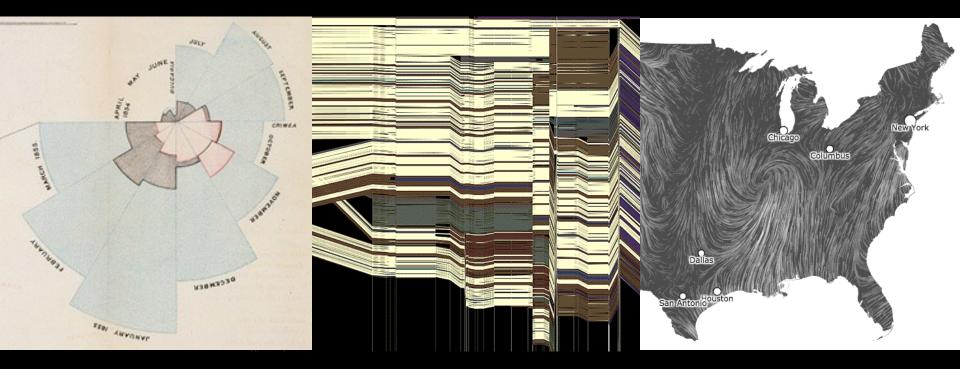
# CSE 442 - Data Visualization Data and Image Models



Jeffrey Heer University of Washington

# Last Time: Value of Visualization

### The Value of Visualization

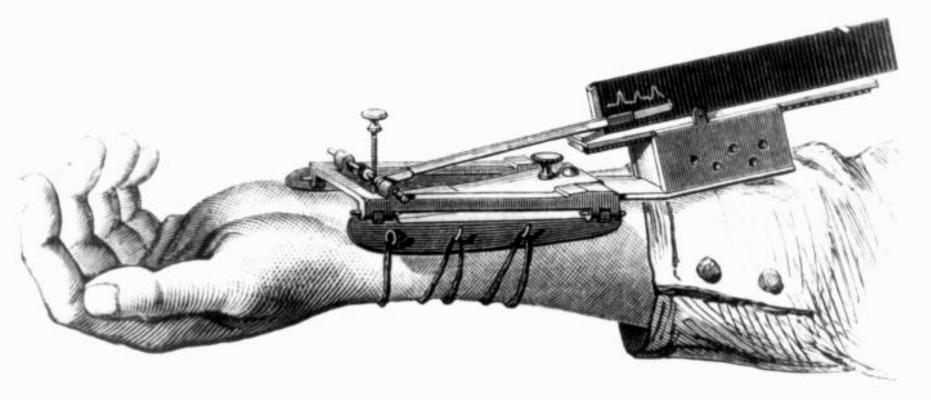
**Record** information

Blueprints, photographs, seismographs, ...

**Analyze** data to support reasoning Develop and assess hypotheses Find patterns / Discover errors in data Expand memory

**Communicate** information to others

- Share and persuade
- Collaborate and revise

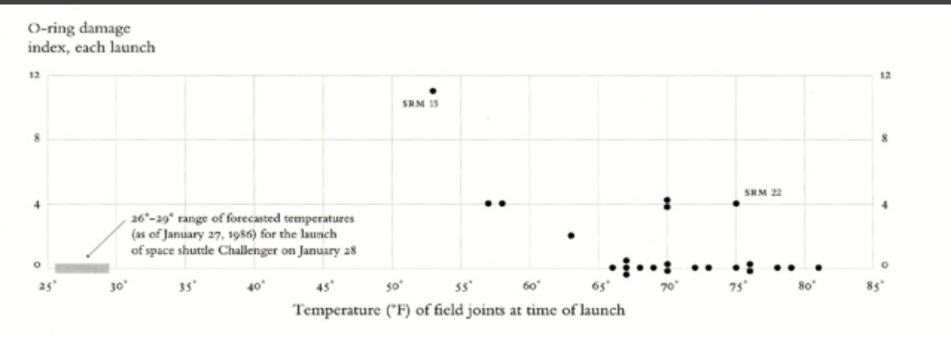


1.

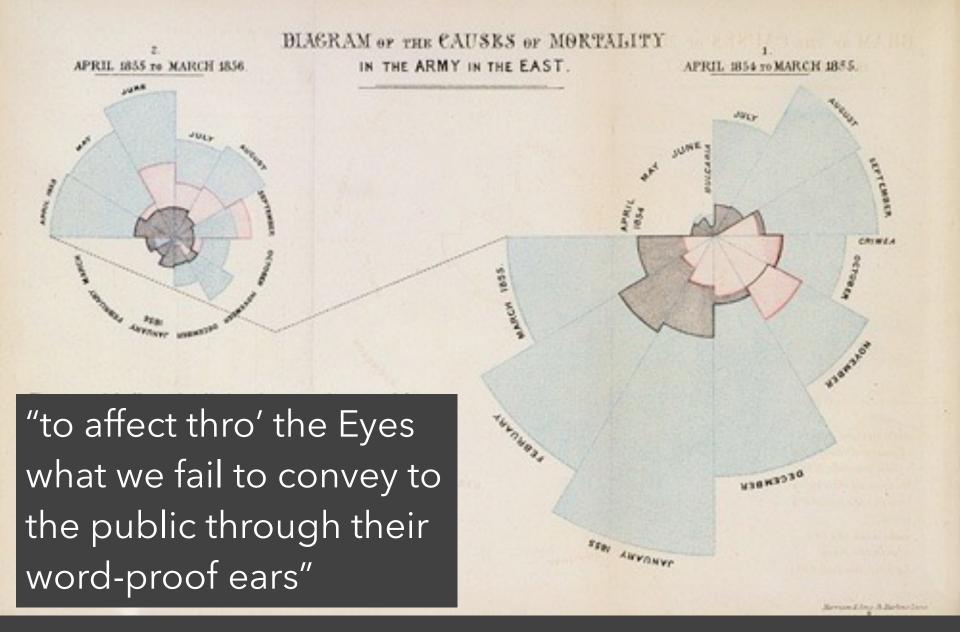
Marey's sphygmograph in use. 1860. La méthode graphique dans les sciences expérimentales et principalement en physiologie et en médecine.

