NREL
Transforming ENERGY

# A Picture is Worth a Thousand Data Points: Introduction to Visualization

**Andy Berres** (they/them)
National Renewable Energy Laboratory

ACM Tapia Conference 2023
September 14, 2023

Photo by Dennis Schroeder, NREL 55200

# Quick Bio

- BS, MS, PhD Computer Science
  - Minors: Mathematics (BS)/Biology (MS)
  - Research: Visualization, Geometry, Topology

- These days: Applied Data Science
  - Climate
  - Traffic
  - Buildings
  - Power Grid

# Contents

**1** **Why should you care**

**2** **What not to do**

**3** **What to do instead**

**4** **Hands-on practice**

**5** **Resources**

# Why should you care?

| x | y |
|---|---|
| 40 | 99.859 |
| 50 | 99.859 |
| 61.2821 | 91.3974 |
| 69.2308 | 79.0897 |
| 76.4103 | 59.859 |
| 75.1282 | 43.3205 |
| 70.5128 | 26.0128 |
| 60.5128 | 13.3205 |
| 50 | 8.7051 |
| 37.4359 | 6.3974 |
| 28.7179 | 14.859 |
| 22.0513 | 26.3974 |
| 16.4103 | 47.9359 |
| 16.6667 | 64.859 |
| 19.7436 | 80.2436 |
| 28.4615 | 92.1667 |
| 31.7949 | 69.859 |
| 54.359 | 70.6282 |
| 29.2308 | 49.859 |
| 31.5385 | 41.7821 |
| 33.8462 | 35.2436 |
| 40.2564 | 28.7051 |
| 49.2308 | 28.7051 |
| 57.4359 | 35.2436 |
| 59.4872 | 45.6282 |
| 53.3333 | 31.7821 |
| 44.8718 | 28.3205 |
| 60.2564 | 52.5513 |
| 17.9487 | 36.3974 |
| 44.1026 | 6.3974 |
| 72.8205 | 71.0128 |

What do you think this might be?



Much nicer!

# Showcase 2: Datasaurus Dozen



Matejka, Justin, and George Fitzmaurice. "Same stats, different graphs: generating datasets with varied appearance and identical statistics through simulated annealing." In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, pp. 1290-1294. ACM, 2017.

# Why Should You Care? – Recap

*Visualization will help you understand your data better.*

*Visualization gives you information at a glance.*

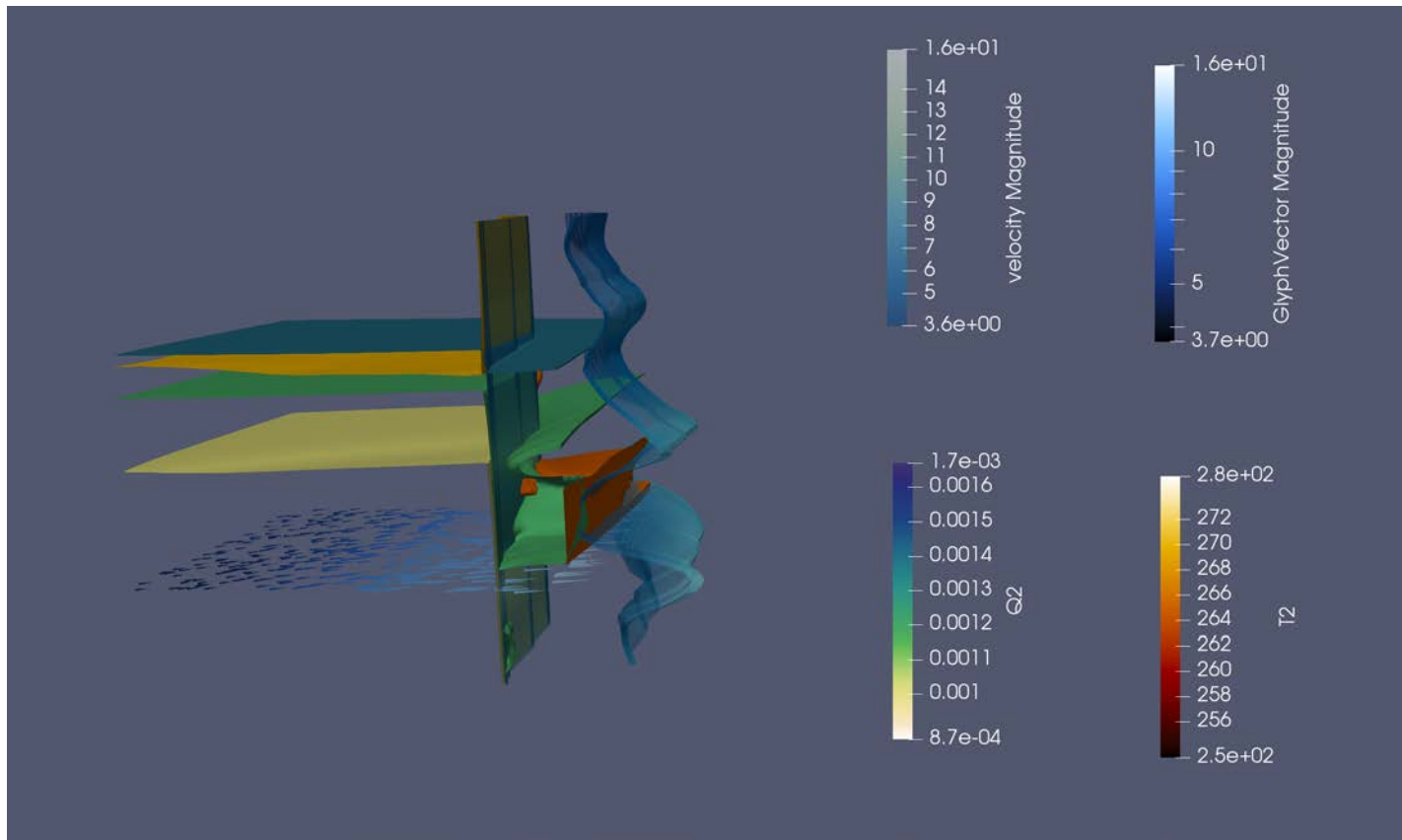*Visualization is a form of communication.*

