

Alexander Lex

@alexander_lex

<http://alexander-lex.net>



THE
UNIVERSITY
OF UTAH

Literate Visualization: Making Visual Analysis

Sessions Reproducible and Reusable



visualization
design lab



www.sci.utah.edu

visualization

**The purpose of computing is insight,
not numbers.**

pictures

[Card, Mackinlay, Shneiderman]

[Richard Wesley Hamming]

Banana *M. acuminata*

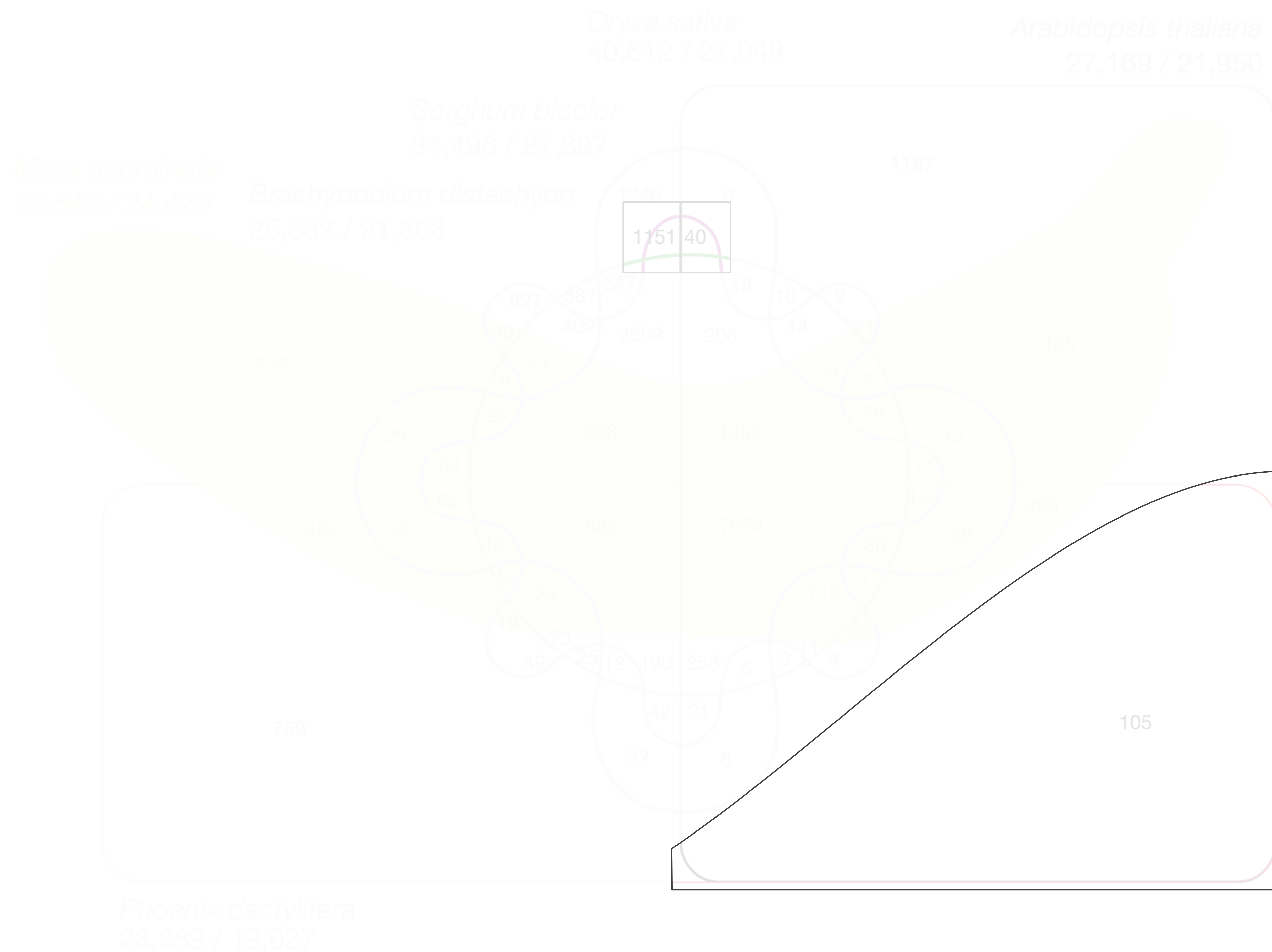
Date *P. dactylifera*

Cress *Arabidopsis thaliana*

Rice *Oryza sativa*

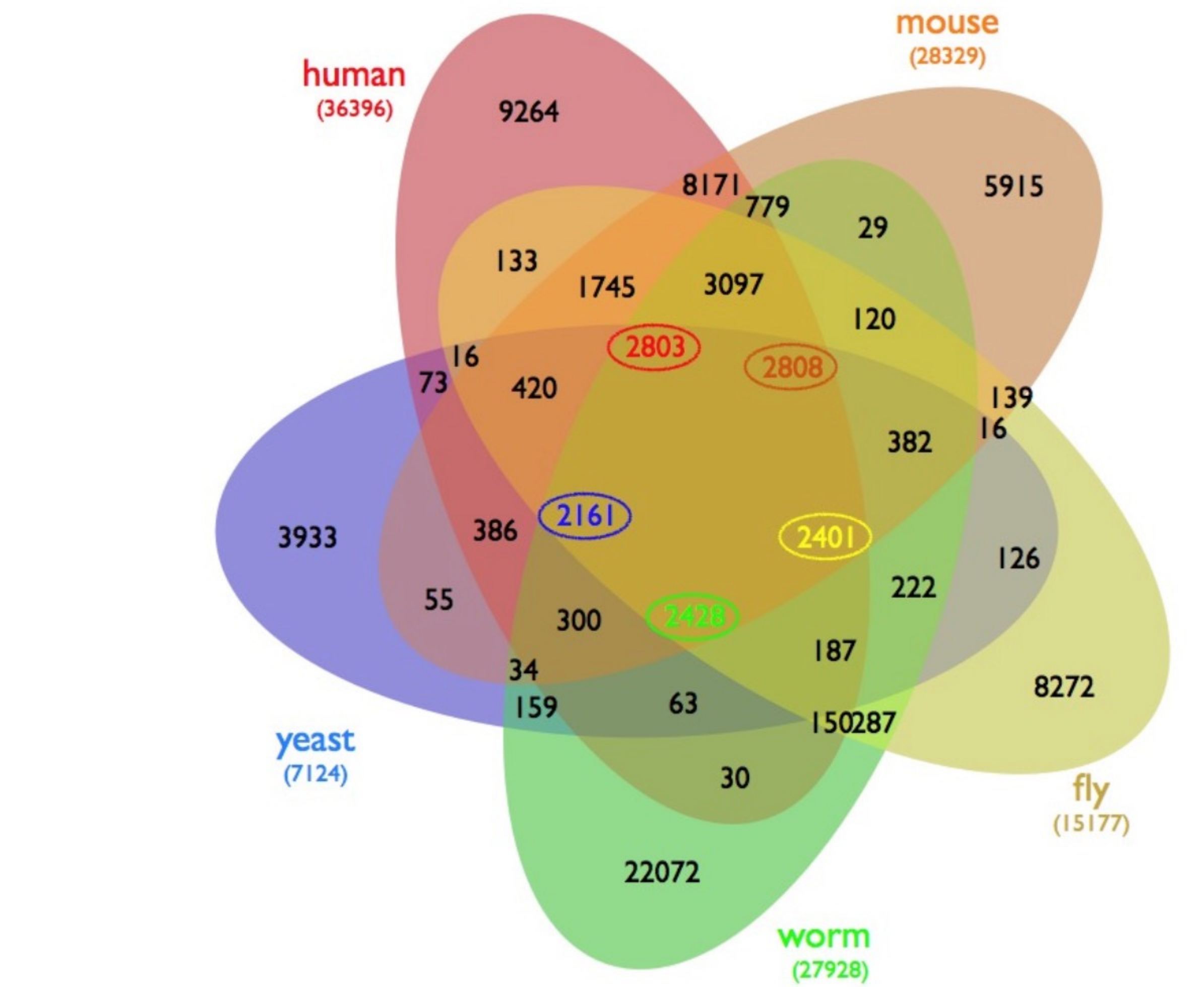
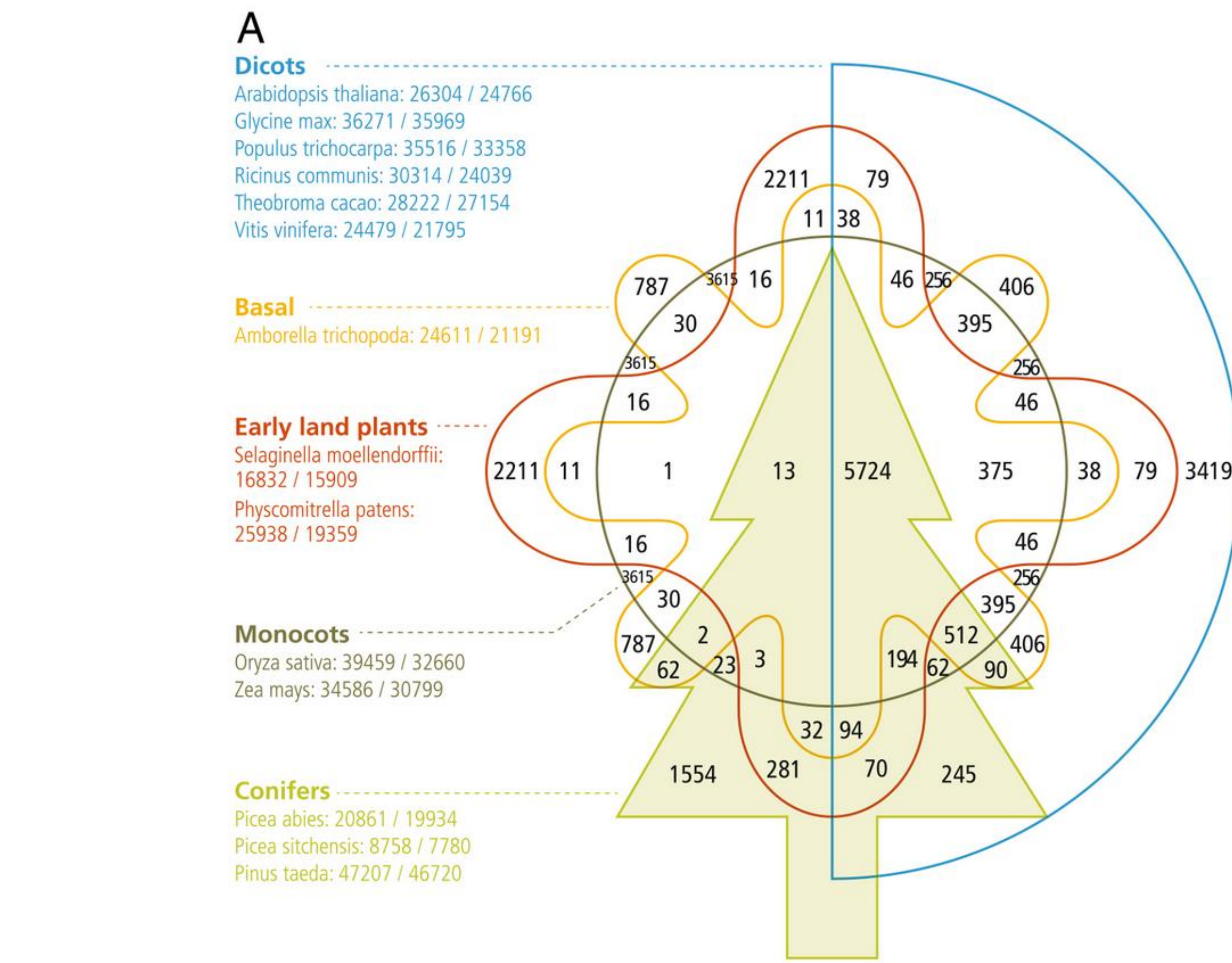
Sorghum *Sorghum bicolor*

Brome *Brachypodium distachyon*

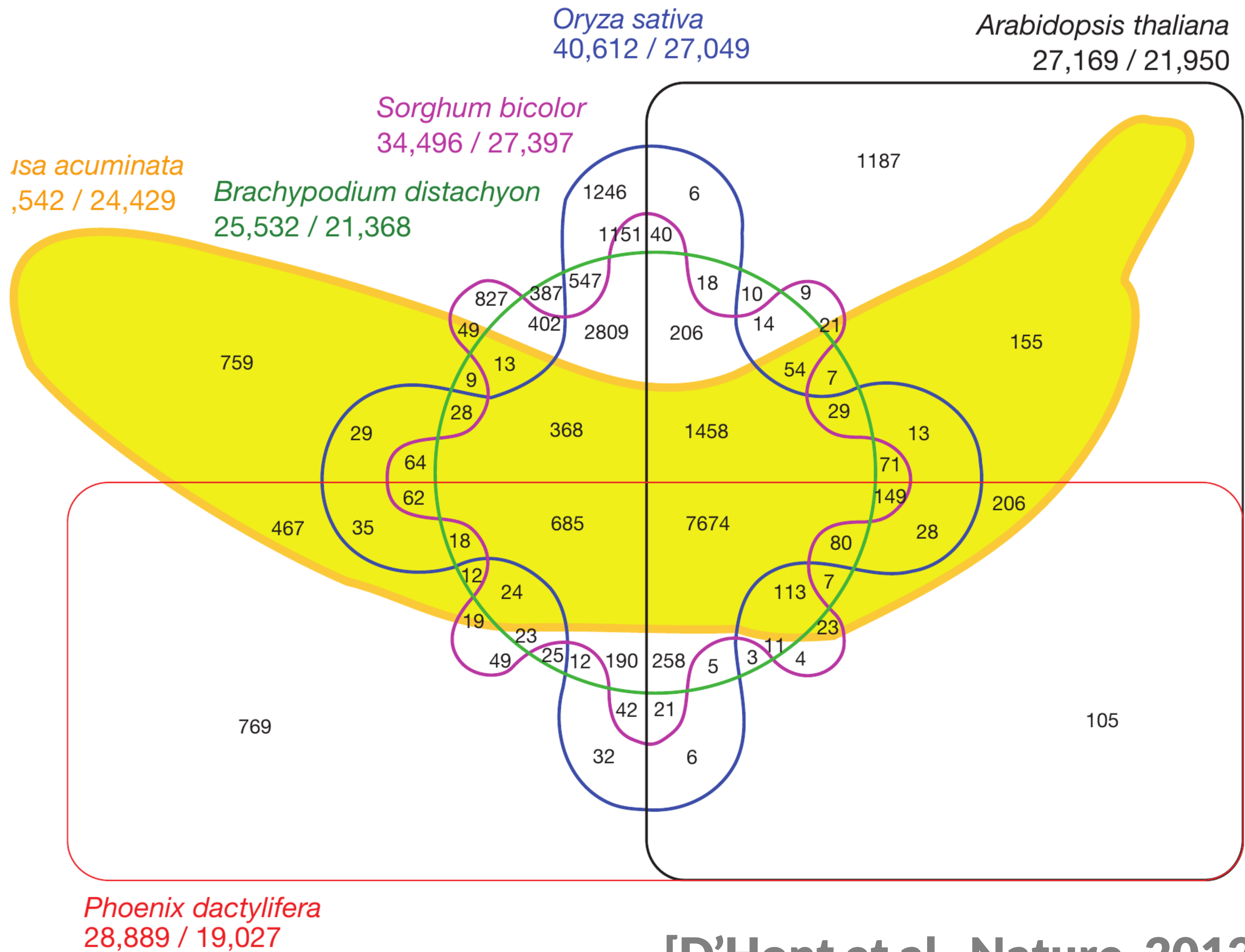


[D'Hont et al., Nature, 2012]

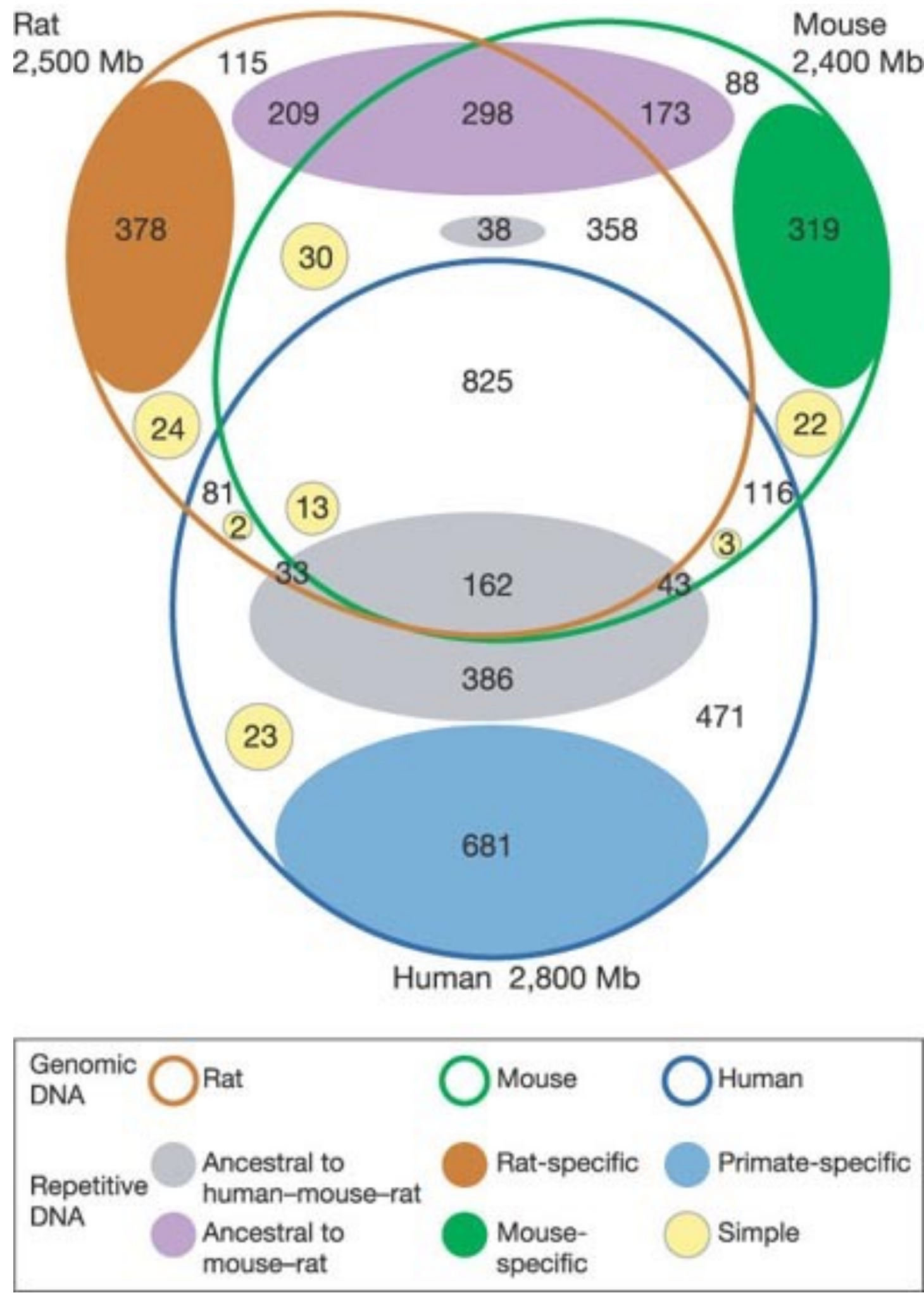
[Neale et al., BMC Genome Biology, 2014]



[Wiles et al., BMC Systems Biology]

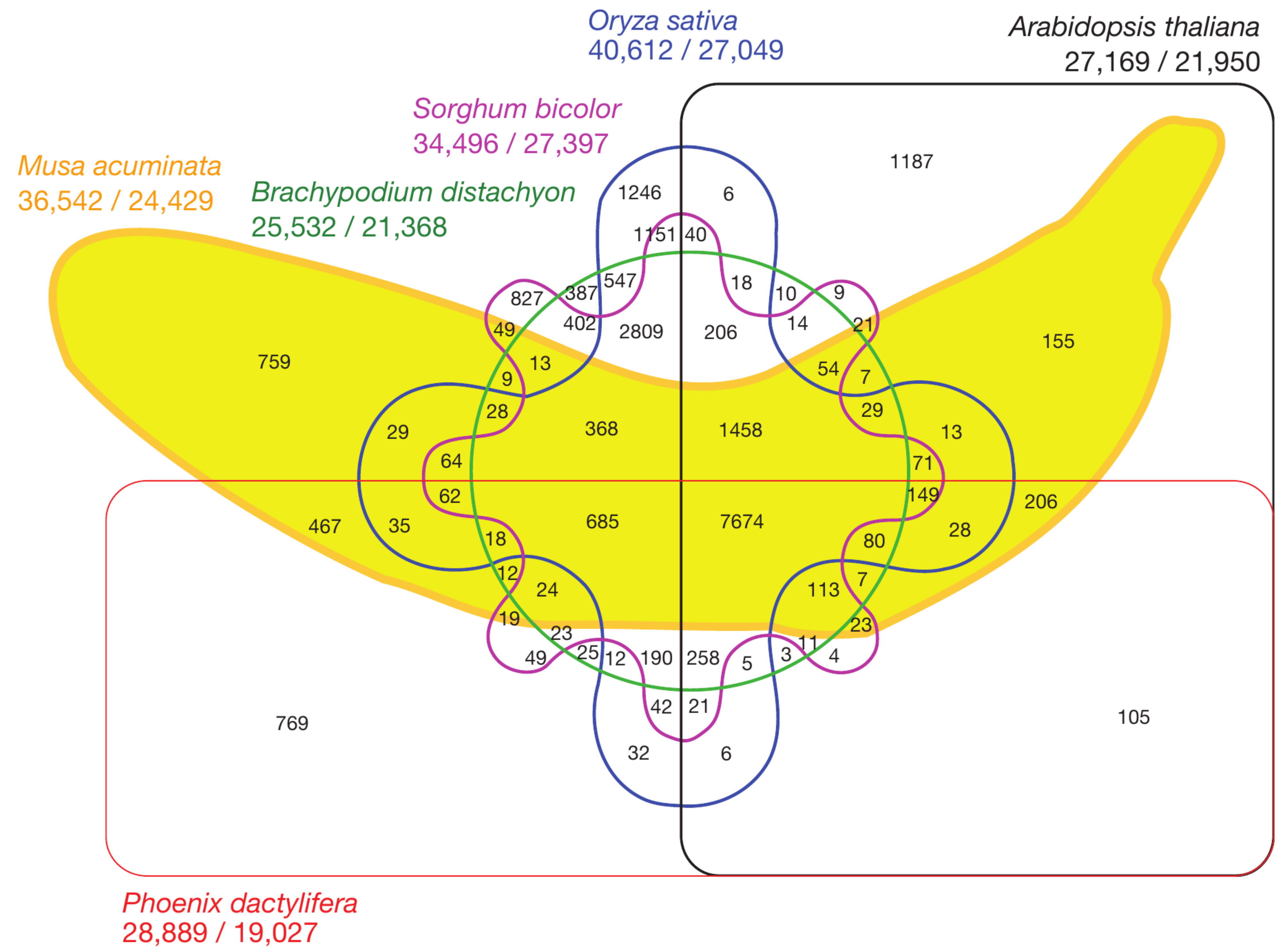


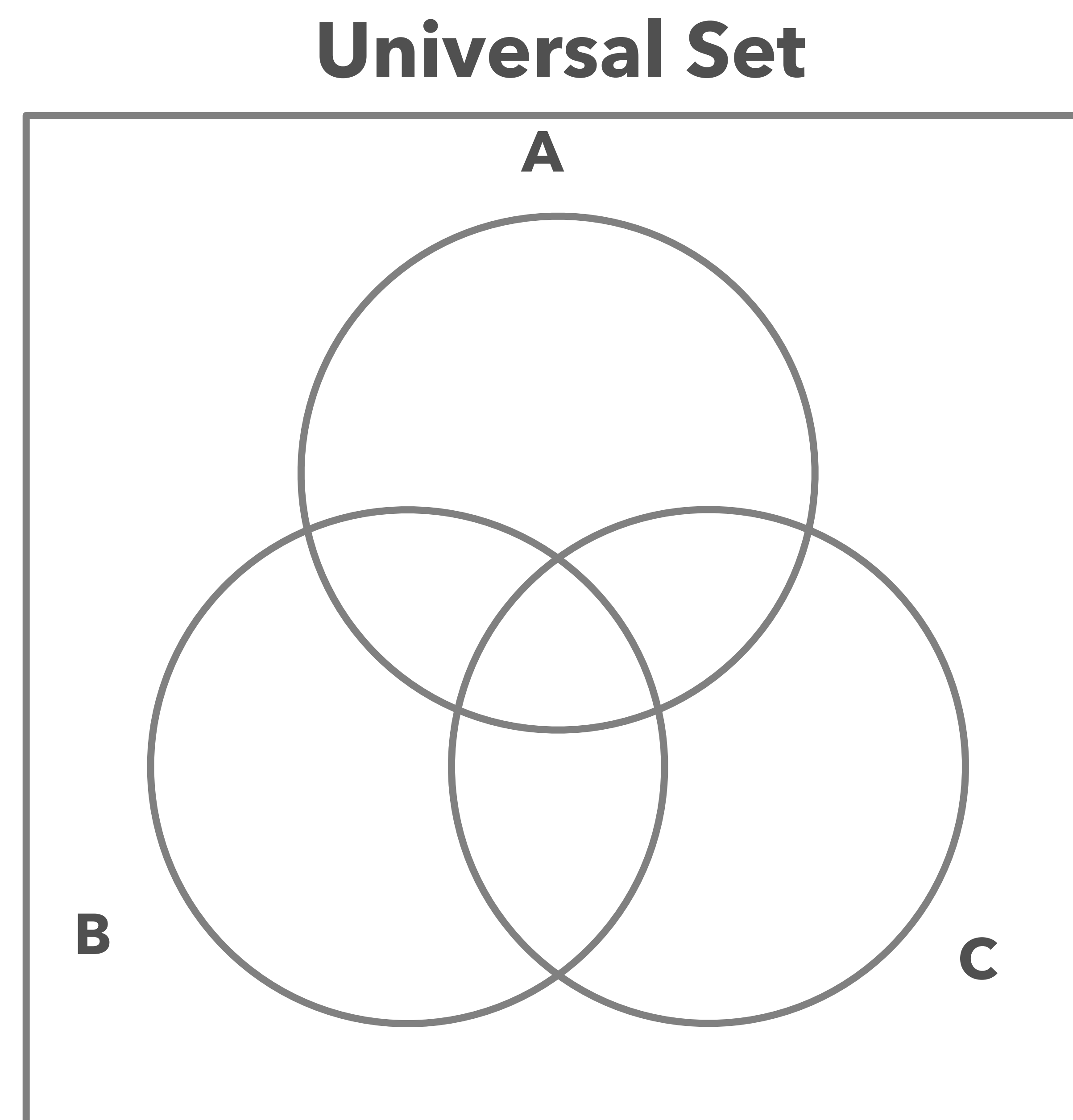
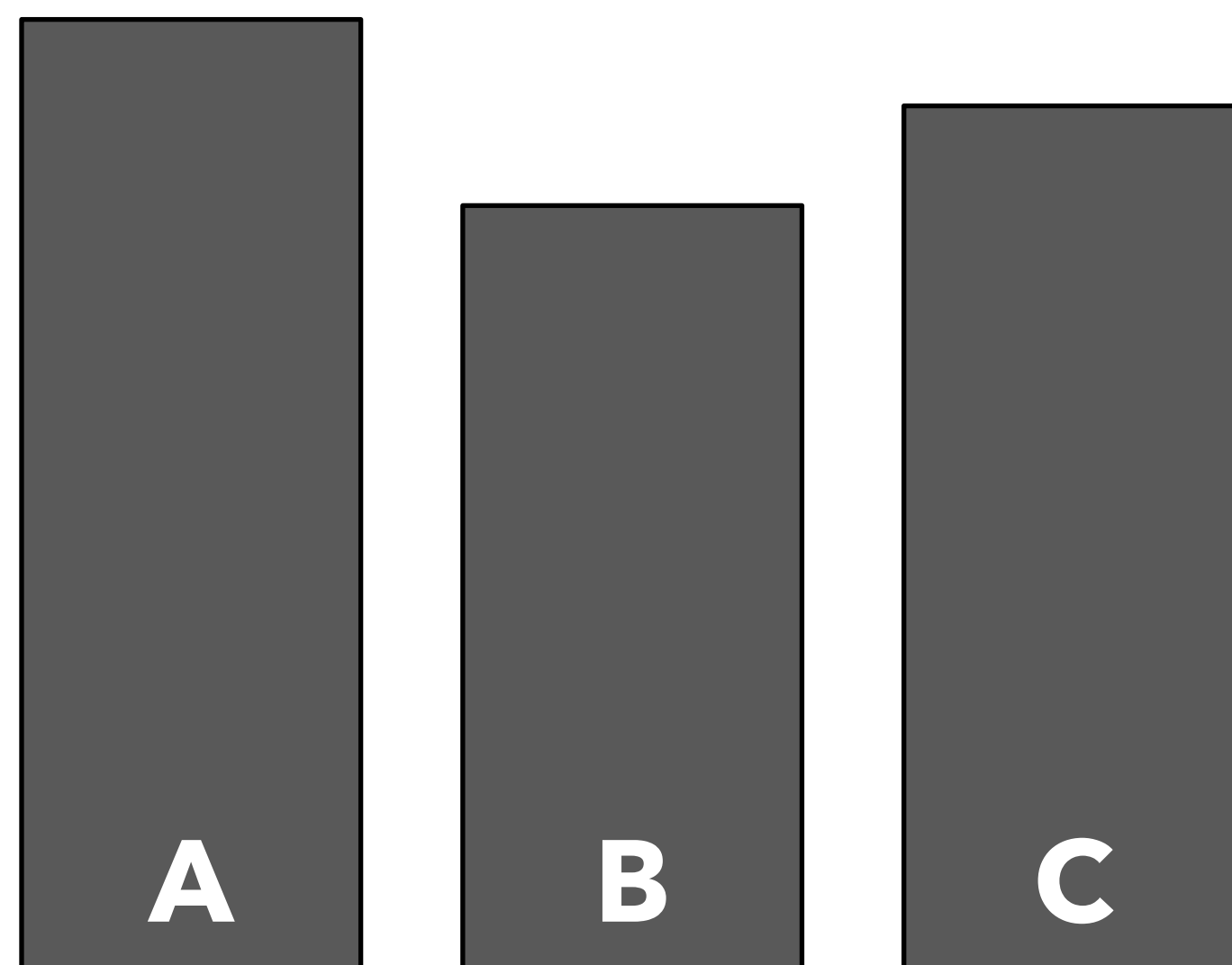
[D'Hont et al., Nature, 2012]

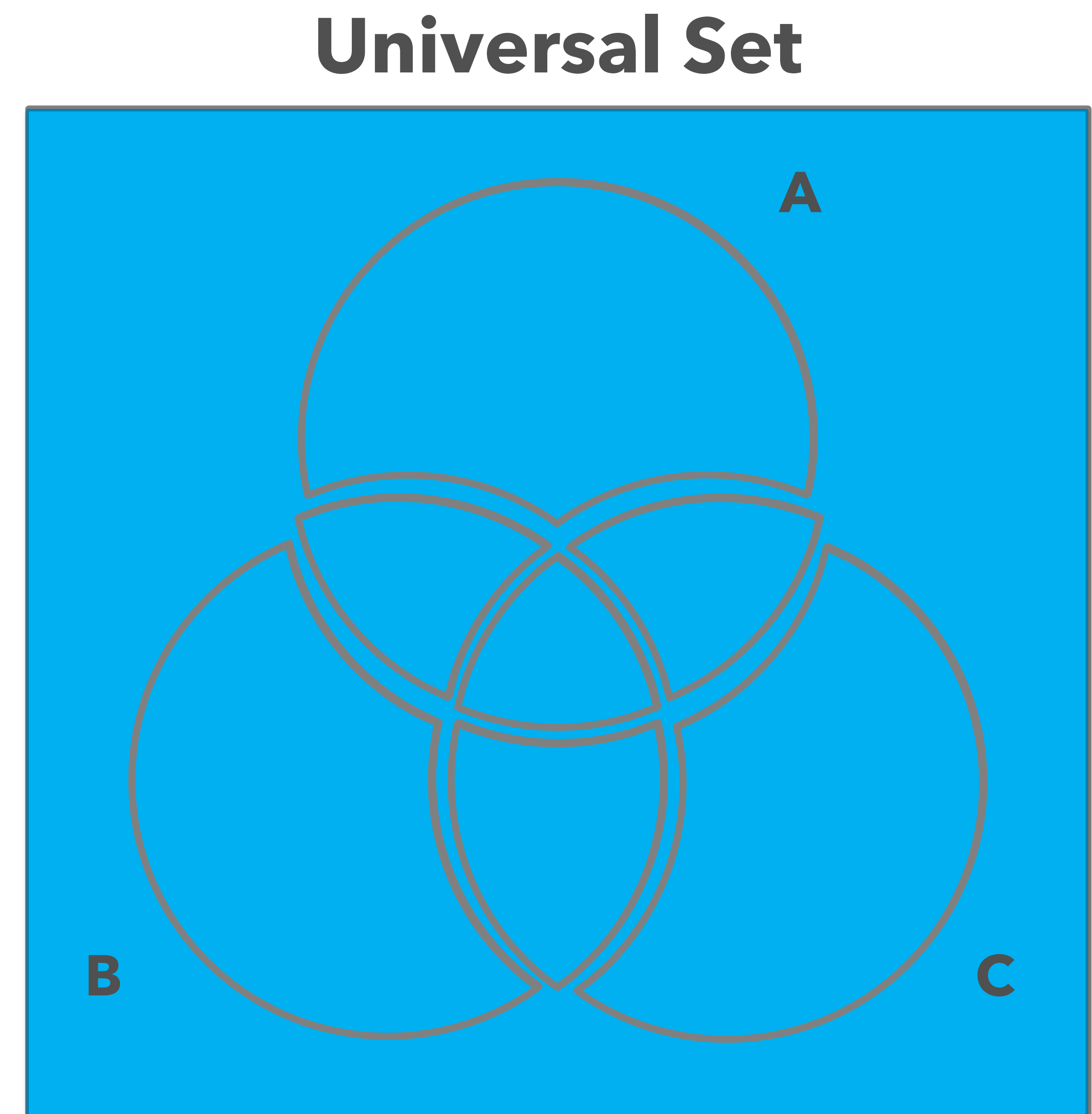
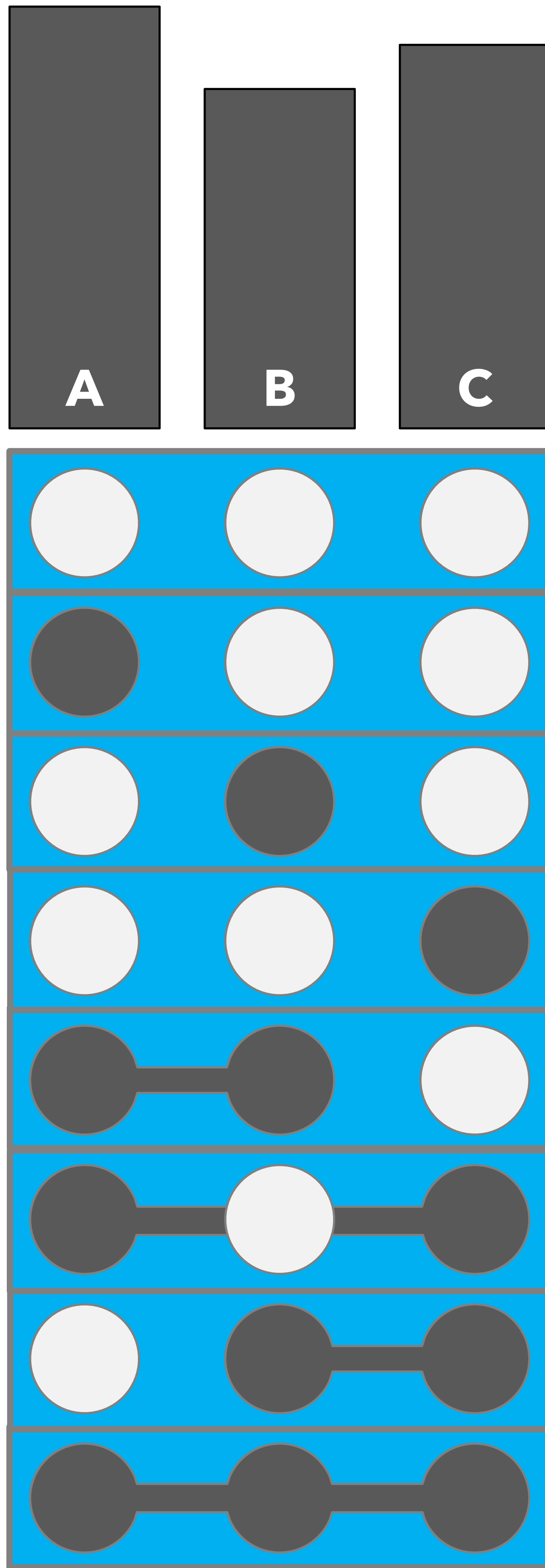


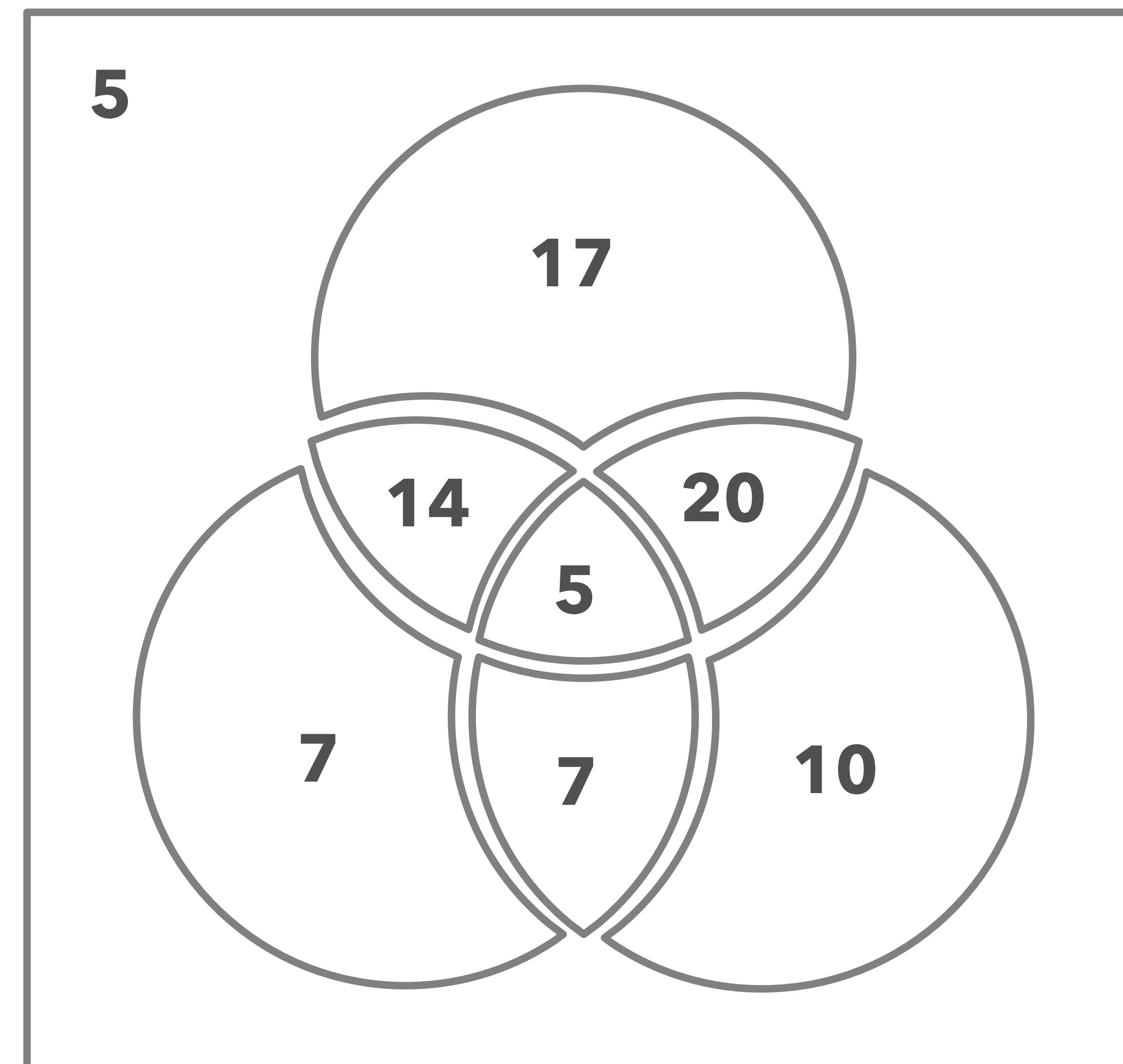
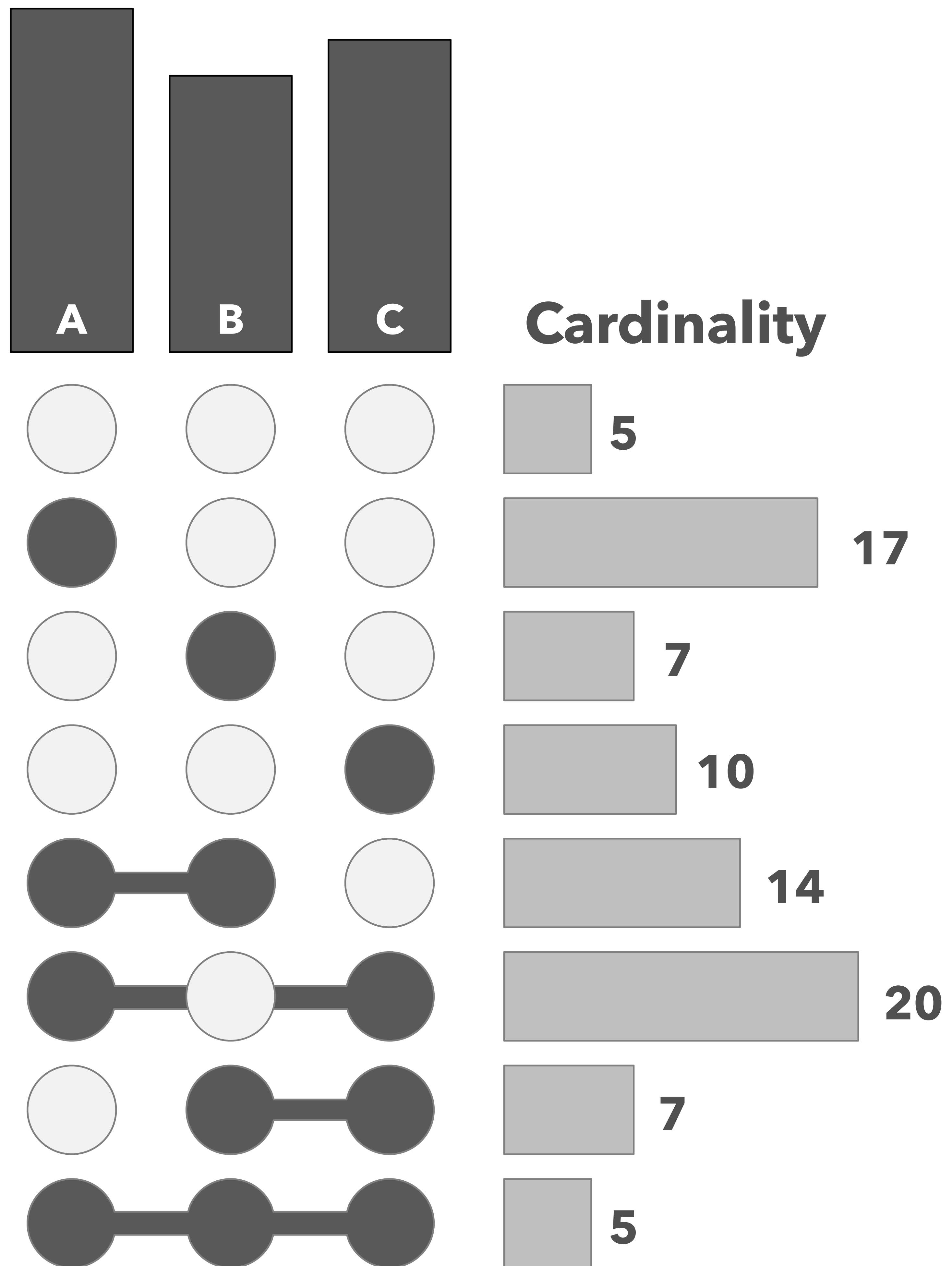
[Gibbs et al., Nature, 2004]

SO CAN WE DO
BETTER?

