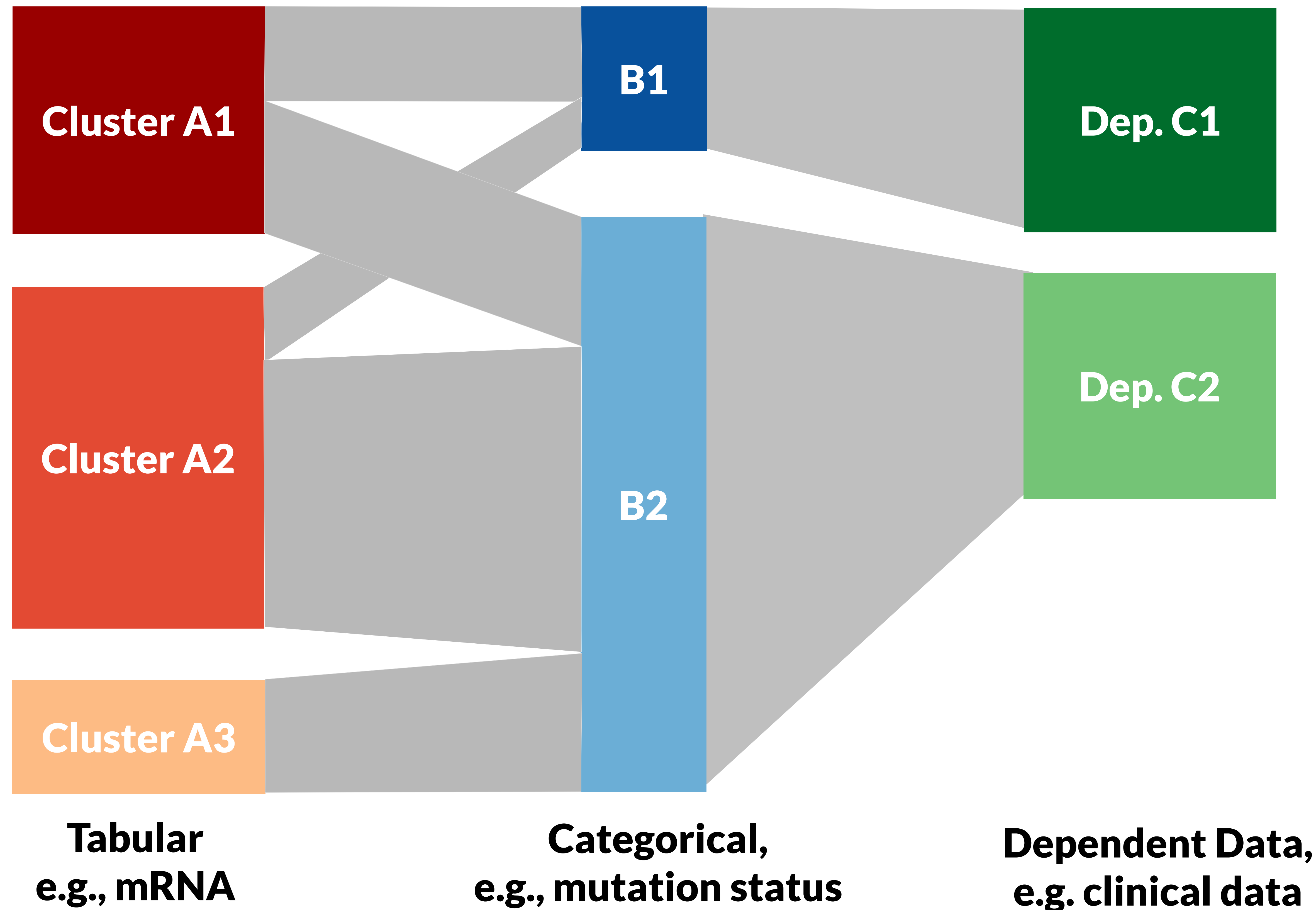
Categorical,
e.g., mutation status

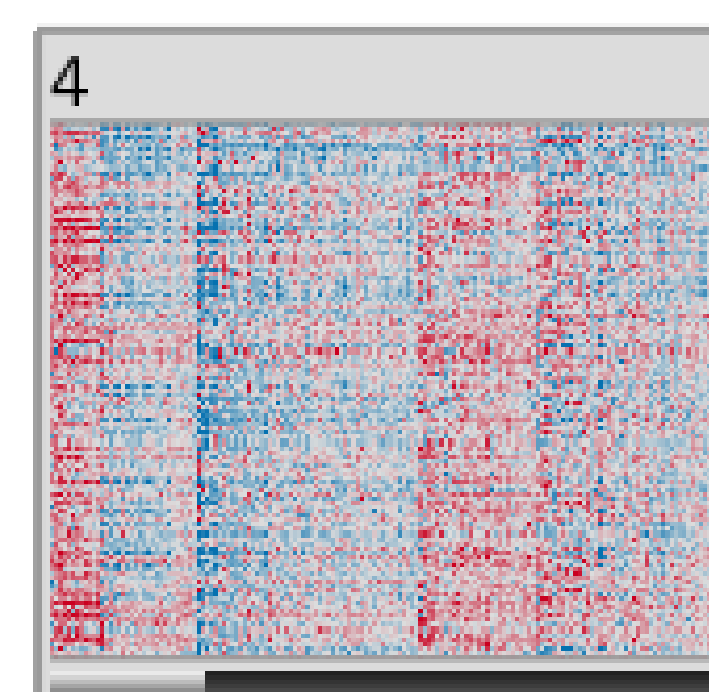
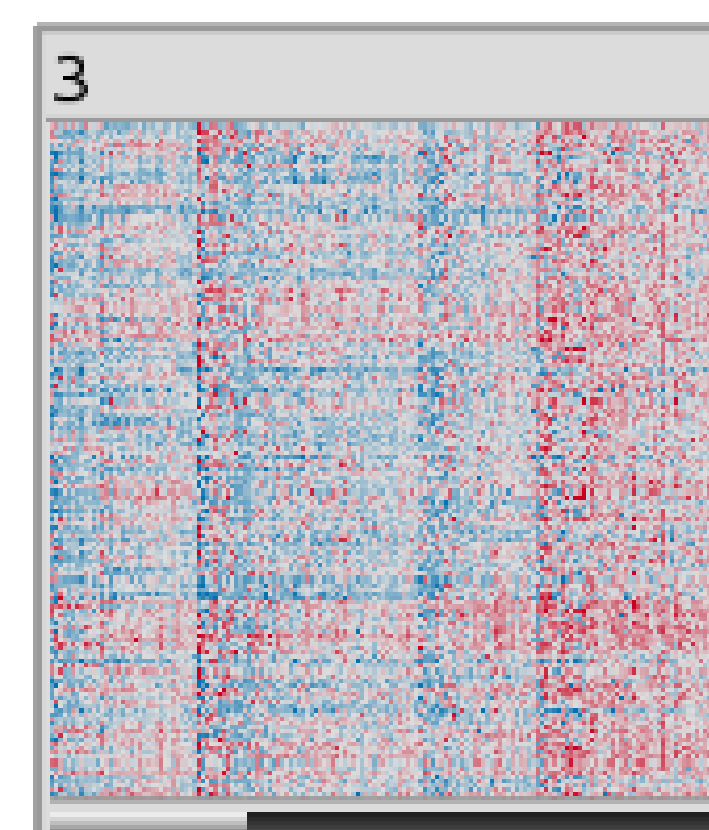
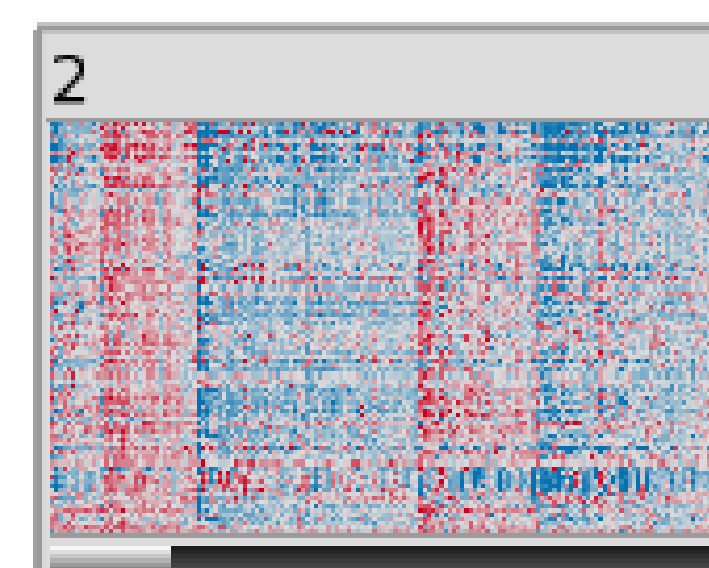
