

Week 2: Tasks & Data, Tables

Tamara Munzner

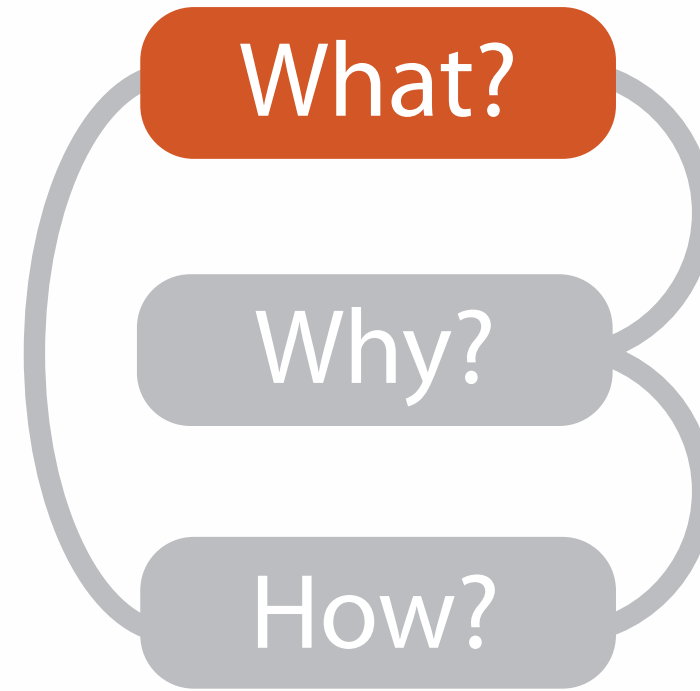
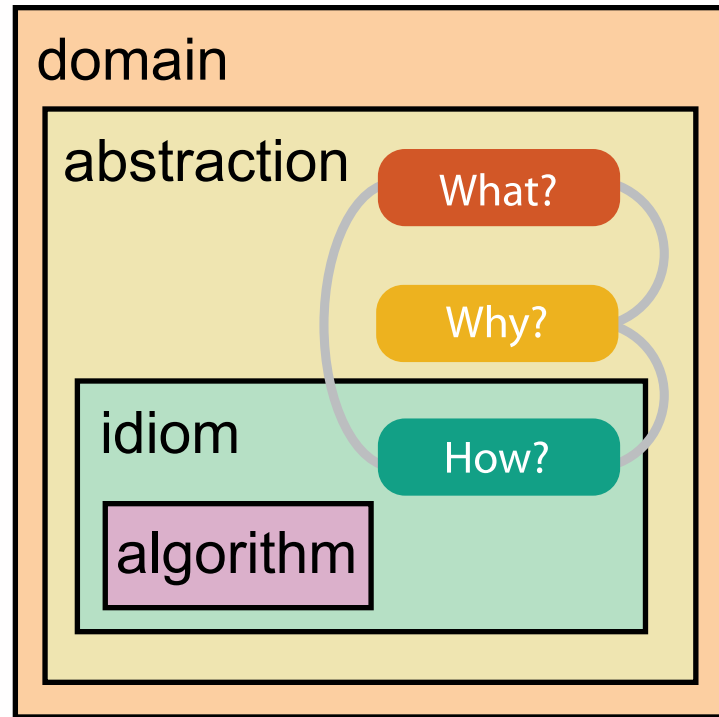
Department of Computer Science
University of British Columbia

JRNL 520M, Special Topics in Contemporary Journalism: Visualization for Journalists

Week 2: 22 September 2015

<http://www.cs.ubc.ca/~tmm/courses/journ15>

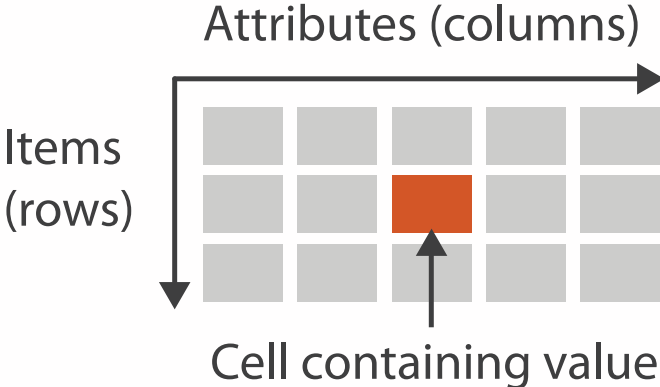
Data abstraction: What



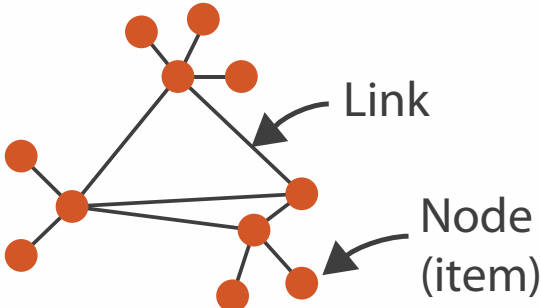
Data Abstraction

➔ Dataset Types

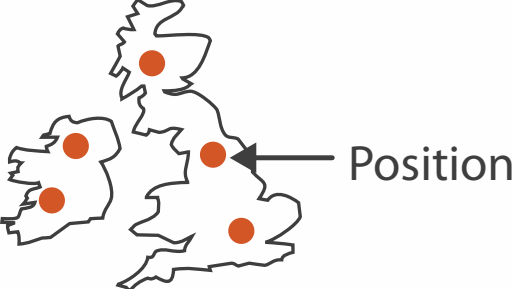
➔ Tables



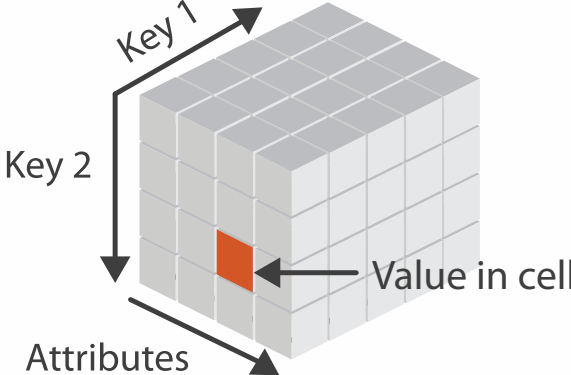
➔ Networks



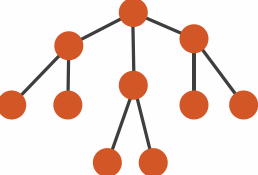
➔ Spatial



➔ *Multidimensional Table*



➔ *Trees*



Attribute types

➔ Attribute Types

➔ Categorical

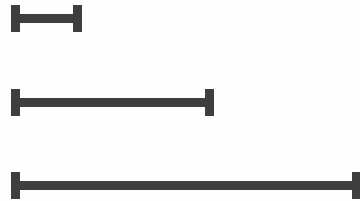


➔ Ordered

➔ *Ordinal*



➔ *Quantitative*



➔ Ordering Direction

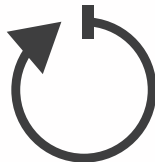
➔ Sequential



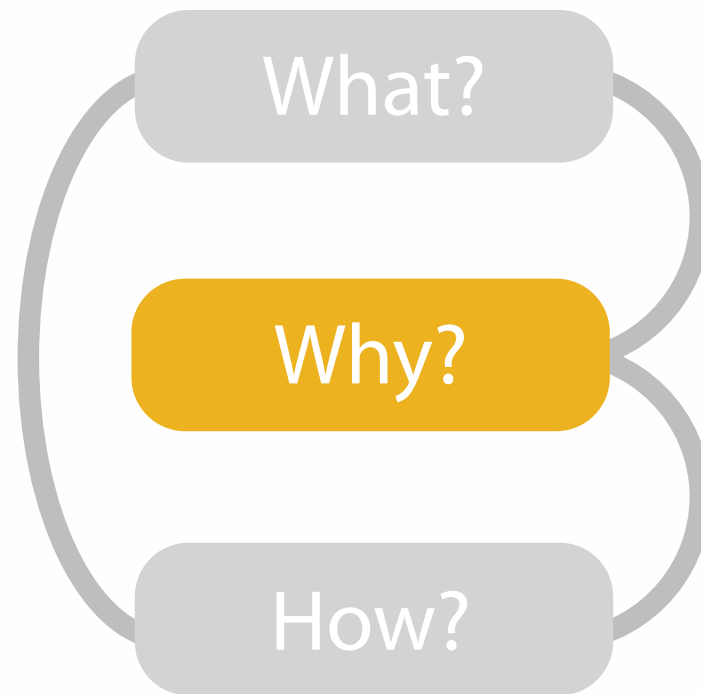
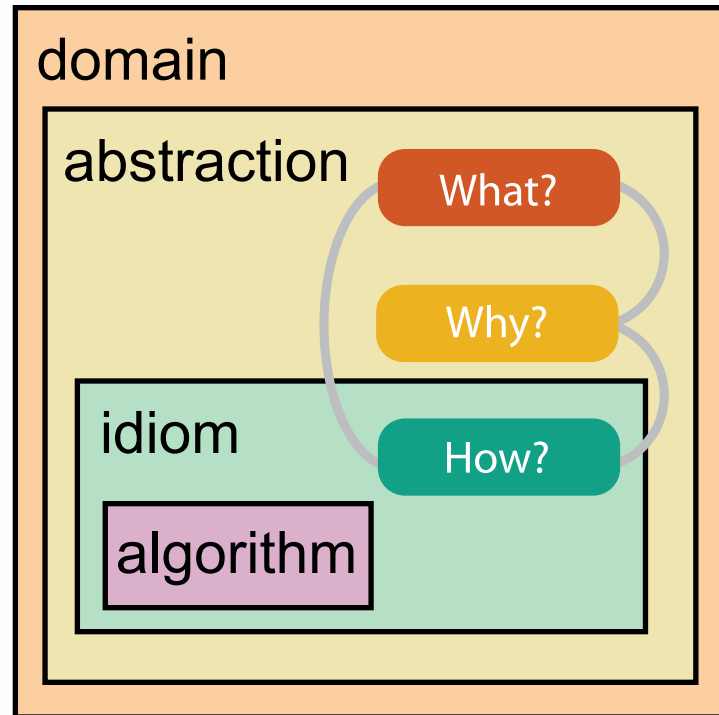
➔ Diverging



➔ Cyclic



Tasks abstraction: Why









Why?




👉 Actions

🎯 Targets




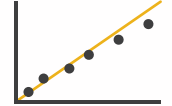
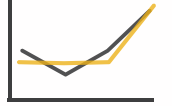
➔ **Analyze**

- ➔ Consume
 - ➔ Discover 
 - ➔ Present 
 - ➔ Enjoy 
- ➔ Produce
 - ➔ Annotate 
 - ➔ Record 
 - ➔ Derive 





➔ **All Data**

- ➔ Trends 
- ➔ Outliers 
- ➔ Features 


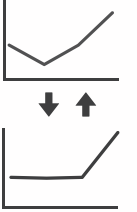

➔ **Attributes**

- ➔ One
 - ➔ Distribution 
 - ➔ Extremes 
- ➔ Many
 - ➔ Dependency 
 - ➔ Correlation 
 - ➔ Similarity 

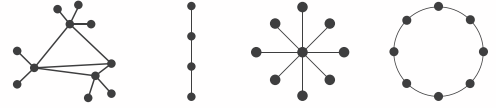

➔ **Search**

	Target known	Target unknown
Location known	 <i>Lookup</i>	 <i>Browse</i>
Location unknown	 <i>Locate</i>	 <i>Explore</i>


➔ **Query**

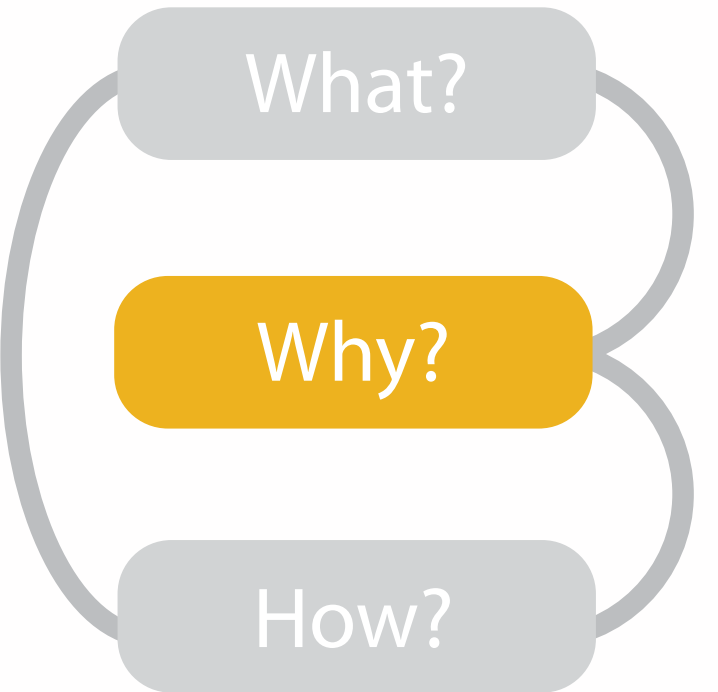
- ➔ Identify 
- ➔ Compare 
- ➔ Summarize 