# What not to do

- The many ways you can mess things up

- Don't try this at home (except do)
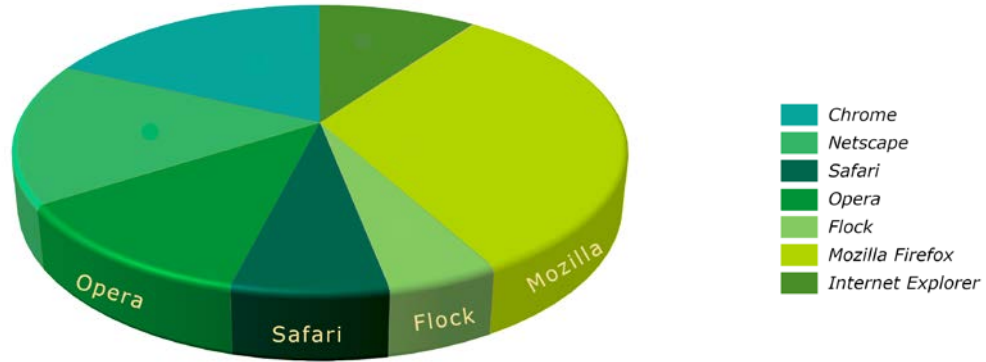
- A game of "What's wrong with this?"

# Disclaimer

- Examples shown serve illustrative purposes.
  - Some of them are actually pretty good
  - Some commit a whole list of sins
  - Some of them are used to advertise tool features
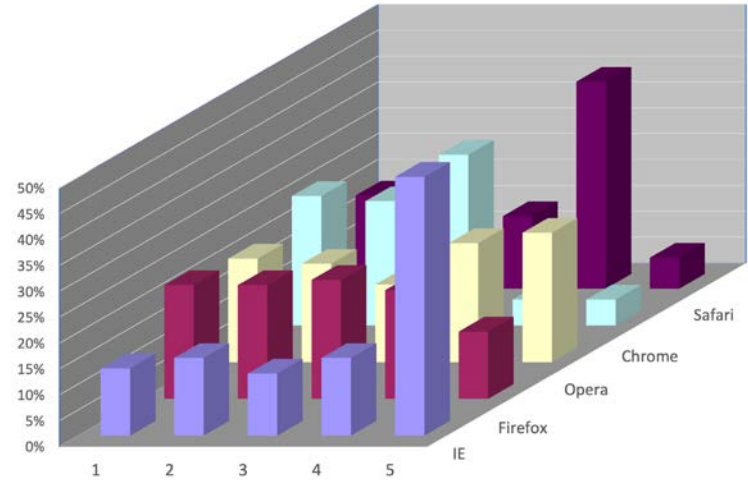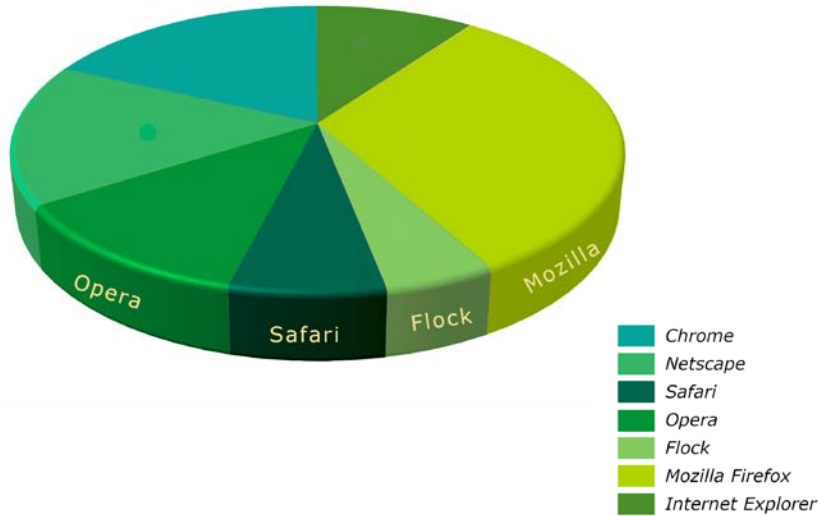
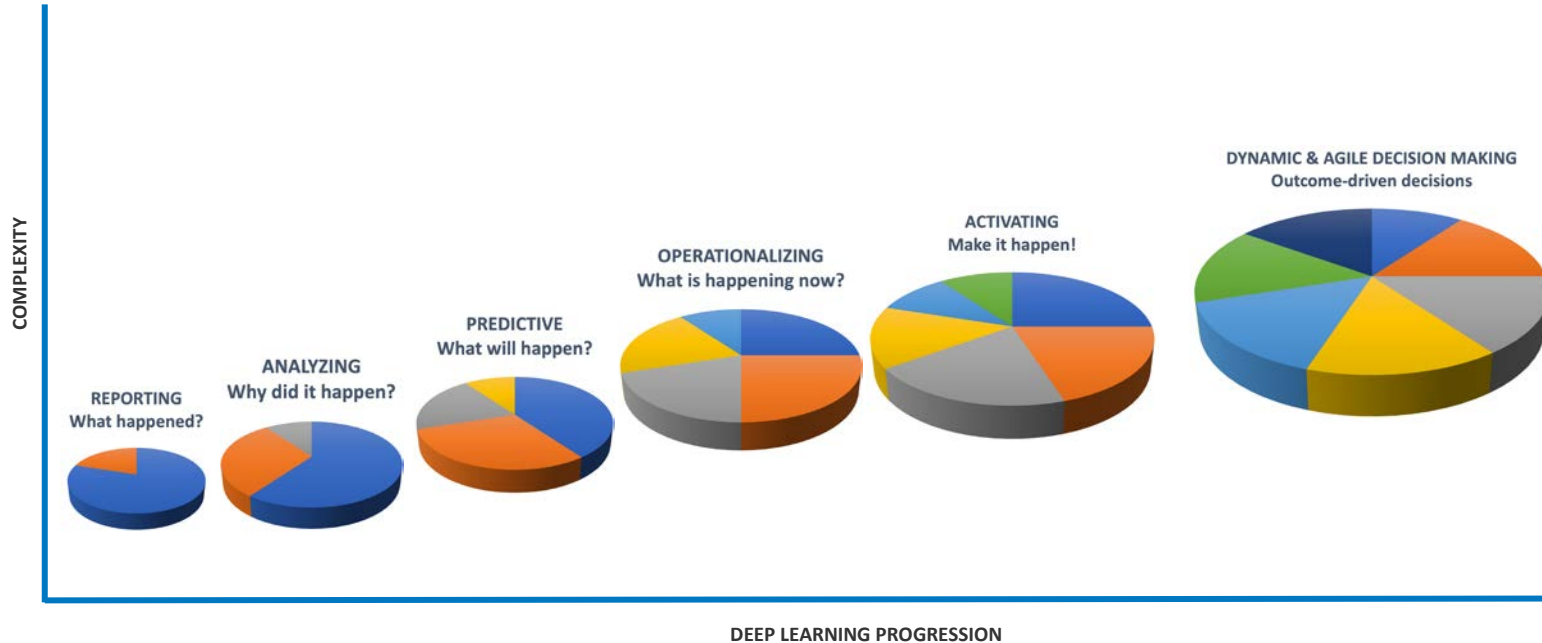# Too much information

