

# G53FIV: Fundamentals of Information Visualization

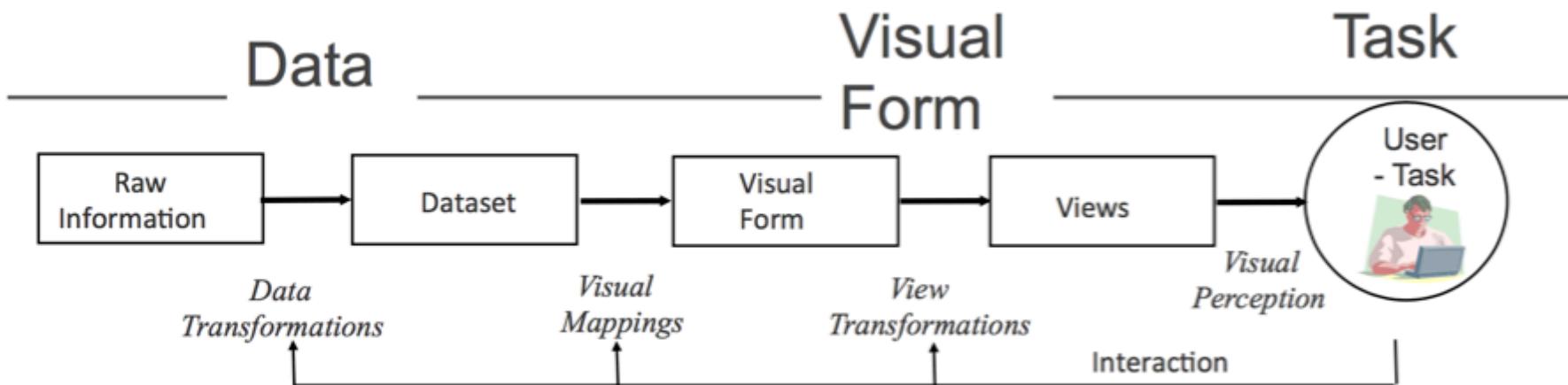
## Lecture 13: Recap of Fundamentals

Ke Zhou  
School of Computer Science  
Ke.Zhou@nottingham.ac.uk

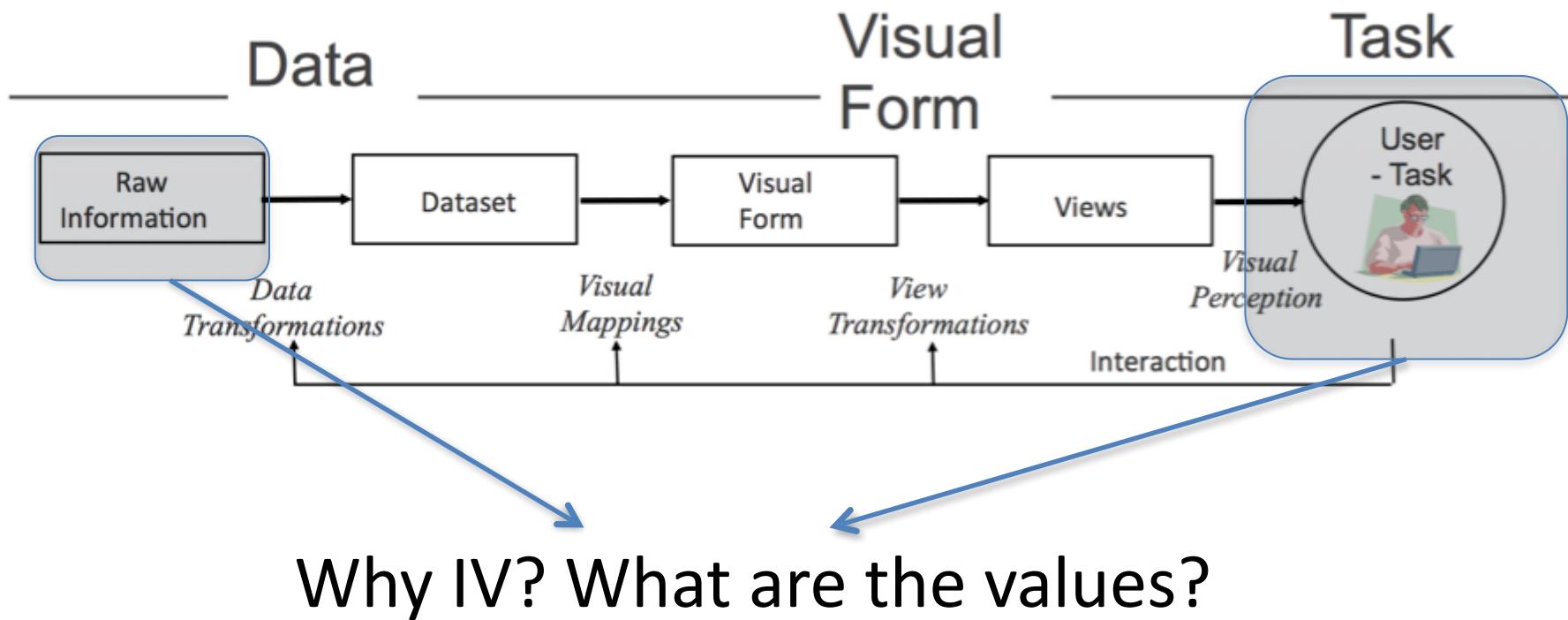
<https://moodle.nottingham.ac.uk/course/view.php?id=68644>

# A Recap of Fundamentals

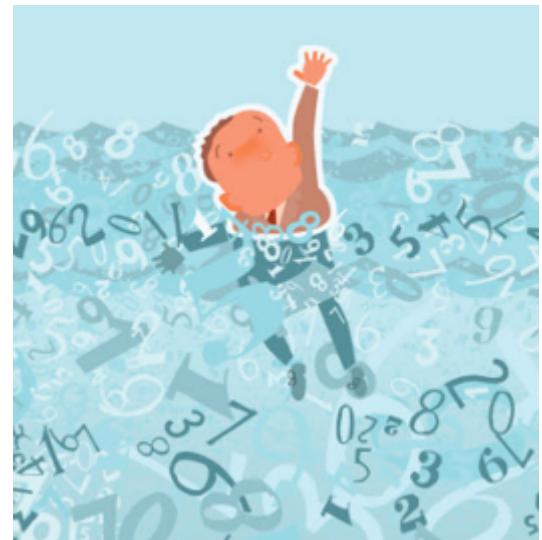
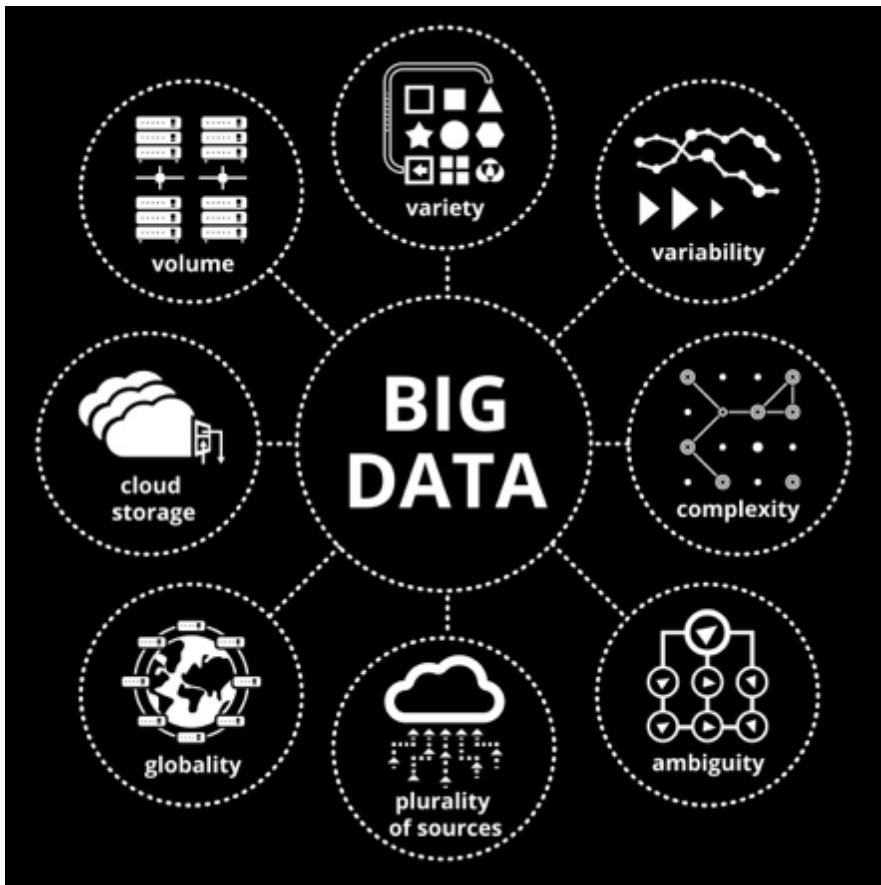
# Information Visualization



# Information Visualization



# Information Overload



# Objective

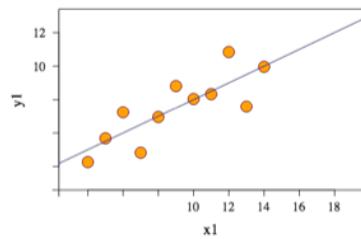
- Transform the data into information (understanding, insight) thus making it useful



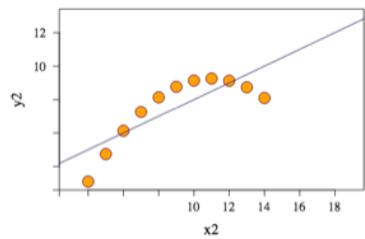
# Anscombe's Quartet

	Set A		Set B		Set C		Set D	
	X	Y	X	Y	X	Y	X	Y
0	10	8.04	10	9.14	10	7.46	8	6.58
1	8	6.95	8	8.14	8	6.77	8	5.76
2	13	7.58	13	8.74	13	12.74	8	7.71
3	9	8.81	9	8.77	9	7.11	8	8.84
4	11	8.33	11	9.26	11	7.81	8	8.47
5	14	9.96	14	8.10	14	8.84	8	7.04
6	6	7.24	6	6.13	6	6.08	8	5.25
7	4	4.26	4	3.10	4	5.39	19	12.50
8	12	10.84	12	9.13	12	8.15	8	5.56
9	7	4.82	7	7.26	7	6.42	8	7.91
10	5	5.68	5	4.74	5	5.73	8	6.89
mean	9.00	7.50	9.00	7.50	9.00	7.50	9.00	7.50
std	3.32	2.03	3.32	2.03	3.32	2.03	3.32	2.03
corr	0.82		0.82		0.82		0.82	
lin. reg.	$y = 3.00 + 0.500x$		$y = 3.00 + 0.500x$		$y = 3.00 + 0.500x$		$y = 3.00 + 0.500x$	

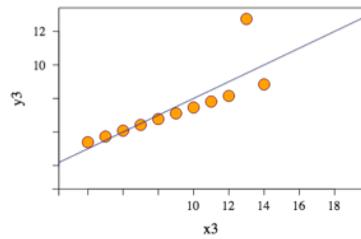
"what most people would see in their mind's eye [for a linear relationship with some unexplained variation]"



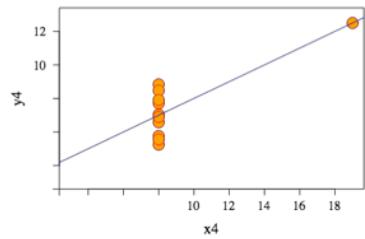
"y has a smooth curved relation with x, possibly quadratic, and there is little residual variability"



"all but one of the observations lie close to a straight line (not the one yielded by the standard regression calculation)"

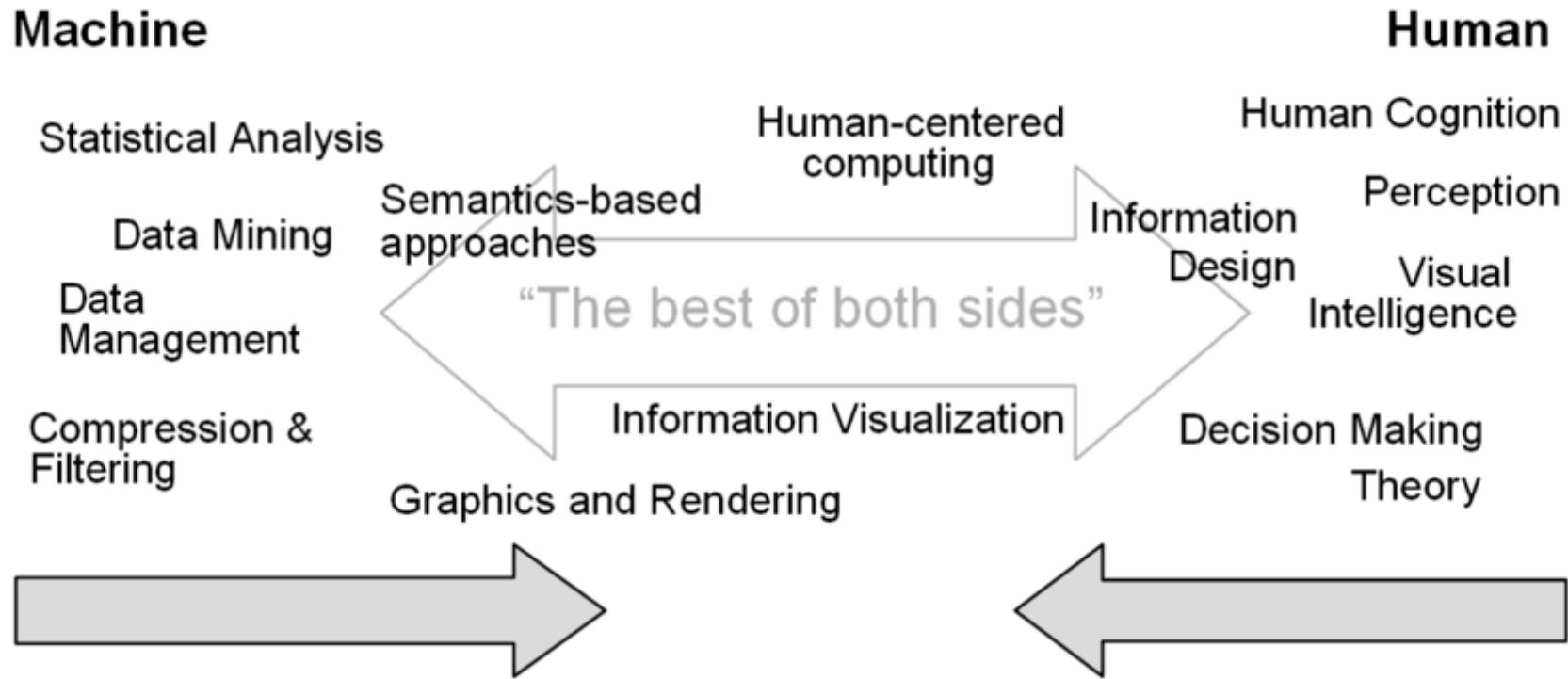


"all the information about the slope of the regression line resides in one observation"



[Anscombe, 1973]

# The Best of Both Sides



**Fig. 2.** Visual analytics integrates scientific disciplines to improve the division of labor between human and machine.

# Key Values of Visualizations

- **Record** information
  - Blueprints, photographs, seismographs, ...
- **Communicate** information to others
  - Share and persuade
  - Collaborate and revise
- Analyze data to **support reasoning**
  - Find patterns / Discover errors in data
  - Expand memory
  - Develop and assess hypotheses

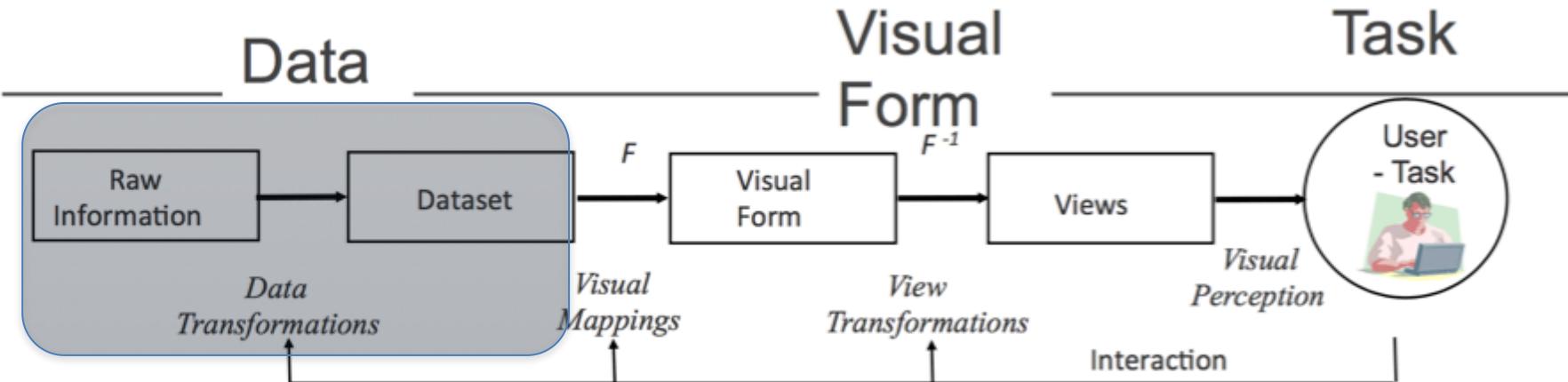


# Data and Initial Tasks: Before the Actual Visualization

- Pick a dataset of your interest
- Pose the initial questions/tasks that you would like to answer/accomplish
- Assess the fitness of the data
- **Visualization**
- Refine your questions/tasks

Remember our course work and the house price case study?

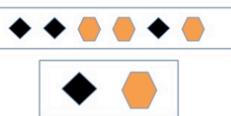
# Information Visualization



00210	+43.005895	-071.013202	U	PORTSMOUTH	33	015
00211	+43.005895	-071.013202	U	PORTSMOUTH	33	015
00212	+43.005895	-071.013202	U	PORTSMOUTH	33	015
00213	+43.005895	-071.013202	U	PORTSMOUTH	33	015
00214	+43.005895	-071.013202	U	PORTSMOUTH	33	015
00215	+43.005895	-071.013202	U	PORTSMOUTH	33	015
00501	+40.922326	-072.637078	U	HOLTSVILLE	36	103
00544	+40.922326	-072.637078	U	HOLTSVILLE	36	103
00601	+18.165273	-066.722583		ADJUNTAS	72	001
00602	+18.493103	-067.180953		AGUADA	72	003
00603	+18.455913	-067.145780		AGUADILLA	72	005
00604	+18.493520	-067.135883		AGUADILLA	72	005
00605	+18.465162	-067.141486	P	AGUADILLA	72	005
00606	+18.172947	-066.944111		MARICAO	72	093
00610	+18.288685	-067.139696		ANASCO	72	011



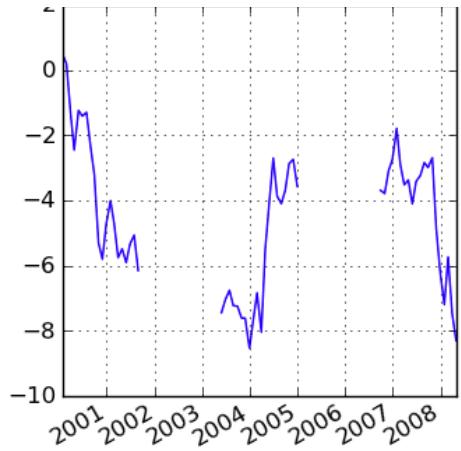
- Data transformation:** create a structural model (schema), mapping raw data into data tables

- FILTER Rows 
- SELECT Column Types 
- ArRANGE Rows (SORT) 
- Mutate (into something new) 
- Summarize by Groups 



# Data Processing

- Data cleaning and filtering
  - for quality control
  - Remove (Outlier, missing data)
  - Modify (conversion of format, etc.)
- Data adjustment
  - Depends on your task and questions to ask
  - Relational algebra:
    - e.g. Aggregation, mean, sort, projection
  - Reformatting and Integration



# R is a tool for...

## Data Manipulation

- connecting to data sources
- slicing & dicing data

## Modeling & Computation

- statistical modeling
- numerical simulation

## Data Visualization

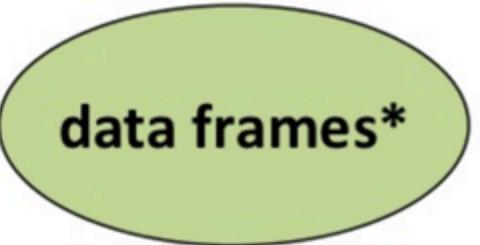
- visualizing fit of models
- composing statistical graphics

munge

model

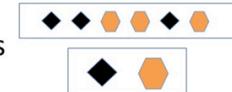
visualize

# R Data Structures

	Linear	Rectangular
Homogeneous	? 	
Heterogeneous		

# dplyr

- dplyr takes the `%>%` operator and uses it to great effect for manipulating data frames
  - Works only with data frames
  - 5 basic “verbs” work for 90% of data manipulations manipulations

Verbs	What does it do?	
<code>filter()</code>	Select a subset of ROWS by conditions	
<code>arrange()</code>	Reorders ROWS in a data frame	
<code>select()</code>	Select the COLUMNS of interest	
<code>mutate()</code>	Create new columns based on existing columns (mutations!)	
<code>summarise()</code>	Aggregate values for each group, reduces to single value	

# Pipe Operator

- **Library(magrittr)**
  - A R package launched on Jan 2014
  - A “magic” operator called the PIPE was introduced
  - `%>%`
  - i.e. “AND THEN”, “PIPE TO”

```
round(sqrt(1000), 3)  
  
library(magrittr)  
1000 %>% sqrt %>% round()  
1000 %>% sqrt %>% round(., 3)
```

