

Administrivia

Visualization Critique

What is this a visualization of?

What visual encodings
(mappings from data -> channel)
are used?

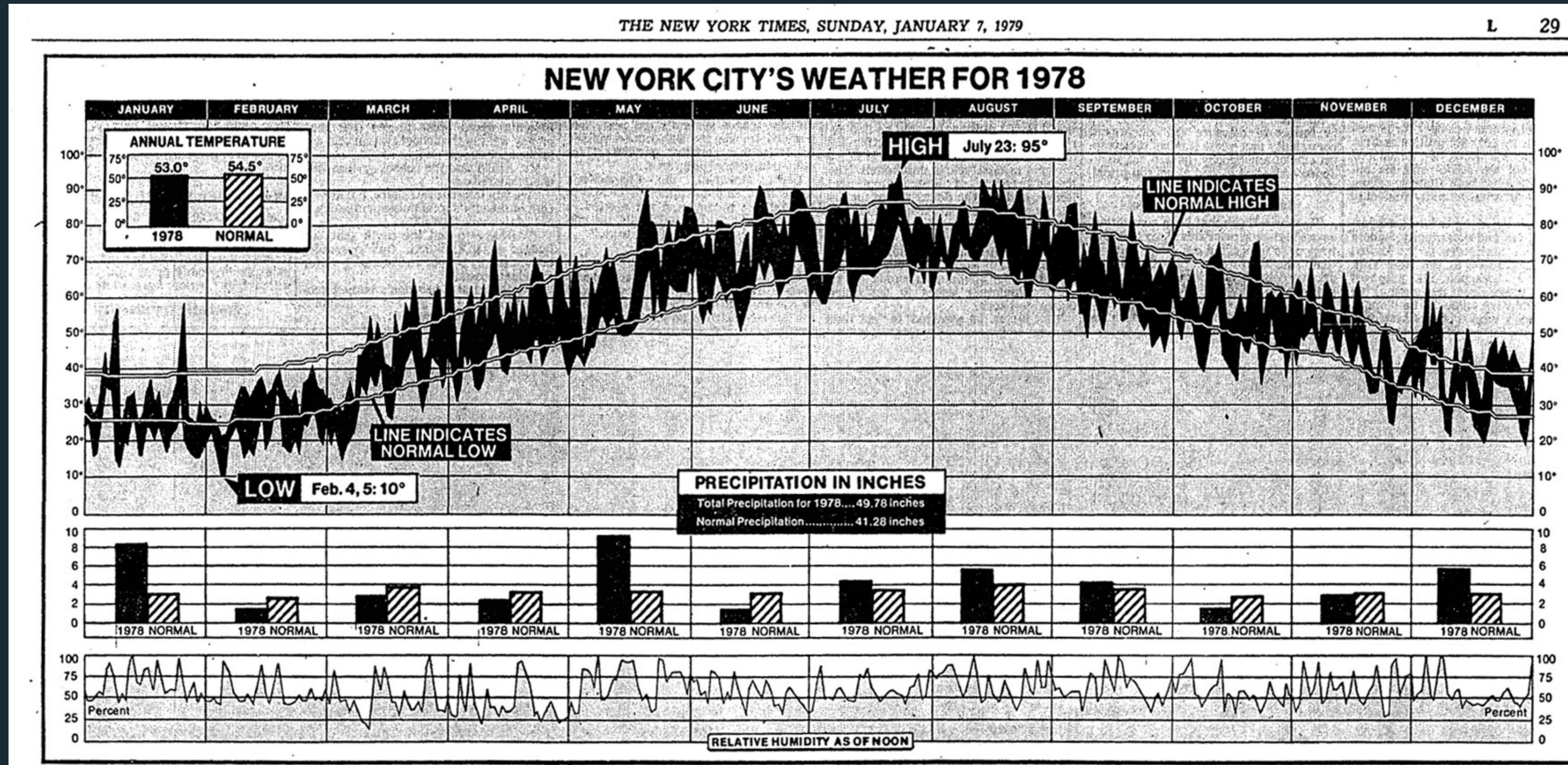
To help structure your critique:

- > "I like..."
- > "I wish..."
- > "What if...?"

Think (~ 3 mins).

Pair (~7 mins).

Share.

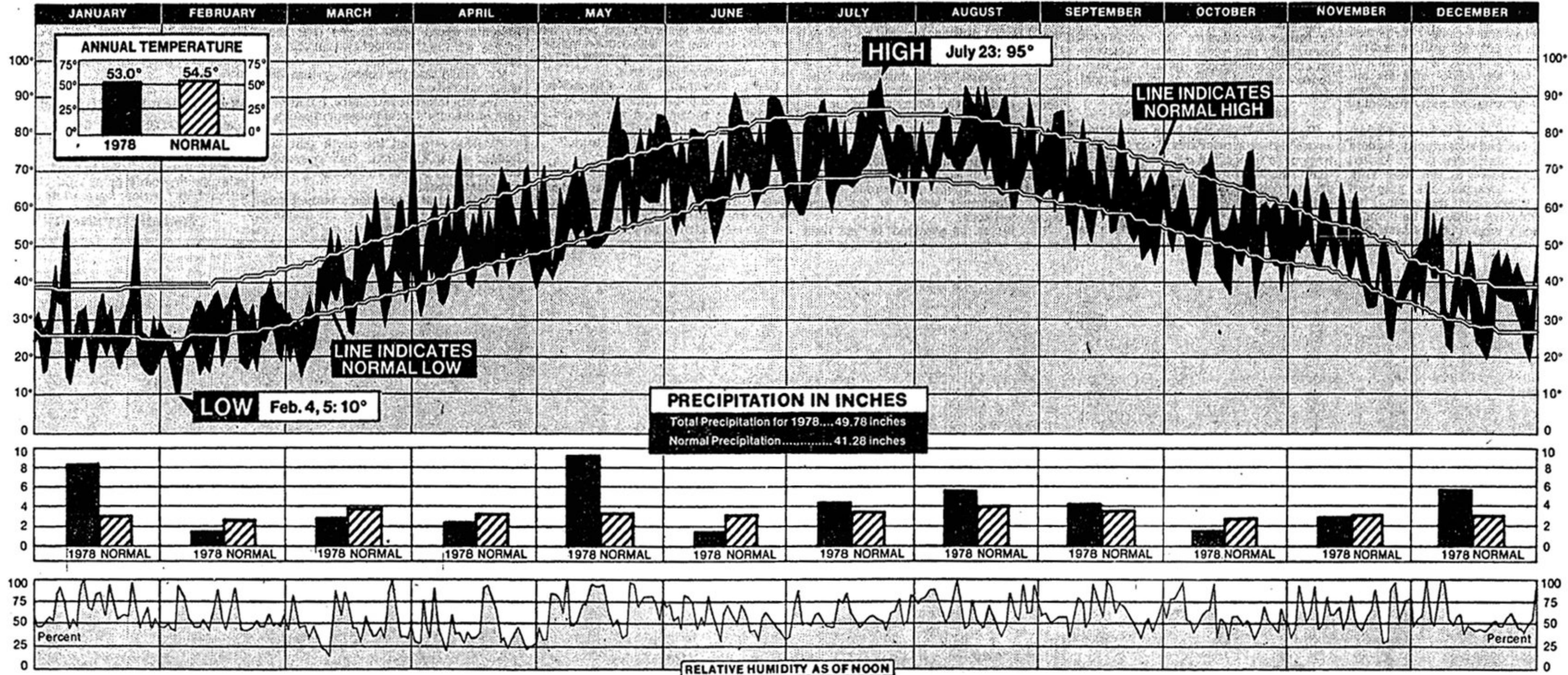


Visualization Critique

THE NEW YORK TIMES, SUNDAY, JANUARY 7, 1979

L 29

NEW YORK CITY'S WEATHER FOR 1978

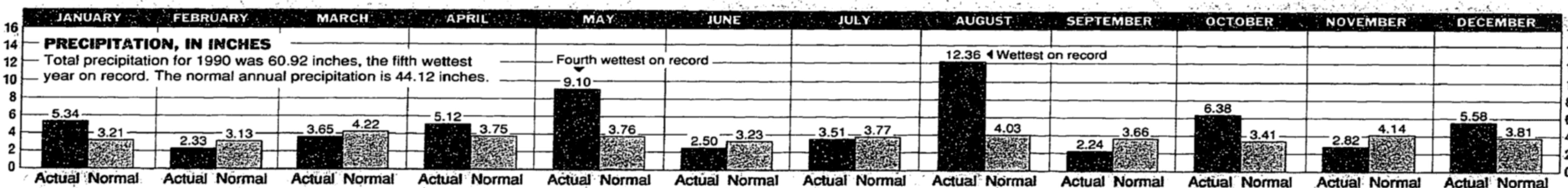
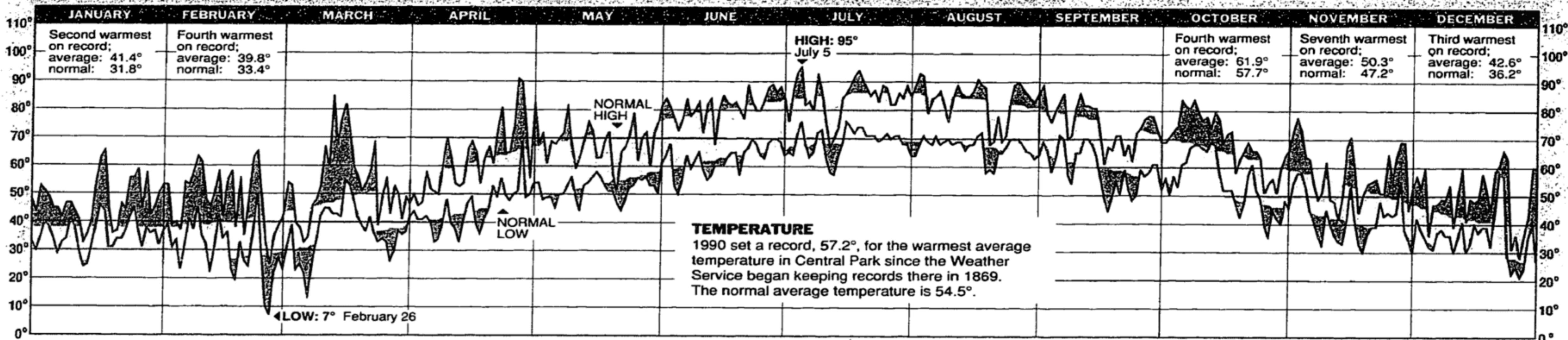


Visualization Critique

THE NEW YORK TIMES WEATHER SUNDAY, JANUARY 6, 1991

L+ 23

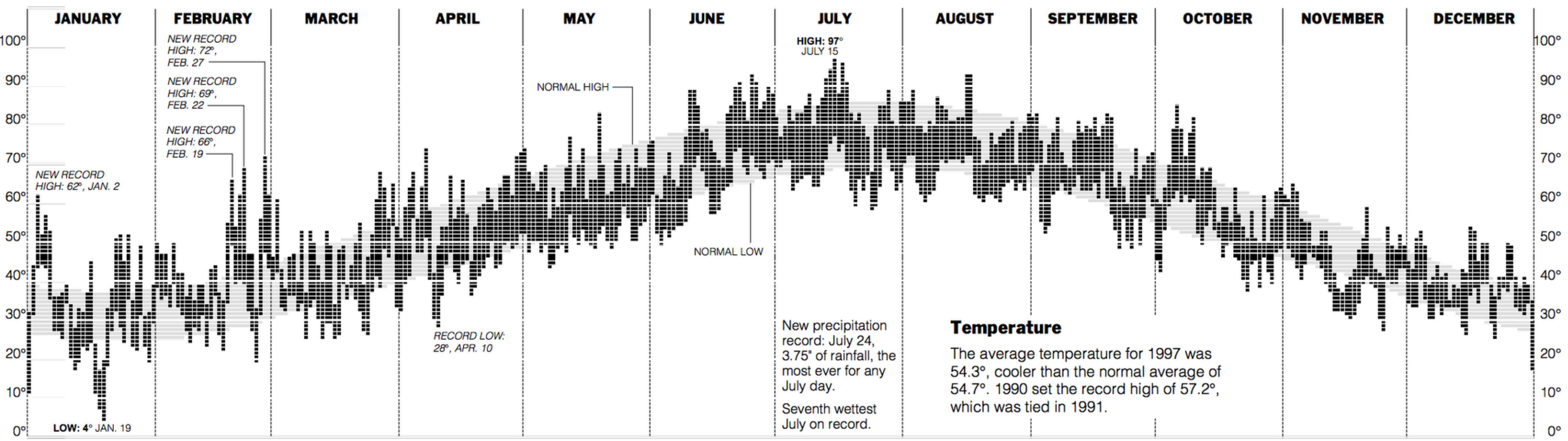
New York's Weather for 1990



Source: National Weather Service

Visualization Critique

New York City's Weather for 1997



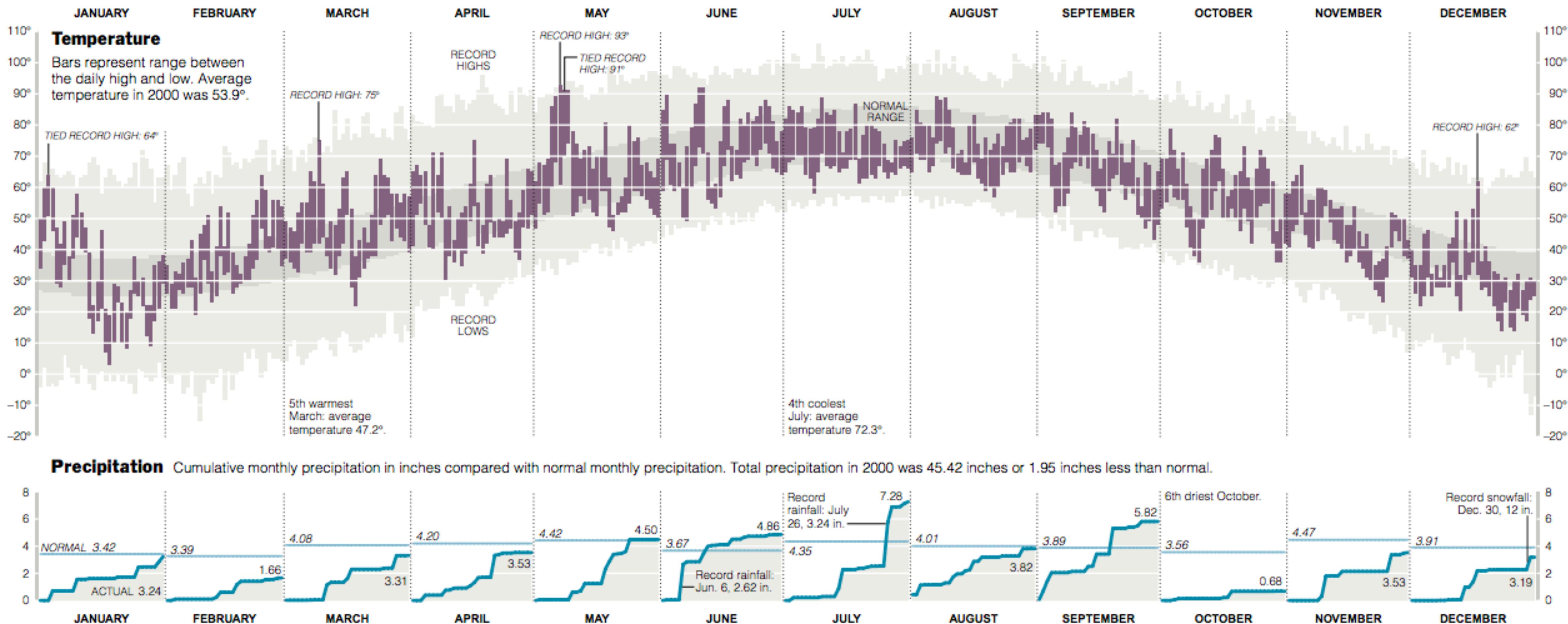
Precipitation, in Inches

Total precipitation for 1997 was 43.93 inches. The normal annual precipitation is 47.25 inches.

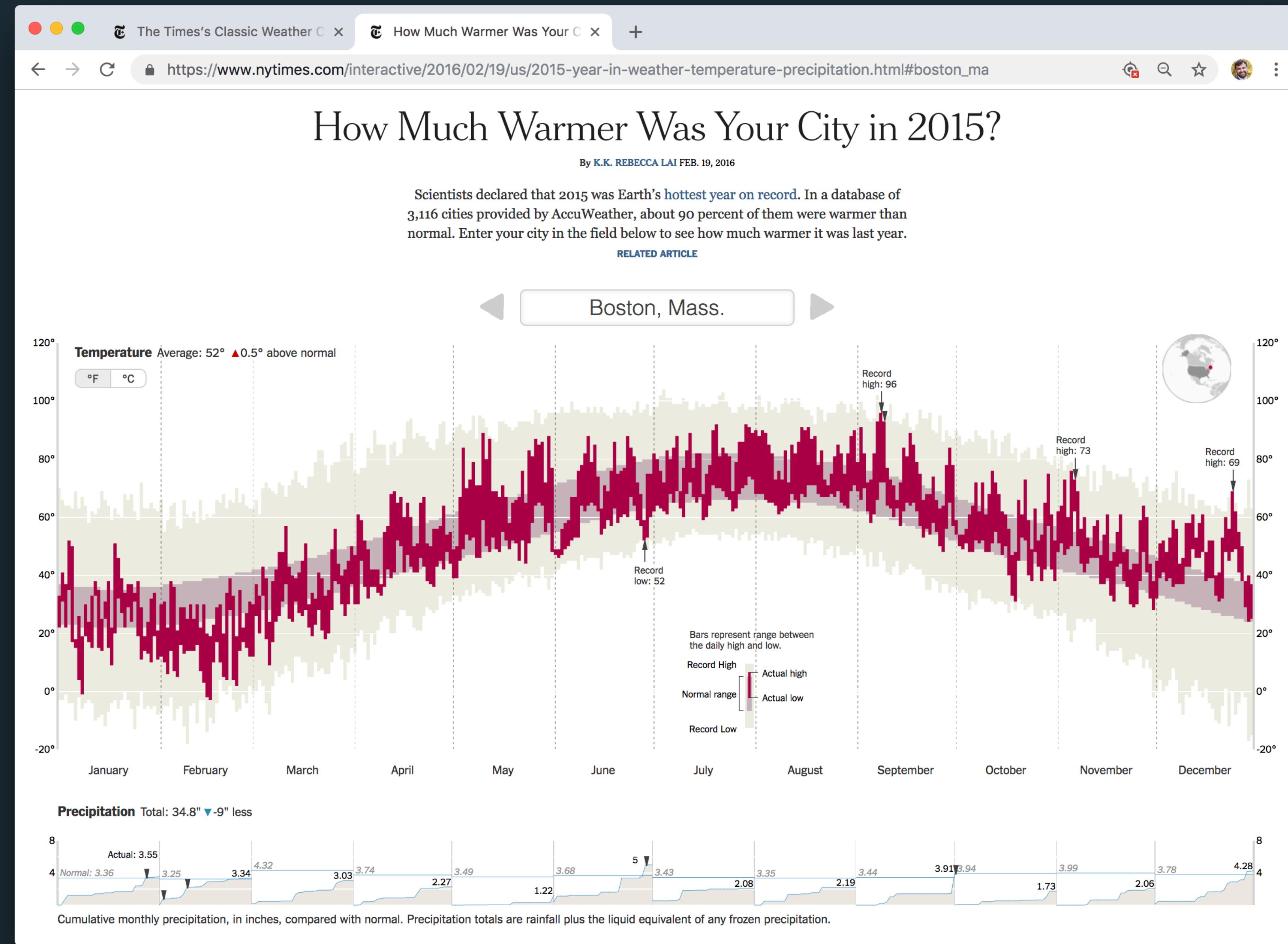


Visualization Critique

New York City's Weather in 2000



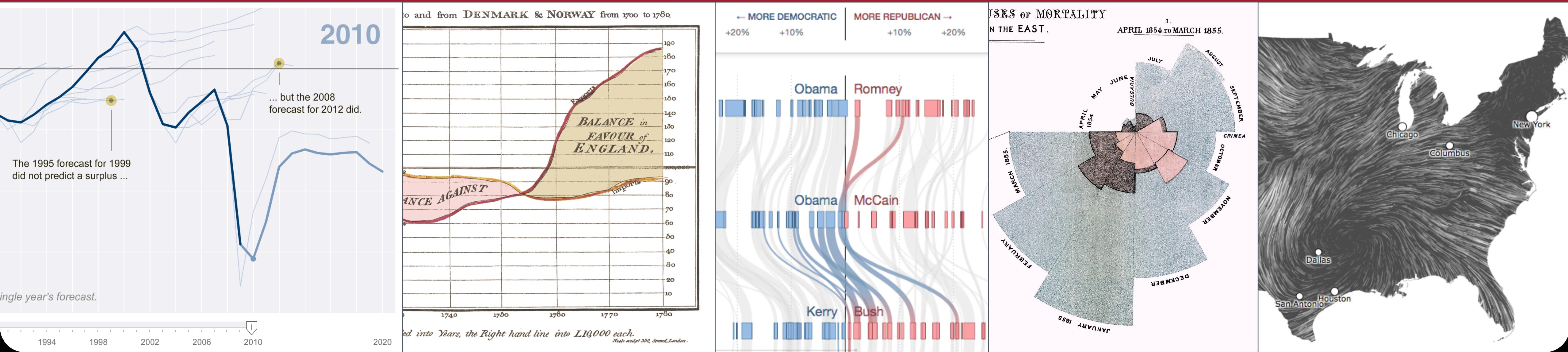
Visualization Critique



6.894: Interactive Data Visualization

Exploratory Data Analysis

Arvind Satyanarayan



Mapping or Visual Encoding

Data → Visual

Physical Data Types

int, float, string

Conceptual Data Types

temperature, location

Attribute Types

nominal, ordinal,

quantitative

(ratio or interval)