E.J. Marey's sphygmograph [from Braun 83]

### Make a Decision: Challenger

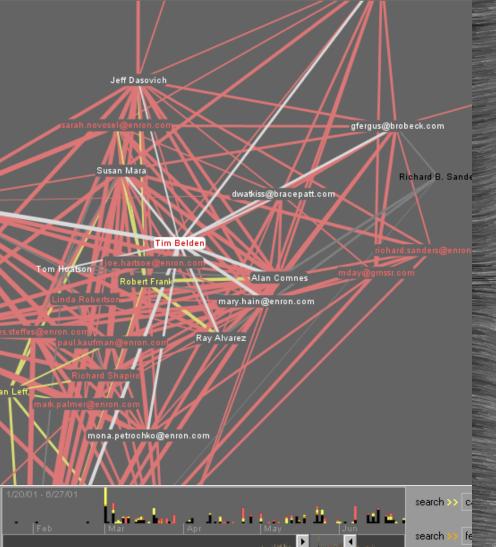


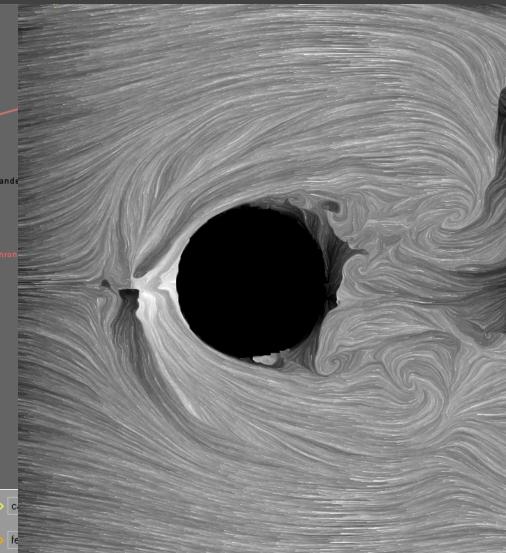
Visualizations drawn by Tufte show how low temperatures damage O-rings [Tufte 97]



1856 "Coxcomb" of Crimean War Deaths, Florence Nightingale

### InfoVis vs. SciVis?





# Data & Image Models

### The Big Picture

task questions, goals assumptions

data physical data type conceptual data type

domain metadata semantics conventions processing algorithms image visual channel graphical marks

### Topics

Properties of Data Properties of Images Mapping Data to Images Data Models

### Data Models / Conceptual Models

**Data models** are formal descriptions Math: sets with operations on them Example: integers with + and x operators

**Conceptual models** are mental constructions Include semantics and support reasoning

Examples (data vs. conceptual)1D floats vs. temperatures3D vector of floats vs. spatial location

### **Taxonomy of Data Types** (?)

1D (sets and sequences) Temporal 2D (maps) 3D (shapes) nD (relational) Trees (hierarchies) Networks (graphs)

Are there others?

The eyes have it: A task by data type taxonomy for information visualization [Shneiderman 96]

- N Nominal (labels or categories)
  - Fruits: apples, oranges, ...

- N Nominal (labels or categories)
  - Fruits: apples, oranges, ...
- O Ordered
  - Quality of meat: Grade A, AA, AAA

- N Nominal (labels or categories)
  - Fruits: apples, oranges, ...
- O Ordered
  - Quality of meat: Grade A, AA, AAA
- Q Interval (location of zero arbitrary)
  - Dates: Jan, 19, 2006; Location: (LAT 33.98, LONG -118.45)
  - Only differences (i.e. intervals) may be compared

- N Nominal (labels or categories)
  - Fruits: apples, oranges, ...
- O Ordered
  - Quality of meat: Grade A, AA, AAA
- Q Interval (location of zero arbitrary)
  - Dates: Jan, 19, 2006; Location: (LAT 33.98, LONG -118.45)
  - Only differences (i.e. intervals) may be compared
- Q Ratio (zero fixed)
  - Physical measurement: Length, Mass, Temp, ...
  - Counts and amounts

- N Nominal (labels or categories)
  - Operations: =, ≠
- O Ordered
  - Operations: =,  $\neq$ , <, >
- Q Interval (location of zero arbitrary)
  - Operations: =, ≠, <, >, -
  - Can measure distances or spans
- Q Ratio (zero fixed)
  - Operations: =,  $\neq$ , <, >, -, %
  - Can measure ratios or proportions

### From Data Model to N, O, Q

**Data Model** 32.5, 54.0, -17.3, ... Floating point numbers

**Conceptual Model** Temperature (°C)

**Data Type** Burned vs. Not-Burned (N) Hot, Warm, Cold (O) Temperature Value (Q)

### **Dimensions & Measures**

**Dimensions** (~ independent variables) Discrete variables describing data (N, O) Categories, dates, binned quantities

**Measures** (~ dependent variables) Data values that can be aggregated (Q) Numbers to be analyzed Aggregate as sum, count, avg, std. dev...

# Example: U.S. Census Data

### **Example: U.S. Census Data**

People Count: # of people in group
Year: 1850 - 2000 (every decade)
Age: 0 - 90+
Sex: Male, Female
Marital Status: Single, Married, Divorced, ...

