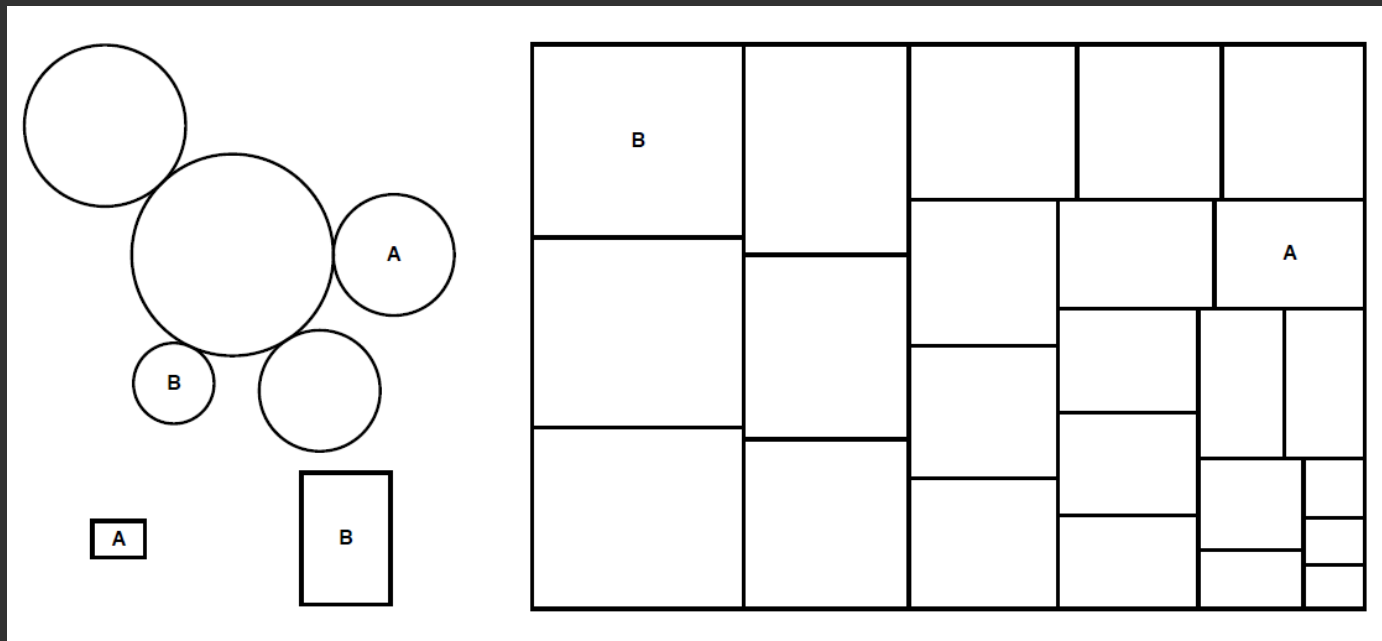# Crowdsourcing Graphical Perception
## Using Mechanical Turk to Assess Visualization Design

Jeffrey Heer & Michael Bostock
Stanford University

| Set A | | Set B | | Set C | | Set D | |
|---|---|---|---|---|---|---|---|
| X | Y | X | Y | X | Y | X | Y |
| 10 | 8.04 | 10 | 9.14 | 10 | 7.46 | 8 | 6.58 |
| 8 | 6.95 | 8 | 8.14 | 8 | 6.77 | 8 | 5.76 |
| 13 | 7.58 | 13 | 8.74 | 13 | 12.74 | 8 | 7.71 |
| 9 | 8.81 | 9 | 8.77 | 9 | 7.11 | 8 | 8.84 |
| 11 | 8.33 | 11 | 9.26 | 11 | 7.81 | 8 | 8.47 |
| 14 | 9.96 | 14 | 8.1 | 14 | 8.84 | 8 | 7.04 |
| 6 | 7.24 | 6 | 6.13 | 6 | 6.08 | 8 | 5.25 |
| 4 | 4.26 | 4 | 3.1 | 4 | 5.39 | 19 | 12.5 |
| 12 | 10.84 | 12 | 9.11 | 12 | 8.15 | 8 | 5.56 |
| 7 | 4.82 | 7 | 7.26 | 7 | 6.42 | 8 | 7.91 |
| 5 | 5.68 | 5 | 4.74 | 5 | 5.73 | 8 | 6.89 |

**Summary Statistics**
$u_X = 9.0$  $\sigma_X = 3.317$
$u_Y = 7.5$  $\sigma_Y = 2.03$

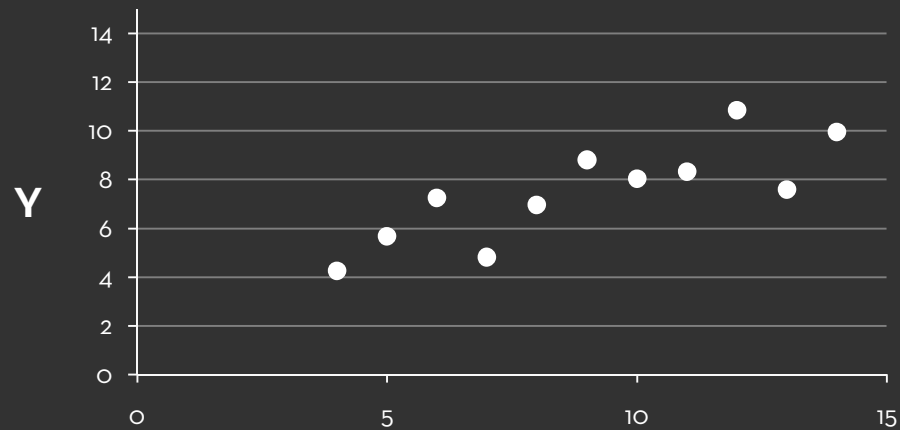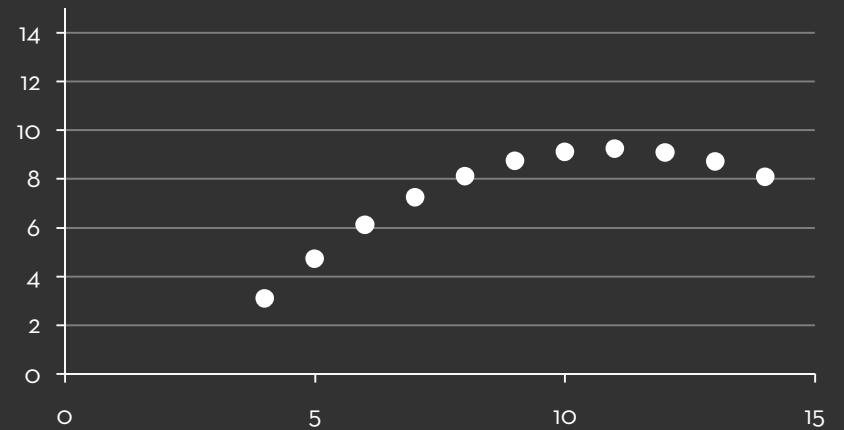**Linear Regression**
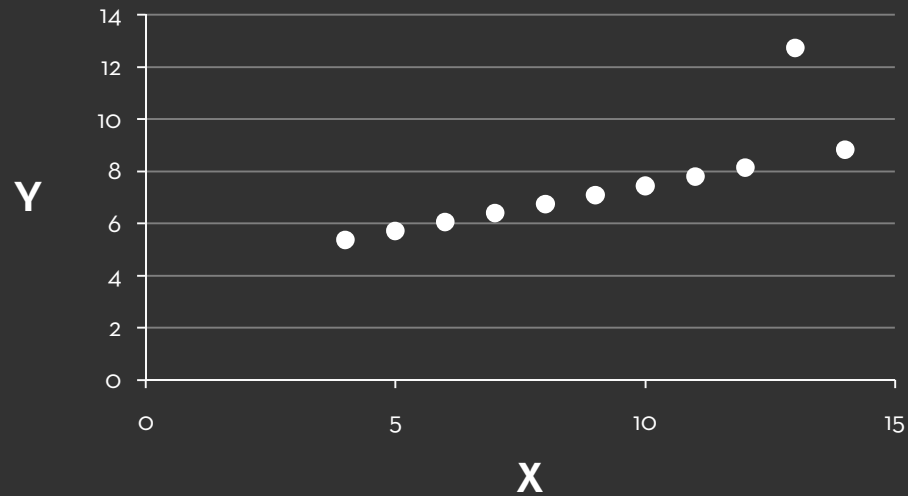$Y^2 = 3 + 0.5\,X$
$R^2 = 0.67$

[Anscombe 73]

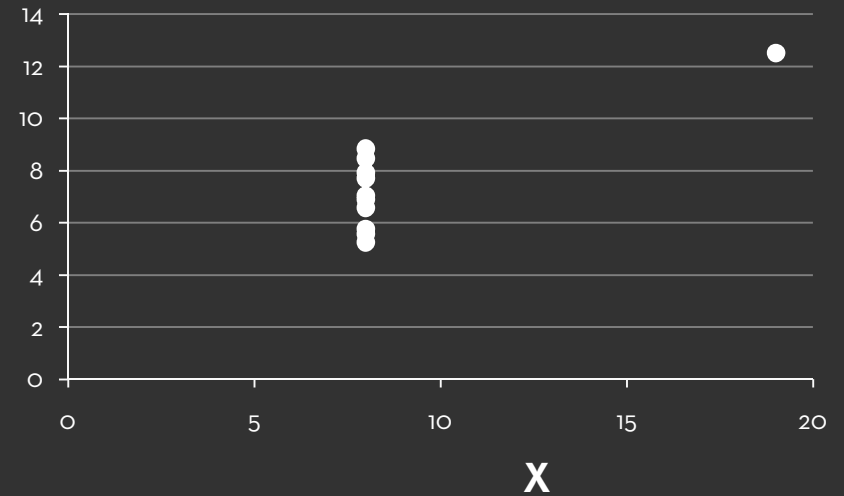hawaii.naist.jp/research/visual_e.html

Introduction | **Dashboard** | **Status** | **Account Settings**

## Mechanical Turk is a marketplace for work.

We give businesses and developers access to an on-demand, scalable workforce.
Workers select from thousands of tasks and work whenever it's convenient.
**127,286 HITs** available. <u>View them now.</u>

