

STANFORD / CS448B



Data and Image Models

Jeffrey Heer

ASSISTANT: Jason Chuang

12 January 2009

<http://cs448b.stanford.edu>

Announcements

Auditors

- Requirements: Come to class and participate (online as well)

Class participation requirements

- Complete readings before class
- In-class discussion
- Post at least 1 discussion substantive comment/question on wiki within a week of each lecture

Class wiki

<https://graphics.stanford.edu/wikis/cs448b-09-winter>

Assignment 1: Visualization Design

Design a static visualization for a given data set.

Deliverables (post to the course wiki)

- Image of your visualization
- Short description and design rationale (≤ 4 paragraphs)

Due by end of day (11:59p) **today**.

Course Assistant

Jason Chuang (jcchuang@cs.stanford.edu)

Office Hours

- Thursday: 9 - 11
- Friday: 10 - 12
- Gates 398

Ph.D. student in Computer Science

Research: Understanding color perception by analyzing linguistic usage of color words.

Last Time: Value of Visualization

Three functions of visualizations

Record: store information

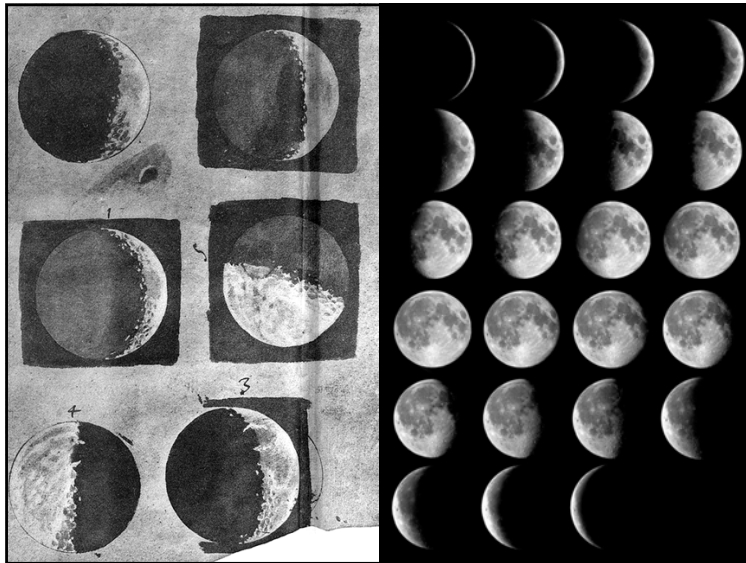
- Photographs, blueprints, ...

Analyze: support reasoning about information

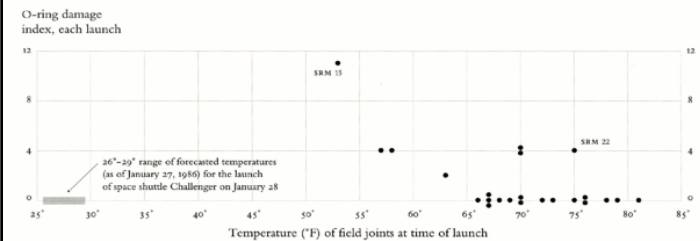
- Process and calculate
- Reason about data
- Feedback and interaction