# Inappropriate Color Choice

# Useless 3D

- Humans are phenomenally bad at comparing volumes

- 3D often leads to occlusion

# Lack of Context

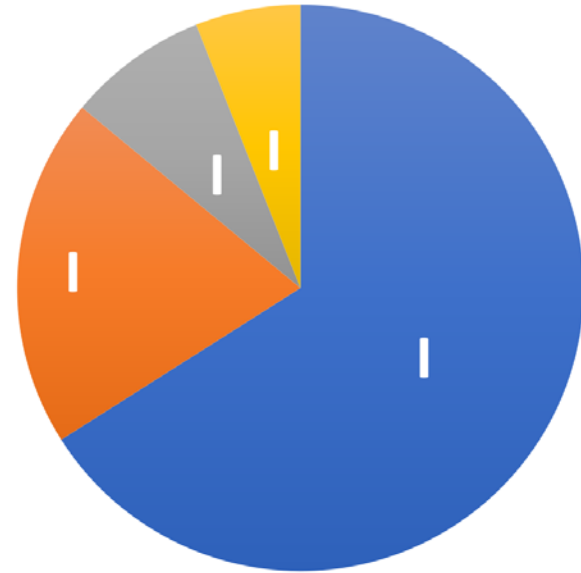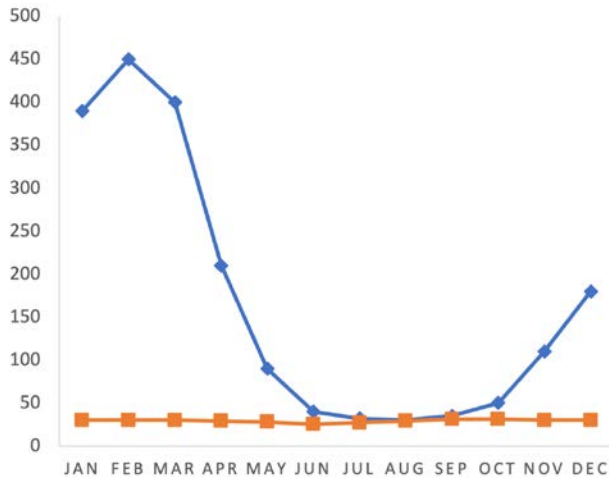- Missing/useless labels

- Missing axis ticks

- Missing legend

- Missing caption
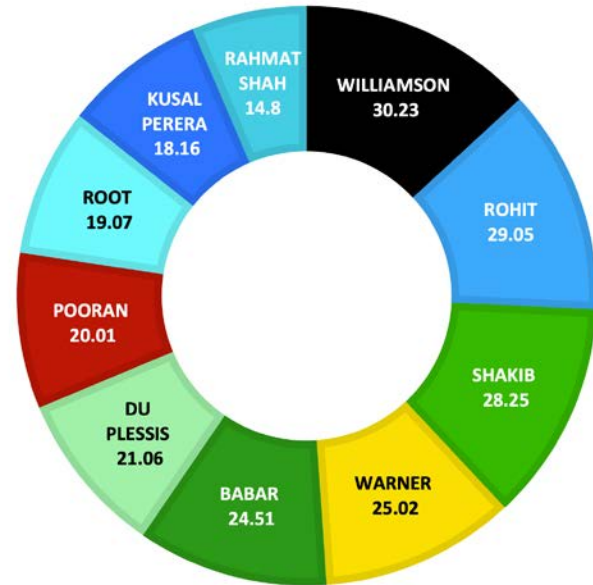
**General US Population**

# Inappropriate Representation



**TEMPERATURE & RAIN CHART**

Can you tell the temperature?