## Make Money
### by working on HITs

HITs - *Human Intelligence Tasks* - are individual tasks that you work on. <u>Find HITs now.</u>

**As a Mechanical Turk Worker you:**

- Can work from home
- Choose your own work hours
- Get paid for doing good work

**Find an interesting task** → **Work** → **Earn money**

[Find HITs Now]

or <u>learn more about being a **Worker**</u>

## Get Results
### from Mechanical Turk Workers
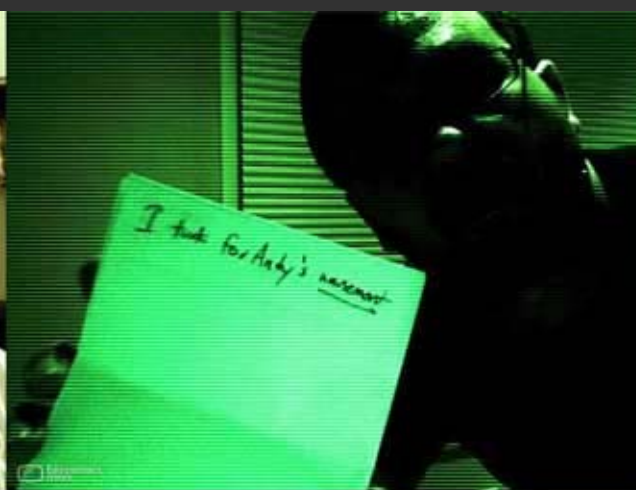
Ask workers to complete HITs - *Human Intelligence Tasks* - and get results using Mechanical Turk. <u>Register Now</u>

**As a Mechanical Turk Requester you:**

- Have access to a global, on-demand, 24 x 7 workforce
- Get thousands of HITs completed in minutes
- Pay only when you're satisfied with the results

**Fund your account** → **Load your tasks** → **Get results**

[Get Started]

# Using MTurk for Research

*Machine Learning, Comp. Vision & Info. Retrieval*
User-Generated Metadata, Labeling Data

*Kittur, Chi & Suh: Wikipedia Article Quality*
Use verifiable questions to reduce gaming
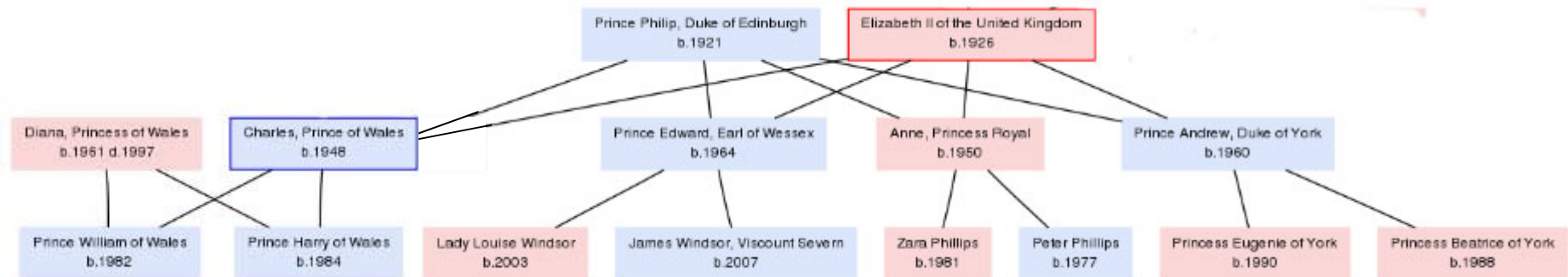Make sincere responses as easy as insincere ones

*Mason & Watts: Financial Incentives*
Higher reward → faster completion, same quality

AN EXAMPLE:

**TimeNets** for **Genealogical Data**

# Visualizing Genealogical Graphs

Prince Philip, Duke of Edinburgh

Elizabeth II of the United ...

Prince Edward, Earl of Wessex

Sophie, The Countess of Wessex

...

Lady ...

Prince Andrew, Duke of York

Sarah, Duchess of York

Princess Eugenie of Y...

Princess Beatrice of York

Anne, Princess Royal

Mark Phillips

Timothy Laurence

Zara Phillips

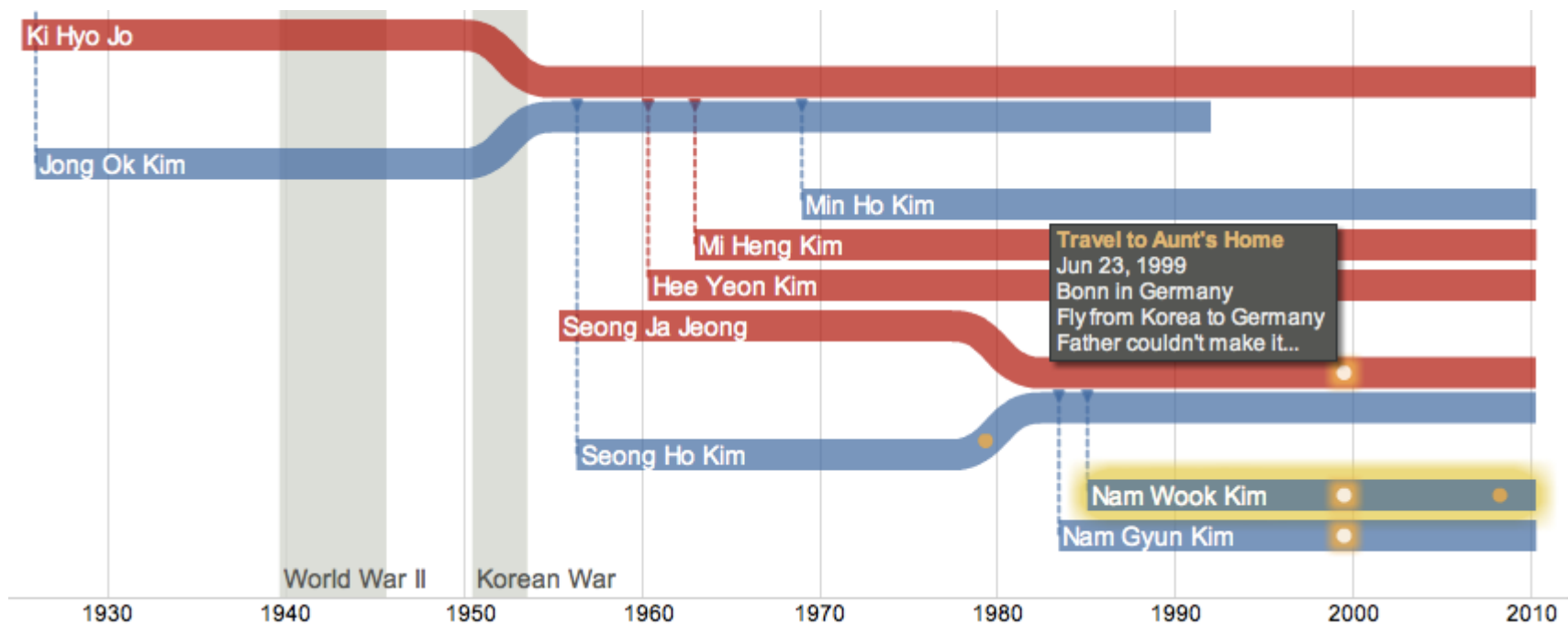Peter Mark Andrew Phi...