Take 1000, and then its sqrt  
And then round it



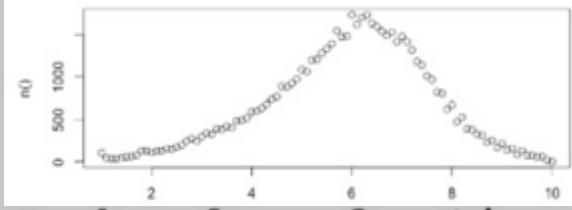
# Chain the “Verbs” Together

- Chain them together

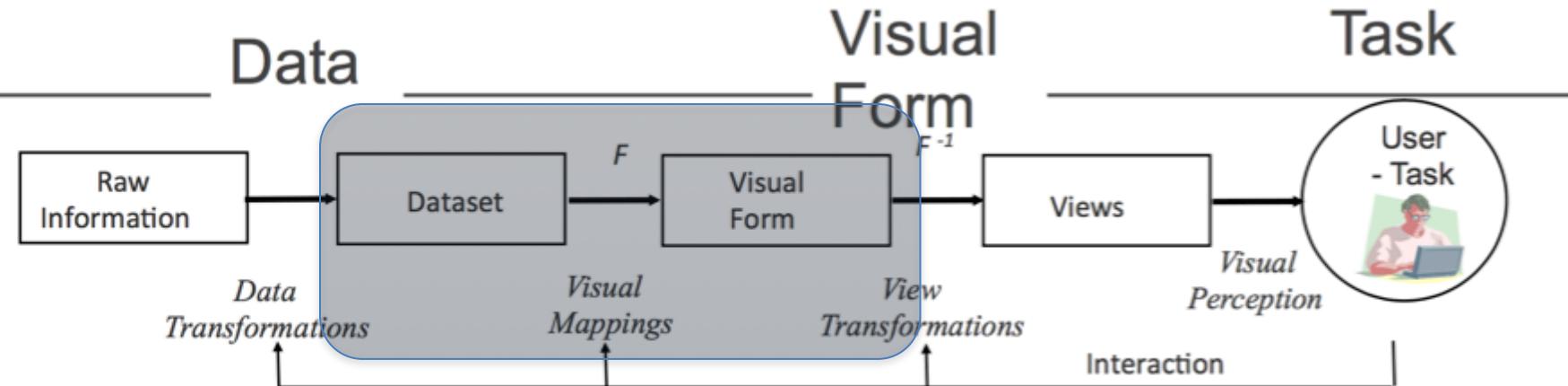
```
producers_nightmare <-  
  filter(movies_df, !is.na(budget)) %>%  
  mutate(costPerMinute = budget/length) %>%  
  arrange(desc(costPerMinute)) %>%  
  select(title, costPerMinute)
```

- Can also be fed to a “plot” command

```
movies %>%  
  group_by(rating) %>%  
  summarize(n()) %>%  
  plot() # plots the histogram of movies by Each value of rating
```



# Information Visualization

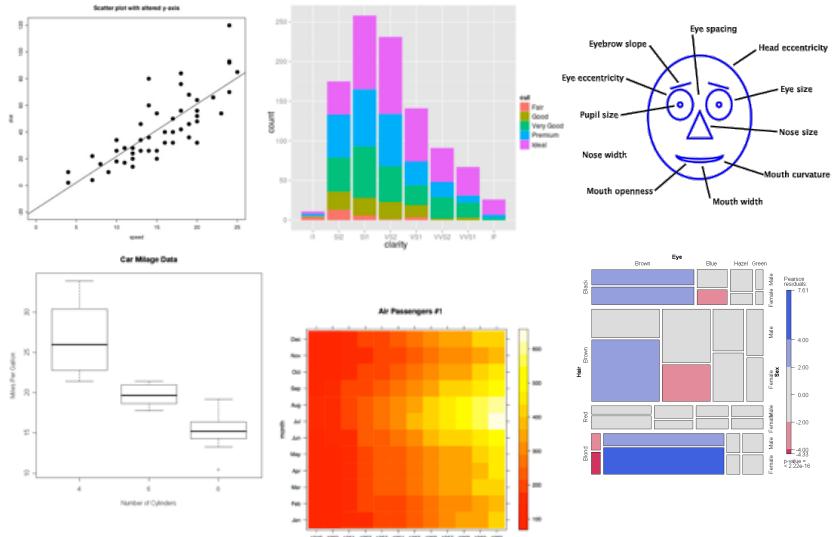


00210	+43.005895	-071.013202	U	PORTRSMOUTH	33	015
00211	+43.005895	-071.013202	U	PORTRSMOUTH	33	015
00212	+43.005895	-071.013202	U	PORTRSMOUTH	33	015
00213	+43.005895	-071.013202	U	PORTRSMOUTH	33	015
00214	+43.005895	-071.013202	U	PORTRSMOUTH	33	015
00215	+43.005895	-071.013202	U	PORTRSMOUTH	33	015

	Nominal	Ordinal	Quantitative
Position	✓	✓	✓
Size	✓	✓	~
(Grey)Value	✓	✓	~
Texture	✓	~	✗
Color	✓	✗	✗
Orientation	✓	✗	✗
Shape	✓	✗	✗

✓ = Good  
~ = OK  
✗ = Bad

**Visual mapping:** create a visual spatial model, transforming data tables into visual structures



# Data, Image and Design

# Data Models

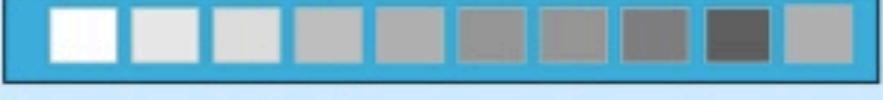
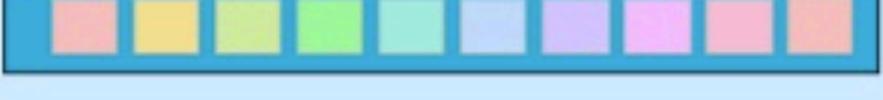
- Nominal, Ordinal, Quantitative?
- Dimension or Measure?

– Year	Q-Internal (O)	Dimension
– Age	Q-Ratio (O)	Depends
– Marital	N	Dimension
– Sex	N	Dimension
– People	Q-Ratio	Measure

	A	B	C	D	E
1	year	age	marst	sex	people
2	1850	0	0	1	1483789
3	1850	0	0	2	1450376
4	1850	5	0	1	1411067
5	1850	5	0	2	1359668
6	1850	10	0	1	1260099
7	1850	10	0	2	1216114
8	1850	15	0	1	1077133
9	1850	15	0	2	1110619
10	1850	20	0	1	1017281
11	1850	20	0	2	1003841
12	1850	25	0	1	862547
13	1850	25	0	2	799482
14	1850	30	0	1	730638
15	1850	30	0	2	639636
16	1850	35	0	1	588487
17	1850	35	0	2	505012
18	1850	40	0	1	475911
19	1850	40	0	2	428185
20	1850	45	0	1	384211
21	1850	45	0	2	341254
22	1850	50	0	1	321343
23	1850	50	0	2	286580
24	1850	55	0	1	194080

# Image: Visual Encoding Variables

## Bertin's Semiology of Graphics (1967)

- **position**
  - changes in the x, y, (z) location
- **size**
  - change in length, area or repetition
- **shape**
  - infinite number of shapes
- **value**
  - changes from light to dark
- **orientation**
  - changes in alignment
- **colour**
  - changes in hue at a given value
- **texture**
  - variation in pattern
- **motion**

Graphic by: Sheelagh Carpendale

Dr. Ke Zhou (<http://www.cs.nott.ac.uk/~pszkz/>)

# Levels of Organization

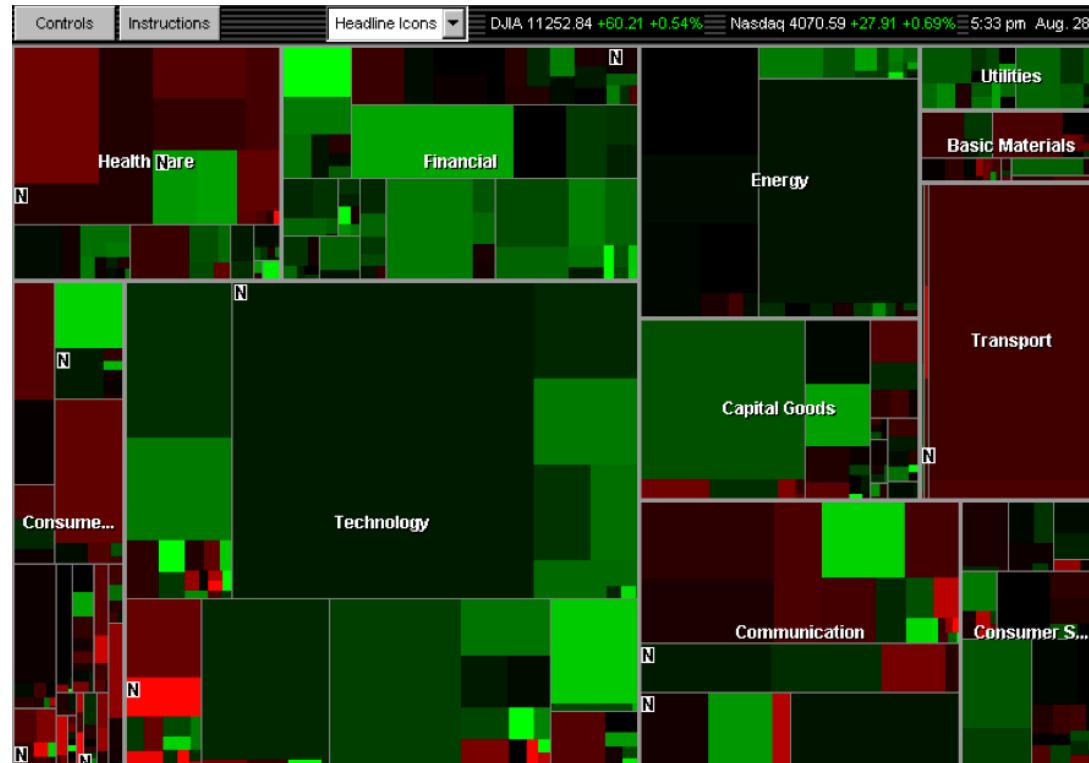
	Nominal	Ordinal	Quantitative
Position	✓	✓	✓
Size	✓	✓	~
(Grey)Value	✓	✓	~
Texture	✓	~	✗
Color	✓	✗	✗
Orientation	✓	✗	✗
Shape	✓	✗	✗

✓ = Good

~ = OK

✗ = Bad

# Example: Map of the Market

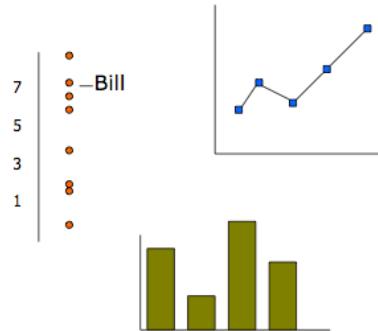


- Rectangle area: market cap (Q);
- Rectangle position: market sector (N)
- Color Hue: loss vs. gain (N, O)
- Color Value: magnitude of loss or gain (Q)

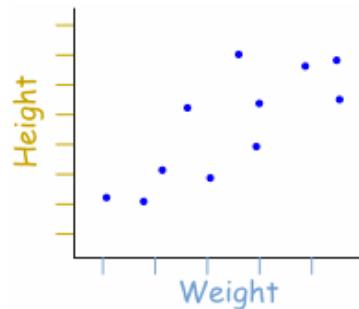
# Graphs and Charts

# Graphs

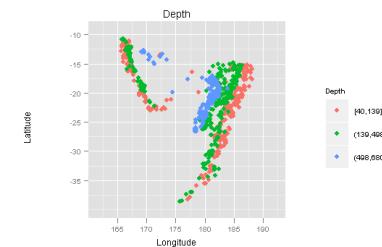
- Data Dimensions
  - 1 - Univariate data
  - 2 - Bivariate data
  - 3 - Trivariate data
  - >3 - Hypervariate data
- Data Types
  - Nominal, Ordinal, Quantitative
- Visualization Representations
  - Points, Lines, Bars, Boxes



Univariate



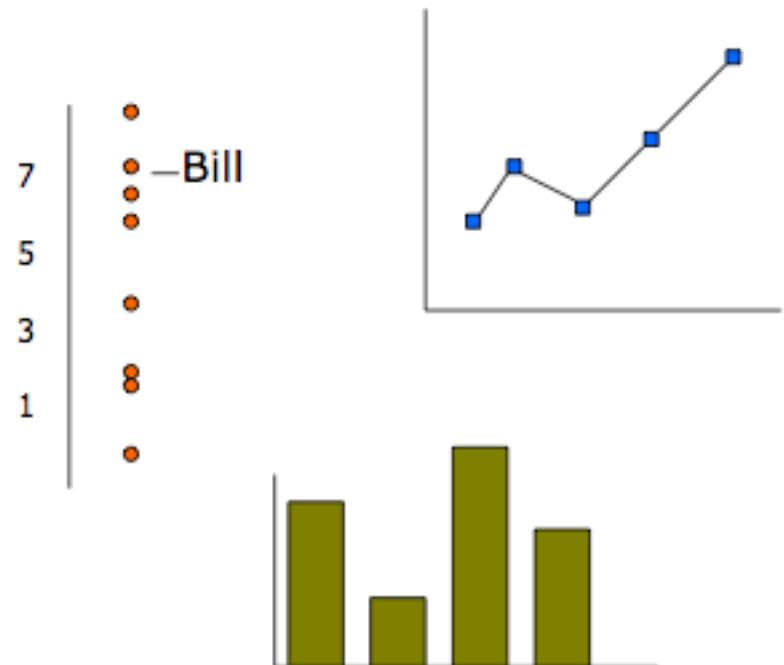
Bivariate



Trivariate

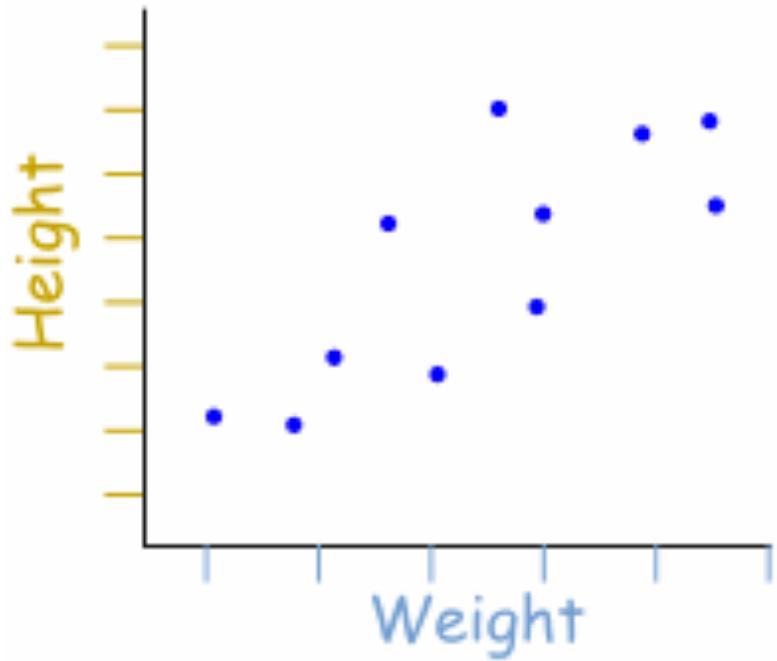
# Univariate Data

- In univariate representations, we often think of the data case as being shown along one dimension, and the value in another.
- Statistical view
  - Independent variable on x-axis (data case)
  - Track dependent variable along y-axis (value)



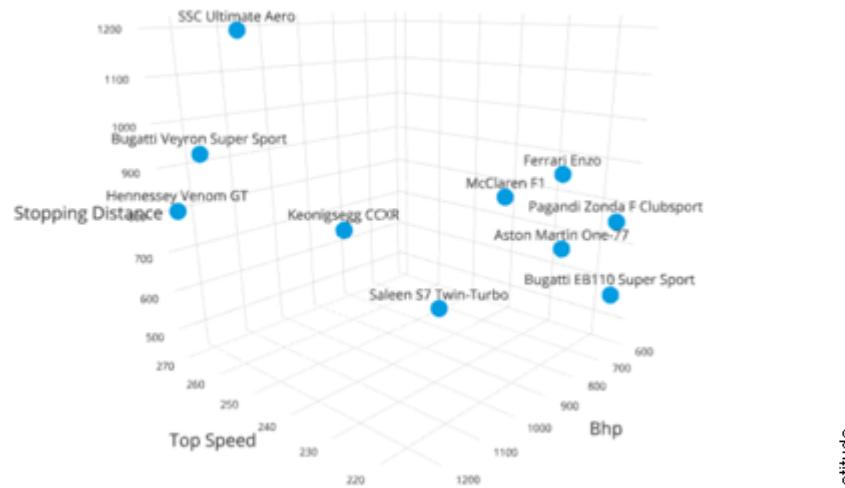
# Bivariate Data

- Scatter plot is commonly used
- Each mark is now a data case
- Objective:
  - Two variables, want to see relationship
  - Is there a linear, curved or random pattern?

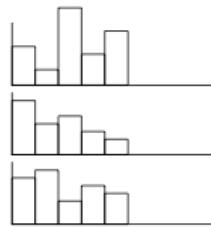


# Trivariate Data

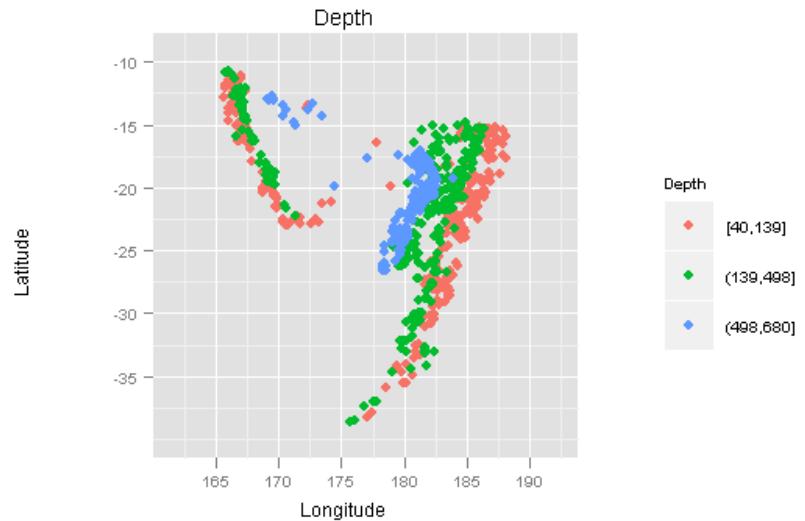
3D scatter plot

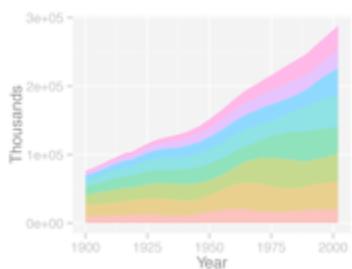
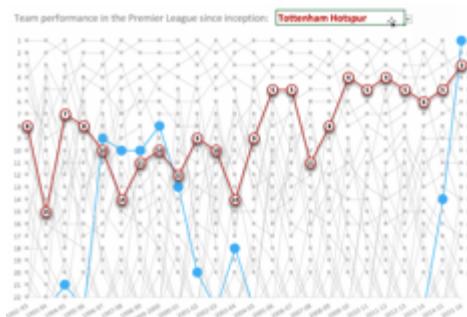
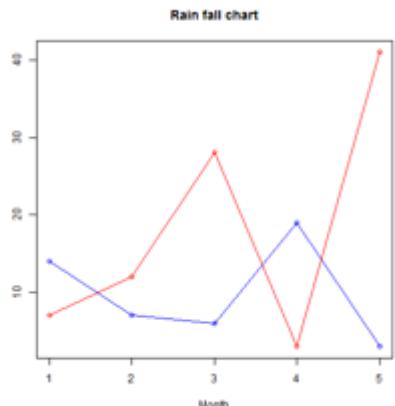
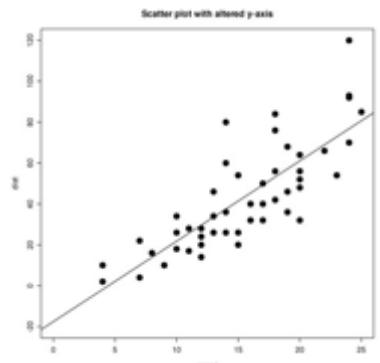


2D + mark  
property



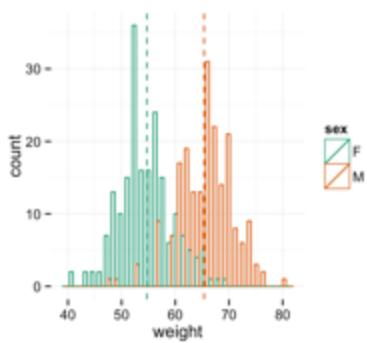
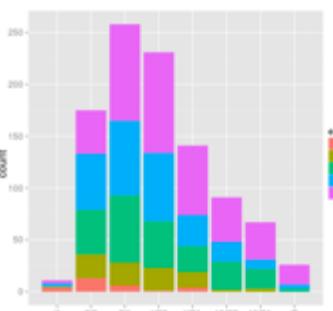
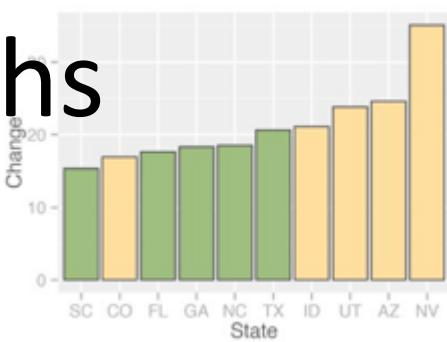
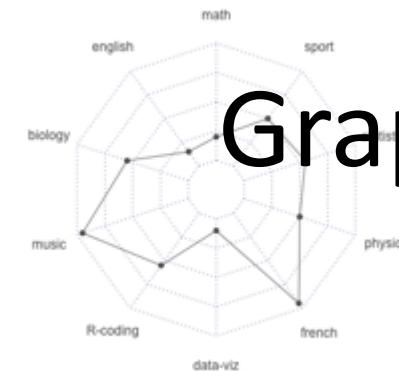
Represent each  
variable in its own  
explicit way





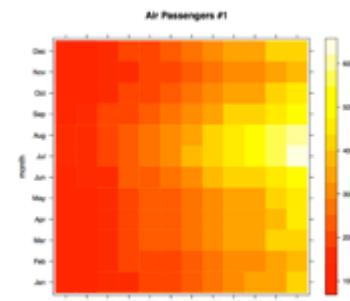
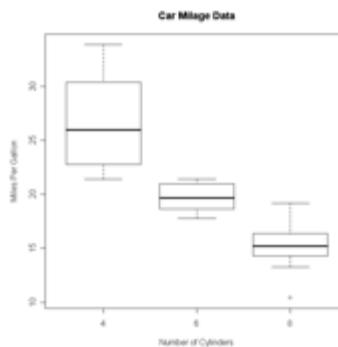
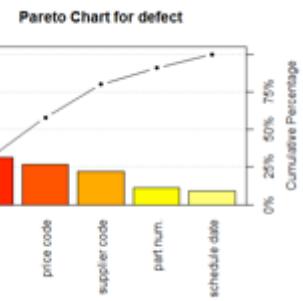
AgeGroup

- <5
- 5-14
- 15-24
- 25-34
- 35-44
- 45-54
- 55-64
- >65



SEX

- M
- F



# Multivariate Data Visualization

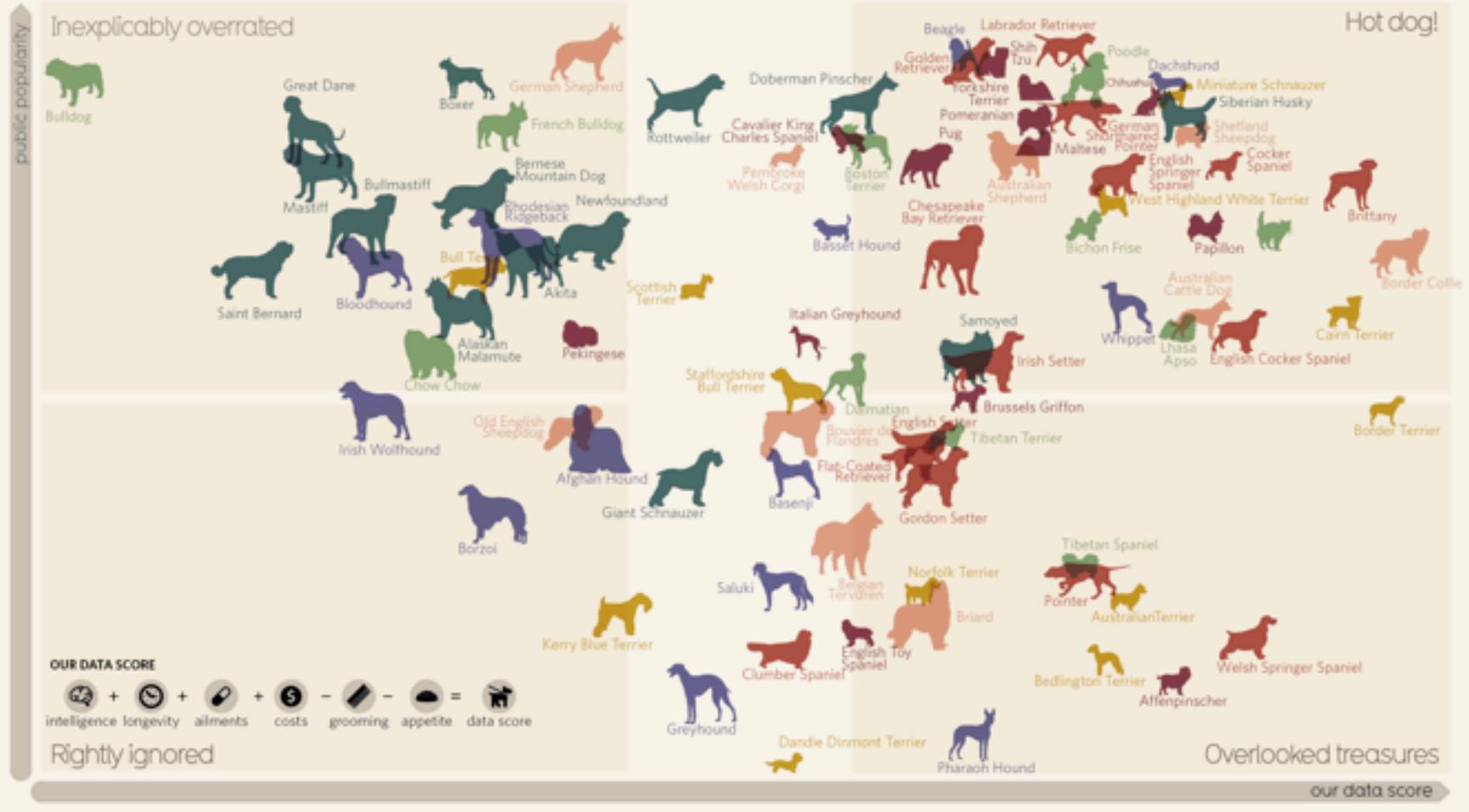
# When to use?