**WORLD'S TOP CRICKET SCORERS**
% OF TEAM'S RUNS SCORED BY TOP SCORER

WILLIAMSON 30.23
ROHIT 29.05
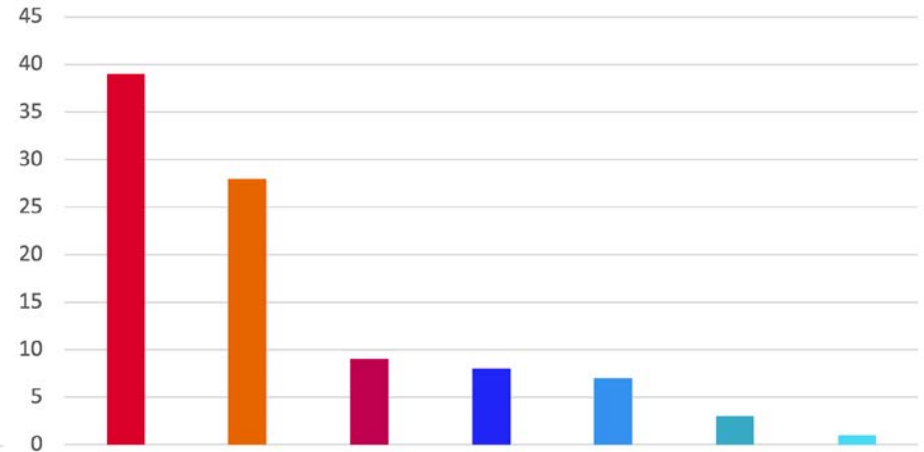SHAKIB 28.25
WARNER 25.02
BABAR 24.51
DU PLESSIS 21.06
POORAN 20.01
ROOT 19.07
KUSAL PERERA 18.16
RAHMAT SHAH 14.8

Numbers add up to more than 100%!

# Misleading Information

Election Results (parties redacted)



39%    28%    9%    8%    7%    3%    1%

Election Results (actual values)



- Axis offset
- Wrong scale

- Choice of projection
- Choice of colormap

- No reproducibility

- Error-prone
  - What settings did you use?
  - Where did you store the data?
  - Which of the 5 versions is it, really?

- What if your data changed? Or you have new data?
  - "Oops, there was an error in my spreadsheet!"
  - "This chart is great! Can you make one for each of the 300 intersections?"
  - "What exact zoom level did I use for my screenshot?"