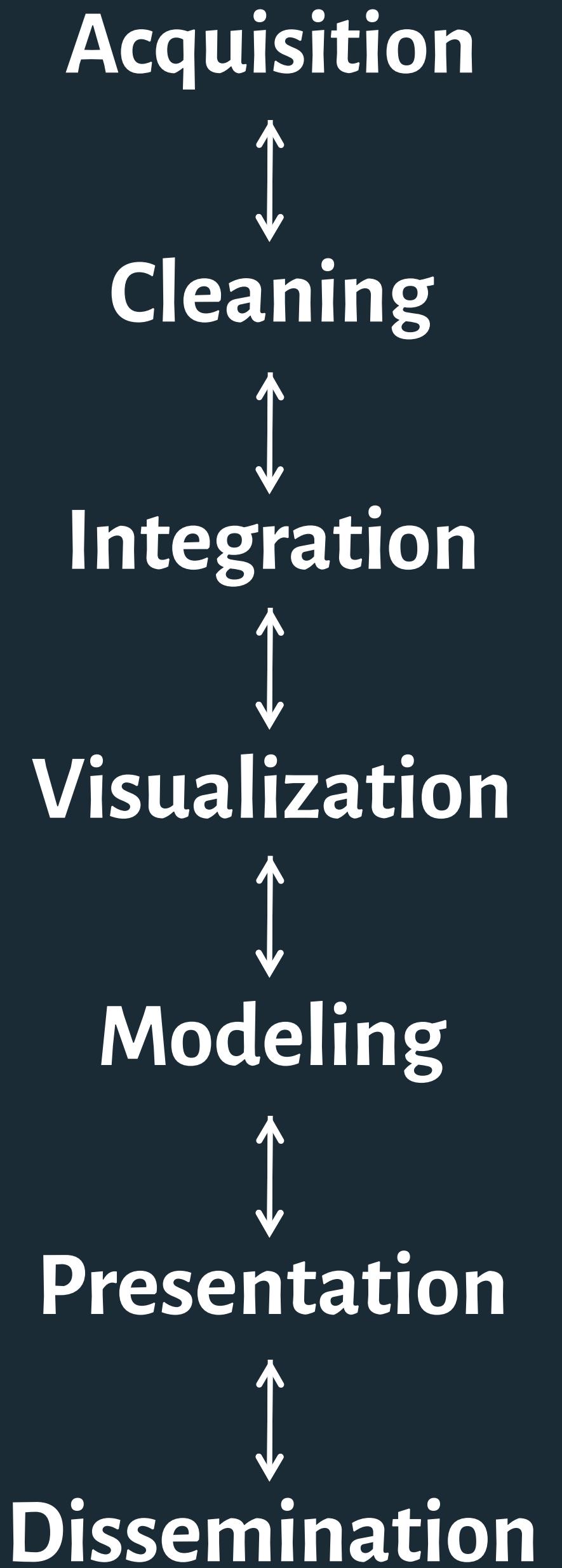
Graphical **Marks**

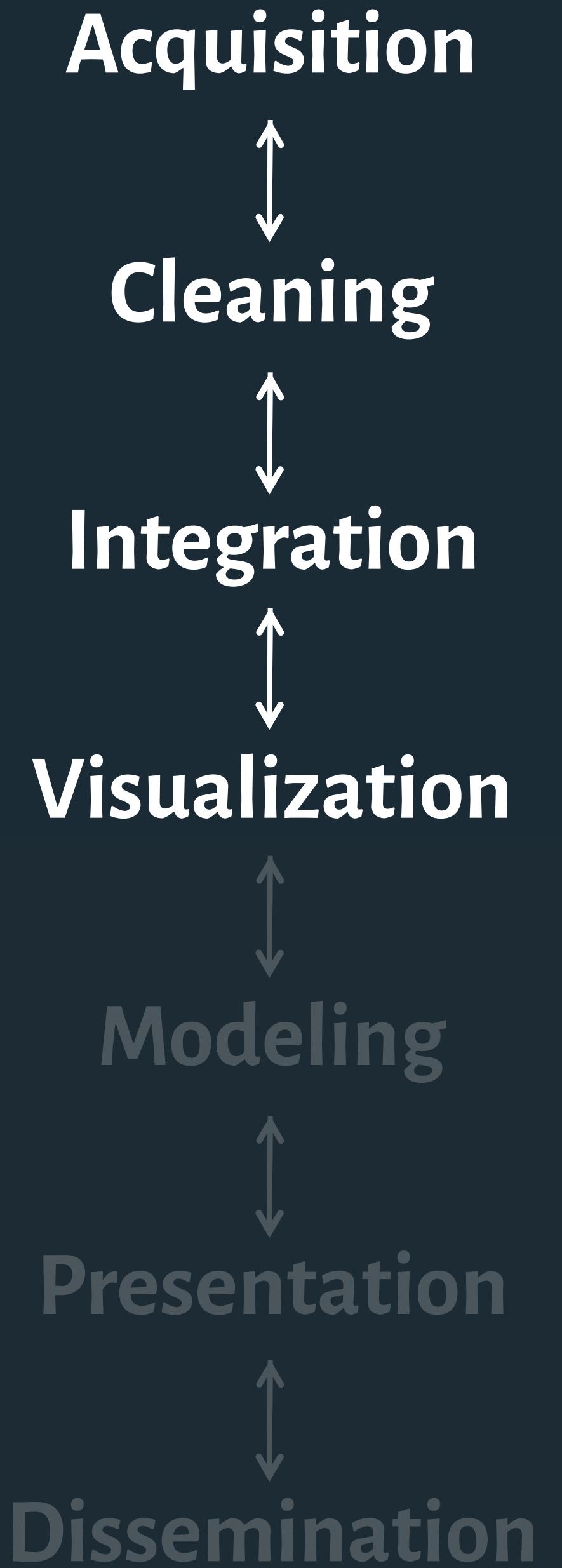
rect, line, point, area

Visual **Channels**

x, y, color, opacity

Data Visualization





Activity: What data quality issues may arise?

Think (~ 2 mins).

Pair (~5 mins).

Share.

Data Skepticism

"The wave of bullshit data is rising, and now it's our turn to figure out how not to get swept away."

– Jacob Harris

Co-founder of NYT Interactive Newsroom Tech Team

Nieman Lab, Dec 2014



Data Provenance

Who collected or produced it?

What was their intent?

Is it a reputable source?

What are the motives of the
data producer (are they an
advocate or lobbyist?)

Data Provenance

Who collected or produced it?

What was their intent?

Is it a reputable source?

What are the motives of the data producer (are they can advocate or lobbyist?)

Rats! D.C. calls pest company about rodents more often than New York

By Julie Zauzmer
October 13, 2014

Pest control company Orkin released a list on Monday of what it termed the “top 20 rattiest cities” — the cities where the company performed the most rodent treatments in 2013.

The Washington region came in third place, behind only Chicago and Los Angeles. (See Orkin’s complete list at the bottom of this post.)

Rikin S. Mehta, a senior deputy director of the city’s Department of Health, told the Post earlier this fall that Washington has “one of the most comprehensive rodent-control programs in the country,” geared toward understanding patterns of rodent behavior, not just exterminating animals when they pop up.

Mehta said that the city has a team of 14 rodent control specialists who are dispatched within two days every time a 3-1-1 call comes in about a rat sighting in a public place or a business. Residents can also report rat sightings to the Department of Health, and see where their neighbors have spotted rodents, by using the Post’s rat tracker.

Over the summer, Mehta’s branch of the Department of Health conducted a rat summit in each ward to ask residents about where they spot the creatures. The results will be published by spring.

Data Provenance

Who collected or produced it?

What was their intent?

Is it a reputable source?

What are the motives of the data producer (are they can advocate or lobbyist?)

Global Terrorism Database <https://www.start.umd.edu/gtd/>

GTD GLOBAL TERRORISM DATABASE

Search the Database

I'm a New User [ADVANCED SEARCH](#)

Browse by: [Go](#)

SEARCH

Information on more than 180,000 Terrorist Attacks

The Global Terrorism Database (GTD) is an open-source database including information on terrorist events around the world from 1970 through 2017 (with annual updates planned for the future). Unlike many other event databases, the GTD includes systematic data on domestic as well as international terrorist incidents that have occurred during this time period and now includes more than 180,000 cases. [Learn more](#)

Read more about [Global Terrorism in 2017](#)

GTD DATA VISUALIZATIONS



The **GTD World Map: 45 Years of Terrorism** displays terrorist violence that occurred worldwide between 1970 and 2015.

[The GTD 2017 World Map is available here.](#)
[The GTD 2016 World Map is available here.](#)
[The GTD 2015 World Map is available here.](#)
[The GTD 2014 World Map is available here.](#)
[The GTD 2013 World Map is available here.](#)
[The GTD 2012 World Map is available here.](#)

THIS DATE IN TERRORISM

February 9

2015 Bantacan, Philippines

02/09/2015: An explosive device was discovered and safely defused in Purok 1A area, Bantacan village, Compostela Valley province, Philippines. No group claimed responsibility for the incident; however, sources attributed the attempted attack to the New People's Army (NPA).

[Learn more](#)

2015 Logo district, Nigeria

02/09/2015: Assailants attacked residents and buildings in Logo district, Benue state, Nigeria. This was one of 24 coordinated raids on villages and communities in this area on February 9, 2015. At least 18 people were killed across attacks. No group claimed responsibility for the incident; however, sources attributed the attack to Fulani militants.

[Learn more](#)

FEATURED

Message from the Global Terrorism Database Manager

For more than a decade, START has compiled and published the Global Terrorism Database (GTD) for use by scholars, analysts, journalists, security professionals, and policy makers. It has been our privilege to work closely with these user communities to continually improve the data and inform stakeholders.

Since 2012, the majority of the costs of collecting the GTD have been funded by the U.S. State Department, for the past year almost exclusively. Our contract with the State Department ended in May 2018 and, although we received only positive feedback from the Bureau of Counterterrorism and our 2018 data collection was well underway, we recently learned that we were not awarded a follow-on contract for base data collection.

At the moment, the loss of the State Department funding means two things: First, we do not currently have funding to complete collection of 2018 data, nor are we able to publish data beyond 2017.

[Continue Reading](#)

Data Provenance

When was the data collected?

Measurements can change over time.

Definitions/interpretations in quantification can change over time.

Is the data recent, and how much does that matter to the insight you wish to convey?

The Three-Year Plunge

To help gauge each city's overall crime level, the FBI tracks eight "index crimes." From 1993 to 2010, Chicago's annual total dropped by 47 percent. But from 2010 to 2013, it dropped a stunning 56 percent, or nearly 19 percent per year, according to data from the Chicago Police Department.

Graph Viewer

Graph Viewer

Roll-up by:

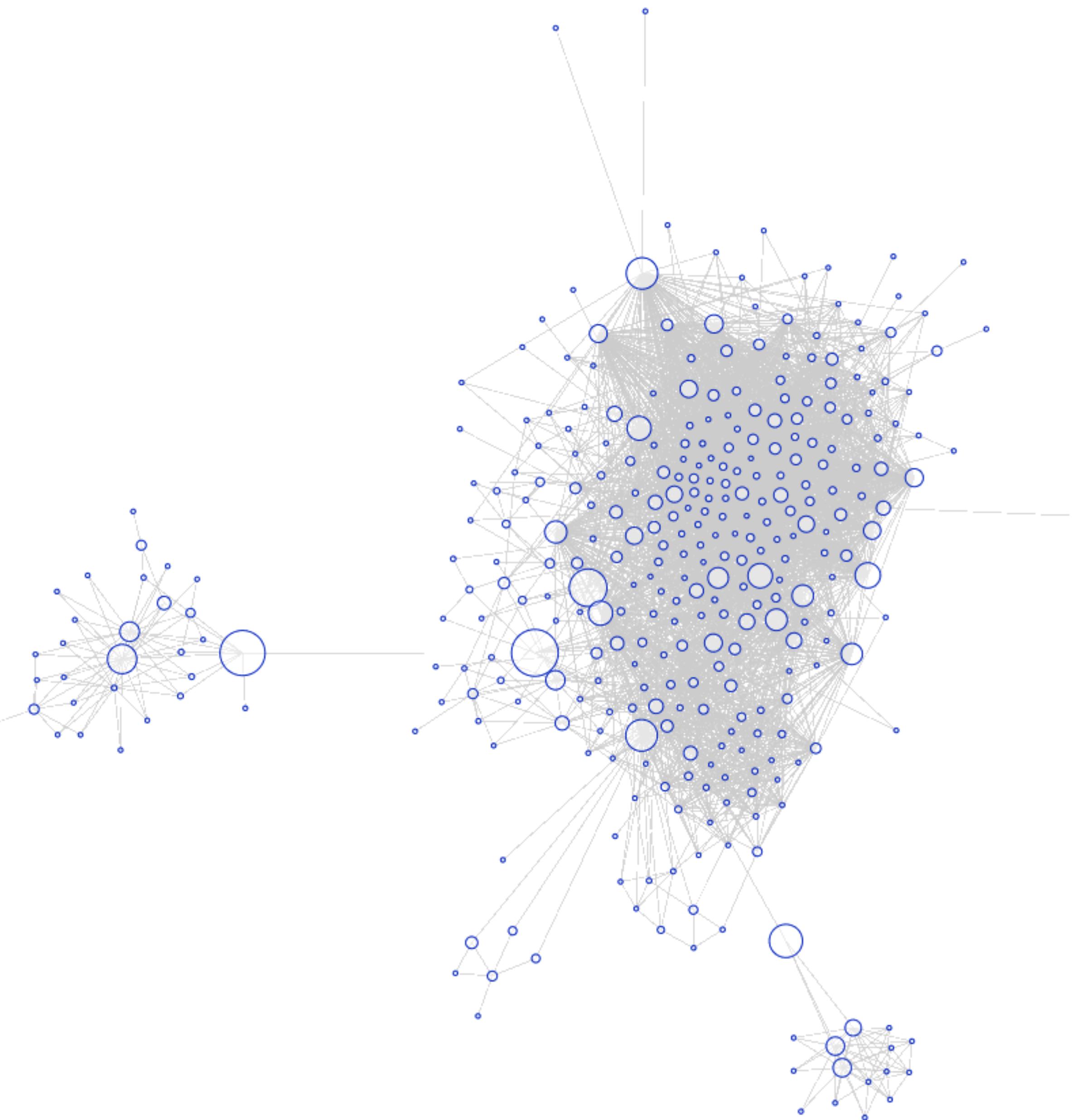
Visualization:

Sort by:

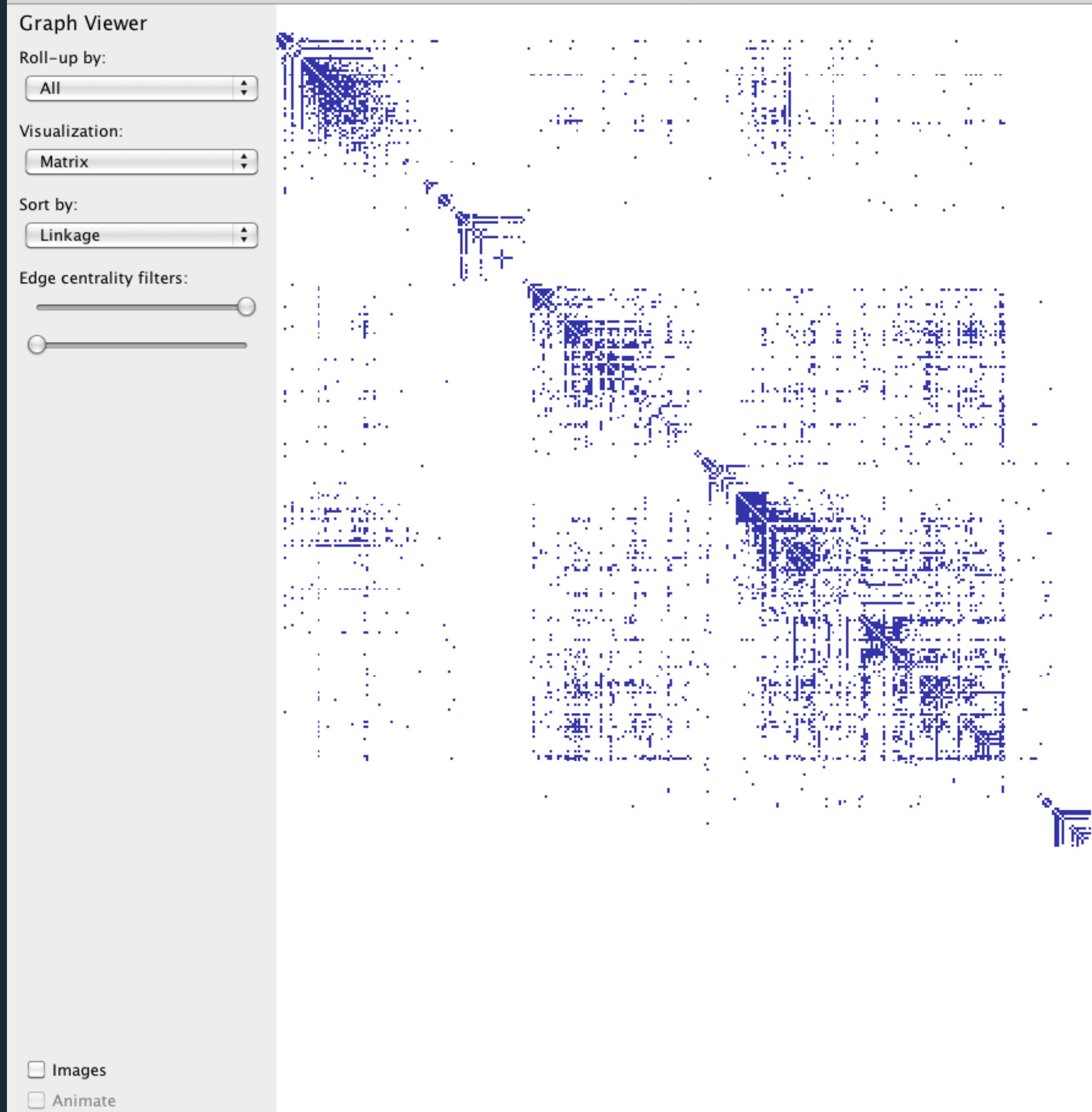
Edge centrality filters:

Images

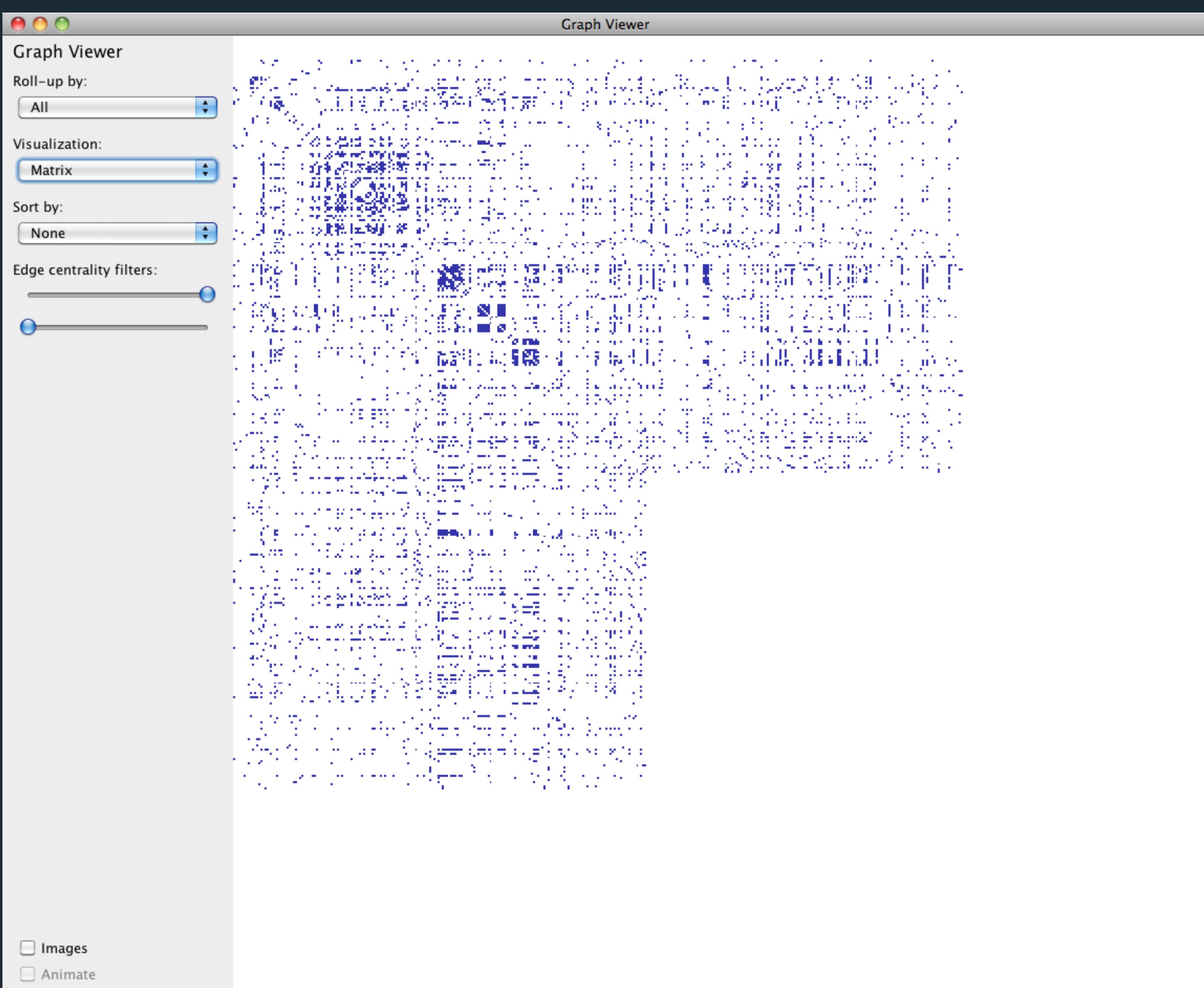
Animate



Graph Viewer



Missing Values



Berkeley

Cornell

Harvard

Harvard University

Stanford

Stanford University

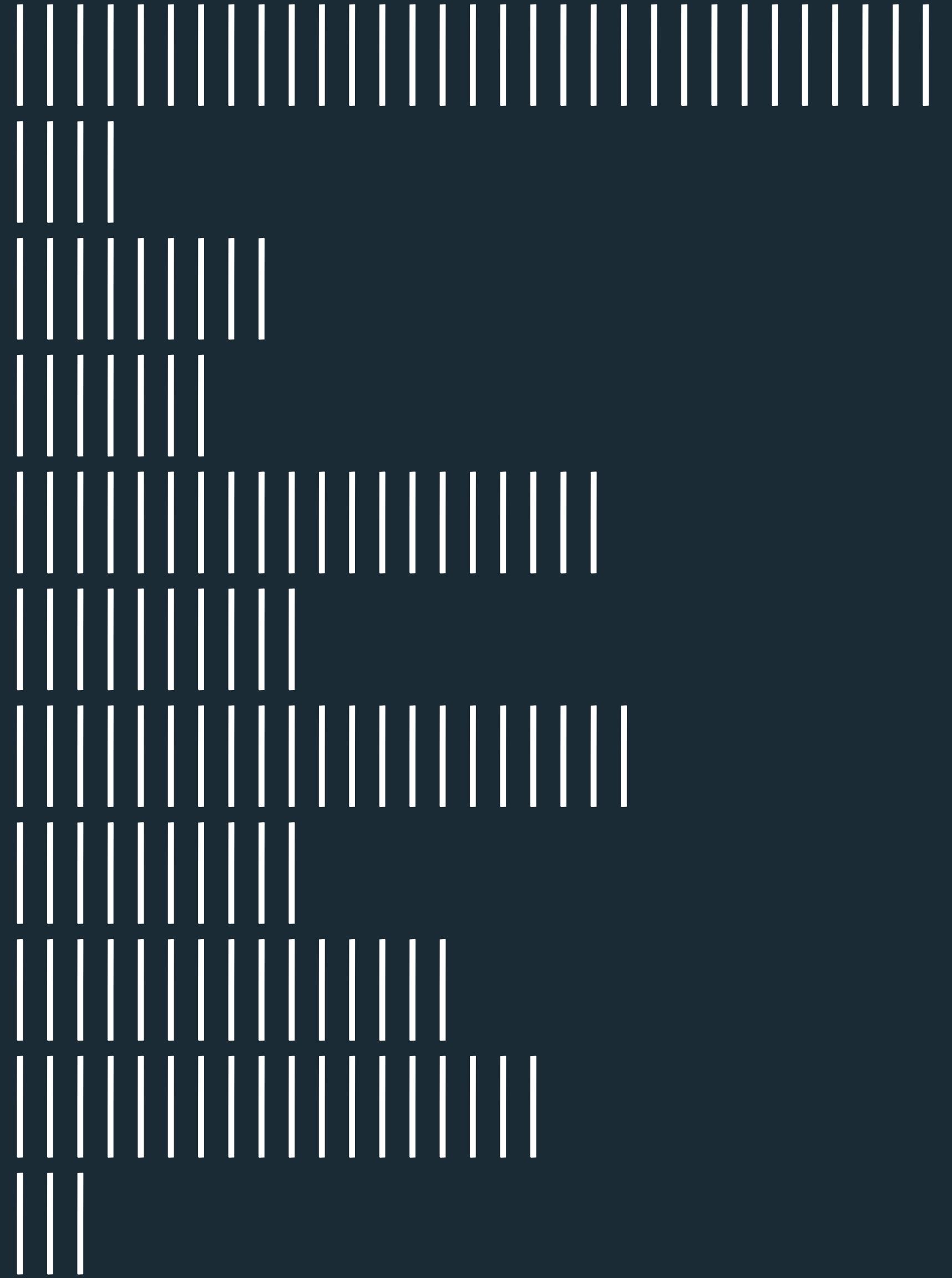
UC Berkeley

UC Davis

University of California at Berkeley

University of California, Berkeley

University of California, Davis



“The first sign that a visualization is good is that **it shows you a problem in your data**. Every successful visualization that I've been involved with has had this stage where you realize, "***Oh my God, this data is not what I thought it would be!***" So already, you've discovered something.”

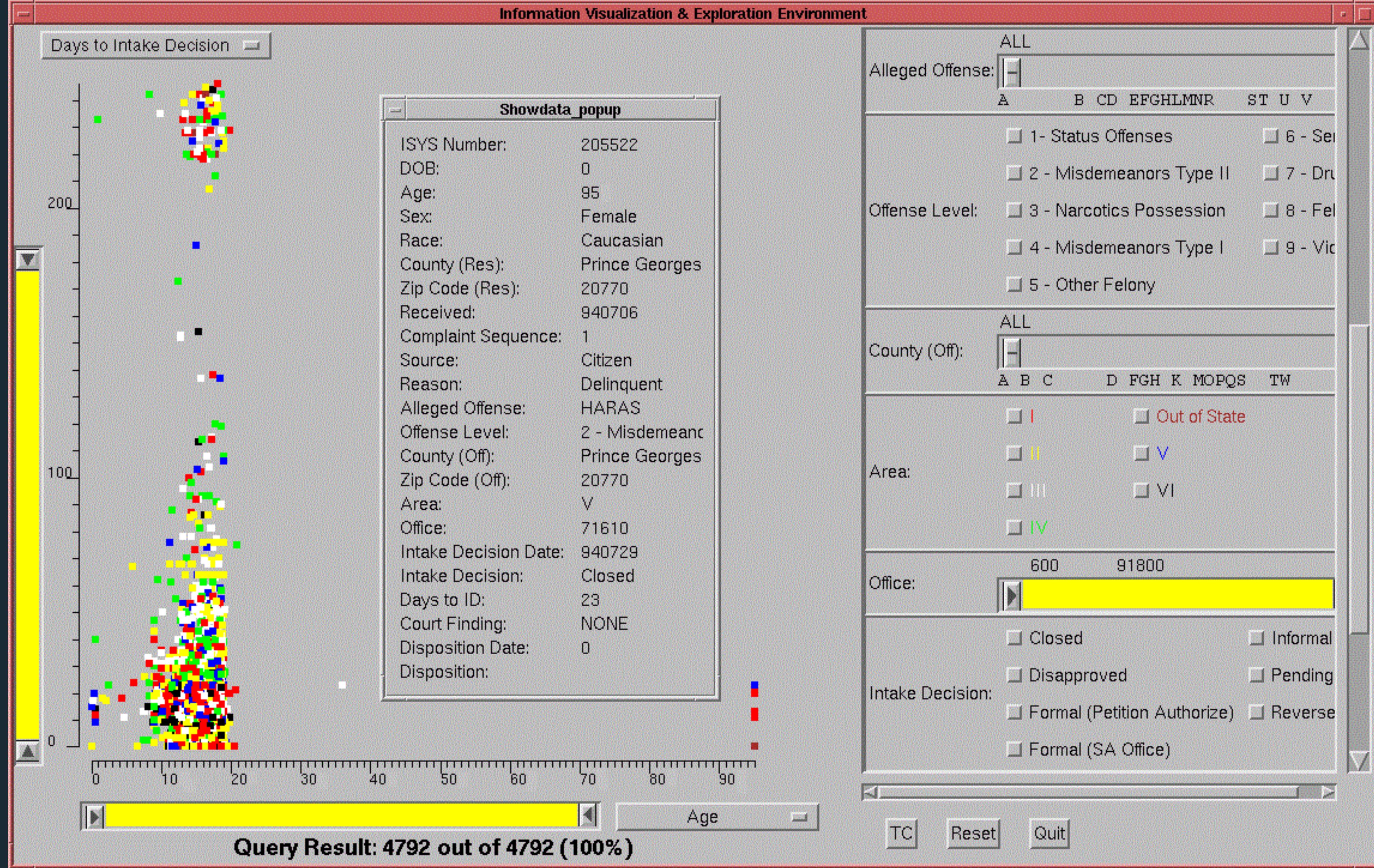
– Martin Wattenberg

Co-lead of Google's People + AI Initiative

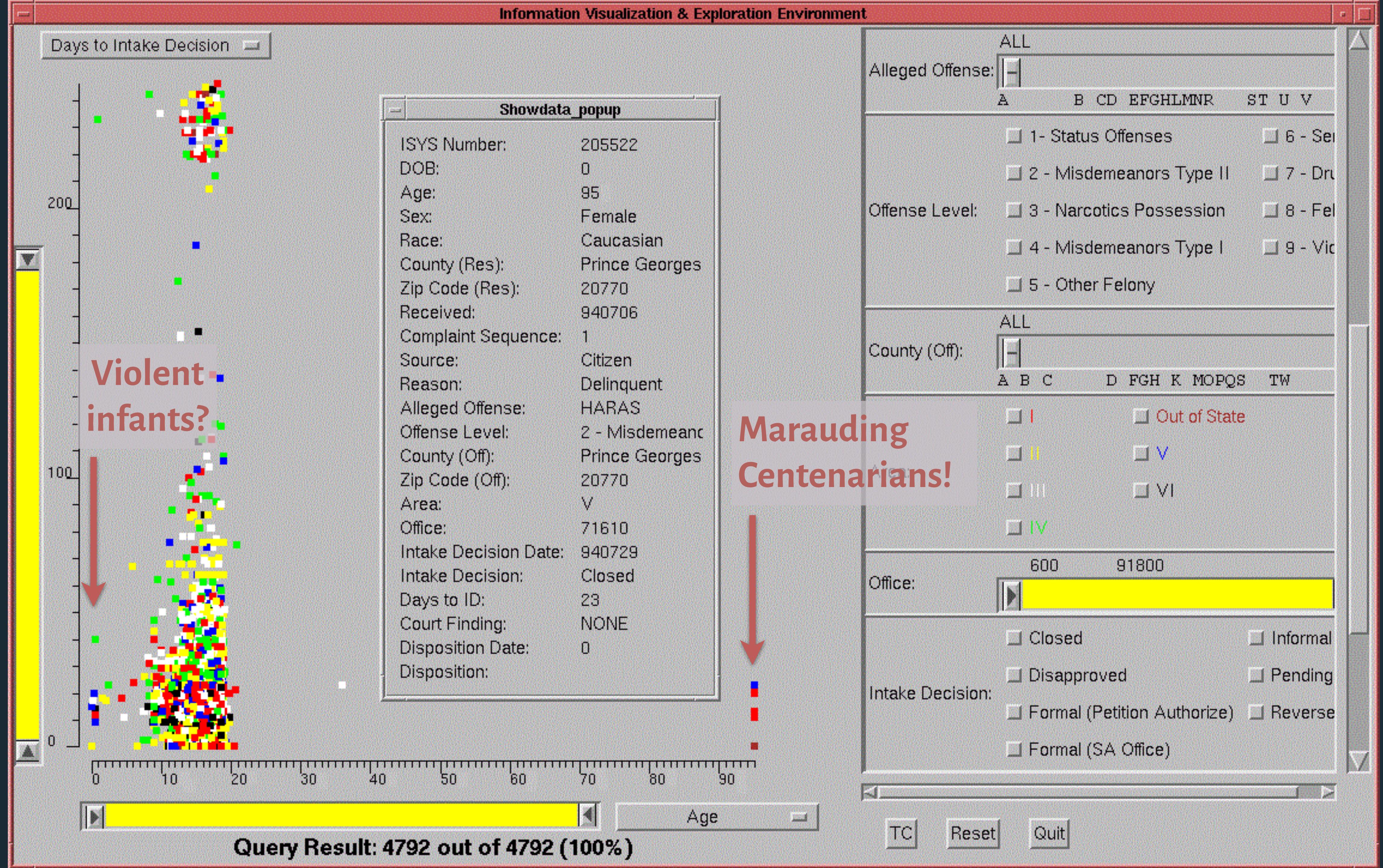
ACM Queue, Mar 2010



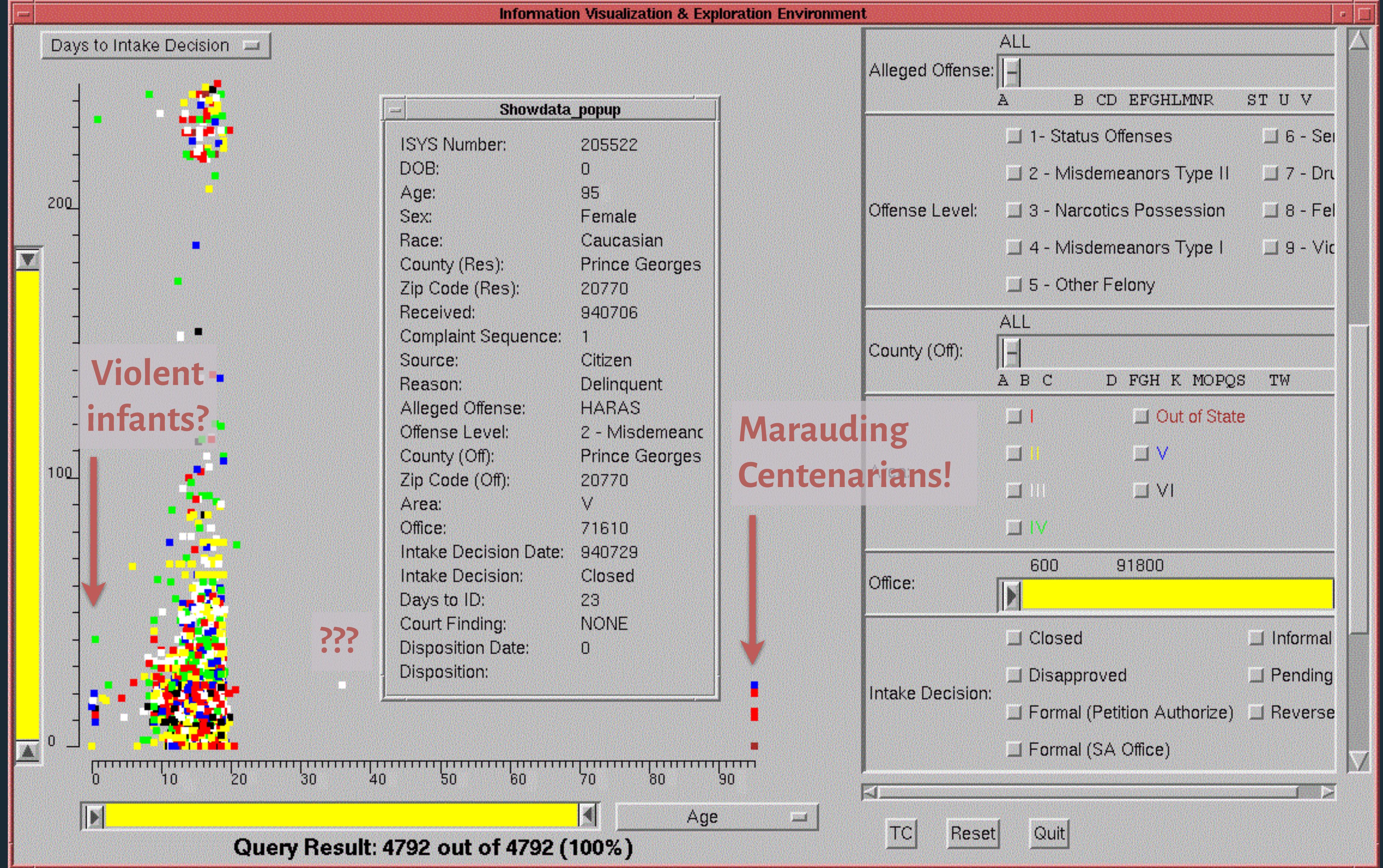
Ch15: Visualization Data Mining Models. Thearling et al.,
Information Visualization in Data Mining & Knowledge Discovery 2002.



Ch15: Visualization Data Mining Models. Thearling et al.,
Information Visualization in Data Mining & Knowledge Discovery 2002.



Ch15: Visualization Data Mining Models. Thearling et al.,
Information Visualization in Data Mining & Knowledge Discovery 2002.





“I spend more than half of my time integrating, cleansing and transforming data without doing any actual analysis. Most of the time I'm lucky if I get to do any “analysis” at all.”

– **Anonymous Data Scientist**

[Kandel et al. VAST 2012]



Big Data Borat

@BigDataBorat

Follow



In Data Science, 80% of time spent prepare data, 20% of time spent complain about need for prepare data.