Communicate: convey information to others

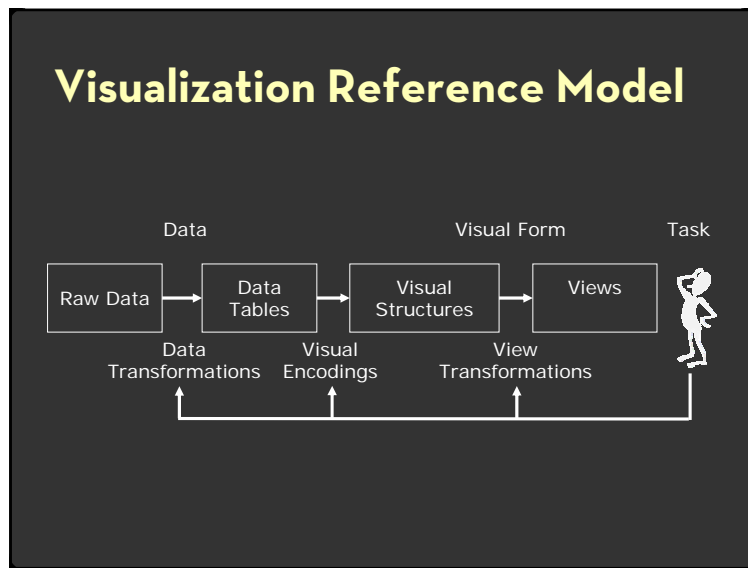
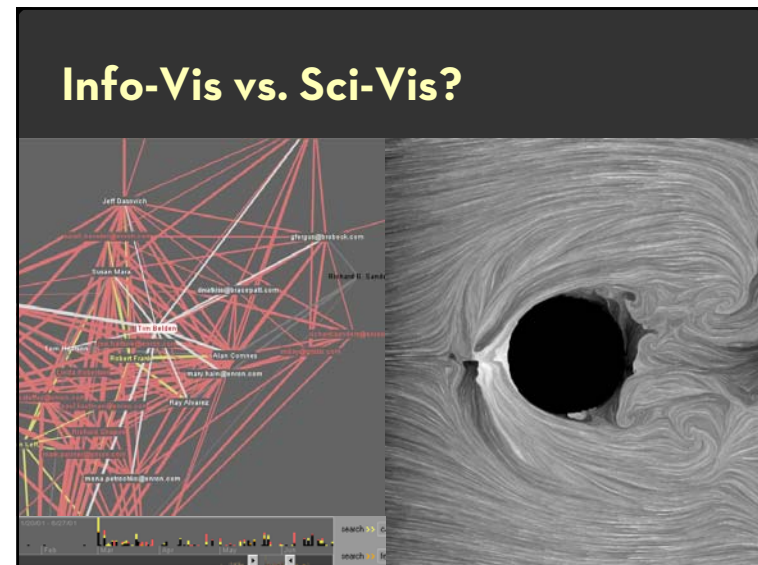
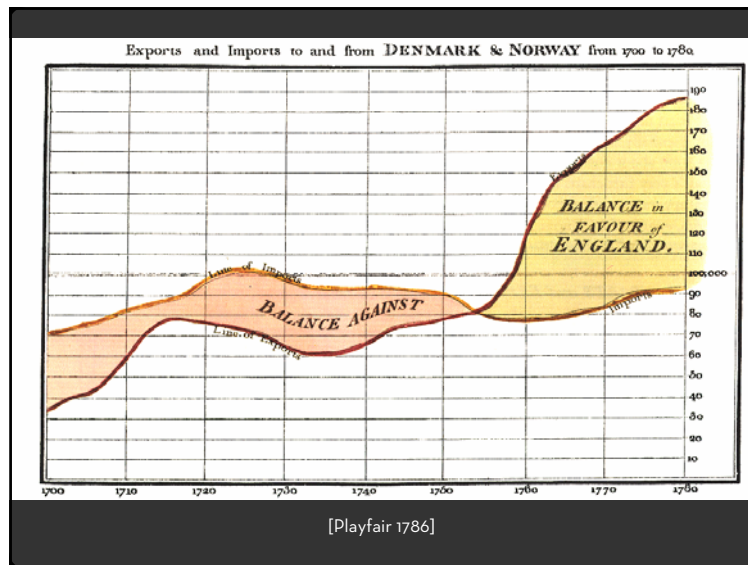
- Share and persuade
- Collaborate and revise
- Emphasize important aspects of data



Make a decision: Challenger



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



Data and Image Models

The Big Picture

task

data

physical type
int, float, etc.
abstract type
nominal, ordinal, etc.

domain

metadata
semantics
conceptual model

processing algorithms

mapping

visual encoding
visual metaphor

image

visual channel
retinal variables

[based on slide from Munzner]

Topics

- Properties of data or information
- Properties of the image
- Mapping data to images

Data

Data models vs. Conceptual models

Data models are low level descriptions of the data

- Math: Sets with operations on them
- Example: integers with + and \times operators

Conceptual models are mental constructions

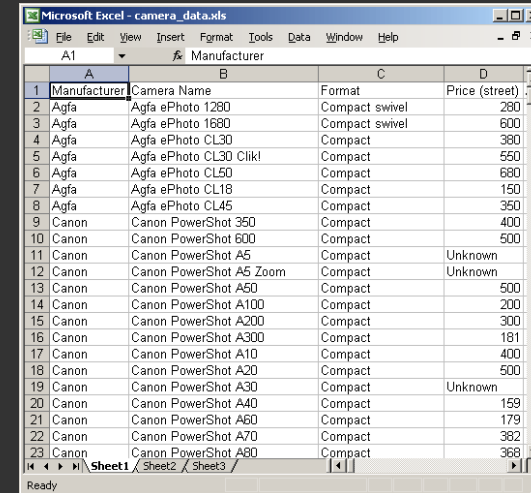
- Include semantics and support reasoning

Examples (data vs. conceptual)

- (1D floats) vs. Temperature
- (3D vector of floats) vs. Space

Relational data model

- Records are fixed-length tuples
- Each column (attribute) of tuple has a domain (type)
- Relation is schema and a table of tuples
- Database is a collection of relations



The screenshot shows a Microsoft Excel spreadsheet titled "camera_data.xls". The data is organized into a table with four columns: Manufacturer, Camera Name, Format, and Price (street). The rows list various camera models from manufacturers like Agfa and Canon, including details like lens type and price.

