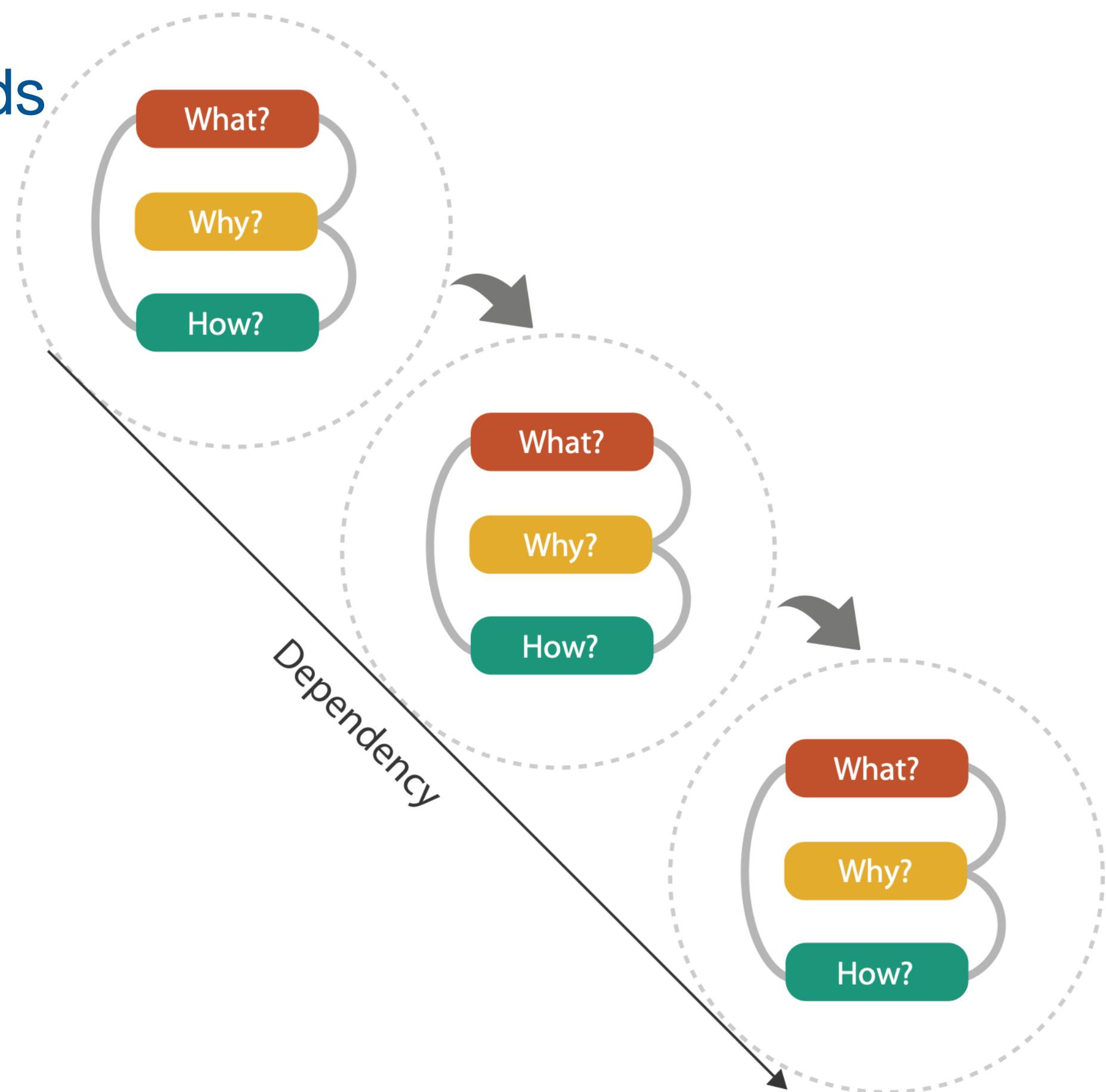


# Visualization for Data Science

## The Zoo II



# Means and ends



# What?

## Datasets

## Attributes

### → Data Types

→ Items    → Attributes    → Links    → Positions    → Grids

### → Attribute Types

→ Categorical



### → Data and Dataset Types

Tables	Networks & Trees	Fields	Geometry	Clusters, Sets, Lists
Items	Items (nodes)	Grids	Items	Clusters, Sets, Lists
Attributes	Links	Positions	Positions	Items

→ Ordered

→ Ordinal

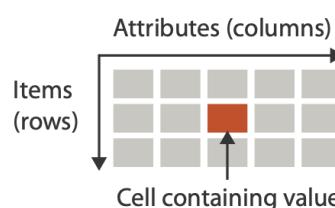


→ Quantitative

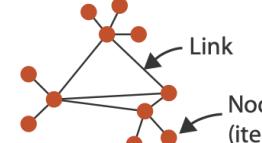


### → Dataset Types

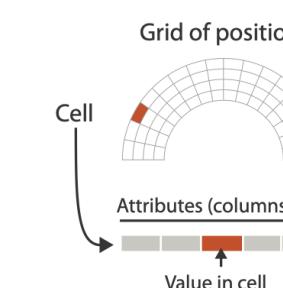
#### → Tables



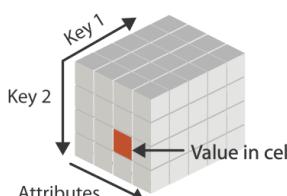
#### → Networks



#### → Fields (Continuous)



#### → Multidimensional Table



#### → Trees



#### → Geometry (Spatial)



### → Dataset Availability

#### → Static



#### → Dynamic



### → Ordering Direction

#### → Sequential

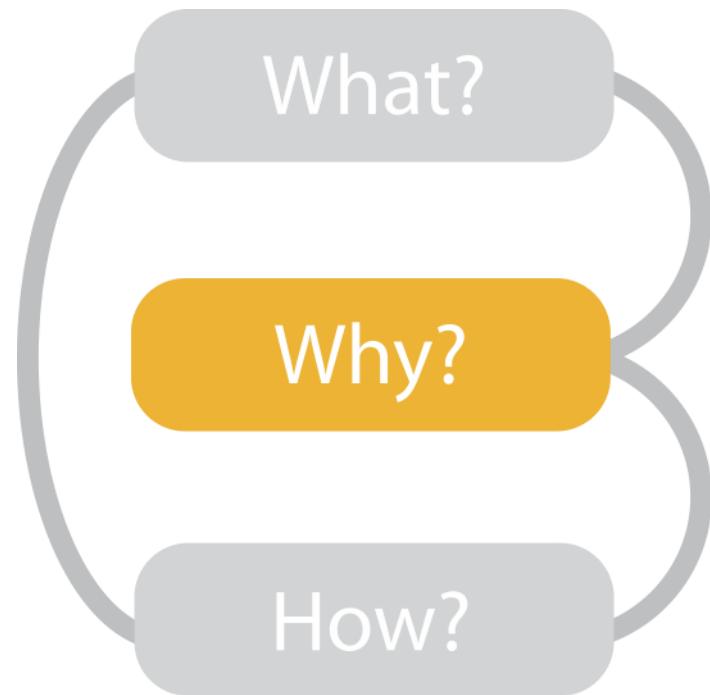


#### → Diverging

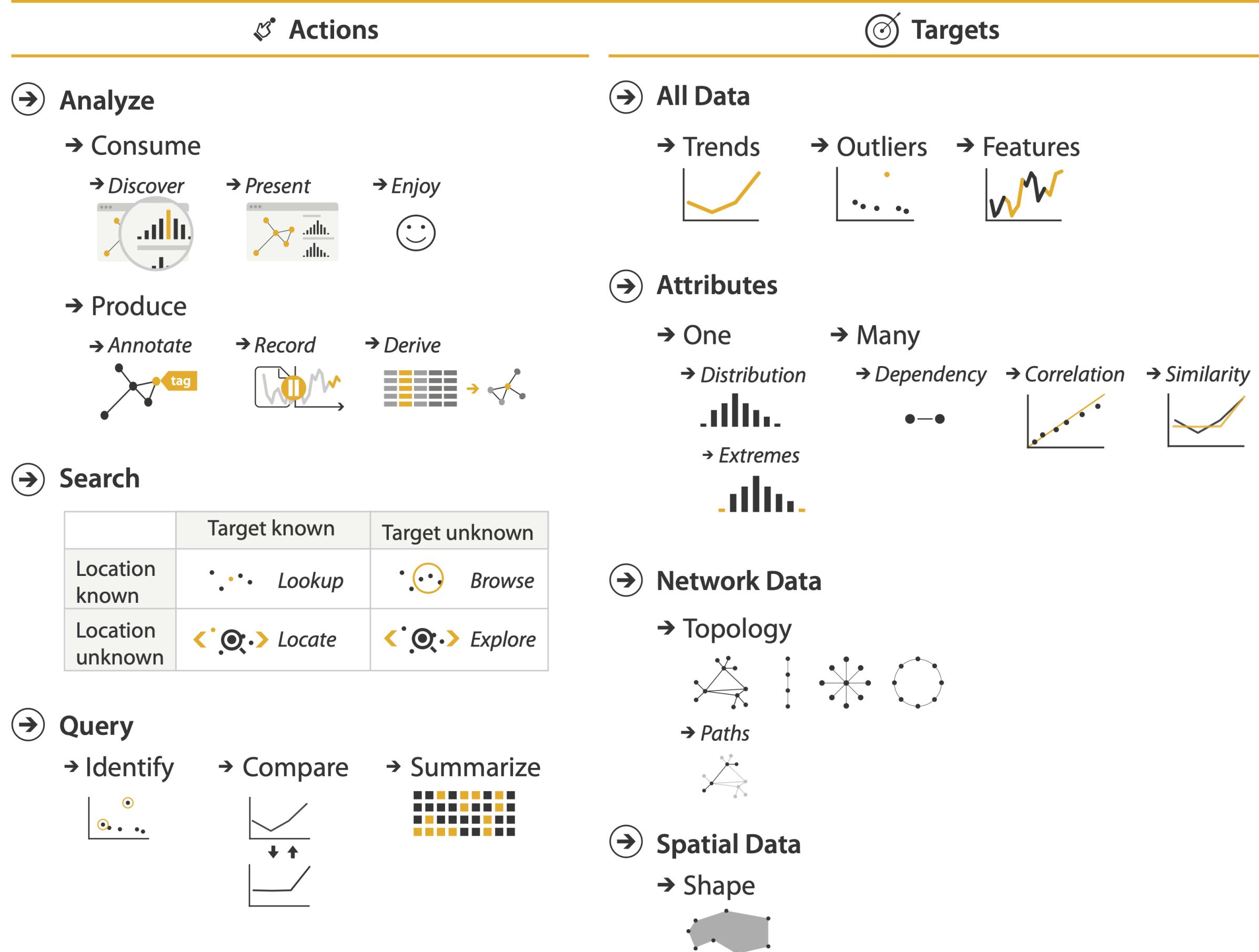


#### → Cyclic

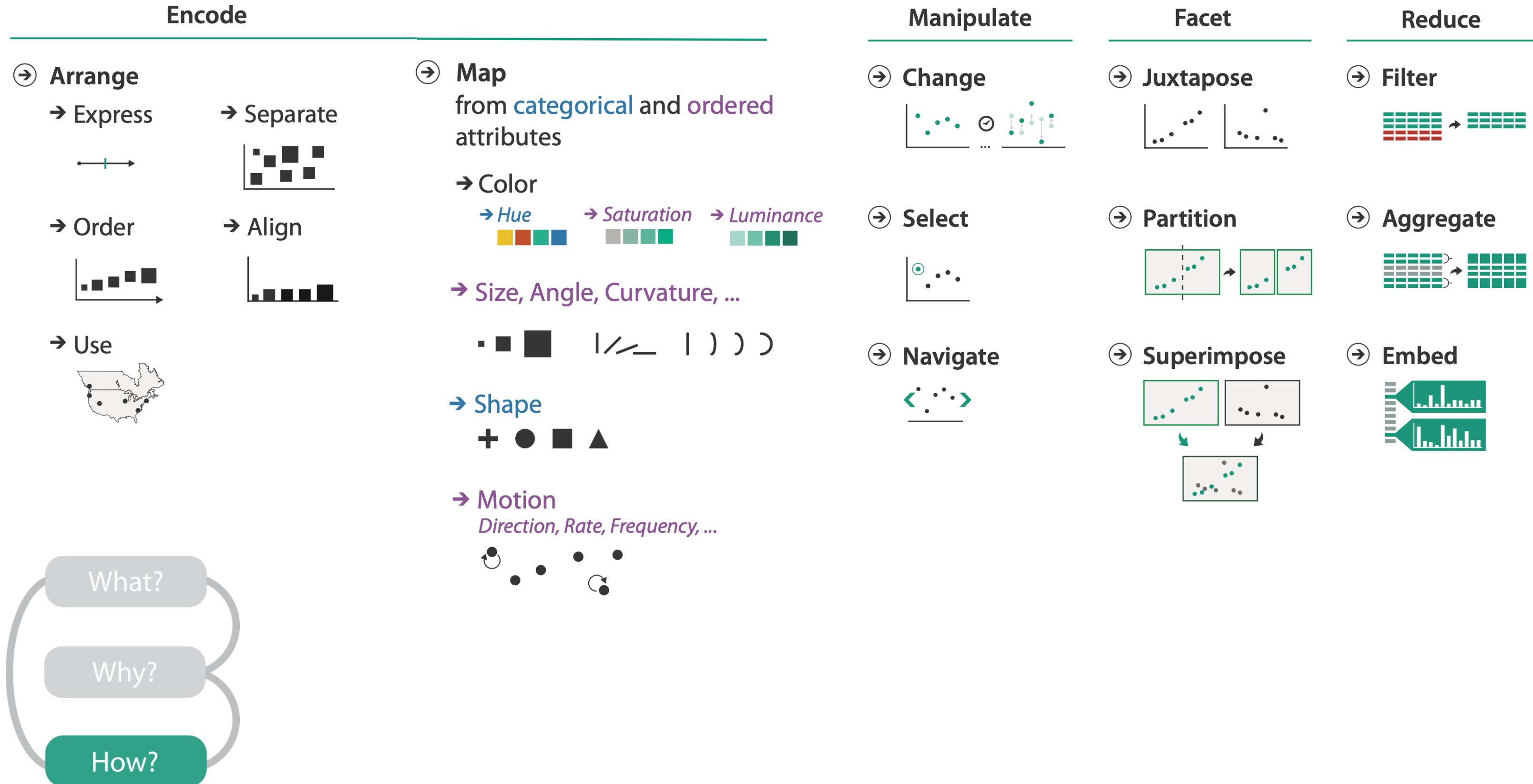




- {action, target} pairs
  - *discover distribution*
  - *compare trends*
  - *locate outliers*
  - *browse topology*



# How?



# Encode

---

→ Arrange

→ Express



→ Separate



→ Order

→ Align



→ Use



# Accuracy: User studies

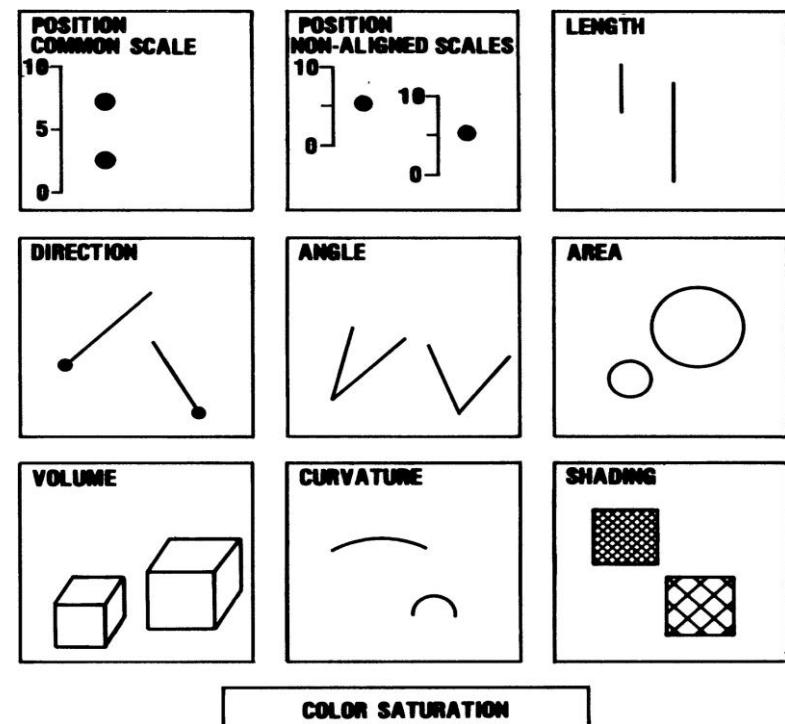


Figure 1. Elementary perceptual tasks.

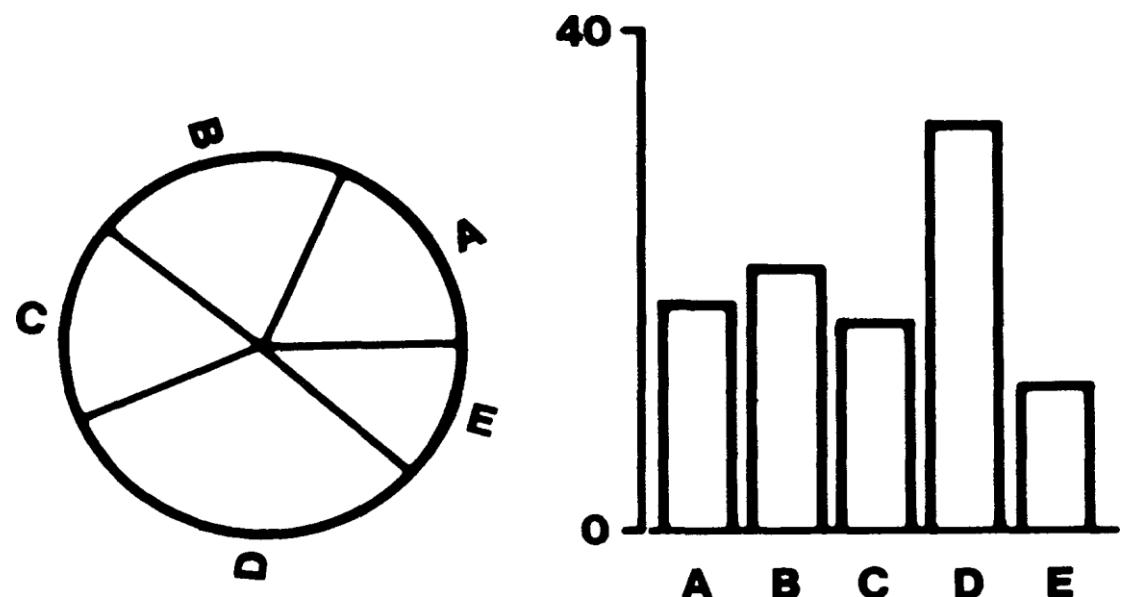


Figure 3. Graphs from position–angle experiment.

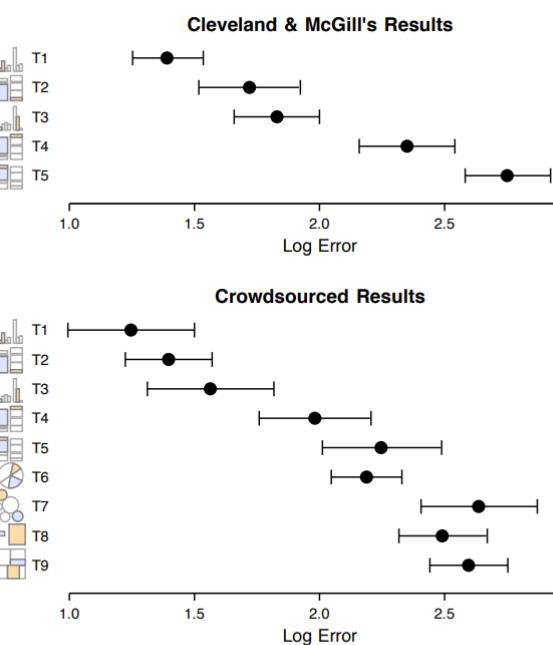
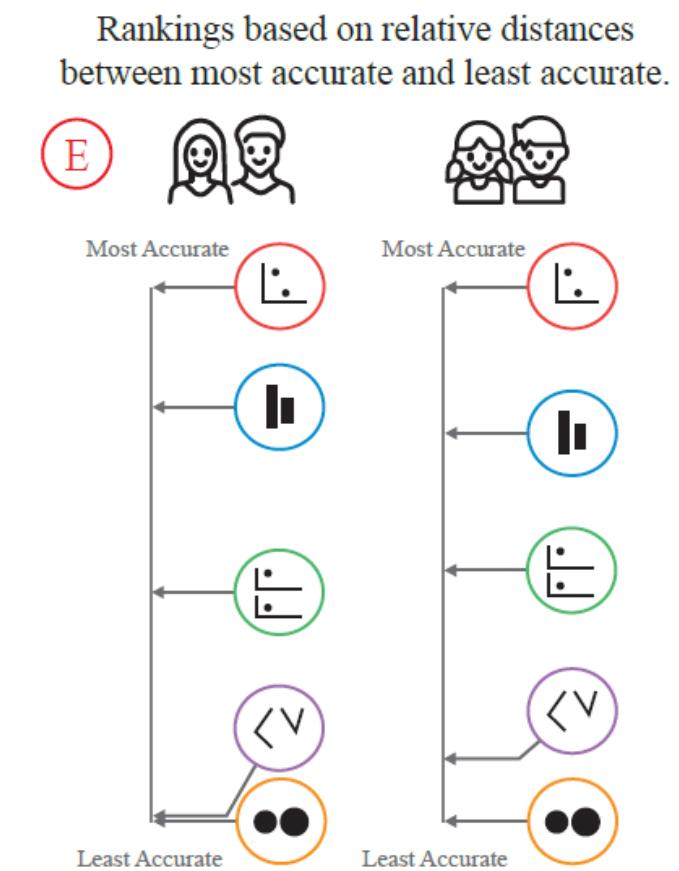
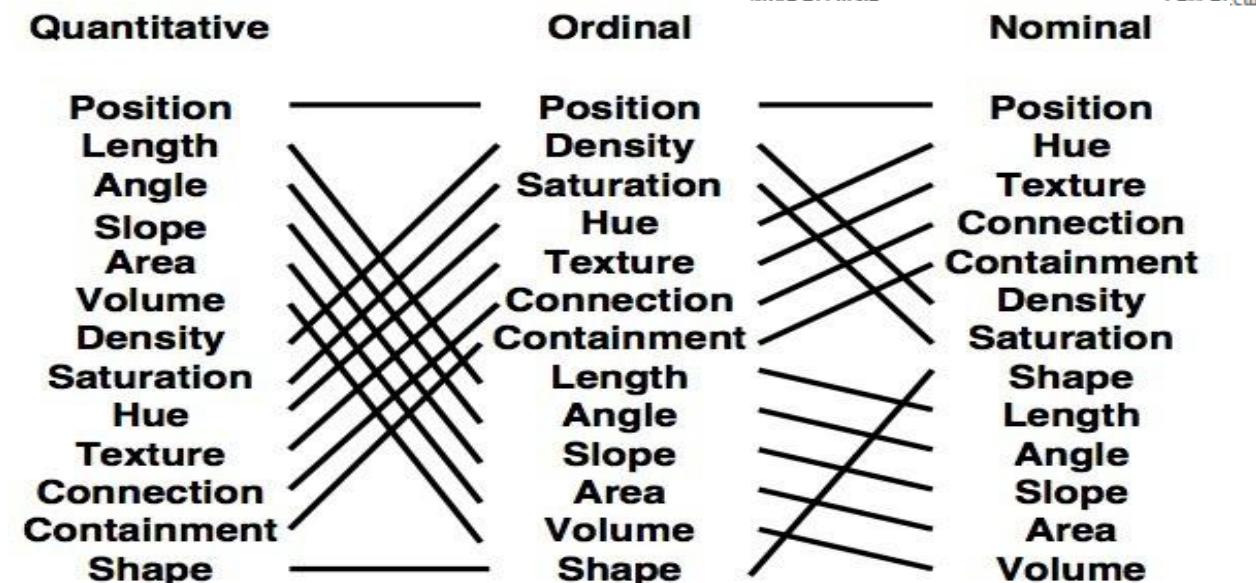


Figure 4: Proportional judgment results (Exp. 1A & B).  
Top: Cleveland & McGill's [7] lab study. Bottom: MTurk studies. Error bars indicate 95% confidence intervals.