### Example: U.S. Census

### People Count

Year

Age

Sex

**Marital Status** 

2,348 data points

	А	В	С	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
25	1850	55	0	2	187208
26	1850	60	0	1	174976
27	1850	60	0	2	162236
28	1850	65	0	1	106827
29	1850	65	0	2	105534
30	1850	70	0	1	73677
31	1850	70	0	2	71762
32	1850	75	0	1	40834
33	1850	75	0	2	40229
34	1850	80	0	1	23449
35	1850	80	0	2	22949
36	1850	85	0	1	8186
37	1850	85	0	2	10511
38	1850	90	0	1	5259
39	1850	90	0	2	6569
40	1860	0	0	1	2120846
41	1860	0	0	2	2092162

## Census: N, O, Q?

People Count Year Age Sex Marital Status Q-Ratio Q-Interval (*O*) Q-Ratio (*O*) N

### **Census: Dimension or Measure?**

People Count Year Age Sex Marital Status Measure Dimension Depends! Dimension Dimension Data Tables & Transformations

### **Relational Data Model**

Represent data as a **table** (*relation*)

Each **row** (or *tuple*) represents a record Each record is a fixed-length tuple

Each **column** (or *field*) represents a variable Each field has a *name* and a *data type* 

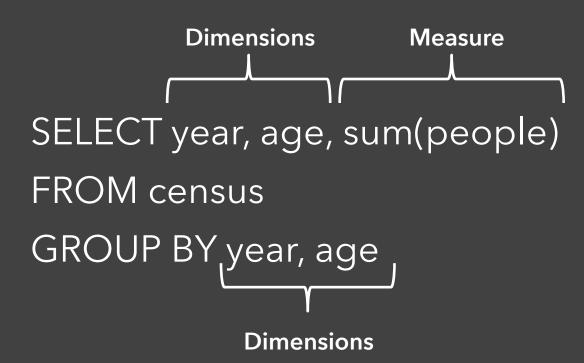
A table's **schema** is the set of names and types A **database** is a collection of tables (relations)

### Relational Algebra [Codd '70] / SQL

**Operations on Data Tables: table(s) in, table out** Projection (select): select a set of columns Selection (where): filter rows Sorting (order by): order records Aggregation (group by, sum, min, max, ...): partition rows into groups + summarize Combination (join, union, ...): integrate data from multiple tables

## **Roll-Up and Drill-Down**

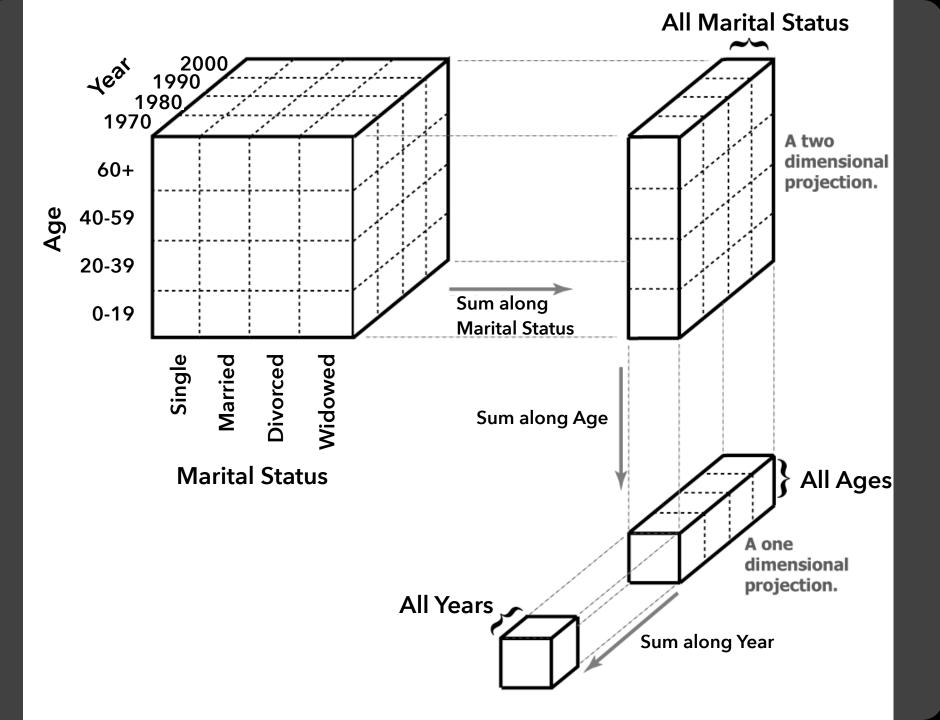
Want to examine population by year and age? **Roll-up** the data along the desired dimensions

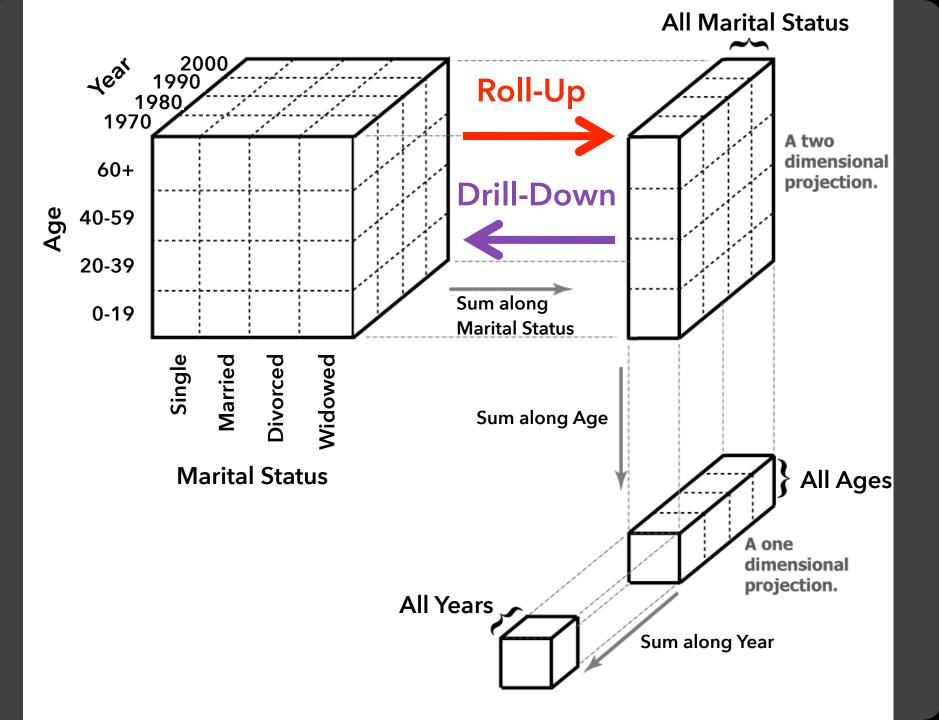


### **Roll-Up and Drill-Down**

Want to see the breakdown by marital status? **Drill-down** into additional dimensions

SELECT year, age, marst, sum(people) FROM census GROUP BY year, age, marst





YEAR	AGE	MARST	SEX	PEOPLE
1850	0	0	1	1,483,789
1850	5	0	1	1,411,067
1860	0	0	1	2,120,846
1860	5	0	1	1,804,467
• • •				

AGE MARSTSEX18501860...0011,483,7892,120,846...5011,411,0671,804,467...