Autumn Kelly

Charles, Prince of Wales
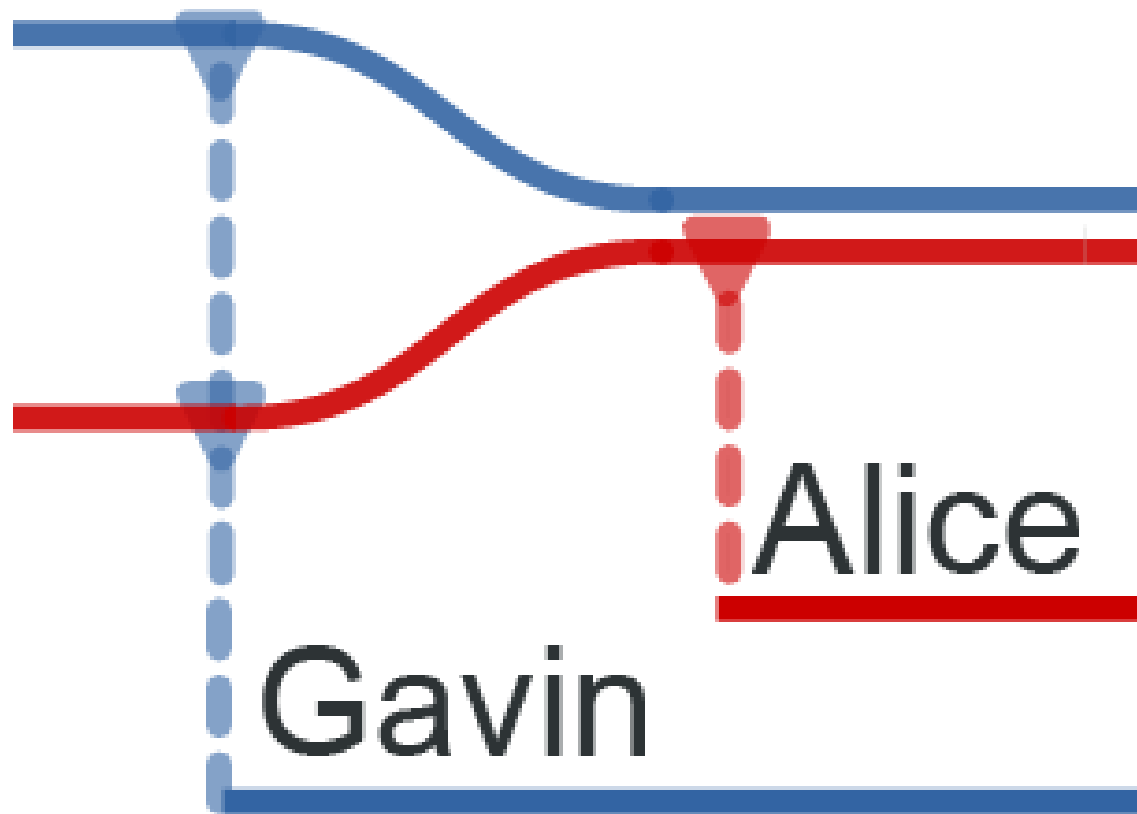
Diana, Princess of W...

Camilla, The Duchess of Cornwall

Prince Harry of Wales

Prince William of Wales

# **TimeNets** = Time x Family Trees



Ki Hyo Jo

Jong Ok Kim

Min Ho Kim

Mi Heng Kim

Hee Yeon Kim

Seong Ja Jeong

Seong Ho Kim

**Travel to Aunt's Home**
Jun 23, 1999
Bonn in Germany
Fly from Korea to Germany
Father couldn't make it...

Nam Wook Kim

Nam Gyun Kim

World War II    Korean War

1930    1940    1950    1960    1970    1980    1990    2000    2010
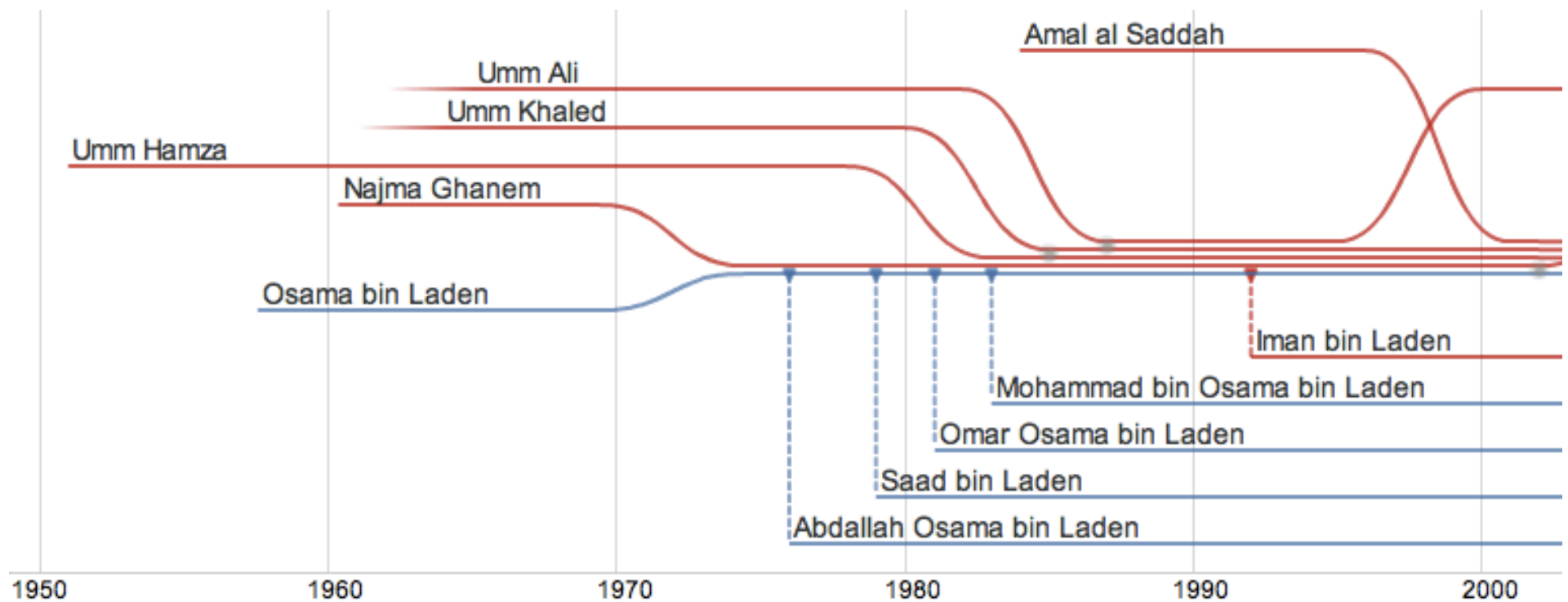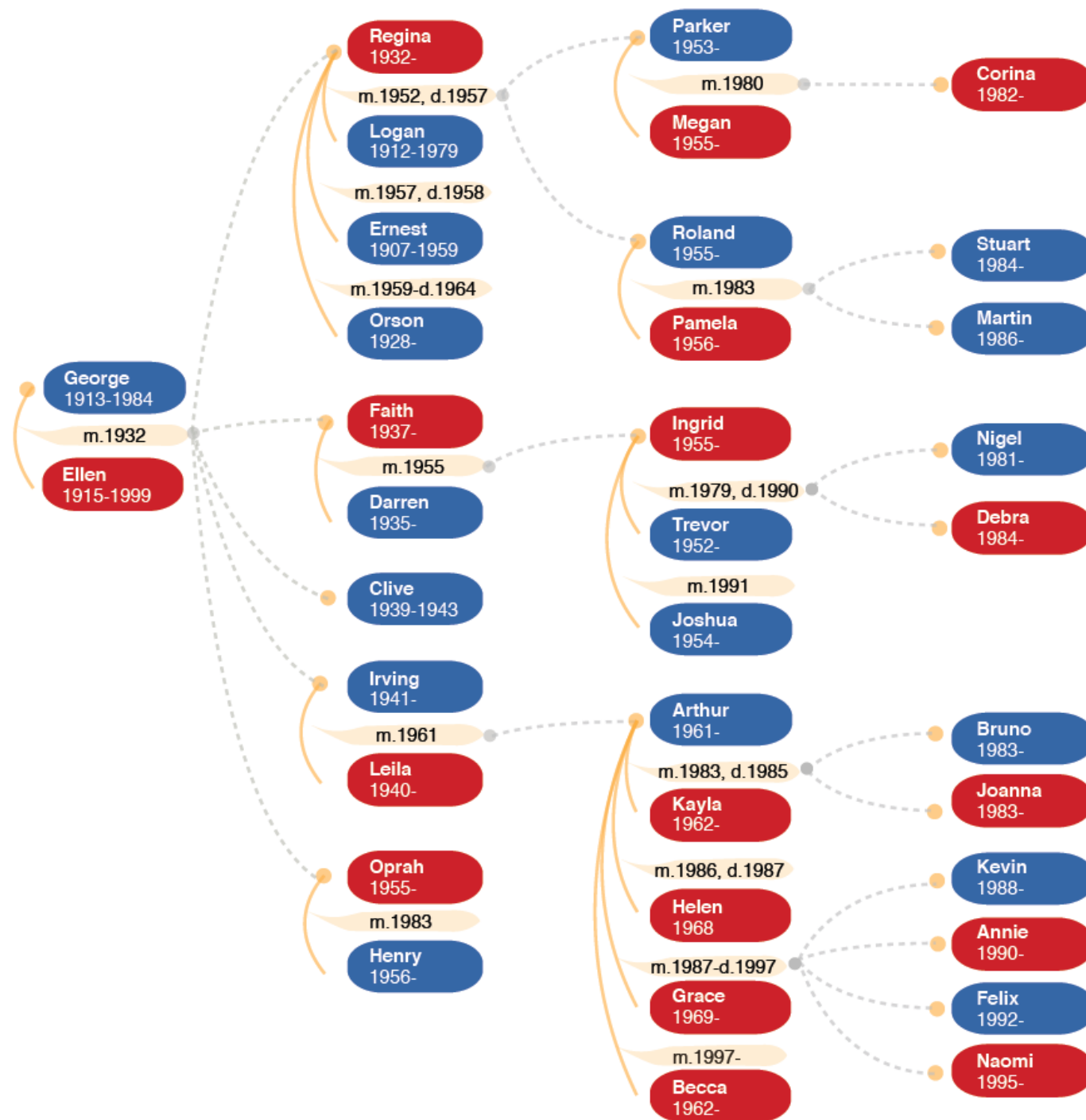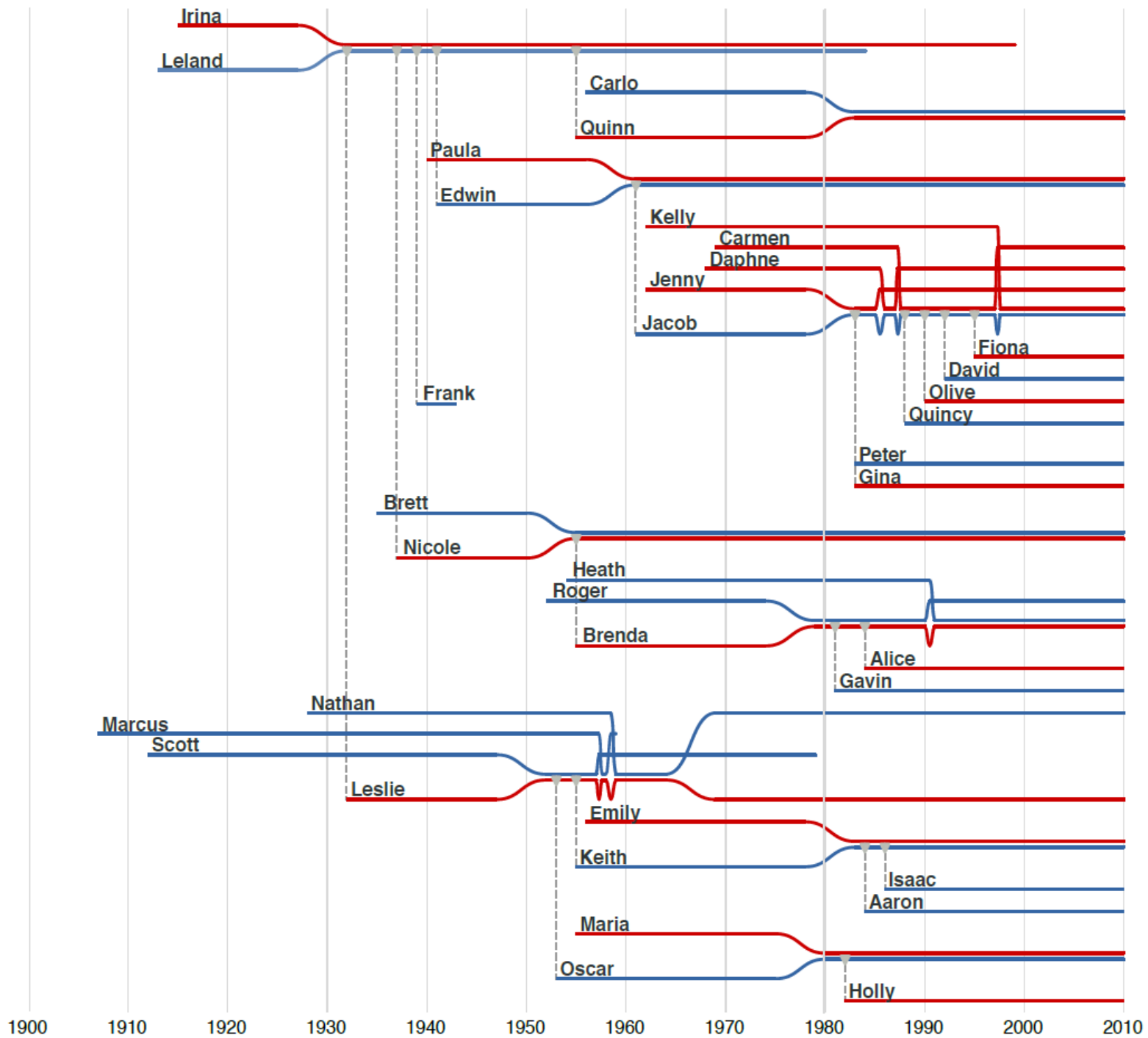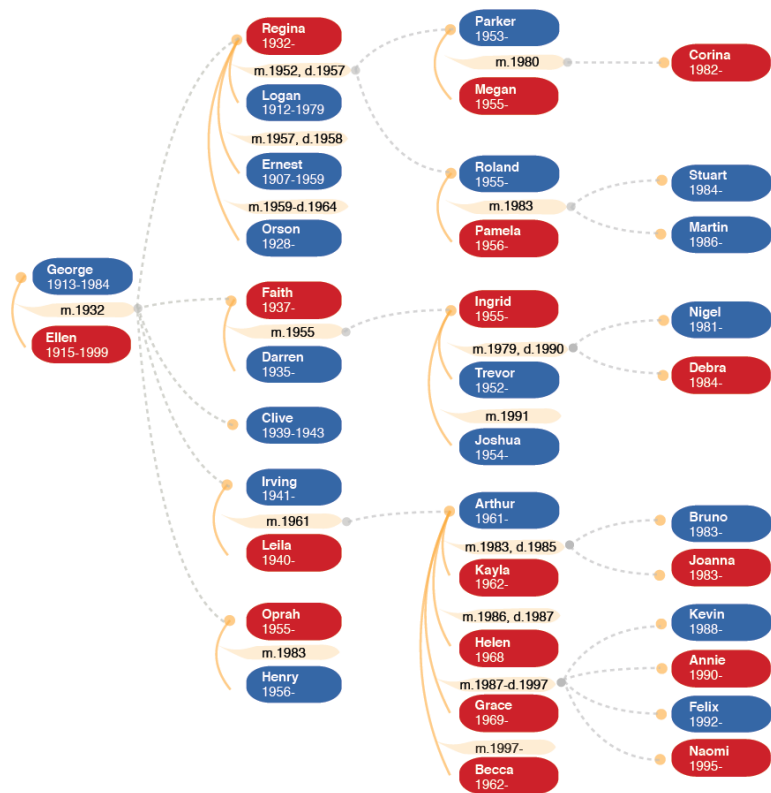
(Out-of-Wedlock Births)

