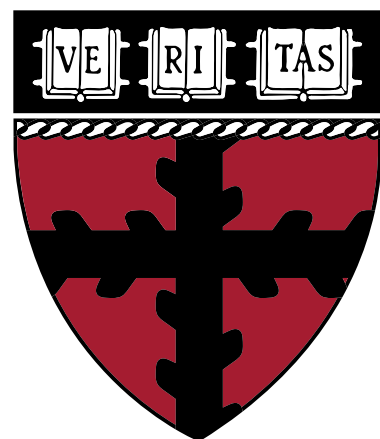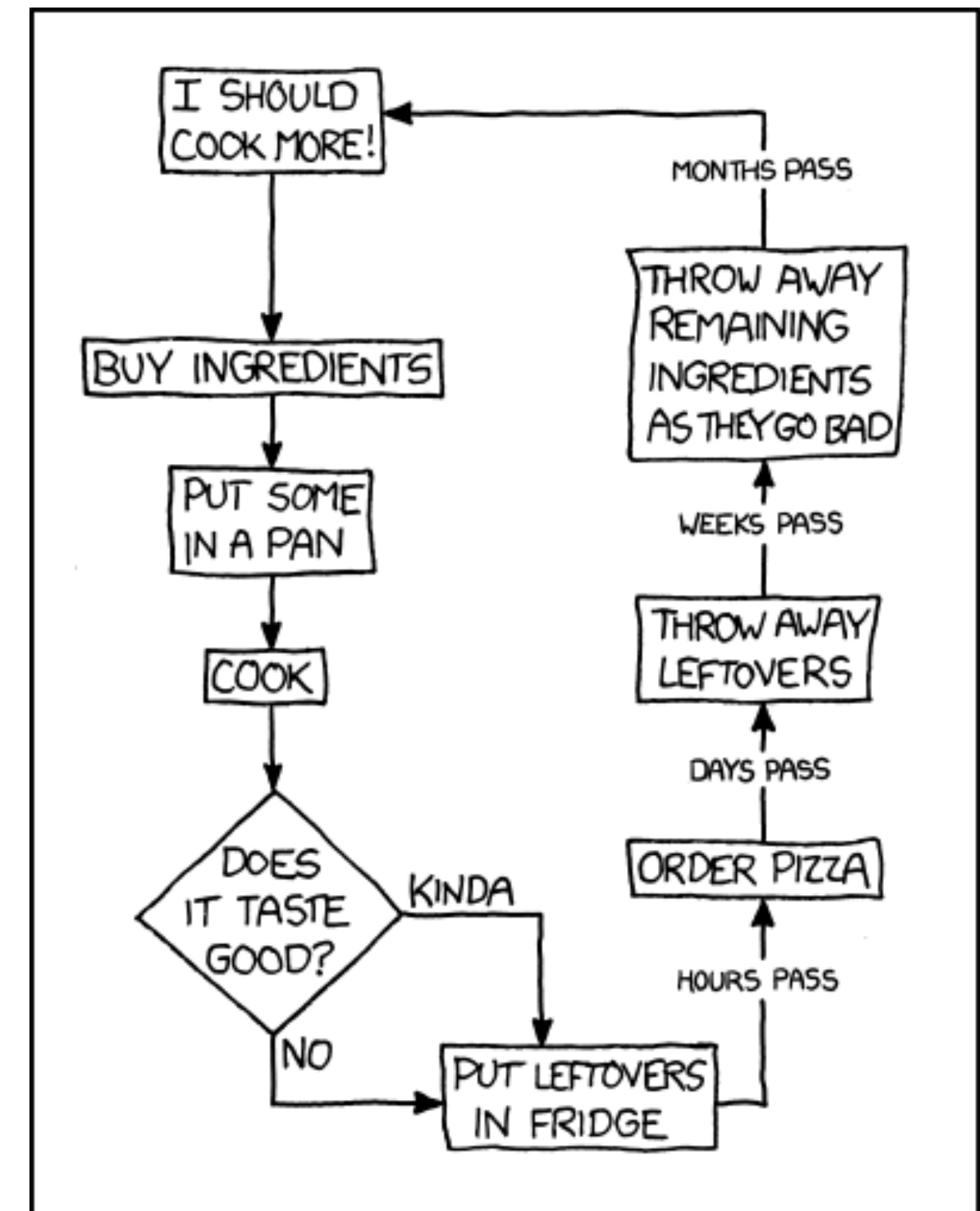# CS171 Visualization

Alexander Lex
alex@seas.harvard.edu

Tables Part II

HARVARD
School of Engineering
and Applied Sciences



[xkcd]

# Next Week

Reading: VAD, Chapters 9

Lecture 11: Text & Documents

Lecture 12: Homework 3 Design Studio

Sections: view coordination, linking & brushing

Updates

Design Studio moved to Thursday
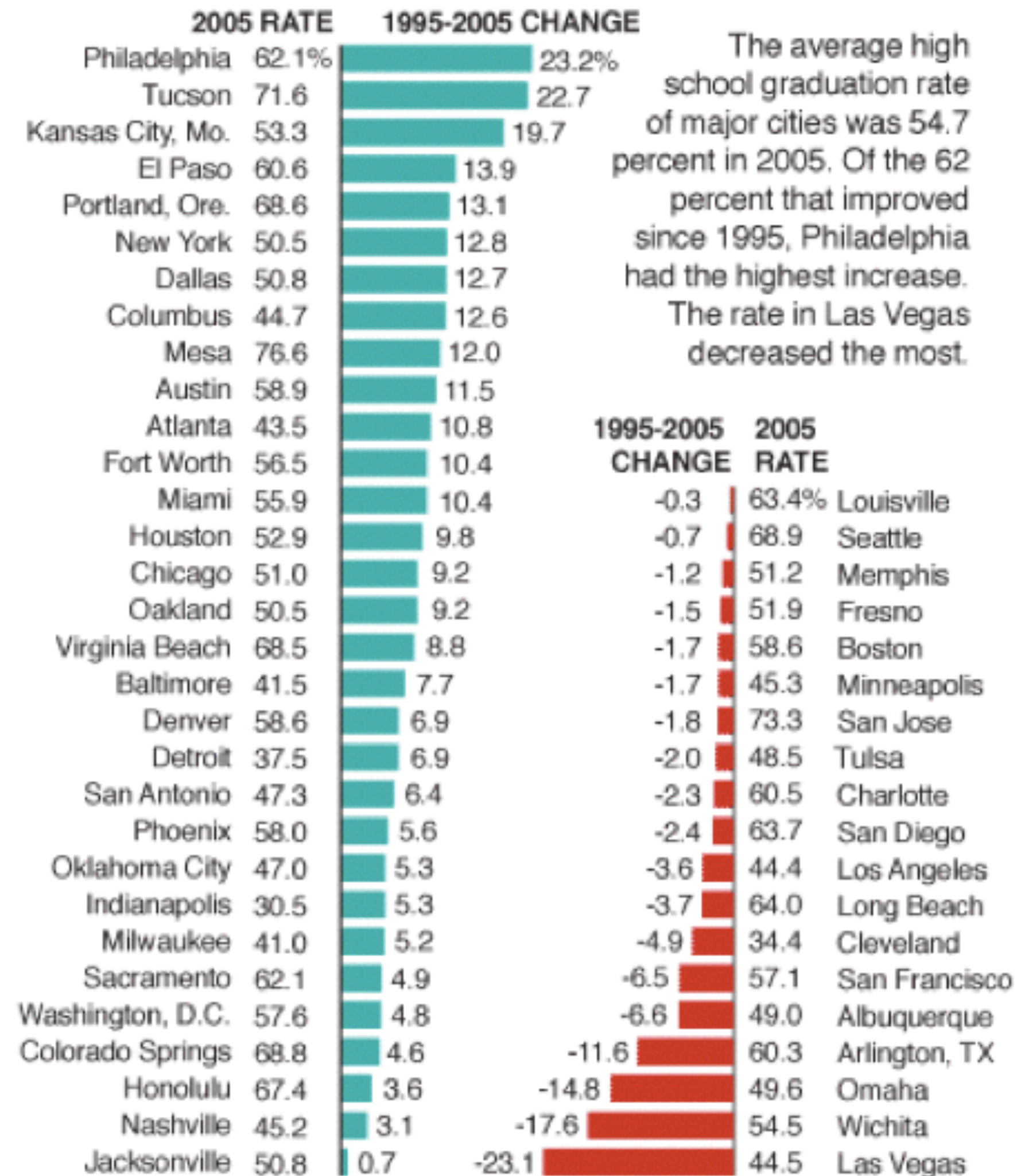
Project Proposal moved to HW 4

# Tables & Multi-Dimensional Data

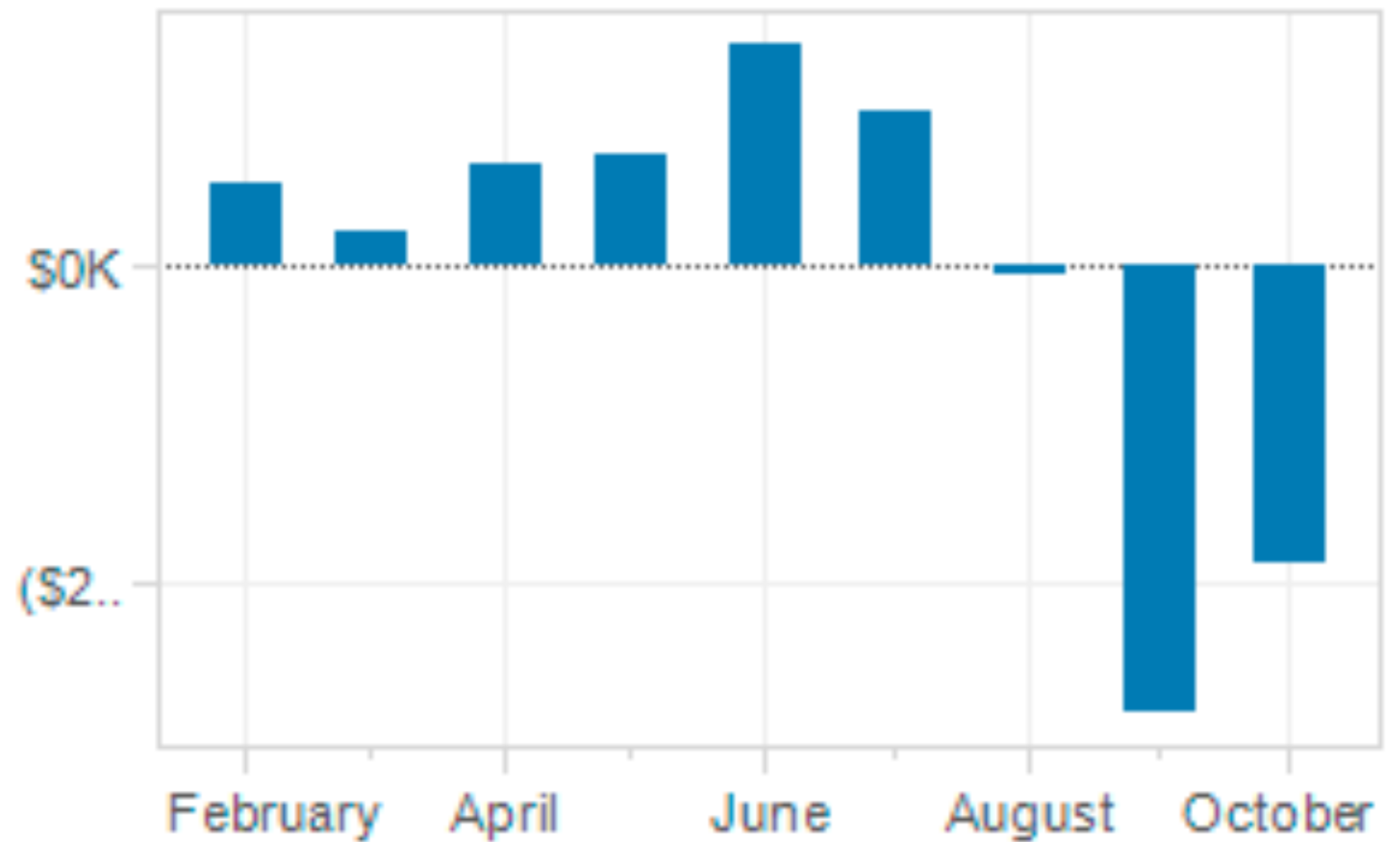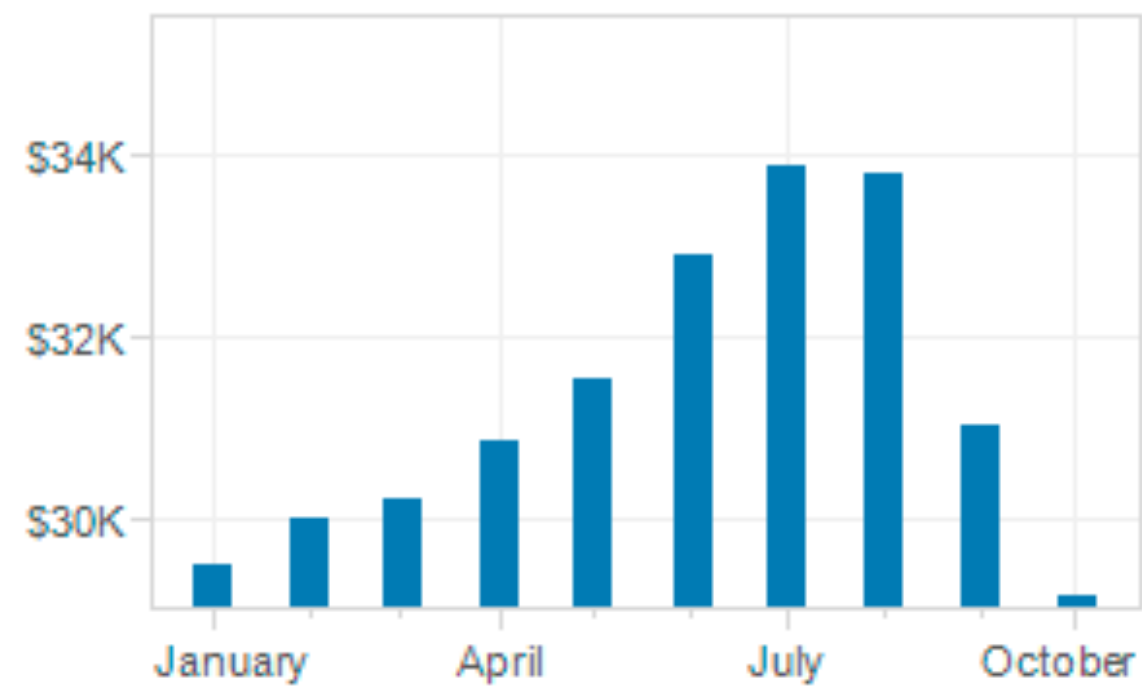# Comparisons

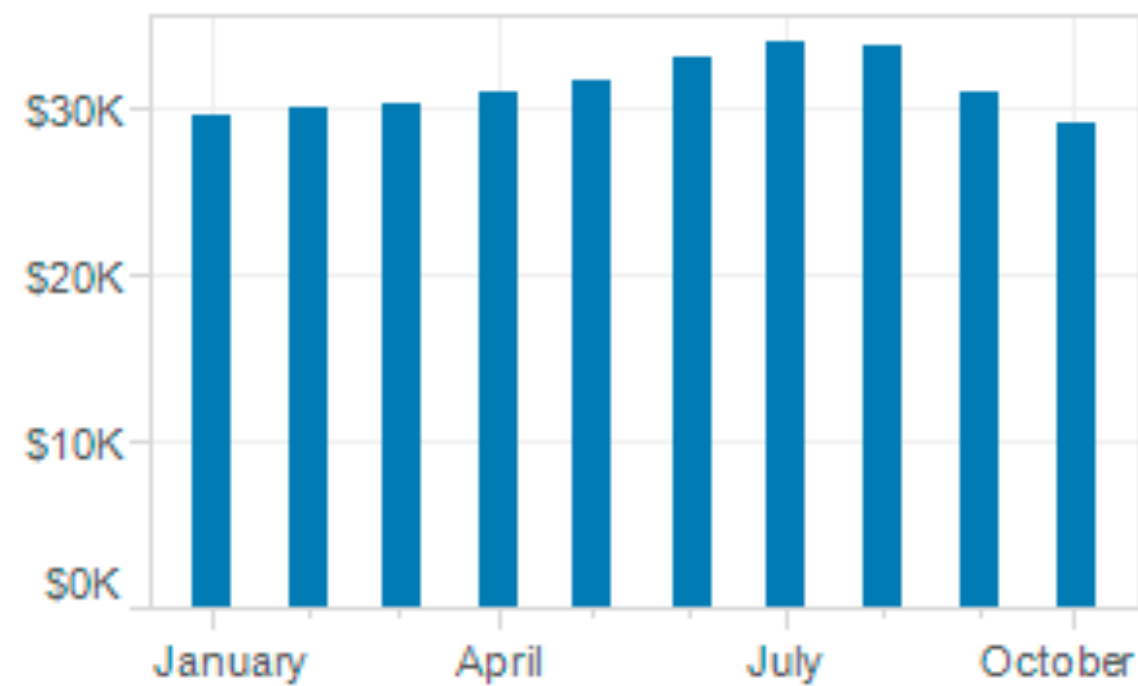# Direction



Graduation rates up in most cities

Graduation rate for principal school district of the largest cities

The average high school graduation rate of major cities was 54.7 percent in 2005. Of the 62 percent that improved since 1995, Philadelphia had the highest increase. The rate in Las Vegas decreased the most.

| | 2005 RATE | 1995-2005 CHANGE |
|---|---|---|
| Philadelphia | 62.1% | 23.2% |
| Tucson | 71.6 | 22.7 |
| Kansas City, Mo. | 53.3 | 19.7 |
| El Paso | 60.6 | 13.9 |
| Portland, Ore. | 68.6 | 13.1 |
| New York | 50.5 | 12.8 |
| Dallas | 50.8 | 12.7 |
| Columbus | 44.7 | 12.6 |
| Mesa | 76.6 | 12.0 |
| Austin | 58.9 | 11.5 |
| Atlanta | 43.5 | 10.8 |
| Fort Worth | 56.5 | 10.4 |
| Miami | 55.9 | 10.4 |
| Houston | 52.9 | 9.8 |
| Chicago | 51.0 | 9.2 |
| Oakland | 50.5 | 9.2 |
| Virginia Beach | 68.5 | 8.8 |
| Baltimore | 41.5 | 7.7 |
| Denver | 58.6 | 6.9 |
| Detroit | 37.5 | 6.9 |
| San Antonio | 47.3 | 6.4 |
| Phoenix | 58.0 | 5.6 |
| Oklahoma City | 47.0 | 5.3 |
| Indianapolis | 30.5 | 5.3 |
| Milwaukee | 41.0 | 5.2 |
| Sacramento | 62.1 | 4.9 |
| Washington, D.C. | 57.6 | 4.8 |
| Colorado Springs | 68.8 | 4.6 |
| Honolulu | 67.4 | 3.6 |
| Nashville | 45.2 | 3.1 |
| Jacksonville | 50.8 | 0.7 |

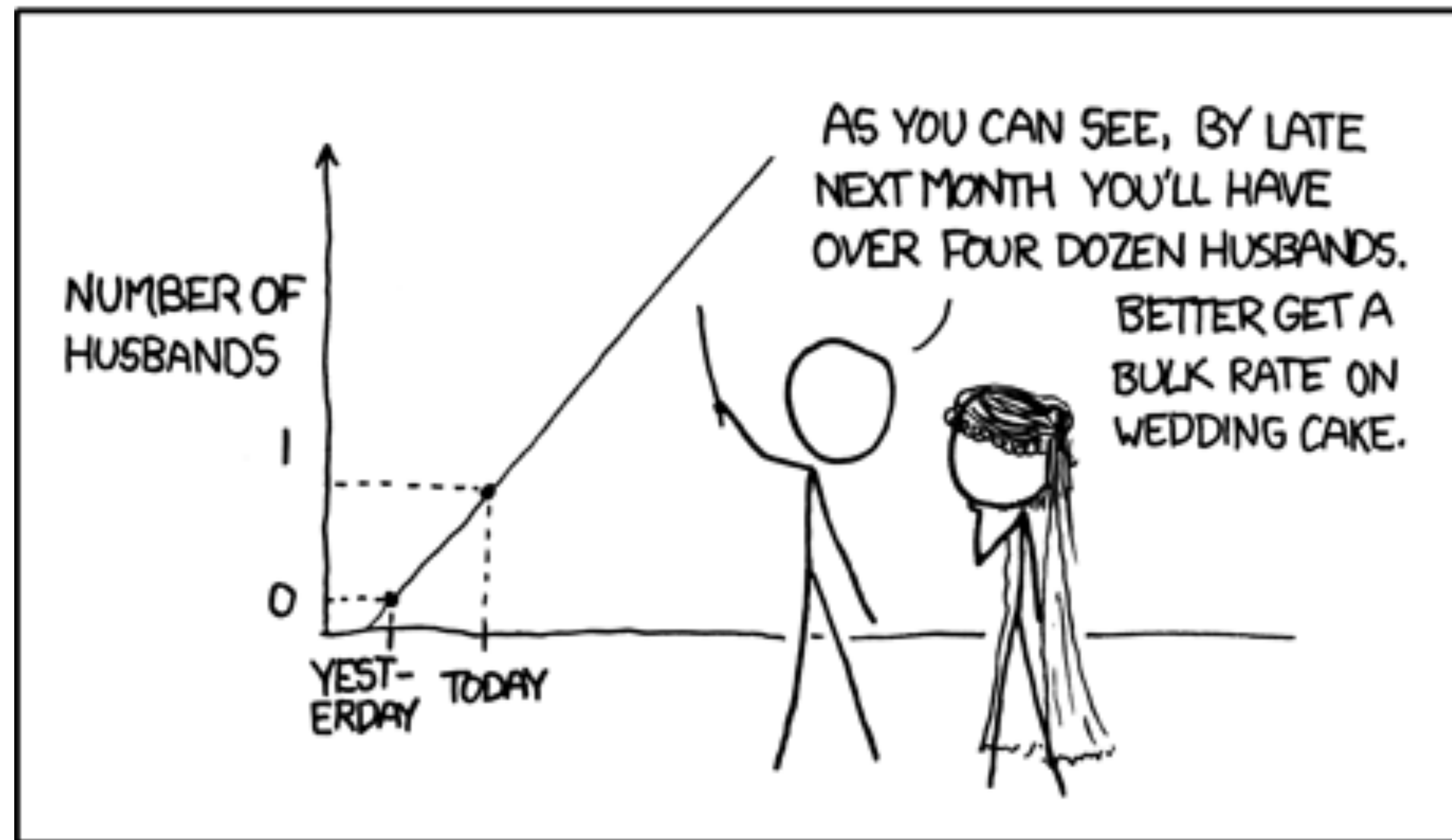| 1995-2005 CHANGE | 2005 RATE | |
|---|---|---|
| -0.3 | 63.4% | Louisville |
| -0.7 | 68.9 | Seattle |
| -1.2 | 51.2 | Memphis |
| -1.5 | 51.9 | Fresno |
| -1.7 | 58.6 | Boston |
| -1.7 | 45.3 | Minneapolis |
| -1.8 | 73.3 | San Jose |
| -2.0 | 48.5 | Tulsa |
| -2.3 | 60.5 | Charlotte |
| -2.4 | 63.7 | San Diego |
| -3.6 | 44.4 | Los Angeles |
| -3.7 | 64.0 | Long Beach |
| -4.9 | 34.4 | Cleveland |
| -6.5 | 57.1 | San Francisco |
| -6.6 | 49.0 | Albuquerque |
| -11.6 | 60.3 | Arlington, TX |
| -14.8 | 49.6 | Omaha |
| -17.6 | 54.5 | Wichita |
| -23.1 | 44.5 | Las Vegas |

SOURCE: EPE Research Center

AP

Nicolas Rapp

# Plot Change Instead

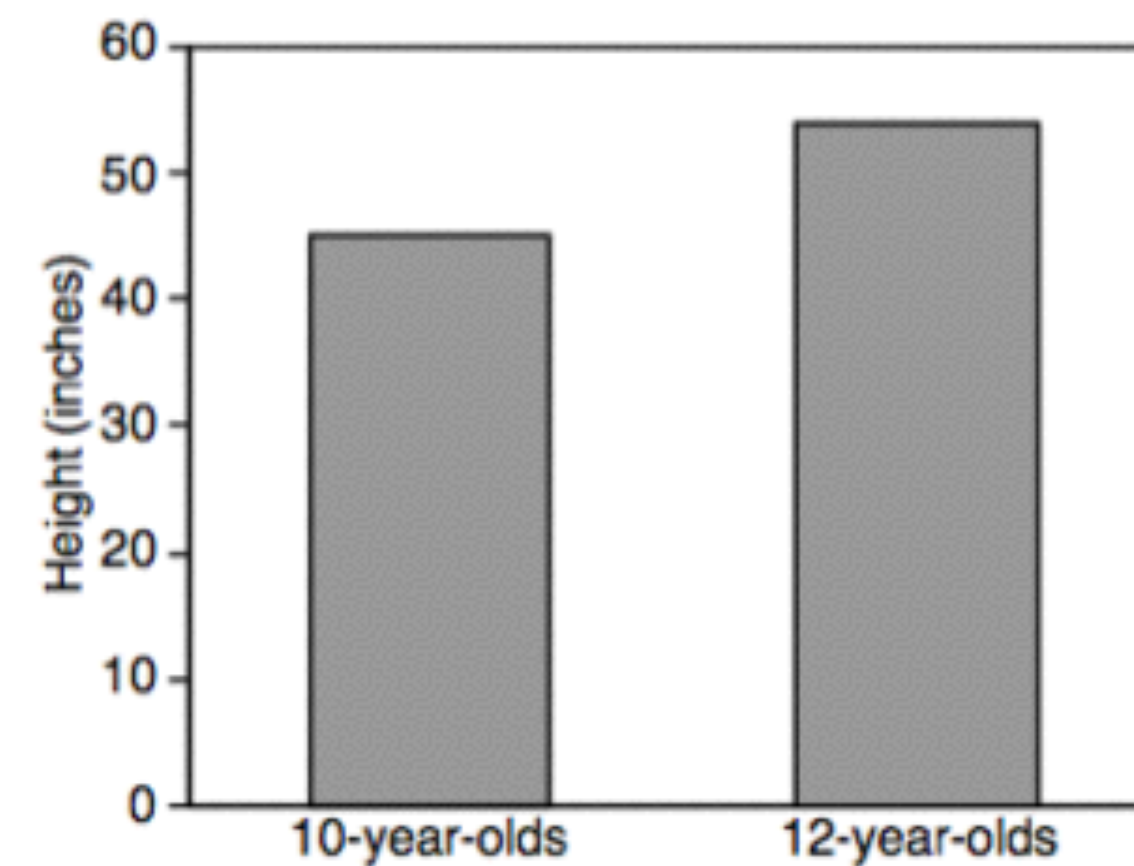# Trends Over Time
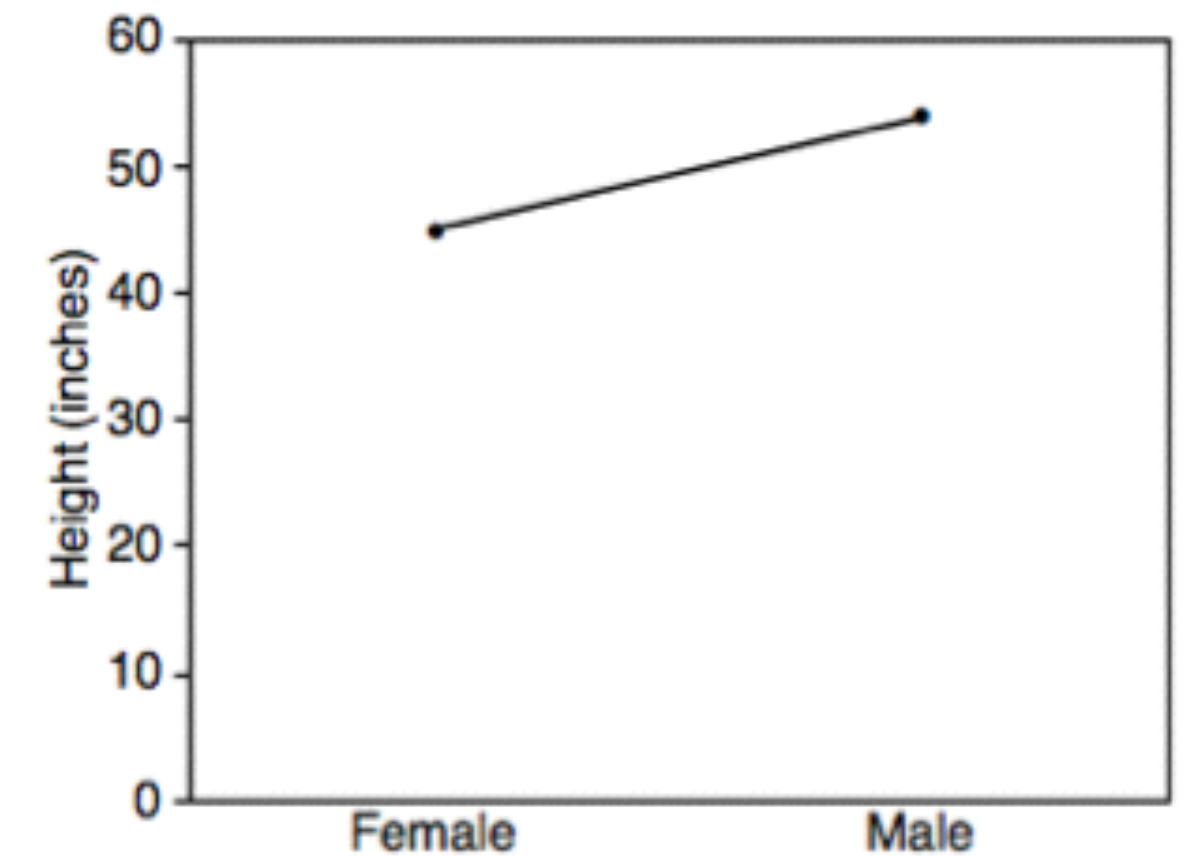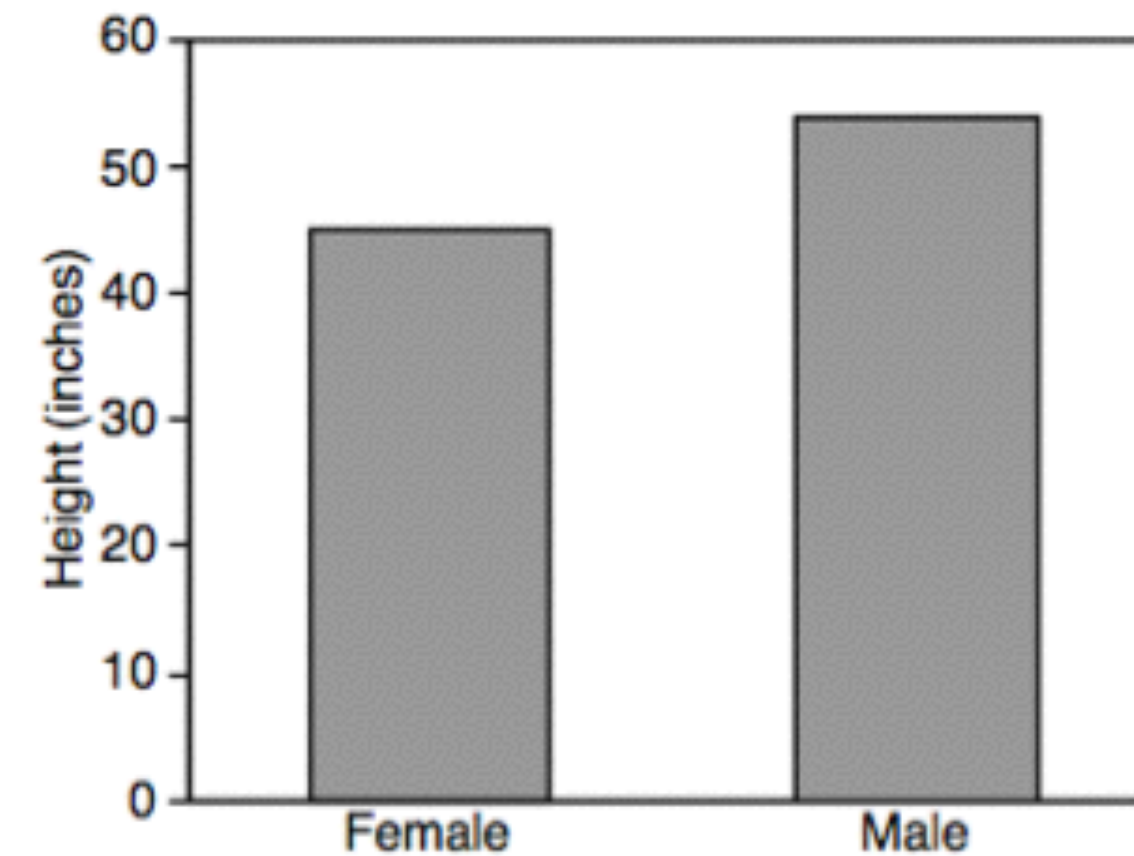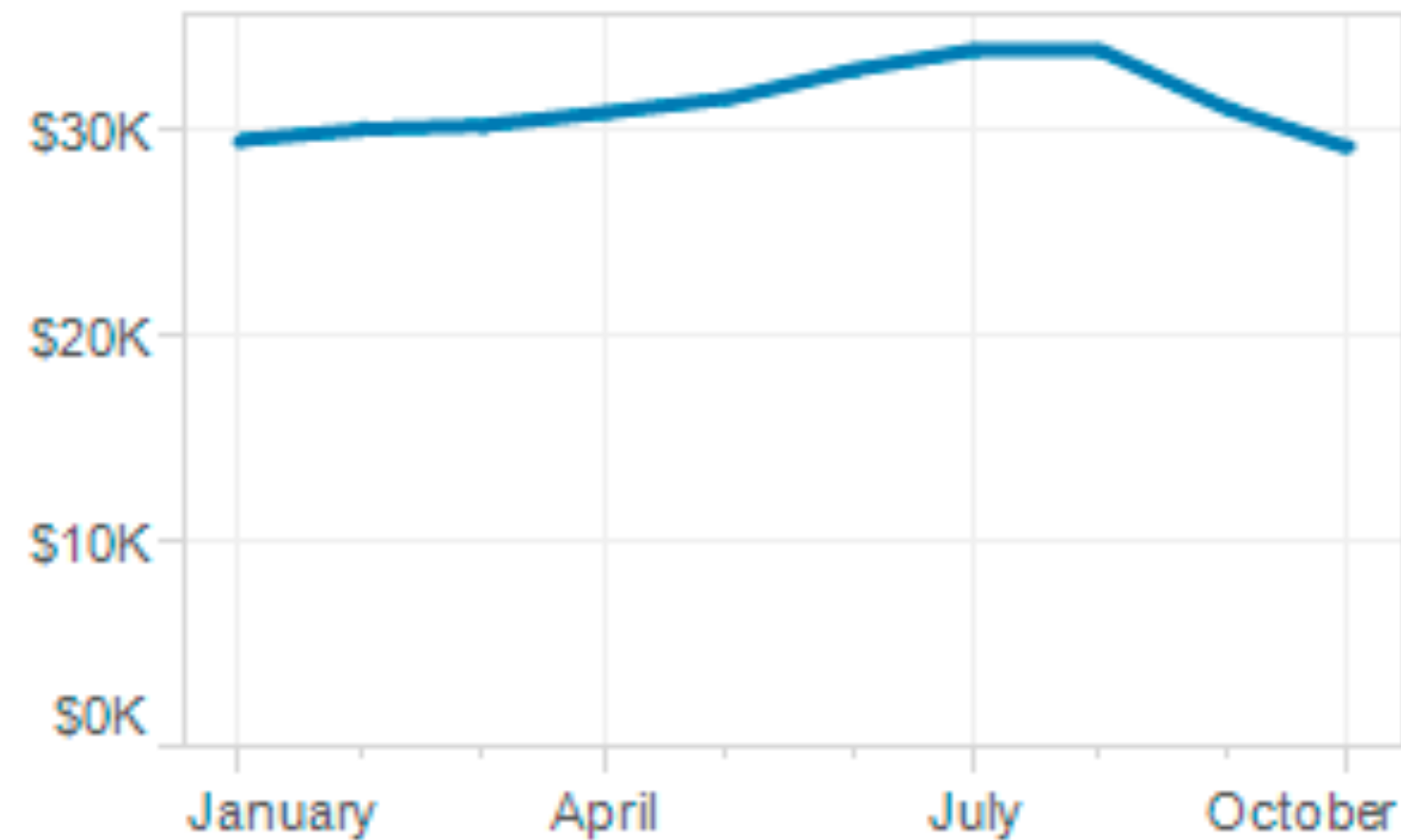
# Bars vs. Lines