Which format might we prefer?

### **Common Data Formats**

#### **CSV: Comma-Separated Values** (d3.csv)

year,age,marst,sex,people
1850,0,0,1,1483789
1850,5,0,1,1411067

• • •

### **Common Data Formats**

#### CSV: Comma-Separated Values (d3.csv)

year,age,marst,sex,people
1850,0,0,1,1483789
1850,5,0,1,1411067

• • •

#### JSON: JavaScript Object Notation (d3.json)

L {"year":1850,"age":0,"marst":0,"sex":1,"people":1483789}, {"year":1850,"age":5,"marst":0,"sex":1,"people":1411067},

### Transformations in JavaScript

#### **Operations on Data Tables: table(s) in, table out**

var array = [ 1, 2, 3, 5, 7, ... ];
// return a new filtered array
array.filter((d) => d > 2);

// sorts an array in-place and return it
array.sort((a, b) => b - a);

// return sum of values in an array
array.reduce((s, d) => s + d, 0);
d3.sum(array);
d3.sum(array, (d) => d.field);

### Aggregation in JavaScript

#### **Aggregation Functions**

d3.sum, d3.mean, d3.median, d3.deviation, ...

#### **Grouping (Nesting) Operations**

var entries = d3.nest()
.key((d) => d.variety)
.rollup((a) => d3.mean(a, (d) => d.yield))
.entries(yields);

For more, see <u>d3-array</u> and <u>d3-collection#nests</u>

# Administrivia

### **Assignment 1: Visualization Design**

#### Design a static visualization for a data set.

College admissions can play a profound role in determining one's future life and career. We've collected admissions data (grouped by gender) for selected departments at a major university.

You must choose the message you want to convey. What question(s) do you want to answer? What insight do you want to communicate?

### **Assignment 1: Visualization Design**

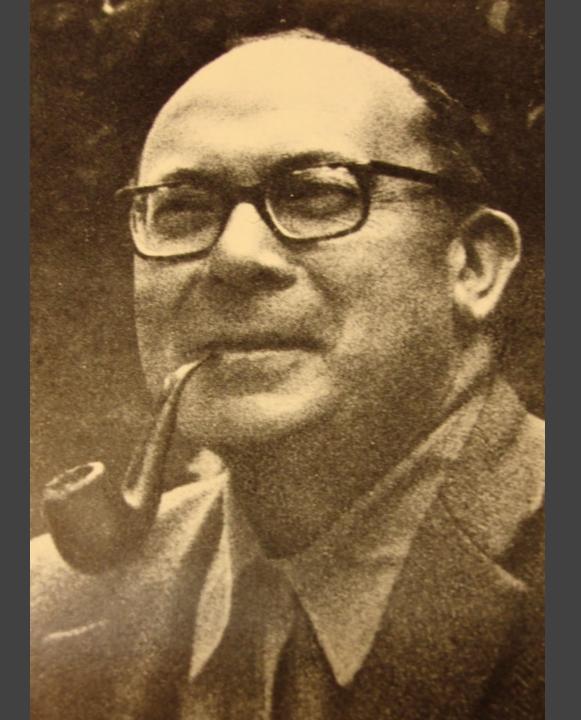
Pick a **guiding question**, use it to title your vis. Design a **static visualization** for that question. You are free to **use any tools** (inc. pen & paper).

Deliverables (upload via Canvas; see A1 page) Image of your visualization (PNG or JPG format) Short description + design rationale (≤ 4 paragraphs)

Due by 5:00 pm, Monday April 3.

### Next Tuesday: Encoding Design

We will **review A1 submissions** So be sure to turn yours in on time! Image Models



# Visual Language is a Sign System

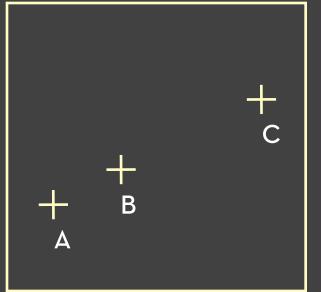


Images perceived as a set of signs Sender encodes information in signs Receiver decodes information from signs

Jacques Bertin

Sémiologie Graphique, 1967

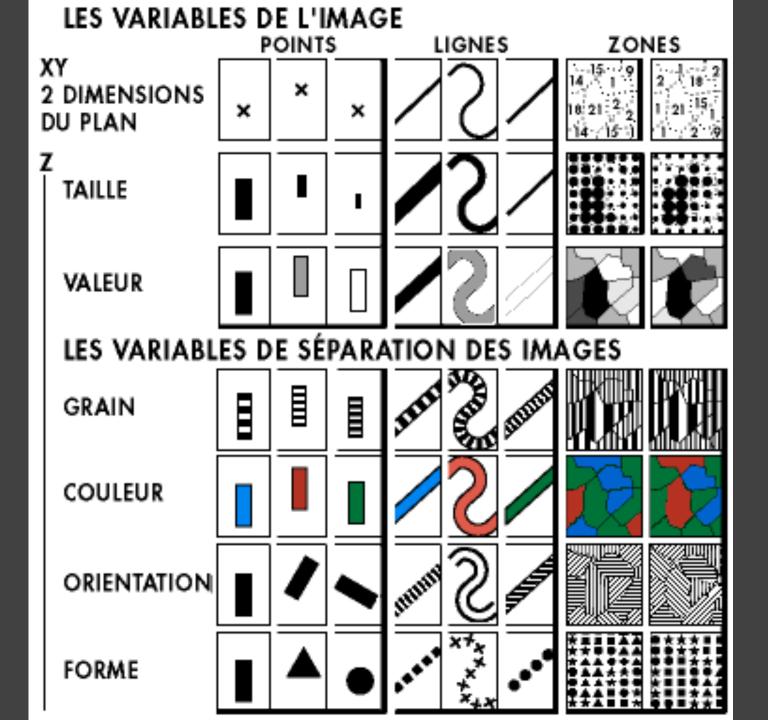
### **Bertin's Semiology of Graphics**



A, B, C are distinguishable
 B is between A and C.
 BC is twice as long as AB.

