

Data visualization

Volodymyr Miz

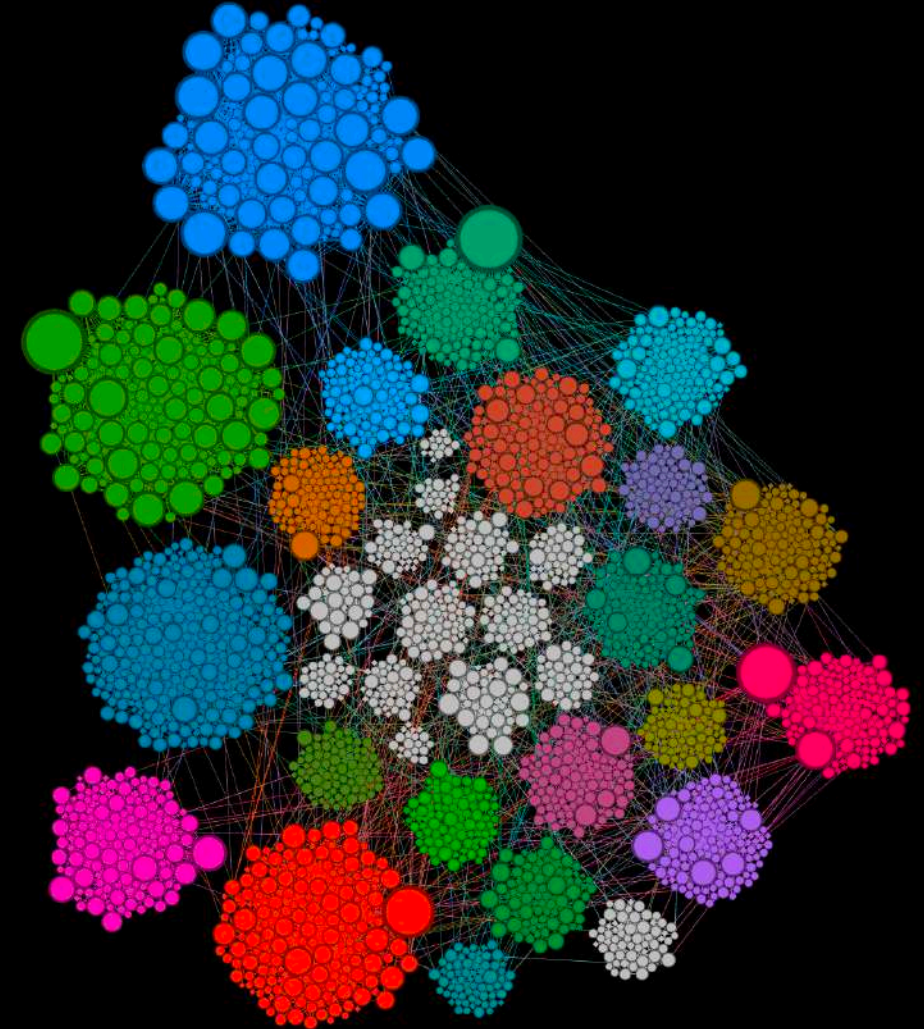
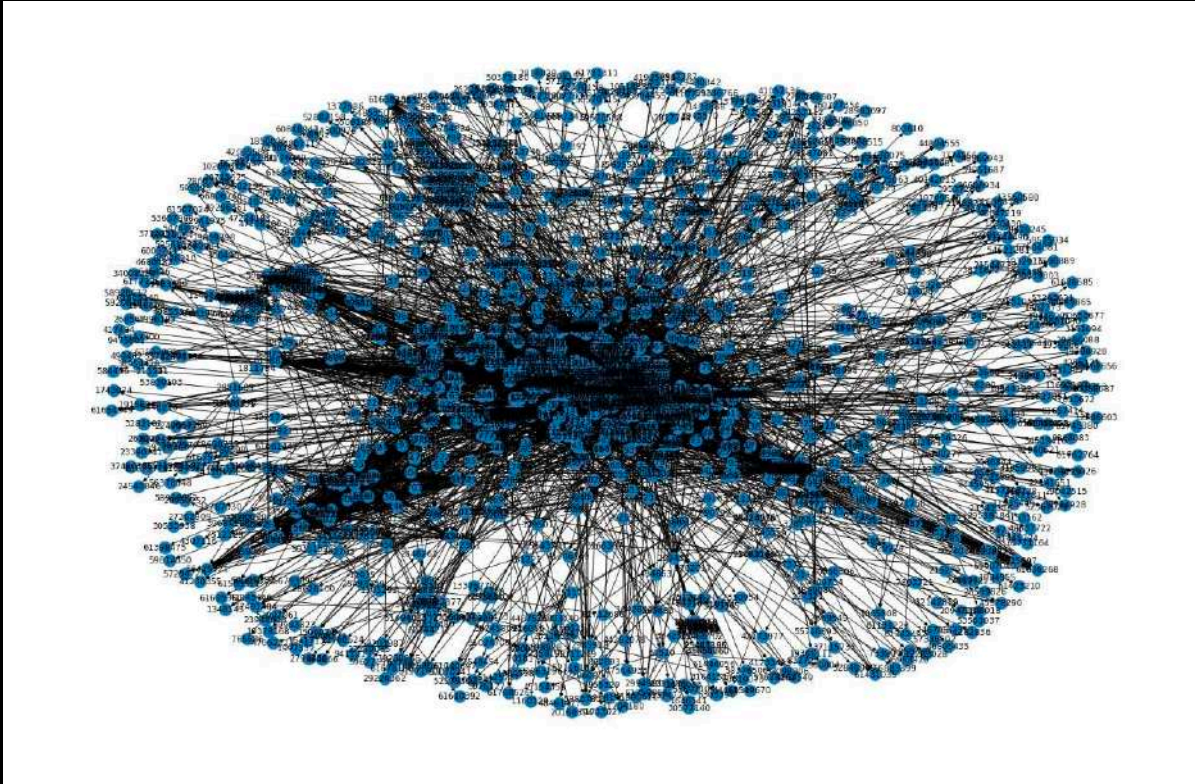
for A Network Tour of Data Science (Fall 2019)
by Prof. Pierre Vandergheynst and Prof. Pascal Frossard