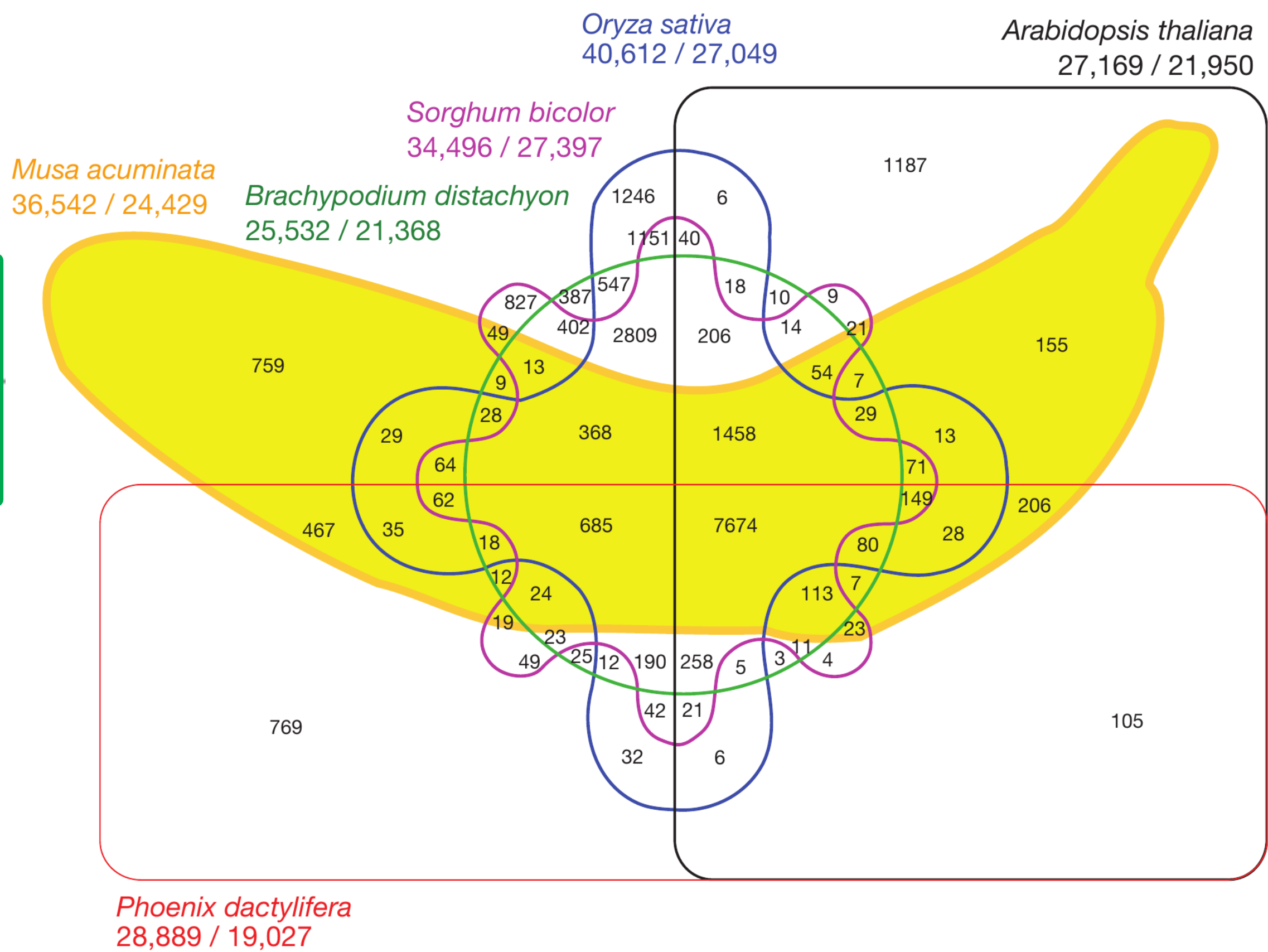
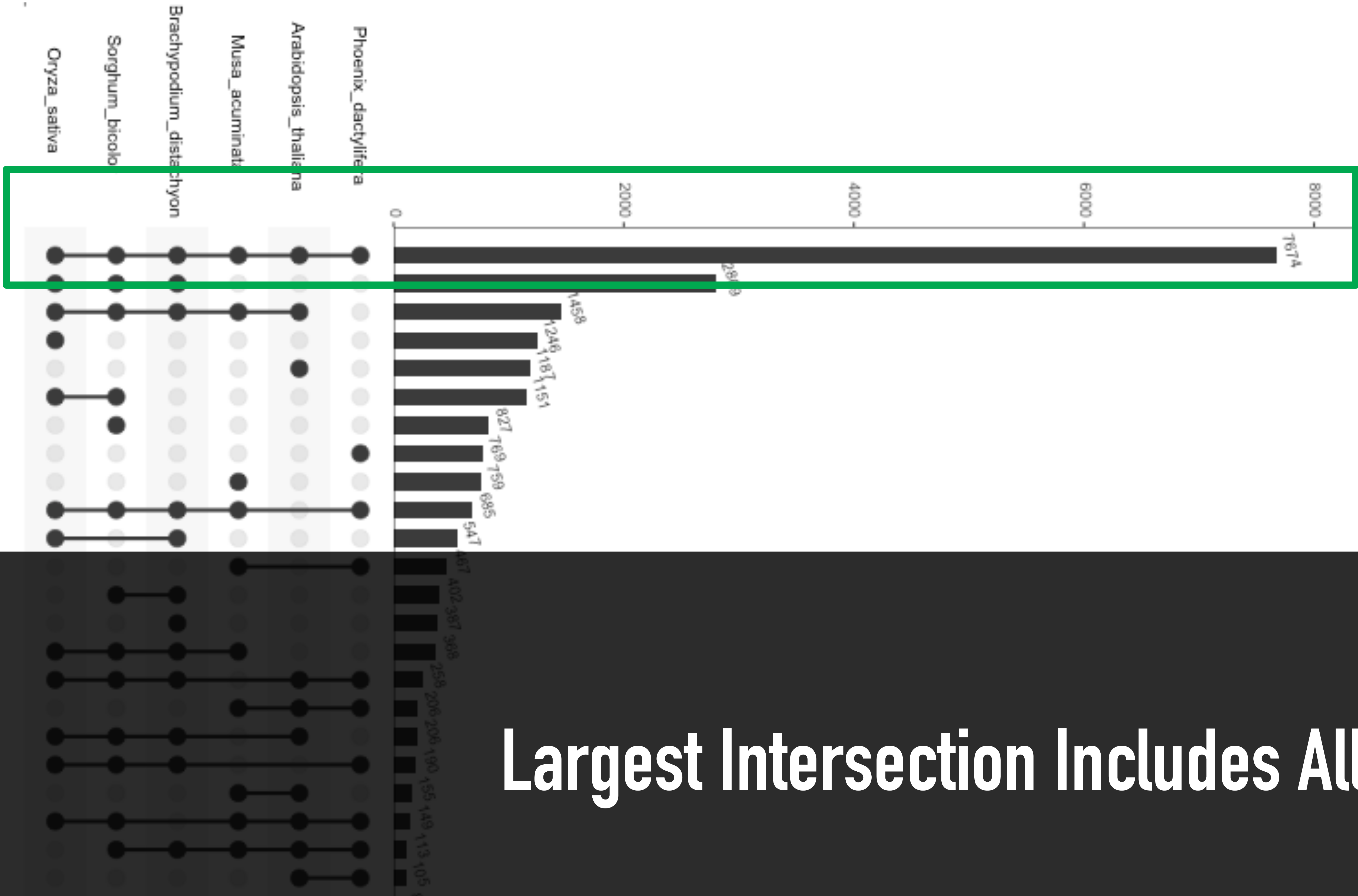






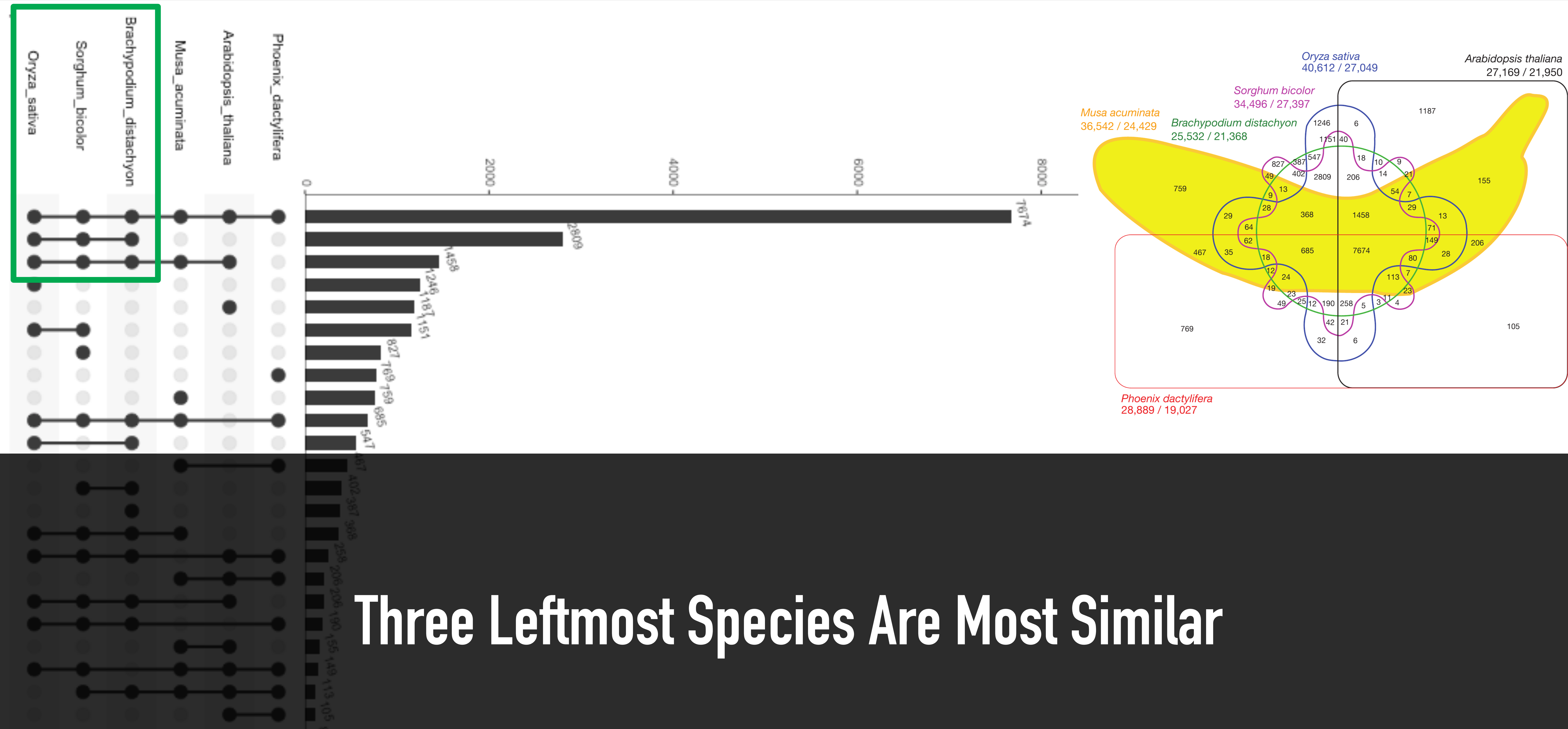


THE BANANA CHART REDESIGNED: UPSET



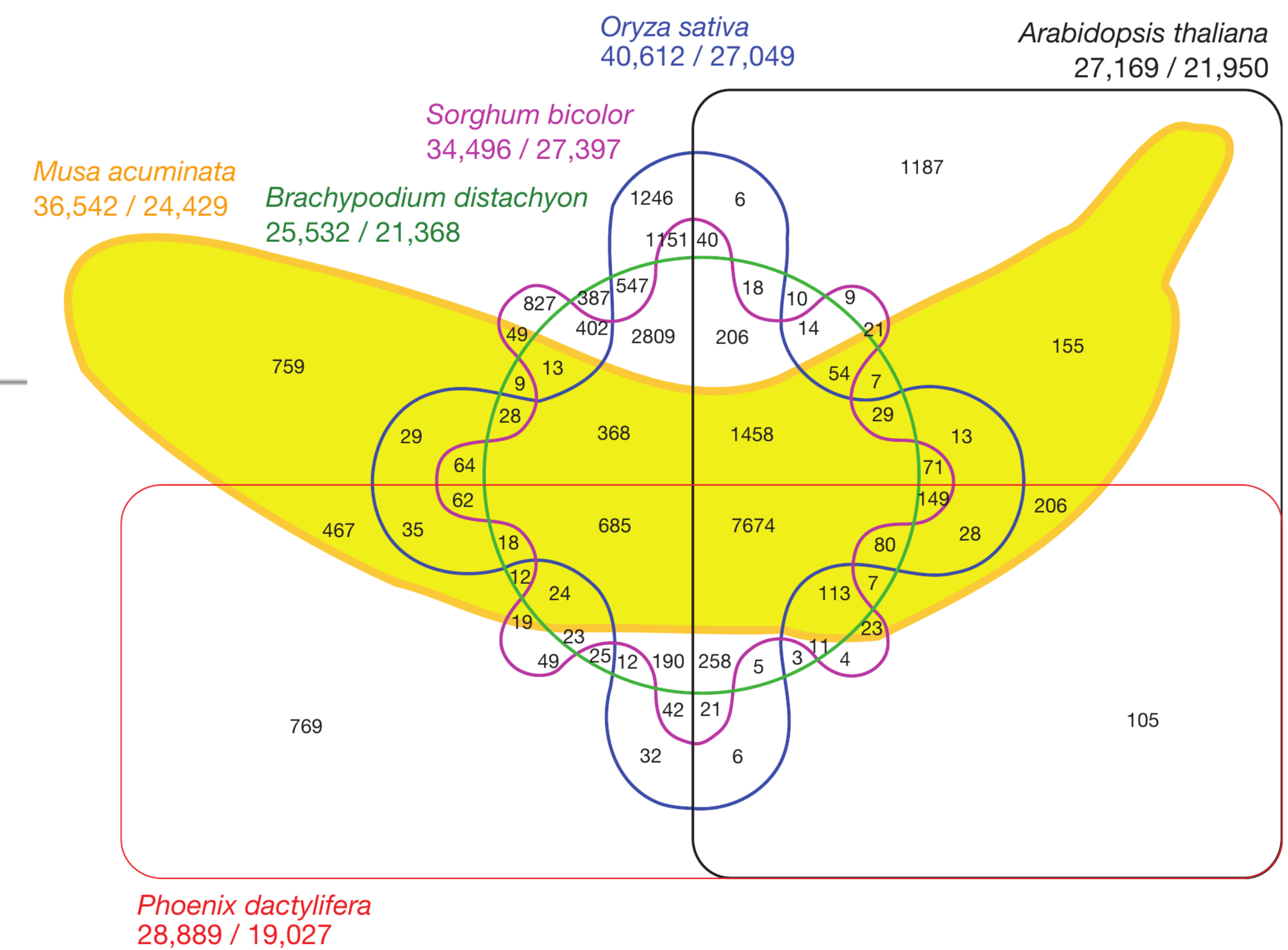
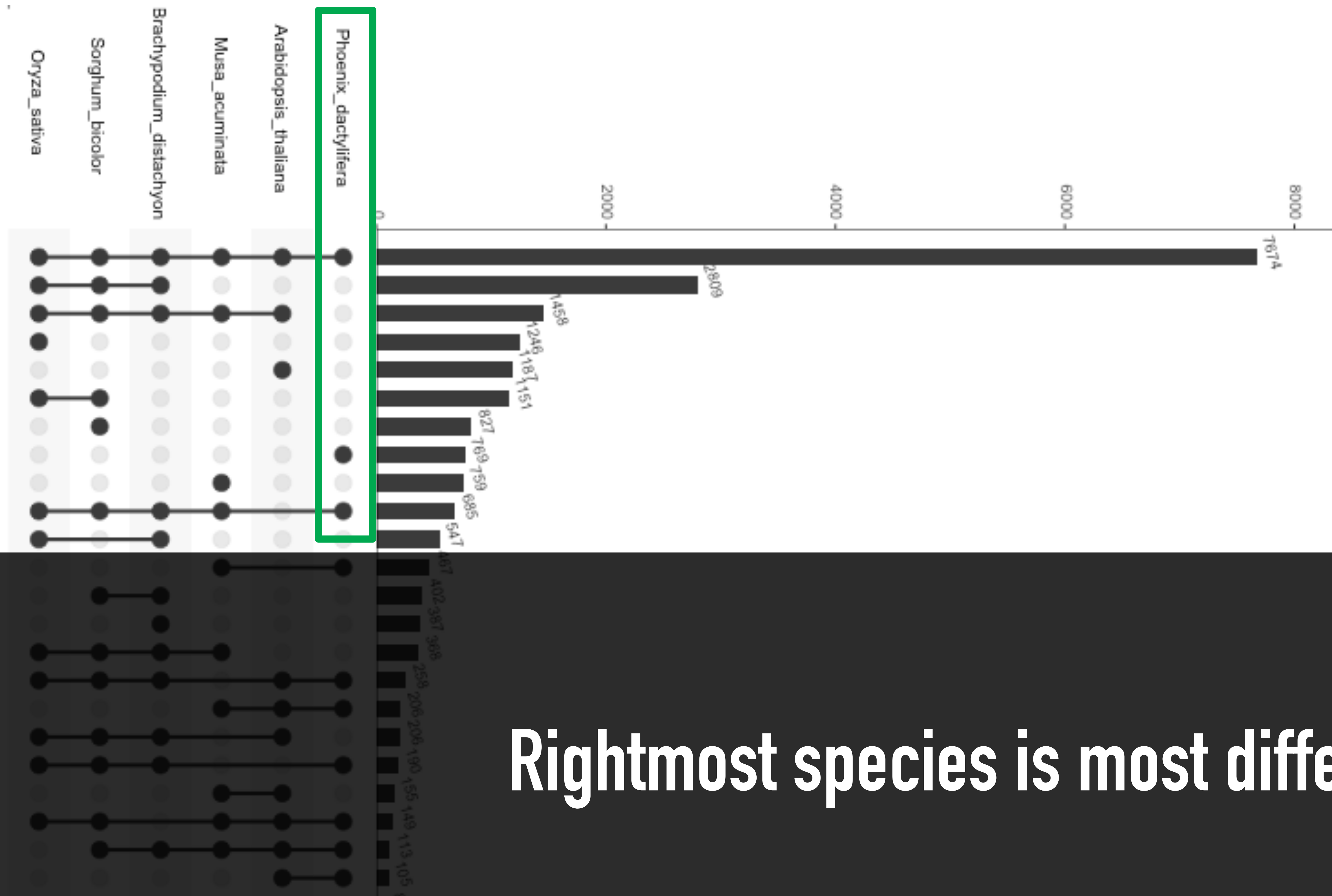
Largest Intersection Includes All Sets

THE BANANA CHART REDESIGNED: UPSET



Three Leftmost Species Are Most Similar

THE BANANA CHART REDESIGNED: UPSET



Rightmost species is most different

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..

Set Selection

0 - 11

Batch Add Sets

Sort Sets

Drama Thriller Romance Horror SciFi War Musical Mystery Fantasy Western Noir

0 1000 2000 3000 3883

Cardinality

Deviation

Release Date

Average Rating

Times Watched

+ Query

Adventure Children Comedy Crime Documentary Action

1777

1003

283

131

123

101

85

78

66

62

38

38

24

13

13

13

9

7

5

4

4

3

2

1

1777

1003

283

131

123

101

85

78

66

62

38

38

24

13

13

13

9

7

5

4

4

3

2

1

1777

1003

283

131

123

101

85

78

66

62

38

38

24

13

13

13

9

7

5

4

4

3

2

1

1777

1003

283

131

123

101

85

78

66

62

38

38

24

13

13

13

9

7

5

4

4

3

2

1

1777

1003

283

131

123

101

85

78

66

62

38

38

24

13

13

13

9

7

5

4

4

3

2

1

1777

1003

283

131

123

101

85

78

66

62

38

38

24

13

13

13

9

7

5

4

4

3

2

1

1777

1003

283

131

123

101

85

78

66

62

38

38

24

13

13

13

9

7

5

4

4

3

2

1

1777

1003

283

131

123

101

85

78

66

62

38

38

24

13

13

13

9

7

5

4

4

3

2

1

1777

1003

283

131

123

101

85

78

66

62

38

38

24

13

13

13

9

7

5

4

4

3

2

1

1777

1003

283

131

123

101

85

78

66

62

38

38

24

13

13

13

9

7

5

4

4

3

2

1

1777

1003

283

131

123

101

85

78

66

62

38

38

24

13

13

13

9

7

5

4

4

3

2

1

1777

1003

283

131

123

101

85

78

66

62

38

38

24

13

13

13

9

7

5

4

4

3

2

1

1777

1003

283

131

123

101

85

78

66

62

38

38

24

13

13

13

9

7

5

4

4

3

2

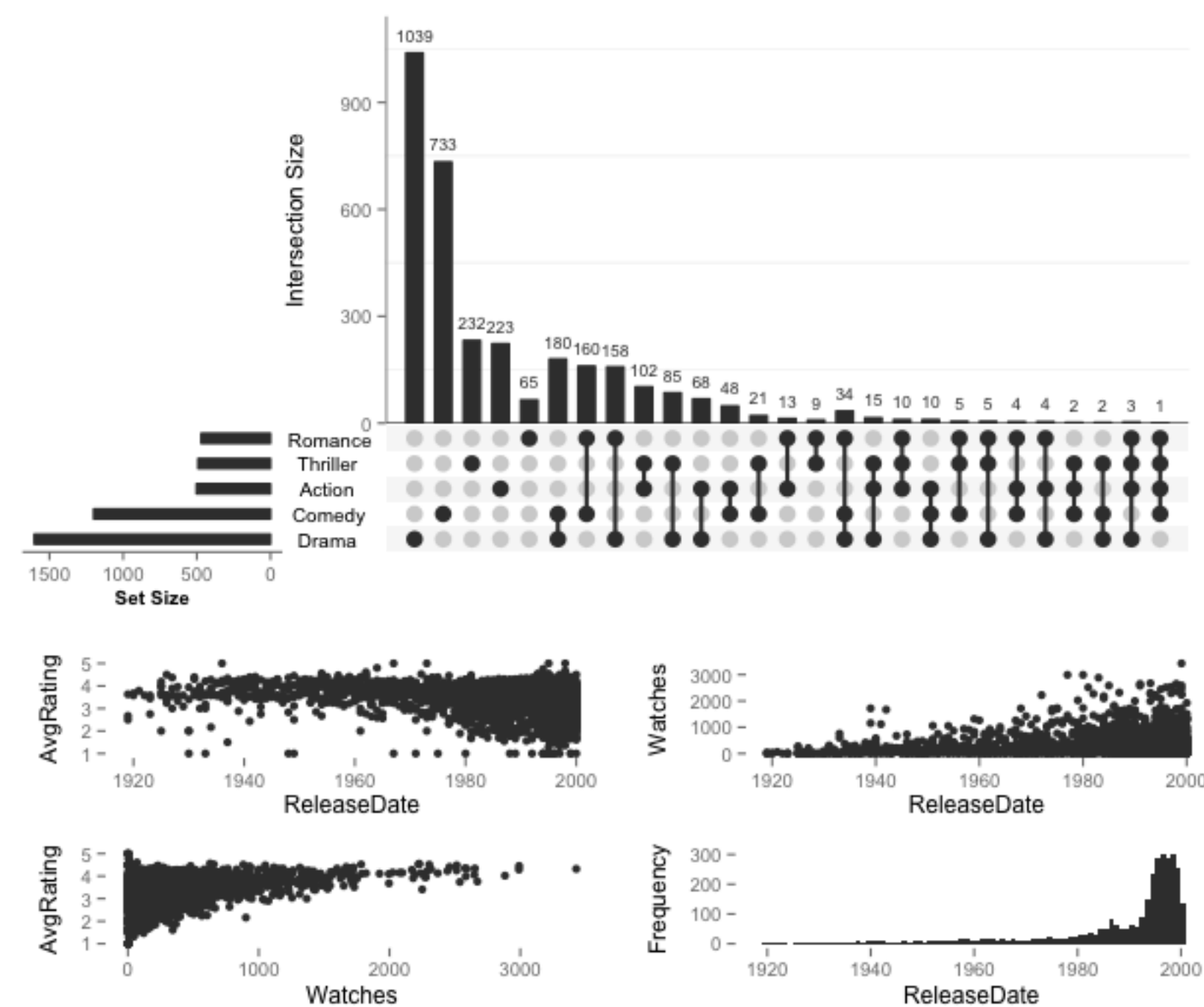
1

UPSET

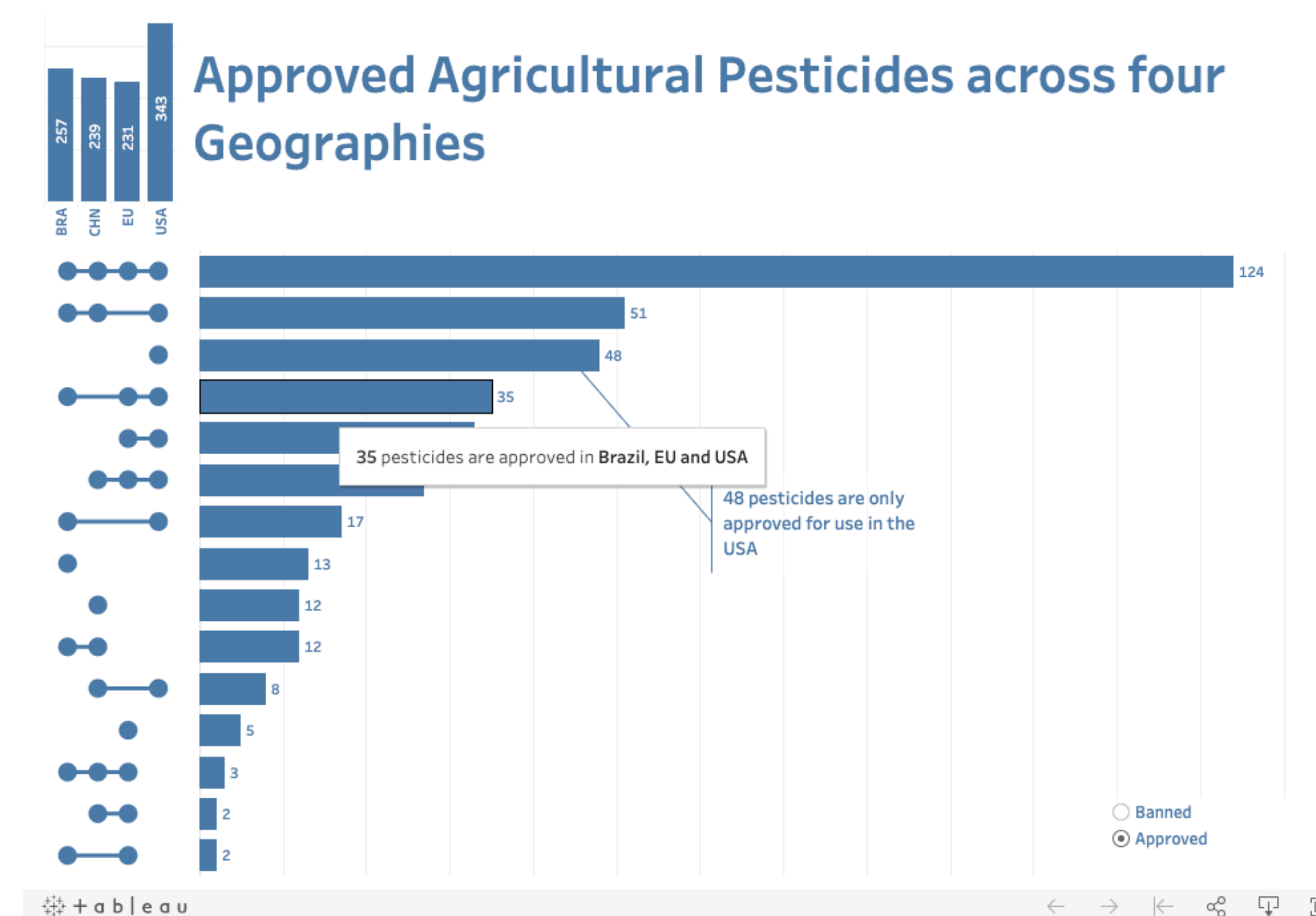
The canonical way to show set data with > 3 sets

Second-most cited VIS paper of the last decade

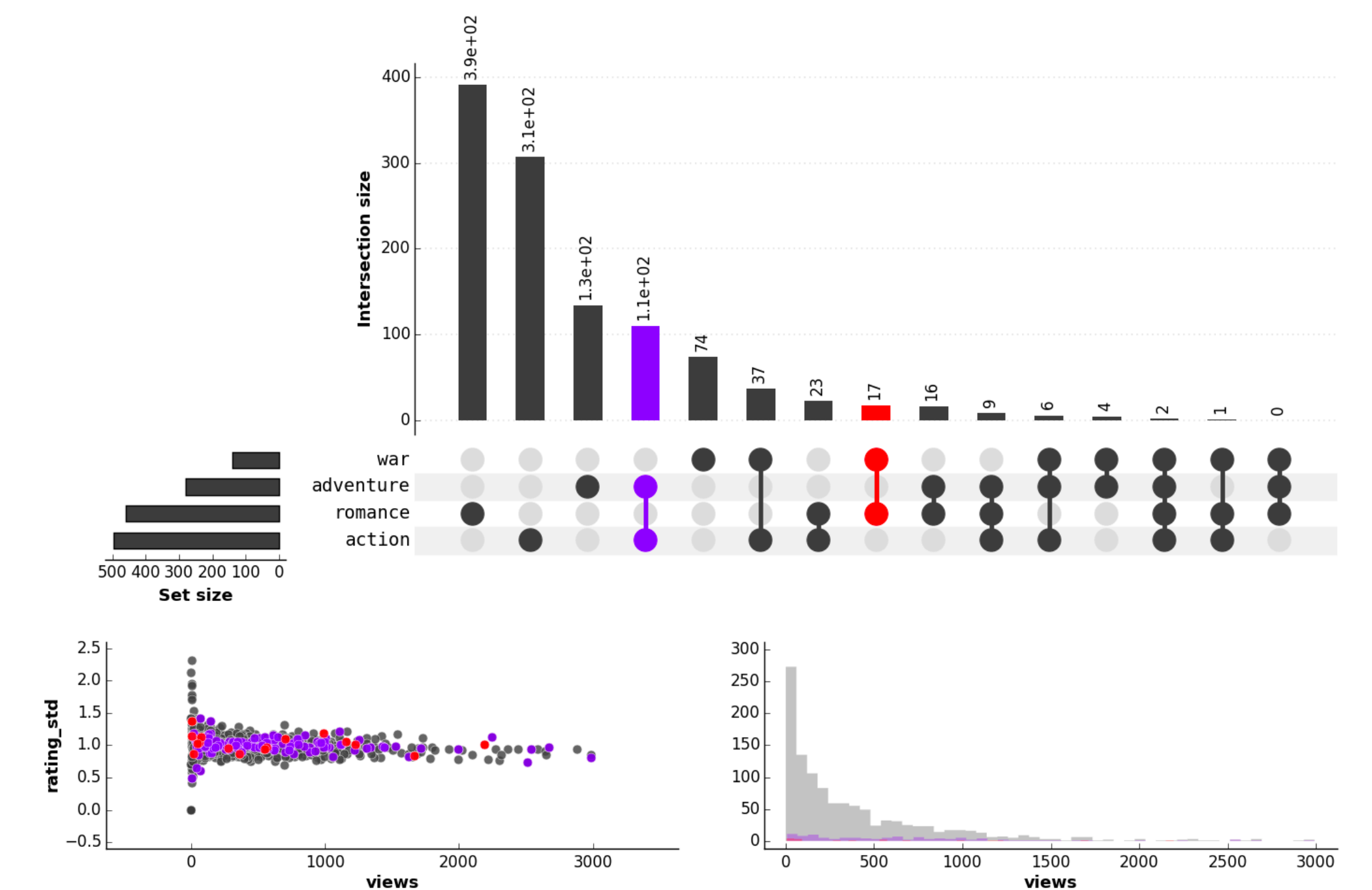
Multiple implementations in various languages



R



Tableau



Python

RESEARCH AREAS

TECHNICAL CONTRIBUTIONS

**Novel Visualization
Techniques**

**Visualization Process
Innovations**

**Data Wrangling
Methods**

DOMAIN DRIVEN TECHNIQUES

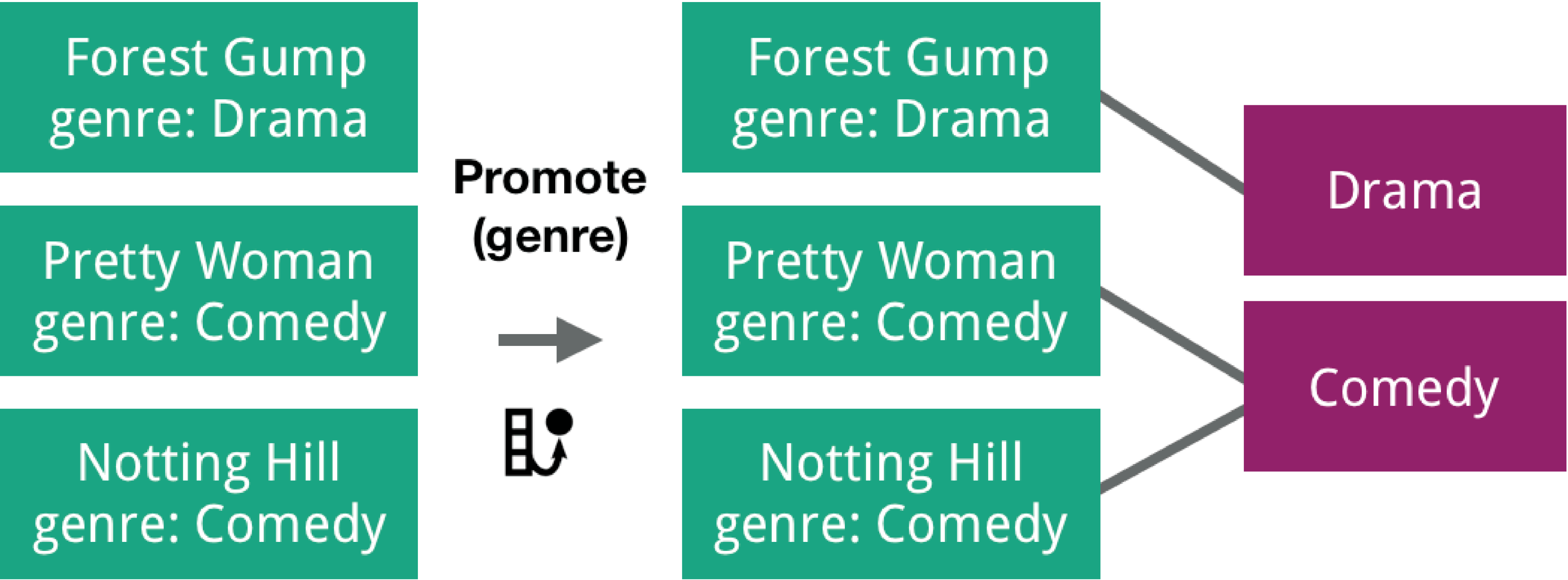
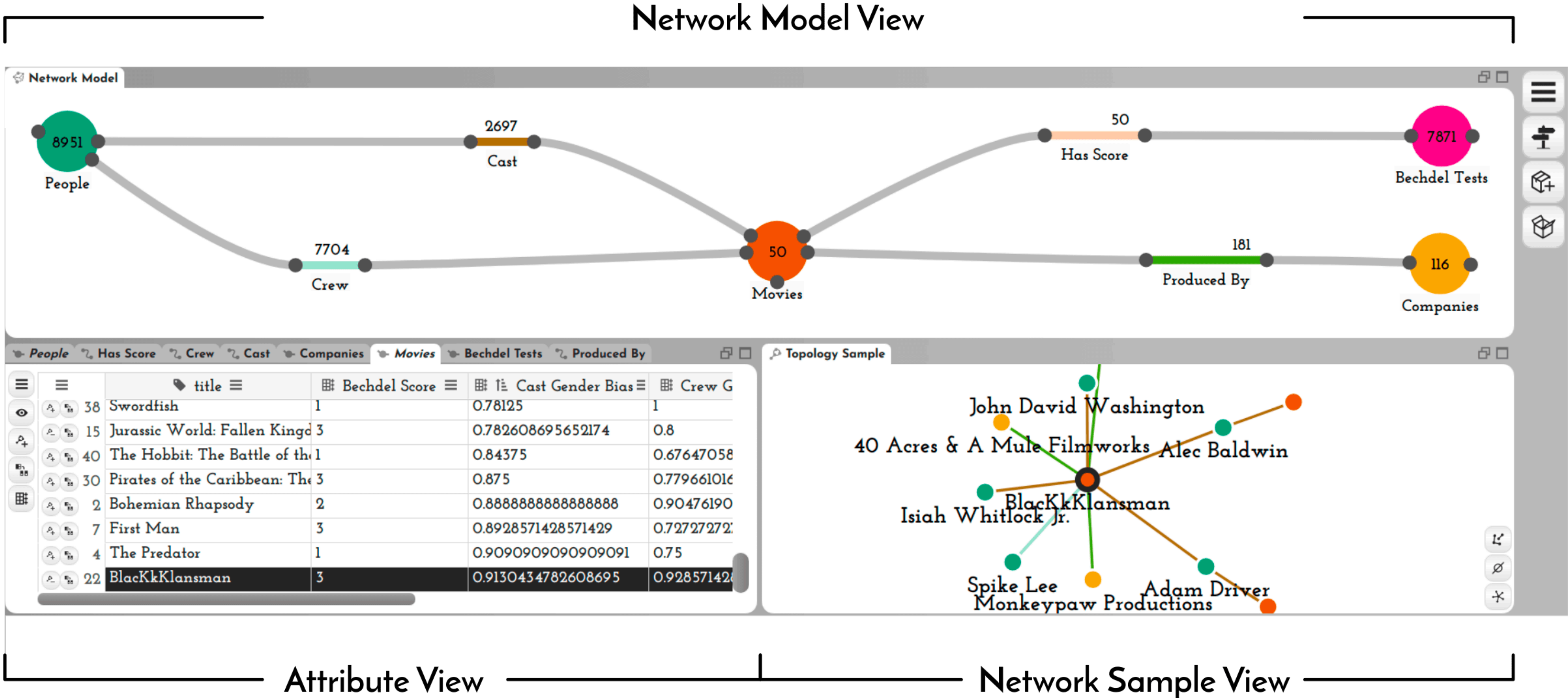
**Tailored Methods
and Systems for High
Impact Science
Problems**