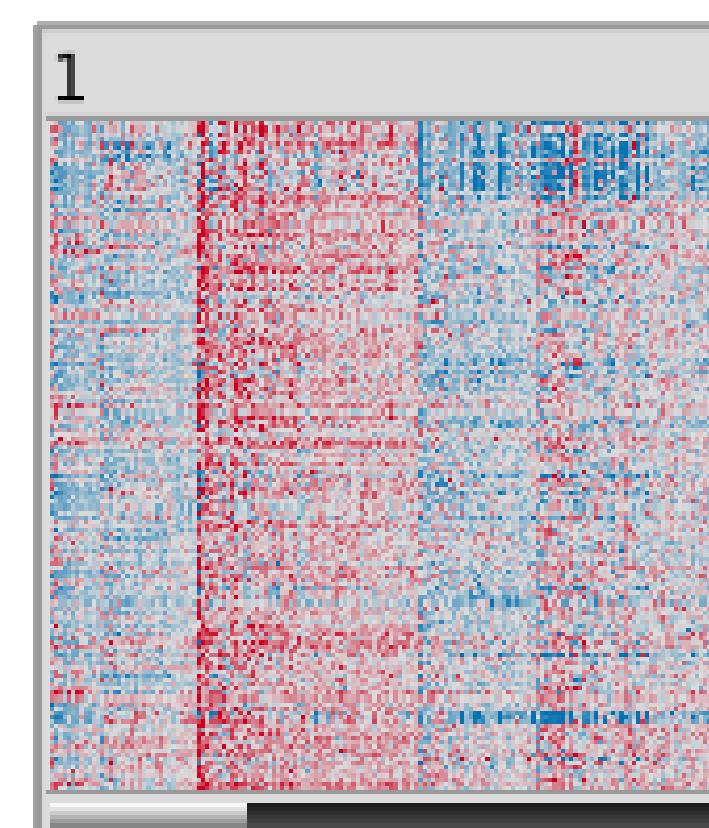
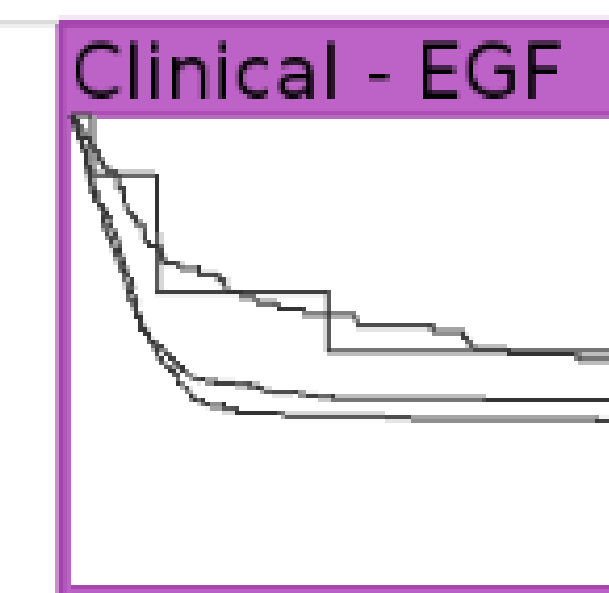
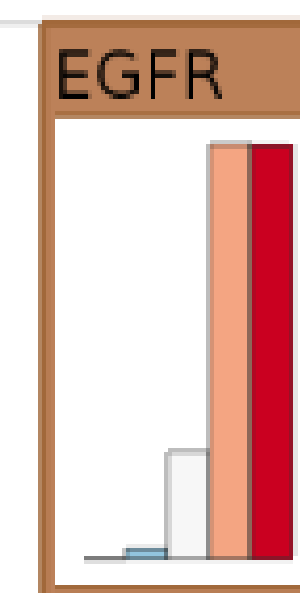
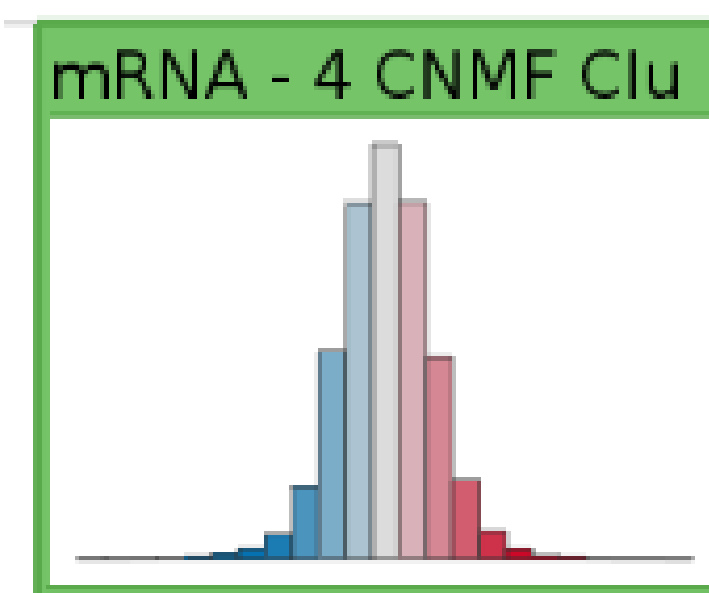