➔ **Network Data**

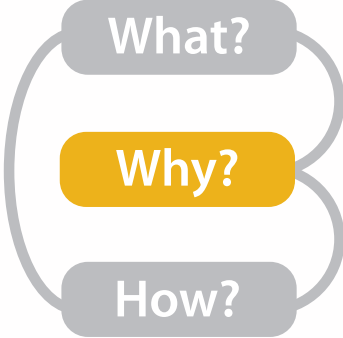
- ➔ Topology 
- ➔ Paths 

➔ **Spatial Data**

- ➔ Shape 



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



Actions: Analyze

- consume
 - discover vs present
 - classic split
 - aka explore vs explain
 - enjoy
- produce
 - newcomer
 - aka casual, social
- produce
 - annotate, record
 - derive
 - crucial design choice

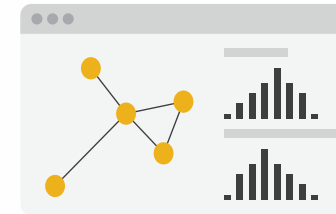
→ Analyze

→ Consume

→ Discover



→ Present

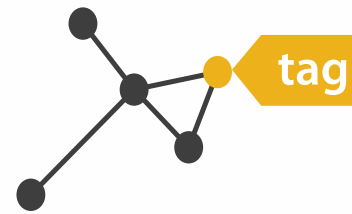


→ Enjoy



→ Produce

→ Annotate



→ Record

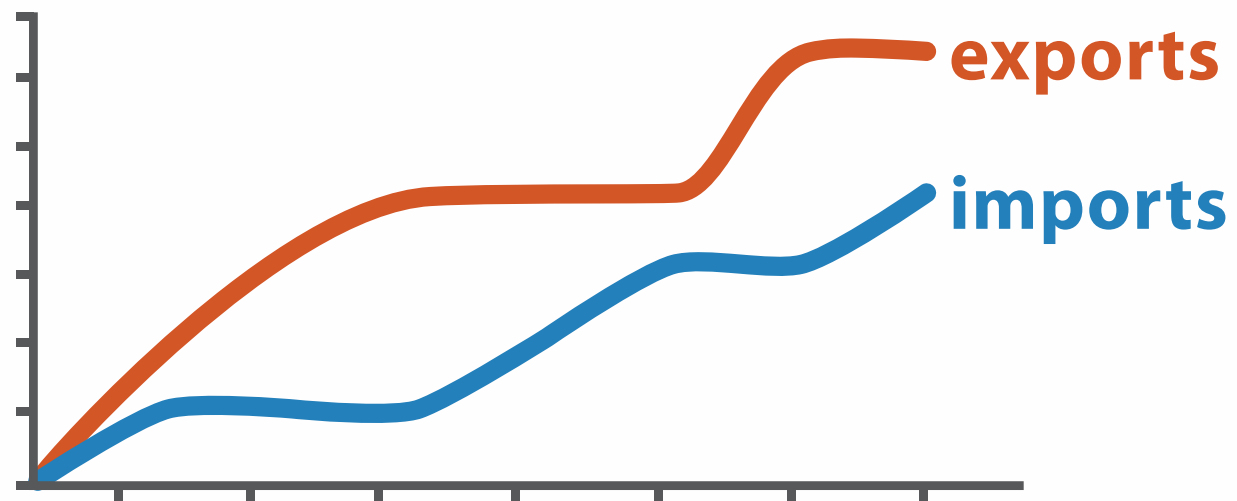


→ Derive

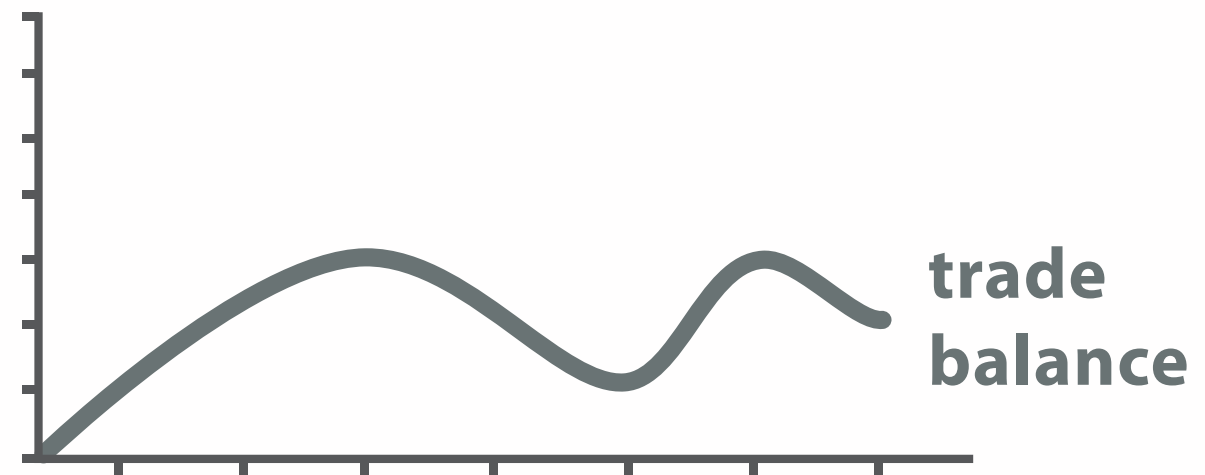


Derive

- don't just draw what you're given!
 - decide what the right thing to show is
 - create it with a series of transformations from the original dataset
 - draw that
- one of the four major strategies for handling complexity



Original Data







$$\text{trade balance} = \text{exports} - \text{imports}$$

Derived Data

Actions: Search, query

- what does user know?
 - target, location
- how much of the data matters?
 - one, some, all

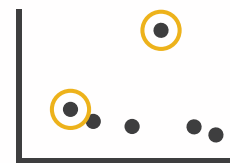
➔ Search

	Target known	Target unknown
Location known	 <i>Lookup</i>	 <i>Browse</i>
Location unknown	 <i>Locate</i>	 <i>Explore</i>

- independent choices for each of these three levels
 - analyze, search, query
 - mix and match

➔ Query

➔ Identify



➔ Compare



↓ ↑



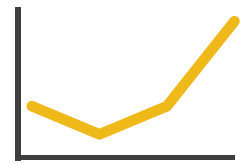
➔ Summarize



Why: Targets

→ ALL DATA

→ Trends



→ Outliers



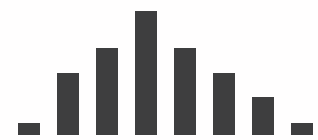
→ Features



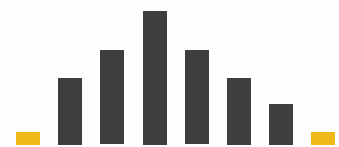
→ ATTRIBUTES

→ One

→ *Distribution*



↓ *Extremes*

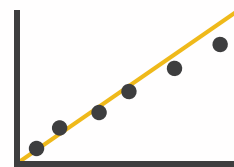


→ Many

→ *Dependency*



→ *Correlation*

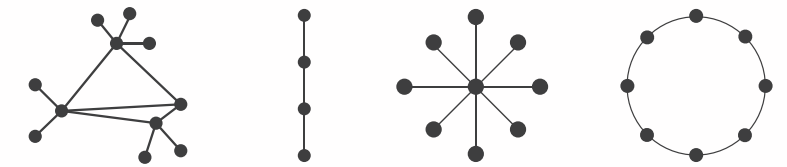


→ *Similarity*

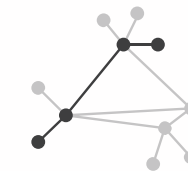


→ NETWORK DATA

→ Topology

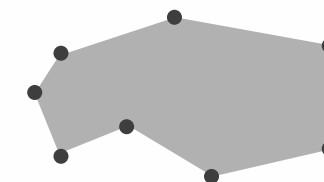


→ *Paths*



→ SPATIAL DATA

→ Shape

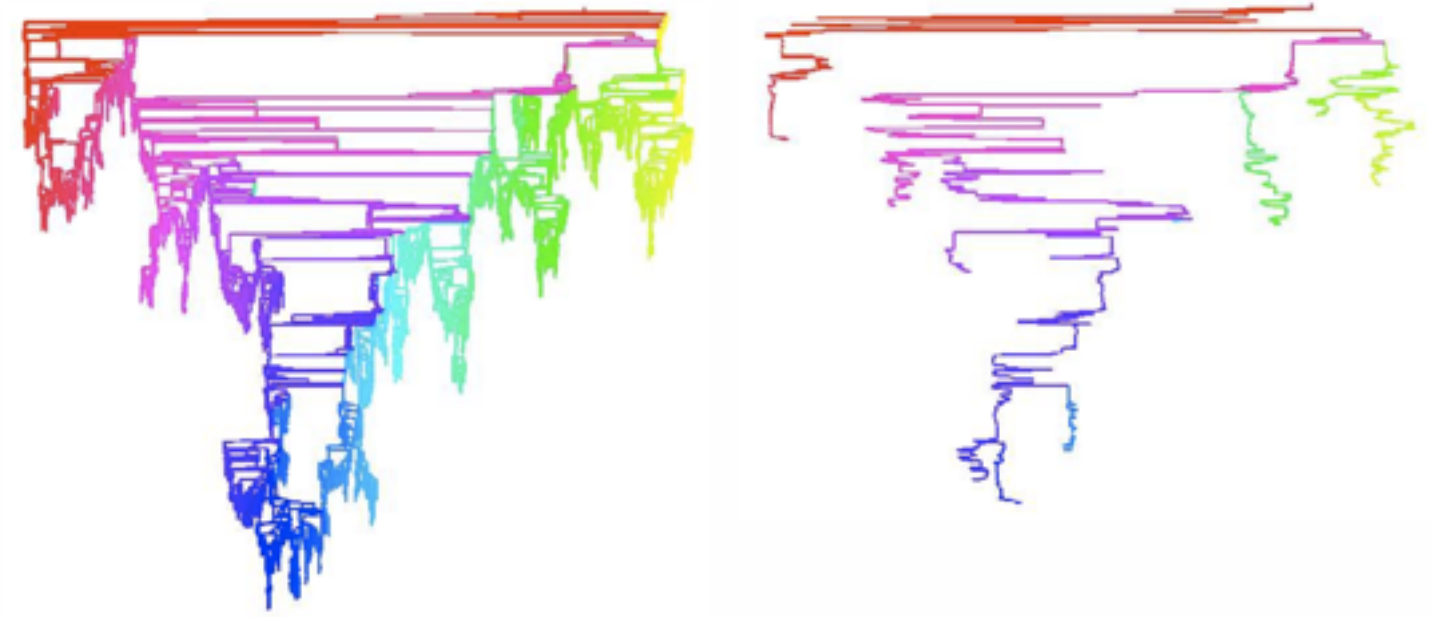


Analysis example: Derive one attribute

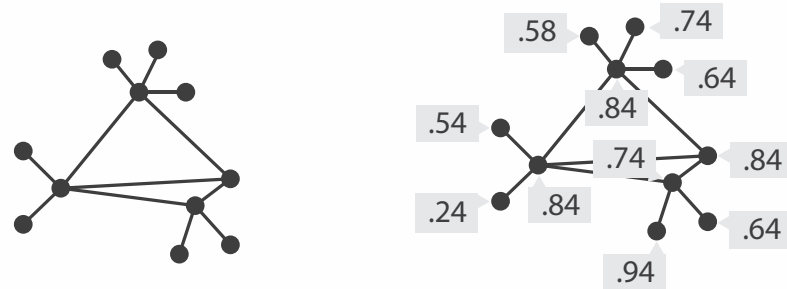
- Strahler number

- centrality metric for trees/networks
- derived quantitative attribute
- draw top 5K of 500K for good skeleton

[Using Strahler numbers for real time visual exploration of huge graphs. Auber. Proc. Intl. Conf. Computer Vision and Graphics, pp. 56–69, 2002.]



Task 1



In
Tree

➔

Out
Quantitative
attribute on nodes

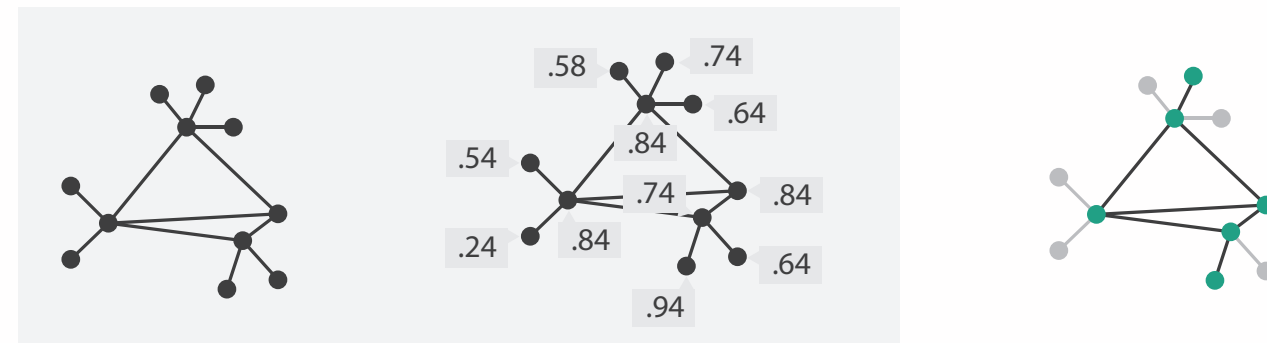
What?

- ➔ In Tree
- ➔ Out Quantitative attribute on nodes

Why?

- ➔ Derive

Task 2



In
Tree

+

In
Quantitative
attribute on nodes

➔

Out
Filtered Tree
Removed
unimportant parts

What?