[Cleveland & McGill, 1984](#)

[Mackinlay, 1986](#)

[Heer & Bostock, 2010](#)

[Panavas et al., 2022](#)

## → Express Values



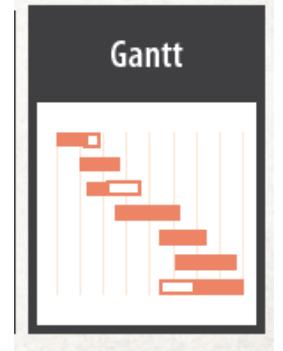
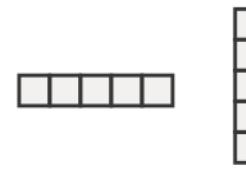
→ 0 Keys



x, y both  
quantity

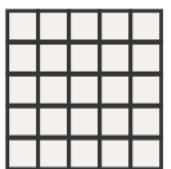
→ 1 Key

*List*



→ 2 Keys

*Matrix*



# Learning Outcomes

- Select a visualization based on the data and task
- Describe visualizations in terms of keys and values that they encode

# Idiom: **heatmap**

- two keys, one value

- data

- 2 categ attrs (gene, experimental condition)
    - 1 quant attrib (expression levels)

- marks: point (rect in Altair)

- separate and align in 2D matrix
      - indexed by 2 categorical attributes

- channels

- color by quant attrib
      - (ordered diverging colormap)

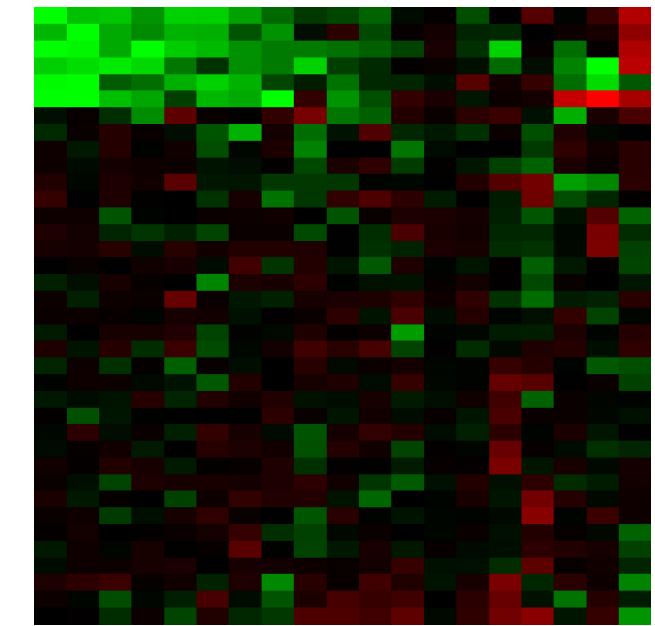
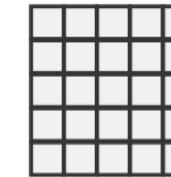
- task

- find clusters, outliers

- scalability

- 1M items, 100s of categ levels, ~10 quant attrib levels

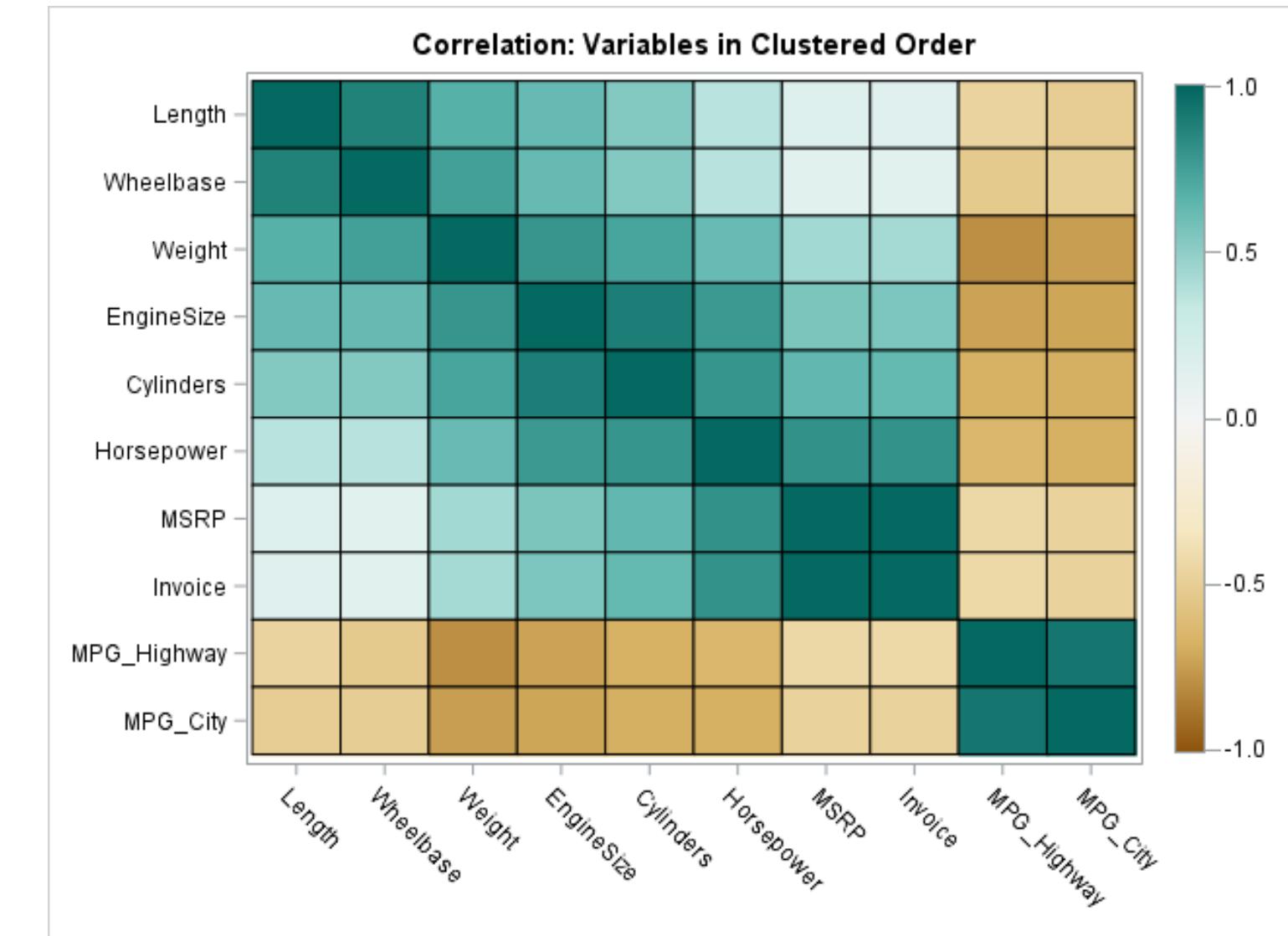
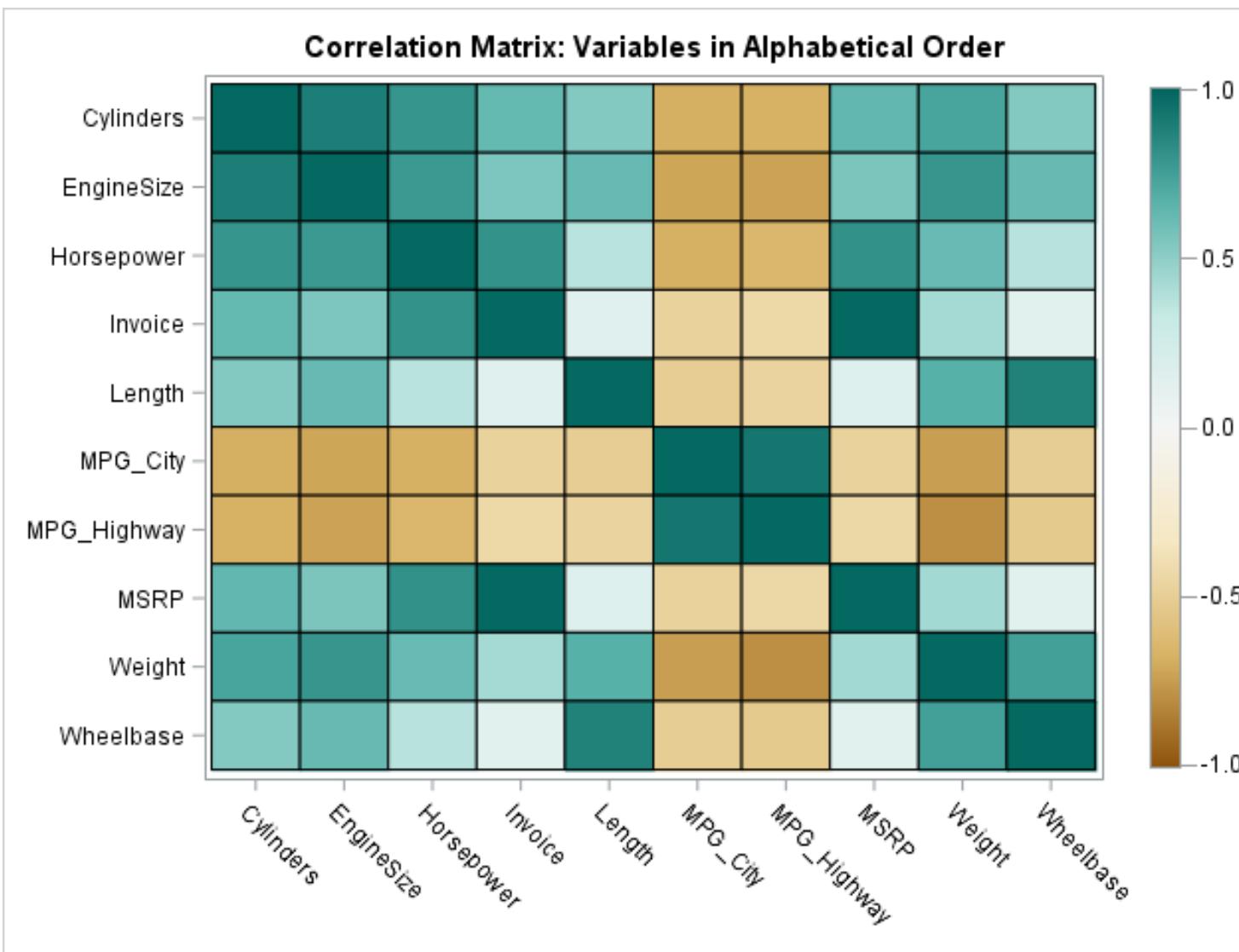
→ 2 Keys  
*Matrix*



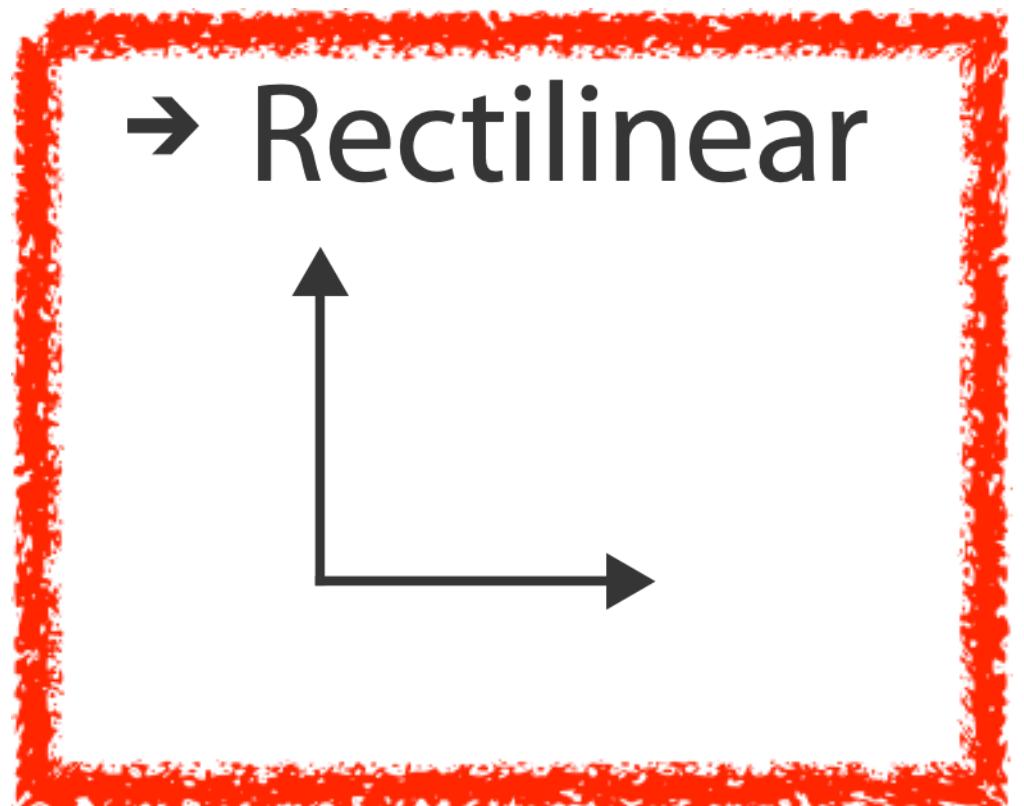
How many shades of red can you see

- A. 2
- B. 3
- C. 4
- D. 5
- E. > 5

# Heatmap reordering

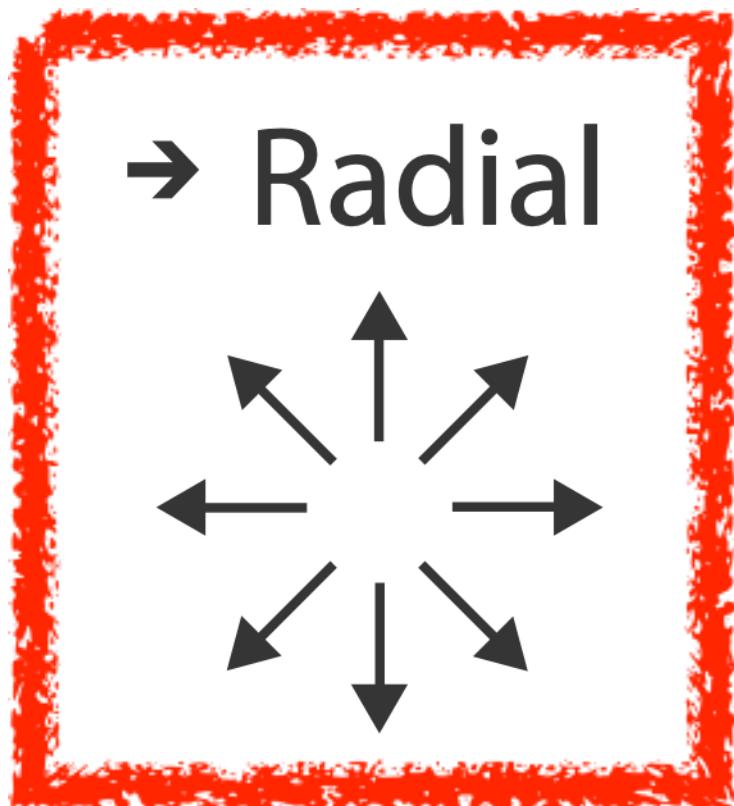


## → Axis Orientation



→ Rectilinear

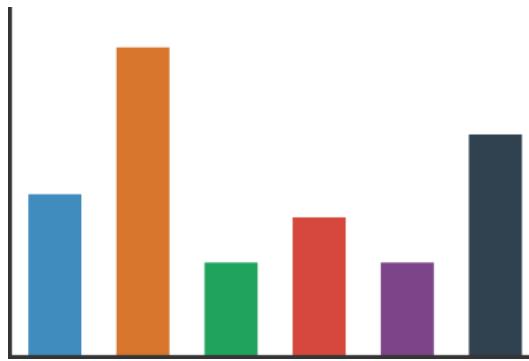
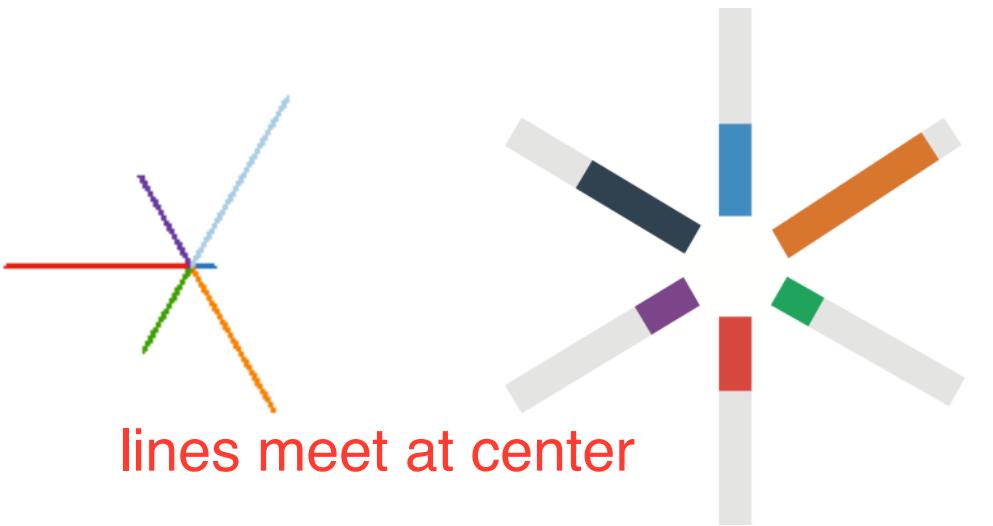
→ Parallel



→ Radial

# Idioms: **radial bar chart, star plot**

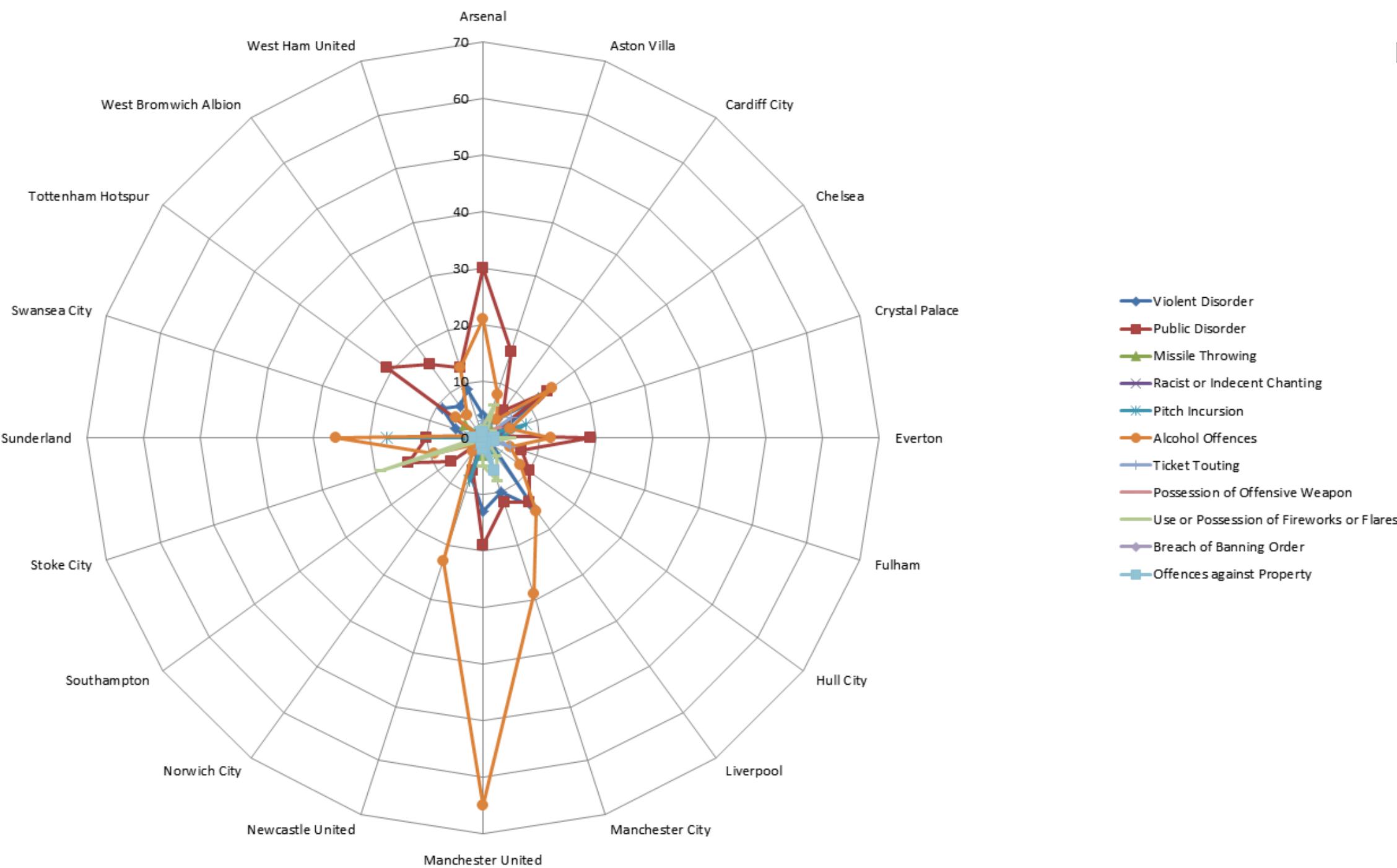
- star plot
  - line mark, radial axes meet at central point
- radial bar chart
  - line mark, radial axes meet at central ring
  - channels: **length**, angle/orientation
- bar chart
  - rectilinear axes, aligned vertically
- accuracy
  - length not aligned with radial layouts
    - less accurately perceived than rectilinear aligned



# Idiom: **radar plot**

- radial line chart
  - point marks,  
radial layout
  - connecting line  
marks
- avoid unless data  
is cyclic

no distinct ordering for the  
category on the circle



shape depicts the meaning

not able to see all data point

# “Radar graphs: Avoid them (99.9% of the time)”



## Os sinais da bússola eleitoral

Disputa de 2010 foi parecida com a de 2006

Alberto Cairo, Alexandre Matos, Carlos Eduardo Cruz-Garcia, Eliseu Barreira Junior, Marco Vergolfi e Ricardo Mendoza

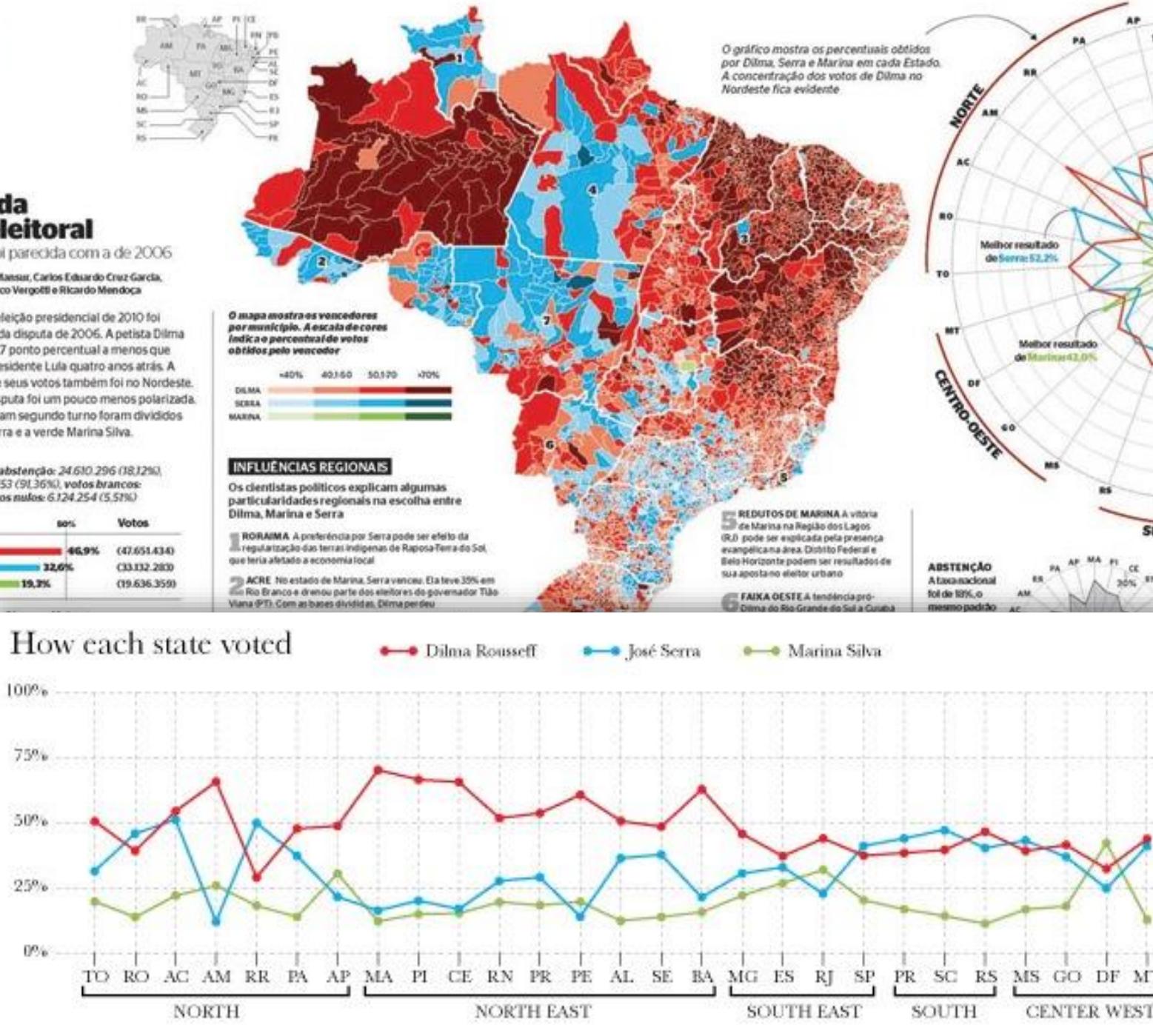
O PRIMEIRO TURNO da eleição presidencial de 2010 foi muito parecido com o da disputa de 2006. A petista Dilma Rousseff teve apenas 1,7 ponto percentual a menos que o índice obtido pelo presidente Lula quatro anos atrás. A concentração maior de seus votos também foi no Nordeste. Dessa vez, porém, a disputa foi um pouco menos polarizada. Os votos que provocaram segundo turno foram divididos entre o tucano José Serra e a verde Marina Silva.

Eleitores: 135.804.433, abstenção: 24.610.296 (18,12%), votos válidos: 110.590.153 (91,36%), votos brancos: 3.479.340 (3,13%) e votos nulos: 6.124.254 (5,51%).