# What not to do-- Recap

- Information Overload

- Inappropriate Color Choice

- Useless 3D

- Lack of Context

- Inappropriate Data Representation

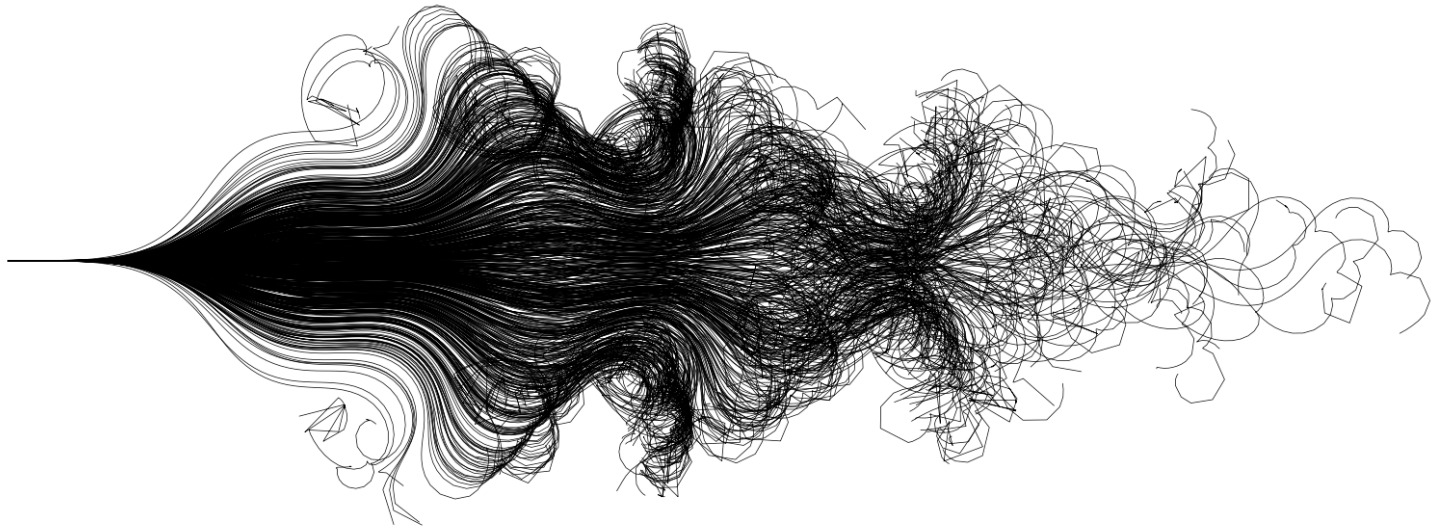- Misleading Information

- Manual Labor

# What to do instead

- Know Your Audience!
- Tell a Story
- Human Factor
- Follow the Rules
- Break the Rules
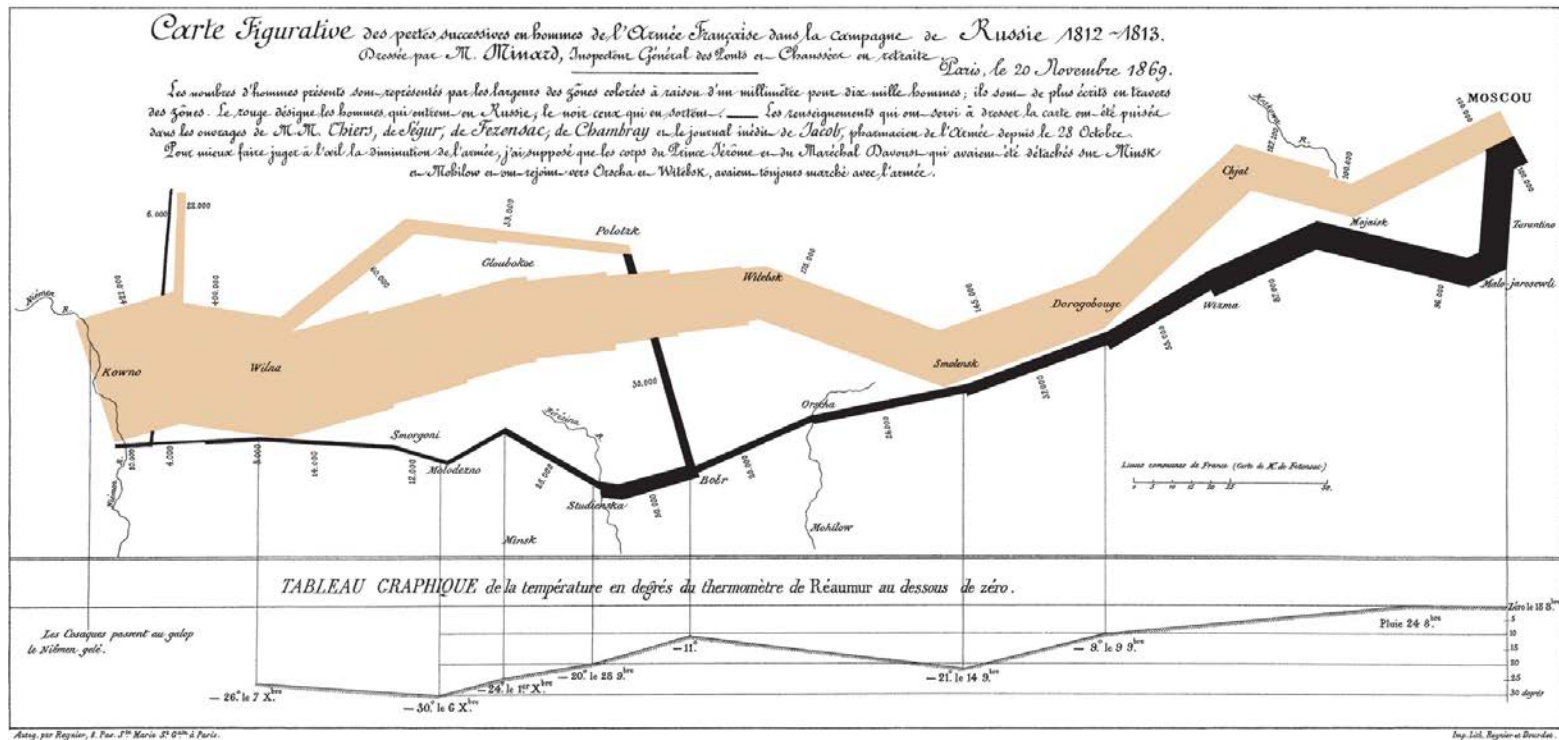
# Know Your Audience

- Objective of the visualization
  - High-level overview (infographic)
  - Expert user (more options, more complexity)

- Visualization literacy
  - Understand which types of visualization will work for your audience – and which types won't
  - For a general audience, aim for simplicity and known concepts
  - For an expert audience (daily interaction), your visualization can be more complex

- Presentation format
  - Paper vs slides vs interactive

Cook, Matthew. "It takes two neurons to ride a bicycle." *Demonstration at NIPS* 4 (2004).
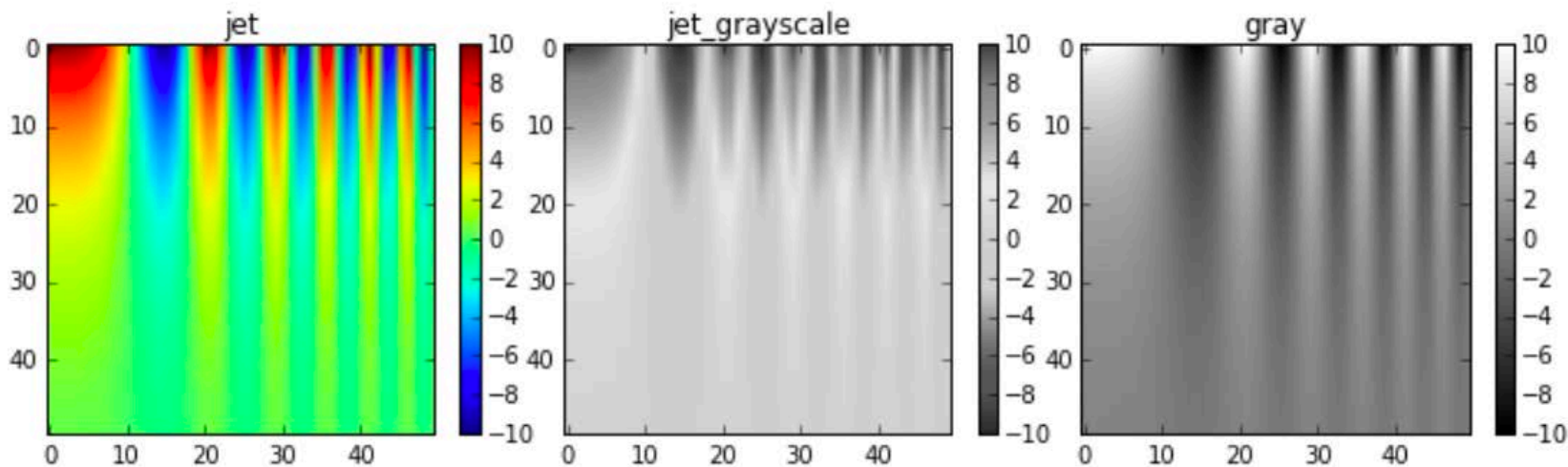
# Tell a Story