- ➔ In Tree
- ➔ In Quantitative attribute on nodes
- ➔ Out Filtered Tree

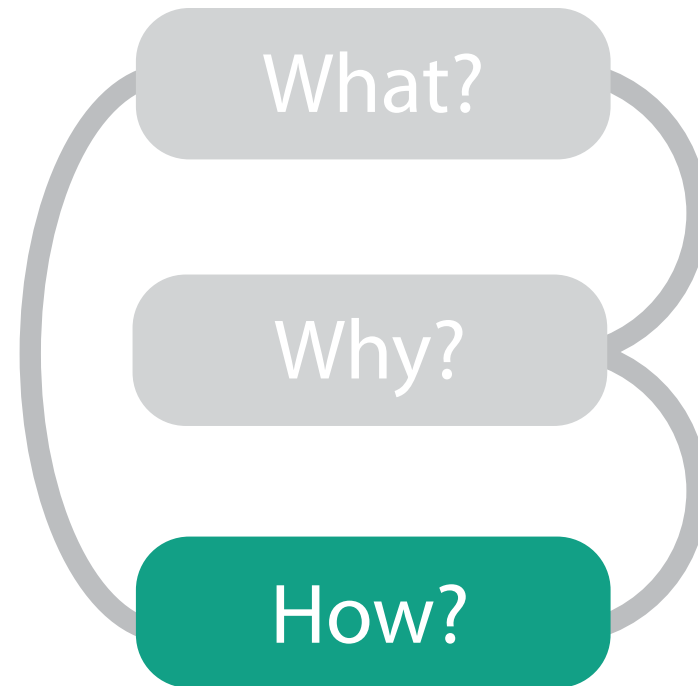
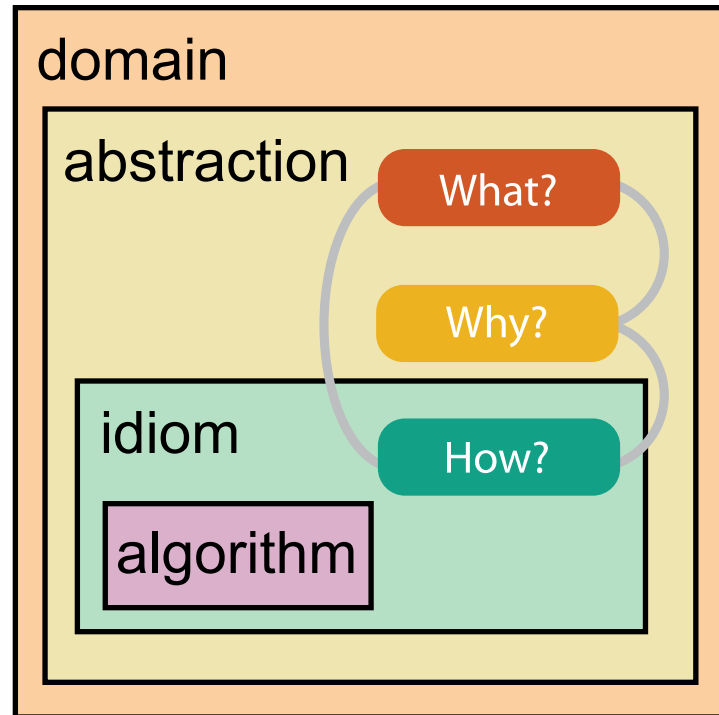
Why?

- ➔ Summarize
- ➔ Topology

How?

- ➔ Reduce
- ➔ Filter

Visual encoding and interaction idiom: How



How?

Encode

→ Arrange

→ Express



→ Order



→ Use



→ Separate



→ Align



→ Map

from **categorical** and **ordered** attributes

→ Color

→ Hue



→ Saturation



→ Luminance



→ Size, Angle, Curvature, ...



→ Shape



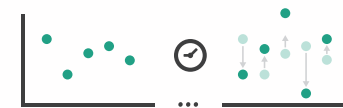
→ Motion

Direction, Rate, Frequency, ...

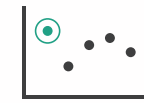


Manipulate

→ Change



→ Select



→ Navigate

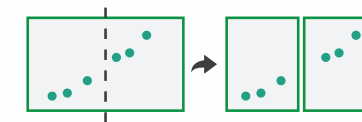


Facet

→ Juxtapose



→ Partition



→ Superimpose



Reduce

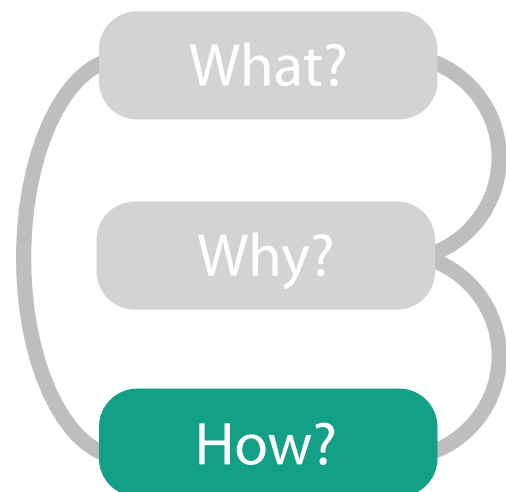
→ Filter



→ Aggregate



→ Embed



How?

Encode

→ Arrange

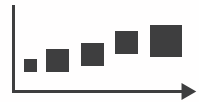
→ Express



→ Separate



→ Order



→ Align



Encode tables: Arrange space

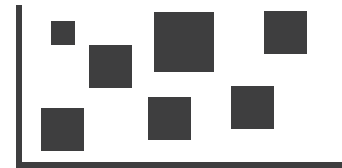
Encode

➔ Arrange

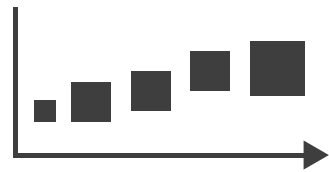
➔ Express



➔ Separate



➔ Order



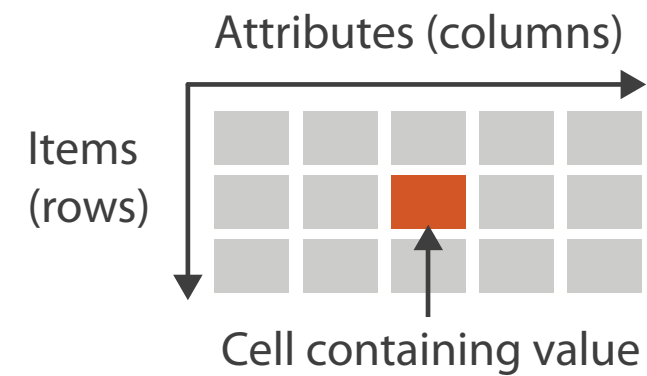
➔ Align



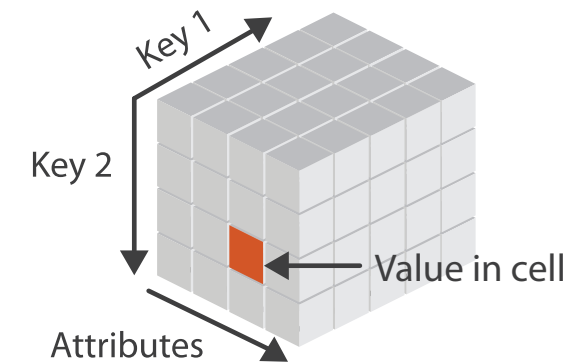
Keys and values

- **key**
 - independent attribute
 - used as unique index to look up items
 - simple tables: 1 key
 - multidimensional tables: multiple keys
- **value**
 - dependent attribute, value of cell
- **classify arrangements by key count**
 - 0, 1, 2, many...

→ Tables



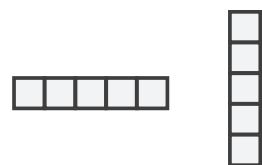
→ *Multidimensional Table*



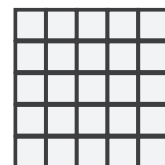
⊕ Express Values



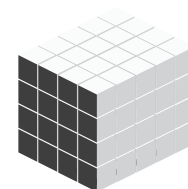
→ 1 Key
List



→ 2 Keys
Matrix



→ 3 Keys
Volume



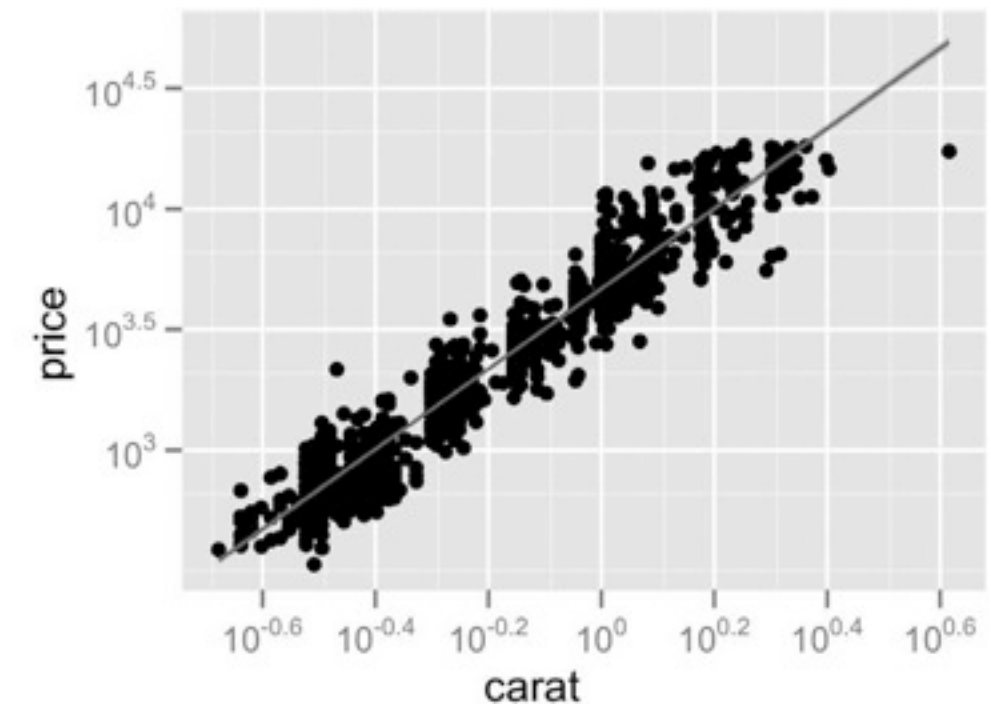
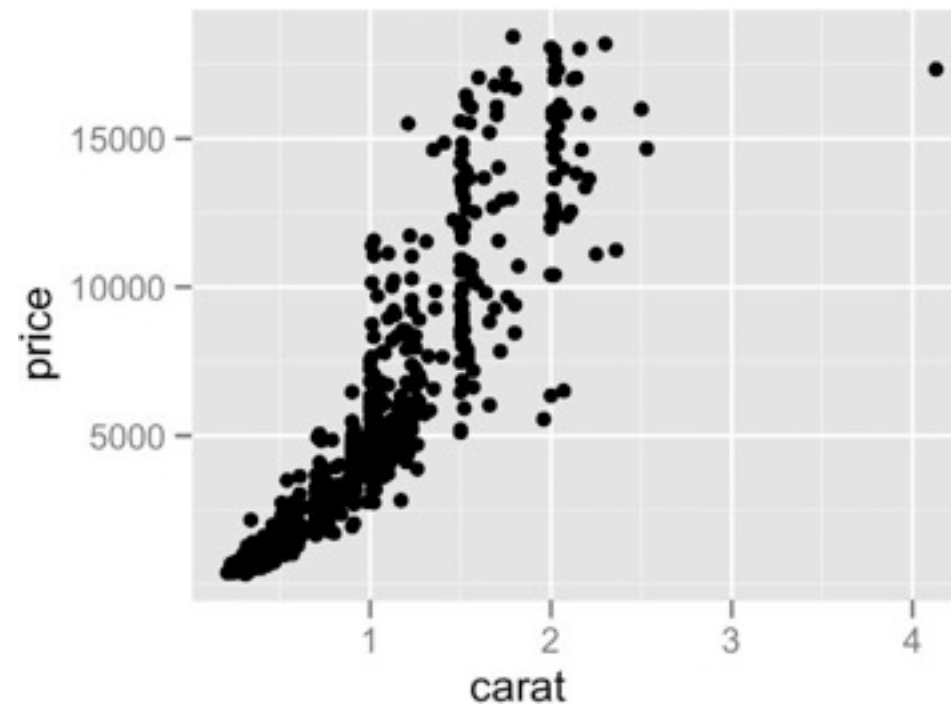
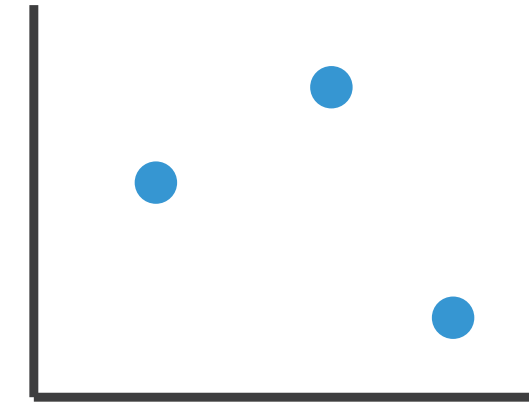
→ Many Keys
Recursive Subdivision