**Clustering of
mRNA Data**

**Stratification on
Copy Number Status**

Other Data – Same Stratification



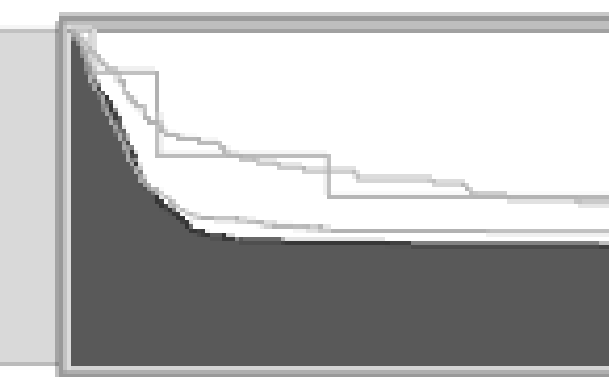
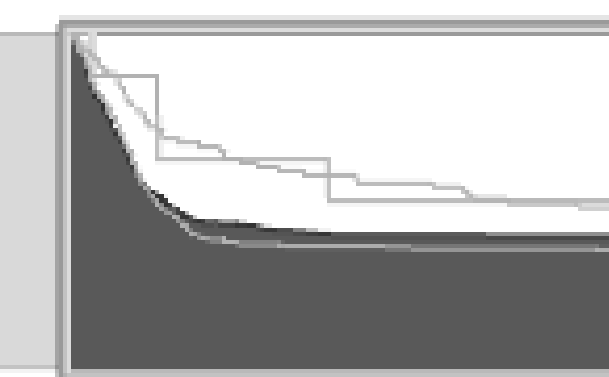
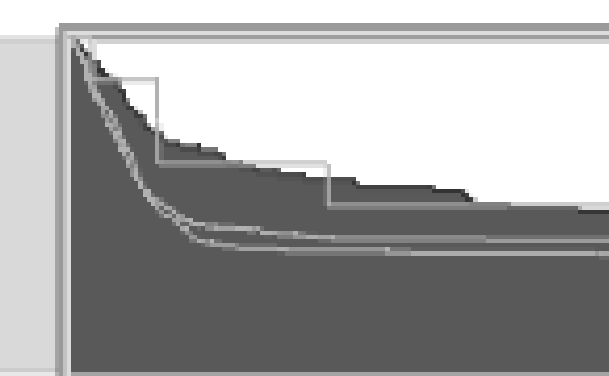
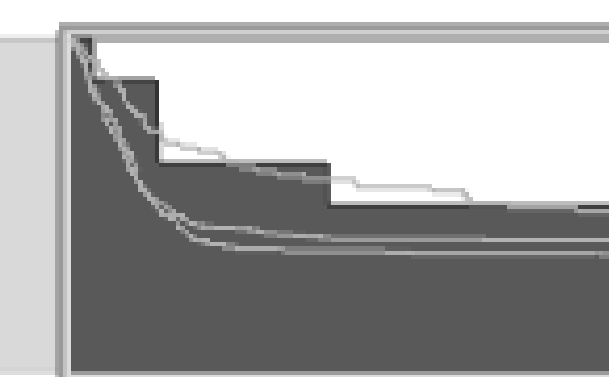