Lines imply connections &
sampling from continuous
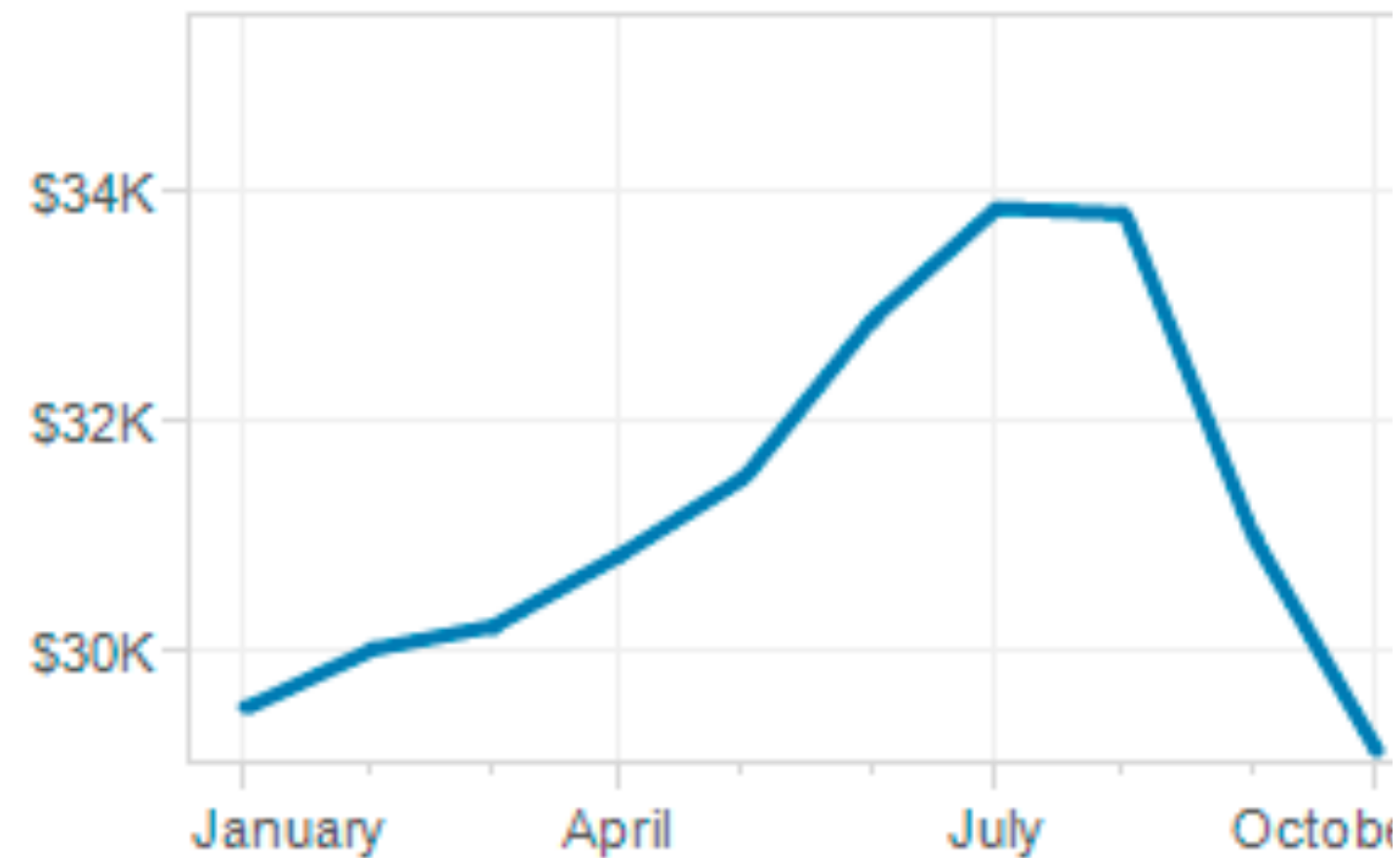data.
Do not use for categorical
data.

# Baseline Problem (again)

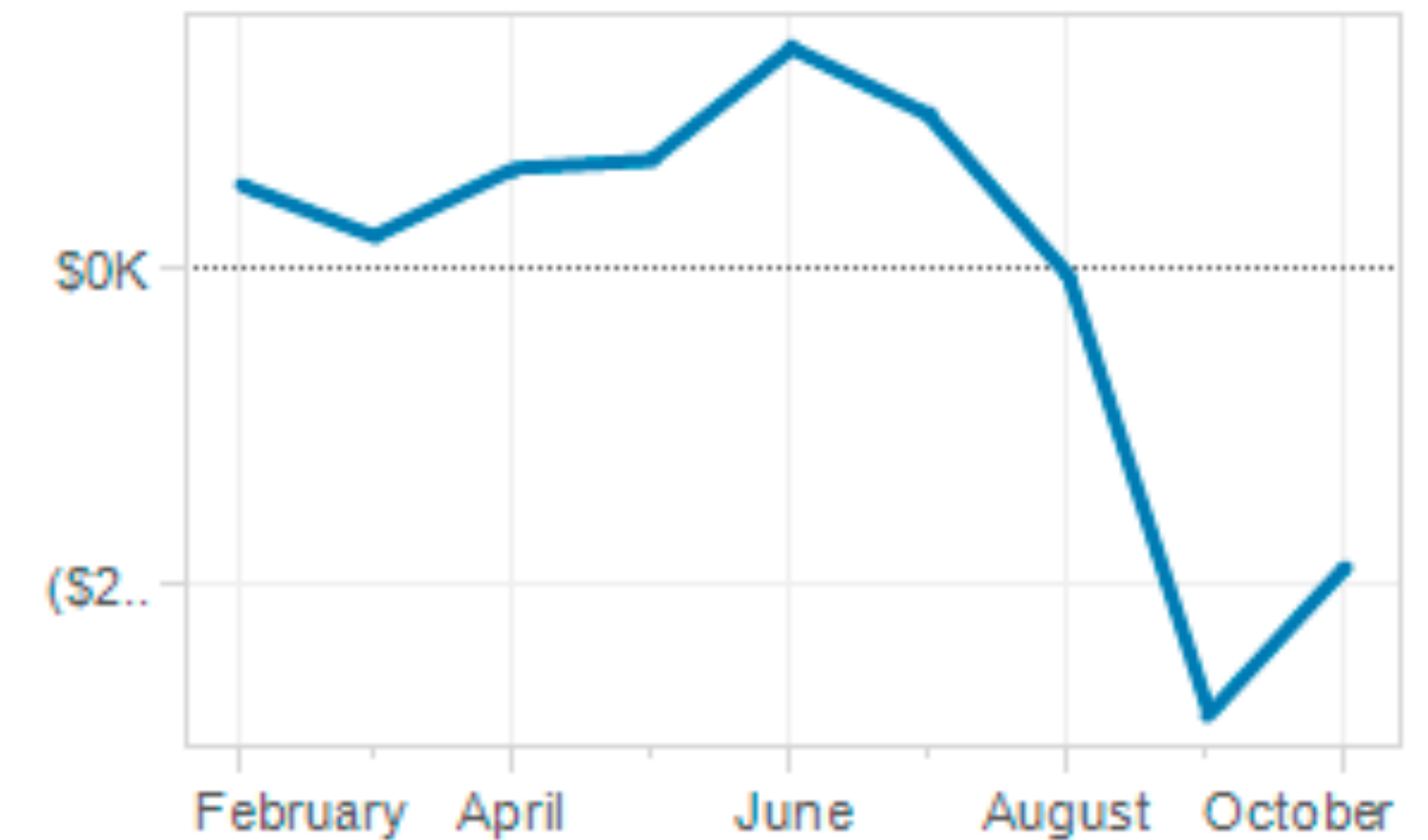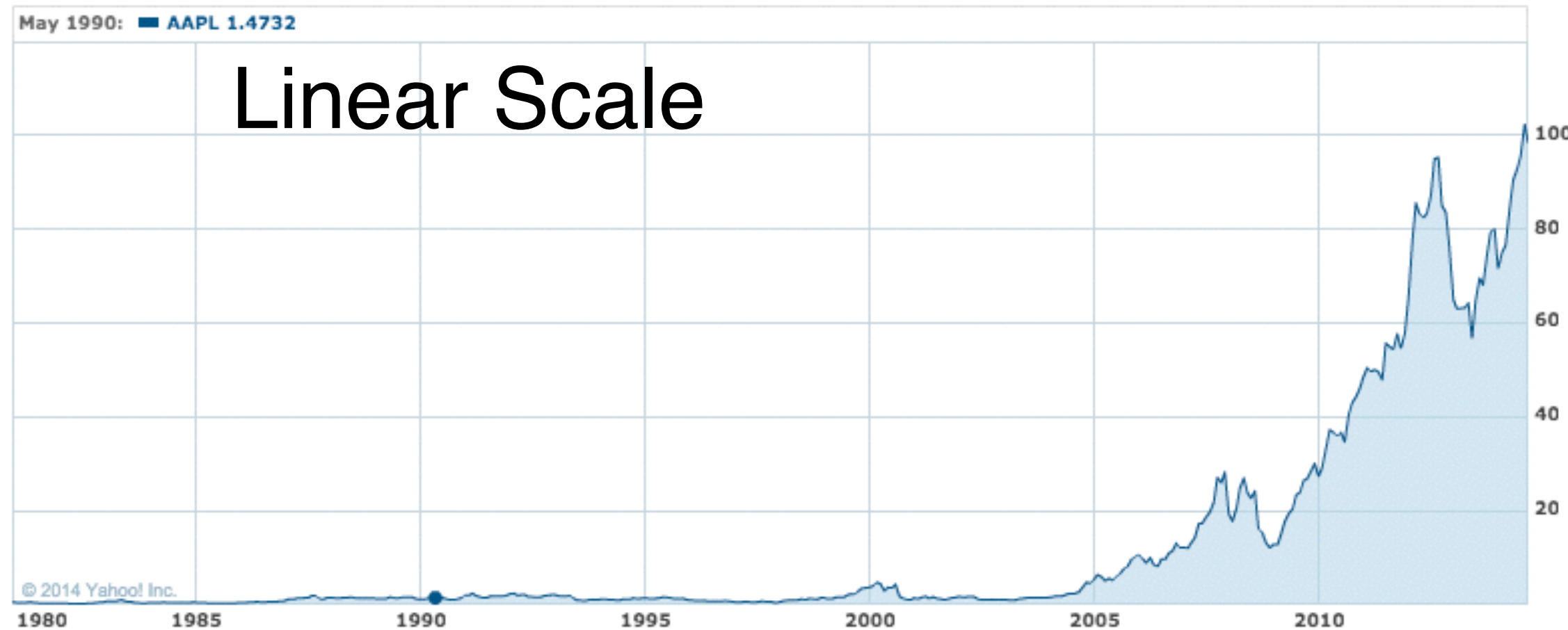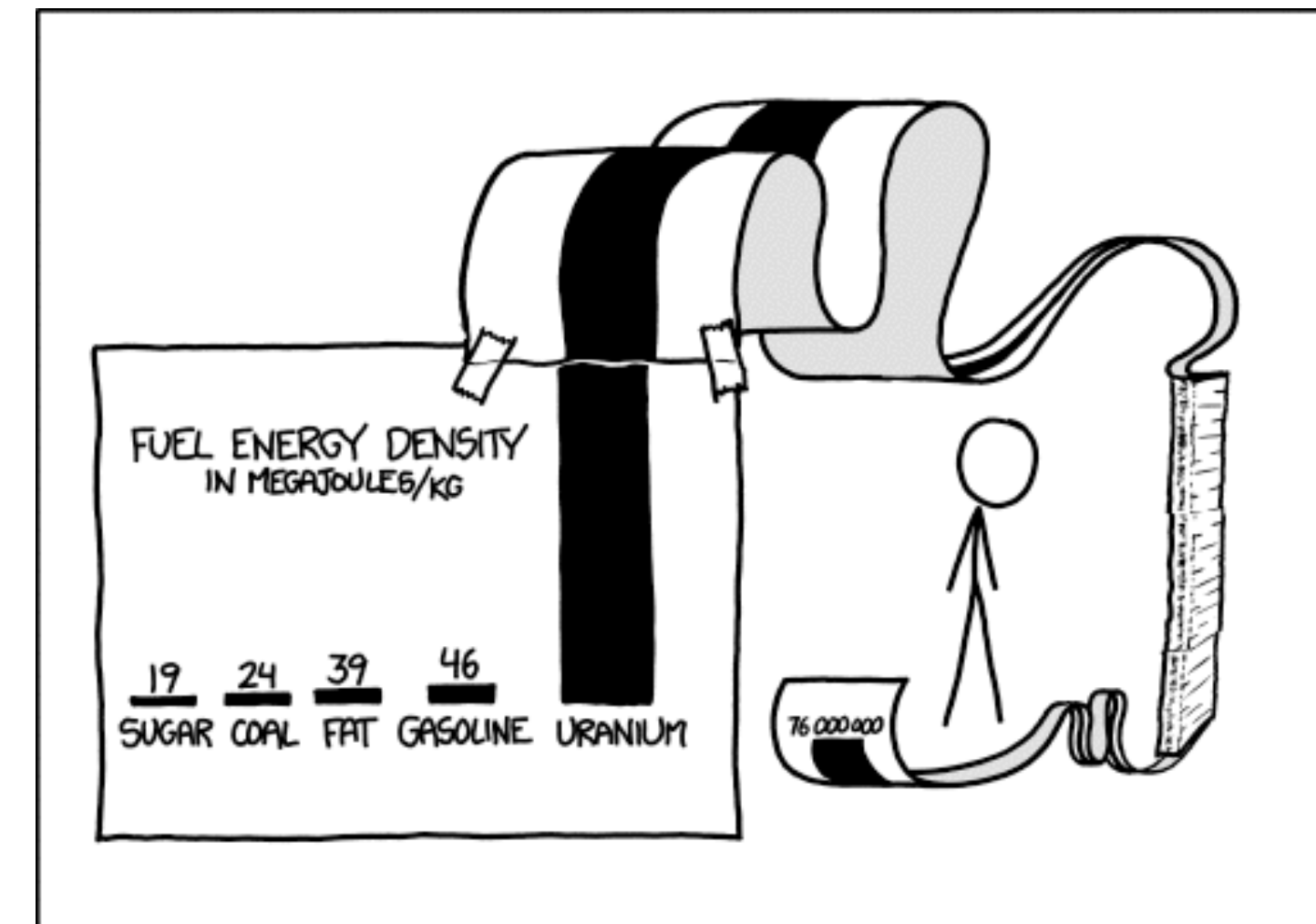True Baseline

Clipped Baseline
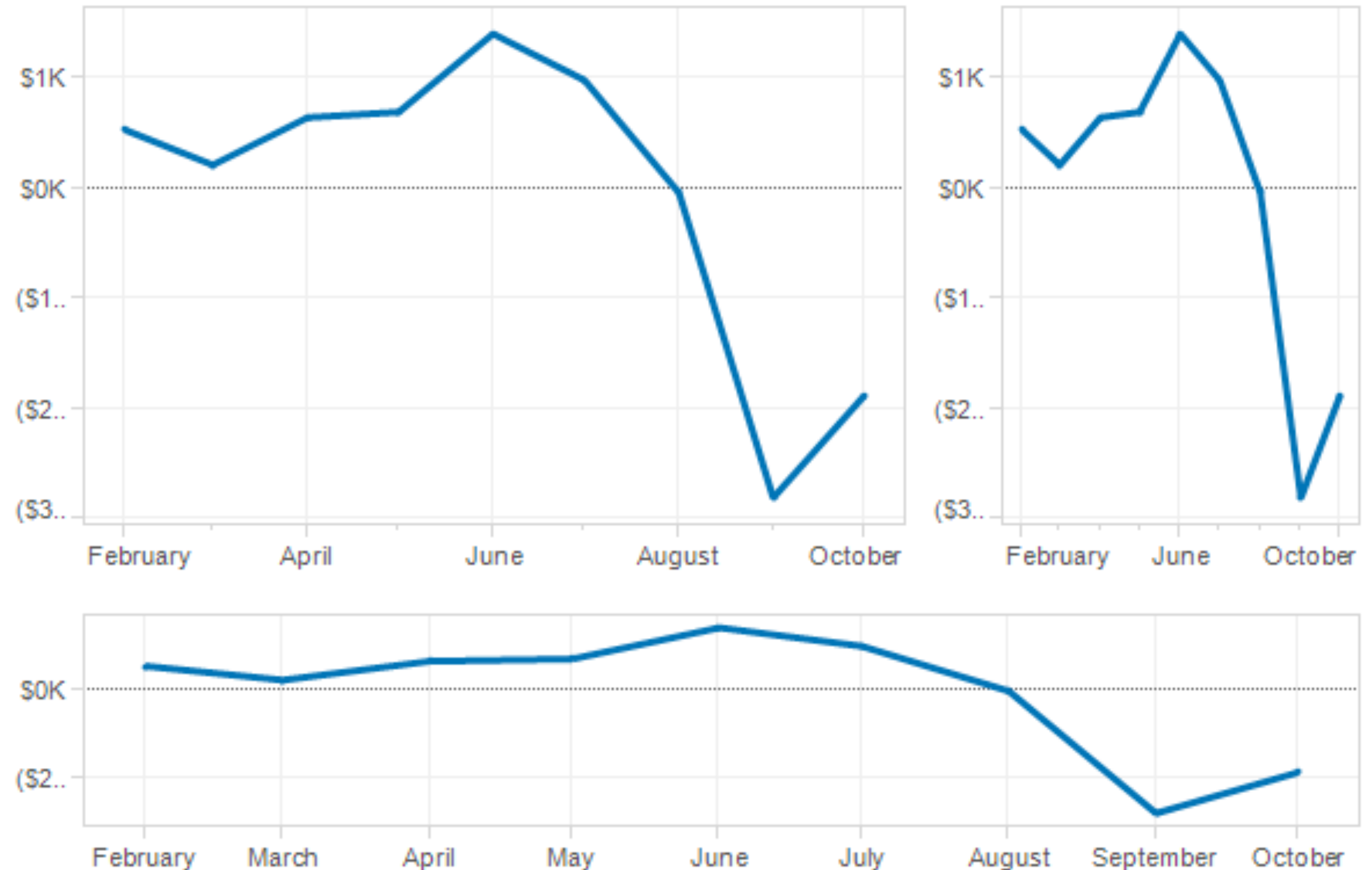
Plotting Change

# Linear vs. Logarithmic Scale



Linear Scale

Log Scale

Apple Stock Price

FUEL ENERGY DENSITY
IN MEGAJOULES/KG

19 SUGAR   24 COAL   39 FAT   46 GASOLINE   URANIUM   76,000,000

SCIENCE TIP: LOG SCALES ARE FOR QUITTERS WHO CAN'T FIND ENOUGH PAPER TO MAKE THEIR POINT *PROPERLY*.

http://xkcd.com/1162/

http://finance.yahoo.com/echarts?s=AAPL

# Aspect Ratios

Rule of Thumb:

Banking to 45°

(average line slope: 45°)

# Correlations

# Scatterplots

# Overplotting



alpha = 1/100

# Compositions

# Stacked Bar Chart

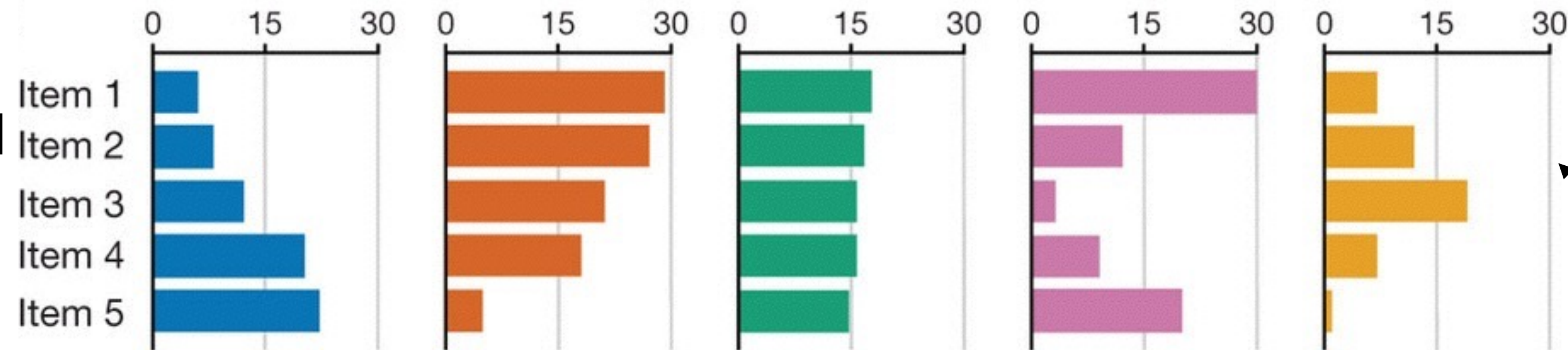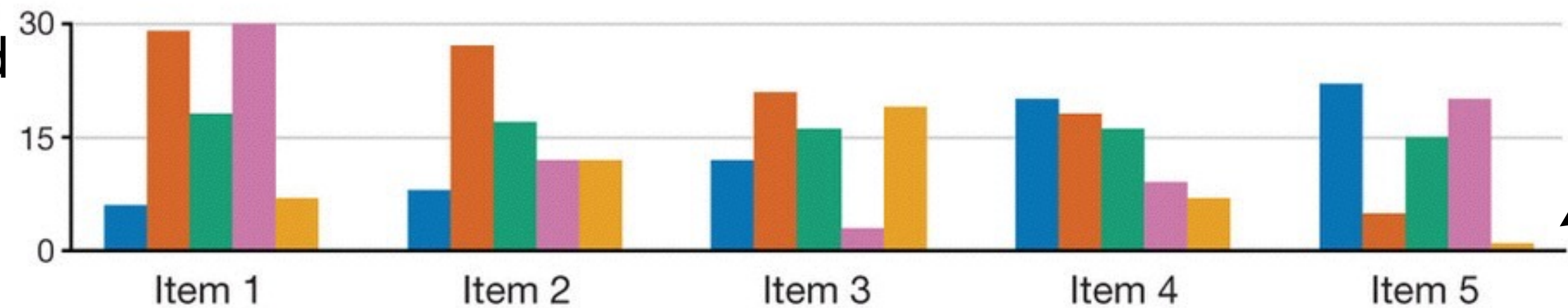# Comparison of bar chart types
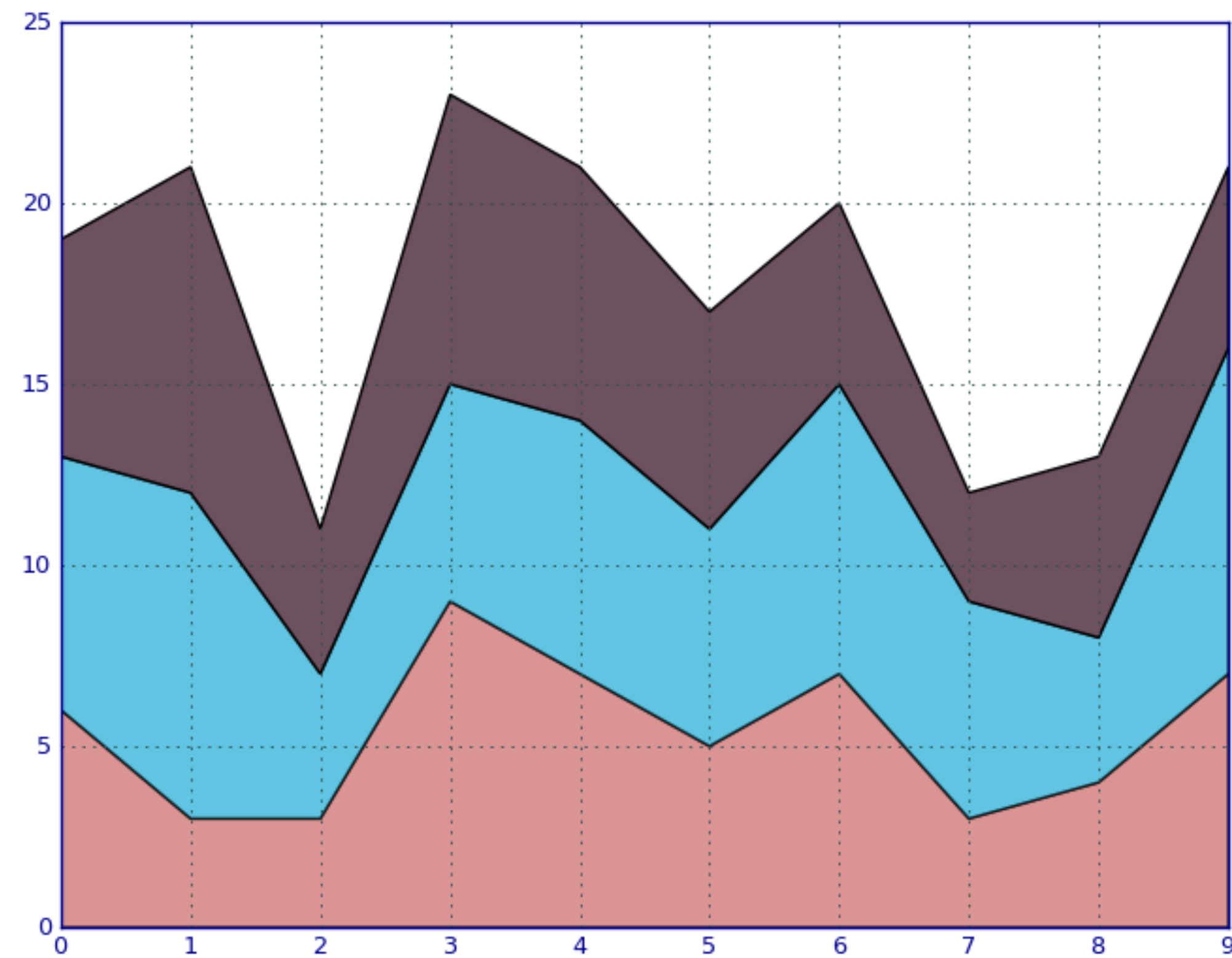


