

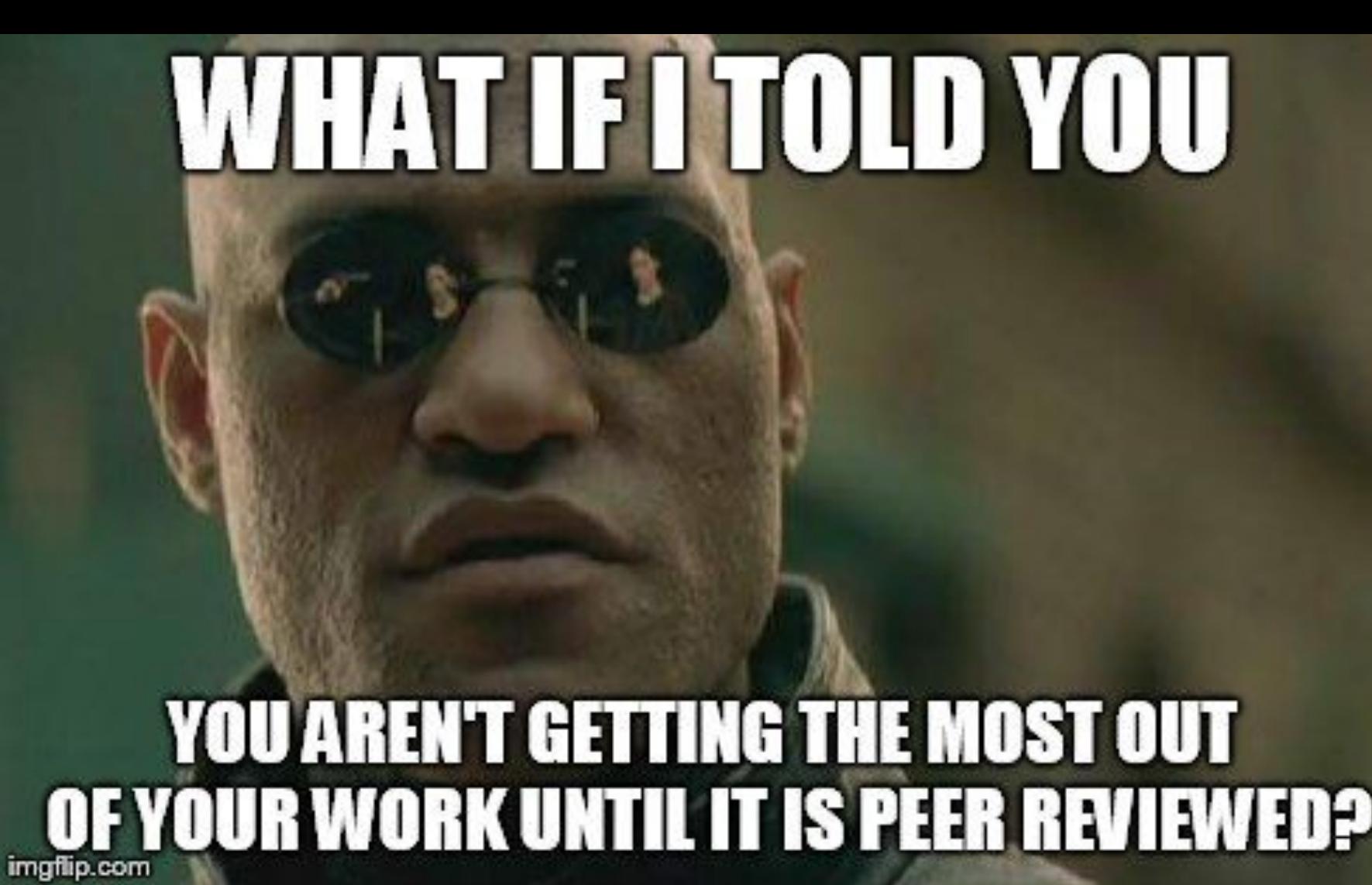
# TABULAR DATA

KIRELL BENZI, PH.D.



Inspired by Lex

[www.kirellbenzi.com](http://www.kirellbenzi.com)



WHAT IF I TOLD YOU

YOU AREN'T GETTING THE MOST OUT  
OF YOUR WORK UNTIL IT IS PEER REVIEWED?

imgflip.com



REVIEW COMPLETED

TOAST TO YOU AND YOUR  
THOUGHTS

memecrunch.com



"ALL I WANTED WAS

SOME REVIEWS...."

imgflip.com



LET'S REVIEW



I WILL FIND  
YOU  
AND I WILL THANK  
YOU

memegenerator.net



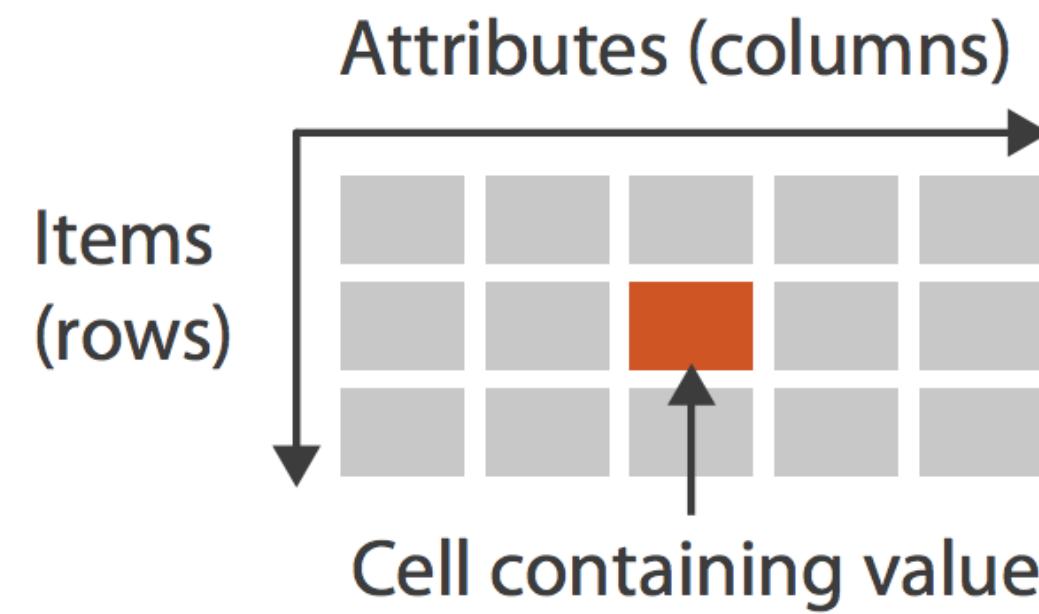
MID YEAR REVIEW...

WHO CARES?

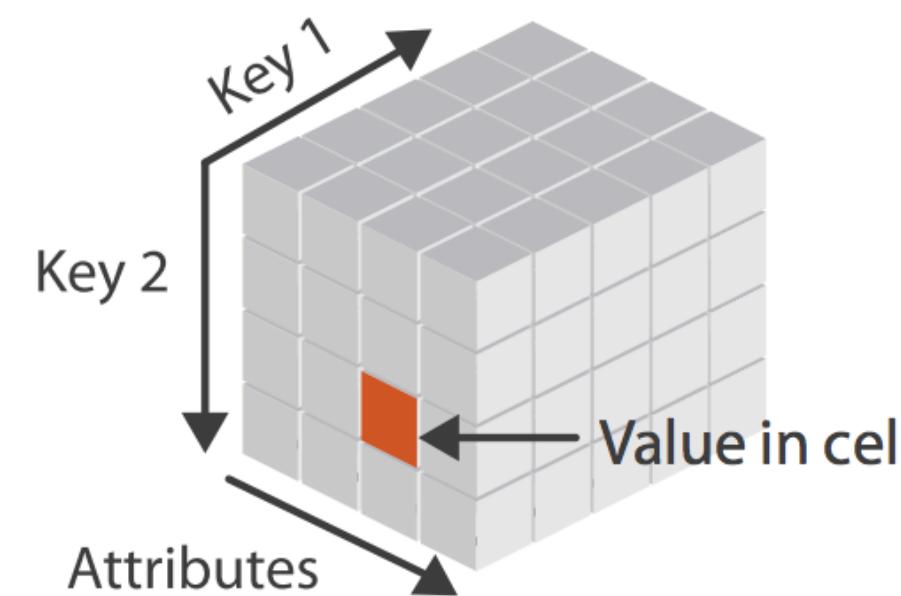
memegenerator.net

# Tabular data

## → Tables



→ *Multidimensional Table*



**19 columns, 500+ rows**

**23.03 MB, Comma-Separated, UTF-8**

album_id	album_comments	album_date_created	album_date_released	album_engineer	albumFavorites	album_handle
1	0	11/26/2008 01:44:45 AM	1/05/2009		4	AWOL_-_A_Way_Of_Life
100	0	11/26/2008 01:55:44 AM	1/09/2009		0	On_Opaque_Things
1000	0	12/04/2008 09:28:49 AM	10/26/2008		0	DMBQ_Live_at_2008_Record
10000	0	9/05/2011 04:42:57 PM			0	Live_at_CKUT_on_Montreal_S
10001	0	9/06/2011 12:02:58 AM	1/01/2006		0	Grounds_Dream_Cosmic_Love
10002	0	9/06/2011 12:53:01 AM			0	Trust_Now
10003	0	9/06/2011 05:46:22 AM	9/06/2011		1	Or_Up_We_Fall

# Arranging tables overview

## Arrange Tables

### → Express Values



### → Separate, Order, Align Regions

→ Separate



→ Order



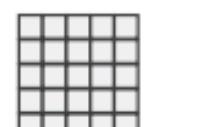
→ Align



→ 1 Key  
*List*



→ 2 Keys  
*Matrix*



→ 3 Keys  
*Volume*



→ Many Keys  
*Recursive Subdivision*



### → Axis Orientation

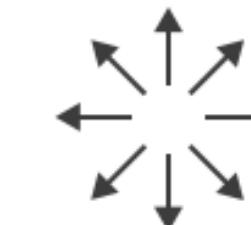
→ Rectilinear



→ Parallel

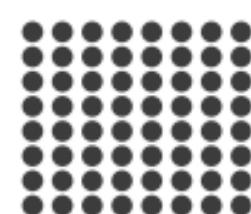


→ Radial



### → Layout Density

→ Dense



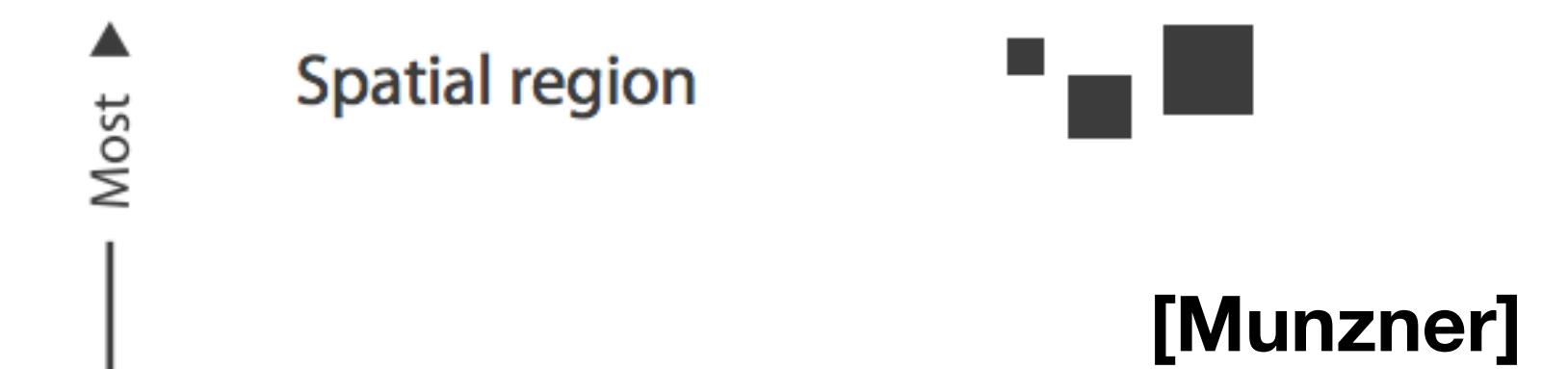
→ Space-Filling



## Channels: Expressiveness Types and Effectiveness Ranks

### → Magnitude Channels: Ordered Attributes

Position on common scale



Position on unaligned scale



[Munzner]

# Table viz scalability

Normal tables can be visualized “as is”

~ 50 attributes ~ 1000 records

High-dim tables need analytical methods

~ 1000 attributes >> 10'000 records

19 columns, 500+ rows  
23.03 MB, Comma-Separated, UTF-8

album_id	album_comments	album_date_created	album_date_released	album_engineer	albumFavorites	album_handle
1	0	11/26/2008 01:44:45 AM	1/05/2009		4	AWOL_-_A_Way_Of_Life
100	0	11/26/2008 01:55:44 AM	1/09/2009		0	On_Opaque_Things
1000	0	12/04/2008 09:28:49 AM	10/26/2008		0	DMBQ_Live_at_2008_Record_
10000	0	9/05/2011 04:42:57 PM			0	Live_at_CKUT_on_Montreal_S
10001	0	9/06/2011 12:02:58 AM	1/01/2006		0	Grounds_Dream_Cosmic_Love
10002	0	9/06/2011 12:53:01 AM			0	Trust_Now
10003	0	9/06/2011 05:46:22 AM	9/06/2011		1	Or_Up_We_Fall

## Table homogeneity?

# Scatterplots

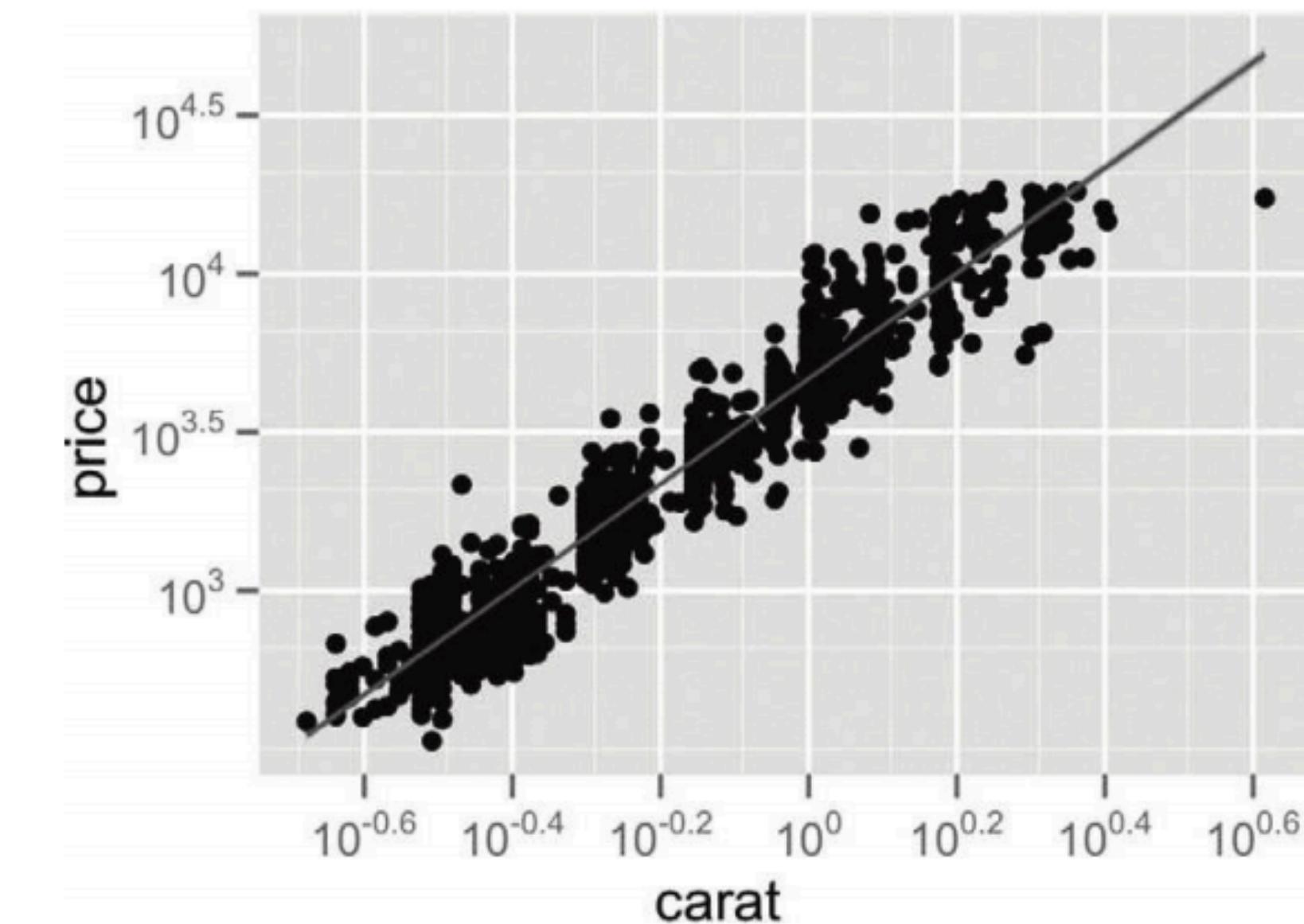
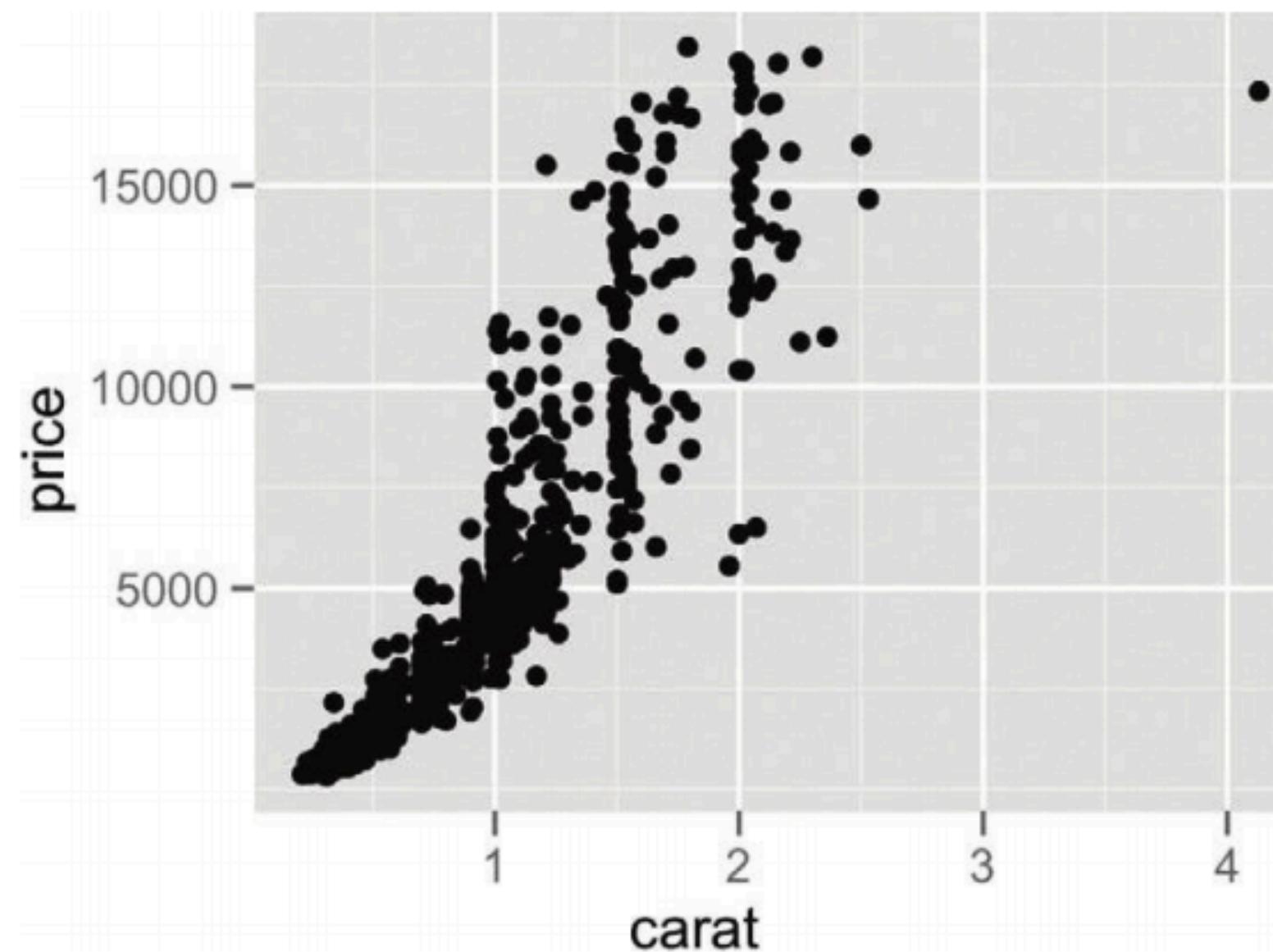
## → Express Values