Heteroz

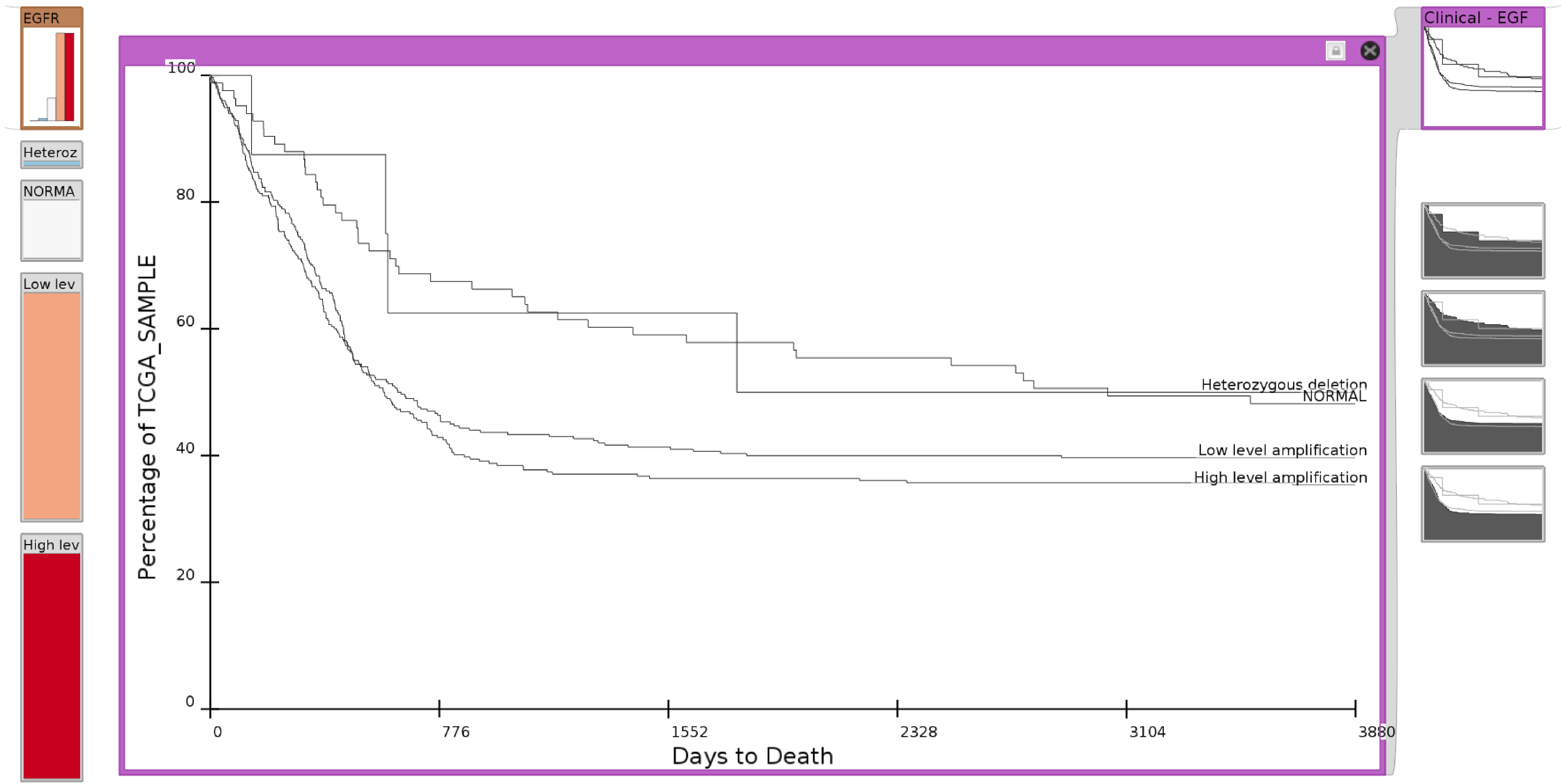
NORMA

Low lev

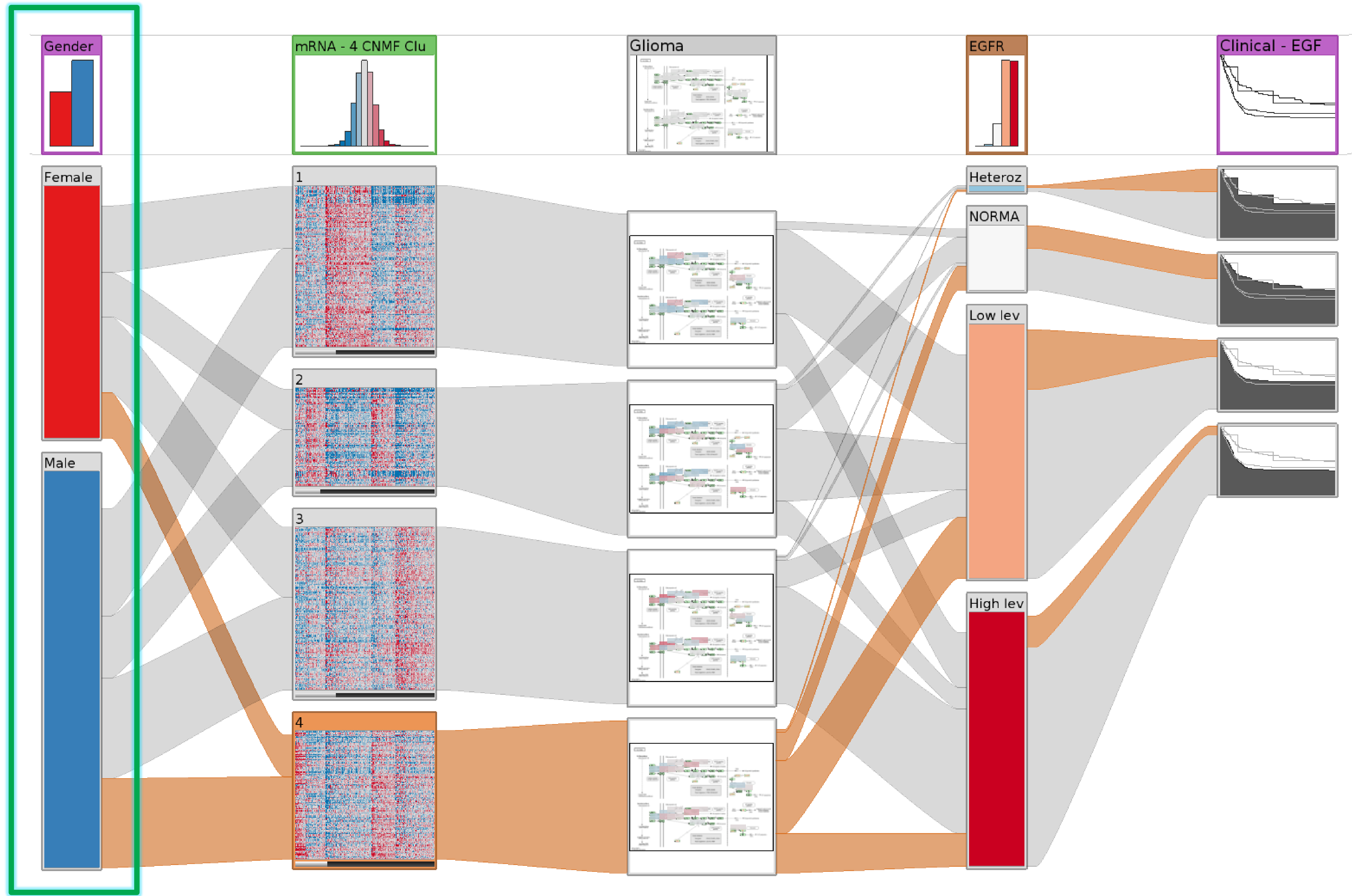
High lev



***Survival data in
Kaplan Meier plots***

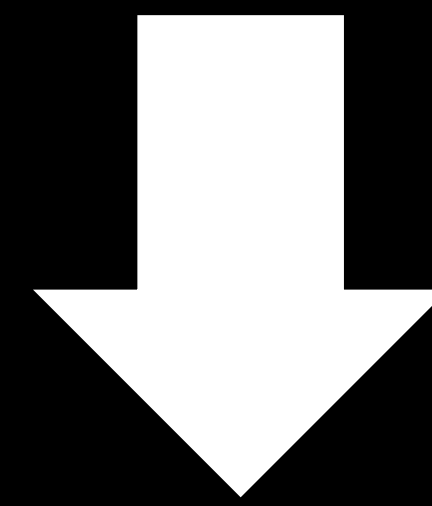


Detail View



***Stratification based on
clinical variable (gender)***