Idiom: scatterplot

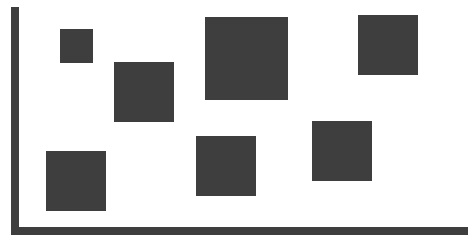
- **express** values
 - quantitative attributes
- no keys, only values
 - data
 - 2 quant attribs
 - mark: points
 - channels
 - horiz + vert position
 - tasks
 - find trends, outliers, distribution, correlation, clusters
 - scalability
 - hundreds of items

⇒ Express Values

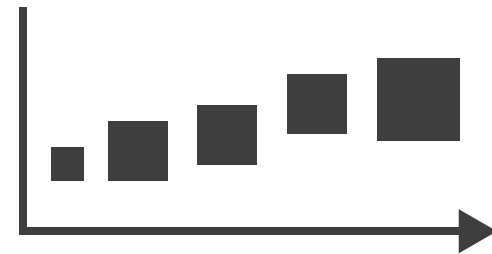


Some keys: Categorical regions

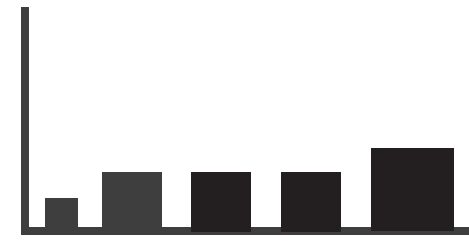
→ Separate



→ Order

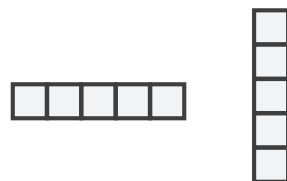


→ Align

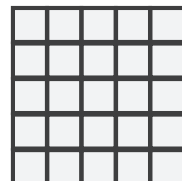


- **regions**: contiguous bounded areas distinct from each other
 - using space to **separate** (proximity)
 - following expressiveness principle for categorical attributes
- use ordered attribute to **order** and **align** regions

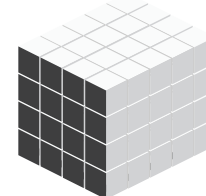
→ 1 Key
List



→ 2 Keys
Matrix



→ 3 Keys
Volume

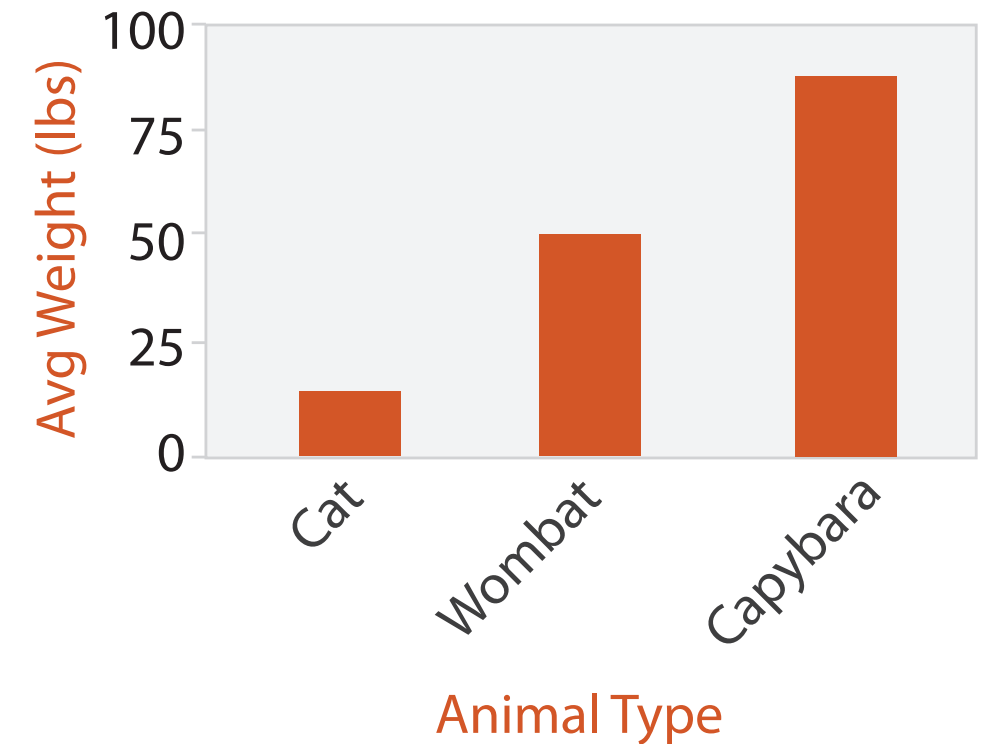
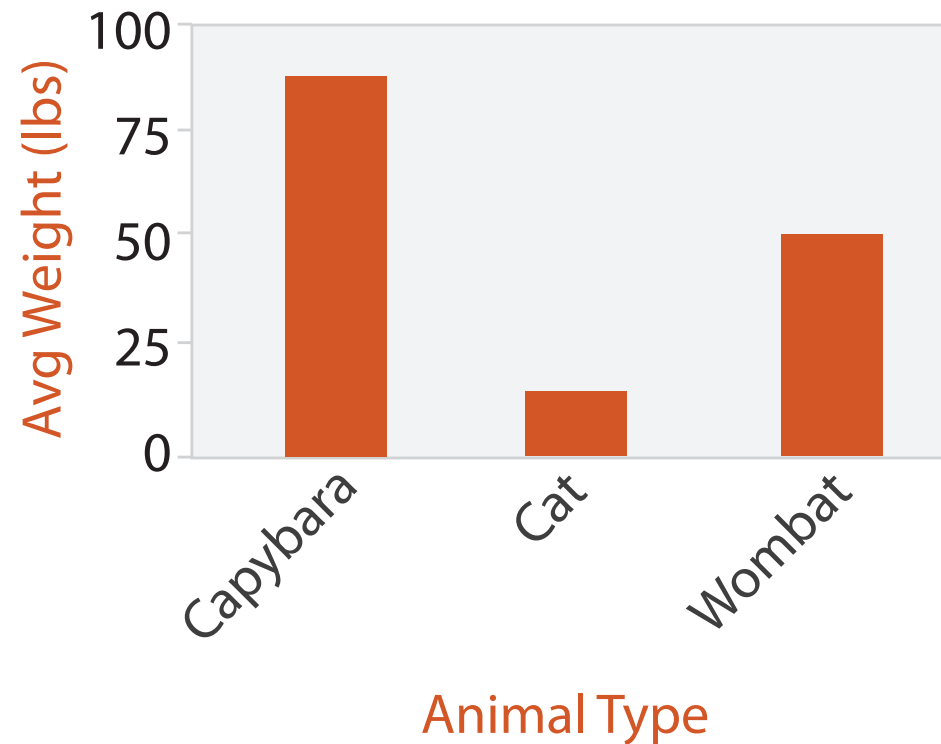


→ Many Keys
Recursive Subdivision



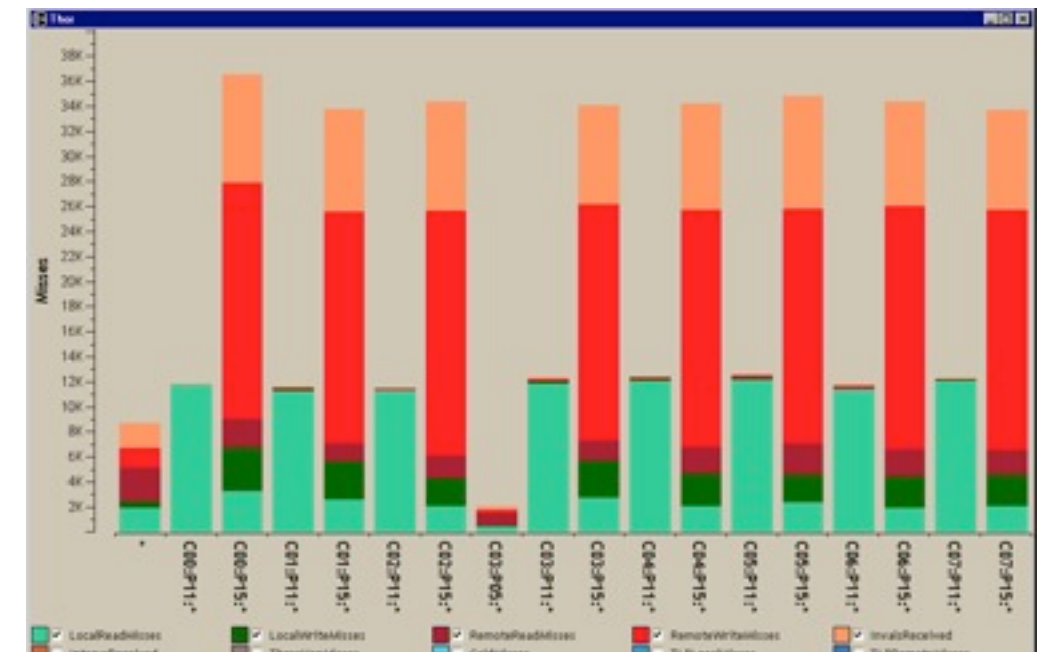
Idiom: bar chart

- one key, one value
 - data
 - 1 categ attrib, 1 quant attrib
 - mark: lines
 - channels
 - length to express quant value
 - spatial regions: one per mark
 - separated horizontally, aligned vertically
 - ordered by quant attrib
 - » by label (alphabetical), by length attrib (data-driven)
 - task
 - compare, lookup values
 - scalability
 - dozens to hundreds of levels for key attrib



Idiom: stacked bar chart

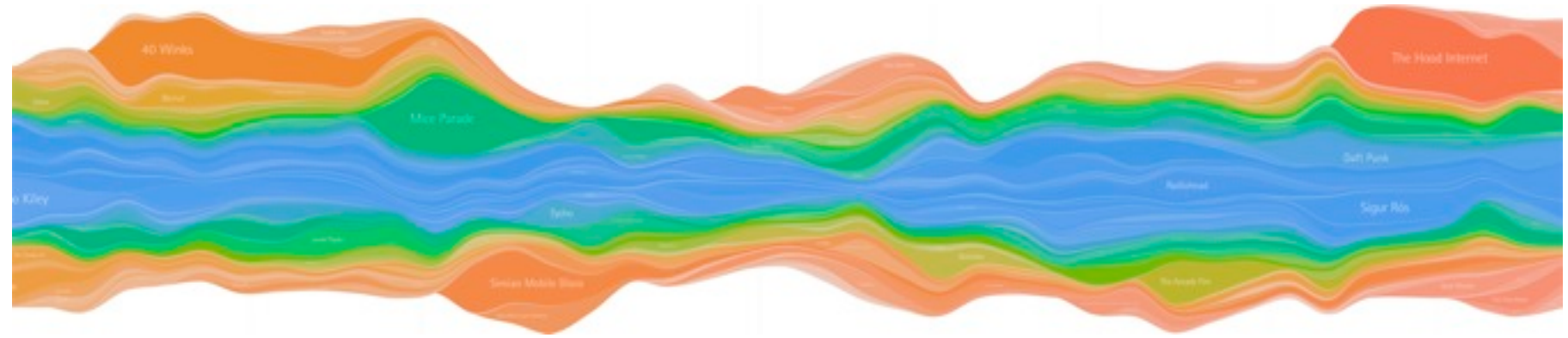
- one more key
 - data
 - 2 categ attrib, 1 quant attrib
 - mark: vertical stack of line marks
 - **glyph**: composite object, internal structure from multiple marks
 - channels
 - length and color hue
 - spatial regions: one per glyph
 - aligned: full glyph, lowest bar component
 - unaligned: other bar components
 - task
 - part-to-whole relationship
 - scalability
 - several to one dozen levels for stacked attrib



[Using Visualization to Understand the Behavior of Computer Systems. Bosch. Ph.D. thesis, Stanford Computer Science, 2001.]