# Elizabeth Taylor (Remarriage)
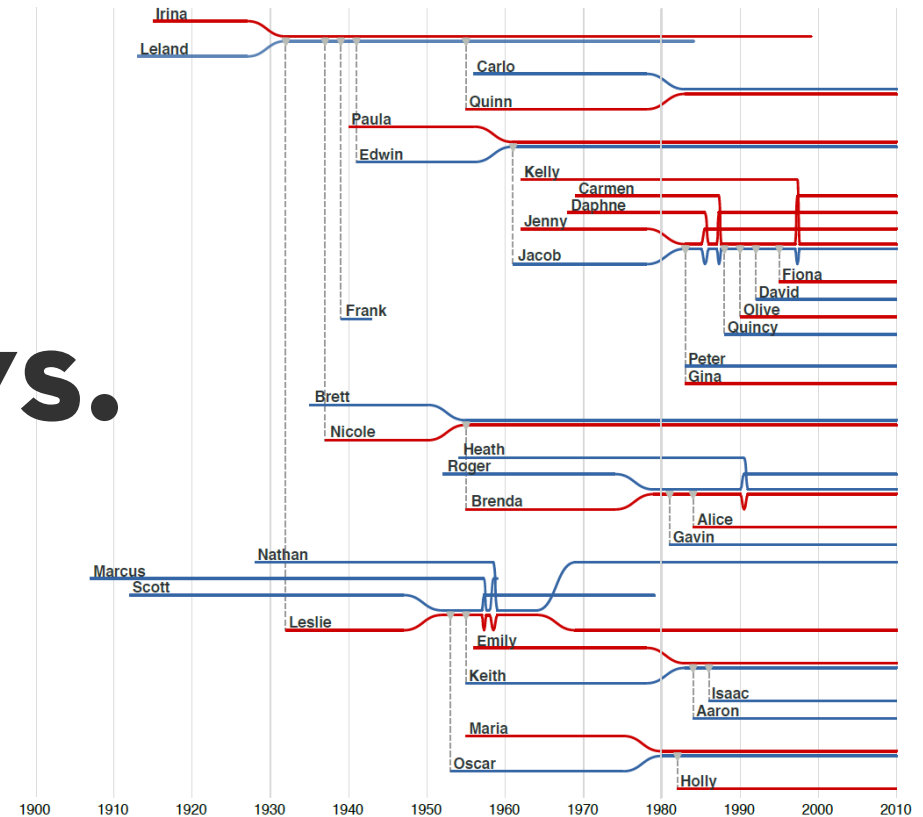
# Osama bin Laden (Polygamy)

**George** 1913-1984 — m.1932 — **Ellen** 1915-1999

**Regina** 1932- — m.1952, d.1957 — **Logan** 1912-1979
— m.1957, d.1958 — **Ernest** 1907-1959
— m.1959-d.1964 — **Orson** 1928-

**Parker** 1953- — m.1980 — **Megan** 1955-
— **Corina** 1982-

**Roland** 1955- — m.1983 — **Pamela** 1956-
— **Stuart** 1984-
— **Martin** 1986-

**Faith** 1937- — m.1955 — **Darren** 1935-

**Clive** 1939-1943

**Ingrid** 1955- — m.1979, d.1990 — **Trevor** 1952-
— m.1991 — **Joshua** 1954-
— **Nigel** 1981-
— **Debra** 1984-

**Irving** 1941- — m.1961 — **Leila** 1940-

**Arthur** 1961- — m.1983, d.1985 — **Kayla** 1962-
— m.1986, d.1987 — **Helen** 1968
— m.1987-d.1997 — **Grace** 1969-
— m.1997- — **Becca** 1962-
— **Bruno** 1983-
— **Joanna** 1983-
— **Kevin** 1988-
— **Annie** 1990-
— **Felix** 1992-
— **Naomi** 1995-