EMPIRICAL & THEORETICAL WORK

**Evaluation
Methodology**

**Design Spaces /
Taxonomies**

Data Wrangling
Methods



TECHNICAL CONTRIBUTIONS

**Data Wrangling
Methods**

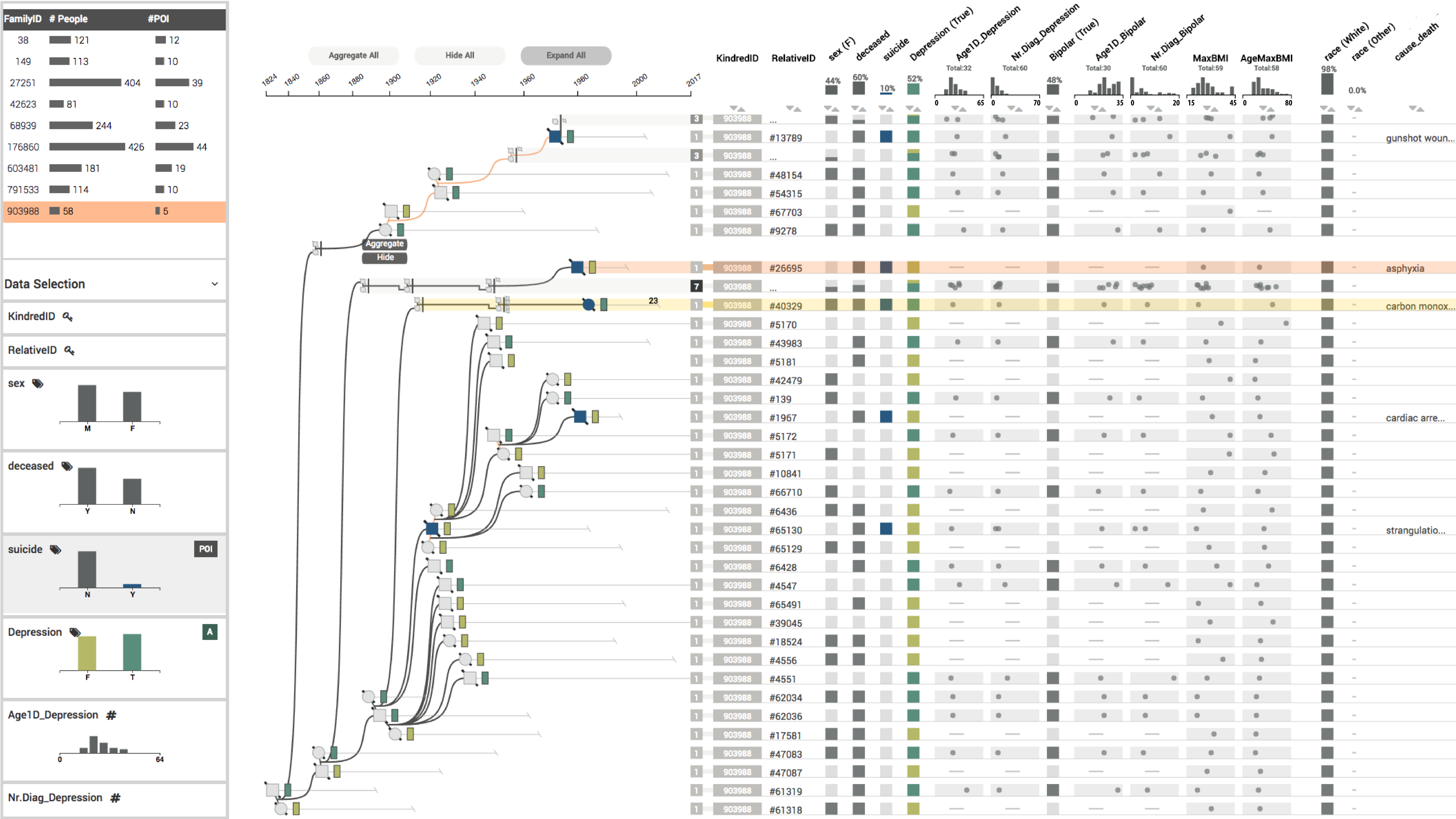
DOMAIN DRIVEN TECHNIQUES

**Tailored Methods
and Systems for High
Impact Science
Problems**

DOMAIN DRIVEN TECHNIQUES

Tailored Methods
and Systems for High
Impact Science
Problems

Genealogies for Clinical Data Analysis



DOMAIN DRIVEN TECHNIQUES

**Tailored Methods
and Systems for High
Impact Science
Problems**

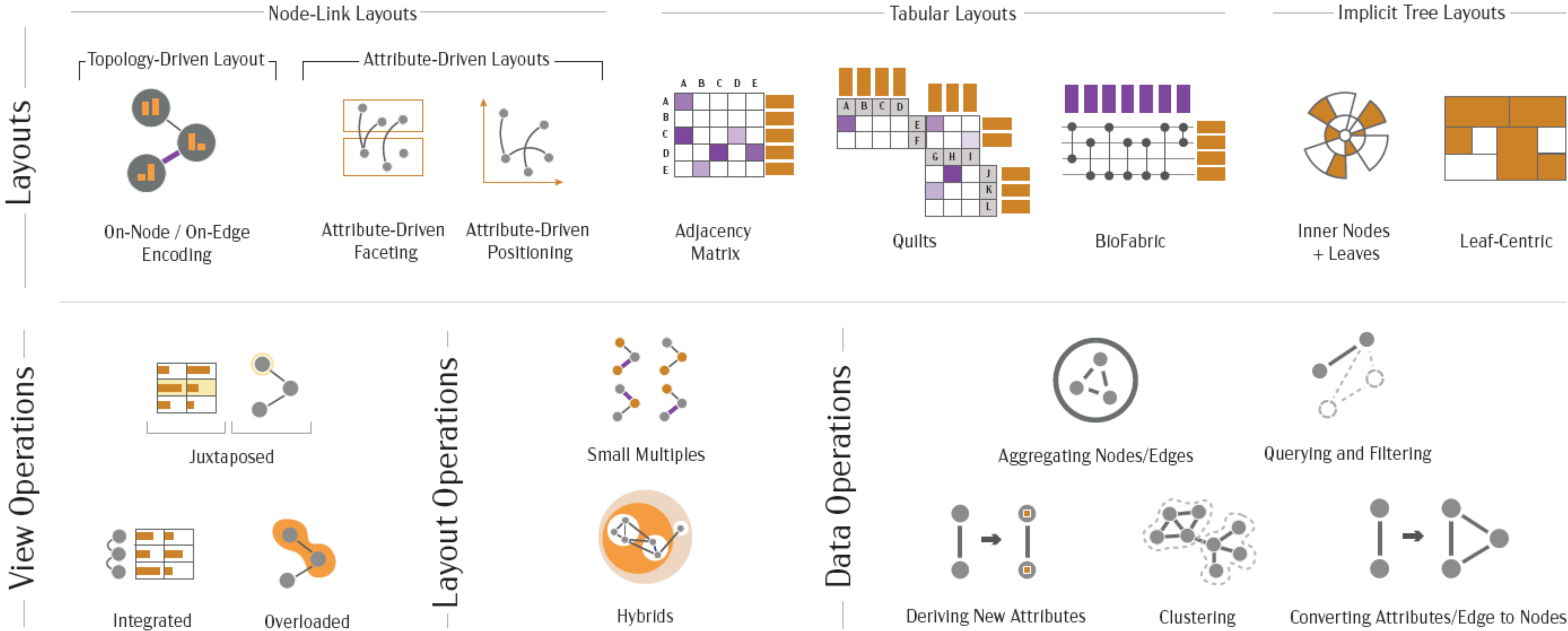
EMPIRICAL & THEORETICAL WORK

**Evaluation
Methodology**

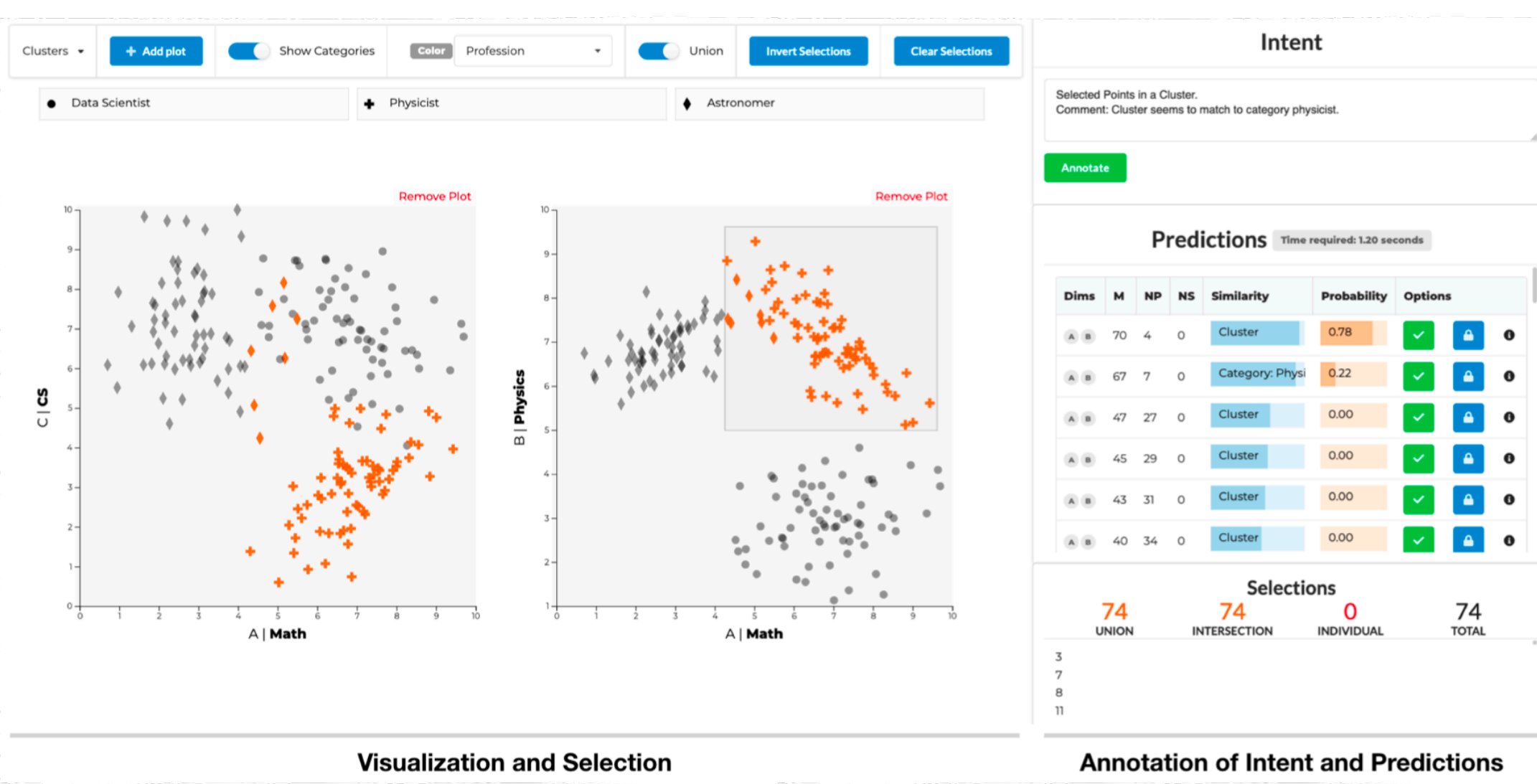
**Design Spaces /
Taxonomies**

Evaluation
Methodology

Design Spaces /
Taxonomies



TECHNICAL CONTRIBUTIONS



LITERATE VISUAL DATA ANALYSIS:

MAKING VISUAL ANALYSIS SESSIONS REPRODUCIBLE AND REUSABLE

REPRODUCIBILITY CRISIS IN SCIENCE

Reproducibility Rates

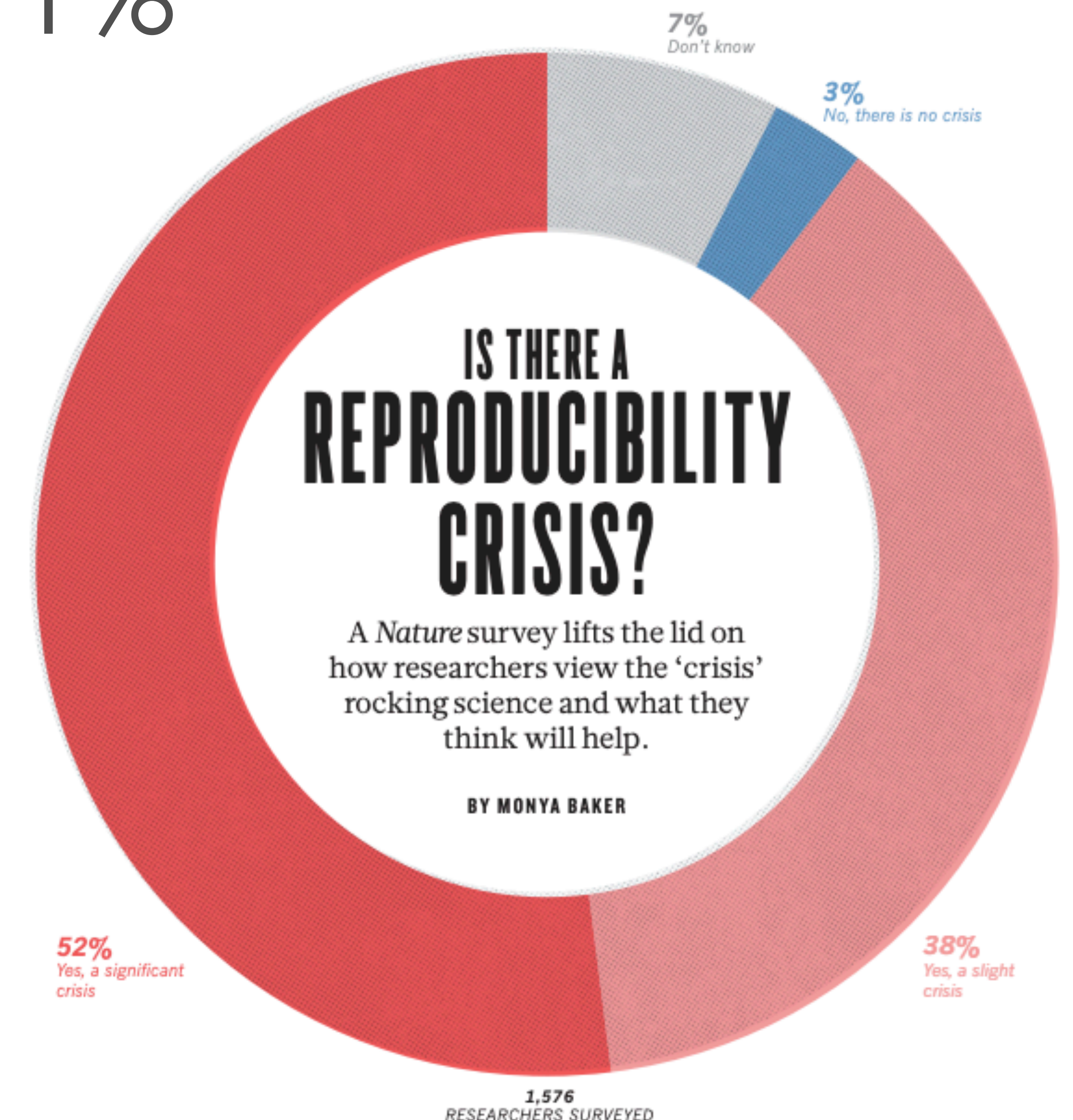
Psychology: 40%

Cancer Biology: 11%

[Baker, Penny, 2016]

[Open Science Collaboration, 2016]

[Begley, Ellis, 2012]



WHY IS THERE A REPRODUCIBILITY CRISIS?

Perverse incentives (publish or perish)

Bias for “flashy” results

Problems in data analysis

Lack of data sharing

Lack of sharing the analysis process

...

**TACKLE PROBLEMS IN DATA ANALYSIS
AND ENHANCE REPRODUCIBILITY:
LITERATE PROGRAMMING**



LITERATE PROGRAMMING

**Explain the why and how
using any means necessary!**

Text

Images / Visualizations

Formulas

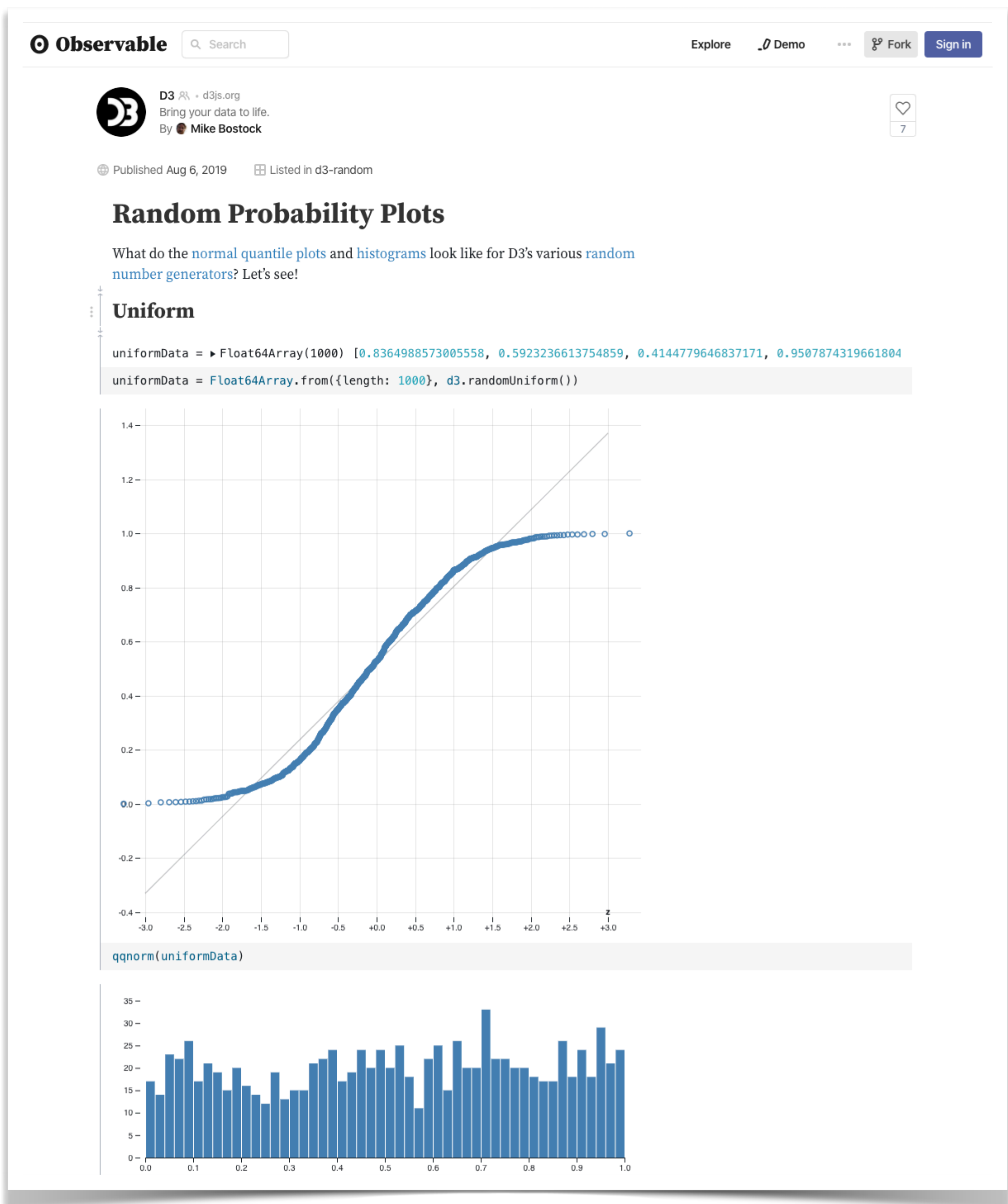
Videos

Links

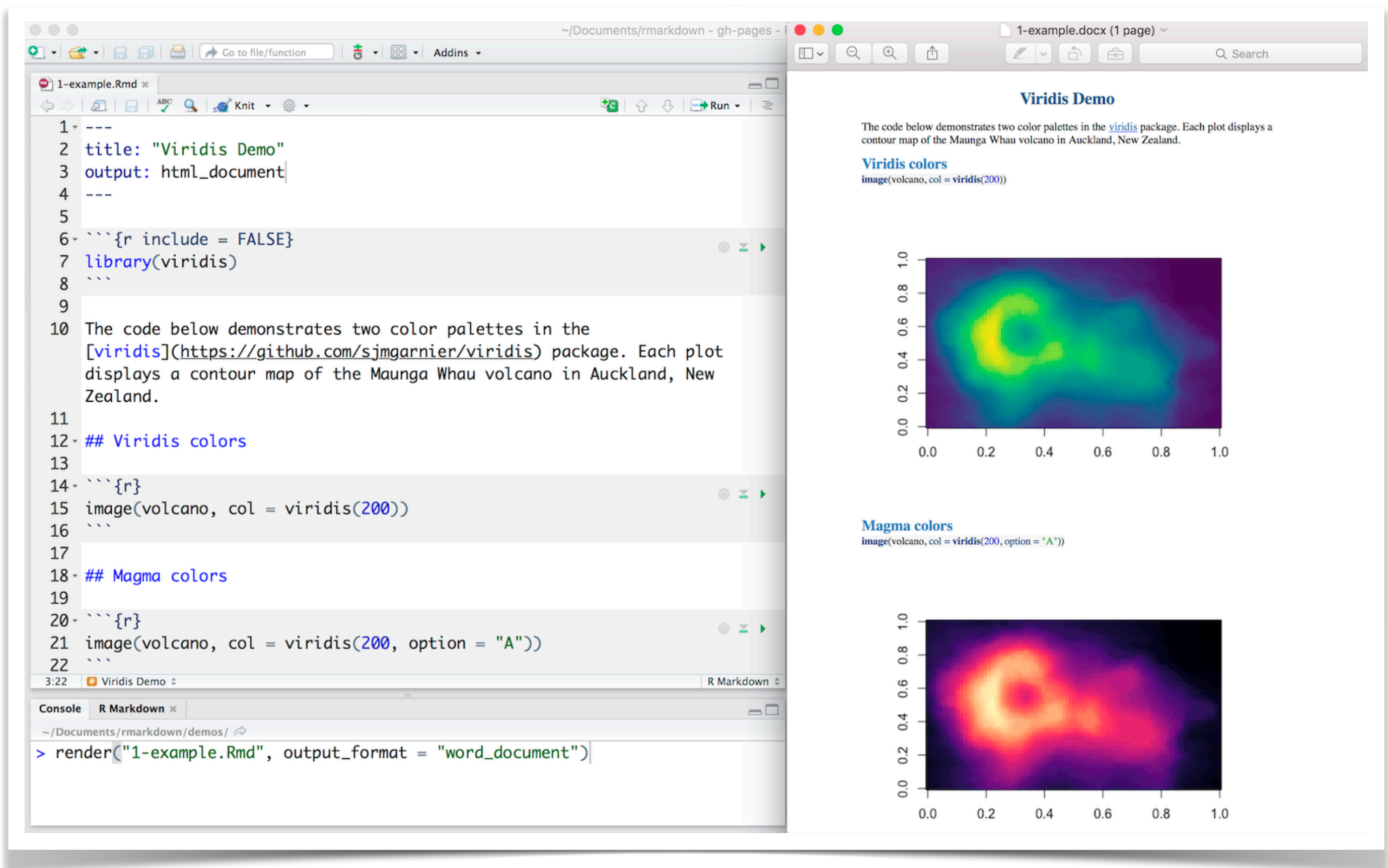
Code

[Donald E. Knuth, 1984]

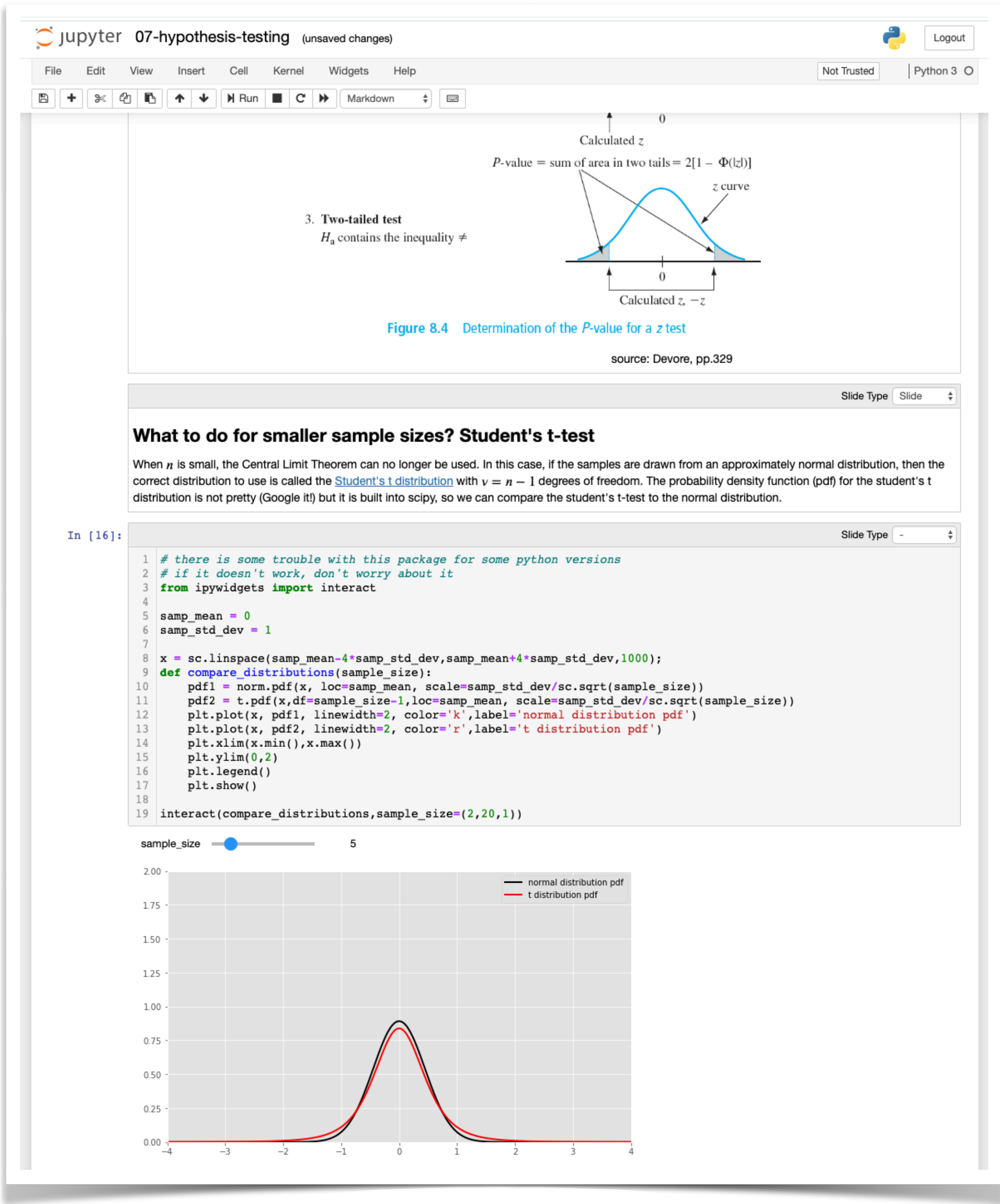
LITERATE PROGRAMMING IN THE WILD



Observable



R Markdown



Jupyter Notebooks

CHARACTERISTICS OF A GOOD NOTEBOOK

Carefully curated, well narrated.

Clean code for readability

Illustrations, Formulas, Visualizations

Complex, multi-stage process

**THERE IS NO STRAIGHTFORWARD WAY TO DO
LITERATE DATA VISUALIZATION**

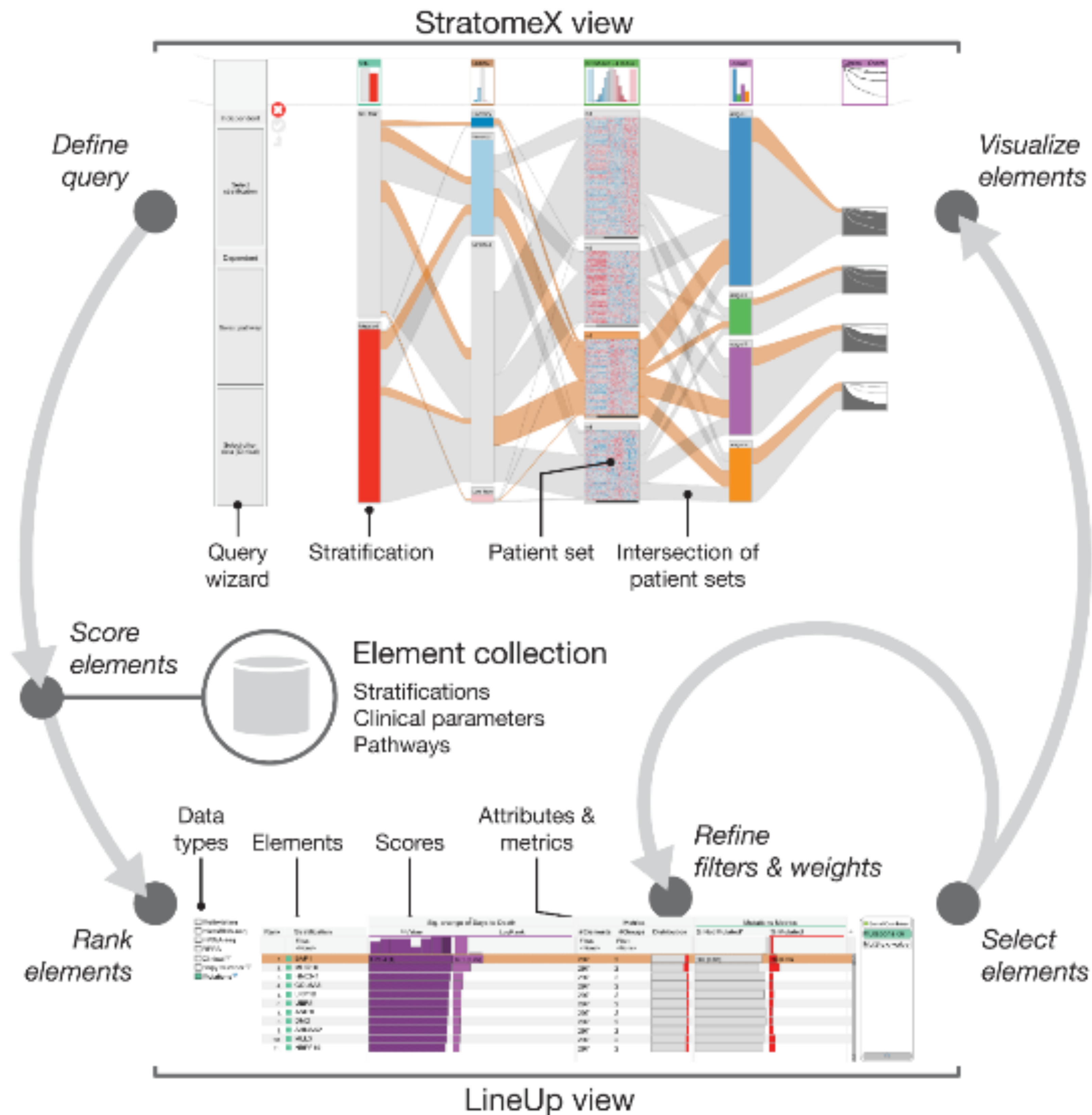
LITERATE VISUAL DATA ANALYSIS

Current State:

no record of **what was done,
let alone **why****

**Idea: make the process of an interactive,
visual analysis session well reasoned and
documented**

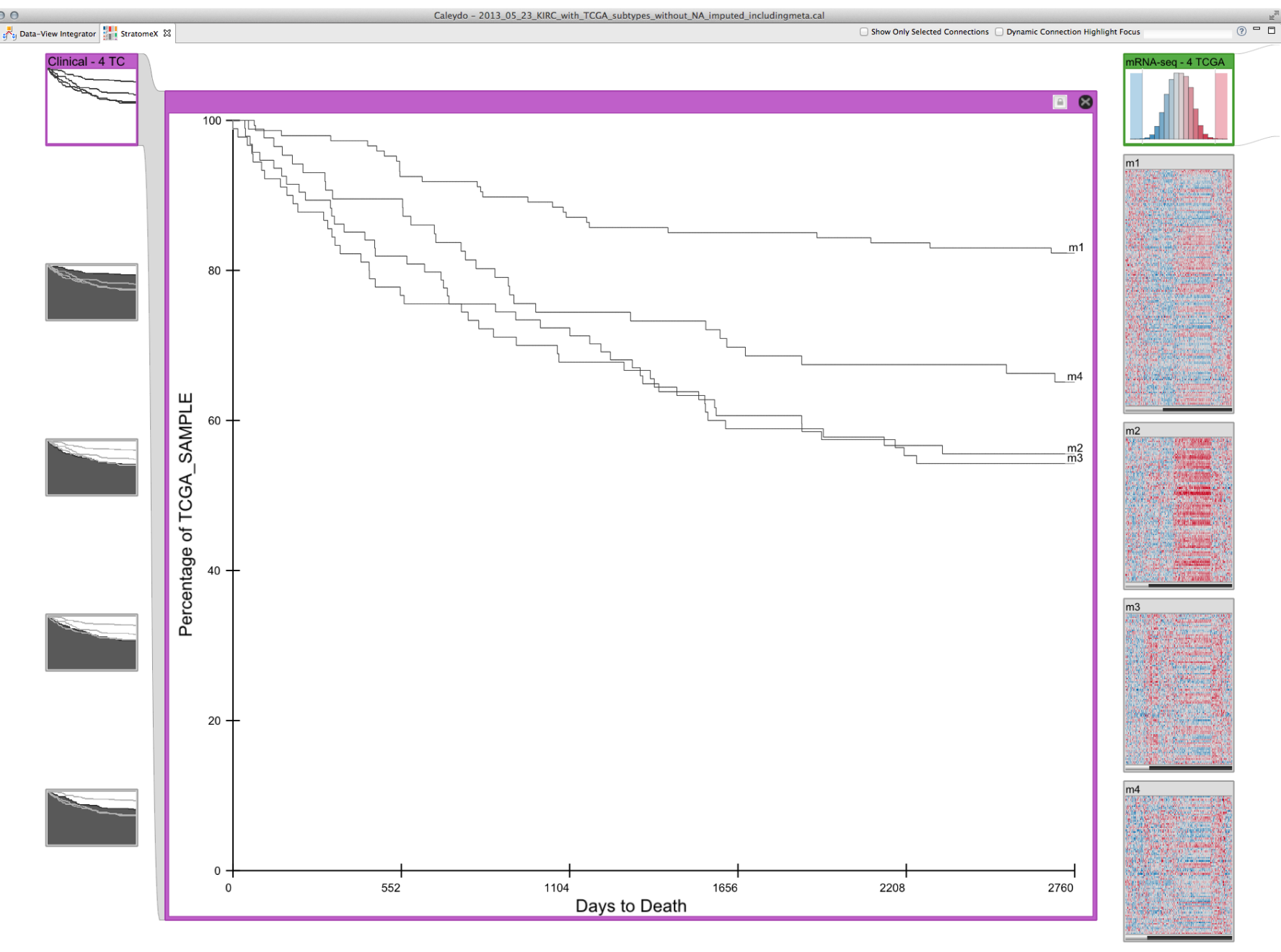
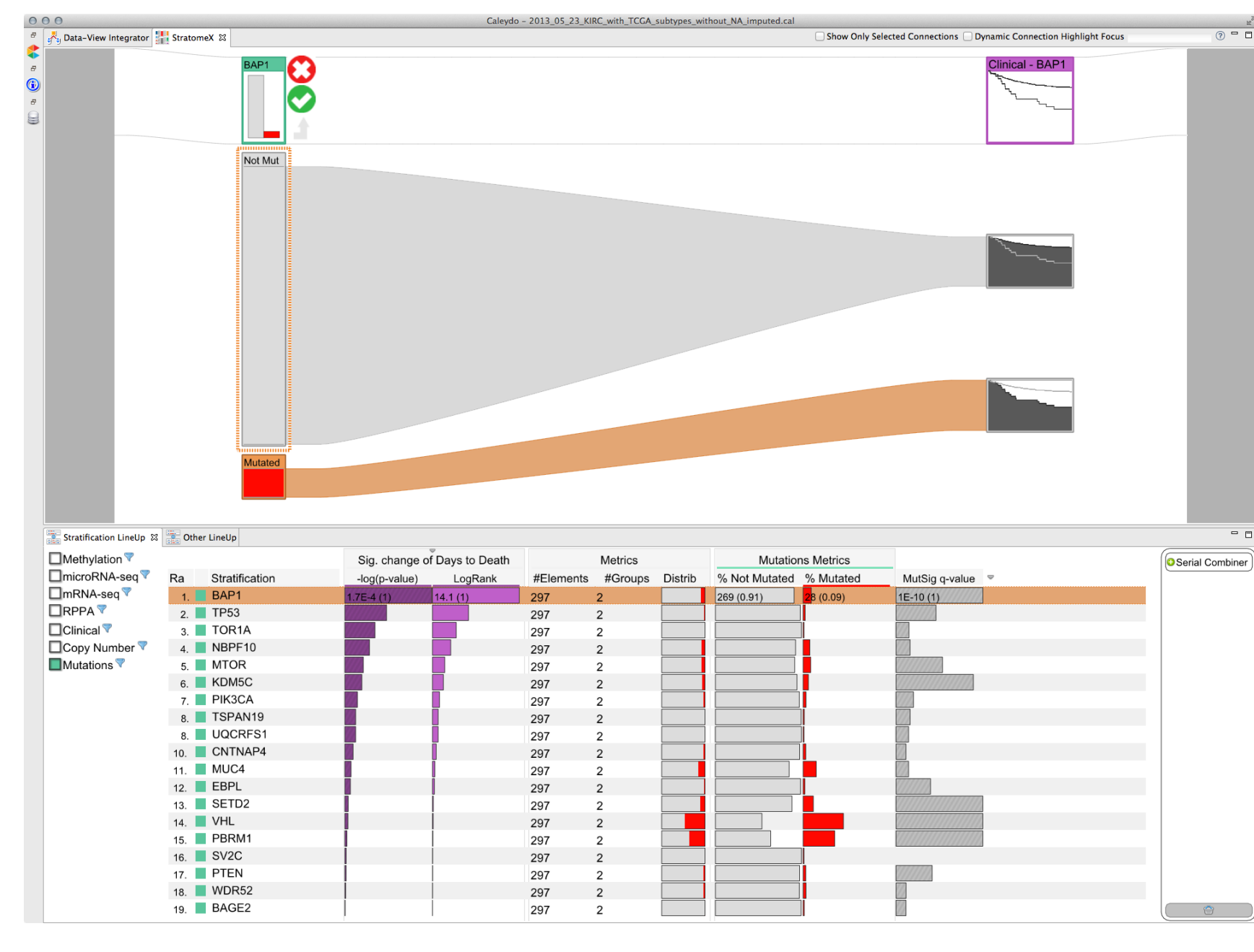
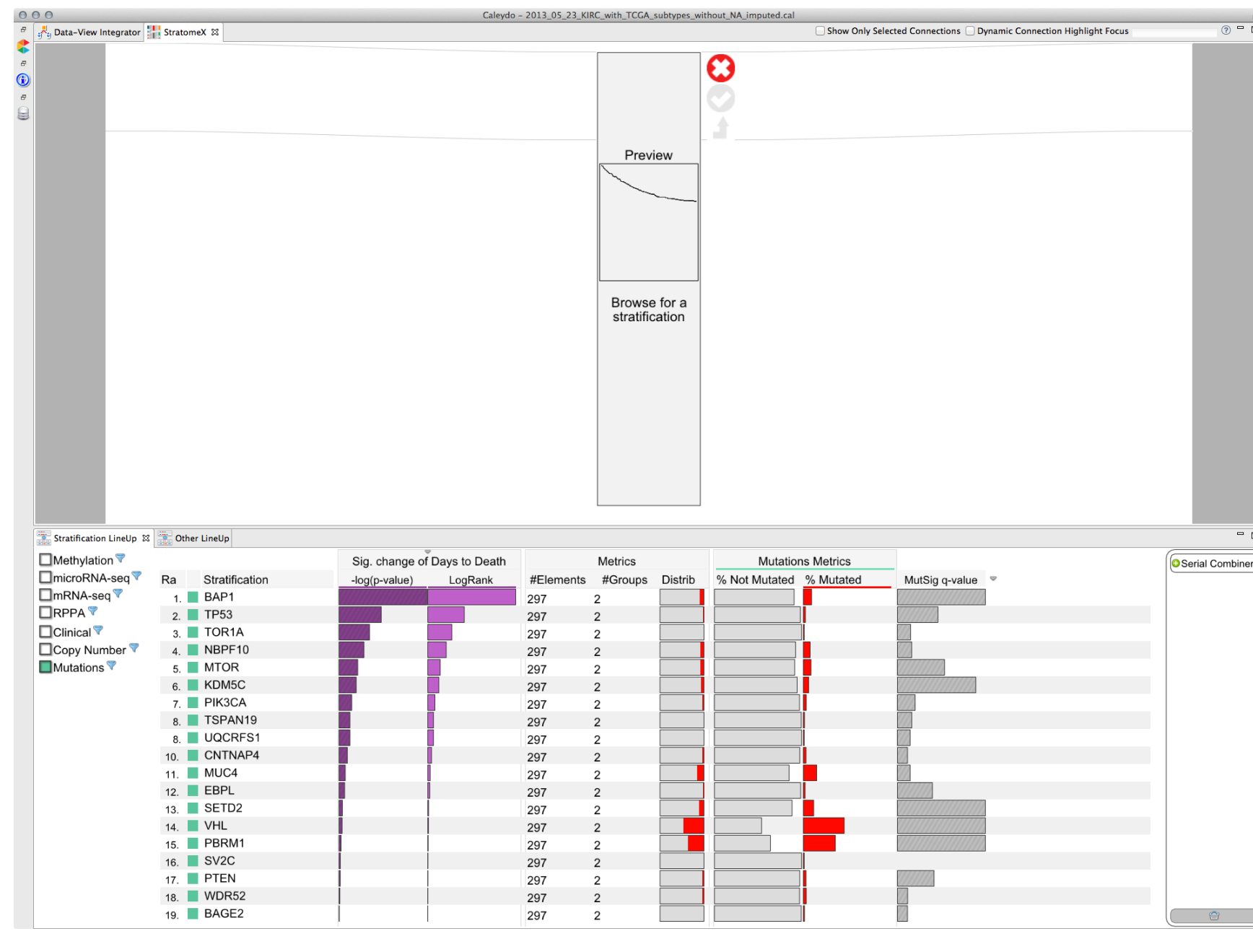
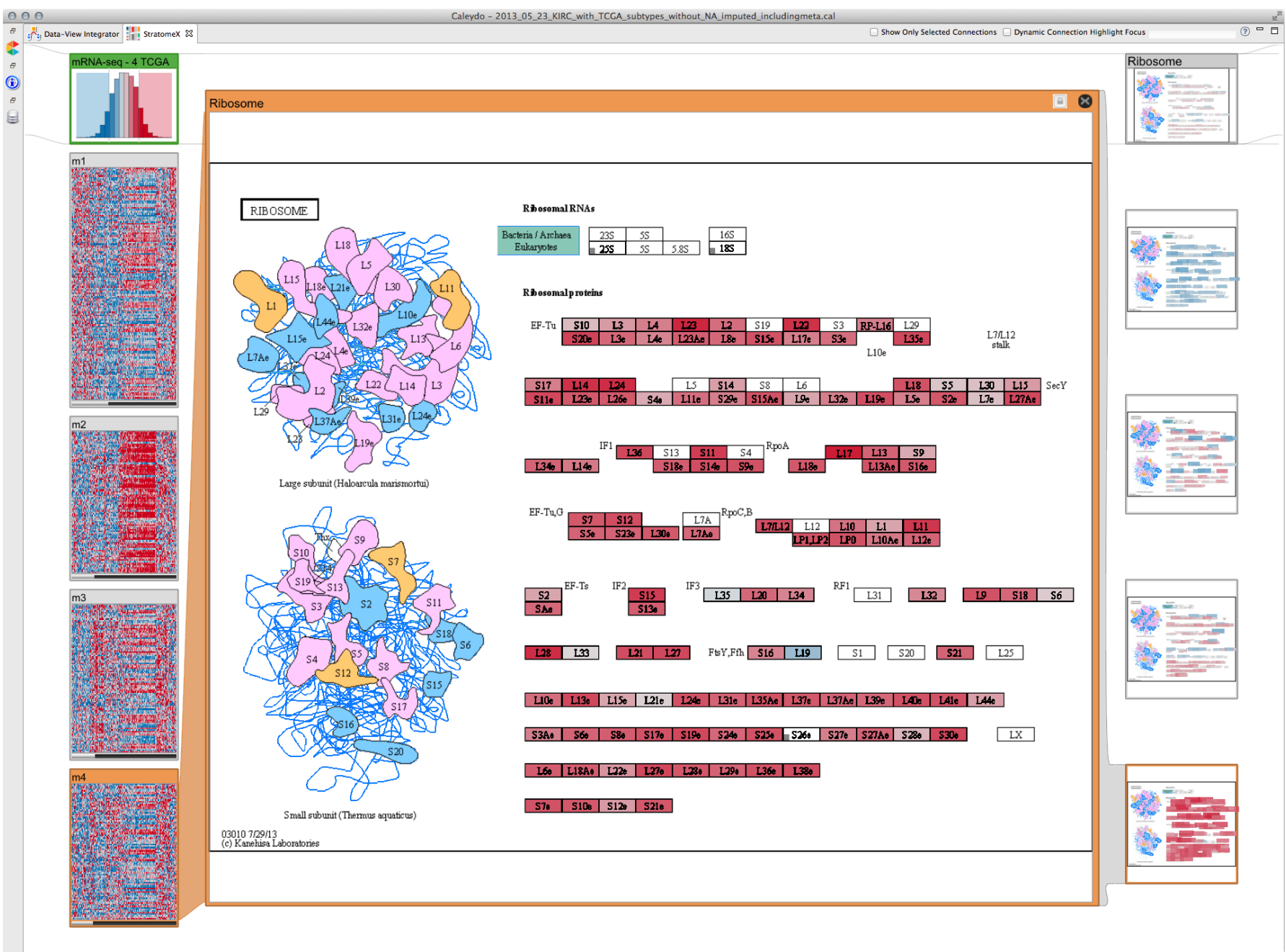
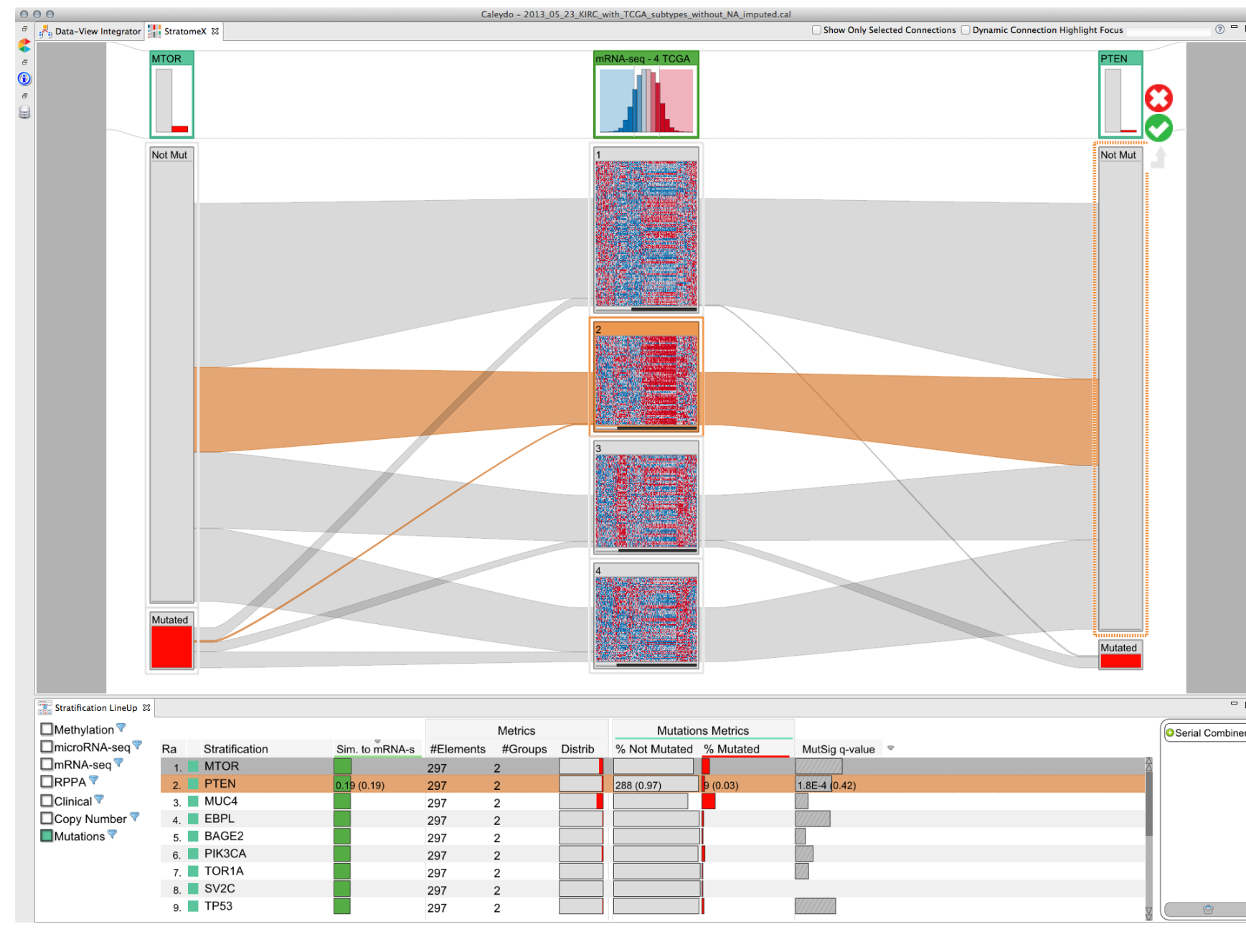
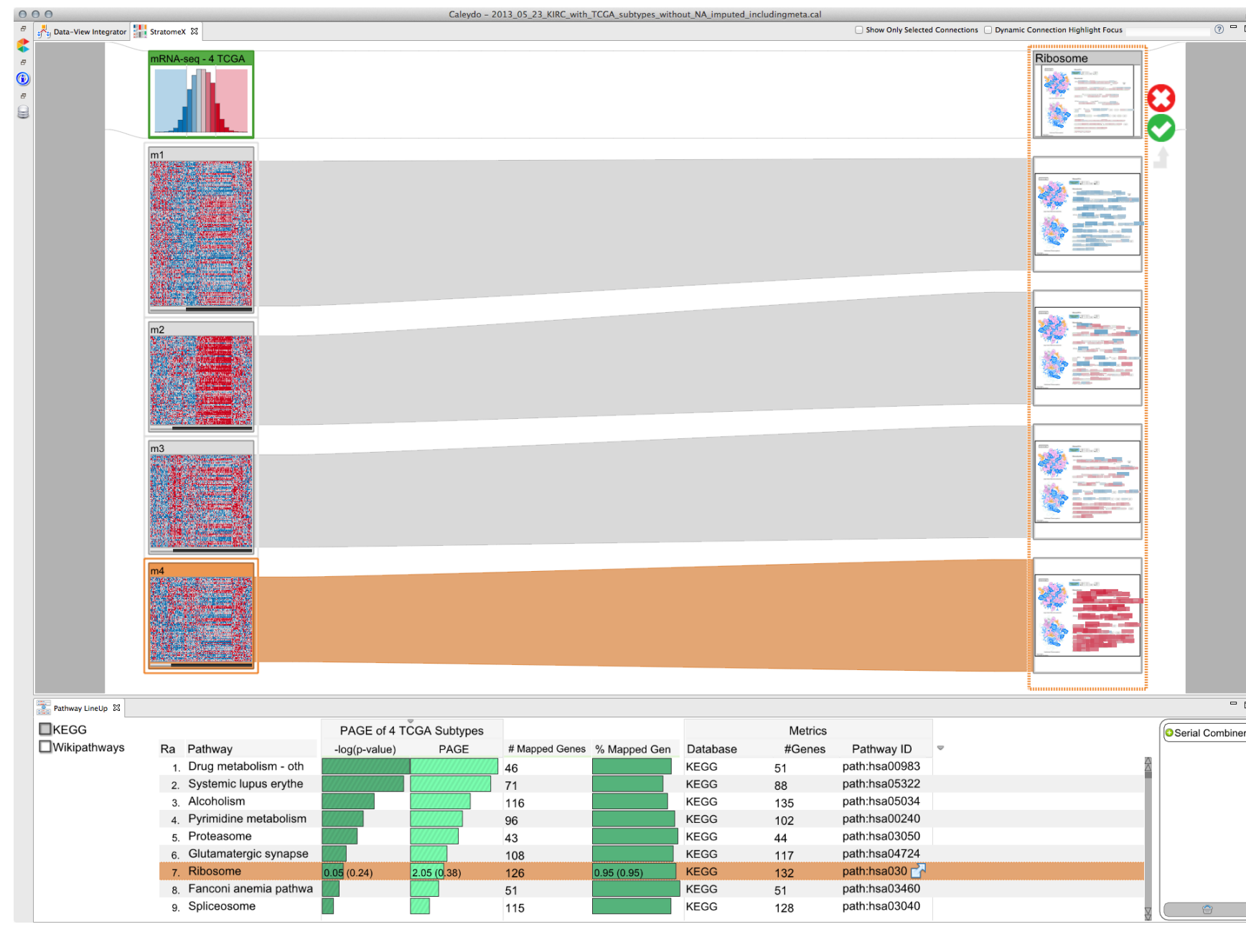
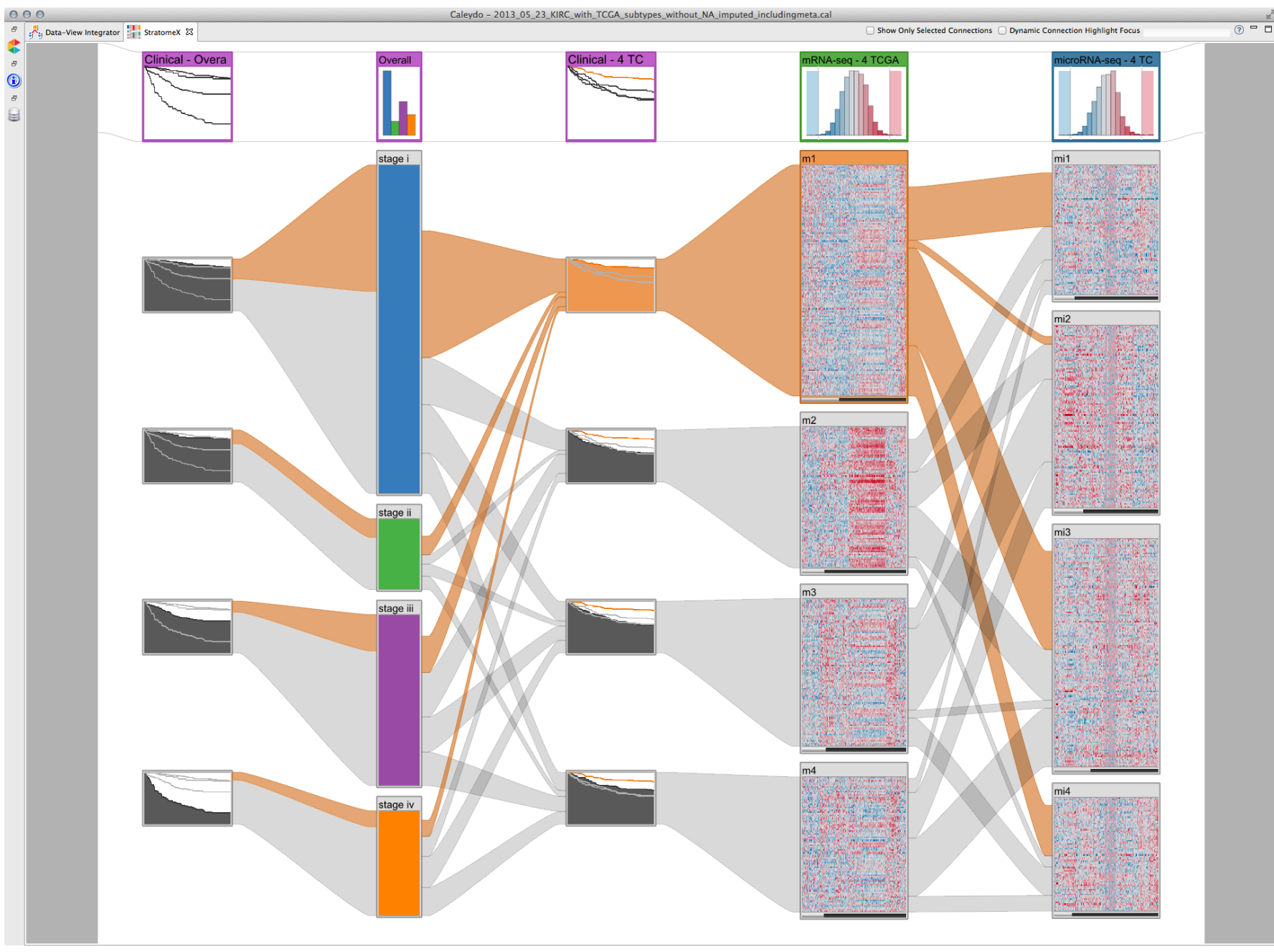
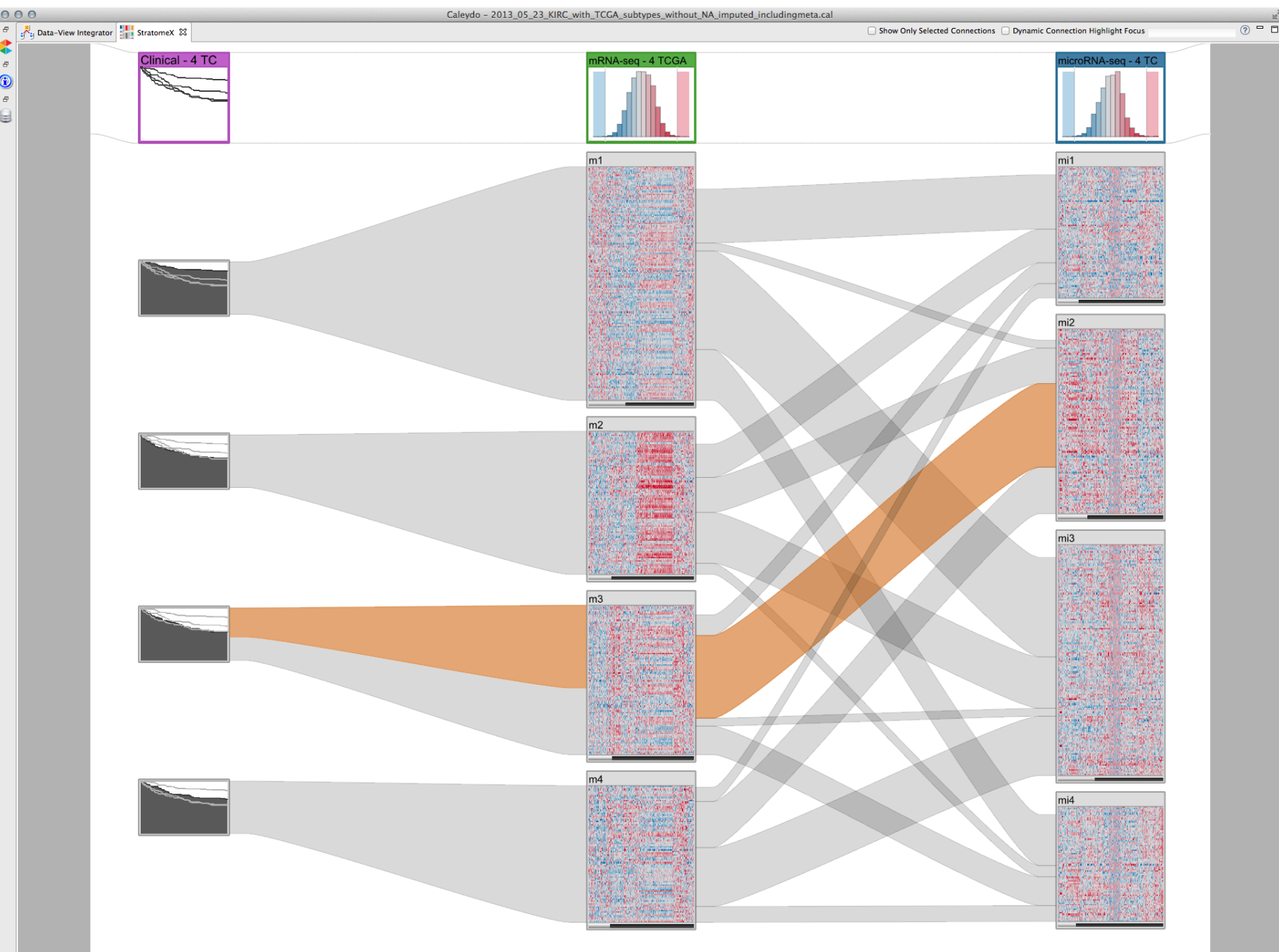
A MANUAL ATTEMPT AT LITERATE VISUALIZATION



Guided Visual Exploration of Genomic Stratifications in Cancer

Nature Methods 11, 9 (2014), 884–885





CASE STUDY

Streit et al., Guided Visual Exploration of Genomic Stratifications in Cancer, Nature Methods, 2014

**IDEA: USE ANALYSIS SESSION
PROVENANCE AS BASIS FOR
A DATA STORY**

PROVENANCE TRACKING

Provenance (from the French provenir, 'to come from/forth') is the chronology of the ownership, custody or location of a historical object. [Wikipedia]

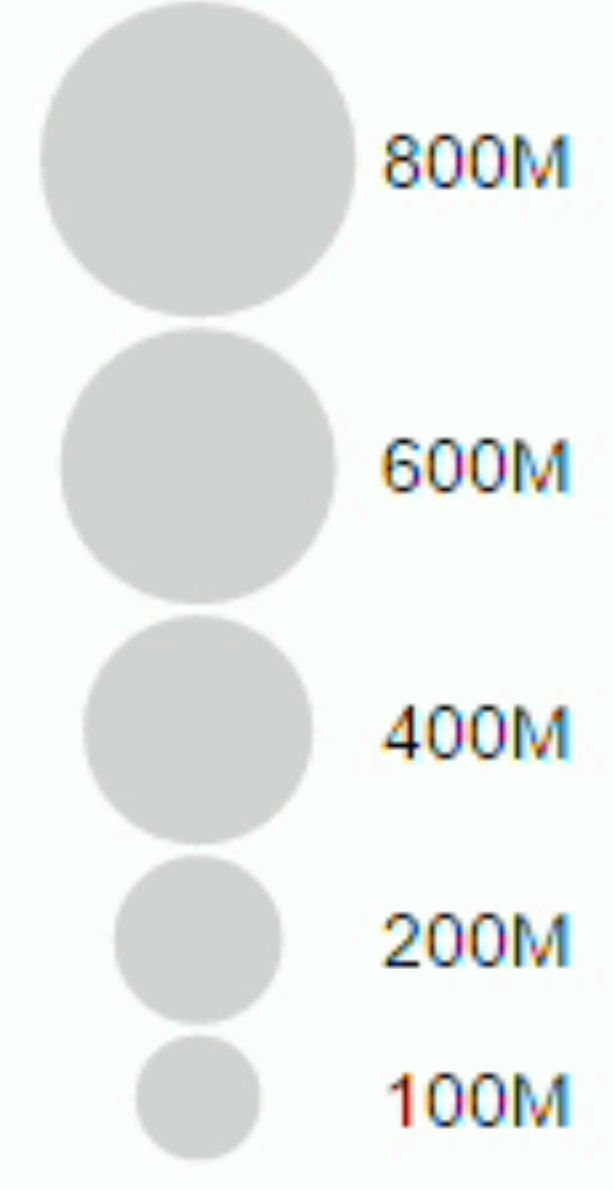
In CS: a log, a record of everything that lead to a state

Color: Continent

- Africa
- America
- Asia
- Europe

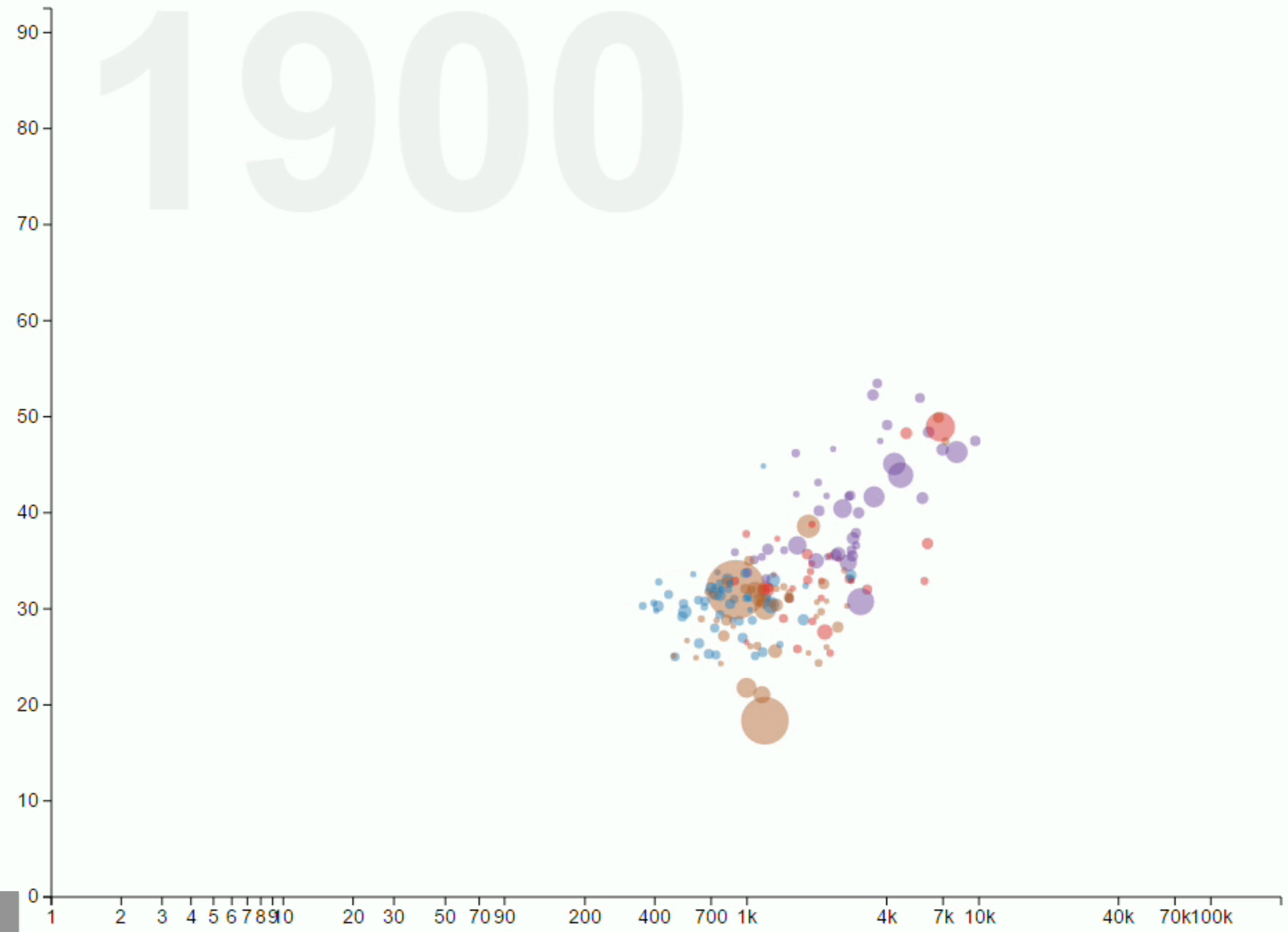
Size: Population

Sqrt



Life Expectancy at Birth

Linear



GDP (PPP\$, inflation-adjusted) Log

Exploration

1800 1820 1840 1860 1880 1900 1920 1940 1960 1980 2000

Provenance

- Start
- X=GDP
- Y=Life Expectancy
- Size=Population
- Color=Continent
- Year 1800
- scale(X)=Log
- Year 1900

- Data
- Visual
- Selections
- Layout
- Analysis

Color:

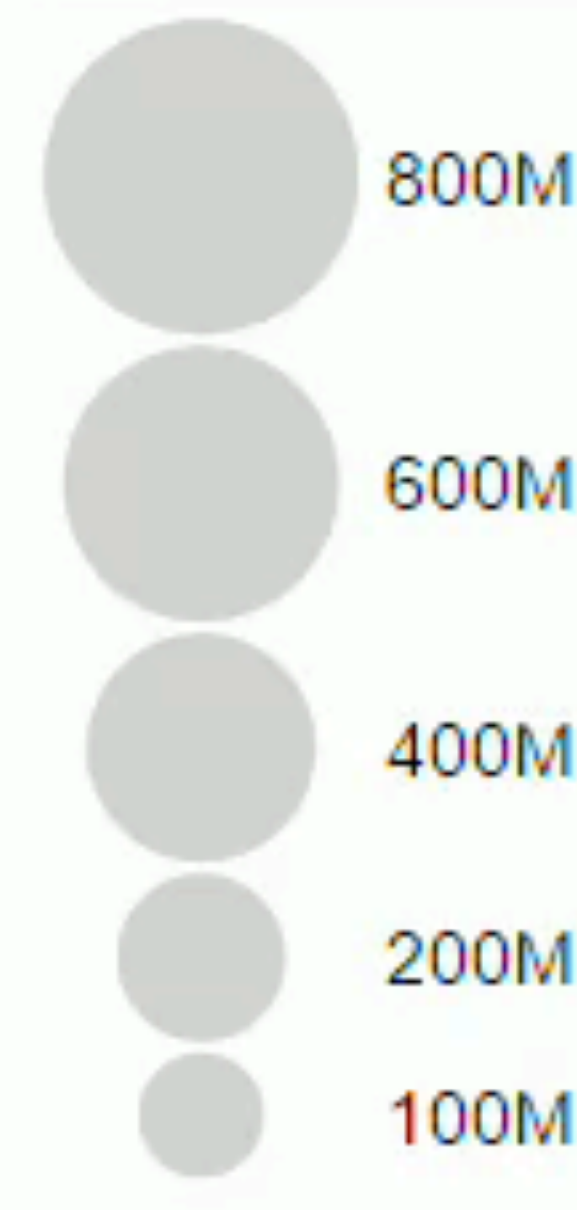
Continent

- Africa
- America
- Asia
- Europe

Size:

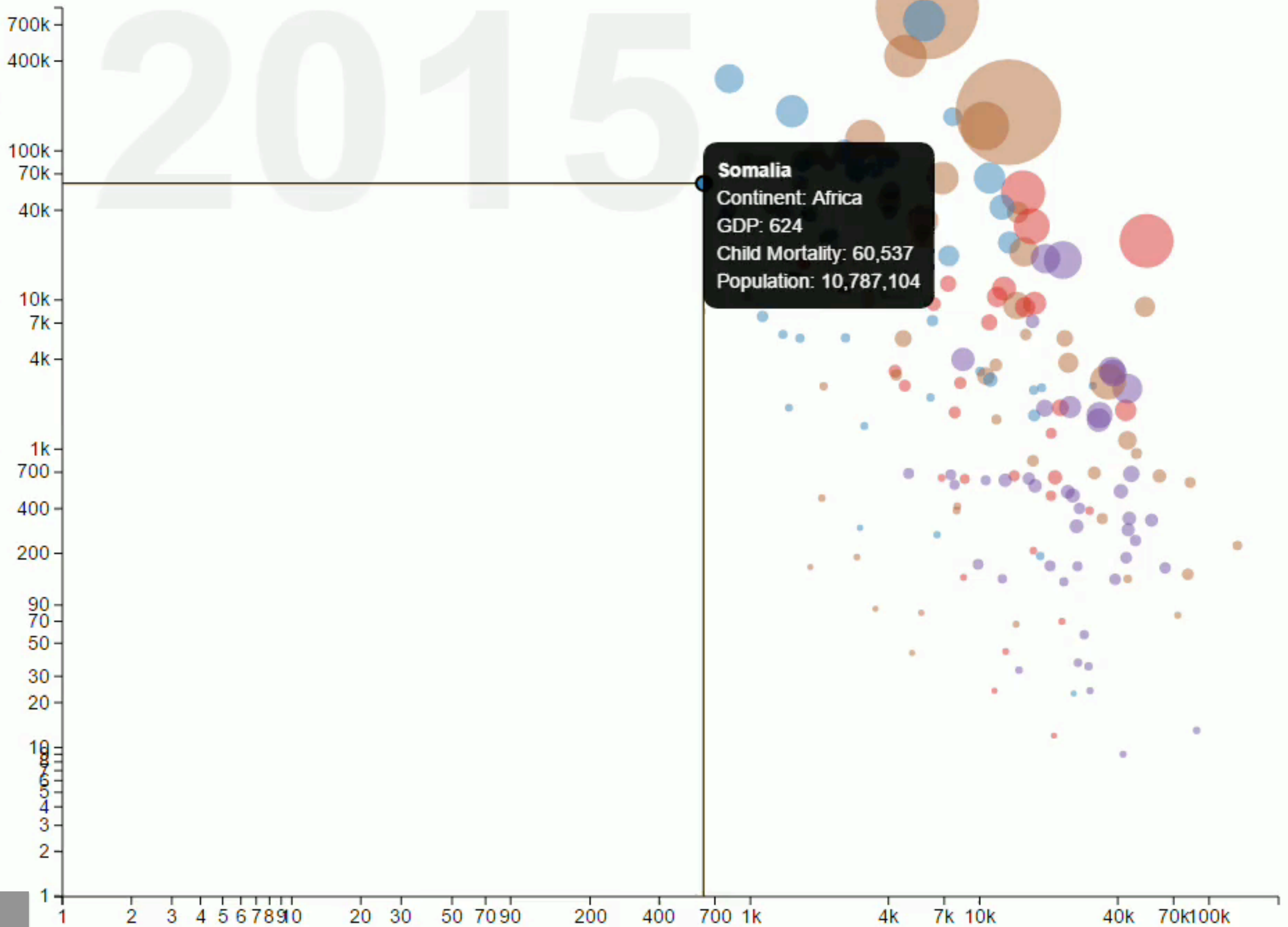
Population

Sqrt



Log

(0-5 year-olds dying per 1,000 born)



Somalia

Continent: Africa

GDP: 624

Child Mortality: 60,537

Population: 10,787,104

GDP

(PPP\$, inflation-adjusted)

Log

Authoring

Provenance

- Year 1800
- Y=Child Mortality
- scale(X)=Log
- scale(Y)=Log
- Year 1860
- Year 1920
- Year 2015
- Country Somalia

- Data
- Visual
- Selections
- Layout
- Analysis



Annotations
A → □



Provenance

- 86 TCGA_SAMPLES
- 62 TCGA_SAMPLES
- 32 TCGA_SAMPLES
- no TCGA_SAMPLES
- 147 TCGA_SAMPLES
- patient.daystodeath (4 TCGA_SAMPLES)
- patient.daystodeath (4 TCGA_SAMPLES)
- patient.daystodeath (4 TCGA_SAMPLES)

- Data
- Visual
- Selections
- Layout
- Analysis

Story

- Compare Stratifications
- microRNA-seq (4 TCGA Subtypes)
- Select Bands 1
- Select Bands 2
- Select Bands 3
- Select Bands 4
- Select Bands 5
- Viewing Stratified Patient Survival
- 147 TCGA_SAMPLES
- Days to Death↔mRNA-seq (4 TCGA Subtypes)
- Kaplan Maier Plot
- Adding a Categorical Parameter
- Days to Death↔microRNA-seq (4 TCGA Subtypes)
- Days to Death↔Clinical (Overall Stage)

⏮

▶

⏭

**PROGRESS;
BUT
PROVENANCE IS
“DUMB”**

**We solved the WHAT,
but not the WHY**

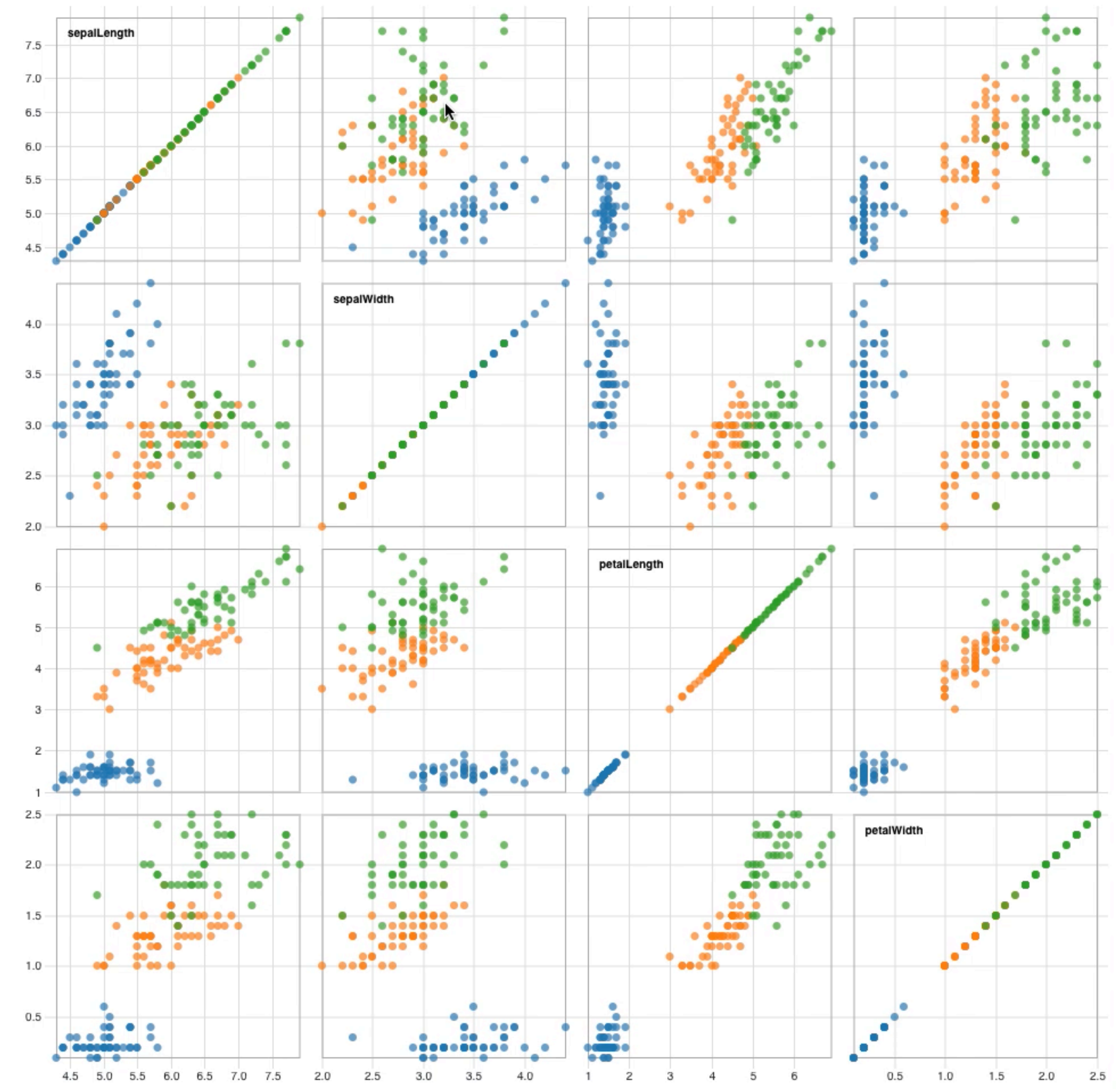
No progress towards reusability..

So, what else can we do?