Charles Joseph Minard (1869): Napoleon's failed attempt at conquering Russia

# Human Factor: Perception

- Not all colormaps are perceptually smooth
- To check how your colormap is doing, convert to grayscale and compare with a grayscale map



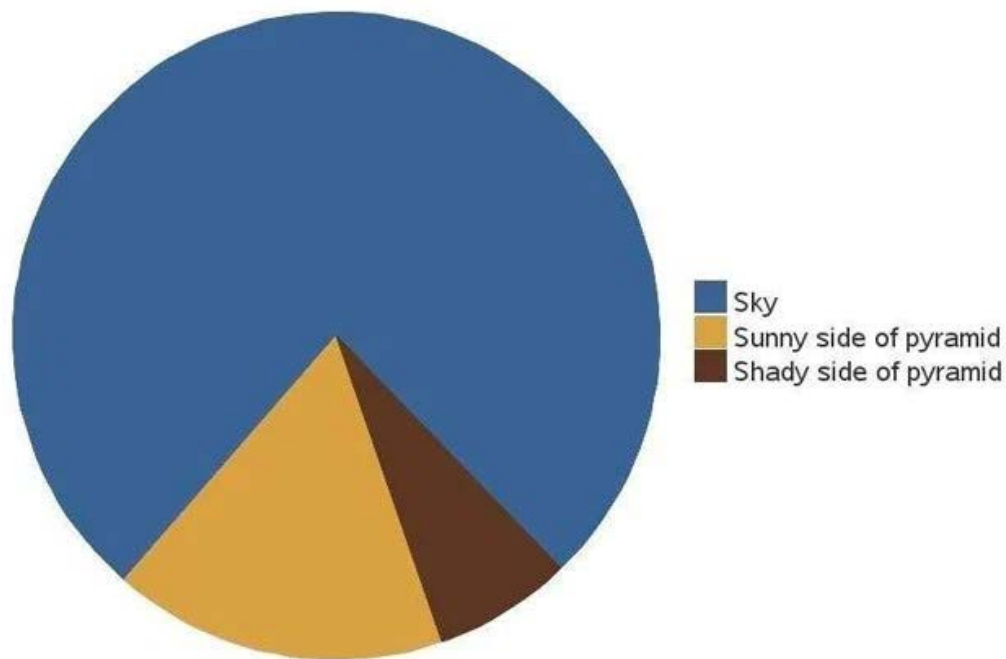https://jakevdp.github.io/blog/2014/10/16/how-bad-is-your-colormap/
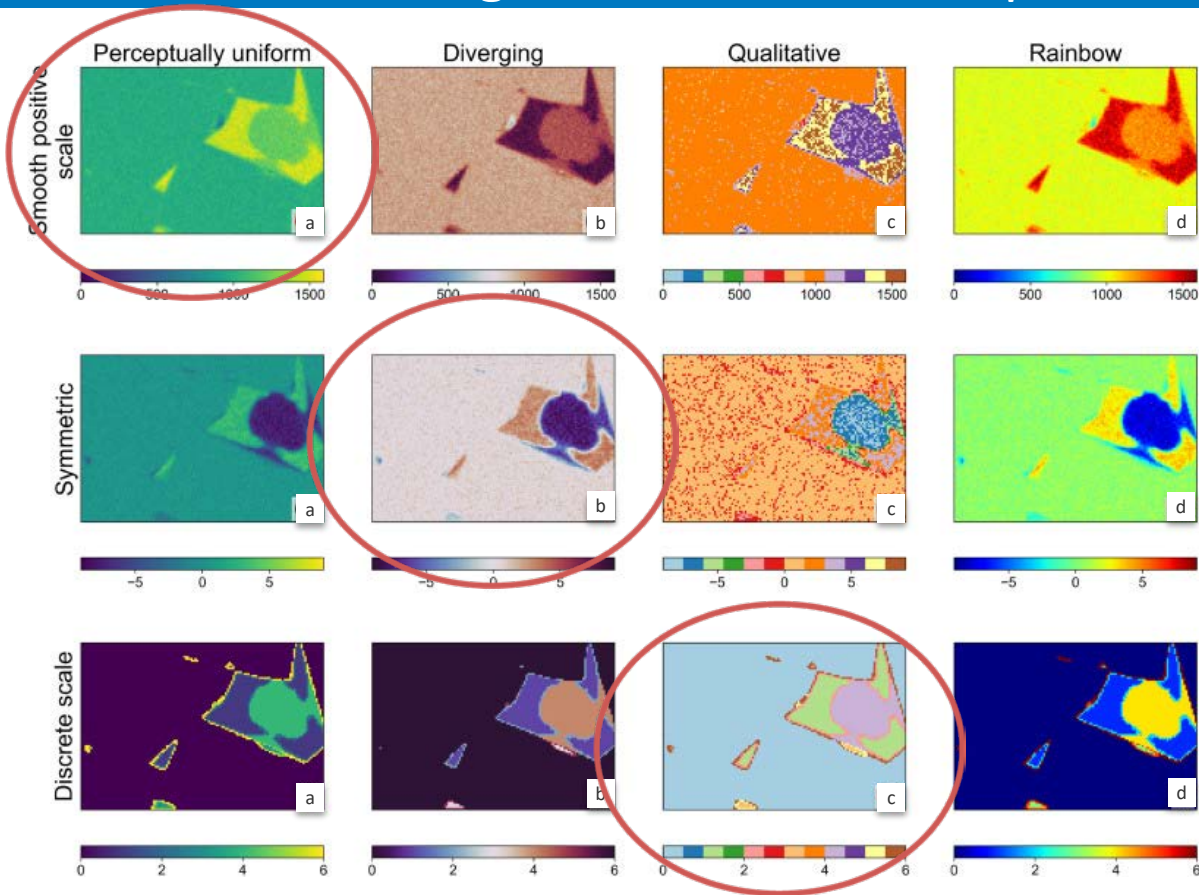
# Human Factor: Colorblindness