Pie Chart

Stacked bar chart

Layered Bar Chart

Grouped Bar Chart

Small Multiples

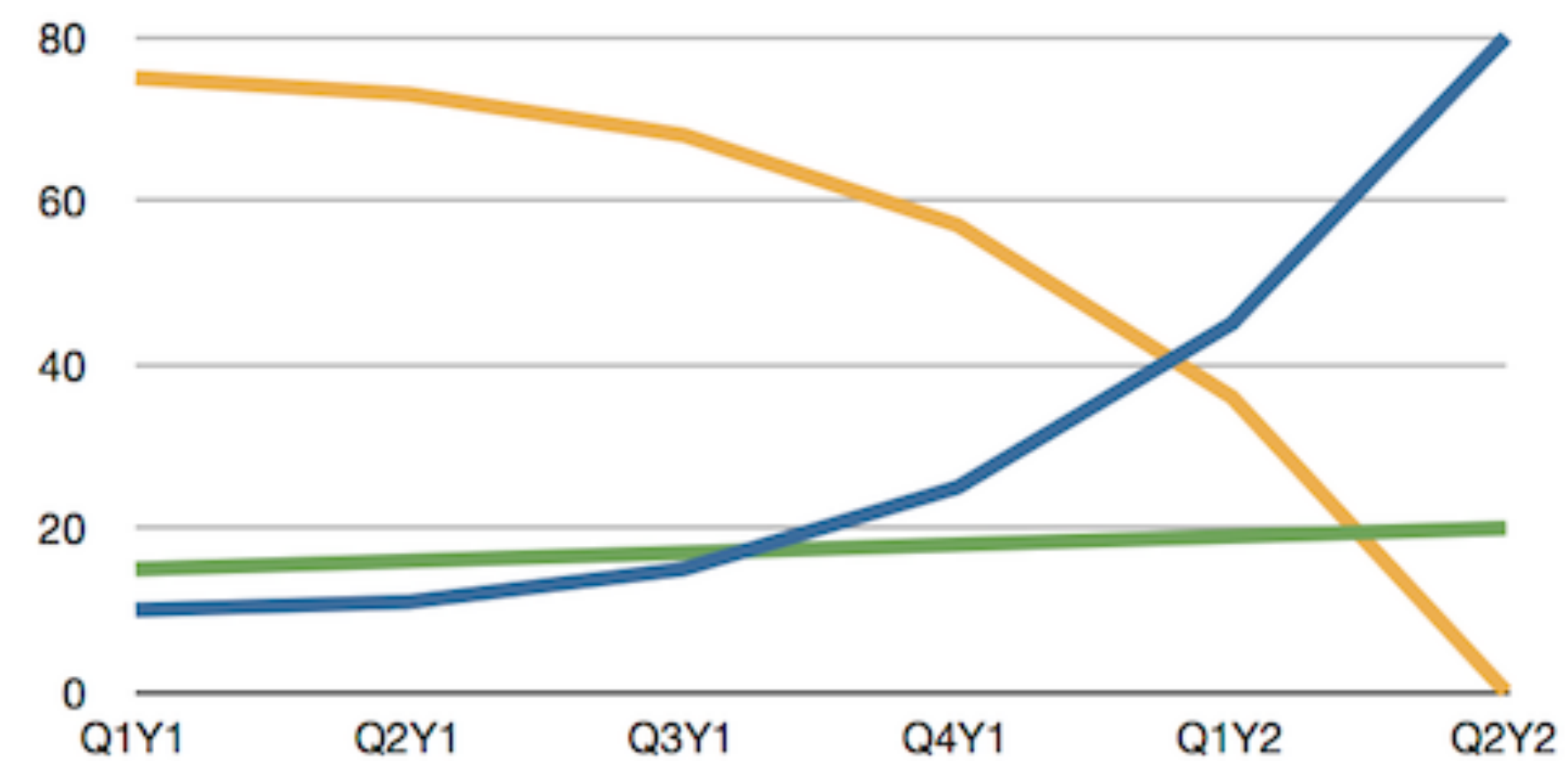# Stacked Area Chart

# 100% Stacked Area Chart



100 % stacked area chart

# Stacked Area vs. Line Graphs
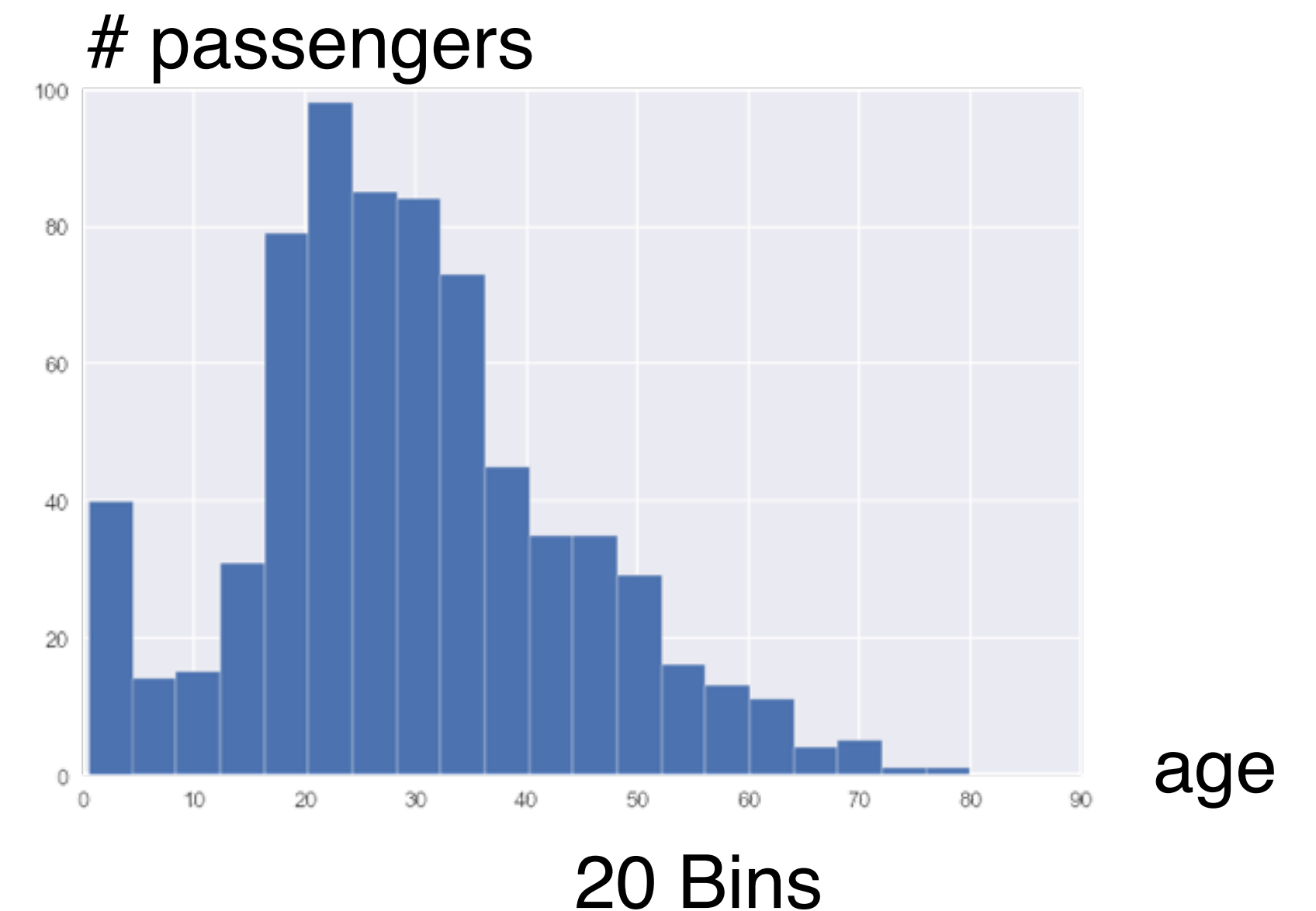


Market Share

leancrew.com &
Practically Efficient

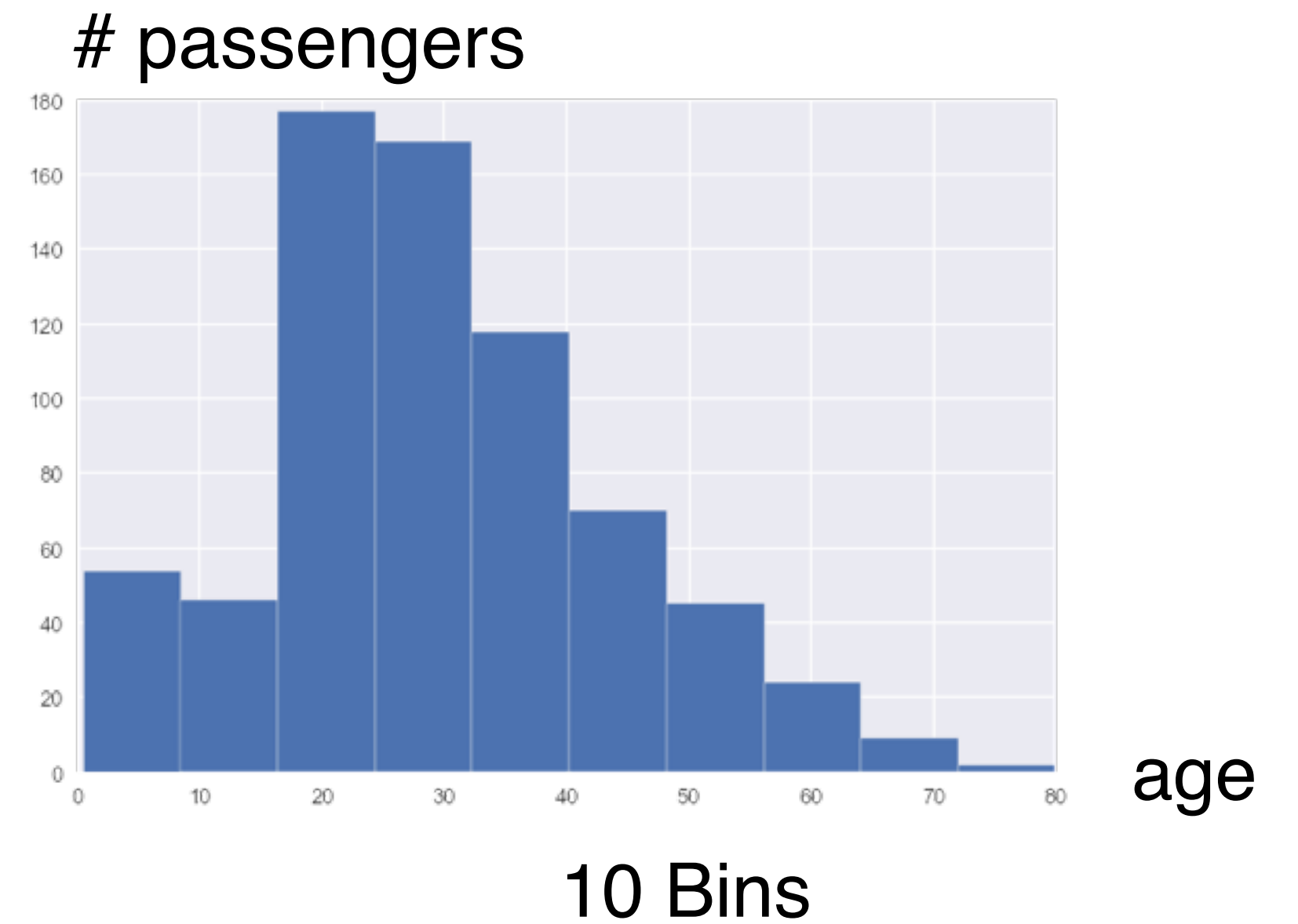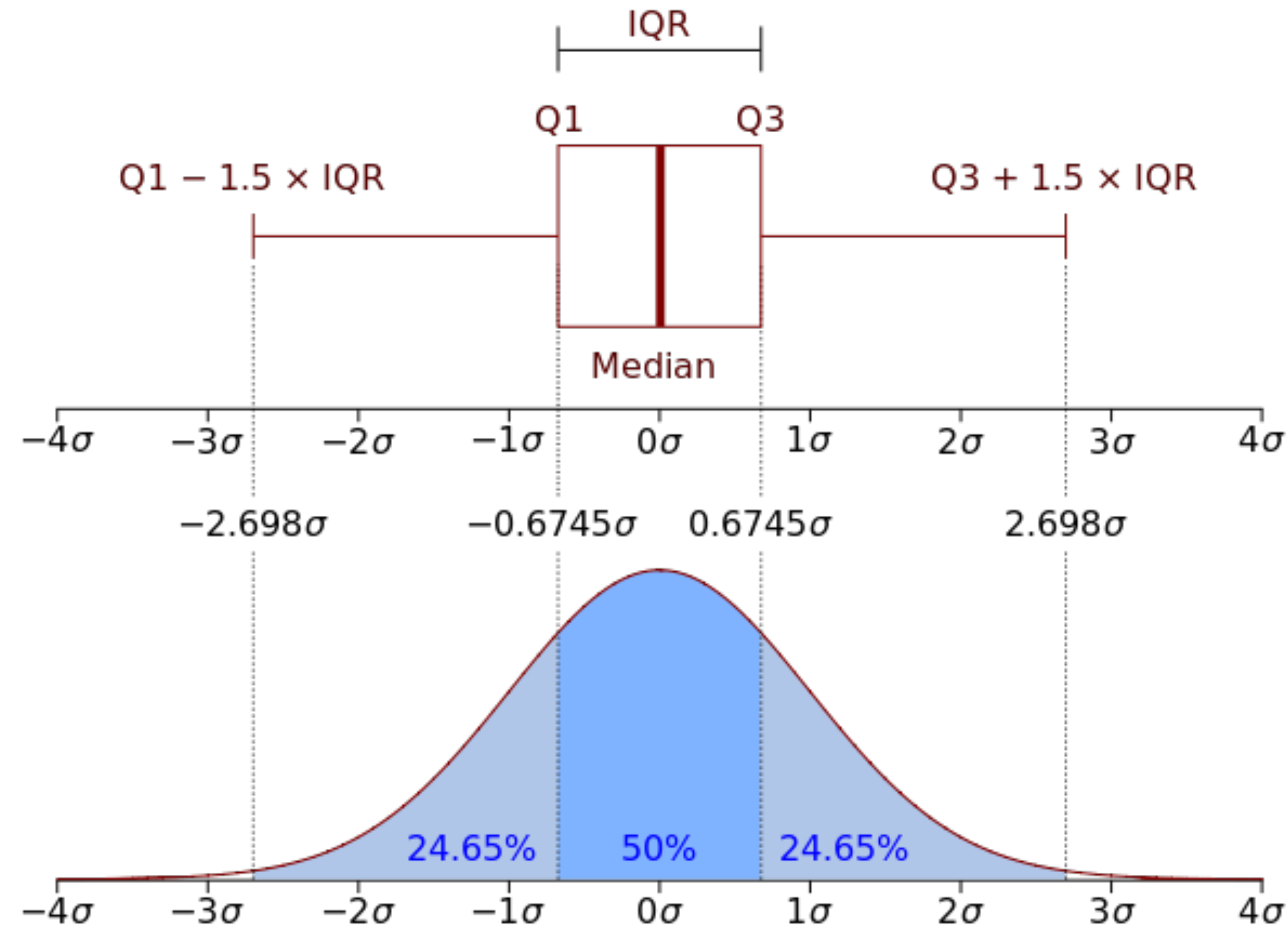# Distributions

# Histogram

#bins hard to predict
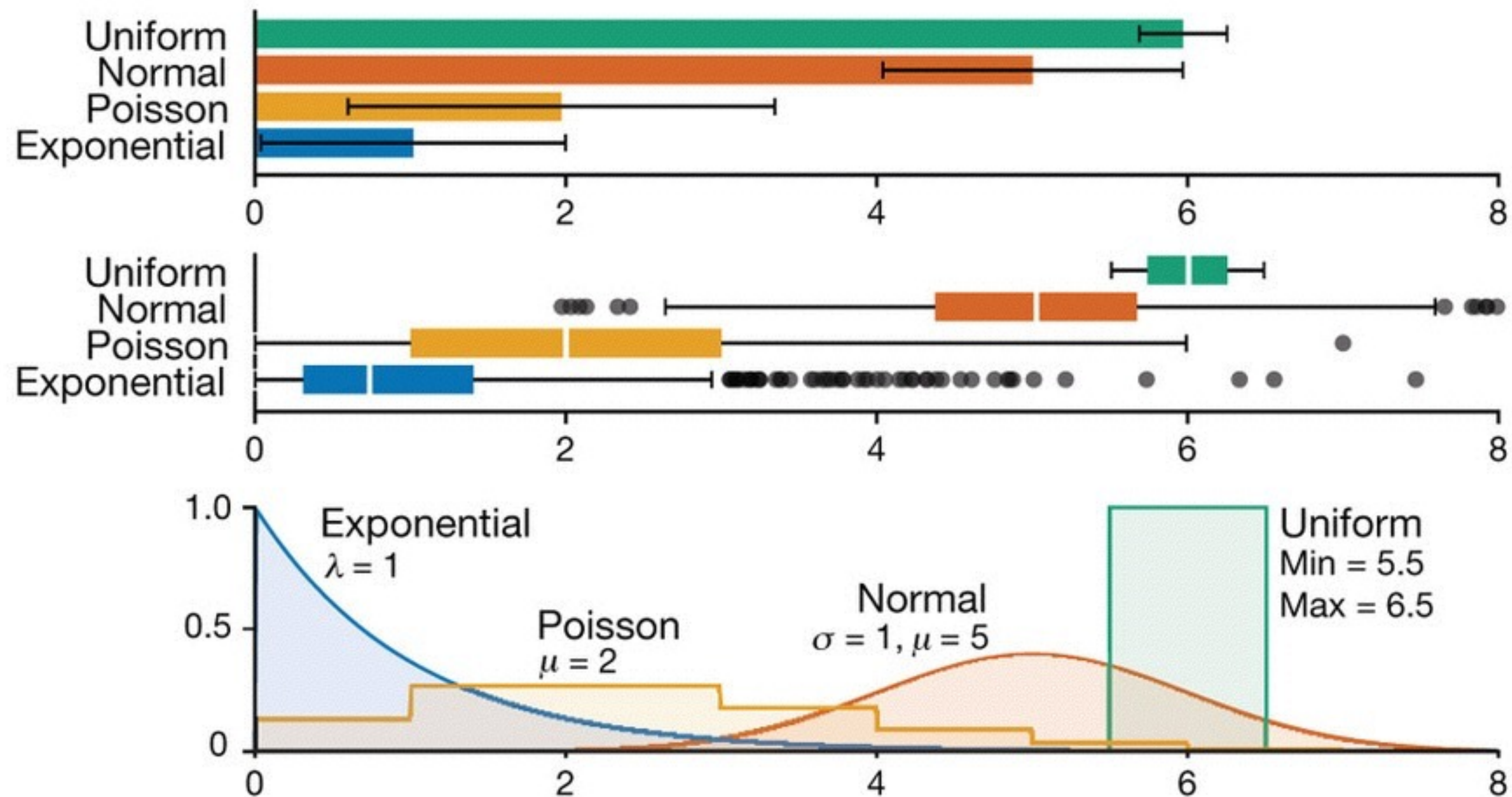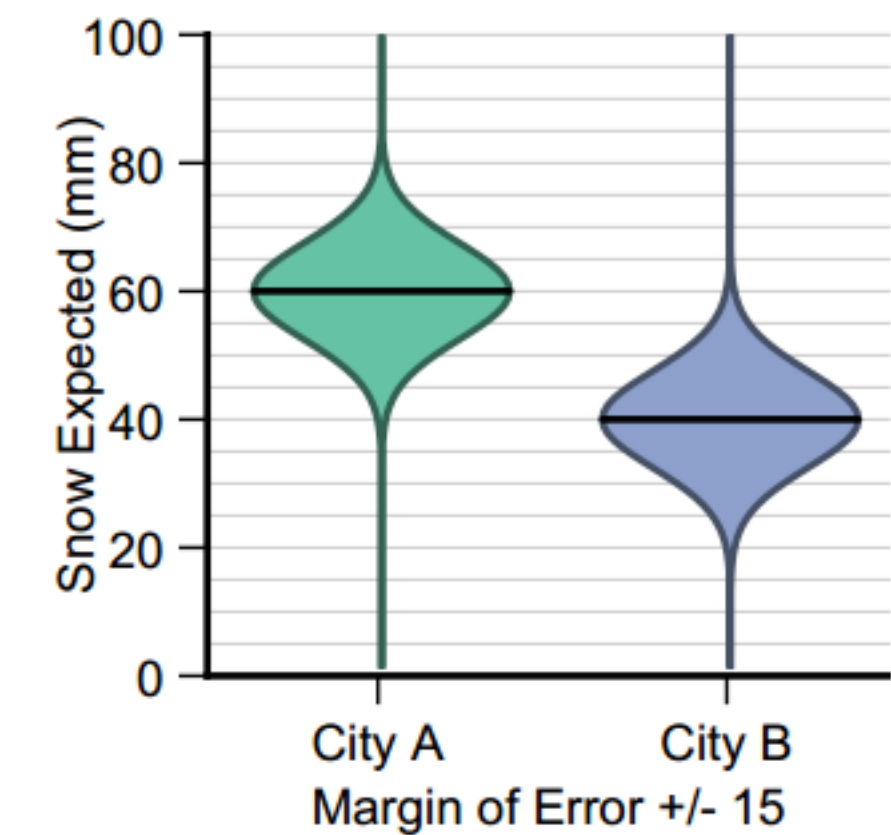
make interactive!

rule of thumb: #bins = sqrt(n)

# passengers



age

10 Bins

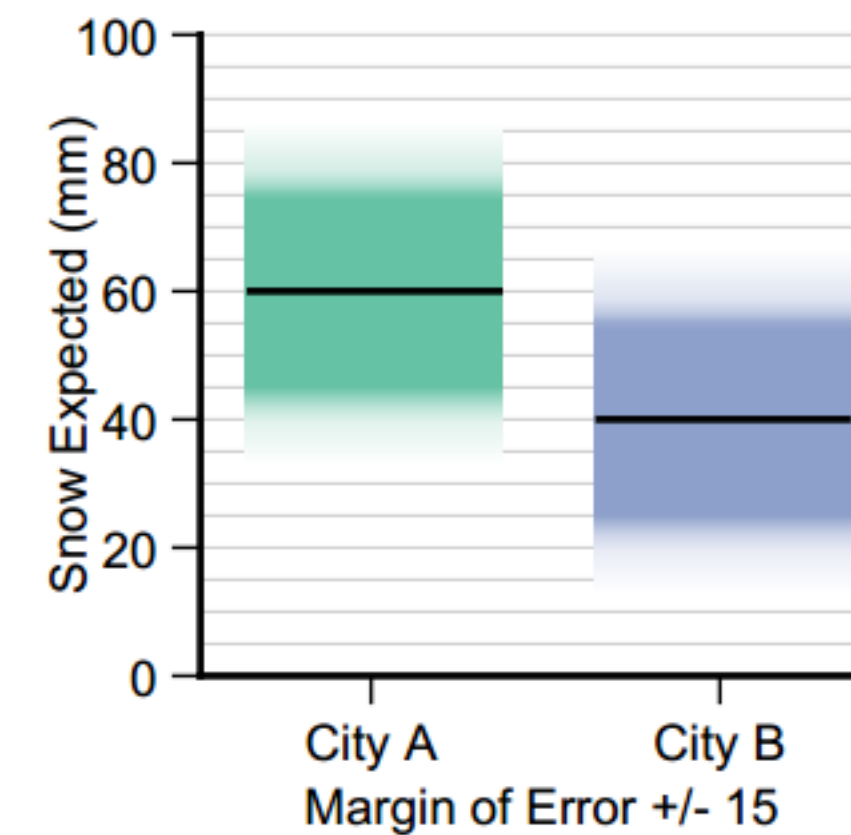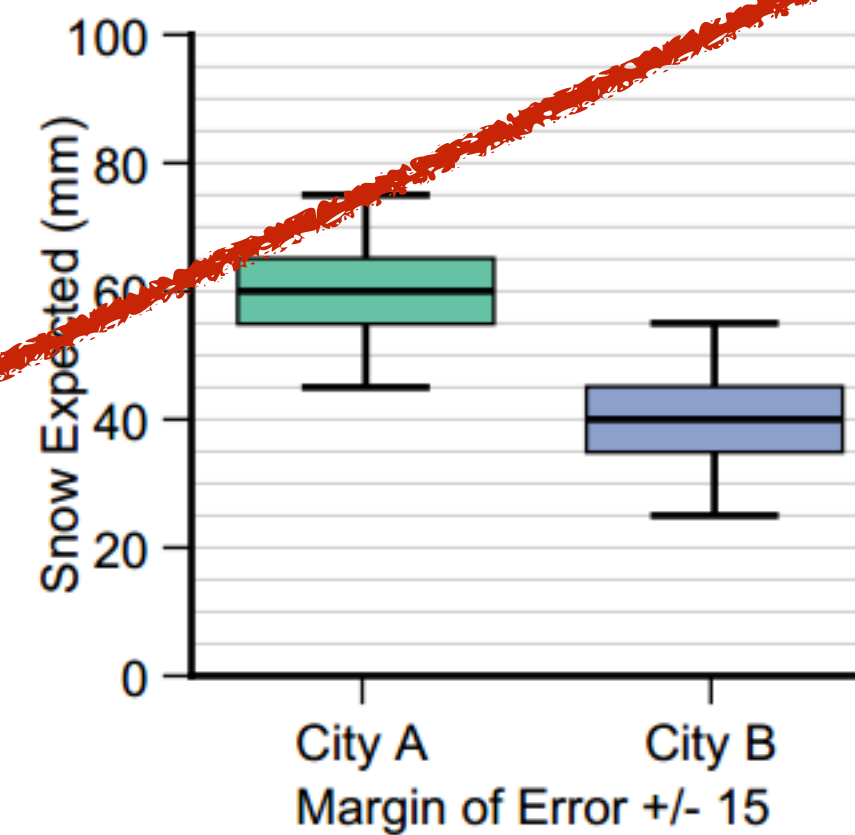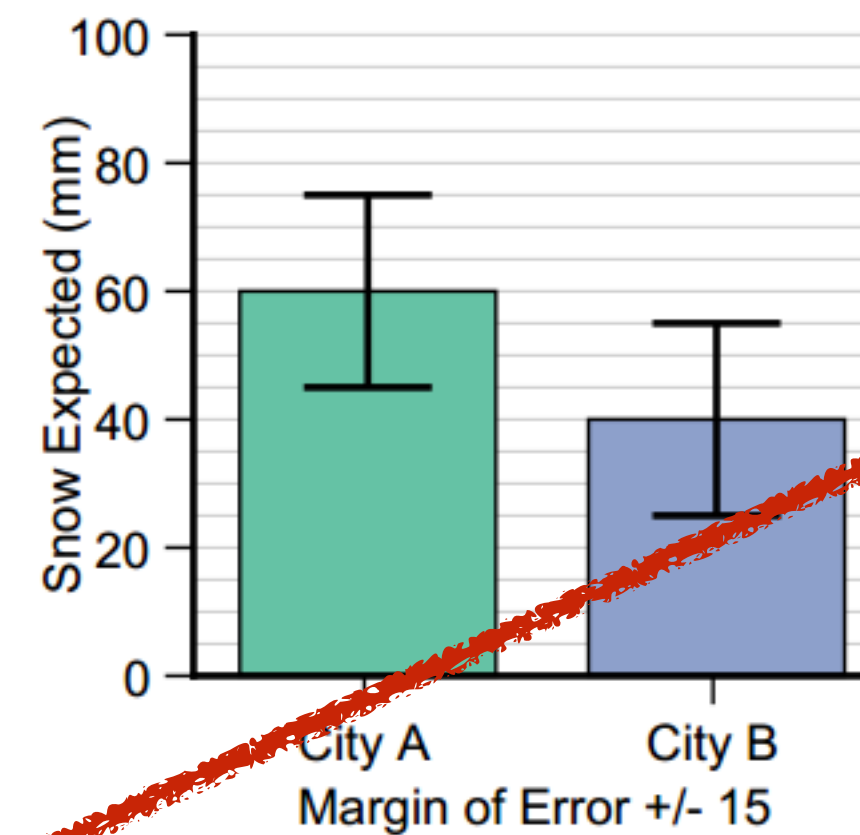# passengers



age

20 Bins

# Box Plots

aka Box-and-Whisker Plot



Wikipedia

# Comparison

# Showing Expected Values & Uncertainty



Error Bars Considered Harmful:
Exploring Alternate Encodings for Mean and Error
Michael Correll, and Michael Gleicher

# Highdimensional Data

# What is High-dimensional Data?

Tabular data, containing

   rows (items)

   columns (attributes or items)

   rows >> columns

|       | Age | Gender | Height |
|-------|-----|--------|--------|
| Bob   | 25  | M      | 181    |
| Alice | 22  | F      | 185    |
| Chris | 19  | M      | 175    |

# High-Dimensional Data Visualization

## How many dimensions?

~50 – tractable with "just" vis

~1000 – need analytical methods

## How many records?

~ 1000 – "just" vis is fine
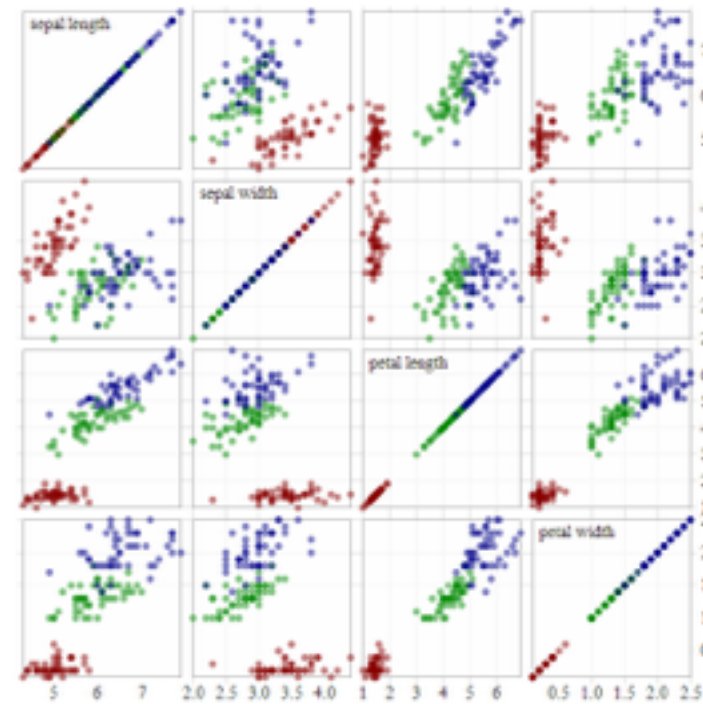
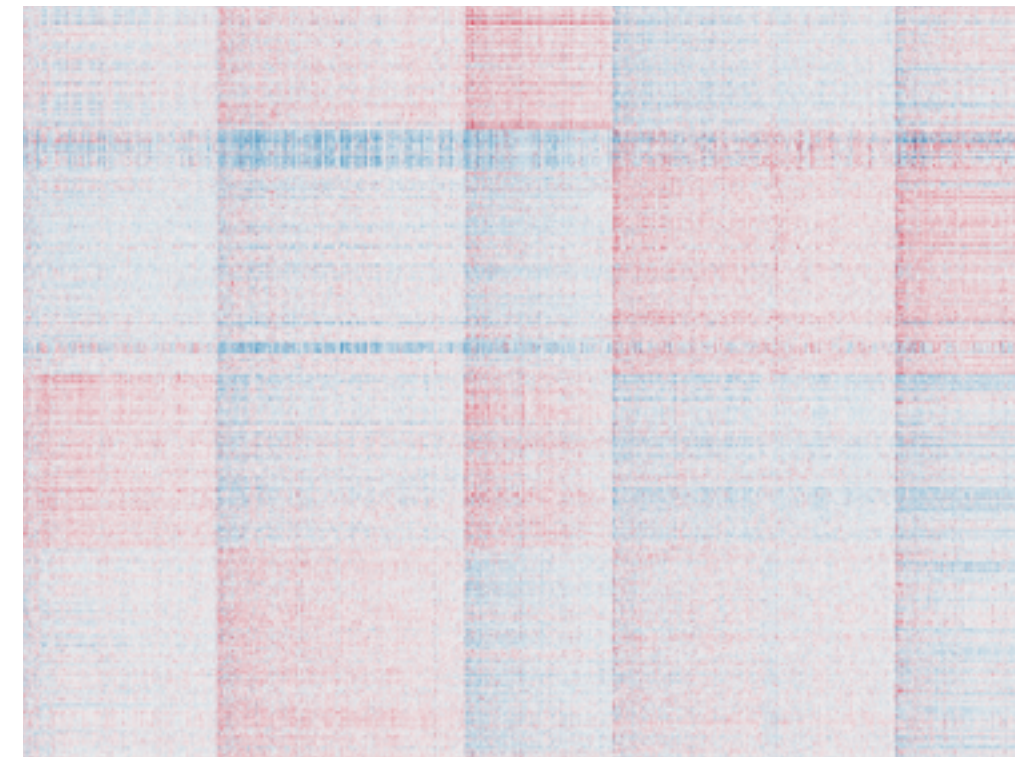>> 10,000 – need analytical methods

**Homogeneity**

**Same data type?**

**Same scales?**

|       | Age | Gender | Height |
|-------|-----|--------|--------|
| *Bob*   | 25  | M      | 181    |
| *Alice* | 22  | F      | 185    |
| *Chris* | 19  | M      | 175    |

|       | BPM 1 | BPM 2 | BPM 3 |
|-------|-------|-------|-------|
| *Bob*   | 65    | 120   | 145   |
| *Alice* | 80    | 135   | 185   |
| *Chris* | 45    | 115   | 135   |

# Analytic Component
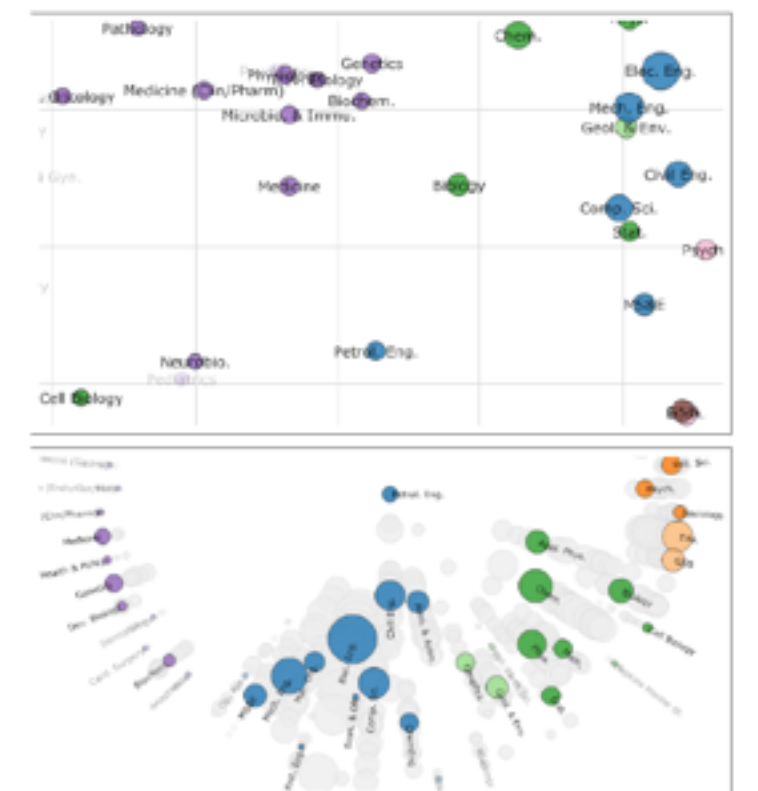

**Scatterplot Matrices**
[Bostock]


**Parallel Coordinates**
[Bostock]


**Pixel-based visualizations /
heat maps**


**Multidimensional Scaling**
[Doerk 2011]


[Chuang 2012]