Manufacturer	Camera Name	Format	Price (street)
Agfa	Agfa ePhoto 1280	Compact swivel	280
Agfa	Agfa ePhoto 1680	Compact swivel	600
Agfa	Agfa ePhoto CL30	Compact	390
Agfa	Agfa ePhoto CL30 Click!	Compact	550
Agfa	Agfa ePhoto CL50	Compact	680
Agfa	Agfa ePhoto CL18	Compact	150
Agfa	Agfa ePhoto CL45	Compact	350
Canon	Canon PowerShot 350	Compact	400
Canon	Canon PowerShot 600	Compact	500
Canon	Canon PowerShot A5	Compact	Unknown
Canon	Canon PowerShot A5 Zoom	Compact	Unknown
Canon	Canon PowerShot A50	Compact	500
Canon	Canon PowerShot A100	Compact	200
Canon	Canon PowerShot A200	Compact	300
Canon	Canon PowerShot A300	Compact	181
Canon	Canon PowerShot A10	Compact	400
Canon	Canon PowerShot A20	Compact	500
Canon	Canon PowerShot A30	Compact	Unknown
Canon	Canon PowerShot A40	Compact	159
Canon	Canon PowerShot A60	Compact	179
Canon	Canon PowerShot A70	Compact	382
Canon	Canon PowerShot A80	Compact	368

Example: Digital cameras

Relational Algebra [Codd]

- Data transformations (SQL)
- Selection (SELECT)
- Projection (WHERE)
- Sorting (ORDER BY)
- Aggregation (GROUP BY, SUM, MIN, ...)
- Set operations (UNION, ...)
- Join (INNER JOIN)

Statistical data model

- Variables or measurements
- Categories or factors or dimensions
- Observations or cases

Statistical data model

- Variables or measurements
- Categories or factors or dimensions
- Observations or cases

Month	Control	Placebo	300 mg	450 mg
March	165	163	166	168
April	162	159	161	163
May	164	158	161	153
June	162	161	158	160
July	166	158	160	148
August	163	158	157	150

Blood Pressure Study (4 treatments, 6 months)

Dimensions and Measures

Independent vs. dependent variables

- Example: $y = f(x, a)$
- Dimensions: $\text{Domain}(x) \times \text{Domain}(a)$
- Measures: $\text{Range}(y)$

Dimensions and Measures

Dimensions: Discrete variables describing data
Dates, categories of values (independent vars)

Measures: Data values that can be aggregated
Numbers to be analyzed (dependent vars)
Aggregate as sum, count, average, std. deviation

Example: U.S. Census Data

People: # 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

Year

Age

Sex

Marital Status

2348 data points

	A	B	C	D	E
1	year	age	marst	sex	people
2	1850	0	0	1	1482148
3	1850	0	0	2	1450376
4	1850	5	0	1	1411067
5	1850	5	0	2	1339668
6	1850	10	0	1	1160099
7	1850	10	0	2	1218114
8	1850	15	0	1	1077113
9	1850	15	0	2	1110619
10	1850	20	0	1	1017381
11	1850	20	0	2	1003841
12	1850	25	0	1	892577
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	438189
20	1850	45	0	1	384211
21	1850	45	0	2	341234
22	1850	50	0	1	321345
23	1850	50	0	2	286580
24	1850	55	0	1	194090
25	1850	55	0	2	187208
26	1850	60	0	1	174976
27	1850	60	0	2	162236
28	1850	65	0	1	108827
29	1850	65	0	2	105834
30	1850	70	0	1	73677
31	1850	70	0	2	71762
32	1850	75	0	1	40834
33	1850	75	0	2	40219
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	5239
39	1850	90	0	2	6569
40	1860	0	0	1	2120846
41	1860	0	0	2	2021103

Census: Dimension or Measure?

People Count

Measure

Year

Dimension

Age

Depends!

Sex (M/F)

Dimension

Marital Status

Dimension

Roll-Up and Drill-Down

Want to examine marital status in each decade?