Show values directly

No keys

Derived attributes can be  
more expressive

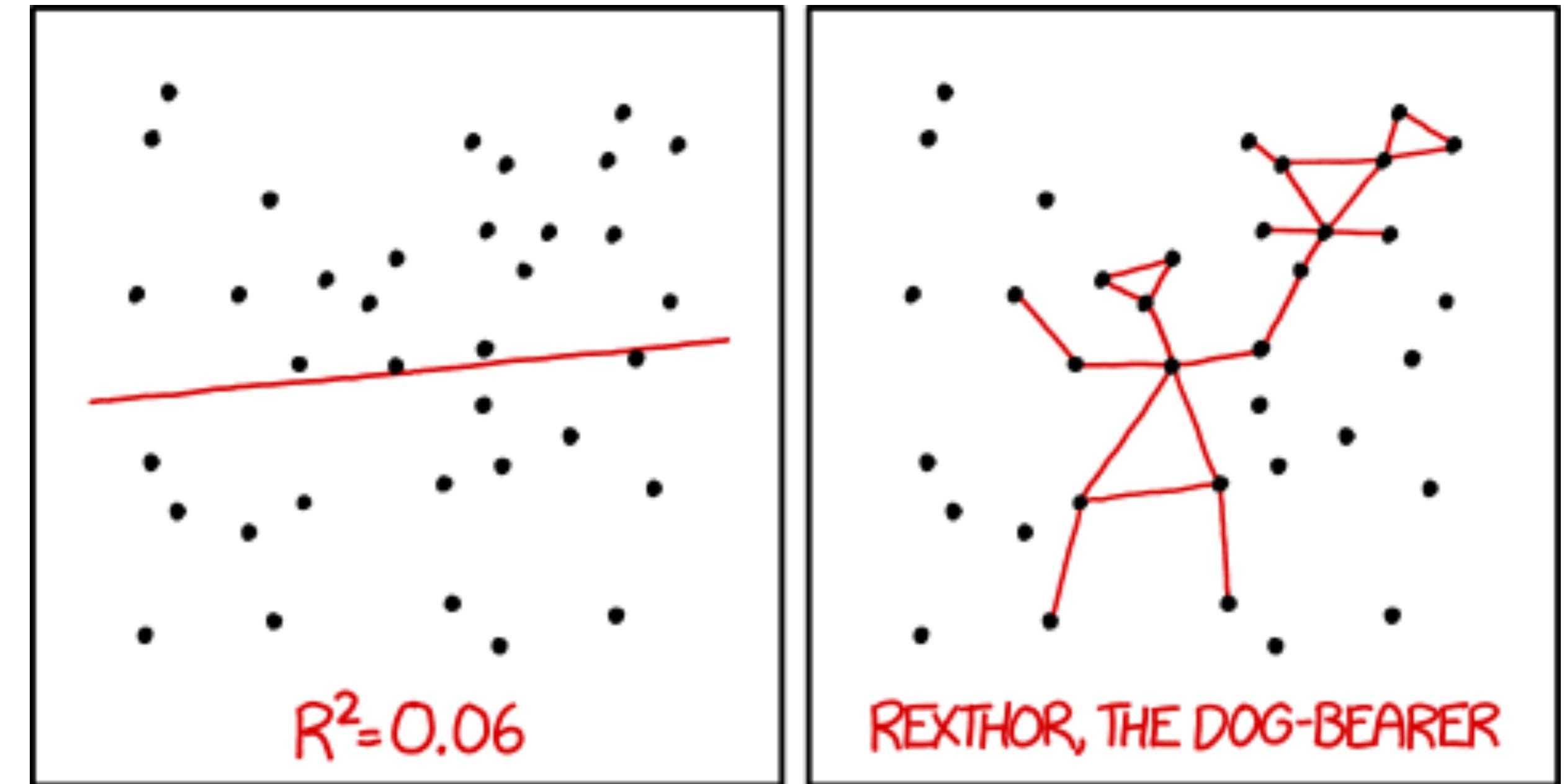
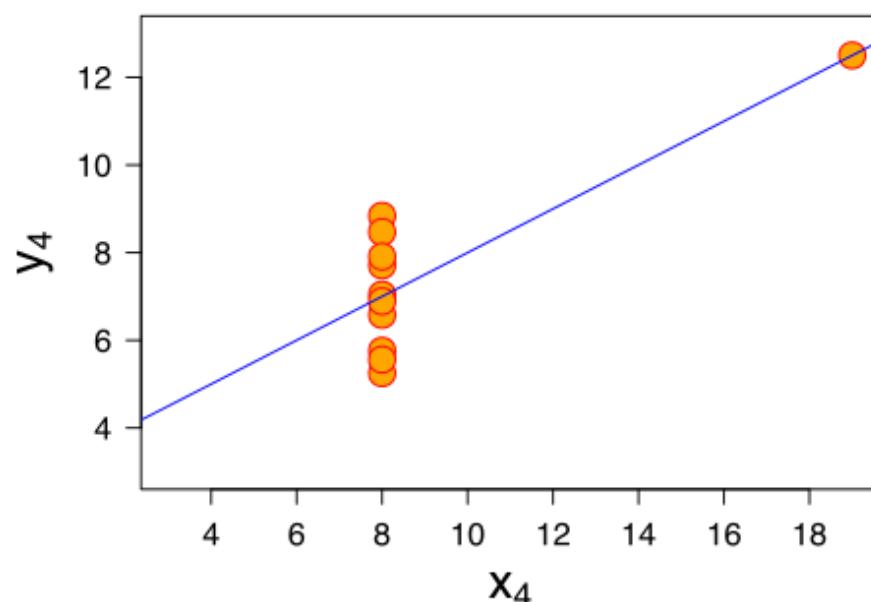
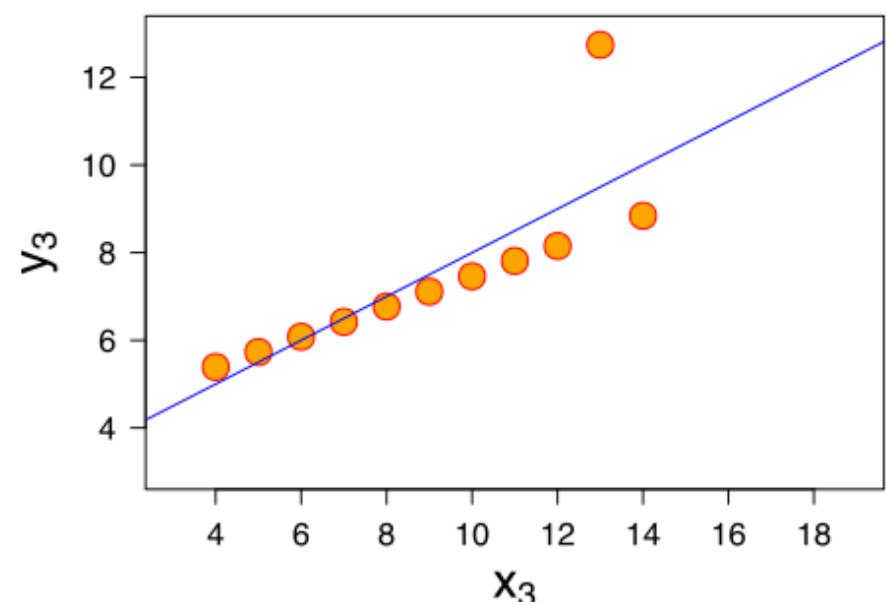
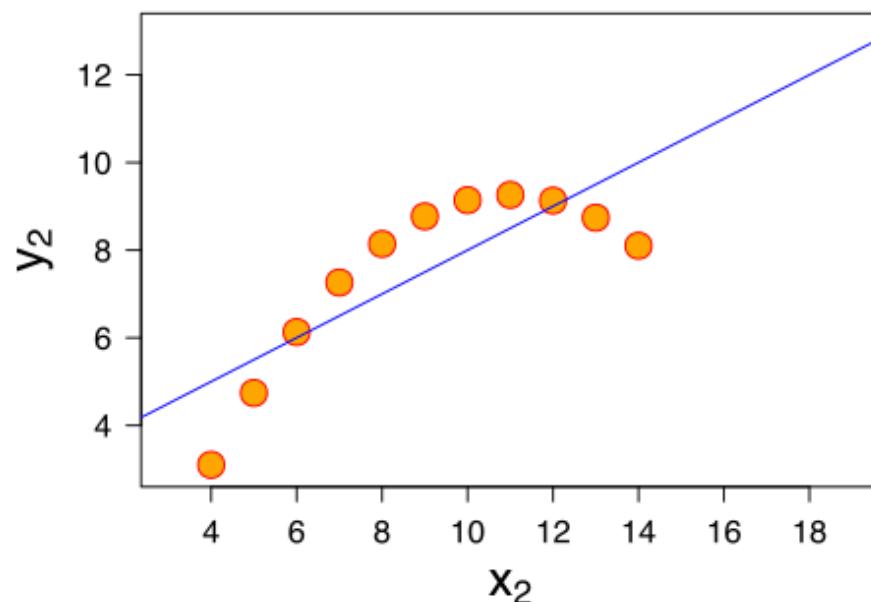
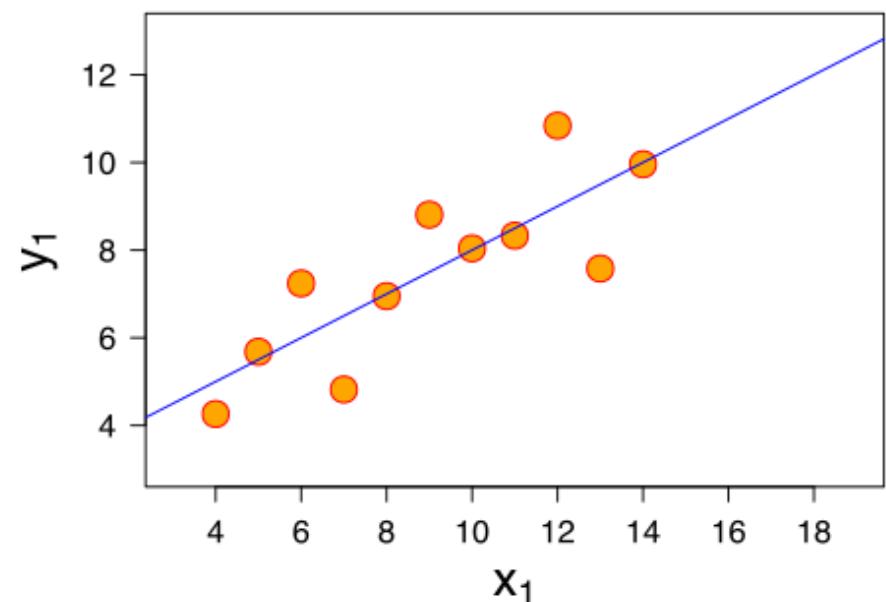


[Wickham 10]

# Showing trends

Helps showing trends at a quick glance

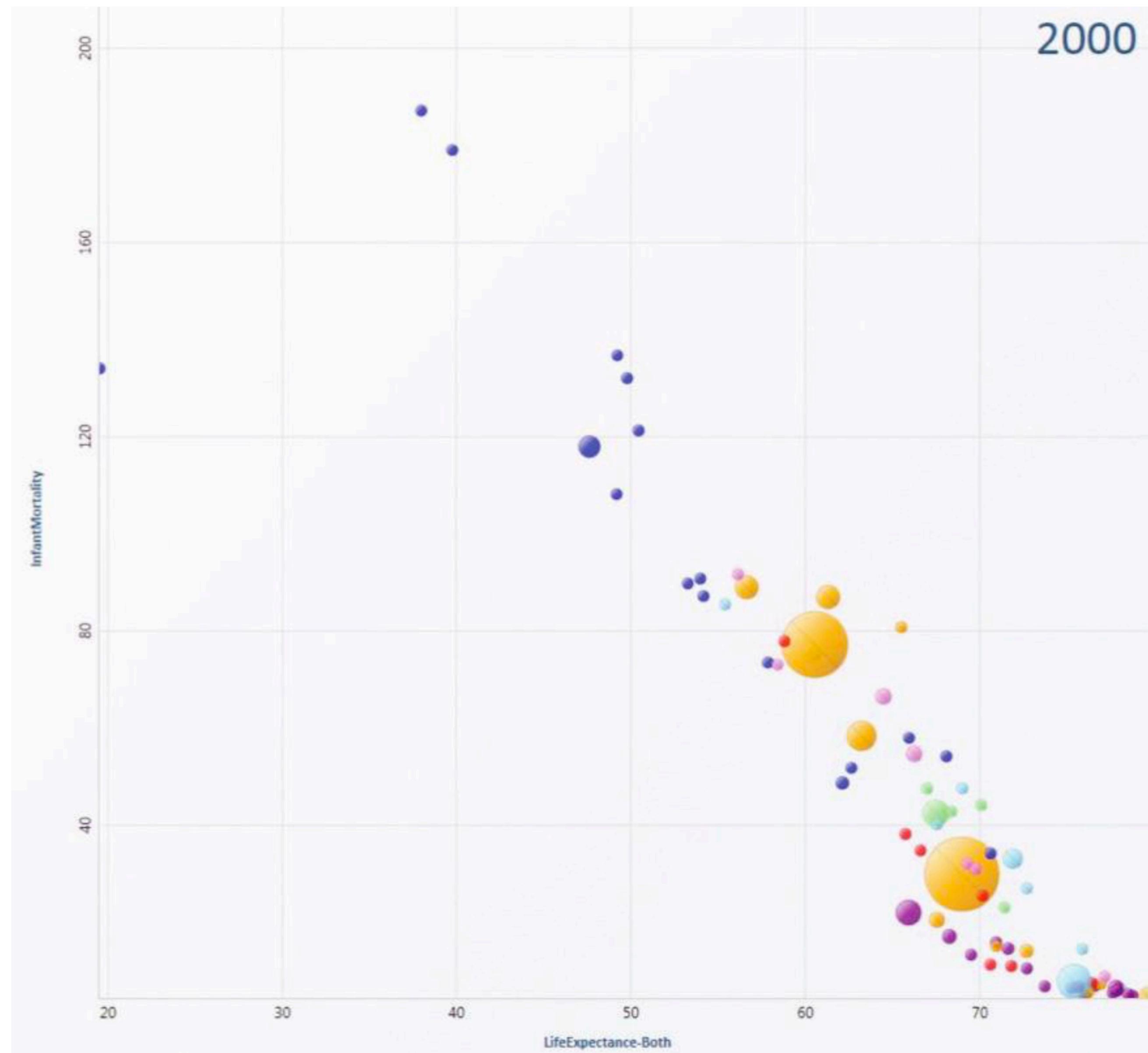
Most often use a linear regression



I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

[xkcd]

