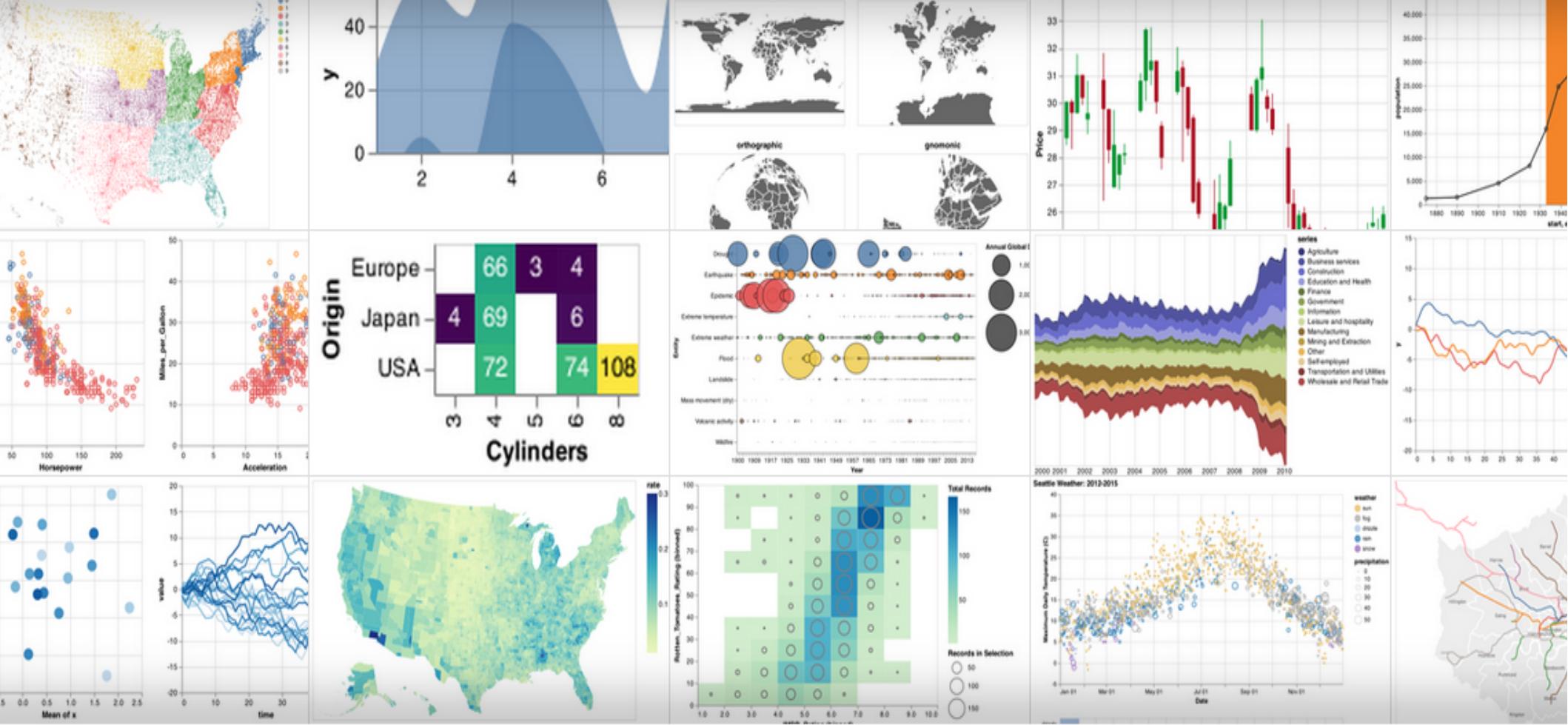




Data 227  
Data Visualization &  
Communication  
Autumn 2023



# Data 227

## Data Visualization & Communication

### Autumn 2023

Just kidding.

# Greetings

- 2023-09-26 Tuesday class
- I am Will Trimble
  - educated at University of Chicago and University of Washington (go Huskies)
- You are ?

# Greetings

- I am Will Trimble
  - educated at University of Chicago and University of Washington (go Huskies)
- You are ?
  - [Undergraduates, mostly 3rd and 4th year, mostly economics and social sciences]
- Have you seen python? R?
  -

# Course Plan

- Plan for the course
  - Half slides/discussion, half interactive programming
  - Five homeworks, two projects, no exams
  - No particular technology required, but I can help you with python.
- This class:
  - Quick overview of topics
  - Software choices for visualization
  - jupyter notebook 10 minutes to pandas
  - HW 1 due Friday Oct 6 - census visualization+ data cleaning

# Week 1 readings

- Hadley Wickham. Tidy Data. Journal of Statistical Software 2014
- The Power of Representation, Chapter 3 in Things That Make Us Smart. Don Norman.
- Chapter 1: Graphical Excellence, In The Visual Display of Quantitative Information. Tufte.

# Encoding channels

VARIABLES OF THE IMAGE		POINT	LINE	AREA (ZONE)
XY 2 DIMENSIONS OF THE PLANE				
Z	SIZE			
DIFFERENTIAL VARIABLES				

# Data types

- **integer**: stored exactly, but with a limit, positive and negative
- **floating point**: stored in scientific notation, not perfectly exact, positive and negative
- **dates**: usually need to be handled separately from numbers and strings (can't add "2021-09-28" and 21:15 otherwise)
- **complex**: encountered usually in physics and engineering; each point has two real numbers with a slightly different algebra
- **vectors and matrices**:  $1 \times n$  or  $m \times n$  objects, usually numbers

Shoutout to ISO-8601. There is a right way to write dates (for science). YYYY-MM-DD always. Why ?

# More data types

- **character:** single letter, number, or punctuation mark, conventionally set aside by quotation marks: “A” “8” 'j' '^M'
- **strings:** sequences of ordered symbols (text)

“This course introduces best practices for presenting and communicating quantitative data.”
- **unicode:** generalization of strings to every human alphabet and writing system “السلام عليكم” “A tsunami (from from Japanese: 津波 , lit. 'harbour wave') is a series of waves”
- **Boolean:** stores True or False, 0 or 1, “binary” data flag

When you import your data, if your computing environment does not guess the data types correctly, you will not be able to use the data.

# Attribute types

- **Nominal :** Labels or categories  
=, ≠ ex: nationality, species, “type”,  
postal code (can be grouped)
- **Ordinal:** Ordered  
=, ≠ <>>  
ex: Likert scale (agree strongly, somewhat agree..)  
Grades A/B/C/D/F
- **Quantitative (interval):**  
=, ≠ <>> +- ex: dates, location  
(ratios meaningless, differences meaningful)
- **Quantitative (ratio):** physical measurement w/ meaningful zero  
=, ≠ <>> + \* / ex: mass, length, counts

# Attribute types

- **Nominal :** Labels or categories  
=, ≠ ex: nationality, species, “type”,  
postal code (can be grouped)
- **Ordinal:** Ordered  
=, ≠ <> ex: Likert scale (agree strongly, somewhat agree..)  
Grades A/B/C/D/F
- **Quantitative (interval):**  
=, ≠ <> +- ex: dates, location

This less-familiar categorization helps us understand what visualization techniques are feasible with which data.

## What?

### Datasets

### Attributes

# Chart types

## → Data Types

→ Items → Attributes → Links → Positions → Grids

## → Attribute Types

→ Categorical



## → Data and Dataset Types

Tables	Networks & Trees	Fields	Geometry	Clusters, Sets, Lists
Items	Items (nodes)	Grids	Items	Items
Attributes	Links	Positions	Attributes	