**no / little analytics**

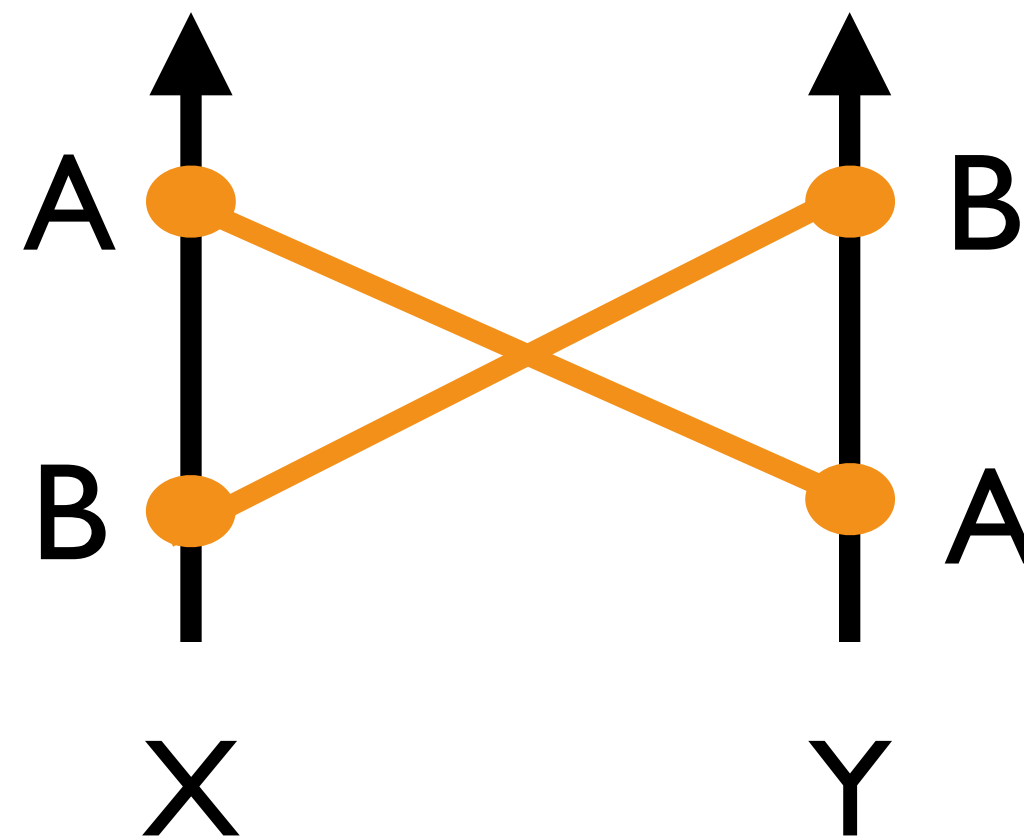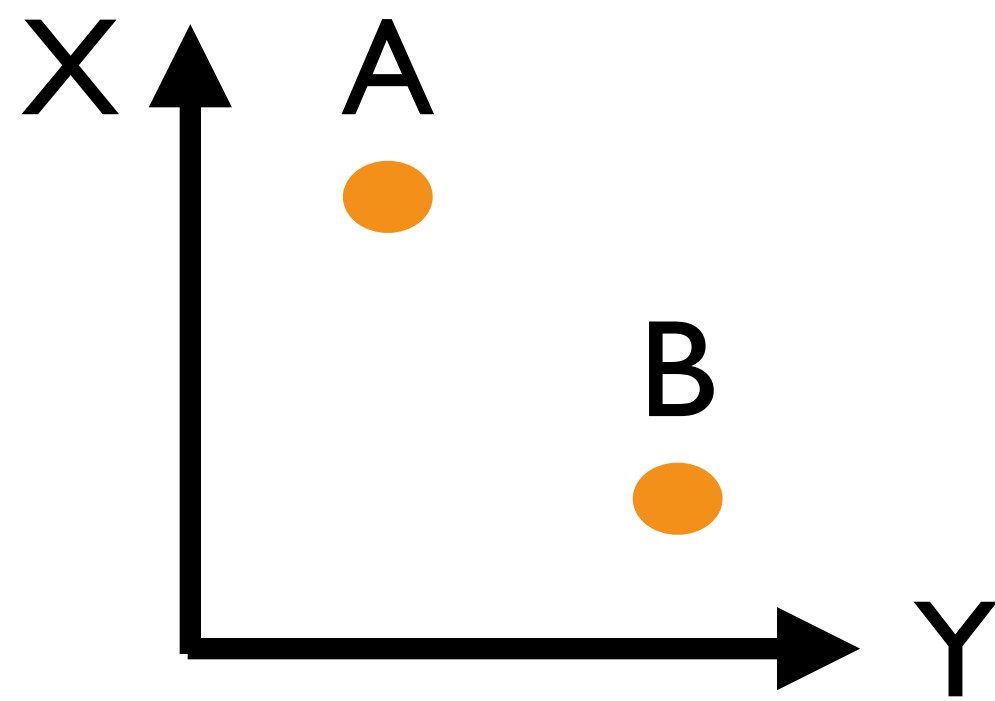**strong analytics
component**

# Geometric Methods

# Parallel Coordinates (PC)

Inselberg 1985

Axes represent attributes

Lines connecting axes represent items
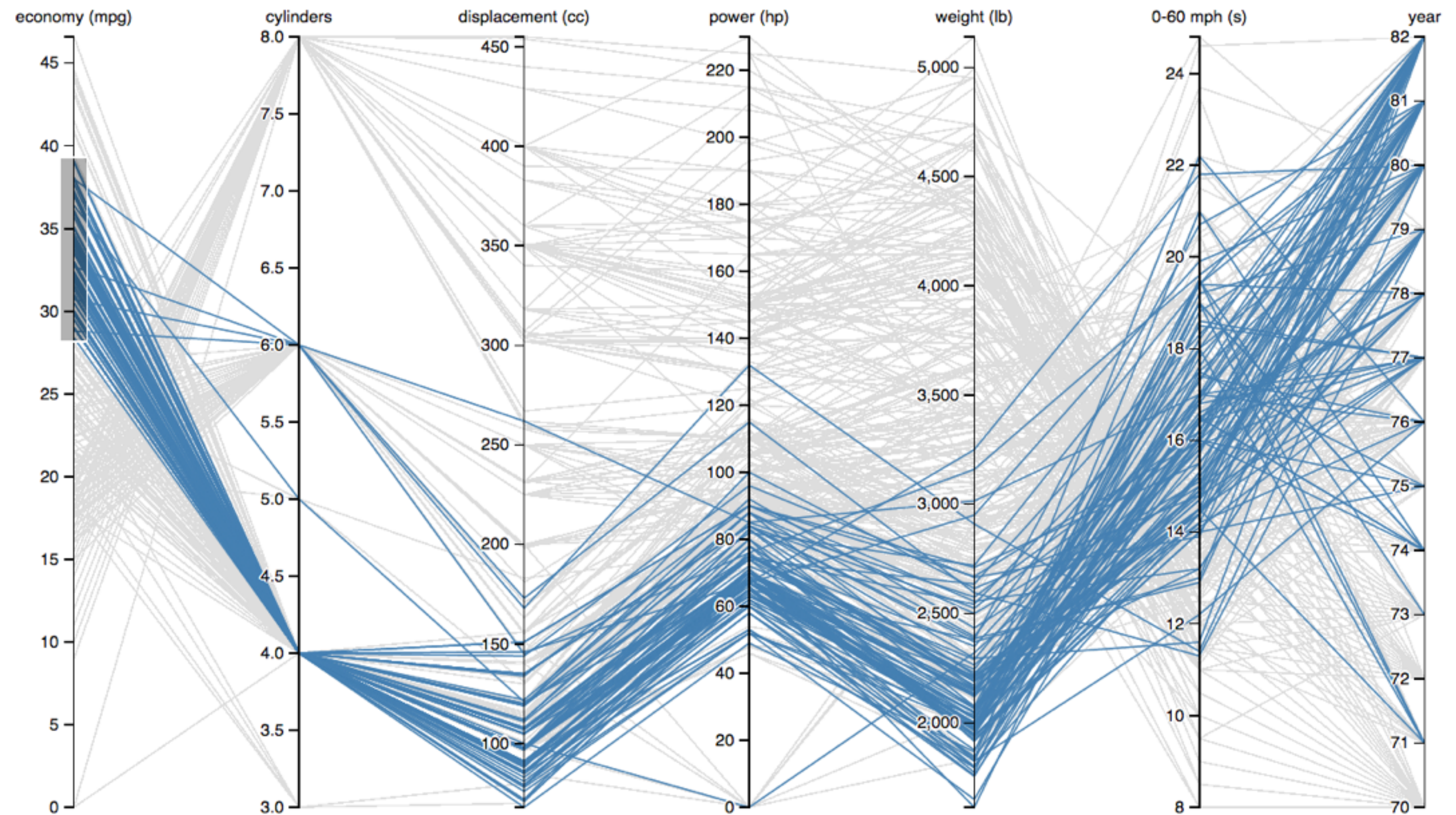
# Parallel Coordinates

Each axis represents dimension

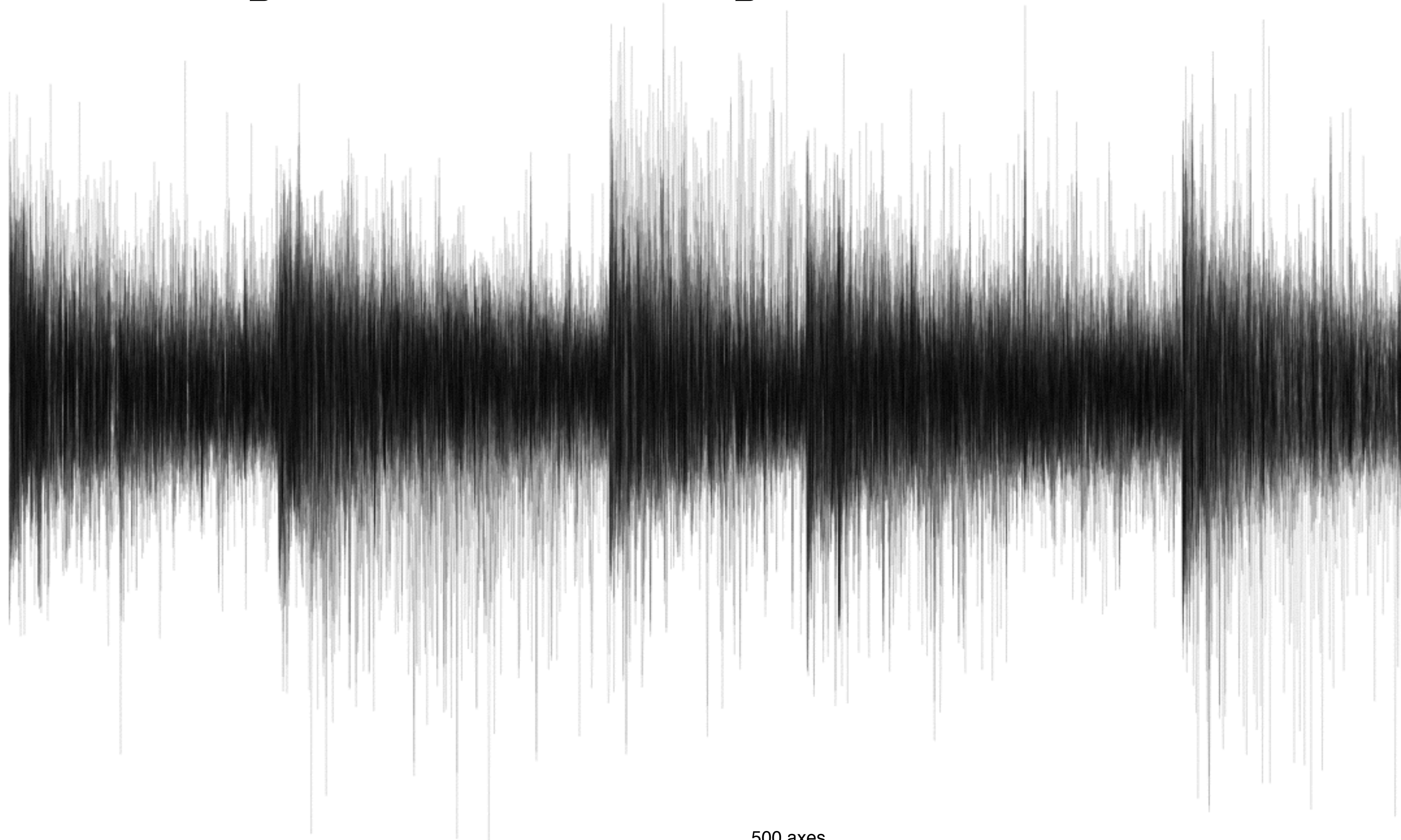Lines connecting axis represent records

Suitable for

all tabular data types

heterogeneous data

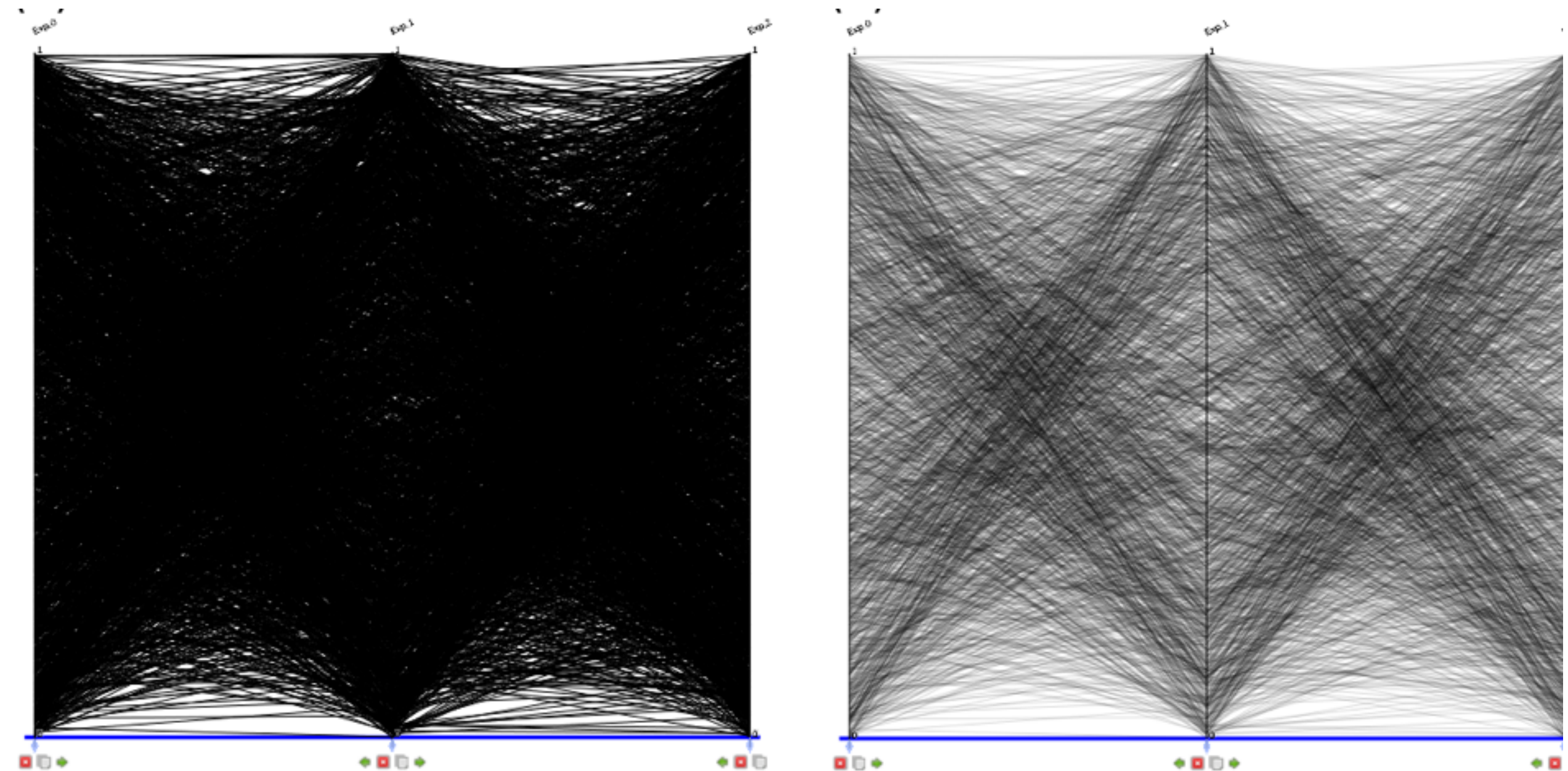# PC Limitation: Scalability to Many Dimensions



500 axes

# PC Limitation: Scalability to Many Items

Solutions:
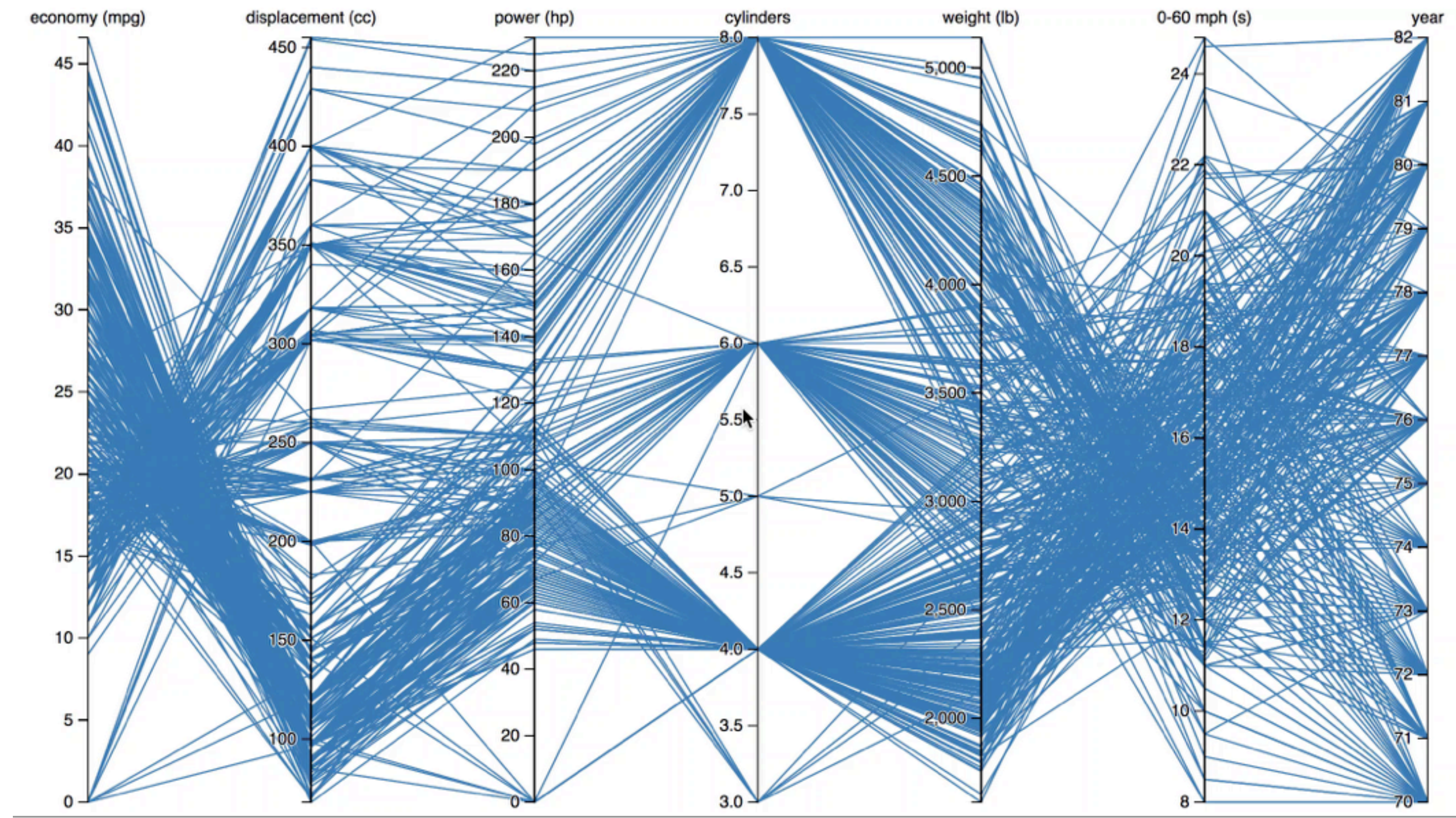
Transparency

Bundling, Clustering

Sampling

# PC Limitations

**Correlations only between adjacent axes**
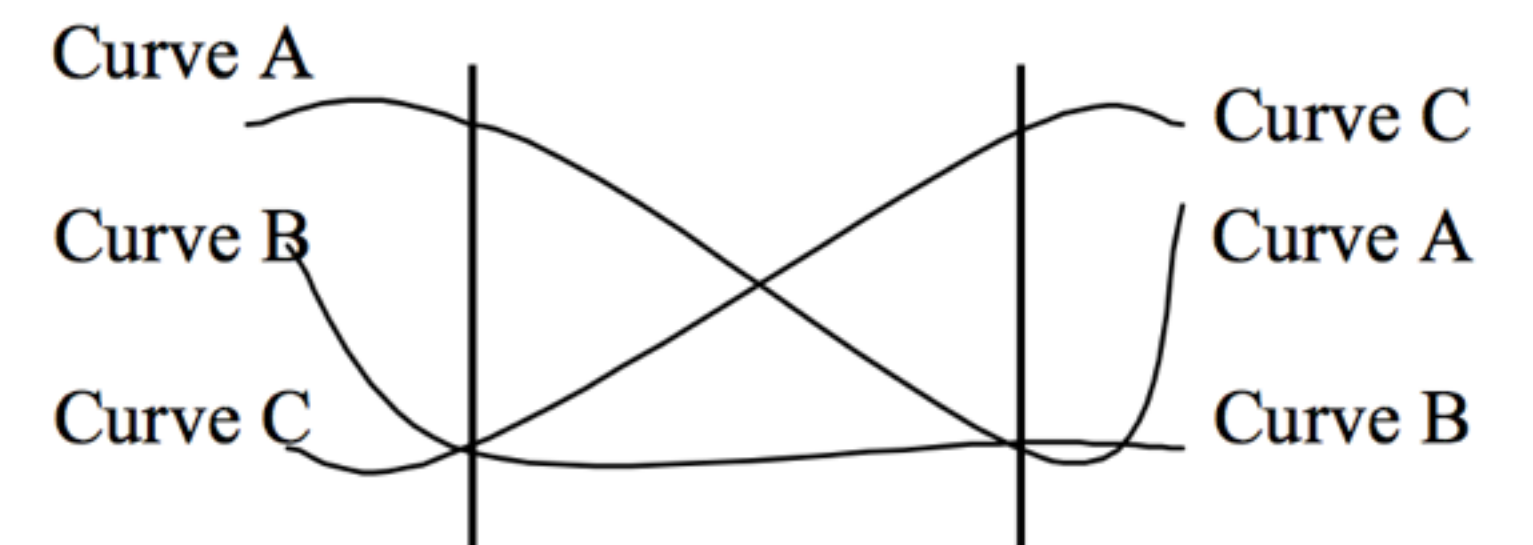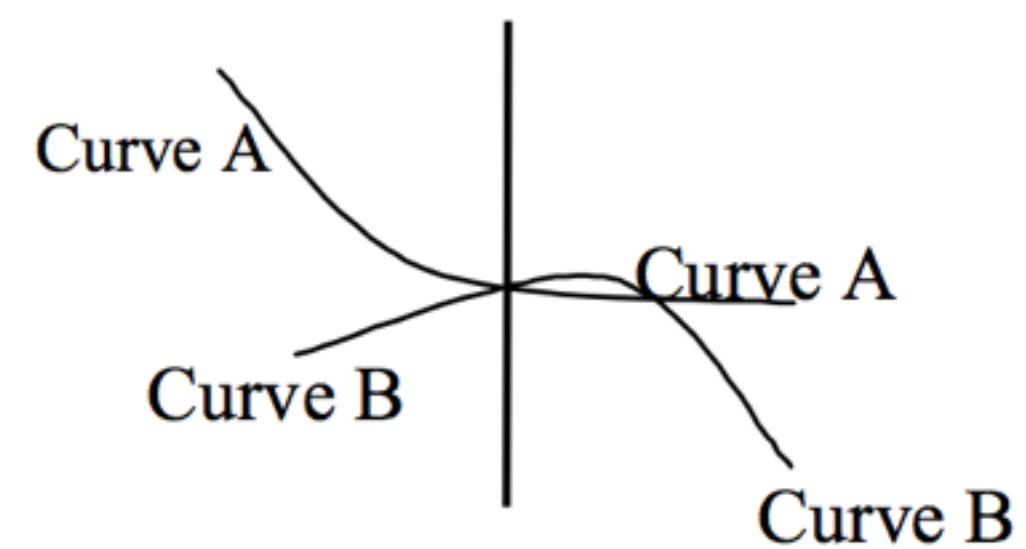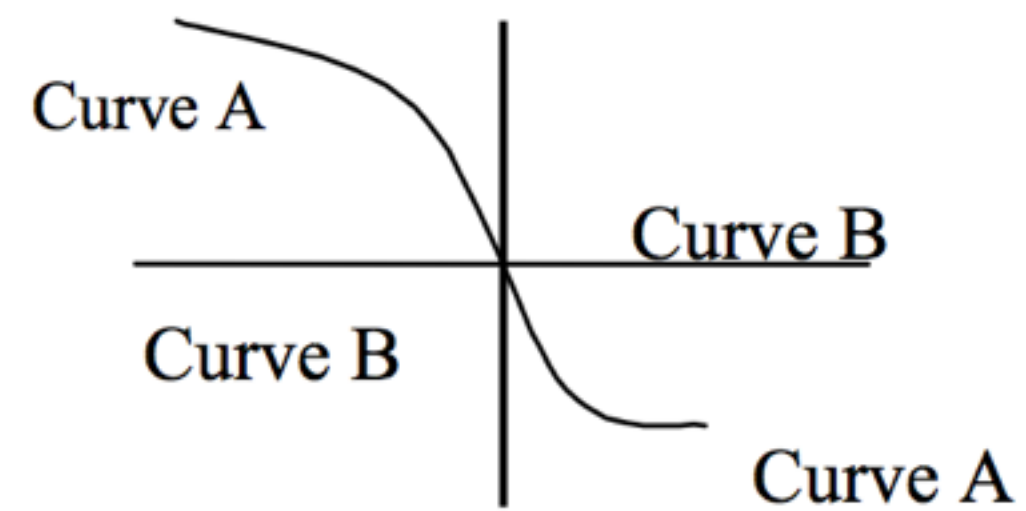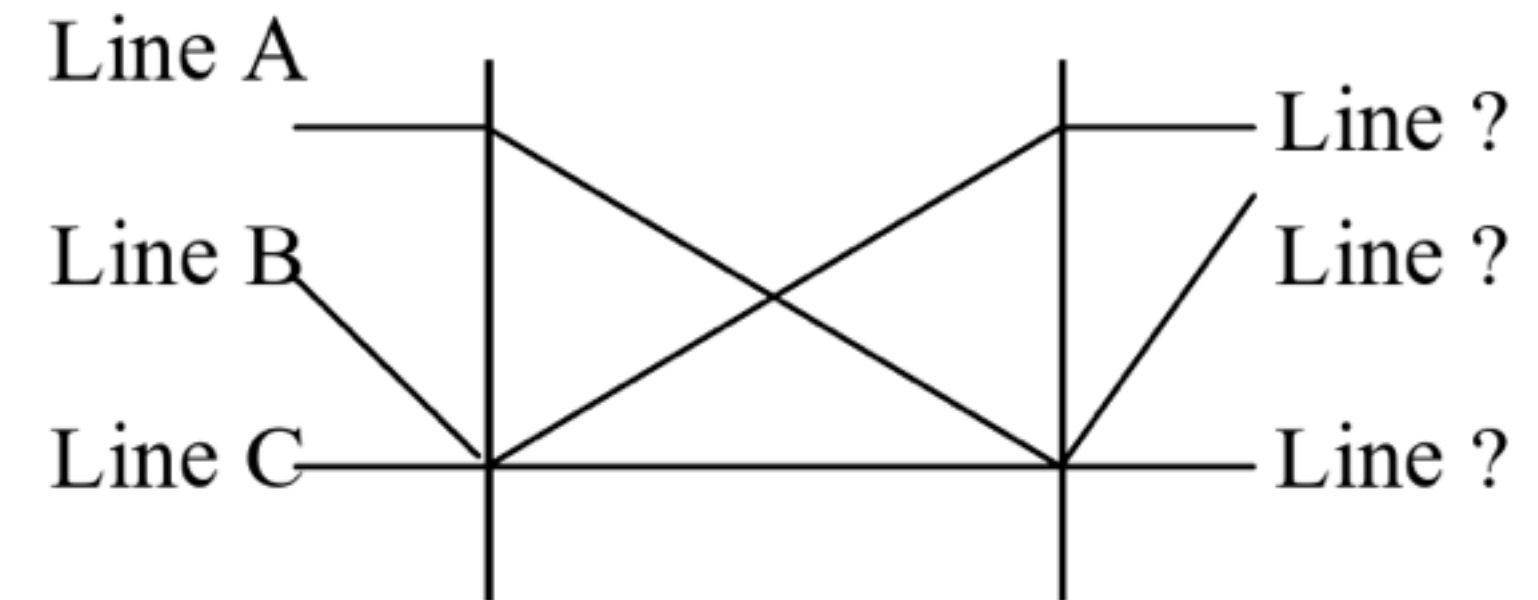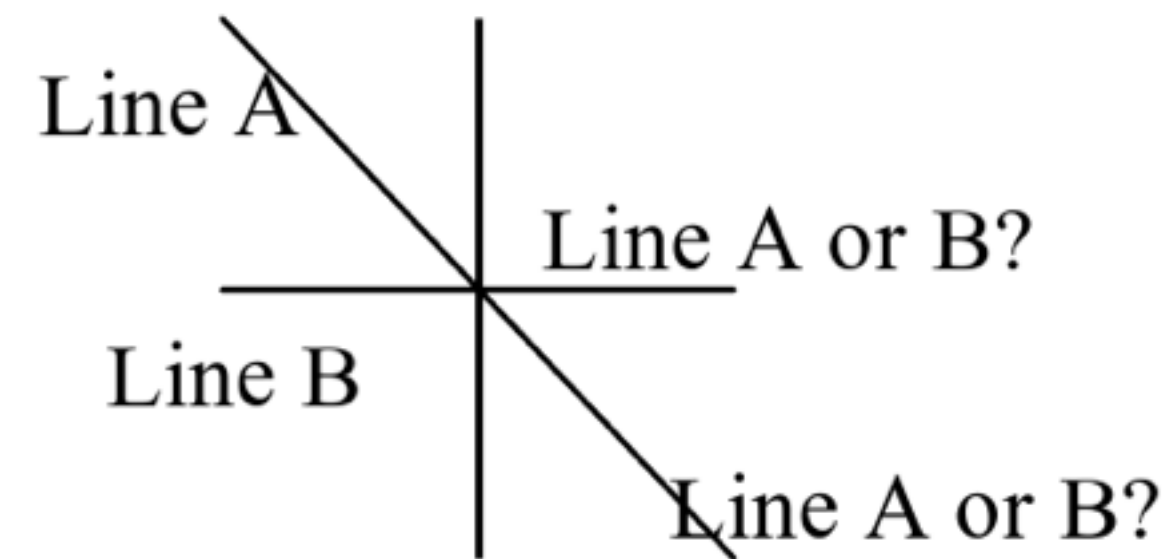
Solution: Interaction

Brushing

Let user change order

# PC Limitation: Ambiguity

Solutions:

Brushing

Curves

# Parallel Coordinates

Shows primarily relationships between adjacent axis

Limited scalability (~50 dimensions, ~1-5k records)