Idiom: streamgraph

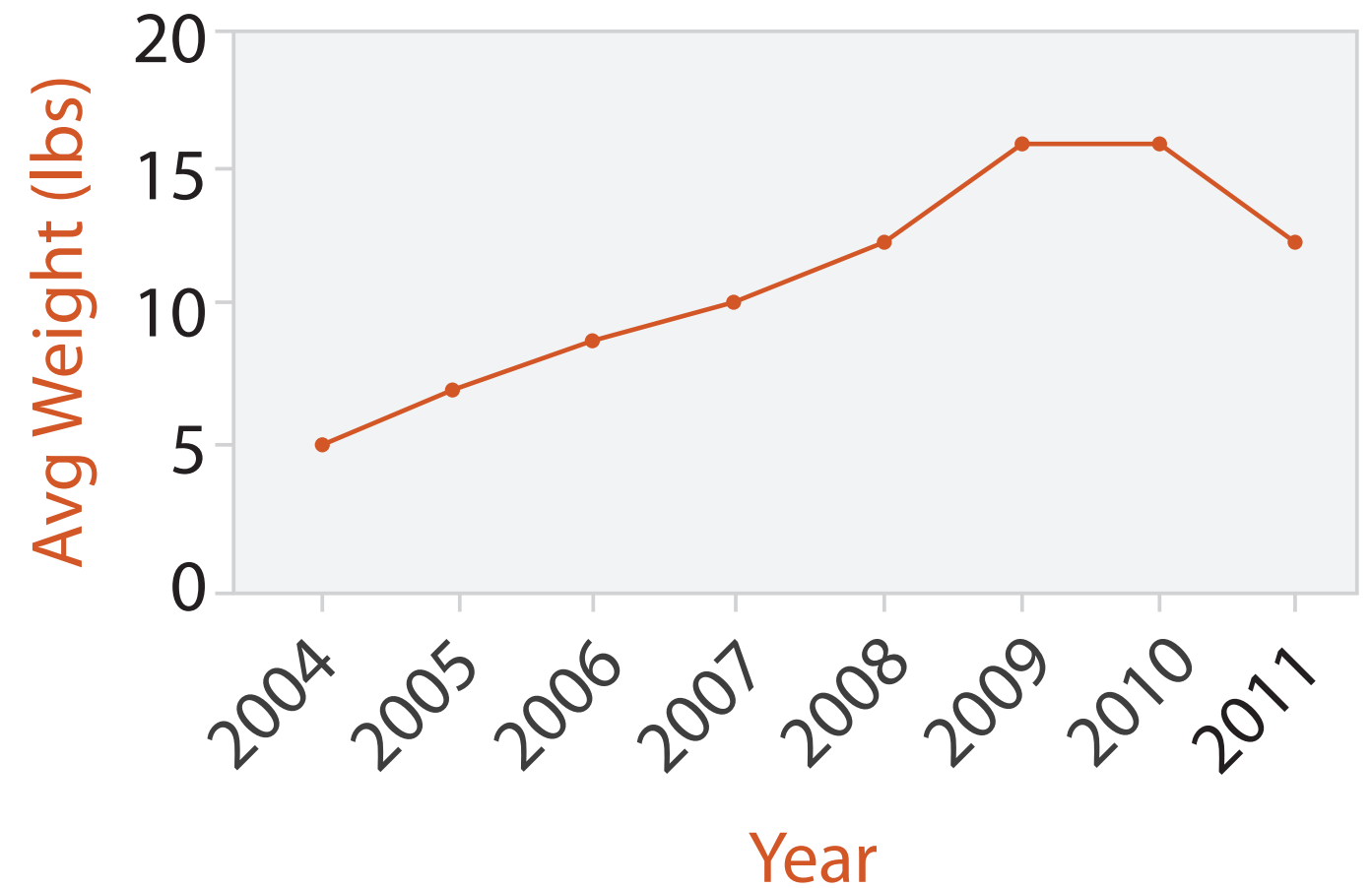
- generalized stacked graph
 - emphasizing horizontal continuity
 - vs vertical items
 - data
 - | categ key attrib (artist)
 - | ordered key attrib (time)
 - | quant value attrib (counts)
 - derived data
 - geometry: layers, where height encodes counts
 - | quant attrib (layer ordering)
 - scalability
 - hundreds of time keys
 - dozens to hundreds of artist keys
 - more than stacked bars, since most layers don't extend across whole chart



[Stacked Graphs Geometry & Aesthetics. Byron and Wattenberg. IEEE Trans. Visualization and Computer Graphics (Proc. InfoVis 2008) 14(6): 1245–1252, (2008).]

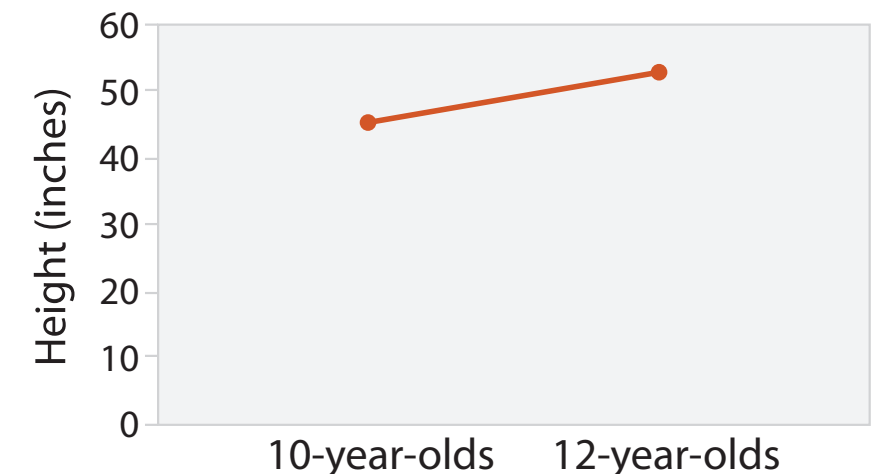
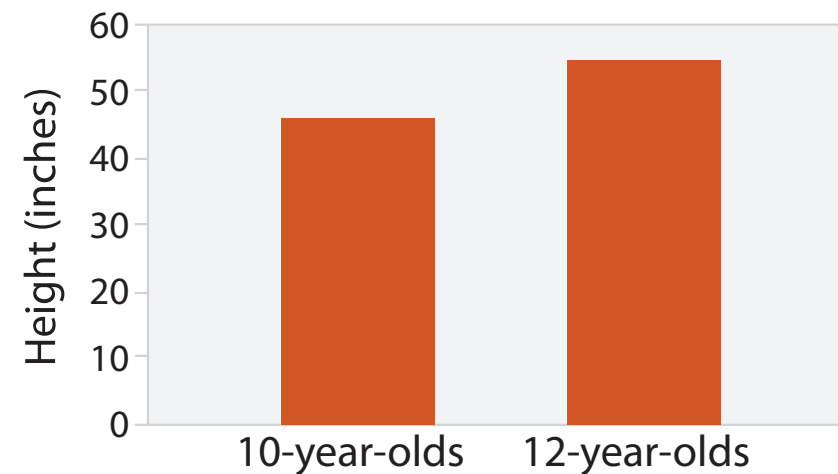
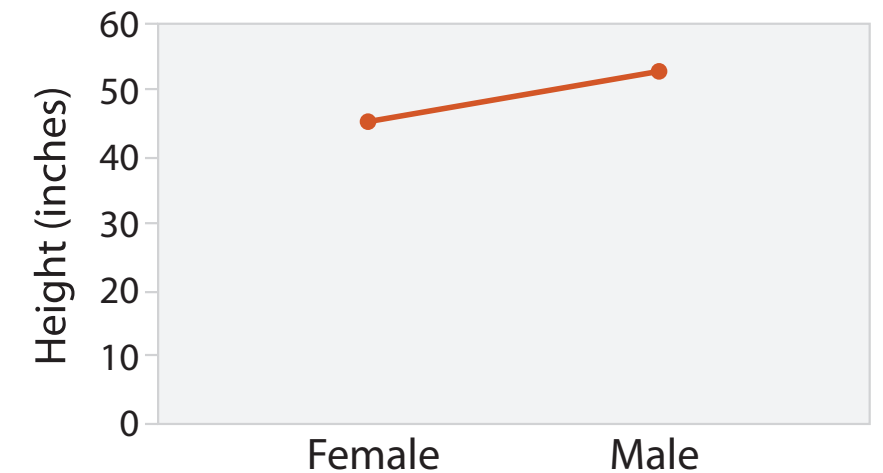
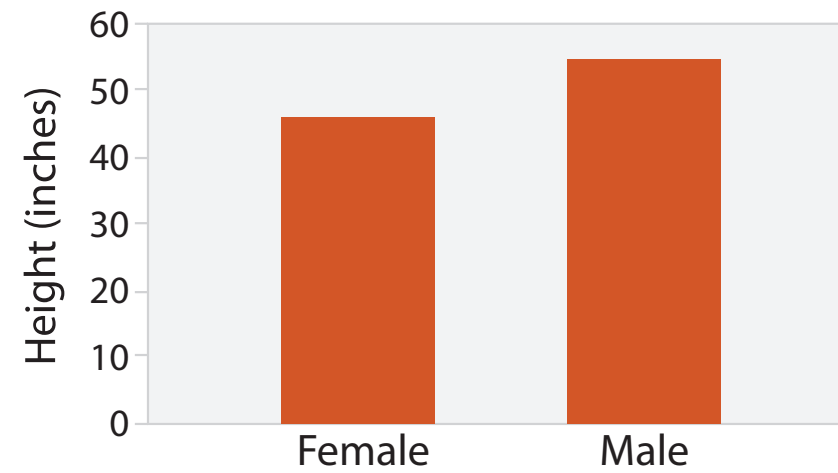
Idiom: **line chart**

- one key, one value
 - data
 - 2 quant attribs
 - mark: points
 - line connection marks between them
 - channels
 - aligned lengths to express quant value
 - separated and ordered by key attrib into horizontal regions
 - task
 - find trend
 - connection marks emphasize ordering of items along key axis by explicitly showing relationship between one item and the next



Choosing bar vs line charts

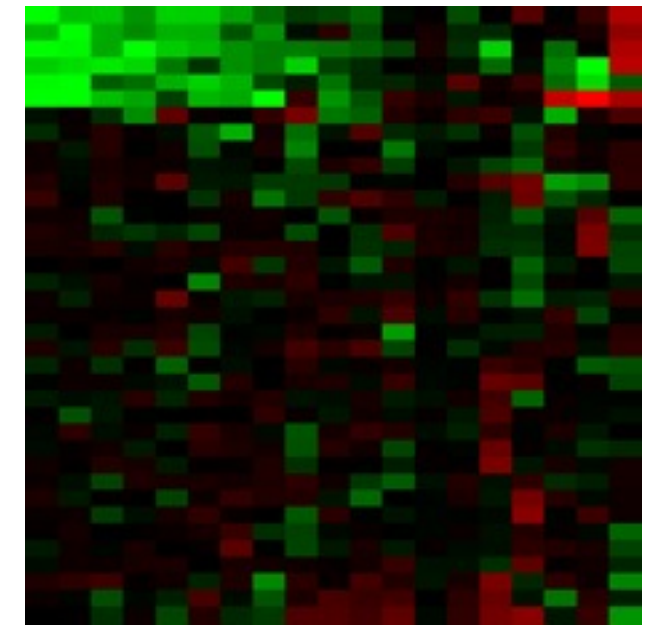
- depends on type of key attrib
 - bar charts if categorical
 - line charts if ordered
- do not use line charts for categorical key attribs
 - violates expressiveness principle
 - implication of trend so strong that it overrides semantics!
 - “The more male a person is, the taller he/she is”



after [Bars and Lines: A Study of Graphic Communication. Zacks and Tversky. Memory and Cognition 27:6 (1999), 1073–1079.]

Idiom: heatmap

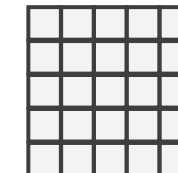
- two keys, one value
 - data
 - 2 categ attribs (gene, experimental condition)
 - 1 quant attrib (expression levels)
 - marks: area
 - separate and align in 2D matrix
 - indexed by 2 categorical attributes
 - channels
 - color by quant attrib
 - (ordered diverging colormap)
 - task
 - find clusters, outliers
 - scalability
 - 1M items, 100s of categ levels, ~10 quant attrib levels



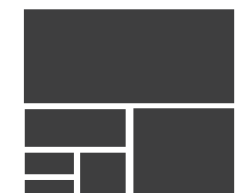
→ 1 Key
List



→ 2 Keys
Matrix

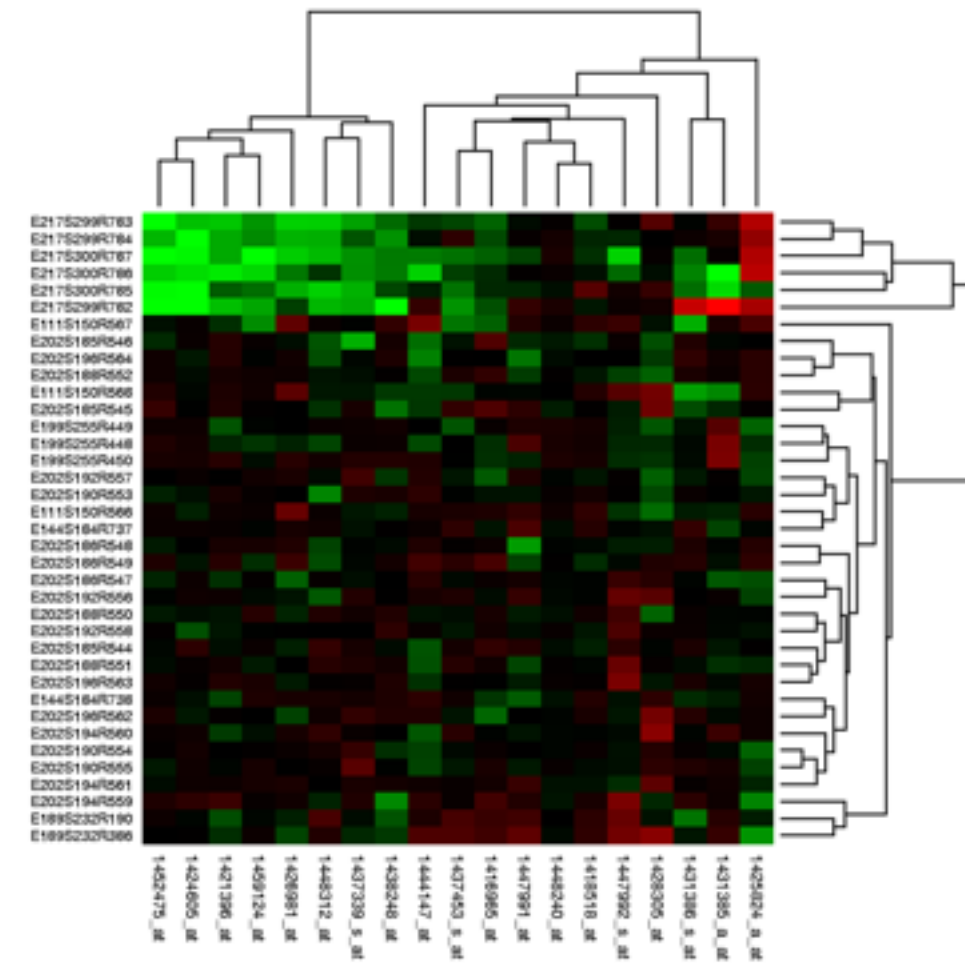


→ Many Keys
Recursive Subdivision



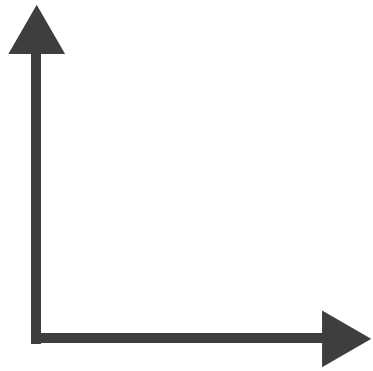
Idiom: cluster heatmap

- in addition
 - derived data
 - 2 cluster hierarchies
 - dendrogram
 - parent-child relationships in tree with connection line marks
 - leaves aligned so interior branch heights easy to compare
 - heatmap
 - marks (re-)ordered by cluster hierarchy traversal