→ Ordered

→ Ordinal

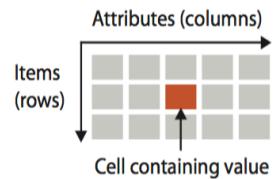


→ Quantitative

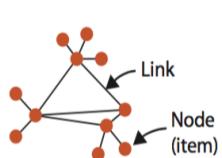


## → Dataset Types

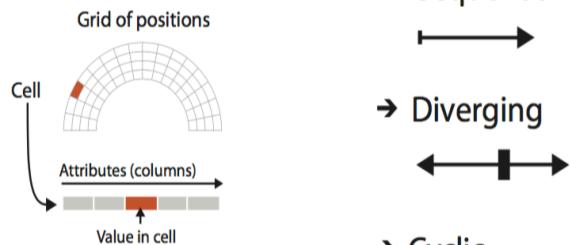
### → Tables



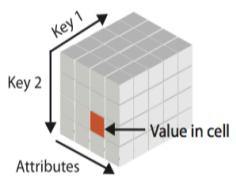
### → Networks



### → Fields (Continuous)



### → Multidimensional Table



### → Trees



### → Geometry (Spatial)

## → Ordering Direction

→ Sequential



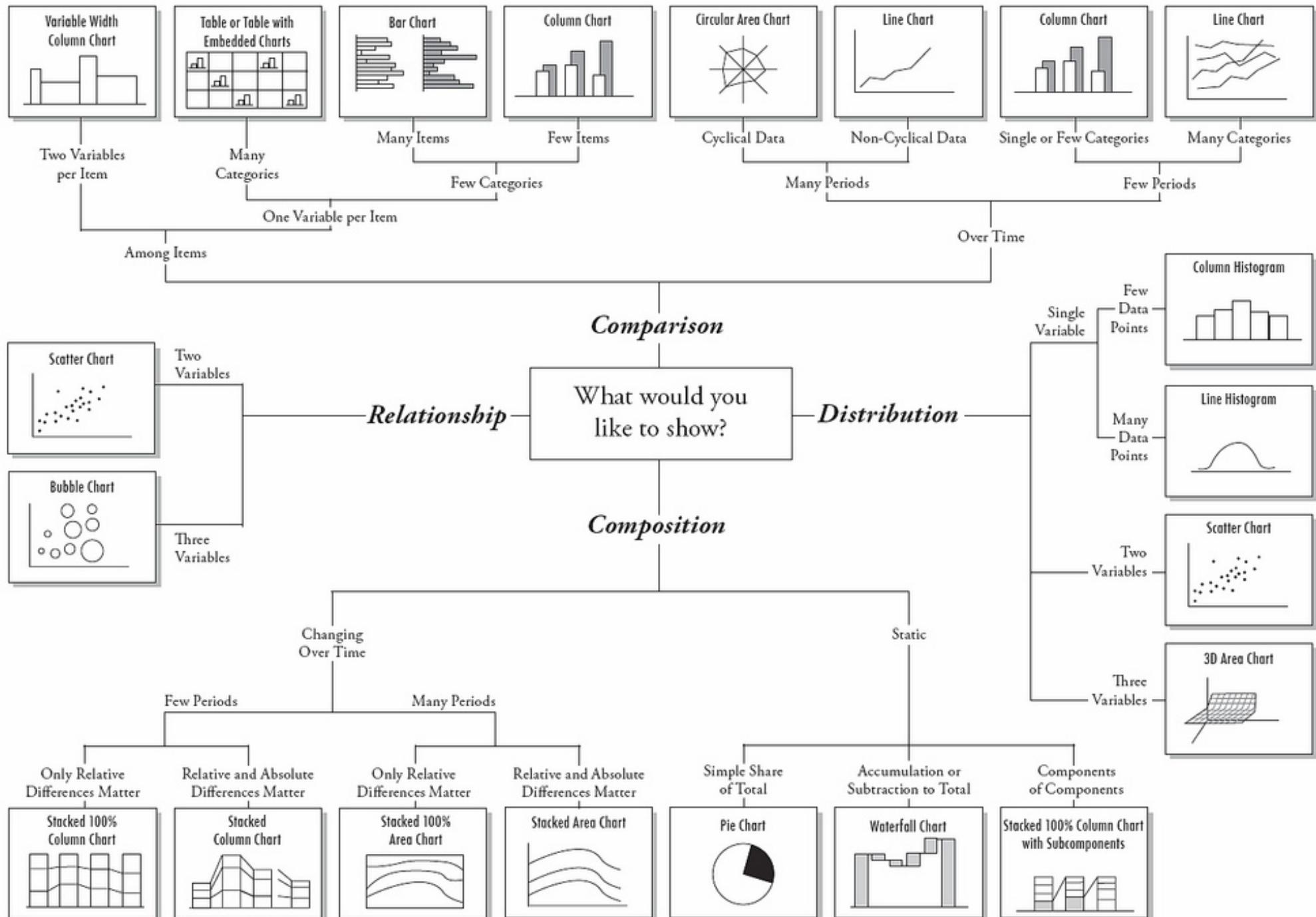
→ Diverging



→ Cyclic



# Chart Suggestions—A Thought-Starter

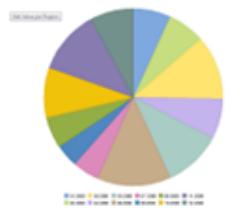


# Chart types

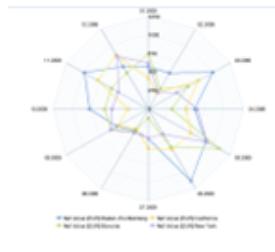
Area chart



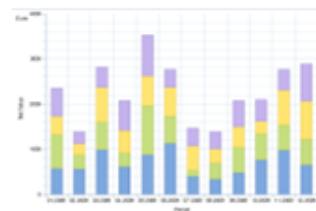
Pie chart



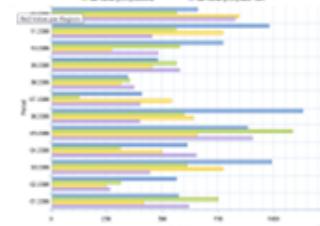
Radar chart



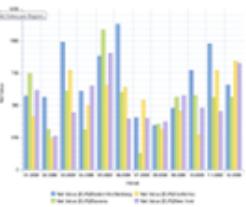
Stacked columns chart



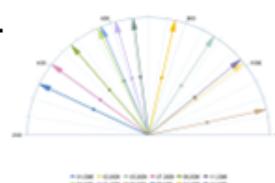
Bar chart



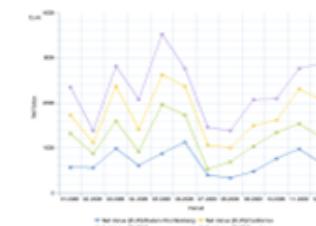
Pipeline chart



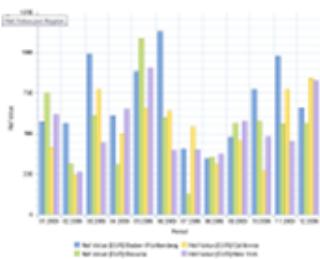
Speedometer chart



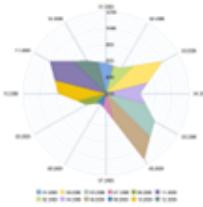
Stacked line chart