Knowledge Driven Approach



Data Driven Approach

Finding Relevant Stratifications

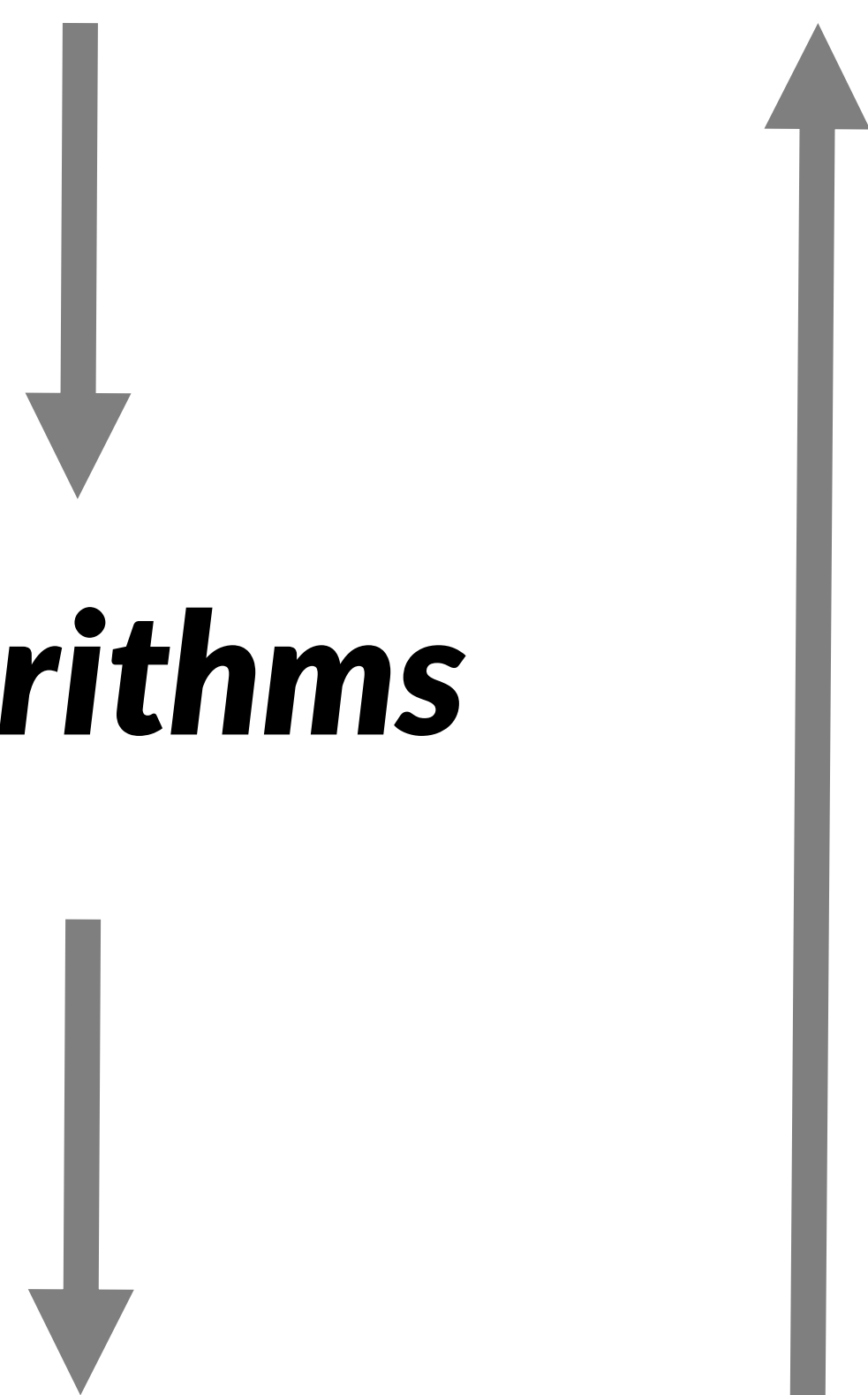
- ~ 10 datasets**
- ~ 15 clusterings per matrix**
- ~ 15,000 stratifications for copy number & mutations**
- ~ 500 pathways**
- ~ 20 clinical variables**

