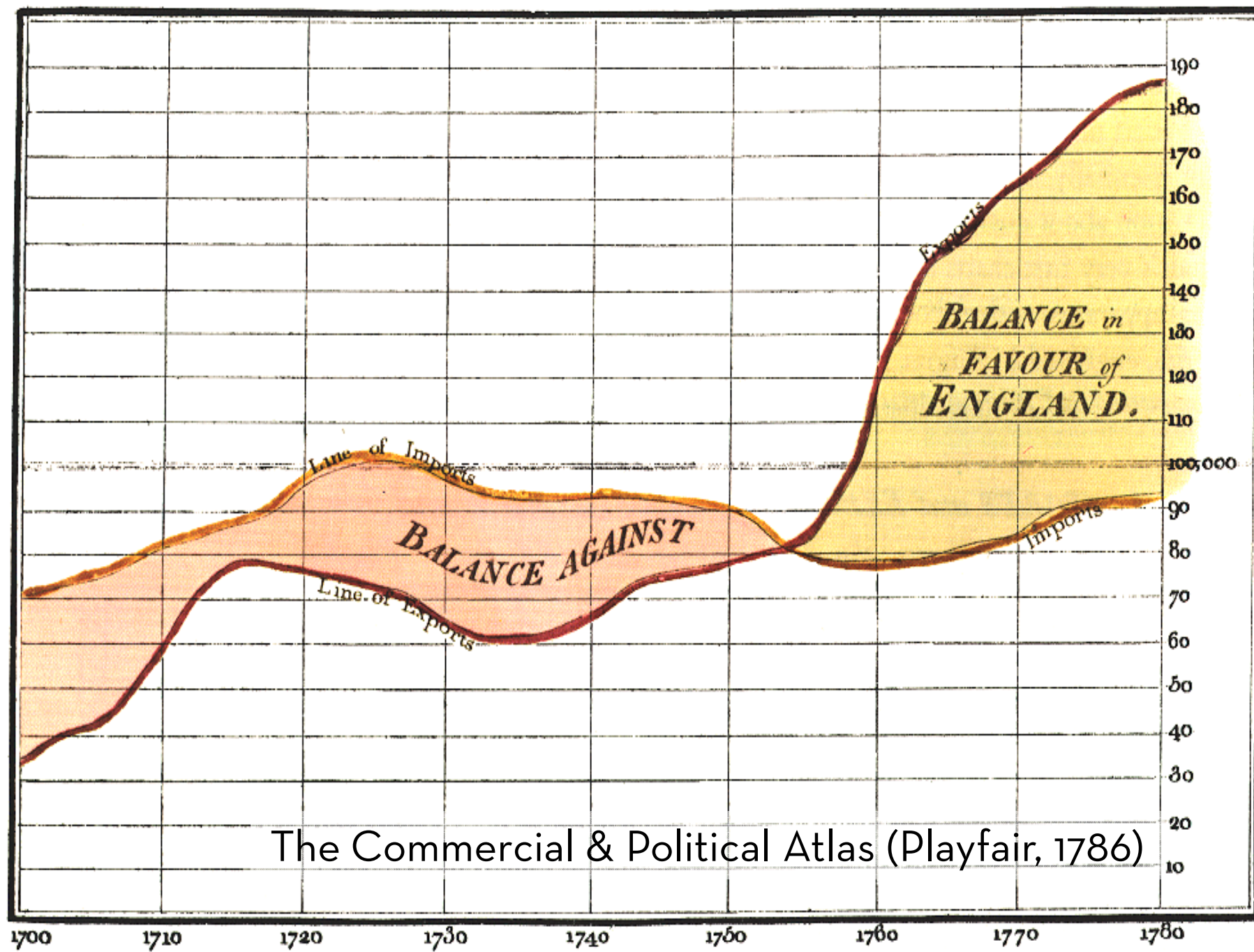# ReVision

Automated Classification, Analysis and Redesign of Chart Images

Jeffrey Heer & Fei-Fei Li
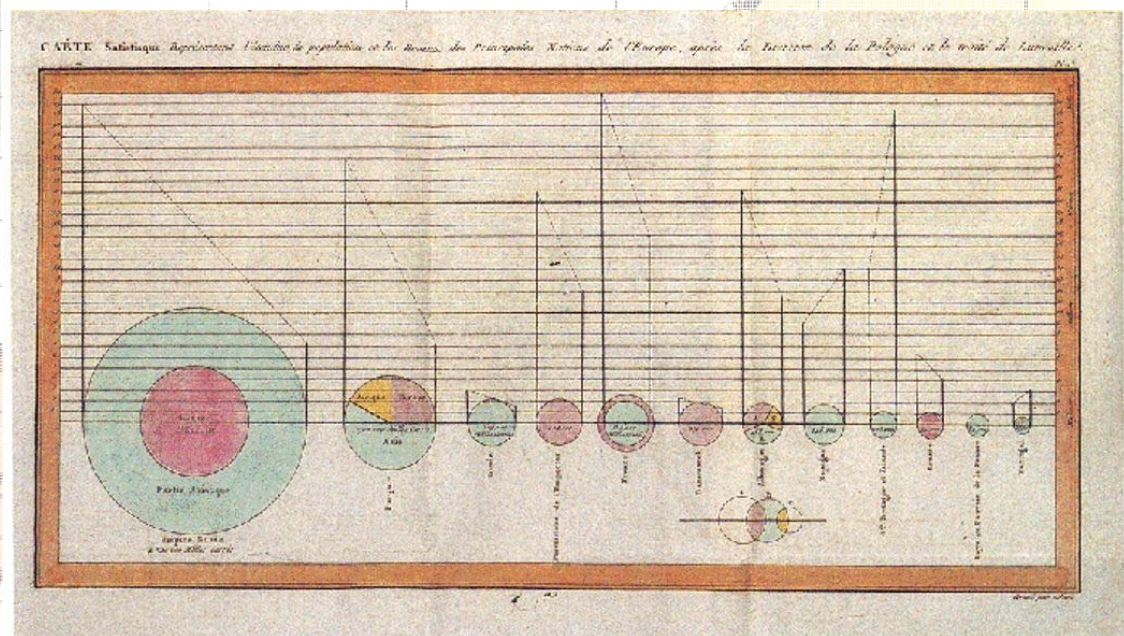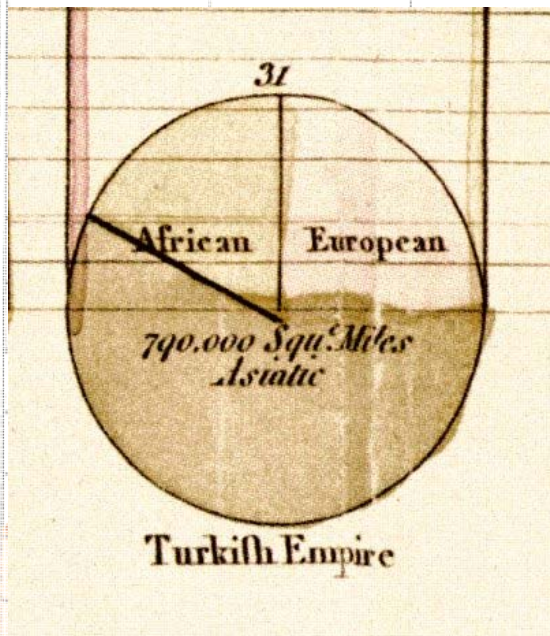
*with Manolis Savva, Nick Kong, Arti Chhajta, & Maneesh Agrawala*

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

BALANCE in FAVOUR of ENGLAND.

BALANCE AGAINST

Line of Imports

Line of Exports

Exports

Imports

The Commercial & Political Atlas (Playfair, 1786)

The Statistical Breviary (Playfair, 1801)

Statistical Atlas of the Eleventh U.S. Census (1890)

# Poor Designs are Prevalent



2005 NIH Research Budget per Death

Cardiovascular (CVD = Heart & Stroke)
Hepatitis C
Diabetes
Hepatitis B
Prostate
Alzheimer's
Parkinson'
AIDS

© Copyright FAIR Foundation 2004

The Issue — In 2005, **$10,833** will be spent per Hepatitis C death, **$8,000** per Hepatitis B death while **$178,975** will be spent per AIDS death



**Percentage of Entries in Consumer Reports' 2006 Used-Cars-to-Avoid List for Model Years 1998 to 2005, by Manufacturer**

| Auto Manufacturer | Percent |
|---|---|