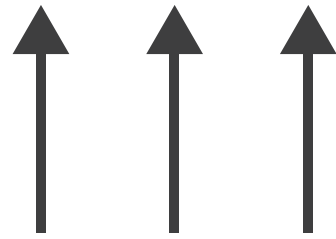


➔ Axis Orientation

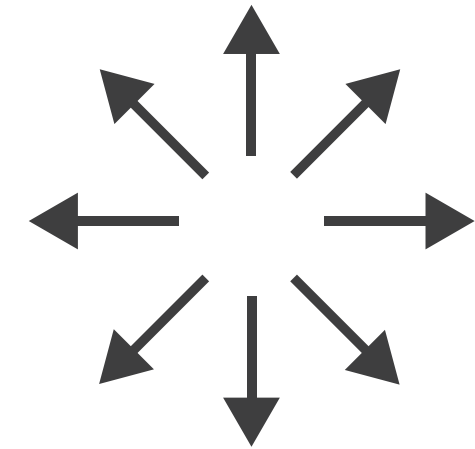
➔ Rectilinear



➔ Parallel



➔ Radial



Idioms: scatterplot matrix, parallel coordinates

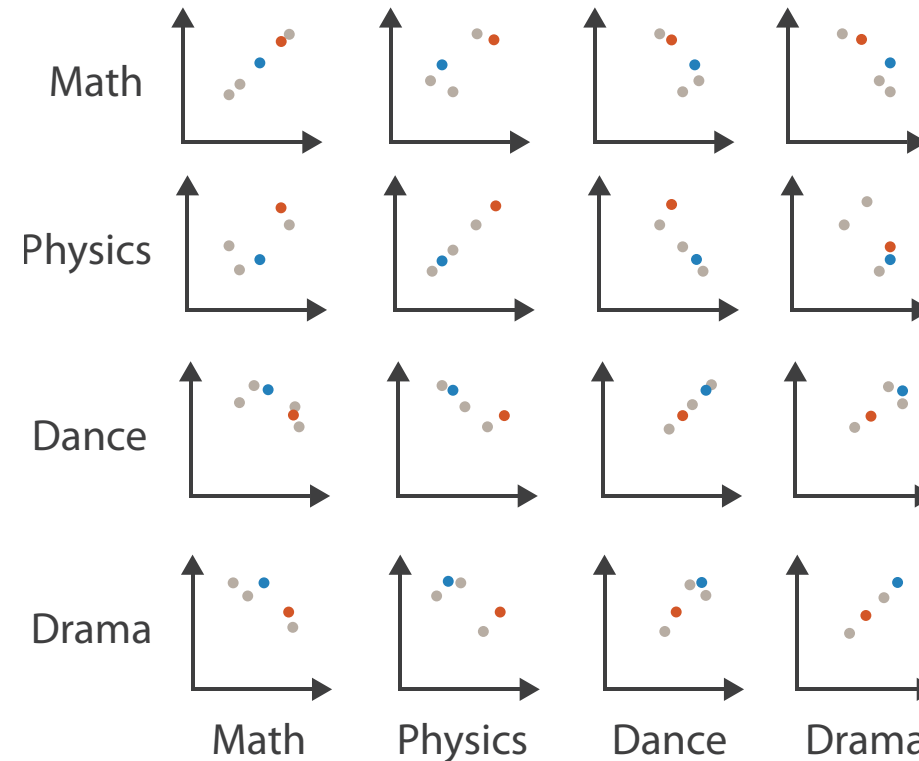
- scatterplot matrix (SPLOM)

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

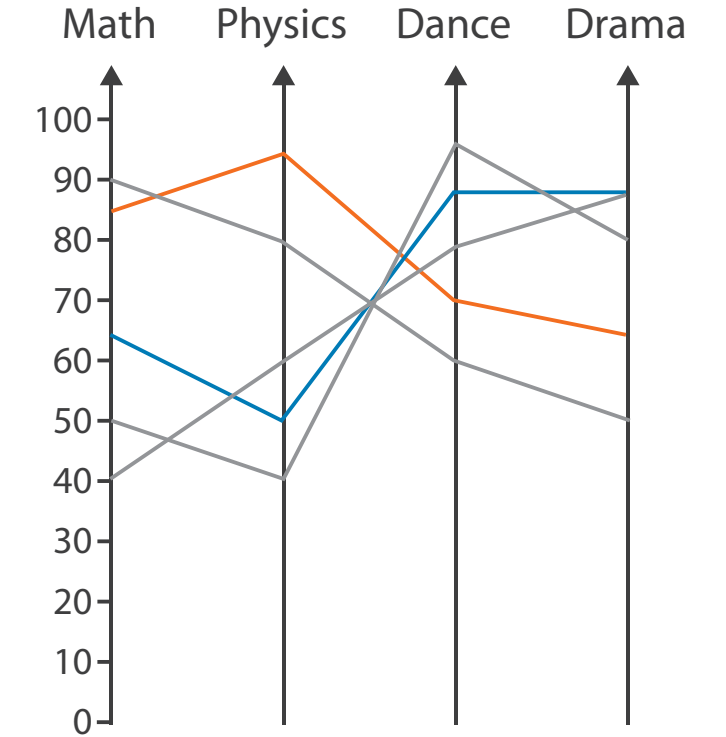
- parallel coordinates

- parallel axes, jagged line representing item
- rectilinear axes, item as point
 - axis ordering is major challenge
- scalability
 - dozens of attribs
 - hundreds of items

Scatterplot Matrix



Parallel Coordinates

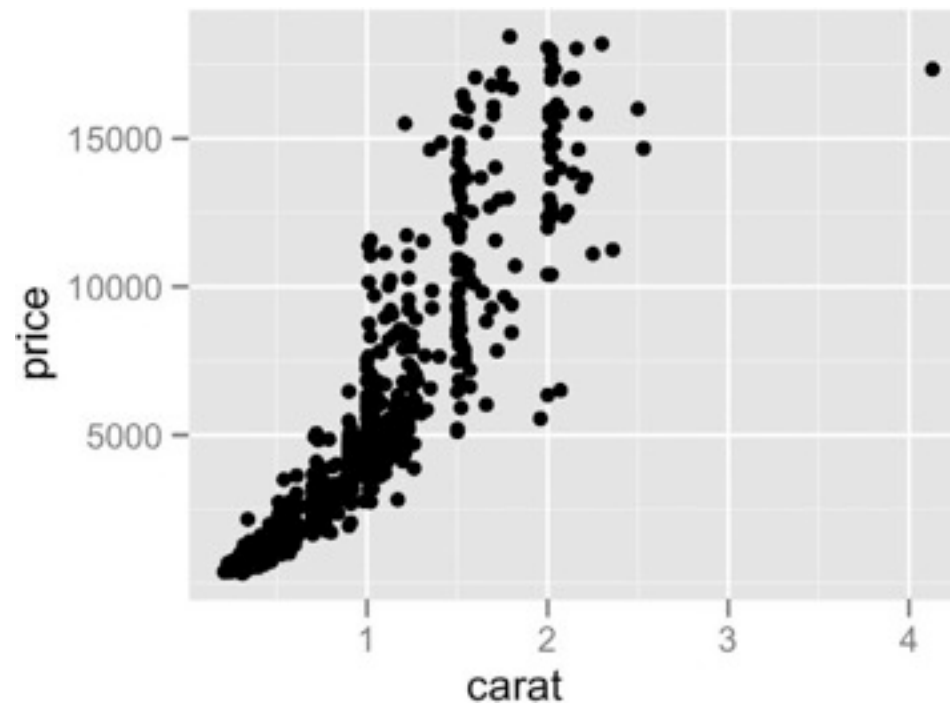


Table

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

Task: Correlation

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



[A layered grammar of graphics. Wickham. *Journ. Computational and Graphical Statistics* 19:1 (2010), 3–28.]

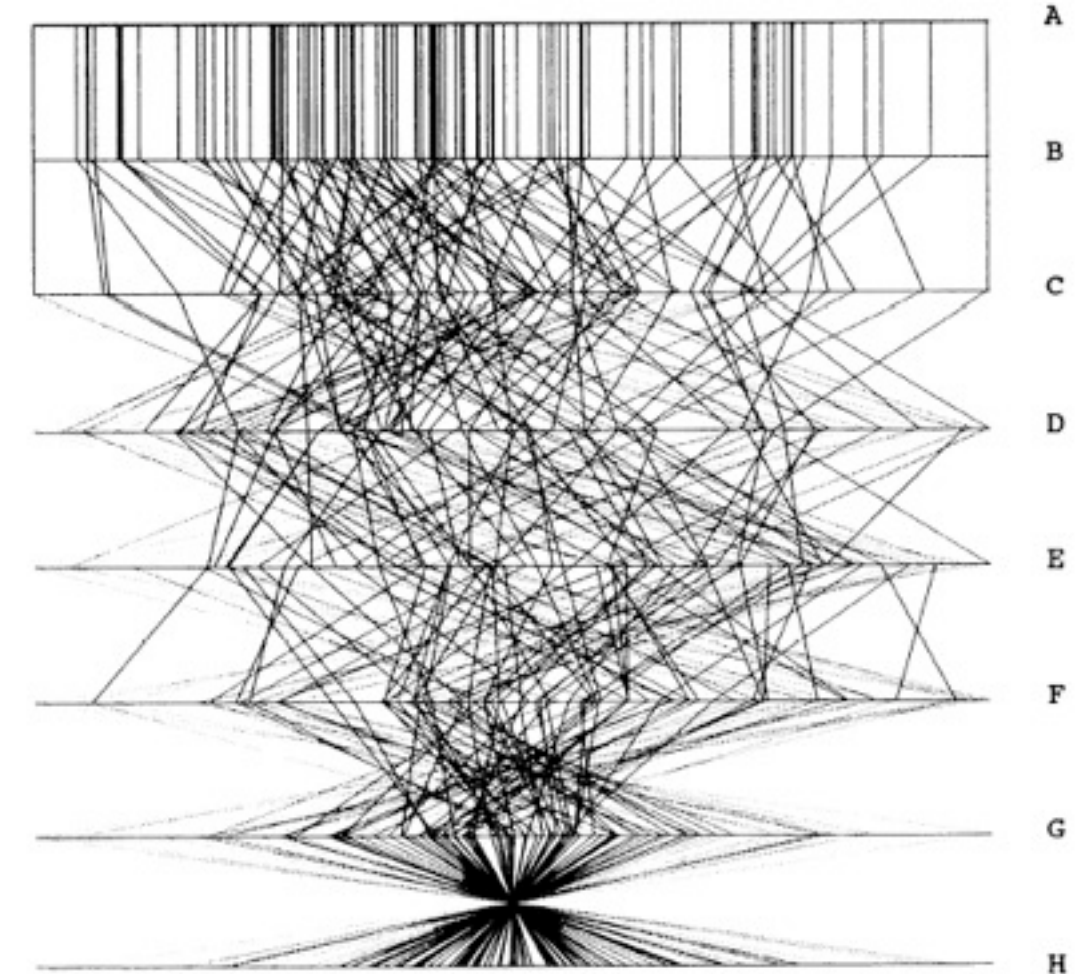
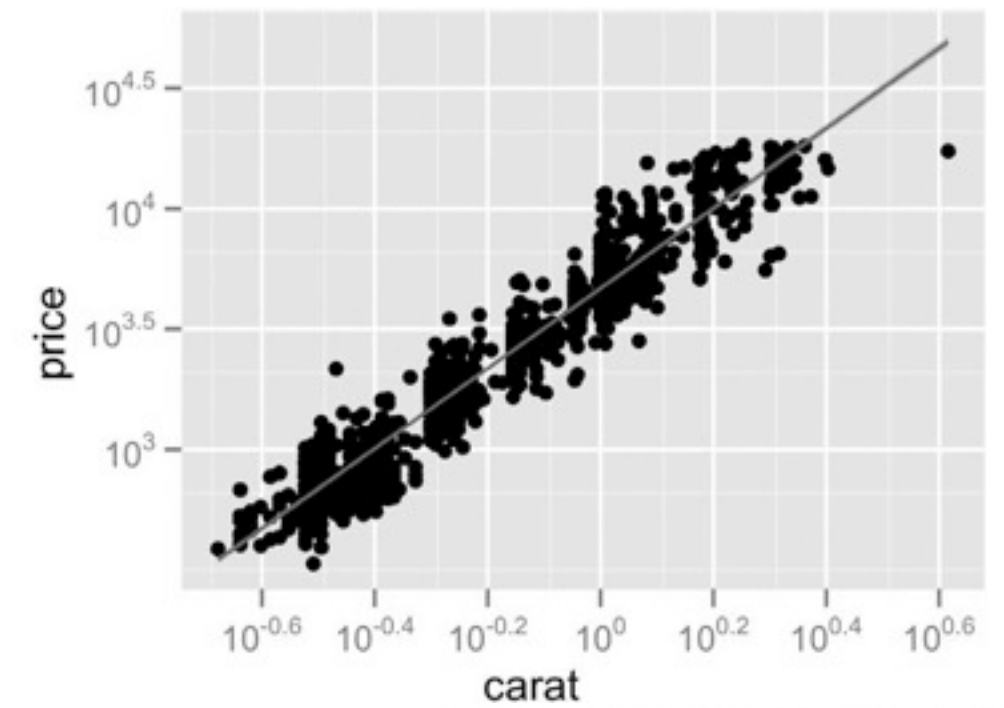
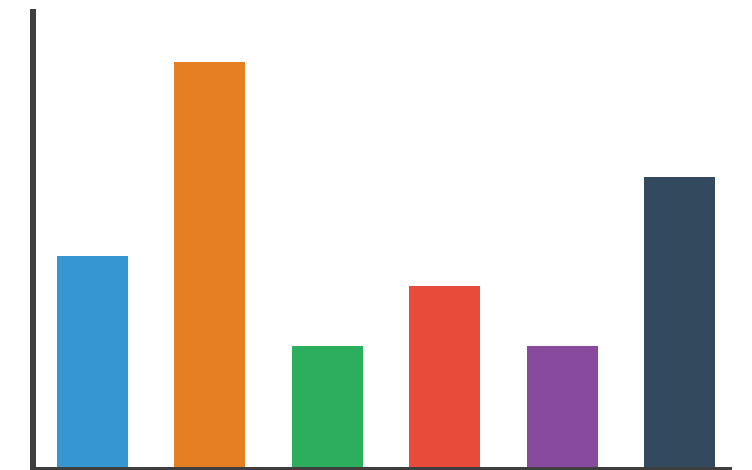
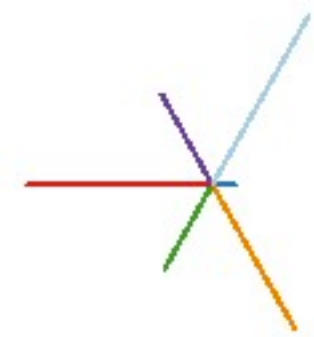