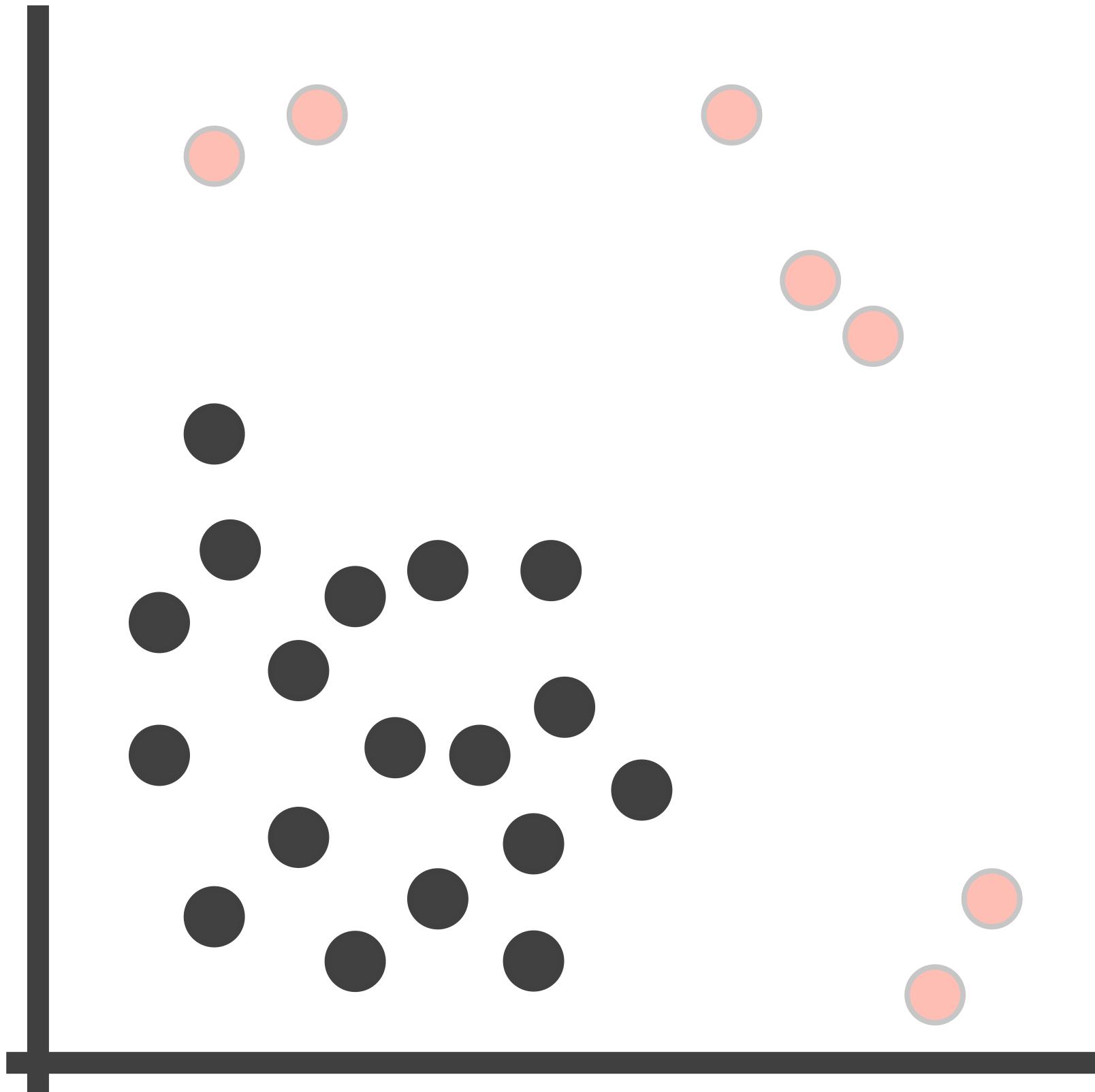
SEMANTIC SELECTIONS

Information Visualization 2021. Kiran Gadhave, Jochen Görtler, Carolina Nobre,
Oliver Deussen, Miriah Meyer, Jeff Phillips, Alexander Lex

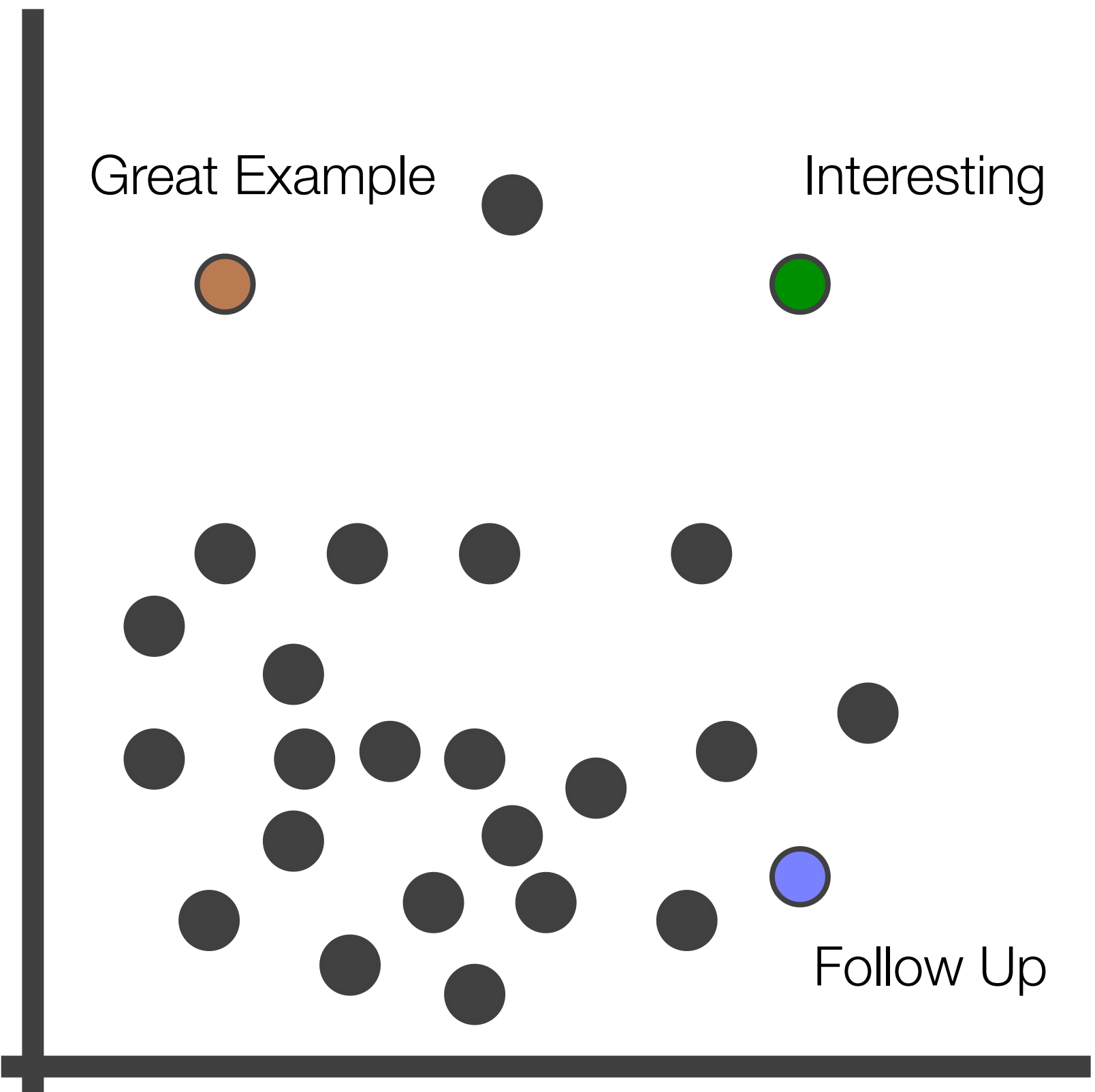
WHAT ARE SELECTIONS?



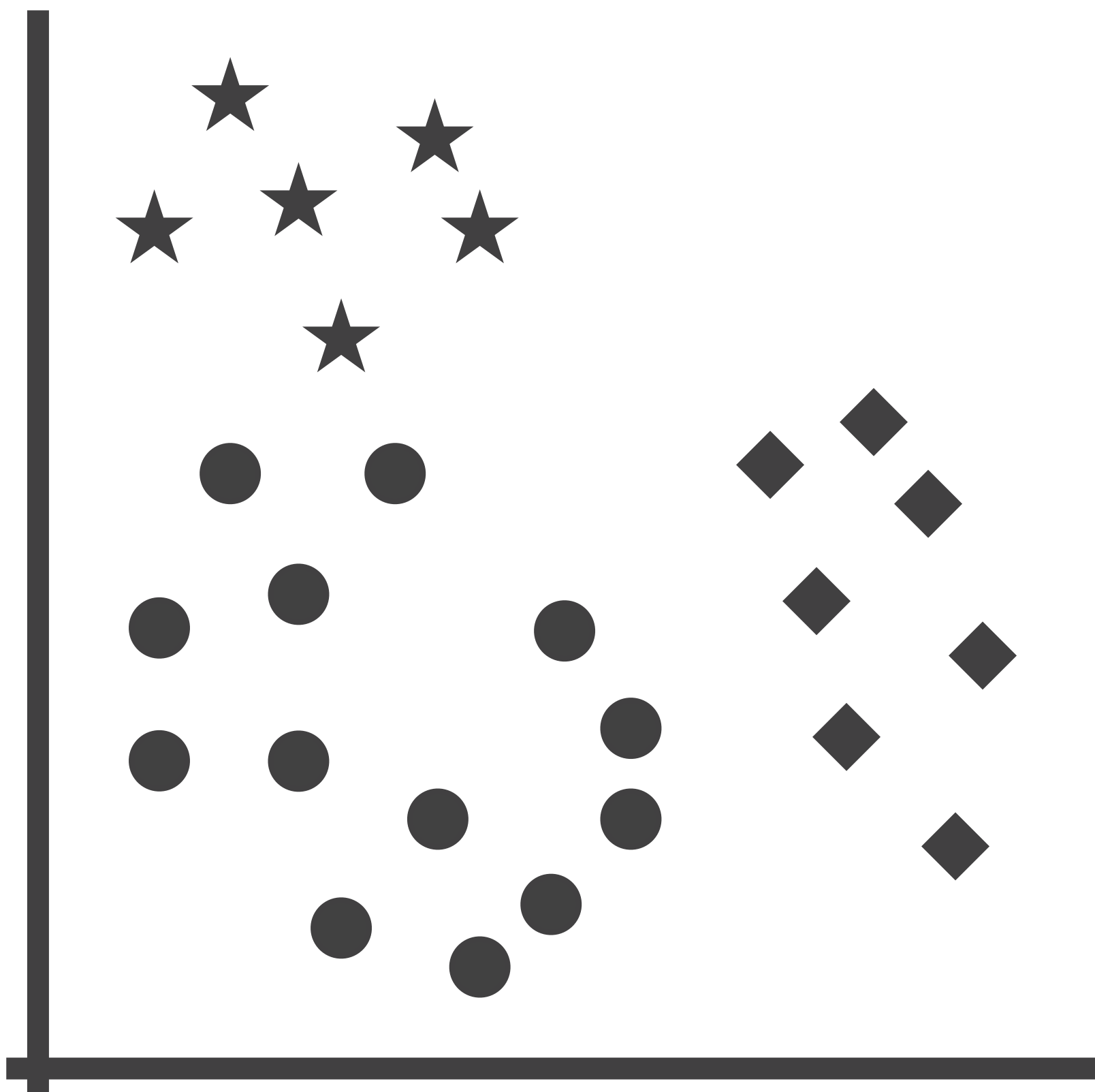
FROM SELECTIONS TO ADVANCED OPERATIONS



Filter

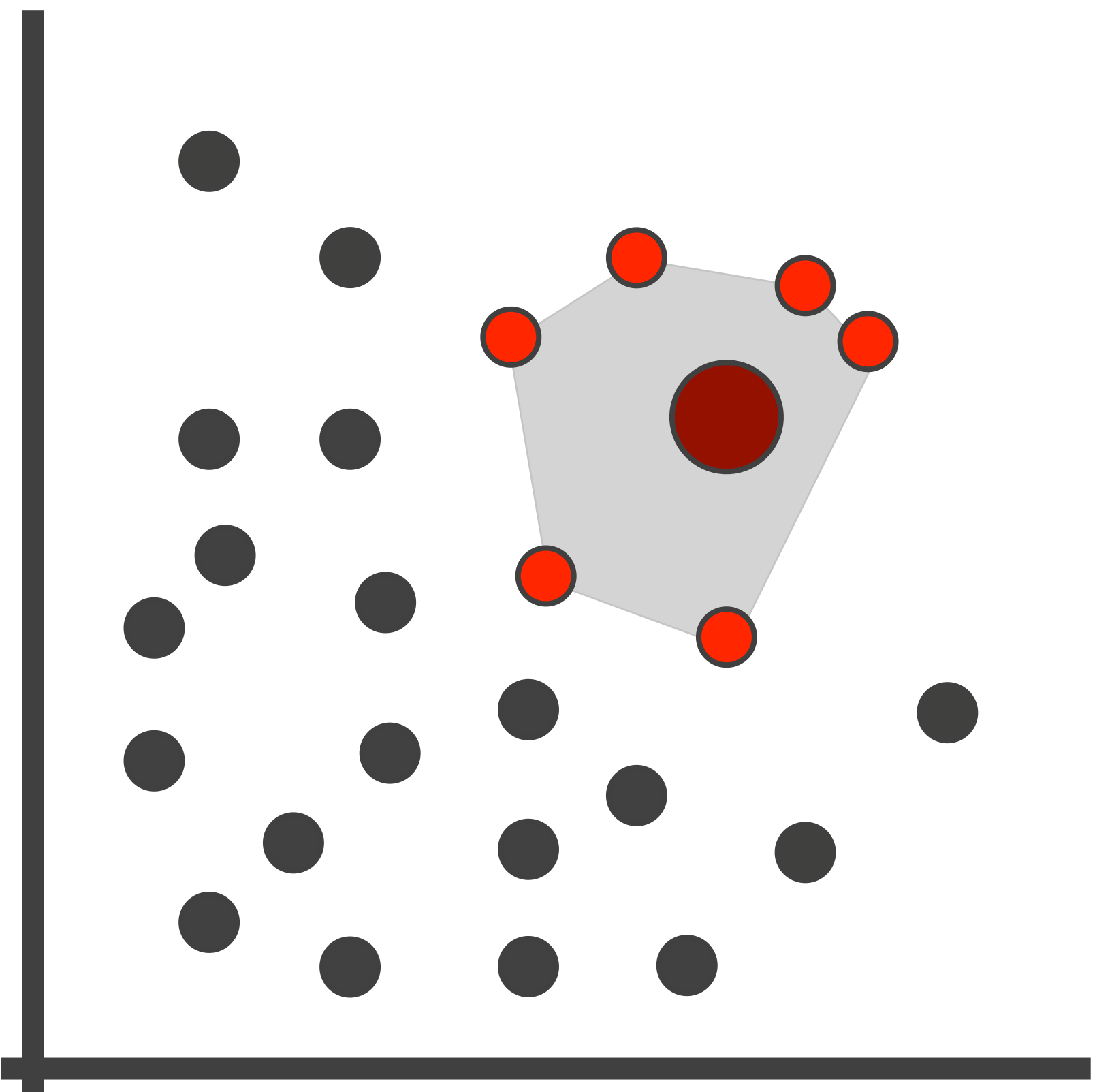


Label



Categorize

- Unassigned
- ◆ Category A
- ★ Category B



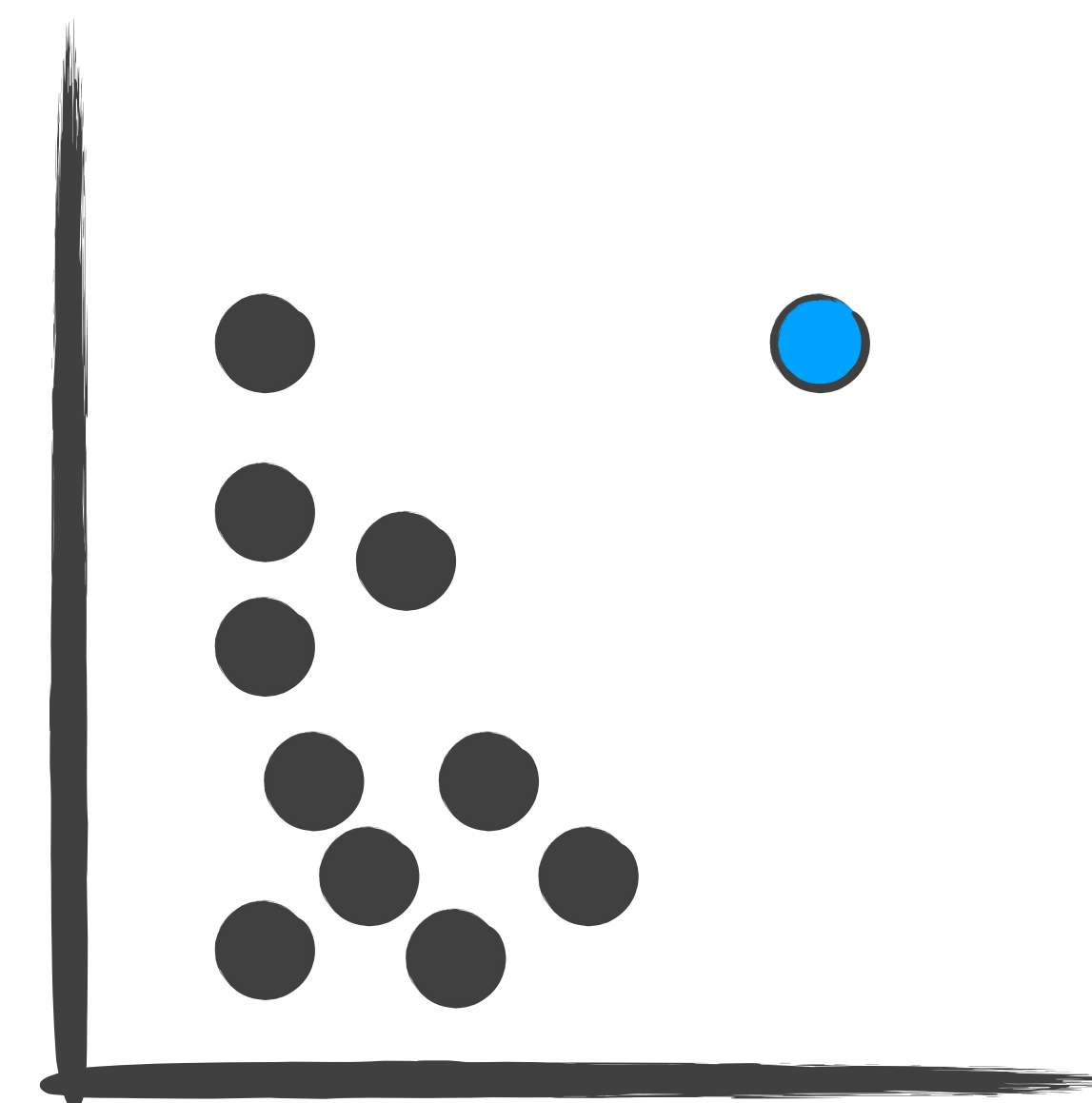
Aggregate

WHAT IS INTENT WHEN SELECTING?

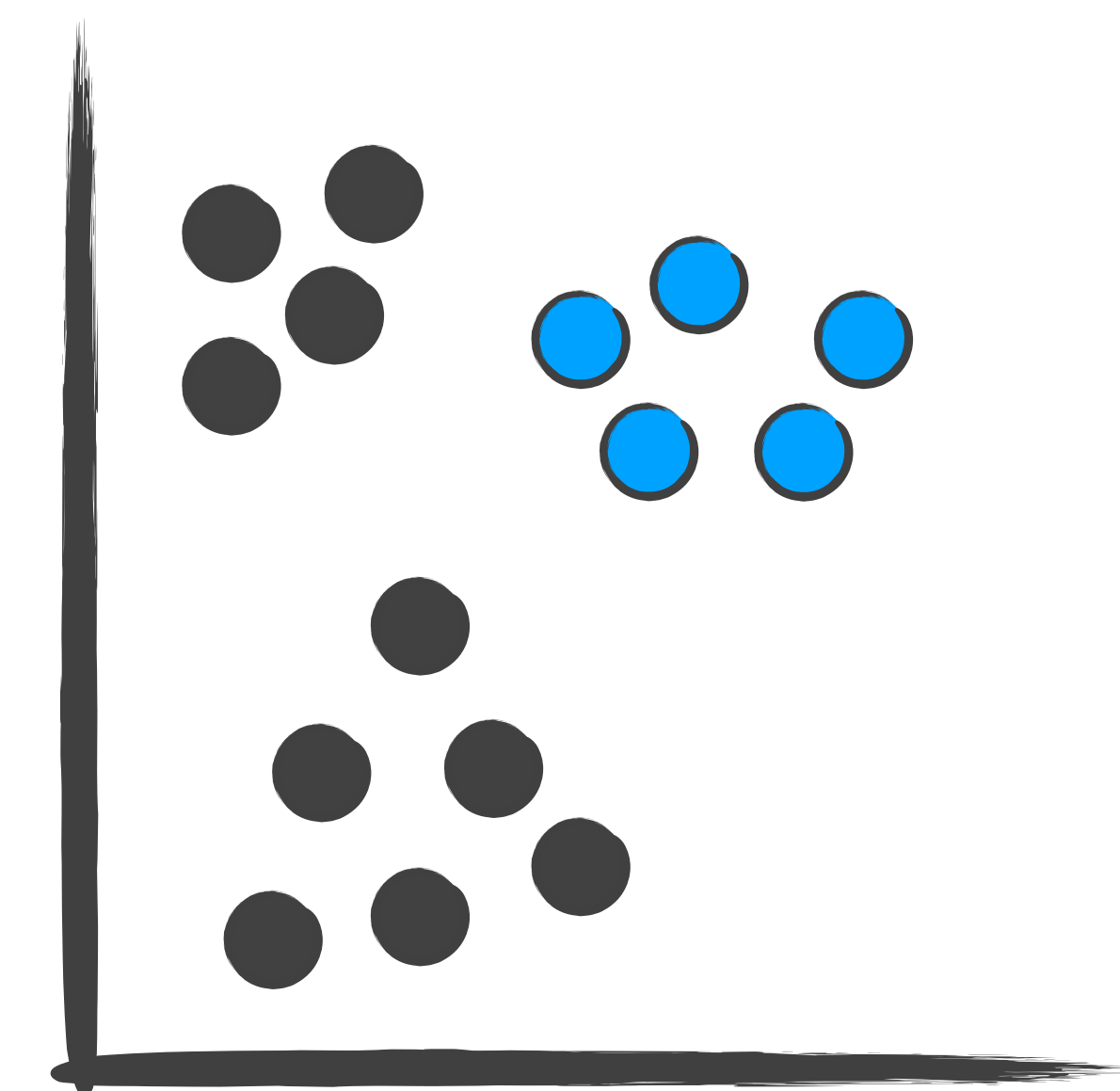
Intent is the user's reason for performing a brush with a visualization.

Domain Specific Intent: Capture through Annotation

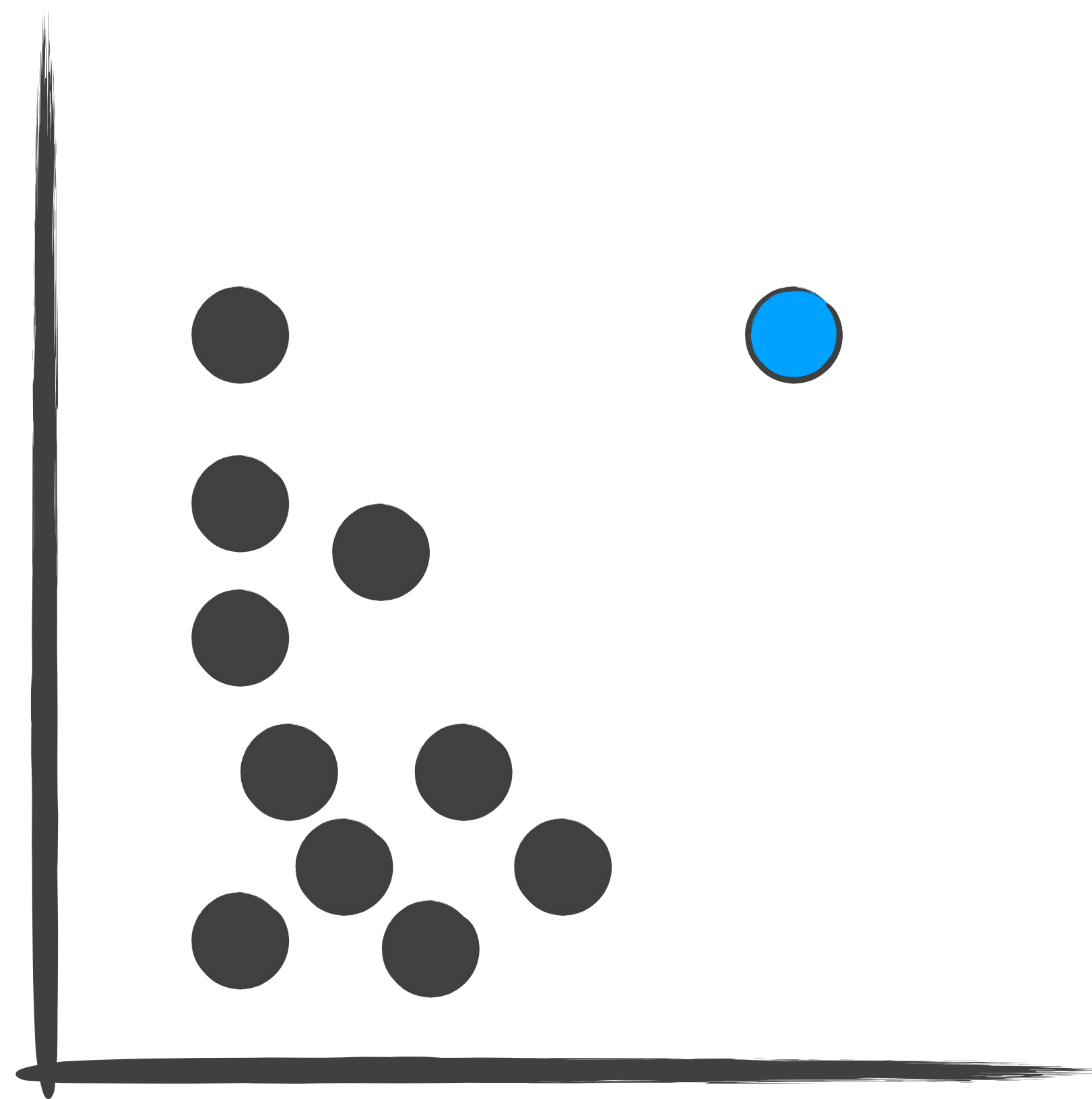
Pattern-Based Intent: Capture Automatically