- Colorblindness is fairly common
  - 8% in people with one X chromosome
  - 0.5% in people with two X chromosomes

- Many colormaps don't consider this!
  - Choose one that was designed with colorblind people in mind
  - Test your visualization:

    https://www.color-blindness.com/coblis-color-blindness-simulator/

# Meaningful Colors



Legend:
- Sky
- Sunny side of pyramid
- Shady side of pyramid

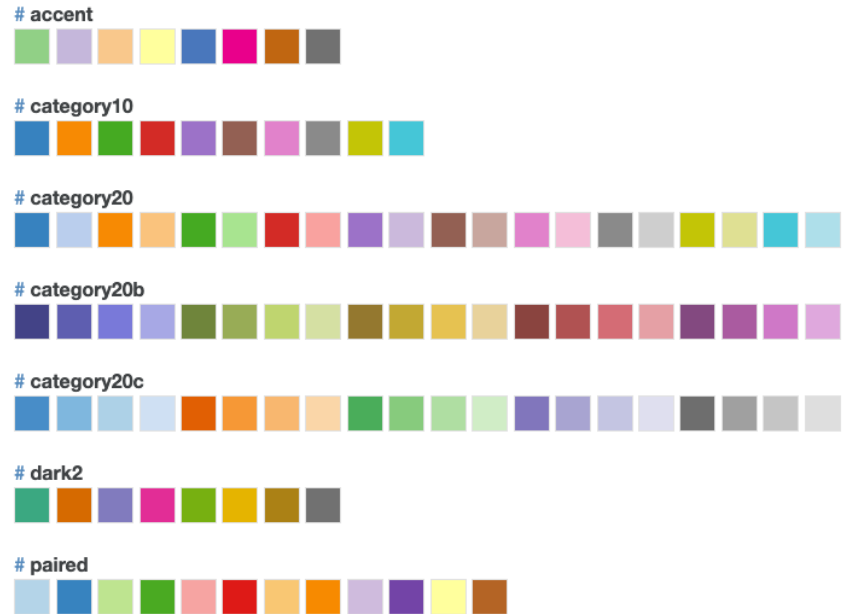# Meaningful Colors: Colormaps That Make Sense



What's the best colormap?

Parish, Chad M., and Philip D. Edmondson. "Data visualization heuristics for the physical sciences." *Materials & Design* 179 (2019): 107868.

# Meaningful Colors: Colormap Choice

***Categorical* or *discrete* colormap**

- Your data consists of distinct categories

- If you have groups of similar meaning, there are colormaps for that, too (category20b, category20c)



# accent

# category10

# category20

# category20b

# category20c

# dark2

# paired

https://vega.github.io/vega/docs/schemes/

# Meaningful Colors: Colormap Choice

***Sequential*** or ***linear*** **colormap**

- Your data has numerical values along one scale

- Your data doesn't have a meaningful midpoint

- Single-hue/multi-hue

https://vega.github.io/vega/docs/schemes/

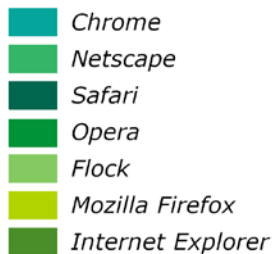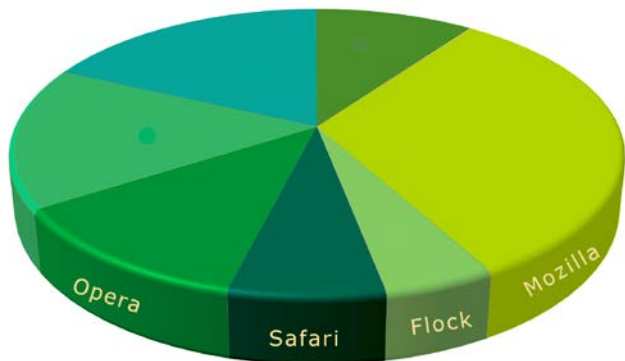# Meaningful Colors: Colormap Choice