Some examples are taken from Dr. Kirell Benzi's class Data Visualization (2018)

Outline

- Visualization and storytelling
 - Designing your visualization
- DOs and DON'Ts
 - Best practices
 - Misinformation
- Tools
- **Graph visualization**
- **Tutorials**
 - Graph layouts
 - Interactive web-based visualization

Motivation



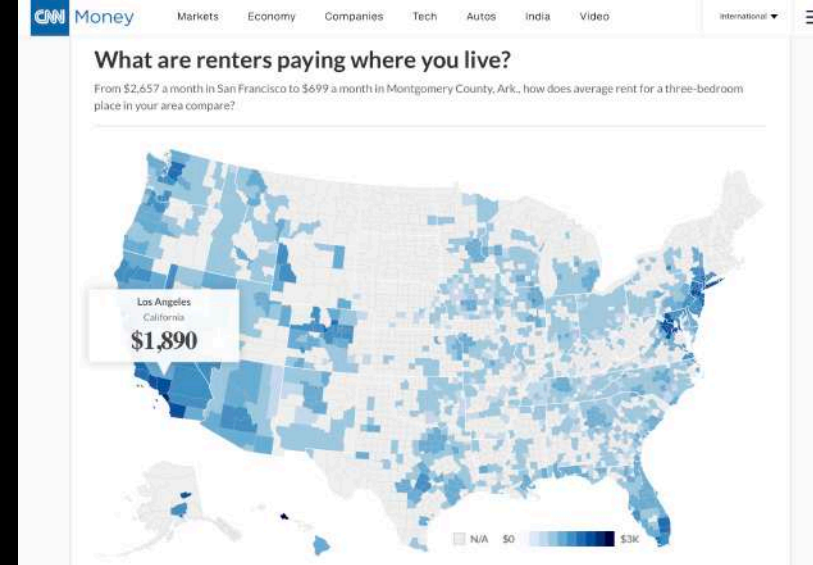
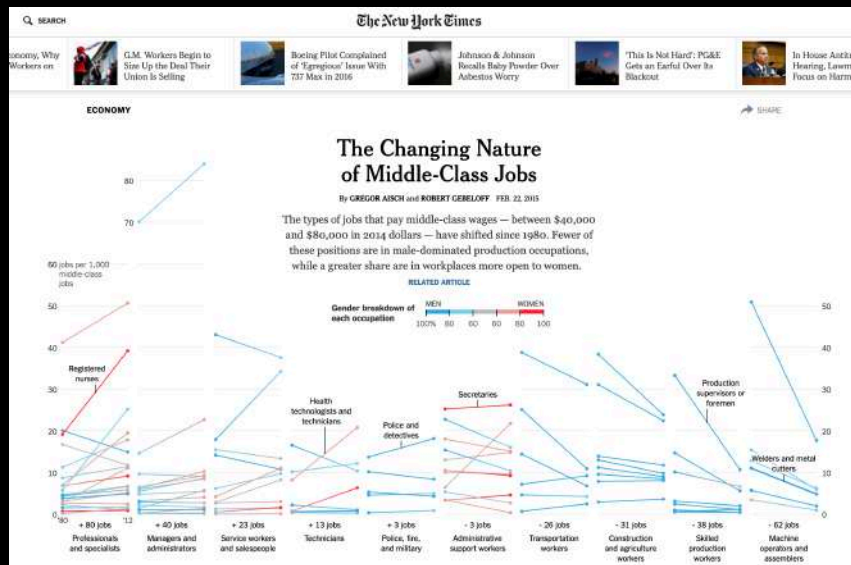
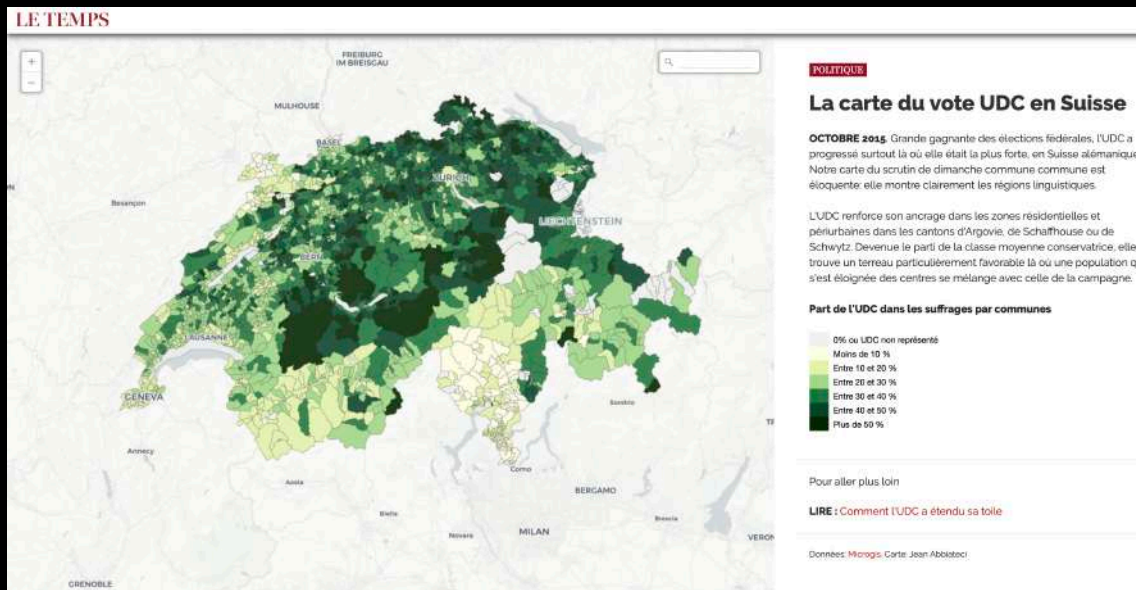
The background of the slide is a dense, intricate network graph. It consists of numerous small, bright green circular nodes connected by thin, dark green lines representing edges. The connections are complex and non-uniform, creating a web-like structure that fills the entire frame. The overall effect is one of a highly interconnected system, possibly representing a social network, a biological system, or a data network.