Candidatos	sos	Votos
Dilma Rousseff	46,9%	(47.651.434)
José Serra	32,6%	(33.332.280)
Marina Silva	19,3%	(19.636.359)

### Outros candidatos

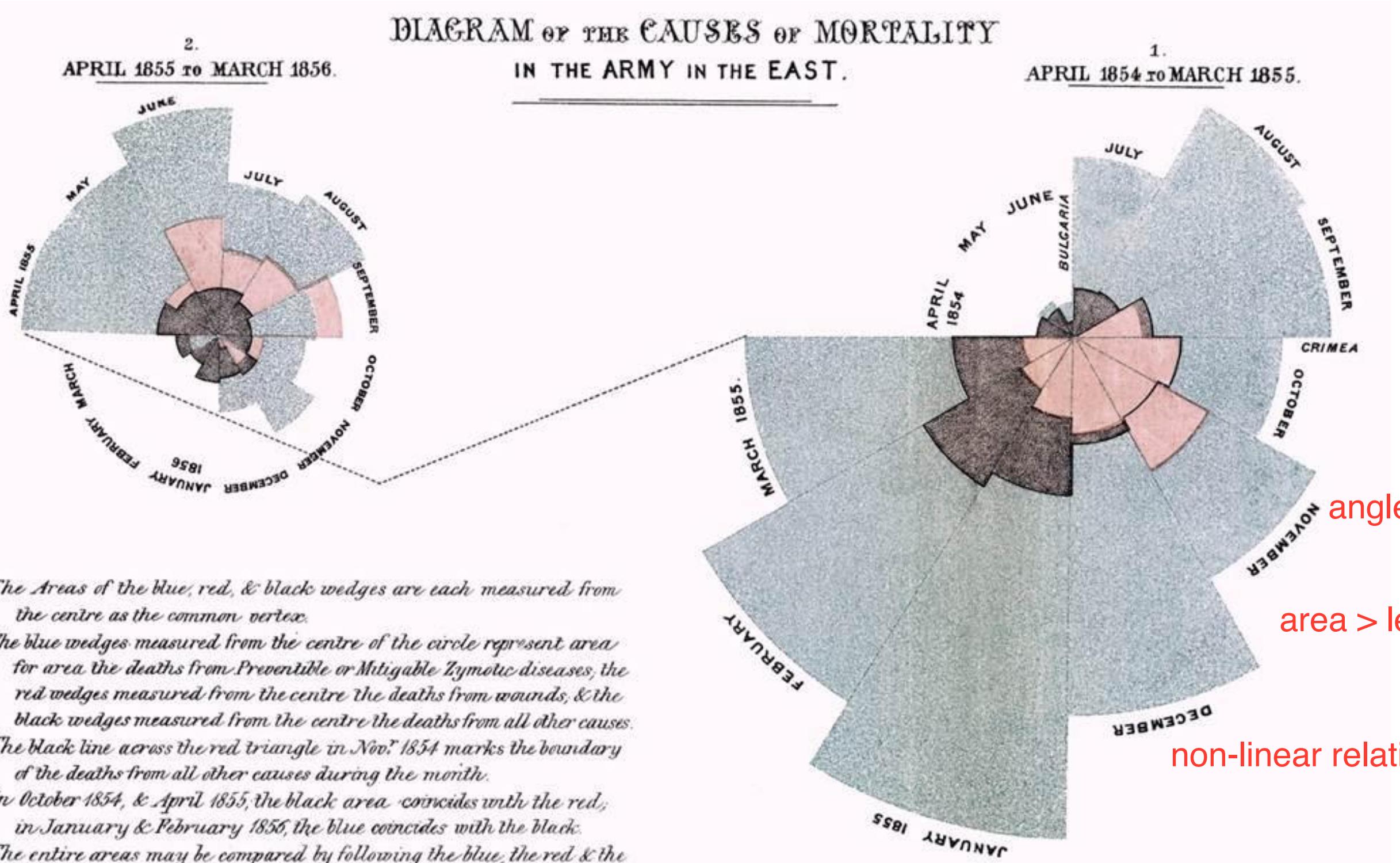
Pônio (PSOL)  
José Maria Eymael (DCS)  
Zé Maria (PSTB)  
Levy Fidelix (PRTB)  
Ivan Pinheiro (PCB)  
Rui Costa Pimenta (PCO)



original  
difficult to interpret

redesign for  
rectilinear

# Coxcomb / nightingale rose / polar area chart



angles are the same

area > length

non-linear relationship

# Florence Nightingale

<https://www.youtube.com/watch?v=u6XqiDccroM>

As we watch the video pay attention to the following

1. The Structure
2. The data represented
3. The case for visualization that Florence Nightingale made
4. The case for sanitation

# Idioms: **pie chart, coxcomb chart**

length (radius) is the same, angles are different

- pie chart

- interlocking area marks with angle channel: 2D area varies

- separated & ordered radially, uniform height

- accuracy: area *less accurate* than rectilinear aligned line length

- task: part-to-whole judgements

- coxcomb chart

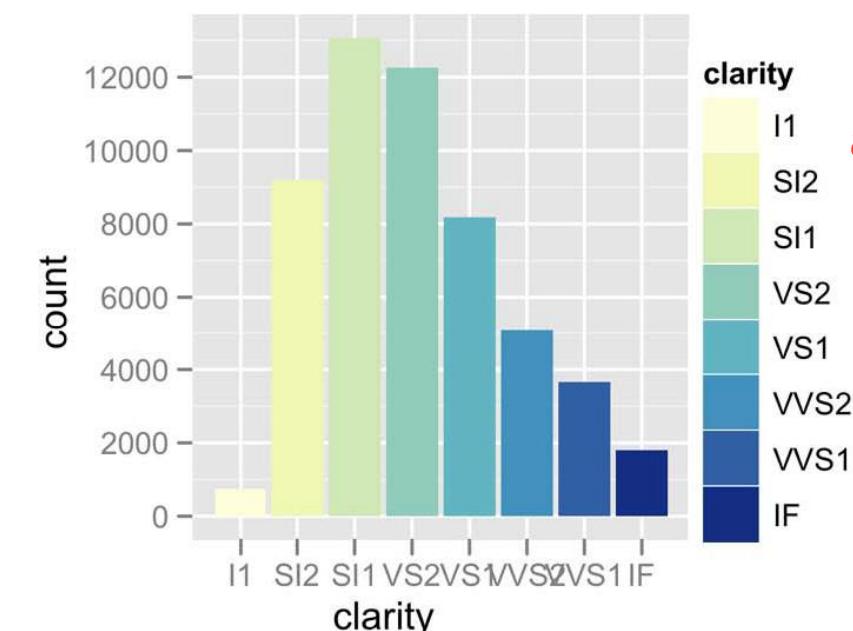
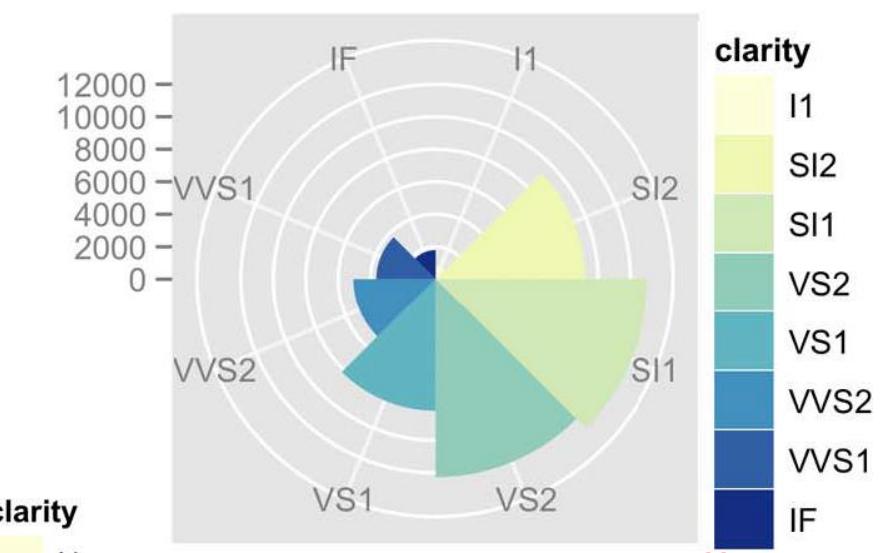
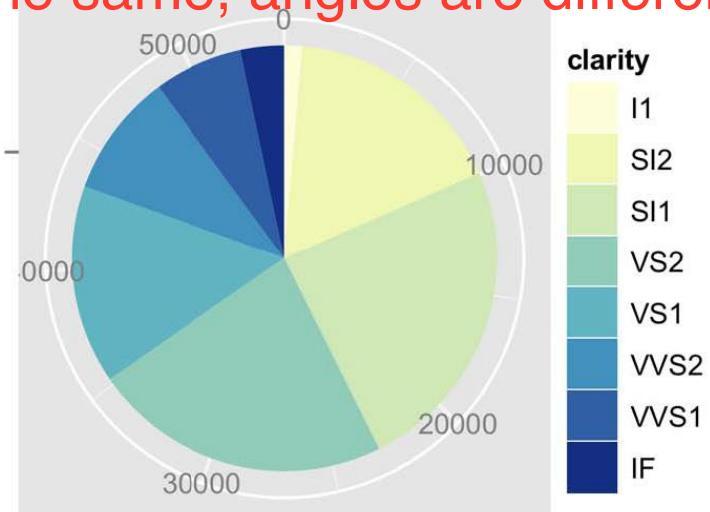
- line marks with length channel: 1D length varies

- separated & ordered radially, uniform width

- direct analog to radial bar charts

- data

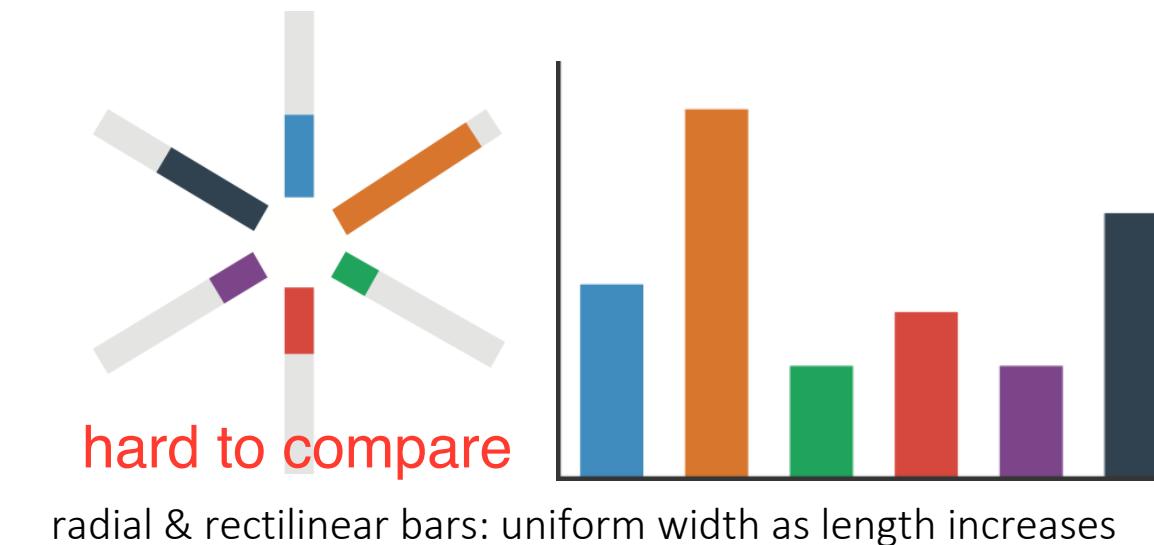
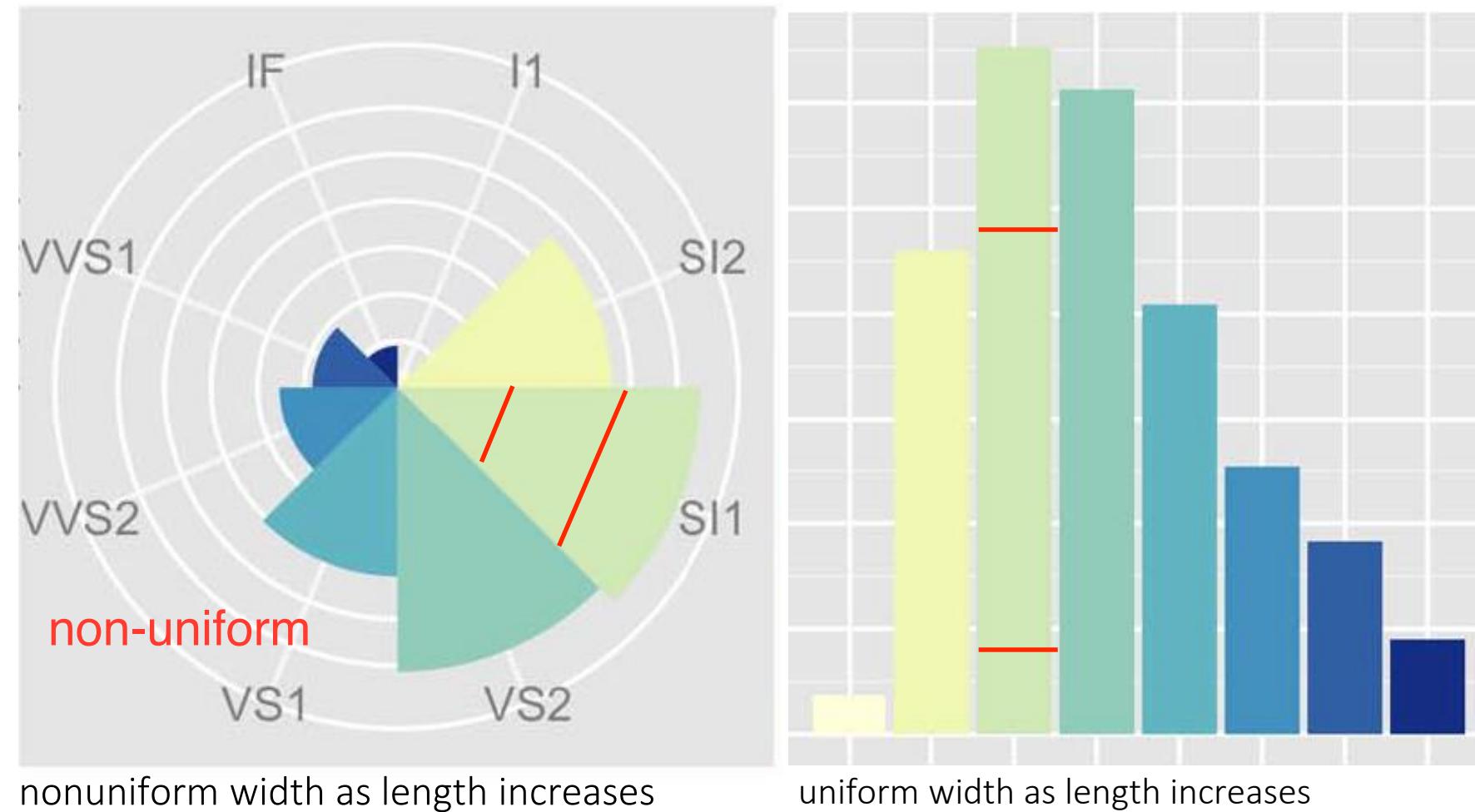
- 1 categ key attrib, 1 quant value attrib



angle same, length different

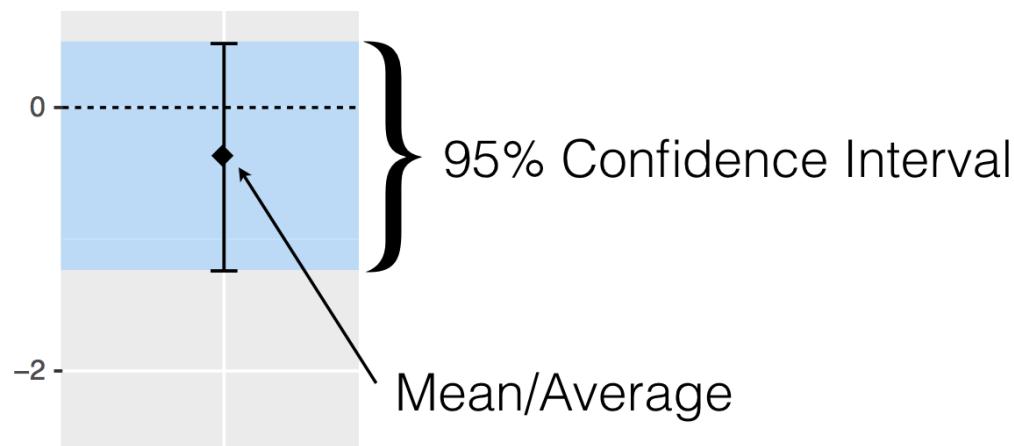
# Coxcomb: perception

- encode: **1D length**
- decode/perceive: **2D area**
- nonuniform line/sector width as length increases
  - so area variation is nonlinear wrt line mark length!
- bar chart safer: uniform width, so area is linear with line mark length
  - both radial & rectilinear cases

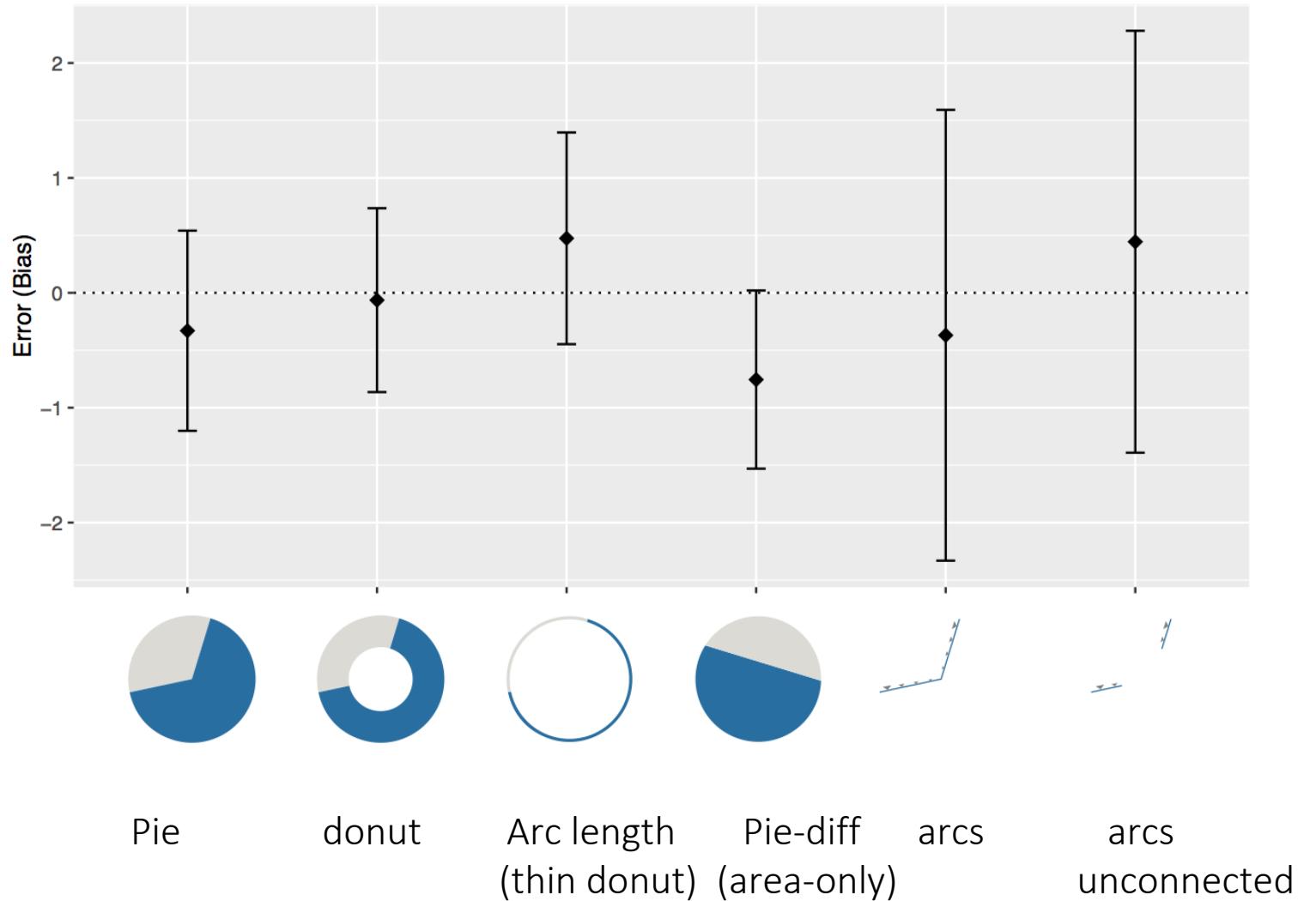


# Pie charts: perception

- some empirical evidence that people respond to arc length
- what we encode is different than what people perceive, most people do not read pie charts by angle
  - It could be areas, arc length alone, or combination
- donut charts no worse than pie charts and this is a direct result from the fact that most people don't read it by angle

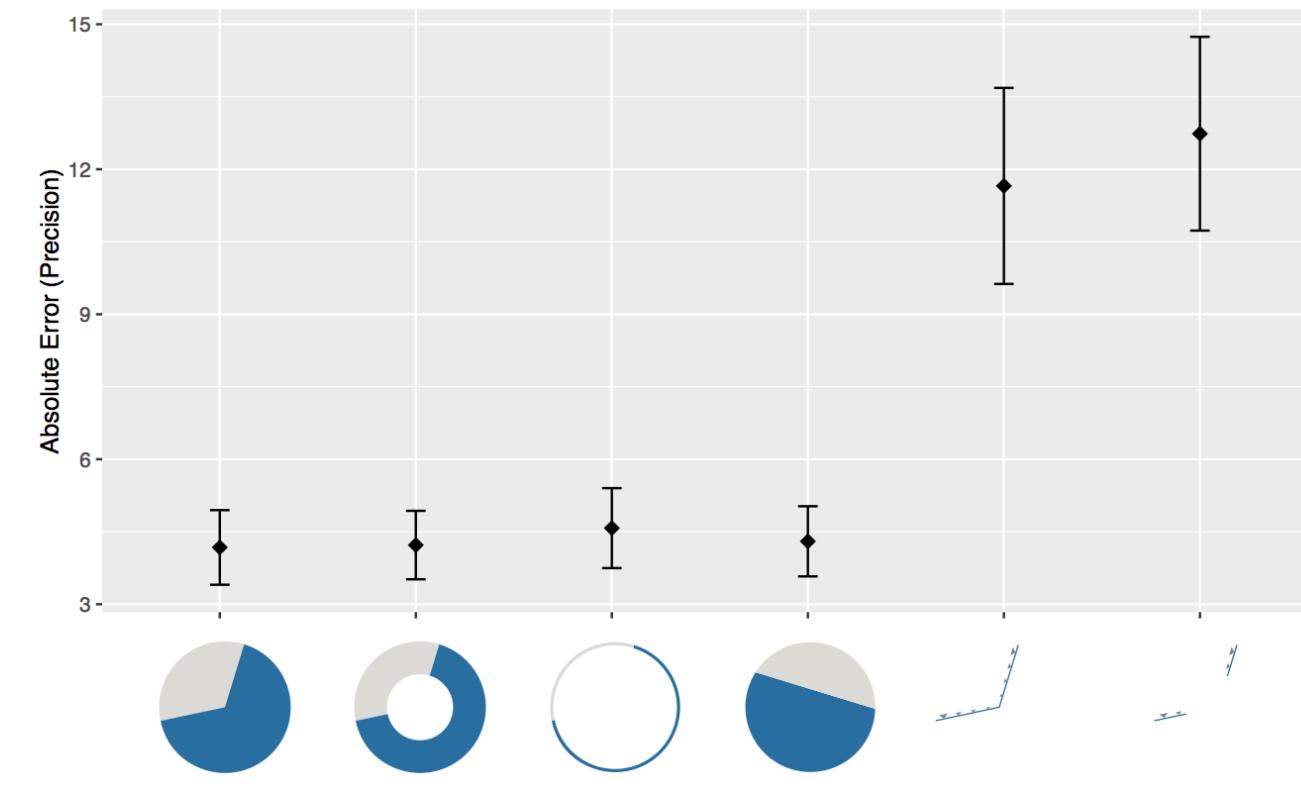


# Pie charts: perception



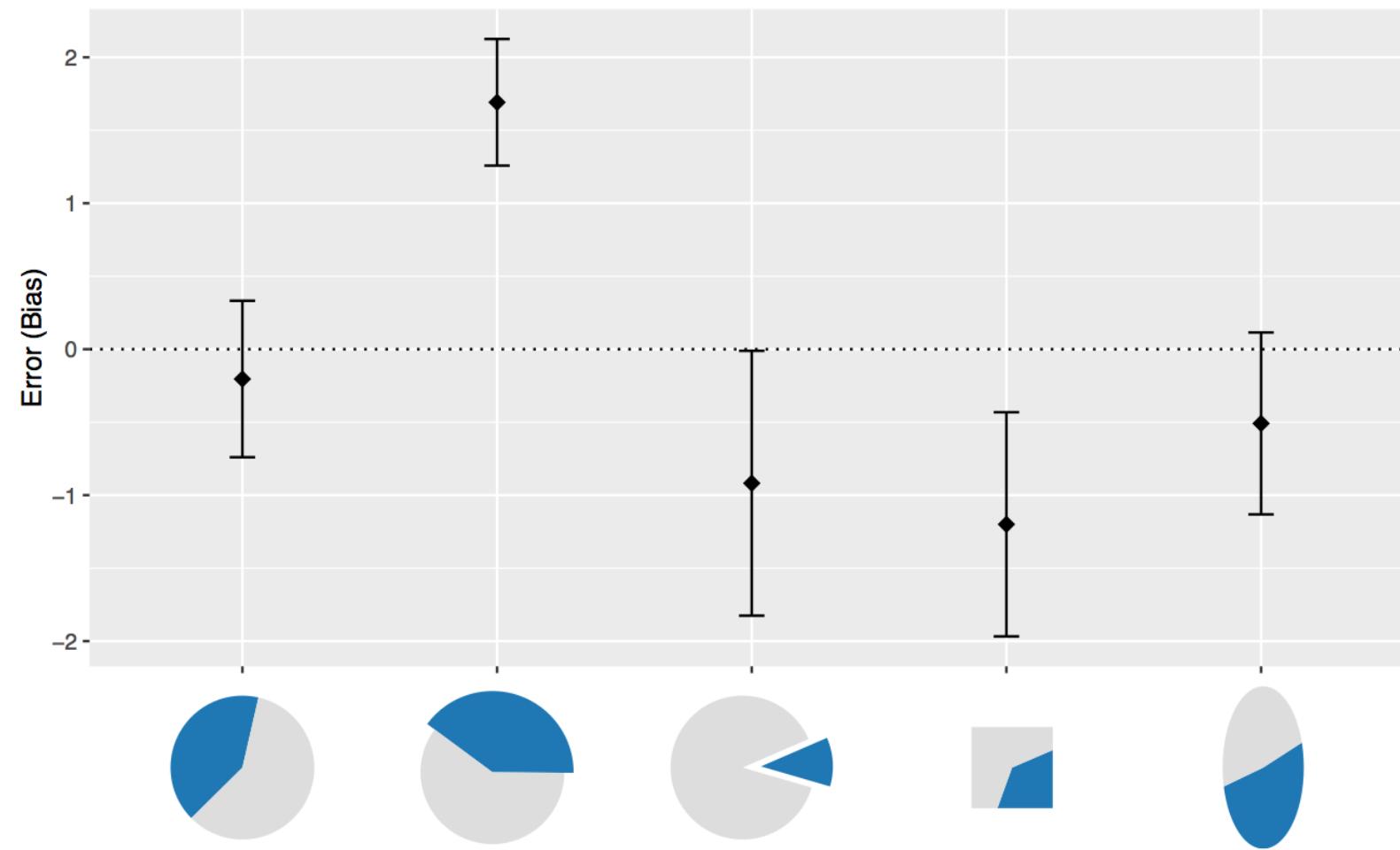
Each of these bars shows us the signed error. That means we can get a sense of the deviation from the real value (the lengths of the bars), but especially the bias: were people systematically over- or underestimating?

arc length/ angle/ area

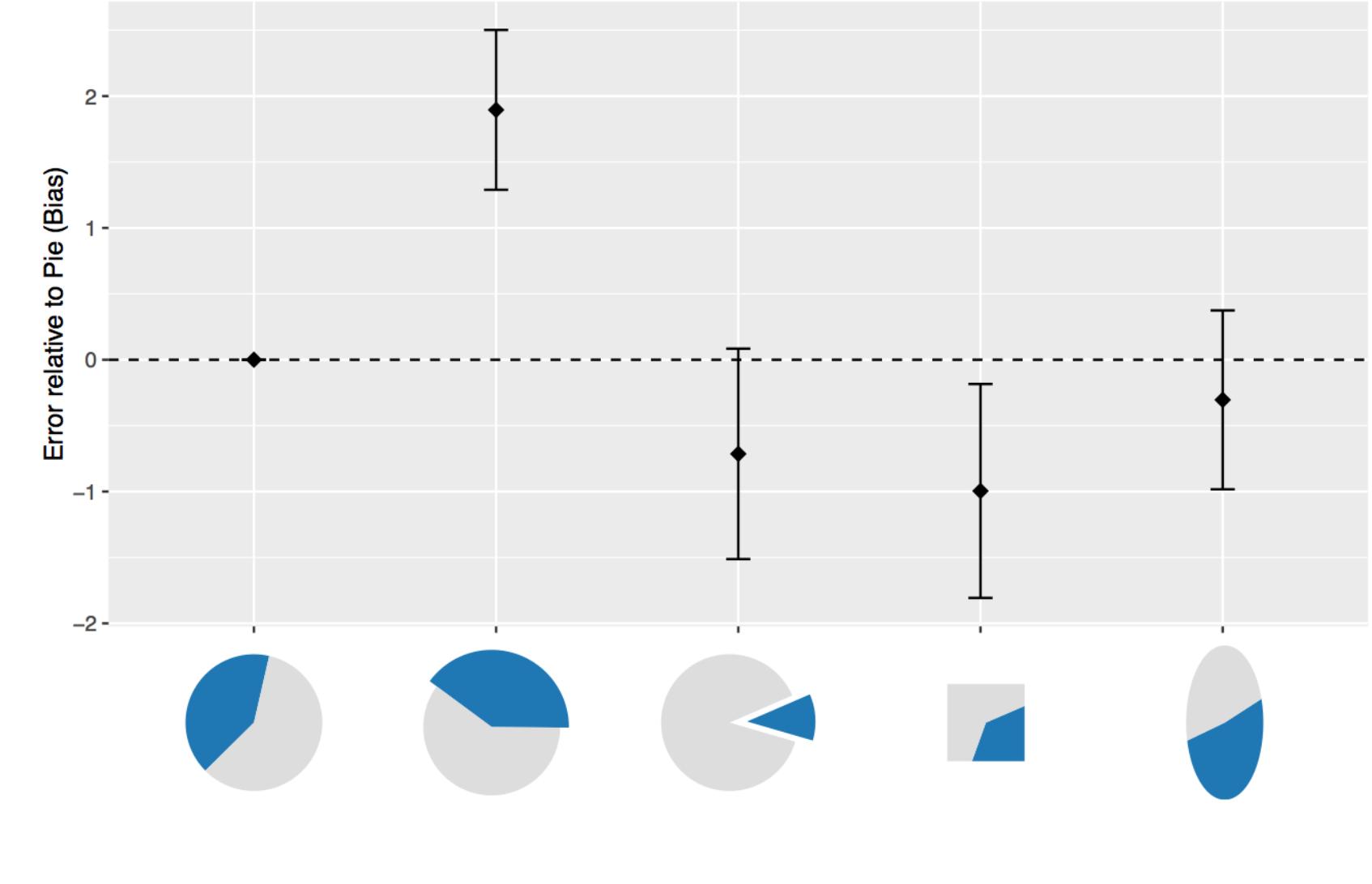


More interesting than signed error (where we take the average of all the errors, both above and below) is absolute error, where we count all deviations in the same direction. That means errors don't even out, so we can see how far off people are, no matter if above or below the correct value.