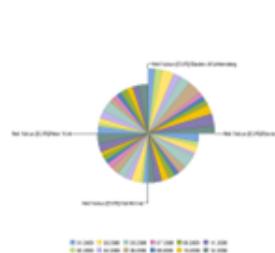
Vertical bar chart



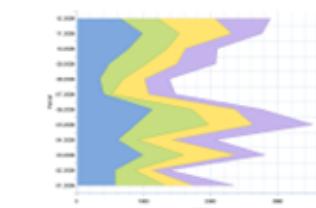
Polar chart



Split pie chart



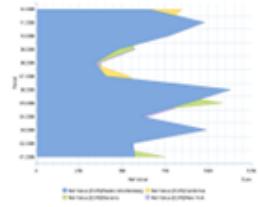
Stacked profile area chart



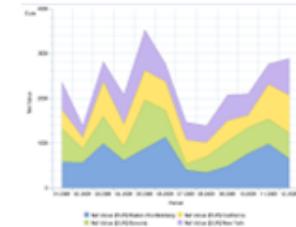
Doughnut chart



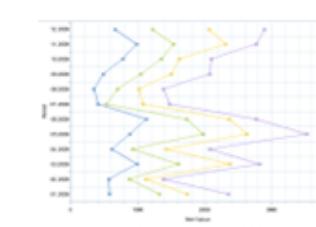
Profile area chart



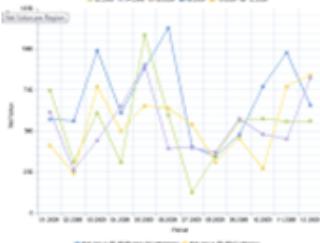
Stacked area chart



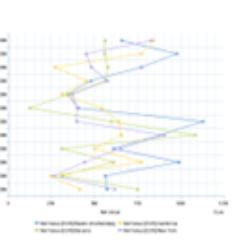
Stacked profile chart



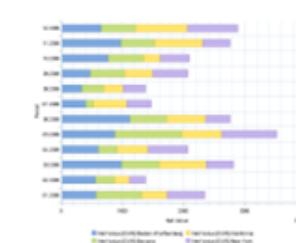
Line chart



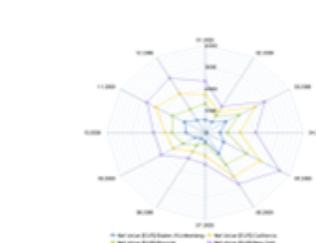
Profile chart



Stacked bar chart



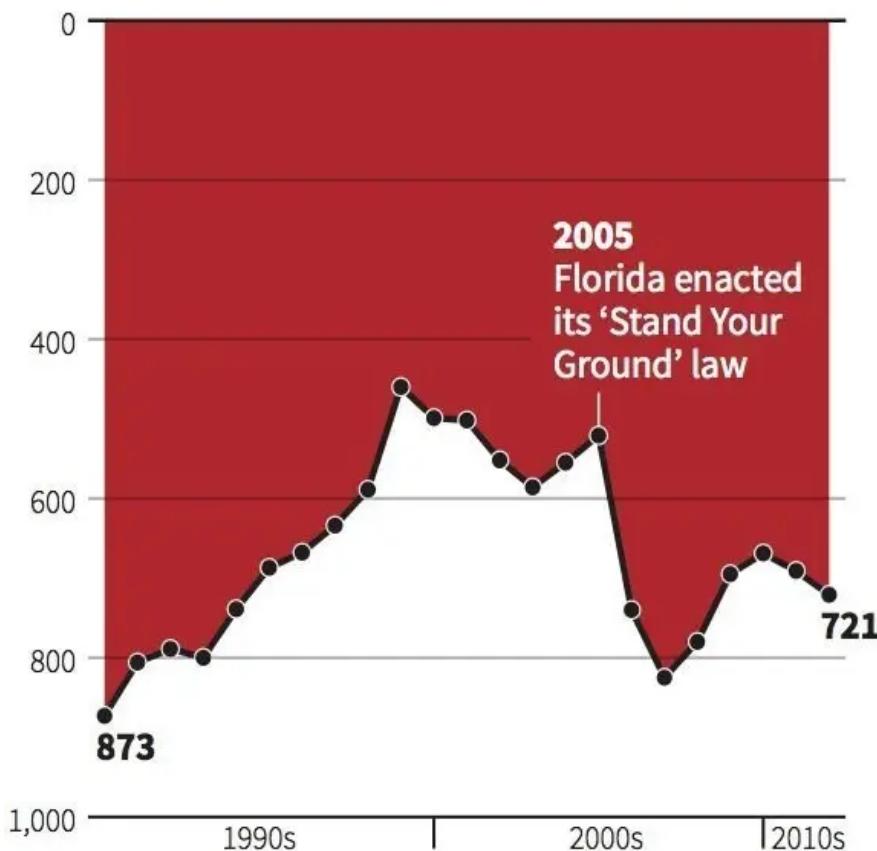
Stacked radar chart



# Ethics

## Gun deaths in Florida

Number of murders committed using firearms



Source: Florida Department of Law Enforcement

C. Chan 16/02/2014

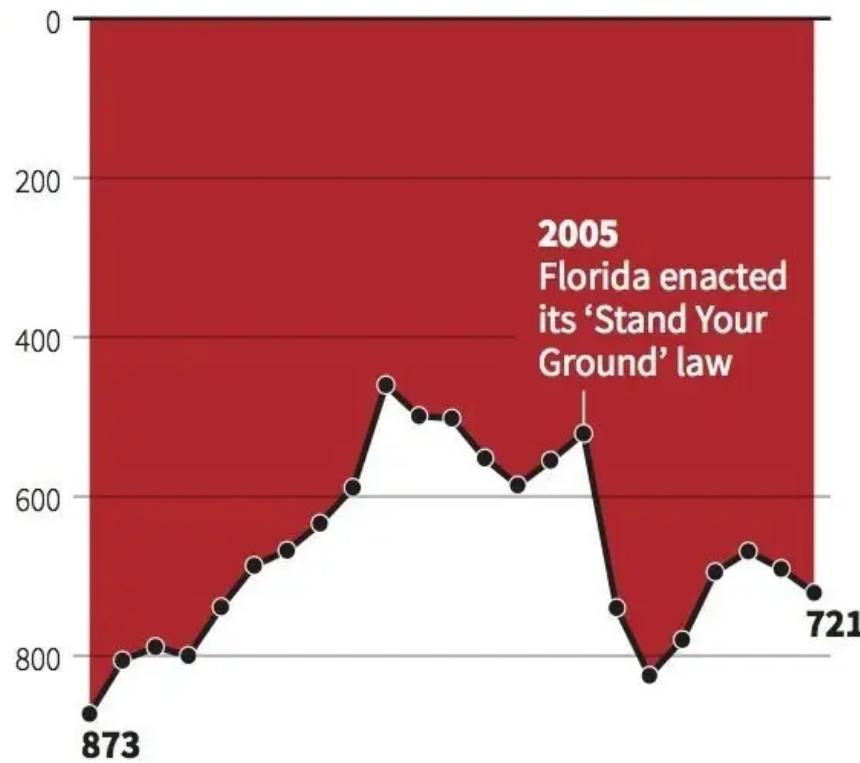
REUTERS

# Ethics

---

## Gun deaths in Florida

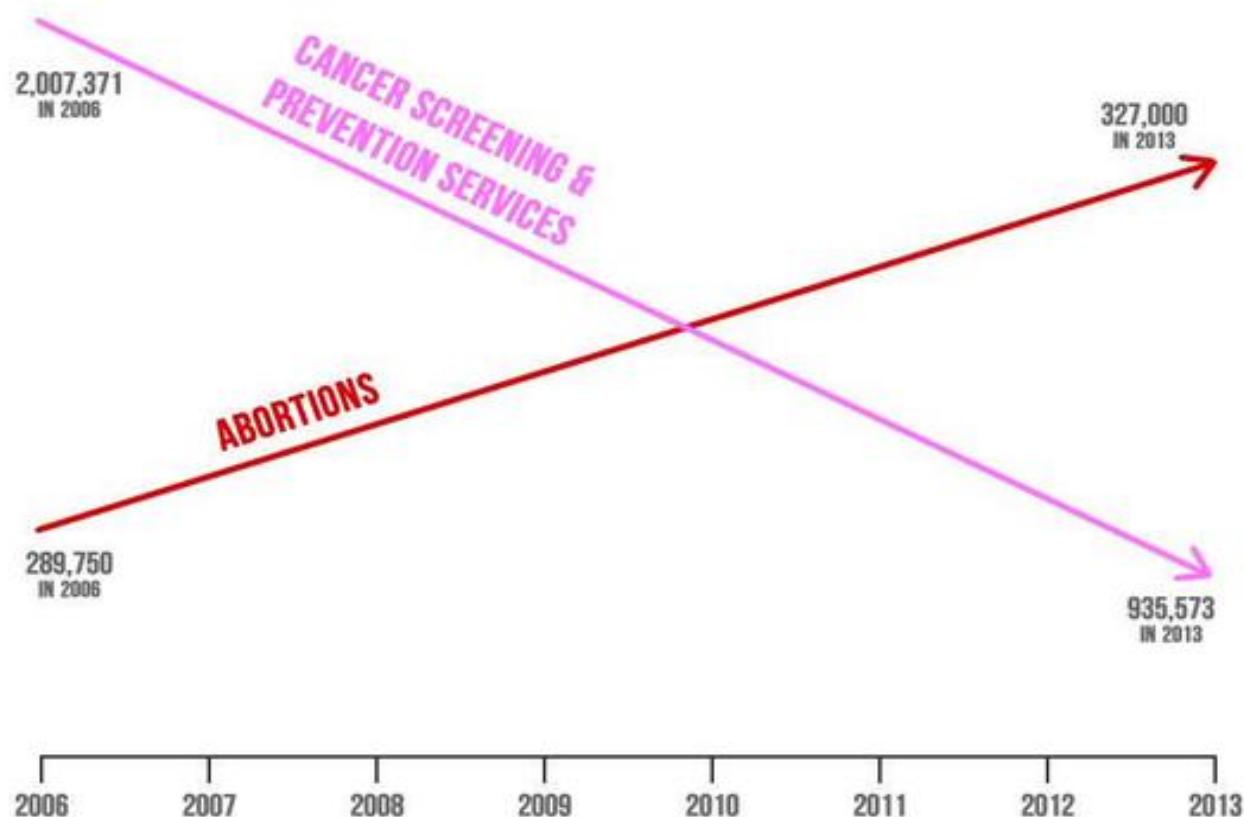
Number of murders committed using firearms



This only has one problem with it..