# Bubble plots



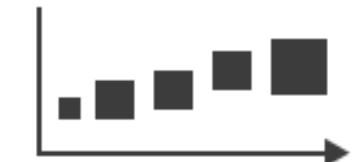
# Simple bar, line, dot charts

④ Separate, Order, Align Regions

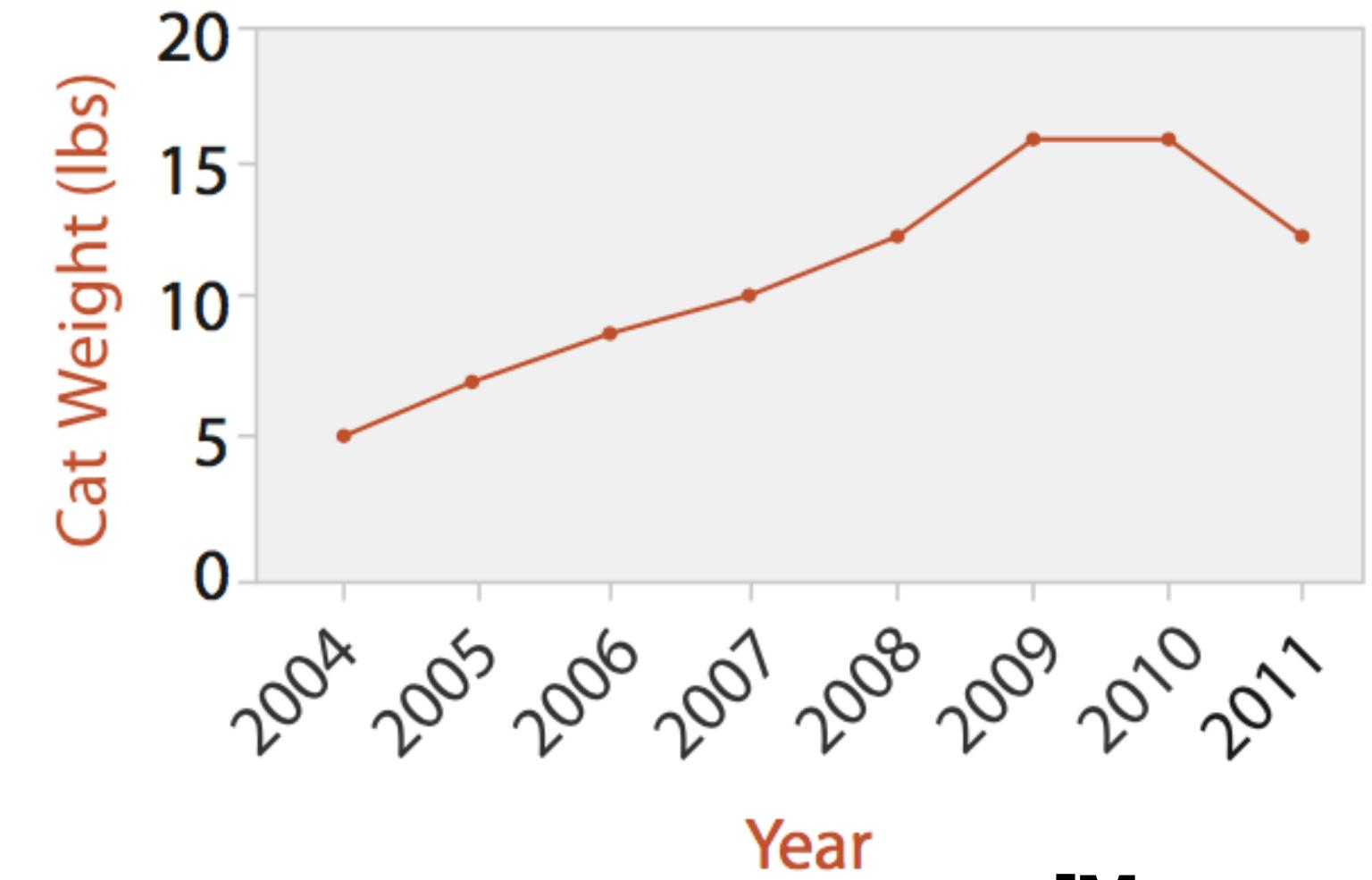
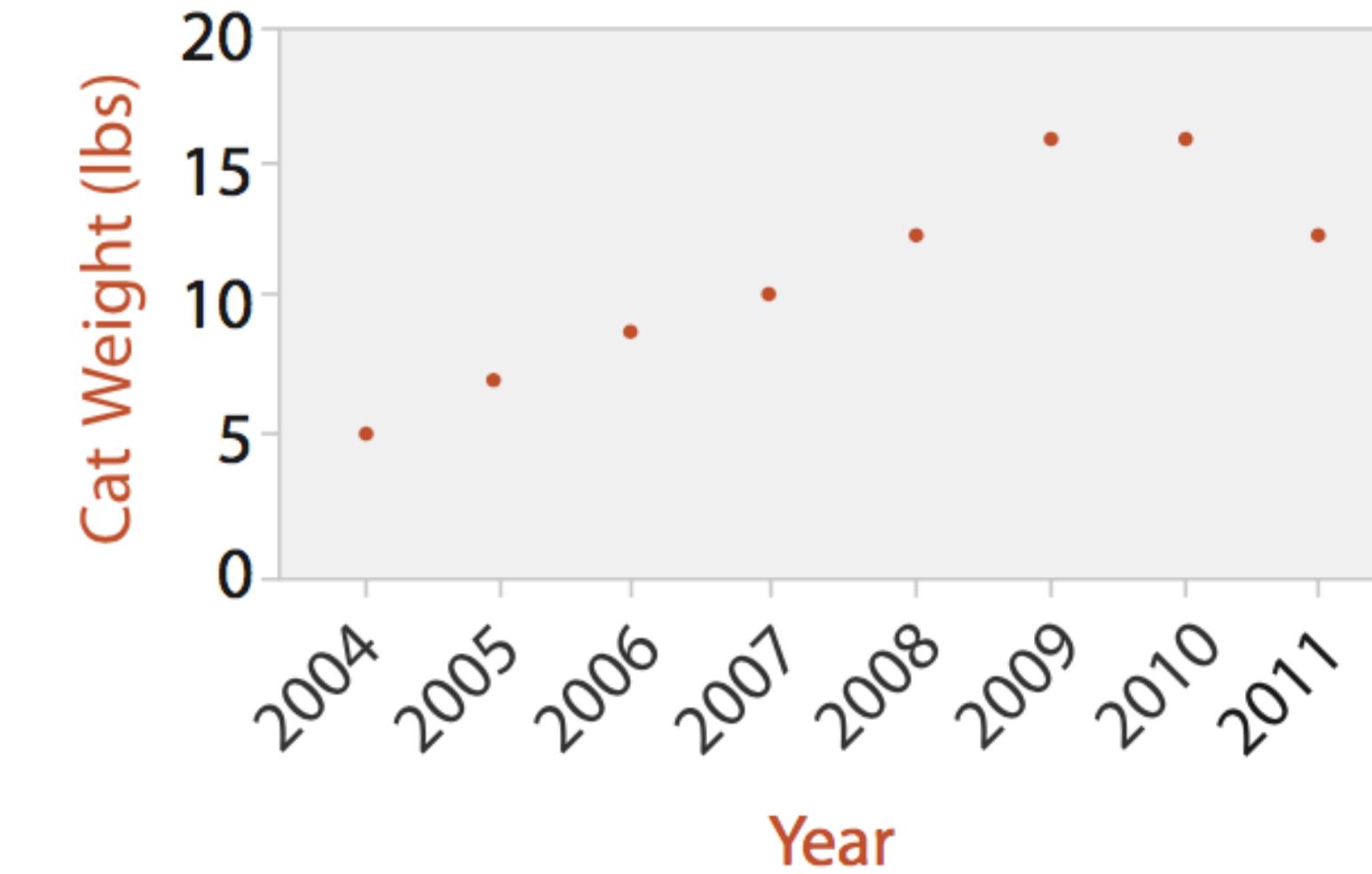
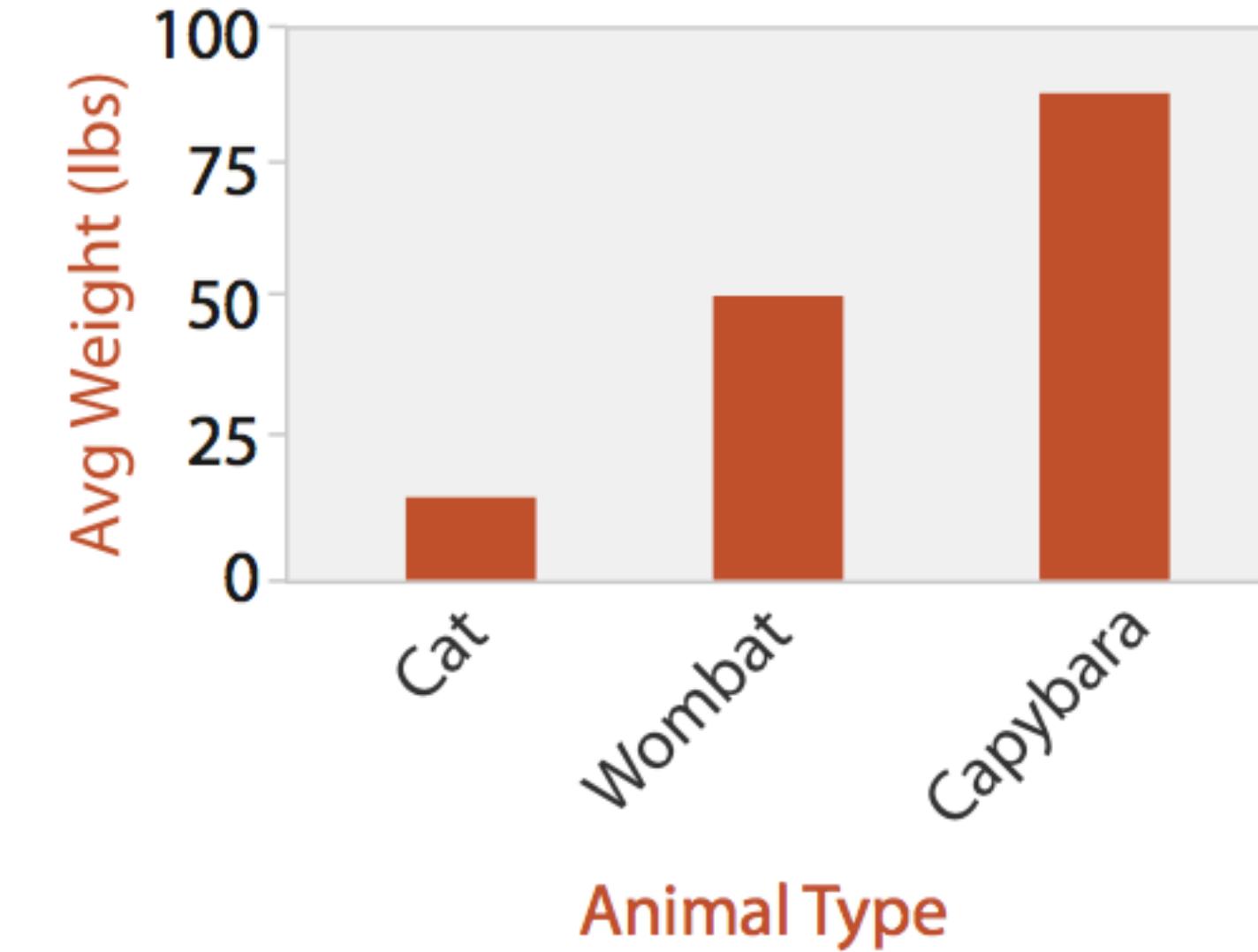
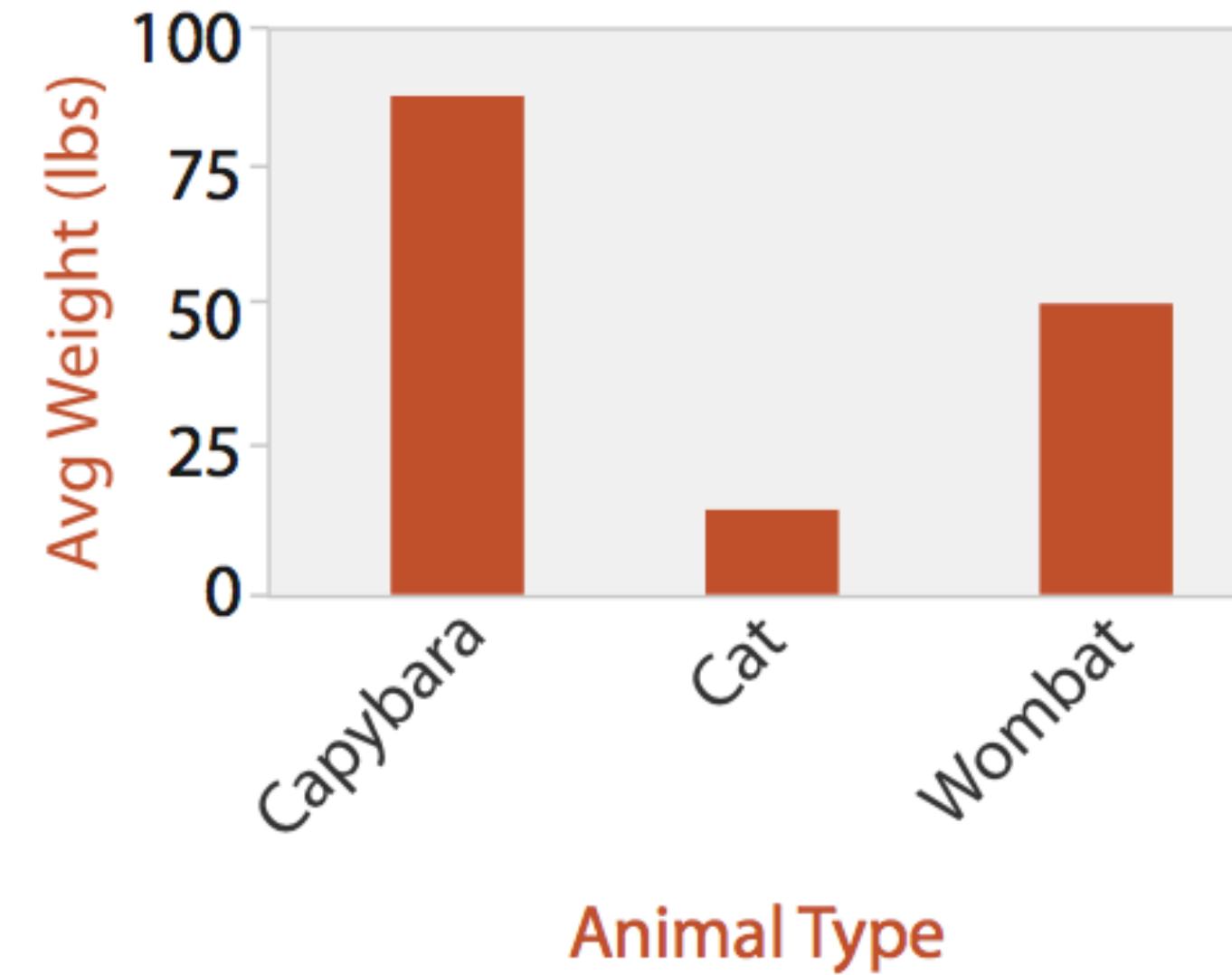
→ Separate



→ Order



→ Align



Encode a single key attribute

# Capybara

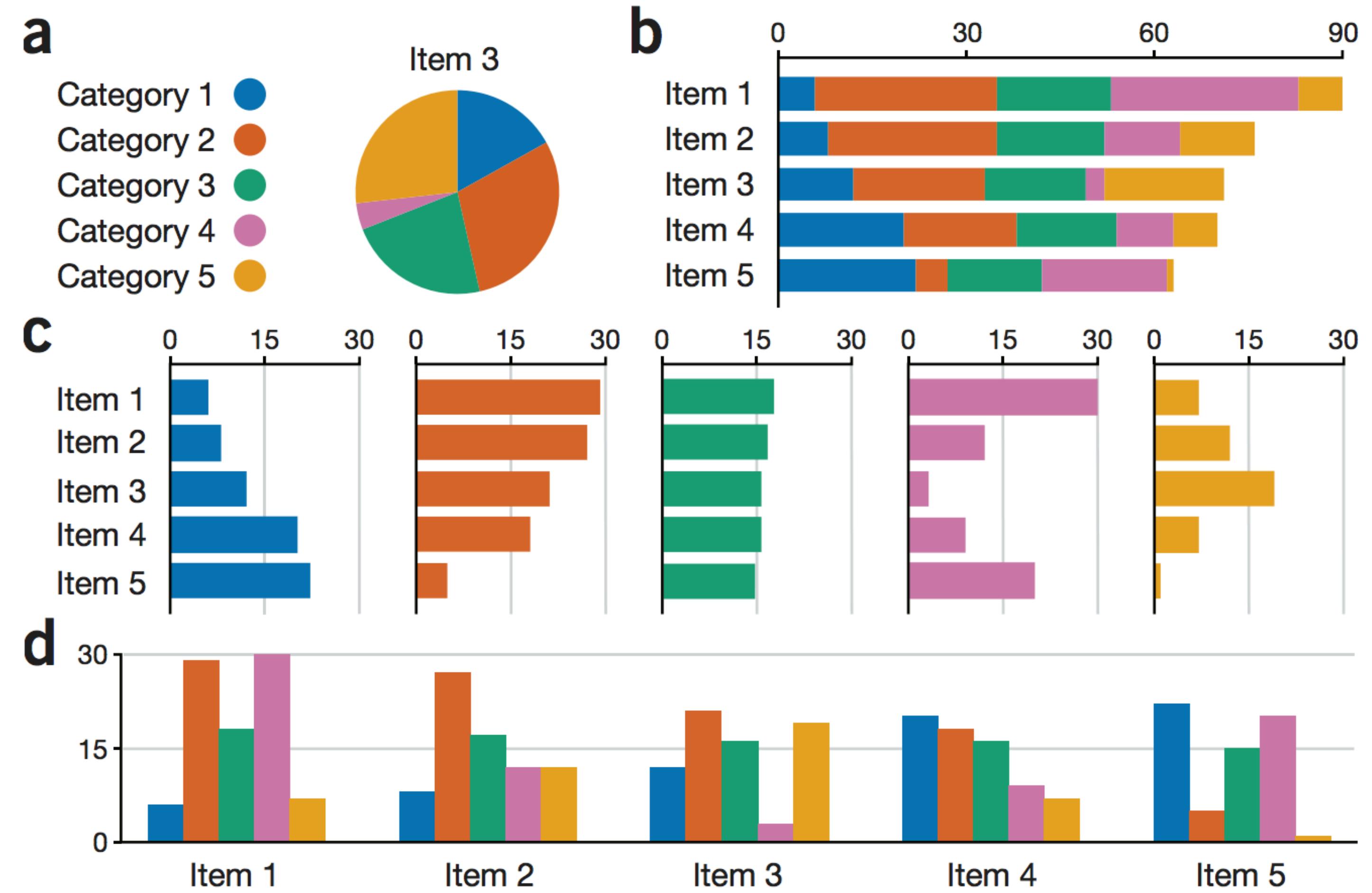
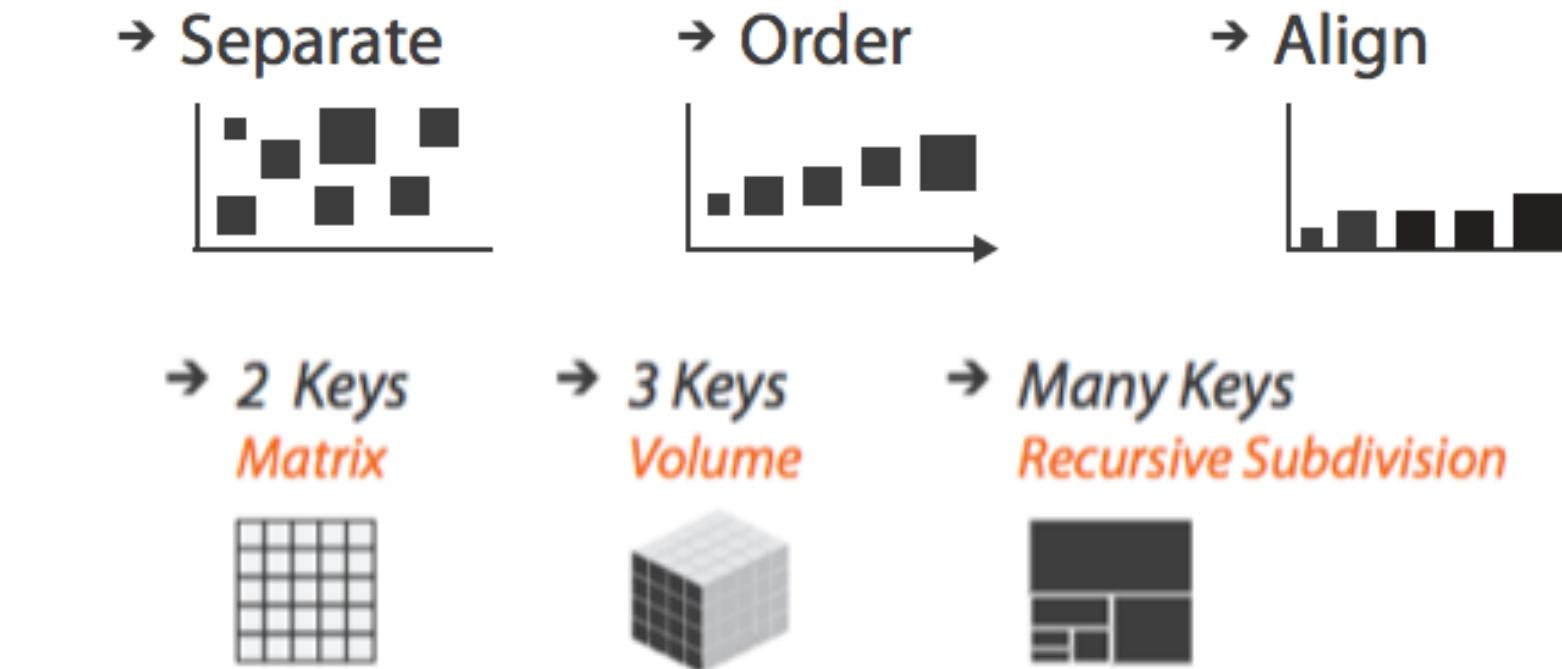