1. Visualization and Storytelling

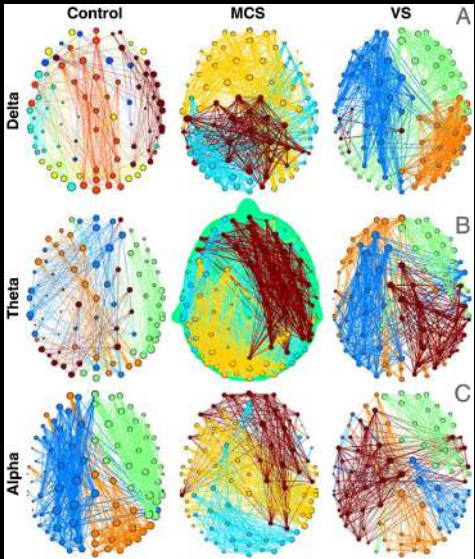
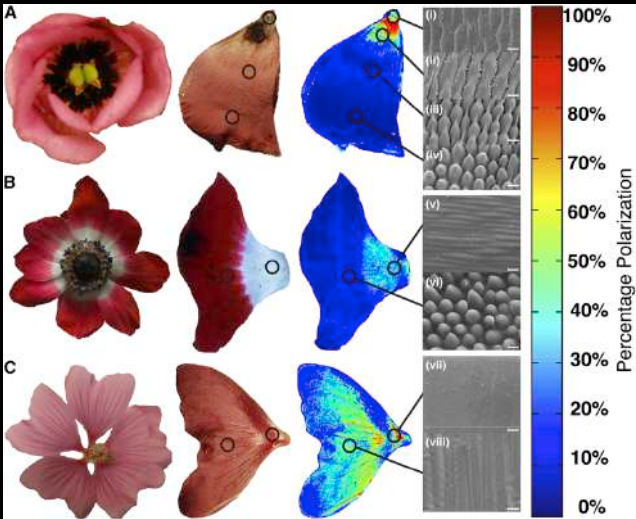
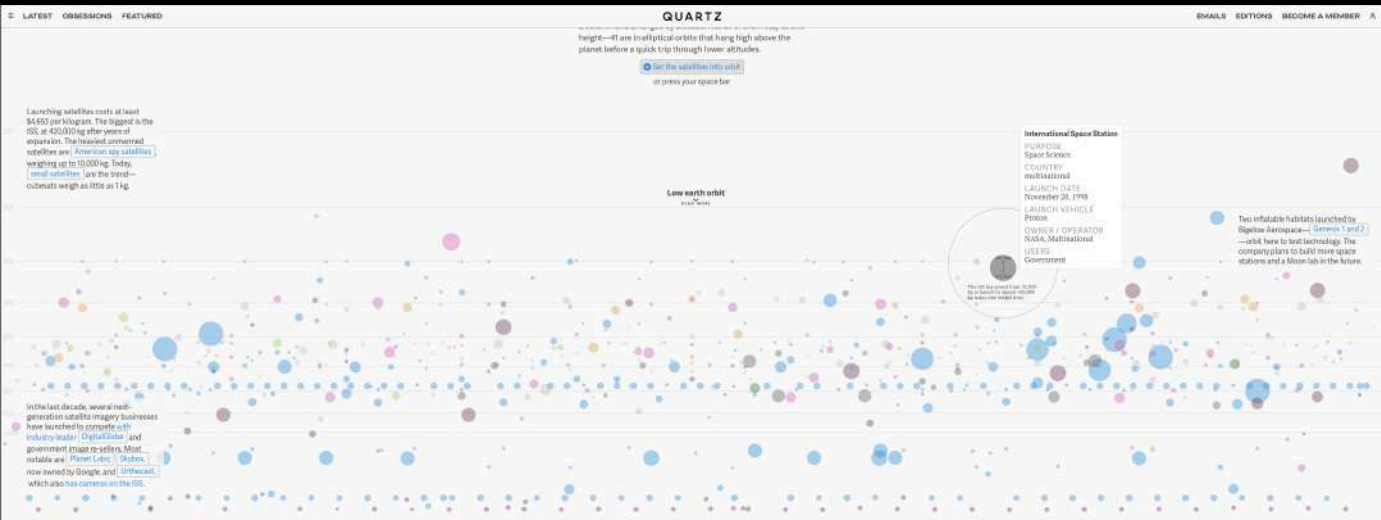
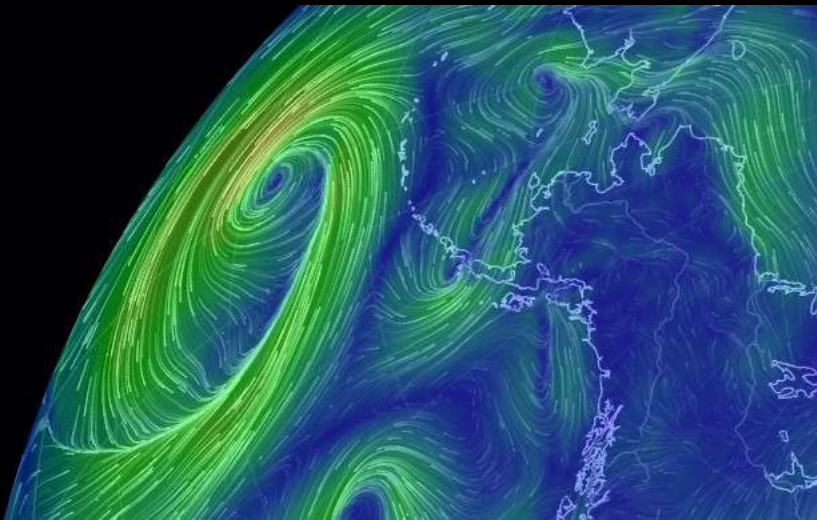
Visualization and Storytelling

- Data journalism
- Designing your visualizations
 - What?
 - Why?
 - How?