Roll-up the data along the desired dimensions

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

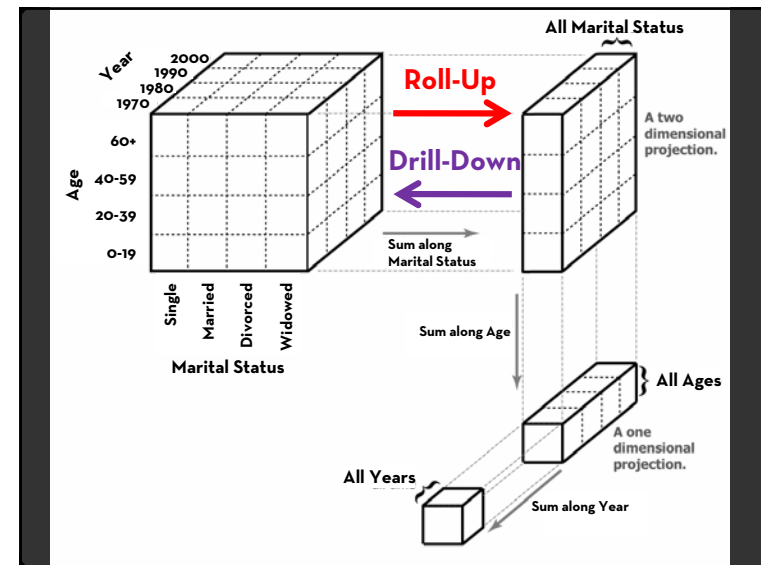
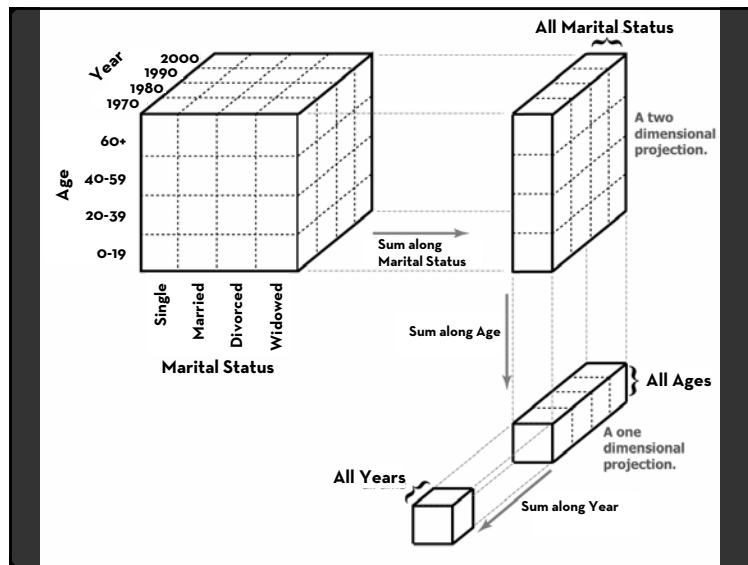
Dimensions: year, marst
Measure: sum(people)

Roll-Up and Drill-Down

Need more detailed information?

Drill-down into additional dimensions

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



Taxonomy

- 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]

Types of variables

Physical types

- Characterized by storage format
- Characterized by machine operations
- Example: bool, short, int32, float, double, string, ...

Abstract types

- Provide descriptions of the data
- May be characterized by methods/attributes
- May be organized into a hierarchy
- Example: plants, animals, metazoans, ...

Nominal, Ordinal and Quantitative

N - Nominal (labels)

- 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)
- Like a geometric point. Cannot compare directly
- Only differences (i.e. intervals) may be compared

Q - Ratio (zero fixed)

- Physical measurement: Length, Mass, Temp, ...
- Counts and amounts
- Like a geometric vector, origin is meaningful

S. S. Stevens, On the theory of scales of measurements, 1946

Nominal, Ordinal and Quantitative

N - Nominal (labels)

- Operations: =, ≠

O - Ordered

- Operations: =, ≠, <, >

Q - Interval (Location of zero arbitrary)

- Operations: =, ≠, <, >, -
- Can measure distances or spans

Q - Ratio (zero fixed)

- Operations: =, ≠, <, >, ·, ÷
- Can measure ratios or proportions

S. S. Stevens, On the theory of scales of measurements, 1946

From data model to N,O,Q data type

Data model

- 32.5, 54.0, -17.3, ...
- floats

Conceptual model

- Temperature (°C)

Data type

- Burned vs. Not burned (N)
- Hot, warm, cold (O)
- Continuous range of values (Q)

[based on slide from Munzner]

ID	Case	Species_No	Species	Organ	Width	Length
2	1	1	1 I. Setosa	Petal	2	14
3	2	1	3 I. Versicolour	Petal	24	56
4	3	1	2 I. Versicolour	Petal	13	45
5	4	1	1 I. Setosa	Sepal	33	50
6	5	1	3 I. Versicolour	Sepal	31	67
7	6	1	2 I. Versicolour	Sepal	26	57
8	7	2	1 I. Setosa	Petal	2	10
9	8	2	3 I. Versicolour	Petal	23	51
10	9	2	2 I. Versicolour	Petal	16	47
11	10	2	1 I. Setosa	Sepal	36	46
12	11	2	3 I. Versicolour	Sepal	31	69
13	12	2	2 I. Versicolour	Sepal	33	63
14	13	3	1 I. Setosa	Petal	2	16
15	14	3	3 I. Versicolour	Petal	20	52
16	15	3	2 I. Versicolour	Petal	14	47
17	16	3	1 I. Setosa	Sepal	31	48
18	17	3	3 I. Versicolour	Sepal	30	65
19	18	3	2 I. Versicolour	Sepal	32	70
20	19	4	1 I. Setosa	Petal	1	14
21	20	4	3 I. Versicolour	Petal	19	51
22	21	4	2 I. Versicolour	Petal	12	40
23	22	4	1 I. Setosa	Sepal	36	49
24	23	4	3 I. Versicolour	Sepal	27	50
25	24	4	2 I. Versicolour	Sepal	26	58
26	25	5	1 I. Setosa	Petal	2	13
27	26	5	3 I. Versicolour	Petal	17	45
28	27	5	2 I. Versicolour	Petal	10	33
29	28	5	1 I. Setosa	Sepal	32	44
30	29	5	3 I. Versicolour	Sepal	25	49
31	30	5	2 I. Versicolour	Sepal	23	50
32	31	6	1 I. Setosa	Petal	2	16

Sepal and petal lengths and widths for three species of iris [Fisher 1936].