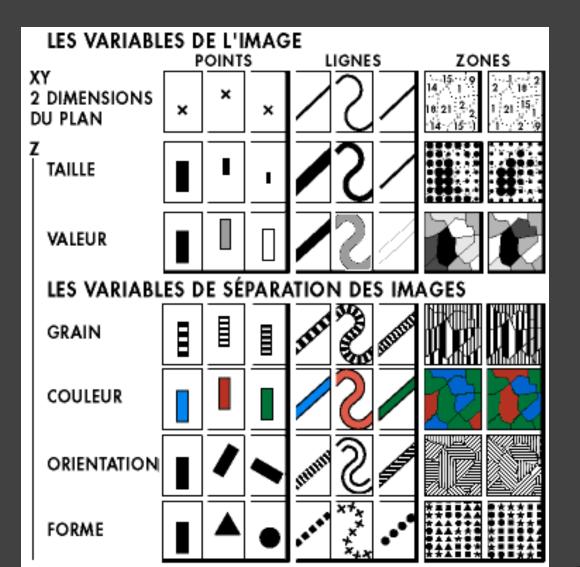
.. Encode quantitative variables

"Resemblance, order and proportional are the three signfields in graphics." - Bertin



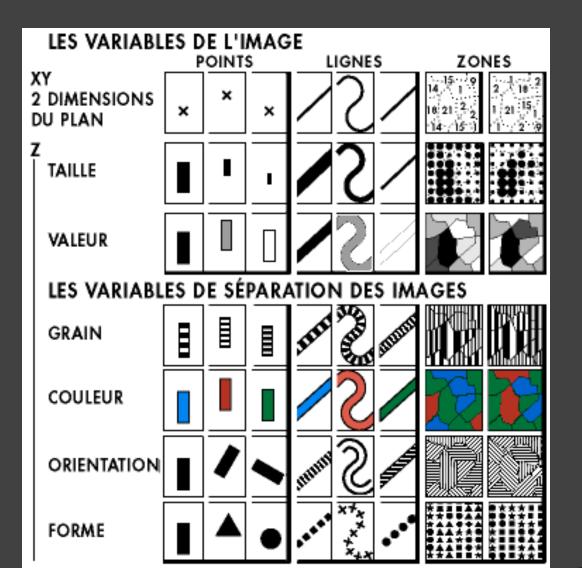
# **Visual Encoding Variables**

Position (x 2) Size Value Texture Color Orientation Shape



# **Visual Encoding Variables**

Position Length Area Volume Value Texture Color Orientation Shape Transparency Blur / Focus ...



## Information in Hue and Value

Value is perceived as ordered

: Encode ordinal variables (O)



 $\therefore$  Encode continuous variables (Q) [not as well]

Hue is normally perceived as unordered

.:. Encode nominal variables (N) using color

### Bertin's "Levels of Organization"

Q

Position

Size

Value

Texture

Color

Orientation

Shape

N	0	٥
Ν	ο	
Ν		
Ν		
Ν		

Ο

 $\mathbf{O}$ 

Ν

Ν

Nominal

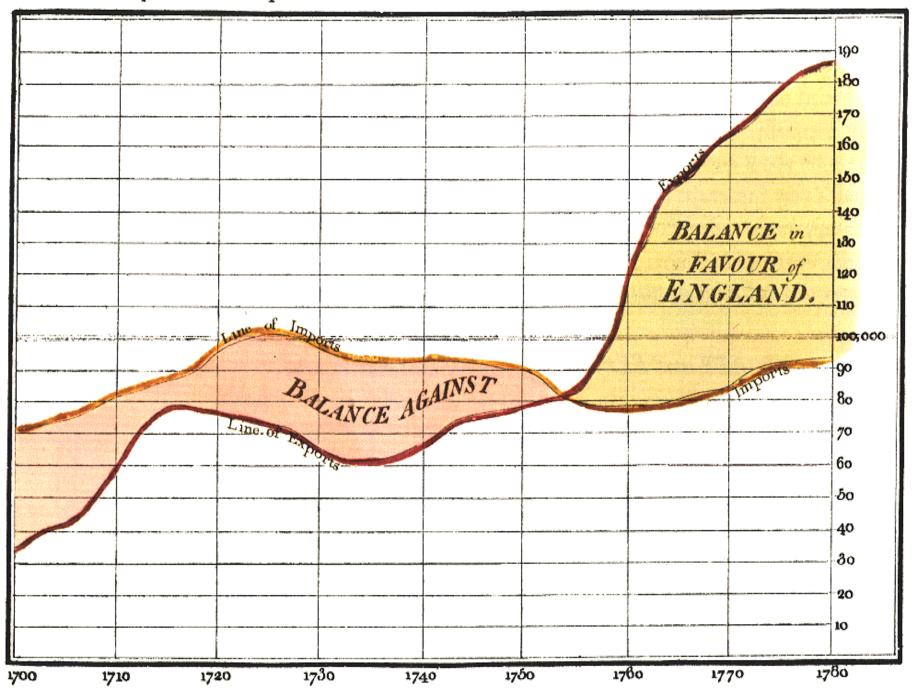
Ordinal

Quantitative

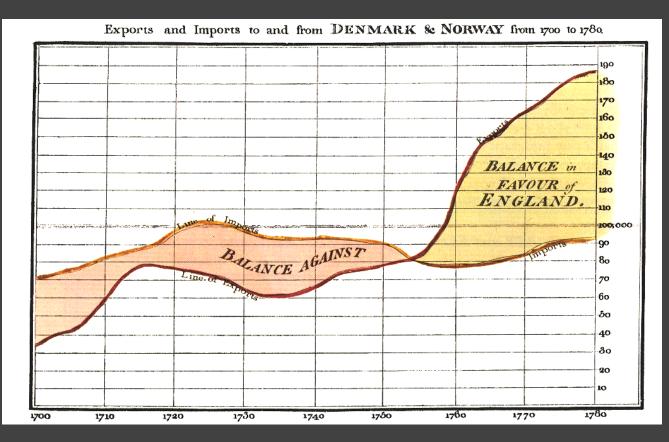
Note:  $\mathbf{Q} \subset \mathbf{O} \subset \mathbf{N}$ 