**Oprah** 1955- — m.1983 — **Henry** 1956-

Irina
Leland
Carlo
Quinn
Paula
Edwin
Kelly
Carmen
Daphne
Jenny
Jacob
Fiona
David
Olive
Quincy
Frank
Peter
Gina
Brett
Nicole
Heath
Roger
Brenda
Alice
Gavin
Nathan
Marcus
Scott
Leslie
Emily
Keith
Isaac
Aaron
Maria
Oscar
Holly

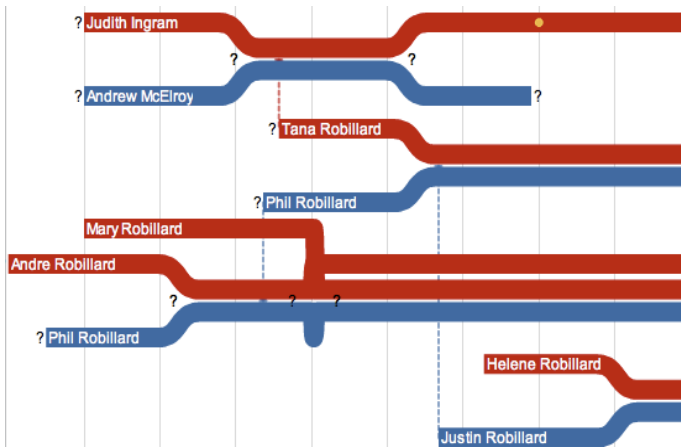1900    1910    1920    1930    1940    1950    1960    1970    1980    1990    2000    2010
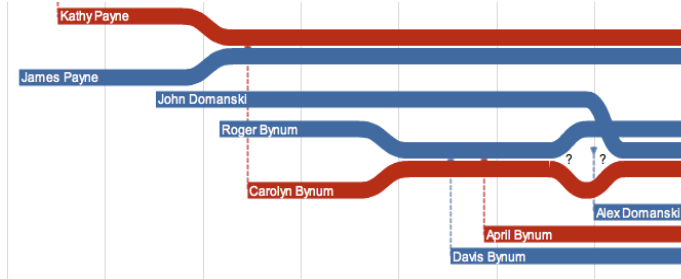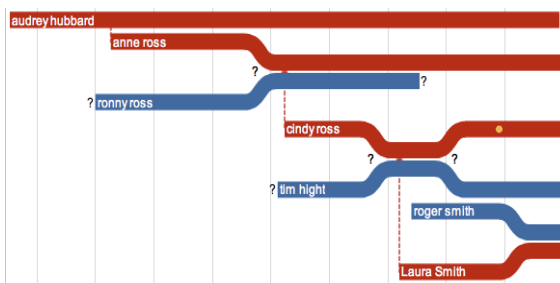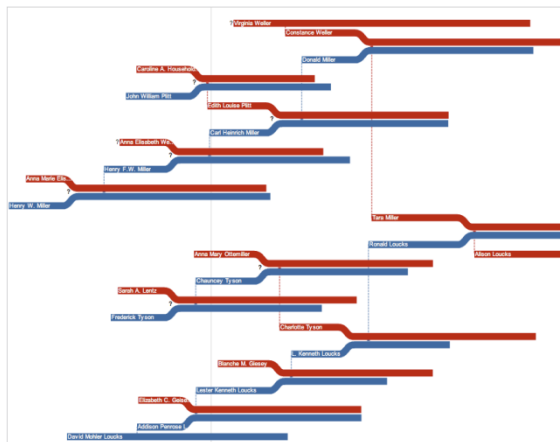
**vs.**

Asked *structural, temporal & struct x temp* tasks

**No accuracy differences** between visuals

**TimeNets were significantly faster** (~25%) for tasks with a *temporal* component

I love the idea of this tool.   I love the look and ease of this program! I

I think that the concept is very good and the effort taken is commendable.

Please don't delete my data! Very cool.   This is a very interesting idea

i was having a lot of fun with this.i love how it shows everything simply

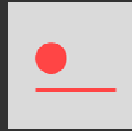hawaii.naist.jp/research/visual_e.html

# Research Goals

1. **Assess the viability** of crowdsourced perception experiments on Mechanical Turk.

2. Demonstrate the use of MTurk to **gain novel insights** for visualization design.

3. Analyze experimental data to **characterize MTurk as an experimental platform**.

# Experiment 1:
# Proportional Judgments of Spatial Data Encodings

Most accurate

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

Length
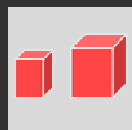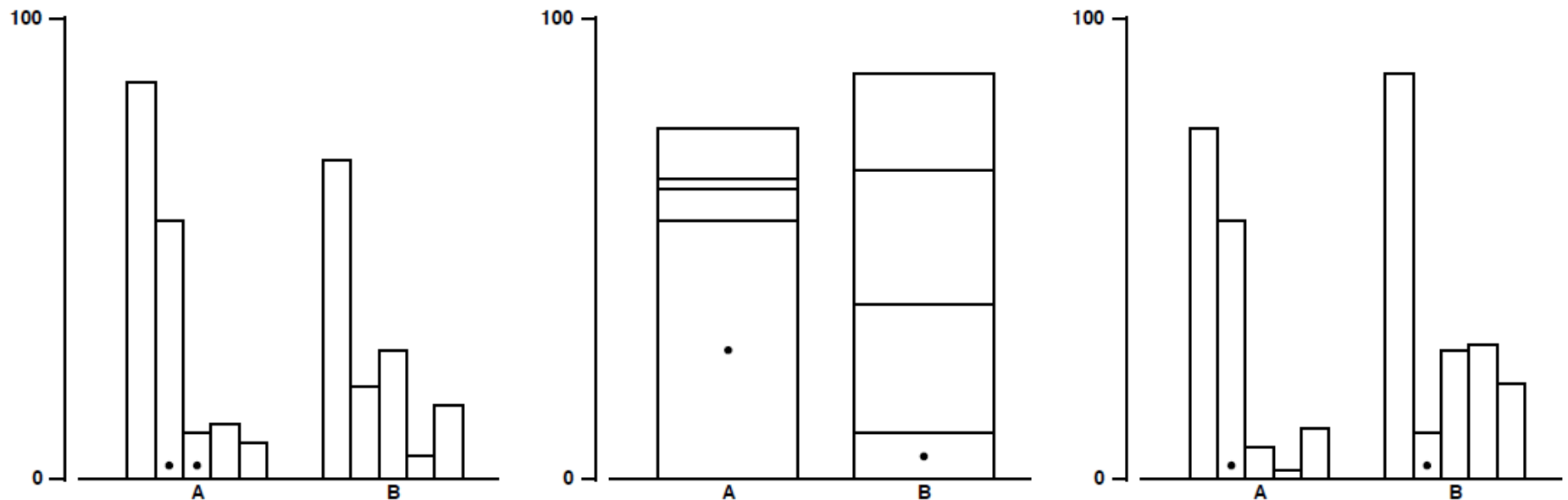
Slope

Angle

Area

Volume

Least accurate

Color hue-saturation-density

Cleveland & McGill '84

# Cleveland & McGill, 1984

Stimuli for position encodings.

Task: estimate **%** smaller element is of the larger

# Experiment 1A: Proportions

Goal: replicate Cleveland & McGill, 1984
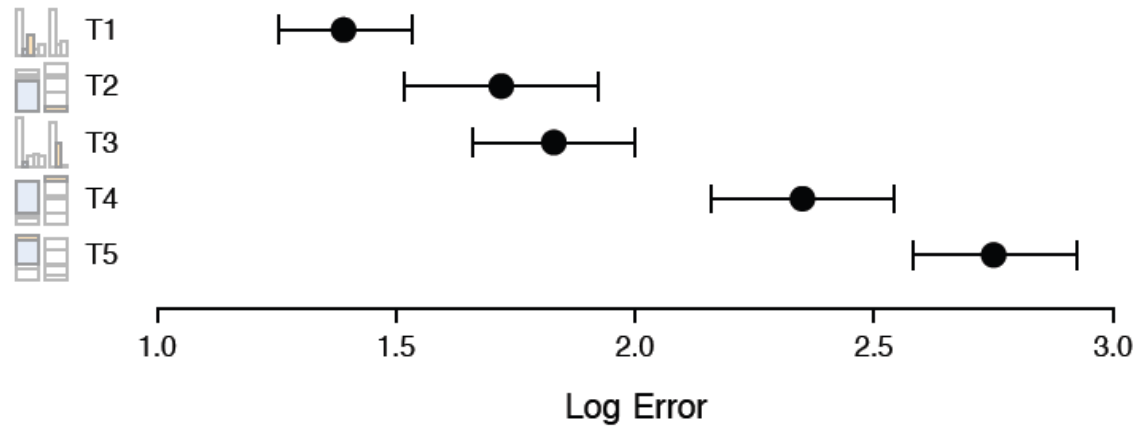
   5    *original types*: position (3) + length (2)