6:47 PM - 26 Feb 2013

540 Retweets **343** Likes



12

540

343



Reported crime in Alabama

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	4525375	4029.3	987	2732.4	309.9
2005	4548327	3900	955.8	2656	289
2006	4599030	3937	968.9	2645.1	322.9
2007	4627851	3974.9	980.2	2687	307.7
2008	4661900	4081.9	1080.7	2712.6	288.6

Reported crime in Alaska

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	657755	3370.9	573.6	2456.7	340.6
2005	663253	3615	622.8	2601	391
2006	670053	3582	615.2	2588.5	378.3
2007	683478	3373.9	538.9	2480	355.1
2008	686293	2928.3	470.9	2219.9	237.5

Reported crime in Arizona

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	5739879	5073.3	991	3118.7	963.5
2005	5953007	4827	946.2	2958	922
2006	6166318	4741.6	953	2874.1	914.4
2007	6338755	4502.6	935.4	2780.5	786.7
2008	6500180	4087.3	894.2	2605.3	587.8

Reported crime in Arkansas

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	2750000	4033.1	1096.4	2699.7	237
2005	2775708	4068	1085.1	2720	262
2006	2810872	4021.6	1154.4	2596.7	270.4
2007	2834797	3945.5	1124.4	2574.6	246.5
2008	2855390	3843.7	1182.7	2433.4	227.6

Reported crime in California

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	35842038	3423.9	686.1	2033.1	704.8
2005	36154147	3321	692.9	1915	712
2006	36457549	3175.2	676.9	1831.5	666.8
2007	36553215	3032.6	648.4	1784.1	600.2
2008	36756666	2940.3	646.8	1769.8	523.8

Reported crime in Colorado

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	4601821	3918.5	717.3	2679.5	521.6

DataWrangler

Suggestions

Delete rows 8,10

Delete empty rows

Delete rows where Property_crime_rate is null

Delete rows where Year is null

Script

Export

► Split data repeatedly on newline into rows

► Split data repeatedly on ;

rows: 408 prev next

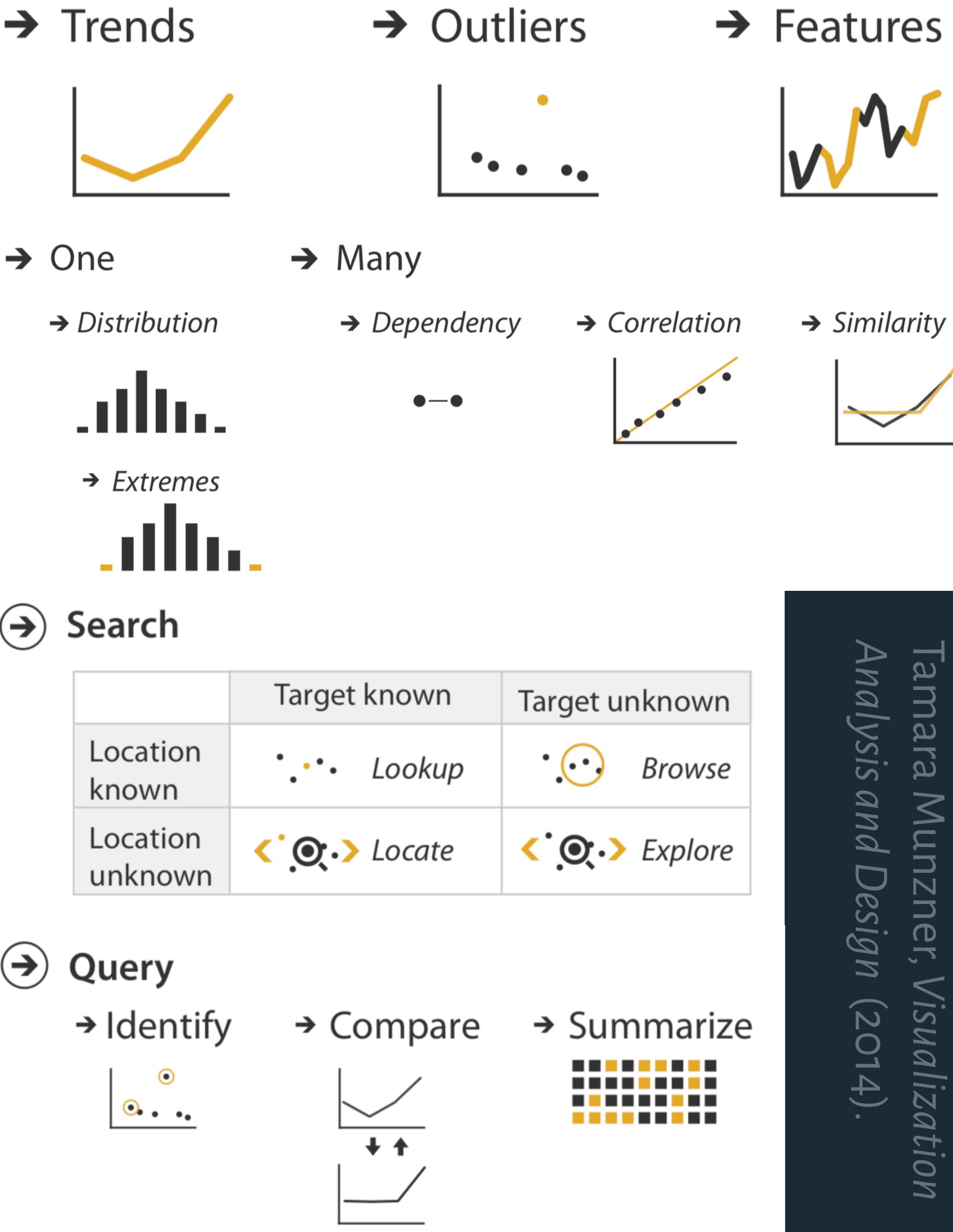
#	Year	#	Property_crime_rate
1	Reported crime in Alabama		
2			
3	2004		4029.3
4	2005		3900
5	2006		3937
6	2007		3974.9
7	2008		4081.9
8			
9	Reported crime in Alaska		
10			
11	2004		3370.9
12	2005		3615
13	2006		3582
14	2007		3373.9

Wrangler: Interactive Visual Specification of Data Transformation Scripts. Sean Kandel et al., ACM CHI 2011.

Exploratory Visual Analysis

Process

1. Construct graphics to address questions.
2. Inspect "answer" and ask new questions.
3. Iterate...



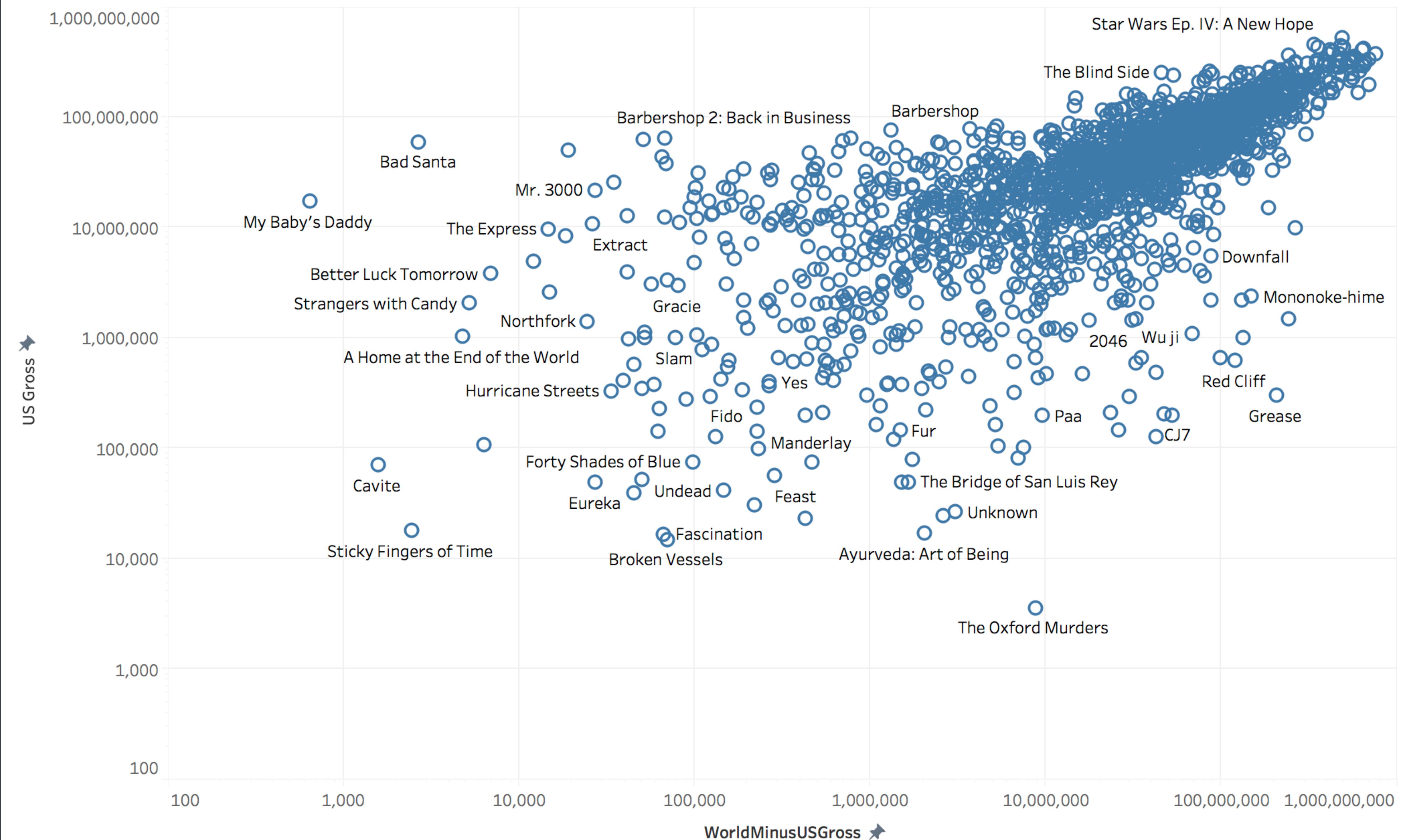
Analysis Example: Motion Pictures Data

Analysis Example: Motion Pictures Data

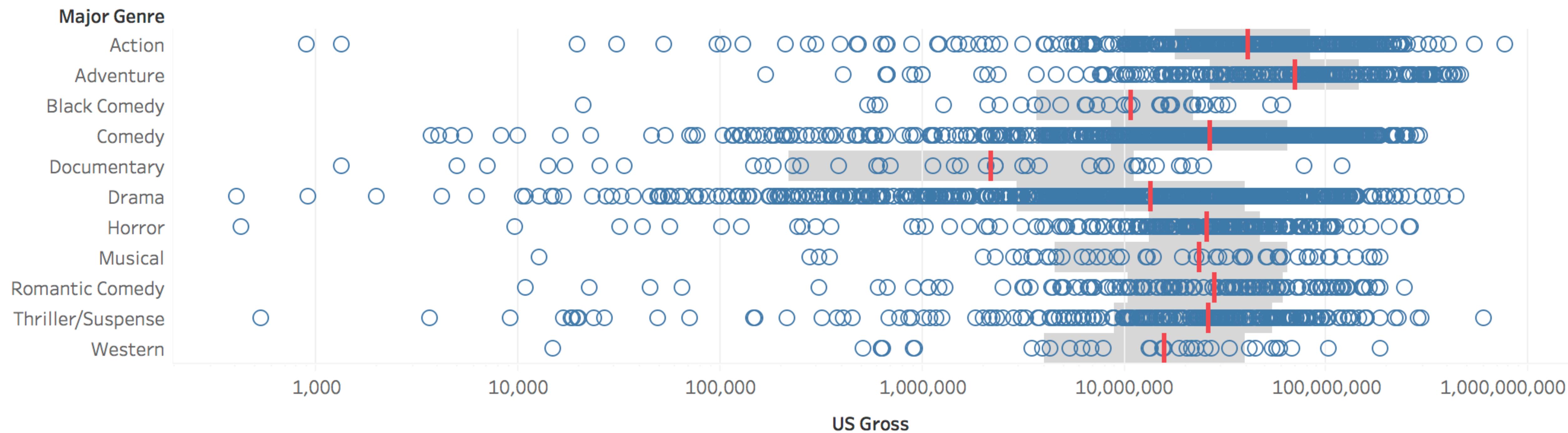
A sample of 3,201 movies collected in 2010.

Title	String (N)
IMDB Rating	Number (Q)
Rotten Tomatoes Rating	Number (Q)
Genre	String (N)
Release Date	Date (T)
US Gross	Number (Q)
Worldwide Gross	Number (Q)