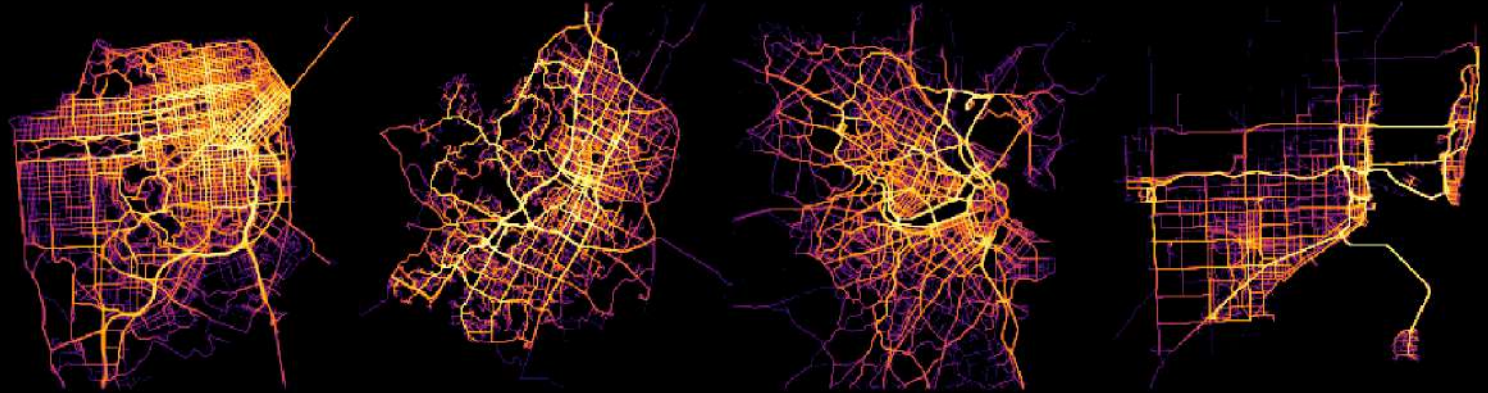
Data journalism. News



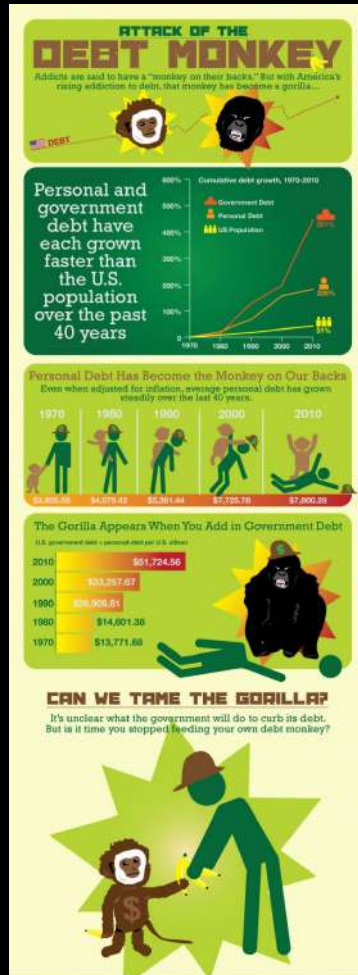
Data journalism. Science Communication



Data journalism. Tech and Marketing

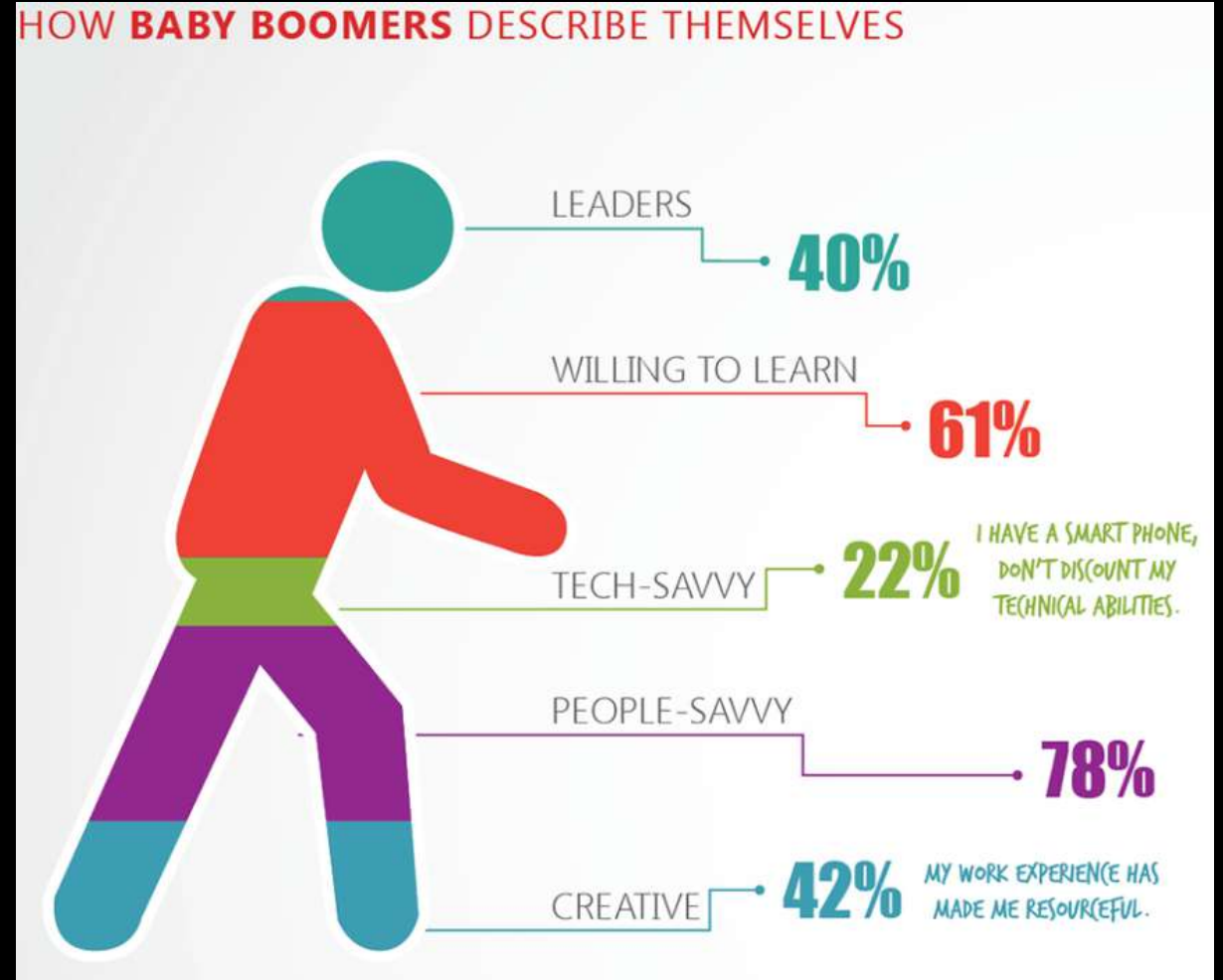


Bad Visualization Kills Good Content

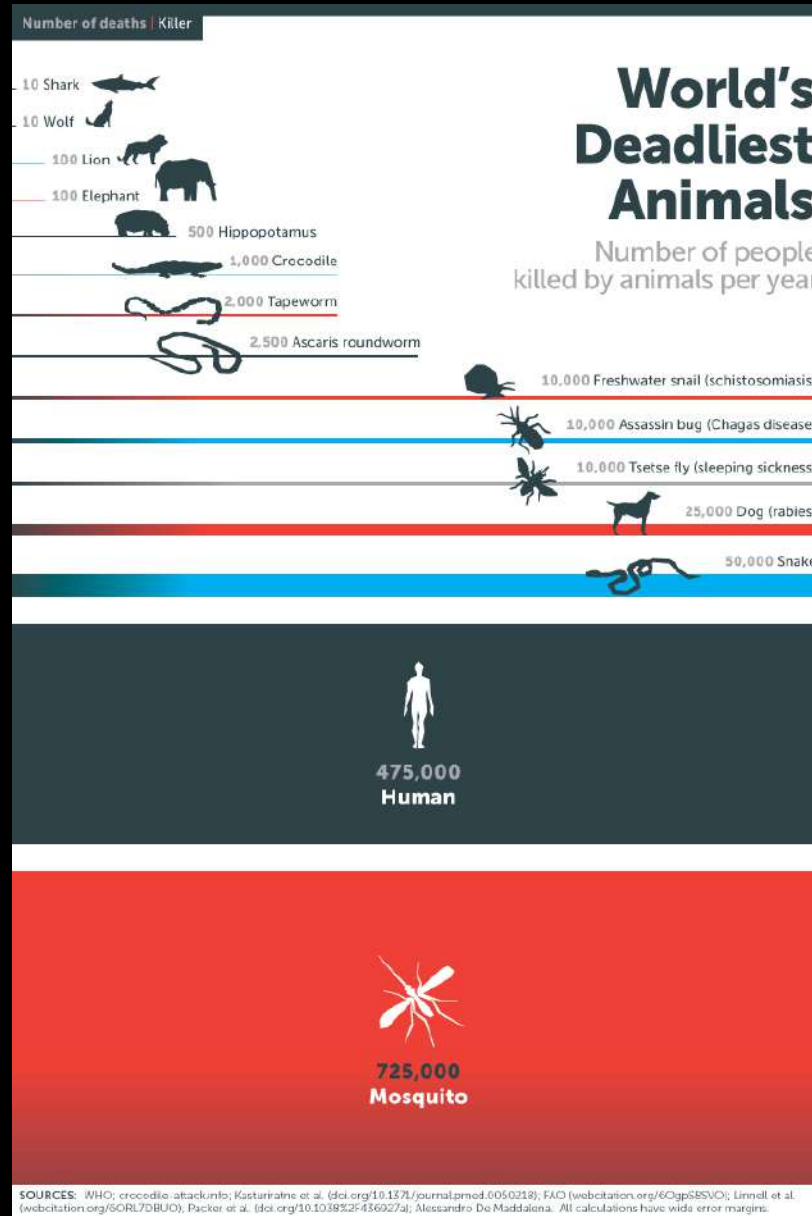


Sources
U.S. Federal Reserve (Consumer Debt Figures, May 7, 2012)
http://www.federalreserve.gov/releases/m3/m3d.htm
U.S. Census Bureau (U.S. population figures, May 2010)
http://www.census.gov/prod/2010/states/totals/2010-02.pdf
U.S. Treasury Department (Government Debt Figures, October 1, 2010)
http://www.treasury.gov/press/releases/2010/02/2010-02-01.pdf
All figures adjusted for inflation into 2010 dollars.