# Pie charts variations

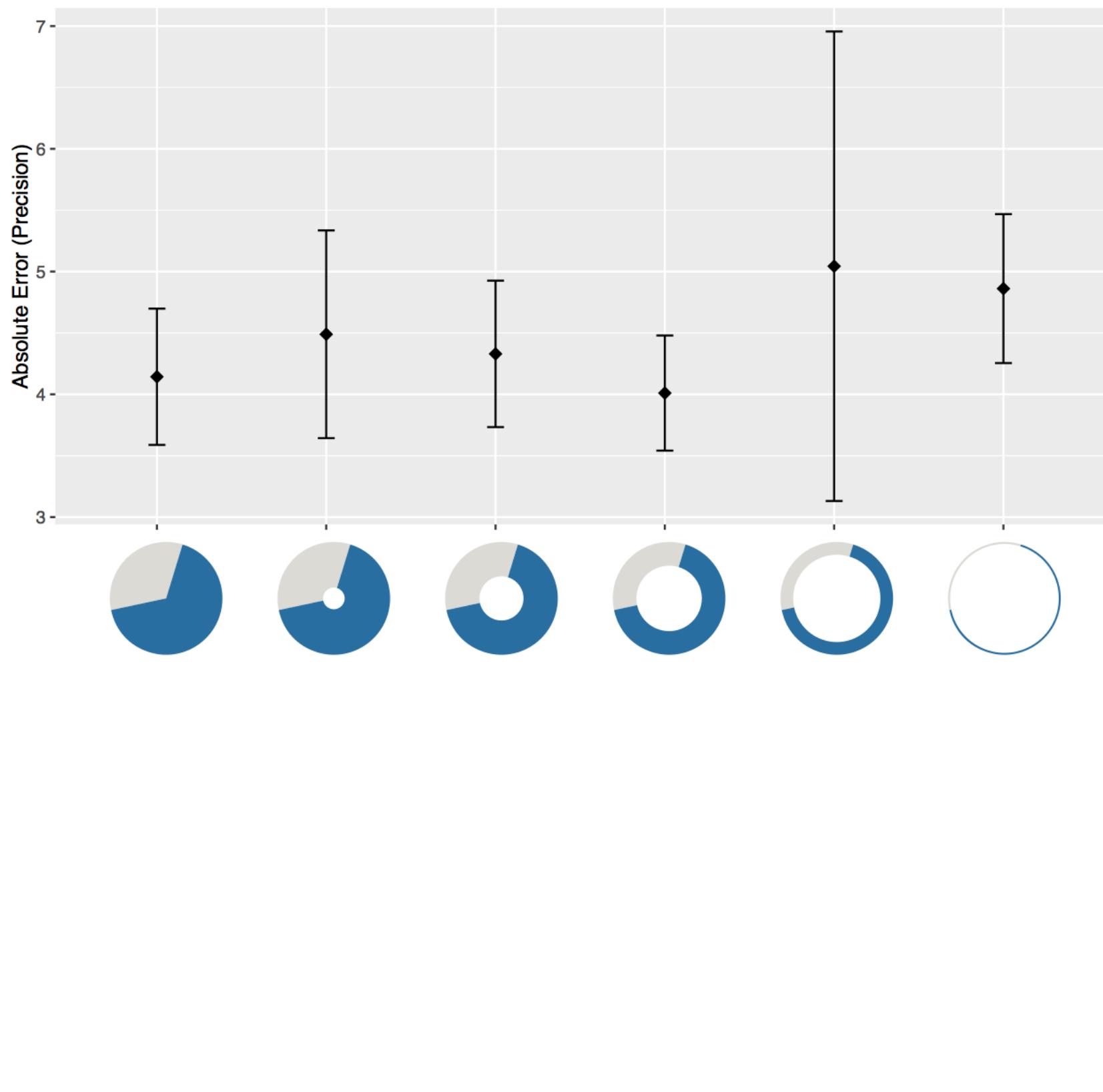
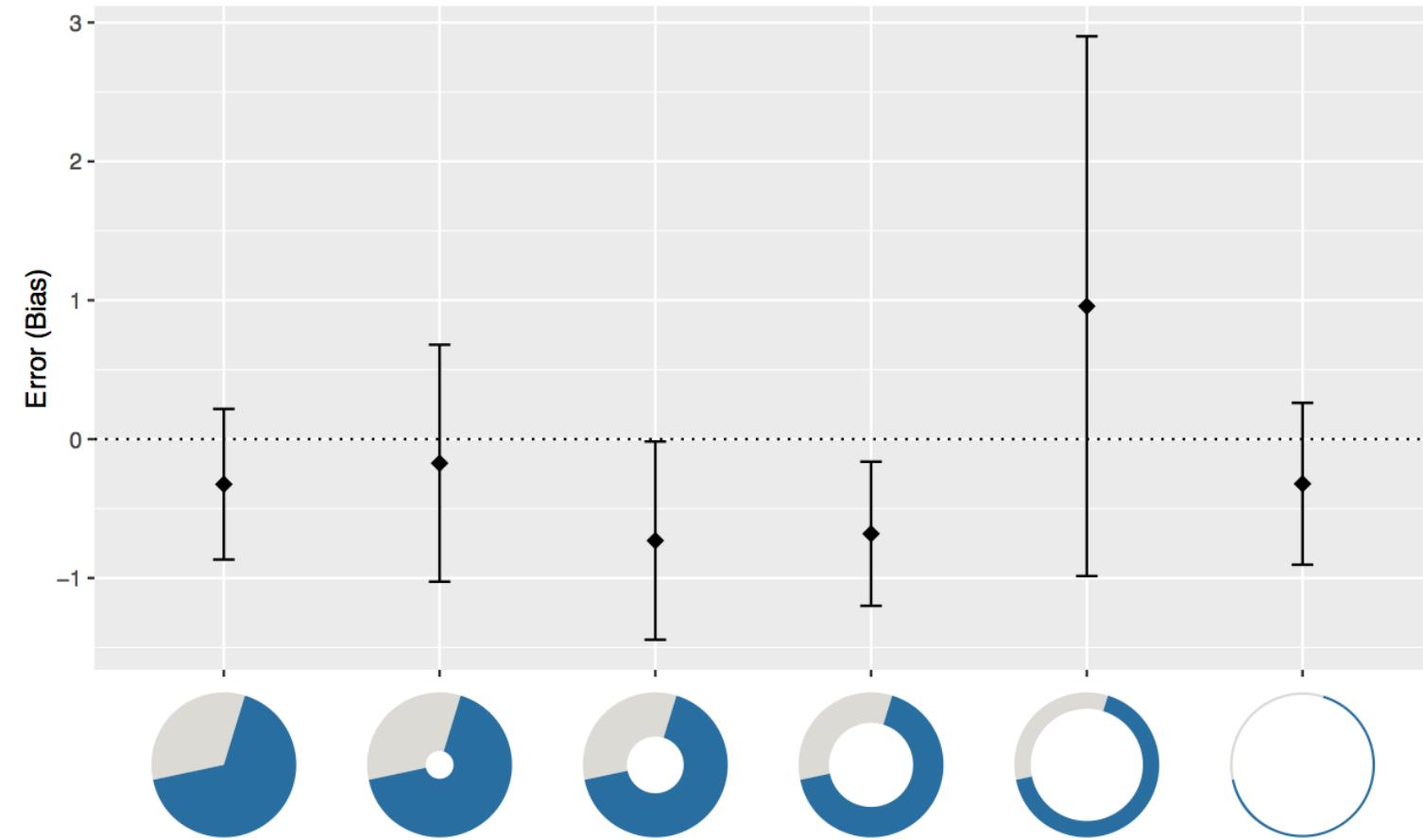


overestimate



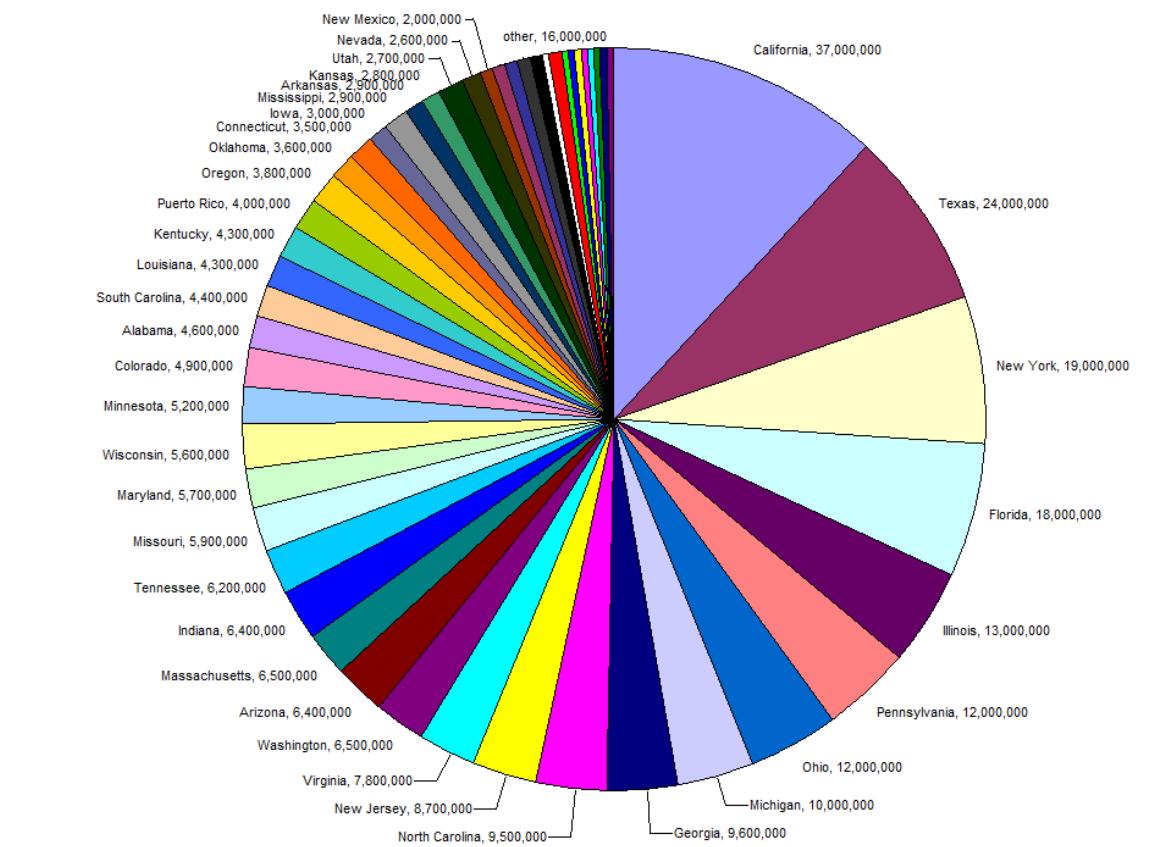
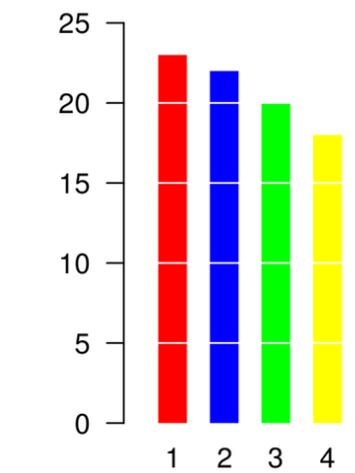
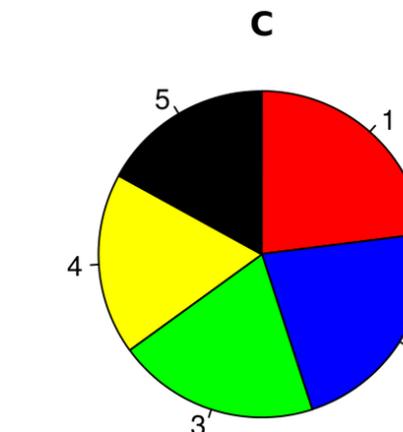
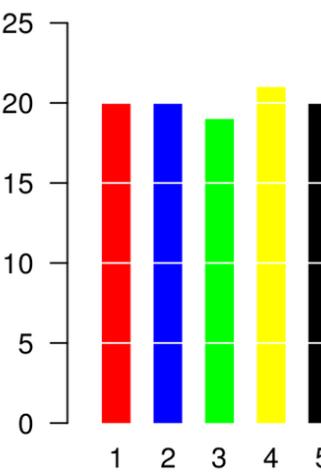
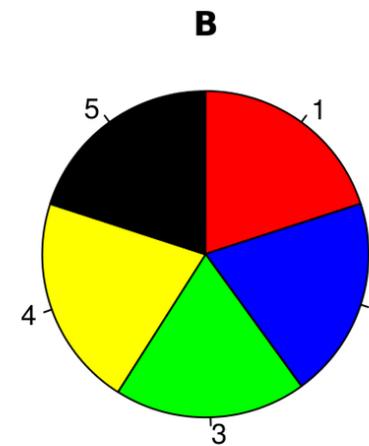
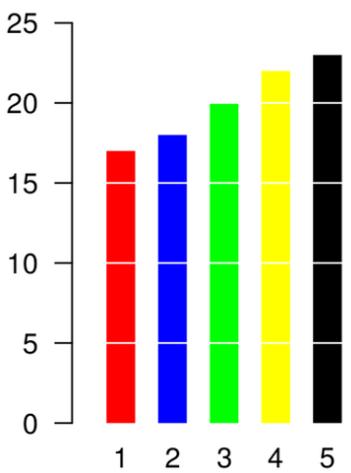
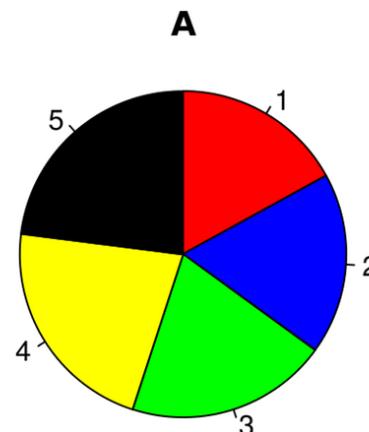
If the decomposed pie charts in the first study didn't convince you, this definitely should: the larger slice gets overestimated systematically. This is exactly what you'd expect if pie charts were read by arc length or area (since those are larger in comparison due to the larger radius), but not if you're in the angle camp. This is the smoking gun, right there!

# Donut chart Radii



# Pie charts: best practices

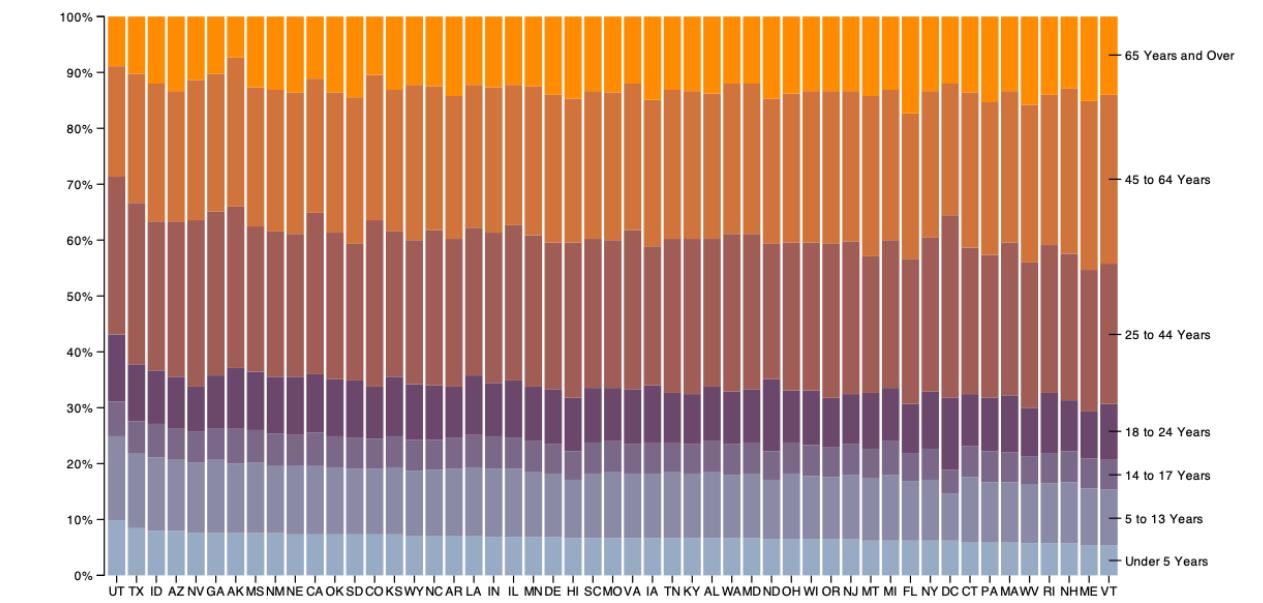
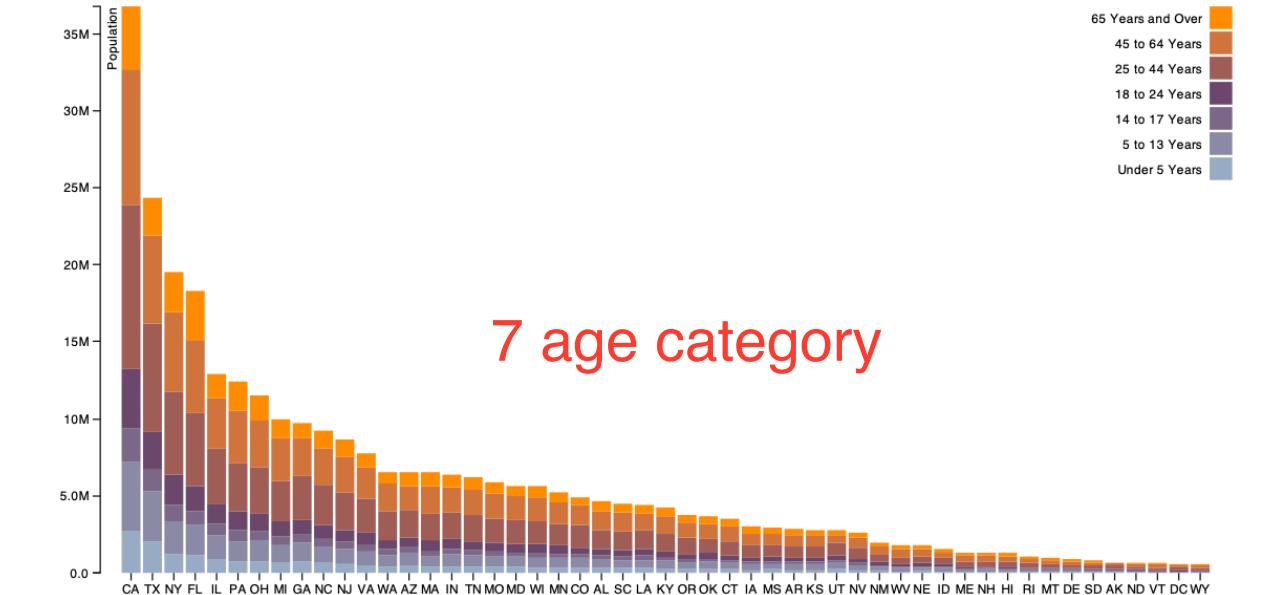
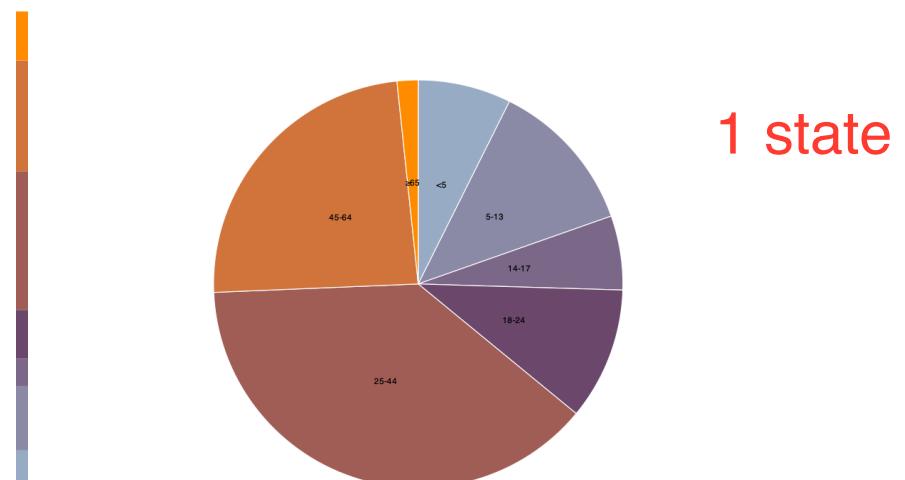
- not so bad for two (or few) levels, for part-to-whole task
- dubious for several levels if details matter
- terrible for many levels

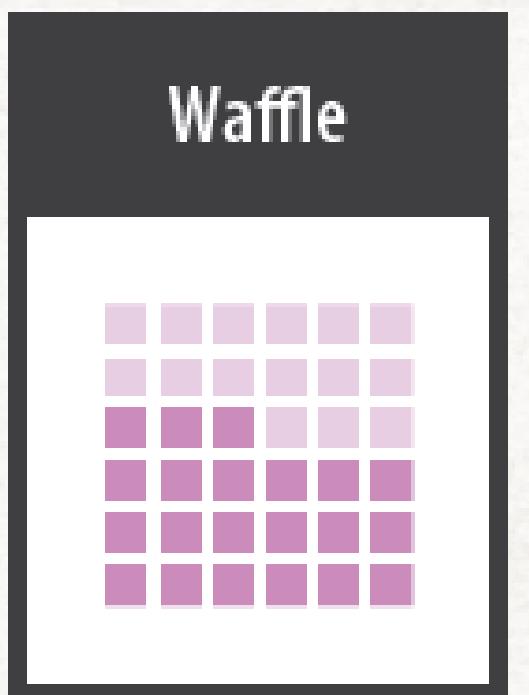
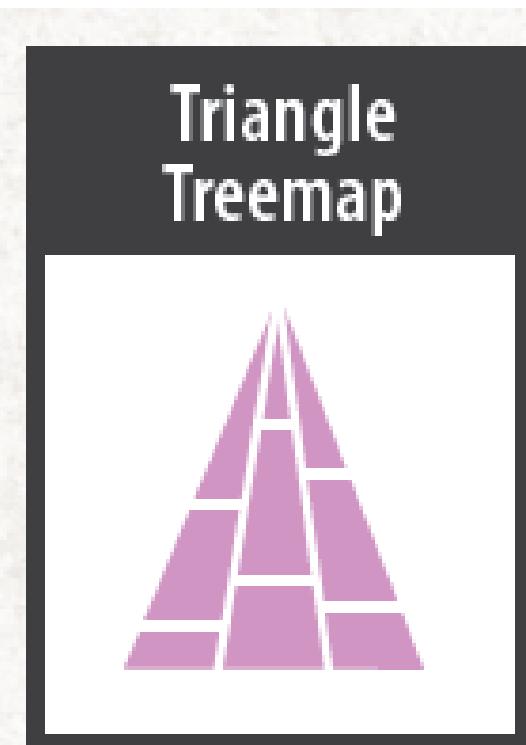
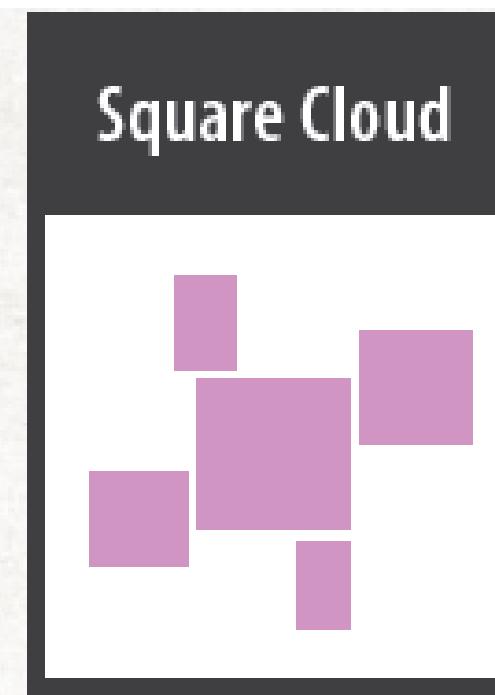
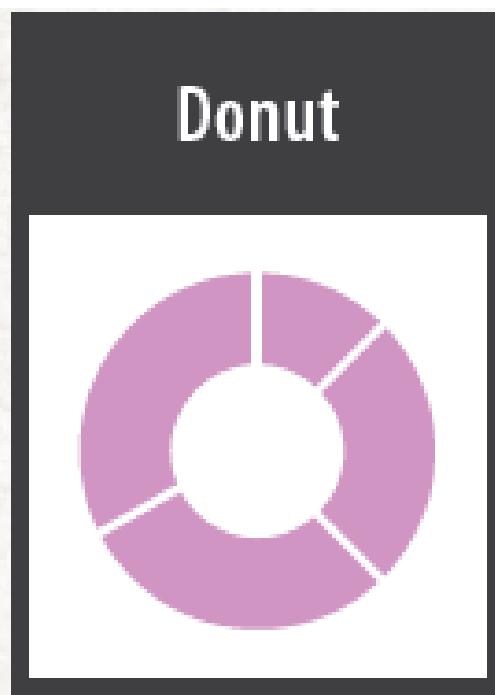
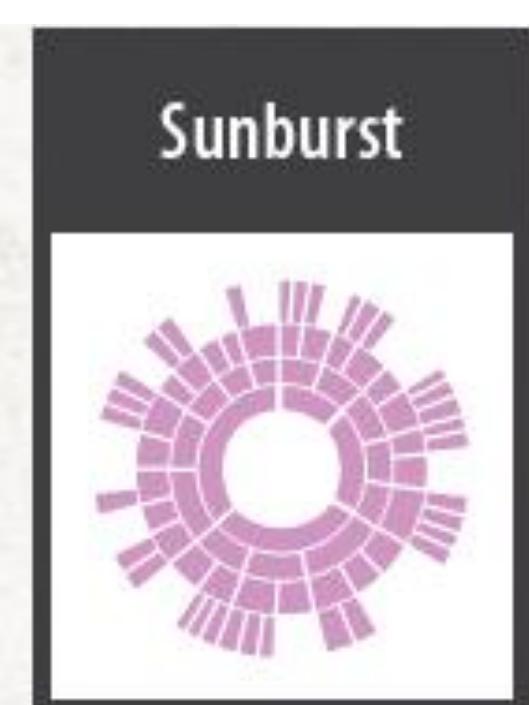
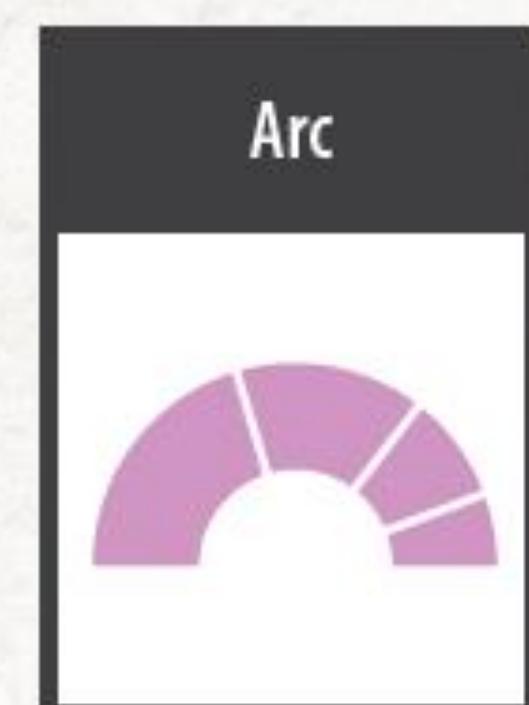
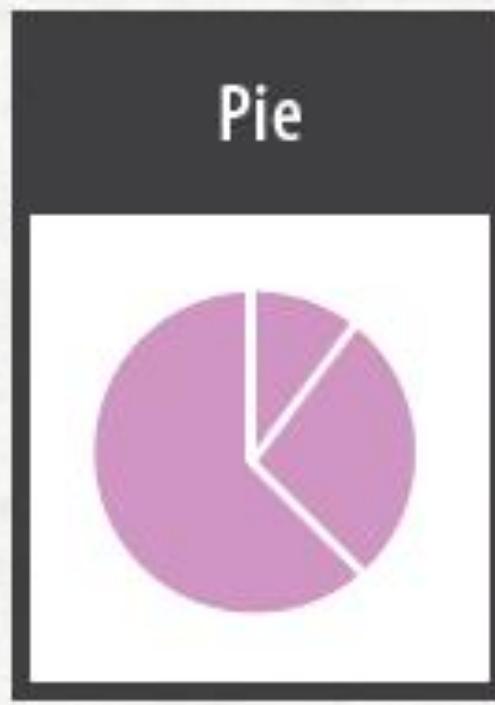