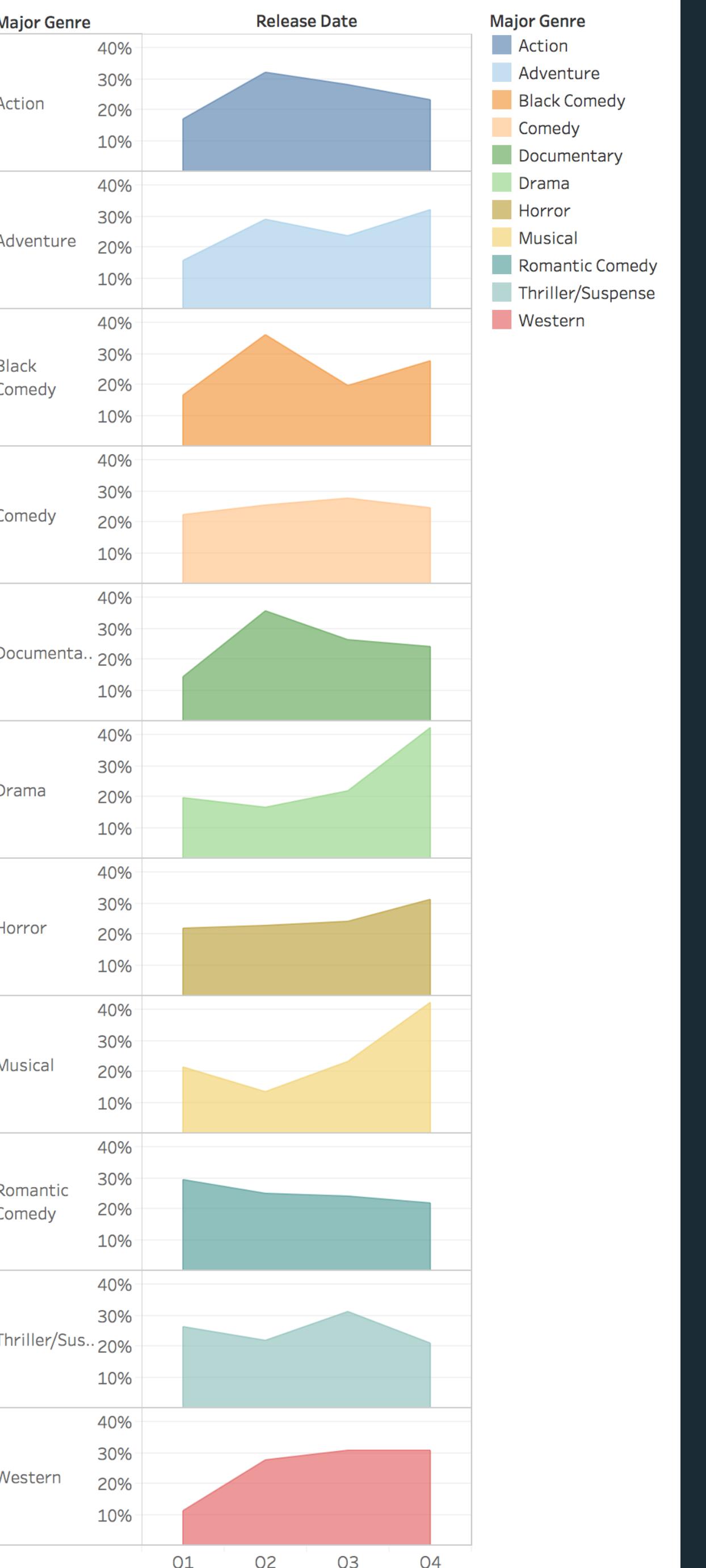
US vs WW



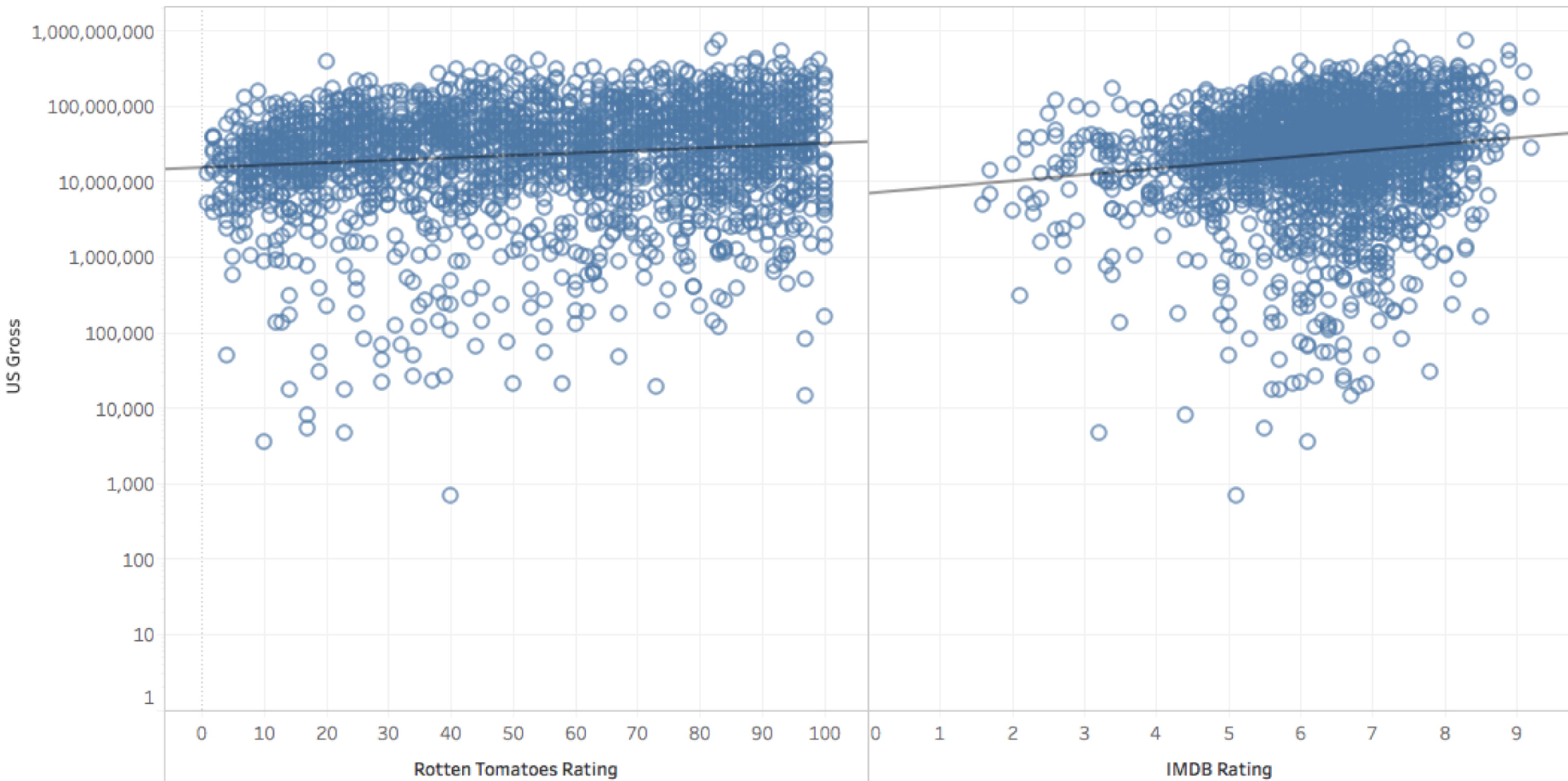
Distribution of US Gross by Genre



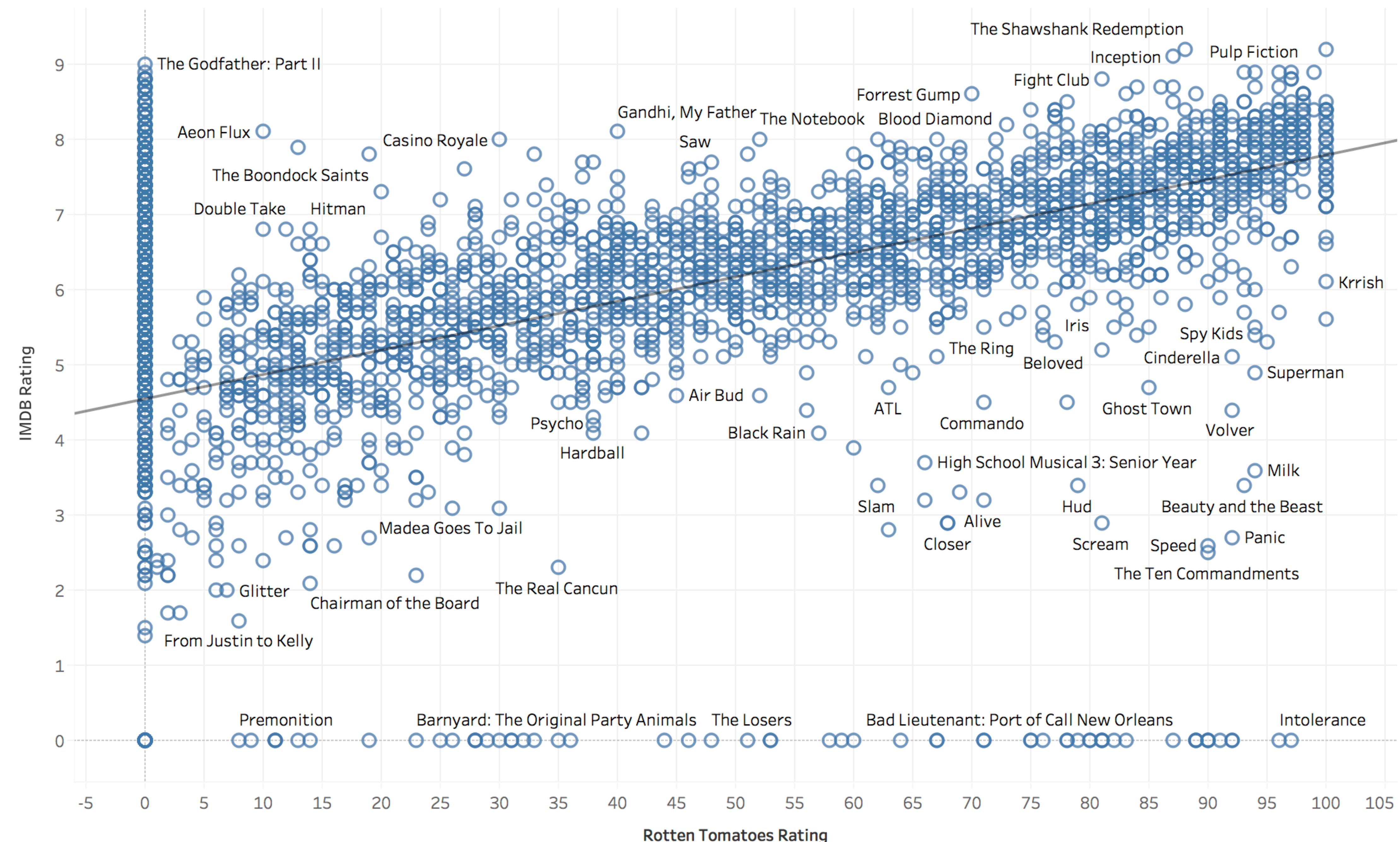
% Releases per Quarter per Genre

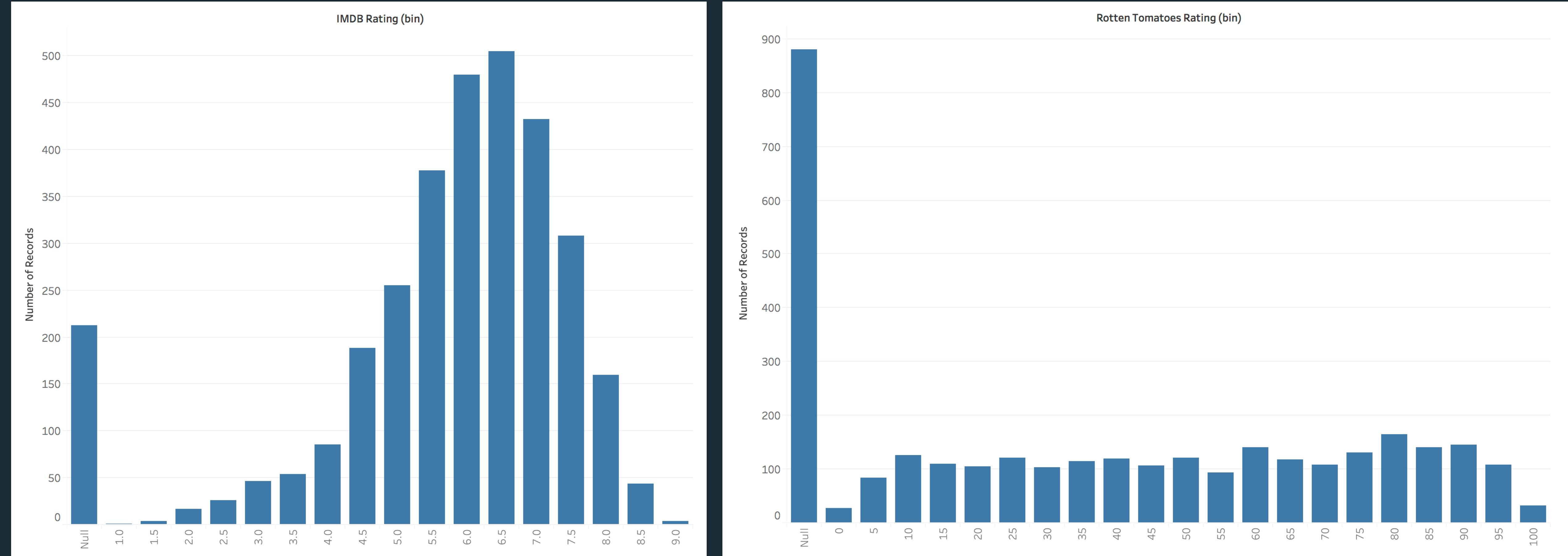


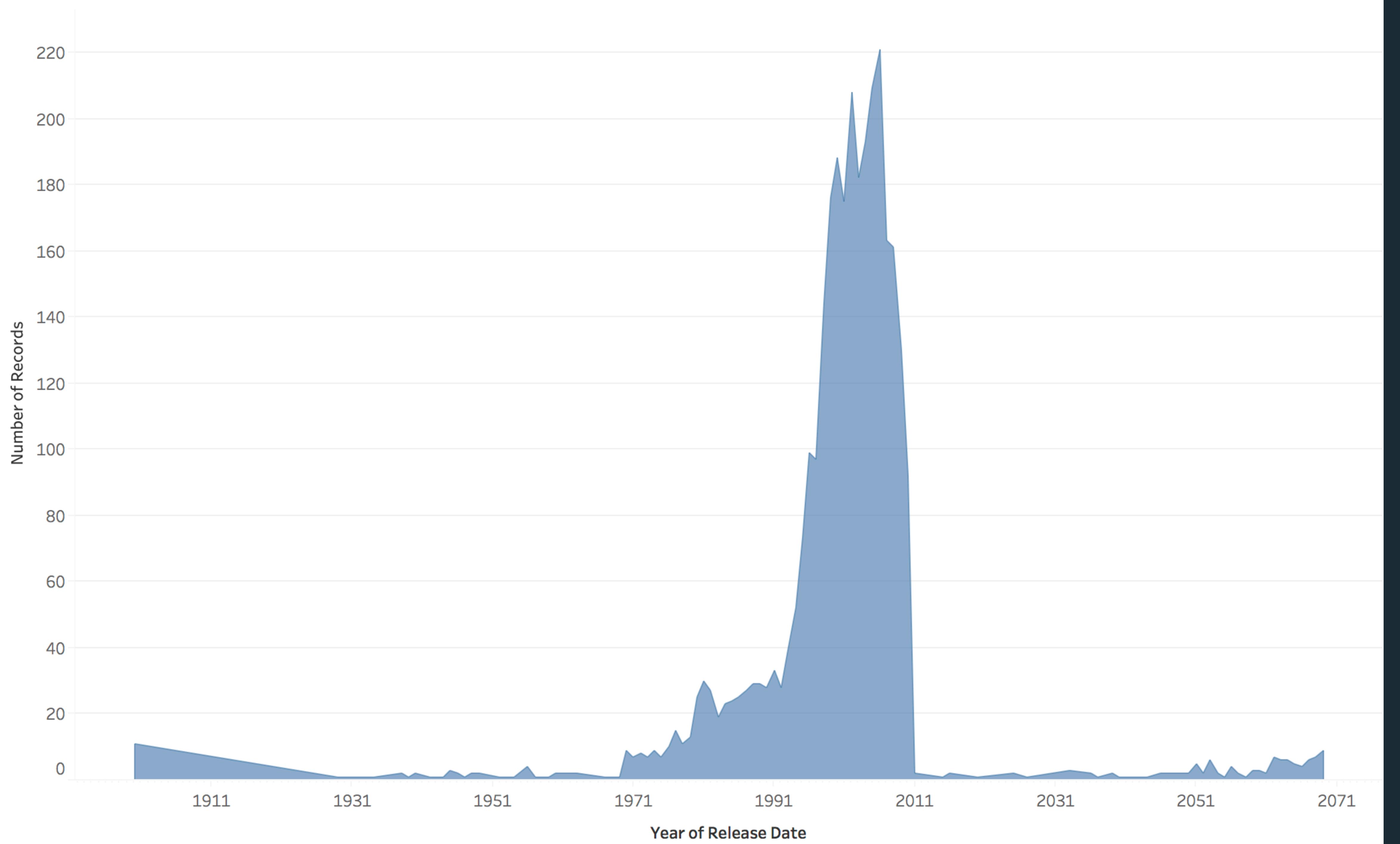
US Gross by Ratings



Audience vs. Critic Ratings







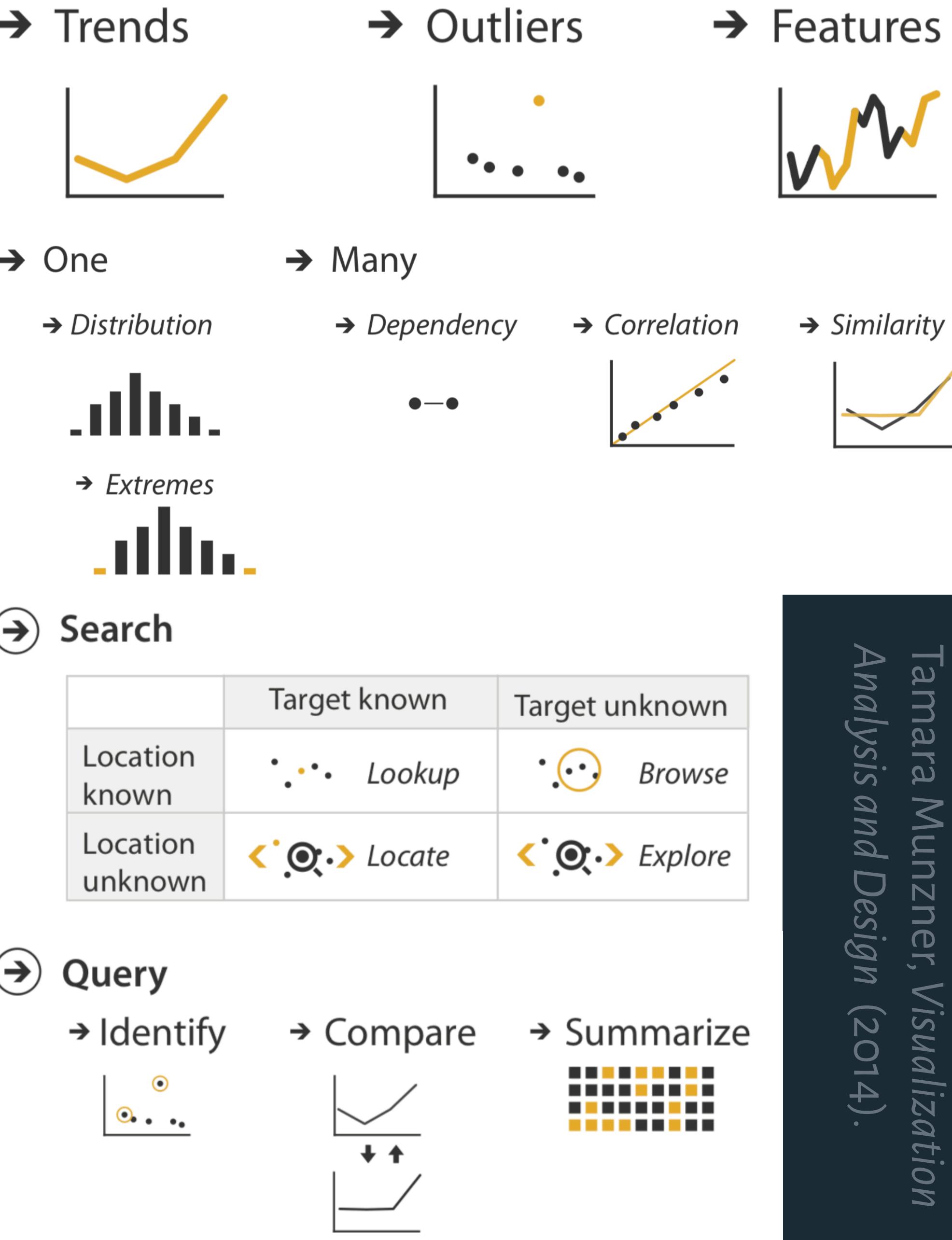
Exploratory Visual Analysis

Process

1. Construct graphics to address questions.
2. Inspect "answer" and ask new questions.
3. Iterate...

Lessons

- ✓ Check **data quality** and your **assumptions**.
- ✓ Start with **univariate summaries**, then consider **relationships between variables**.



Analysis Example: Antibiotic Effectiveness

What questions might we ask?

Collected prior to 1951

Genus of Bacteria

String (N)

Species of Bacteria

String (N)

Antibiotic Applied

String (N)

Gram-Staining?

Pos / Neg (N)

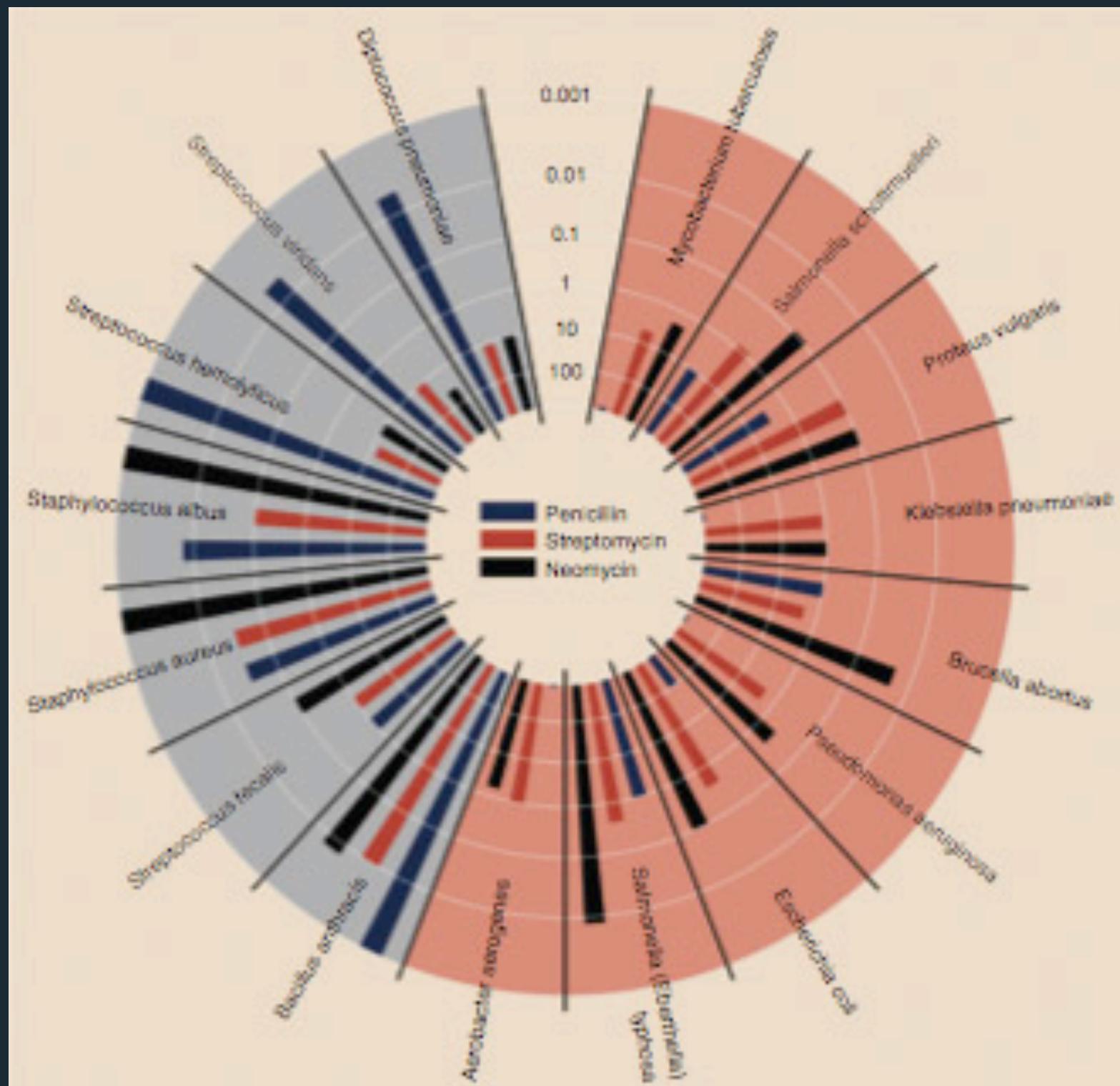
Min. Inhibitory Con. (g)

Number (Q)

Table 1—Burton's Data

Bacteria	Antibiotic				Gram Staining
	Penicillin	Streptomycin	Neomycin		
<i>Aerobacter aerogenes</i>	870	1	1.6		negative
<i>Brucella abortus</i>	1	2	0.02		negative
<i>Brucella anthracis</i>	0.001	0.01	0.007		positive
<i>Diplococcus pneumoniae</i>	0.005	11	10		positive
<i>Escherichia coli</i>	100	0.4	0.1		negative
<i>Klebsiella pneumoniae</i>	850	1.2	1		negative
<i>Mycobacterium tuberculosis</i>	800	5	2		negative
<i>Proteus vulgaris</i>	3	0.1	0.1		negative
<i>Pseudomonas aeruginosa</i>	850	2	0.4		negative
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008		negative
<i>Salmonella schottmuelleri</i>	10	0.8	0.09		negative
<i>Staphylococcus albus</i>	0.007	0.1	0.001		positive
<i>Staphylococcus aureus</i>	0.03	0.03	0.001		positive
<i>Streptococcus fecalis</i>	1	1	0.1		positive
<i>Streptococcus hemolyticus</i>	0.001	14	10		positive
<i>Streptococcus viridans</i>	0.005	10	40		positive

How do the drugs compare?



Bacteria	Penicillin	Antibiotic Streptomycin	Neomycin	Gram stain
<i>Aerobacter aerogenes</i>	870	1	1.6	-
<i>Brucella abortus</i>	1	2	0.02	-
<i>Bacillus anthracis</i>	0.001	0.01	0.007	+
<i>Diplococcus pneumoniae</i>	0.005	11	10	+
<i>Escherichia coli</i>	100	0.4	0.1	-
<i>Klebsiella pneumoniae</i>	850	1.2	1	-
<i>Mycobacterium tuberculosis</i>	800	5	2	-
<i>Proteus vulgaris</i>	3	0.1	0.1	-
<i>Pseudomonas aeruginosa</i>	850	2	0.4	-
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	-
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	-
<i>Staphylococcus albus</i>	0.007	0.1	0.001	+
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	+
<i>Streptococcus fecalis</i>	1	1	0.1	+
<i>Streptococcus hemolyticus</i>	0.001	14	10	+
<i>Streptococcus viridans</i>	0.005	10	40	+

Encodings

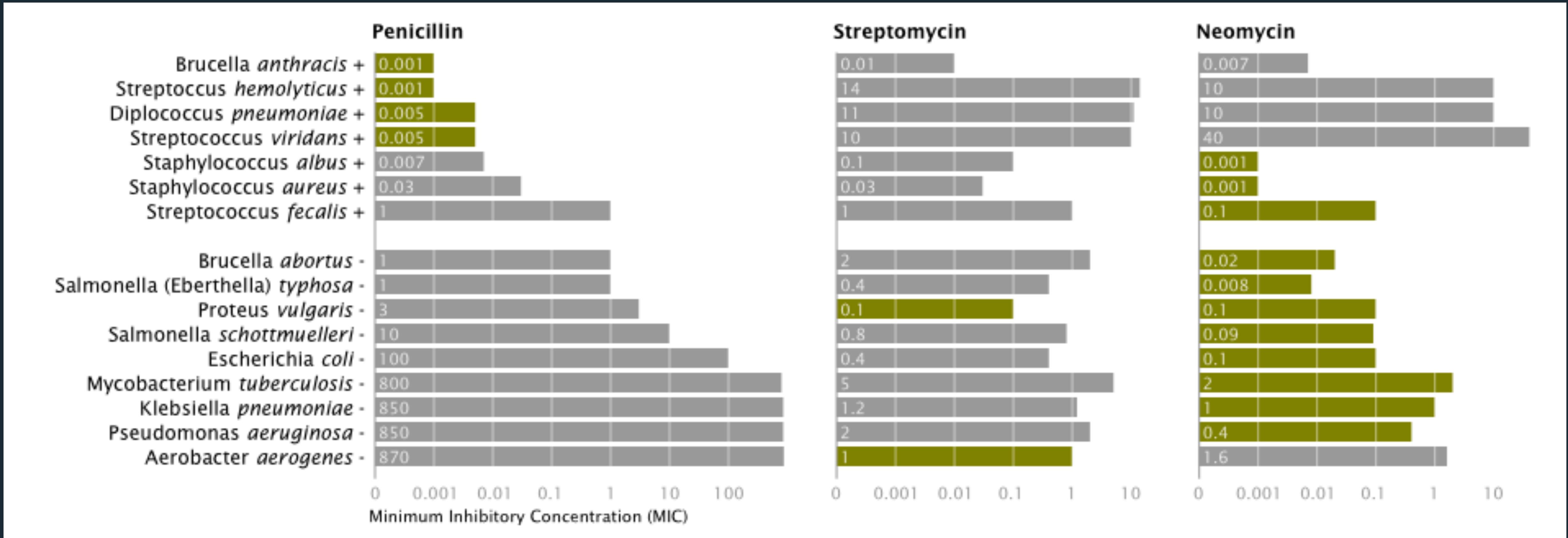
Radius: $1 / \log(\text{MIC})$

Bar Color: Antibiotic

Background Color:
Gram Staining

Original graphic by Will Burtin, 1951.