- Use tables when
  - The document will be used to **look up individual values**
  - The document will be used to **compare individual values**
  - **Precise values** are required
  - The quantitative info to be communicated involves **more than one unit of measure**
- Use graphs when
  - The message is contained in the **shape** of the values
  - The document will be used to **reveal relationships** among values
  - Especially useful when the number of data points is huge

(Optional Reading) Stephen Few. 2012. Show Me the Numbers: Designing Tables and Graphs to Enlighten (2nd ed.). Analytics Press, , USA.

# Best in Show

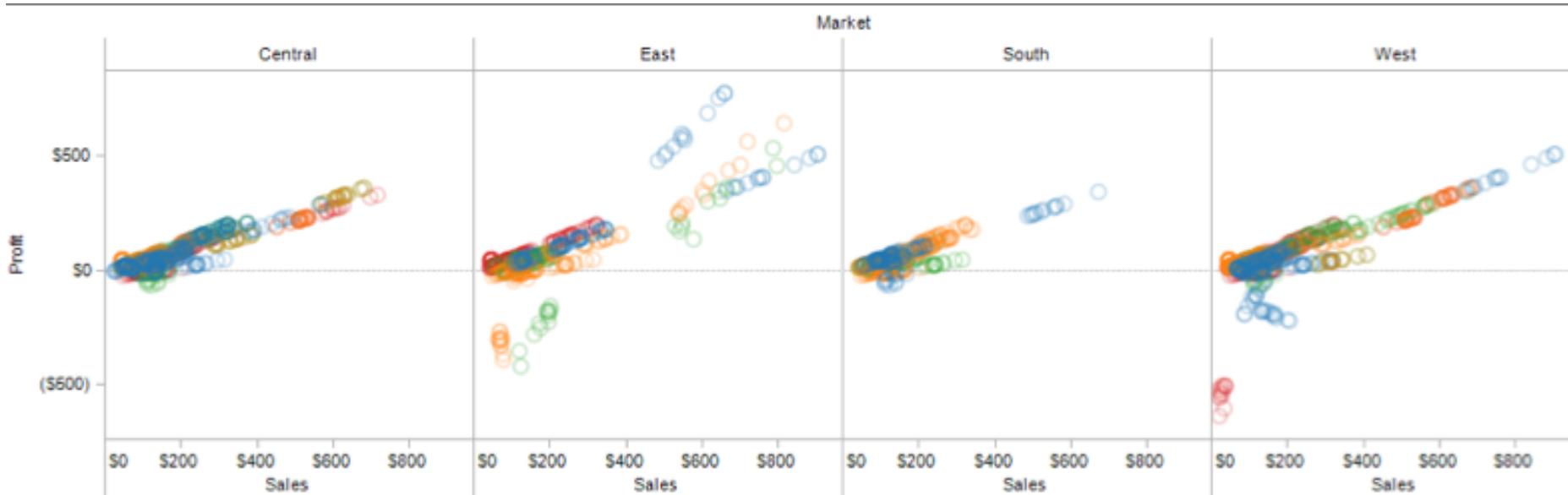
The ultimate data-dog



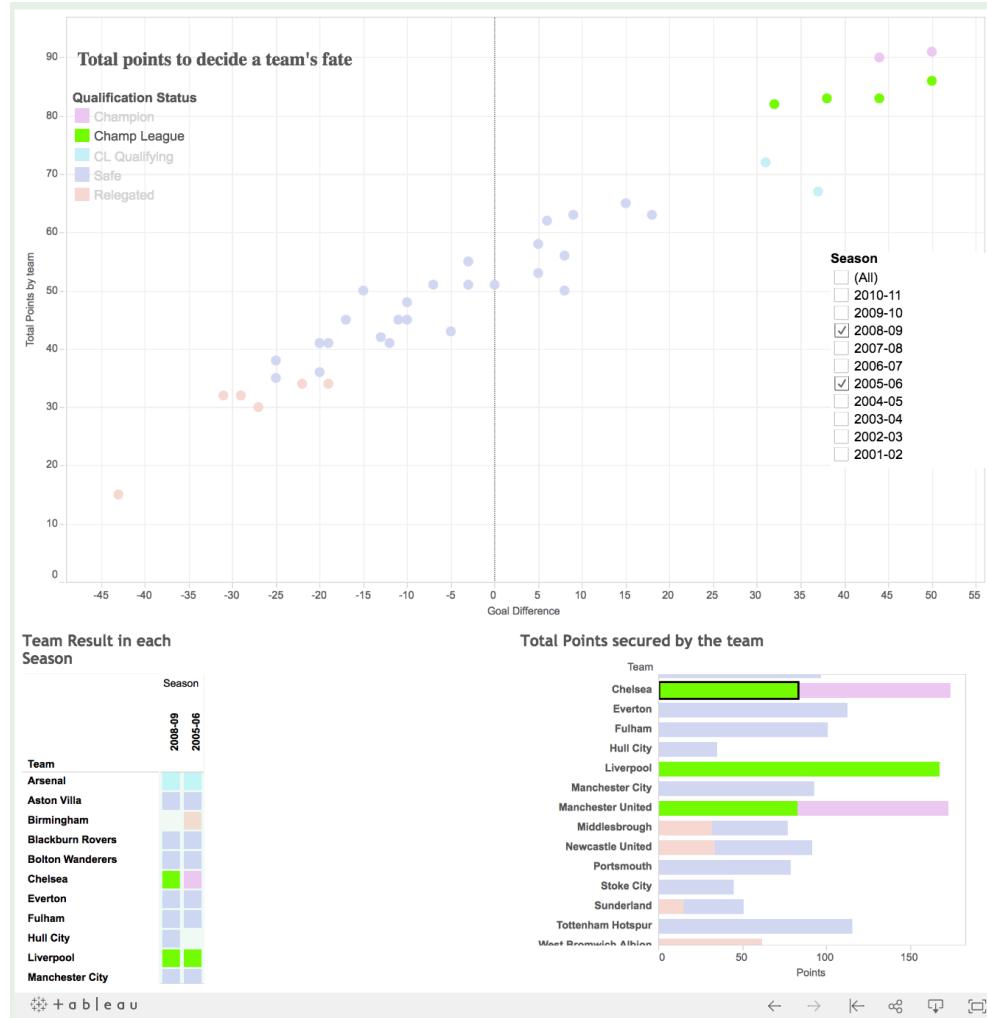
- Iconic Representations: Glyph (graphical object) represents a data case
- Visual properties of glyph represent different variables

# Trellis Display (Small Multiples)

- It subdivides space to enable comparison across multiple plots.
- Typically nominal or ordinal variables are used as dimensions for subdivision.



# Multiple Coordinated Views



# Minard 1869: Napoleon's March

*Figurative Map of the successive losses in men of the French Army in the Russian campaign 1812 ~ 1813*

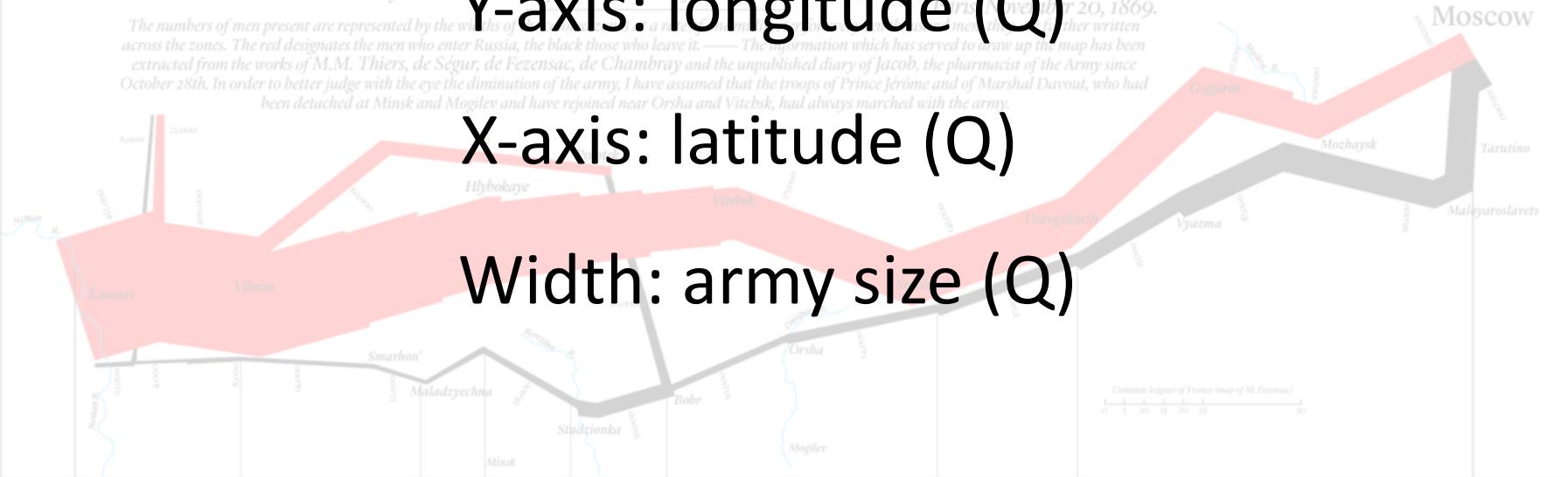
Drawn by M. Minard, Inspector General of Bridges and Roads (retired), Paris November 20, 1869.

The numbers of men present are represented by the widths of the lines which form a ribbon, the dimensions of which have been measured in the map itself after written across the zones. The red designates the men who enter Russia, the black those who leave it. — The information which has served to draw up the map has been extracted from the works of M.M. Thiers, de Ségrur, de Fezensac, de Chambray and the unpublished diary of Jacob, the pharmacist of the Army since October 28th. In order to better judge with the eye the diminution of the army, I have assumed that the troops of Prince Jérôme and of Marshal Davout, who had been detached at Minsk and Mogilev and have rejoined near Orsha and Vitebsk, had always marched with the army.

Y-axis: longitude (Q)

X-axis: latitude (Q)

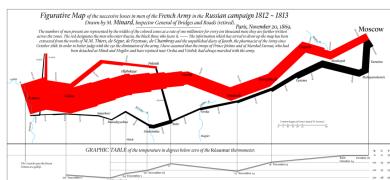
Width: army size (Q)



Y-axis: temperature (Q)

X-axis: longitude (Q) / time (O)

Depicts at least 5 quantitative variables.



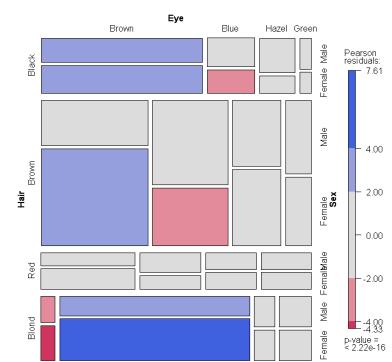
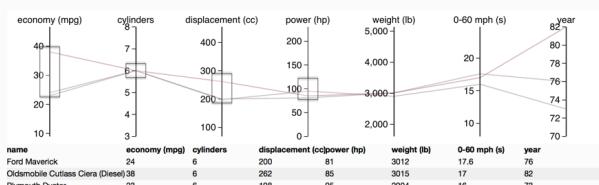
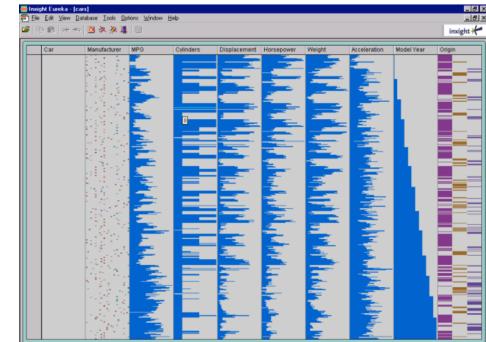
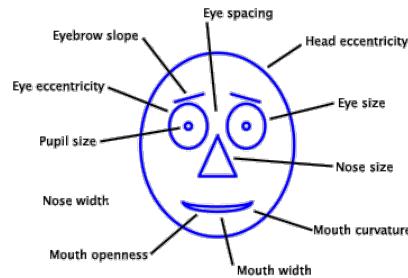
# Multivariate Data Visualization

- Visual Encodings: 8 dimensions?
- Focus: techniques can generally handle all data sets

	Characteristics				
	Selective	Associative	Quantitative	Order	Length
Position	• •	•• ••	↑ . . . . .	↑ . . . . .	Theoretically Infinite
Size	• ●	••●●•●		●>●>●>●	Selection: ~5 Distinction: ~20
Shape					Theoretically Infinite
Value	○●○○○○	●○○●○○●		○<○<○<●<●<●	Selection: <7 Distinction: ~10
Color	• ○	○○●●○●			Selection: <7 Distinction: ~10
Orientation	/ \   /				Theoretically Infinite
Texture	○○○○○○				Theoretically Infinite

# Common Multivariate Data Visualization Techniques

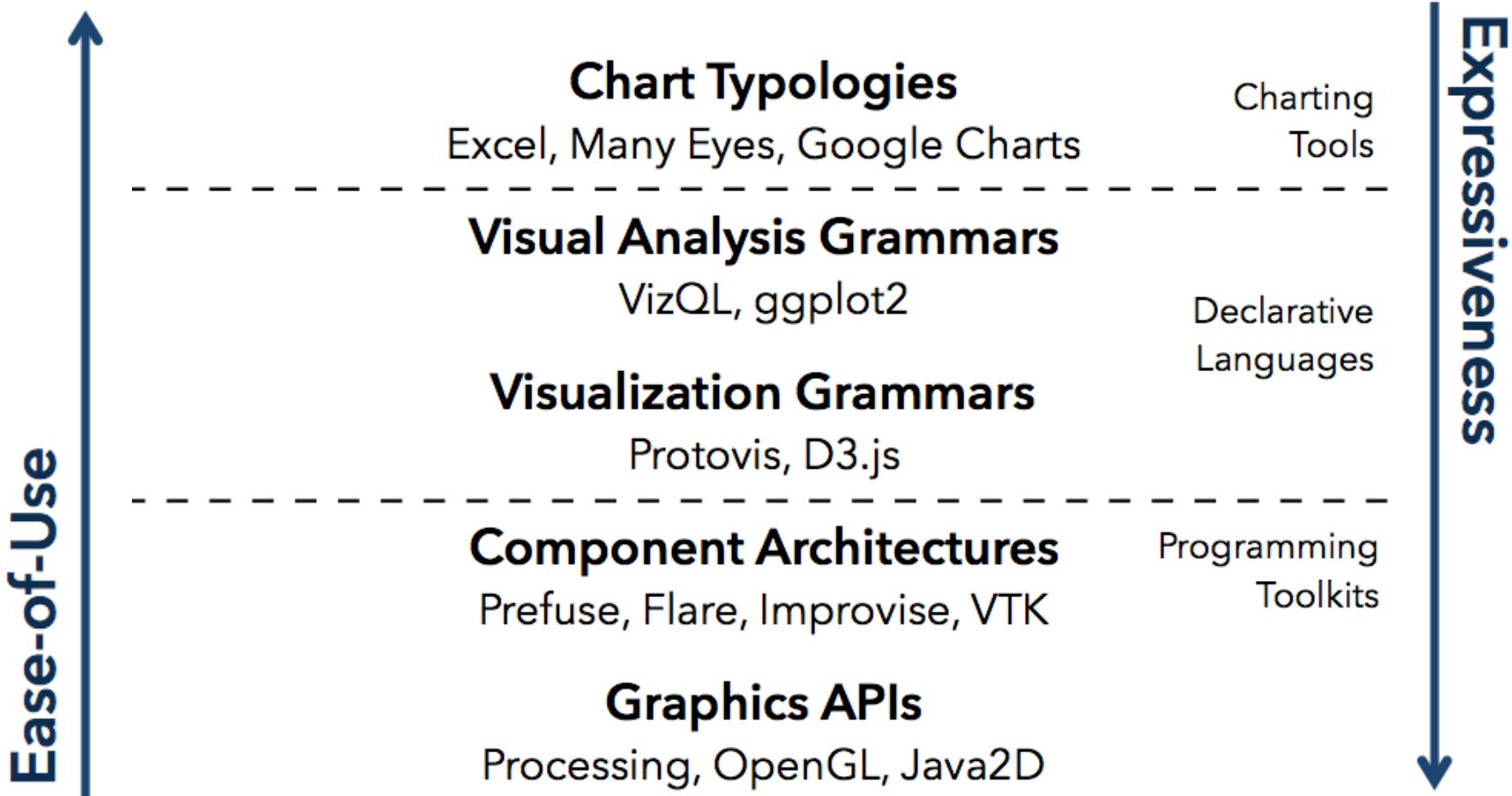
- Chernoff Faces
- Table Lens
- Parallel Coordinates
- Mosaic Plot



# Multivariate Data Visualization

- Strategies:
  - Avoid “over-encoding”
  - Use space and small multiples intelligently
  - Reduce the problem space
  - Use interaction to generate relevant views
- Rarely does a single visualization answer all questions. Instead, the ability to generate appropriate visualizations quickly is key.

# Visualization Tools



# The Advantages of Declarative Languages

- **Faster iteration.** Less code. Larger user base.
- **Better visualization.** Smart defaults.
- **Reuse.** Write-once, then re-apply.
- **Performance.** Optimization, scalability.
- **Portability.** Multiple devices, renderers, inputs.
- **Programmatic generation.** Write programs which output visualizations. Automated search & recommendation.

# Building a Plot in ggplot2

**data** to visualize (a data frame)

map variables to **aes**thetic attributes

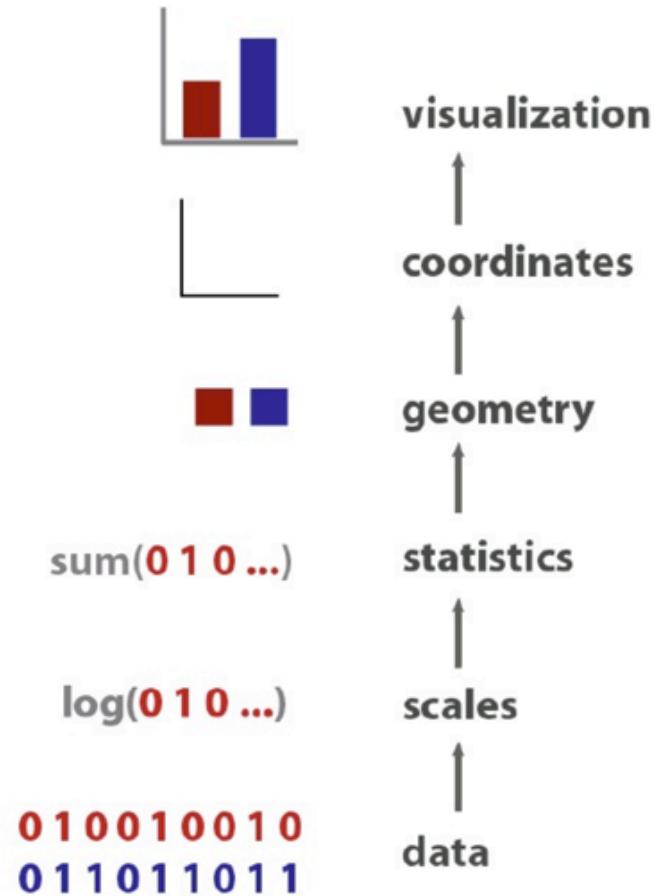
**geometric** objects – what you see (points, bars, etc)

**scales** map values from data to aesthetic space

**facet**ing subsets the data to show multiple plots

**stat**istical transformations – summarize data

**coord**inate systems put data on plane of graphic



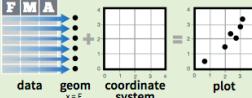
# Data Visualization with ggplot2

## Cheat Sheet

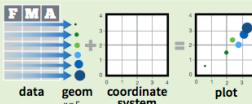


### Basics

**ggplot2** is based on the **grammar of graphics**, the idea that you can build every graph from the same few components: a **data** set, a set of **geoms**—visual marks that represent data points, and a **coordinate system**.



To display data values, map variables in the data set to aesthetic properties of the geom like **size**, **color**, and **x** and **y** locations.