# Wombat



# Stacked bar charts

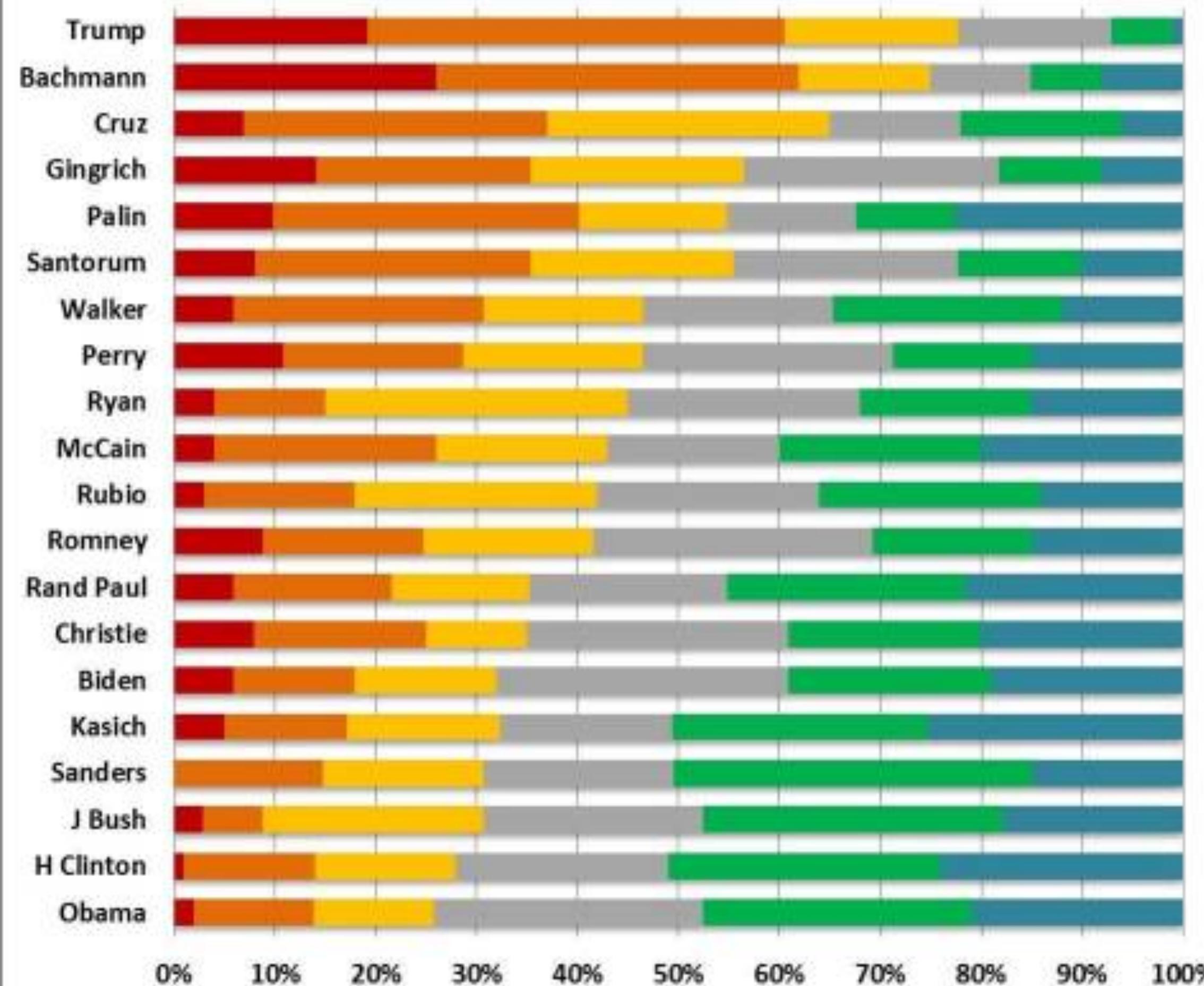
## Separate, Order, Align Regions



# Who Lies More: A Comparison

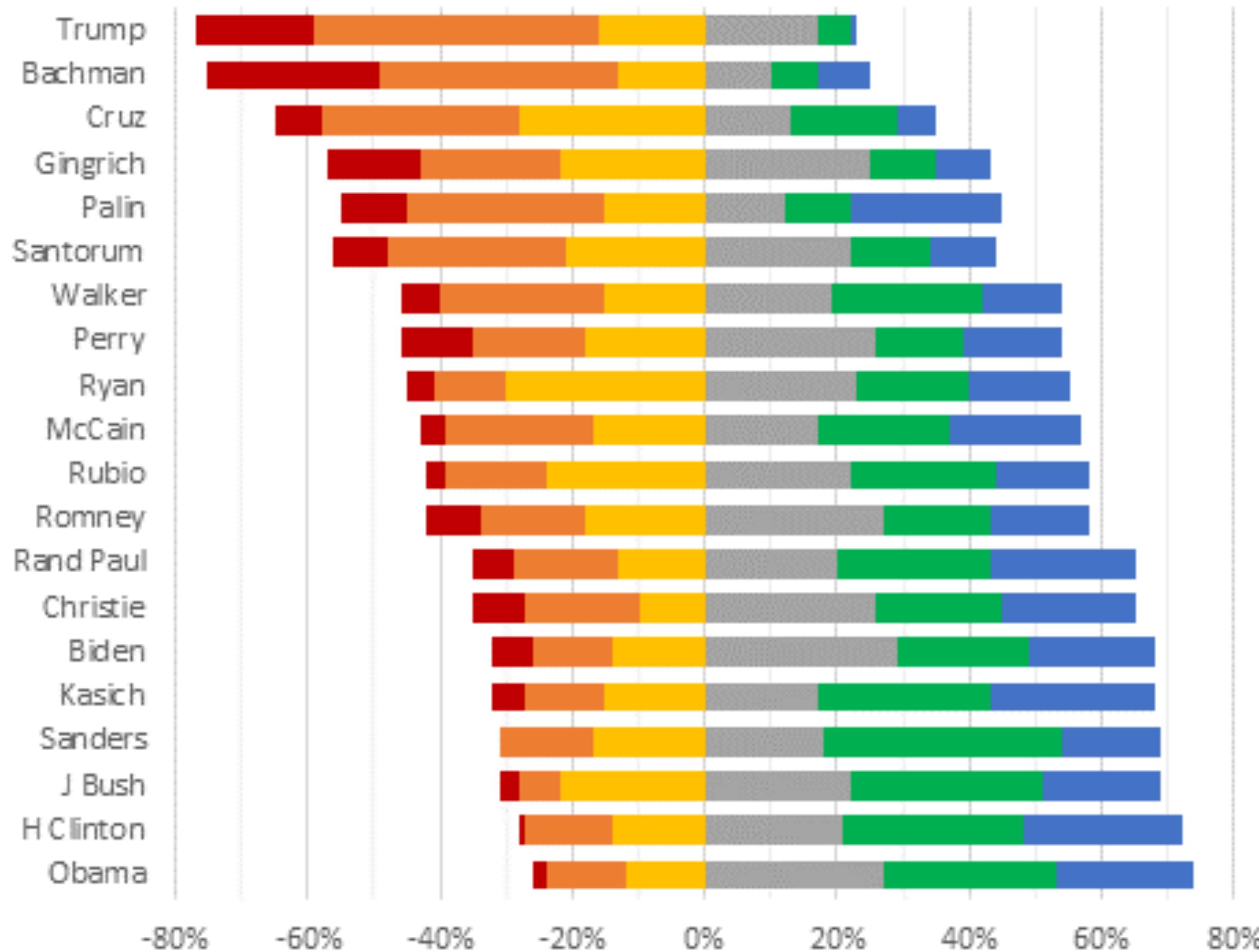
PolitiFact, an independent fact-checking website, has graded more than 50 statements since 2007 from each of these candidates. Here is how they rank.

■ Pants On Fire: ■ False: ■ Mostly False: ■ Half True: ■ Mostly True: ■ True:



## Who Lies More – Diverging Stacked Bar

■ Pants on Fire ■ False ■ Mostly False ■ Half True ■ Mostly True ■ True

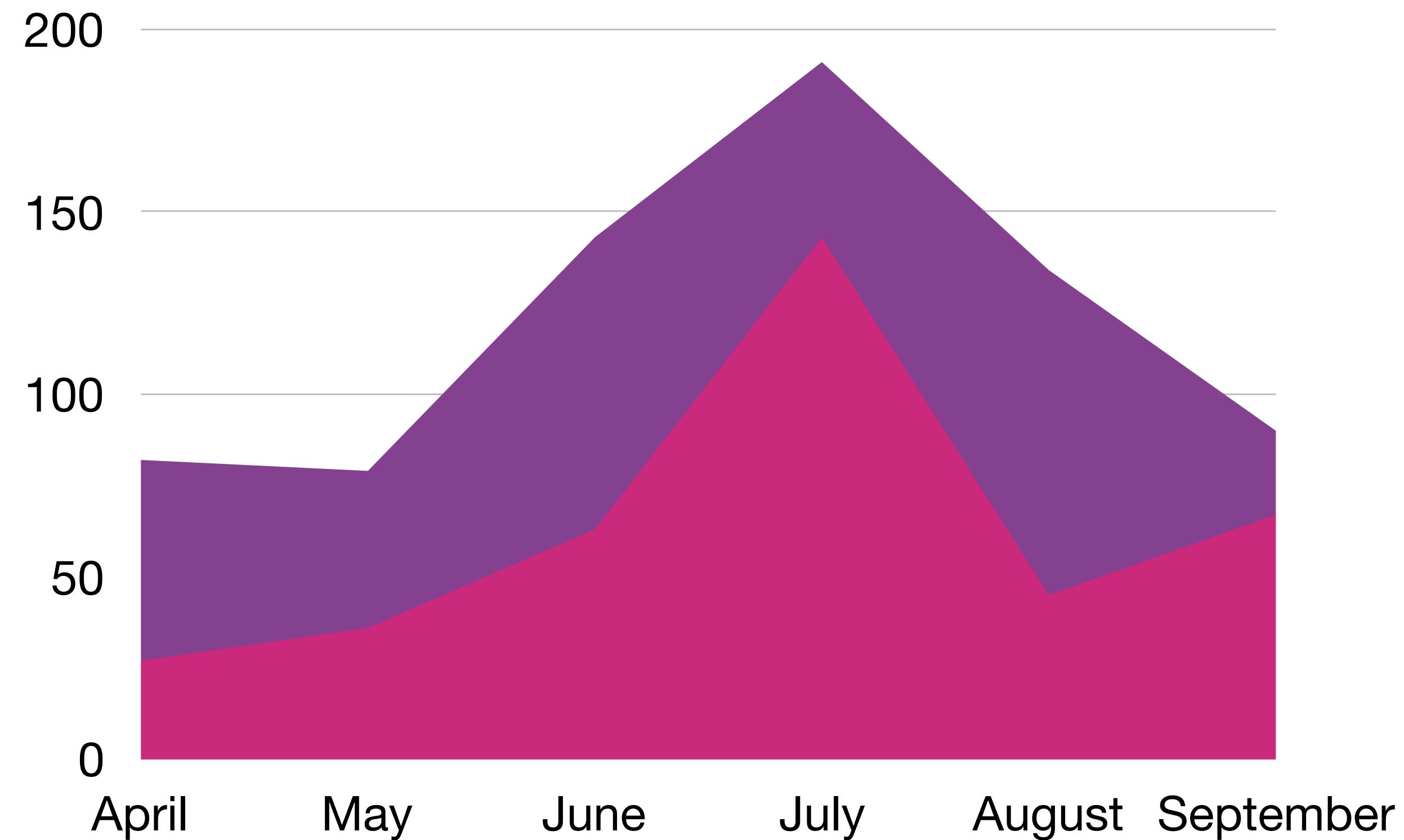


# Stacked area chart

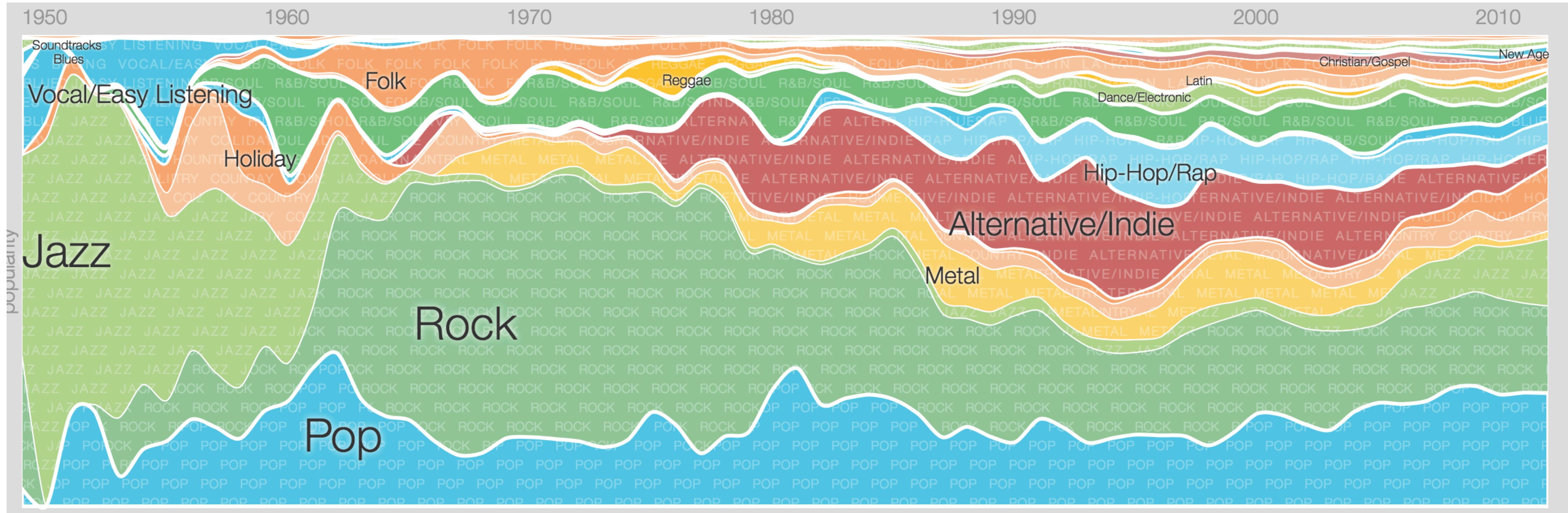
Color the area below the line that connects the values from the dataset

Display the development of quantitative values over an interval

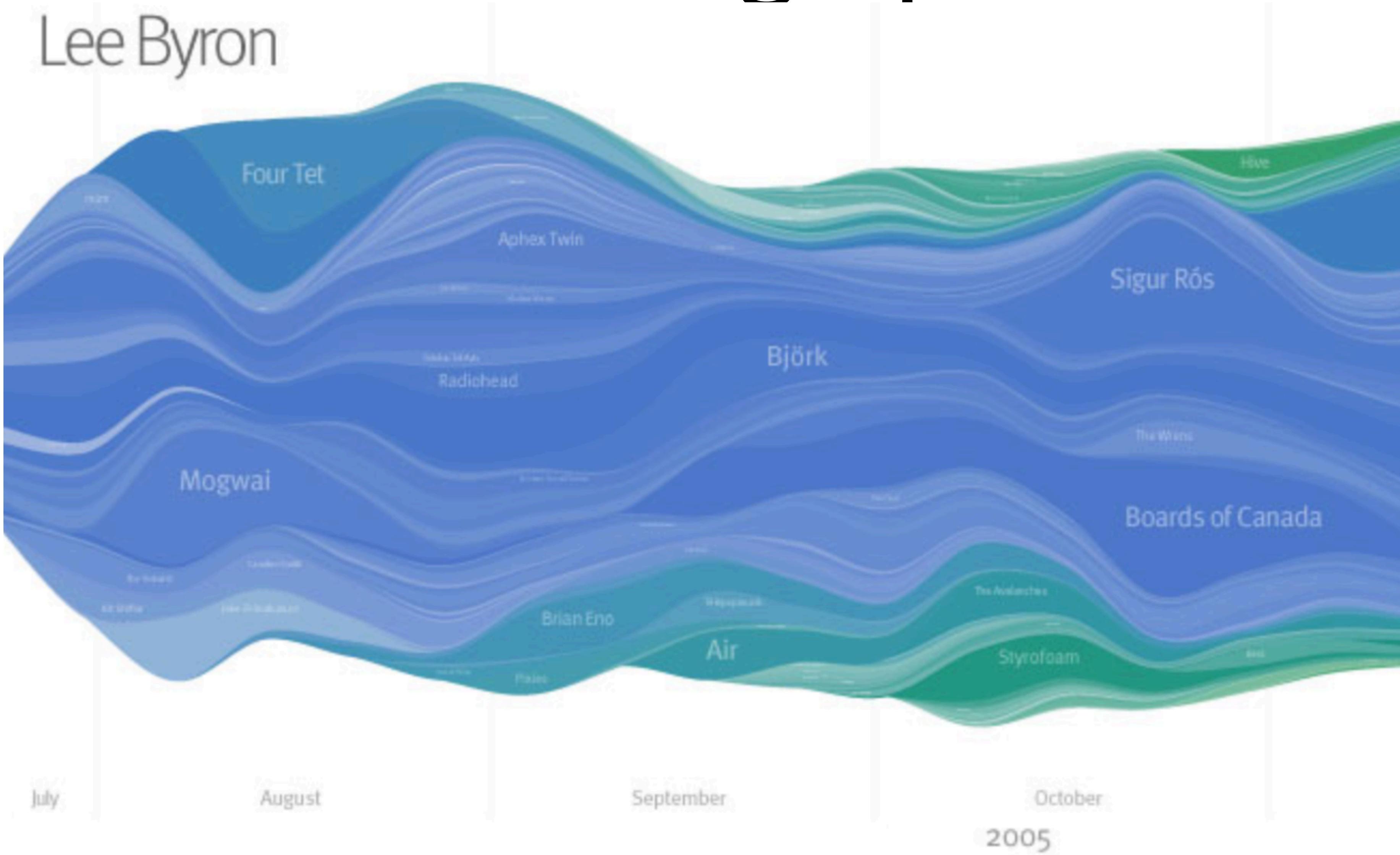
Ordering is paramount



# 100% stacked area chart



# Streamgraph



[Link to original paper](#)