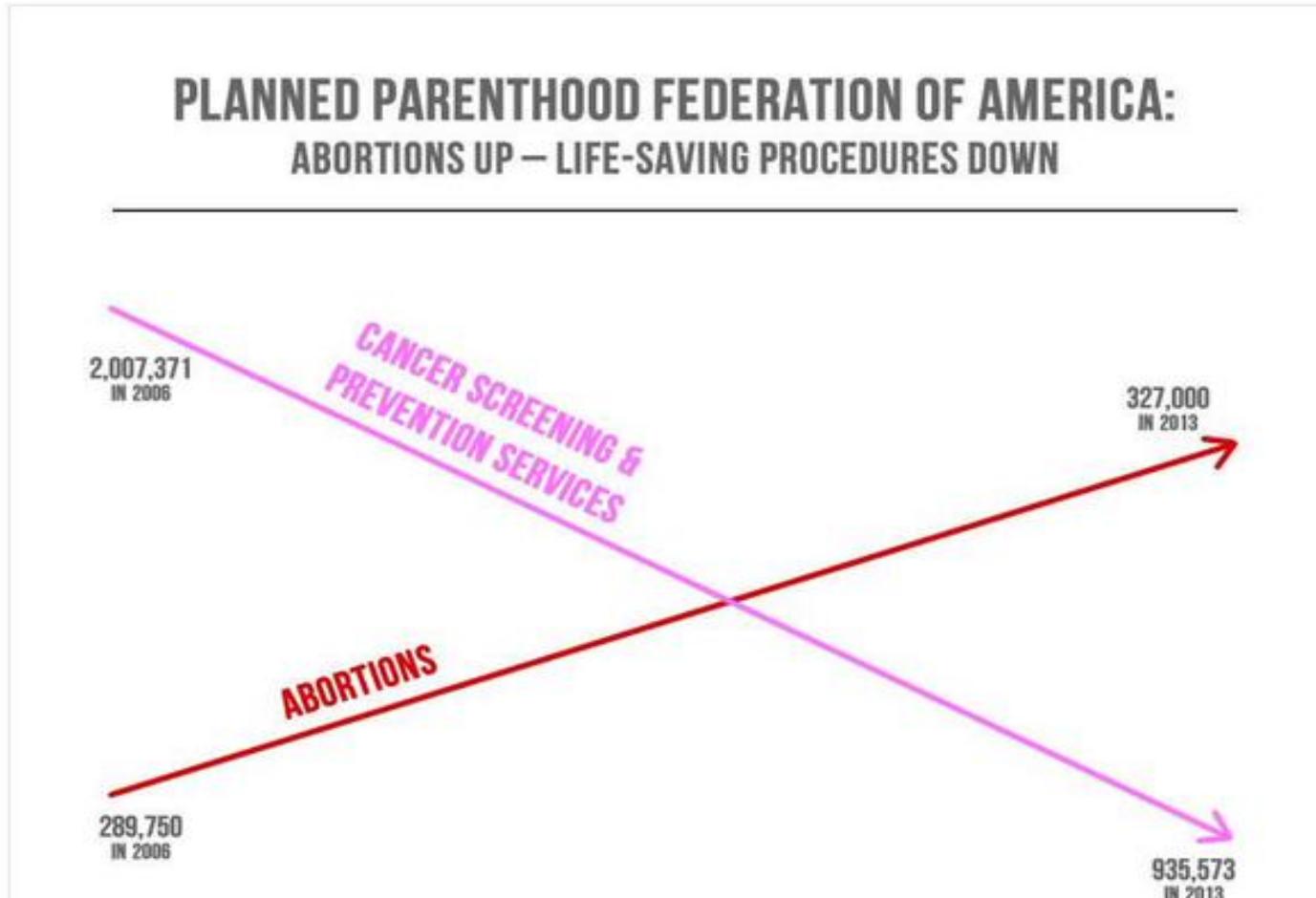
# Ethics

## PLANNED PARENTHOOD FEDERATION OF AMERICA: ABORTIONS UP – LIFE-SAVING PROCEDURES DOWN



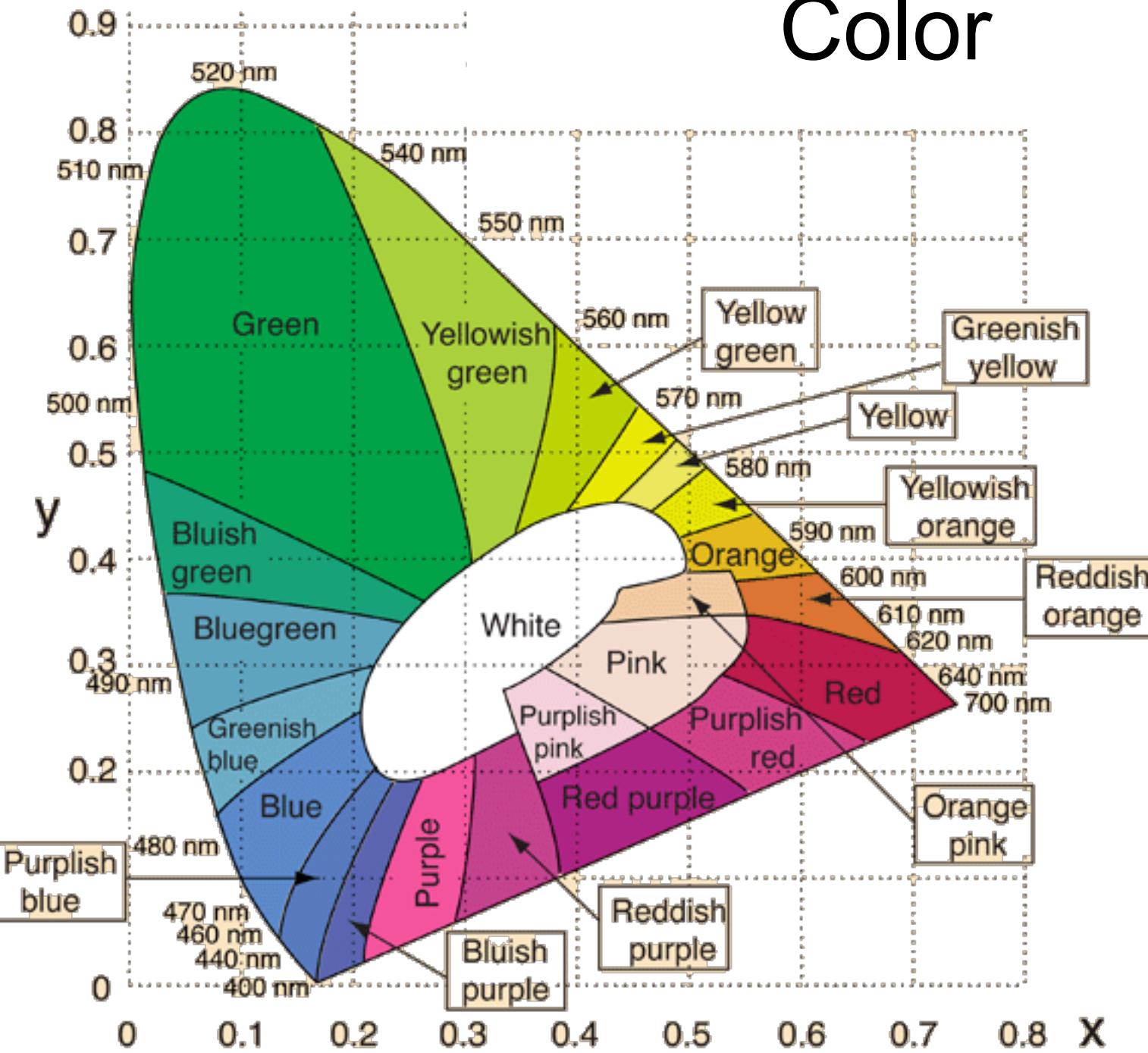
SOURCE: AMERICANS UNITED FOR LIFE

# Ethics

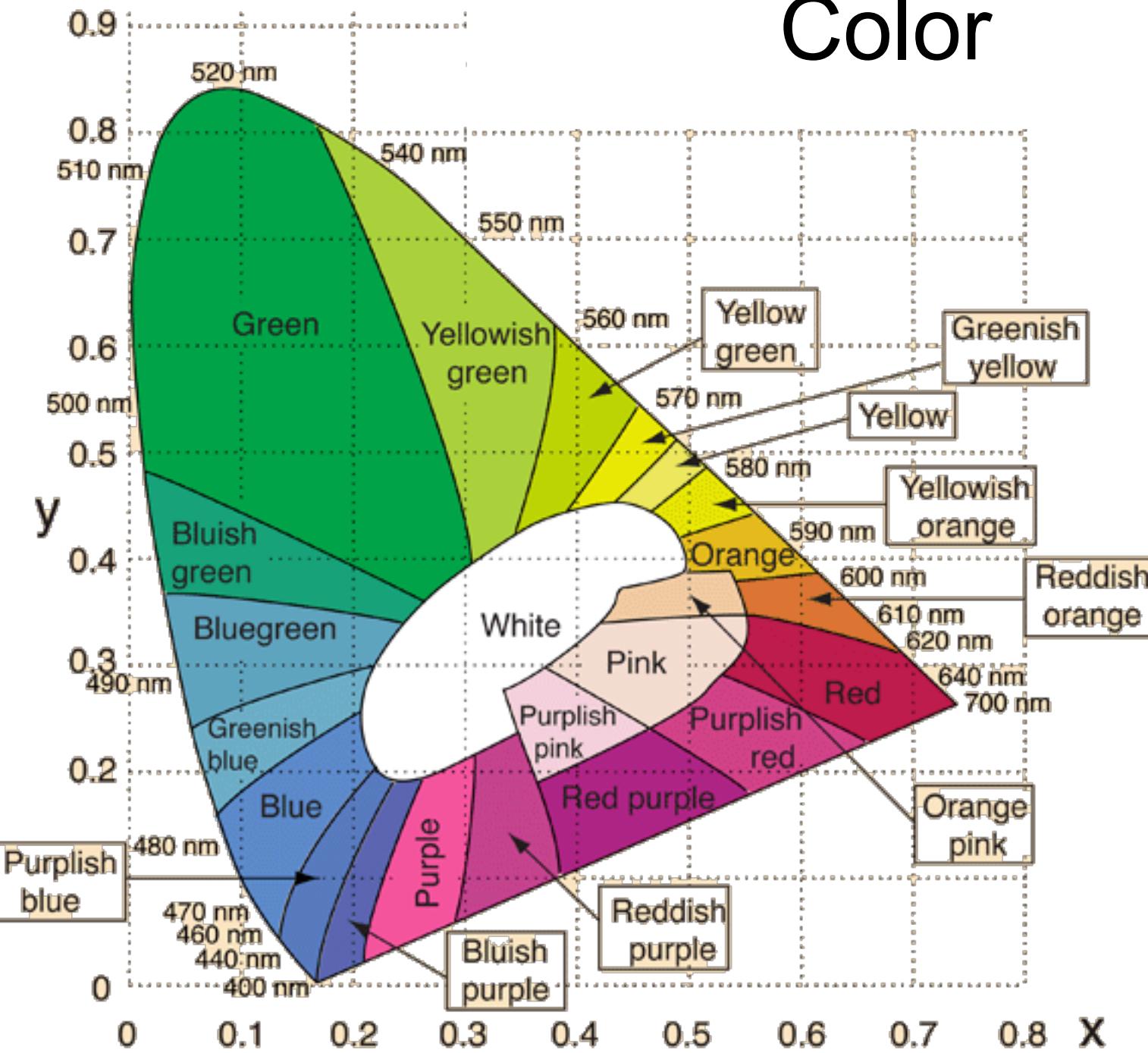


This one has more than one problem –  
using visualization to obscure the truth rather than reveal it