Build a graph with **qplot()** or **ggplot()**

**aesthetic mappings**    **data**    **geom**

```
qplot(x = cty, y = hwy, color = cyl, data = mpg, geom = "point")
```

Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

**ggplot(data = mpg, aes(x = cty, y = hwy))**

Begins a plot that you finish by adding layers to. No defaults, but provides more control than **qplot()**.

**data**    **add layers, elements with +**    **layer = geom + default stat + layer specific mappings**    **additional elements**

```
ggplot(mpg, aes(hwy, cty)) +  
  geom_point(aes(color = cyl)) +  
  geom_smooth(method = "lm") +  
  coord_cartesian() +  
  scale_color_gradient() +  
  theme_bw()
```

Add a new layer to a plot with a **geom\_\***() or **stat\_\***() function. Each provides a geom, a set of aesthetic mappings, and a default stat and position adjustment.

**last\_plot()**

Returns the last plot

**ggsave("plot.png", width = 5, height = 5)**

Saves last plot as 5'x5' file named "plot.png" in working directory. Matches file type to file extension.

**Geoms** - Use a geom to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

### One Variable

#### Continuous

```
a <- ggplot(mpg, aes(hwy))
```



**a + geom\_area(stat = "bin")**  
x, y, alpha, color, fill, linetype, size  
b + geom\_area(aes(y = ..density..), stat = "bin")

**a + geom\_density(kernel = "gaussian")**  
x, y, alpha, color, fill, linetype, size, weight  
b + geom\_density(aes(y = ..count..))

**a + geom\_dotplot()**  
x, y, alpha, color, fill

**a + geom\_freqpoly()**  
x, y, alpha, color, linetype, size  
b + geom\_freqpoly(aes(y = ..density..))

**a + geom\_histogram(binwidth = 5)**  
x, y, alpha, color, fill, linetype, size, weight  
b + geom\_histogram(aes(y = ..density..))

#### Discrete

```
b <- ggplot(mpg, aes(fl))
```



**b + geom\_bar()**  
x, alpha, color, fill, linetype, size, weight

### Graphical Primitives

```
c <- ggplot(map, aes(long, lat))
```



**c + geom\_polygon(aes(group = group))**  
x, y, alpha, color, fill, linetype, size

```
d <- ggplot(economics, aes(date, unemploy))
```



**d + geom\_path(lineend = "butt", linejoin = "round", linemetre = 1)**  
x, y, alpha, color, linetype, size  
**d + geom\_ribbon(aes(ymin = unemploy - 900, ymax = unemploy + 900))**  
x, ymax, ymin, alpha, color, fill, linetype, size

```
e <- ggplot(seals, aes(x = long, y = lat))
```



**e + geom\_segment(aes(xend = long + delta\_long, yend = lat + delta\_lat))**  
x, xend, y, yend, alpha, color, linetype, size



**e + geom\_rect(aes(xmin = long, ymin = lat, xmax = long + delta\_long, ymax = lat + delta\_lat))**  
xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size

### Two Variables

#### Continuous X, Continuous Y

```
f <- ggplot(mpg, aes(cty, hwy))
```



**f + geom\_blank()**



**f + geom\_jitter()**  
x, y, alpha, color, fill, shape, size



**f + geom\_point()**  
x, y, alpha, color, fill, shape, size



**f + geom\_quantile()**  
x, y, alpha, color, linetype, size, weight



**f + geom\_rug(sides = "bl")**  
alpha, color, linetype, size



**f + geom\_smooth(model = lm)**  
x, y, alpha, color, fill, linetype, size, weight



**f + geom\_text(aes(label = cty))**  
x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust



**Discrete X, Continuous Y**  
g <- ggplot(mpg, aes(class, hwy))



**g + geom\_bar(stat = "identity")**  
x, y, alpha, color, fill, linetype, size, weight



**g + geom\_boxplot()**  
lower, middle, upper, x, ymax, ymin, alpha, color, fill, linetype, shape, size, weight



**g + geom\_dotplot(binaxis = "y", stackdir = "center")**  
x, y, alpha, color, fill



**g + geom\_violin(scale = "area")**  
x, y, alpha, color, fill, linetype, size, weight



**Discrete X, Discrete Y**  
h <- ggplot(diamonds, aes(cut, color))



**h + geom\_jitter()**  
x, y, alpha, color, fill, shape, size



**seals\$z <- with(seals, sqrt(delta\_long^2 + delta\_lat^2))**  
m <- ggplot(seals, aes(long, lat))



**m + geom\_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)**  
x, y, alpha, fill



**m + geom\_hex()**  
x, y, alpha, colour, fill size

### Continuous Bivariate Distribution

```
i <- ggplot(movies, aes(year, rating))
```



**i + geom\_bin2d(binwidth = c(5, 0.5))**  
xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size, weight



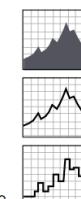
**i + geom\_density2d()**  
x, y, alpha, colour, linetype, size



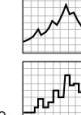
**i + geom\_hex()**  
x, y, alpha, colour, fill size

### Continuous Function

```
j <- ggplot(economics, aes(date, unemploy))
```



**j + geom\_area()**  
x, y, alpha, color, fill, linetype, size



**j + geom\_line()**  
x, y, alpha, color, linetype, size



**j + geom\_step(direction = "hv")**  
x, y, alpha, color, linetype, size

### Visualizing error

```
df <- data.frame(grp = c("A", "B"), fit = 4.5, se = 1.2)
```

```
k <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```



**k + geom\_crossbar(fatten = 2)**  
x, y, ymax, ymin, alpha, color, fill, linetype, size



**k + geom\_errorbar()**  
x, ymax, ymin, alpha, color, linetype, size, width (also **geom\_errorbarh()**)



**k + geom\_linerange()**  
x, ymin, ymax, alpha, color, linetype, size



**k + geom\_pointrange()**  
x, y, ymin, ymax, alpha, color, fill, linetype, shape, size

### Maps

```
data <- data.frame(murder = USArrests$Murder, state = tolower(rownames(USArrests)))
```



**map <- map\_data("state")**



**l <- ggplot(data, aes(fill = murder))**



**l + geom\_map(aes(map\_id = state), map = map) + expand\_limits(x = map\$long, y = map\$lat)**  
map\_id, alpha, color, fill, linetype, size

### Three Variables