Microsoft Excel - fischer.iris2.colored.xls

ID	Case	Species_No	Species	Organ	Width	Length
1	1	1	1. I. Setosa	Petal	2	14
2	1	3	3. I. Versicolour	Petal	24	56
3	1	2	2. I. Versicolour	Petal	13	45
4	3	1	1. I. Setosa	Sepal	33	50
5	4	1	3. I. Versicolour	Sepal	31	67
6	5	1	2. I. Versicolour	Sepal	26	57
7	6	1	1. I. Setosa	Petal	2	10
8	7	2	3. I. Versicolour	Petal	23	51
9	8	2	2. I. Versicolour	Petal	16	47
10	9	2	1. I. Setosa	Sepal	36	46
11	10	2	3. I. Versicolour	Sepal	31	69
12	11	2	2. I. Versicolour	Sepal	33	63
13	12	2	1. I. Setosa	Petal	2	16
14	13	3	3. I. Versicolour	Petal	20	62
15	14	3	3. I. Versicolour	Petal	14	47
16	15	3	2. I. Versicolour	Petal	31	48
17	16	3	3. I. Versicolour	Sepal	30	65
18	17	3	2. I. Versicolour	Sepal	32	70
19	18	3	1. I. Setosa	Petal	1	14
20	19	4	3. I. Versicolour	Petal	19	51
21	20	4	2. I. Versicolour	Petal	12	40
22	21	4	1. I. Setosa	Sepal	36	49
23	22	4	3. I. Versicolour	Sepal	27	59
24	23	4	2. I. Versicolour	Sepal	26	59
25	24	4	1. I. Setosa	Petal	2	13
26	25	5	3. I. Versicolour	Petal	17	45
27	26	5	2. I. Versicolour	Petal	10	33
28	27	5	1. I. Setosa	Sepal	32	44
29	28	5	3. I. Versicolour	Sepal	25	49
30	29	5	2. I. Versicolour	Sepal	23	60
31	30	5	1. I. Setosa	Petal	2	16
32	31	5	1. I. Setosa	Petal	2	10

Ready

Z O P

Image

Visual language is a sign system



Jacques Bertin

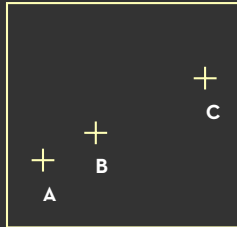
Images perceived as a set of signs
 Sender encodes information in signs
 Receiver decodes information from signs

Sémiologie Graphique, 1967

LES VARIABLES DE L'IMAGE			
	POINTS	LIGNES	ZONES
XY 2 DIMENSIONS DU PLAN	x x x	/ / /	16 15 9 21 21 9 14 15 9
Z TAILLE	■ ■ ■	/ / /	■ ■ ■
VALEUR	■ ■ ■	/ / /	■ ■ ■
LES VARIABLES DE SÉPARATION DES IMAGES			
GRAIN	■ ■ ■	/ / /	■ ■ ■
COULEUR	■ ■ ■	/ / /	■ ■ ■
ORIENTATION	■ ■ ■	/ / /	■ ■ ■
FORME	■ ■ ■	/ / /	■ ■ ■

[Bertin, Sémiologie Graphique, 1967]

Information in position



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

∴ Encode quantitative variables (Q)

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

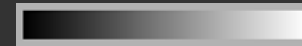
Information in color and value

Value is perceived as ordered

∴ Encode ordinal variables (O)



∴ Encode continuous variables (Q) [not as well]



Hue is normally perceived as unordered

∴ Encode nominal variables (N) using color



Visual encoding variables

Position
Size
Value
Texture
Color
Orientation
Shape

LES VARIABLES DE L'IMAGE			
	POINTS	LIGNES	ZONES
XY 2 DIMENSIONS DU PLAN	x x x	~ ~ ~	~ ~ ~
Z TAILLE	▬ ▬ ▬	~ ~ ~	~ ~ ~
VALEUR	▬ ▬ ▬	~ ~ ~	~ ~ ~
LES VARIABLES DE SÉPARATION DES IMAGES			
GRAIN	▬ ▬ ▬	~ ~ ~	~ ~ ~
COULEUR	▬ ▬ ▬	~ ~ ~	~ ~ ~
ORIENTATION	▬ ▬ ▬	~ ~ ~	~ ~ ~
FORME	▬ ▬ ▬	~ ~ ~	~ ~ ~

- Note: Bertin does not consider 3D or time
- Note: Card and Mackinlay extend the number of vars.