| General Motors Corporation | 36 |
| DaimlerChrysler AG | 23 |
| Ford Motor Company | 16 |
| Volkswagen AG | 14 |
| BMW AG | 4 |
| Hyundai Motor Company | 2 |
| Nissan Motor Company | 2 |
| Mazda Motor Corporation | 1 |
| Isuzu Motors Ltd. | 1 |
| Porsche AG | 1 |
| Honda Motor Company | 1 |
| Mitsubishi Motors Corporation | 0 |
| Subaru Division of Fuji Heavy Industries, Ltd. | 0 |
| Toyota Motor Corporation | 0 |

100 Most Active Tweeters

download11
suhd
iggym
paviles
Systim
silverfighter
saurabhshah
giographix
DianaKhalil
dotnetshoutout
jeffsand
LukCAD
inkhead
alexpuig
phpcamp
MSExpression
hashajax

y's / Christie's

Worldwide Sales
et Share Analysis

1988    1989    1990

Planning Board approves
Widewaters development

IC students occupy Job Hall

Cayuga Vocal Ensemble
ushers in the holidays with
"Judas Maccabeus"

www.lifeslittlemysteries.com

Pie Chart

sumer Brands North America, the makers of the

ost favorite types of pie?

Apple 47%

Pumpkin 37%

Chocolate creme 32%

Apple crumb 25%    Cherry 27%

*Total adds up to more than 100 percent because people were asked to rank their three favorite types of pie.