```
seals$z <- with(seals, sqrt(delta_long^2 + delta_lat^2))
```



**m + geom\_contour(aes(z = z))**  
x, y, z, alpha, colour, linetype, size, weight



**m + geom\_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)**  
x, y, alpha, fill



**m + geom\_tile(aes(fill = z))**  
x, y, alpha, color, fill, linetype, size



# Visualizing Text

# Challenges

- High Dimensionality
  - Where possible use text to represent text...
  - ... which terms are the most descriptive?
- Context and Semantics
  - Provide relevant context to aid understanding.
  - Show (or provide access to) the source text.
- Modeling Abstraction
  - Determine your analysis task.
  - Understand abstraction of your language models.
  - Match analysis task with appropriate tools and models.

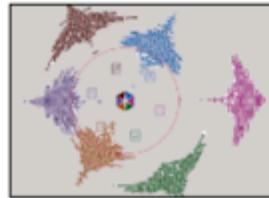
# Text Visualization



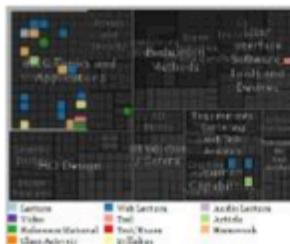
Visualizing text  
Showing words,  
phrases, and  
sentences



- Content
- Context
- Relationship to others



Visualization for IR  
Helping search



# Text Processing Pipeline

- Tokenization
  - Segment text into terms.
  - Remove stop words? [a](#), [an](#), [the](#), [of](#), [to](#), [be](#)
  - Numbers and symbols? [#gocard](#), [@nottinghamforestfbball](#)
  - Entities? [Nottingham](#), [Trump](#).
- Stemming
  - Group together different forms of a word.
  - Porter stemmer? [visualization\(s\)](#), [visualize\(s\)](#), [visually](#) -> [visual](#)
  - Lemmatization? [goes](#), [went](#), [gone](#) -> [go](#)
- Ordered list of terms

# Bag of Words Model

- Ignore ordering relationships within the text
- A document ≈ vector of term weights
  - Each dimension corresponds to a term (10,000+)
  - Each value represents the relevance
    - For example, simple term counts
- Aggregate into a document-term matrix
  - Document vector space model

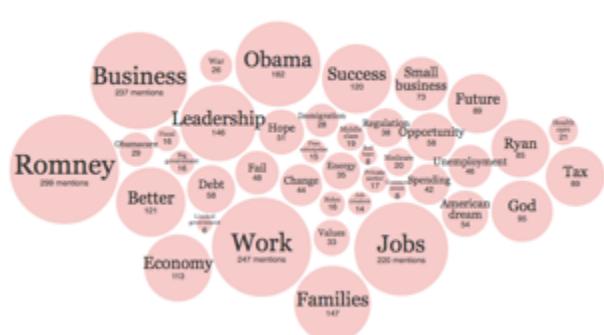
# Keyword Weighting

- Term Frequency
  - $tf_{td} = \text{count}(t) \text{ in } d$
  - Can take log frequency:  $\log(1 + tf_{td})$
  - Can normalize to show proportion  $(tf_{td} / \sum_t tf_{td})$
- TF.IDF: Term Freq by Inverse Document Freq
  - $tf.idf_{td} = \log(1 + tf_{td}) \times \log(N/df_t)$ 
    - $df_t$  = # docs containing t;
    - N = # of docs

# Word Counts and Tag Cloud

#### At the Republican Convention, the Words Being Used

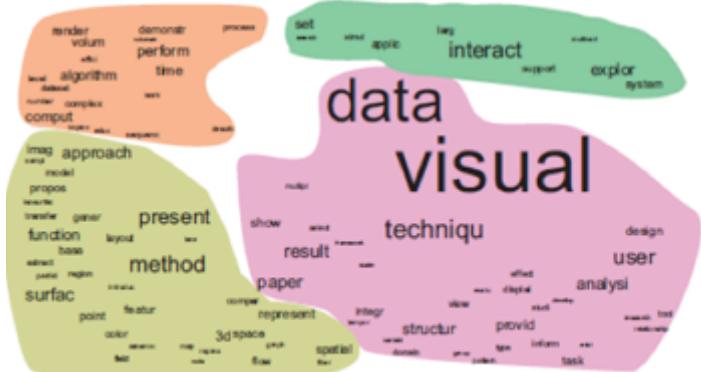
A look at how often speakers at the Republican National Convention have used certain words and phrases so far, based on an analysis of transcripts from the Federal News Service.



<http://www.nytimes.com/interactive/2012/08/28/us/politics/convention-word-counts.html>

Dr. Ke Zhou (<http://www.cs.nott.ac.uk/~pszkz/>)

# Context and Semantics

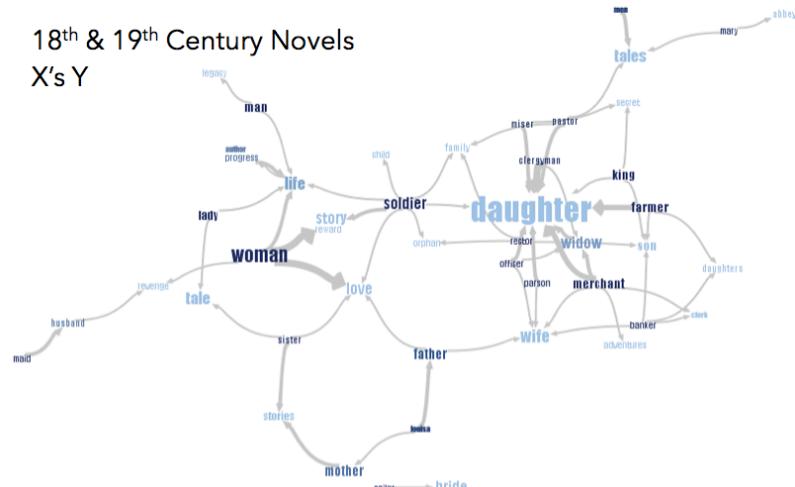


The screenshot shows the Concordance - Larkin.Concordance application window. The menu bar includes File, Text, Search, Edit, Headwords, Contexts, View, Tools, and Help. Below the menu is a toolbar with icons for opening files, saving, printing, and other functions. A vertical sidebar on the right contains buttons for 'Context', 'Text', 'Table', 'List', 'Index', and 'Note'. The main area displays a table of search results:

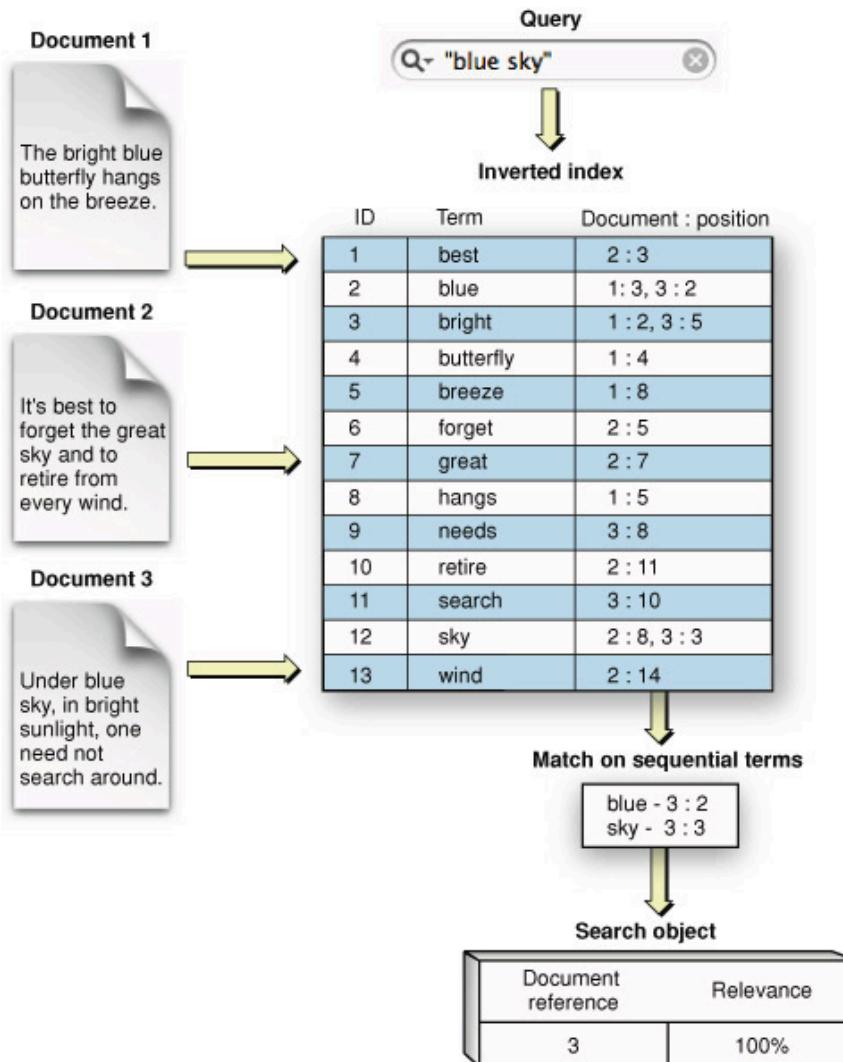
Headword	No.	Context...	Word	...Context	Reference
HEAR	15		That my own	heart	dritts and cries, having no...
HEARD	9		By the shout of the	heart	continually at work
HEARING	7	Nothing to adapt the skill of the	heart	to, still	And the wave And the wave
HEARS	3	The tread, the beat of it, it is my own	heart	.	Träumerlei
HEARSE	1	Because I follow it to my own	heart	.	many famous
<b>HEART</b>	<b>25</b>		My	is tickling like the sun:	I was washed i...
HEART'S	2		heart	sharpened to a candid co...	The March Pa...
HEART-SHAPED	1		Contract my	heart	Lines on a Yo
HEARTH	1		Having no	heart	Home is so Se...
HEARTS	7	And the boy pulling his	heart	out in the Gents	Essential Bea...
HEARTY	1	A harbour for the	heart	against distress.	Bridge for the...
HEAT	6	These I would choose my	heart	to lead	After-Dinner F...
HEAT-HAZE	1	Time in his little cinema of the	heart	.	Time and Spac...
HEATH	1	This petrified	heart	has taken,	A Stone Ch...
HEATS	1	How should they sweep the girl clean	heart	.	I see a girl dra...
HEAVE	1	Hands that the	heart	can govern	Heaviest of th...
HEAVEN	4	For the	heart	to be loveless, and as col...	Dawn
HEAVEN-HOLDING	1	With the unguessable	heart	riding	One man walk...
HEAVIER-THAN...	1	If hands could free you,	heart	.	If hands could
HEAVIEST	2	That overflows the	heart	.	Pour away thi...

At the bottom, there are buttons for 'Words' (7318), 'Tokens' (37070), 'At word' (2990), 'Deleted lines' (1 [24]), 'Word sort' (Asc alpha [string]), 'Context sort' (Asc occurrence order), and a large 'Next' button.

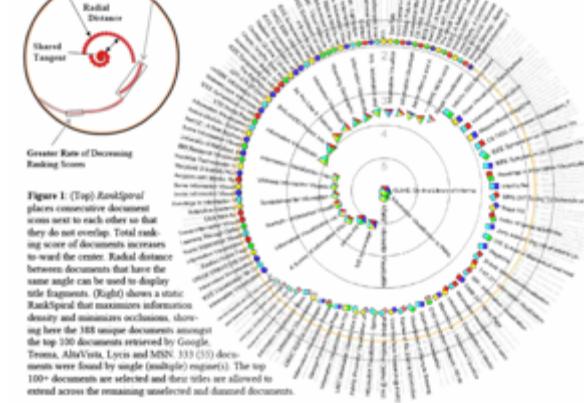
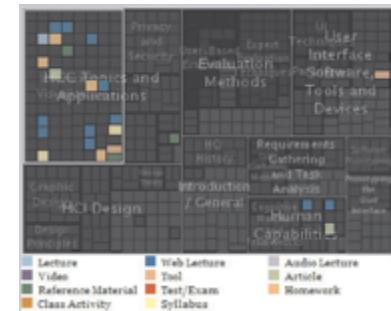
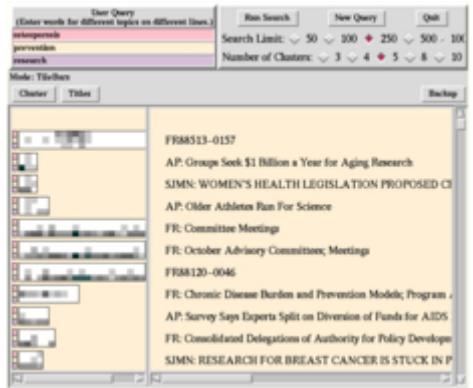
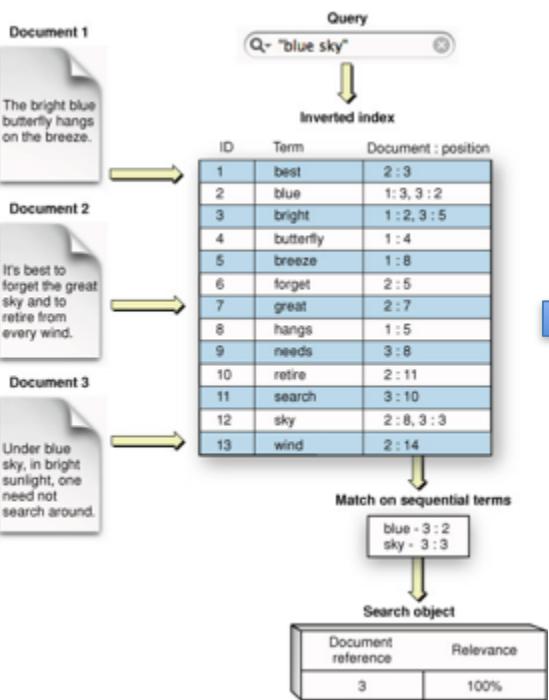
## 18<sup>th</sup> & 19<sup>th</sup> Century Novels



# Information Retrieval (IR)



# Visualizing for IR



# Visualizing Time Series Data

# Tasks

- Often asked questions:
  - when was something greatest/least?
  - is there a pattern? are two series similar?
  - does a data element exist at time t, and when?