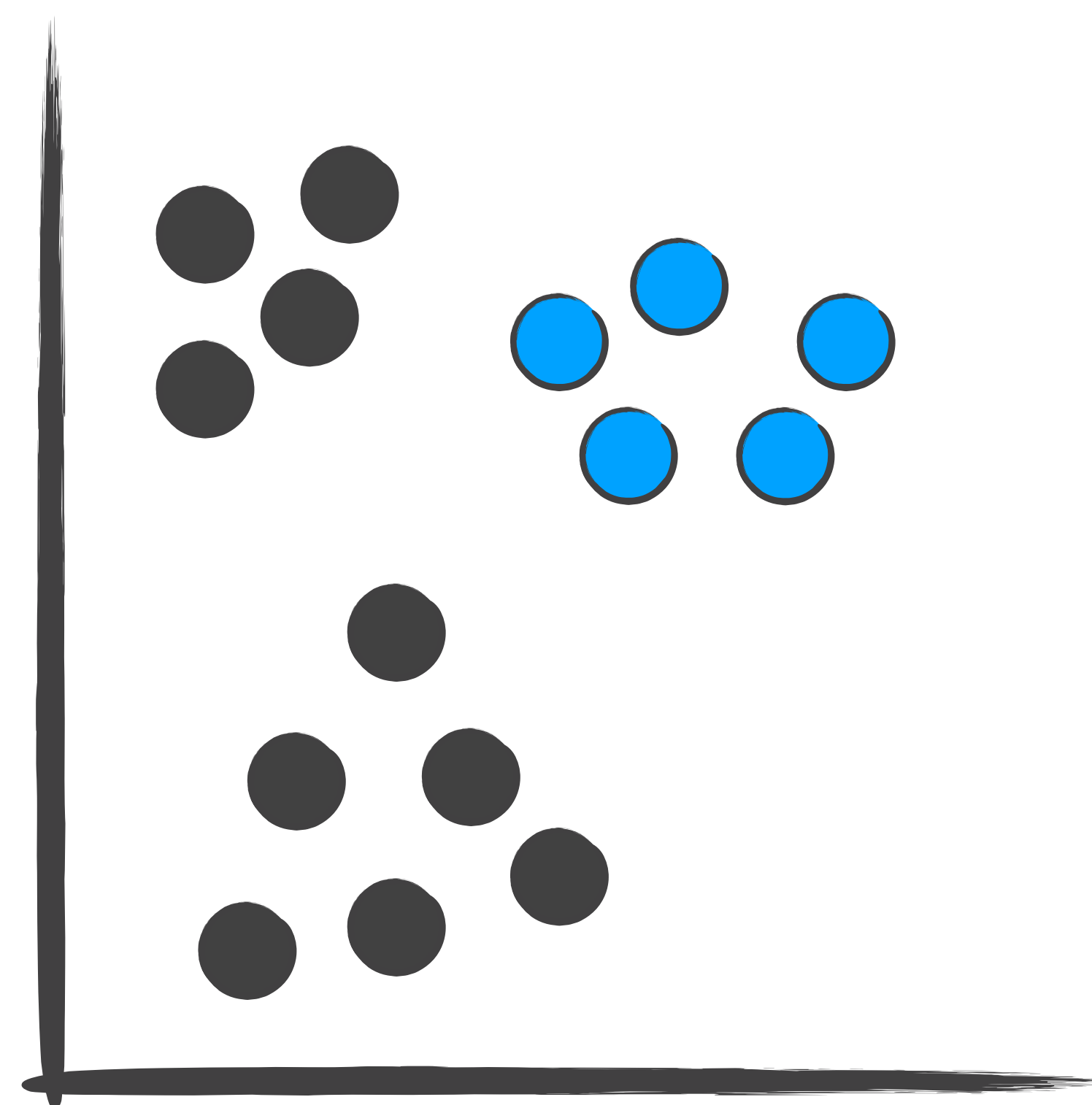
Outlier



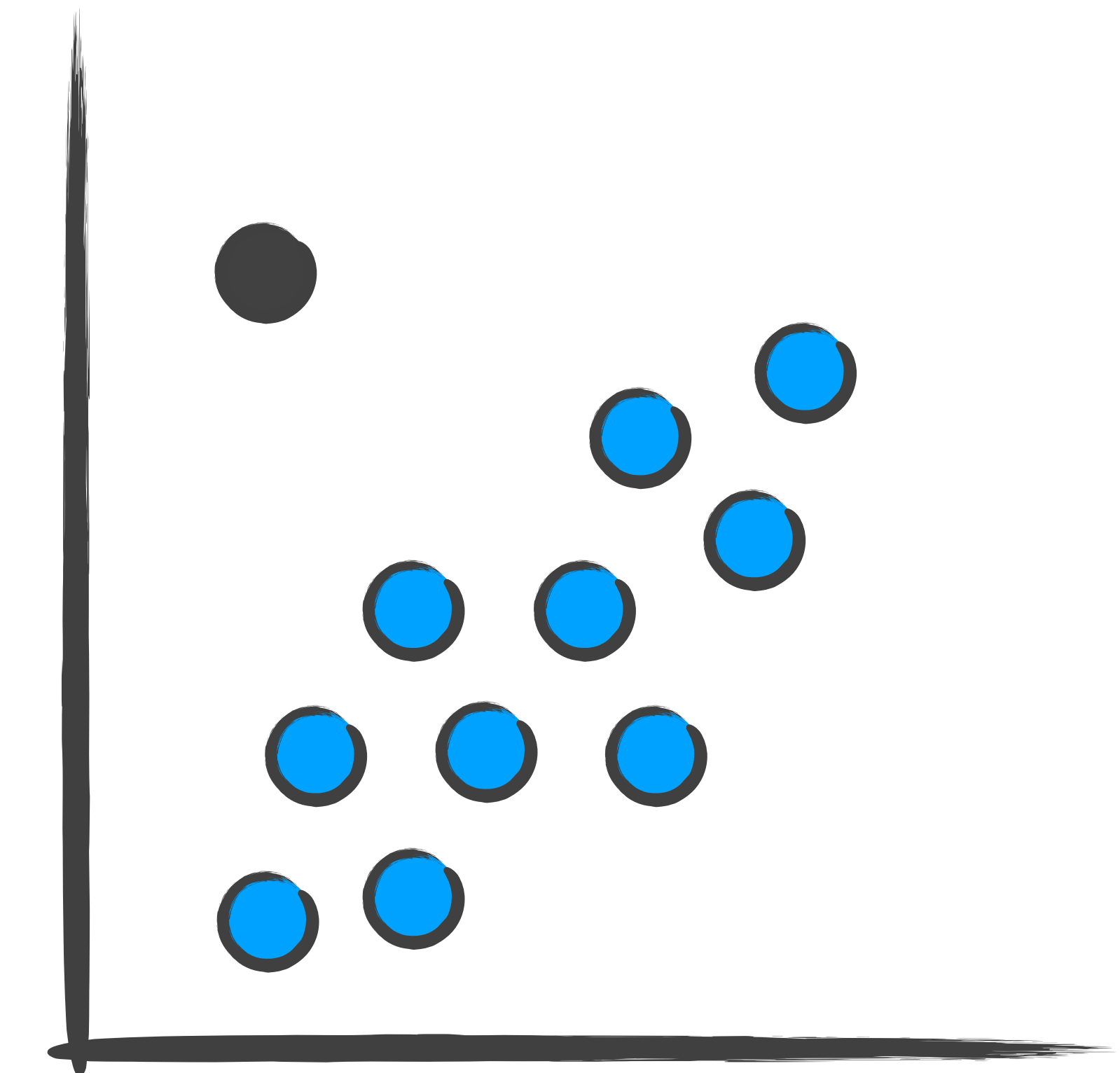
Clusters



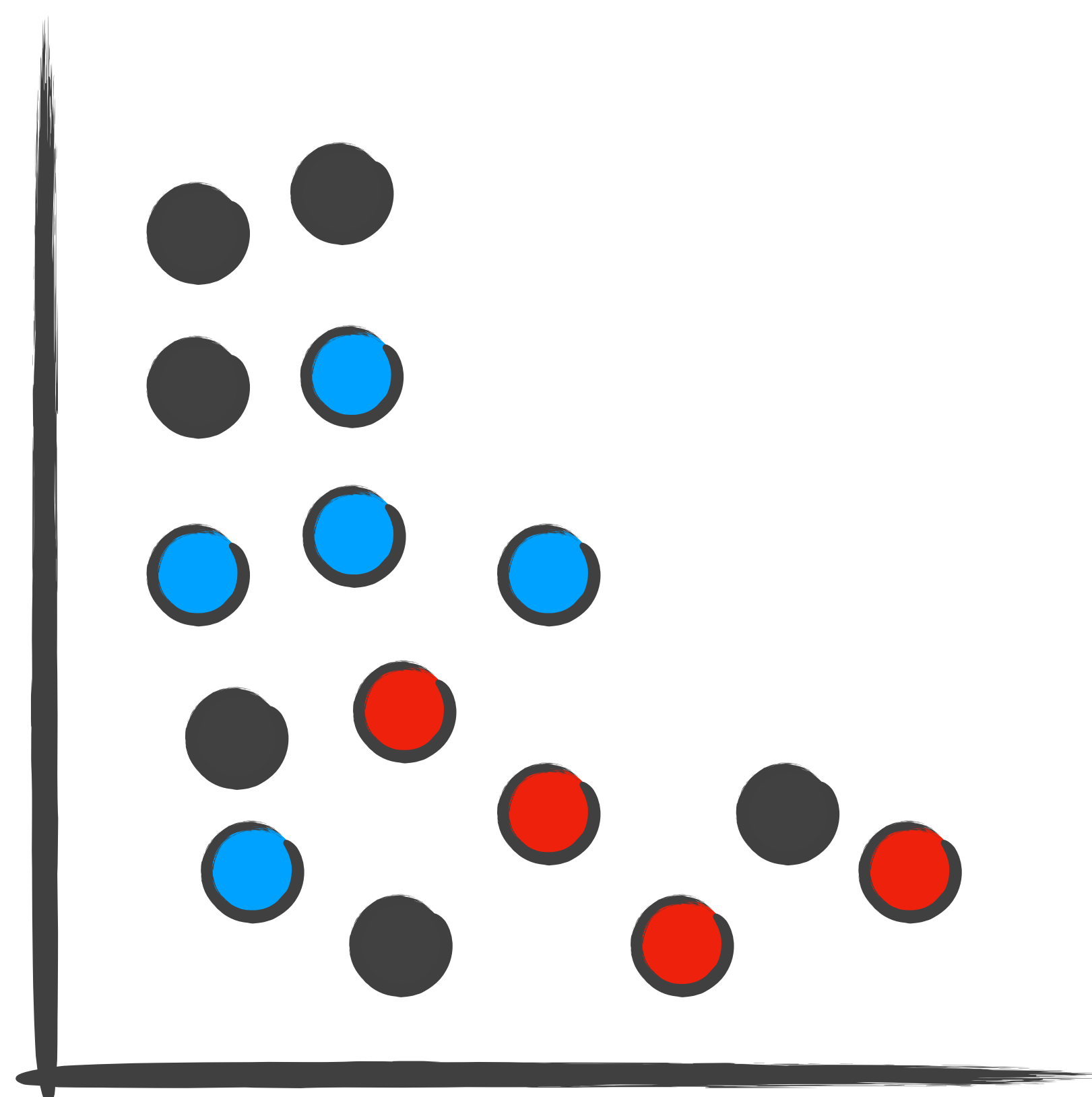
Outlier



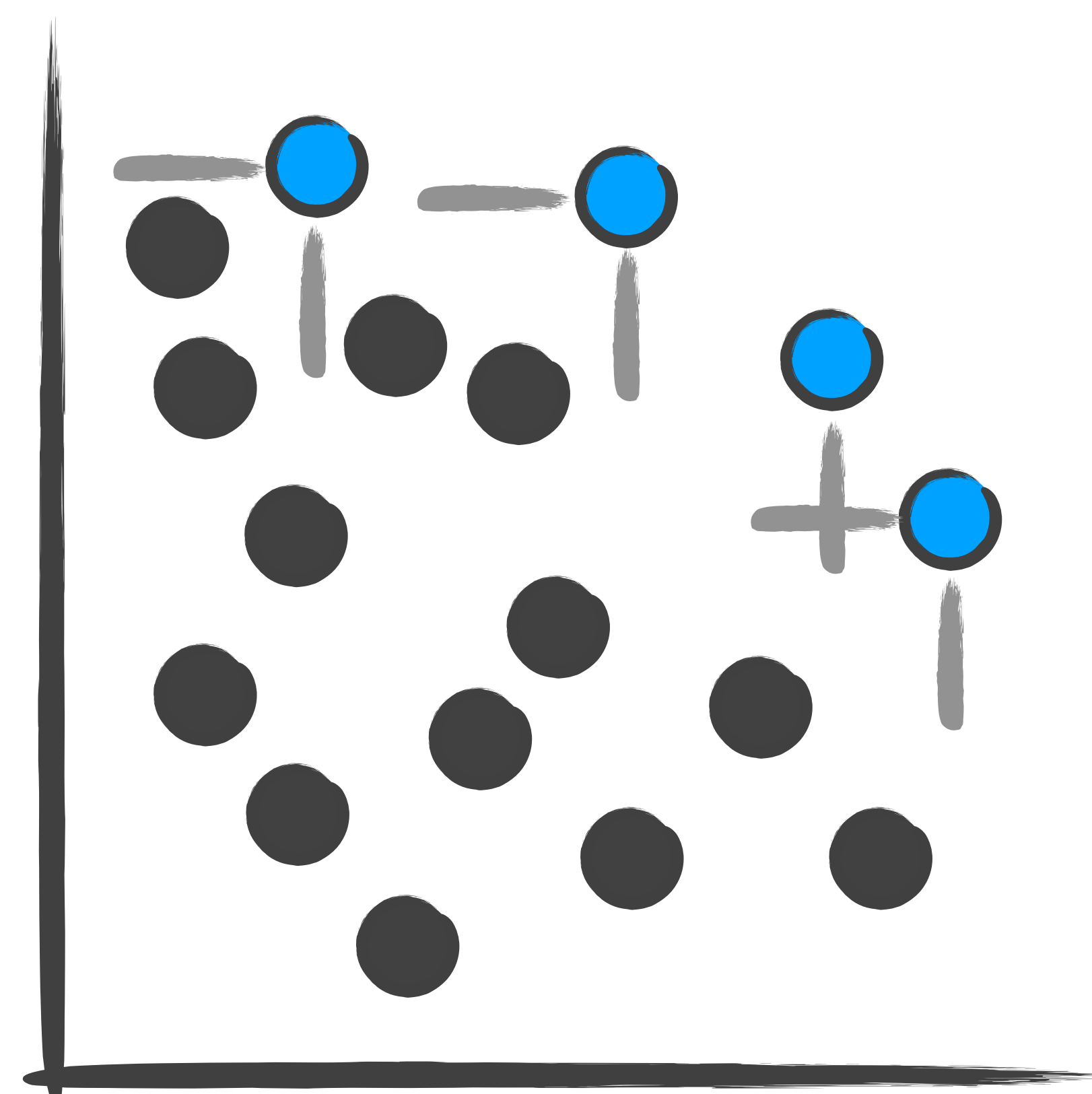
Clusters



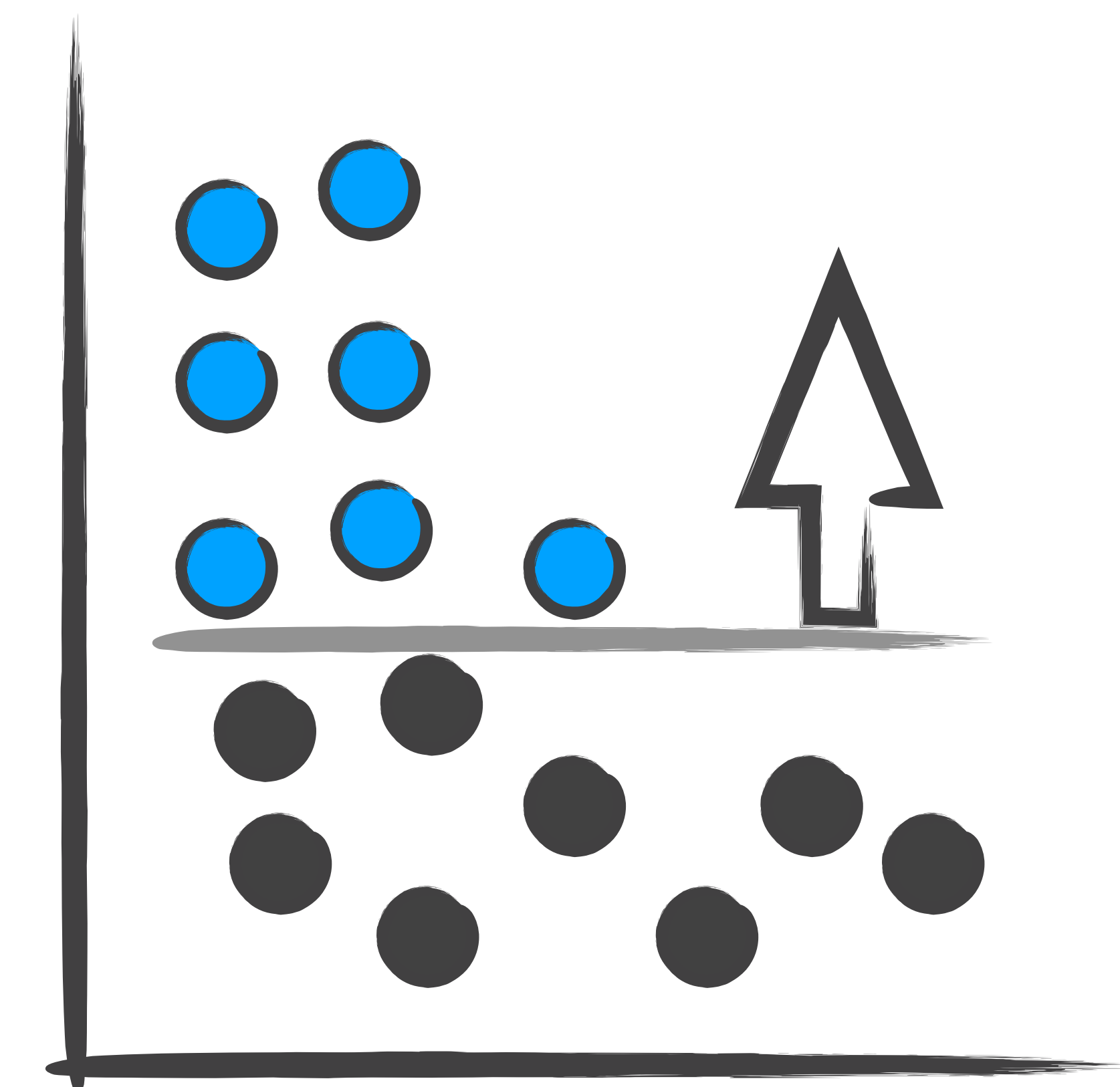
Correlation



Categories



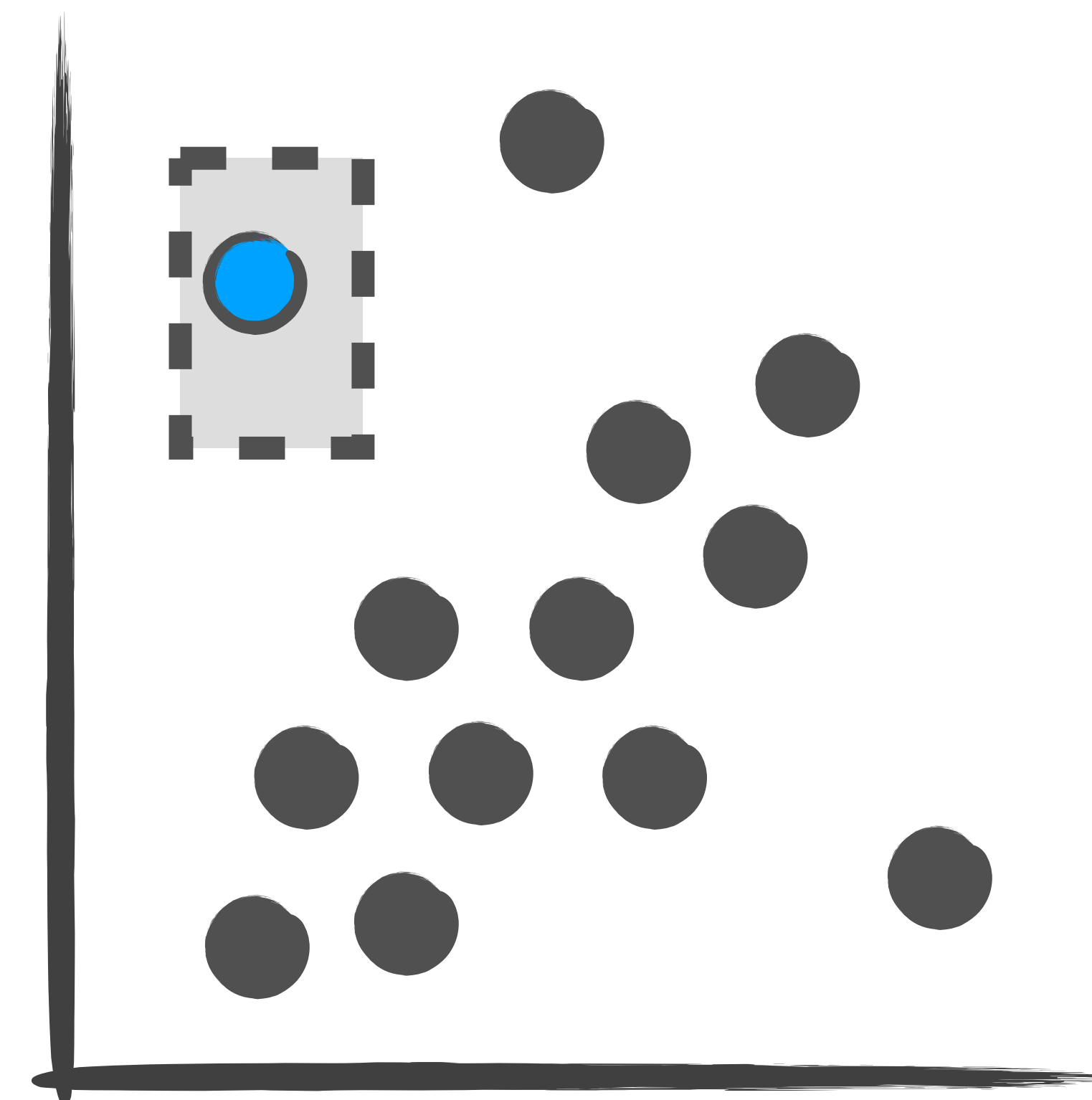
Multivariate Optimization



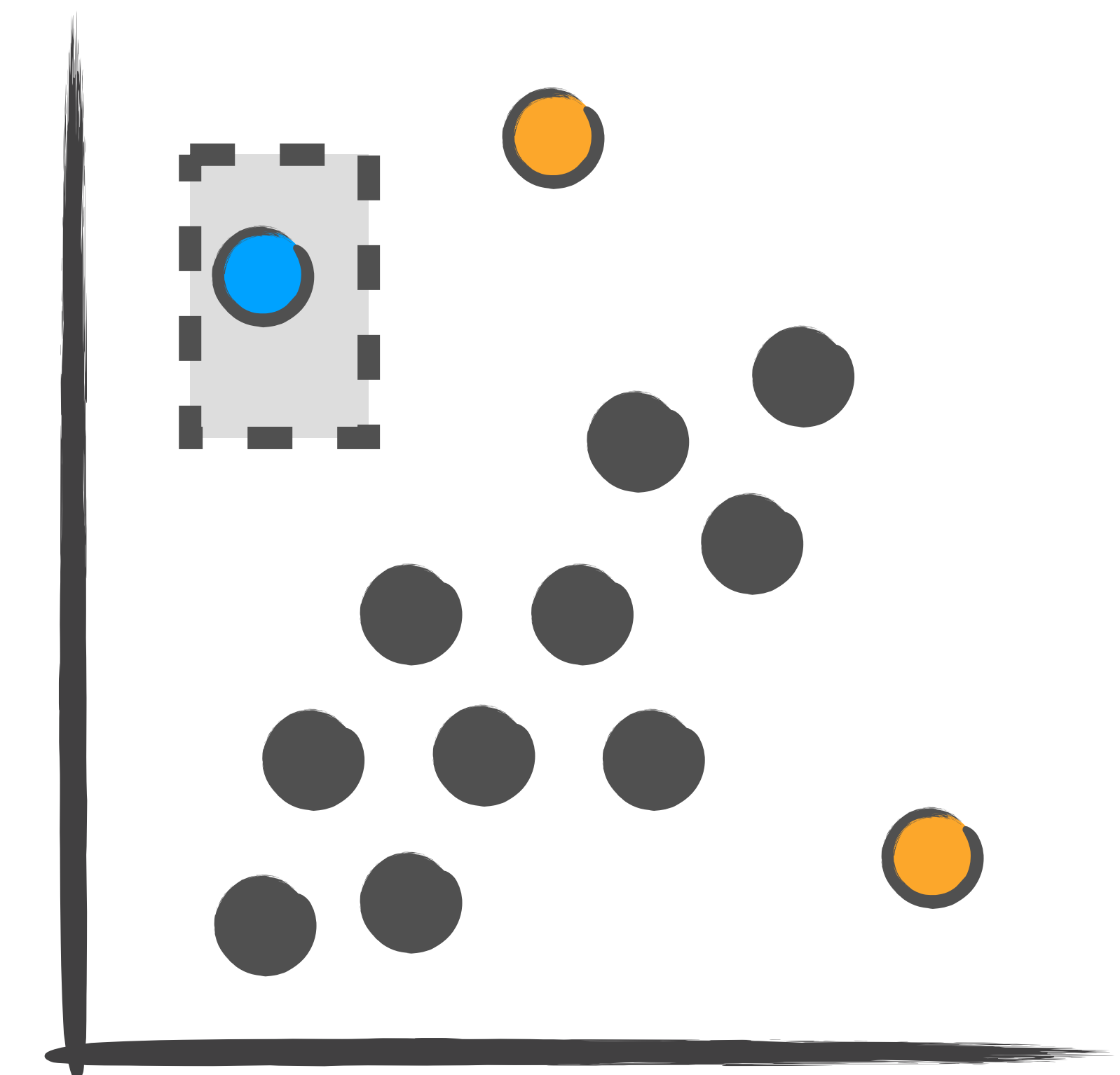
Ranges

WHY DO WE CARE?

Speed up complex
selections



Selection

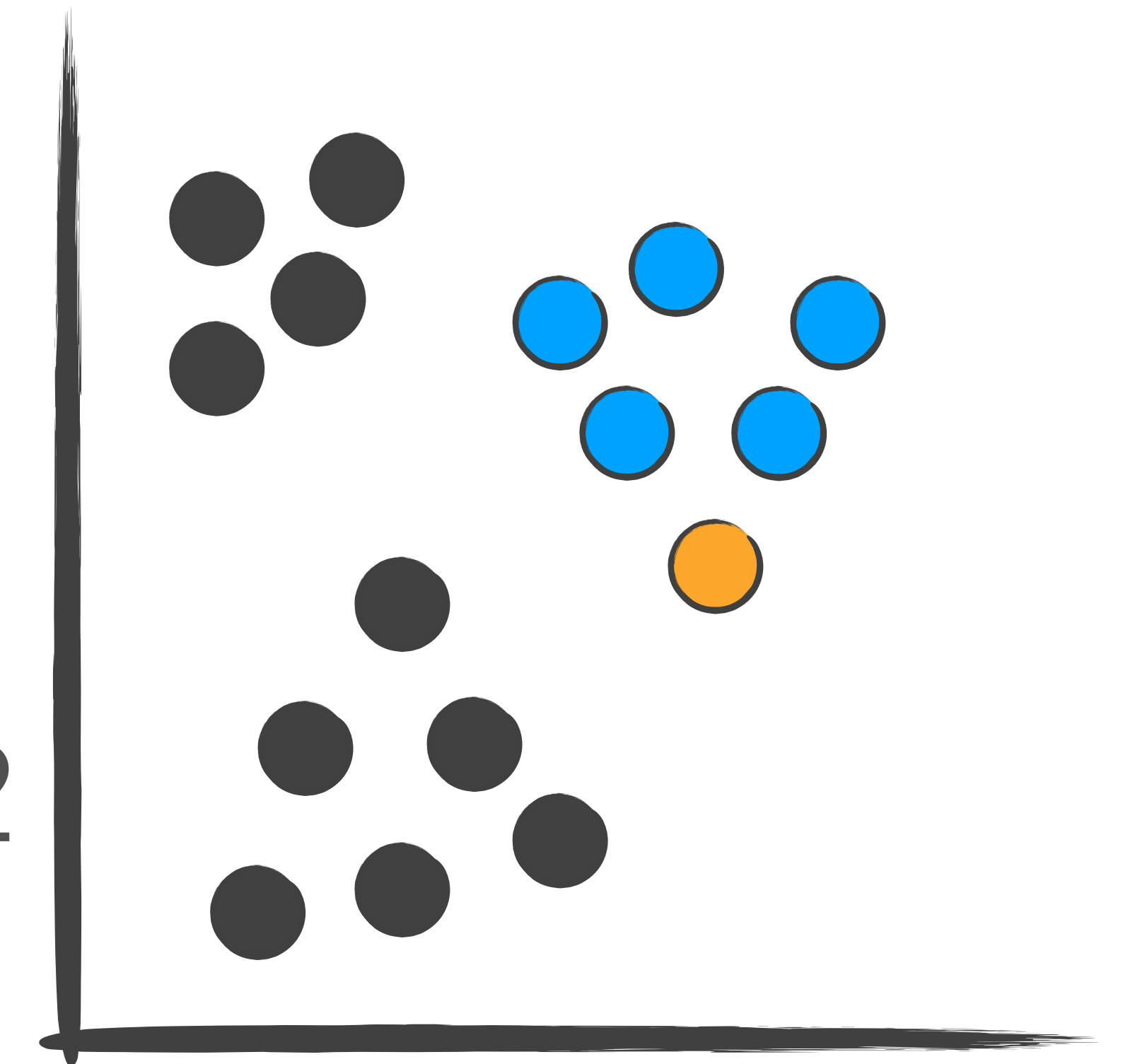


Outliers?

WHY DO WE CARE?

ID Based Selection:

Selected Elements: 7, 9, 13, 18, 22



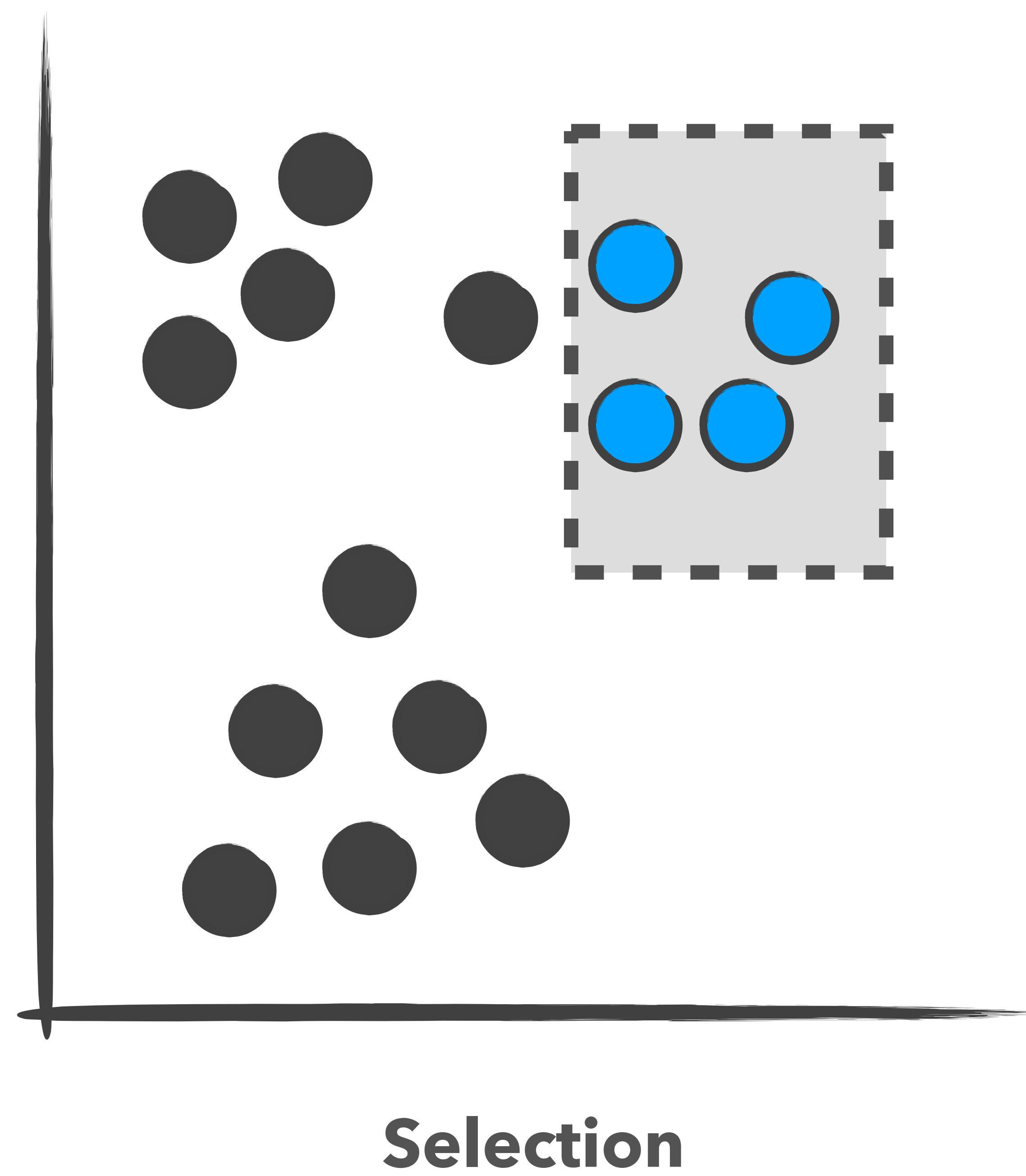
Semantic Selection:

Elements in K-Means cluster centered at $[2, 3]$

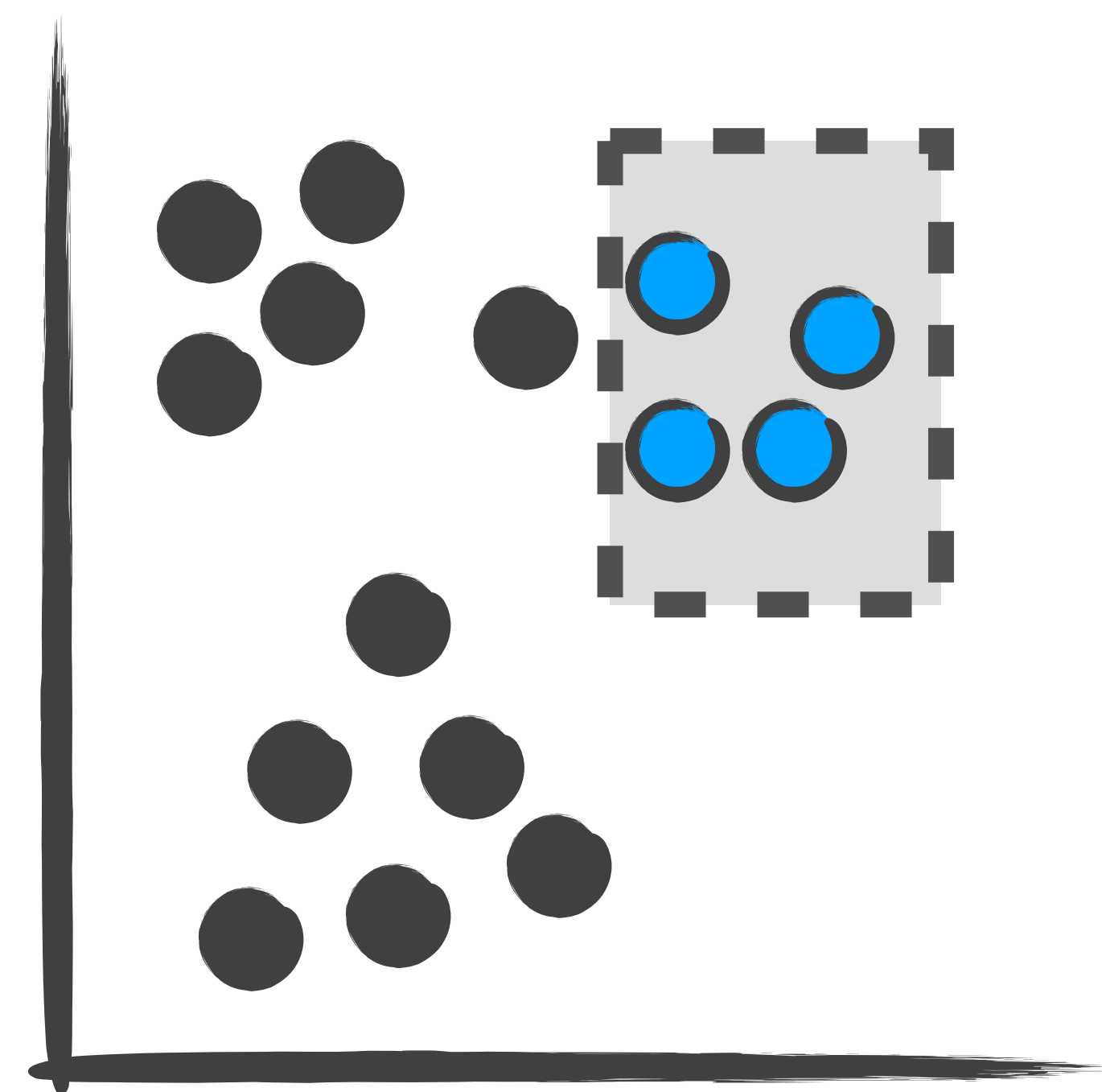
Meaningful, higher level concept:
improves reproducibility

Robust to changes and updates in dataset:
enables re-usability

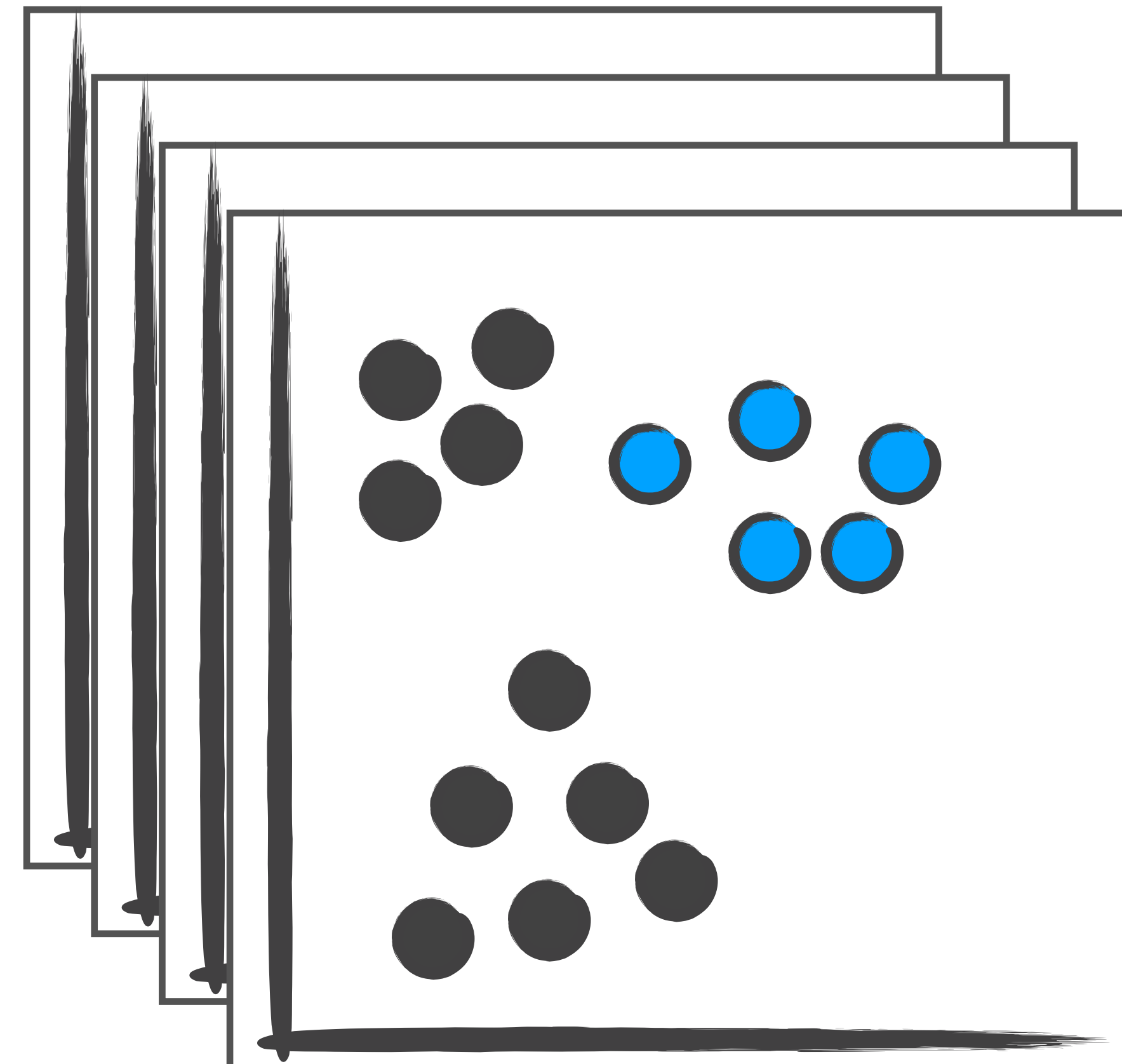
HOW DO WE INFER INTENT?



HOW DO WE INFER INTENT?



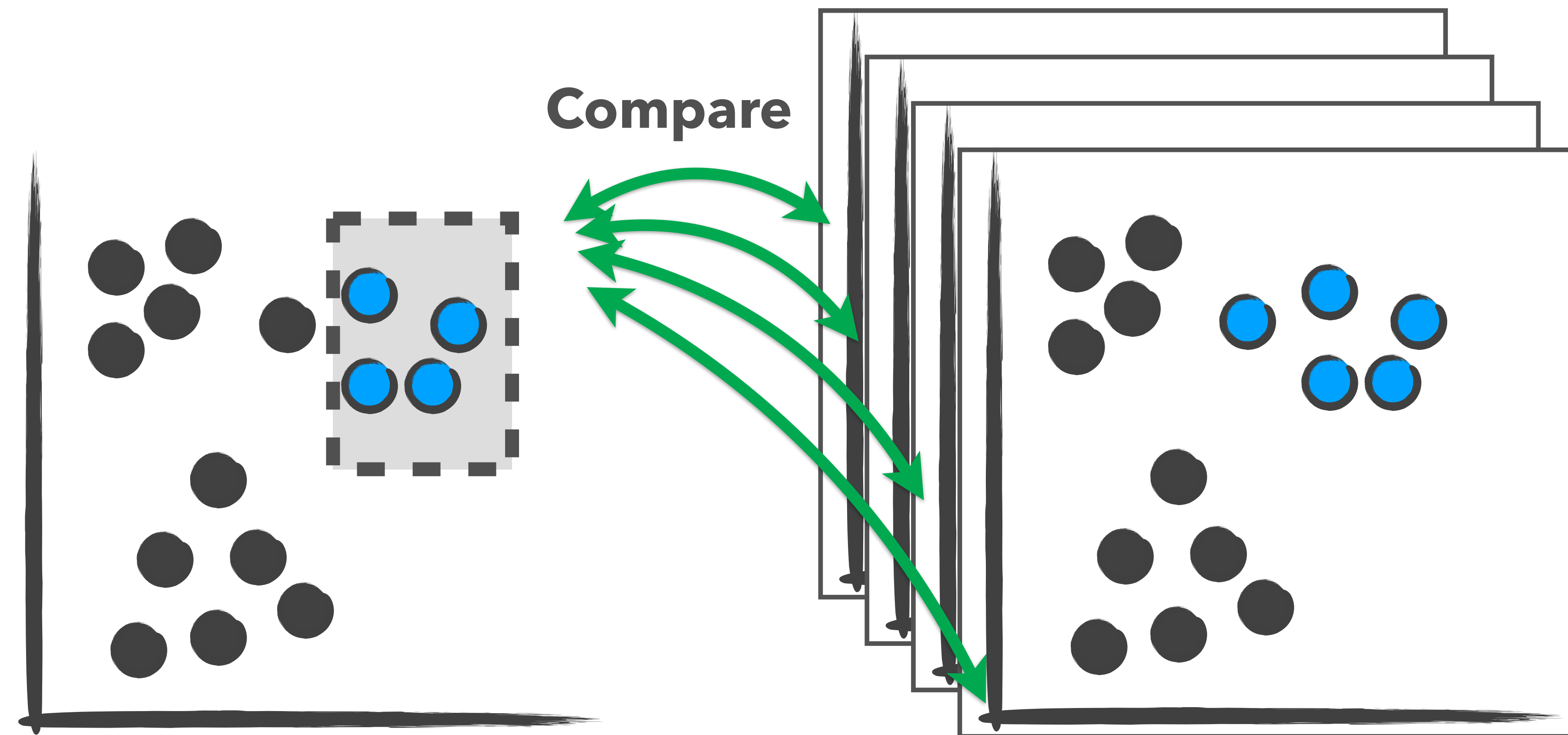
Selection



Predictions

K-Means
DBScan
Regression
Outlier Detection
Skyline
Decision Trees / Ranges
Categories

HOW DO WE INFER INTENT?



Selection

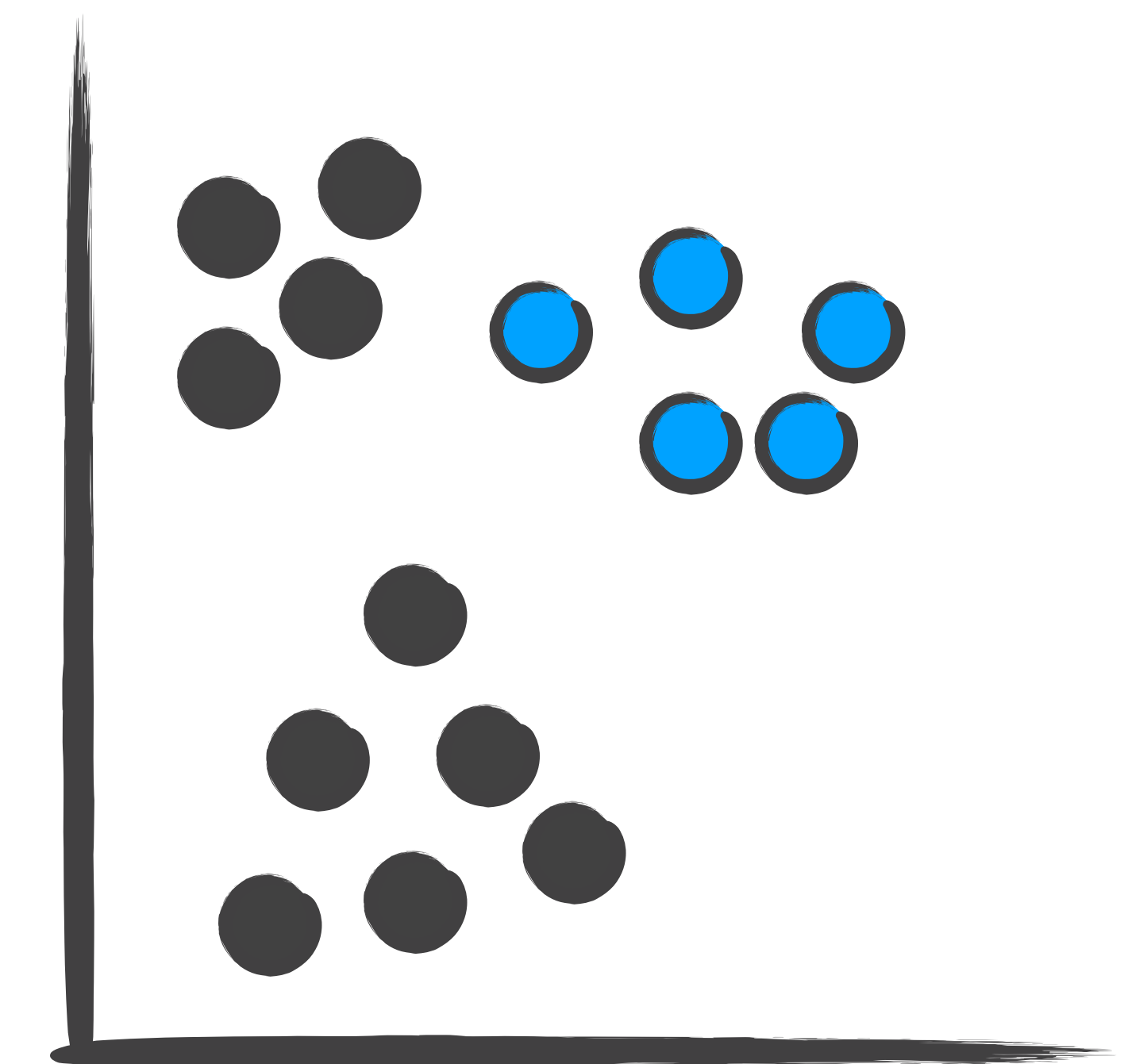
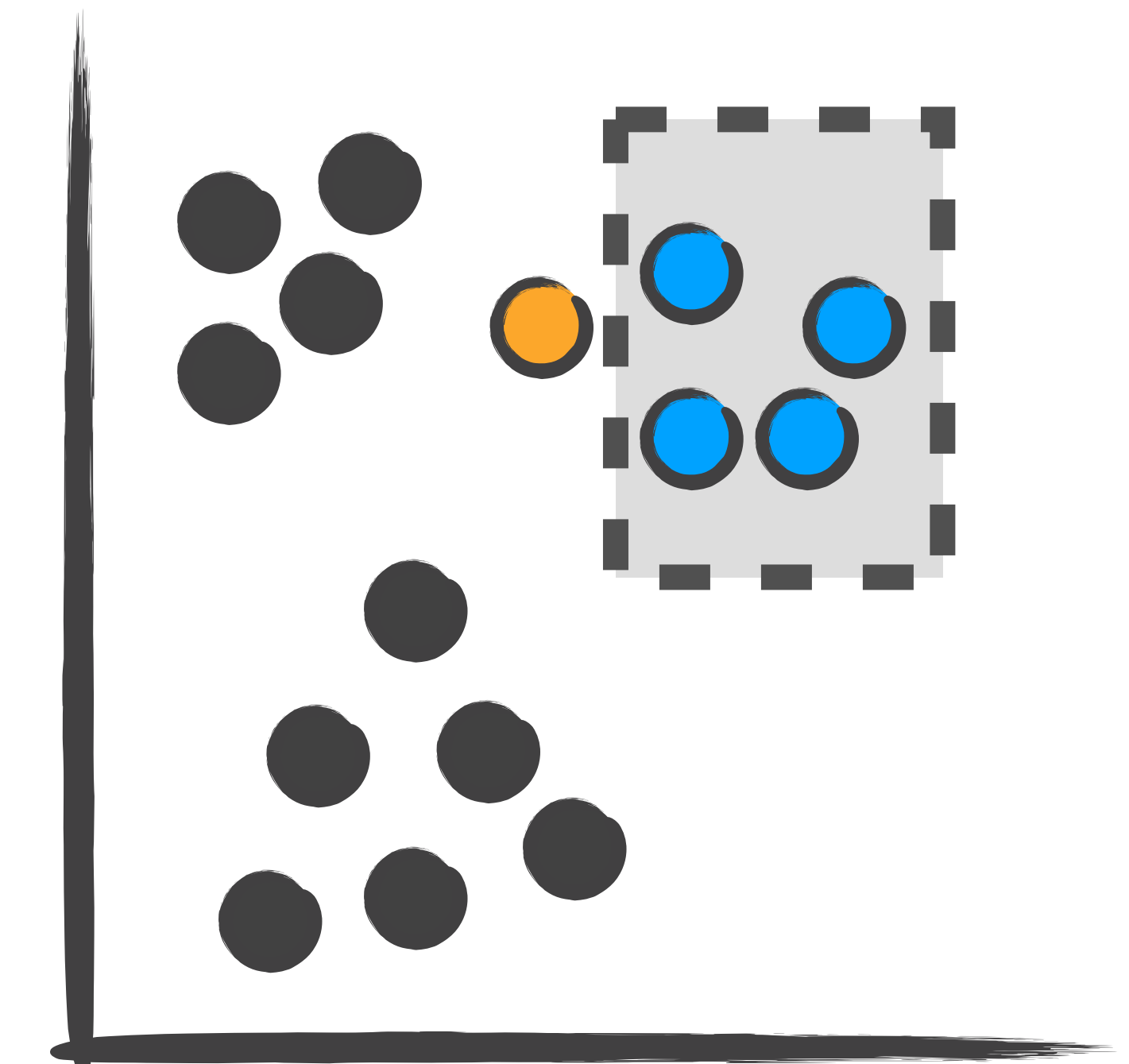
Predictions

K-Means
DBScan
Regression
Outlier Detection
Skyline
Decision Trees / Ranges
Categories

1. Range 
2. Cluster 
3. Outlier 

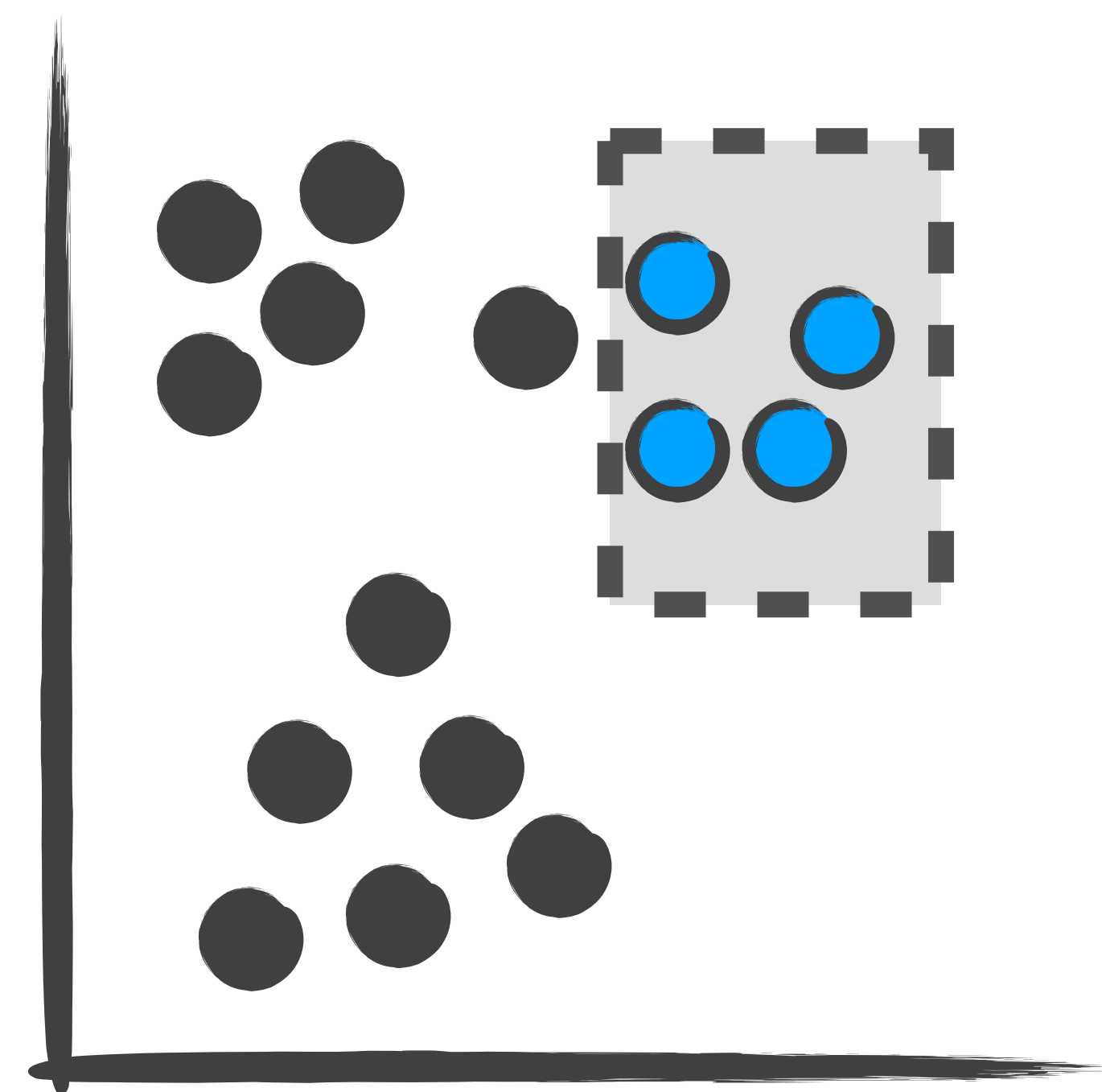
Ranking

Jaccard Distance
Naive Bayes
Classifier
Heuristic
Measures

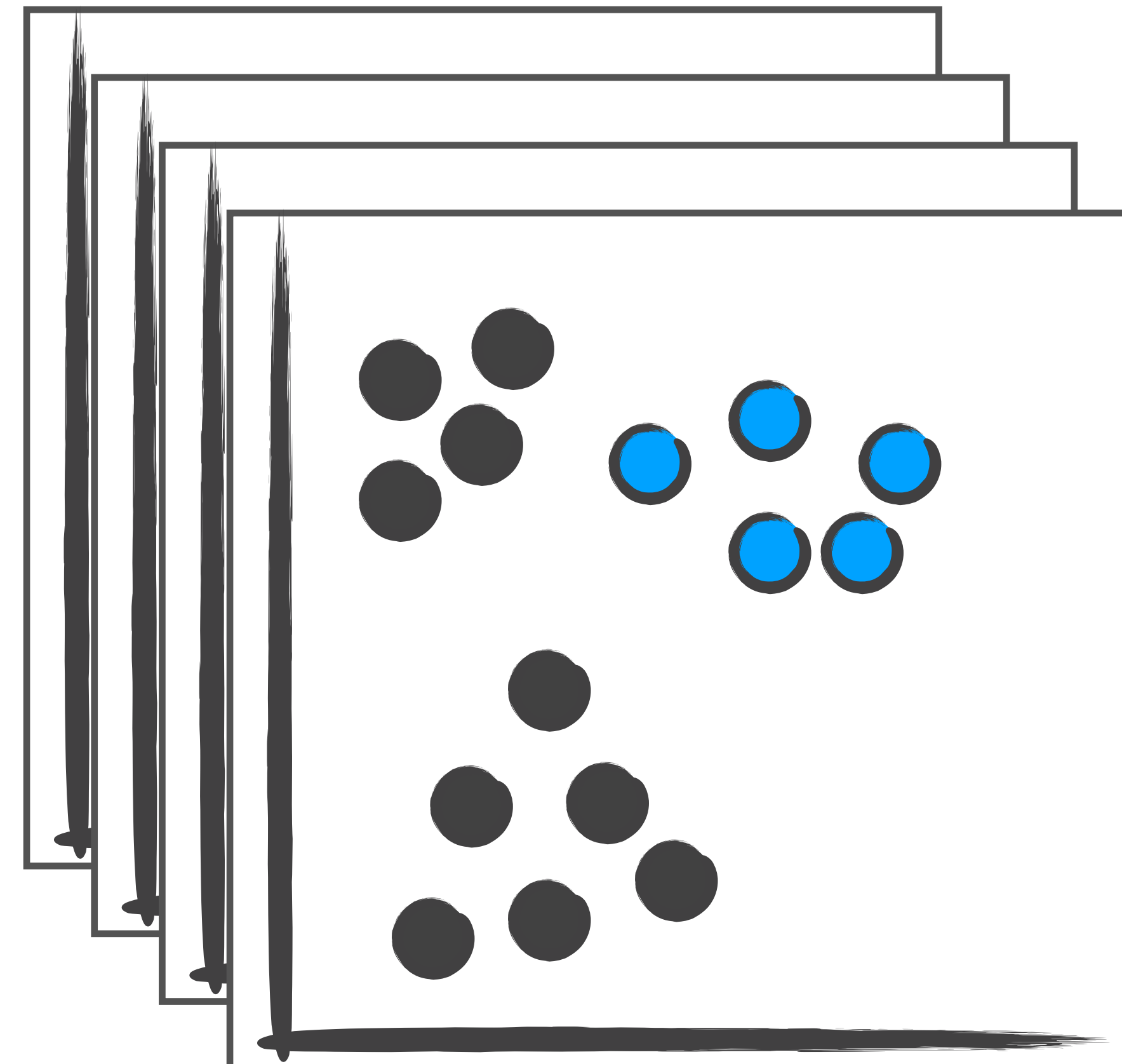


$$J(S, C) = \frac{|S \cap C|}{|S \cup C|}$$

HOW DO WE INFER INTENT?