How do the drugs compare?



X-Axis: Antibiotic | log(MIC)

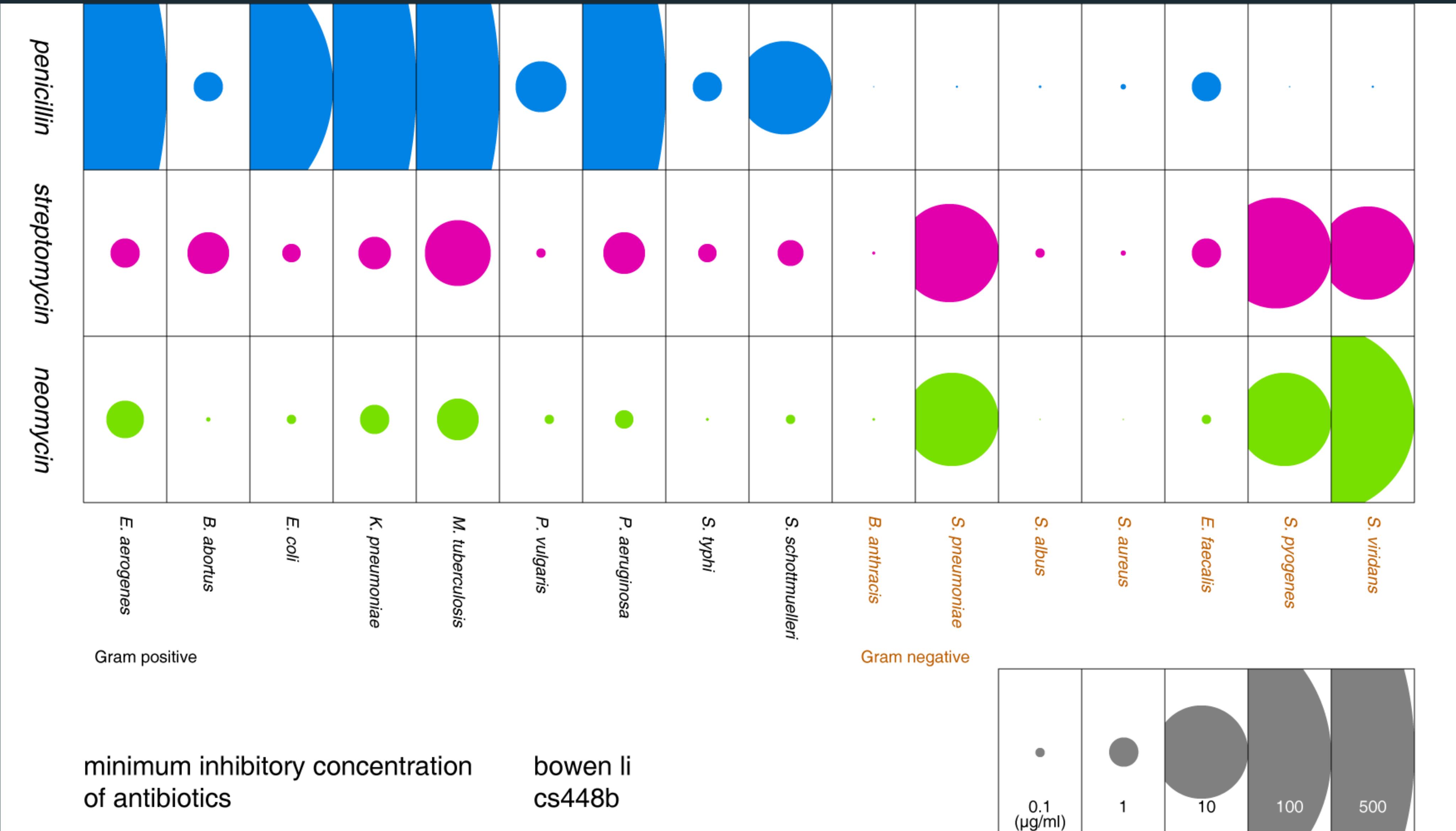
Y-Axis: Gram-Staining | Species

Color: Most Effective?

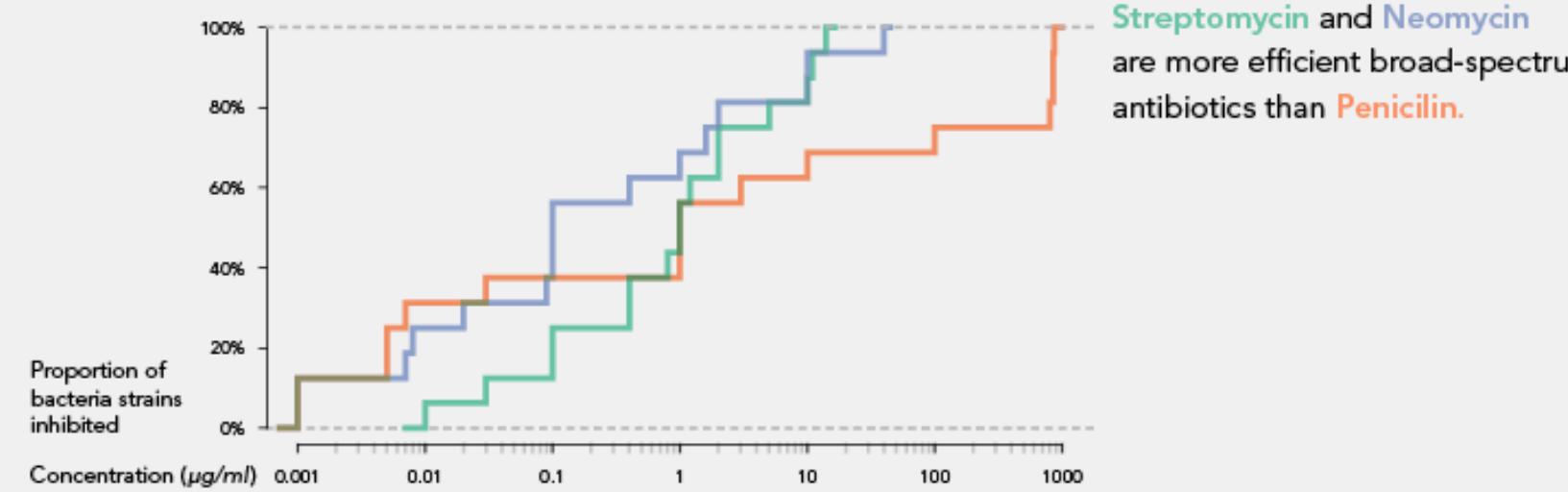
Mike Bostock, Stanford CS448b (Winter 2009).

How do the drugs compare?

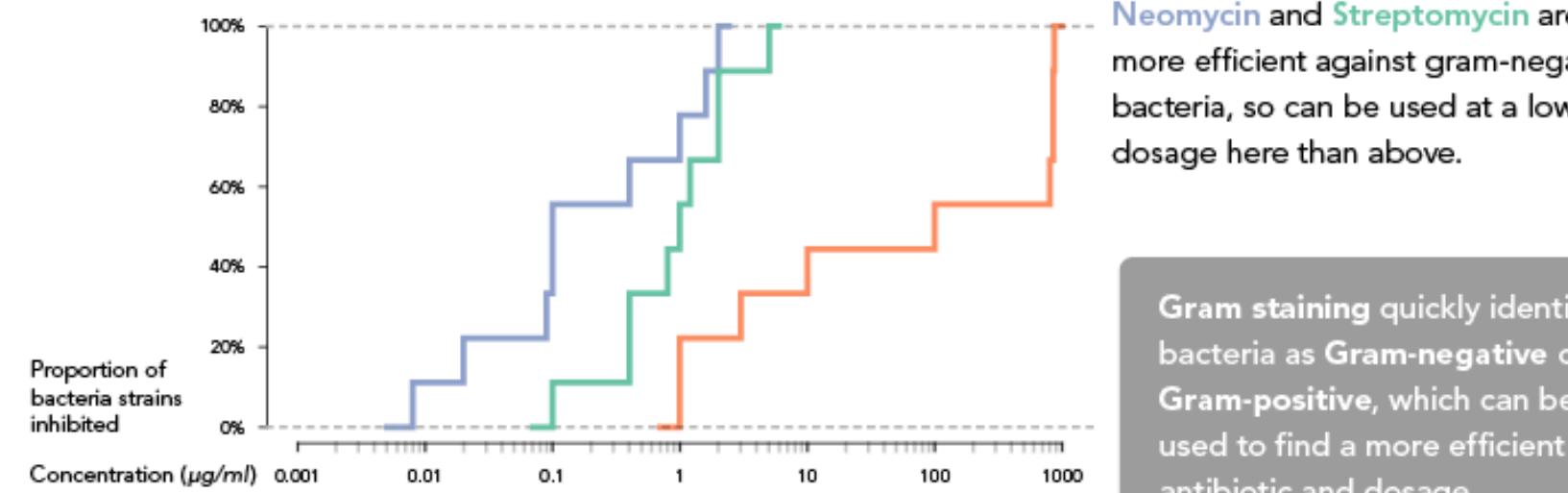
Bowen Li, Stanford CS448b (Fall 2009).



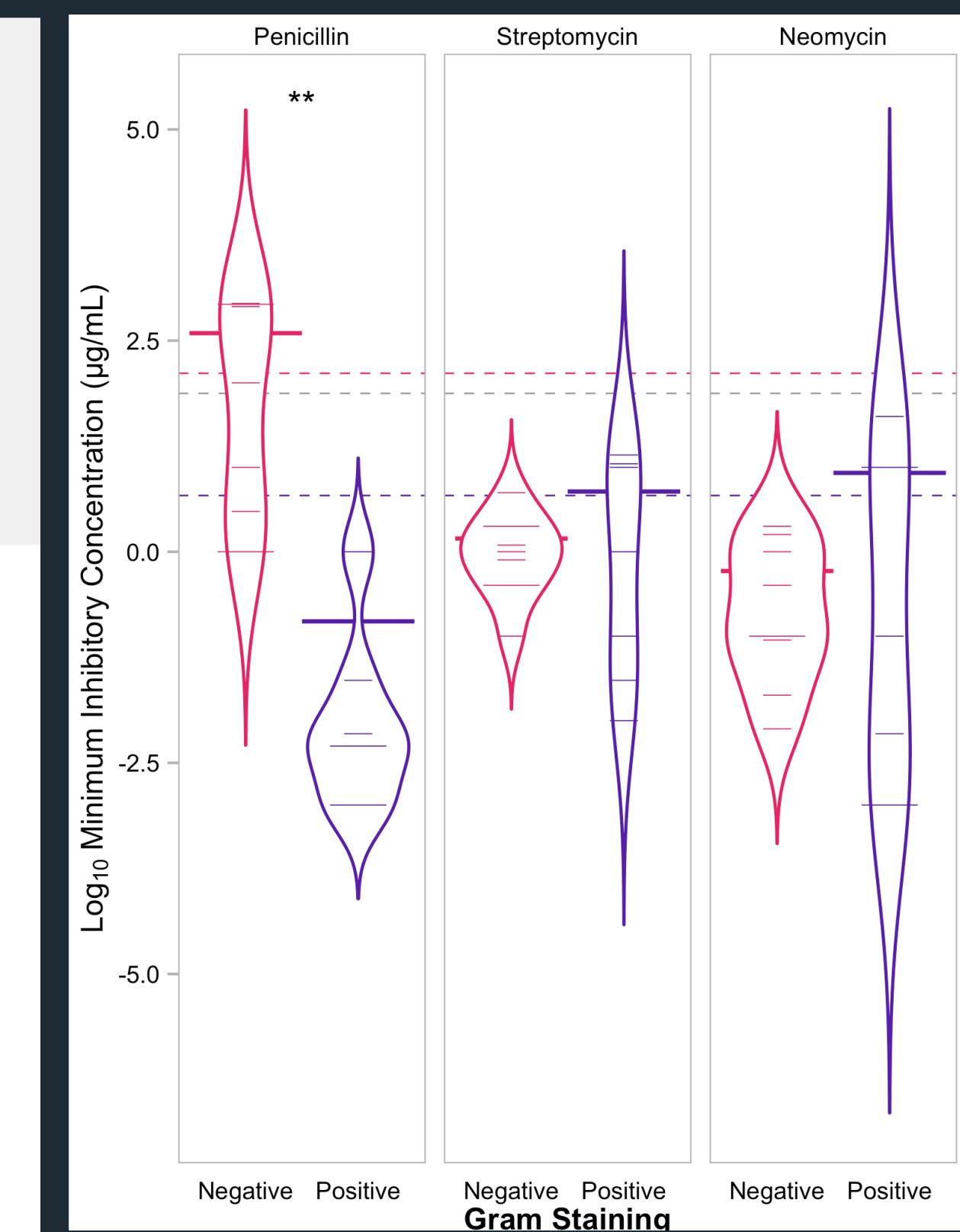
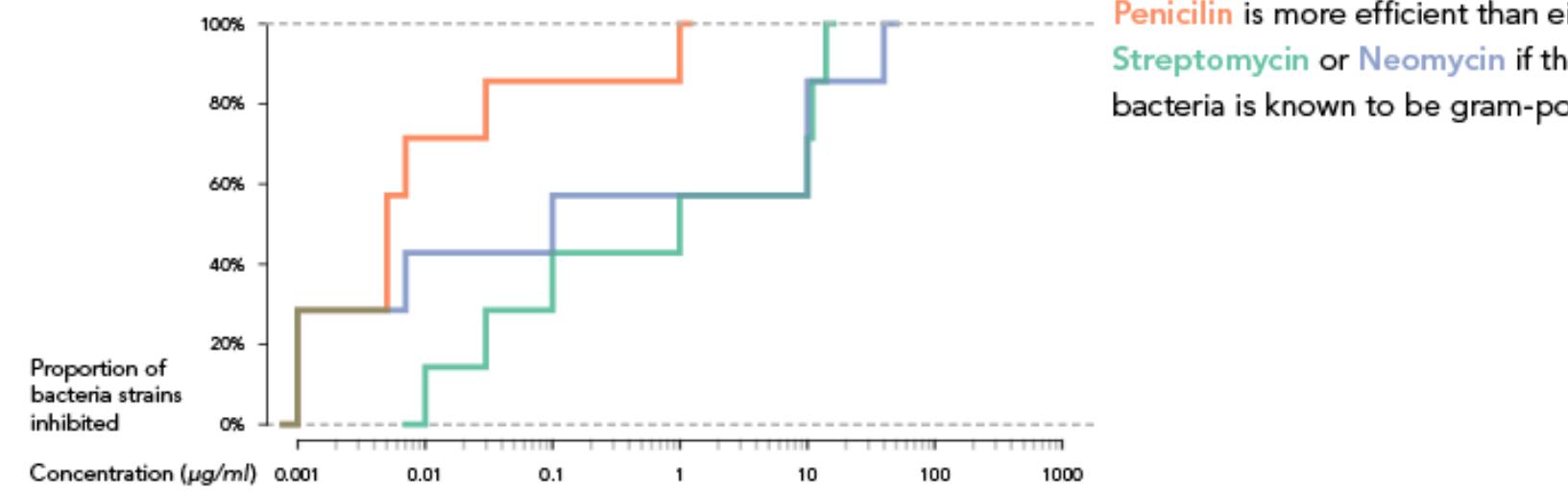
All bacteria