Calculate scores for matches

Rank the results

Integrating Computation

Data Visualization

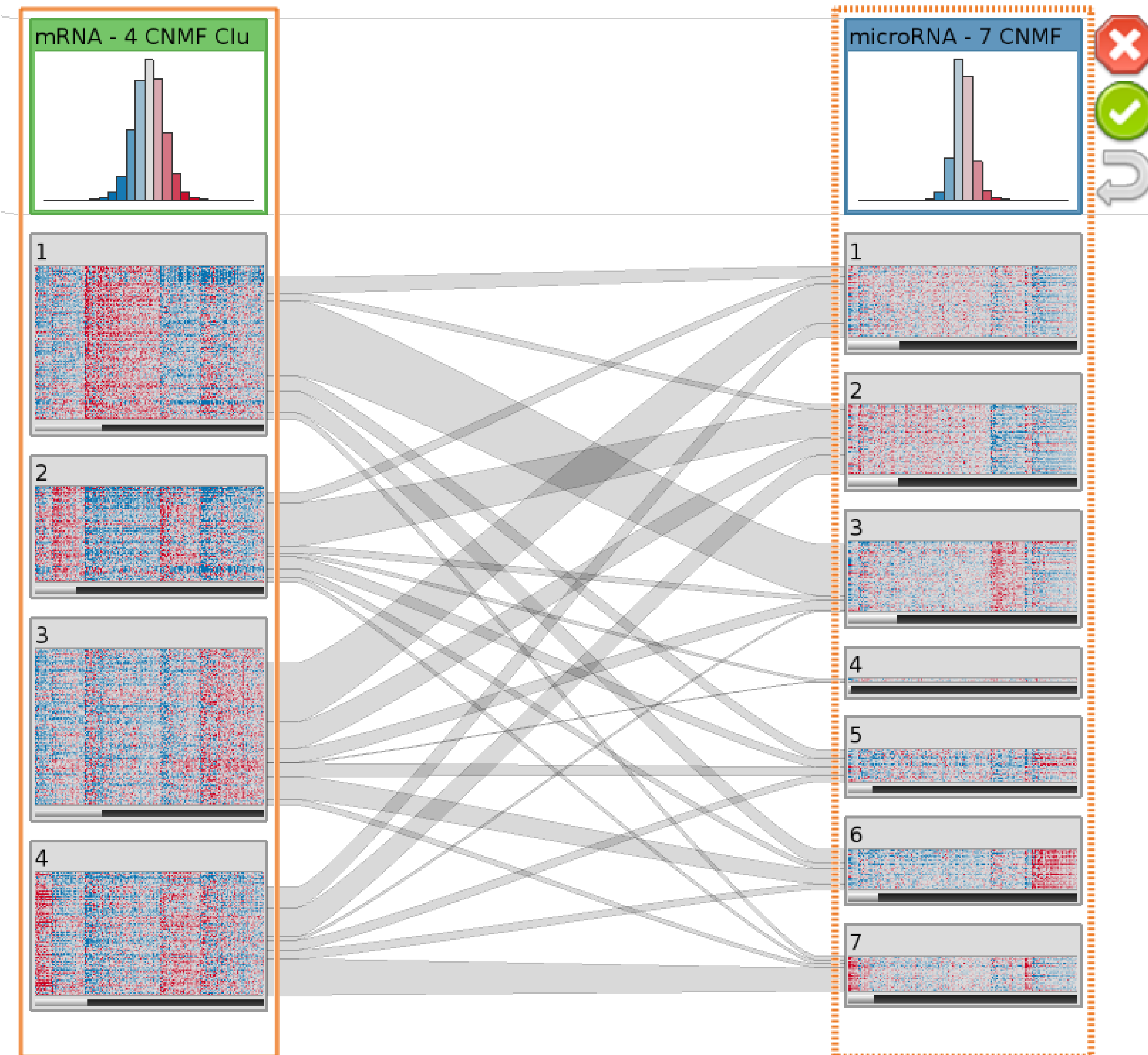


Algorithms

Algorithm Output Visualization

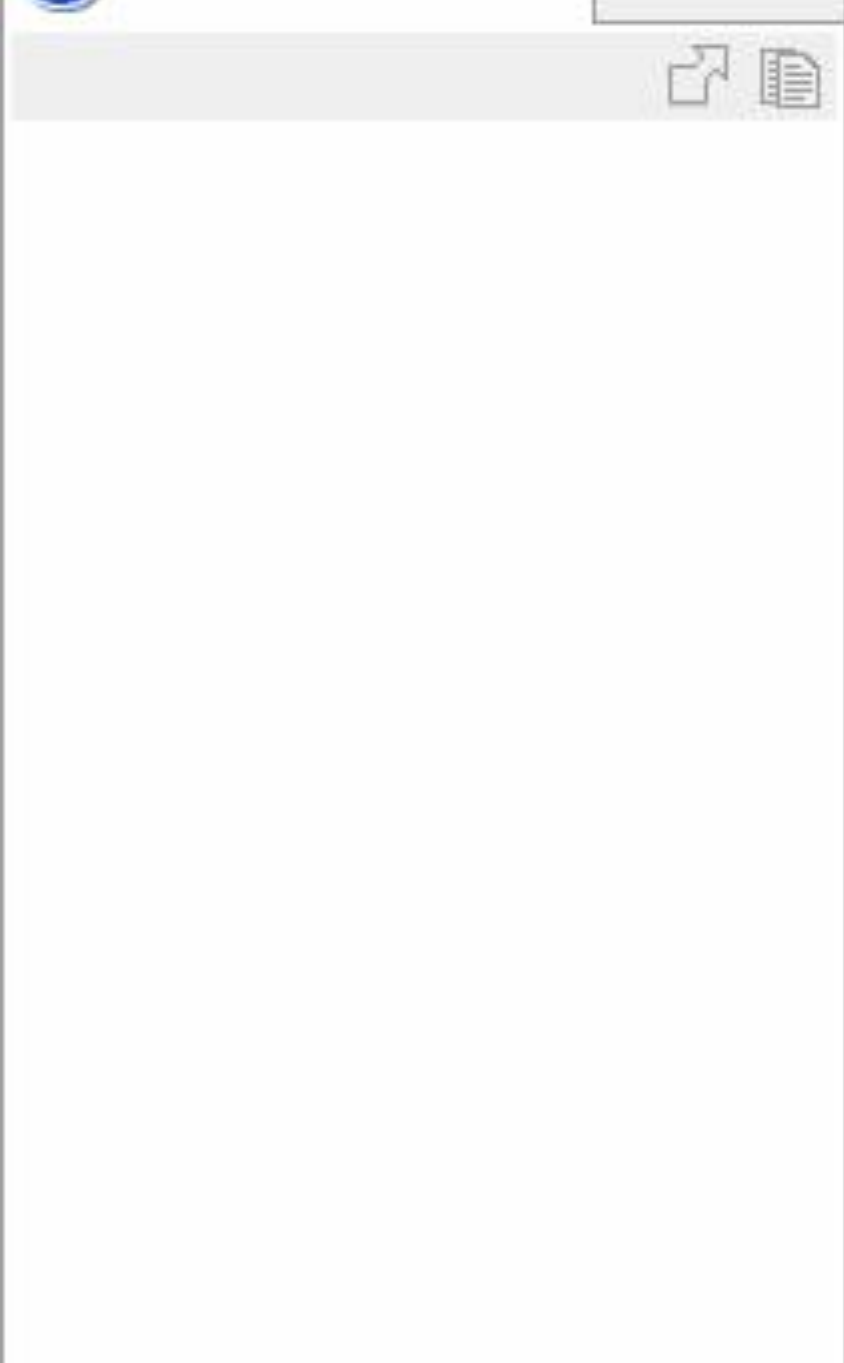
Query column

Result column



Stratification LineUp						
Other LineUp						
Pathway LineUp						
Metrics						
Rank	Stratification	Sim. to mRNA	#Element	#Groups	Distrib	
10.	4 Hierarchical Cluster		529	4		
11.	3 Hierarchical Cluster		529	3		
12.	8 CNMF Clustering		112	8		