Selection



Predictions

K-Means
DBScan
Regression
Outlier Detection
Skyline
Decision Trees / Ranges
Categories

1. Range 
2. Cluster 
3. Outlier 

I think this cluster...

Ranking

Jaccard Distance
Naive Bayes
Classifier
Heuristic
Measures

**Confirming Intent
& Annotation**

Clusters ▾

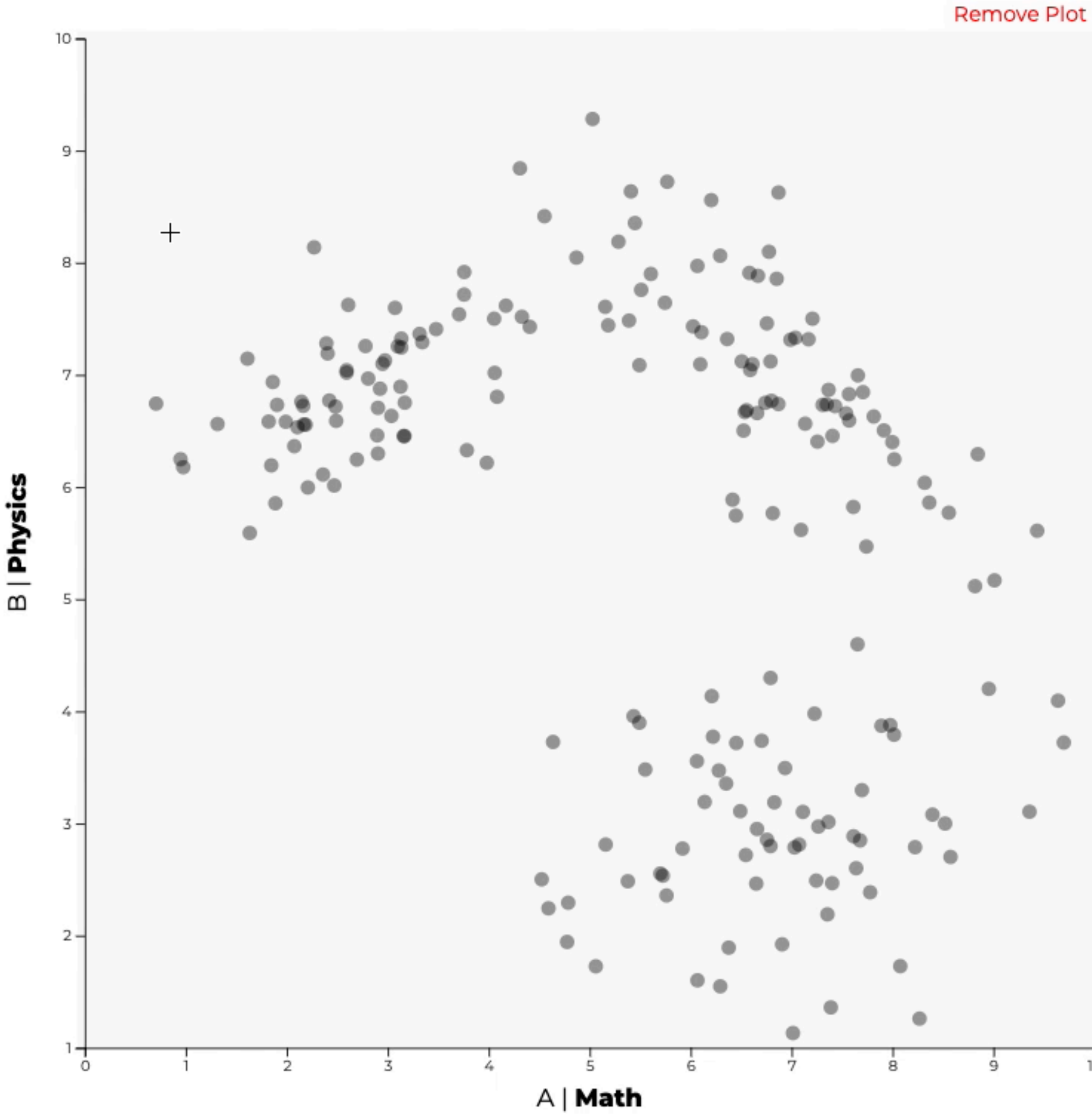
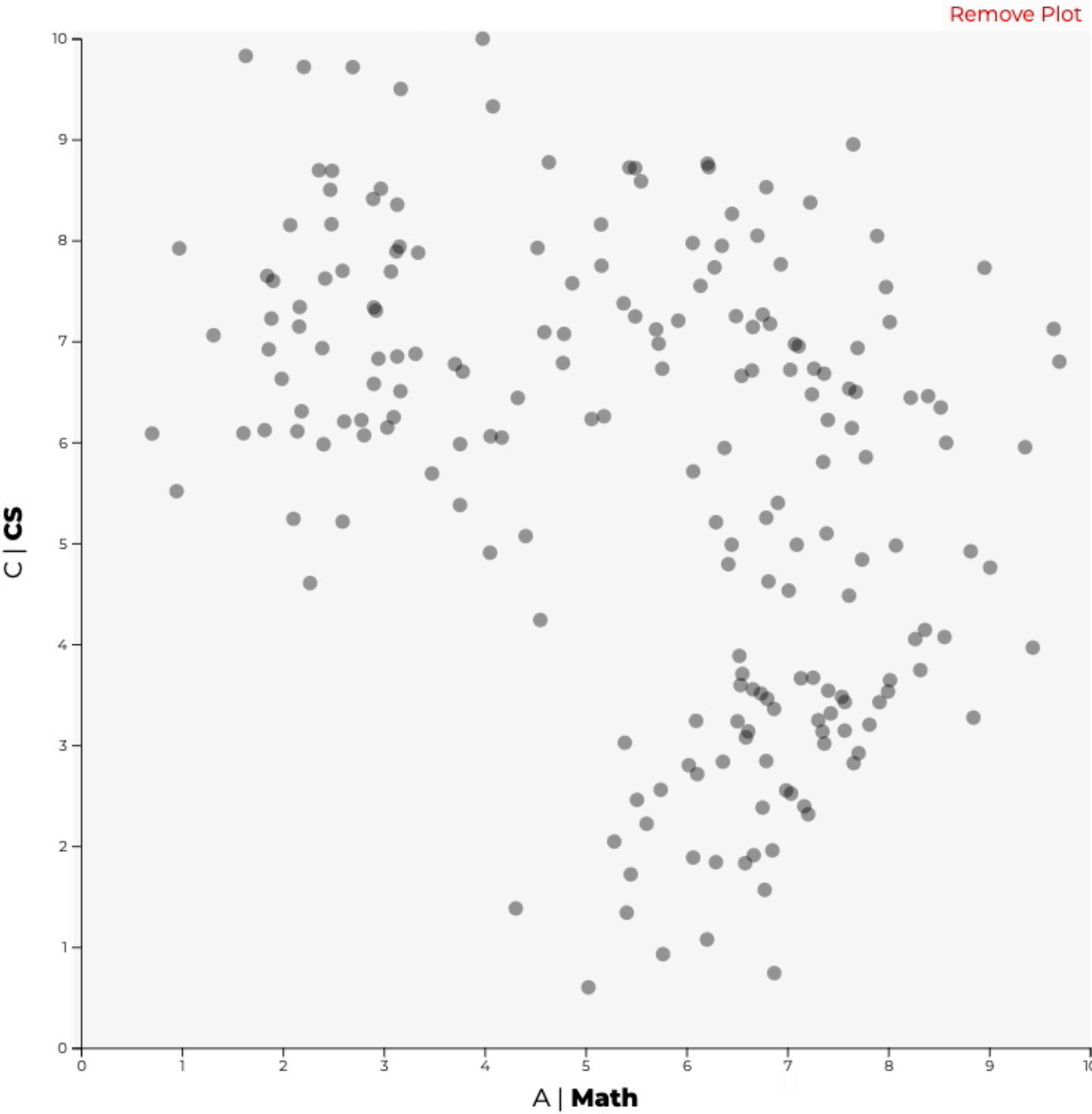
+ Add plot

☐ Show Categories

☒ Union

Invert Selections

Clear Selections



Visualization and Selection

Intent

Please interact

Annotate

Predictions

Time required: 0.01 seconds

Selections

0
UNION

0
INTERSECTION

0
INDIVIDUAL

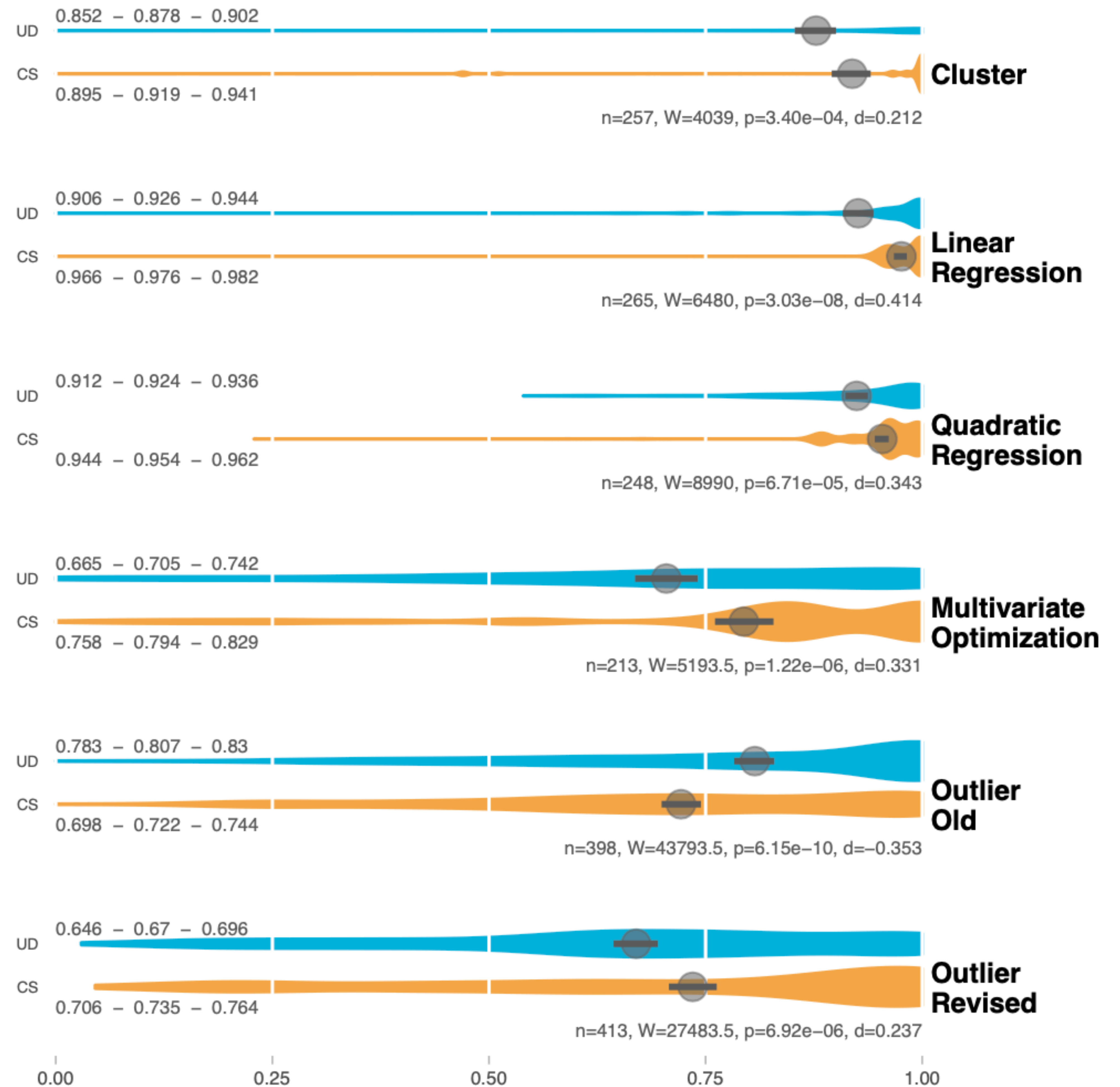
0
TOTAL

Annotation of Intent and Predictions

VALIDATION

User Driven (UD)

Computer Supported
(CD)

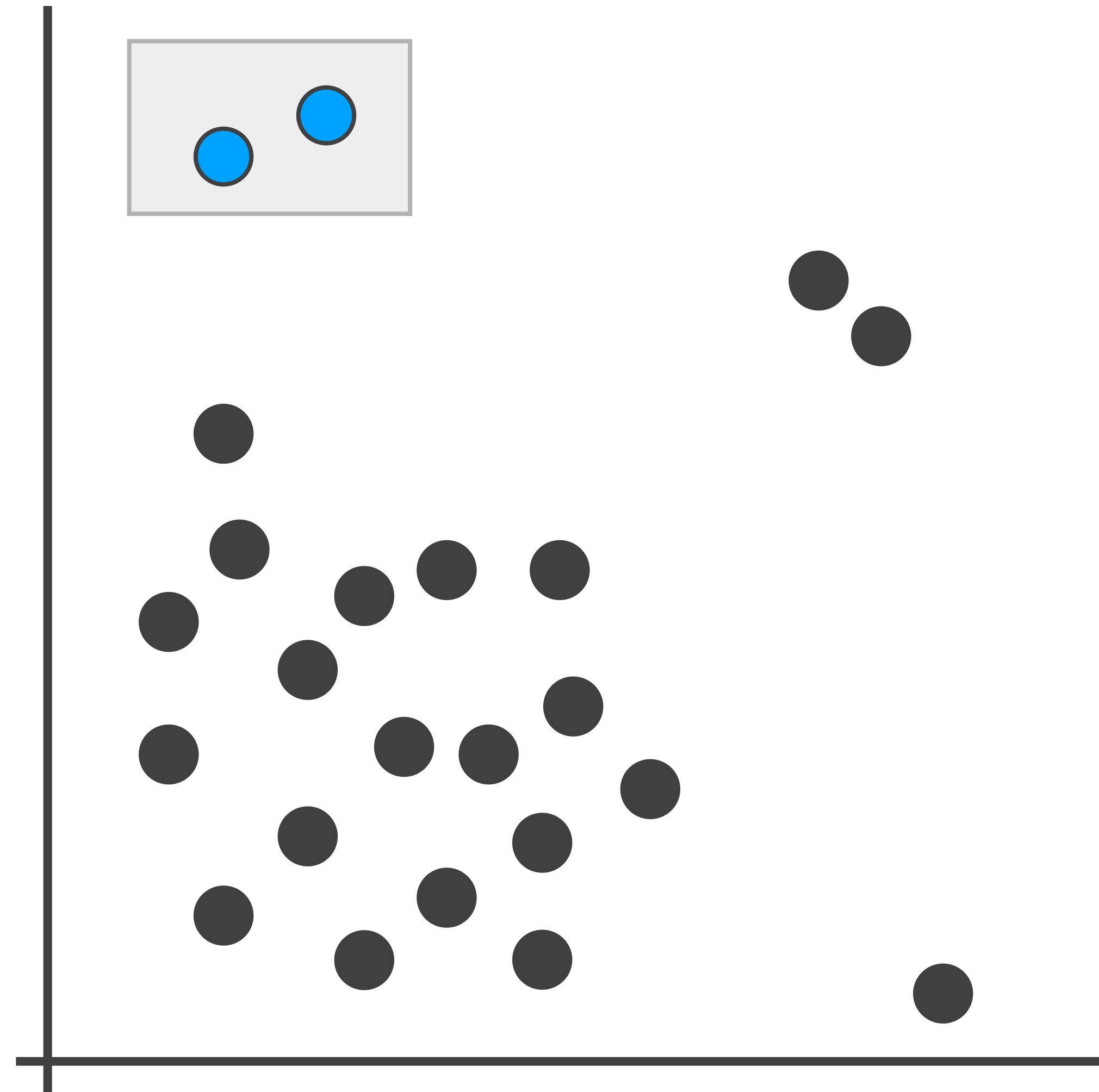


WORKFLOWS

**Based on semantic selections, we
can create reusable workflows!**

Pros and Cons

Interactive Visual Analysis

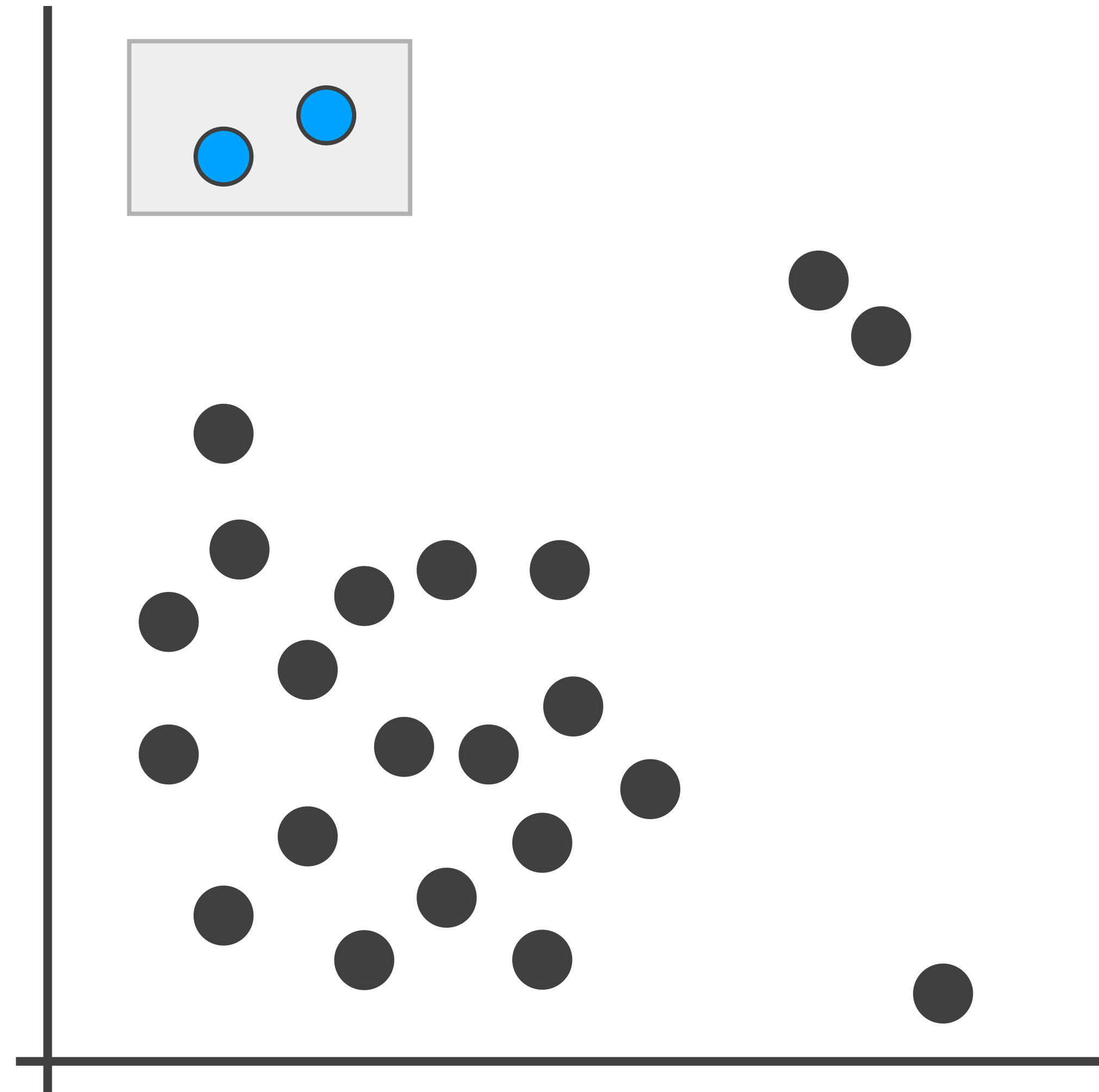


Computational Analysis



Pros and Cons

Interactive Visual Analysis



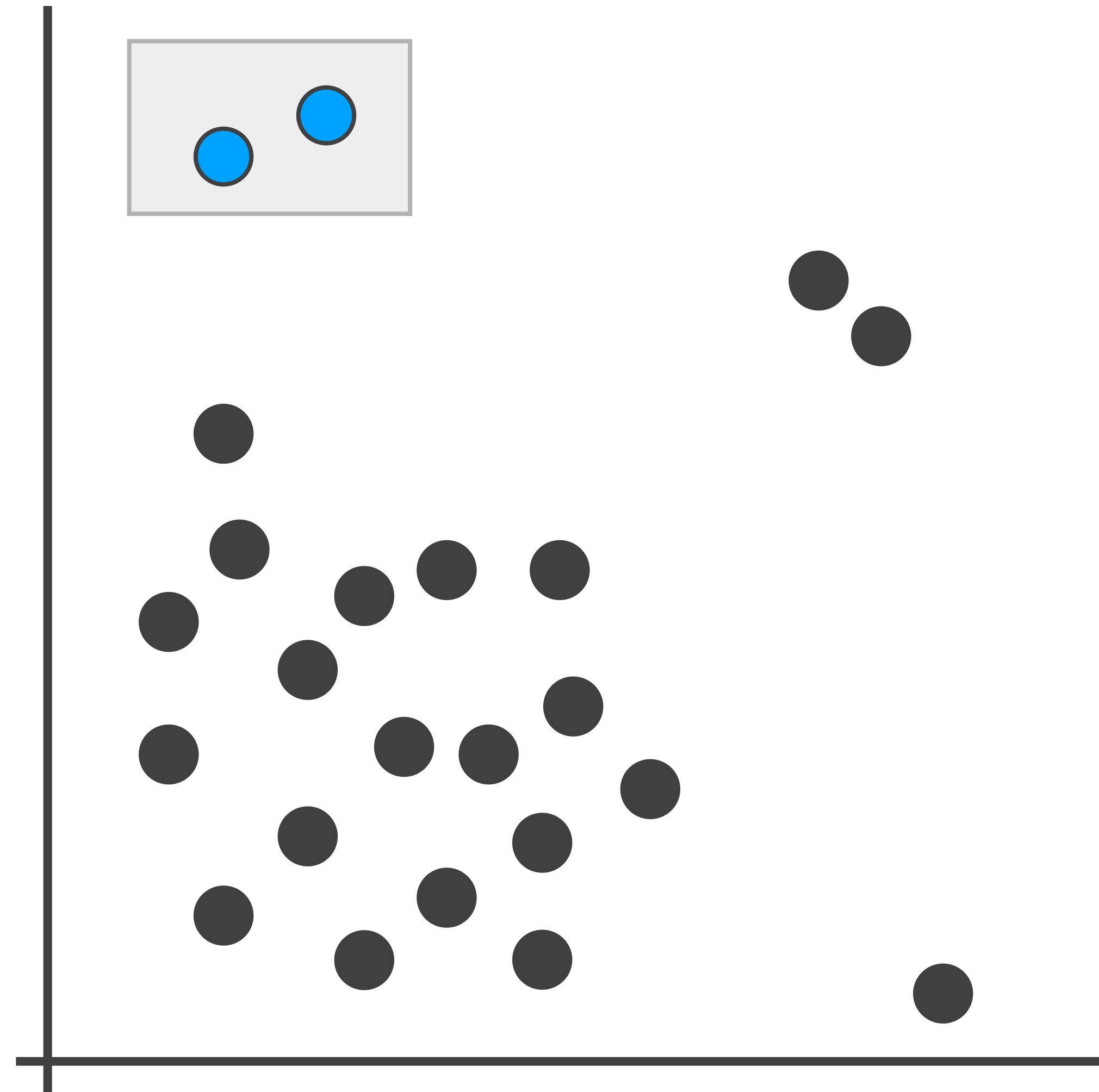
Intuitive and Fast
Uses Human Perceptual Capabilities

Computational Analysis



Pros and Cons

Interactive Visual Analysis



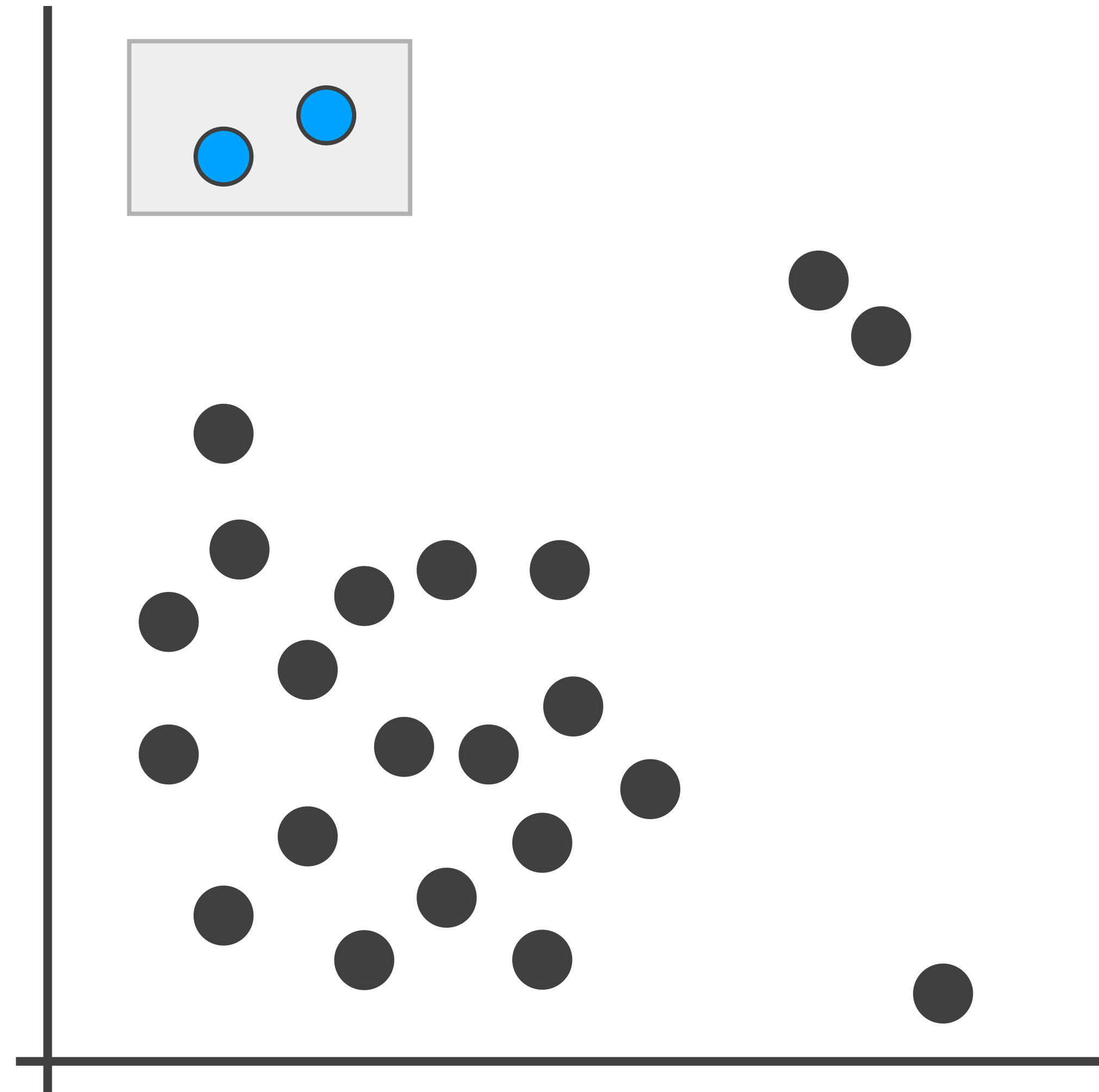
**Need to redo the analysis
when datasets update**