Transparency of lines

Interaction is crucial

Axis reordering

Brushing

Filtering
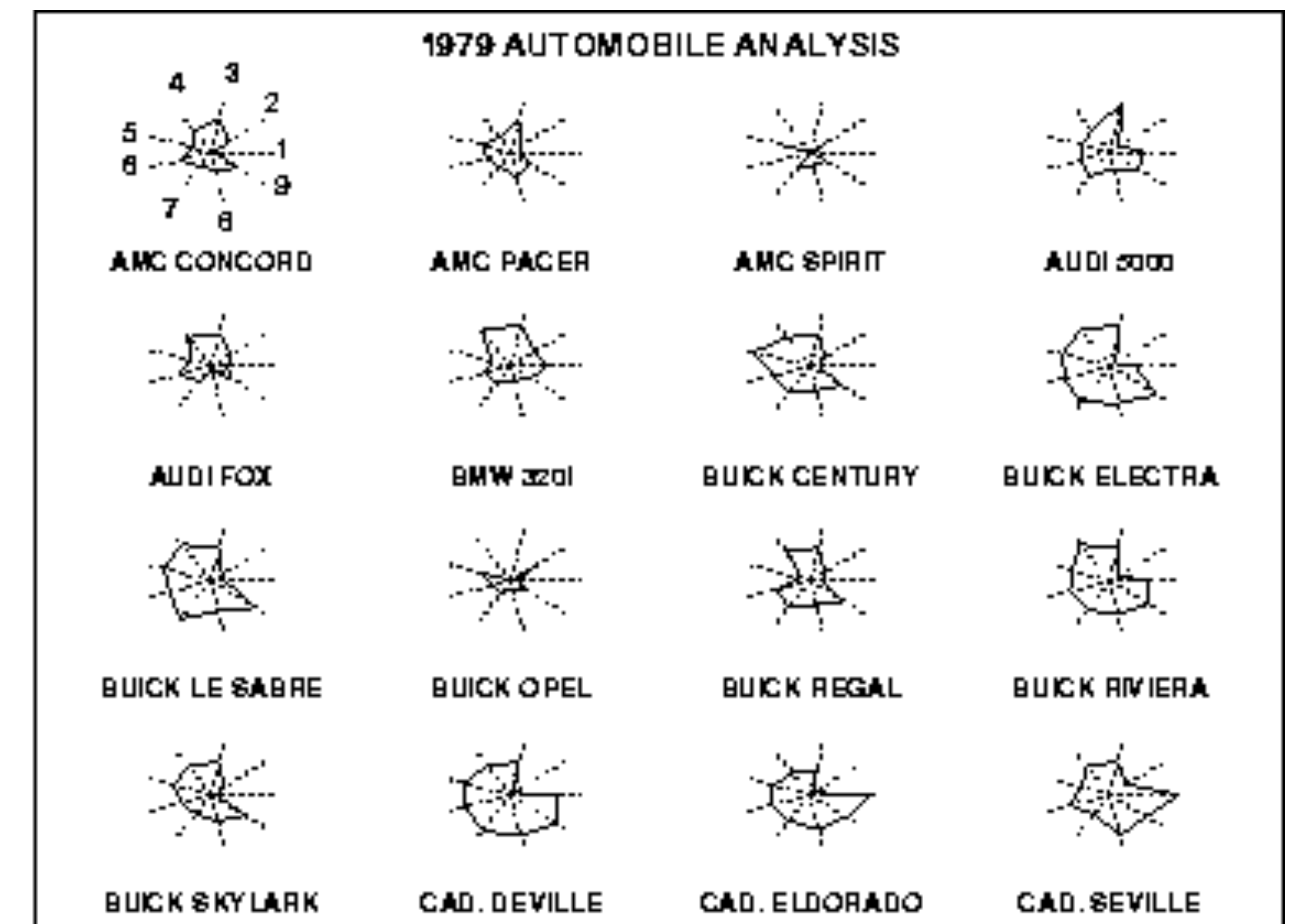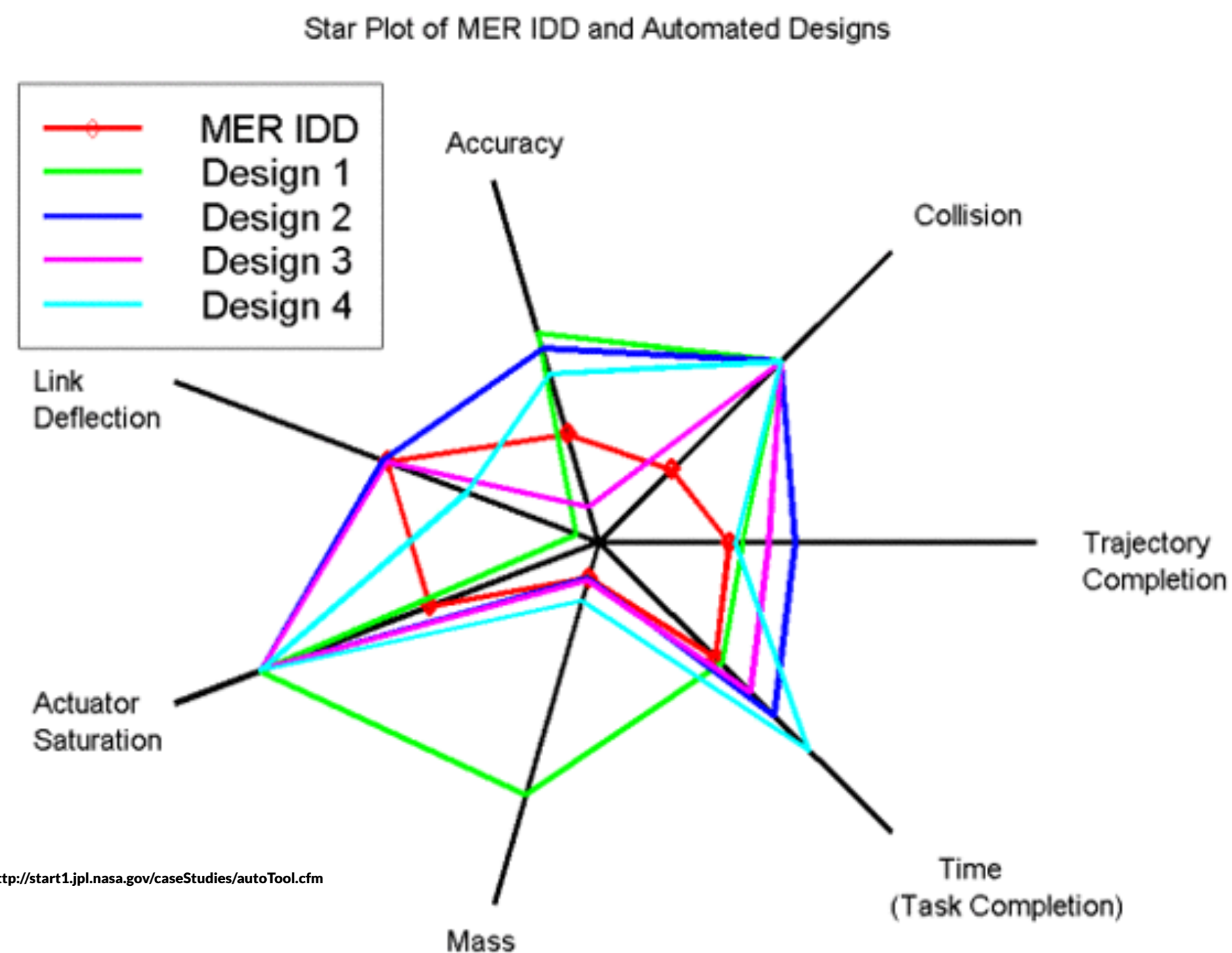
**Algorithmic support:**

Choosing dimensions

Choosing order
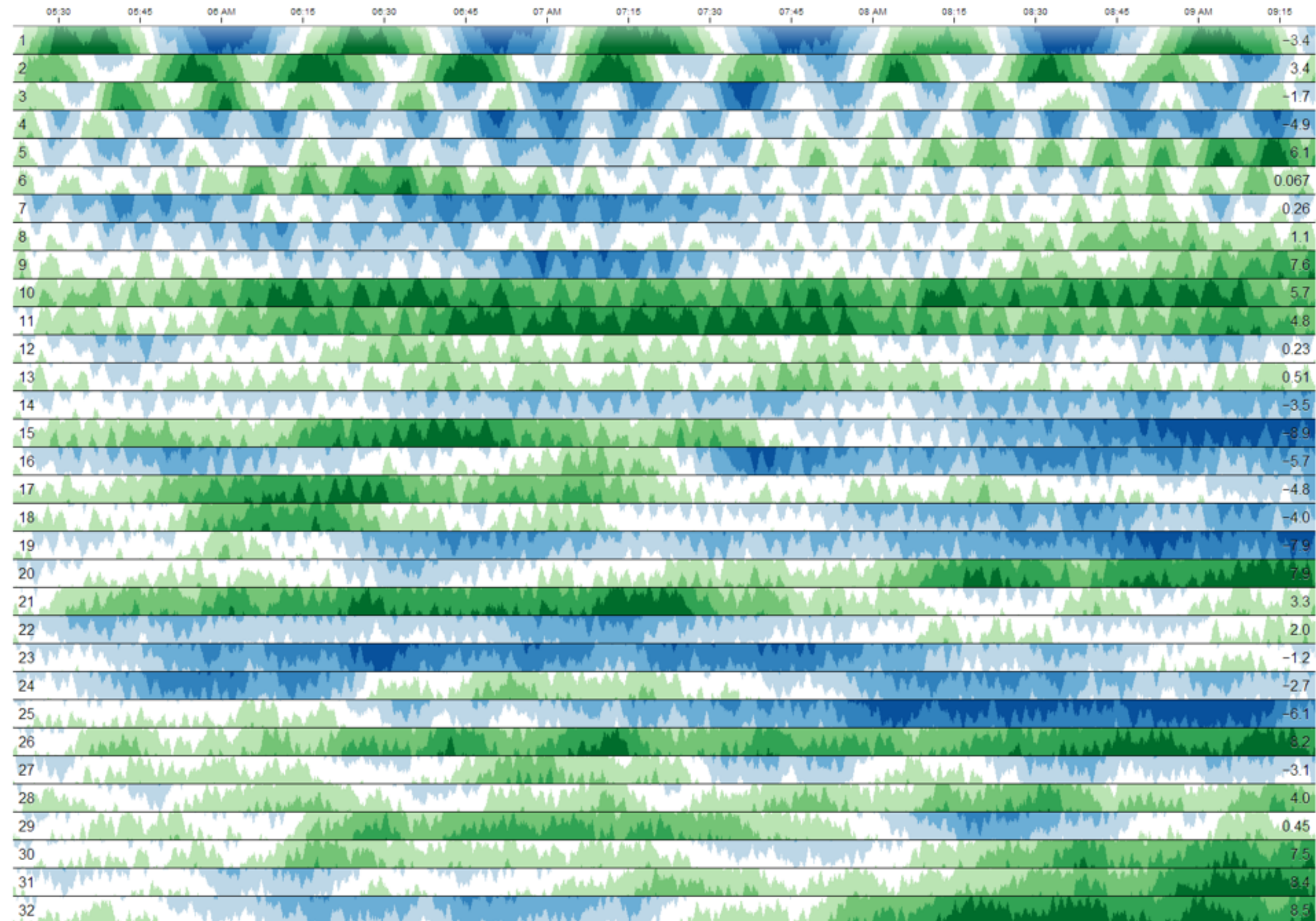
Clustering & aggregating records

# Star Plot

[Coekin1969]

Similar to parallel coordinates

Radiate from a common origin



Star Plot of MER IDD and Automated Designs

MER IDD
Design 1
Design 2
Design 3
Design 4

Accuracy
Collision
Trajectory Completion
Time (Task Completion)
Mass
Actuator Saturation
Link Deflection

http://start1.jpl.nasa.gov/caseStudies/autoTool.cfm



1979 AUTOMOBILE ANALYSIS

AMC CONCORD    AMC PACER    AMC SPIRIT    AUDI 5000
AUDI FOX    BMW 320i    BUICK CENTURY    BUICK ELECTRA
BUICK LE SABRE    BUICK OPEL    BUICK REGAL    BUICK RIVIERA
BUICK SKYLARK    CAD. DEVILLE    CAD. ELDORADO    CAD. SEVILLE

http://www.itl.nist.gov/div898/handbook/eda/section3/starplot.htm

http://bl.ocks.org/kevinschaul/raw/8833989/

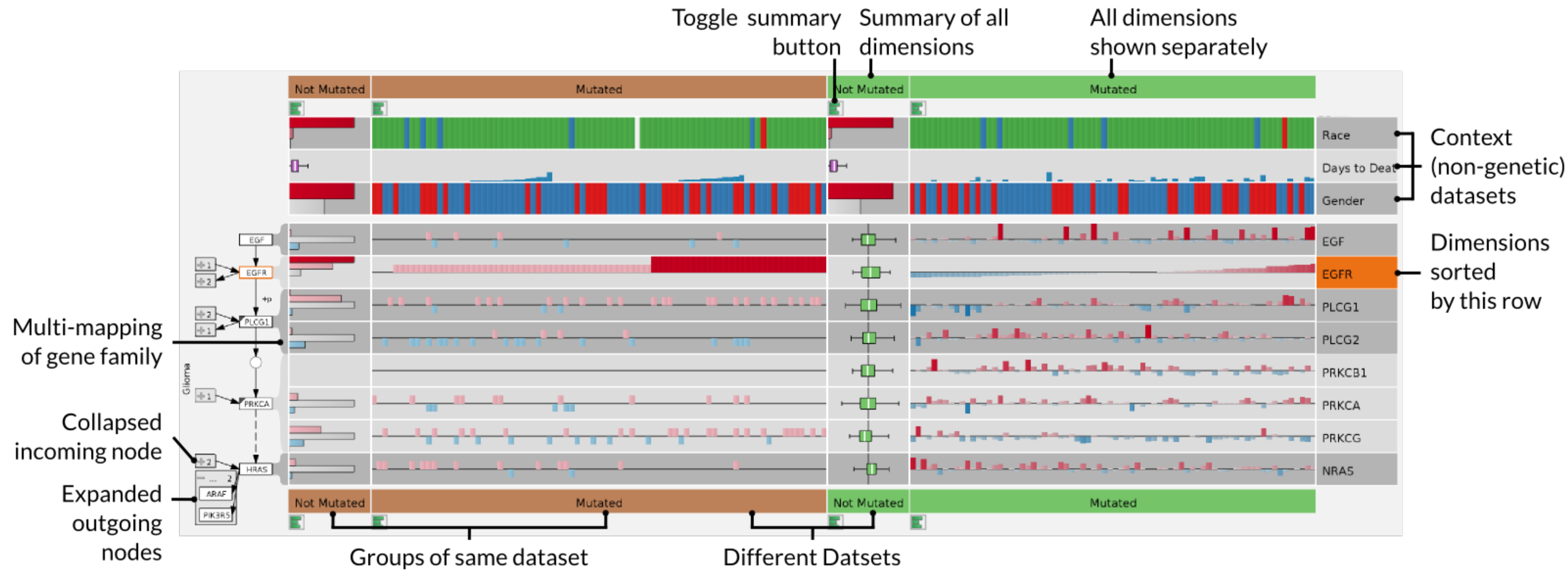# Multiple Line Charts

# Combining Various Charts

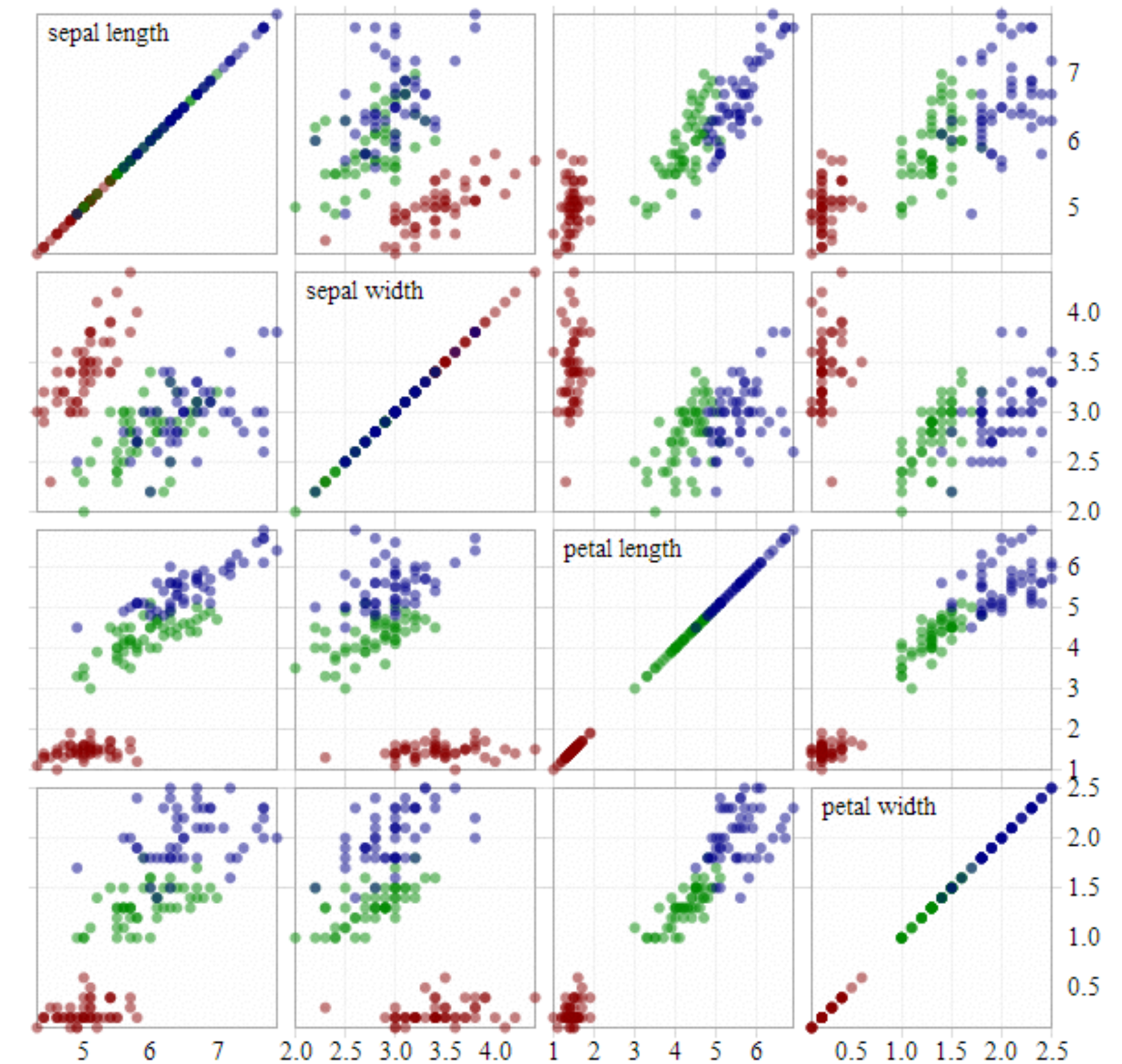# Scatterplot Matrices (SPLOM)

Matrix of size d*d

Each row/column is one dimension

Each cell plots a scatterplot of two dimensions

# Scatterplot Matrices

Limited scalability (~20 dimensions, ~500-1k records)

Brushing is important

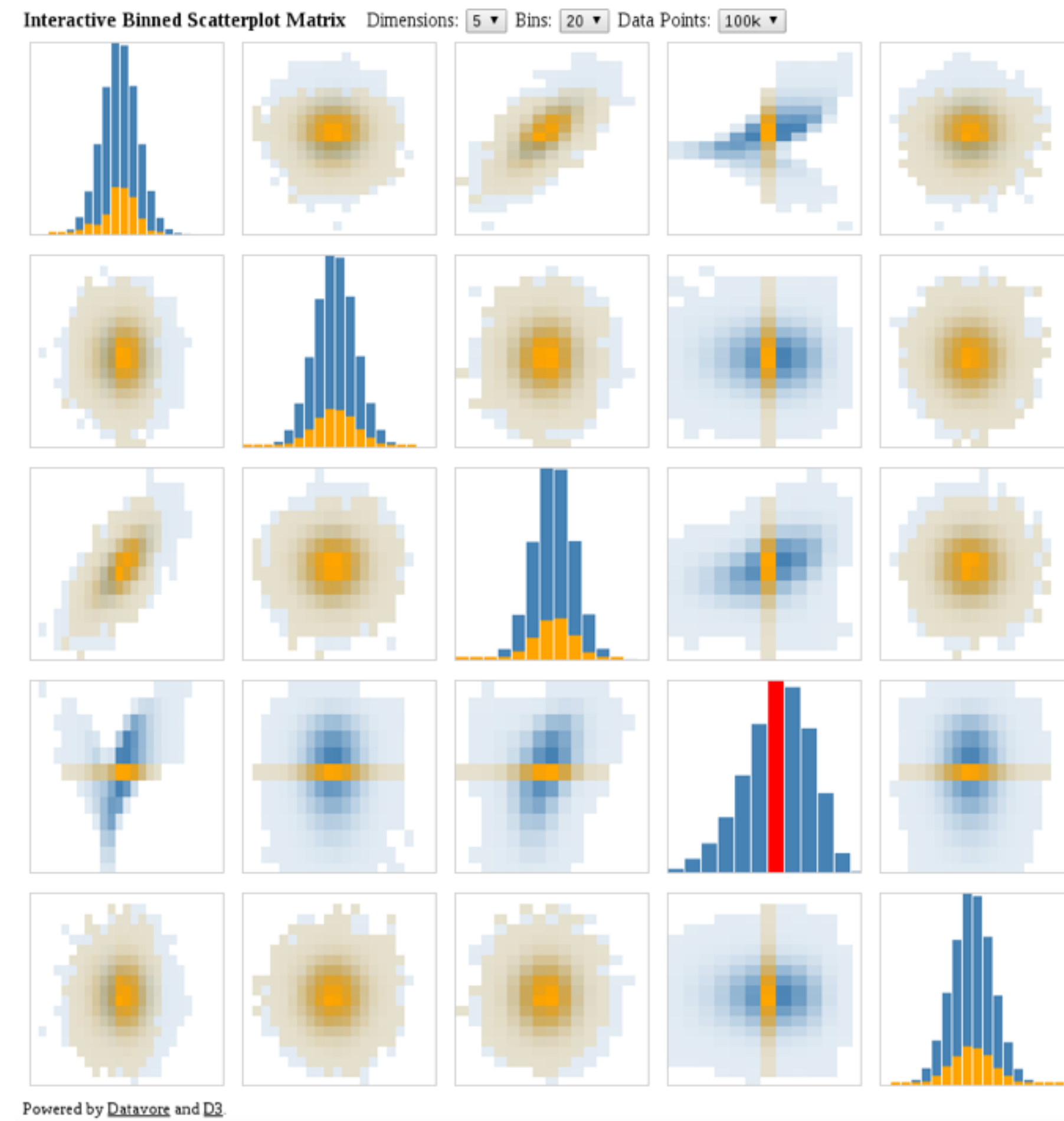Often combined with "Focus Scatterplot" as F+C technique

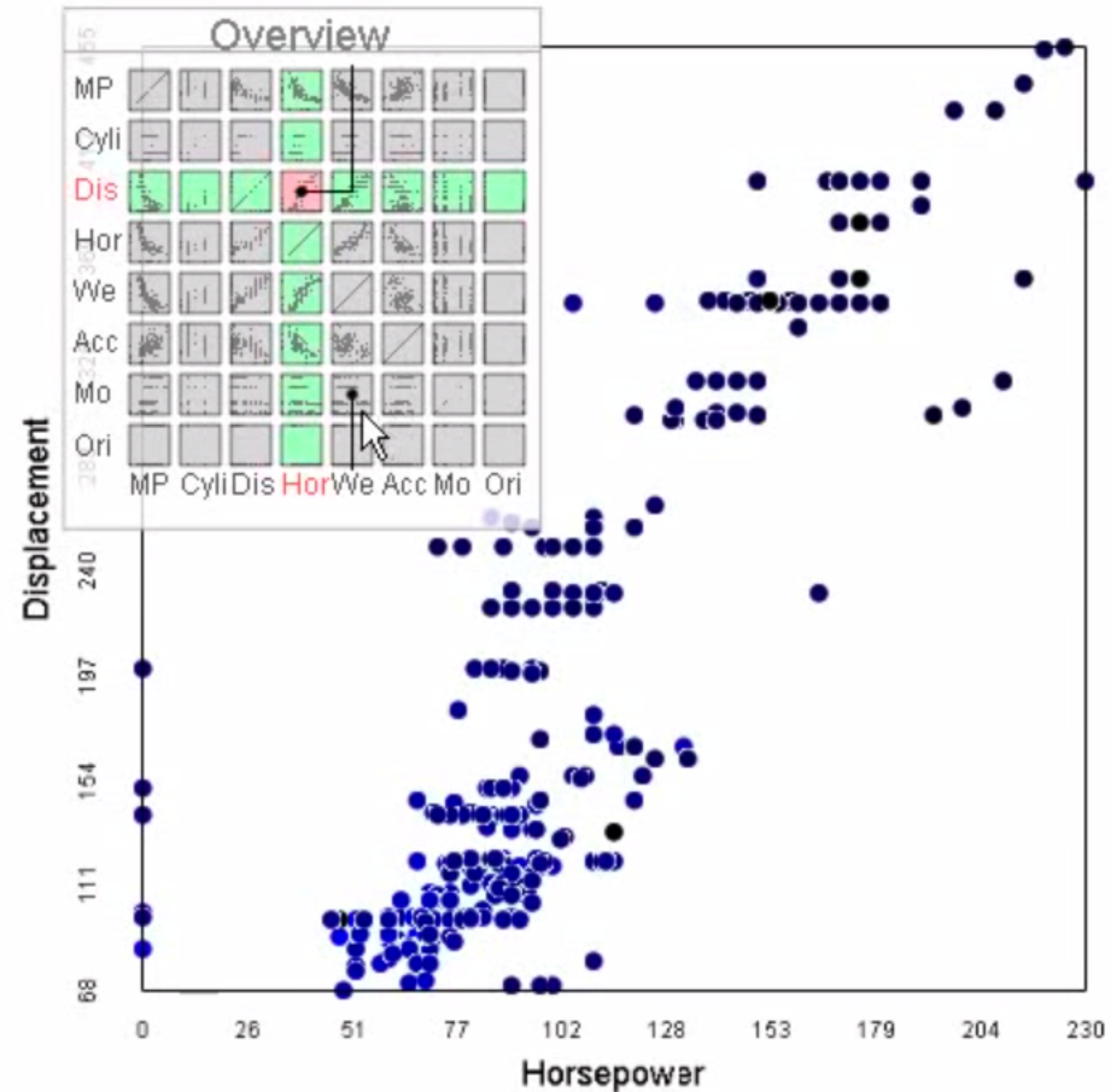**Algorithmic approaches:**

Clustering & aggregating records

Choosing dimensions
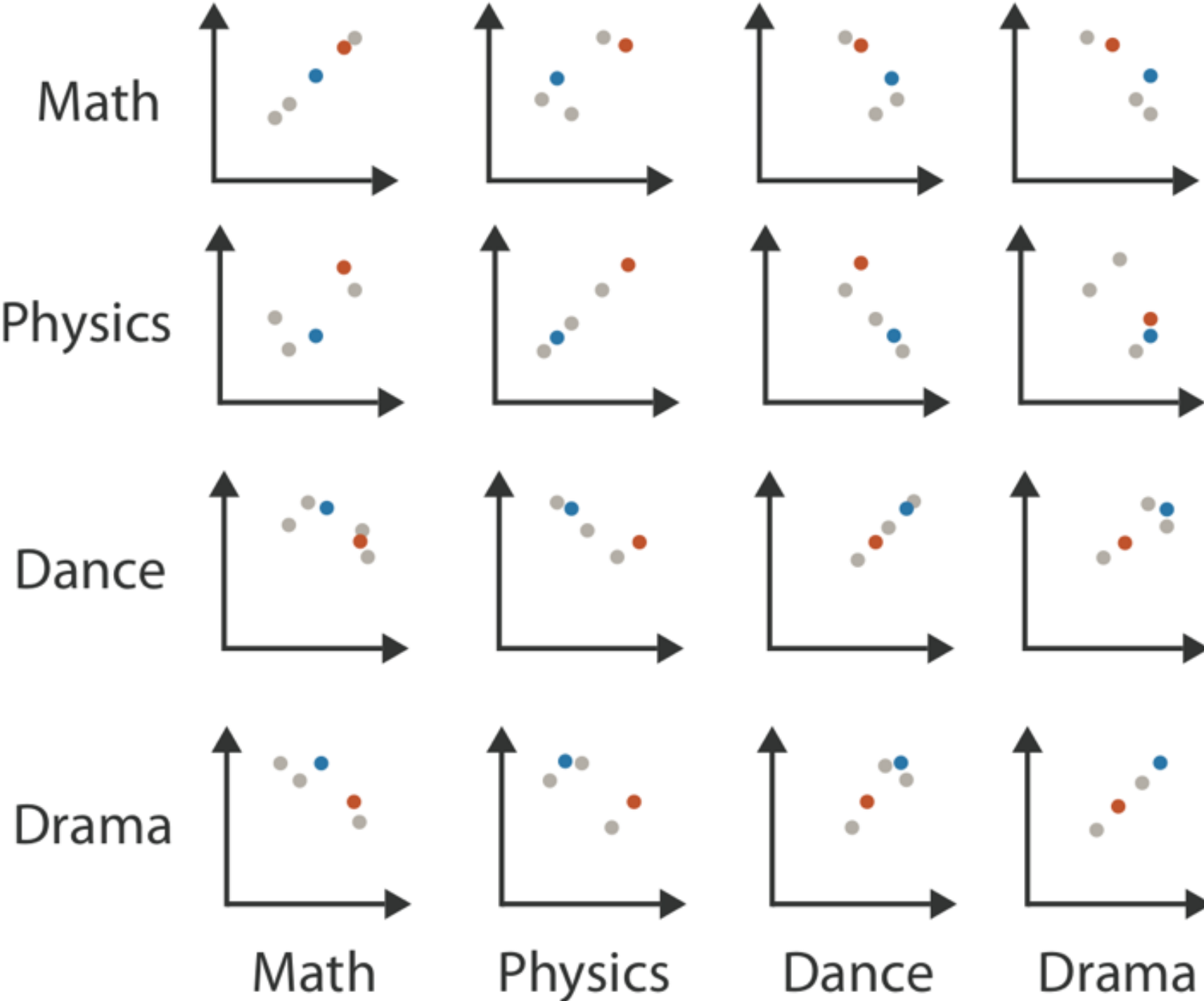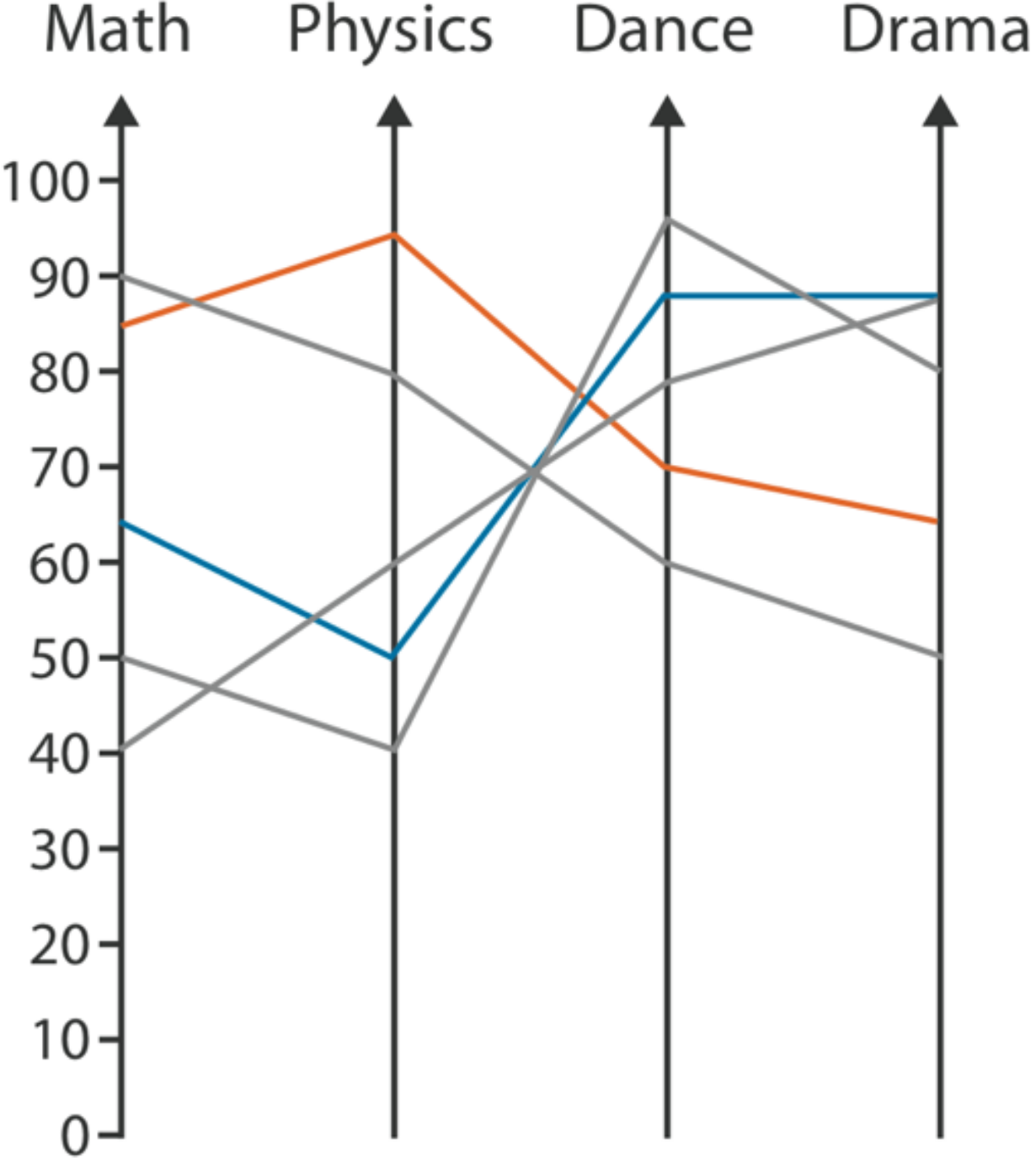
Choosing order