# Deconstructions

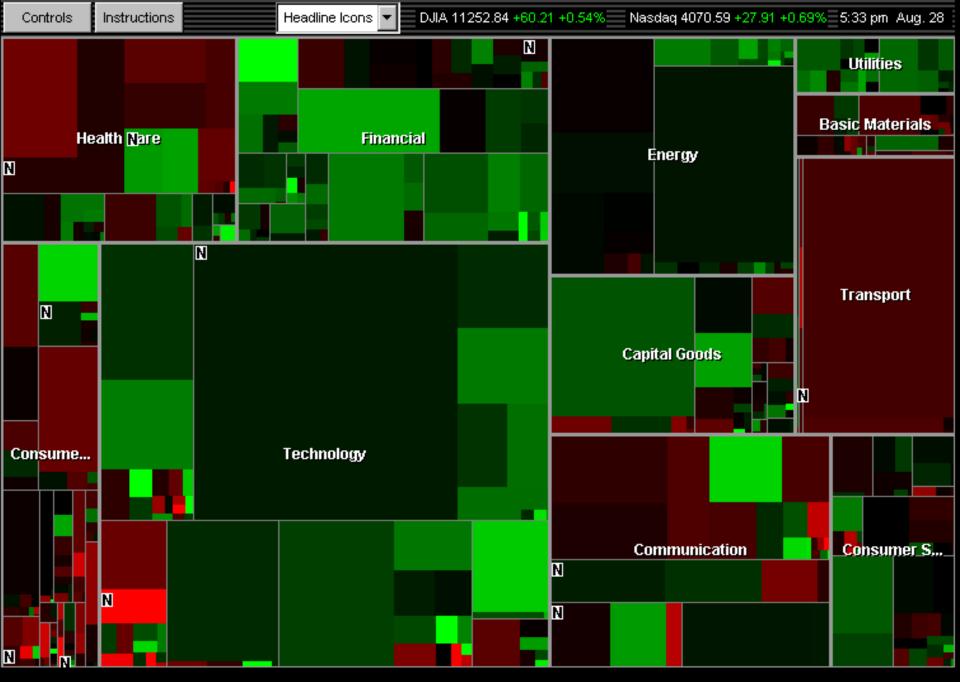
Exports and Imports to and from DENMARK & NORWAY from 1700 to 1780.



# William Playfair, 1786

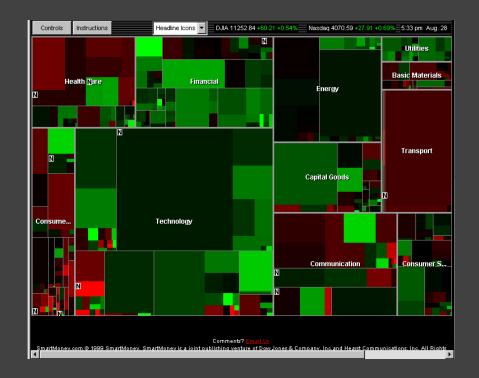


X-axis: year (Q) Y-axis: currency (Q) Color: imports/exports (N, O)



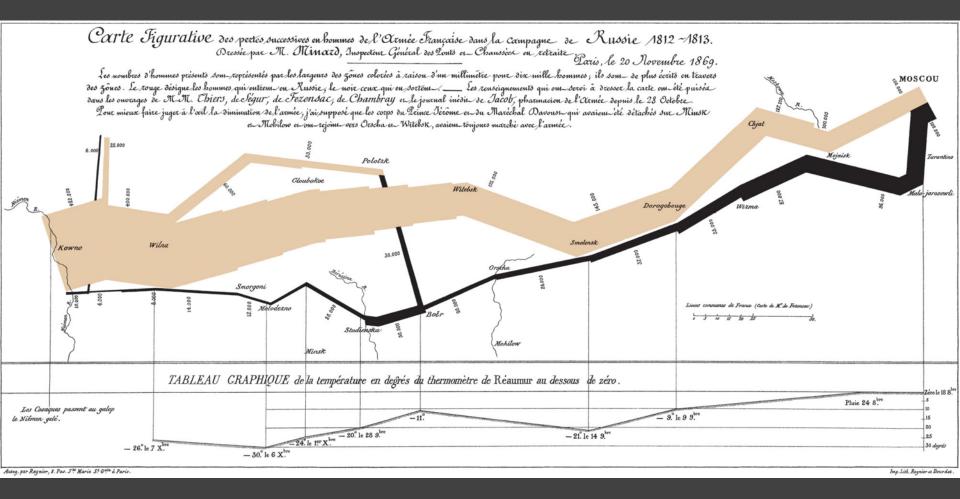
http://www.smartmoney.com/marketmap/

### Wattenberg's Map of the Market

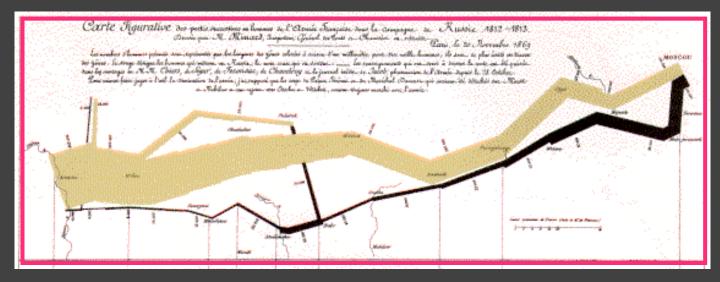


Rectangle Area: market cap (Q) Rectangle Position: market sector (N), market cap (Q) Color Hue: loss vs. gain (N, O) Color Value: magnitude of loss or gain (Q)