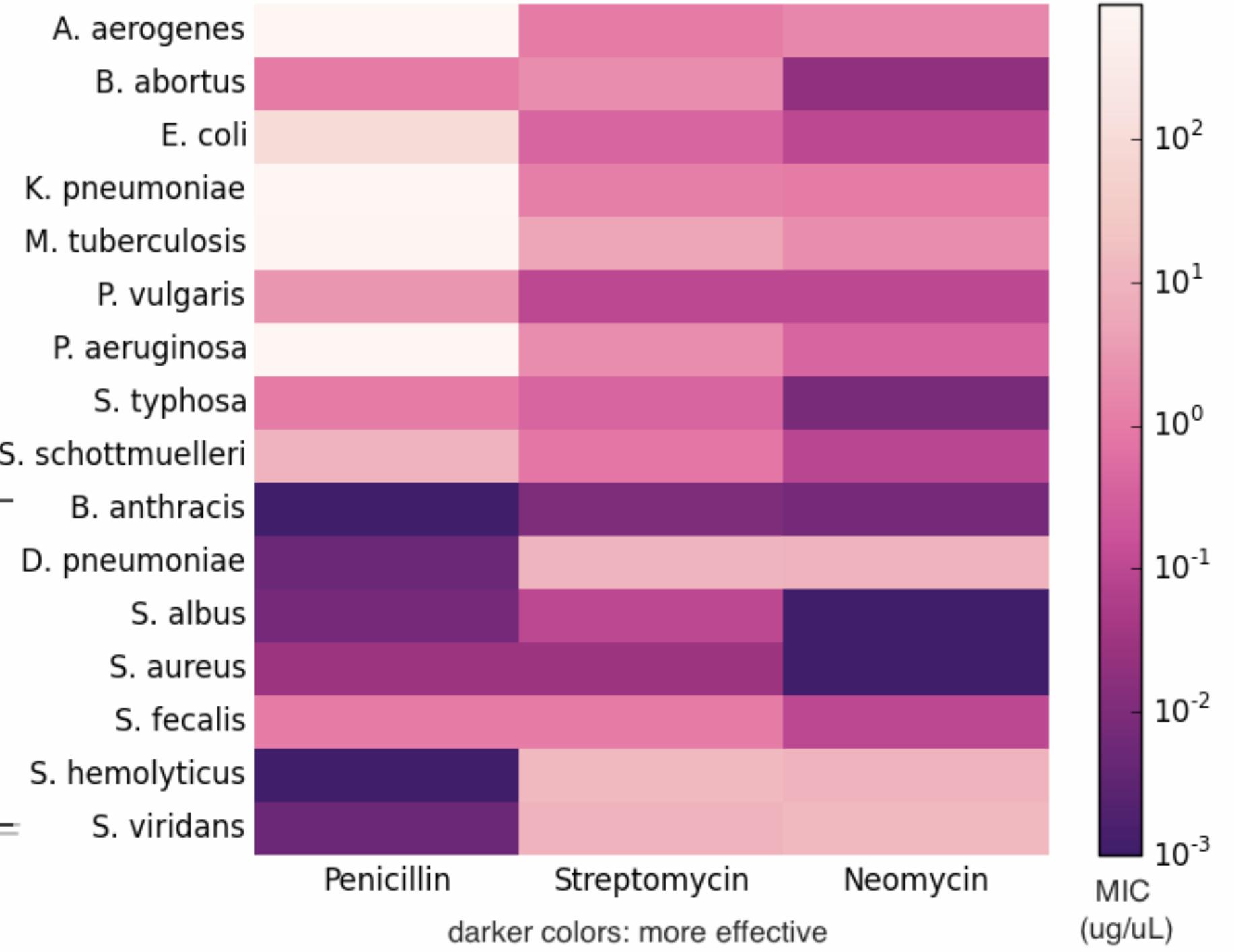
Gram-negative bacteria only



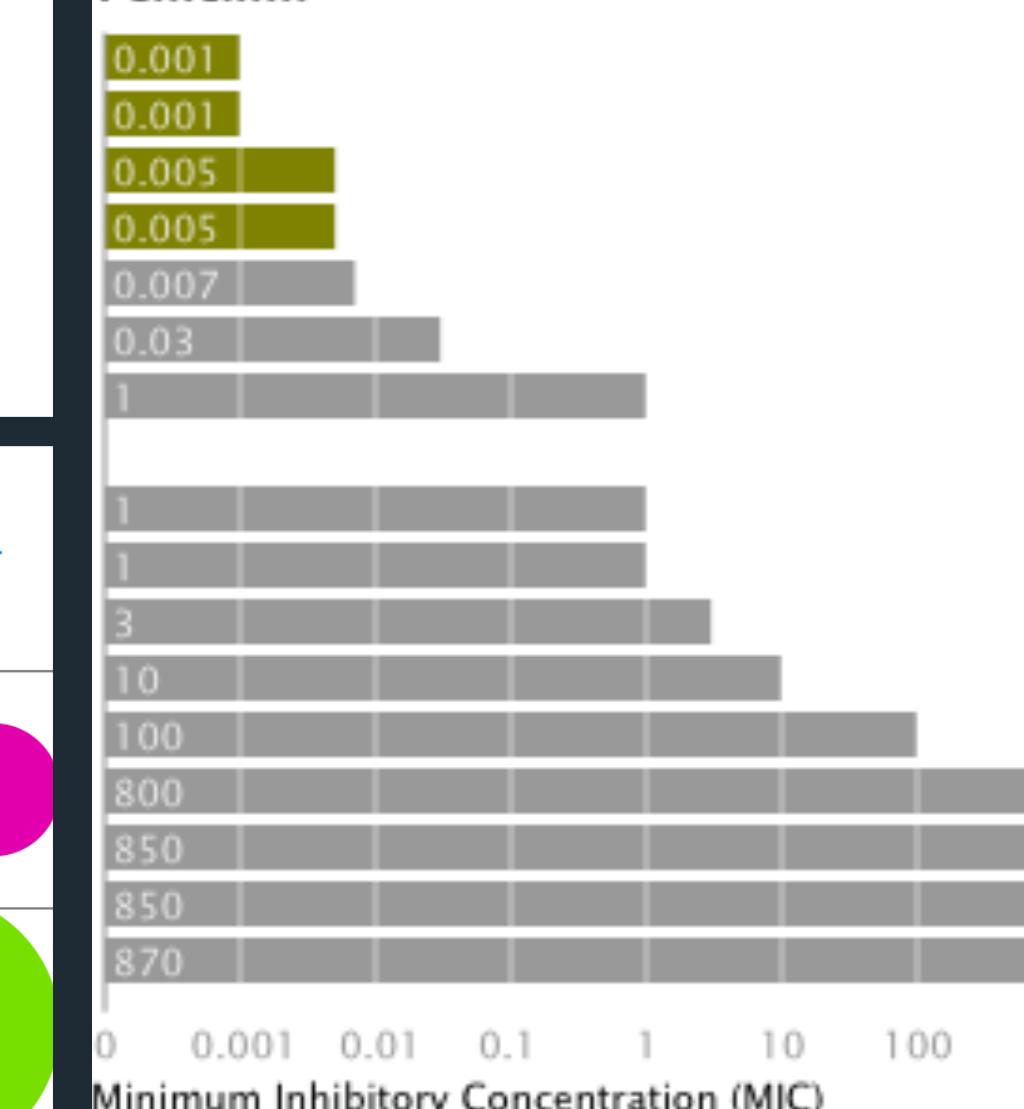
Gram-positive bacteria only



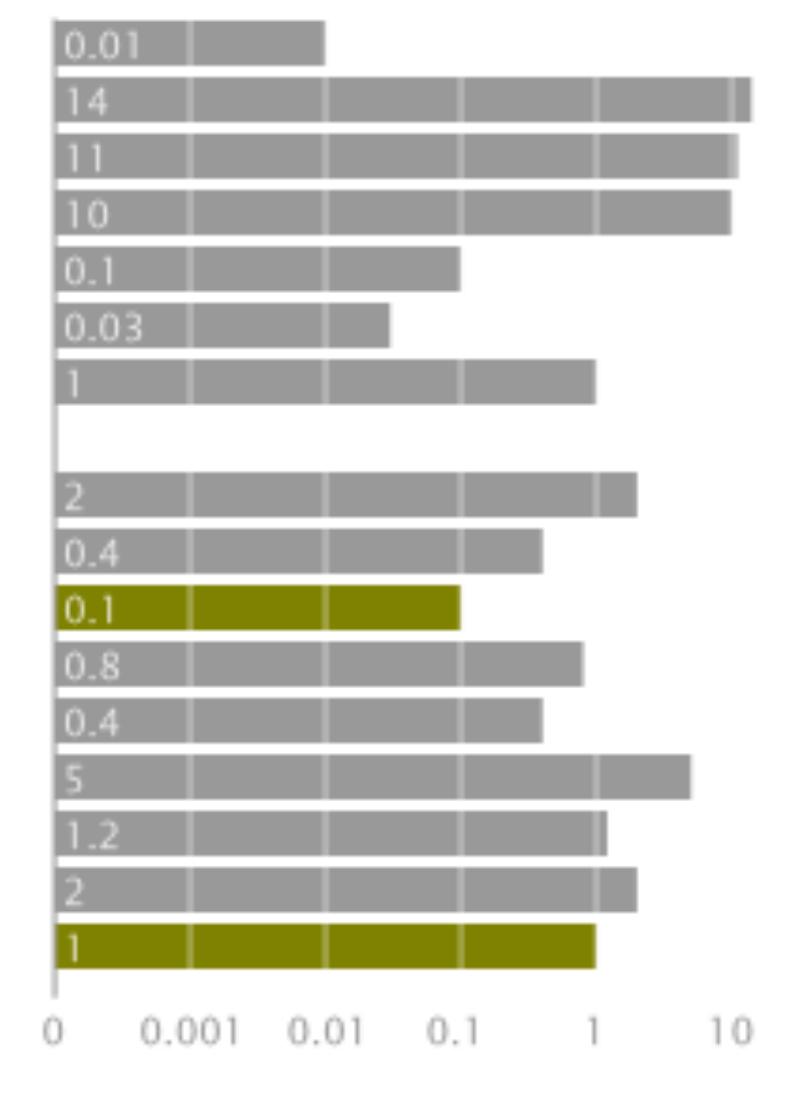
Effectiveness of Antibiotics



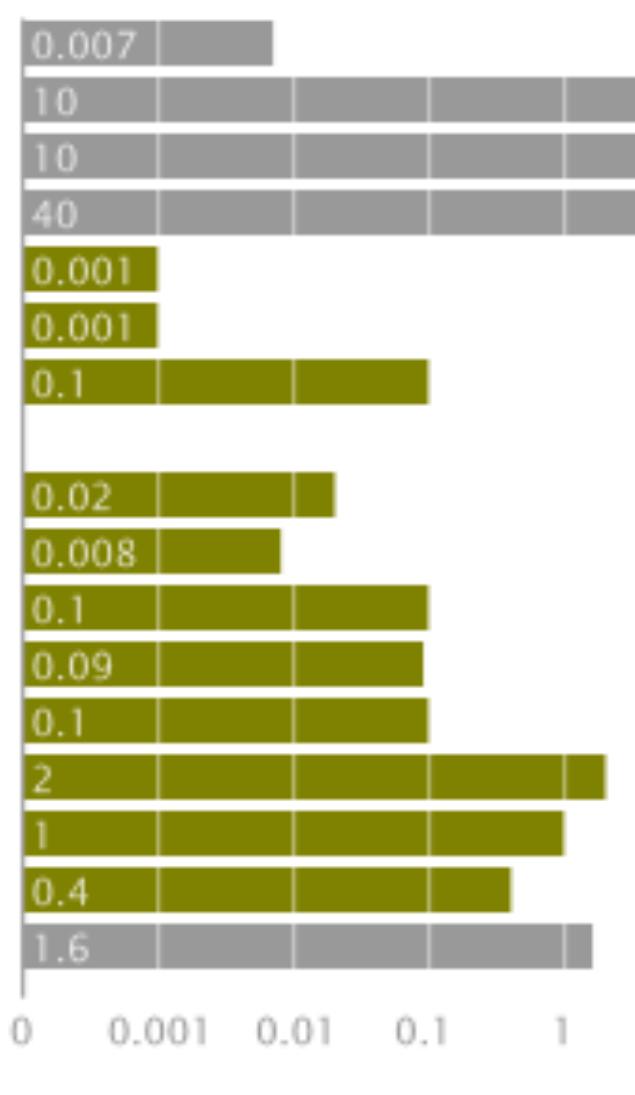
Penicillin



Streptomycin

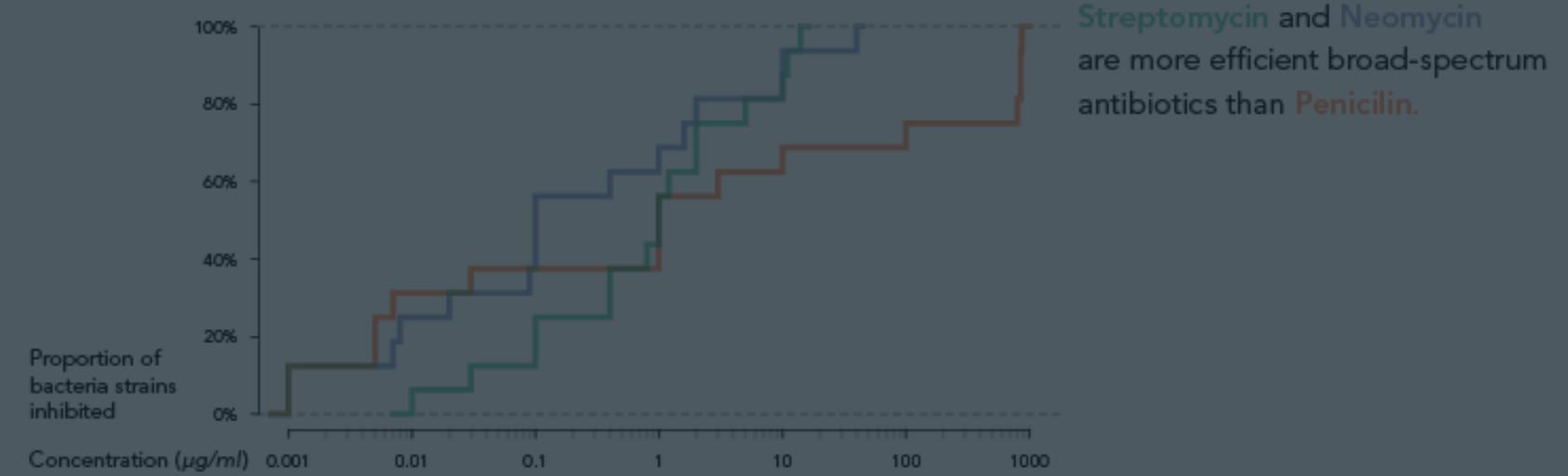


Neomycin

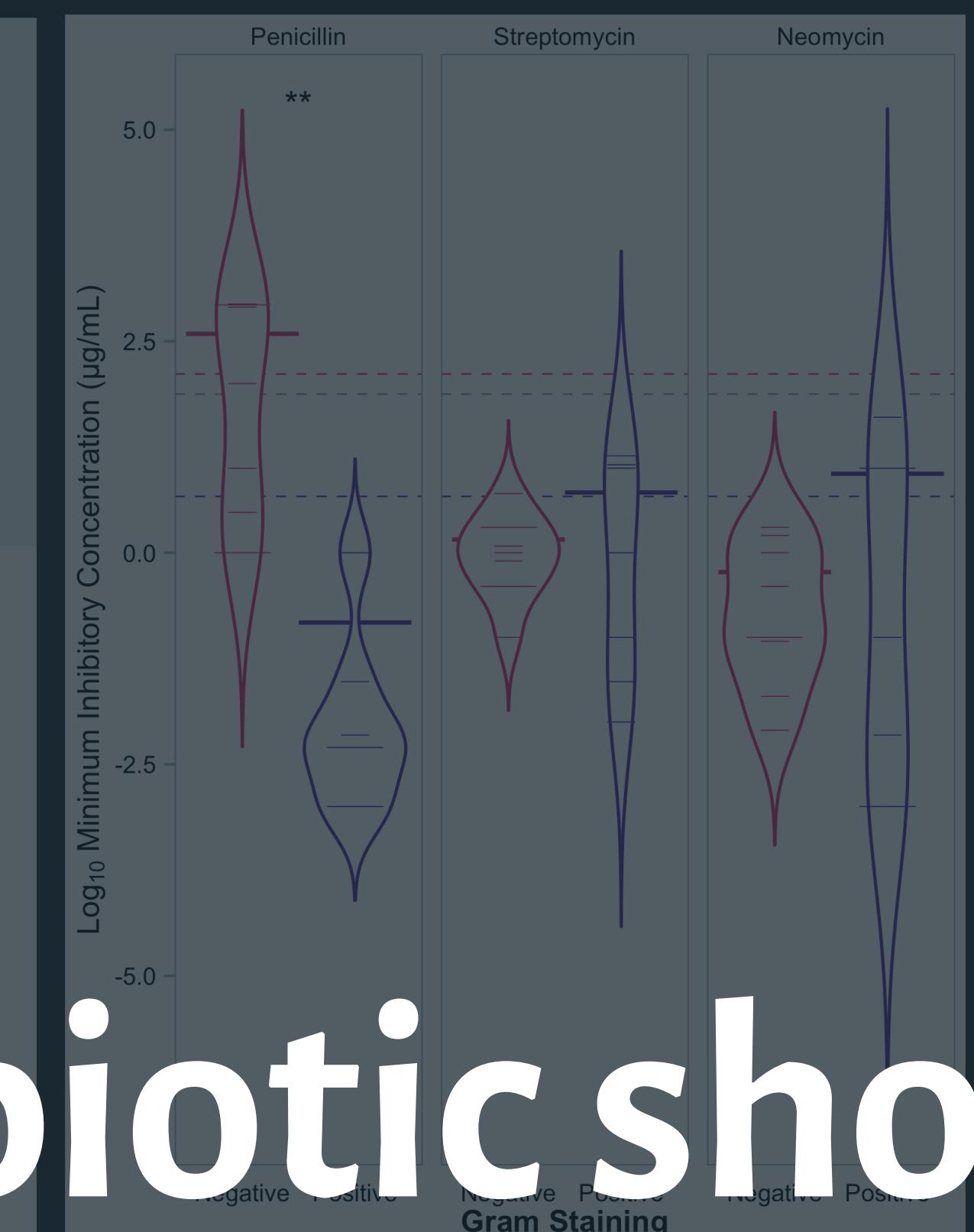
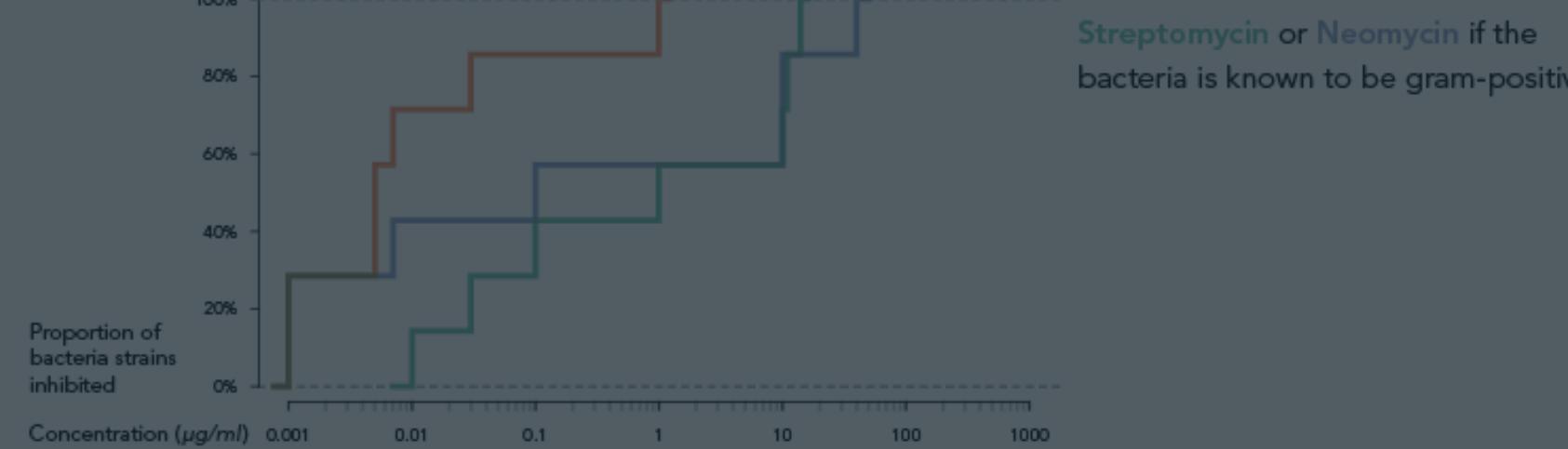
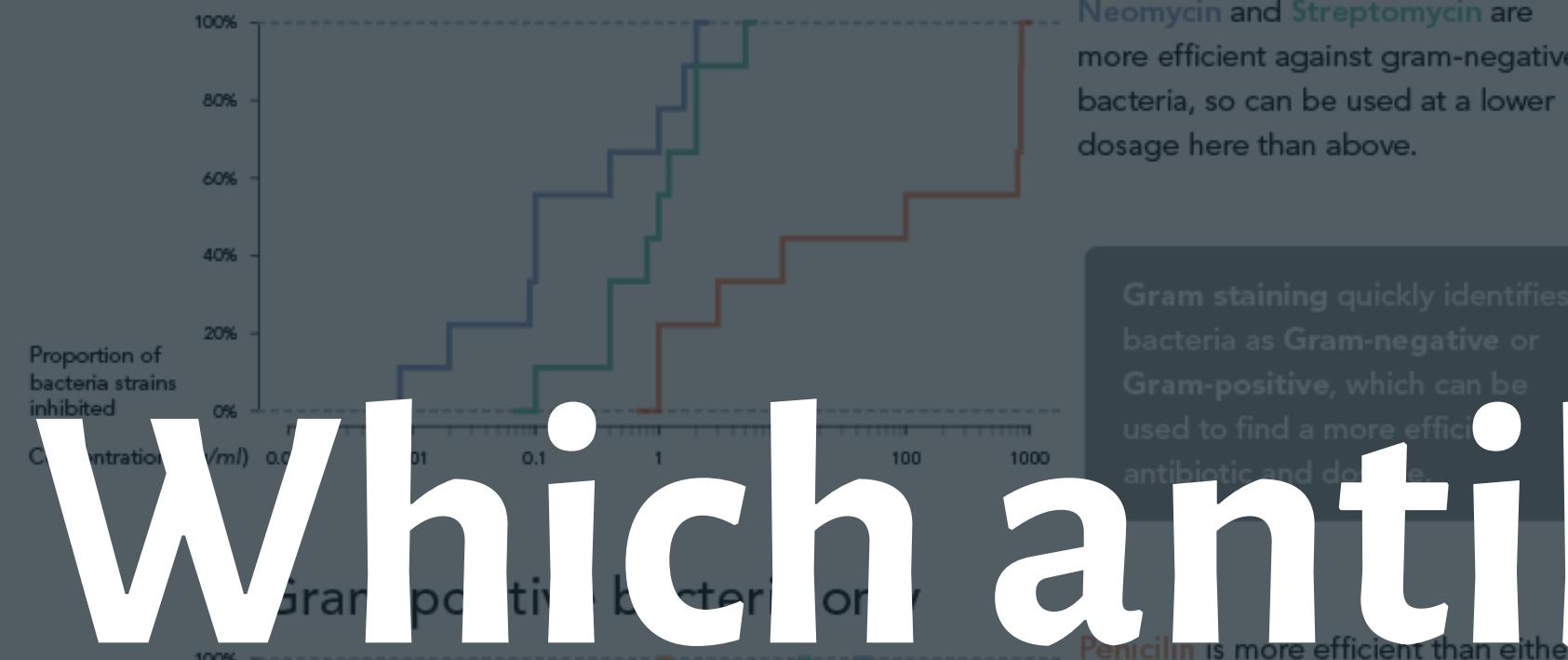


Which antibiotic should one use?