# Color



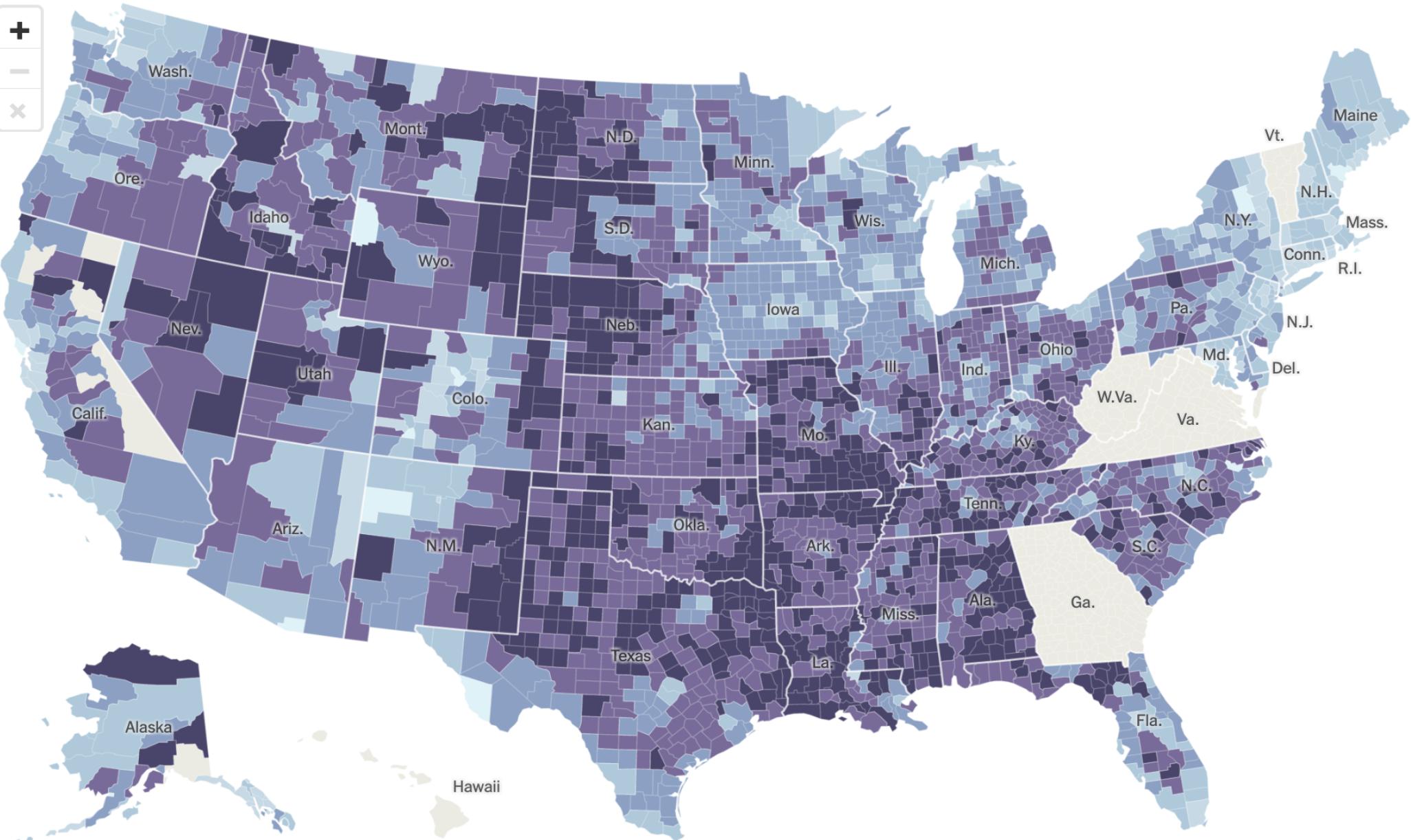
# Color



### Percent of residents who are not fully vaccinated

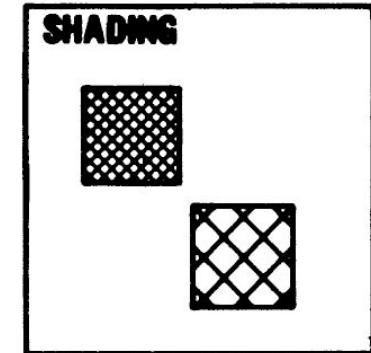
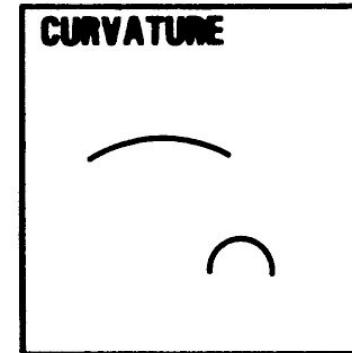
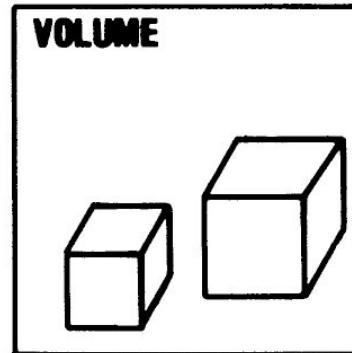
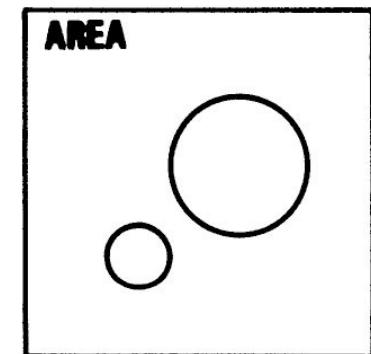
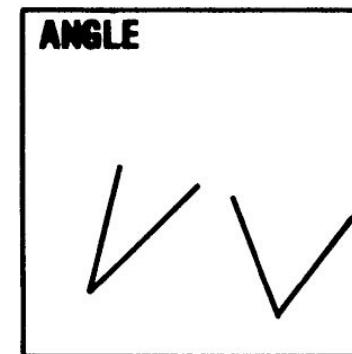
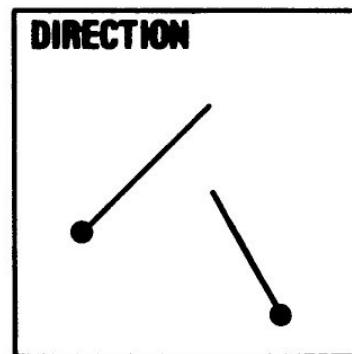
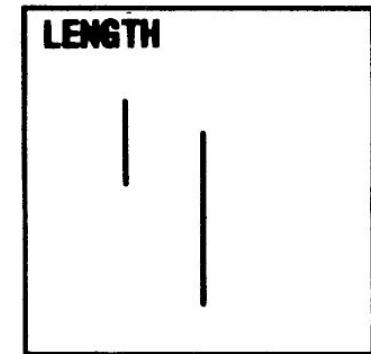
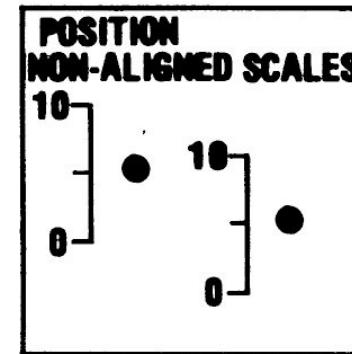
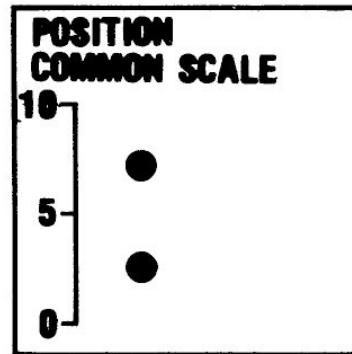
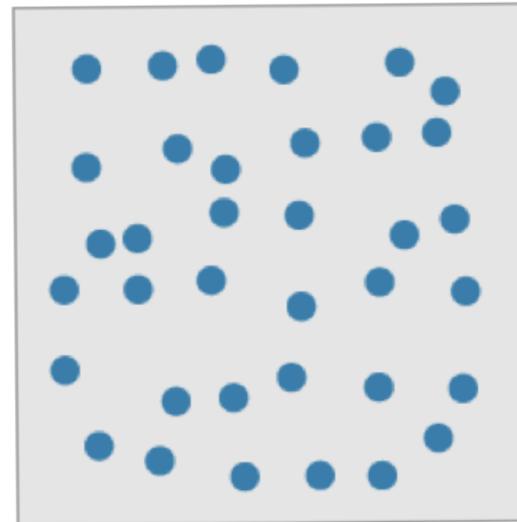
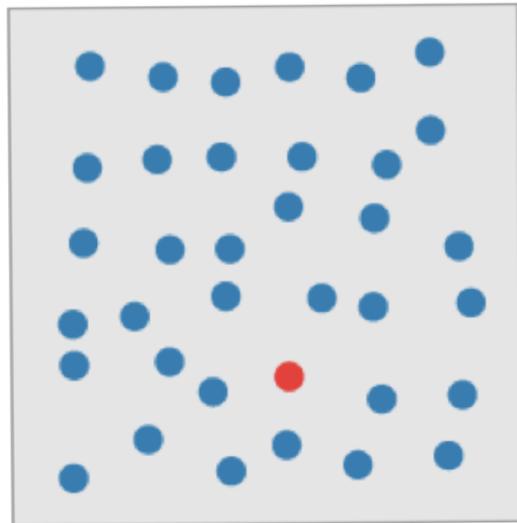
30 40 50 60 70%

Insufficient data



Note: No C.D.C. data available for Hawaii, Texas and some counties. Four additional states were excluded because more than a quarter of the data was missing. Data from Texas and Colorado excludes shots given by most federal agencies. Data is as of July 29. • Sources: Centers for Disease Control and Prevention; Texas Department of State Health Services; Colorado Department of Public Health & Environment; Massachusetts Department of Public Health; U.S. Census Bureau • By Ashley Wu and Albert Sun

# Perception



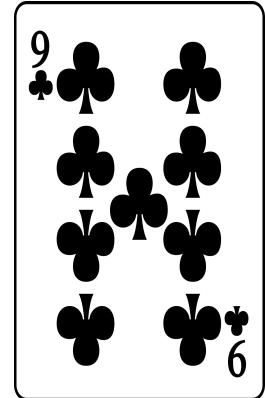
**COLOR SATURATION**

# What is Visualization for?

- Document, record, preserve
- Analyze
  - Tool to help us understand
    - enhance memory
    - take advantage of perceptual preprocessing
  - Reveal peculiarities, limitations of data
- Communicate
  - Inform, persuade
  - Make decisions

# Activity: Number Scrabble

- Two-player game played with cards A-9
- Players alternately choose one card
- First player to collect any three cards which sum to 15 wins
- If all cards are exhausted and no three cards held by one player sum to 15, game ends in a draw.