Figure 3. Parallel Coordinate Plot of Six-Dimensional Data Illustrating Correlations of $\rho = 1, .8, .2, 0, -.2, -.8, \text{ and } -1$.

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

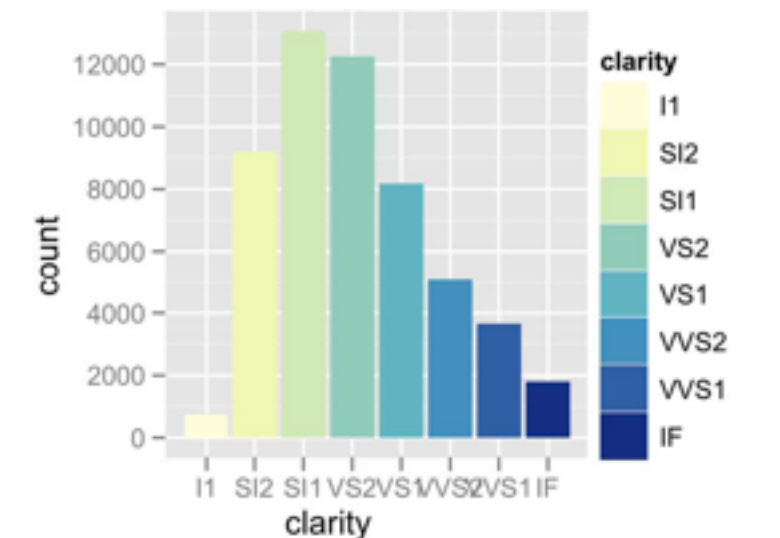
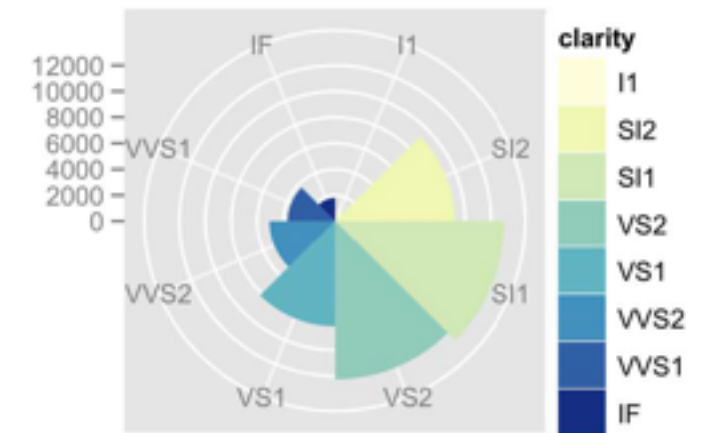
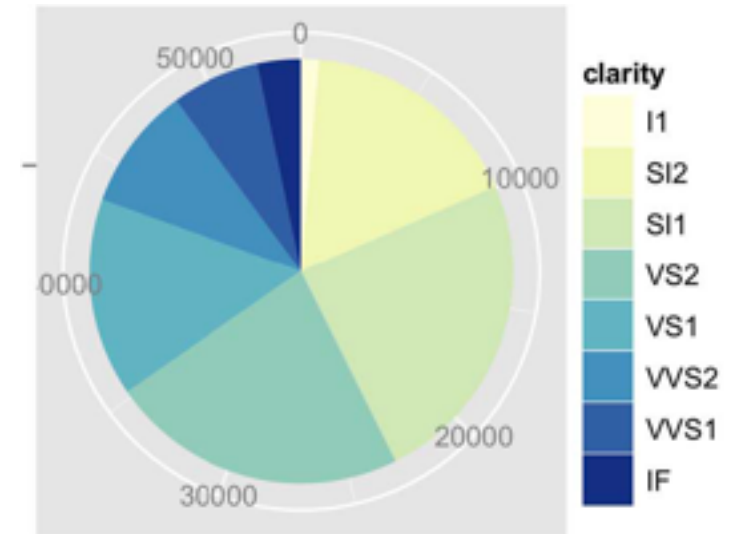
Idioms: radial bar chart, star plot

- radial bar chart
 - radial axes meet at central ring, line mark
- star plot
 - radial axes, meet at central point, line mark
- bar chart
 - rectilinear axes, aligned vertically
- accuracy
 - length unaligned with radial
 - less accurate than aligned with rectilinear



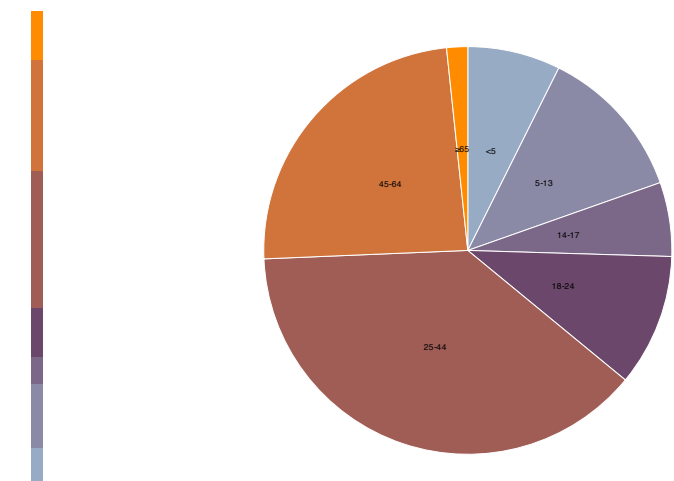
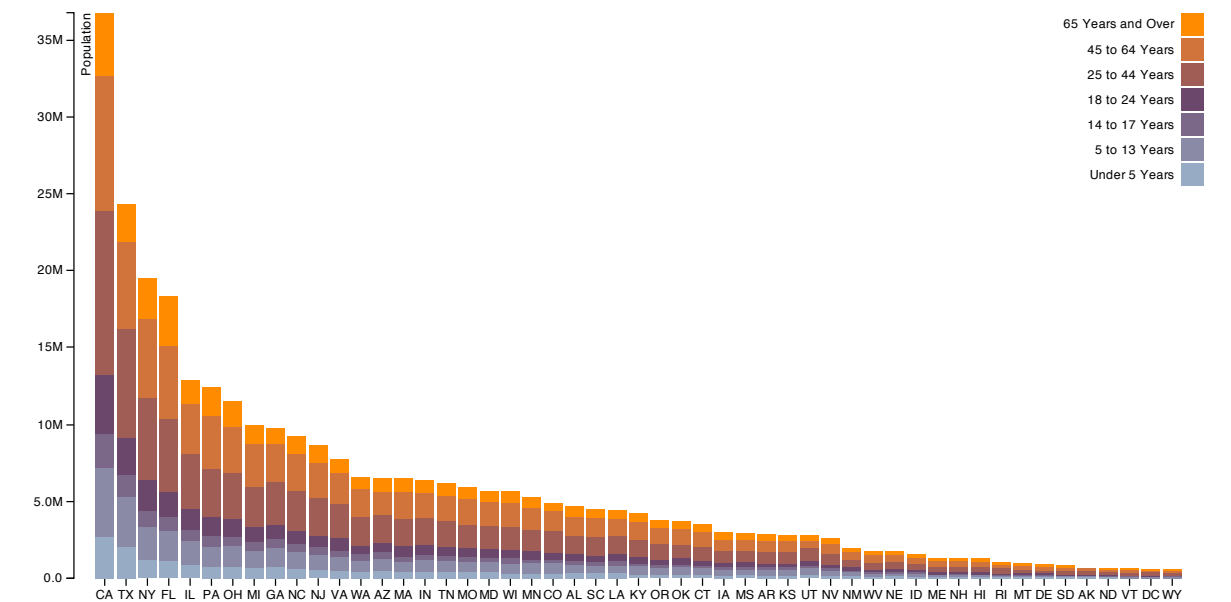
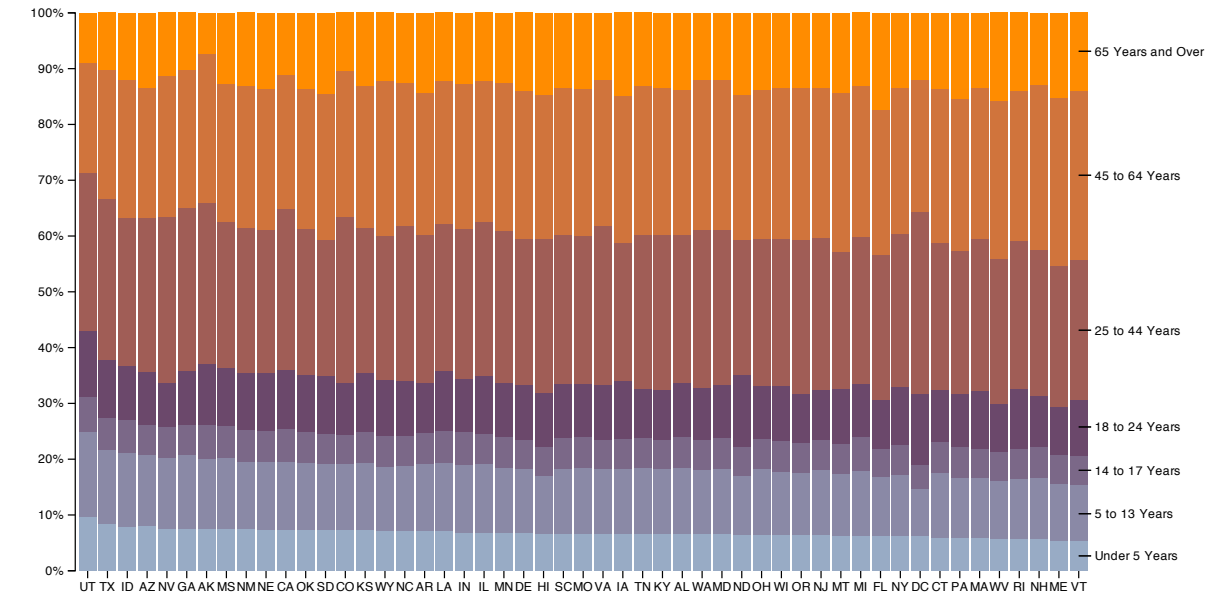
Idioms: pie chart, polar area chart

- pie chart
 - area marks with angle channel
 - accuracy: angle/area much less accurate than line length
- polar area chart
 - area marks with length channel
 - more direct analog to bar charts
- data
 - 1 categ key attrib, 1 quant value attrib
- task
 - part-to-whole judgements