  - how long does a data element exist and how often?
  - how fast are data elements changing
  - in what order do data elements appear?
  - do data elements exist together?

(Optional Reading) Müller, Wolfgang, and Heidrun Schumann. "Visualization for modeling and simulation: visualization methods for time-dependent data—an overview." Proceedings of the 35th conference on Winter simulation: driving innovation. Winter Simulation Conference, 2003.

# Visualizing Time-oriented Data

## A Systematic View

**The TimeViz Browser**  
A Visual Survey of Visualization Techniques for Time-Oriented Data  
by Christian Tominski and Wolfgang Aigner

# of Techniques: 115

Search:

How to use filters:

- Want: Show me!
- Indifferent: I don't care.
- Hide: I'm not interested!

Data

Frame of Reference

- Abstract
- Spatial

Number of Variables

- Univariate
- Multivariate

Time

Arrangement

- Linear
- Cyclic

Time Primitives

- Instant
- Interval

Visualization

Mapping

- Static
- Dynamic



Aigner, Wolfgang, et al. "Visualizing time-oriented data—a systematic view." Computers & Graphics 31.3 (2007): 401-409.

<http://www.timeviz.net/>

Dr. Ke Zhou (<http://www.cs.nott.ac.uk/~pszkz/>)

# Taxonomy

Time	Temporal primitives	time points (a) (b) (c) (d) (e) (f) (g) (i)									time intervals (g) (h)			
	Structure of time	linear (a) (b) (c) (d) (f) (g) (h) (i)					cyclic (e)			branching (h)				
Data	Frame of reference	abstract (c) (d) (f) (g) (h) (i)							spatial (a) (b) (e) (i)					
	Number of variables	univariate (a) (b) (f) (g) (h)							multivariate (c) (d) (e) (i)					
	Level of abstraction	data (a) (b) (c) (d) (e) (f) (g) (h) (i)							data abstractions (b) (g) (i)					
Representation	Time dependency	static (c) (d) (e) (g) (h) (i)							dynamic (a) (b) (f) (i)					
	Dimensionality	2D (a) (c) (d) (g) (h) (i)							3D (b) (e) (f) (i)					

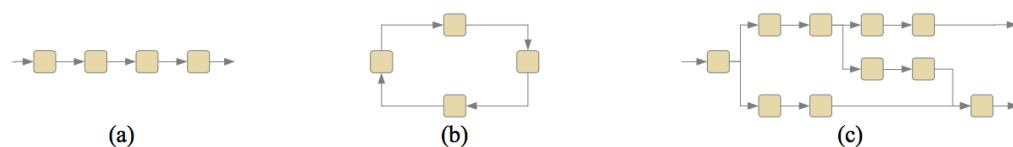
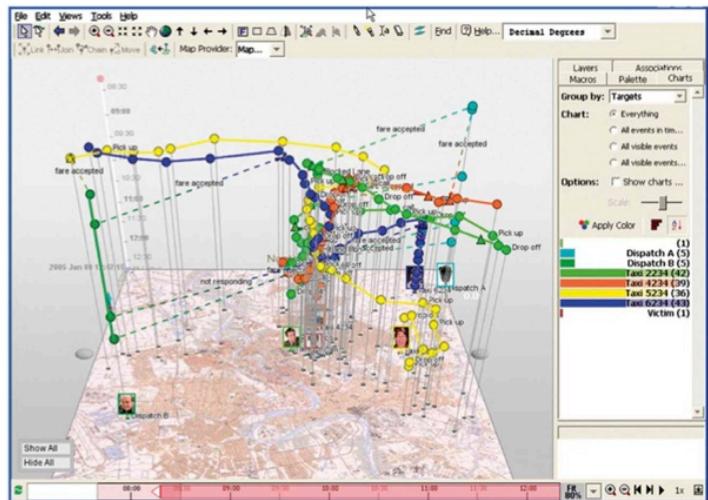
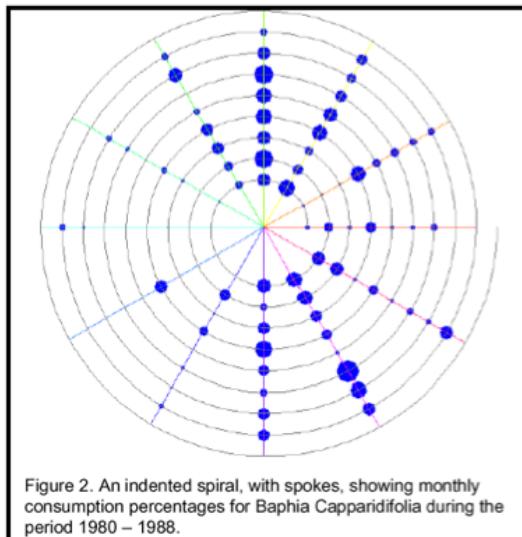
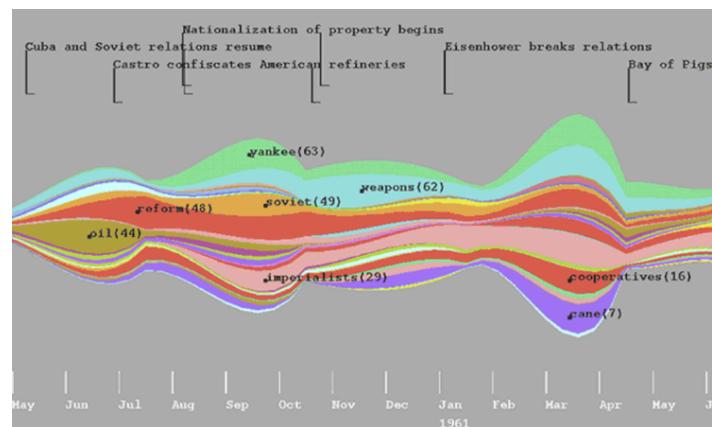
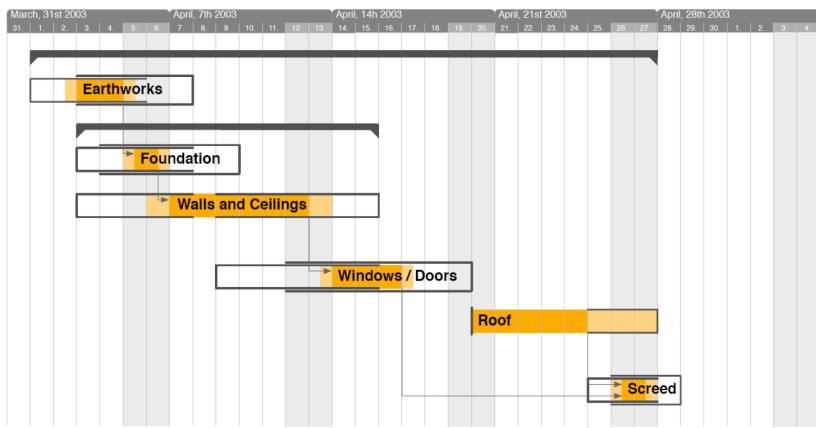


Fig. 2. Structure of time: (a) Linear time; (b) Cyclic time; (c) Branching time.

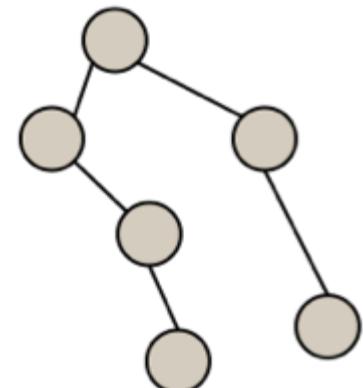
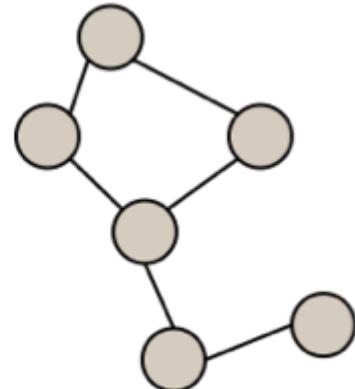
# Time Series Visualization



# Visualizing Trees and Graphs

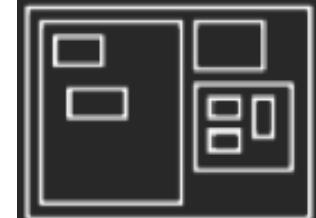
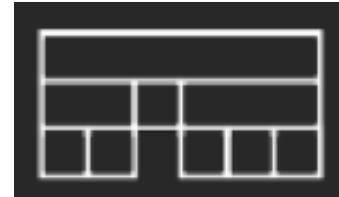
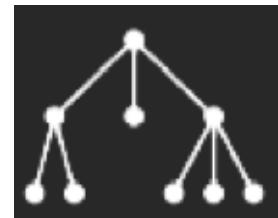
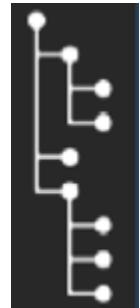
# Graphs and Trees

- Graphs
  - Model relations among data
  - Nodes and edges
- Trees
  - Graphs with hierarchical structure
    - Connected graph with  $N-1$  edges
  - Nodes as parents and children



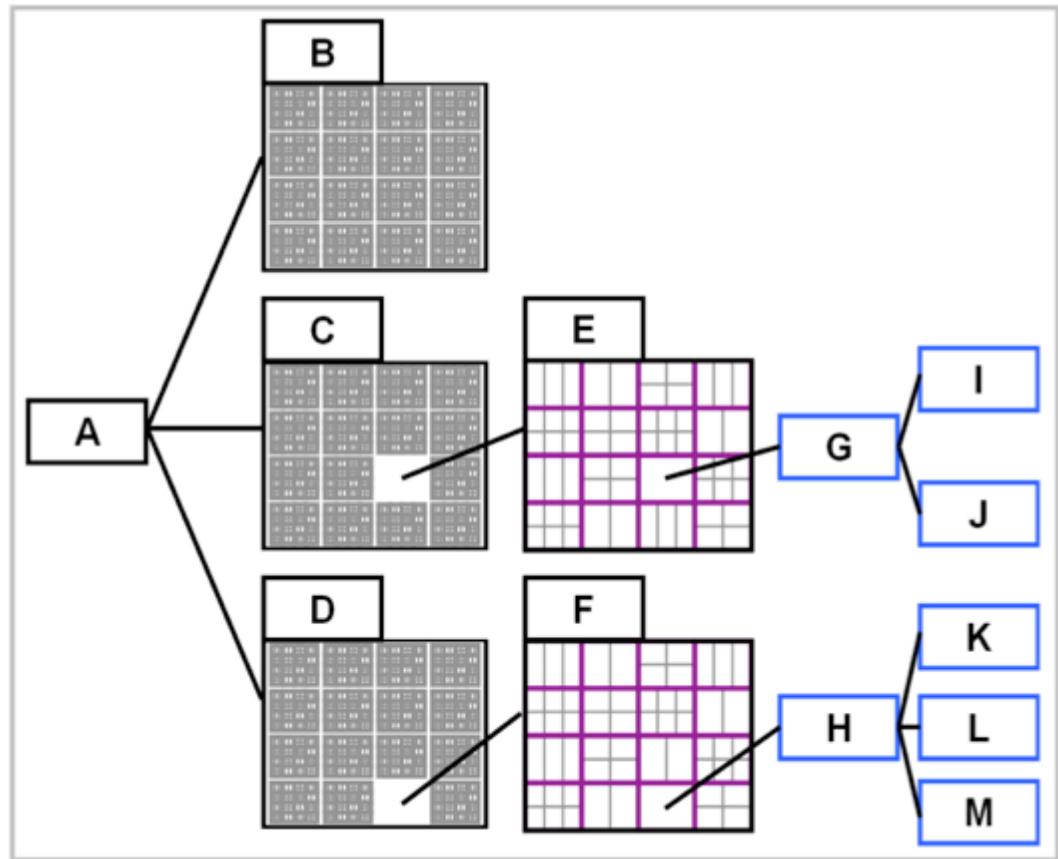
# Tree Visualizations

- Indented lists
  - Linear list, indentation encodes depth
- Node-link trees
  - Nodes connected by lines/curves
- Layered diagrams
  - Relative position and alignment
- Treemaps
  - Represent hierarchy by enclosure

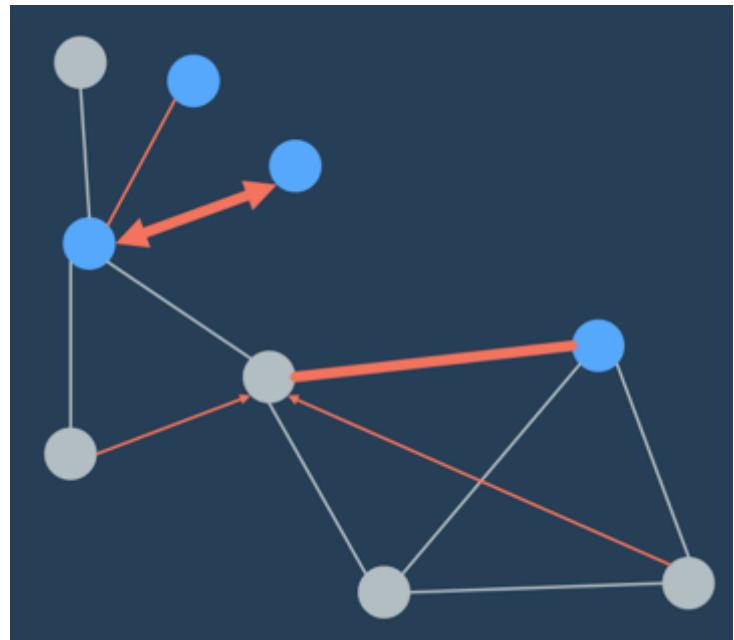
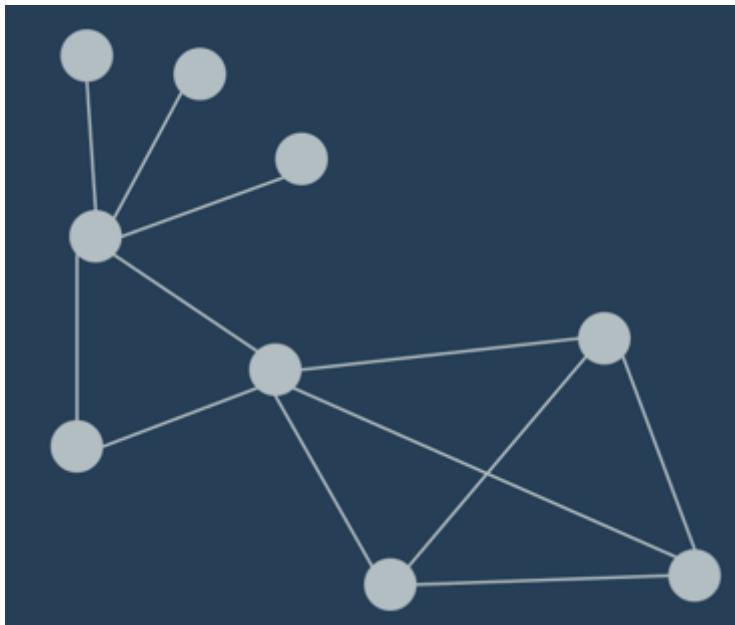


# Hybrids

- Elastic Hierarchies
  - Node-link diagram with treemap nodes.
- Video:  
<https://www.youtube.com/watch?v=nvslqYQ75yA>

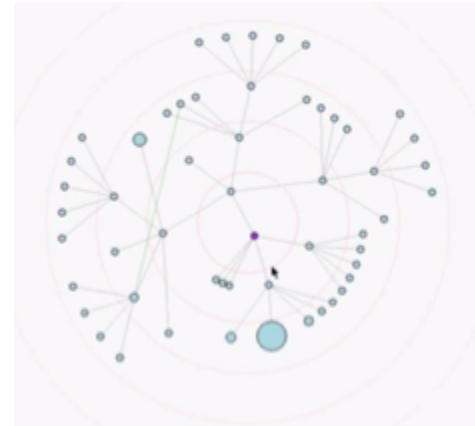


# What's in a Graph?



# Graph Visualization

- Two representations:
  - Node-link diagrams
  - Matrices



- Major Node-Link Layouts

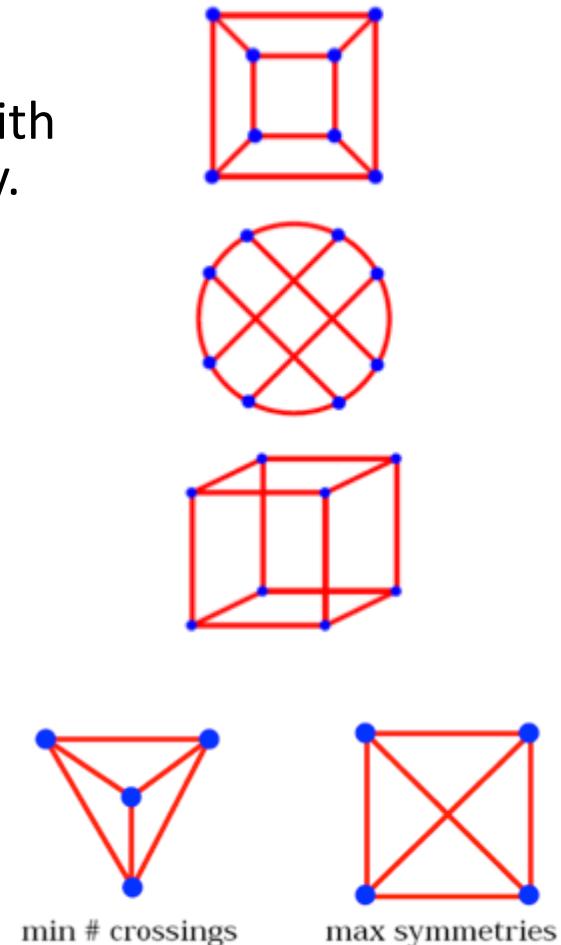
Tree in the Graph



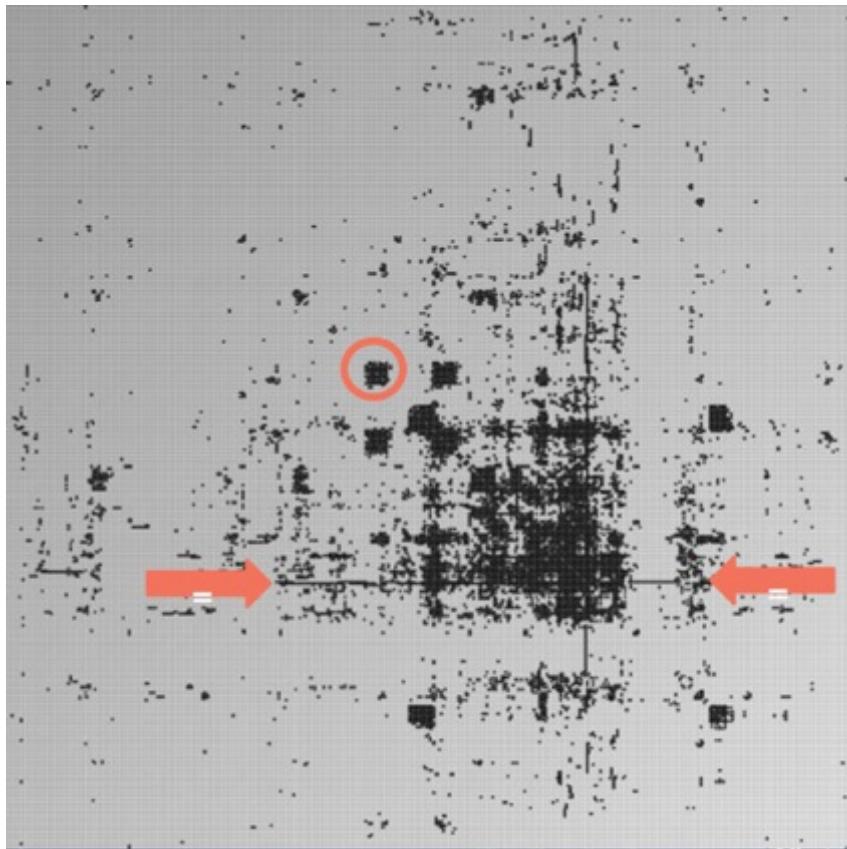
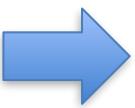
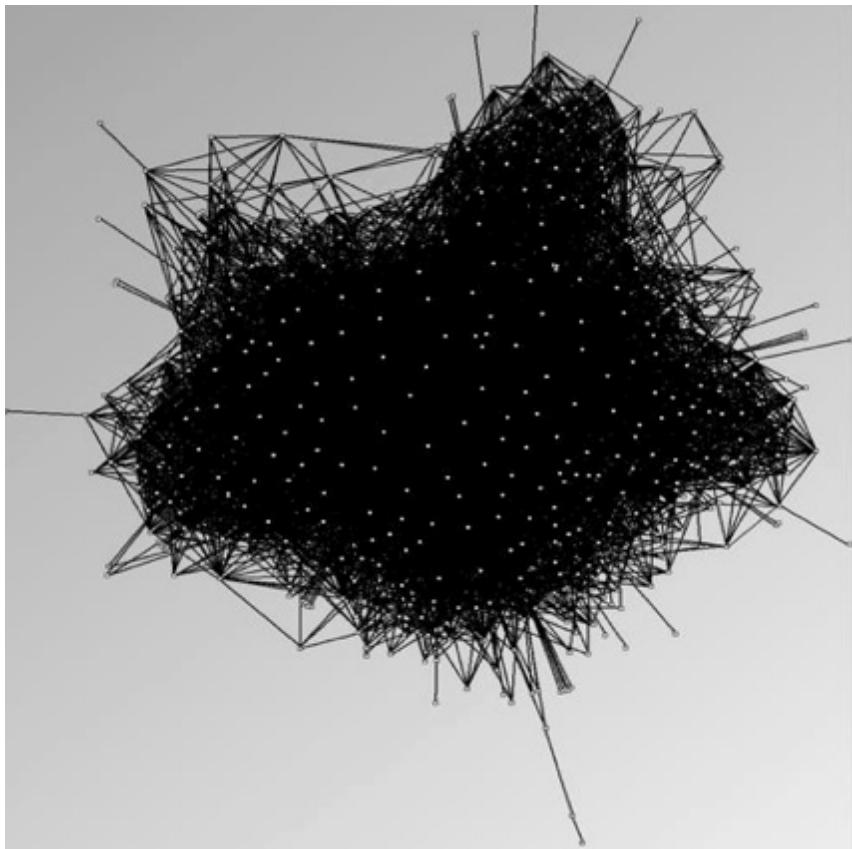
Hierarchical Graph Layout

# Optimization Techniques

- Treat layout as an optimization problem
  - Define layout using an energy model along with constraints: equations the layout should obey.
  - Use optimization algorithms to solve
- Commonly posed as a physical system
  - Charged particles, springs, drag force, ...
- Different constraints can be introduced
  - Minimize edge crossings
  - Minimize area
  - Minimize line bends
  - Minimize line slopes
  - Maximize smallest angle between edges
  - Maximize symmetry



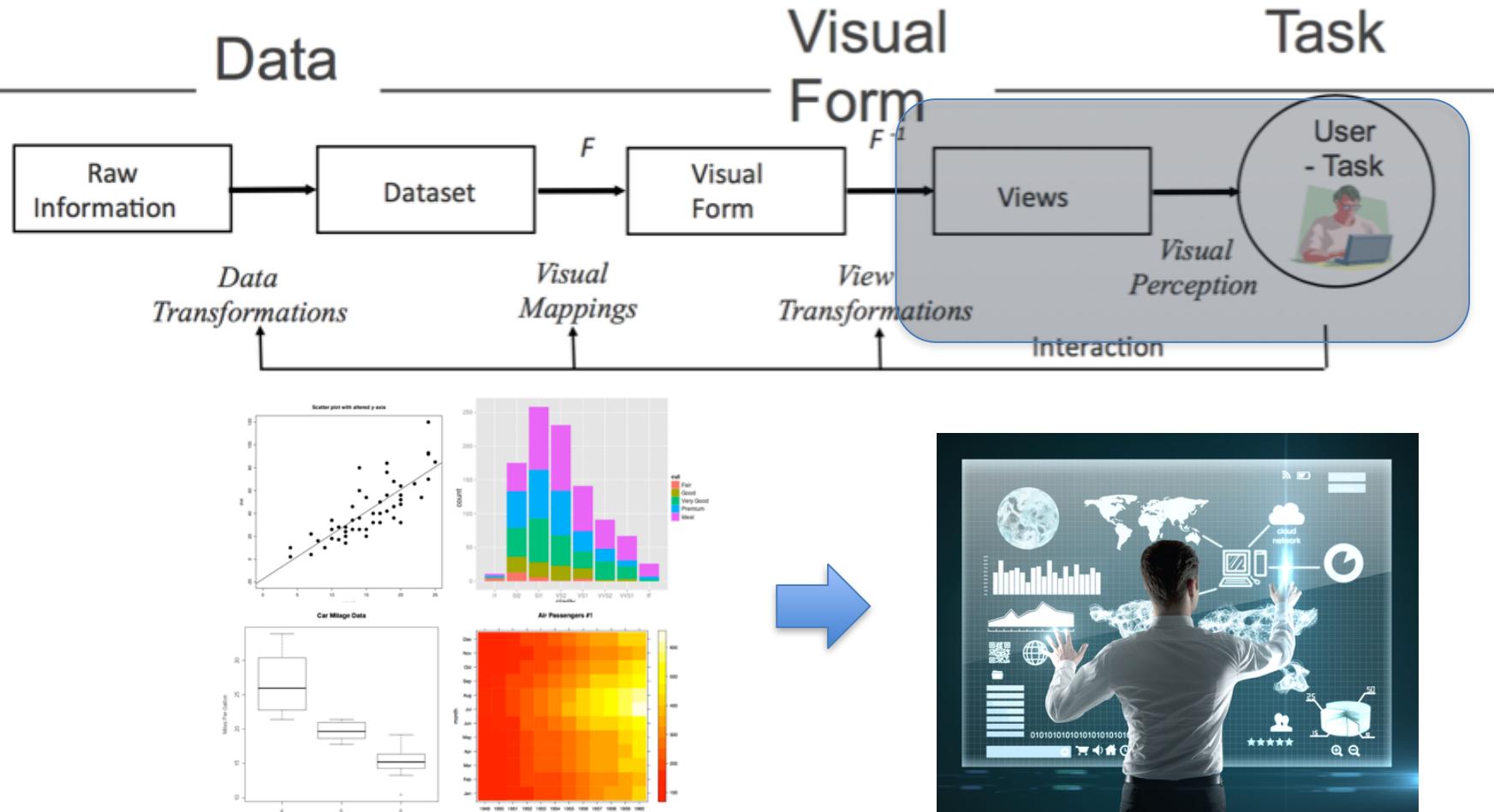
# Matrices



# Matrix vs. Node-link

Matrix	Node-link
Require learning	<b>Familiar</b>
<b>No overlap</b>	Node overlap
<b>No crossings</b>	Link crossing
Use a lot of space	<b>More compact</b>
<b>Dense graphs</b>	<b>Sparse graphs</b>

# Information Visualization



**Design:** how can we design the visualization best for human to accomplish their tasks?

# What design criteria should we follow?

# What Design Criteria to Follow?

- **Expressiveness**

- A set of facts is expressible in a visual language if the sentences (i.e. the visualizations) in the language (1) **express all the facts** in the set of data, and (2) **only the facts** in the data.

***Tell the truth***

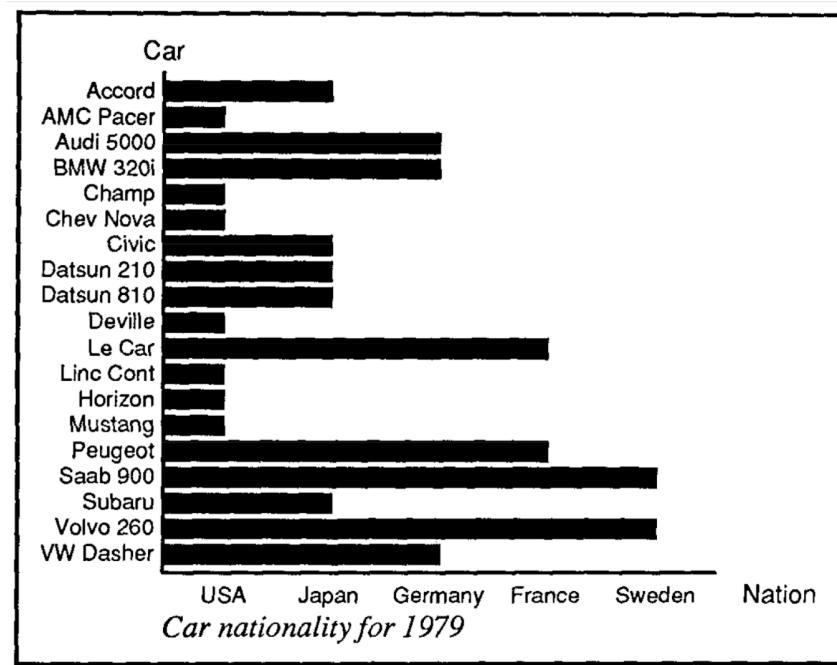
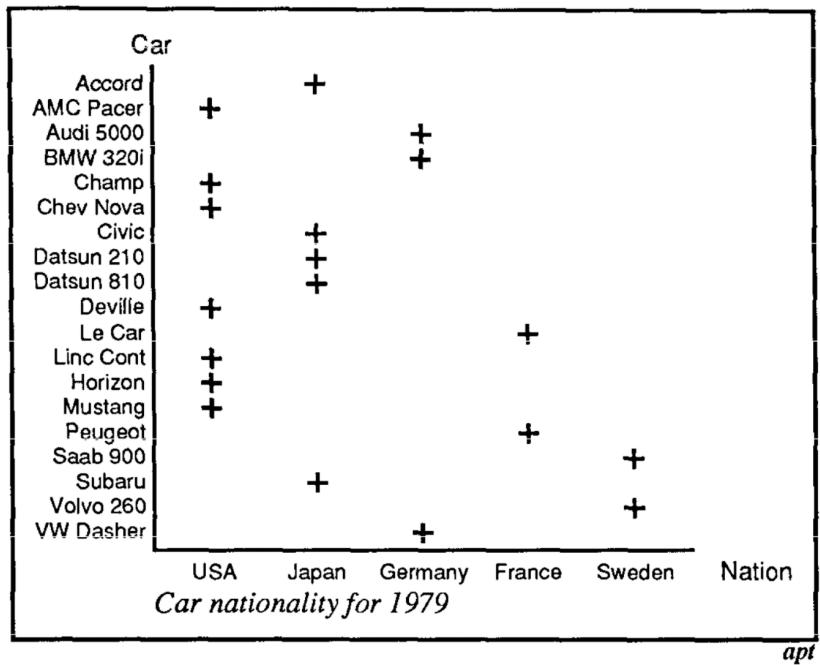
- **Effectiveness**

- A visualization is more effective than another visualization if the information conveyed by one visualization is more readily perceived than the information in the other visualization.

***Use proper encoding***

Mackinlay, Automating the design of graphical presentations of relational information, 1986.

# Expressiveness



An Alternative