SOURCES: SCHWAN'S CONSUMER BRANDS N.A. PIE PREFERENCE SURVEY, 2008; DREAMSTIME

KARL TATE, lifeslittlemysteries.com

**2012 PRESIDENTIAL RUN**

GOP CANDIDATES

BACK PALIN

70%

63%

60%

BACK HUCKABEE    BACK ROMNEY

FOX
47°

SOURCE: OPINIONS
DYNAMIC

In 2005,

THE SHRINKING FAMILY DOC
In California

Percentage of Doctors Devoted Solely to Family

1964     1975
27%      16.0%

1: 4,232
6,212

1: 3,167
6,694

1: 2,247 RATIO TO POPULATION
8,023 Doctors

Los Angeles Times, August 5, 1979, p. 3.

# Bitmap Images are Common



Hard to extract data

Hard to index and search

Hard to redesign

# The Challenge

**Computational Visualization Interpretation**
Given a bitmap image of a chart:
   Determine the chart type (bar, pie, *etc*)
   Infer a model of the underlying data

**Applications**
Retarget to new charts or visual styles
Data extraction, indexing & retrieval
Evaluate visualization designs

**2005 NIH Research Budget per Death**

Cardiovascular (CVD = Heart & Stroke)

Hepatitis C

Hepatitis B

Diabetes

Prostate

Alzheimer's

Parkinson'

AIDS

© Copyright FAIR Foundation 2004

The Issue → In 2005, **$10,833** will be spent per Hepatitis C death, **$8,000** per Hepatitis B death while **$178,975** will be spent per AIDS death

## Type: Pie

↓

## Data Table

| | |
|---|---|
| AIDS | 70.0% |
| Parkinson' | 6.0% |
| Prostate | 5.2% |
| Alzheimer's | 5.1% |
| Diabetes | 5.1% |
| Hepatitis C | 4.0% |
| Hepatitis B | 3.5% |
| Cardiovasc... | 1.1% |

# ReVision

1.  Classification
2.  Mark & Data Extraction
3.  Automated Redesign

1. **Classification**
2. Mark & Data Extraction
3. Automated Redesign

# Asset allocation by type



**"Pie"**

**"Bar"**

Edad

| | |
|---|---|
| 90-99 | |
| 80-89 | |
| 70-79 | |
| 60-69 | |
| 50-59 | |
| 40-49 | |
| 30-39 | |
| 20-29 | |
| 10-19 | |
| 0-9 | |

0    5,000    10,000    15,000    20,000    25,000

**Miles de habitantes**

# Prior Work

State-of-the-Art: **Prasad et al. 2007**

Combines many vision features: Image Segmentation, HOG, SIFT, and more

Input to a Multi-Class Support Vector Machine

Achieves **80%** classification accuracy across 5 classes: bar, pie, line, scatter & surface charts

Our Approach: **Unsupervised Feature Learning**

Image features at the level of *graphical marks*

**Edad**

90-99
80-89
70-79
60-69
50-59
40-49
30-39
20-29
10-19
0-9

0    5,000    10,000    15,000    20,000    25,000

**Miles de habitantes**

Perform K-Means clustering...

Cluster centroids are treated as a codebook.
These patches serve as our image features.

Bar Charts

Pie Charts

Scatter Plots

# Asset allocation by type

# Chart Image Classification

# Chart Image Classification

Asset allocation by type

Platinum
Silver
Gold
Cash
Bonds
Stocks

# Chart Image Classification



Asset allocation by type

Platinum
Silver
Gold
Cash
Bonds
Stocks

# Chart Image Classification

# Chart Image Classification

# Chart Image Classification



Asset allocation by type

Platinum
Silver
Gold
Cash
Bonds
Stocks

# Chart Image Classification



Asset allocation by type

Platinum
Silver
Gold
Cash
Bonds
Stocks

# Chart Image Classification

# Chart Image Classification



```
SVM
Classifier
```
→ **"Pie"** →
```
Mark & Data
Extraction
```

# Classification Results

Prasad Corpus: 667 images, 5 chart types

| 5-way Classifier | Accuracy |
| --- | --- |
| Prasad et al. | 84% |
| ReVision Image Features | 88% |

*What about text features?*

# Implied Volatility Surface



**Strike**

**Expiry**

**Implied Volatility**

Strike axis values: 5, 10, 12.5, 17, 19.5, 22, 24.5, 27, 29.5, 32, 34.5, 37, 40, 45

Expiry axis values: Feb 06, Apr 06, Jul 06, Nov 06, Jan 07, Jan 08

Implied Volatility axis values: 5%, 25%, 45%, 65%, 85%, 105%, 125%, 145%, 165%, 185%, 205%

# Textual Features

Histograms of:

  Text pixels per image region (8 x 8 grid)

  Text label geometry (position, width, height)

  Angles between label pairs

  Distances between label pairs

Each histogram is normalized and then
  concatenated to the final feature vector.