**LARGE AMOUNT AND SIMLIAR  
use pie chart with ordering**

# Idioms: **normalized stacked bar chart**

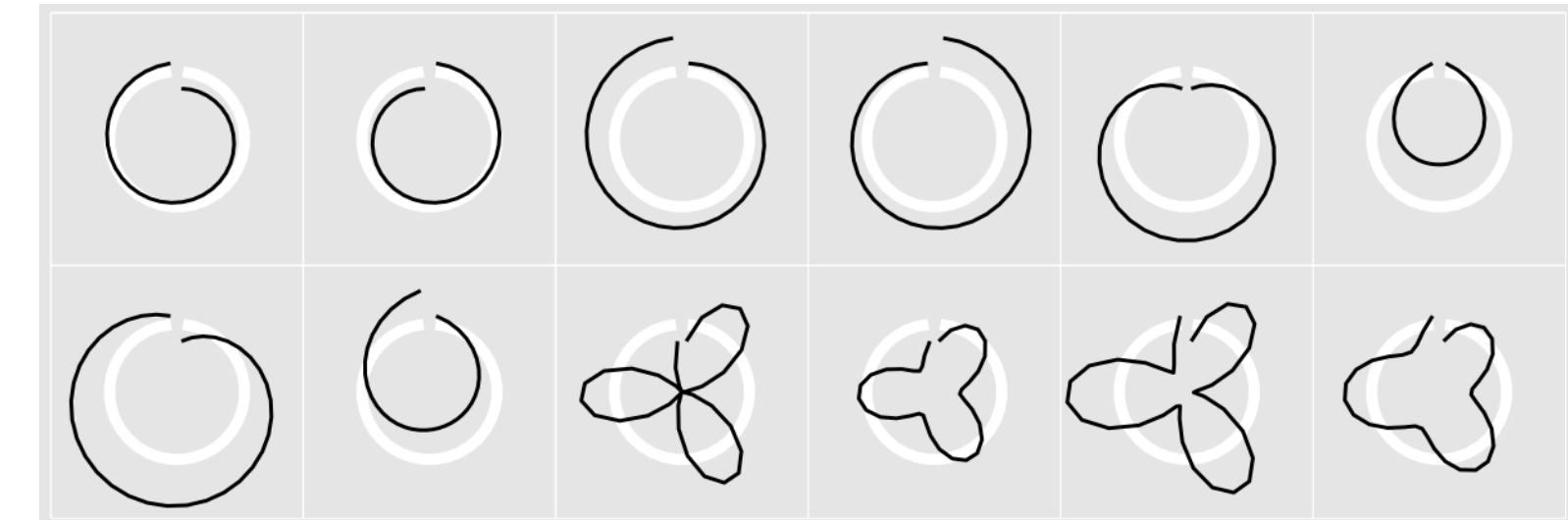
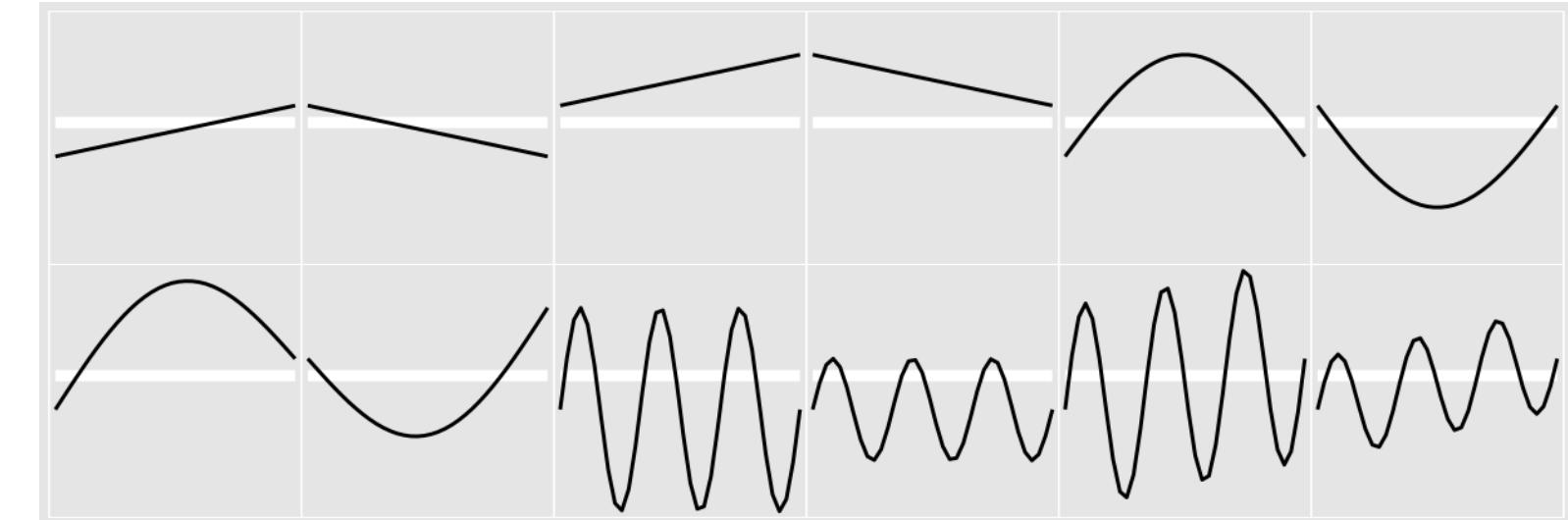
- task
  - part-to-whole judgements
- normalized stacked bar chart
  - stacked bar chart, normalized to full vert height
  - single stacked bar equivalent to full pie
    - high information density: requires narrow rectangle
- pie chart
  - information density: requires large circle





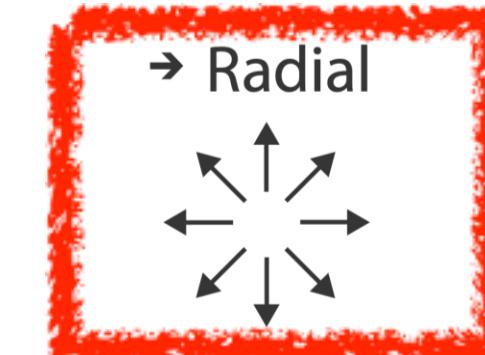
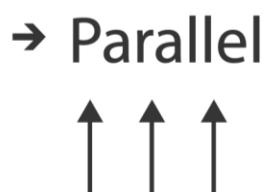
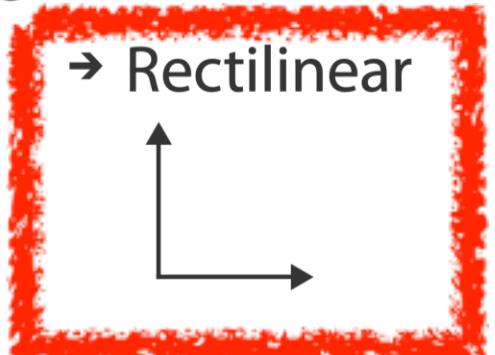
# Idiom: **glyphmaps**

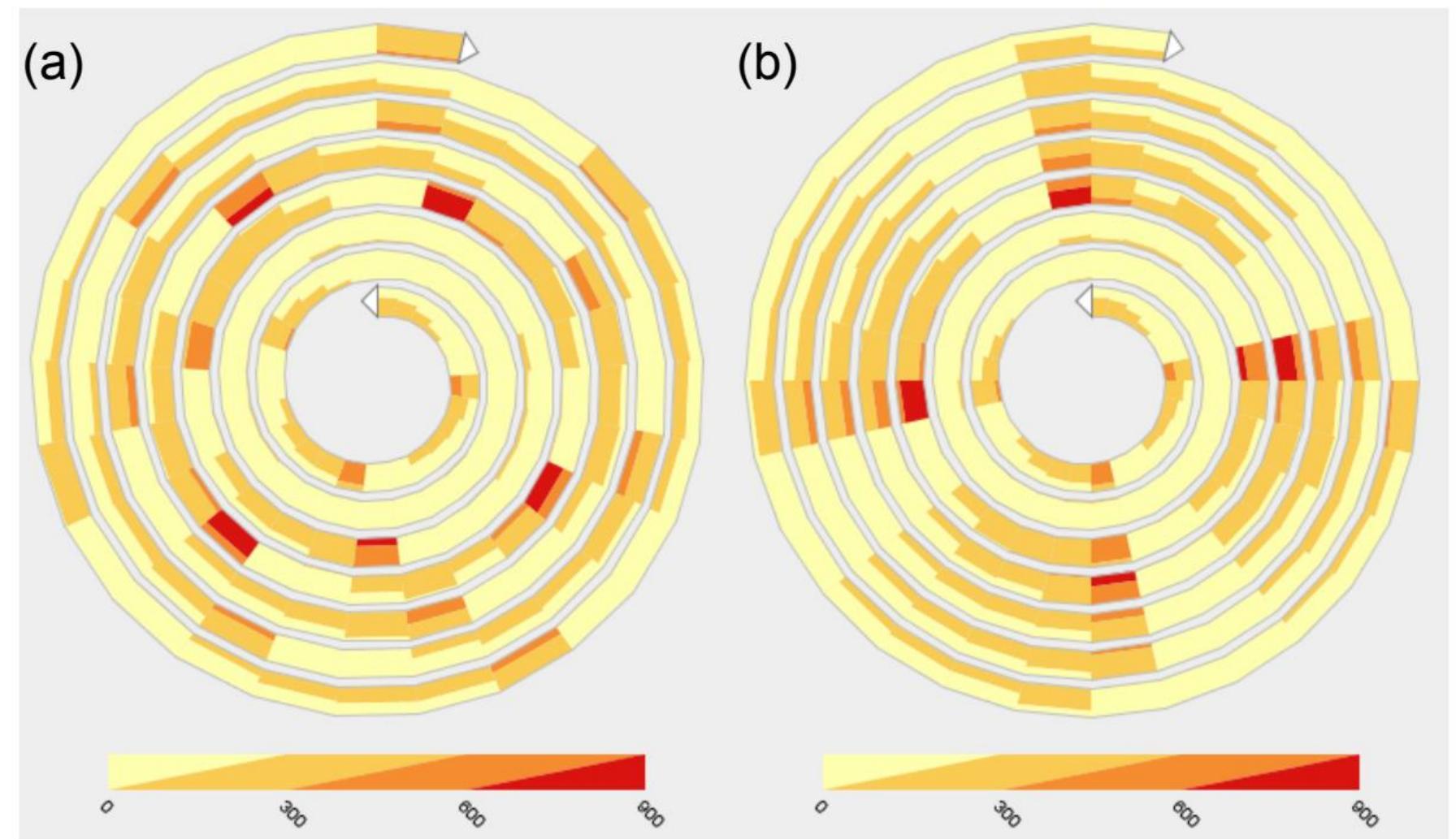
- rectilinear good for linear vs nonlinear trends
- radial good for cyclic patterns
  - evaluating periodicity



[*Glyph-maps for Visually Exploring Temporal Patterns in Climate Data and Models.*  
Wickham, Hofmann, Wickham, and Cook. *Environmetrics* 23:5 (2012), 382–393.]

## → Axis Orientation

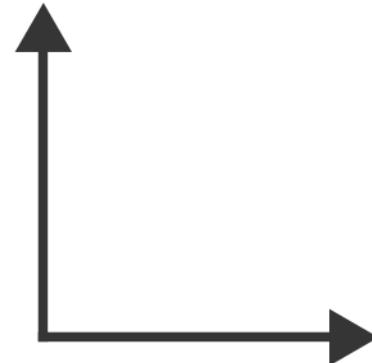




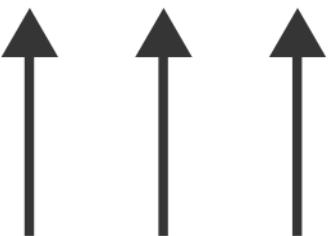
**Figure 6:** Finding a pattern – (a) Cycle length = 25; (b) Cycle length = 28

## → Axis Orientation

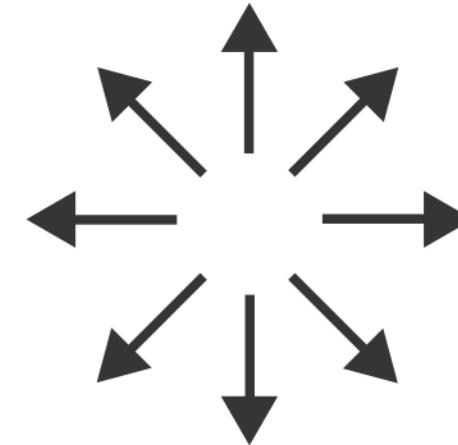
→ Rectilinear



→ Parallel



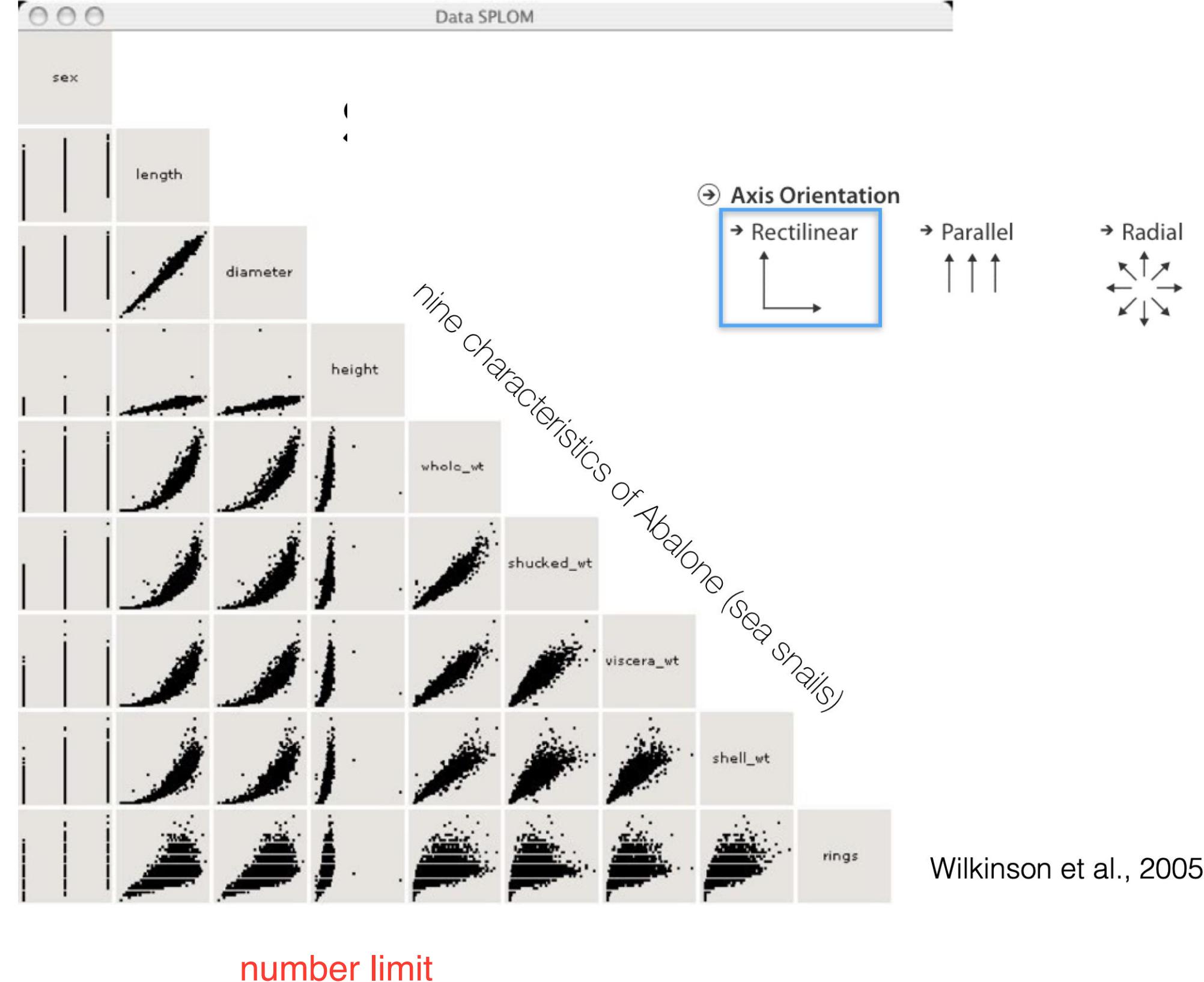
→ Radial



# Idiom: **SPLOM**

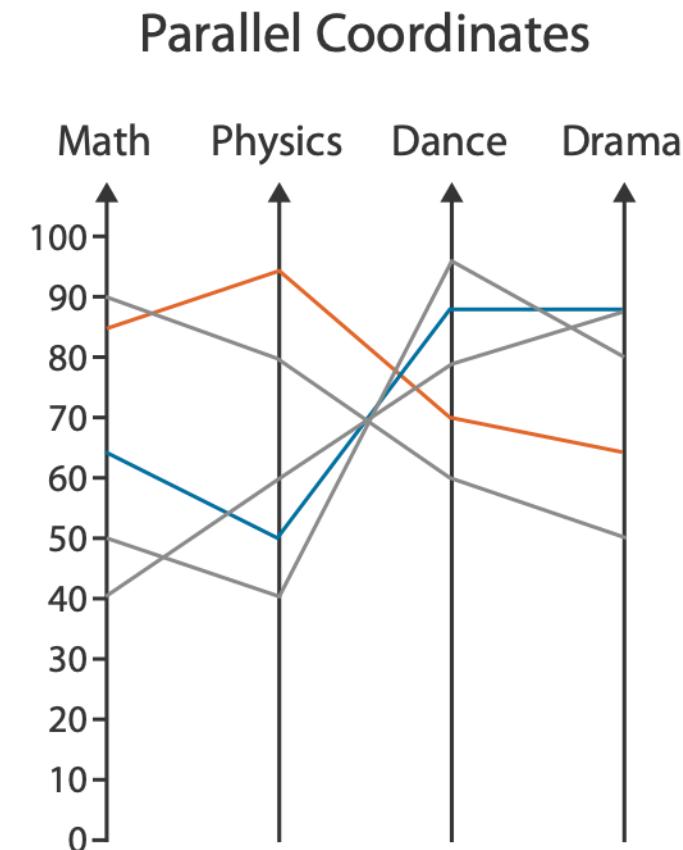
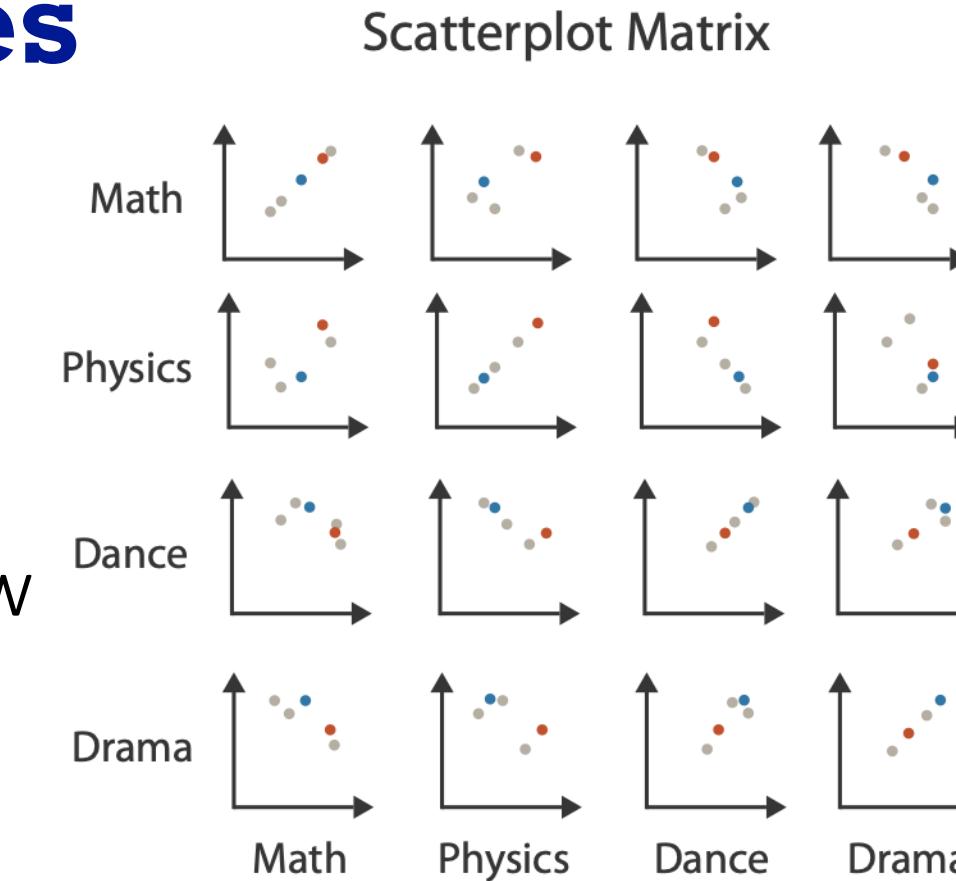
scatterplot matrix  
(SPLOM)

- rectilinear axes,  
point mark
- all possible pairs of axes
- scalability
  - one dozen attrs
  - dozens to hundreds of items



# Idioms: parallel coordinates

- scatterplot limitation
  - visual representation with orthogonal axes
  - can show only two attributes with spatial position channel
- alternative: line up axes in parallel to show many attributes with position
  - item encoded with a line with n segments
  - n is the number of attributes shown
- parallel coordinates
  - parallel axes, jagged line for item
  - rectilinear axes, item as point
    - axis ordering is major challenge
  - scalability
    - dozens of attrs
    - hundreds of items



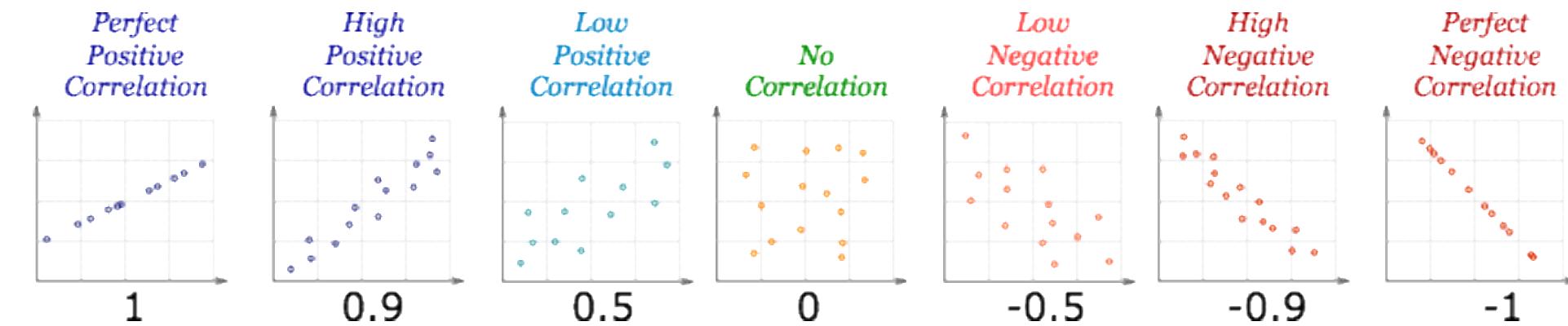
Table

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

after [Visualization Course Figures. McGuffin, 2014.  
<http://www.michaelmcguffin.com/courses/vis/>]