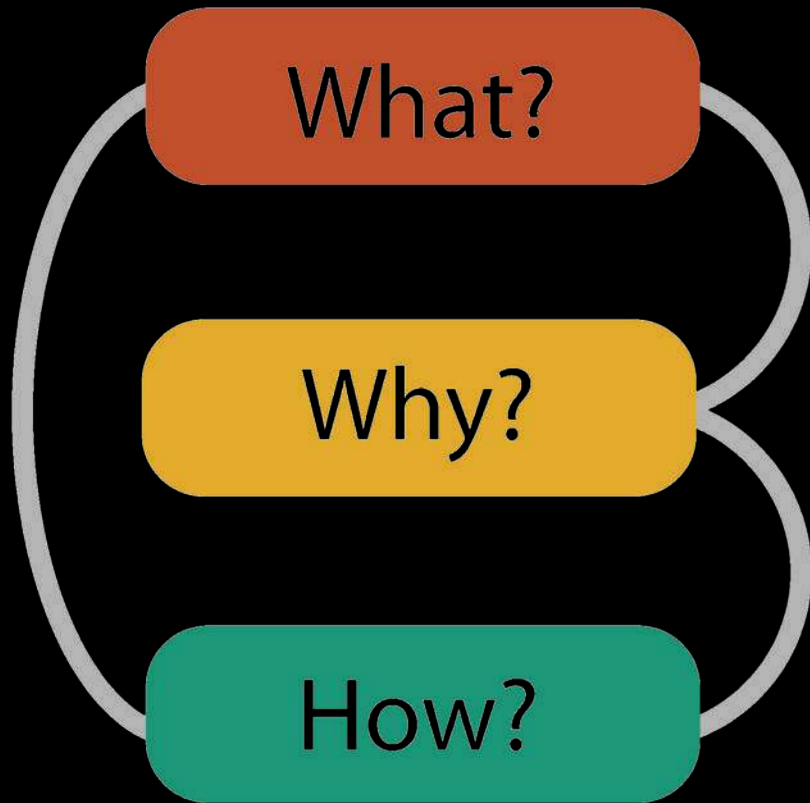
Created for
MoneyRates
Make the most of your money
© SunBurst, Inc. 2012



Good Visualization Breaks Stereotypes



What? Why? How?