\+ 2   *new types*: angle + circular area

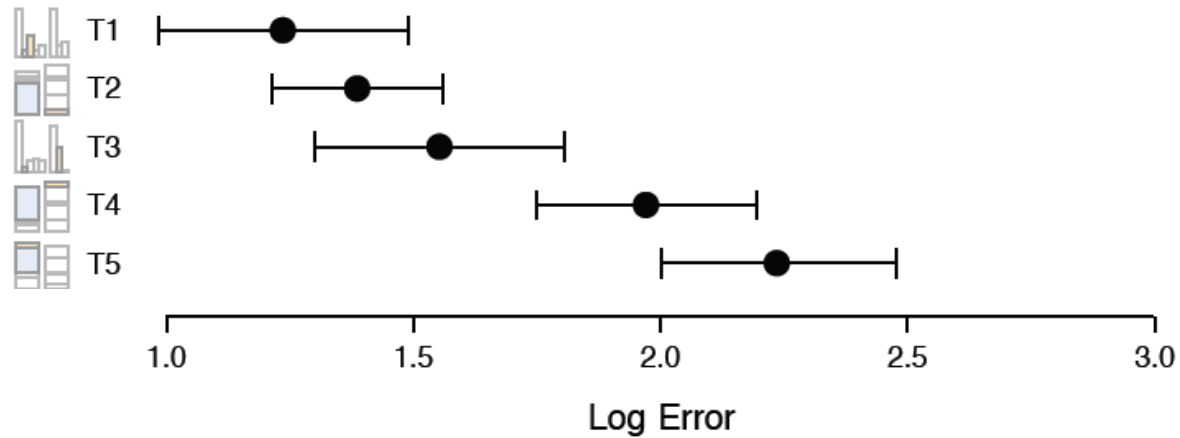x 10  *proportional differences*

N=50 assignments, \$0.05 per HIT

**Task:** estimate **%** smaller element is of the larger

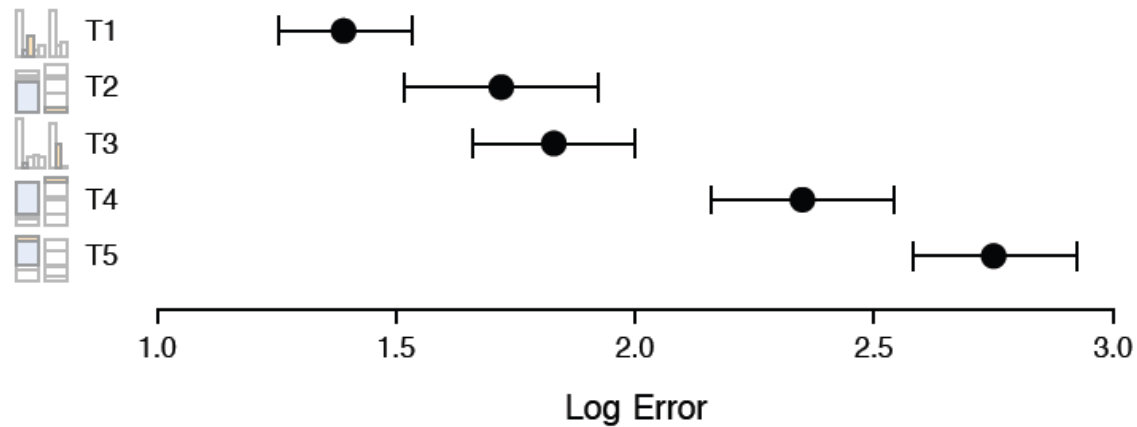**Error = $\log_2(\,|\,\text{true\%} - \text{estimated\%}\,| + 1/8)$**
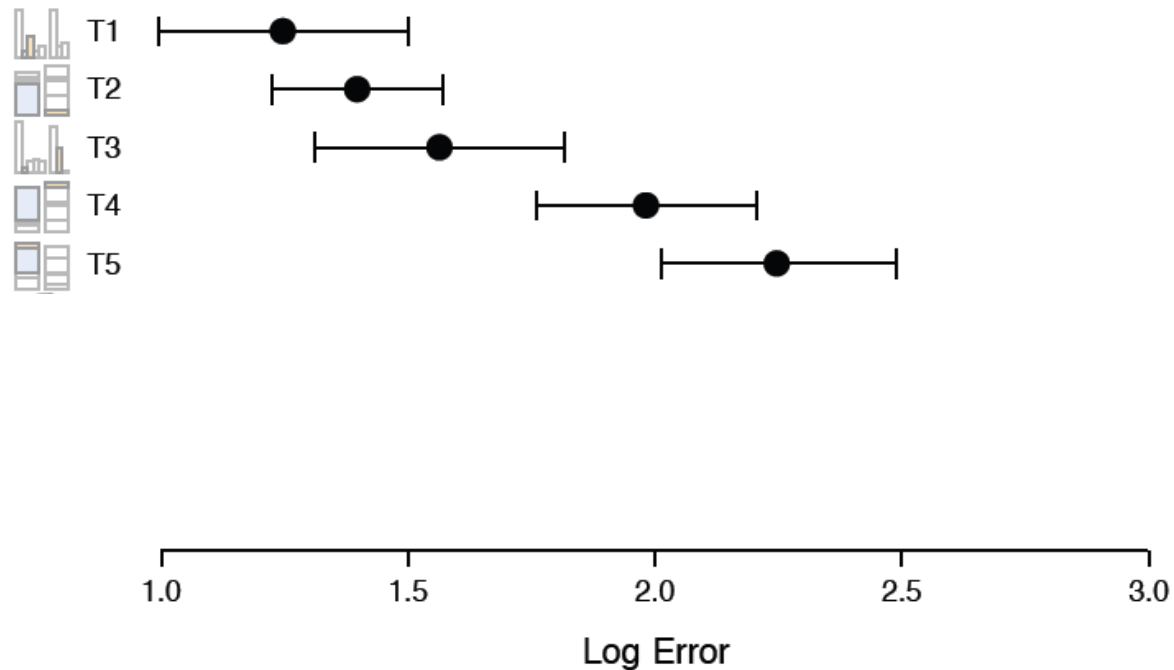
# Cleveland & McGill, 1984 (Lab Study)



# Our Crowdsourced Study

Cleveland & McGill, 1984 (Lab Study)

T1
T2
T3
T4
T5

Log Error

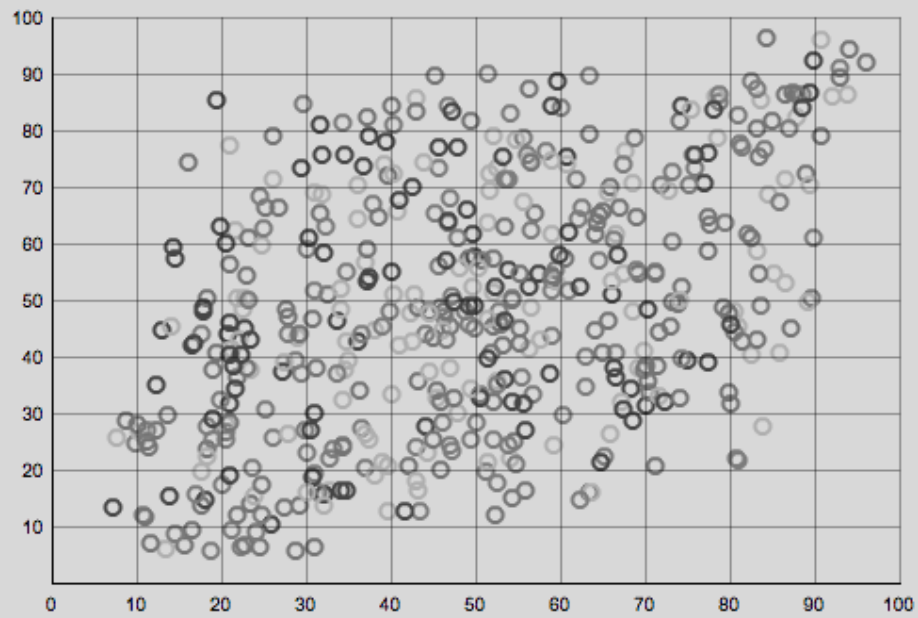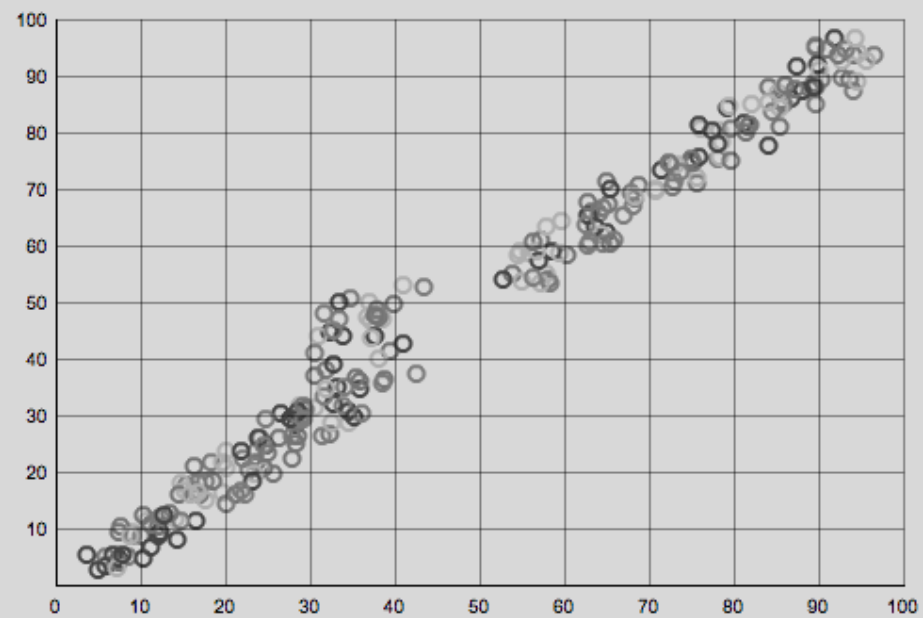Our Crowdsourced Study

T1
T2
T3
T4
T5

Log Error

# Experiment 2:
## Gridline Alpha Contrast

# Experiment 2 Tasks

**2L**: Adjust the grid so that it is as light as possible while still being usably perceptible.

**2D**: Adjust the grid strength to meet your best judgment of how obvious it can be before it becomes too intrusive and sits in front of the image; some users have called this a 'fence'.