## 5 ALGORITHMS FOR STACKED GRAPH DESIGN

There are four main ingredients that determine a generalized stacked graph. The shape of the overall silhouette is the first ingredient; this shape is critical since it determines the overall slopes and curvature of the individual layers. The second important parameter is the ordering of the layers, which may be chosen to conform to different aesthetic criteria. As with any visualization, labels are important as well—but the organic forms of a Streamgraph mean that labels

following figures, we use a synthetic data set with randomly assigned

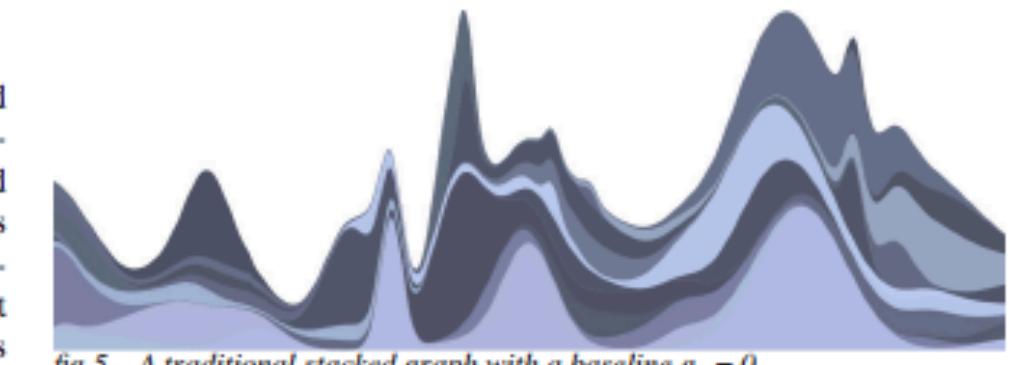


fig 5 – A traditional stacked graph with a baseline  $g_0 = 0$

colors to distinguish layers.) In this layout the size of the sum of the series is easy to read, potentially at the expense of the legibility of the individual layers as mentioned in design issue (D).

A simple alternative layout was suggested by Havre et al in the ThemeRiver system. They used a layout symmetric around the x-axis. Mathematically, this can be expressed as:

$$g_0 + g_n = 0$$

or, from the definition of  $g_o$ ,

$$2g_0 + \sum_{i=1}^n f_i = 0$$

which yields the ThemeRiver (fig 6) solution for  $g_o$ :

$$g_0 = -\frac{1}{2} \sum_{i=1}^n f_i$$

This is a simple enough definition, but there is another way to look at this formula which suggests some generalizations. What is special about the symmetric layout? Aside from a certain aesthetic quality, this layout has the effect of minimizing some important quantities. In particular, at each point, the silhouette is as close as possible to the x-axis, and in addition the slopes of the top and bottom of the silhouette are in a sense as small as possible (in the sense of total sums of squares). This directly addresses design issue (C) by making the overall graph much less “spiky” thus greatly reducing the horizontal space needed to satisfy Cleveland’s principal. To see this, recall the following fact about averages. For any set of real numbers  $\{a_1, \dots, a_n\}$ , the value of  $x$  that minimizes

$$\sum_{i=1}^n (x + a_i)^2$$

is

$$x = -\frac{1}{n} \sum_{i=1}^n a_i$$

From this fact it follows that the value of  $g_o$  that gives the symmetric ThemeRiver layout minimizes the sum of squares of the top and bottom of the silhouette of the graph (at each point in  $[0,1]$ ):

$$\text{silhouette}(g_0) = g_0^2 + g_n^2$$

since

$$g_0^2 + g_n^2 = g_0^2 + (g_0 + \sum_{i=1}^n f_i)^2$$

bottom, as in the case  $i=0$ , we take the sum to be empty and equal to zero.)

We might also be interested in minimizing the sum of squares of the slopes at each value of  $x$ :

$$\text{wiggle}(g_0) = \sum_{i=0}^n g_i'^2 = \sum_{i=0}^n (g_0' + \sum_{j=1}^i f_j')^2$$

By the same logic, the deviation quantity is minimized when

$$g_0 = -\frac{1}{n+1} \sum_{i=0}^n \sum_{j=1}^i f_j = -\frac{1}{n+1} \sum_{i=1}^n (n-i+1)f_i$$

Moreover, differentiating this yields:

$$g_0' = -\frac{1}{n+1} \sum_{i=0}^n \sum_{j=1}^i f_j'$$

which minimizes the *wiggle* measure as well. In other words, this choice of  $g_o$  has the effect of simultaneously minimizing both distance from the x-axis and variation in slope. Moreover, the formula is extremely efficient to calculate. As fig 7 shows, it creates a more “even” layout compared to the ThemeRiver graph. This minimization directly addresses design issue (A) by attempting to reduce the effect of middle “wiggles” on the surrounding layers.

Nonetheless, the layout can be improved further. In particular, it is potentially problematic that very thin layers receive the same treatment as very thick ones. After all, the thick layers are visually more important. Thus we might want instead to optimize the following quantity, which represents an average of the squares of slopes between the midpoints of each layer, weighted by layer thickness:

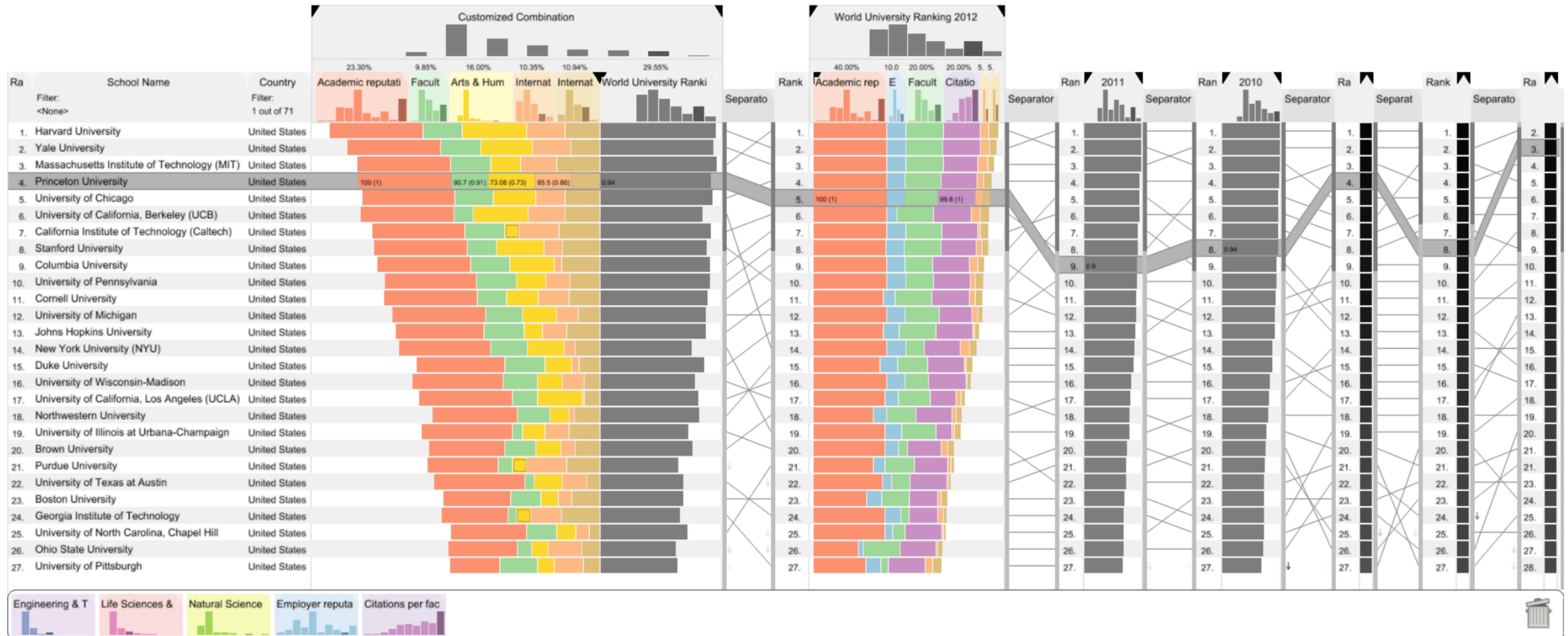
$$\begin{aligned} \text{weighted\_wiggle}(g_o) &= \\ \sum_{i=1}^n \left( \frac{1}{2}(g_i' + g_{i-1}') \right)^2 f_i &= \sum_{i=1}^n \left( g_0' + \frac{1}{2}f_i' + \sum_{j=1}^{i-1} f_j' \right)^2 f_i \end{aligned}$$

By the properties of weighted averages, this is minimized at each point in  $[0,1]$  when

$$g_0' = -\frac{1}{\sum f_i} \sum_{i=0}^n \left( \frac{1}{2}f_i' + \sum_{j=1}^{i-1} f_j' \right) f_i$$

which can be integrated numerically to yield a solution for  $g_o$ . Indeed, it turns out to be equivalent to the algorithm used in the Streamgraph method, as portrayed in Listening History and Box Office Revenue

# LineUp: Rankings





# Pixel-based display (dense)

Each pixel encodes a value as a color

## Pros

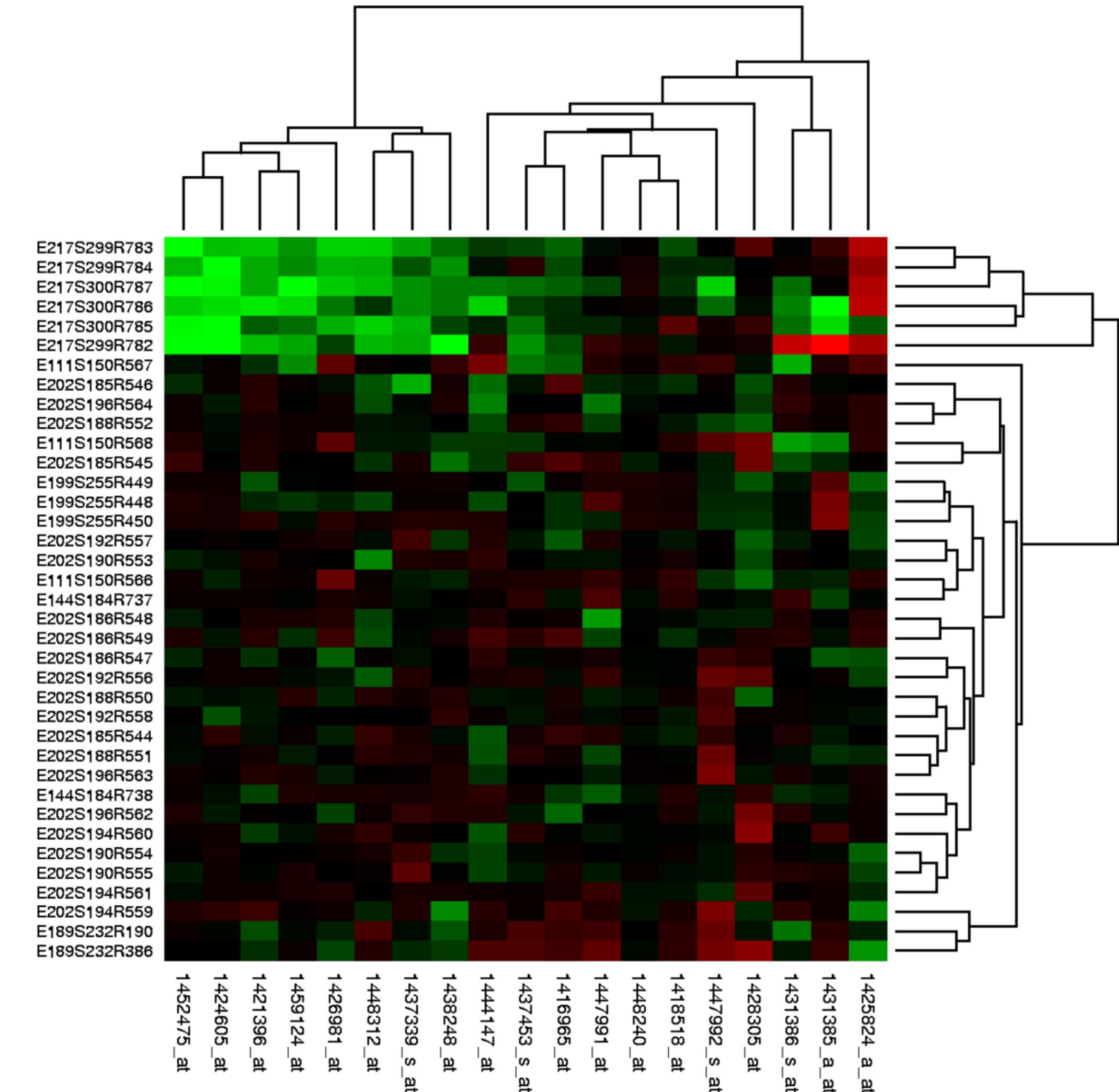
Scalability

Good for same type of data

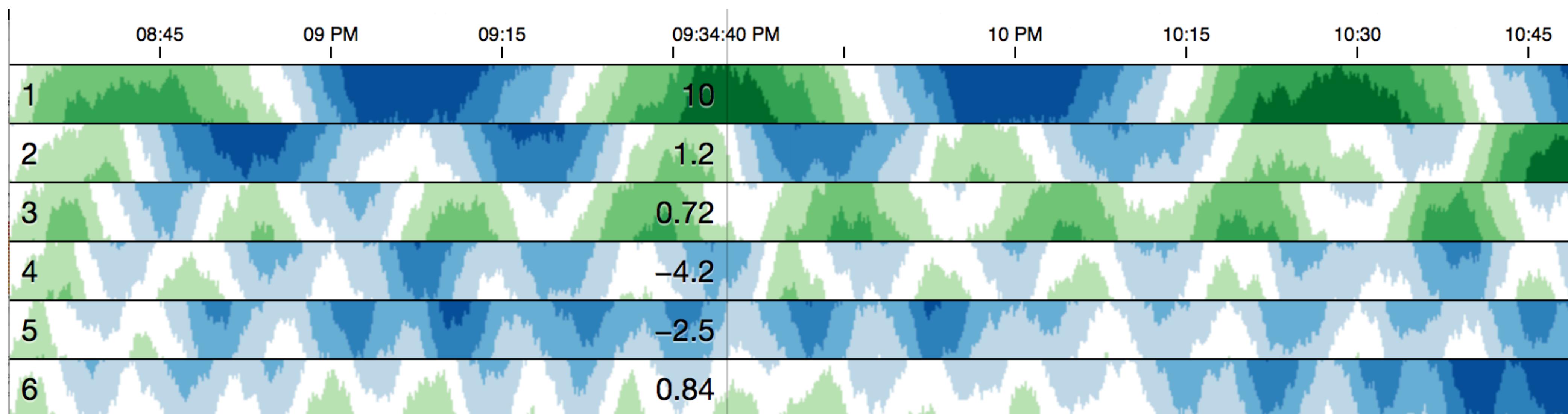
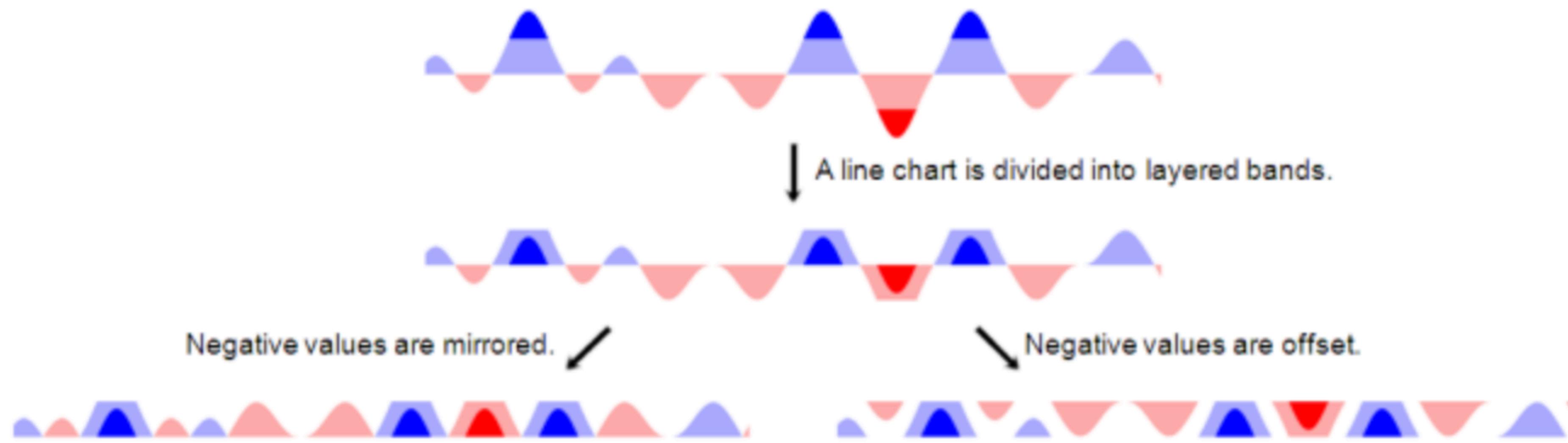
## Cons