13.	6 CNMF Clustering		112	6		
14.	7 CNMF Clustering	0.707 (0.71)	491	7		

Ranked Stratifications



Project: TCGA GBM Package

Dataset: mRNA

Tcga_samples: 528

Genes: 12042

Processing Info

Dataset Stats

Meta Data

Persp.: 3

Tcga_samples: 156 (29.55%)

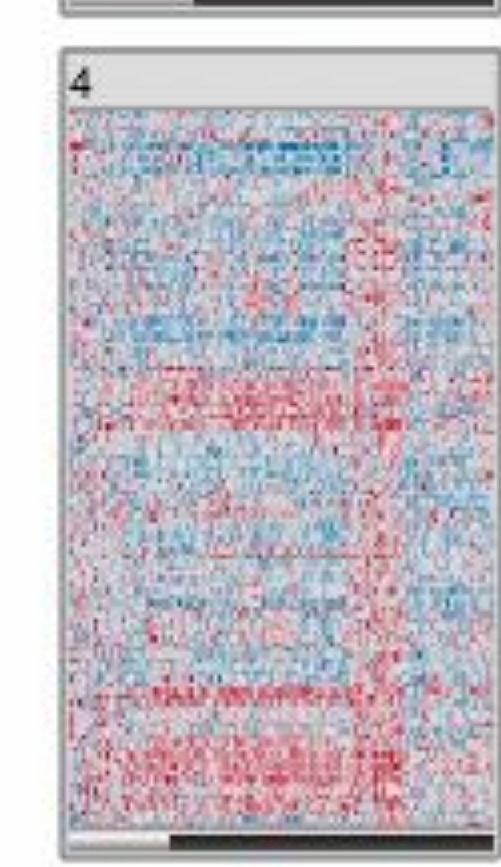
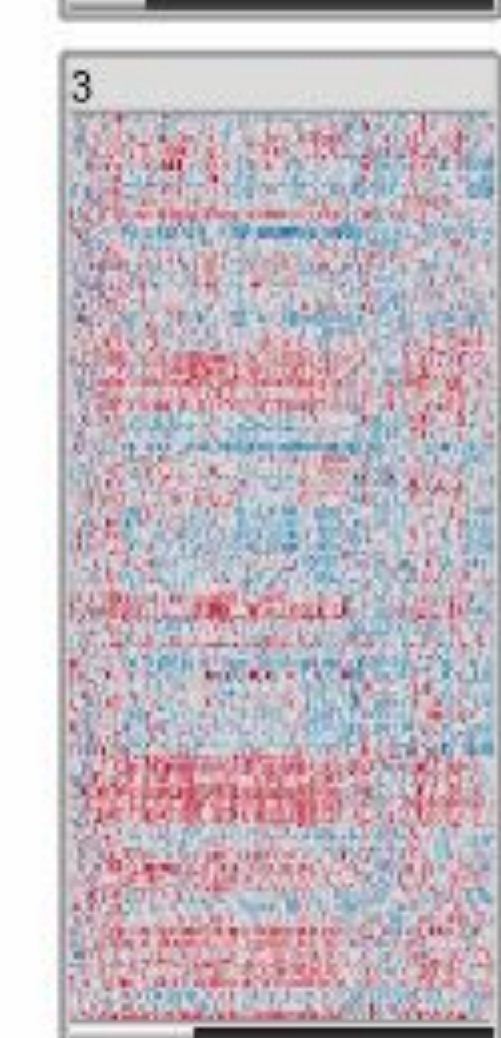
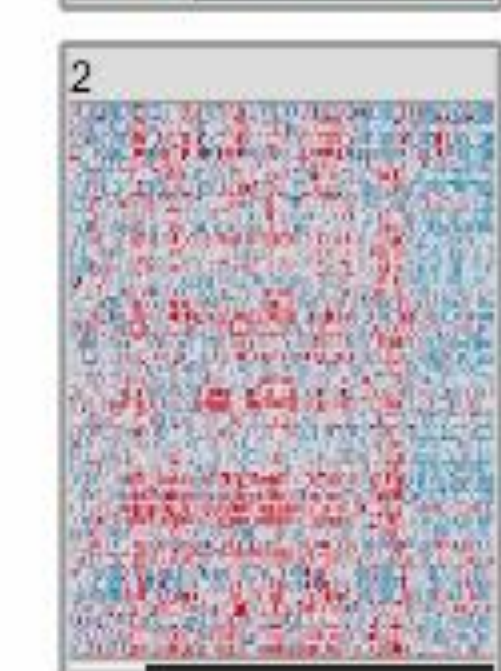
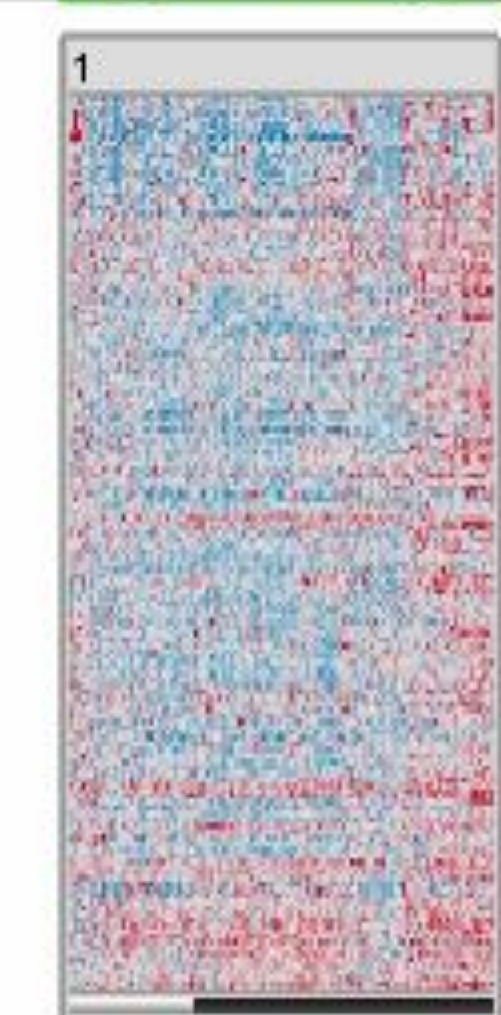
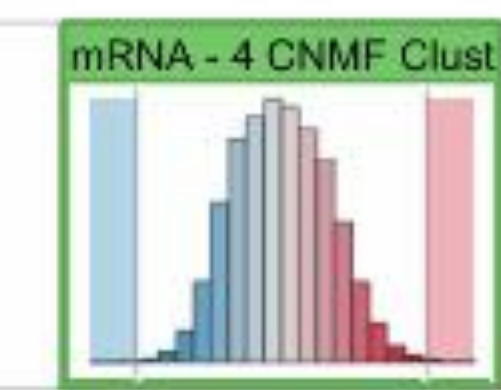
Genes: 1500 (12.46%)