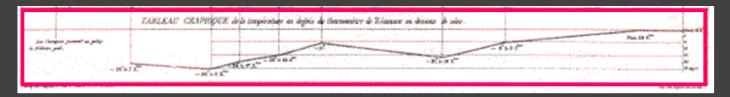
### Minard 1869: Napoleon's March



### Single-Axis Composition





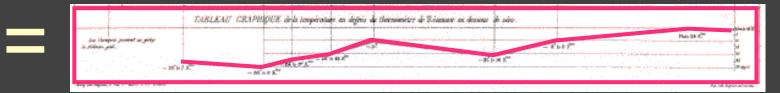




### **Mark Composition**

Y-axis: temperature (Q)

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



Temp over space/time (Q x Q)

### Mark Composition

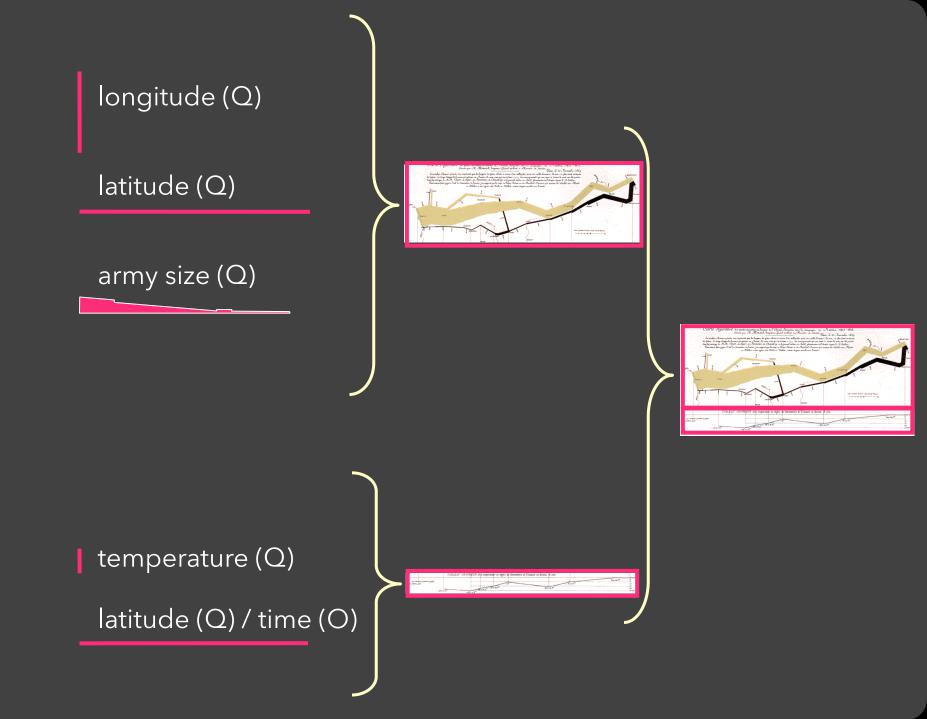
Y-axis: longitude (Q)



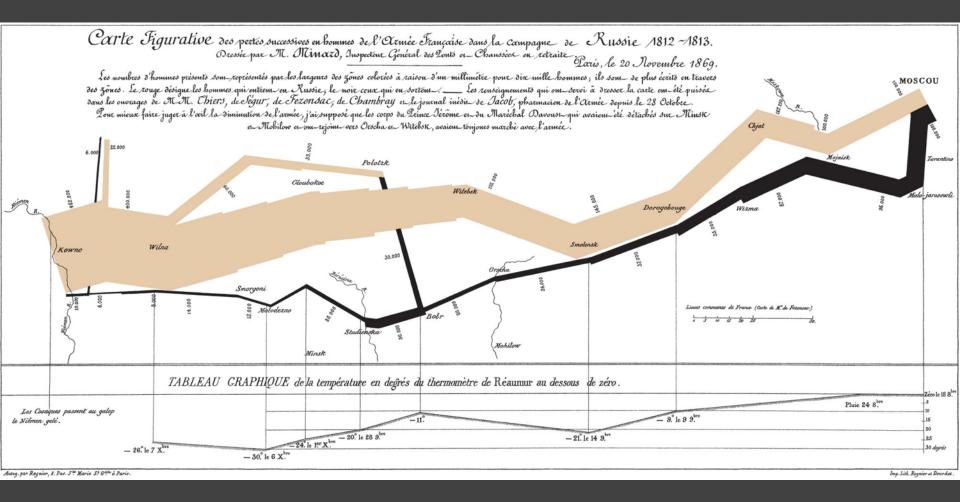




Army position  $(Q \times Q)$  and army size (Q)



### Minard 1869: Napoleon's March



#### Depicts at least 5 quantitative variables. Any others?

# Formalizing Design

# **Choosing Visual Encodings**

Assume k visual encodings and n data attributes. We would like to pick the "best" encoding among a combinatorial set of possibilities of size  $(n+1)^k$ 

#### **Principle of Consistency**

The properties of the image (visual variables) should match the properties of the data.

#### **Principle of Importance Ordering**

Encode the most important information in the most effective way.

### Design Criteria [Mackinlay 86]

#### Expressiveness

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

#### 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.

### Design Criteria [Mackinlay 86]

#### Expressiveness

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

#### 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.

### Can not express the facts

A multivariate relation may be *inexpressive* in a single horizontal dot plot because multiple records are mapped to the same position.