# SPLOM Aggregation - Heat Map



**Datavore:** http://vis.stanford.edu/projects/datavore/splom/

# SPLOM F+C, Navigation

# Table

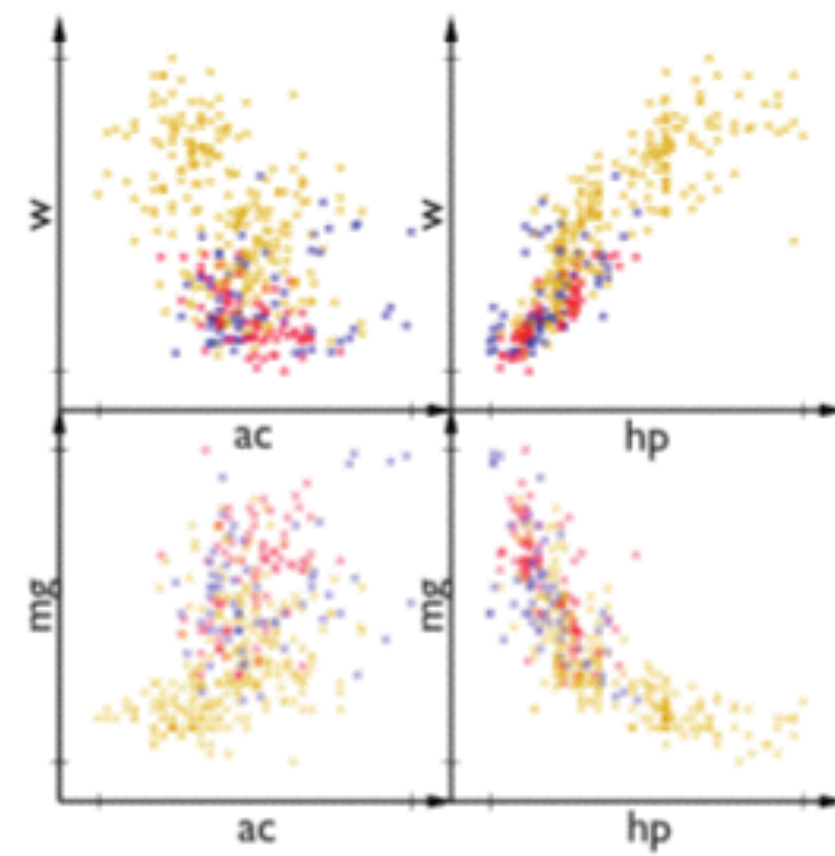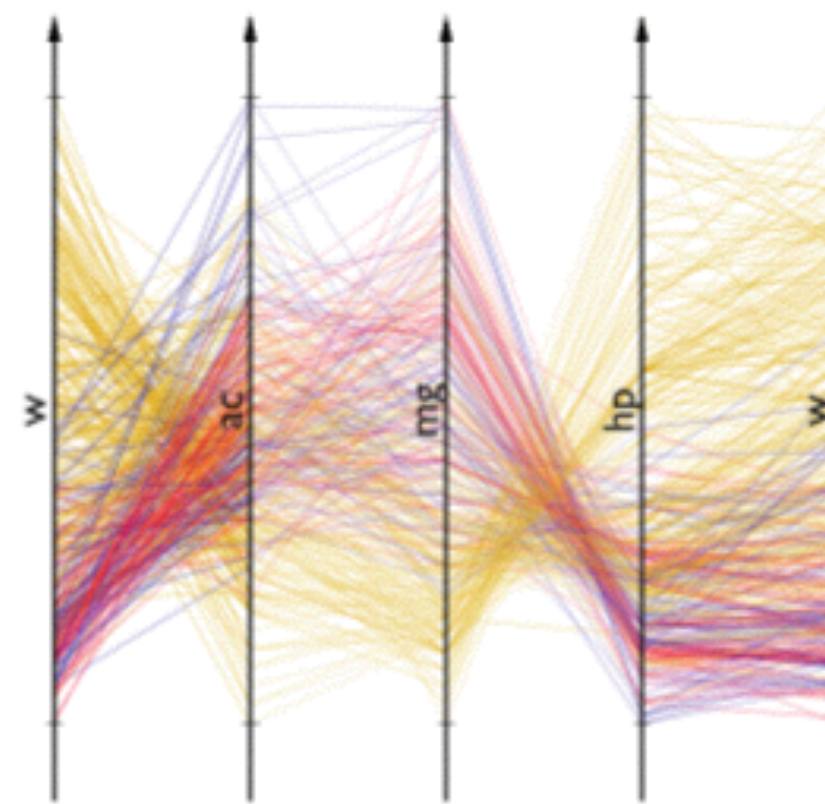| Math | Physics | Dance | Drama |
|------|---------|-------|-------|
| 85 | 95 | 70 | 65 |
| 90 | 80 | 60 | 50 |
| 65 | 50 | 90 | 90 |
| 50 | 40 | 95 | 80 |
| 40 | 60 | 80 | 90 |

# Scatterplot Matrix
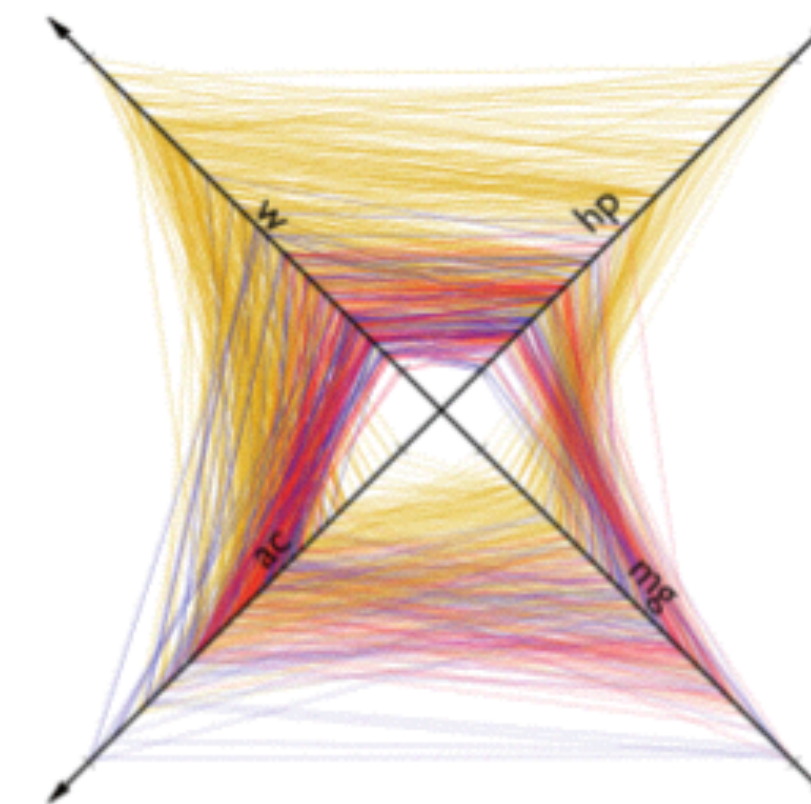
# Parallel Coordinates
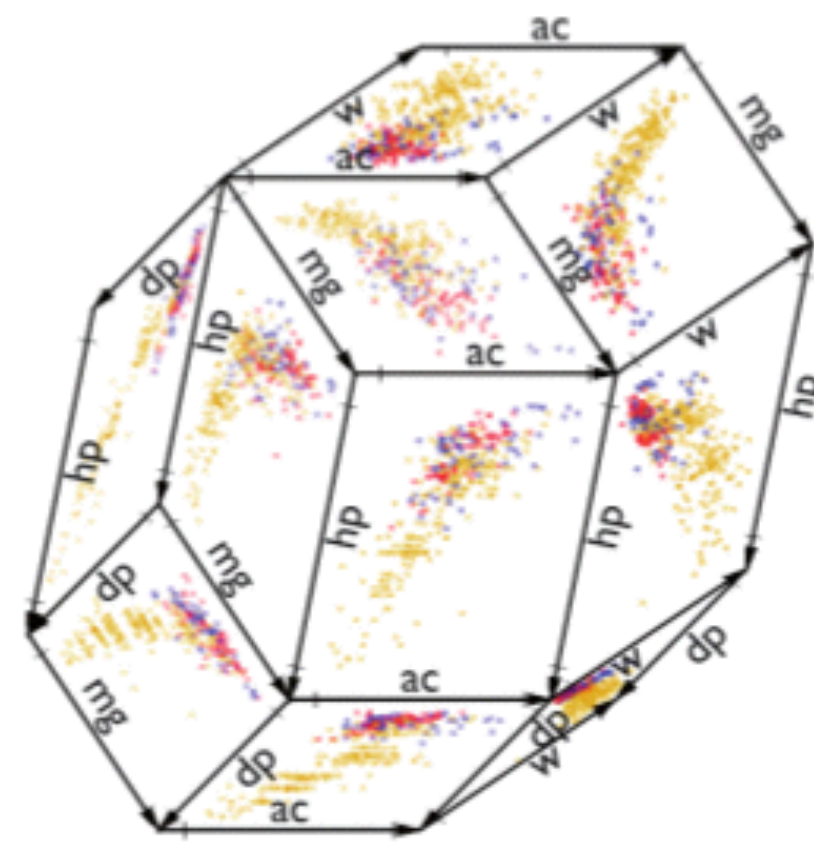
# Flexible Linked Axes (FLINA)
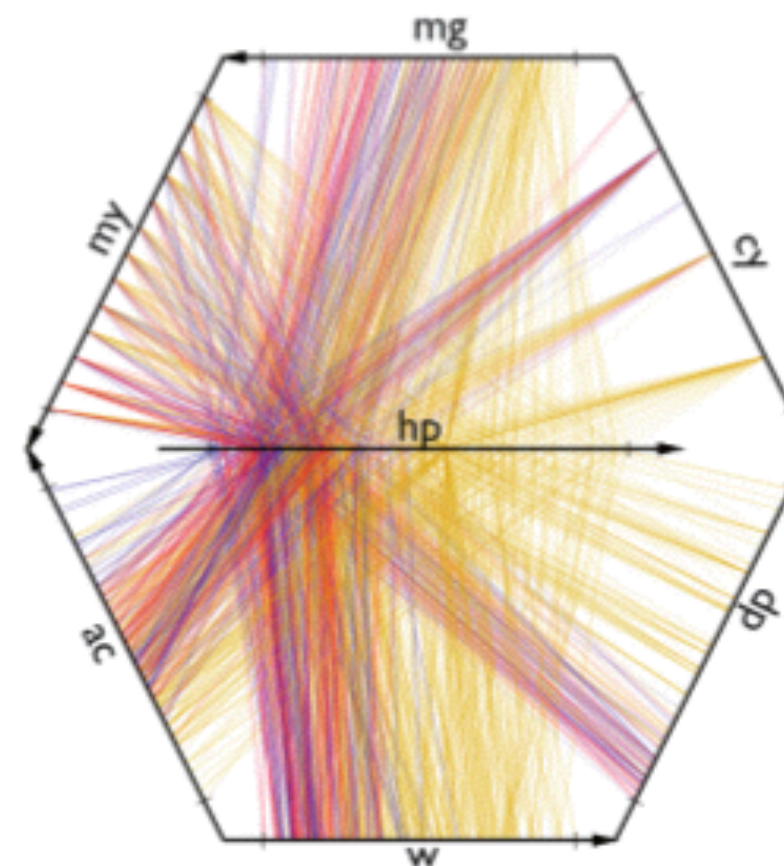


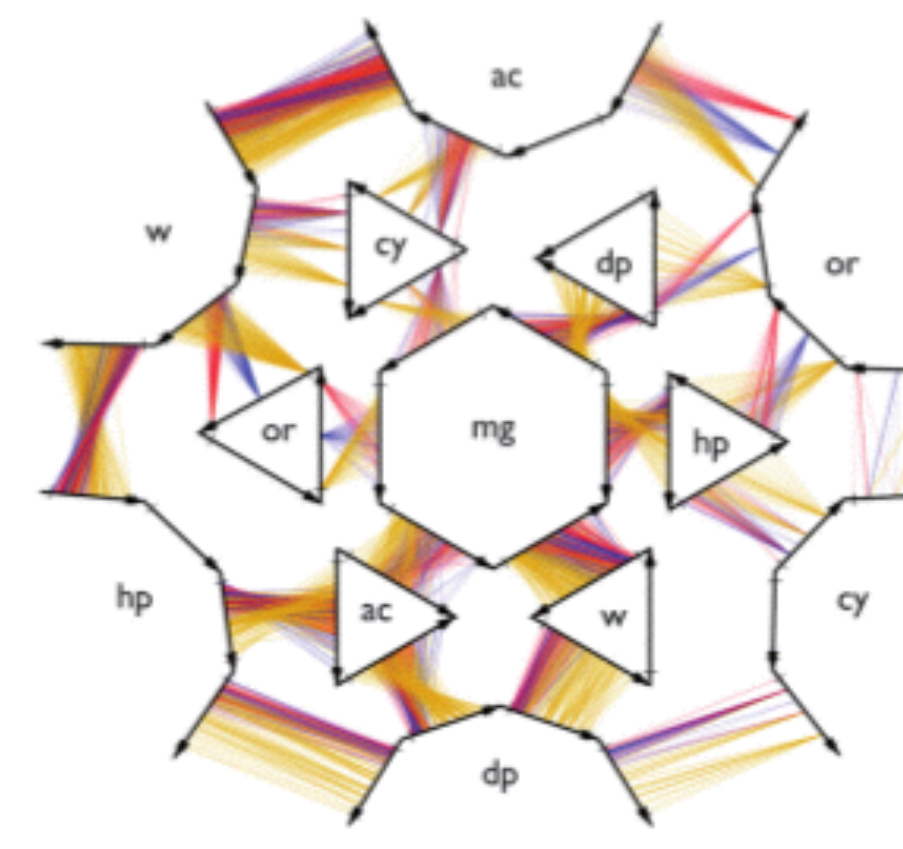(a) scatterplots
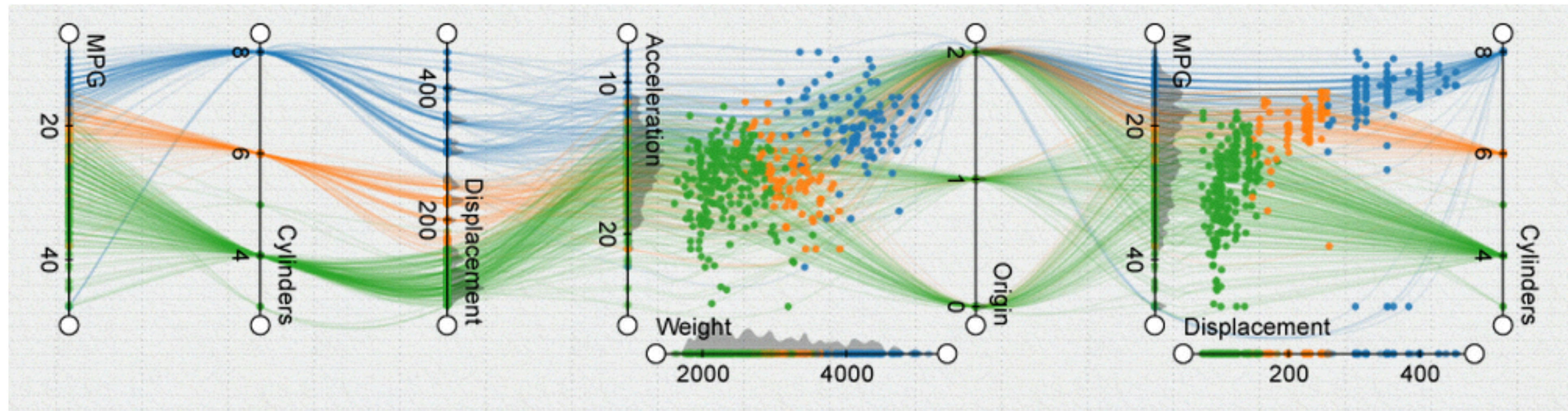(b) Parallel Coordinates Plot
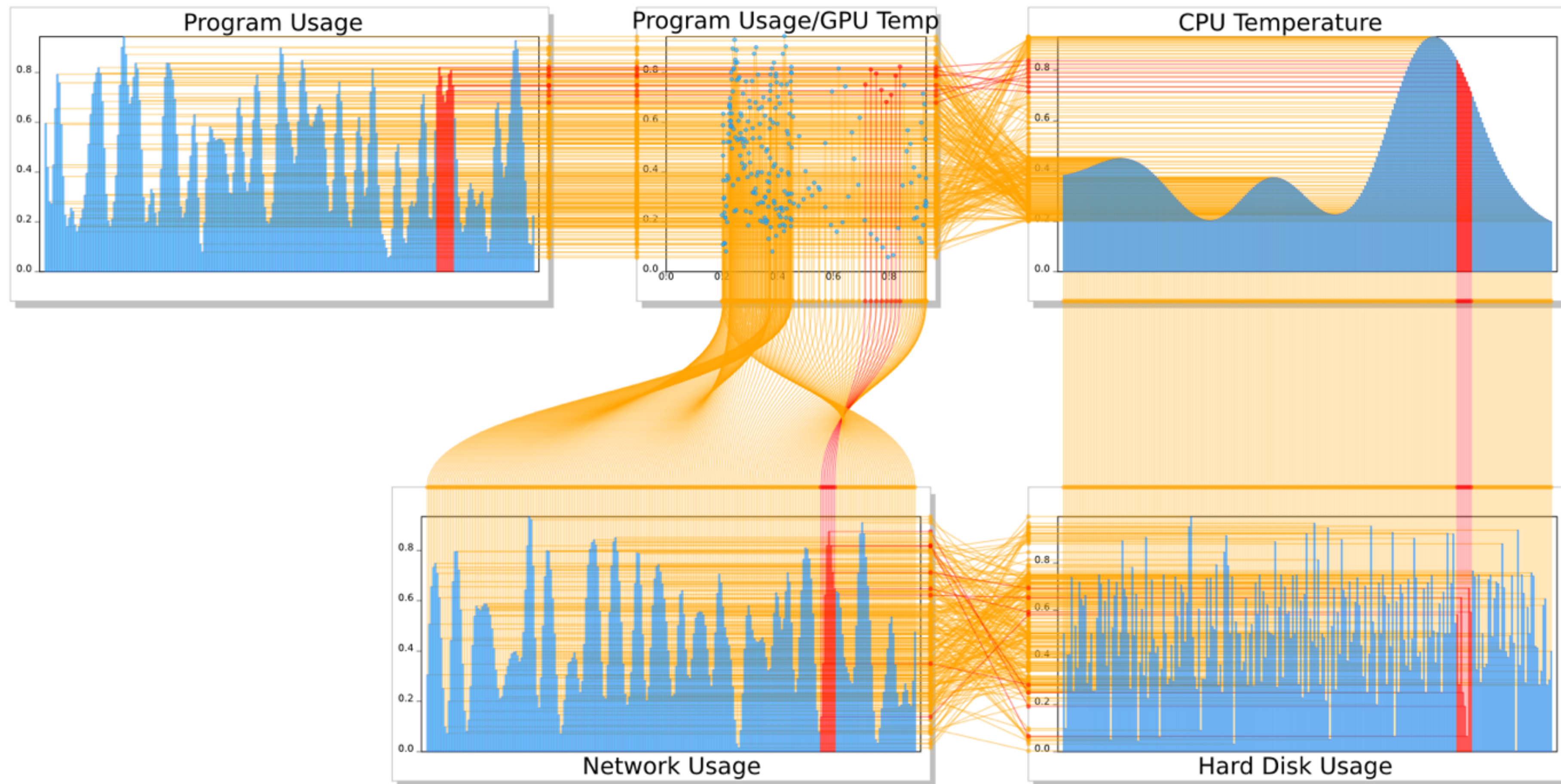(c) radar chart
(d) Hyperbox
(e) Time Wheel
(f) Many-to-many PCP

Claessen & van Wijk 2011

# Web-based implementation of FLINA concept



http://vis.pku.edu.cn/mddv/val/
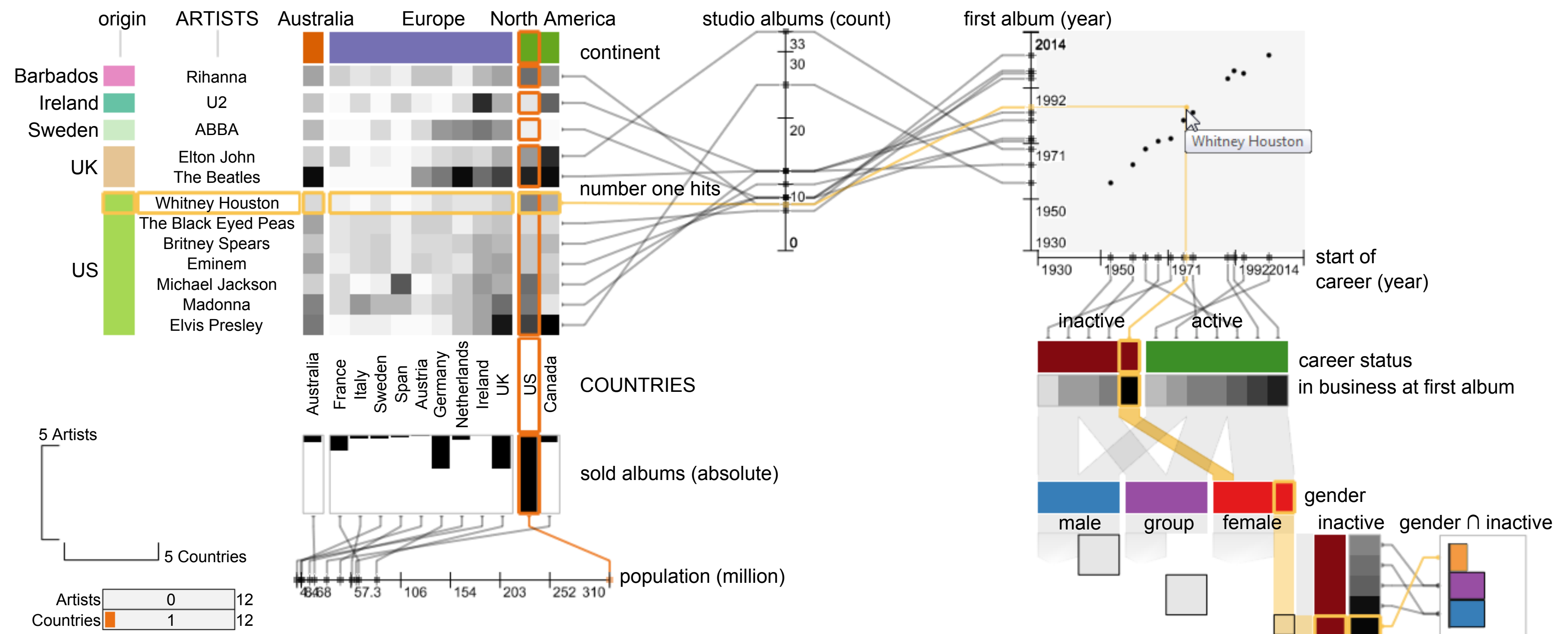
# Connected Charts

# Domino



Gratzl et al. 2014

# Data Reduction

## Sampling

**Don't show every element, show a (random) subset**

**Efficient for large dataset**

**Apply only for display purposes**

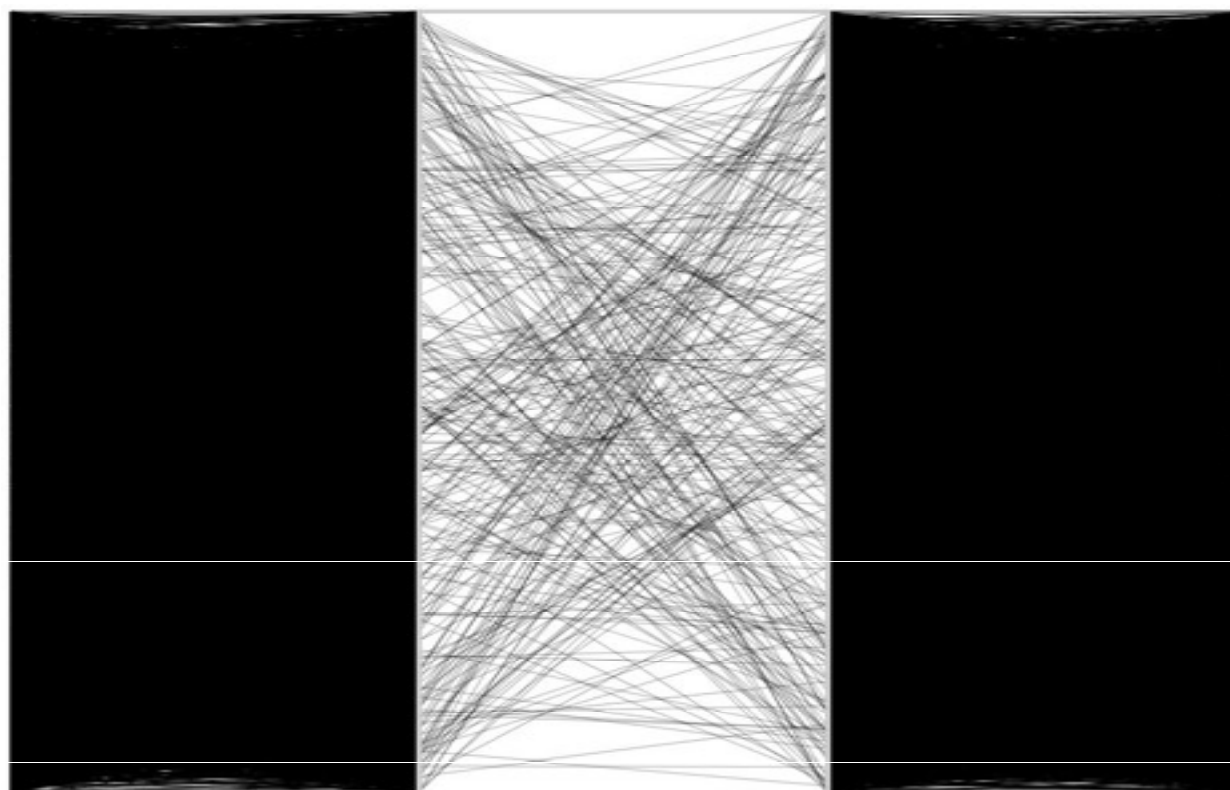**Outlier-preserving approaches**



[Ellis & Dix, 2006]

## Filtering

**Define criteria to remove data, e.g.,**

minimum variability

> / < / = specific value for one dimension

consistency in replicates, …

**Can be interactive, combined with sampling**

# Filter Example



Time of Day

Arrival Delay (min.)    reset

Distance (mi.)

Date

March 31, 2001

37,336 of 231,083 flights selected.

| 10:57 PM | MSY | HOU | 303 mi. | +29 min. |
| 09:33 PM | BWI | PVD | 328 mi. | +64 min. |
| 09:30 PM | TPA | PBI | 174 mi. | +30 min. |
| 09:29 PM | BWI | PVD | 328 mi. | +20 min. |
| 09:10 PM | LAS | PHX | 256 mi. | +53 min. |
| 09:02 PM | MSY | HOU | 303 mi. | +26 min. |

http://square.github.io/crossfilter/

# Pixel Based Methods

# Pixel Based Displays

Each cell is a "pixel", value encoded in color / value
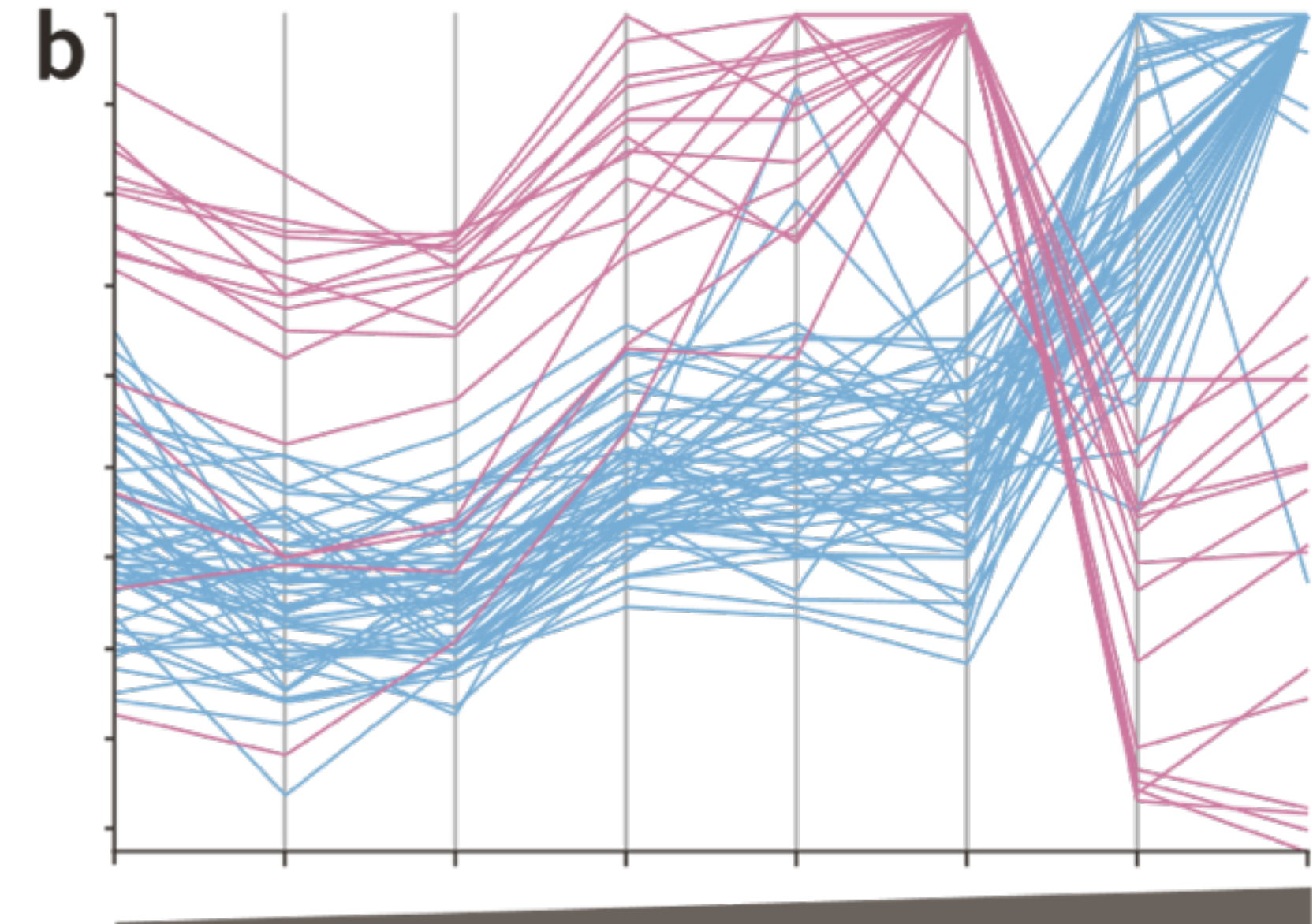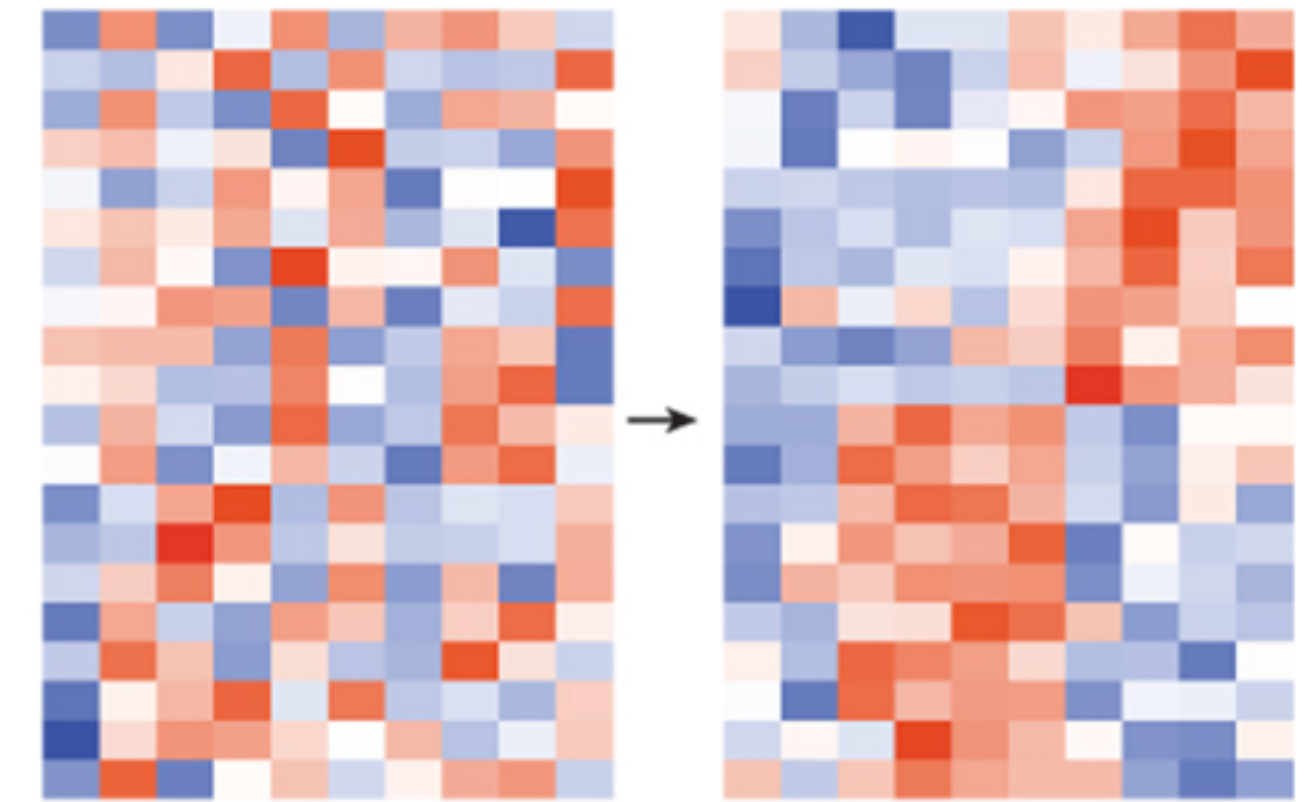
Meaning derived from ordering