Bertin's "Levels of Organization"

Position

N	O	Q
---	---	---

Nominal

Size

N	O	Q
---	---	---

Ordered

Value

N	O	Q
---	---	---

Quantitative

Note: Q < O < N

Texture

N	O	Q
---	---	---

Color

N	O	Q
---	---	---

Orientation

N	O	Q
---	---	---

Shape

N	O	Q
---	---	---

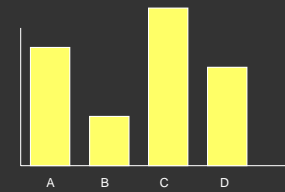
Note: Bertin actually breaks visual variables down into differentiating (≠) and associating (≠)

Encoding Rules

Univariate data

	factors		
	A	B	C
1			

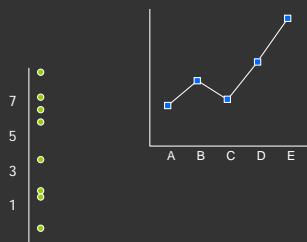
variable



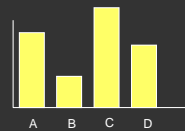
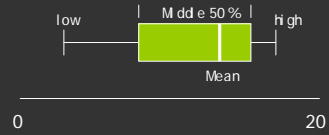
Univariate data

	factors		
	A	B	C
1			

variable



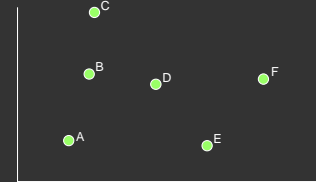
Tuke ybo x plot



[based on slide from Stasko]

Bivariate data

	factors		
	A	B	C
1			
2			

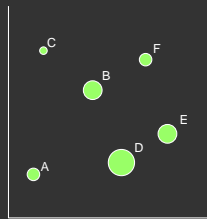


Scatter plot is common

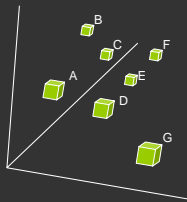
[based on slide from Stasko]

Trivariate data

	A	B	C
1			
2			
3			



3D scatter plot is possible



[based on slide from Stasko]

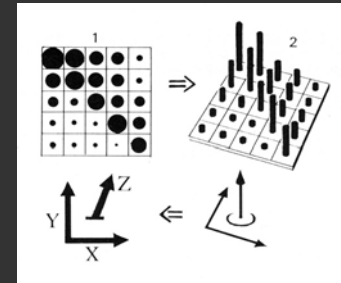
Three variables

Two variables $[x,y]$ can map to points

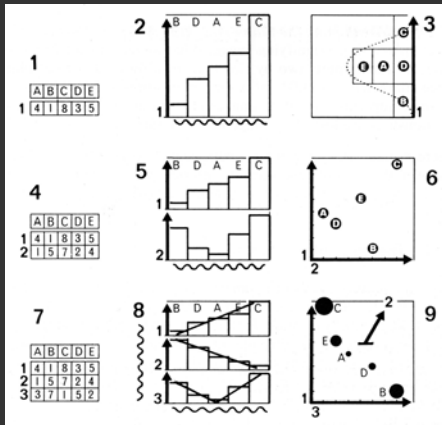
- Scatterplots, maps, ...

Third variable $[z]$ must use

- Color, size, shape, ...



Large design space (visual metaphors)



[Bertin, Graphics and Graphic Info. Processing, 1981]

Multidimensional data

How many variables can be depicted in an image?

	A	B	C
1			
2			
3			
4			
5			
6			
7			
8			

Multidimensional data

How many variables can be depicted in an image?

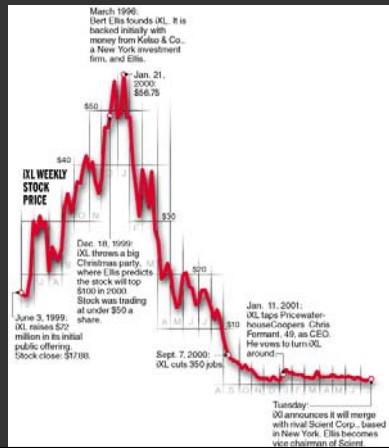
	A	B	C
1			
2			
3			
4			
5			
6			
7			
8			

"With up to three rows, a data table can be constructed directly as a single image ... However, an image has only three dimensions. And this barrier is impassible."