Computational Analysis



Pros and Cons

Interactive Visual Analysis



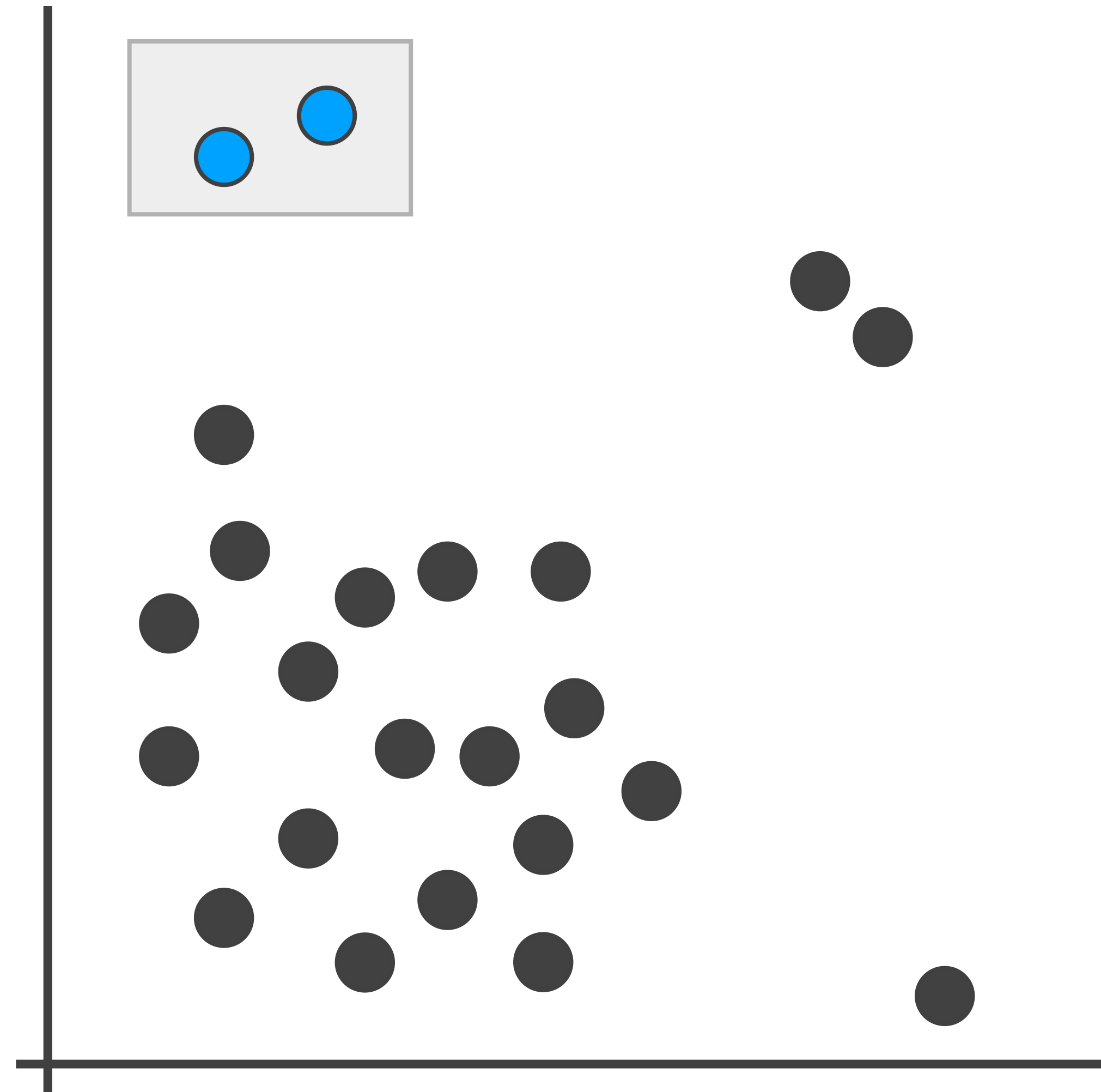
Computational Analysis



**Flexible, powerful,
reusable.**

Pros and Cons

Interactive Visual Analysis



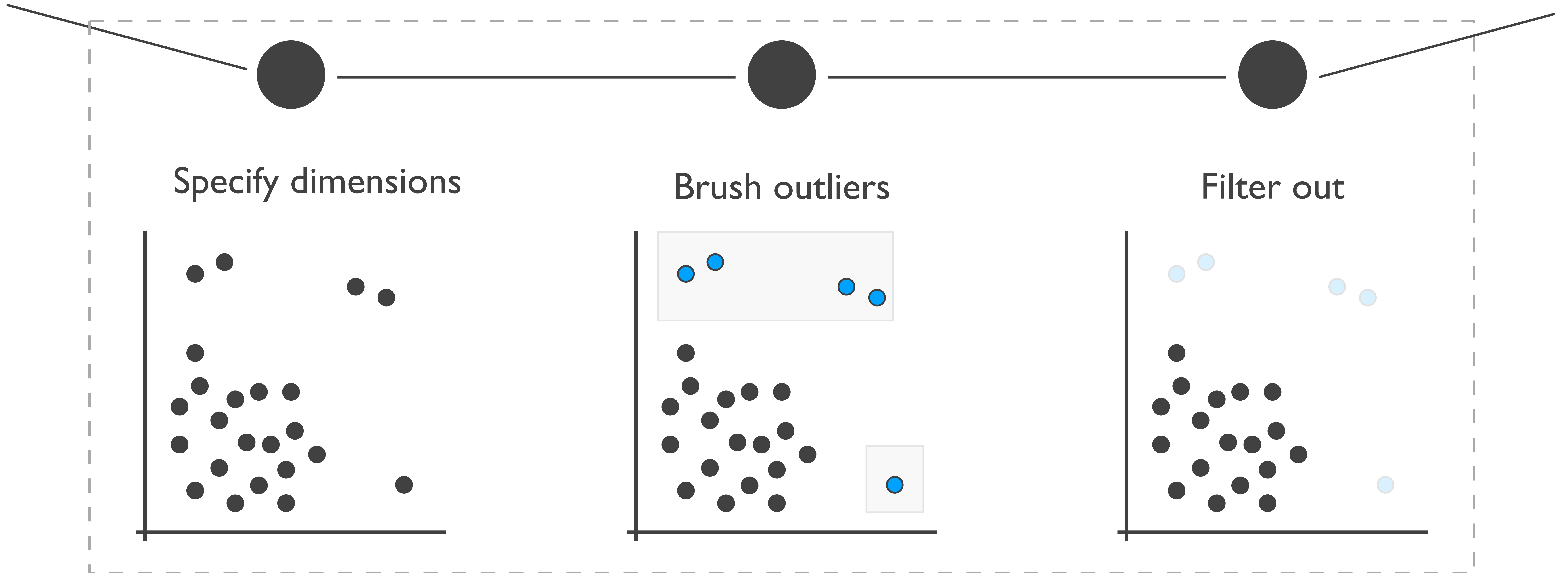
Computational Analysis



Time Consuming
Require analyst can and
want to program

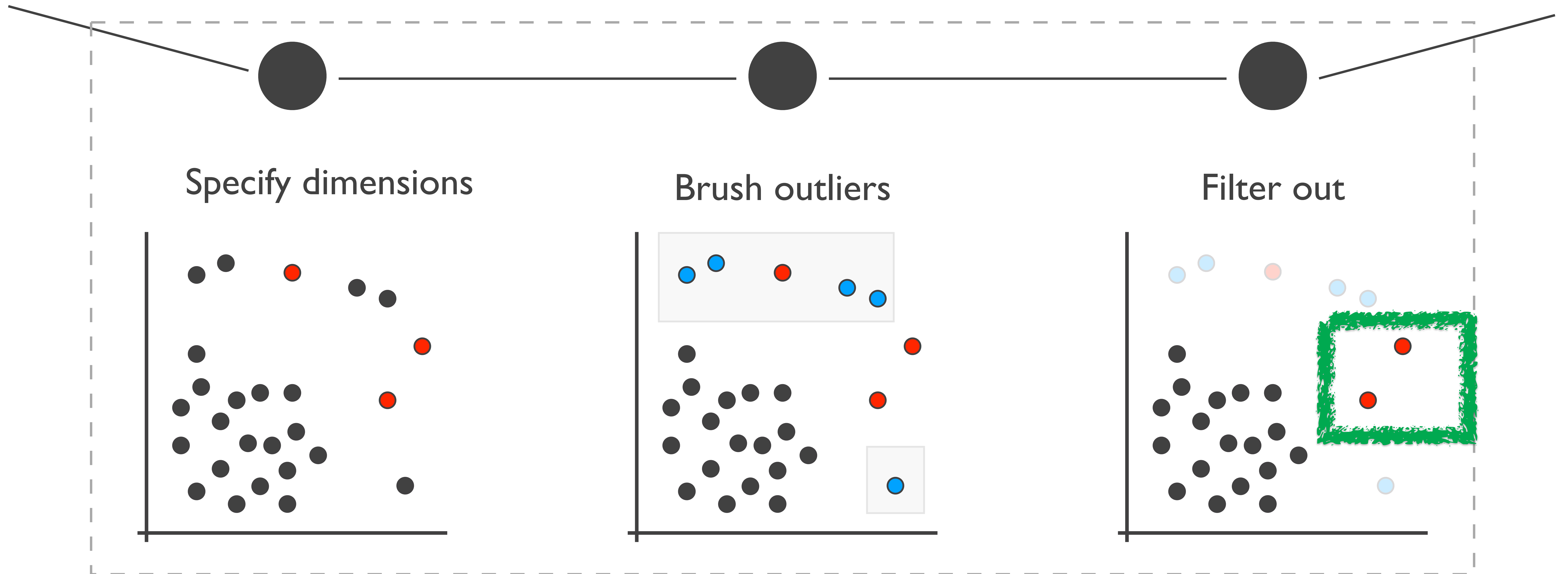
Difficult to see what's
in the data

IDEA: CAPTURING WORKFLOWS

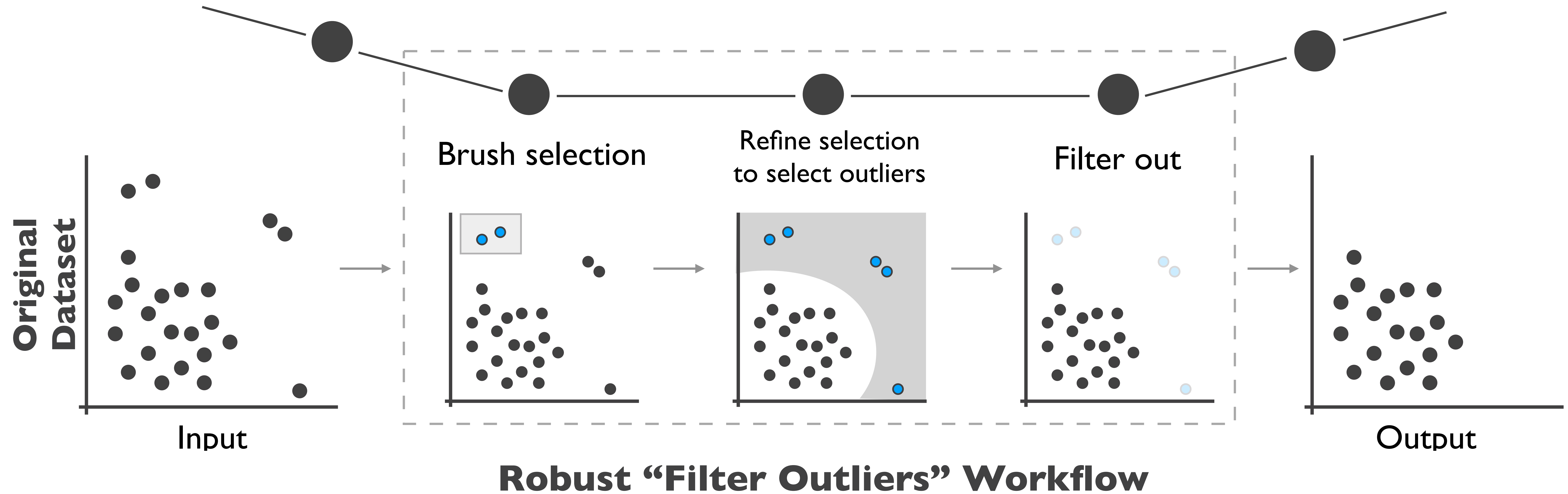


“Filter Outliers” Workflow

USING WORKFLOW ON UPDATED DATA



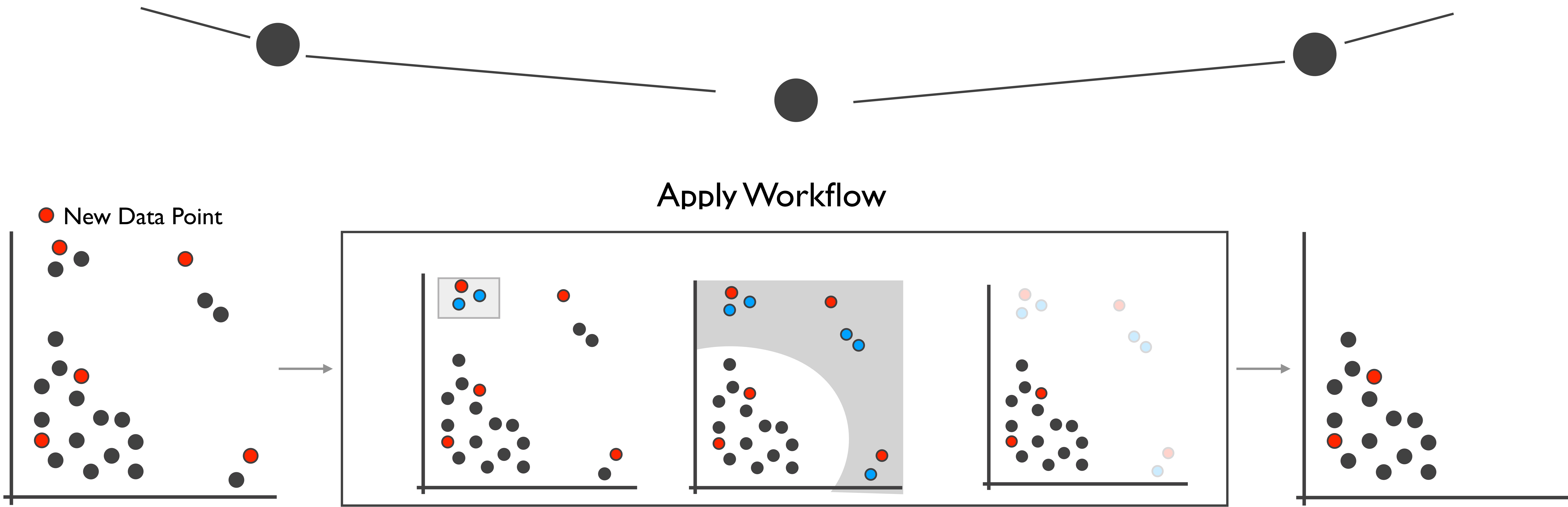
CAPTURING SEMANTICS OF WORKFLOWS



REUSING WORKFLOWS ON **UPDATED** DATA

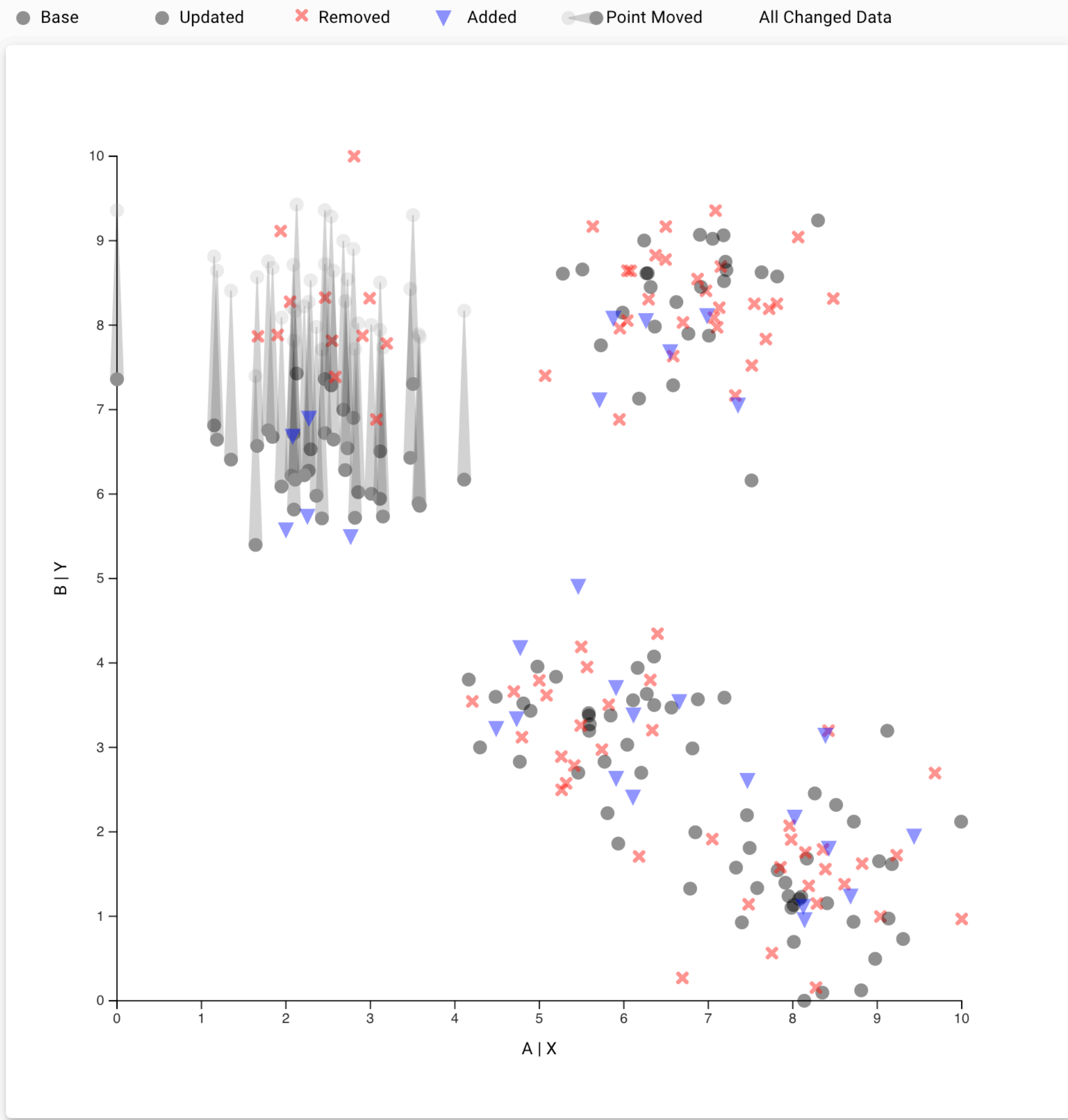
Updated Dataset

Apply Workflow



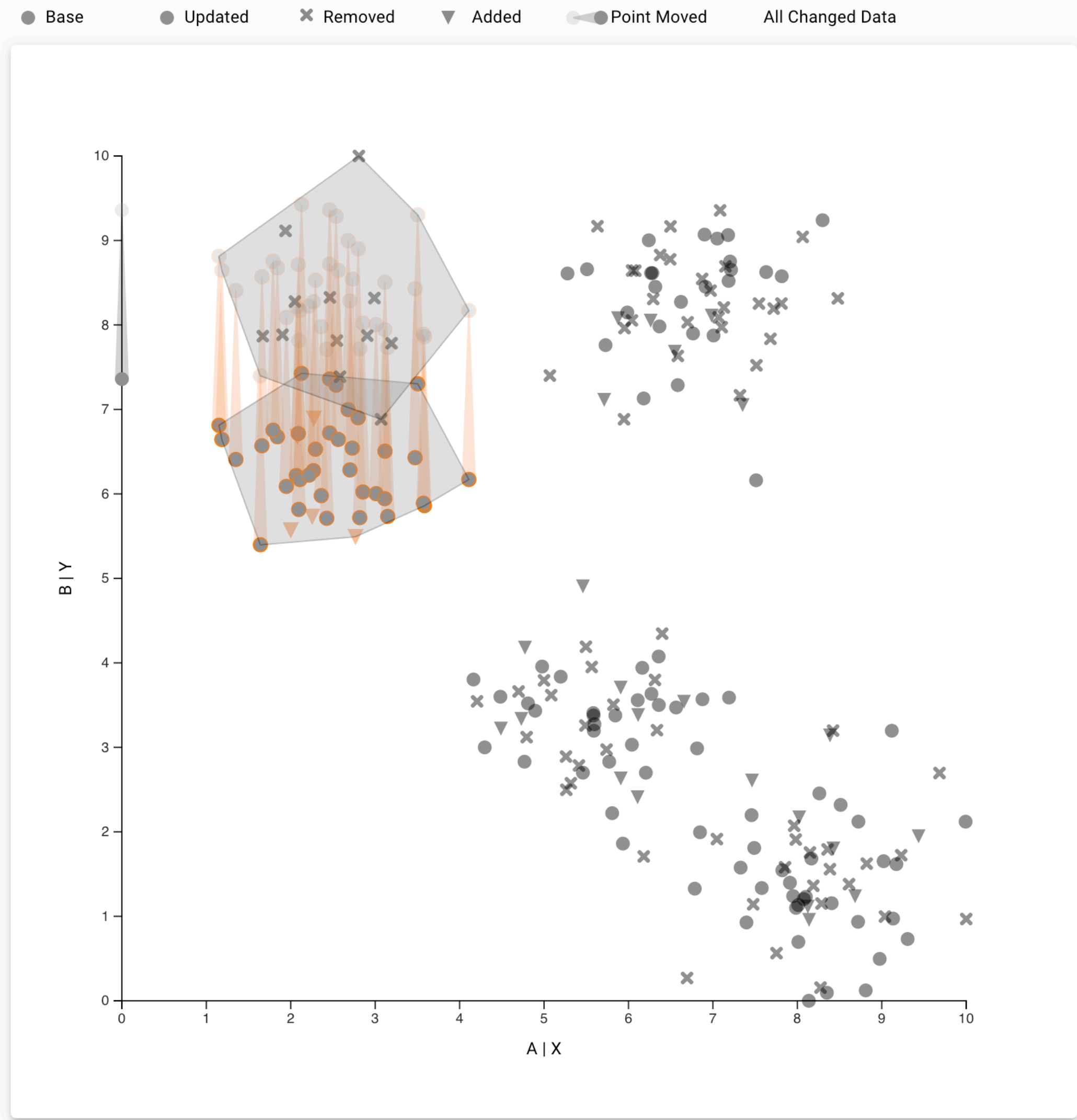
REUSING SELECTIONS ON UPDATED DATASETS

cluster_simple_v1 cluster_simple_v2 cluster_simple_v3 cluster_simple_v6 cluster_simple_v4 cluster_simple_v5 ☒ Compare Datasets



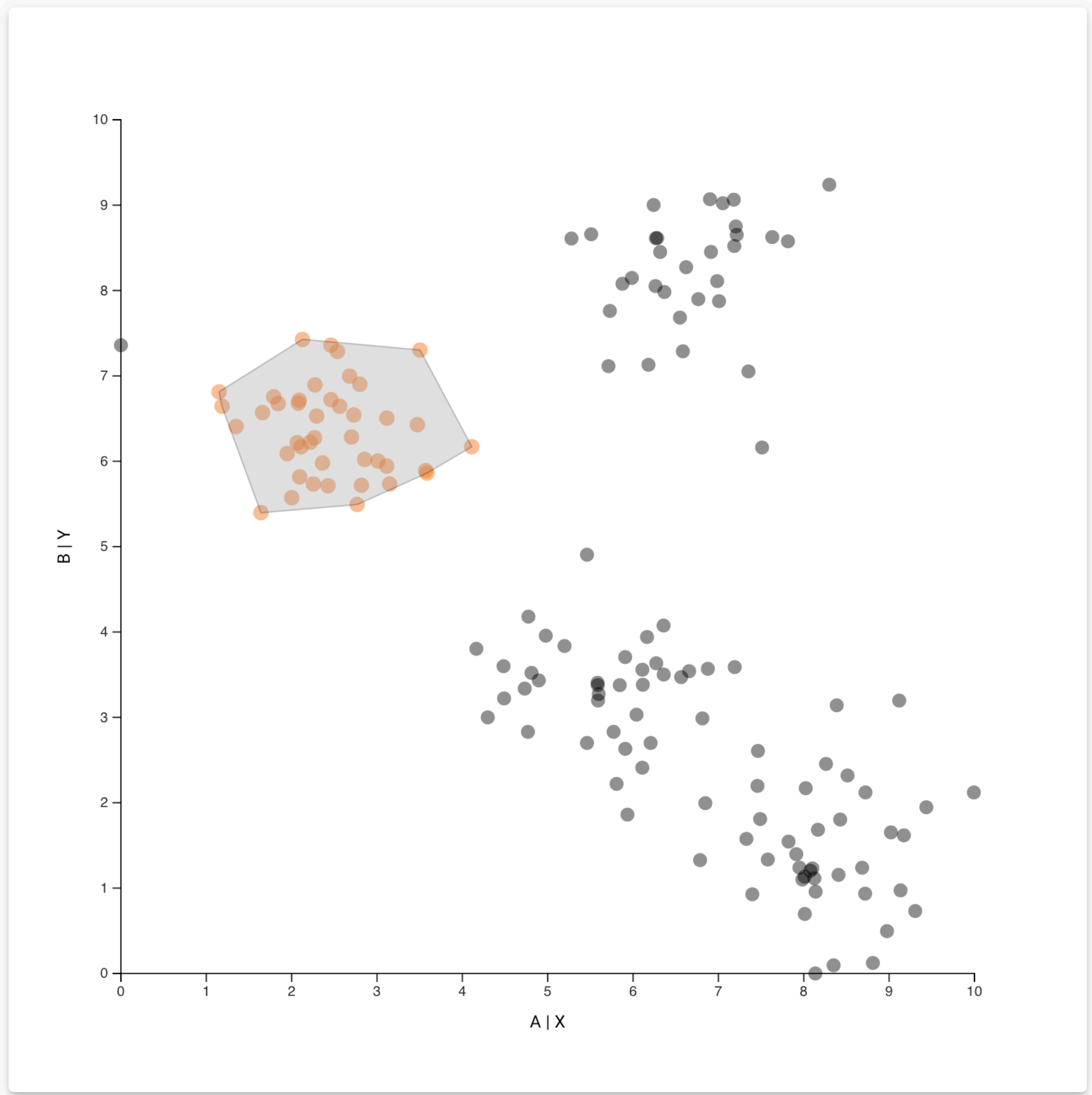
Changed Dataset

cluster_simple_v1 cluster_simple_v2 cluster_simple_v3 cluster_simple_v6 cluster_simple_v4 cluster_simple_v5 ☒ Compare Datasets



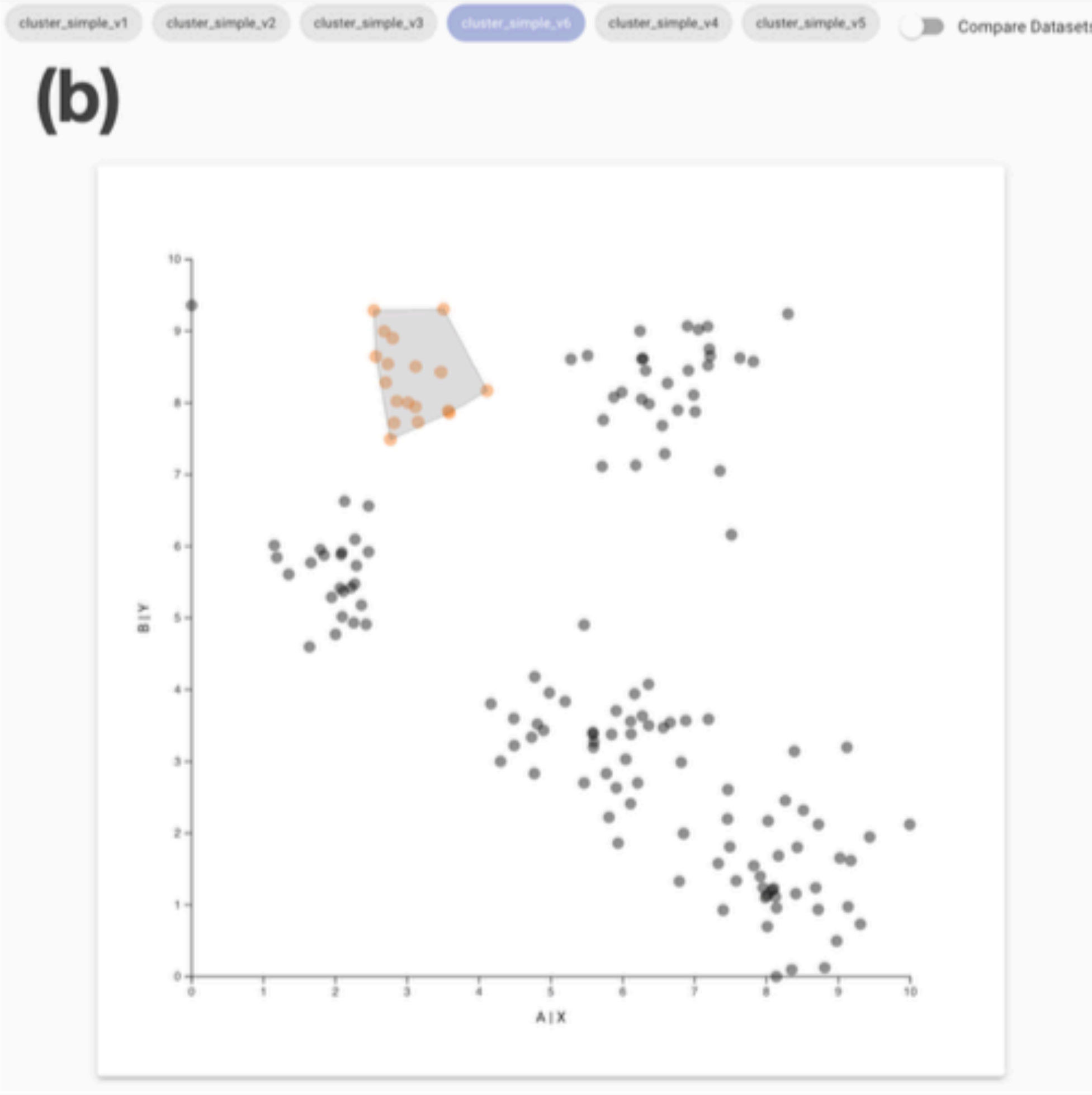
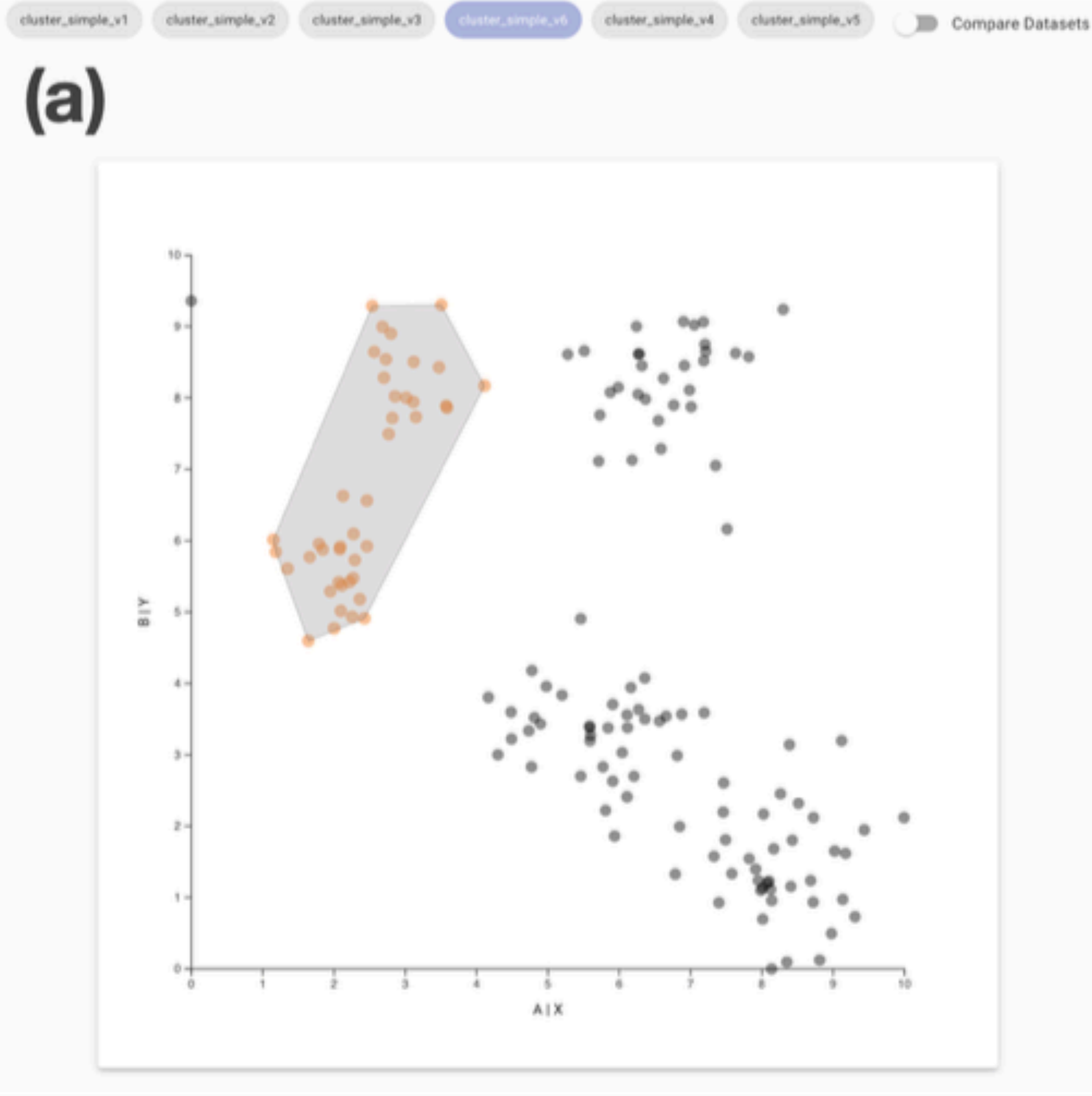
Tracking A Selected Cluster

cluster_simple_v1 cluster_simple_v2 cluster_simple_v3 cluster_simple_v6 cluster_simple_v4 cluster_simple_v5 ☐ Compare Datasets

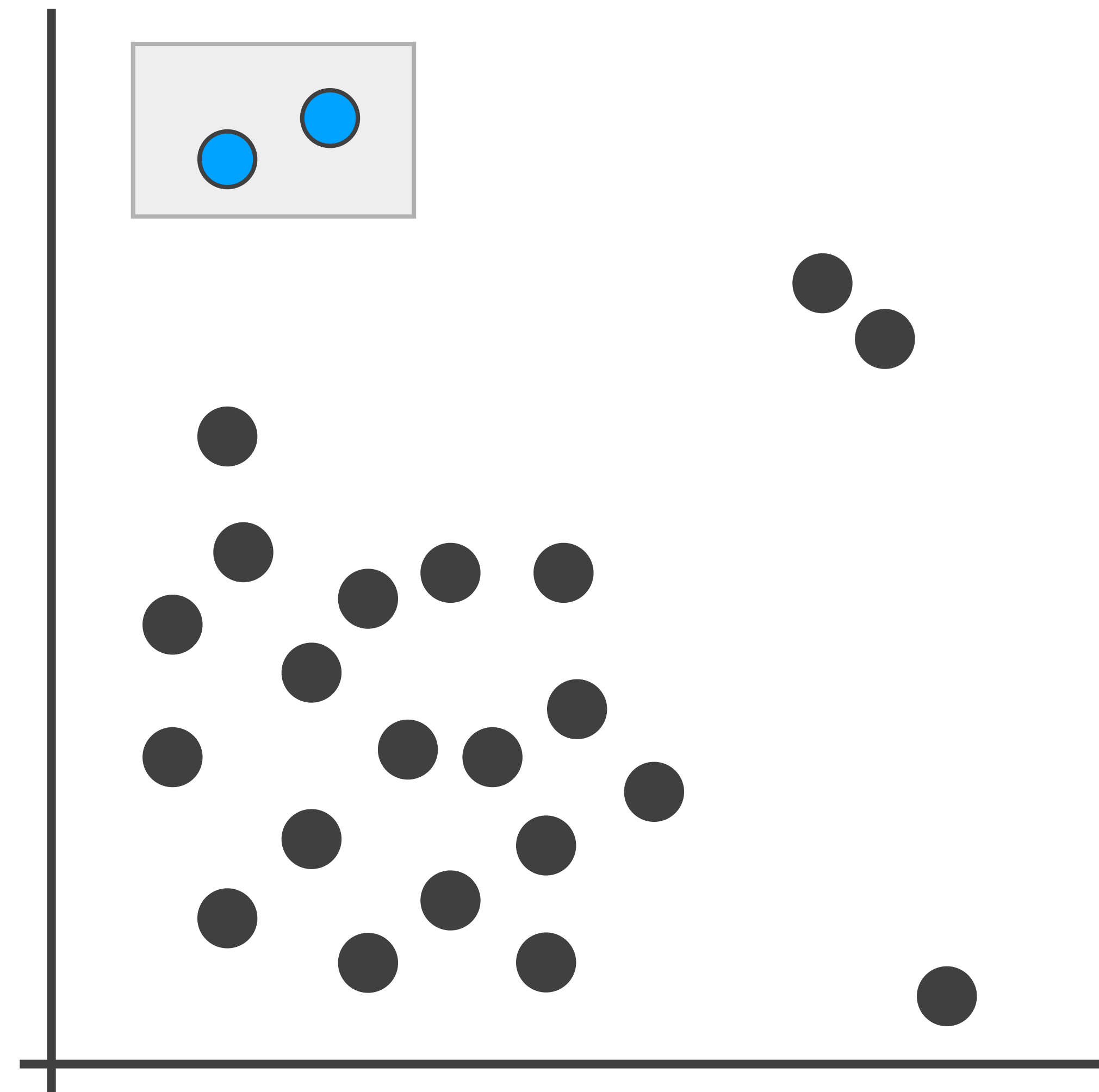


Selected Cluster on Changed Dataset

HUMAN REVIEWS



Interactive Visual Analysis



Workflow
Database

Computational Analysis



USING WORKFLOW IN A COMPUTATIONAL NOTEBOOK

```
# Installing the reapply-workflows adds a module called backend
# This module exposes the Reapply class which initializes the library

from backend import Reapply
```

```
# Here we load the reapply_workflows library.
r = Reapply()

# We add a workflow from our workflow database.
workflow = r.load_workflow("workflow1617808681620")
```

```
# Print the workflow name
print("Workflow: ", workflow.name, "\n")

# Description of the workflow and the operations in it
workflow.describe
```

Workflow: Deleting Cluster

```
| Root
+--| Add Plot
    +--| Added brush to: X-Y
        +--| Cluster Selection
            +--| Filter: Out
```

```
# Prints the reapply results for all interactions, along with review status.
```

```
# Apply the workflow to target dataset.
# apply function requires the target dataset
# and the label column as arguments.
res = workflow.apply(target, "Label")

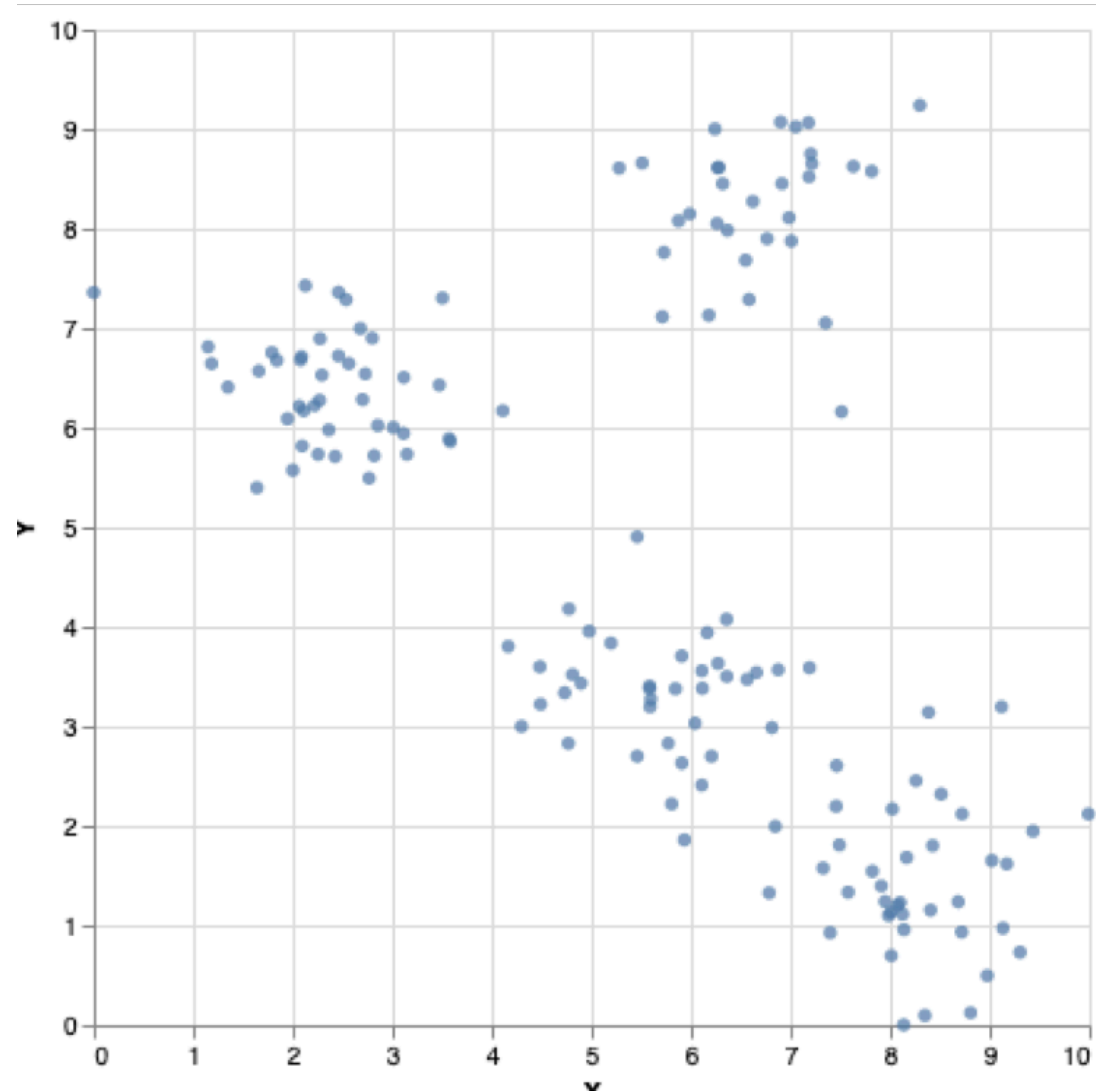
# Results is an array of datasets for each interaction
# we grab the final one.
result_dataset = res.results[-1]['data']
result_dataset
```

This workflow has not been reviewed for all interactions.
Please go to following url: <https://reapply-workflows.git/>

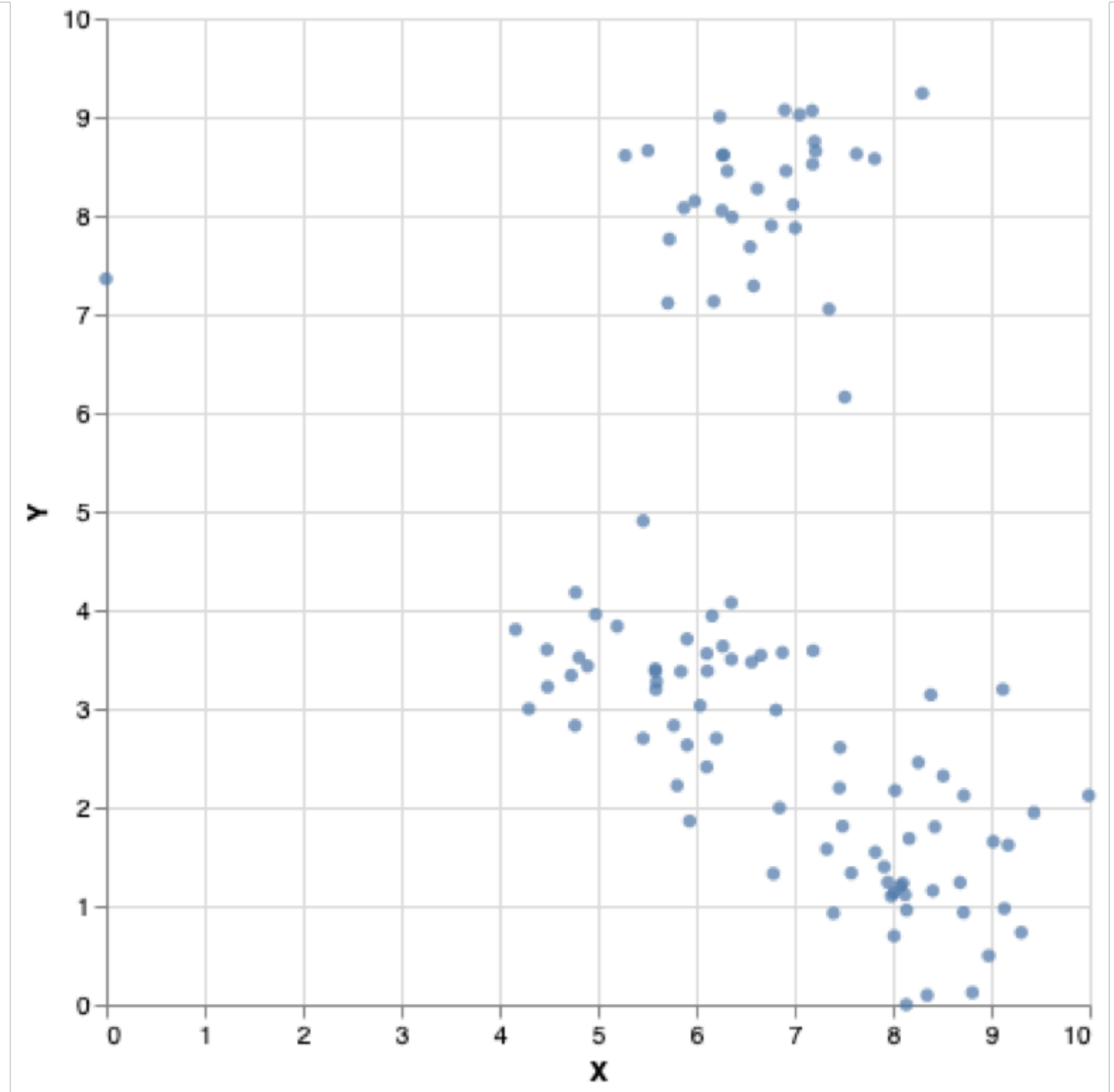
	Label	X	Y
3	P52	6.58351	7.28796
5	P171	4.77421	4.17980
8	P199	8.34966	0.09550
9	P183	8.42670	1.80299
10	P61	4.29760	2.99981
...
141	P138	7.35179	7.05215
142	P46	6.62171	8.27311

BEFORE AND AFTER

150 rows × 3 columns



108 rows × 3 columns



Brush Type



Transforms

FILTER

LABEL

AGGREGATE

Categories

Enable category encoding

Showing categories

Case Category

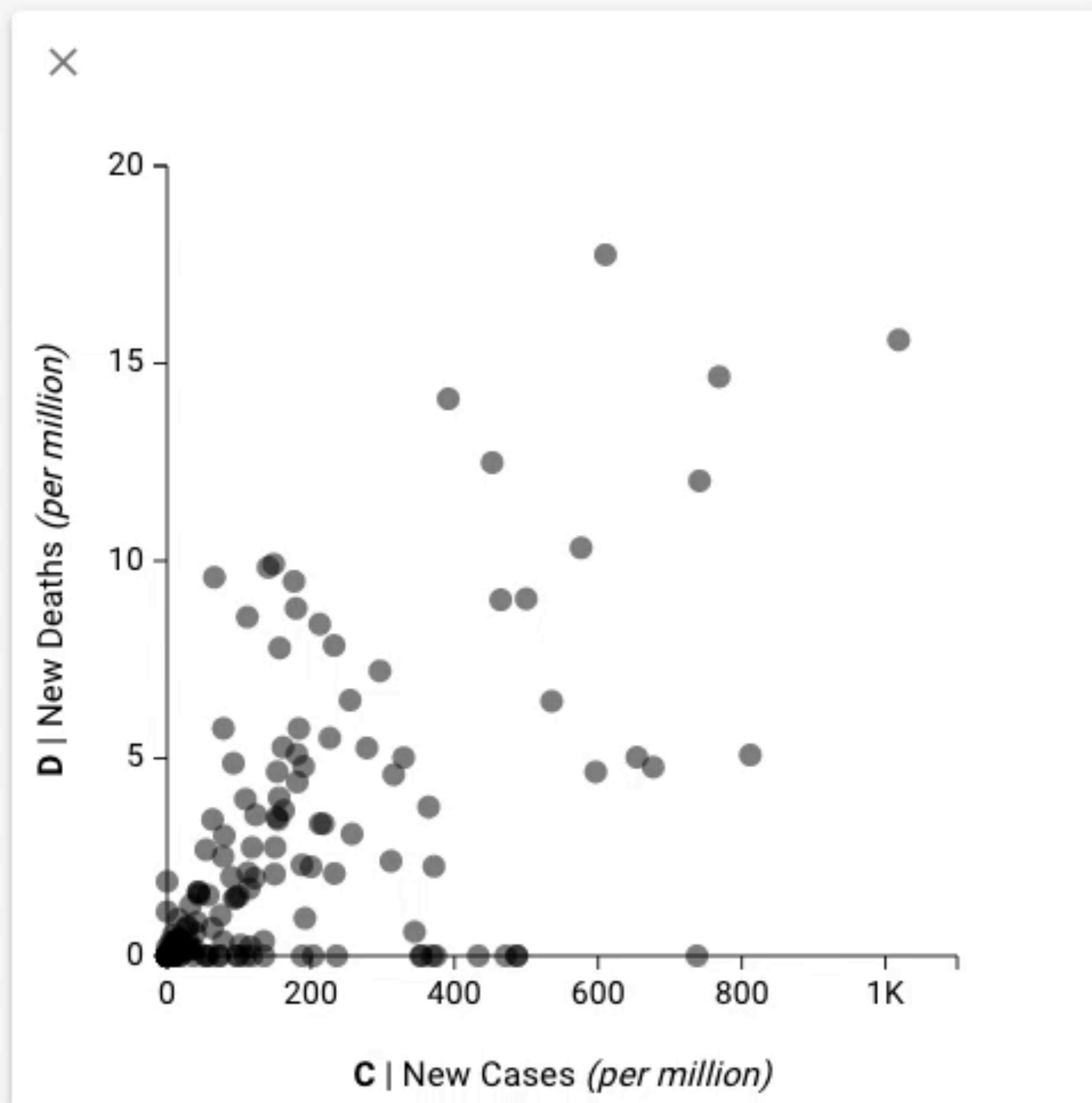
- UNASSIGNED
- + HIGH DEATHS - HIGH CASES
- ◆ LOW DEATH - HIGH CASES
- HIGH DEATHS - LOW CASES

jan_2021

june_2021

dec_2020

Compare



NO VALID SELECTIONS

Graph

Bookmarks/Annotations

Undo

Redo

Root

- jan_2021 Adding scatterplot f..
- jan_2021 Adding scatterplot f..
- jan_2021 Remove Scatterplot

COVID Case Study

+

<>

TOWARDS LITERATE & REUSABLE VISUALIZATION

Semantic selections & annotations allow us to have **higher-level information** about events in an analysis process.

Analysts have the **means to justify** their choices.

Makes it possible to **reuse interactive analysis processes** on updated datasets.

Alexander Lex

@alexander_lex

<http://alexander-lex.net>



Thanks to: **Kiran Gadhave, Zach Cutler**, Marc Streit, Jochen Görtler, Oliver Deussen, Miriah Meyer, Jeff Phillips, Samuel Gratzl, Holger Stitz,, Nils Gehlenborg, Hendrik Strobelt, Romain Vuillemot, Hanspeter Pfister, and many others!



visualization
design lab