# Experiment 2: Gridline Alpha

4 *plot density*: none, sparse, medium, dense

x 5 *background*: #f3, #d8, #be, #a5, #8e

x 3 *replications*

N=24 assignments, $0.02 per HIT

Record alpha value, User-Agent, JS "screen" info

## Stone & Bartram, 2009 (Lab Study)

Data Density — None, Sparse, Medium, Dense

Alpha: 0.05, 0.10, 0.15, 0.20, 0.25, 0.30, 0.35, 0.40, 0.45, 0.50

O Light ● Dark

## Our Crowdsourced Study

Data Density — None, Sparse, Medium, Dense

Alpha: 0.05, 0.10, 0.15, 0.20, 0.25, 0.30, 0.35, 0.40, 0.45, 0.50

O Light ● Dark

# Inferred Display Configuration

Operating System (& gamma?) from User-Agent



MacOS < 10.6: $\gamma = 1.8$  vs.  PC: $\gamma = 2.2$

**Alpha x pixel resolution**: $r = 0.07$, $p < 0.01$

**Alpha x color depth**: $r = -0.18$, $p < 0.01$

# Experiment 3:
## Chart Size & Gridline Spacing

# Mechanical Turk:
Performance and Cost

# Turkers Overlap Across Studies



31% (51/186) Turkers participated in 2 or more

Only 7% (13) from Exp. 1A participated later

# *Samplers* and *Streakers*

# Quality with Qualification

**High quality results**: Only 0.75% of responses were rejected outliers.

Removing qualification tasks resulted in **over 10% unusable responses.**

**Verifiable answers** reduce gaming incentive and insincere responses.

# Standard HITs Frustrate Timing

Expected time per HIT: 10s
Observed time per HIT: 42s ($\mu$=54s, $\sigma$=41s)
→ **Timing data is not reliable.**

Strategies for Fine-Grained Timing
- *Macro-Task* (batch of micro-tasks)
- *Ready-Set-Go* HIT interface
  → Successful in subsequent studies.

# HIT Completion Rates



Orange ≥ 4¢   Blue = 2¢

Raise reward → faster results; ≅ quality

# Crowdsourcing Reduces Costs

**6x** cost savings (vs. $15/subject lab rate)

**9x** savings possible (using $0.02 rewards)

Study time drops **from 2 weeks to 1-3 days**

→ **Crowdsourcing provides up to an order of magnitude $$ and time savings**

→ With constant cost, it enables **more studies, more variables, more subjects**

# Future Work

**Multiple methods studies**: how to best balance the laboratory with online crowdsourcing?

**Better tools** for crowdsourced experimentation. Facilitate experimental control and adaptation.

**Community resources for evaluation**: share "market" data, share experimental designs, facilitate replication and meta-analysis.

**Extend crowdsourcing methods**
to an even greater diversity
of experimental designs.

# Color Naming Experiment

## Instructions

In each task, enter a specific color name that you believe best describes the color shown in the center rectangle. Use as many words as you need. For example, specific names might range from "dark red" to "crimson" to "scarlet". Next, from the provided list of basic color names, select the name that you believe best matches the center rectangle's color.

0

**What is the most specific (exact) color name you would use to name this color?** (required)

**What is the most general (basic) color name you would use to name this color?** (required)

Select one ▾

# Color Naming Experiment

## Instructions

In each task, enter a specific color name that you believe best describes the color shown in the center rectangle. Use as many words as you need. For example, specific names might range from "dark red" to "crimson" to "scarlet". Next, from the provided list of basic color names, select the name that you believe best matches the center rectangle's color.

0

**What is the most specific (exact) color name you would use to name this color?** (required)

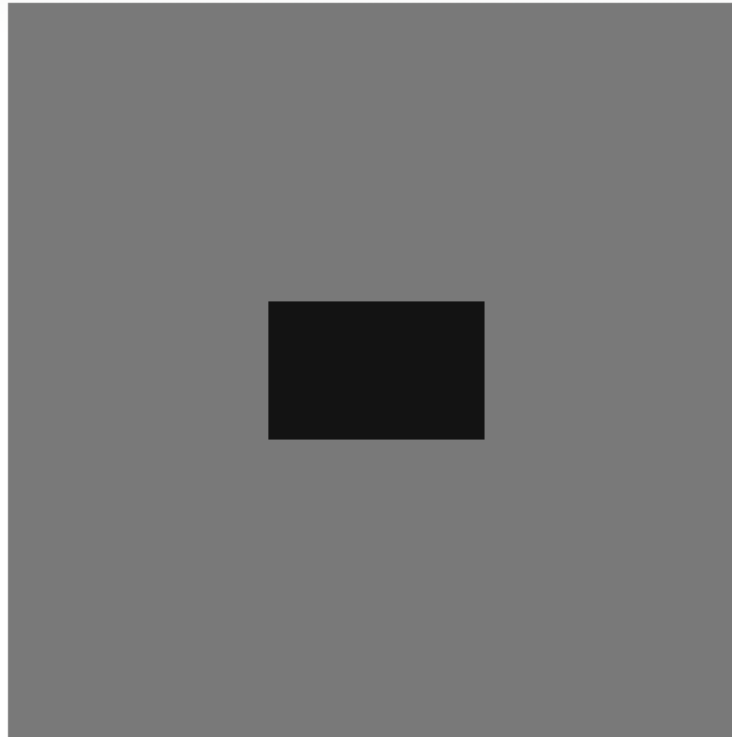**What is the most general (basic) color name you would use to name this color?** (required)