# Number Scrabble example game

- Player A takes 8
- Player B takes 2
- Player A takes 4
- Player B takes 3
- Player A takes 5
- B to move?

# Number Scrabble example game

- Player A takes 8
- Player B takes 2
- Player A takes 4
- Player B takes 3
- Player A takes 5
- B to move?

How much do you know about arithmetic?

# Number Scrabble example game

- Player A takes 8
- Player B takes 2
- Player A takes 4
- Player B takes 3
- Player A takes 5
- B to move?

A: 8, 4, 5

B: 2, 3

Unclaimed:  
1,6,7,9

# Number Scrabble example game

A: 8, 4, 5

B: 2, 3

Unclaimed:  
1,6,7,9

3 choose 2 = 3 pairs

A can try to play next

# Number Scrabble example game

A: 8, 4, 5

B: 2, 3

Unclaimed:  
1,6,7,9

$$8 + 4 + 3 = 15$$

$$8 + 5 + 2 = 15$$

$$4 + 5 + 6 = 15$$

# Number Scrabble example game

A: 8, 4, 5

B: 2, 3

Unclaimed:  
1,6,7,9

$$8 + 4 + 3 = 15$$

$$8 + 5 + 2 = 15$$

$$4 + 5 + 6 = 15$$

# Number Scrabble example game

A: 8, 4, 5

B: 2, 3

Unclaimed:  
1,6,7,9

$$8 + 4 + 3 = 15$$

$$8 + 5 + 2 = 15$$

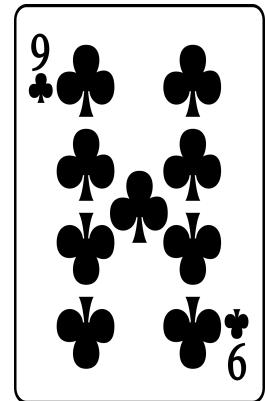
$$4 + 5 + 6 = 15$$

If I color them, B's only non-losing move is evident.

# Activity: Number Scrabble

- There are eight triplets of distinct natural numbers between 1 and 9 that sum to 15.

1	5	9
1	6	8
2	4	9
2	5	8
2	6	7
3	4	8
3	5	7
4	5	6

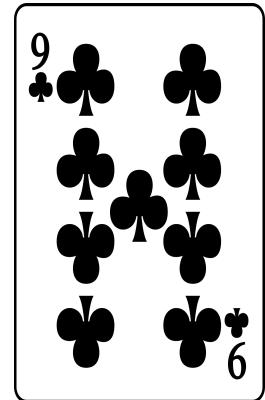


- Four solutions involve 5, and four solutions are imbalanced around 5

# Activity: Number Scrabble

- There are eight triplets of distinct natural numbers between 1 and 9...

1	5	9
1	6	8
2	4	9
2	5	8
2	6	7
3	4	8
3	5	7
4	5	6

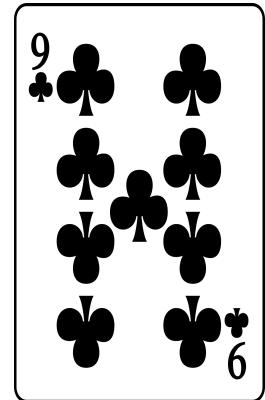


- 1, 3, 9, 7 have only two paths to victory
- 2, 4, 6, 8 have three each
- 5 is in four of the solutions

# Activity: Number Scrabble

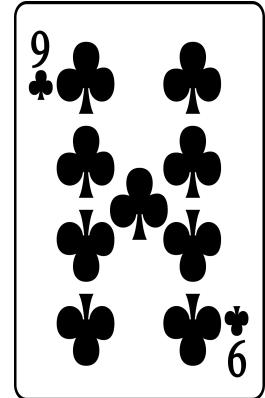
- There are eight triplets of distinct natural numbers between 1 and 9...

1	5	9
1	6	8
2	4	9
2	5	8
2	6	7
3	4	8
3	5	7
4	5	6



- **1, 3, 9, 7** have only two paths to victory
- **2, 4, 6, 8** have three each
- **5** is in four of the solutions

# Activity: Number Scrabble



- 1, 3, 9, 7 have only two paths to victory
- 2, 4, 6, 8 have three each
- 5 is in four of the solutions

# Number Scrabble example game

A: 8, 4, 5

B: 2, 3

Unclaimed:  
1,6,7,9

6	7	2
1	5	9
8	3	4

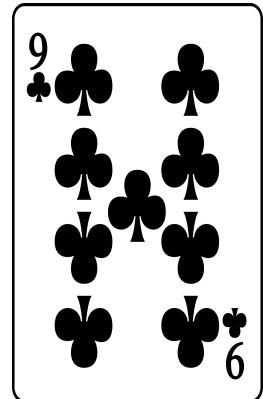
$$8 + 4 + 3 = 15$$

$$8 + 5 + 2 = 15$$

$$4 + 5 + 6 = 15$$

# Activity: Number Scrabble

6	7	2
1	5	9
8	3	4



- Writing down the cards helped us see what was done, what was missing
- Coloring the numbers helped us see some group theory properties
- The tic-tac-toe representation turns the problem into a solved one.

# Visualization tools



- Ubiquitous, limited to keypresses and mouse clicks
- Excellent for interacting with primates / getting data and judgements out of your head



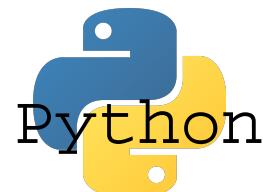
- Commercial, free license for students, entirely graphical and purported to be easy to use.



Google refine

- Browser-based free data wrangling power tool; find outliers, perform transforms on columns. Won't need for this class

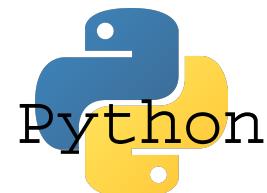
# The stack



L A P A C K  
L -A P -A C -K  
L A P A -C -K  
L -A P -A -C K

Stratigraphic layers -- lower layers are older, will likely last longer into the future.

# The stack



L A P A C K  
L -A P -A C -K  
L A P A -C -K  
L -A P -A -C K

# The stack



Vega-Lite

matplotlib

pandas

Vega

NumPy



Python

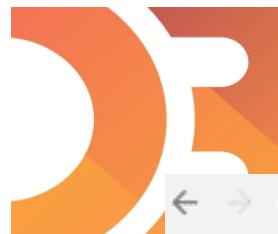
L A P A C K  
L -A P -A C -K  
L A P A -C -K  
L -A P -A -C K

JS



```
<!doctype html>
<html>
<head>
  <style>
    .bar {
      fill: steelblue;
    }
  </style>
  <script src="https://d3js.org/d3.v4.min.js"></script>
<body>
<svg width="600" height="500"></svg>
<script>
  var svg = d3.select("svg"),
    margin = 200,
    width = svg.attr("width") - margin,
    height = svg.attr("height") - margin
  svg.append("text")
    .attr("transform", "translate(100,0)")
    .attr("x", 50)
    .attr("y", 50)
    .attr("font-size", "24px")
    .text("XYZ Foods Stock Price")
  var xScale = d3.scaleBand().range([0,
width]).padding(0.4),
    yScale = d3.scaleLinear().range([height, 0]);
  var g = svg.append("g")
    .attr("transform", "translate(" + 100 + "," +
+ 100 + ")");
  d3.csv("XYZ.csv", function(error, data) {
    if (error) {
      throw error;
    }
    xScale.domain(data.map(function(d) { return
d.year; }));
    yScale.domain([0, d3.max(data, function(d)
{ return d.value; })]);
    g.append("g")
      .attr("transform", "translate(0," + height + ")")
      .call(d3.axisBottom(xScale))
      .append("text")
        .attr("y", height - 250)
        .attr("x", width - 100)
        .attr("text-anchor", "end")
        .attr("stroke", "black")
        .text("Year");
      g.append("g")
        .call(d3.axisLeft(yScale).tickFormat(function(d) {
          return "$" + d;
        }))
        .ticks(10)
        .append("text")
        .attr("dx", -5)
        .attr("dy", 5)
        .attr("text-anchor", "right")
        .text("Value");
  });
</script>
</body>
</html>
```