# Classification Results

Prasad Corpus: 667 images, 5 chart types

| **5-way Classifier** | **Accuracy** |
|---|---|
| Prasad et al. | 84% |
| ReVision Image Features | 88% |
| ReVision Text Features | 66% |
| ReVision All Features | 90% |

| **Binary Classifiers** | **Accuracy** |
|---|---|
| ReVision All Features | 96% |

# Classification Results



Extended Corpus: 2,601 images, 10 chart types

| Classifier | Accuracy |
| --- | --- |
| ReVision Image Feat., 10-way | 80% |
| ReVision Image Feat., Binary | 96% |

1. **Classification**
2. Mark & Data Extraction
3. Automated Redesign

1. Classification
2. **Mark & Data Extraction**
3. Automated Redesign

# Initial Assumptions

No 3D effects, stacked bars, or exploded pies

Marks are shaded with (nearly) constant colors

# Bar chart extraction



| Connected Components | Find Foreground Rectangles | Infer Chart Orientation | Extract Baseline Axis | Data extraction |
|---|---|---|---|---|

| Y-value | X-value |
|---|---|
| 50 | A |
| 25 | B |
| 4 | C |
| 75 | D |

# Bar charts: Bar extraction

# Bar charts: Bar extraction

# Bar charts: Bar extraction

Connected Components → Find Foreground Rectangles → Infer Chart Orientation → Extract Baseline Axis → Data extraction

# Bar charts: Bar extraction

# Bar charts: Bar extraction

# Bar charts: Bar extraction

# Bar charts: Axis extraction

# Bar charts: Axis extraction

# Bar charts: Axis extraction

# Bar charts: Axis extraction

# Bar charts: Data extraction

Connected Components → Find Foreground Rectangles → Infer Chart Orientation → Extract Baseline Axis → Data extraction



| Y-value | X-value |
|---------|---------|
|         |         |
|         |         |
|         |         |
|         |         |

# Bar charts: Data extraction



| Connected Components | Find Foreground Rectangles | Infer Chart Orientation | Extract Baseline Axis | Data extraction |



d=60

| Y-value | X-value |
| --- | --- |
|  |  |
|  |  |
|  |  |
|  |  |

Scale: 60 / (70-40) = 2 pixels/unit

# Bar charts: Data extraction

| Connected Components | Find Foreground Rectangles | | Infer Chart Orientation | Extract Baseline Axis | | Data extraction |
|---|---|---|---|---|---|---|



| Y-value | X-value |
|---|---|
| 50 | |
| 25 | |
| 4 | |
| 75 | |

Scale: 60 / (70-40) = 2 pixels/unit

# Bar charts: Data extraction

| Connected Components | Find Foreground Rectangles | | Infer Chart Orientation | Extract Baseline Axis | | Data extraction |
|---|---|---|---|---|---|---|

| Y-value | X-value |
|---|---|
| 50 | |
| 25 | |
| 4 | |
| 75 | |

# Bar charts: Data extraction

Connected Components → Find Foreground Rectangles → Infer Chart Orientation → Extract Baseline Axis → Data extraction

| Y-value | X-value |
|---------|---------|
| 50 | A |
| 25 | B |
| 4 | C |
| 75 | D |

# Pie chart extraction



Compute gradient → Fit ellipse with RANSAC → Unroll the pie → Peaks in horizontal derivative → Data extraction

| Percentage | Category |
|------------|----------|
| 22.3 | A |
| 22.4 | B |
| 10.8 | C |
| 5.6 | D |
| 5.6 | E |
| 33.3 | F |

# Extraction Results



Bar chart titled "Extraction Results" showing Number of Charts for Bar and Pie categories.

**Bar:** Total charts 52, Mark extractions 41 (79%), Data extractions 29 (56%)

**Pie:** Total charts 53, Mark extractions 33 (62%), Data extractions 21 (40%)

Legend: ■ Total charts ■ Mark extractions ■ Data extractions

# Average Data Extraction Error

Error = | true_value – est_value | / max_value

Bar chart:
- Less than $10,000: 562
- $10,000 - $14,999: 434
- $15,000 - $24,999: 966
- $25,000 - $34,999: 830
- $35,000 - $49,999: 925
- $50,000 - $74,999: 891
- $75,000 - $99,999: 238
- $100,000 - $149,999: 163
- $150,000 - $199,999: 27
- $200,000 or more: 25

Pie chart:
- Not immediate deaths: 10%
- Immediate deaths: 44%
- Diseases: 29%
- 100% Pensions: 17%

Bar Charts: 0.0093

Pie Charts: 0.0034

Average chart size: 452 x 342 pixels

1. Classification
2. **Mark & Data Extraction**
3. Automated Redesign

1. Classification
2. Mark & Data Extraction
3. **Automated Redesign**

Compare area of circles

Compare length of bars

# Graphical Perception Experiments

Absolute Log Estimation Error

Heer & Bostock, 2010

Most accurate

Position (common) scale
Position (non-aligned) scale

Length

Slope

Angle

Area

Volume

Least accurate

Color hue-saturation-density

# Asset allocation by type

Platinum

Silver

Gold

Cash

Bonds

Stocks

# ReVision Gallery

## Input Image (upload)



Asset allocation by type

## Data Table (export)

| Label | % of Total |
| --- | --- |
| Cash | 32.7% |
| Bonds | 22.4% |
| Stocks | 22.3% |
| Gold | 10.8% |
| Platinum | 6.1% |
| Silver | 5.7% |