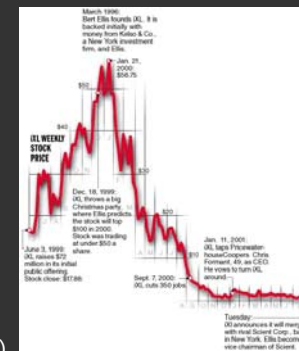
Bertin

Deconstructions

Stock chart from the late 90s

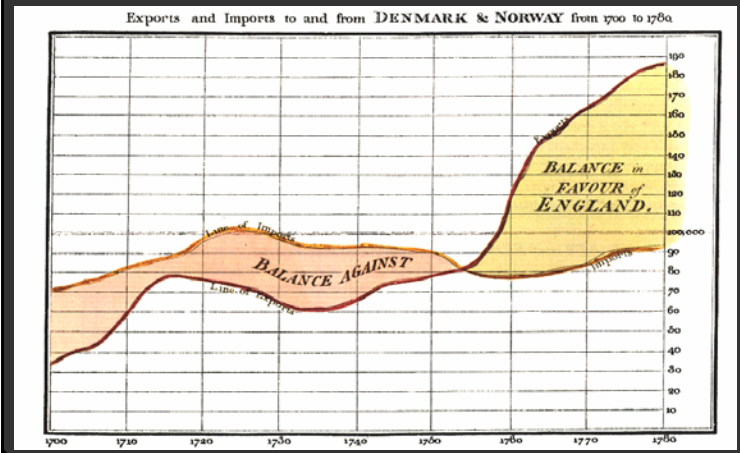


Stock chart from the late 90s

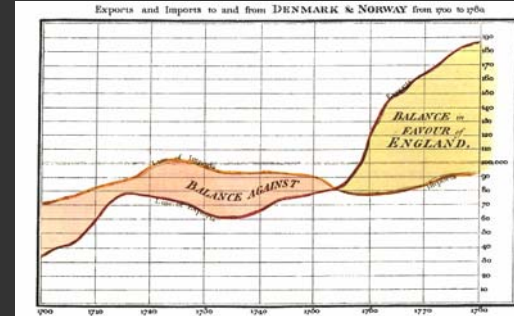


n x-axis: time (Q)
n y-axis: price (Q)

Playfair 1786

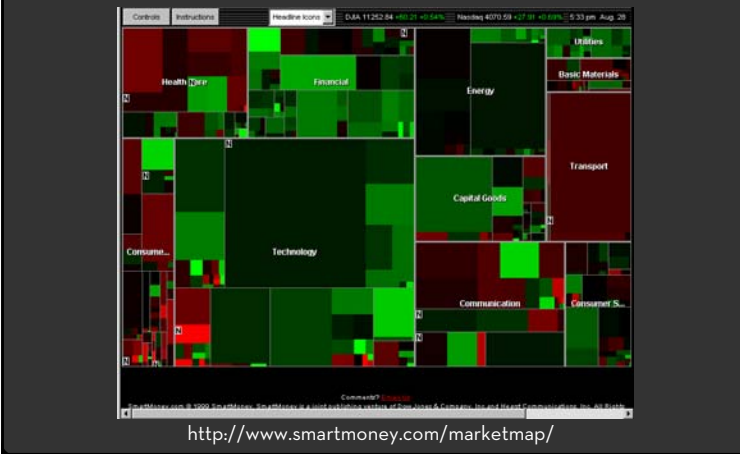


Playfair 1786

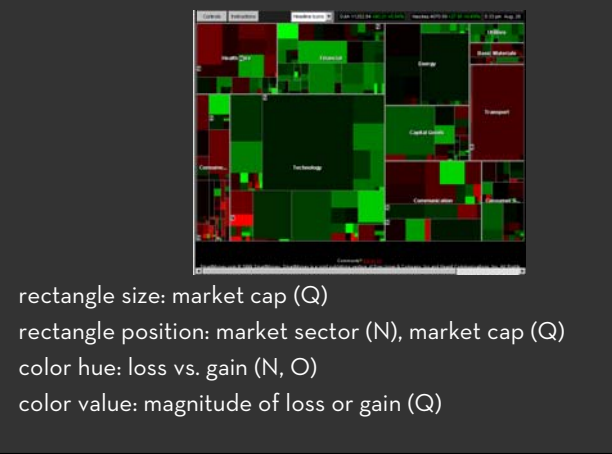


x-axis: year (Q)
 y-axis: currency (Q)
 color: imports/exports (N, O)