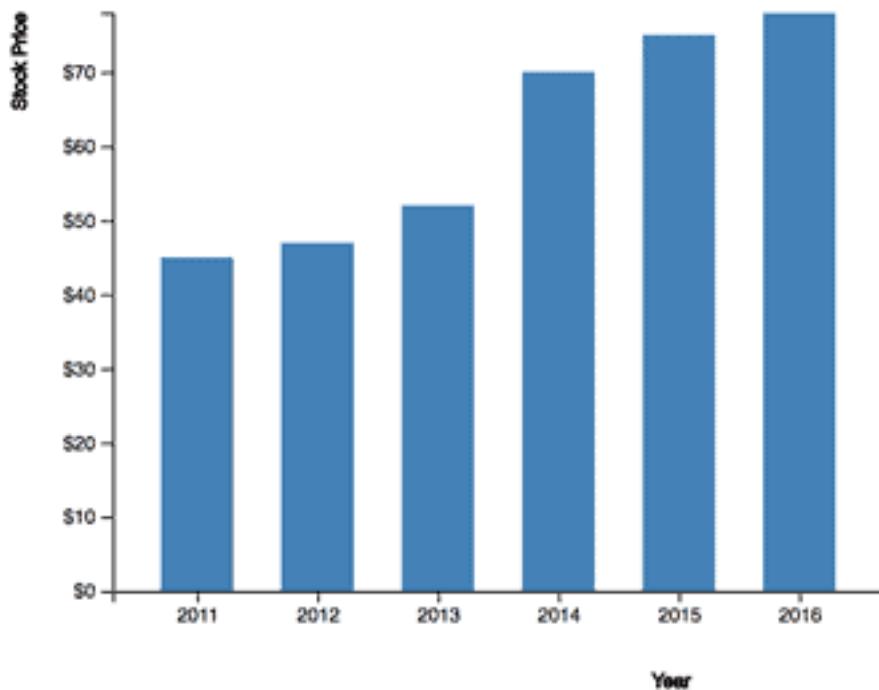
This is the sample code for a bar graph



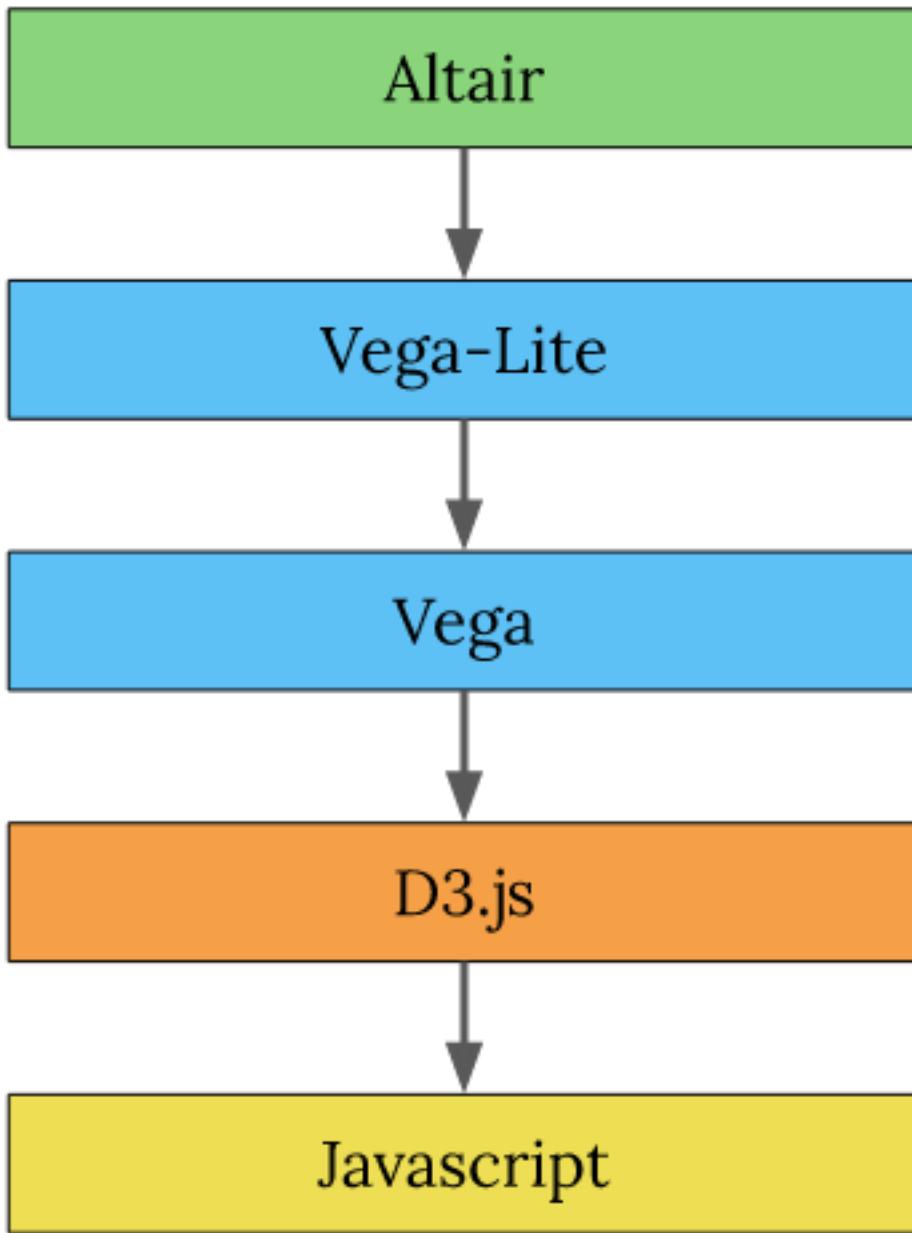
localhost:63342/d3\_tutorial/stockprice.html



## XYZ Foods Stock Price

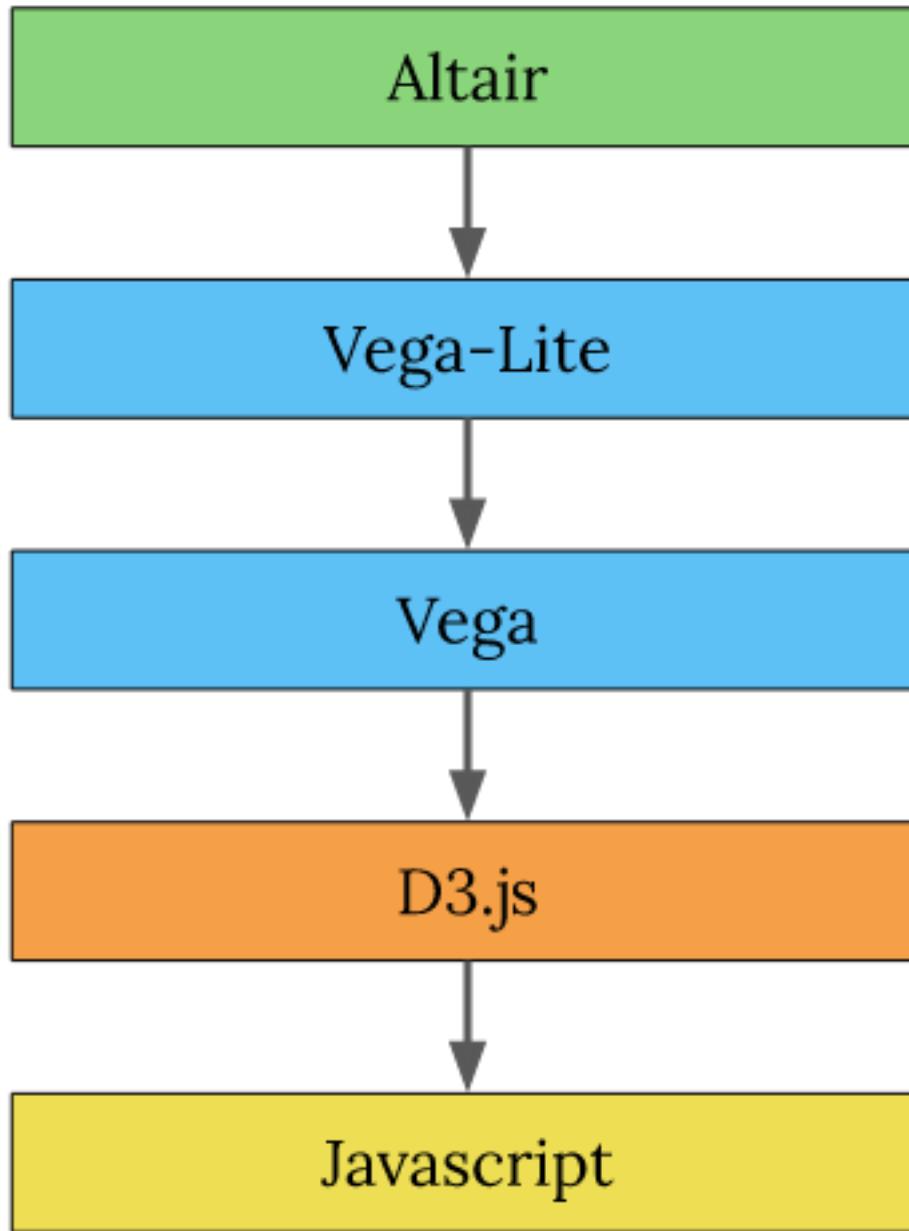


```
<!doctype html>
<html>
<head>
  <style>
    .bar {
      width: 100px;
      height: 100px;
      background-color: #3366CC;
      margin: 10px auto;
    }
  </style>
  <script src="https://d3js.org/d3.v3.min.js"></script>
<body>
<div style="text-align: center;">
  <h2>XYZ Foods Stock Price</h2>
  <div>
    <img alt="Bar chart showing XYZ Foods Stock Price from 2011 to 2016." data-bbox="147 294 562 729"/>
    <table border="1" data-bbox="147 729 562 942">
      <thead>
        <tr>
          <th>Year</th>
          <th>Stock Price</th>
        </tr>
      <tbody>
        <tr>
          <td>2011</td>
          <td>$45</td>
        </tr>
        <tr>
          <td>2012</td>
          <td>$48</td>
        </tr>
        <tr>
          <td>2013</td>
          <td>$52</td>
        </tr>
        <tr>
          <td>2014</td>
          <td>$68</td>
        </tr>
        <tr>
          <td>2015</td>
          <td>$72</td>
        </tr>
        <tr>
          <td>2016</td>
          <td>$75</td>
        </tr>
      </tbody>
    </table>
  </div>
  <script>
    var data = [
      {"year": "2011", "value": 45}, {"year": "2012", "value": 48}, {"year": "2013", "value": 52}, {"year": "2014", "value": 68}, {"year": "2015", "value": 72}, {"year": "2016", "value": 75}
    ];
    var width = 100;
    var height = 100;
    var xScale = d3.scale.ordinal()
      .domain(data.map(function(d) { return d.year; }));
    var yScale = d3.scale.linear()
      .domain([0, d3.max(data, function(d) { return d.value; })])
      .range([0, height]);
    var g = d3.select("div").append("div")
      .attr("width", width)
      .attr("height", height);
    g.append("p").text("Stock Price");
    g.append("p").text("Year");
    g.selectAll("bar")
      .data(data)
      .enter().append("div")
      .attr("width", width)
      .attr("height", function(d) { return height - yScale(d.value); })
      .attr("x", function(d) { return xScale(d.year); })
      .attr("y", function(d) { return yScale(d.value); })
      .attr("fill", "#3366CC");
  </script>
</div>
</body>
</html>
```