## Font

Lucida Grande

## Color

ManyEyes Red

# ReVision Gallery

## Input Image (upload)



Asset allocation by type

## Data Table (export)

| Label | % of Total |
|---|---|
| Cash | 32.7% |
| Bonds | 22.4% |
| Stocks | 22.3% |
| Gold | 10.8% |
| Platinum | 6.1% |
| Silver | 5.7% |

## Font

Lucida Grande

## Color

ManyEyes Red

**Individual heights**

Height/cm

175
170
165
160
155
150
145
140
135

Melissa, Jayne, Hayley, Suzanne, Kate, Cara, Clare, Kirsty, Devinder, Caroline, Sandra, Kay, Sarah, Amanda, Lisa, Carol, Maria, Samantha, Joanne

Name
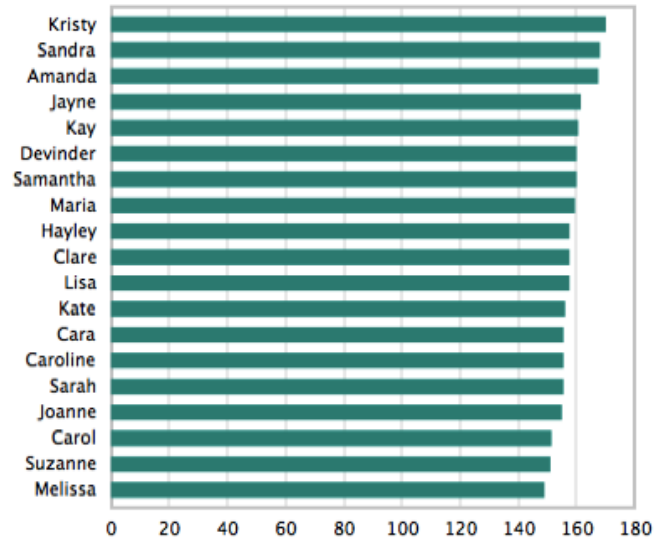
# ReVision Gallery

## Input Image (upload)



## Data Table (export)

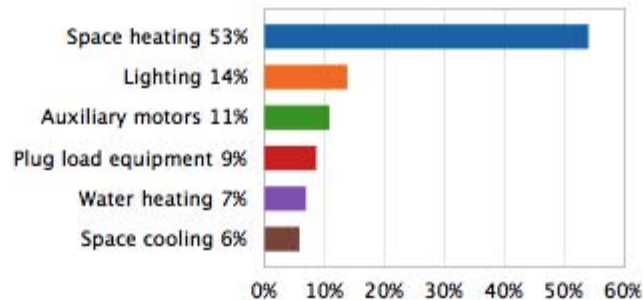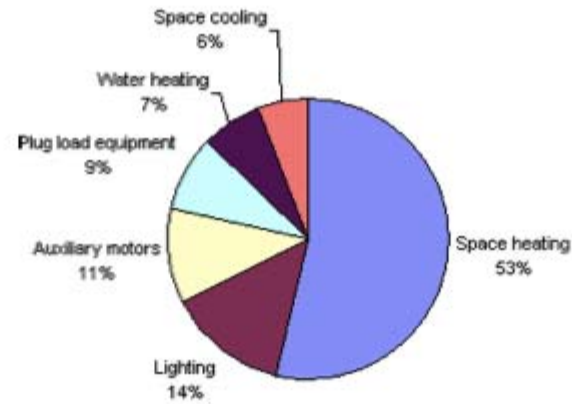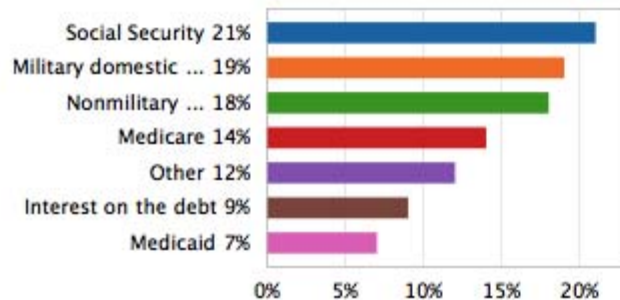| Label | Value |
|---|---|
| Kristy | 170 |
| Sandra | 168 |
| Amanda | 167 |
| Jayne | 161 |
| Kay | 161 |
| Devinder | 160 |
| Samantha | 160 |
| Maria | 159 |
| Hayley | 158 |
| Clare | 158 |
| Lisa | 158 |
| Kate | 156 |
| Cara | 155 |
| Caroline | 155 |
| Sarah | 155 |
| Joanne | 155 |
| Carol | 151 |
| Suzanne | 151 |
| Melissa | 149 |