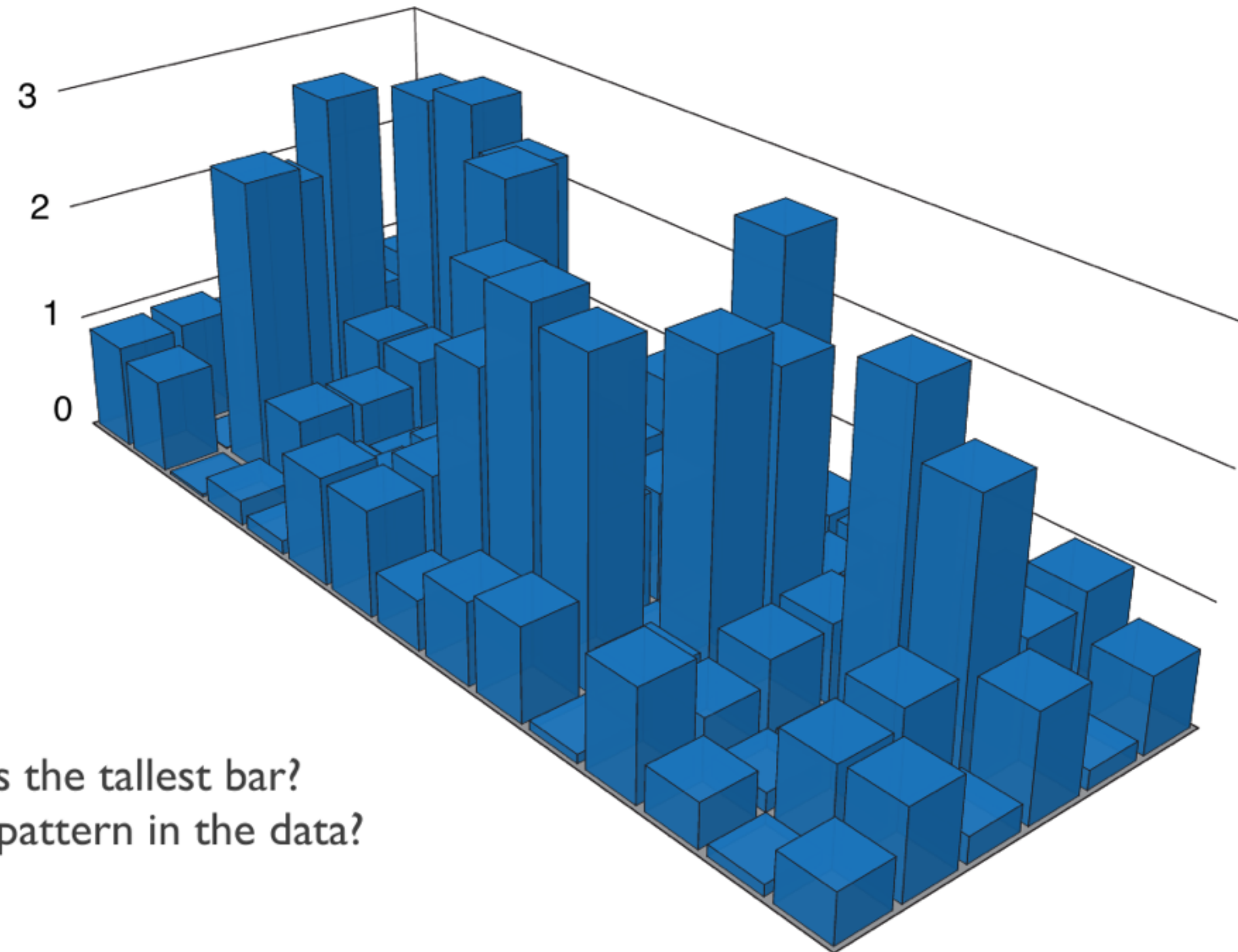If no ordering inherent, clustering is used

Scalable – 1 px per item
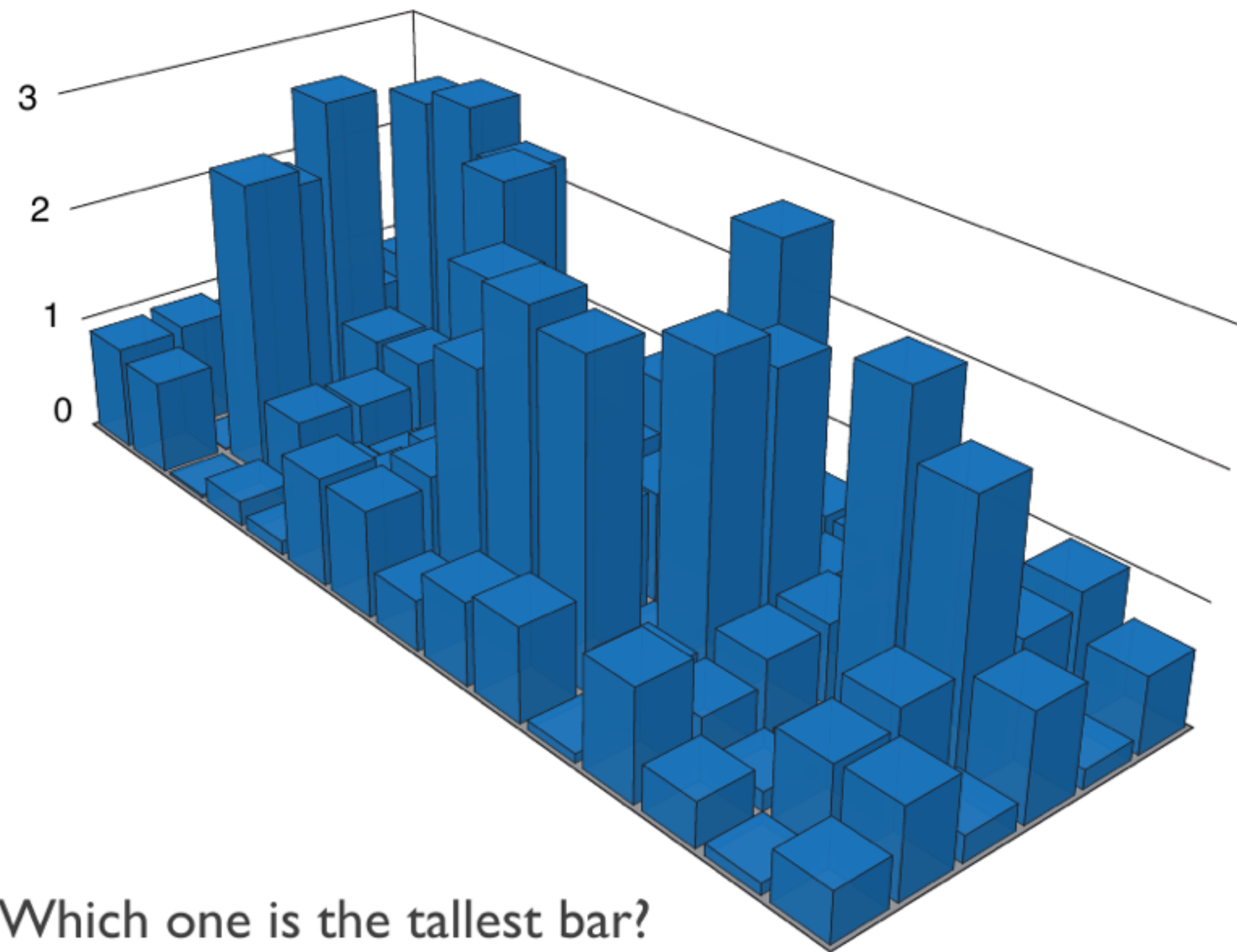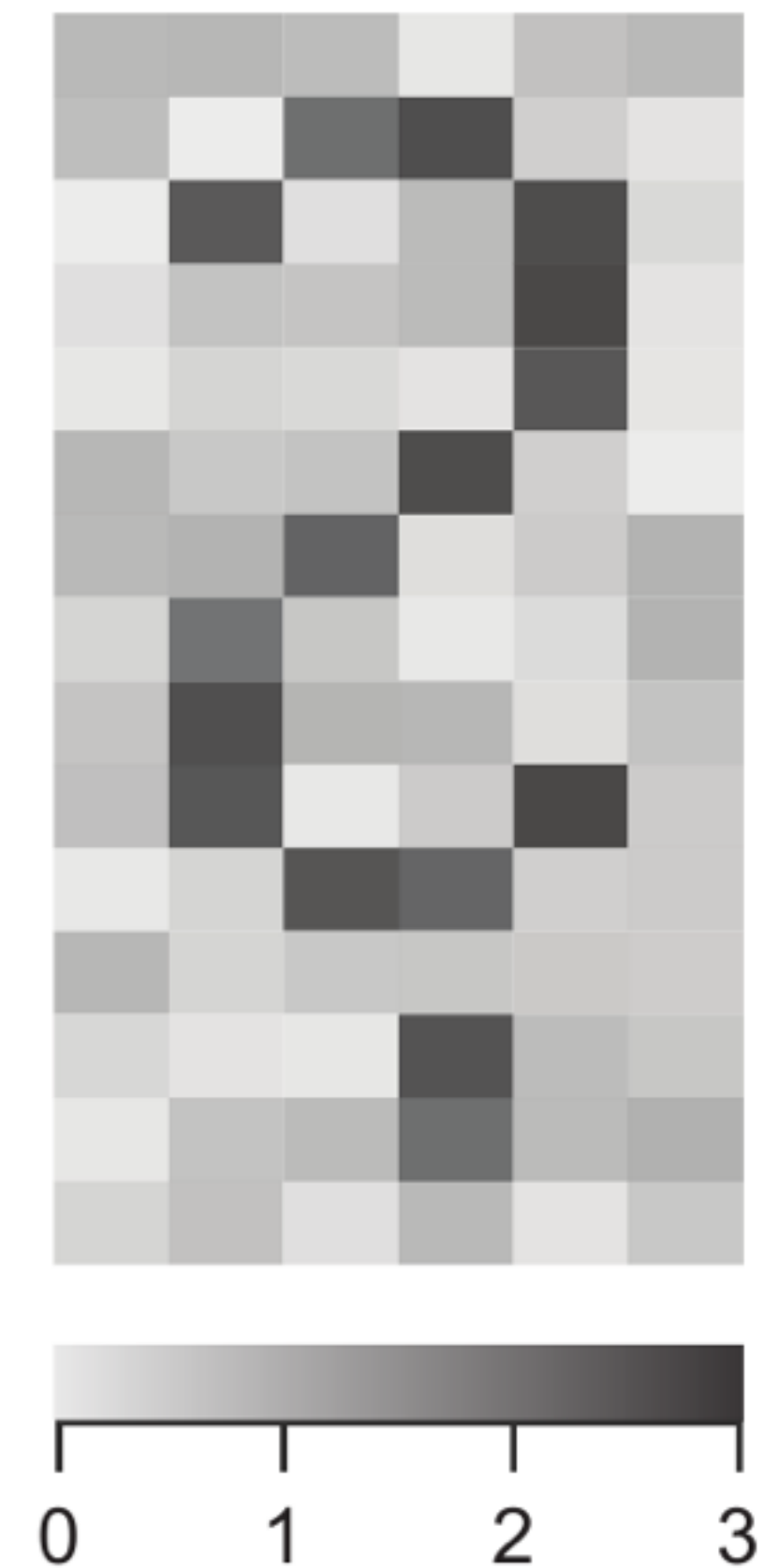
Good for homogeneous data

same scale & type



[Gehlenborg & Wong 2012]

# 3D Pitfall: Occlusion & Perspective



Which one is the tallest bar?
What is the pattern in the data?

[Gehlenborg and Wong, Nature Methods, 2012]

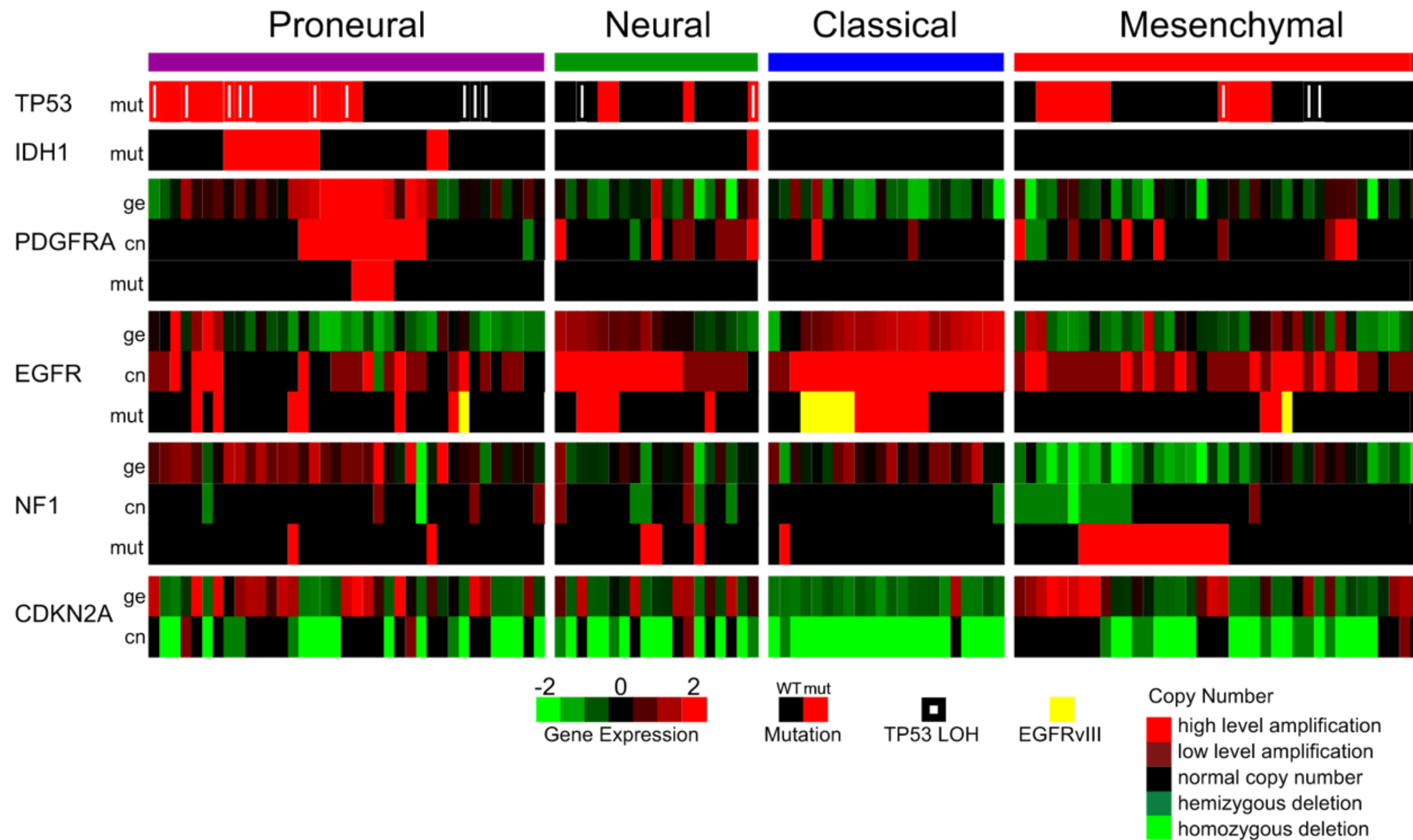# 3D Pitfall: Occlusion & Perspective



Which one is the tallest bar?
What is the pattern in the data?



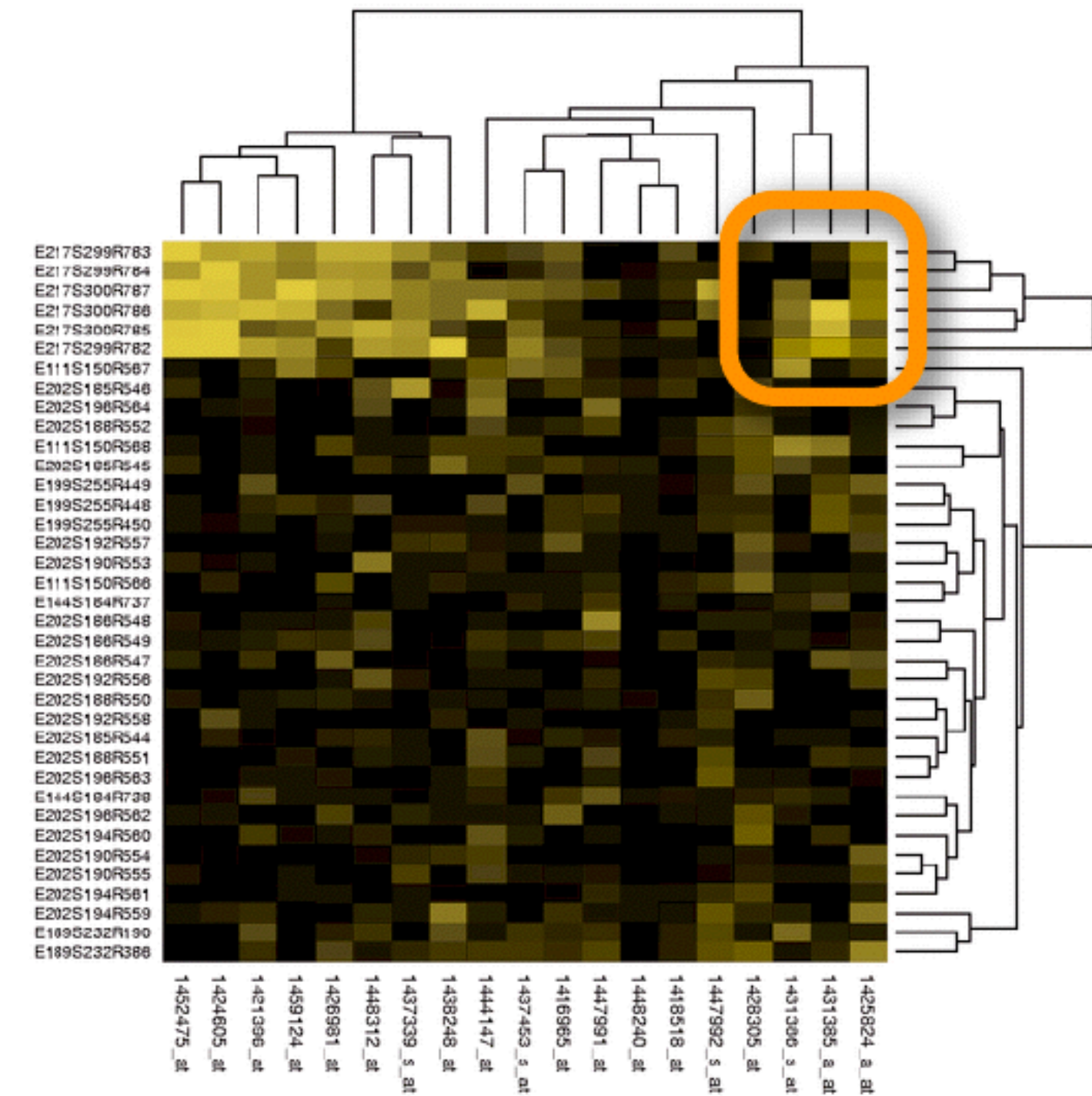[Gehlenborg and Wong, Nature Methods, 2012]

# Heterogeneous Data?


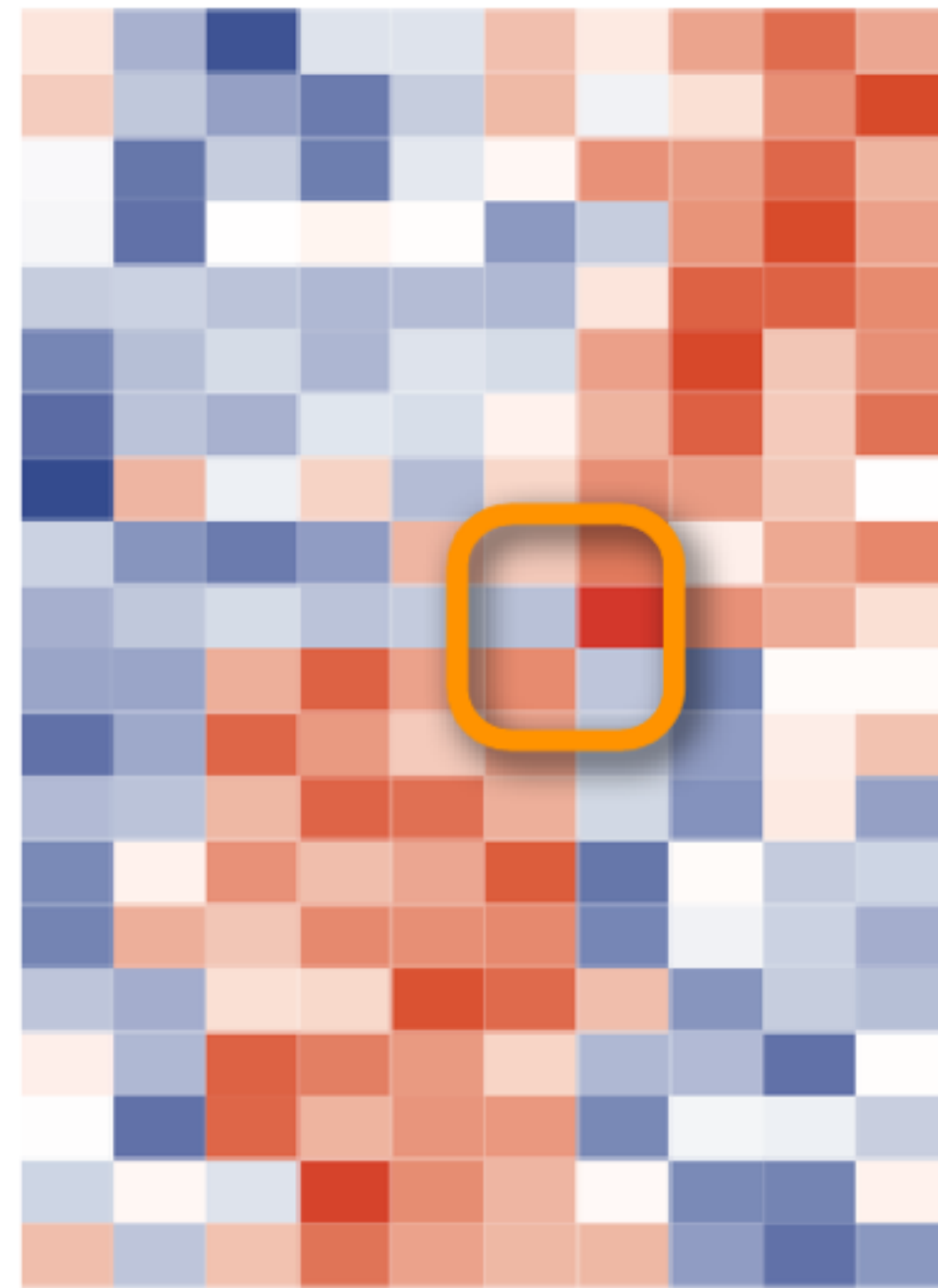
[Verhaak 2012]
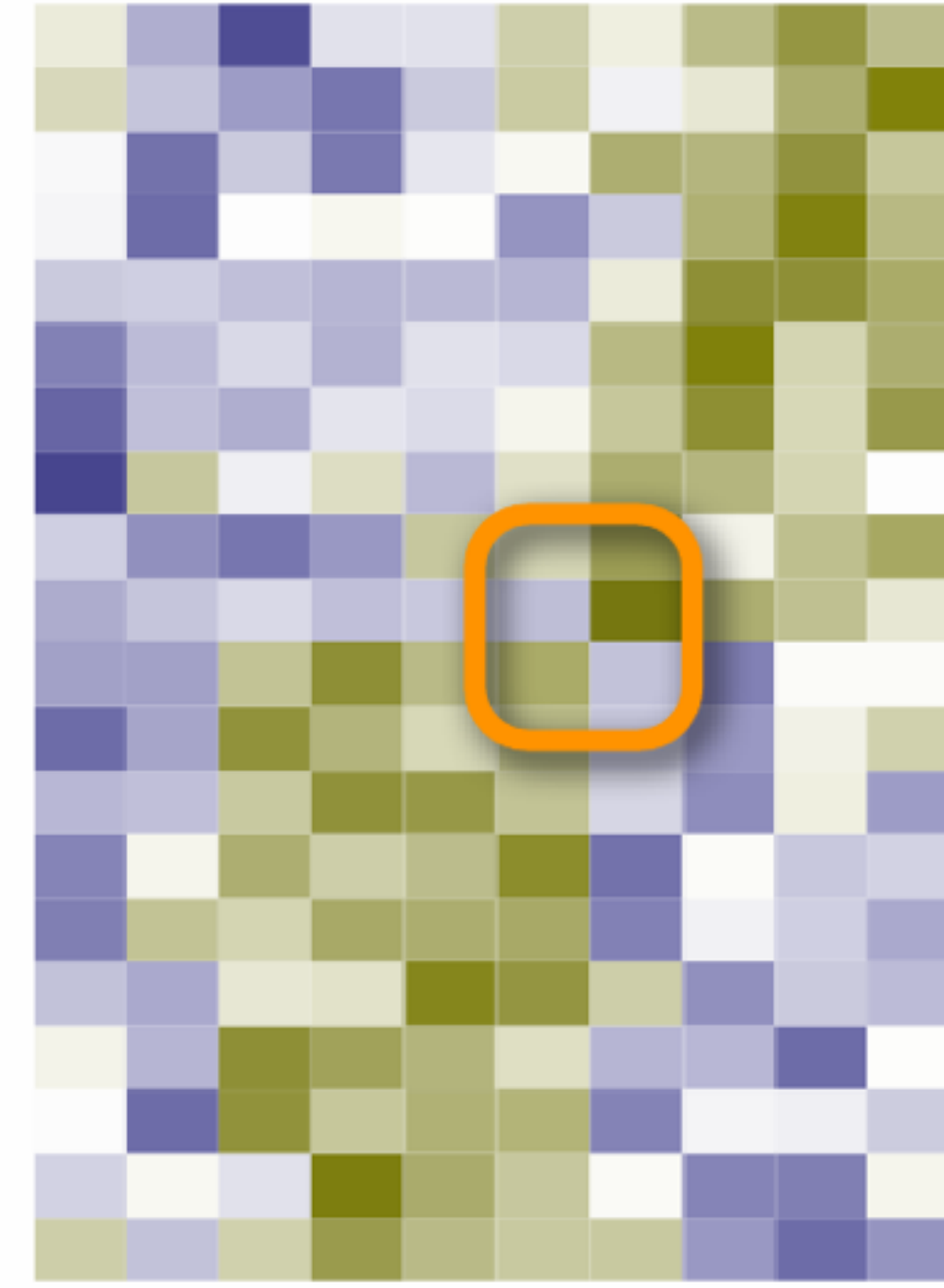
# Bad Color Mapping



Normal Vision

Deuteranope Vision
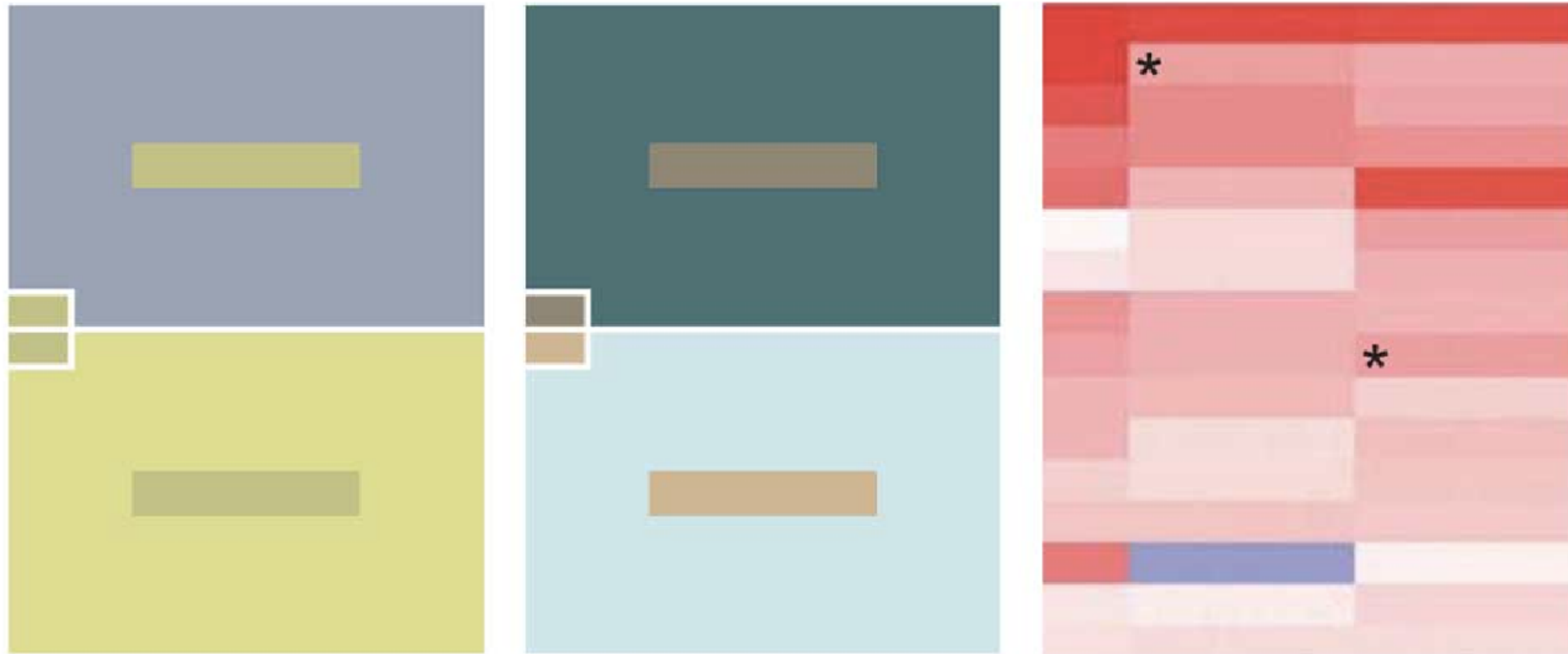("Red-Green Blindness")

# Good Color Mapping



Normal Vision

Deuteranope Vision
("Red-Green Blindness")

# Color is relative!

# Clustering

Classification of items into "similar" bins

Based on similarity measures

Euclidean distance, Pearson correlation, ...