# What Design Criteria to Follow?

- **Expressiveness**

- A set of facts is expressible in a visual language if the sentences (i.e. the visualizations) in the language express (1) all the facts in the set of data, and (2) only the facts in the data.

***Tell the truth***

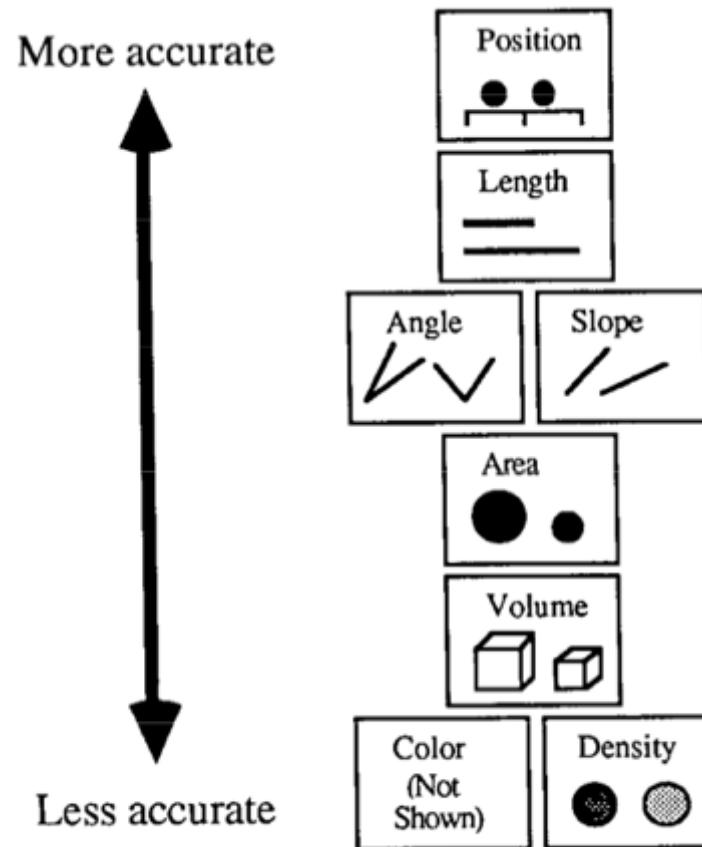
- **Effectiveness**

- A visualization is more effective than another visualization if the information conveyed by one visualization **is more readily perceived** than the information in the other visualization.

***Use proper encoding***

Mackinlay, Automating the design of graphical presentations of relational information, 1986.

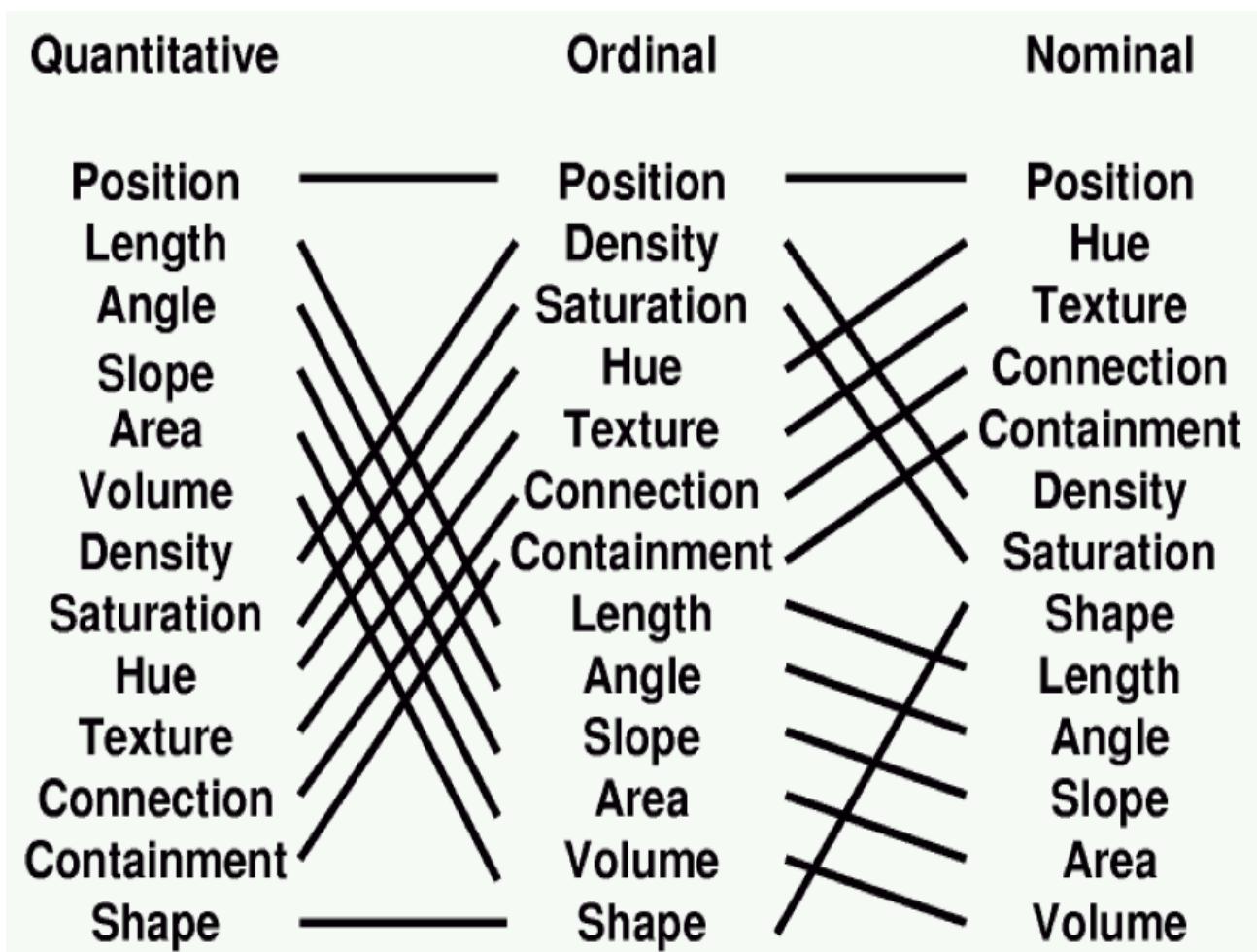
# Effectiveness: Accuracy Ranking for Quantitative Information



Mackinlay, Automating the design of graphical presentations of relational information, 1986.

# Conjectured Effectiveness of Encodings by Data Type

- Nominal/  
Ordinal  
variables:  
detect  
differences
- Quantitative  
variables:  
estimate  
magnitudes

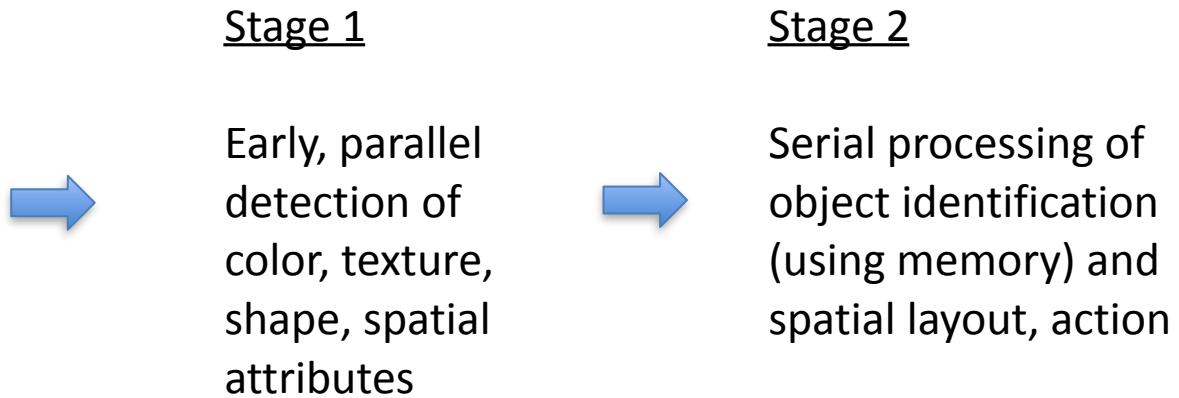
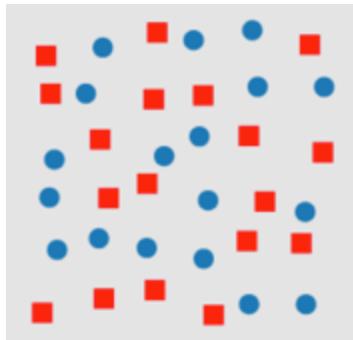


Mackinlay, Automating the design of graphical presentations of relational information, 1986.

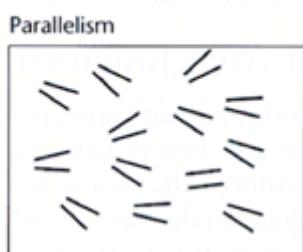
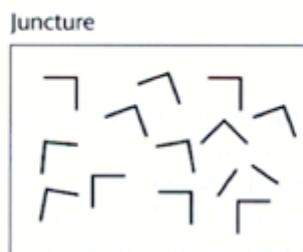
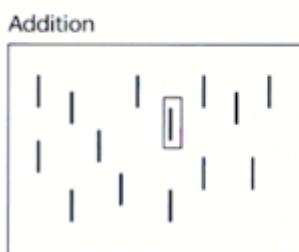
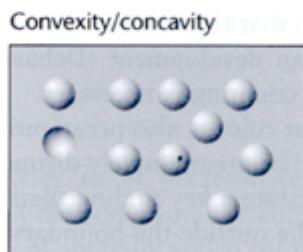
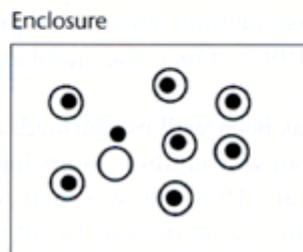
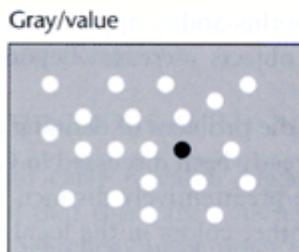
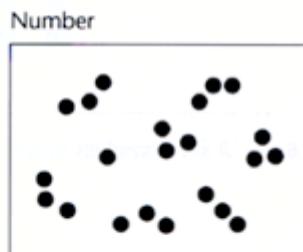
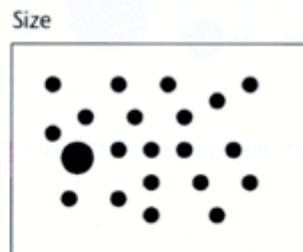
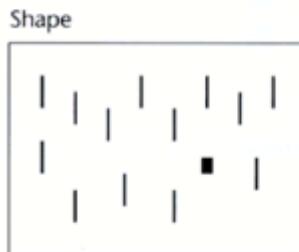
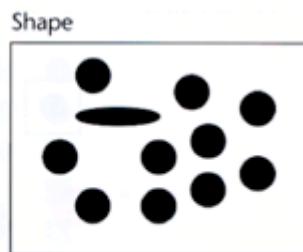
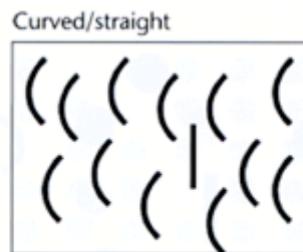
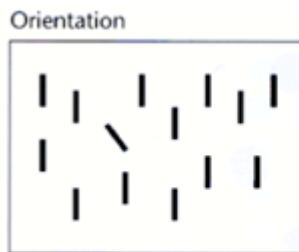
# Visual Perception

# Perceptual Processing Model

- Two stage process
  - Parallel extraction of low-level properties of scene
  - Sequential goal-directed processing



# Pre-Attentive Features



- length
- width
- size
- curvature
- number
- terminators
- intersection
- closure
- hue
- intensity
- flicker
- direction of motion
- binocular lustre
- stereoscopic depth
- 3-D depth cues
- lighting direction

# Pre-Attentive Feature Conjunctions

- Spatial conjunctions are often pre-attentive
- Motion and 3D disparity
- Motion and color
- Motion and shape
- 3D disparity and color
- 3D disparity and shape
- Most conjunctions are not pre-attentive

# Gestalt Grouping Principles

“All else being equal, elements that are related by X tend to be grouped perceptually into higher-order units.”

— Stephen Palmer

- Proximity
- Similarity
- Connectedness
- Continuity
- Symmetry
- Closure
- Figure/Ground
- Common Fate

# Change Blindness

- We don't always see everything that is there!
- Is the viewer able to perceive changes between two scenes?
  - If so, may be distracting
  - Can do things to minimize noticing changes
- Video: <http://www.simonslab.com/videos.html>

# Summary of Design Criteria

- Choose expressive and effective encodings
  - Rule-based tests of expressiveness
  - Perceptual effectiveness rankings
    - Prioritizes encodings that are most easily/accurately interpreted
    - Principle of Importance Ordering: Encode more important information more effectively (Mackinlay)

# Interaction

# Representation and Interaction

- Two main components of information visualization
- Very challenging to come up with innovative, new visual representations
- But can do interesting work with how user interacts with the view or views
  - Analysis is a process, often iterative with different interactions

“The effectiveness of information visualization hinges on two things: its ability to clearly and accurately represent information and our ability to interact with it to figure out what the information means.”

S. Few, <Now you see it>

# Taxonomy of Interactions

- Dix and Ellis (1998)
  - Highlighting and focus;
  - accessing extra info;
  - overview and context;
  - same representation, changing parameters;
  - Linking representations
- Keim (2002)
  - Projection
  - Filtering
  - Zooming
  - Distortion
  - Linking and brushing
- Few's Principles
  - Comparing
  - Sorting
  - Adding variables
  - Filtering
  - Highlighting
  - Aggregating
  - Re-expressing
  - Re-visualizing
  - Zooming and panning
  - Re-scaling
  - Accessing details on demand
  - Annotating
  - Bookmarking

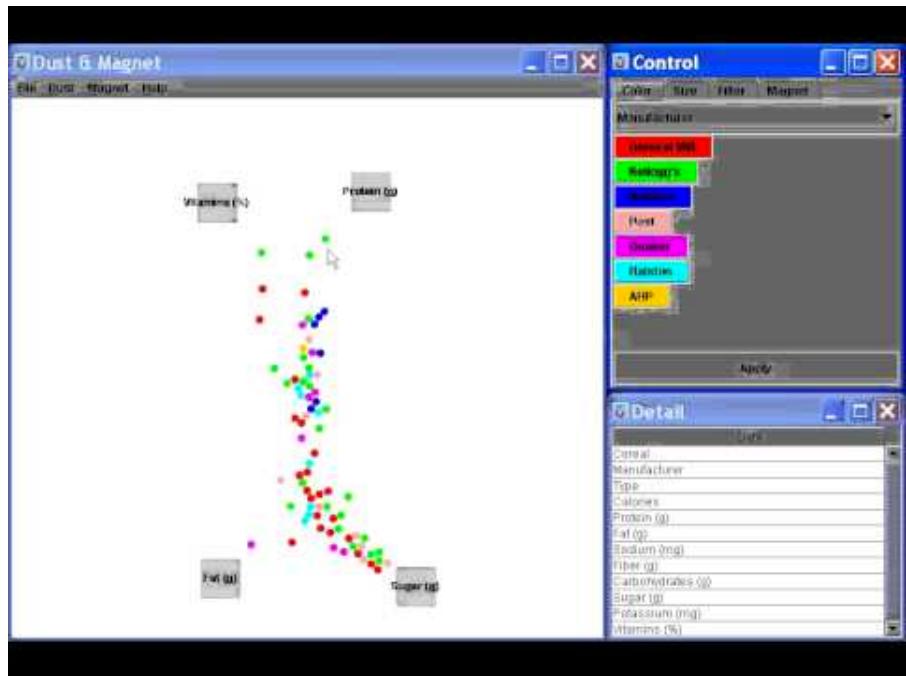
# A Summary of Existing Taxonomy

- Survey
  - 59 papers
    - Papers introducing new interaction systems
    - Well-known papers in subareas of information visualization
  - 51 systems
    - Commercial Infovis Systems (SeeIT, Spotfire, TableLens, InfoZoom, etc.)
  - Collected 311 individual interaction techniques
- Affinity Diagram Method

Yi, Ji Soo, Youn ah Kang, and John Stasko. "Toward a deeper understanding of the role of interaction in information visualization." IEEE transactions on visualization and computer graphics 13.6 (2007): 1224-1231.

# Taxonomy of Interactions based on User Intent

- 7 Categories
- Select
- Explore
- Reconfigure
- Encode
- Abstract/Elaborate
- Filter
- Connect



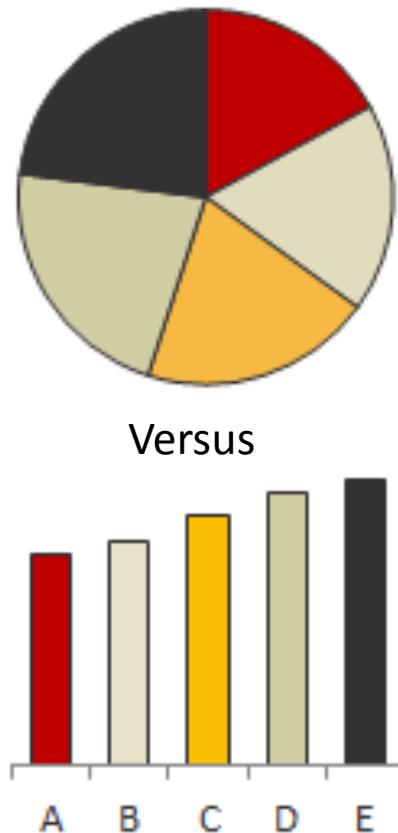
# Interaction is Vital for Exploration

- Interaction facilitates a dialog between the user and the visualization system