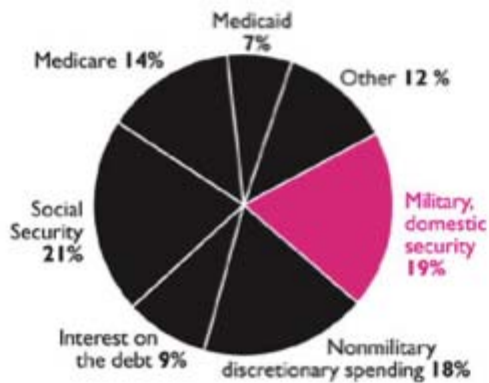## Font

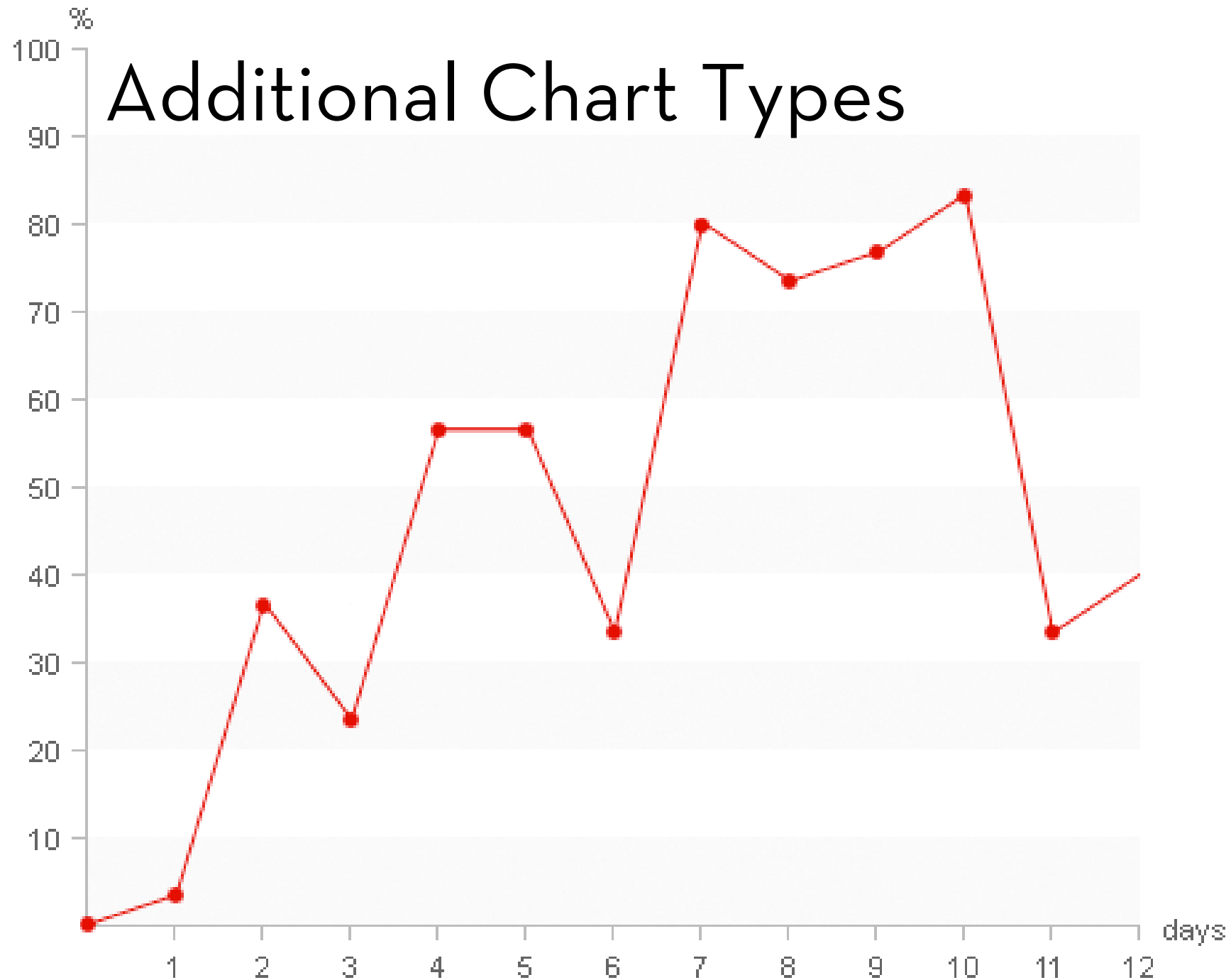Lucida Grande

## Color

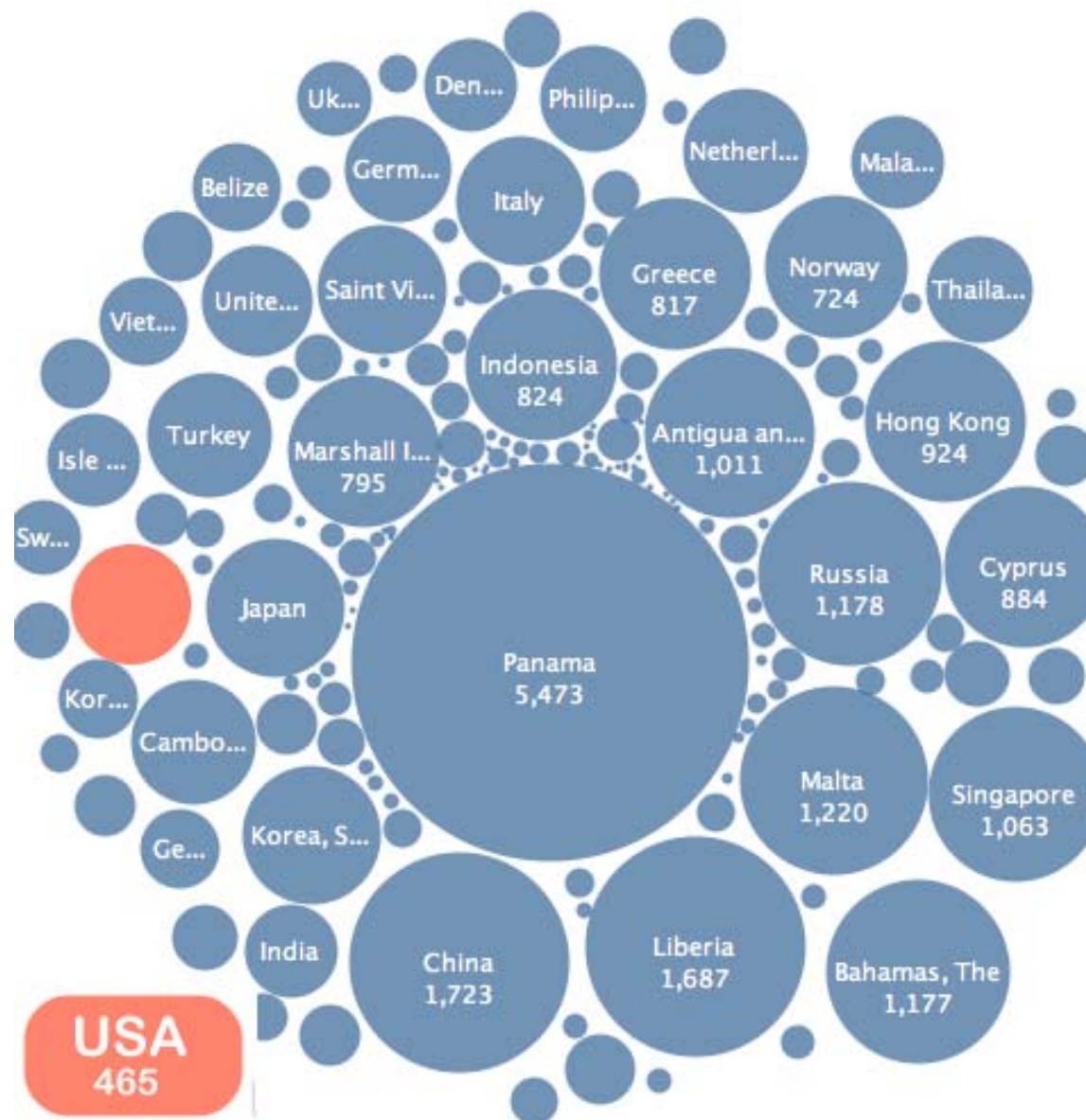Economist

# Bar Chart Redesigns

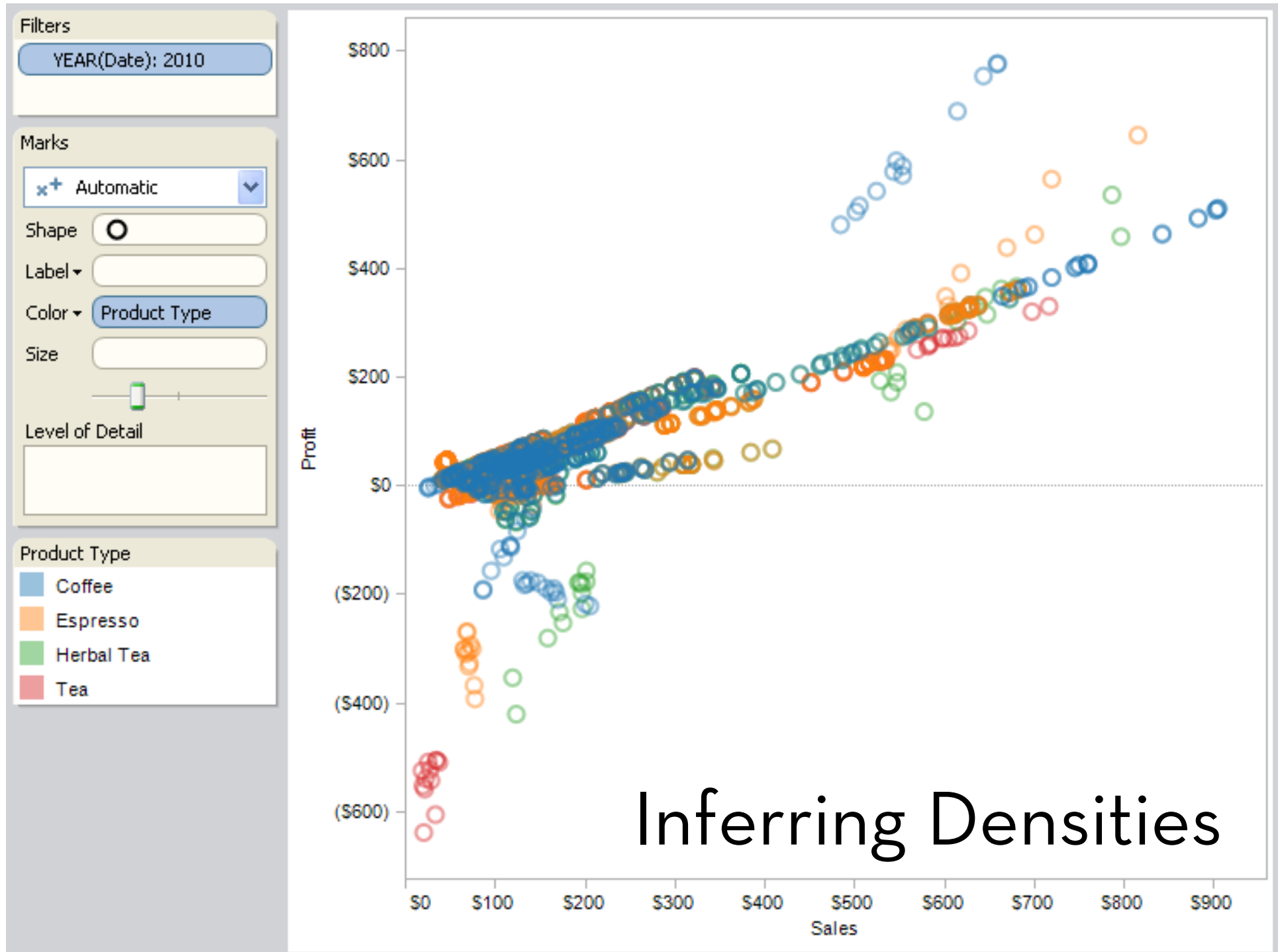# Pie Chart Redesigns

# Future Work

# Additional Chart Types

# Additional Chart Types

# Extracting Legends

Inferring Densities