# Task: Correlation

- scatterplot matrix
  - positive correlation
    - diagonal low-to-high
  - negative correlation
    - diagonal high-to-low
  - uncorrelated: spread out
- parallel coordinates
  - positive correlation
    - parallel line segments
  - negative correlation
    - all segments cross at halfway point
  - uncorrelated
    - scattered crossings



<https://www.mathsisfun.com/data/scatter-xy-plots.html>

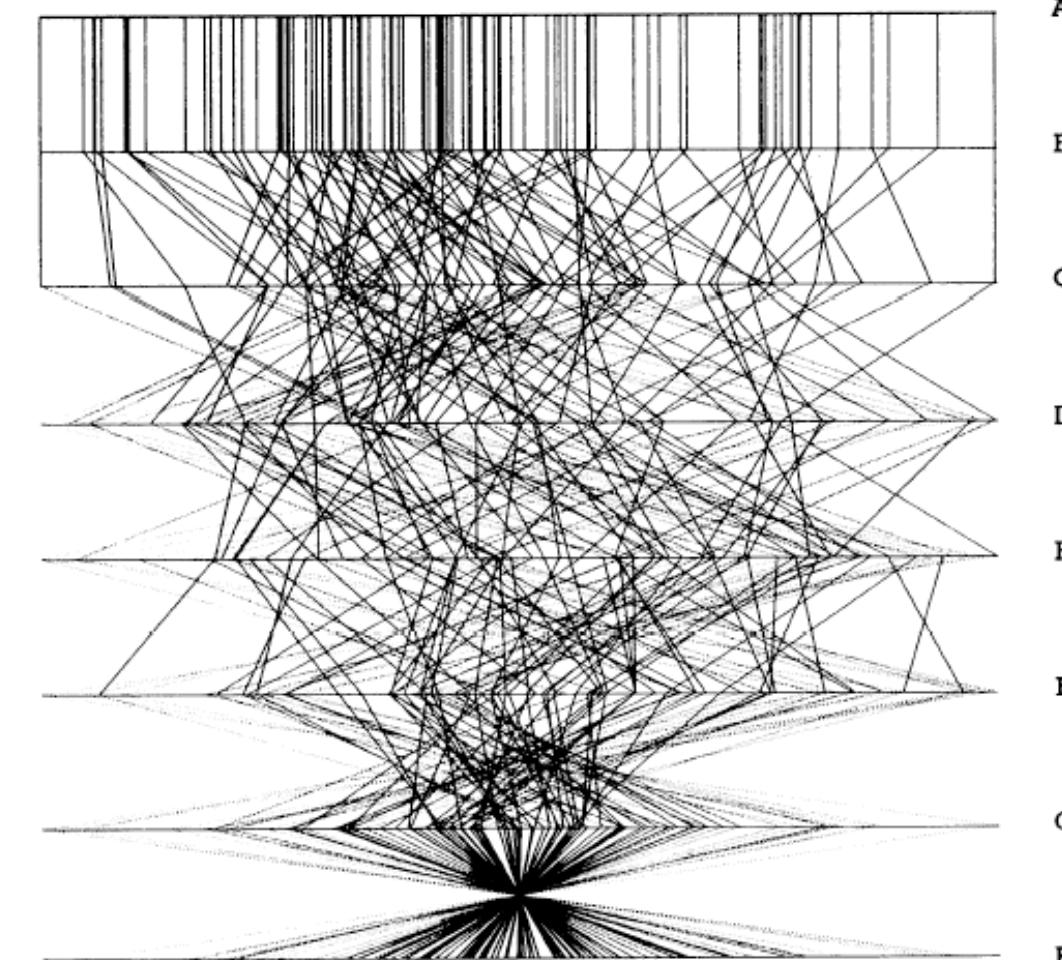
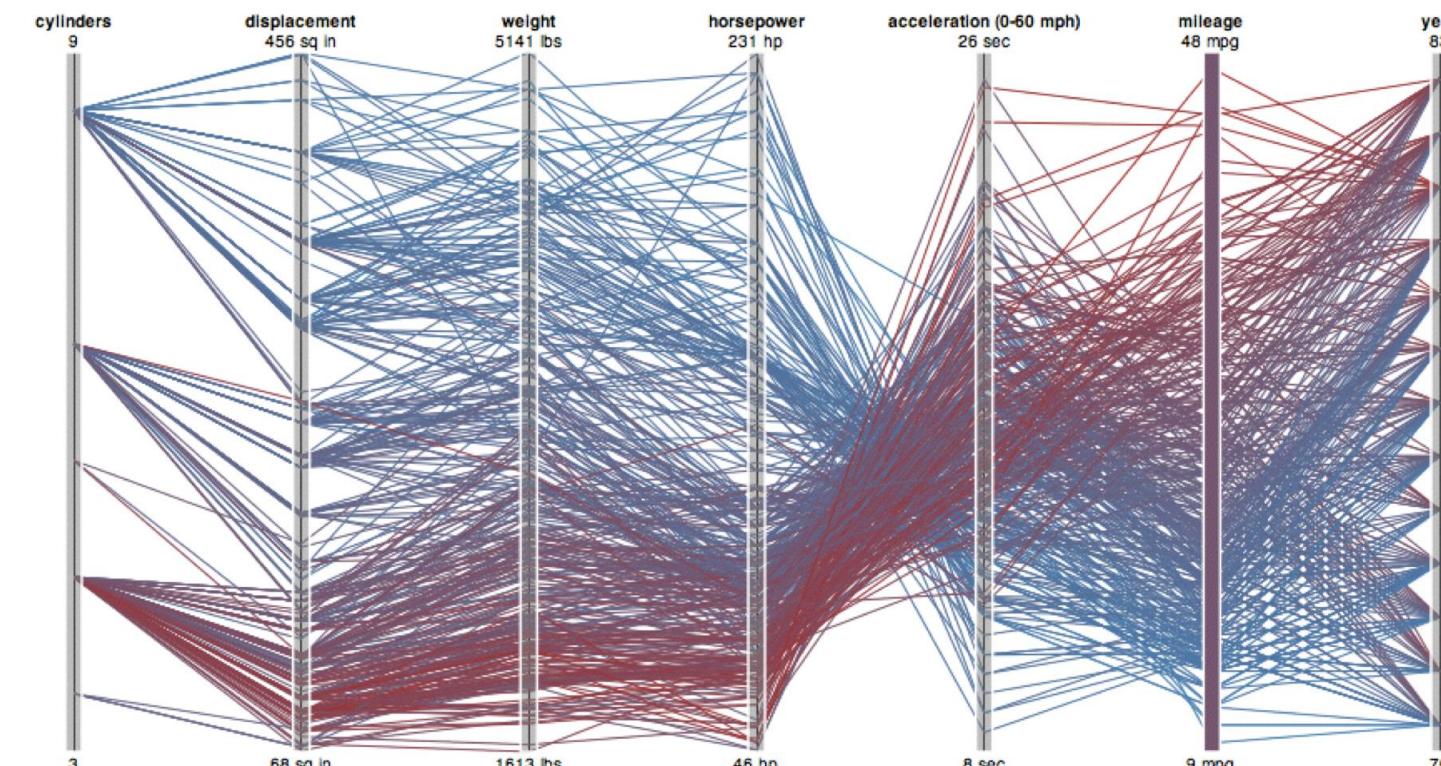


Figure 3. Parallel Coordinate Plot of Six-Dimensional Data Illustrating Correlations of  $\rho = 1, .8, .2, 0, -.2, -.8, \text{ and } -1$ .

[Hyperdimensional Data Analysis Using Parallel Coordinates. Wegman. Journ. American Statistical Association 85:411 (1990), 664–675.]

# Parallel coordinates, limitations

- visible patterns only between neighboring axis pairs
- how to pick axis order?
  - usual solution: reorderable axes, **interactive** exploration
  - same weakness as many other techniques
    - downside of interaction: human-powered search
  - some algorithms proposed, none fully solve



# Orientation limitations

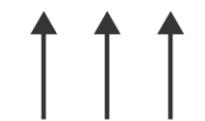
- rectilinear: scalability wrt #axes
  - 2 axes best, 3 problematic, 4+ impossible
- parallel: unfamiliarity, training time
- radial: perceptual limits
  - polar coordinate asymmetry
    - angles lower precision than length
    - **nonuniform** sector width/size depending on radial distance
  - frequently problematic
    - but sometimes can be deliberately exploited!
      - for 2 attrs of very unequal importance

## → Axis Orientation

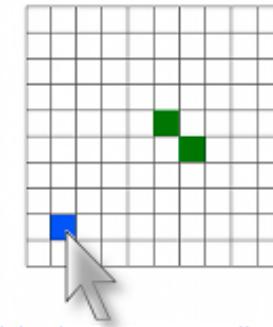
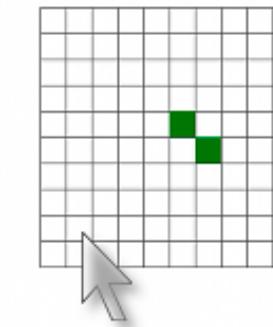
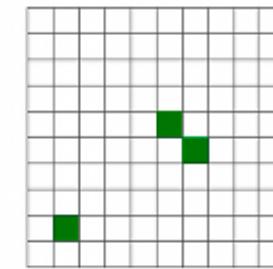
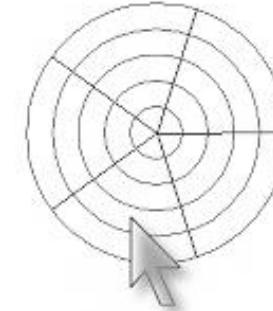
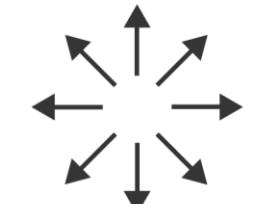
→ Rectilinear



→ Parallel



→ Radial



# Arrange tables

## → Express Values

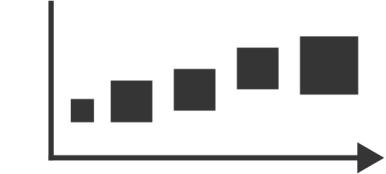


## → Separate, Order, Align Regions

→ Separate



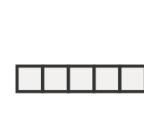
→ Order



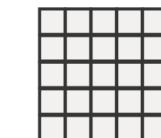
→ Align



→ 1 Key  
*List*

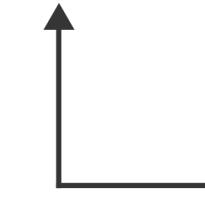


→ 2 Keys  
*Matrix*



## → Axis Orientation

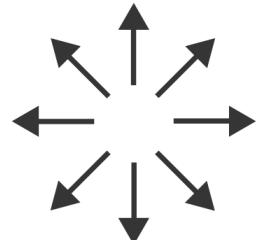
→ Rectilinear



→ Parallel

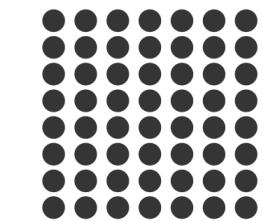


→ Radial



## → Layout Density

→ Dense

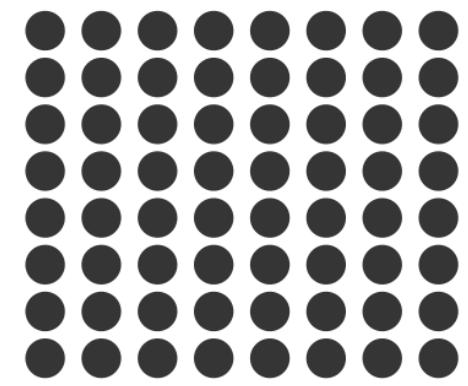


# Layout density



## Layout Density

→ Dense



→ Space-Filling

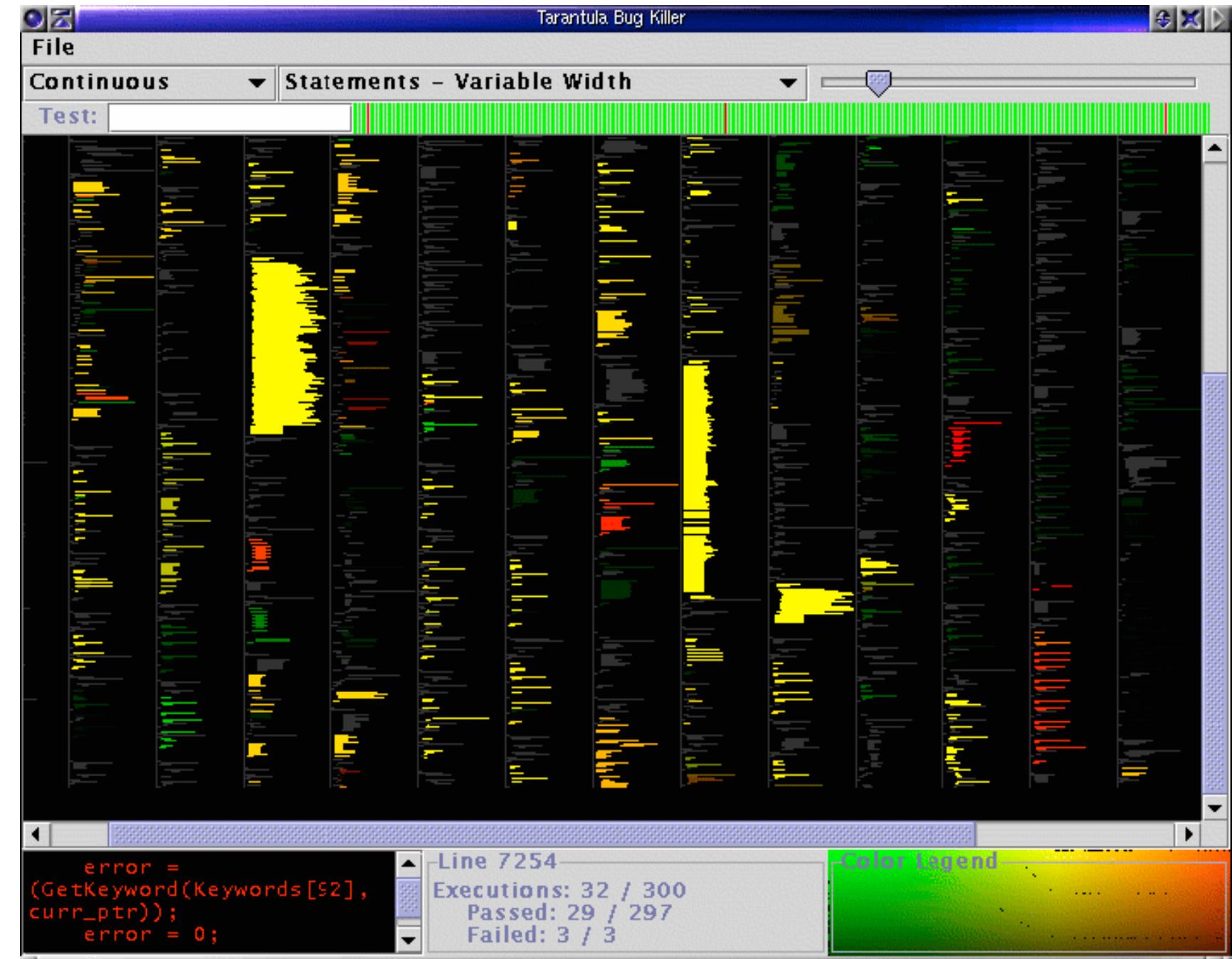
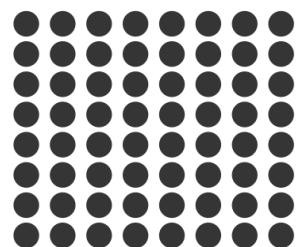


# Idiom: Dense software overviews

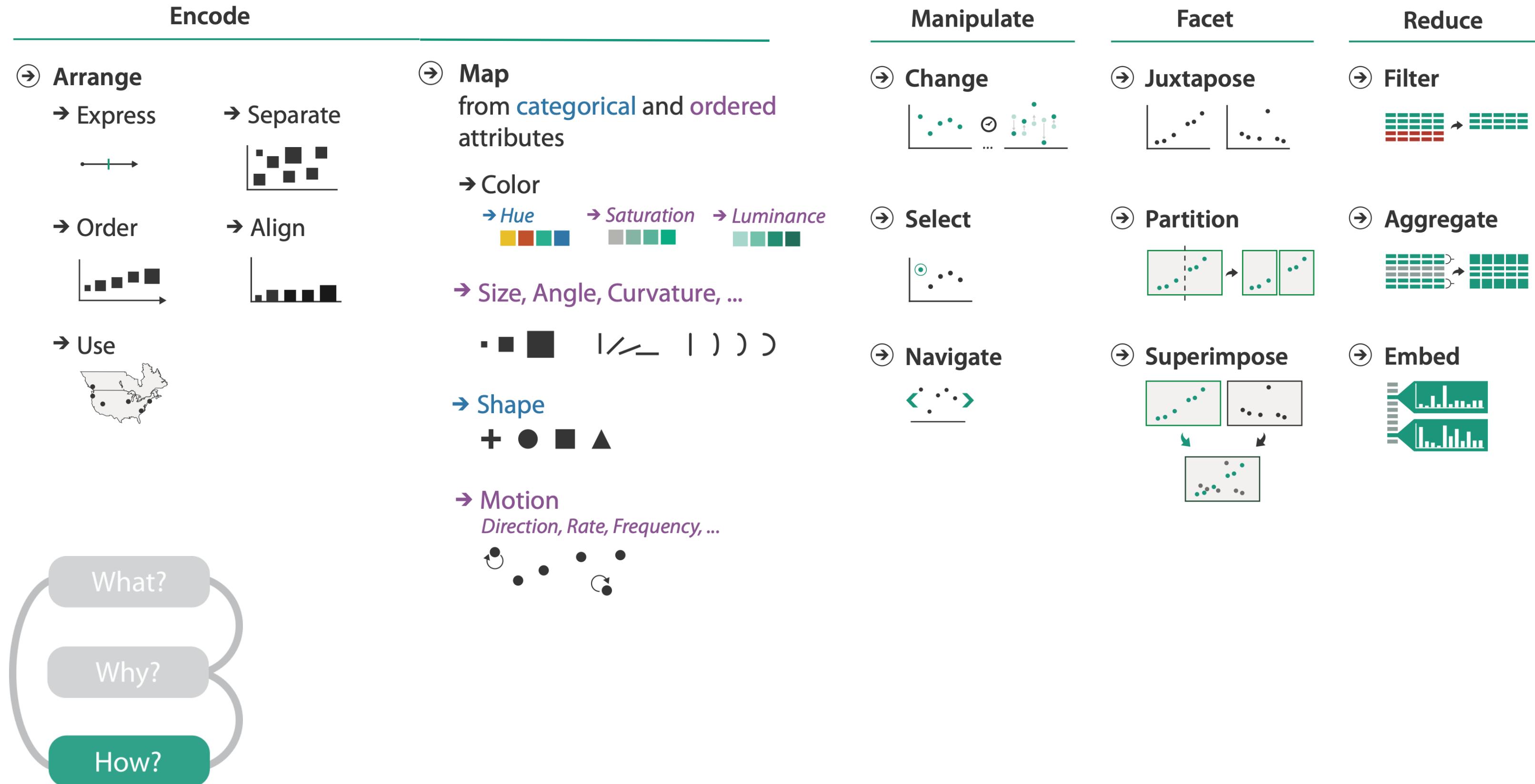
- data: text
  - text + 1 quant attrib per line
- derived data:
  - one pixel high line
  - length according to original
- color line by attrib
- scalability
  - 10K+ lines

➔ Layout Density

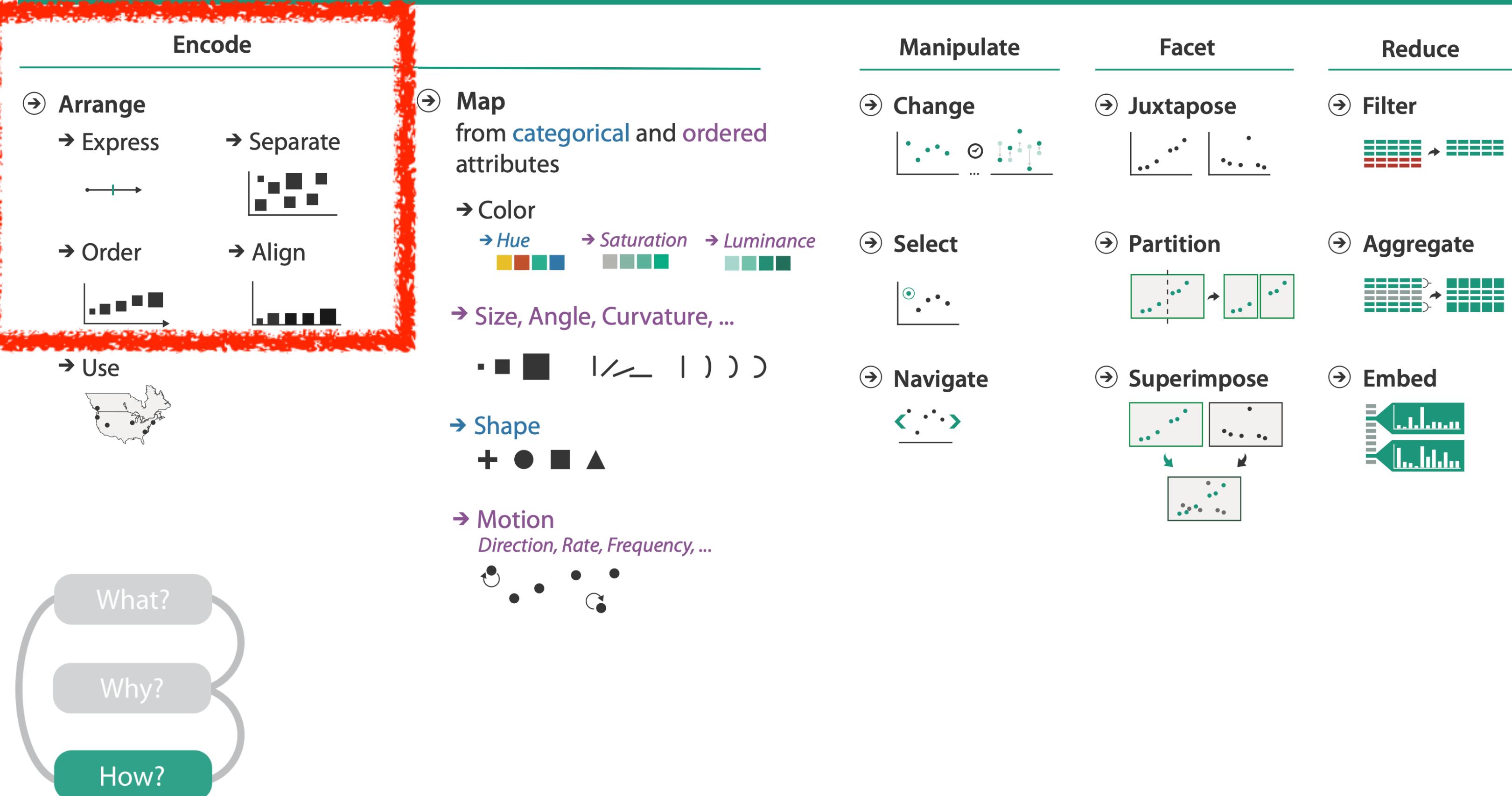
➔ Dense



# How?

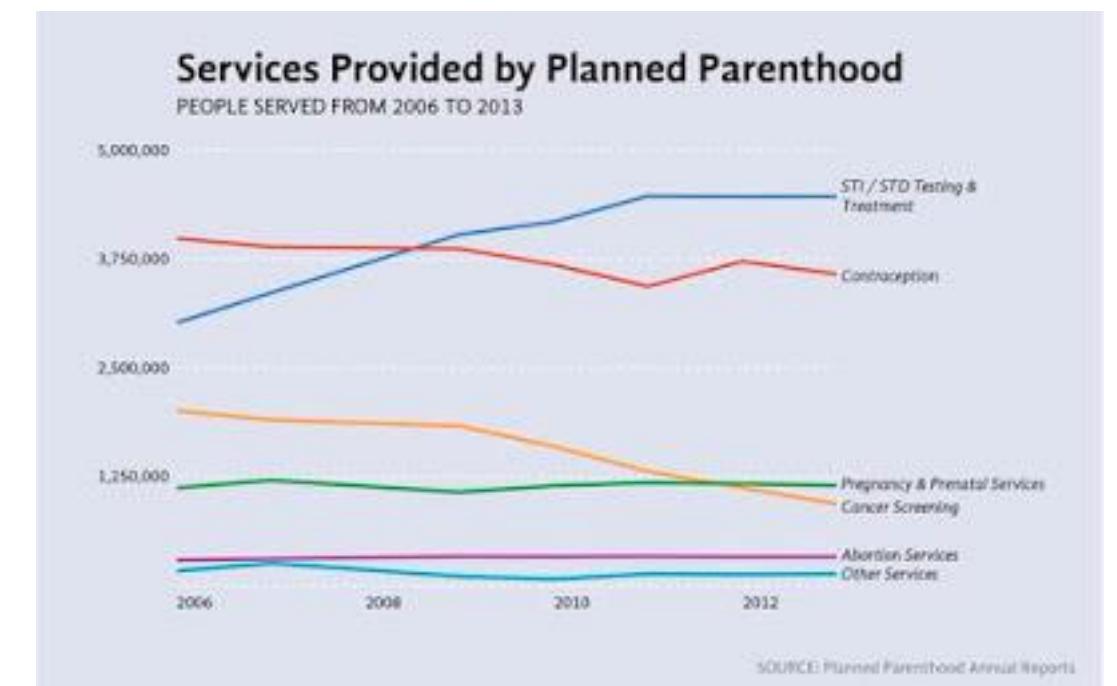
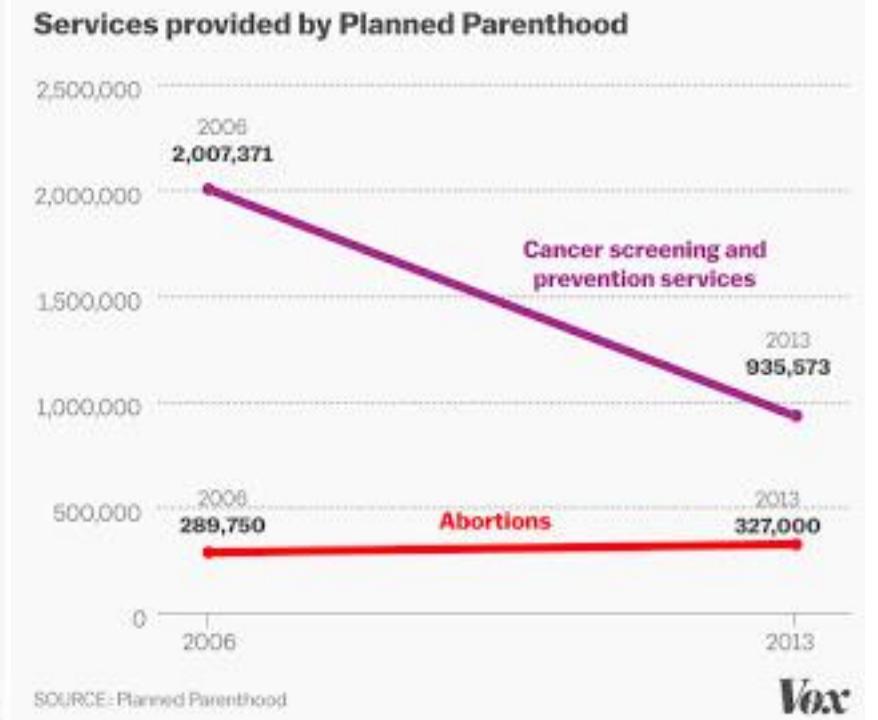
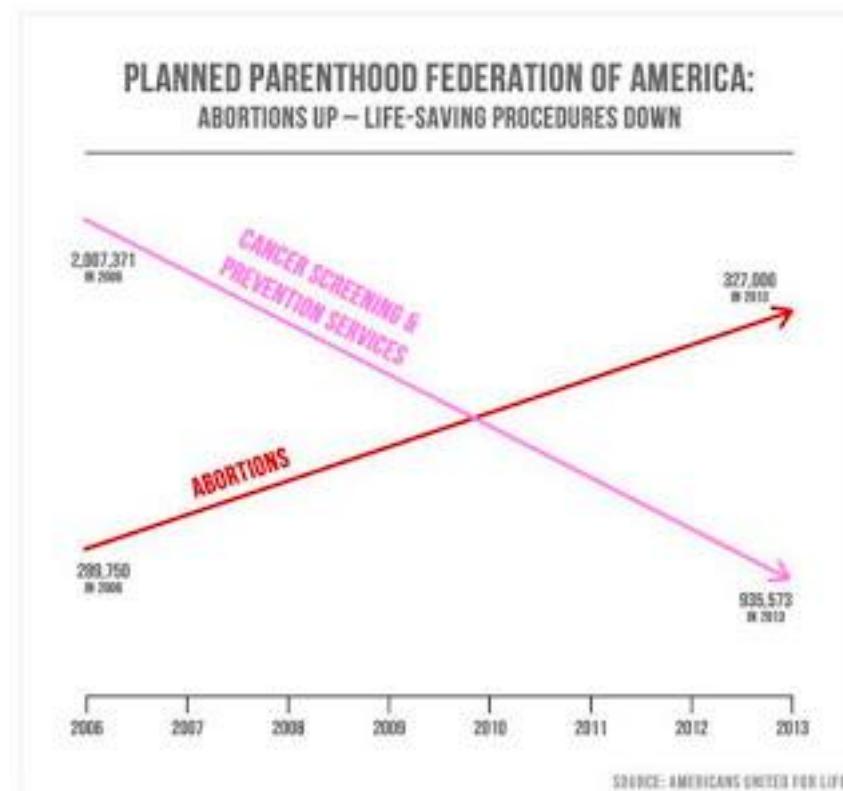