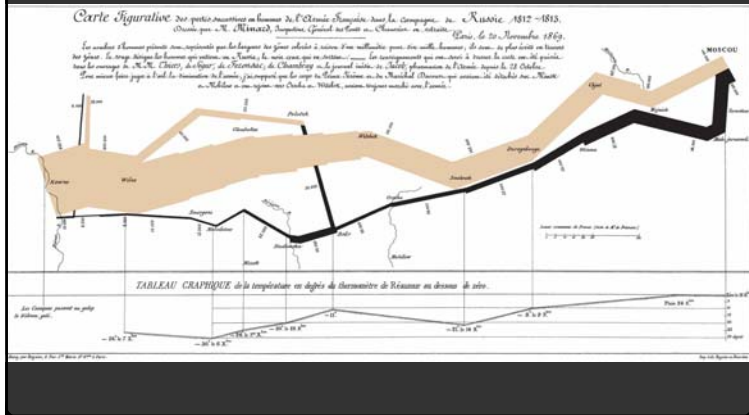
Wattenberg 1998



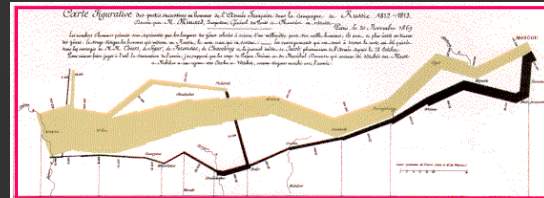
Wattenberg 1998



Minard 1869: Napoleon's march



Single axis composition



+

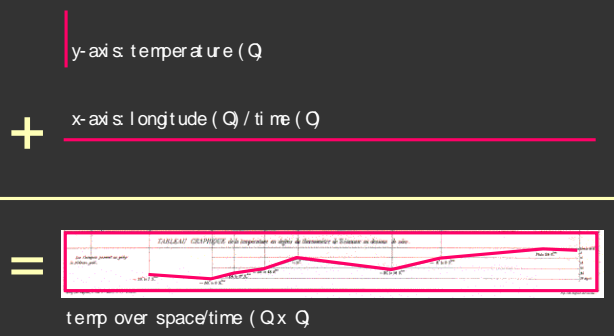


=



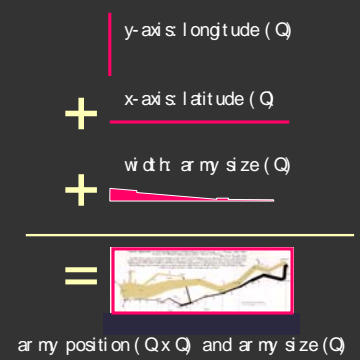
[based on slide from Mackinlay]

Mark composition



[based on slide from Mackinlay]

Mark composition



[based on slide from Mackinlay]

Longitude (Q)

Latitude (Q)

army size (Q)

temperature (Q)

Latitude (Q) / time (O)

[based on slide from Mackinlay]

Minard 1869: Napoleon's march

Depicts at least 5 quantitative variables. Any others?

Automated design

Jock Mackinlay's APT 86

Combinatorics of encodings

Challenge:
Pick the best encoding from the exponential number of possibilities $(n+1)^8$

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.

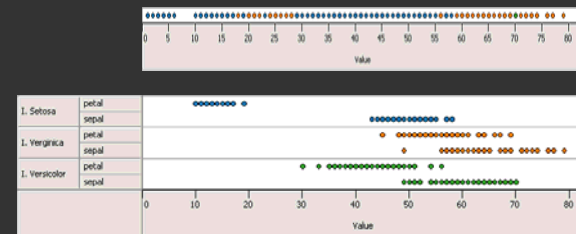
Mackinlay's expressiveness criteria

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.

Cannot express the facts

A one-to-many ($1 \rightarrow N$) relation cannot be expressed in a single horizontal dot plot because multiple tuples are mapped to the same position



Expresses facts not in the data

A length is interpreted as a quantitative value;
 \therefore Length of bar says something untrue about N data

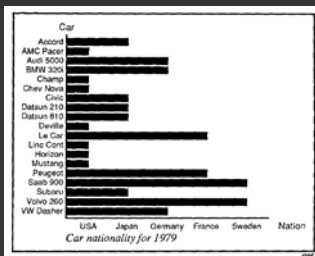


Fig. 11. Incorrect use of a bar chart for the Nation relation. The lengths of the bars suggest an ordering on the vertical axis, as if the USA cars were longer or better than the other cars, which is not true for the Nation relation.

[Mackinlay, APT, 1986]

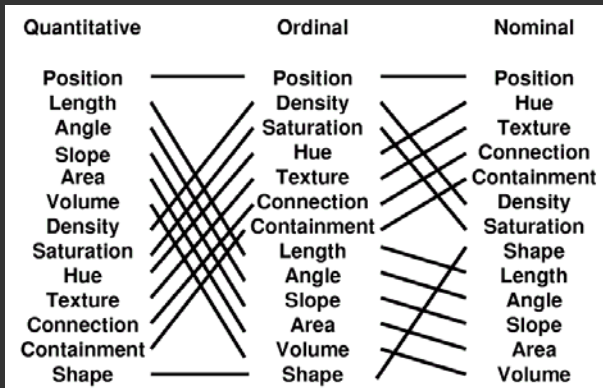
Mackinlay's effectiveness criteria

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.

Subject of *Graphical Perception* lecture

Mackinlay's ranking



Conjectured *effectiveness* of the encoding

Mackinlay's design algorithm

User formally specifies data model and type

- Additional input: ordered list of data variables to show

APT searches over design space

- Tests expressiveness of each visual encoding
- Generates image for encodings that pass test
- Tests perceptual effectiveness of resulting image

Outputs the “most effective” visualization

Limitations

Does not cover many visualization techniques

- Bertin and others discuss networks, maps, diagrams
- Does not consider 3D, animation, illustration, photography, ...

Does not model interaction

Summary

Formal specification

- Data model
- Image model
- Encodings mapping data to image

Choose expressive and effective encodings

- Formal test of expressiveness
- Experimental tests of perceptual effectiveness