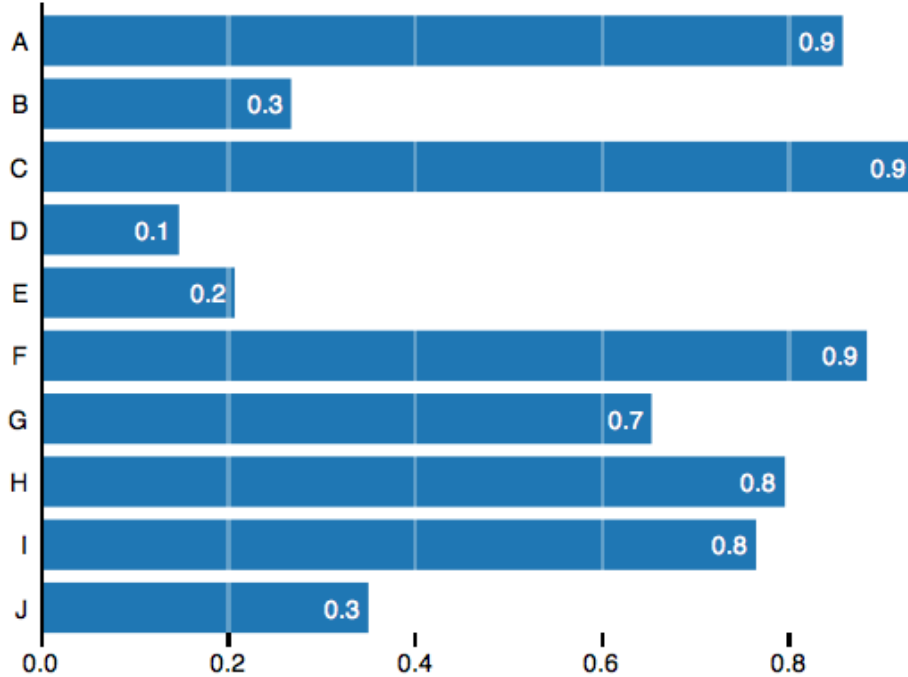
# Generative Model
Infer chart specifications, not just backing data.

```
/* Sizing and scales. */
var w = 400,
    h = 250,
    x = pv.linear(0, 1.1)
        .range(0, w),
    y = pv.ordinal(pv.range(10))
        .splitBanded(0, h, 4/5);

/* The root panel. */
var vis = new pv.Panel()
    .width(w)
    .height(h)
    .bottom(20)
    .left(20)
    .right(10)
    .top(5);

/* The bars. */
var bar = vis.add(pv.Bar)
    .data(data)
    .top(function() y(this.index))
    .height(y.range().band)
    .left(0)
    .width(x);
```

# ReVision
## Automated Classification, Analysis and Redesign of Chart Images

Jeffrey Heer & Fei-Fei Li

*with Manolis Savva, Nick Kong, Arti Chhajta, & Maneesh Agrawala*