Select one

# Experiment zur Benennung von Farben

## Instructions Hide

Anleitung: So funktioniert diese Umfrage

Bitte geben Sie bei jeder der folgenden Aufgaben den Farbnamen in das mittlere Feld ein, die Ihrer Meinung nach die angezeigte Farbe am besten beschreibt. Sie konnen die Farbe dabei so frei oder spezifisch beschreiben wie Sie mochten, z.B. "dunkelrot", "purpurrot", "scharlachrot". Nun wahlen Sie aus der Liste der Grundfarben den Namen aus, der Ihrer Meinung nach die Farbe des mittleren Rechtecks am besten beschreibt.
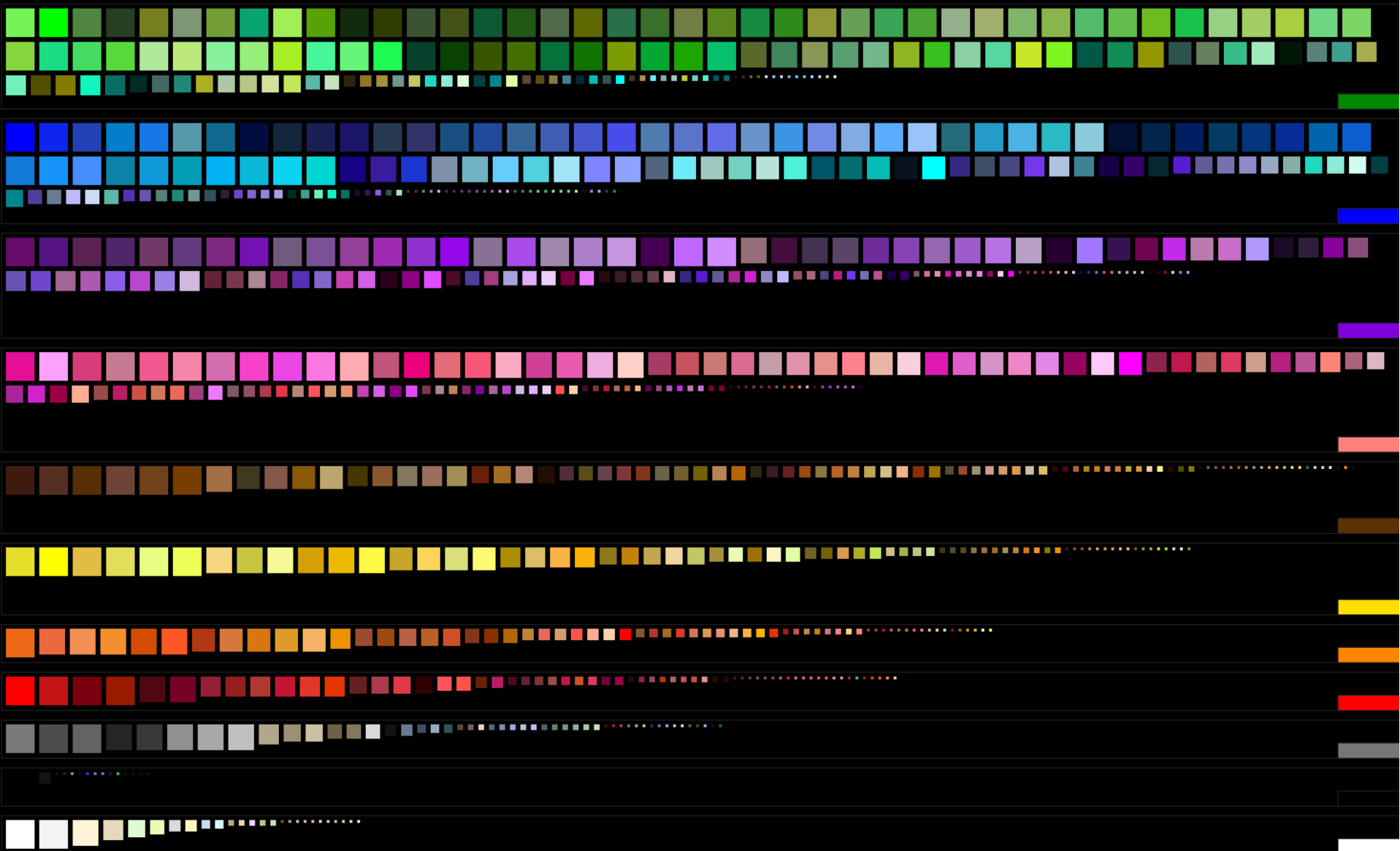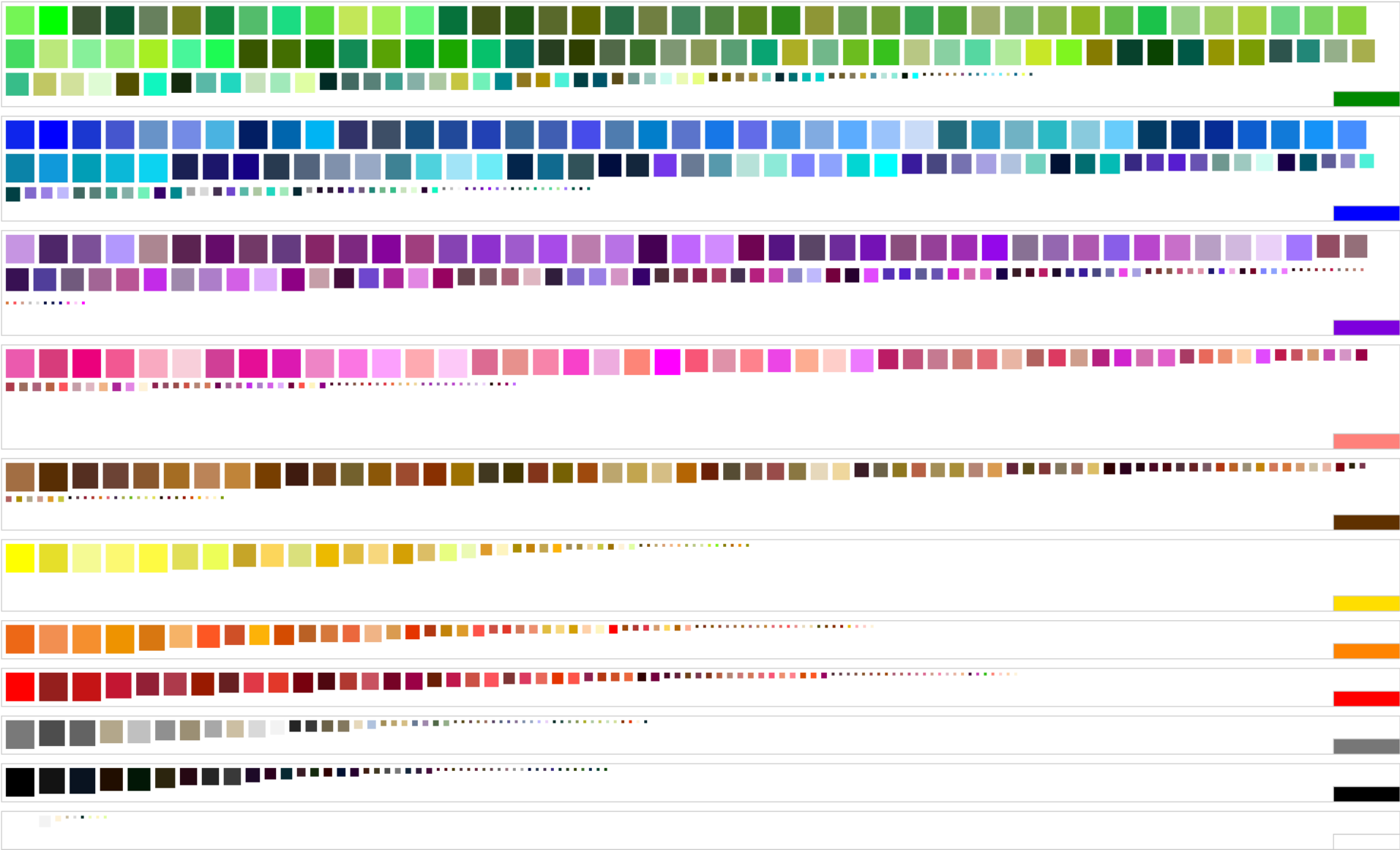
0

Welchen spezifischen Namen wurden Sie dieser Farbe am ehesten zuordnen? (required)

Welchen allgemeinen Namen wurden Sie dieser Farbe am ehesten zuordnen? (required)

Visualizing common color terms
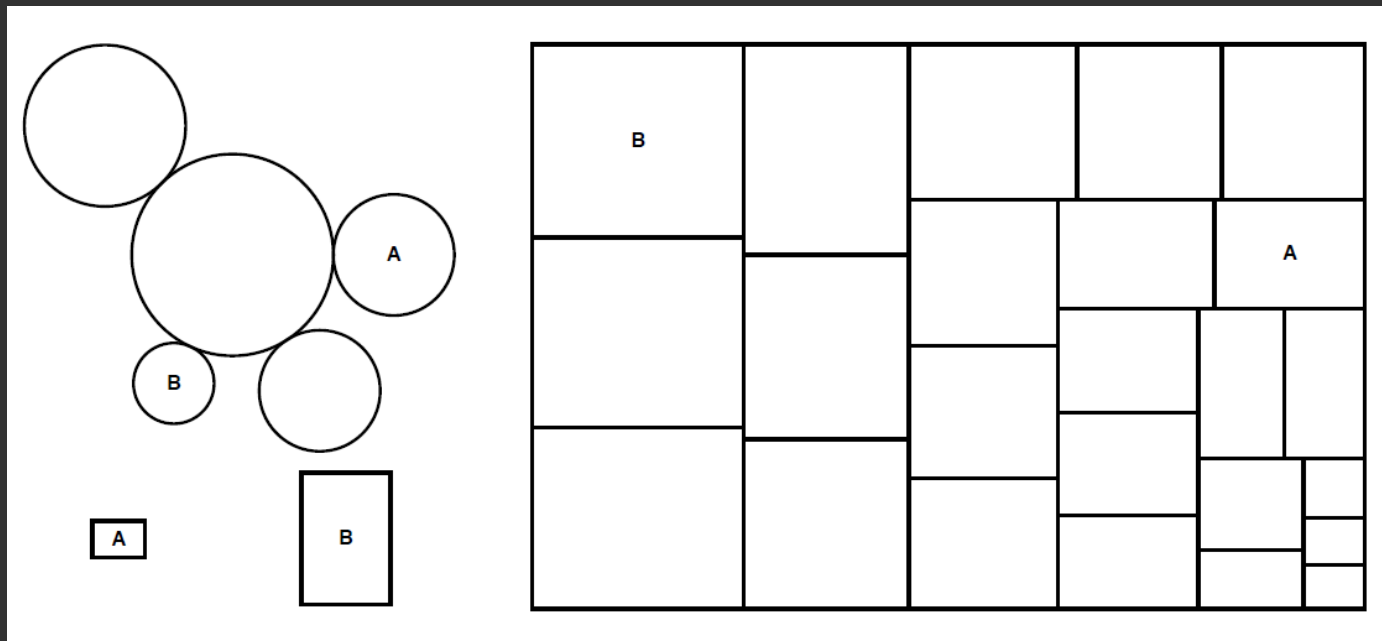
Visualizing common color terms

Your Jobs

# Your Jobs   Show active jobs   Show completed jobs

**Create New Job**

⭐ 11920   **Experience d'appellation de couleurs** (tag)
0 judgments, 454 units, 10 golds, created on **May 23, 2010**
Running ▼

⭐ 11919   **Experience d'appellation de couleurs** (tag)
192 judgments, 454 units, 10 golds, created on **May 23, 2010**
Running ▼

⭐ 11918   **Experience d'appellation de couleurs** (tag)
348 judgments, 454 units, 10 golds, created on **May 23, 2010**
Running ▼

⭐ 10458   **Experiment zur Benennung von Farben** (tag)
252 judgments, 454 units, 10 golds, created on **May 10, 2010**
Running ▼

⭐ 10452   **Experiment zur Benennung von Farben** (tag)
144 judgments, 454 units, 10 golds, created on **May 10, 2010**
Running ▼

⭐ 10418   **Experiment zur Benennung von Farben** (tag)
624 judgments, 454 units, 10 golds, created on **May 10, 2010**
Running ▼

⭐ 6829   **Color Naming Experiment** (tag)
11,256 judgments, 454 units, 10 golds, created on **Mar 04, 2010**
Finished ▼

⭐ 6828   **Color Naming Experiment** (tag)
10,788 judgments, 454 units, 10 golds, created on **Mar 04, 2010**
Finished ▼

⭐ 6827   **Color Naming Experiment** (tag)
11,304 judgments, 454 units, 10 golds, created on **Mar 04, 2010**
Finished ▼

# Crowdsourcing Graphical Perception
## Using Mechanical Turk to Assess Visualization Design



Jeffrey Heer & Michael Bostock
http://hci.stanford.edu/jheer