Partitional Algorithms

divide data into set of bins

# bins either manually set (e.g., k-means) or automatically determined (e.g., affinity propagation)

Hierarchical Algorithms

Produce "similarity tree" – dendrogram

Bi-Clustering

Clusters dimensions & records

Fuzzy clustering

allows occurrence of elements in multiples clusters

# Clustering Applications

Clusters can be used to

   order (pixel based techniques)

   brush (geometric techniques)
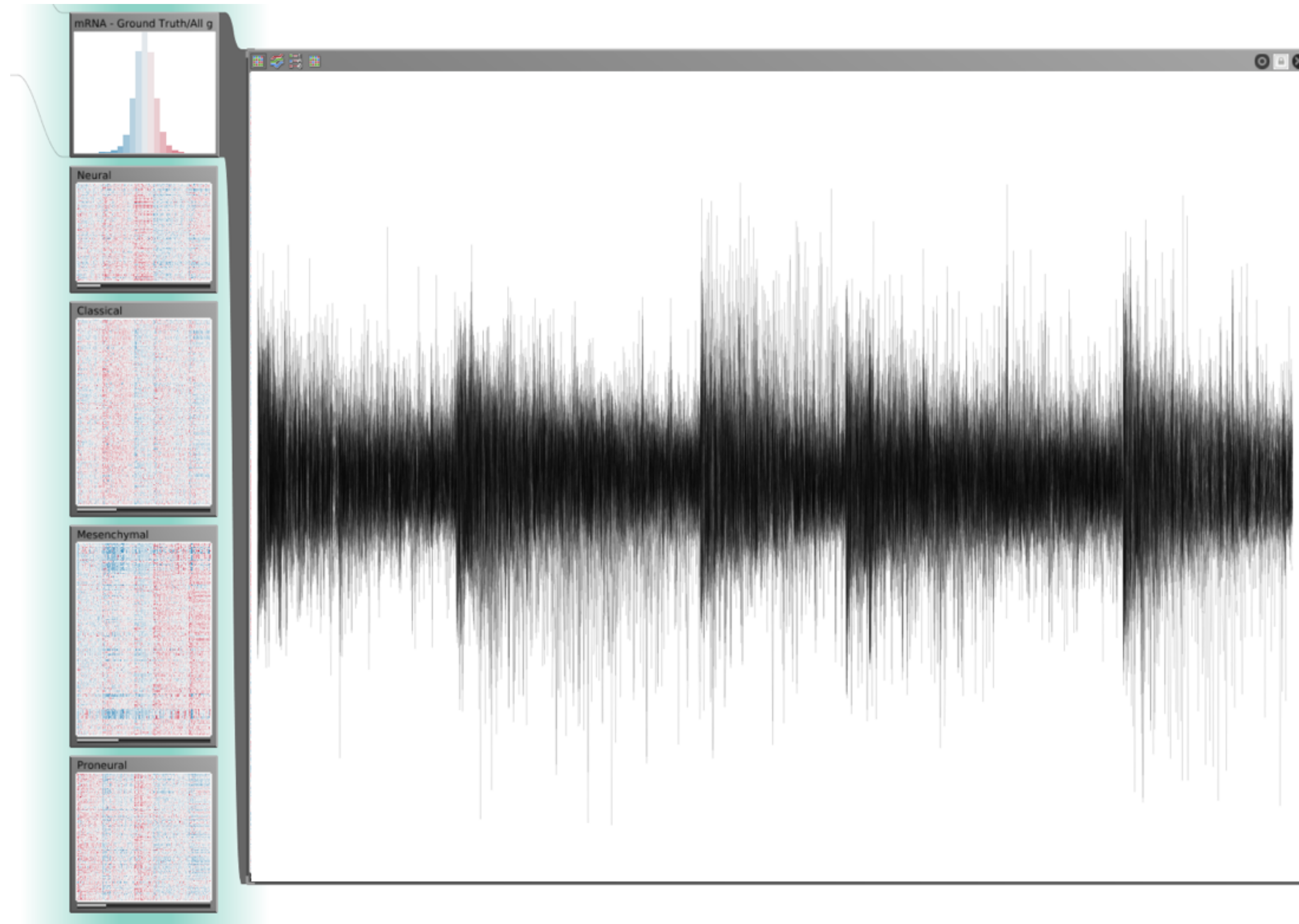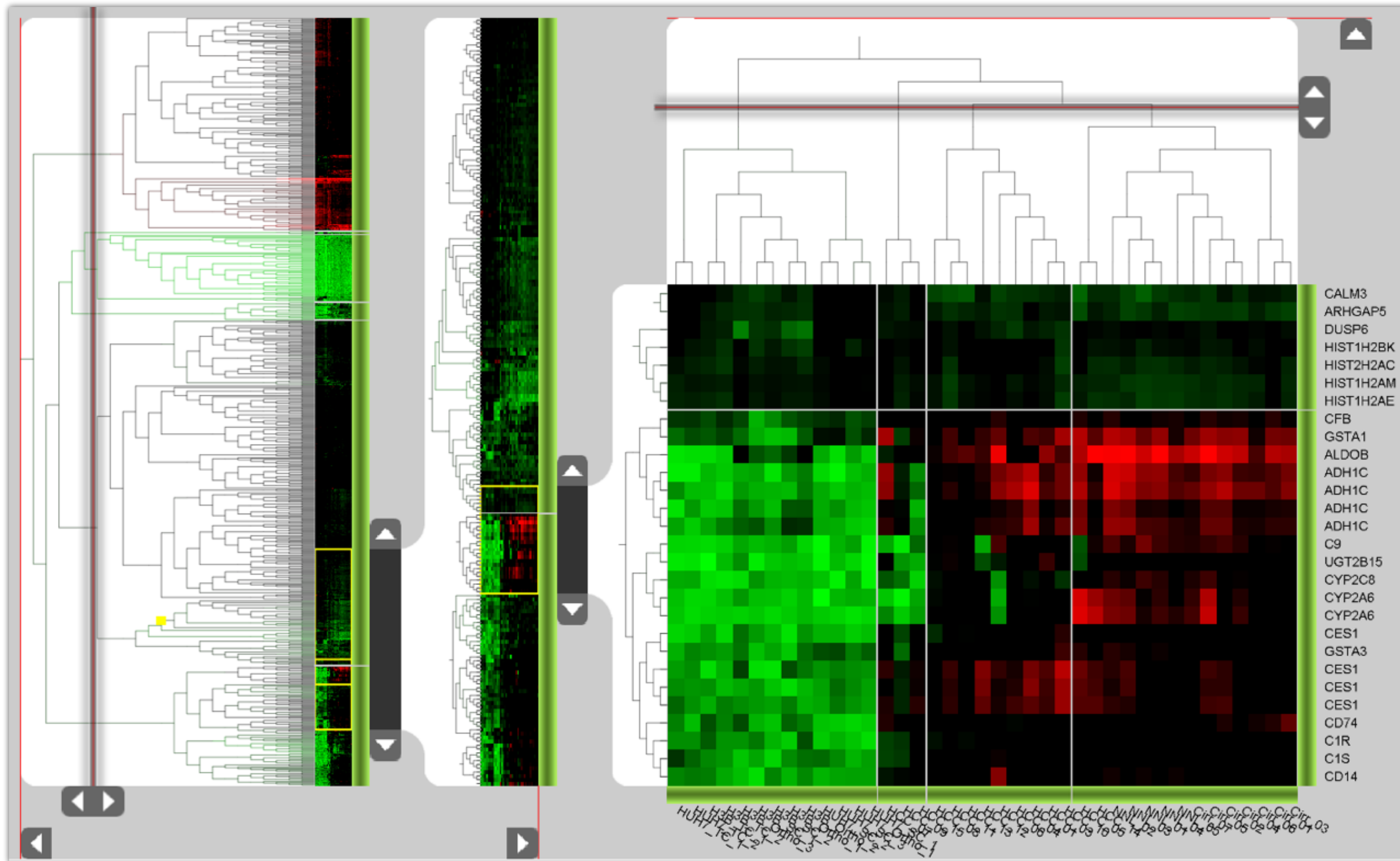
   aggregate

Aggregation

   cluster more homogeneous than whole dataset

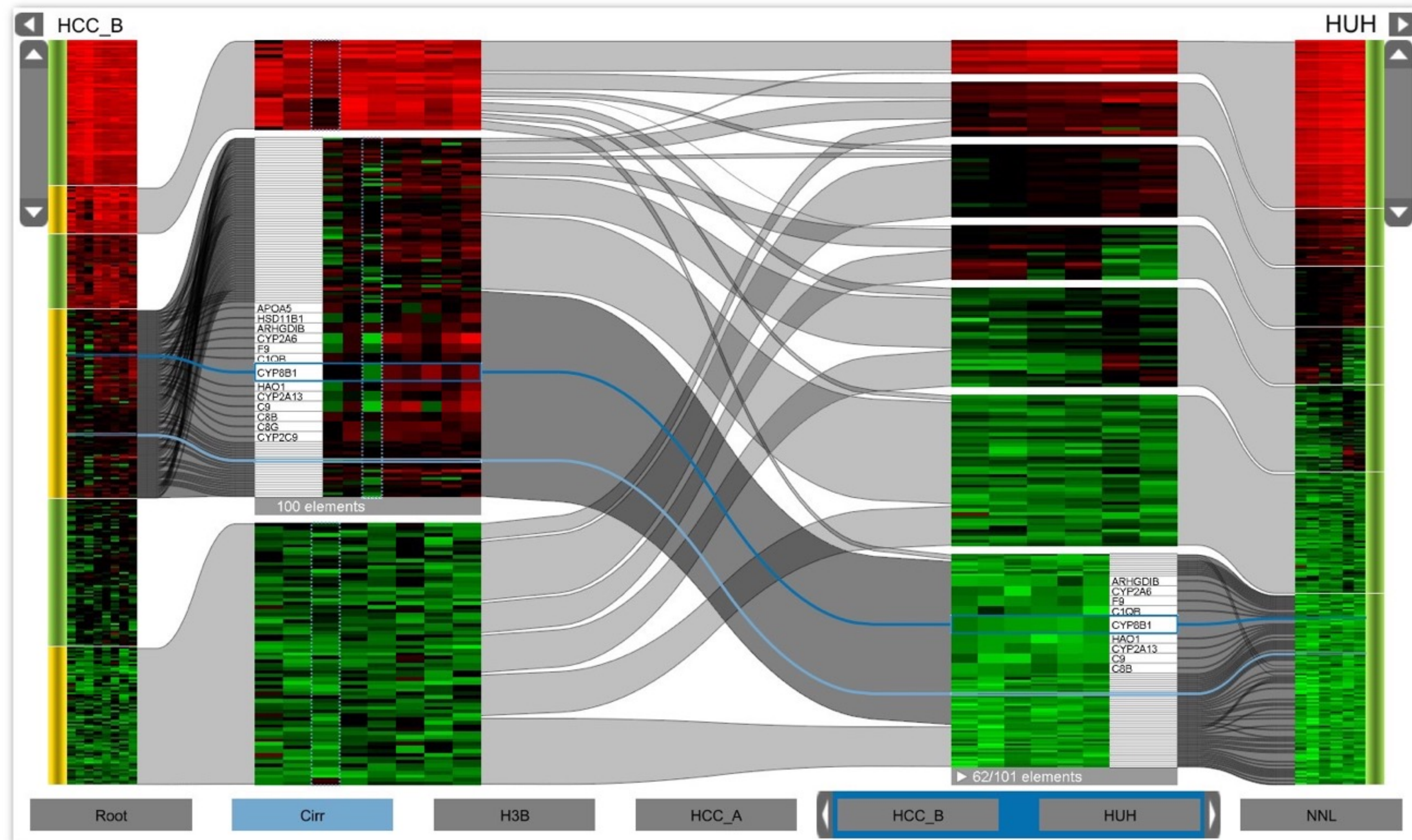   statistical measures, distributions, etc. more meaningful

# Clustered Heat Map
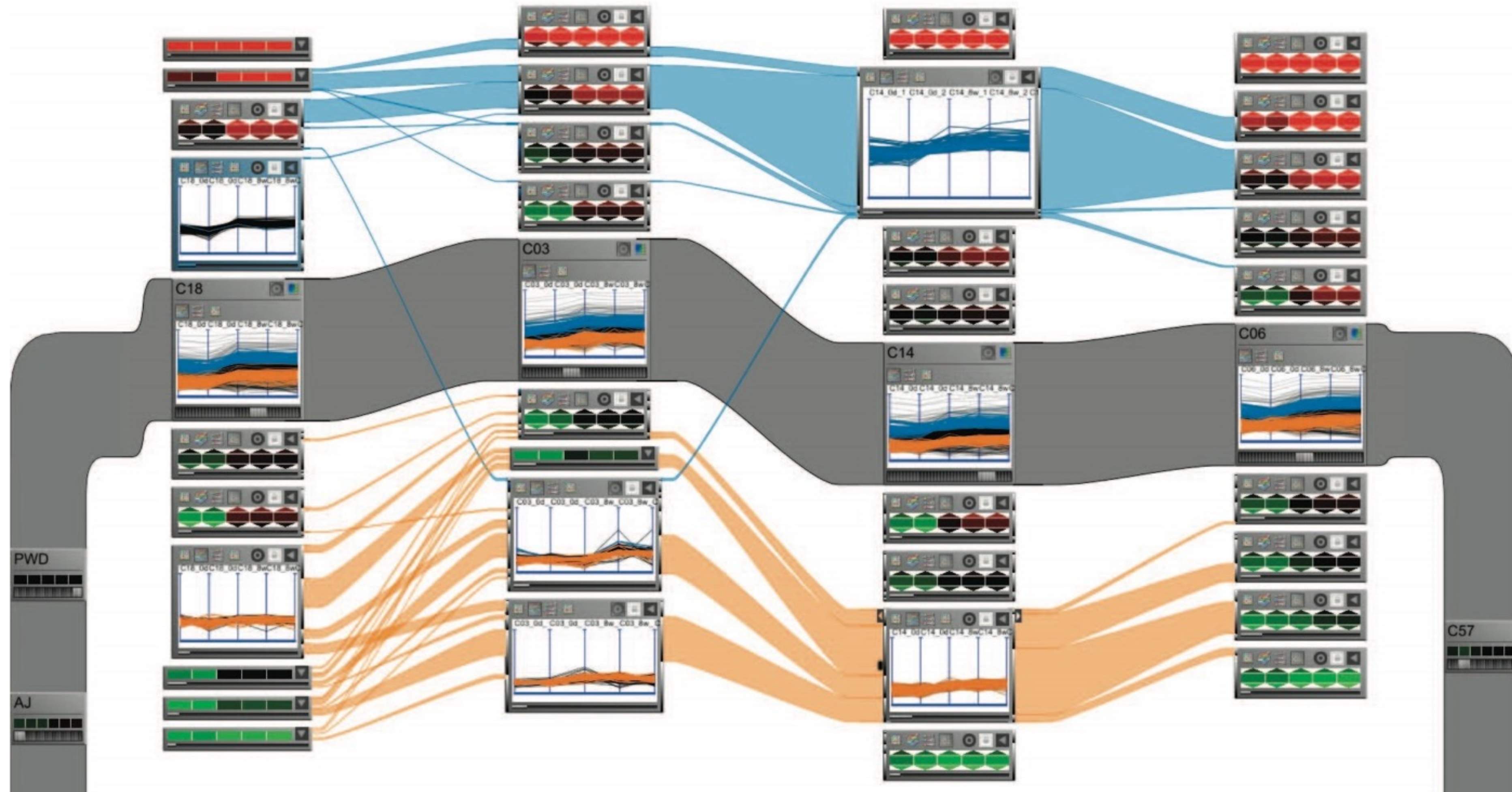
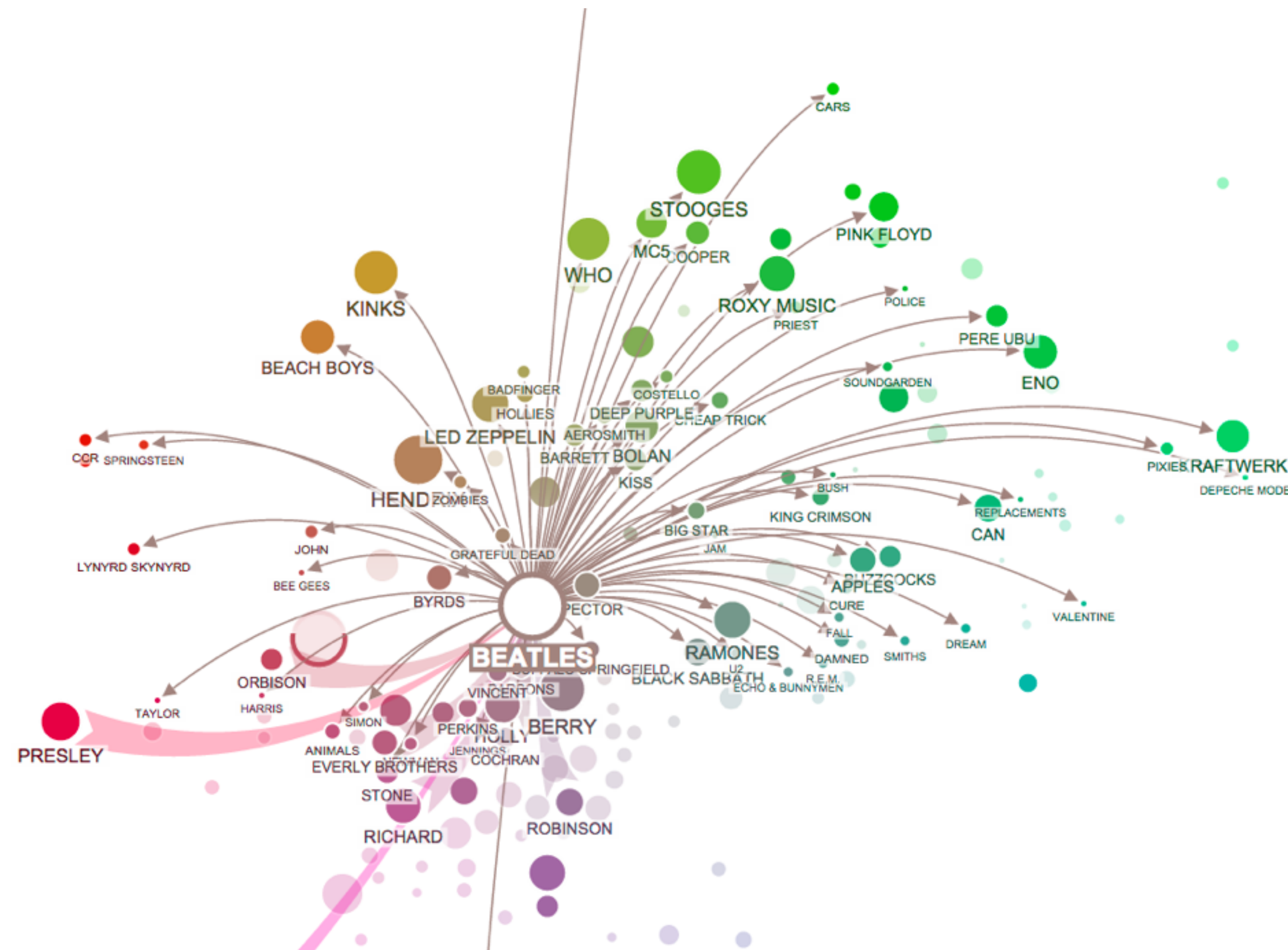# F+C Approach, with Dendrograms



CALM3
ARHGAP5
DUSP6
HIST1H2BK
HIST2H2AC
HIST1H2AM
HIST1H2AE
CFB
GSTA1
ALDOB
ADH1C
ADH1C
ADH1C
ADH1C
C9
UGT2B15
CYP2C8
CYP2A6
CYP2A6
CES1
GSTA3
CES1
CES1
CES1
CD74
C1R
C1S
CD14

# Cluster Comparison

# Aggregation

# Design Critique

# EdgeMaps: http://goo.gl/q8Cv7t



http://mariandoerk.de/edgemaps/demo/#music

# Dimensionality Reduction
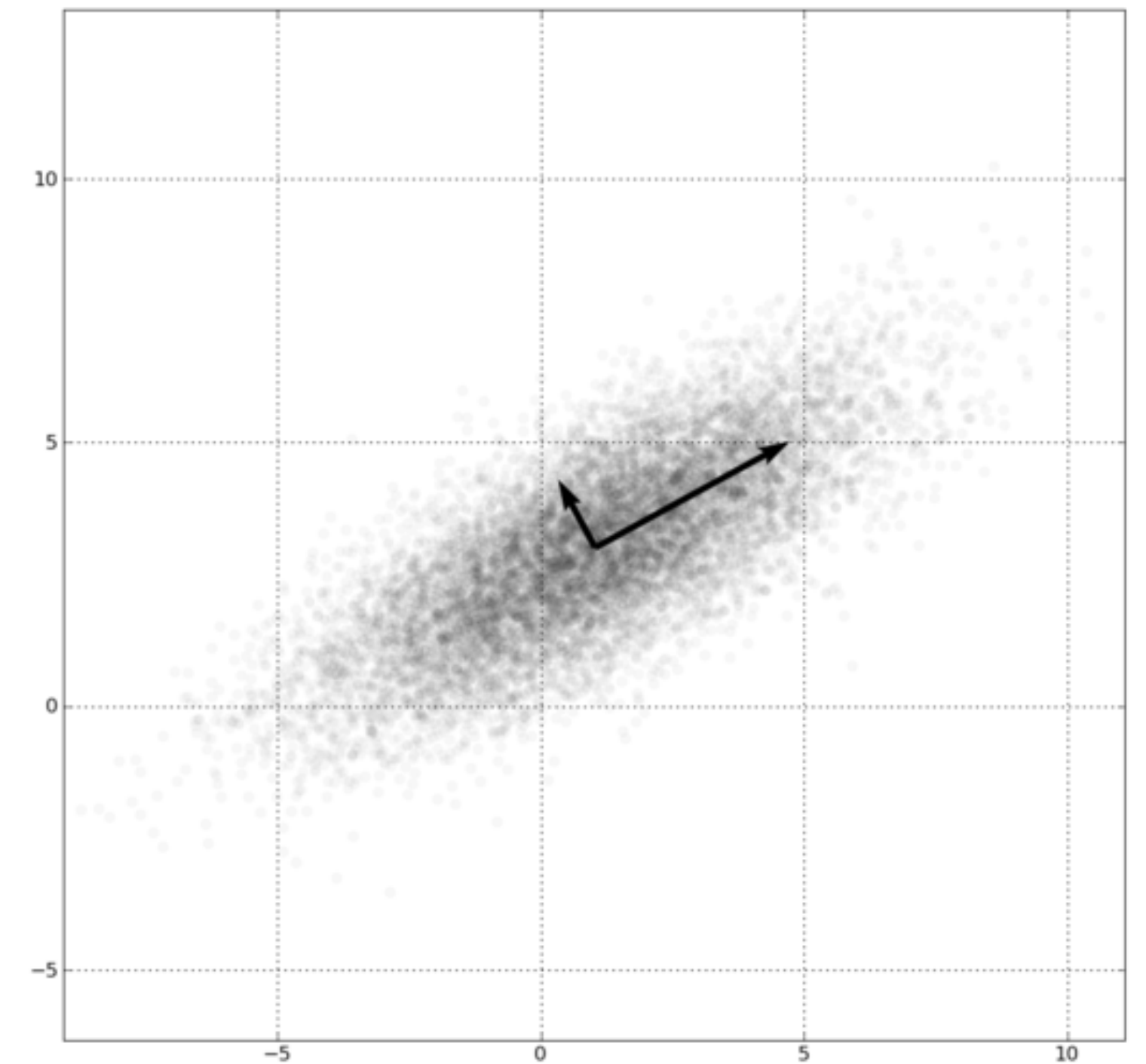
# Dimensionality Reduction

Reduce high dimensional to lower dimensional space
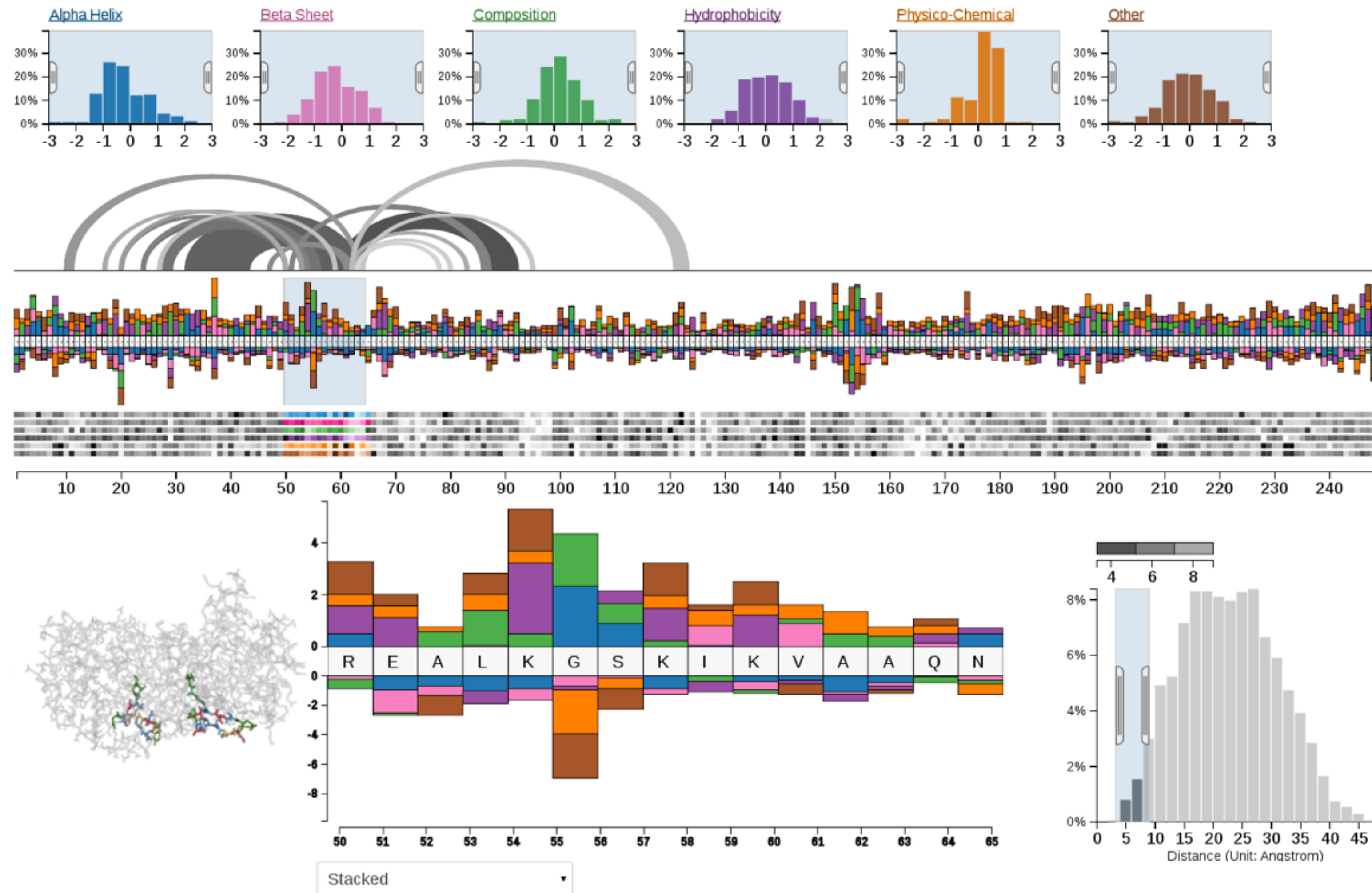
Preserve as much of variation as possible

Plot lower dimensional space

*Principal Component Analysis (PCA)*

linear mapping, by order of variance

# PCA Example – CS 171 Project 2013

[Mercer & Pandian]

# Multidimensional Scaling
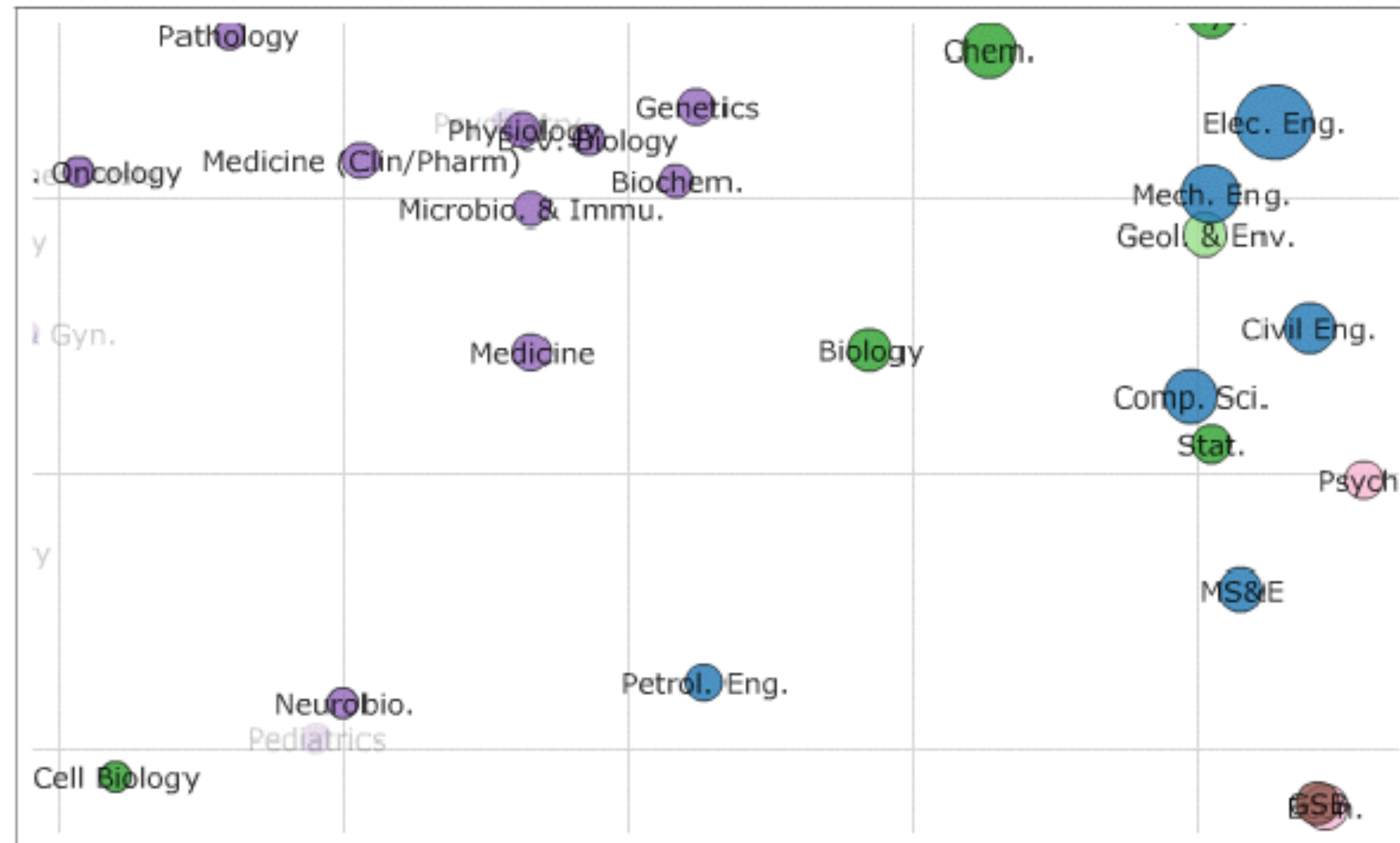
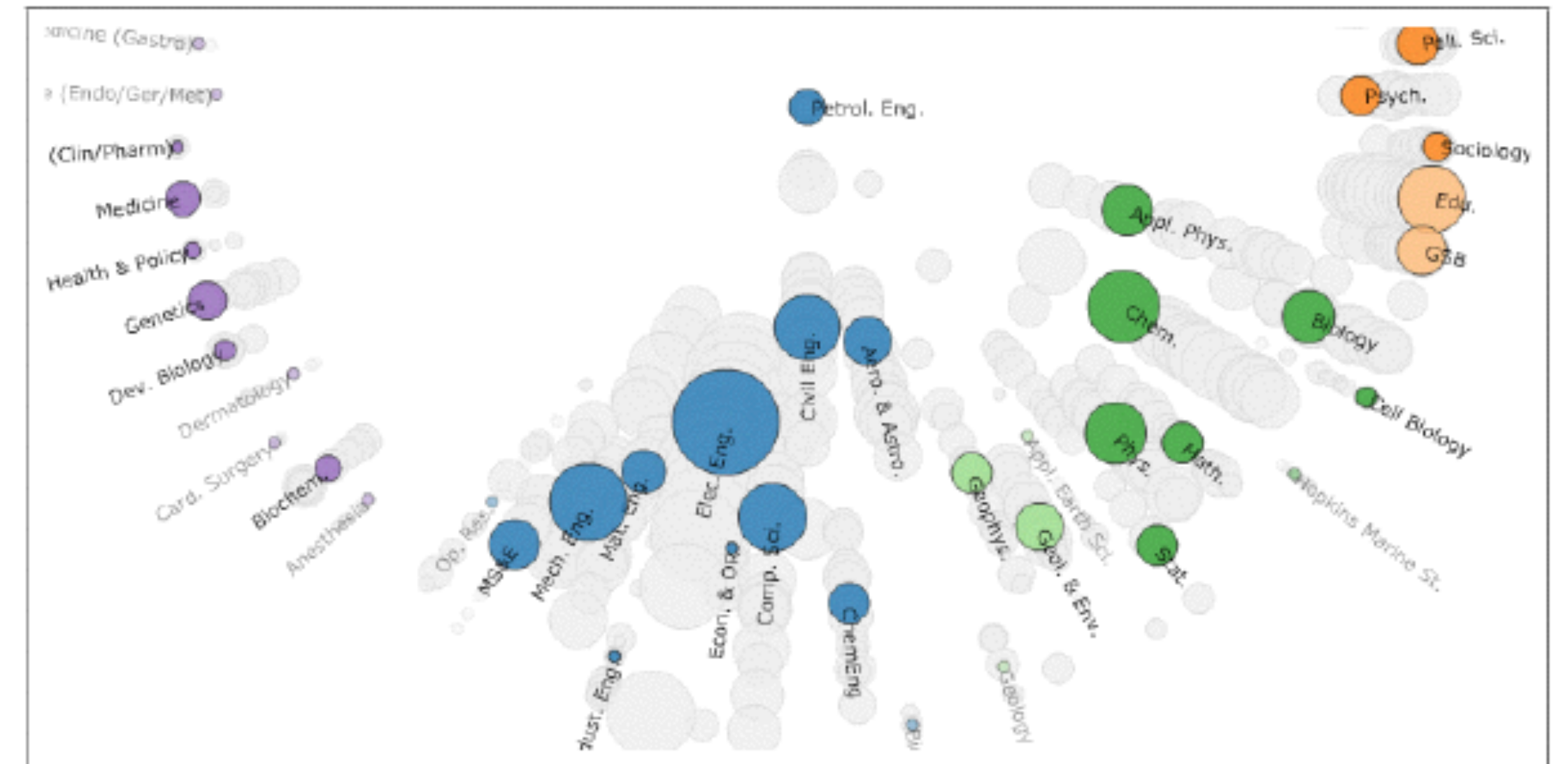Nonlinear, better suited for some DS

Popular for text analysis

[Doerk 2011]

# Can we Trust Dimensionality Reduction?

**Topical distances between departments in a 2D projection**

**Topical distances between the selected Petroleum Engineering and the others.**



[Chuang et al., 2012]

http://www-nlp.stanford.edu/projects/dissertations/browser.html