**_Divergent_ colormap**

- Quantitative data
- Meaningful mid-point
  - Zero
  - Average value
- *Usually* refers to white in the middle
- *Convergent* colormap *usually* refers to black in the middle
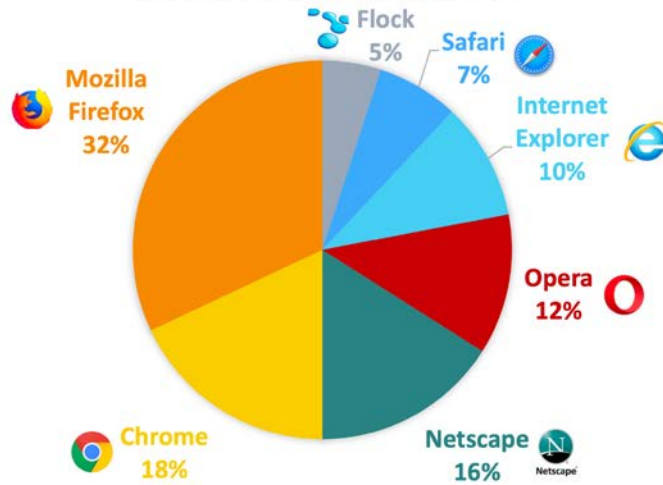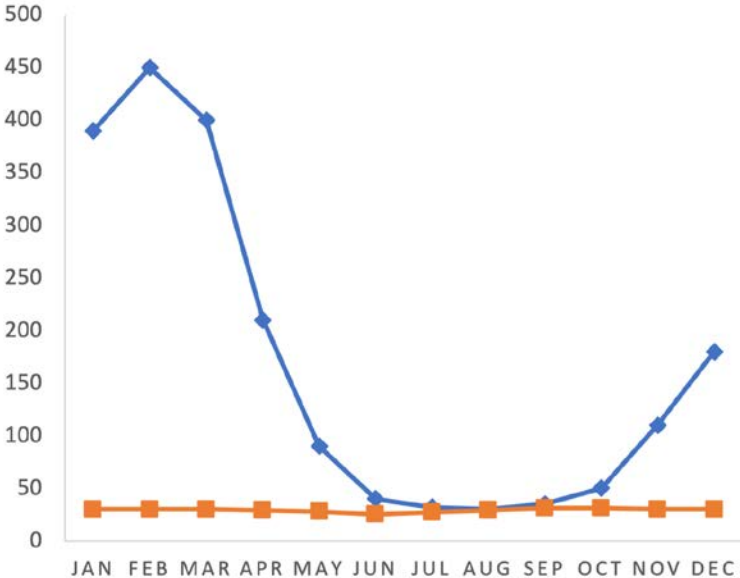- Sometimes you'll see other colors in the middle



https://vega.github.io/vega/docs/schemes/

# Demo: How to Fix a Bad Chart



Legend:
- Chrome
- Netscape
- Safari
- Opera
- Flock
- Mozilla Firefox
- Internet Explorer

**STEP 4: MEANINGFUL COLORS**

- Flock 5%
- Safari 7%
- Internet Explorer 10%
- Mozilla Firefox 32%
- Opera 12%
- Chrome 18%
- Netscape 16%

# Key to Better Visualizations

- Faithful representation of numbers

- Appropriate colormap
  - Correct type for task
  - Colorblind-friendly

- Be consistent

- Focus on data

- Provide context

- Be consistent

- Be minimal

# Sometimes It's Okay to Break the Rules

- Axis scaling – make it abundantly clear and be consistent

- Domain-specific conventions
  - Habits are hard to break (people love their colorful visualizations)
  - Different colors/representations can be confusing, e.g., traffic lights
  - Incremental improvement is key!

- If your data is 3D, it's ok to make a 3D visualization!

- Manual labor – if it's a one-off or illustrative, doing it manually is fine

# Hand-on practice

https://tinyurl.com/TapiaVis2023

# Resources

# Reading

- Parish, Chad M., and Philip D. Edmondson. "Data visualization heuristics for the physical sciences." *Materials & Design* 179 (2019): 107868.

- Shneiderman, Ben, Catherine Plaisant, Maxine Cohen, Steven Jacobs, Niklas Elmqvist, and Nicholas Diakopoulos. *Designing the user interface: strategies for effective human-computer interaction*. Pearson, 2016.
  - Search for "Ben Shneiderman 8 golden principles"

- Tufte, Edward R. *The visual display of quantitative information*. Vol. 2. Cheshire, CT: Graphics press, 2001.
  - Search for "Edward Tufte visualization principles"

- Moreland, Kenneth. "Why we use bad color maps and what you can do about it." *Electronic Imaging* 2016, no. 16 (2016): 1-6.

# Colormaps

- Free E-Book on Color Blindness Essentials:
  http://www.color-blindness.com/2010/02/23/color-blind-essentials/

- Finding appropriate colormaps:
  https://colorbrewer2.org/#type=sequential&scheme=BuGn&n=3

- Many other color and visualization tools: https://sciviscolor.org/tools/

# Data Visualization Tools: Python

- Matplotlib
  - Essentially trying to recreate Matlab's plotting
  - Lots of options but steep learning curve

- ggplot
  - Essentially trying to recreate R's plotting

- Bokeh
  - Interactive plots but steep learning curve

- Seaborn
  - Built on top of Matplotlib, uses pandas input
  - Simple to use
  - Caution: changed dramatically between versions

- PySAL
  - Geospatial data, great for maps

- PyLeaflet
  - Python library to include LeafletJS
  - Geospatial data

...and many more!

# Data Visualization Tools: JavaScript

- LeafletJS
  - Maps
  - Easy to use

- OpenLayers
  - Maps
  - Medium difficulty

- WebWorldWind
  - 3D globe
  - Medium difficulty
  - Comes with layer switching

- Cesium
  - Fancy-looking 3D visualizations
  - Advanced

- WebGL
  - Fast compute in browser

- D3.js
  - Charts

- Highcharts
  - Charts

*...and many more!*

# Thank you

**www.nrel.gov**  Andy.Berres@nrel.gov

NREL/PR-5700-87388

Photo from iStock-627281636

NREL
*Transforming* ENERGY