# How?



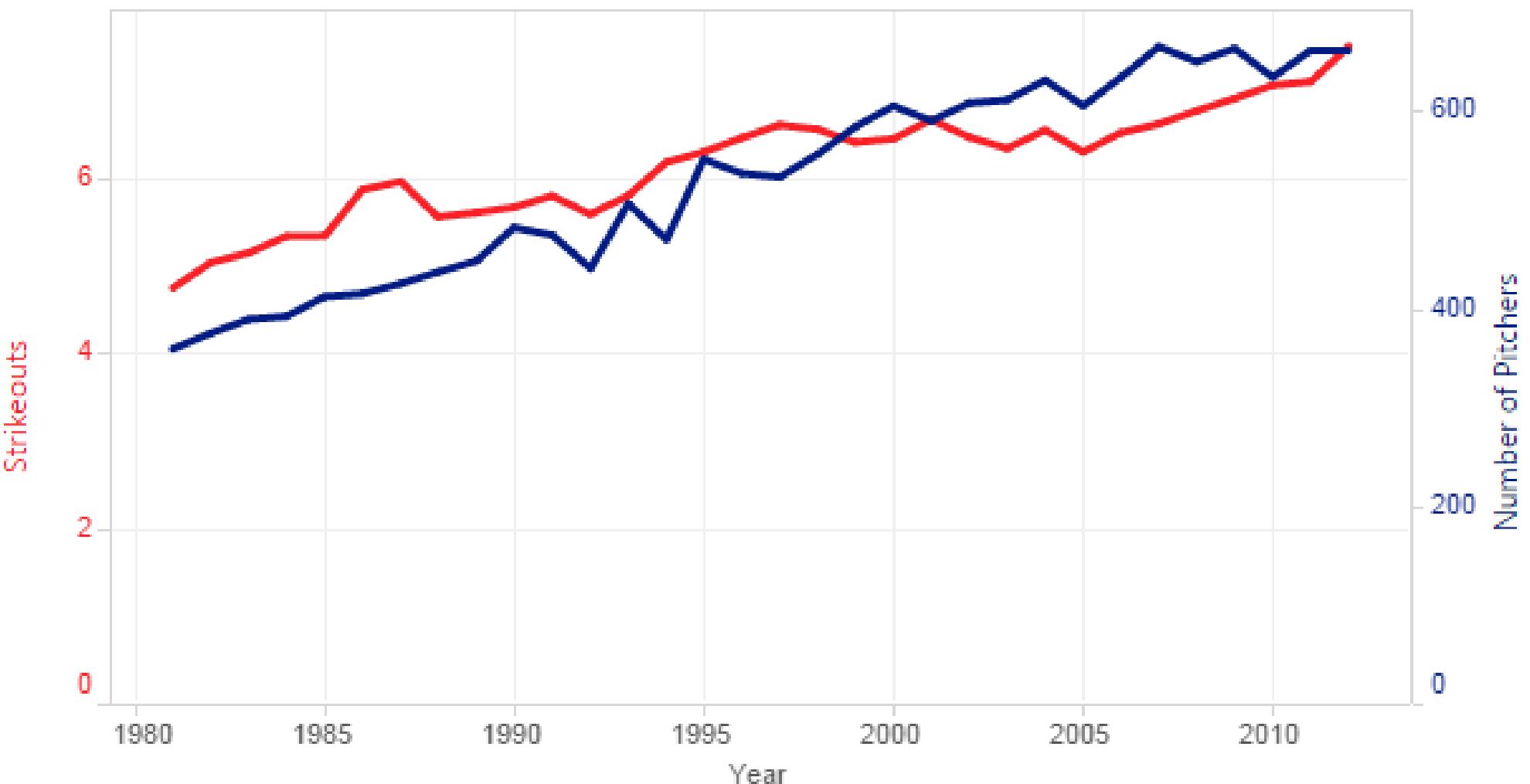
# Chart axes

- labelled axis is critical
- avoid cropping y-axis
  - include 0 at bottom left
  - or slope misleads



# Idiom: **dual-axis line charts**

- controversial
  - acceptable if commensurate
  - beware, very easy to mislead!

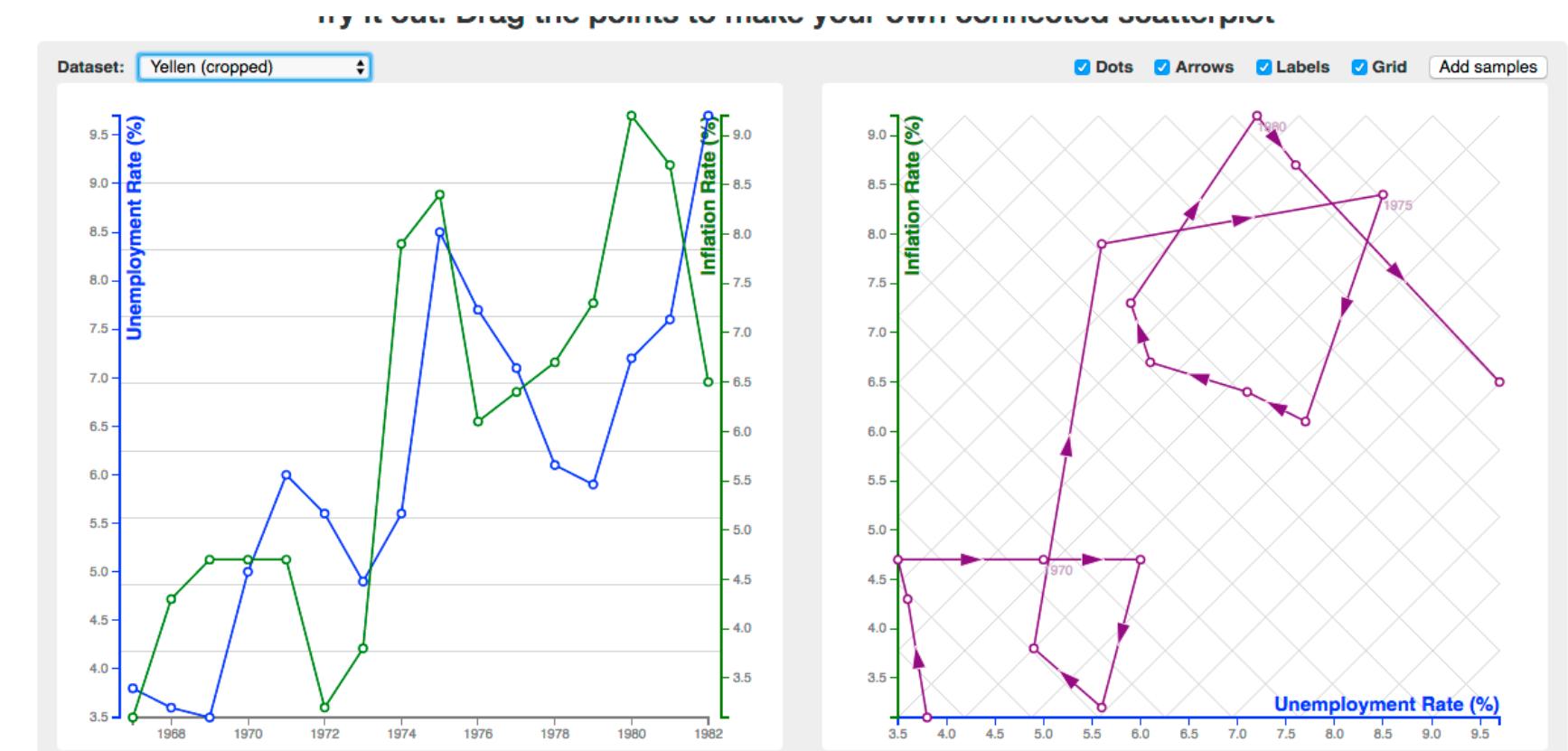
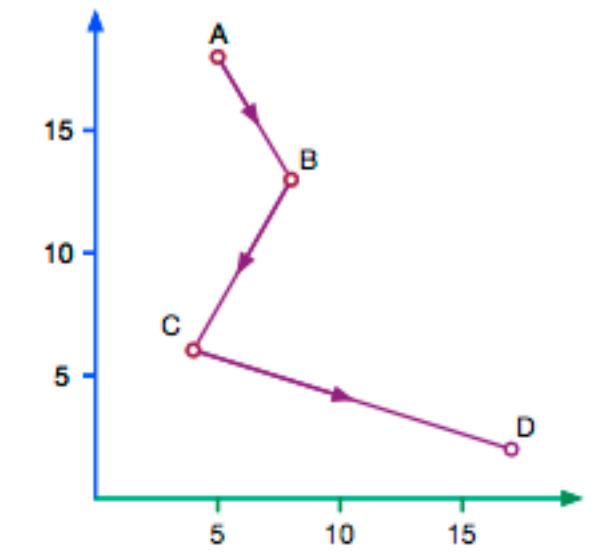
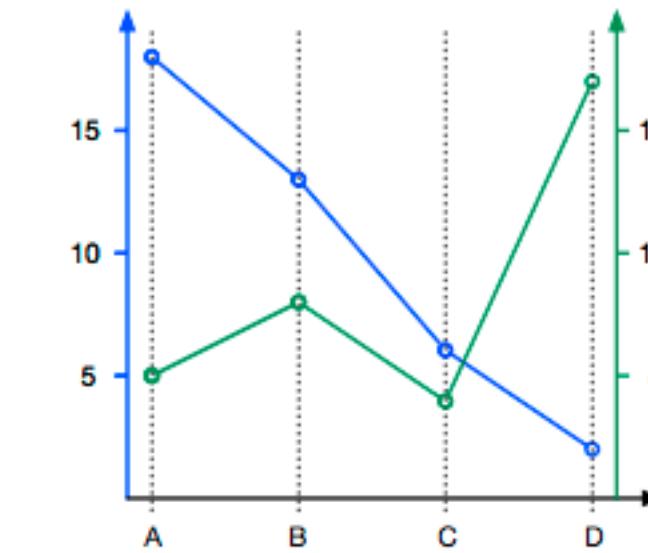
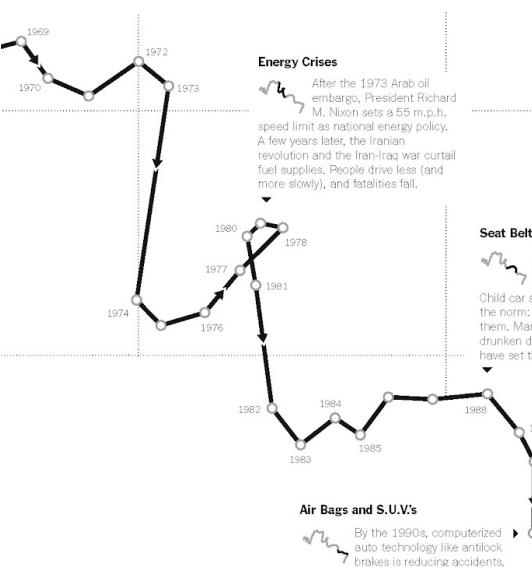


Source | <http://www.baseball-reference.com/leagues/MLB/pitch.shtml>

Ben Jones (@DataRemixed) | 5/4/2013

# Idiom: connected scatterplots

- scatterplot with line connection marks
  - popular in journalism
  - horiz + vert axes: value attrs
  - line connection marks:
    - temporal order
  - alternative to dual-axis charts
    - horiz: time
    - vert: two value attrs
- empirical study
  - engaging, but correlation unclear



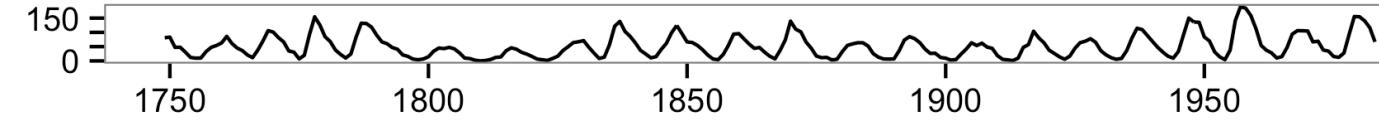
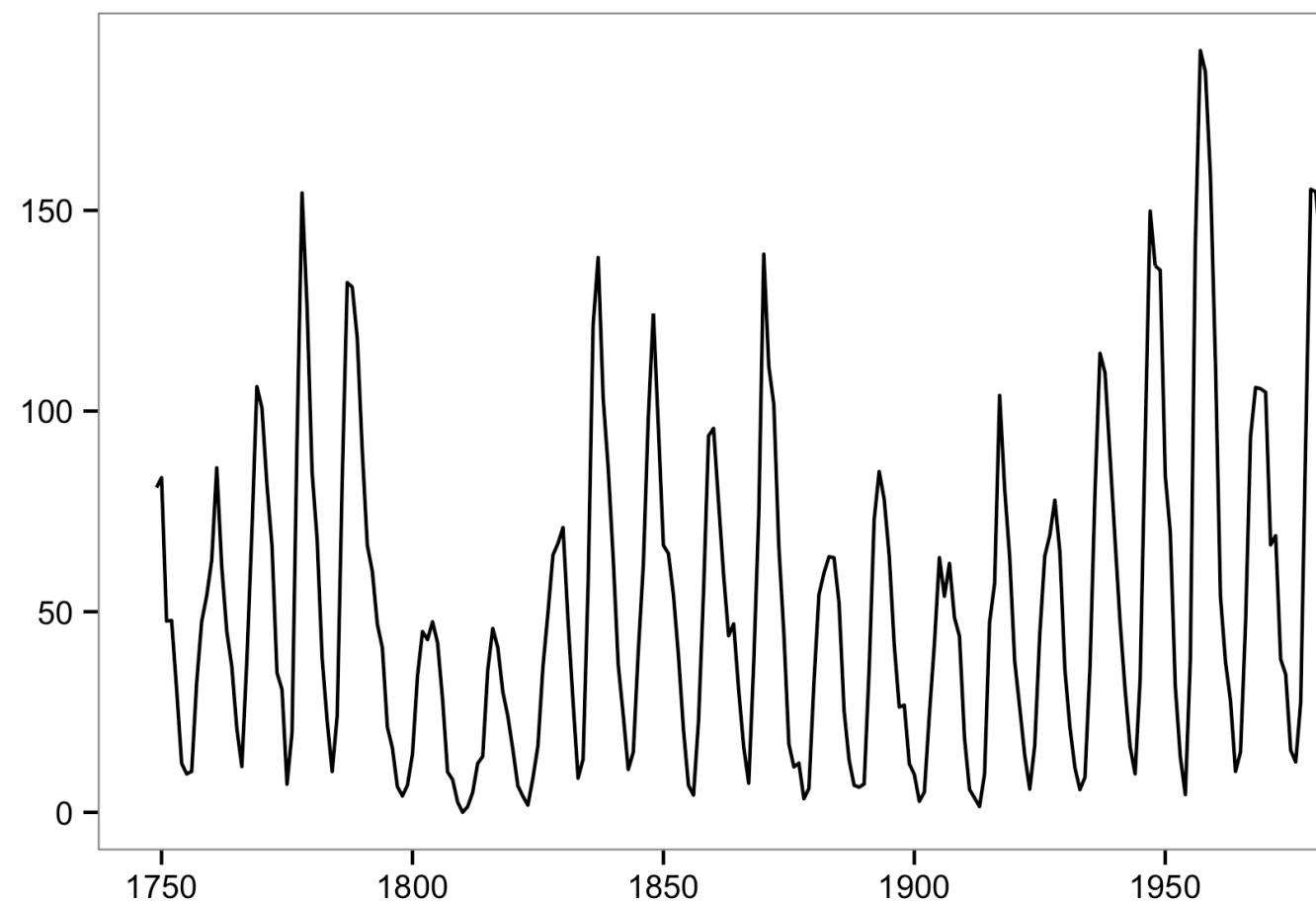
[The Connected Scatterplot for Presenting Paired Time Series.  
Haroz, Kosara and Franconeri. IEEE TVCG 22(9):2174-86,  
2016.]

[http://steveharoz.com/research/connected\\_scatterplot/](http://steveharoz.com/research/connected_scatterplot/)

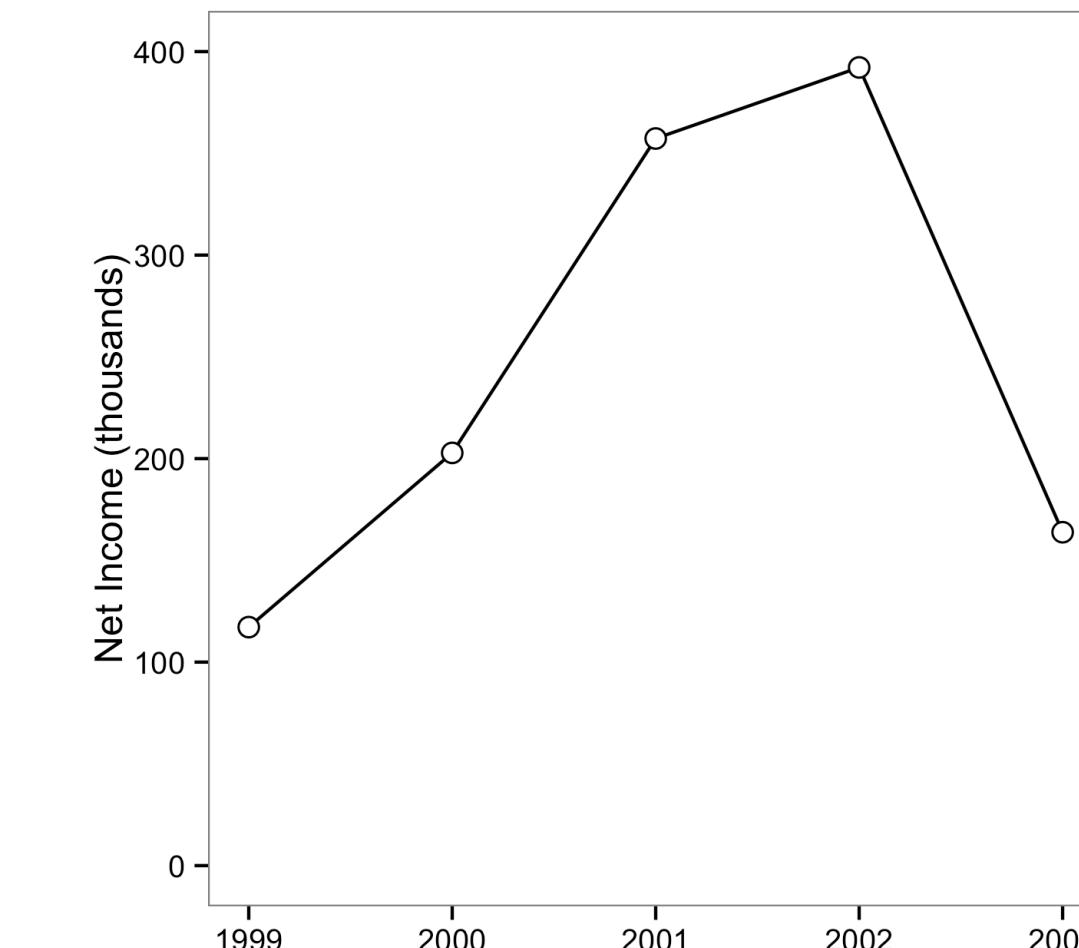
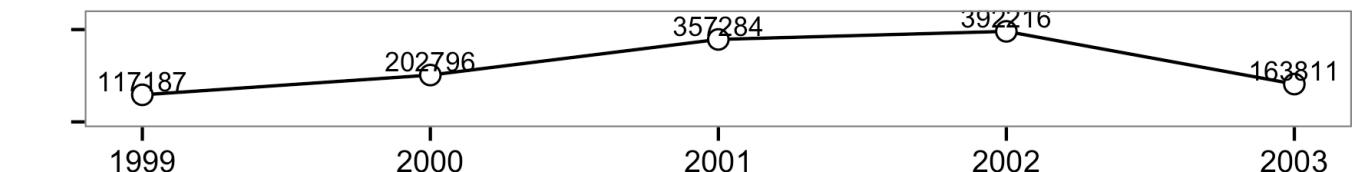
# Choosing line chart aspect ratios

- 1: banking to 45 (1980s)
  - Cleveland perceptual argument: most accurate angle judgement at 45

**Fig 7.1 Sunspot Data: Aspect Ratio 1**

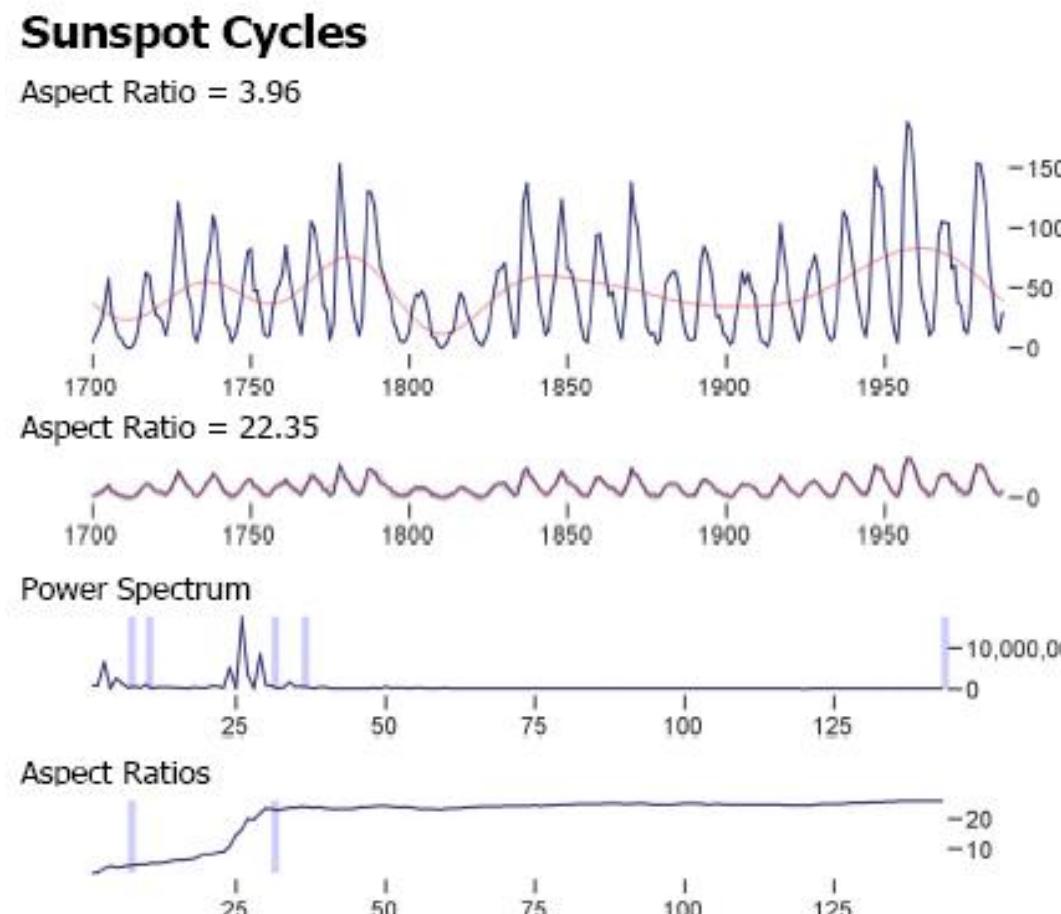


**Fig 7.2 Annual Report: Aspect Ratio 2**



# Choosing line chart aspect ratios

- 2: multi scale banking to 45 (2006)
  - frequency domain analysis to find ratios
    - FFT the data, convolve with Gaussian to smooth
  - find interesting spikes/ranges in power spectrum
    - cull nearby regions if similar, ensure overview
  - create trend curves (red) for each aspect ratio



[\[Multi-Scale Banking to 45 Degrees.  
Heer and Agrawala, Proc InfoVis  
2006\]](#)