•		000											•••••	•••••	•• ••	•
0	 5	10	 15	 20	 25	 30	 35	 40	 45	 50	 55	 60	65	 70	 75	80
	Value															

I. Setosa	petal											
I. Decosa	sepal					00000		• ••				
I. Verginica	petal					•			• •			
i. verginica	sepal						•	******	•••• •••• •	••		
I. Versicolor	petal				• • • • • • • • • • • • • • • • • • • •		• • • • •	•				
	sepal	****										
		0 1	l .0 ;	1 20	30	 40	 50	 60	 70	80		
						Value						

### Expresses facts not in the data

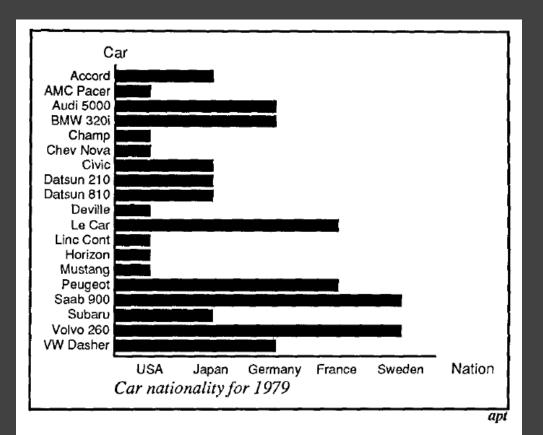


Fig. 11. Incorrect use of a bar chart for the *Nation* relation. The lengths of the bars suggest an ordering on the vertical axis, as if the USA cars were longer or better than the other cars, which is not true for the *Nation* relation.

# A length is interpreted as a quantitative value.

### Design Criteria [Mackinlay 86]

#### Expressiveness

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

#### 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.

## **Design Criteria** [Mackinlay 86]

#### Expressiveness

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

#### 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.

### Design Criteria [Tversky 02]

#### Congruence

The structure and content of the external representation should correspond to the desired structure and content of the internal representation.

#### Apprehension

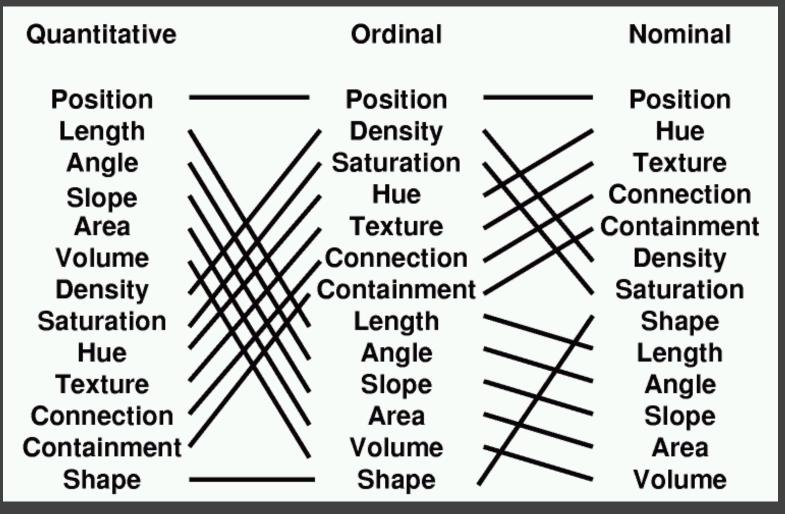
The structure and content of the external representation should be readily and accurately perceived and comprehended.

### Design Criteria Translated

**Tell the truth and nothing but the truth** (don't lie, and don't lie by omission)

Use encodings that people decode better (where better = faster and/or more accurate)

# Mackinlay's Ranking



Conjectured *effectiveness* of encodings by data type

### Mackinlay's Design Algorithm

APT - "A Presentation Tool", 1986

User formally specifies data model and type Input: ordered list of data variables to show

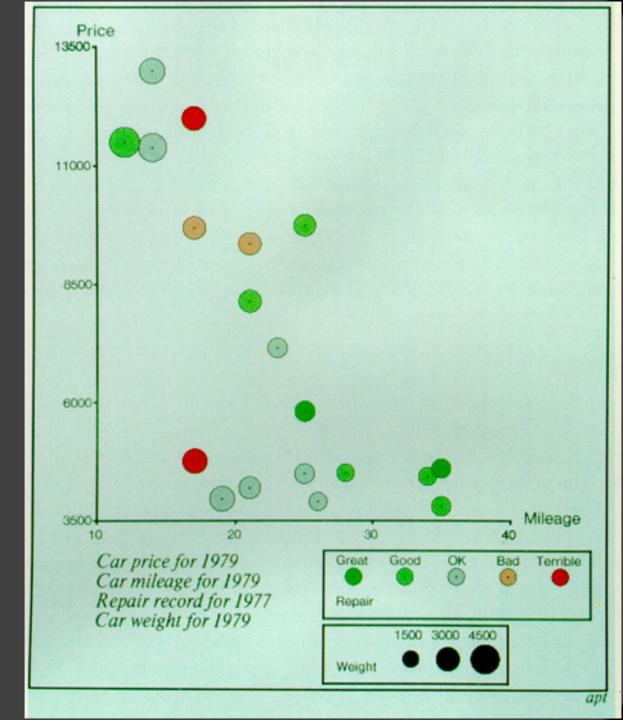
**APT searches over design space** Test expressiveness of each visual encoding Generate encodings that pass test Rank by perceptual effectiveness criteria

Output the "most effective" visualization

APT

Automatically generate chart for car data

Input variables:1. Price2. Mileage3. Repair4. Weight



### Limitations of APT?

### Limitations of APT

**Does not cover many visualization techniques** Networks, hierarchies, maps, diagrams Also: 3D structure, animation, illustration, ...

**Does not consider interaction** 

Does not consider semantics / conventions

Assumes single visualization as output

### Summary: Data & Image Models

#### **Formal specification**

Data model: relational data; N,O,Q types Image model: visual encoding channels Encodings map data to visual variables

**Choose expressive and effective encodings** Rule-based tests of expressiveness Perceptual effectiveness rankings

**Question**: how do we establish effectiveness criteria? *Subject of perception lectures*...

### **Assignment 1: Visualization Design**

Pick a **guiding question**, use it to title your vis. Design a **static visualization** for that question. You are free to **use any tools** (inc. pen & paper).

Deliverables (upload via Canvas; see A1 page) Image of your visualization (PNG or JPG format) Short description + design rationale (≤ 4 paragraphs)

Due by 5:00 pm, Monday April 3.