ordering is critical

a good color map is important



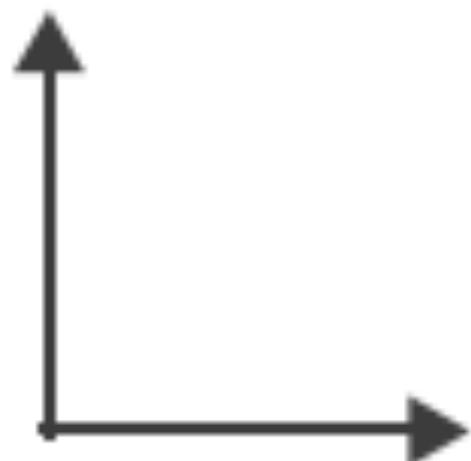
# Multiple Horizon Charts



# Spatial axis orientation

## → Axis Orientation

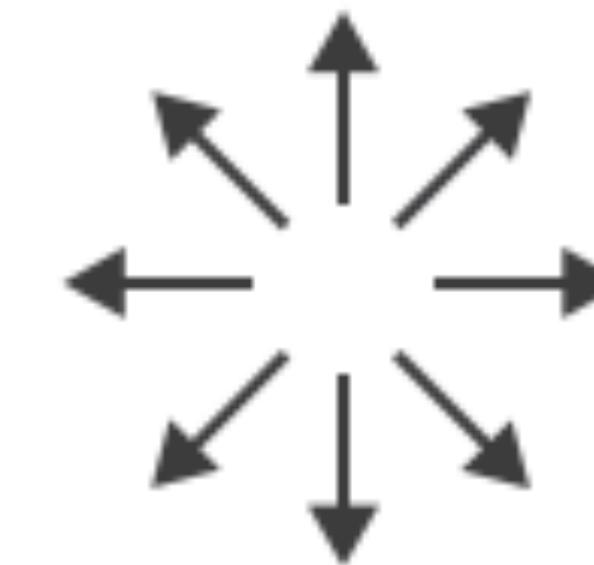
→ Rectilinear



→ Parallel



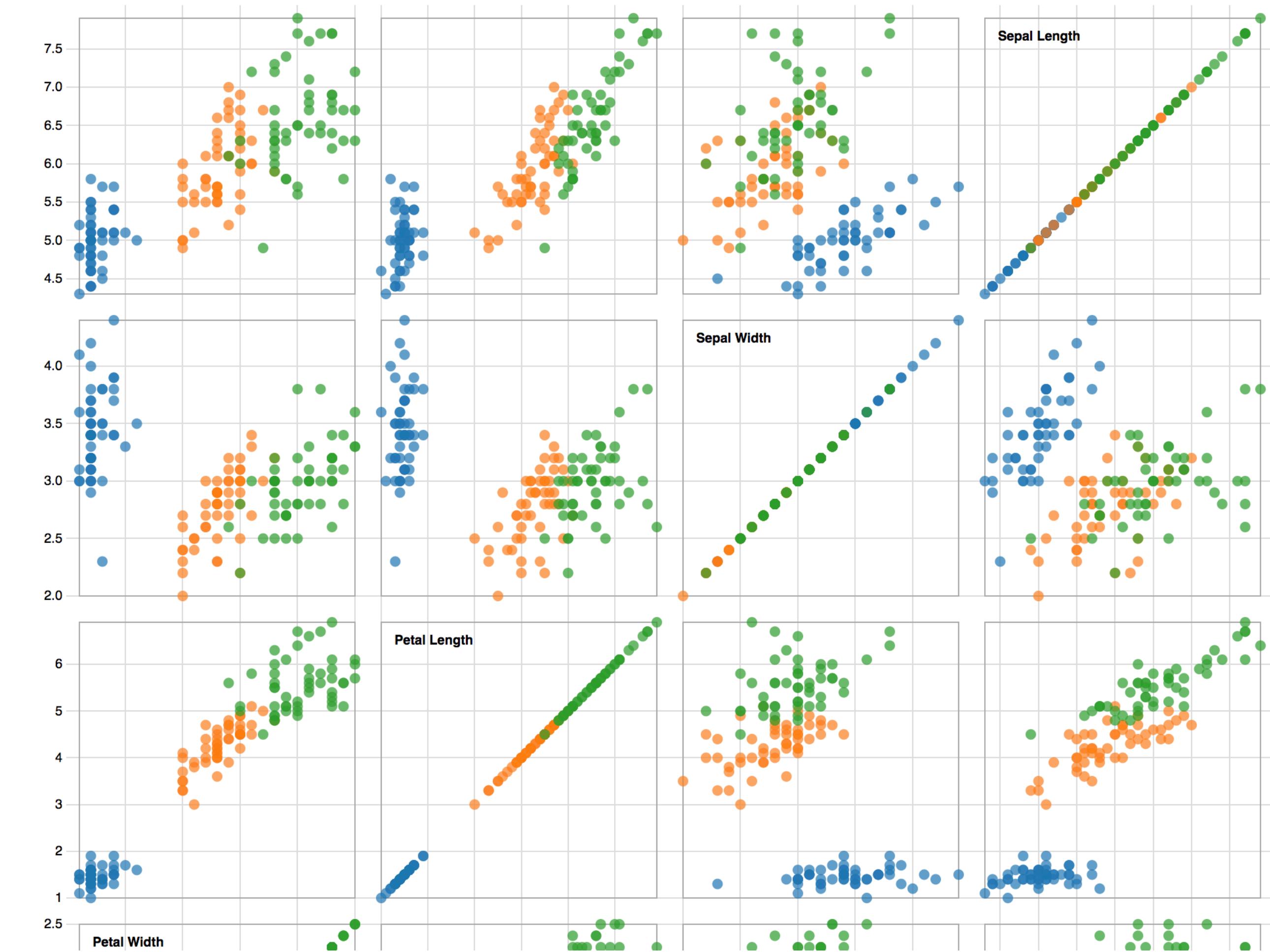
→ Radial



# Scatterplot matrix (SPLOM)

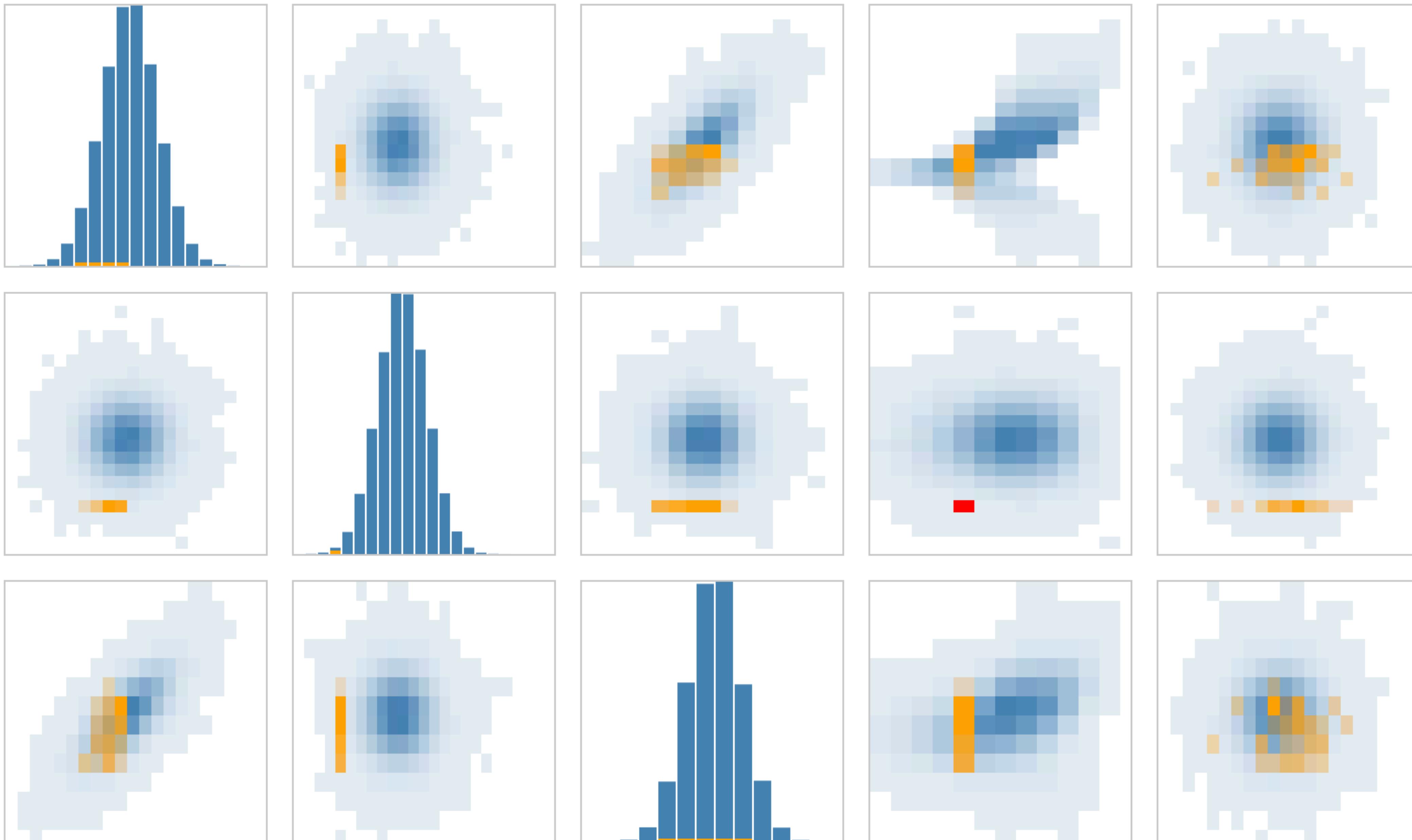
Shows all possible pairwise combinations of attributes

Original attributes are used as rows and columns.



## Interactive Binned Scatterplot Matrix

Dimensions: 5 Bins: 20 Data Points: 100k



# SPLOM

**To be effective SPLOM needs to:**

- have not too many dimensions ~ 20 dim

- have brushing capabilities

- use a F+C technique

**If the dataset is too complex:**

- Use clustering and aggregation

- Choose interesting dimension and order (via exploratory analysis)

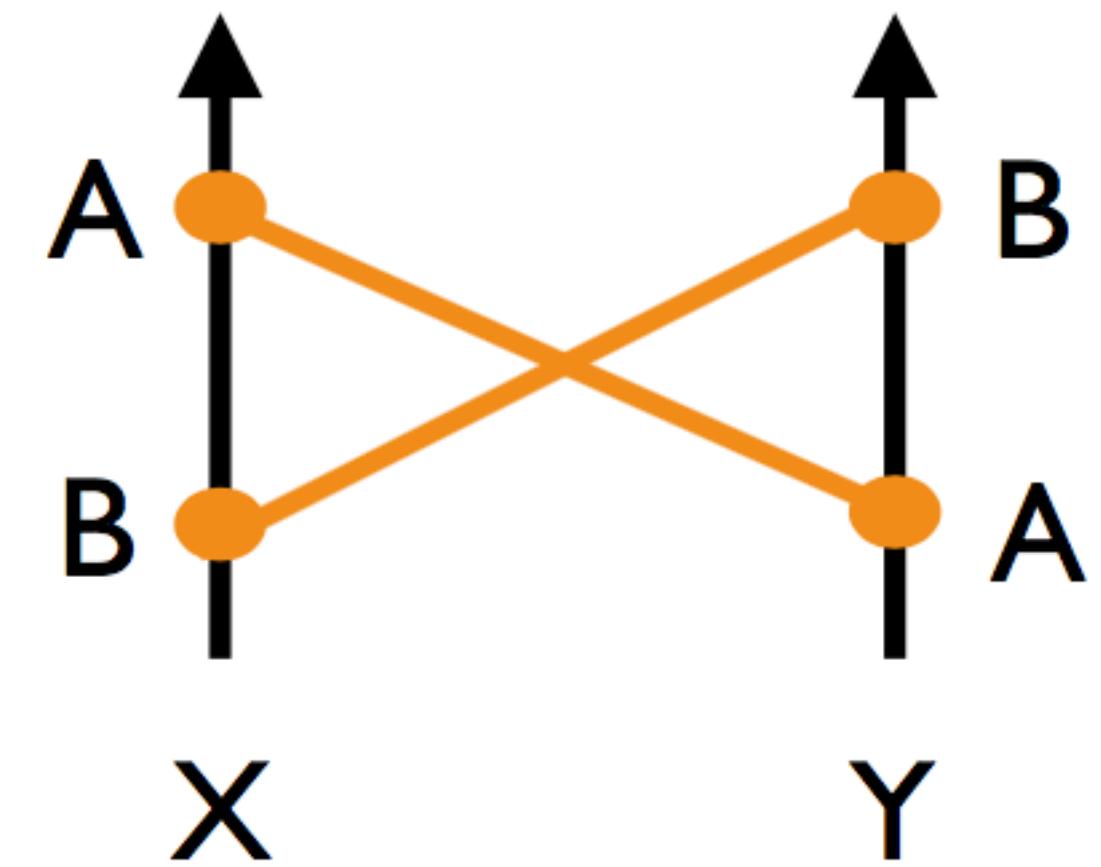
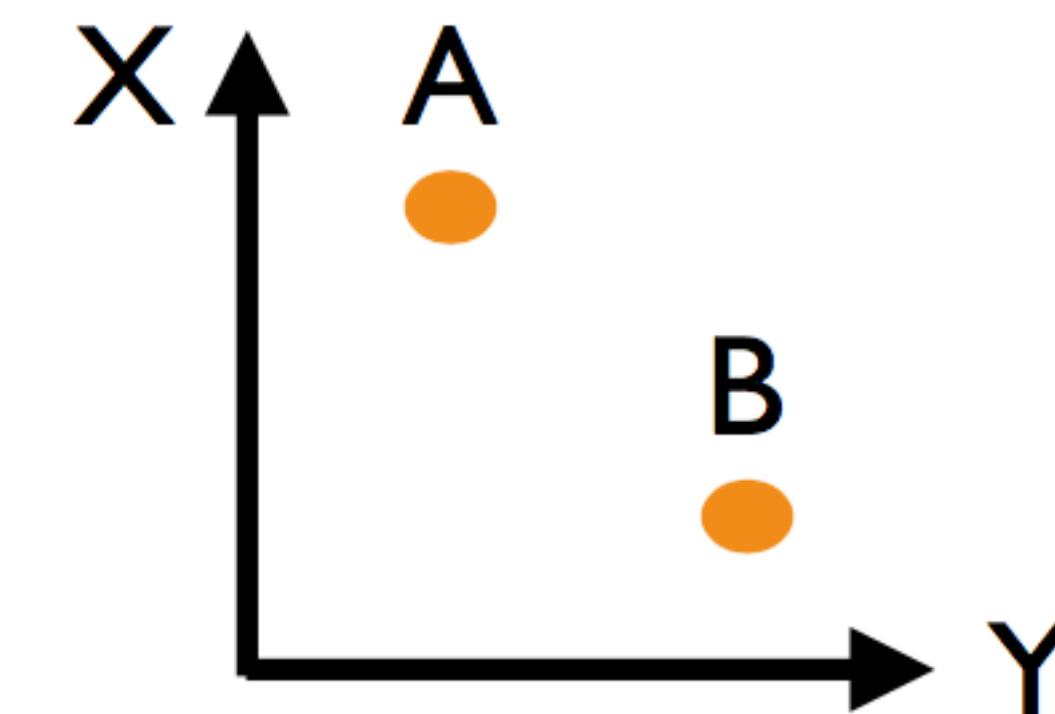
# Parallel Coordinates (PC)

Each axe represents an attribute

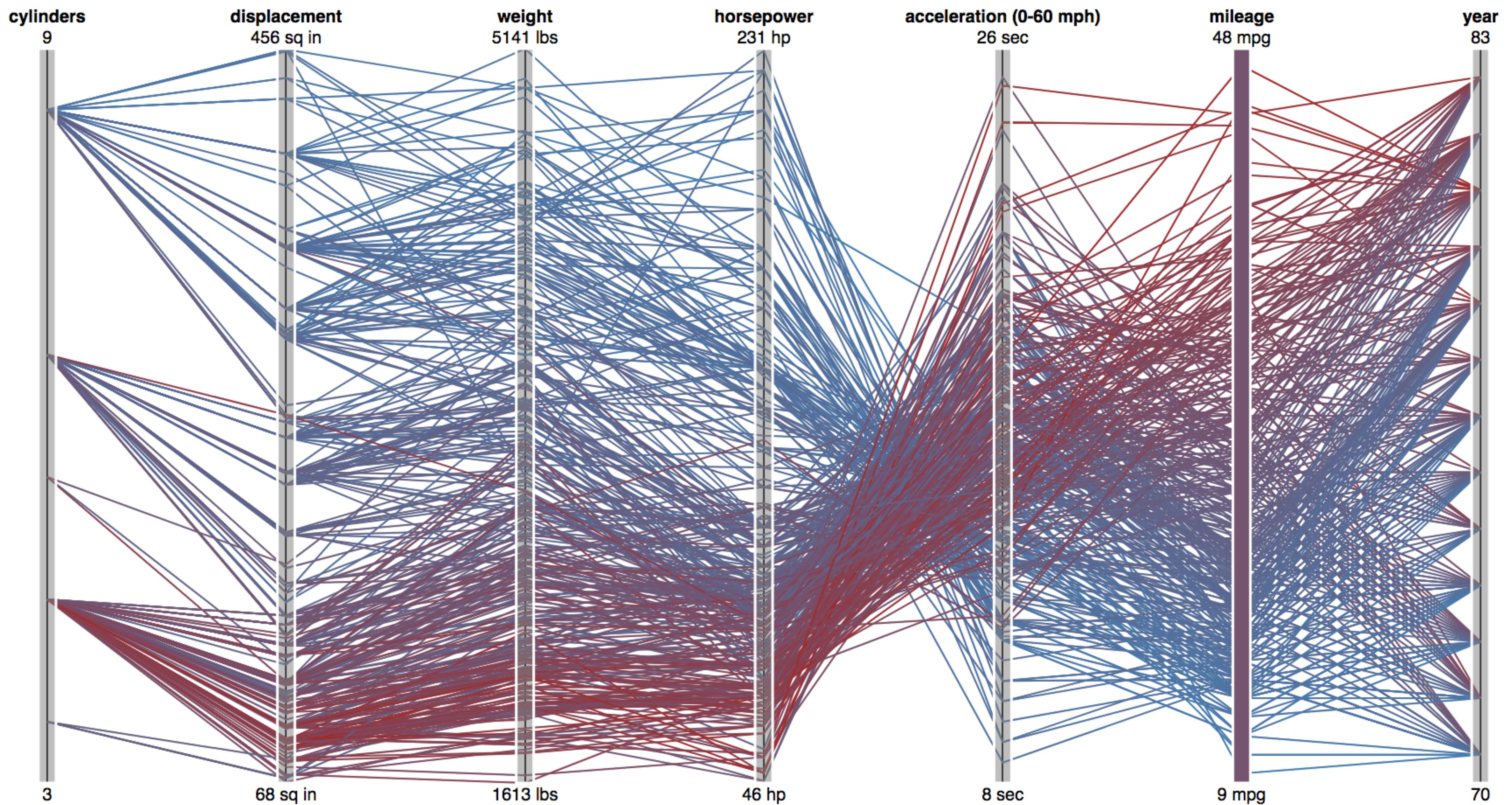
Values are connected if them belong  
to the same record