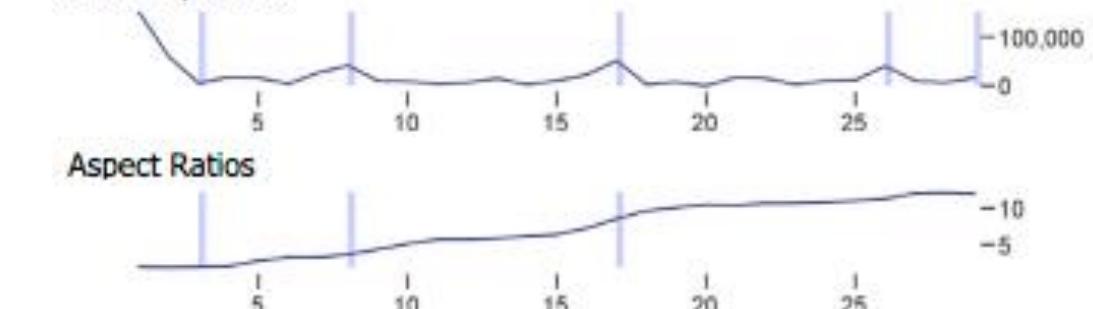
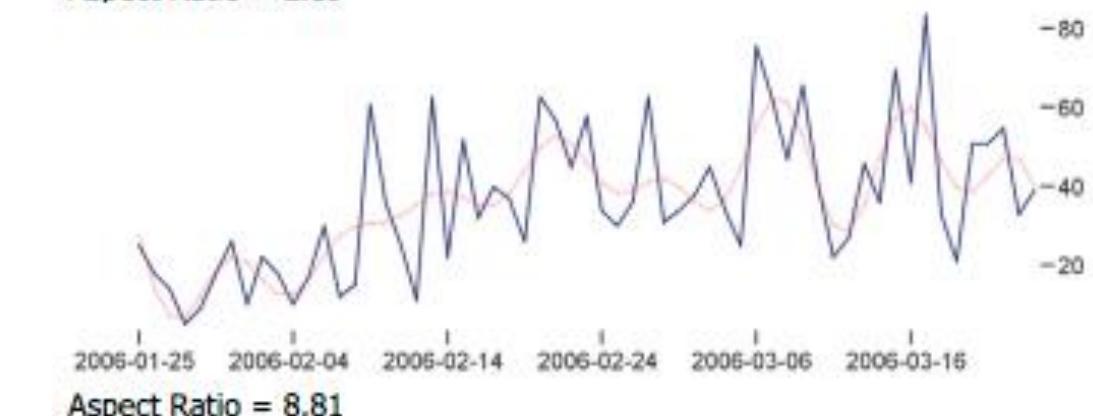
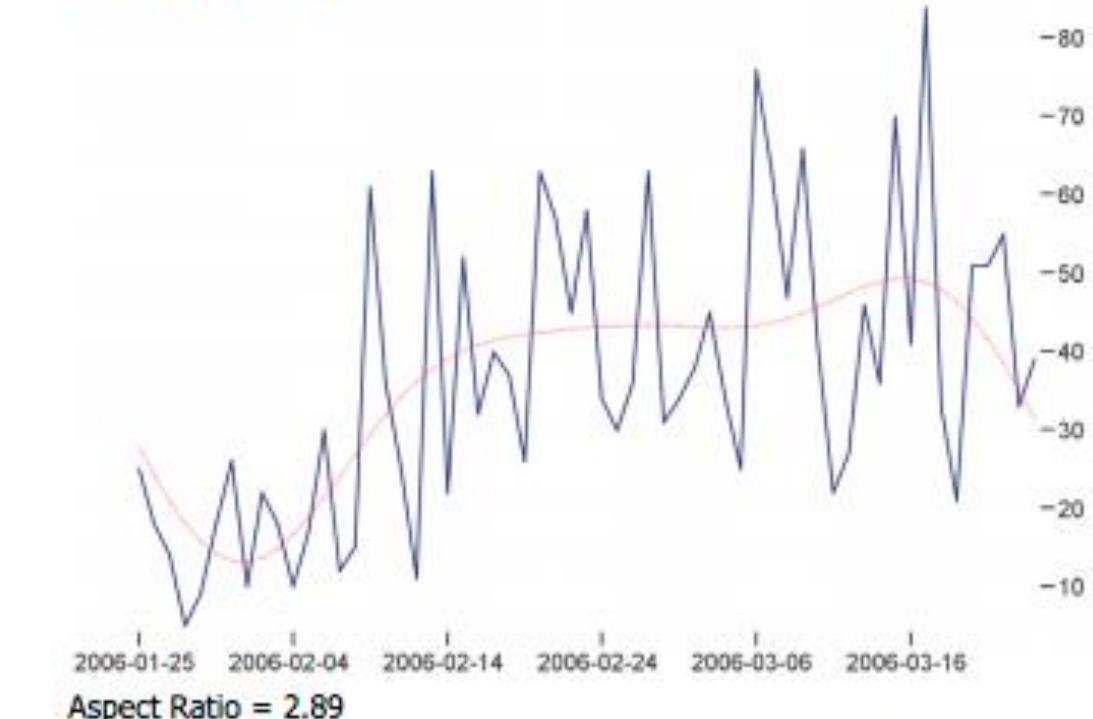
overall

weekly

daily

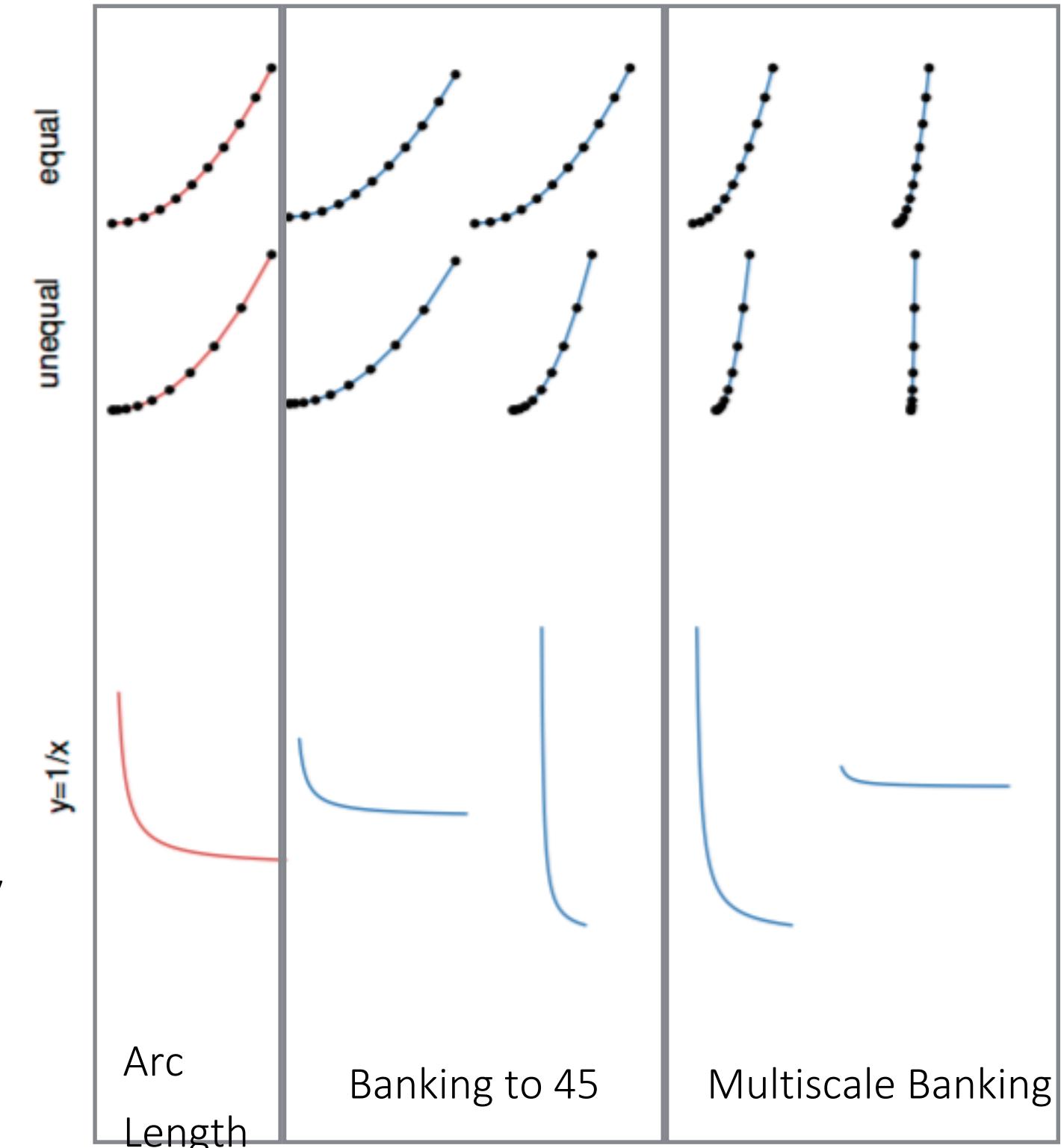
## Downloads of the prefuse toolkit

Aspect Ratio = 1.44



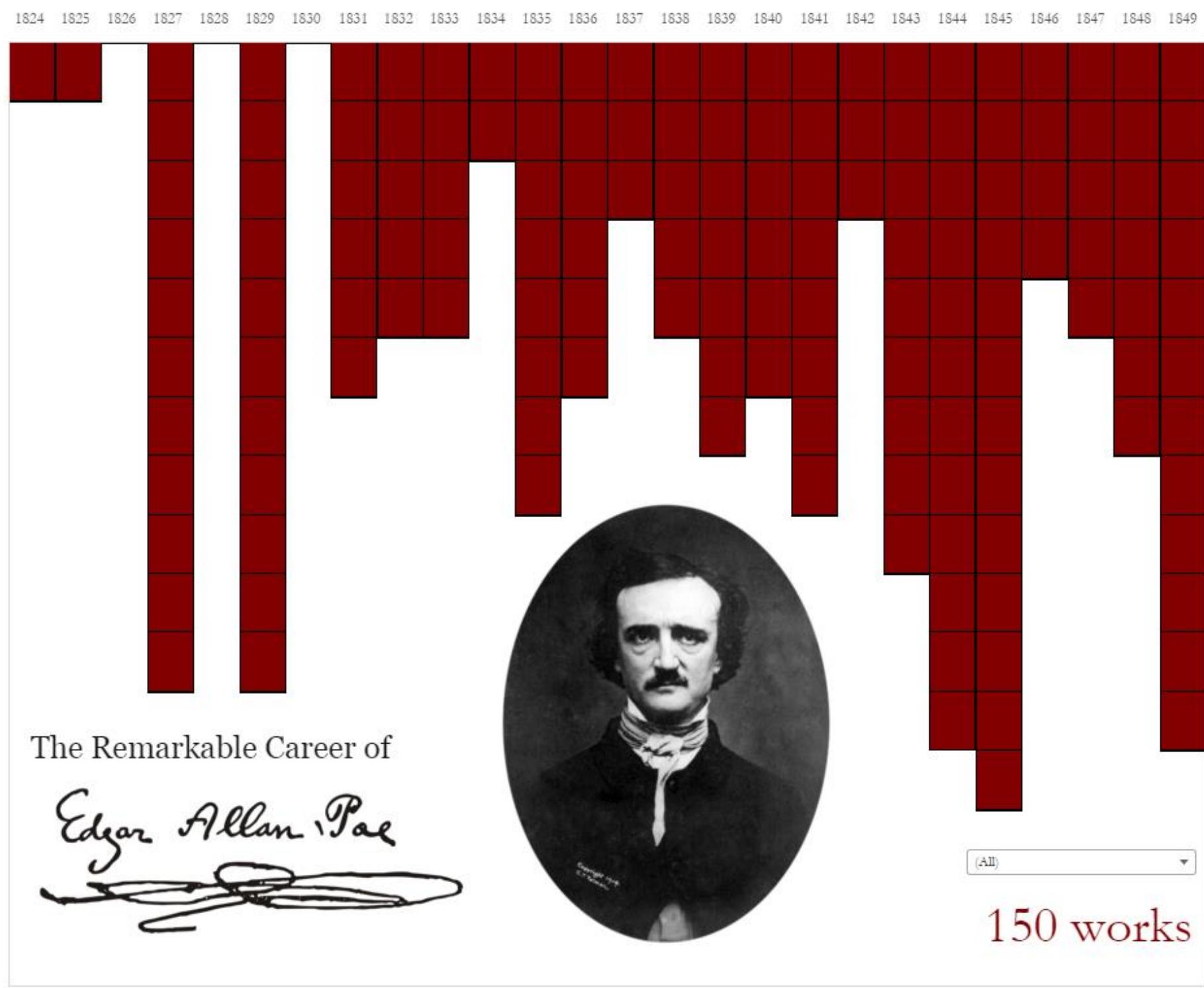
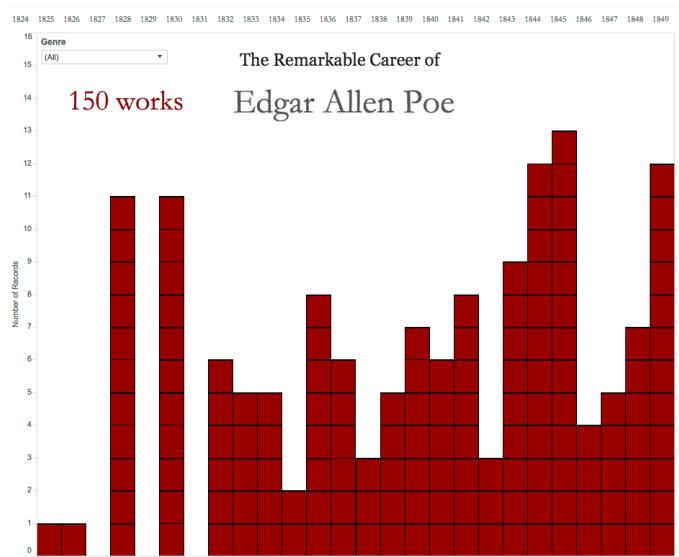
# Choosing line chart aspect ratios

- 3: arc length based aspect ratio (2011)
  - minimize the arc length of curve while keeping the area of the plot constant
  - parametrization and scale invariant
  - symmetry preserving
  - robust & fast to compute
- meta-points from this progression
  - young field; prescriptive advice changes rapidly
  - reasonable defaults required deep dive into perception meets math



# Breaking conventions

- presentation vs exploration
  - engaging/evocative
  - inverted y axis
    - blood drips down on Poe



[https://public.tableau.com/profile/ben.jones#/!](https://public.tableau.com/profile/ben.jones#/!/)

vizhome/EdgarAllanPoeBoring/EdgarAllenPoeBoring

[Slide inspired by Ben Jones]

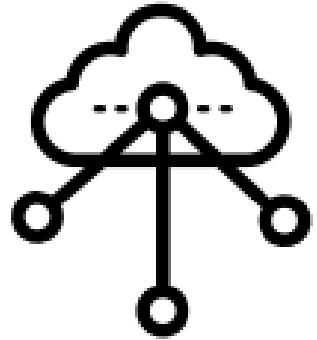
# Data and Task Abstraction Group Exercise – 20 minutes

- Dataset
- <https://www.thesquirrelcensus.com/data>

# Squirrel Design Sprint (20 minutes)

- Group size 3 – 4
- Choose a CEO (aka the Decider – has the final say on decisions)
- Choose a Manager (aka the Facilitator – responsible for keeping track of time and organizing the sprint process)
- Main Task
  - Using the provided dataset identify
  - Who the audience is
  - What questions you want to answer (identify tasks that users wish to perform)
  - Make sure that there is variability in your tasks, use Lec09 to make your questions as specific as possible
- Submit to Gradescope <https://help.gradescope.com/article/m5qz2xsnjy-student-add-group-members>
  - Your submission must include who the CEO is, who the Manager is and a grouping of your questions based on the Task Abstraction Model (Analyze, Search, Query)
  - Also post to jamboard (see EdStem for link) don't add names

Map



## Visualization Theory:

- User-Centered Design
- Data Types
- What is the question?
- Who is the audience?
- What is the data?

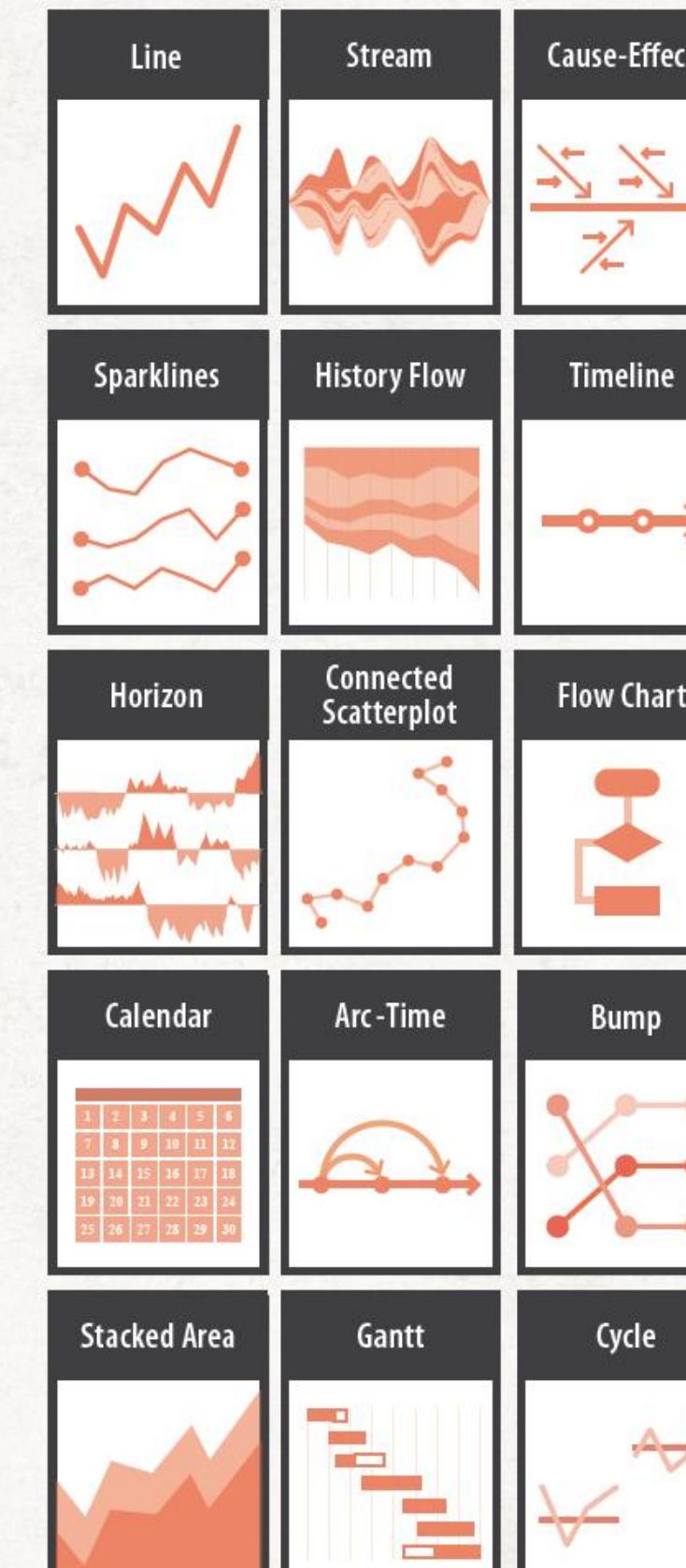
## COMPARING CATEGORIES

Compare values across categories



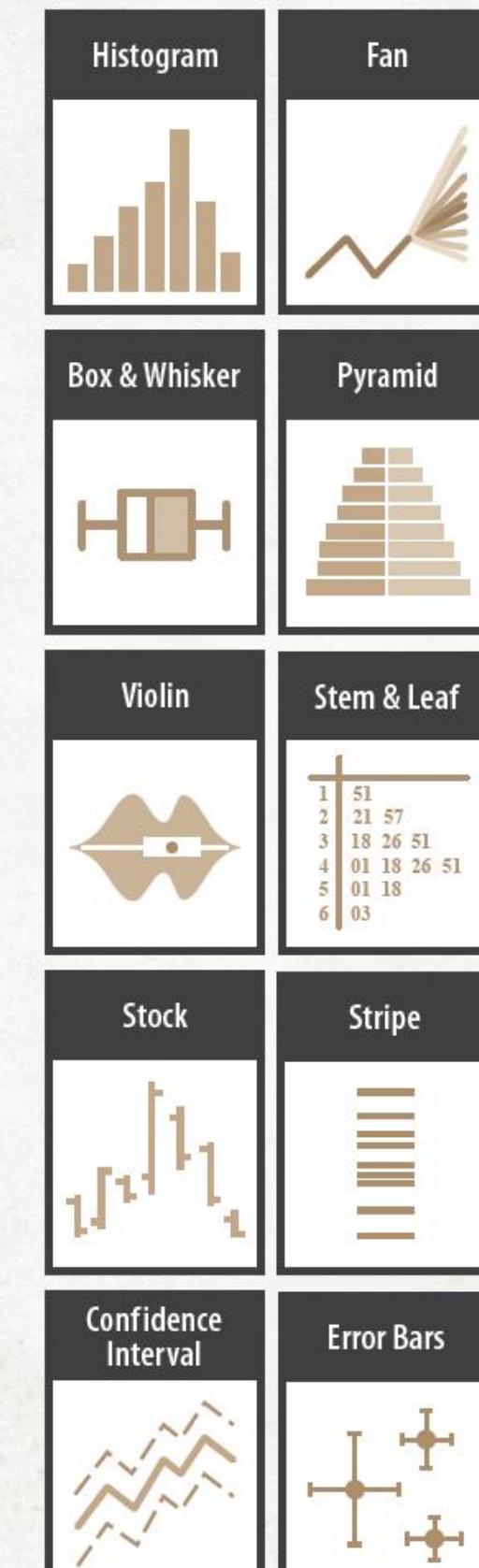
## TIME

Track changes over time



## DISTRIBUTION

Representation of the distribution of data



## GEOSPATIAL

Relates data to its geography

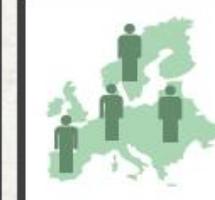
Map



Flow Map



Icon Map



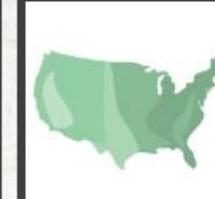
Choropleth



Map with Columns



Isopleth



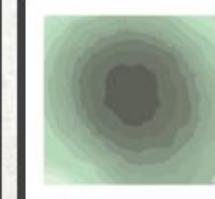
Cartogram



Map with Pie Charts



Contour



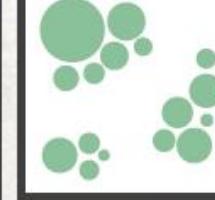
Non-Contiguous Cartogram



Bubble Map



Dorling Map



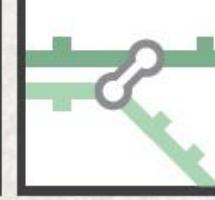
Connection Map



Point Map



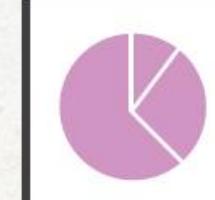
Subway Map



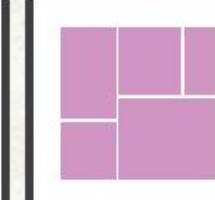
## PART-TO-WHOLE

Relates the part of a variable to its total

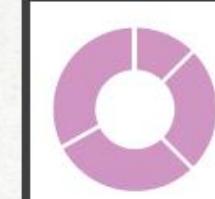
Pie



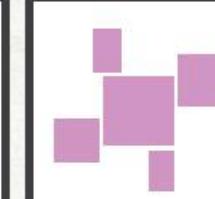
Treemap



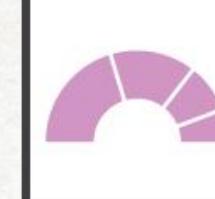
Donut



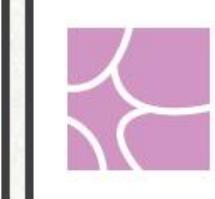
Square Cloud



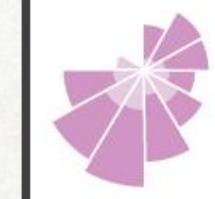
Arc



Voronoi



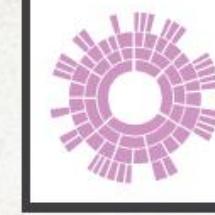
Nightingale



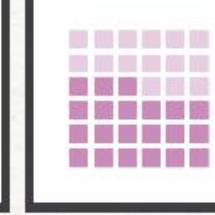
Triangle Treemap



Sunburst



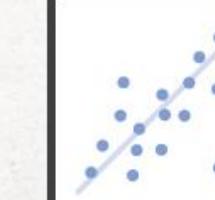
Waffle



## RELATIONSHIP

Illustrates correlations or relationships between variables

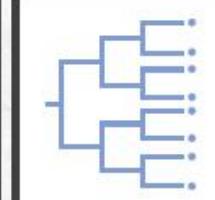
Scatterplot



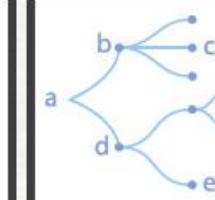
Arc-Connection



Dendrogram



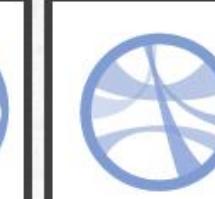
Word Tree



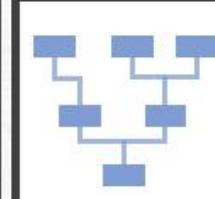
Circle Packing



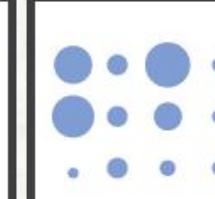
Chord



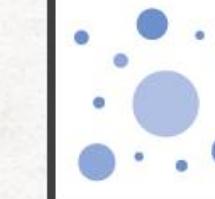
Tree



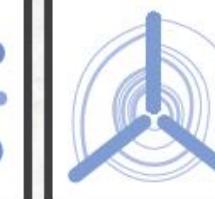
Correlation Matrix



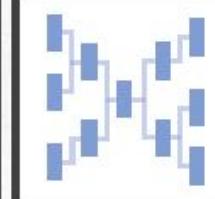
Bubble



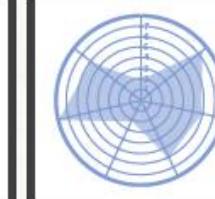
Hive



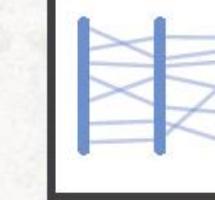
Double Tree



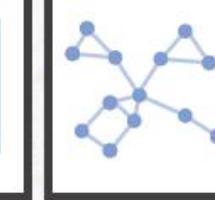
Radar



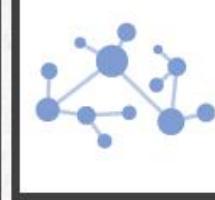
Parallel Coordinates



Force-Directed



Network



Venn Diagram