- Multiple views amplify importance of interaction
- Interaction often helps when you just can't show everything you want

# Evaluation

# How do We Evaluate Visualizations?

- How do we evaluate visualizations?
  - Usability vs. Utility
- What evaluation techniques should we use?
- What do we measure?
  - What data do we gather?
  - What metrics do we use?



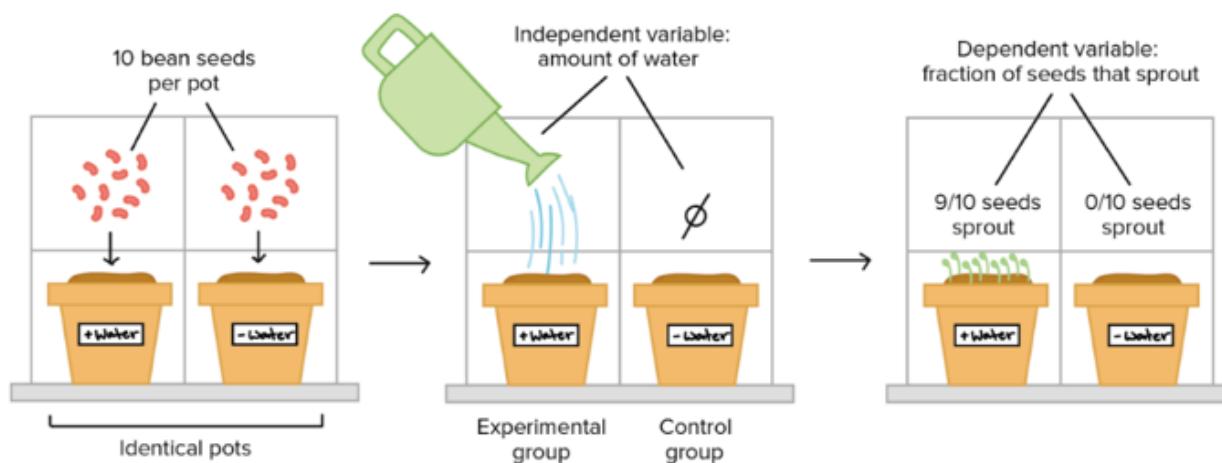
# Evaluation Approaches

- Many different Forms
  - Qualitative, quantitative, objective, subjective, controlled experiments, interpretive observations, ...
- Two popular methodologies
  - Controlled experiments (Quantitative)
  - Subjective assessments (Qualitative)

# Quantitative Methods: Controlled Experiments

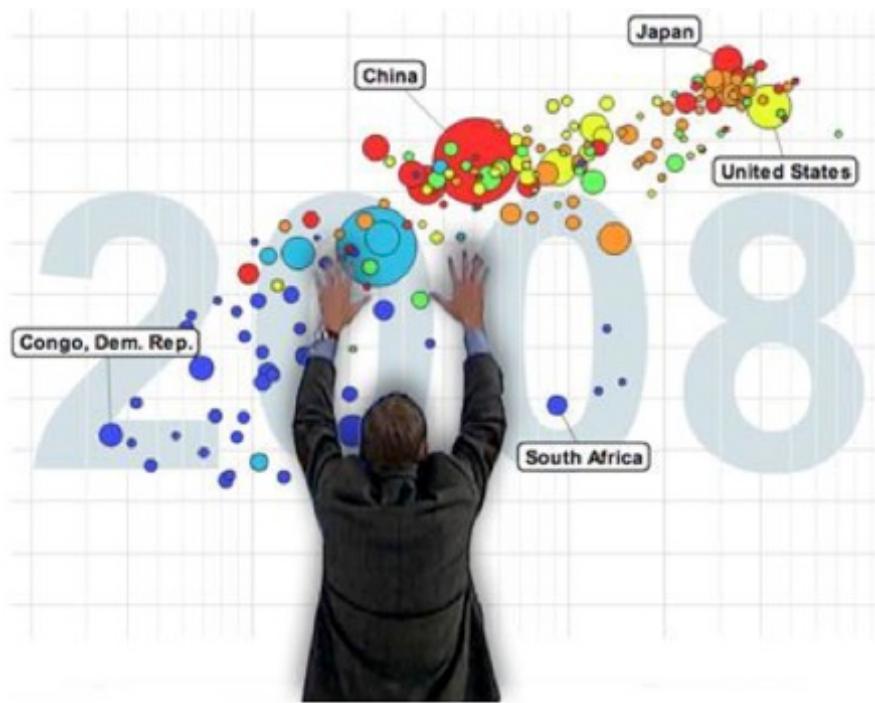
- Good for measuring performance or comparing multiple techniques
- What do we measure?
  - Performance, time, errors,

...



# An Example: Controlled Experiment

- Run an experiment to evaluate three visualization strategies
  - Animation
  - Small multiples
  - Traces
- Especially interested in examining whether animated bubble charts are beneficial for analysis and presentation



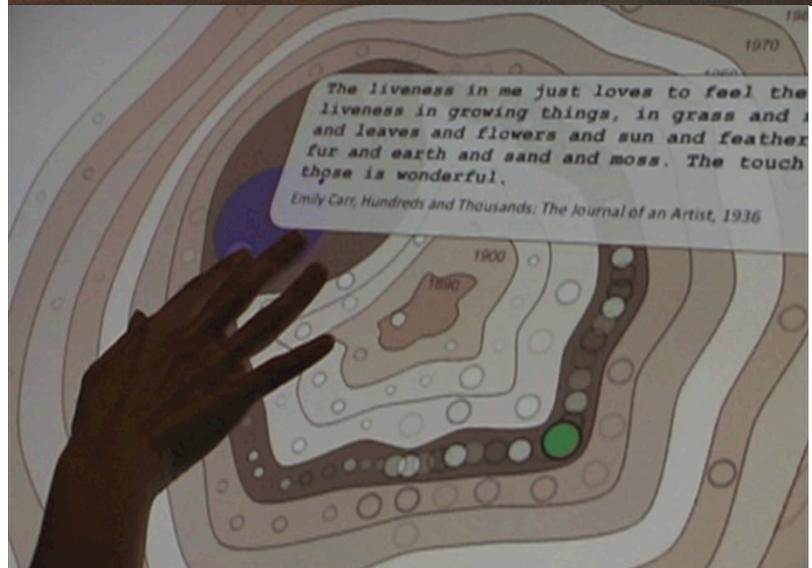
*Do you remember Hans Rosling's TED talk?  
(Lecture 2)*

# Qualitative Methods

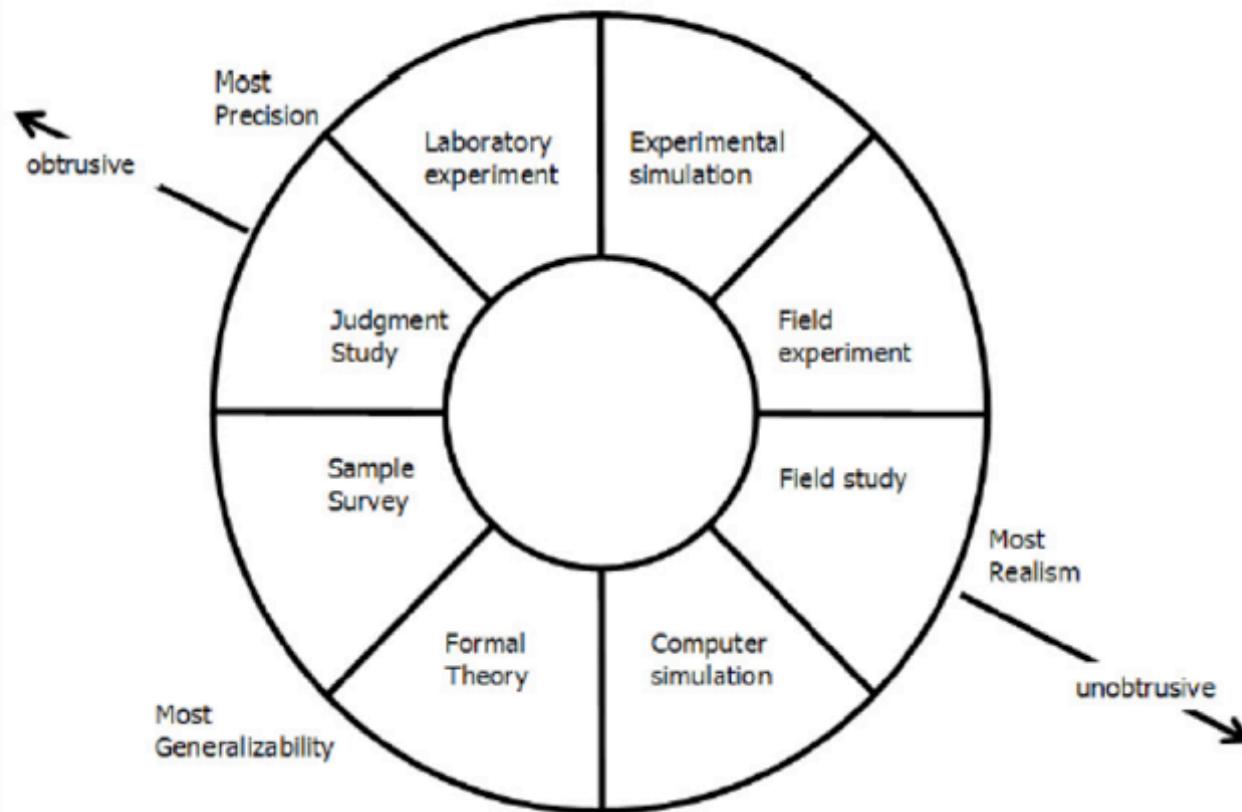
- Types
  - Nested methods
    - Experimenter observation, think-aloud protocol, collecting participant opinions
  - Inspection evaluation methods
    - Heuristics to judge
- Observational context
  - In situ, laboratory, participatory
  - Contextual interviews is important

# An Example: Subjective Assessments

- Evaluating a newly developed visualization system (EMDialog) at the museum (Emily Carr exhibit)
- Discourse visualization
  - Time
  - Context



# Methodology vs. Desirable Features

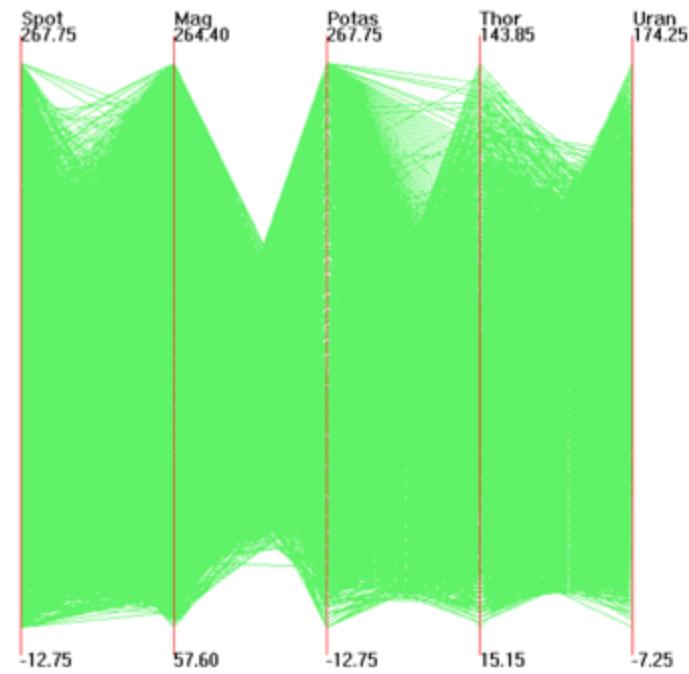


**Fig. 1.** Types of methodologies organized to show relationships to precision, generalizability and realism. (adapted, simplified from McGrath 1995)

# Data Overload

# Data Overload

- Most of the techniques we've examined work for a modest number of data cases or variables
- What happens when you have lots and lots of data cases and/or variables?



Out5d dataset(5 dimensions, 16384 items)

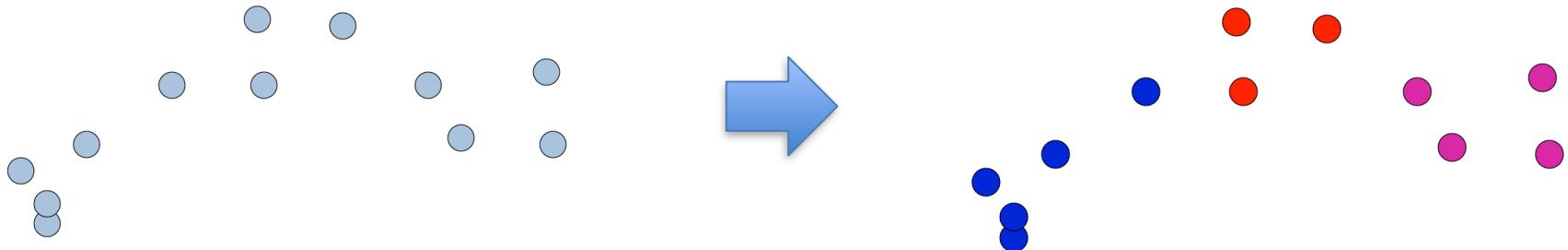
# General Solution

- Data that is similar in most dimensions ought to be drawn together
  - Cluster at high dimensions
- Need to project the data down into the plane and give it some ultra-simplified representation
- Or perhaps only look at certain aspects of the data at any one time

# Clustering and Dimensionality Reduction

- There exist many techniques for clustering high-dimensional data with respect to all those dimensions (too many data points)
  - Affinity propagation
  - k-means
  - Expectation maximization
  - Hierarchical clustering
- There exist many techniques for projecting n-dimensions down to 2-D (dimensionality reduction, too many variables)
  - Multi-dimensional scaling (MDS)
  - Principal component analysis (PCA)
  - Linear discriminant analysis
  - Factor analysis

# K-means Clustering

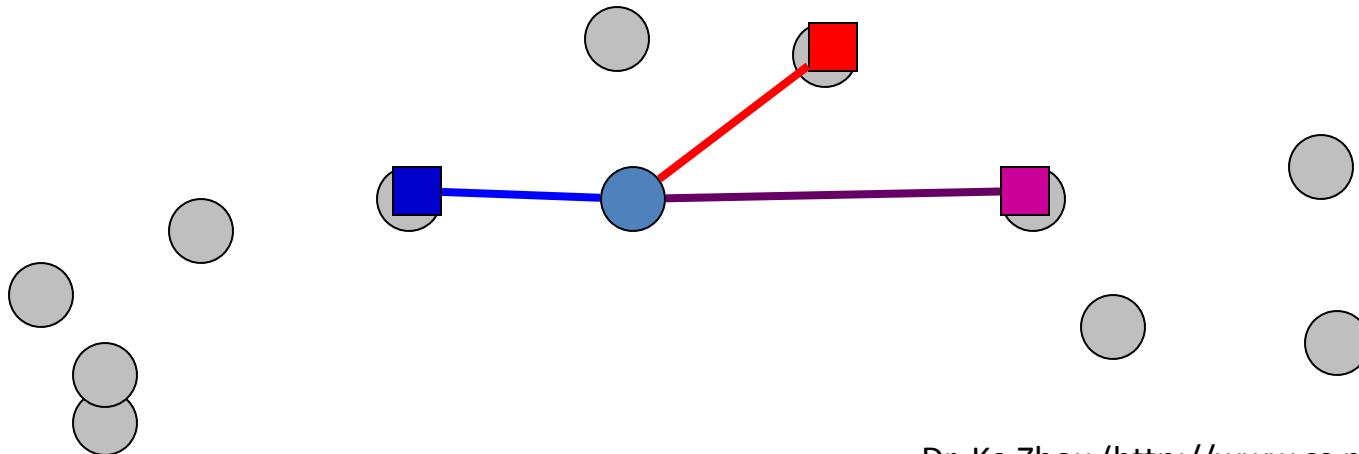


- The most well-known and popular clustering algorithm

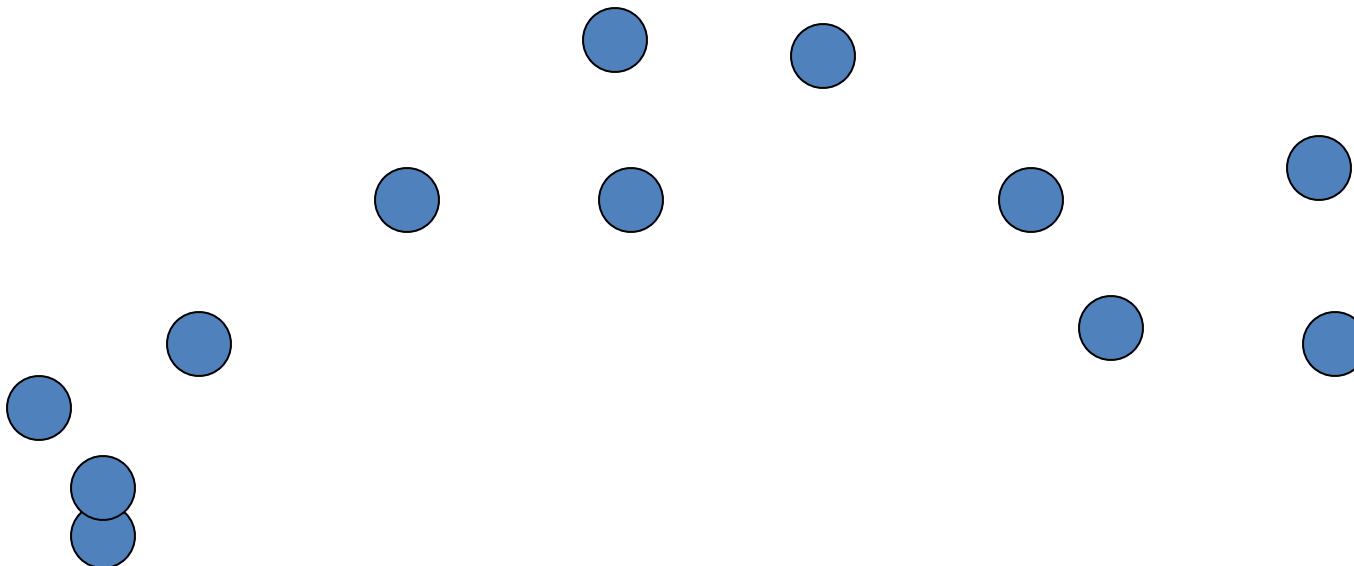
# K-means Algorithm

Iterate:

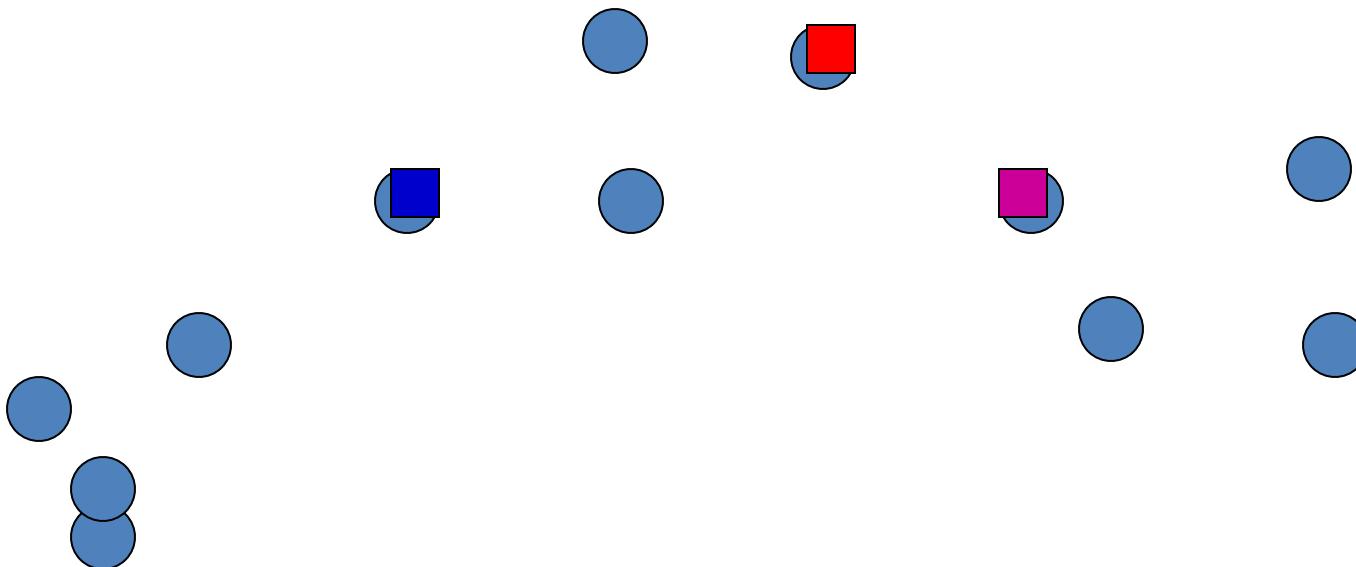
- Start with some initial cluster centers
- Assign/cluster each example to closest center
  - iterate over each point:
    - get distance to each cluster center
    - assign to closest center (hard cluster)
- Recalculate centers as the mean of the points in a cluster



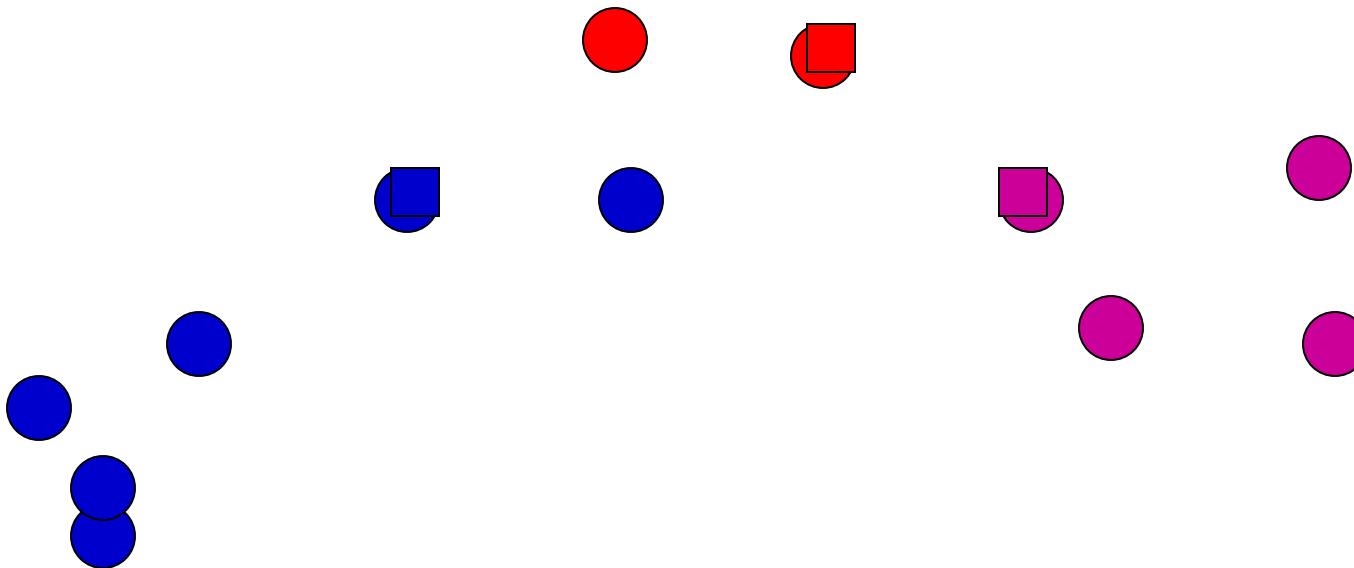
# K-means: an example



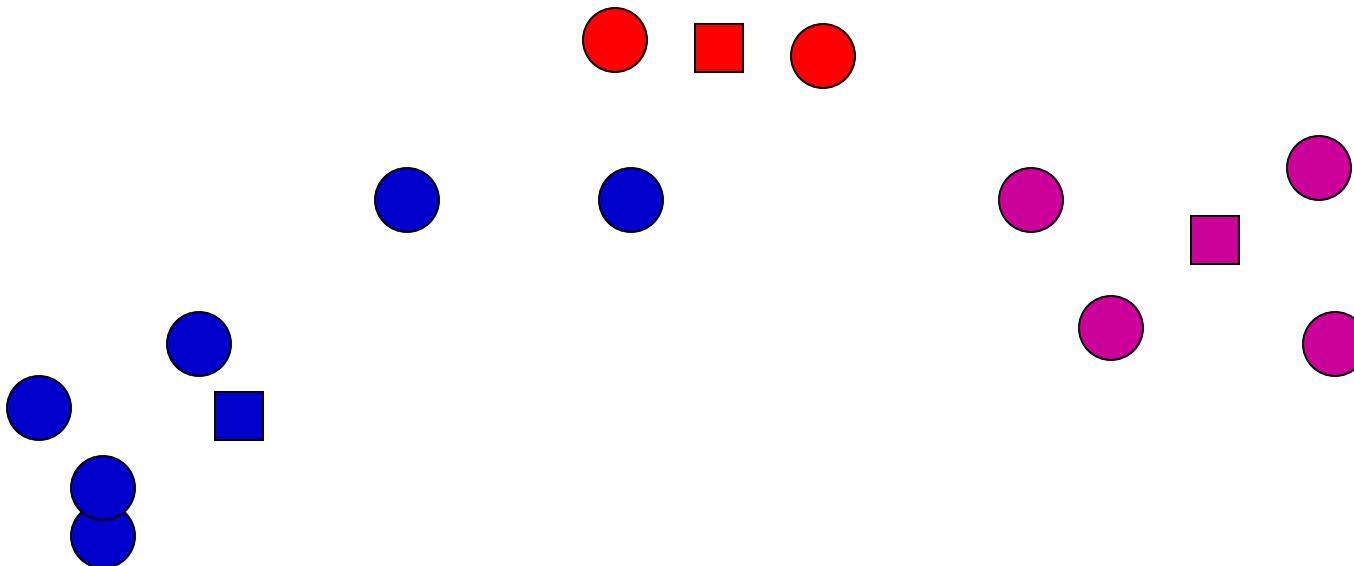
# K-means: Initialize centers randomly



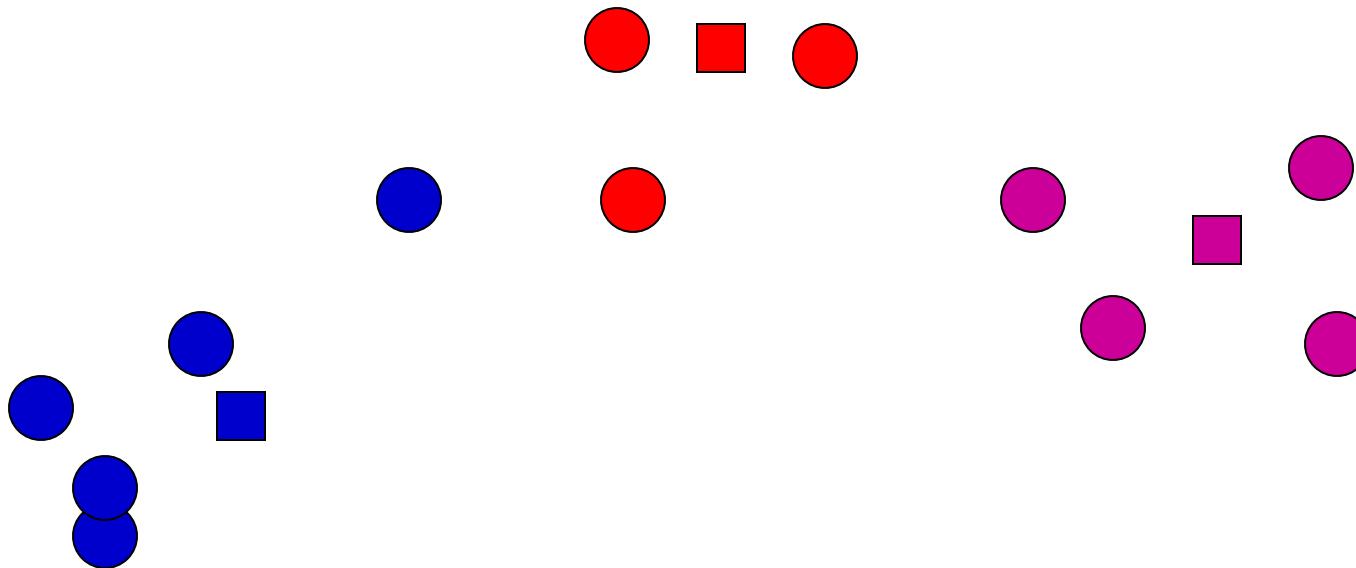
# K-means: assign points to nearest center



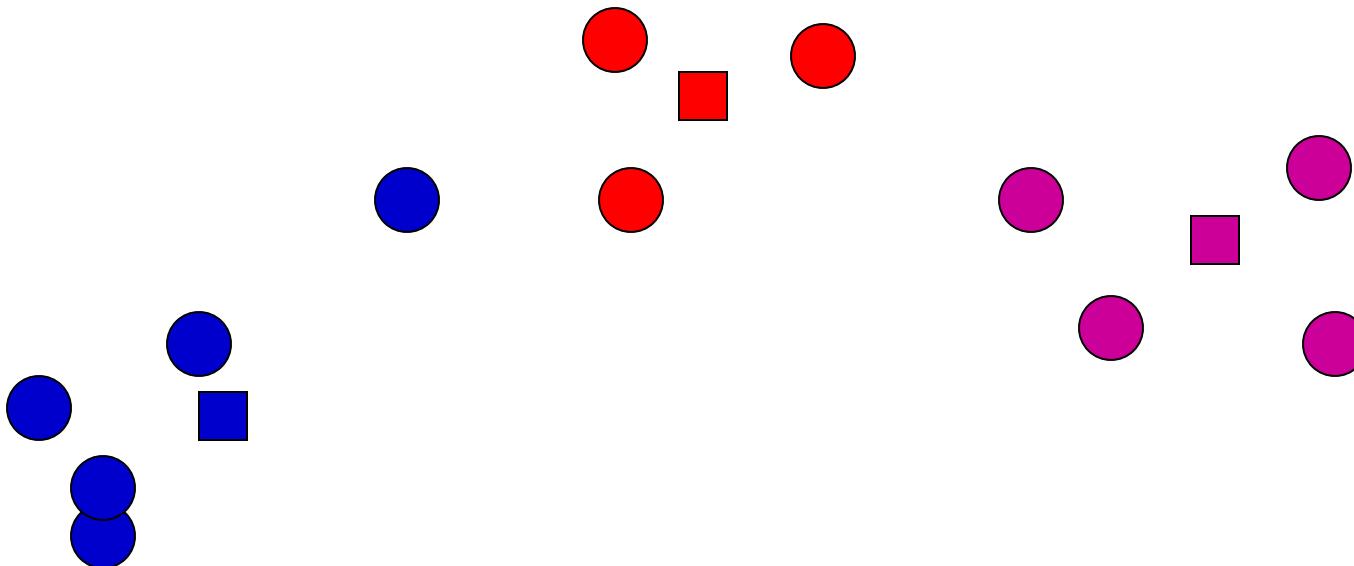
# K-means: readjust centers



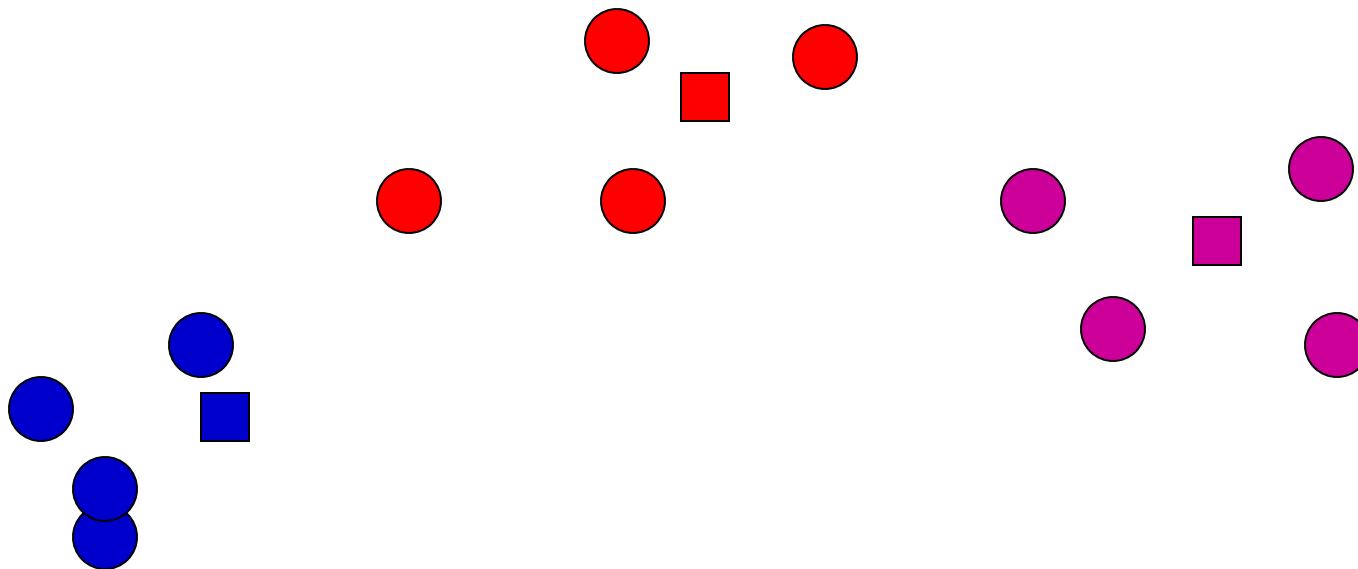
# K-means: assign points to nearest center



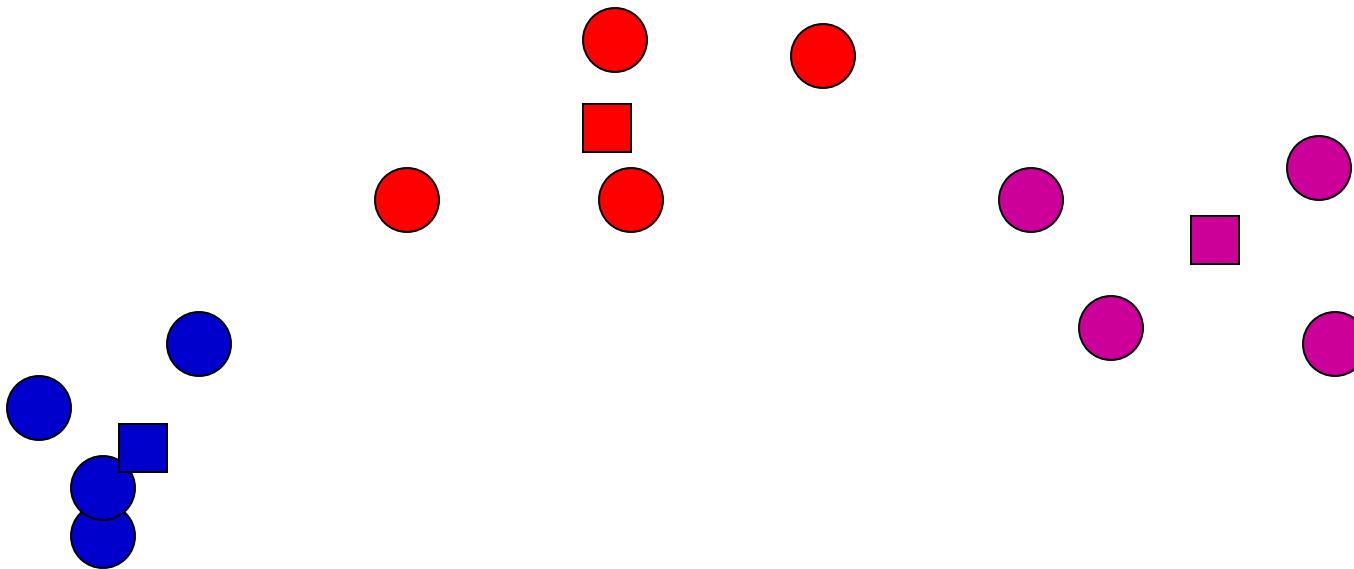
# K-means: readjust centers



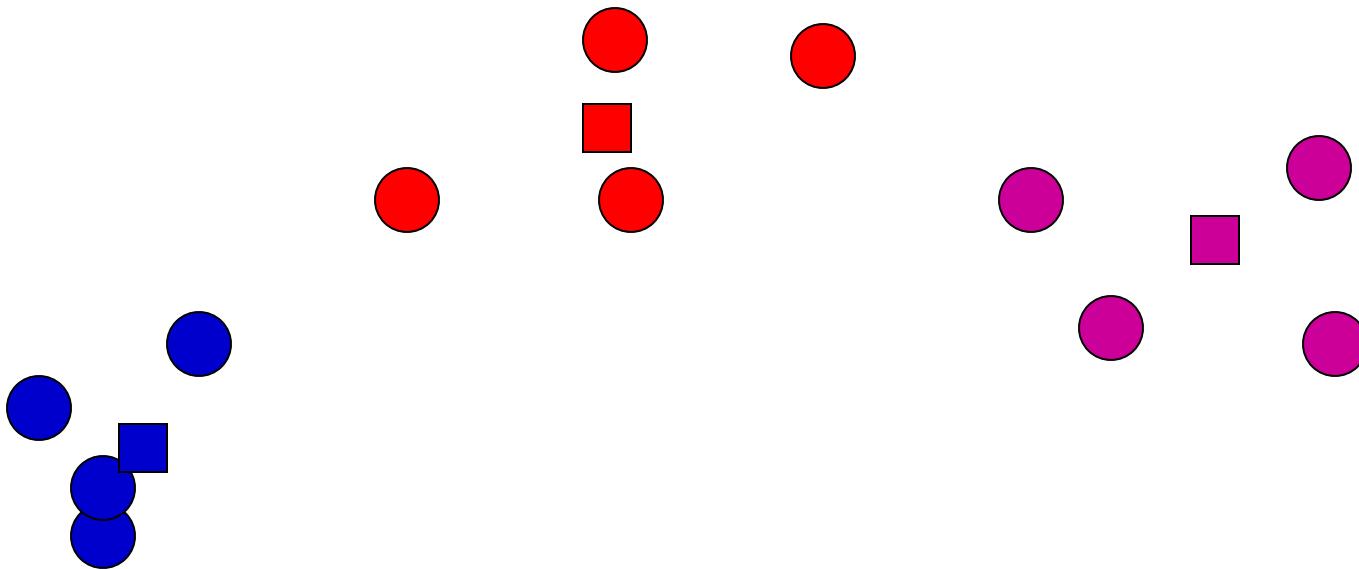
# K-means: assign points to nearest center



# K-means: readjust centers



# K-means: assign points to nearest center



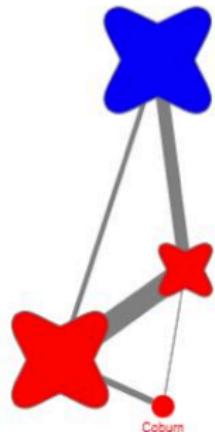
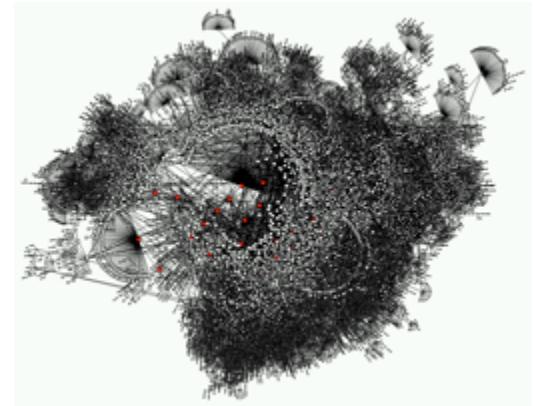
No changes: Done

# Other Reduction Techniques

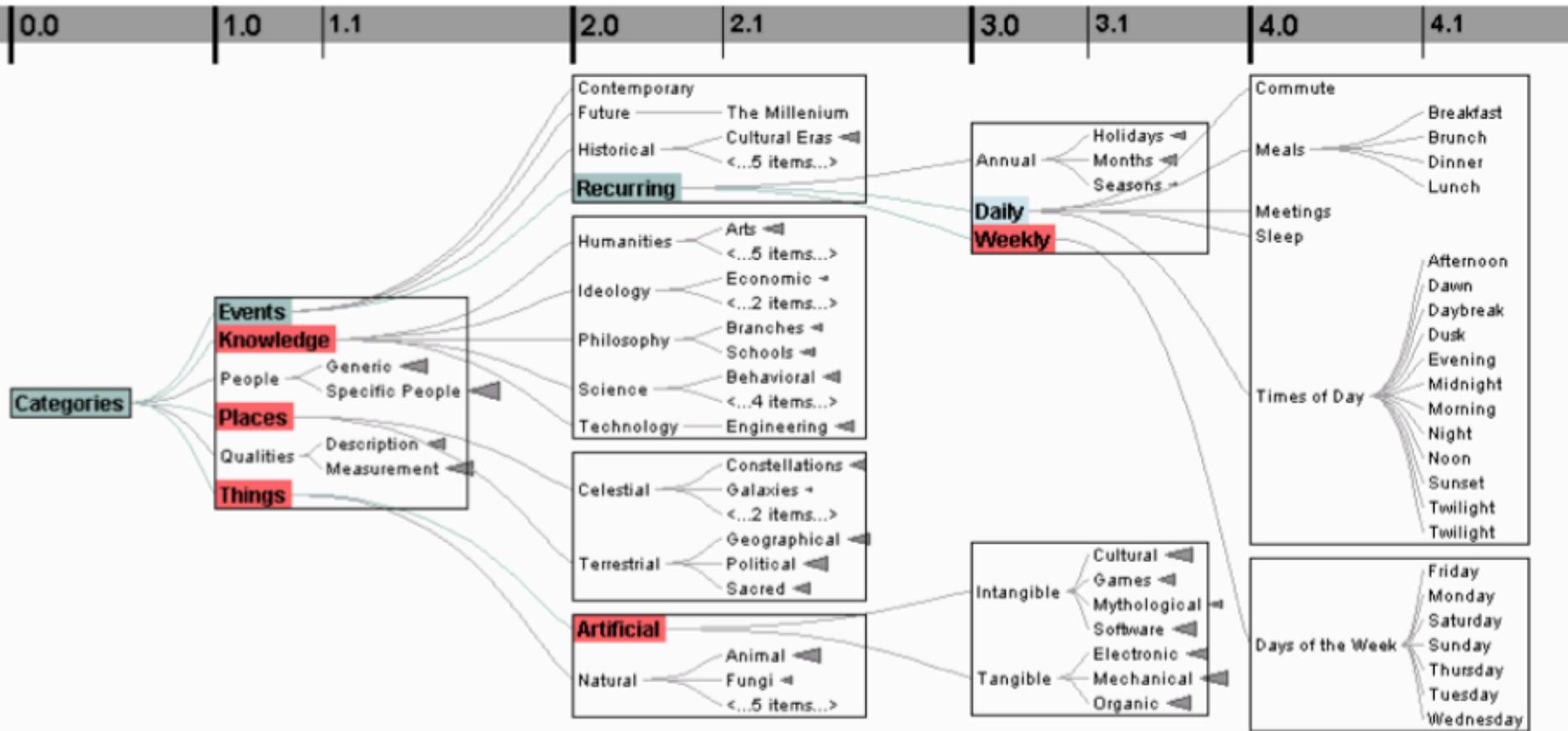
- Other techniques exist to manage scale
  - Sampling: only including every so many data cases or variables
  - Aggregation: combining many data cases or variables
- Interaction
  - Employ user interaction rather than special renderings to help manage scale

# Application: Scalability Issue of Graphs

- Need to cope with messiness
- Solutions
  - Extracting network motifs
  - Taking advantage of node attributes
  - Degree-of-Interest graphs
  - Use the alternative representation: matrix

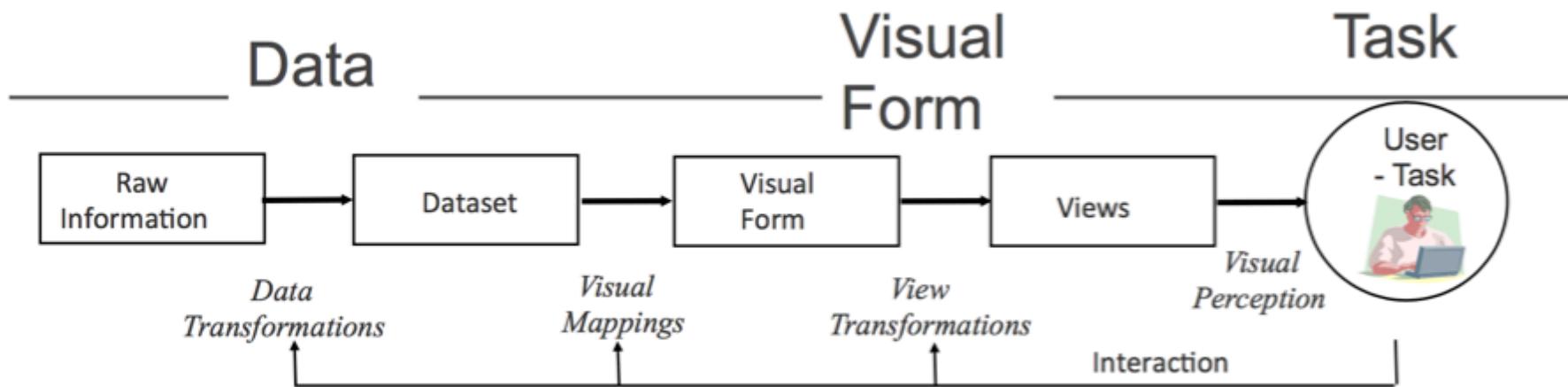


# Interactive: Degree-of-interest Trees/Graphs



- Cull “un-interesting” nodes on a per block basis until all blocks on a level fit within bounds.
- Attempt to center child blocks beneath parents.

# Information Visualization



# Next Lecture

- Topic:
  - Review