All bacteria



Gram-negative bacteria only



Effectiveness of Antibiotics

A. aerogenes

B. abortus

E. coli

K. pneumoniae

M. tuberculosis

P. vulgaris

P. aeruginosa

S. typhosa

S. schottmuelleri

B. anthracis

D. pneumoniae

S. albus

S. aureus

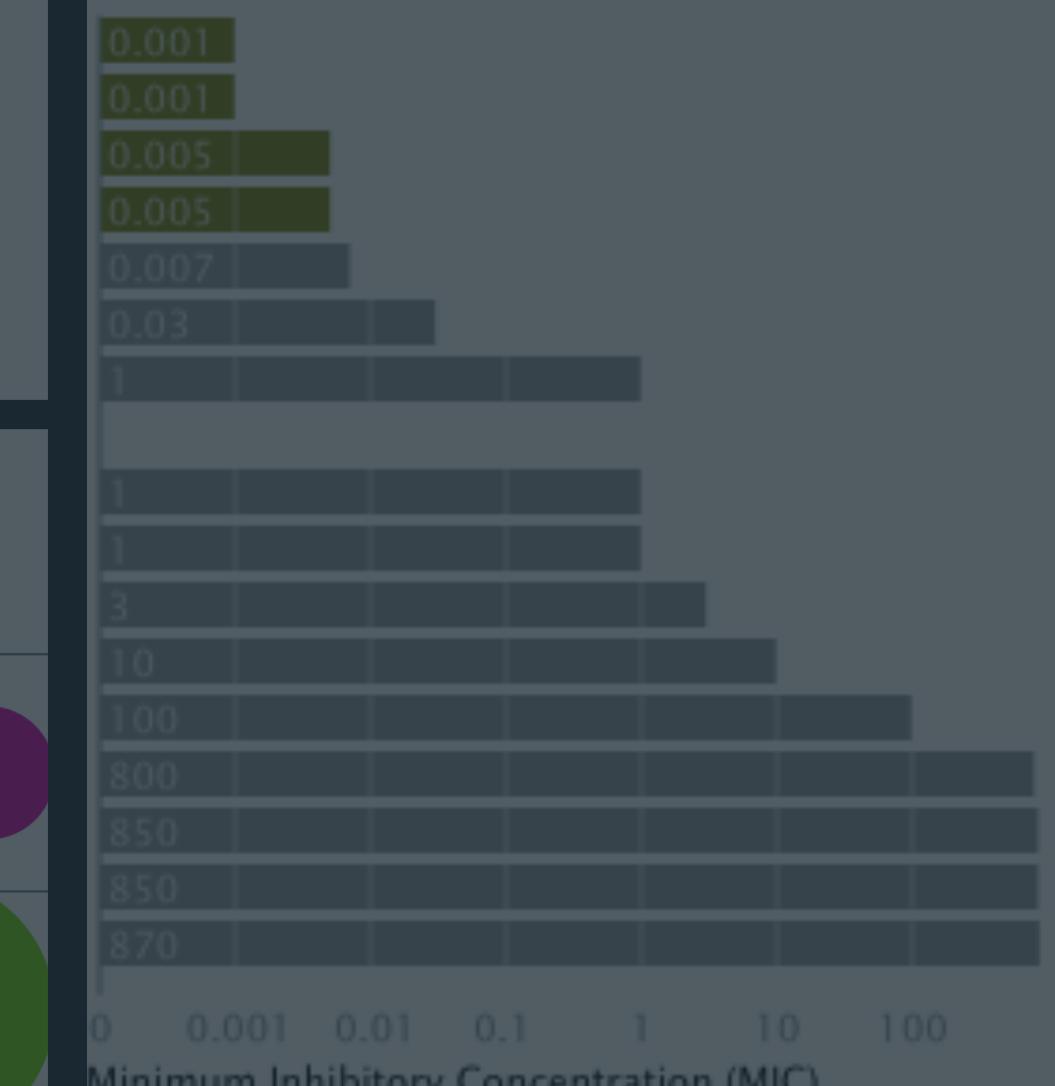
S. fecalis

S. hemolyticus

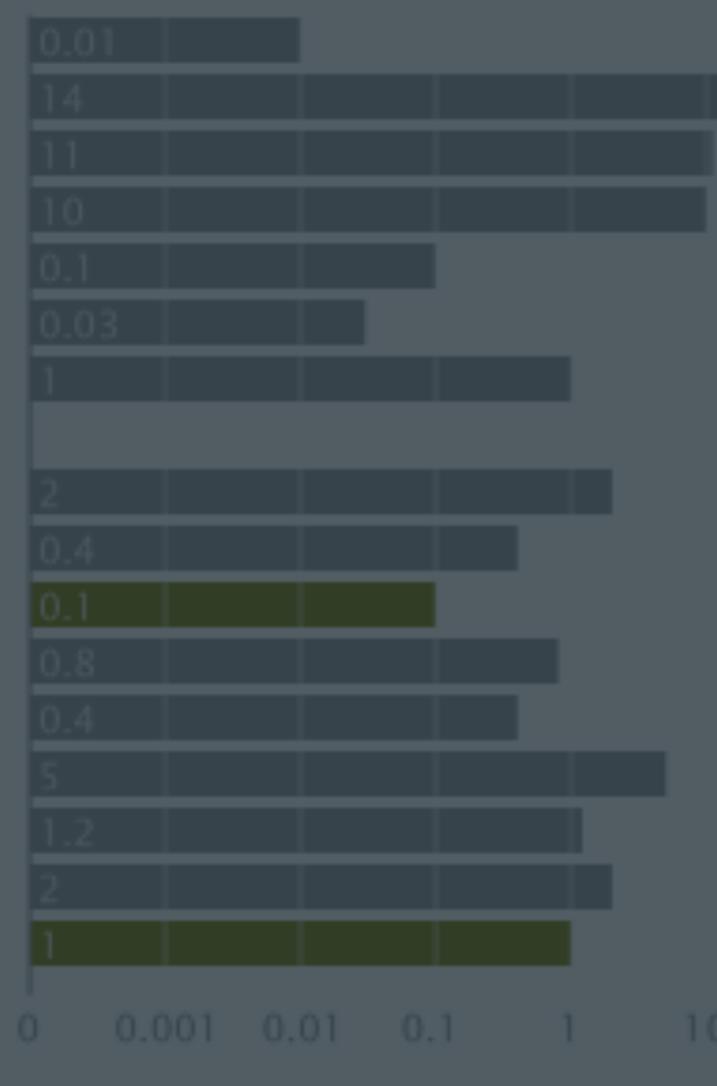
S. virid



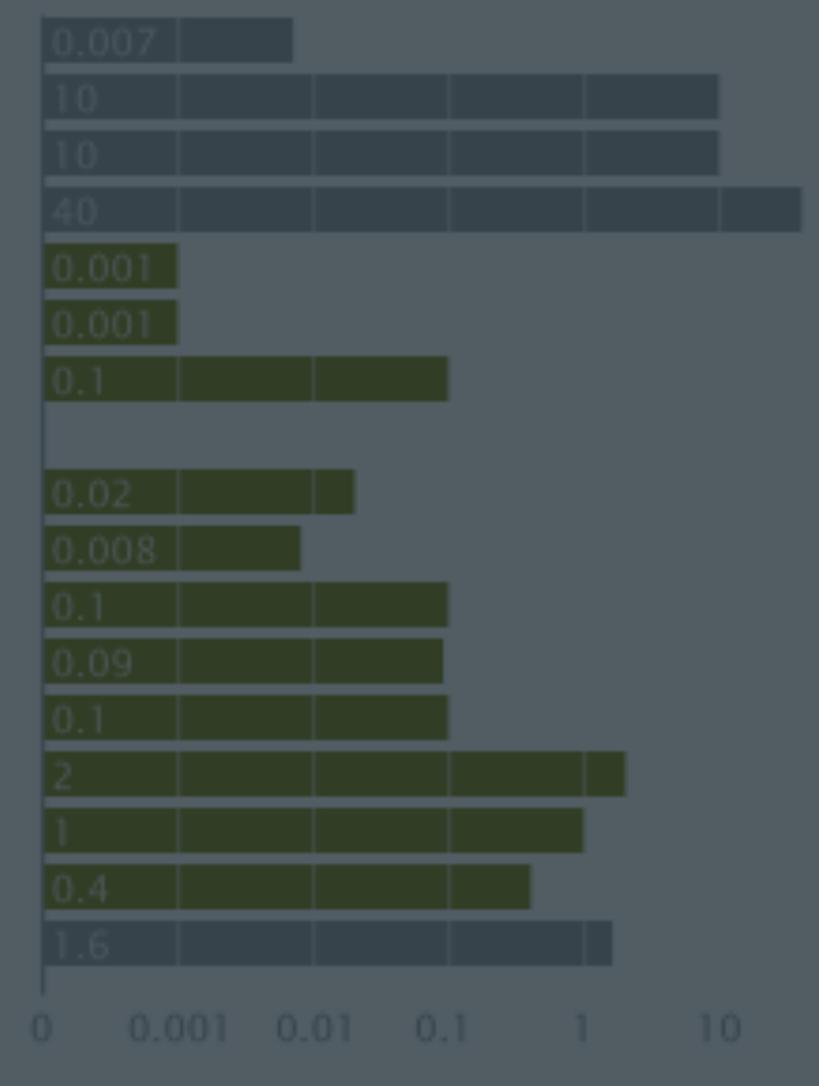
Penicillin

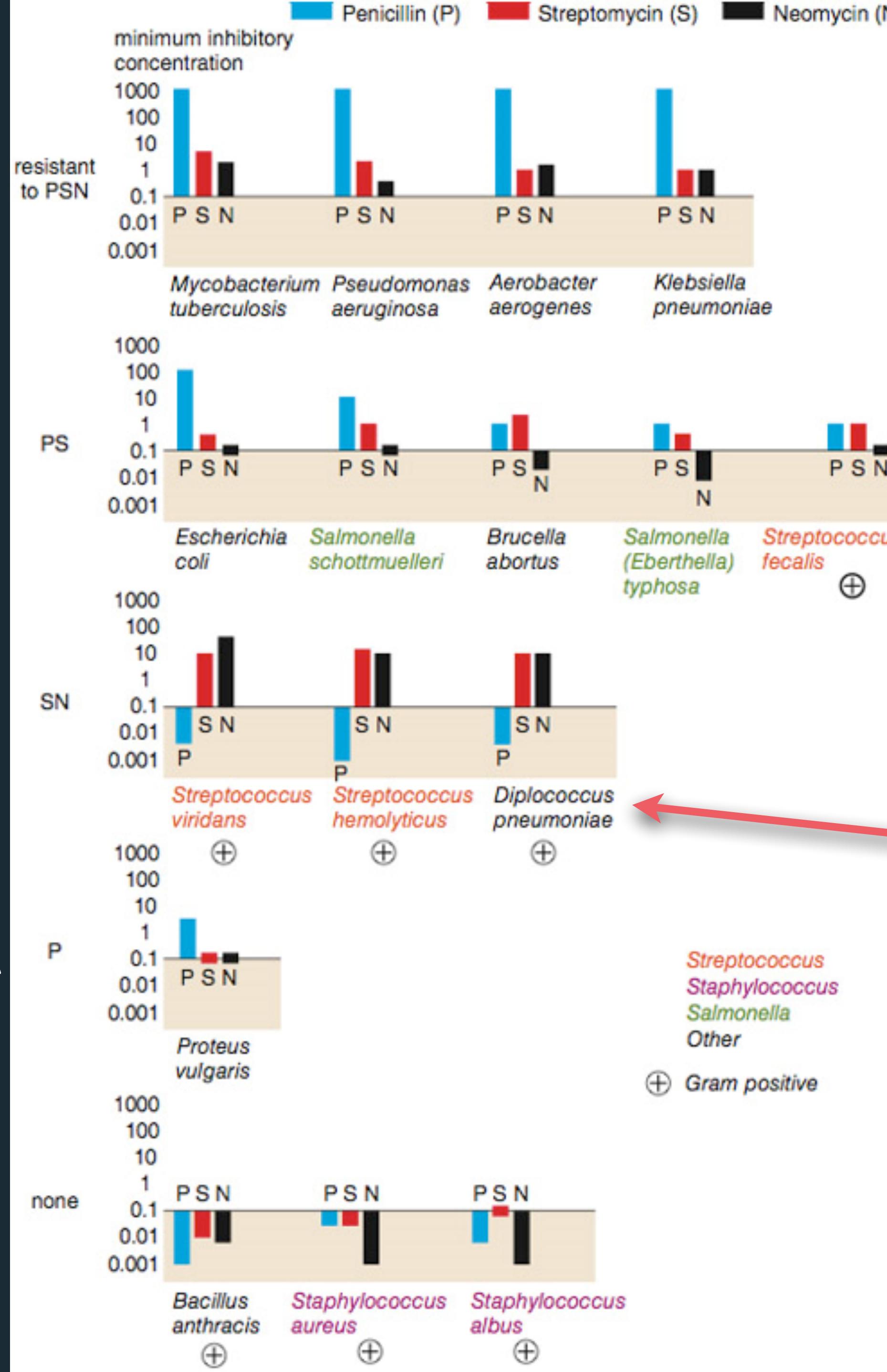


Streptomycin



Neomycin





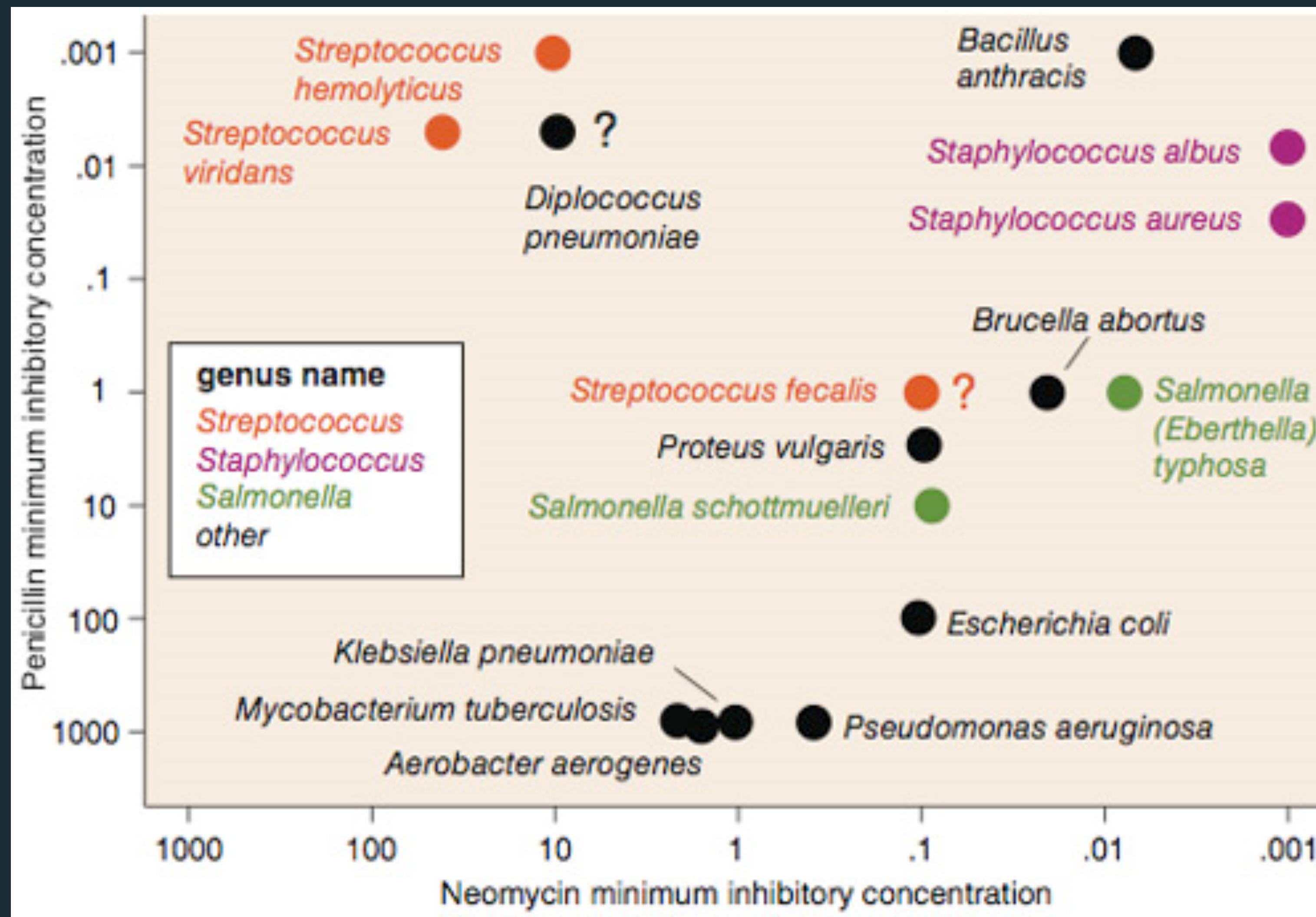
Do the bacteria group by antibiotic resistance?

Not a streptococcus!
(realized ~30 yrs later)

Really a streptococcus!
(realized ~20 yrs later)

Do the bacteria group by resistance? Do different drugs correlate?

Wainer & Lysen. American Scientist, 2009



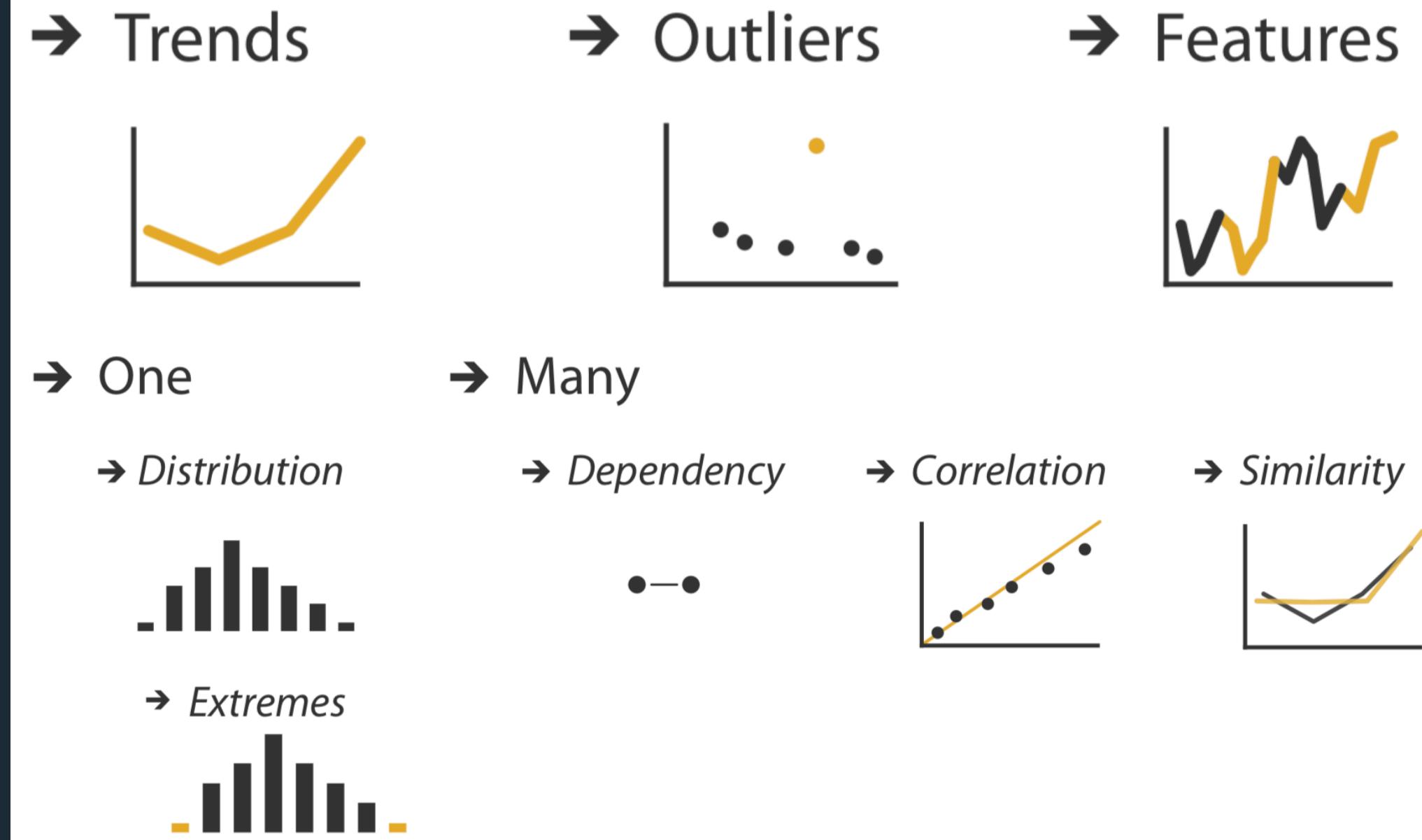
Exploratory Visual Analysis

Process

1. Construct graphics to address questions.
2. Inspect "answer" and ask new questions.
3. Iterate...

Lessons

- ✓ Check **data quality** and your **assumptions**.
- ✓ Start with **univariate summaries**, then consider **relationships between variables**.
- ✓ Avoid **premature fixation**: balance **data variation** and **design variation**.



Search

	Target known	Target unknown
Location known	•.. •.. <i>Lookup</i>	•.. <i>Browse</i>
Location unknown	<i>Locate</i>	<i>Explore</i>

Query