Histogram

A histogram showing the distribution of gene expression values. The x-axis ranges from -4 to 4, and the y-axis represents frequency. The distribution is centered around 0, with a slight skew towards higher values.

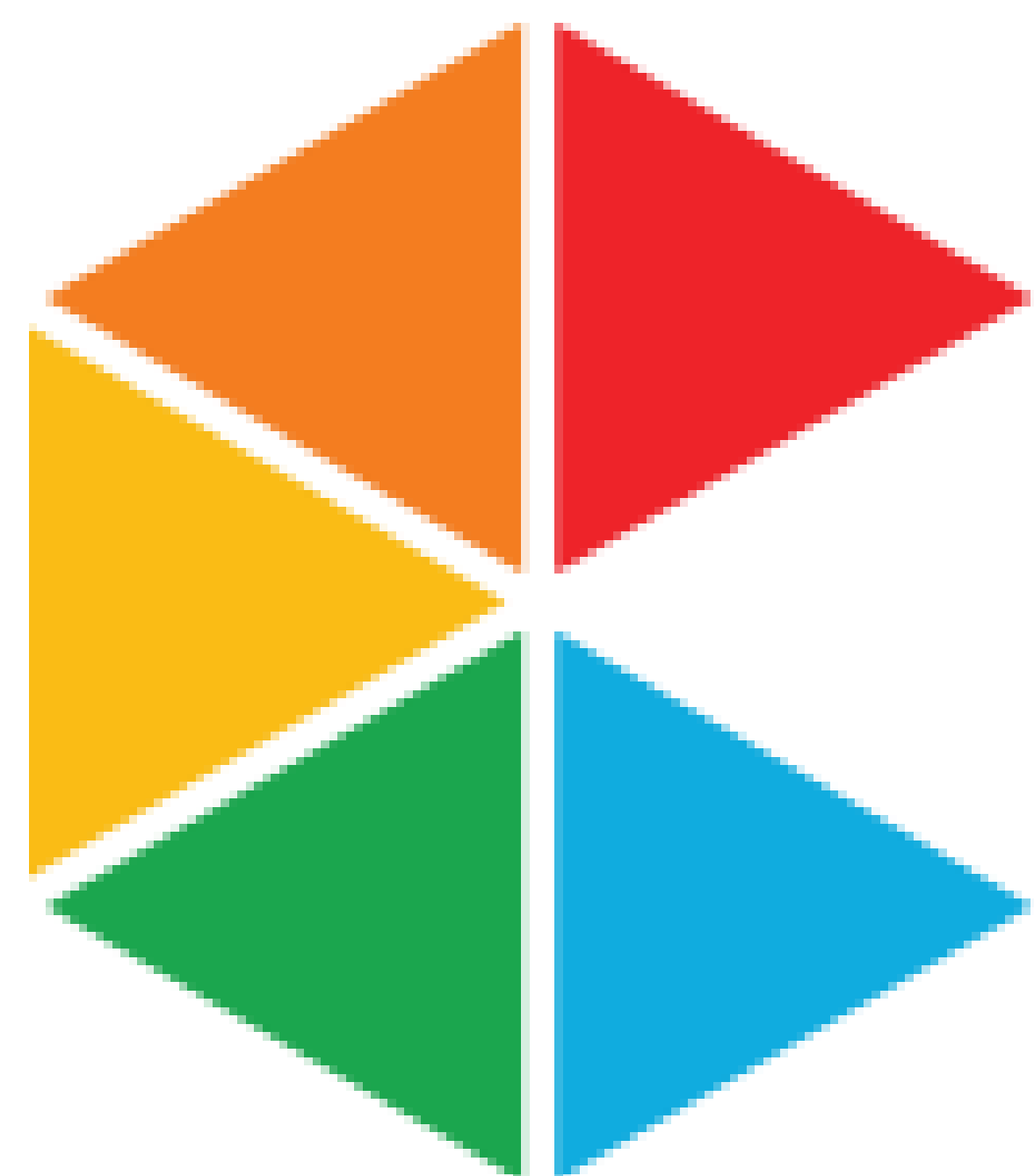
-4 0 4

Colormap



More Information:

<http://stratomex.caleydo.org>



CALEYDO

<http://caleydo.org>

Alexander Lex
@alexander_lex
<http://alexander-lex.com>



HARVARD
School of Engineering
and Applied Sciences

Interactive Visual Data Analysis



Credits:
Marc Streit, Nils Gehlenborg,
Christian Partl, , Samuel Gratzl,
Markus Steinberger, Manuelea
Waldner, Hendrik Strobel, Romain
Vuillemot, Dieter Schmalstieg, Denis
Kalkofen, Mark Borowsky, Anne Mai
Wasserman, Hanspeter Pfister