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



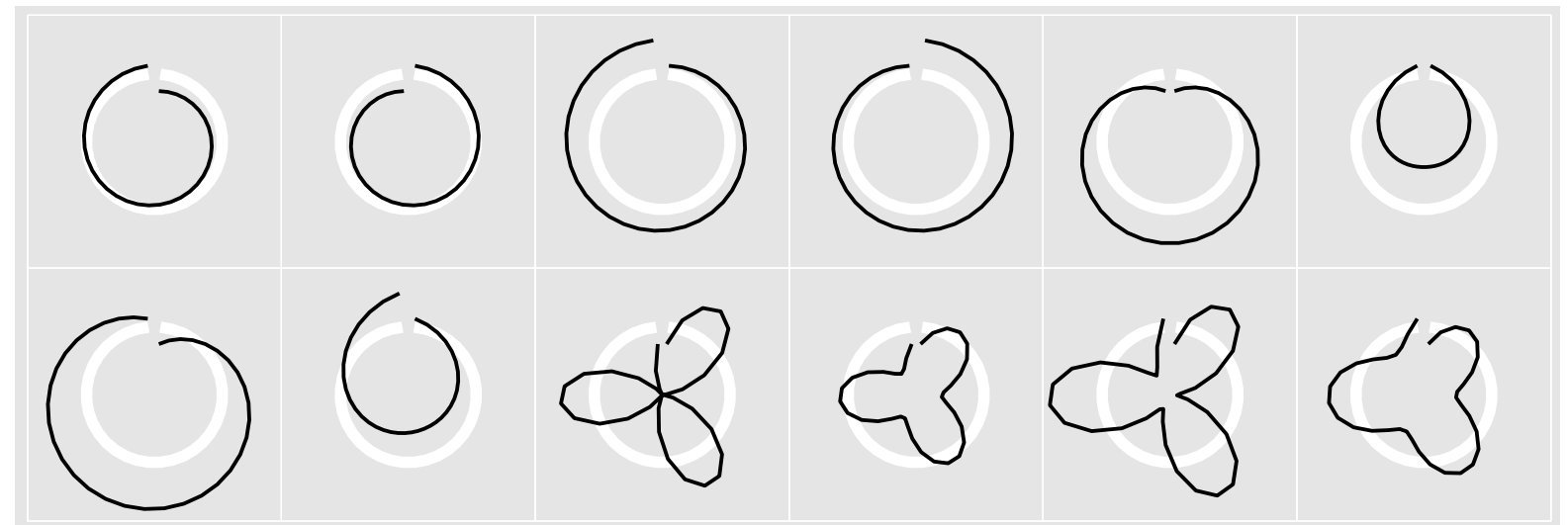
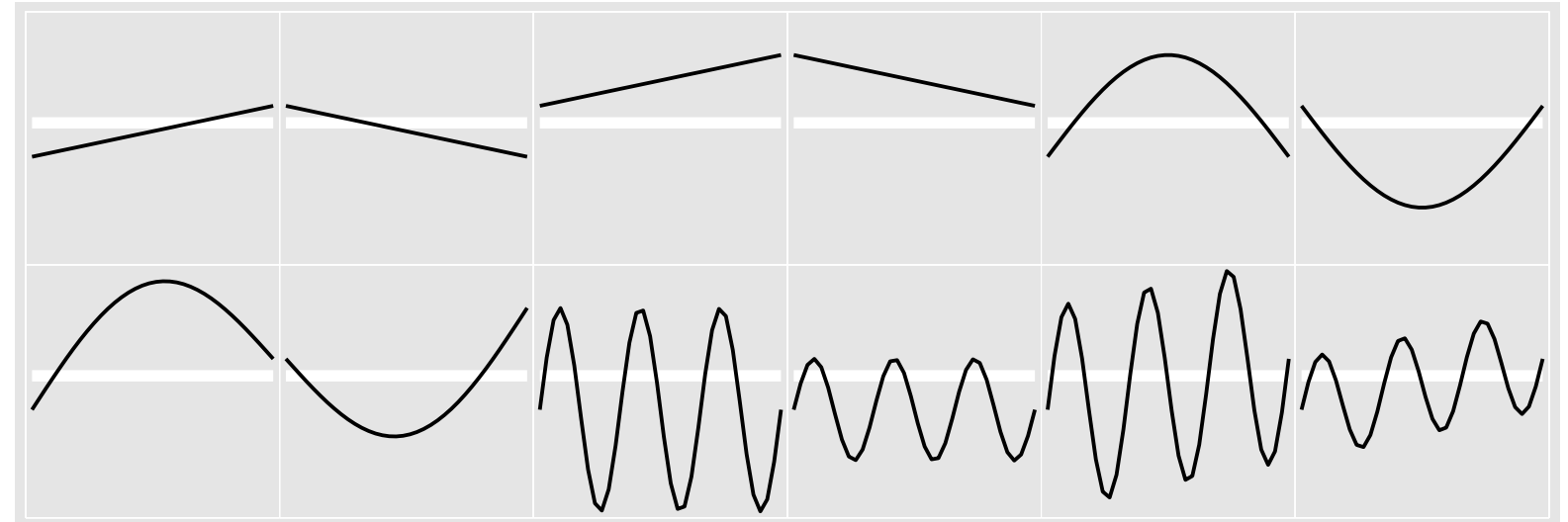
<http://bl.ocks.org/mbostock/3887235>,

<http://bl.ocks.org/mbostock/3886208>,

<http://bl.ocks.org/mbostock/3886394>.

Idiom: **glyphmaps**

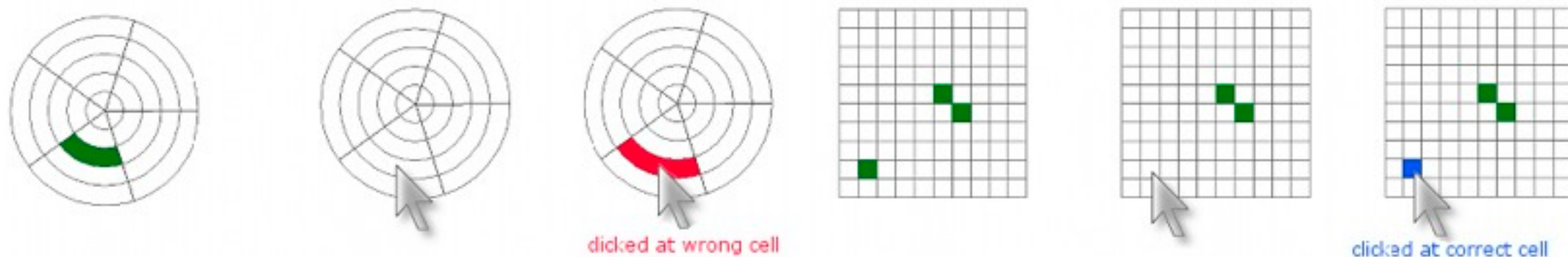
- rectilinear good for linear vs nonlinear trends
- radial good for cyclic patterns



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

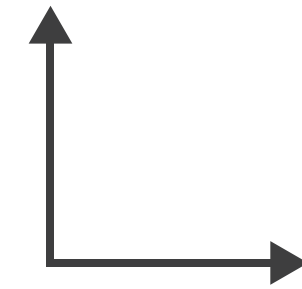
Orientation limitations

- **rectilinear: scalability wrt #axes**
 - 2 axes best
 - 3 problematic
 - more in afternoon
 - 4+ impossible
- **parallel: unfamiliarity, training time**
- **radial: perceptual limits**
 - asymmetry: angles lower precision than lengths
 - sometimes can be exploited!

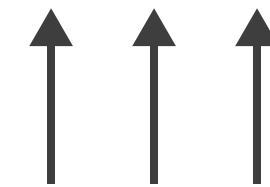


➔ Axis Orientation

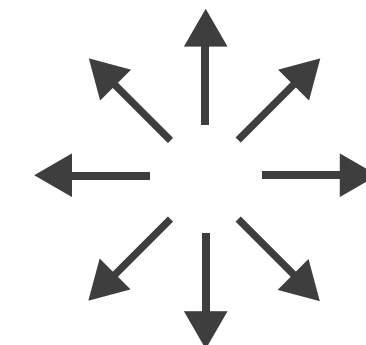
➔ Rectilinear



➔ Parallel

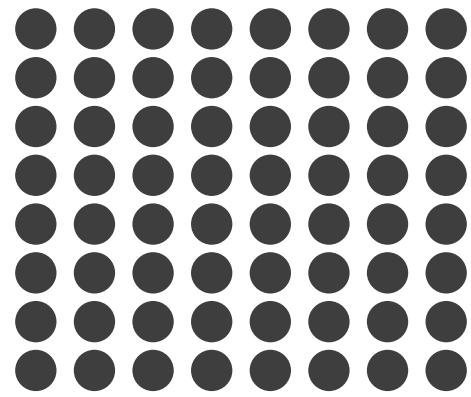


➔ Radial

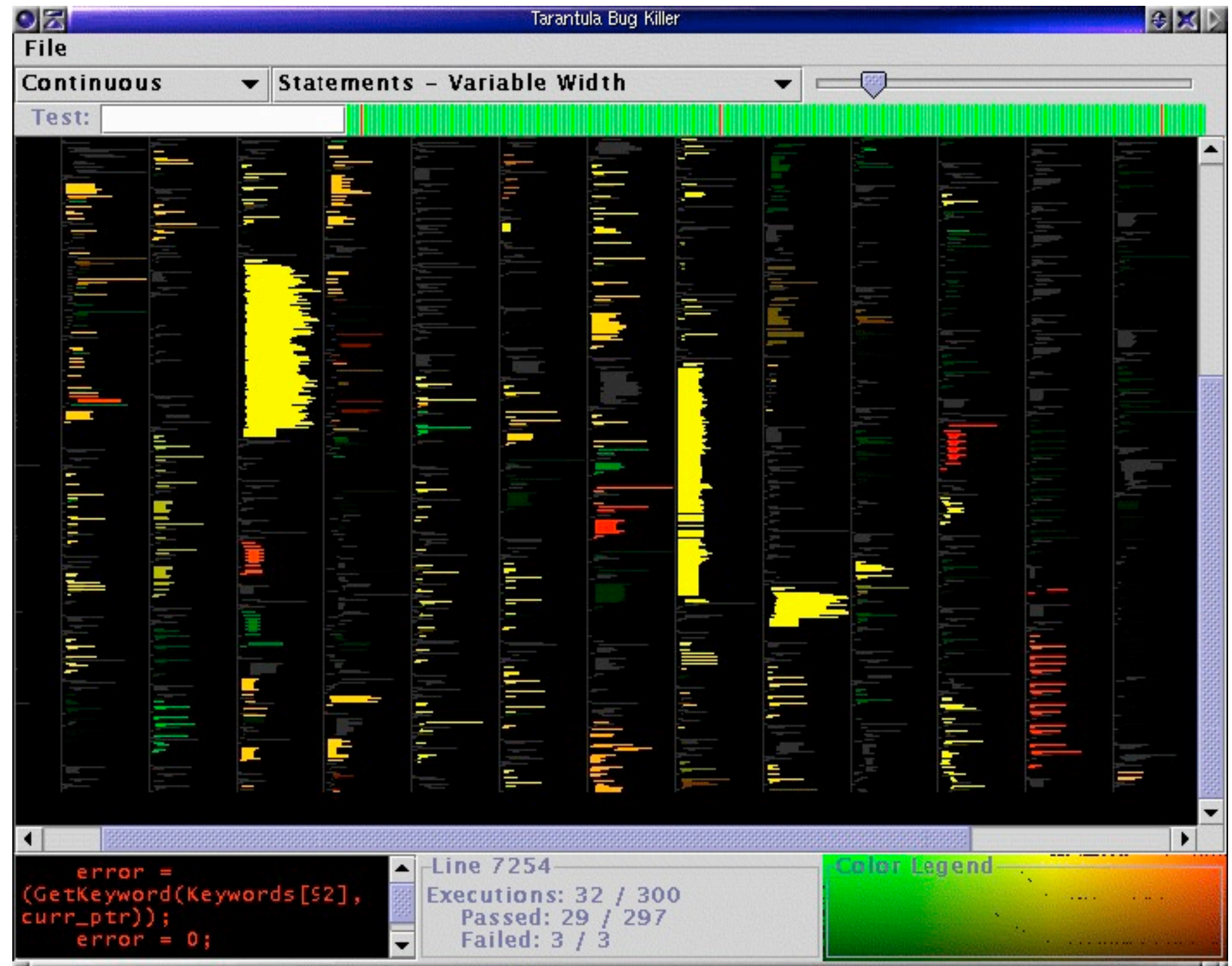


→ Layout Density

→ Dense



dense software overviews



Further reading

- Visualization Analysis and Design. Tamara Munzner. CRC Press, 2014.
 - *Chap 2: Data Abstraction*
 - *Chap 3: Task Abstraction*
 - *Chap 7: Tables*
- *A Brief History of Data Visualization*. Friendly. 2008.
<http://www.datavis.ca/milestones>

Now

- Break (15 min)
- Demo (30 min)
 - Guest lecture/demo from Robert Kosara on data wrangling
 - Tableau and Wrangler
- Lab 2 (45 min)

Lab/Assignment 2

- two main datasets
 - development aid from Guardian Datablog
 - your choice from small set
- focus on tasks and spatial layout as discussed in class for your exploration, story discovery, and writeup
 - provide rationale justifying your design decisions
- submit next week
 - by 9am Tue, email tmm@cs.ubc.ca with subject JOURN Week 2