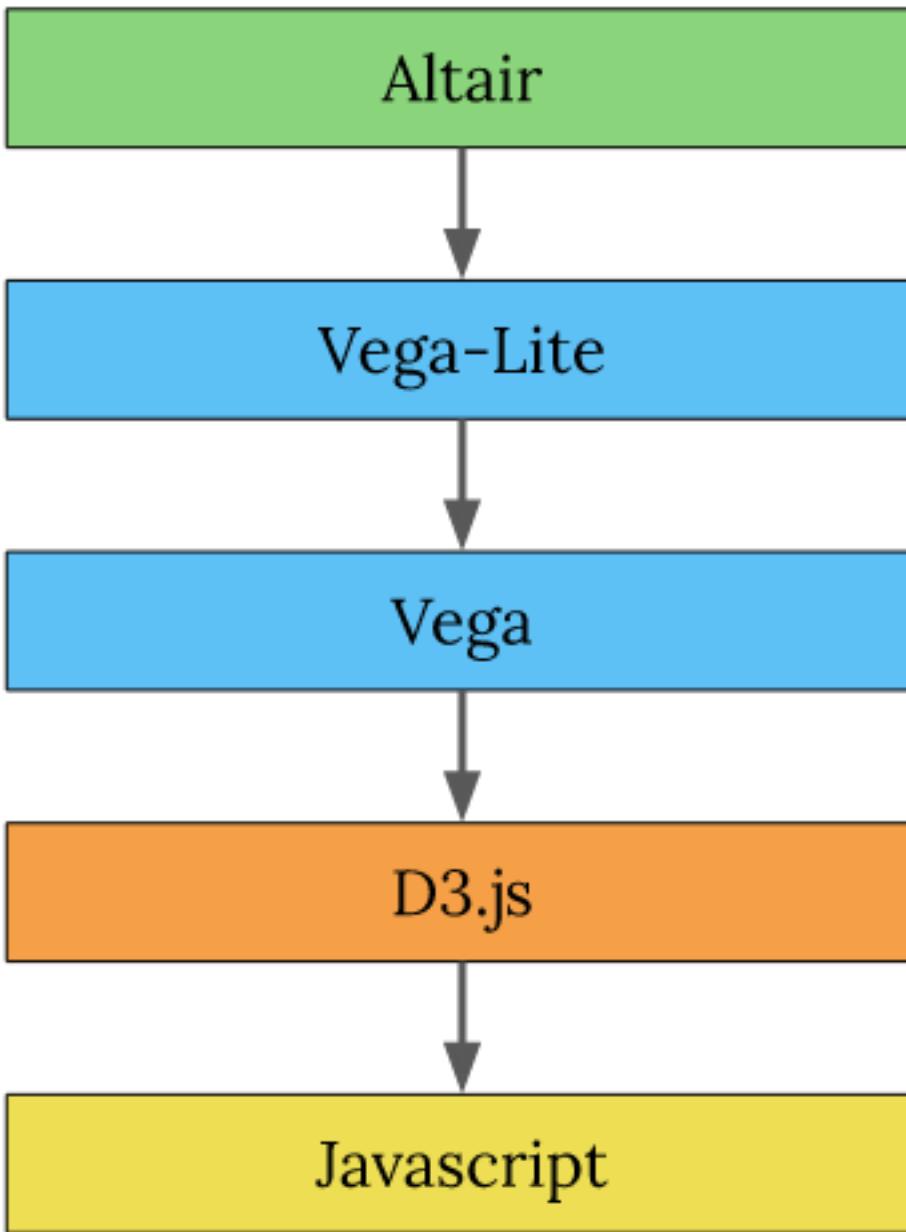
```
FacetedEncoding({  
    x: X({  
        field: 'Year',  
        type: 'temporal'  
    }),  
    y: Y({  
        aggregate: 'count',  
        type: 'quantitative'  
    })  
})
```

This is a specification for a bar graph in Vega-lite

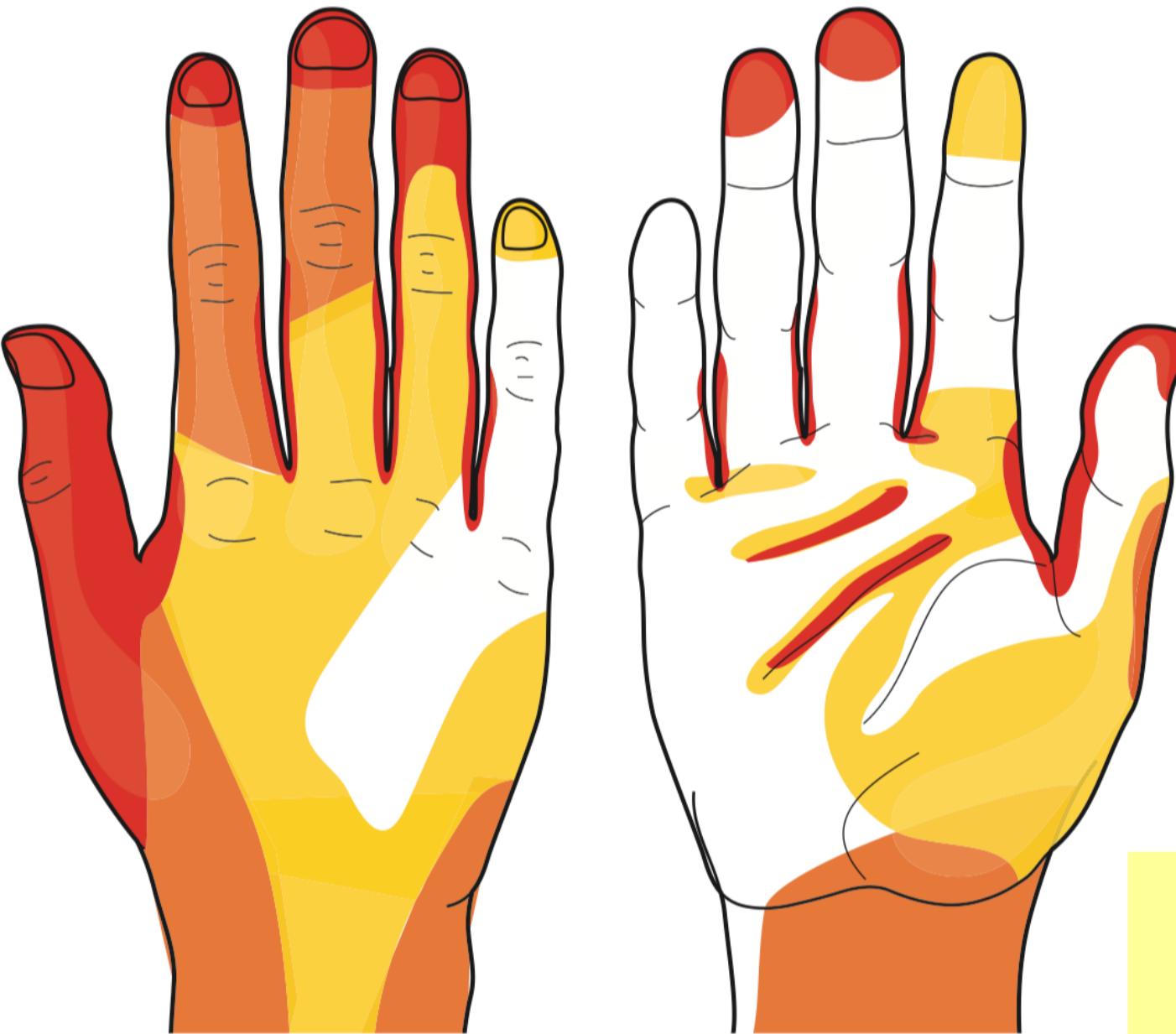


```
alt.Chart(cars).mark_bar()  
.encode(  
    alt.X('Year'),  
    alt.Y('count()'))
```

And here is the specification in altair (python wrapper for vega-lite)



If I find out that altair/vega-lite really isn't working out by 2nd week we can move back to pure python.



MOST OFTEN  
MISSED AREAS

OFTEN MISSED  
AREAS

LESS OFTEN  
MISSED AREAS

Beautiful, easy-to-understand graphic for these interesting times.

# Administrivia

- Office hours
  - WT: 1-2:30 Wednesday,
  - WT: 11-12 Thursday, Friday; other times by appointment
  - AM:
  - AS:
- 7 homeworks, due on Friday 11:59 pm. First one friday of 2nd week.

# Teaching Staff

Teaching Assistants



Anna Moise



Abby Star

Professor



Will Trimble  
Instructor, Data Science  
W 1-2:30  
Th 11-212  
F 11-12

# This and that

TAs:

Abby Starr <[abigailstarr@uchicago.edu](mailto:abigailstarr@uchicago.edu)>

Anna Moise <[amoise@uchicago.edu](mailto:amoise@uchicago.edu)>

- 10 minutes to pandas which takes way more than 10 minutes:
- [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/10min.html](https://pandas.pydata.org/pandas-docs/stable/user_guide/10min.html)