Works with heterogeneous data

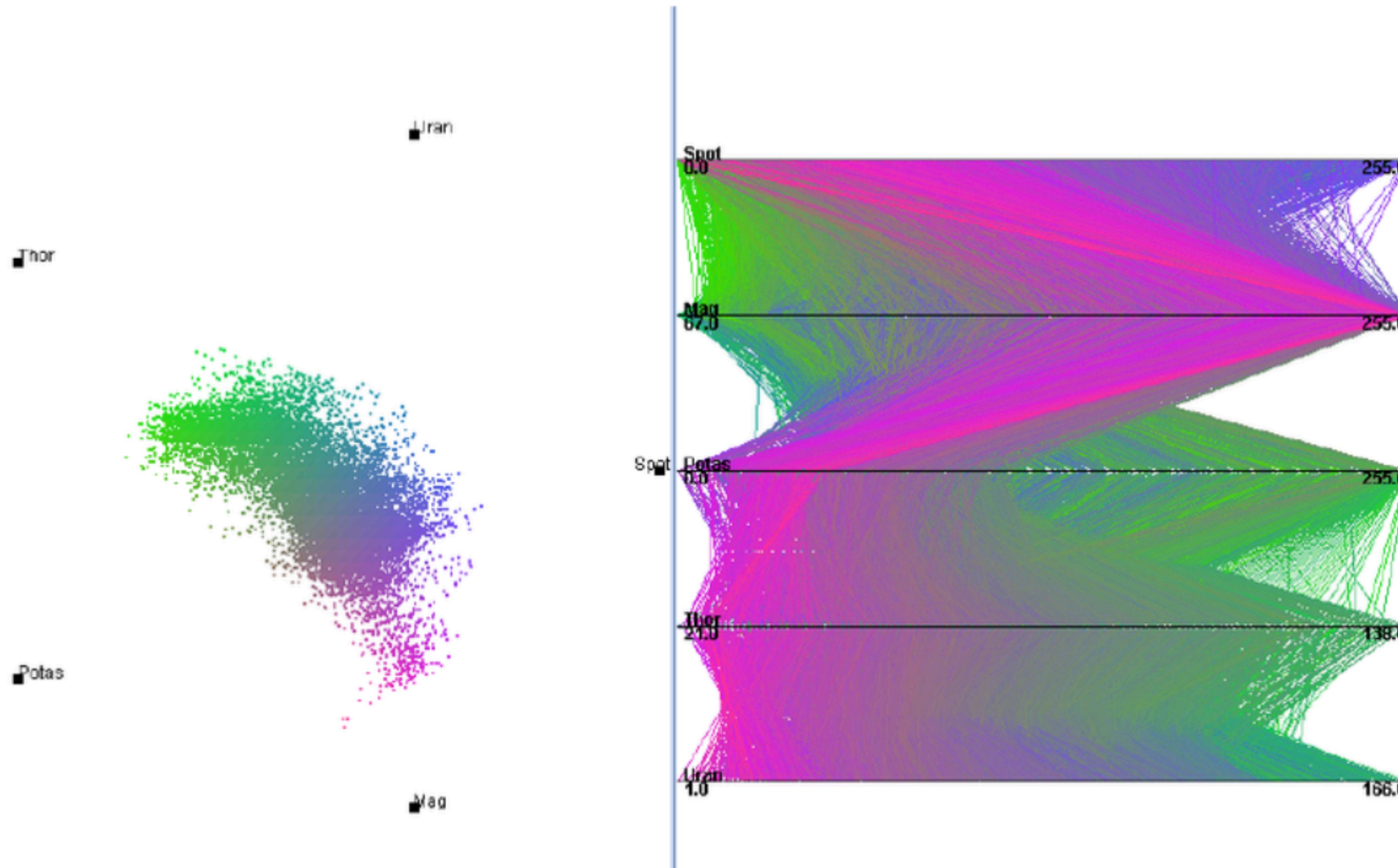
Correlations only between adjacent  
axes



# Parallel Coordinates



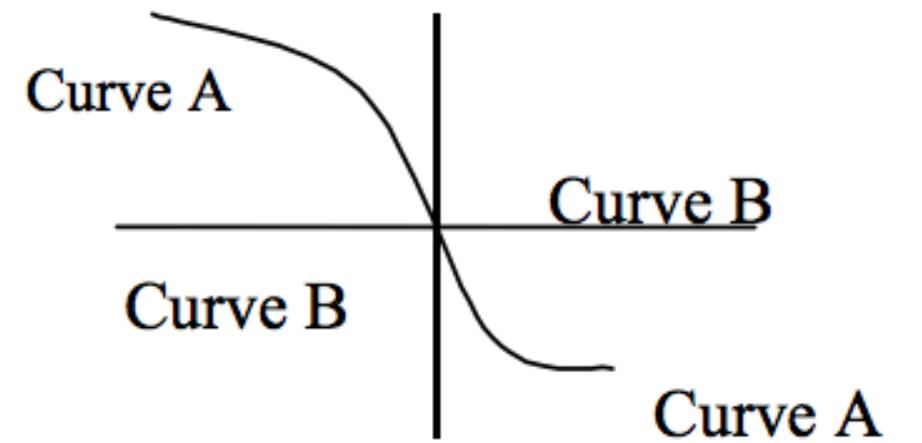
# Issues?



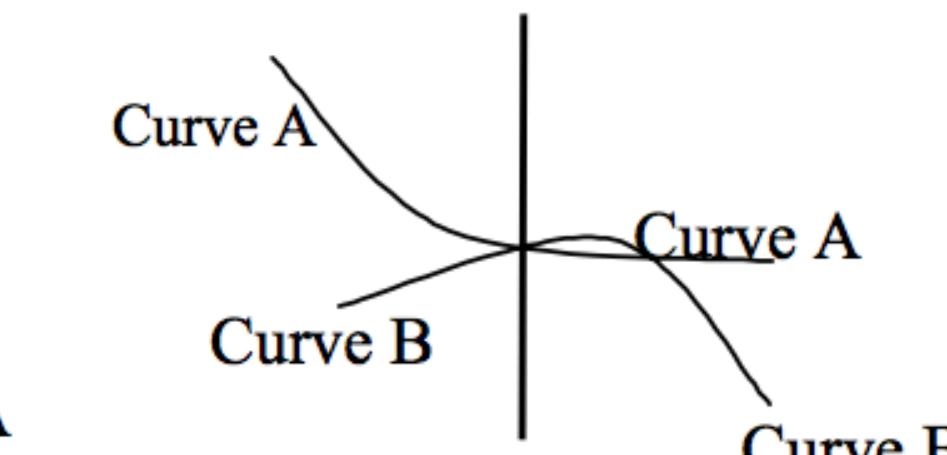
# PC scalability

To scale PC we can use:

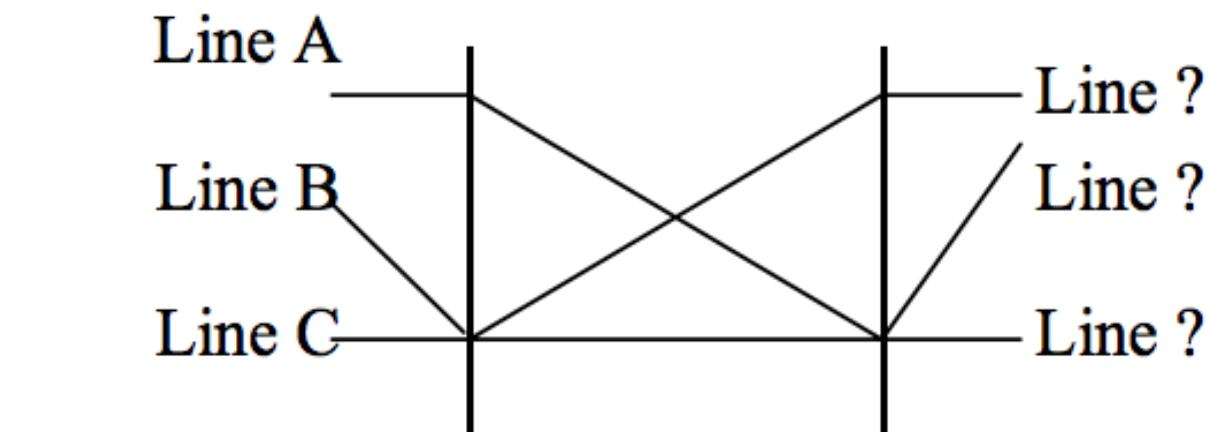
Transparency



Bundling



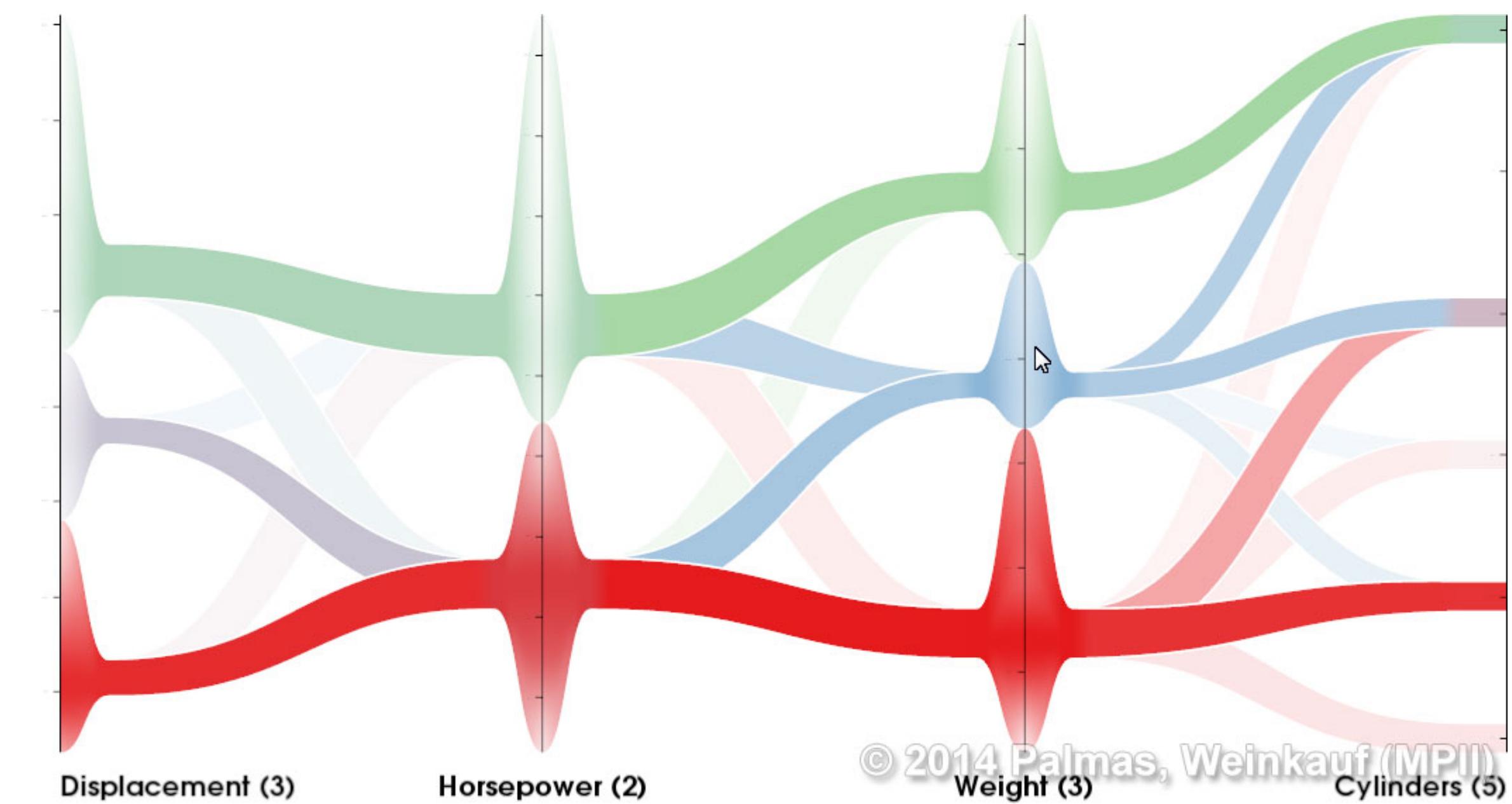
Clustering

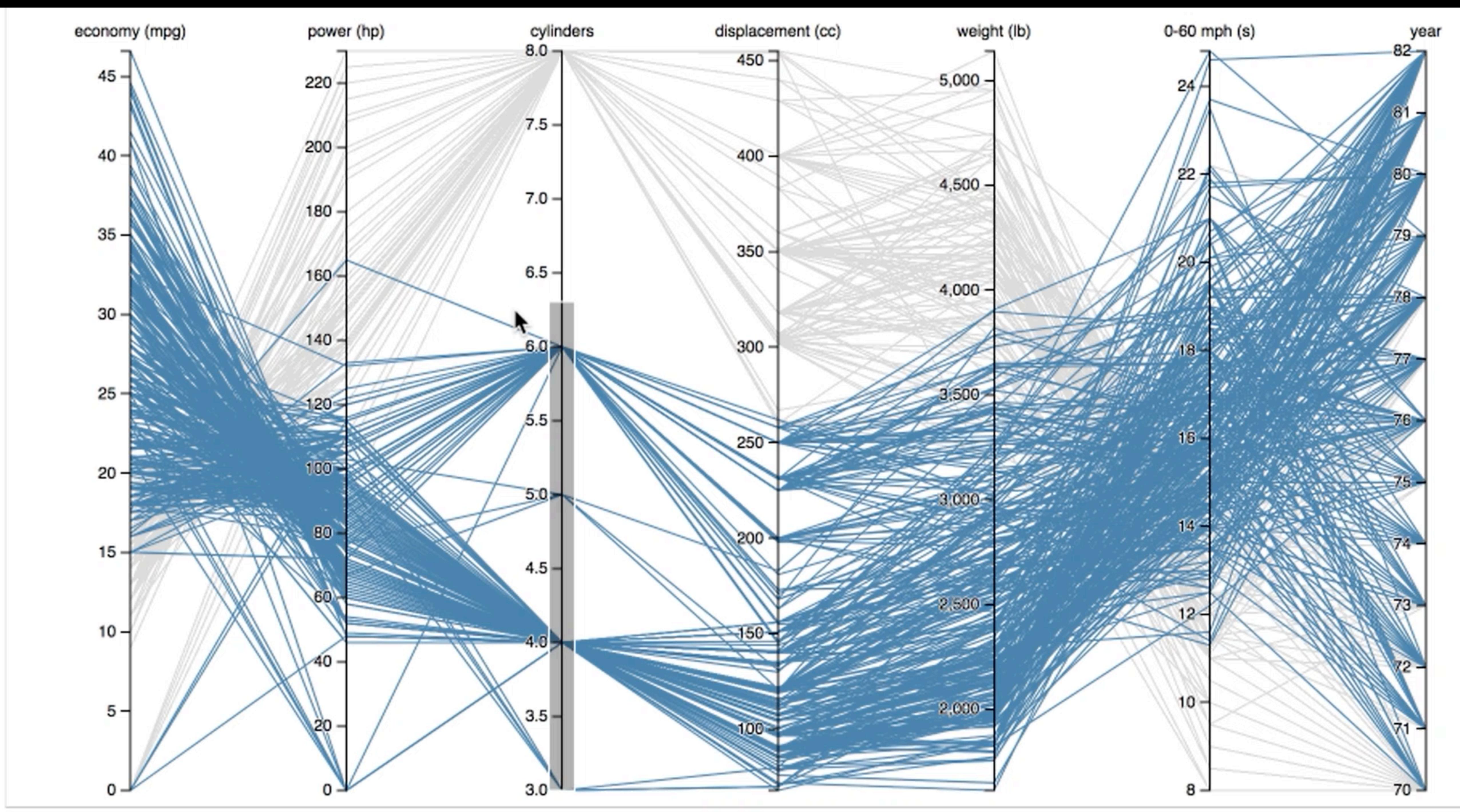


Sampling

Interaction: brushing,  
reorder axis

Filtering

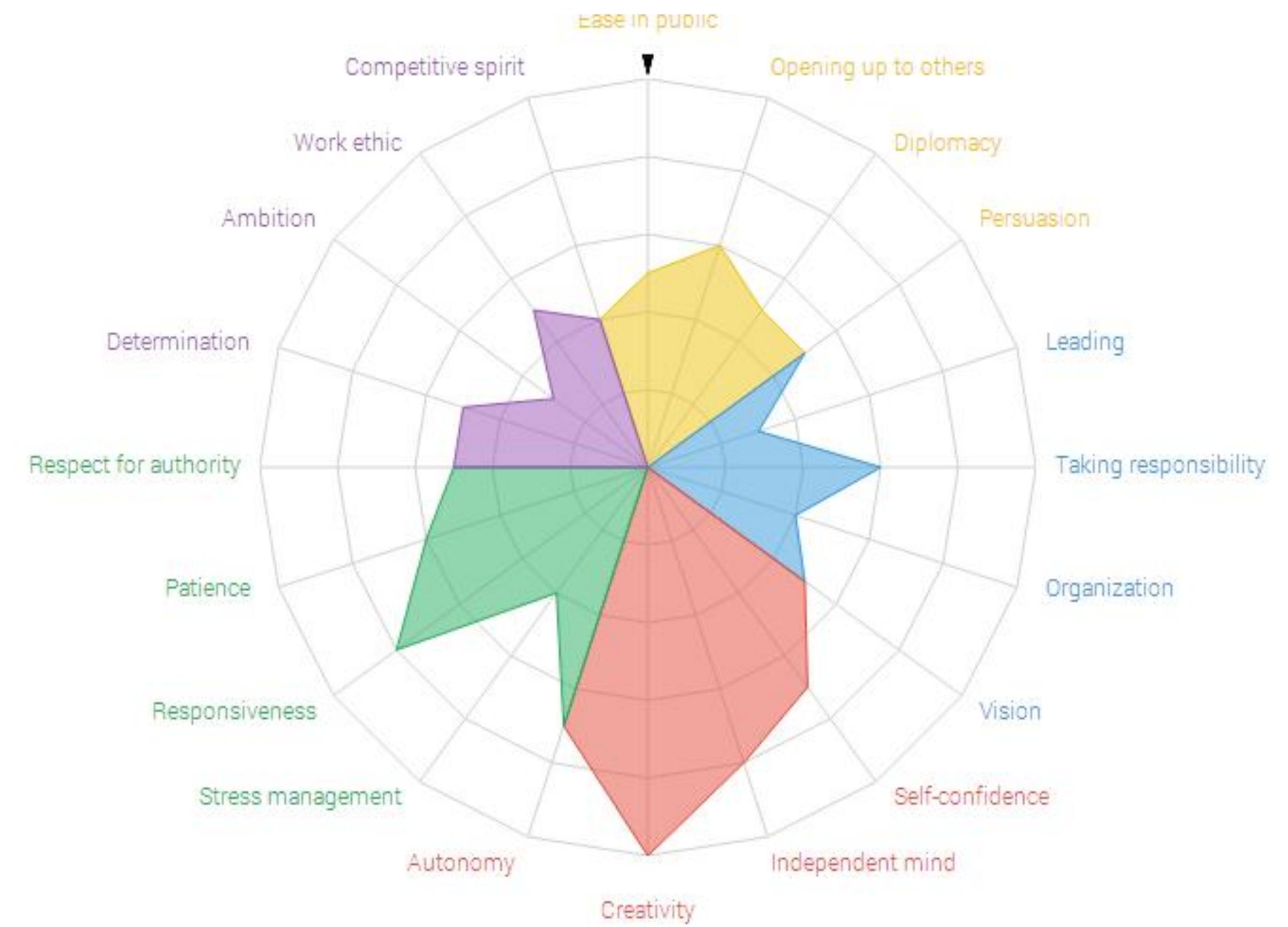
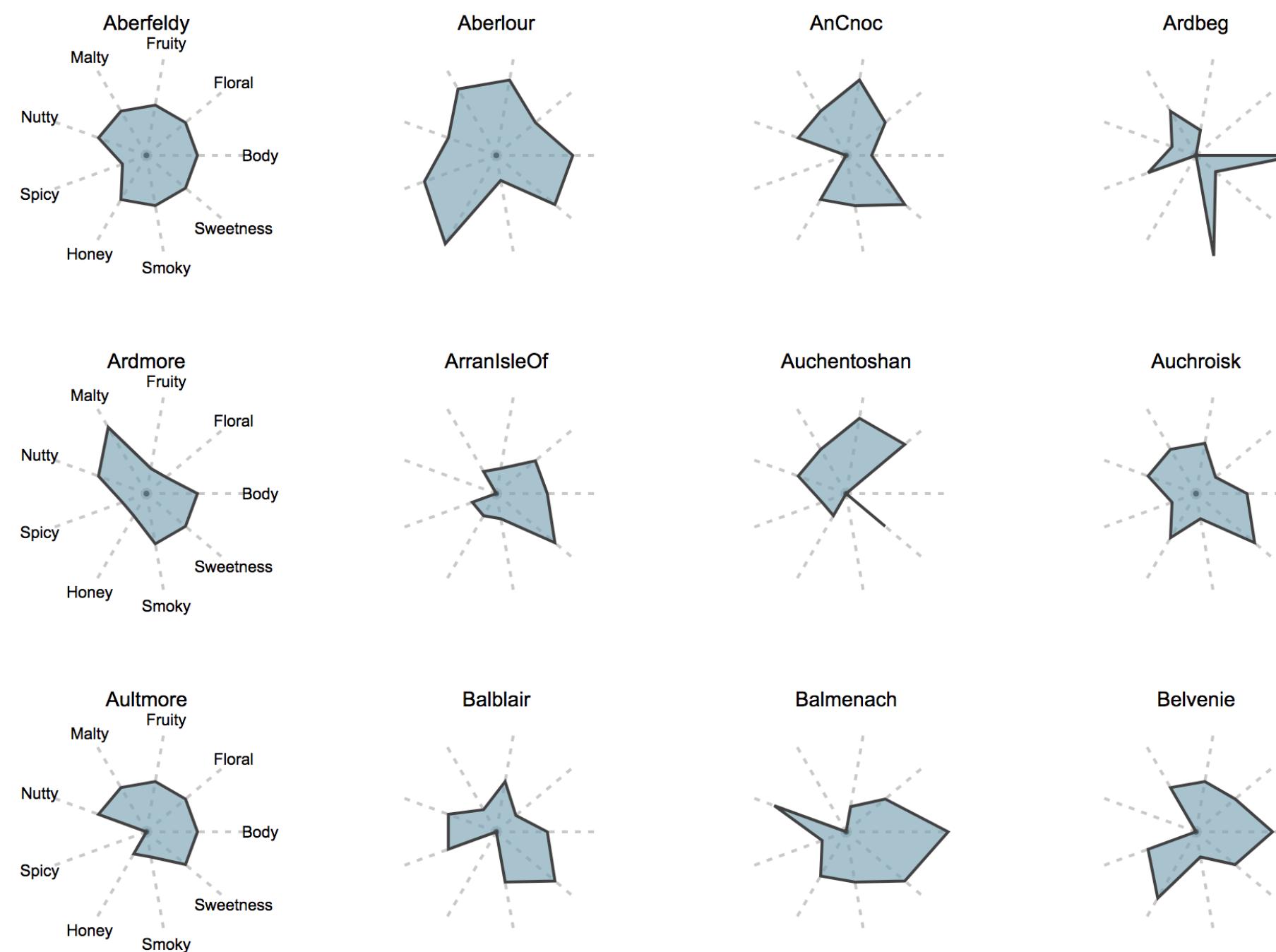




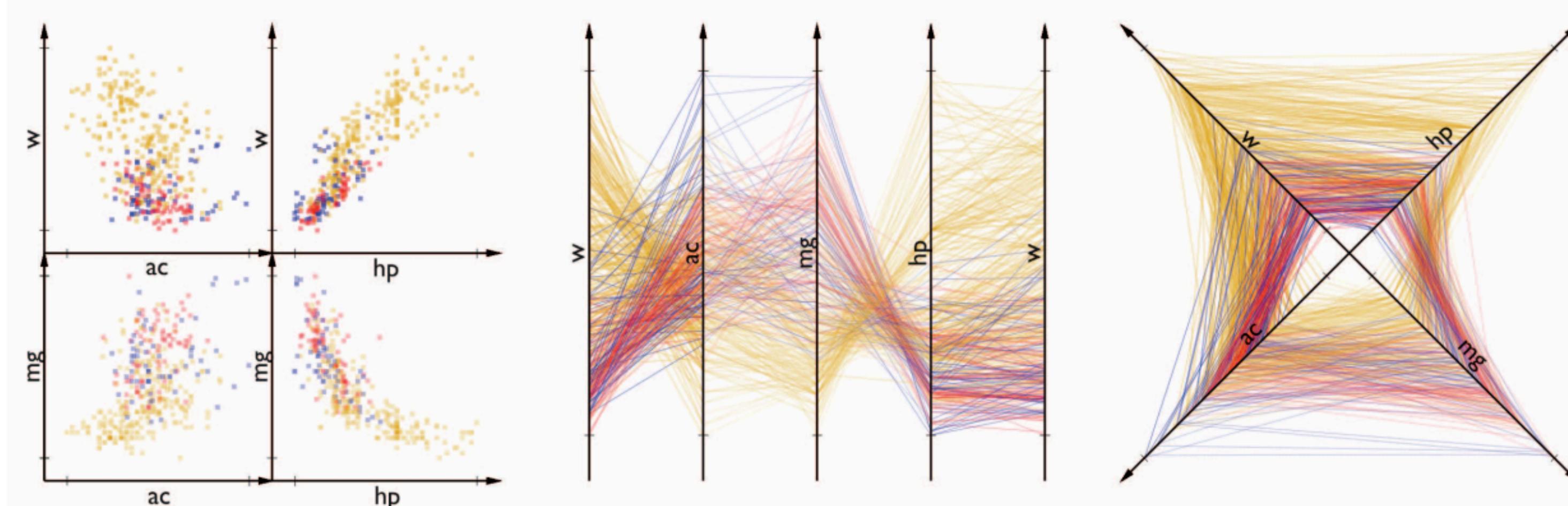
# Star plot

Similar to parallel coordinates

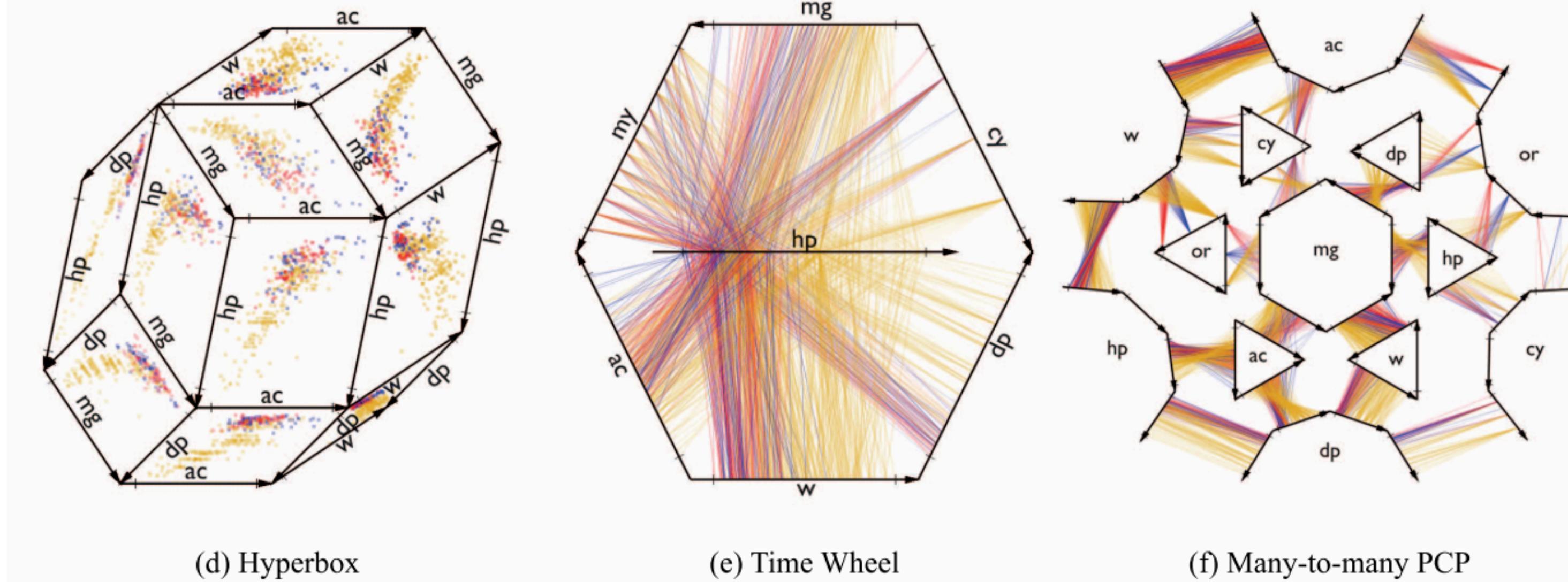
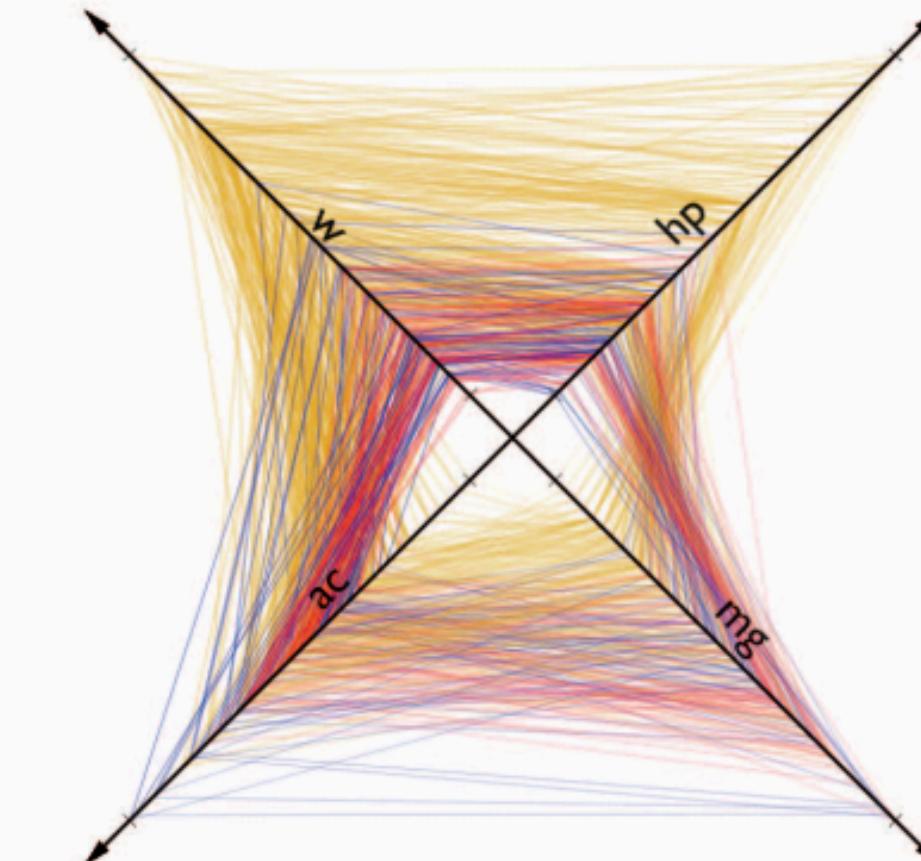
Radial layout from the center



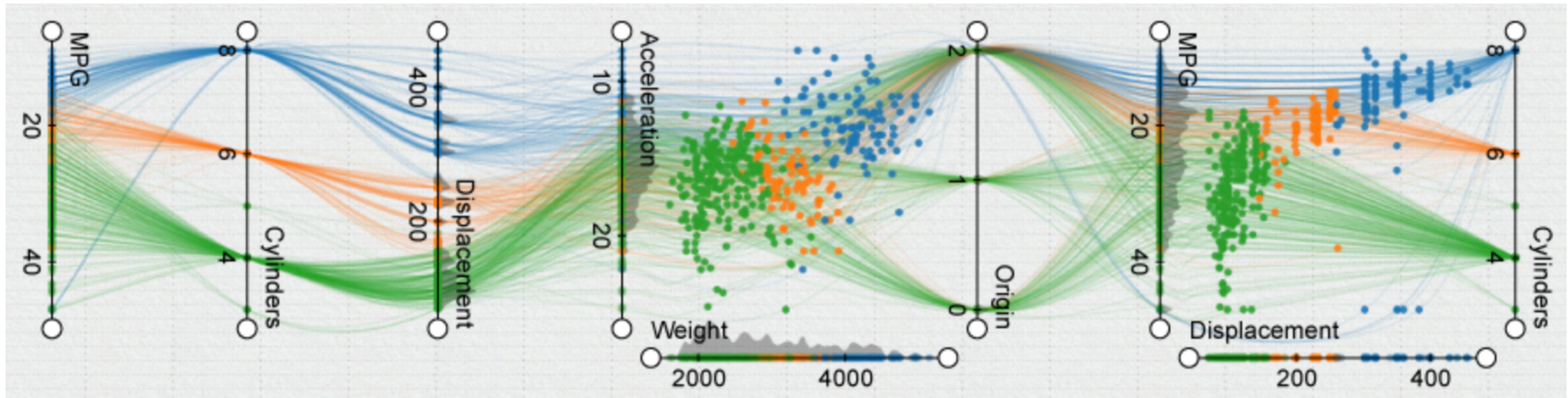
# Flexible Linked Axes (FLINA)



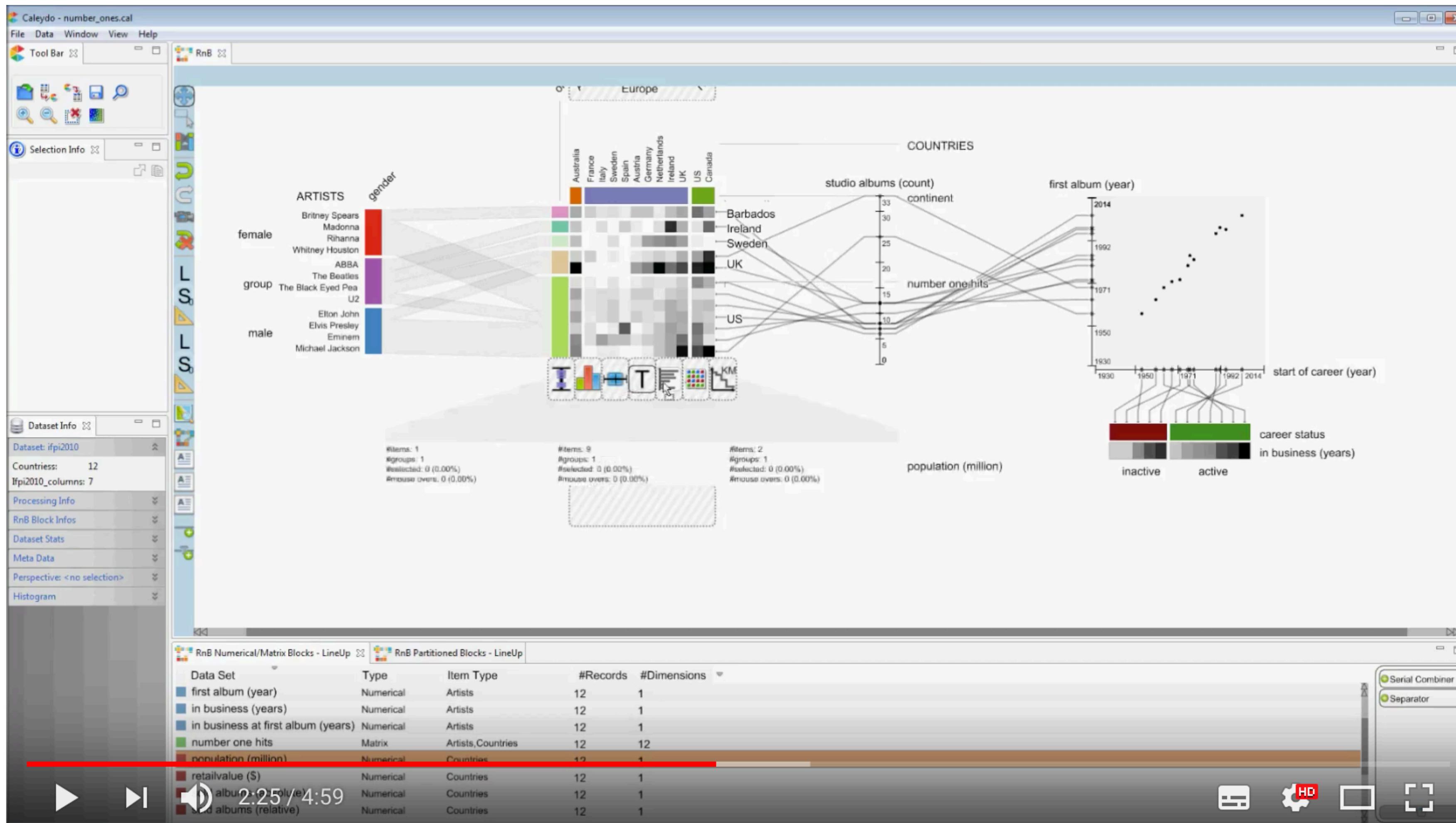
(b) Parallel Coordinates Plot



# Visualization Assembly Line



# Domino



# Parallel set