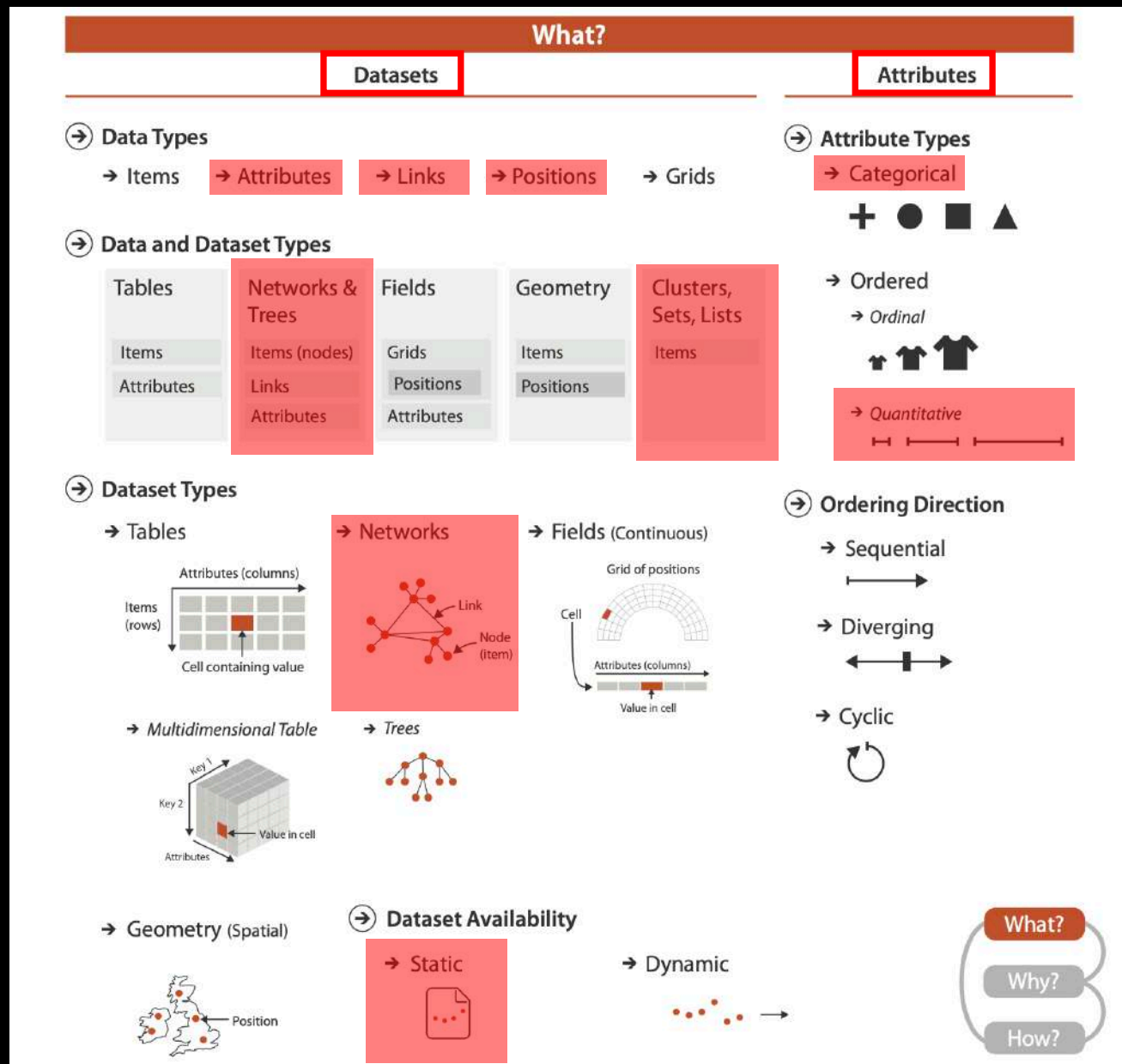


What data are you going to show?

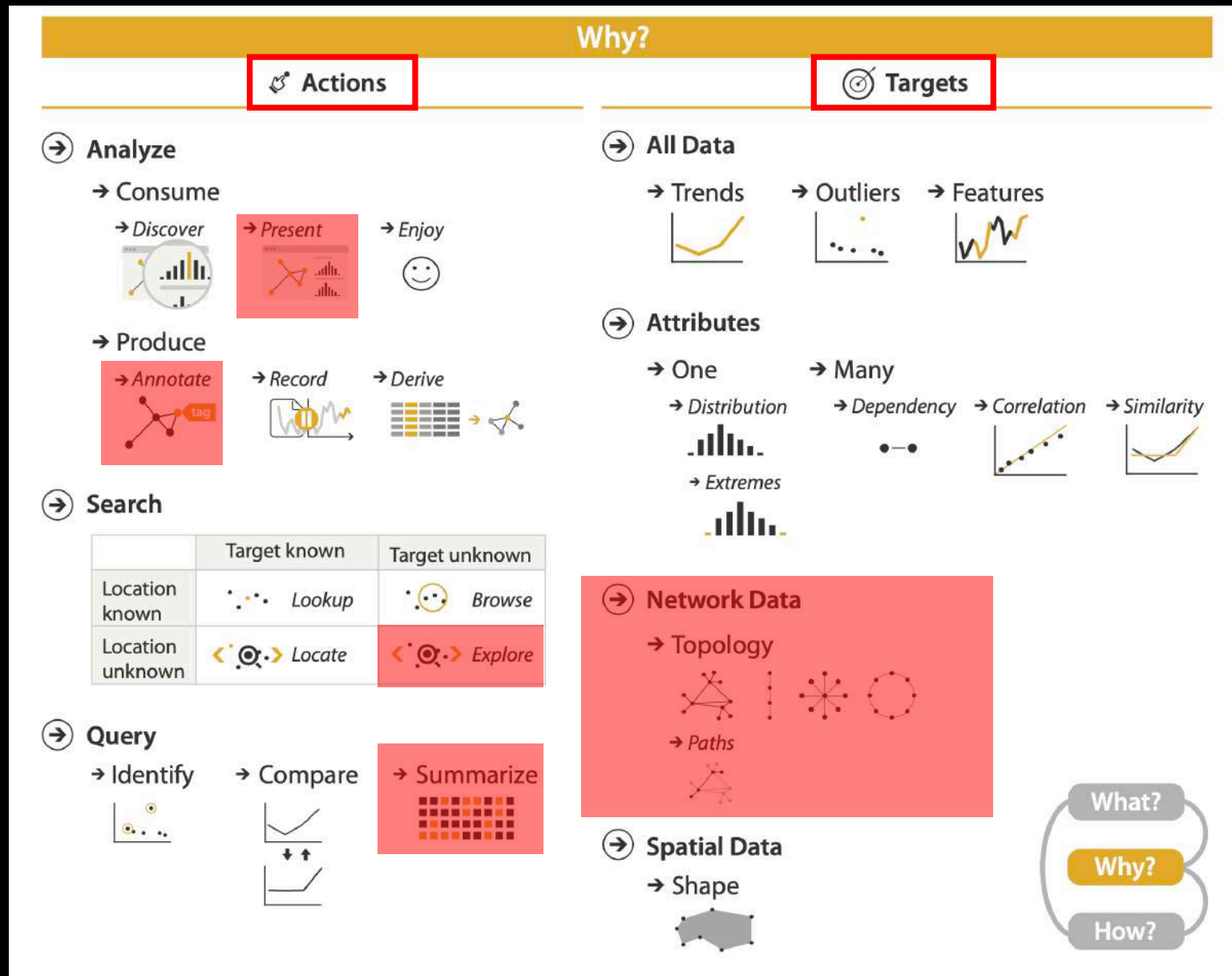
Why are you showing visualization?

How are you going to construct visualization?

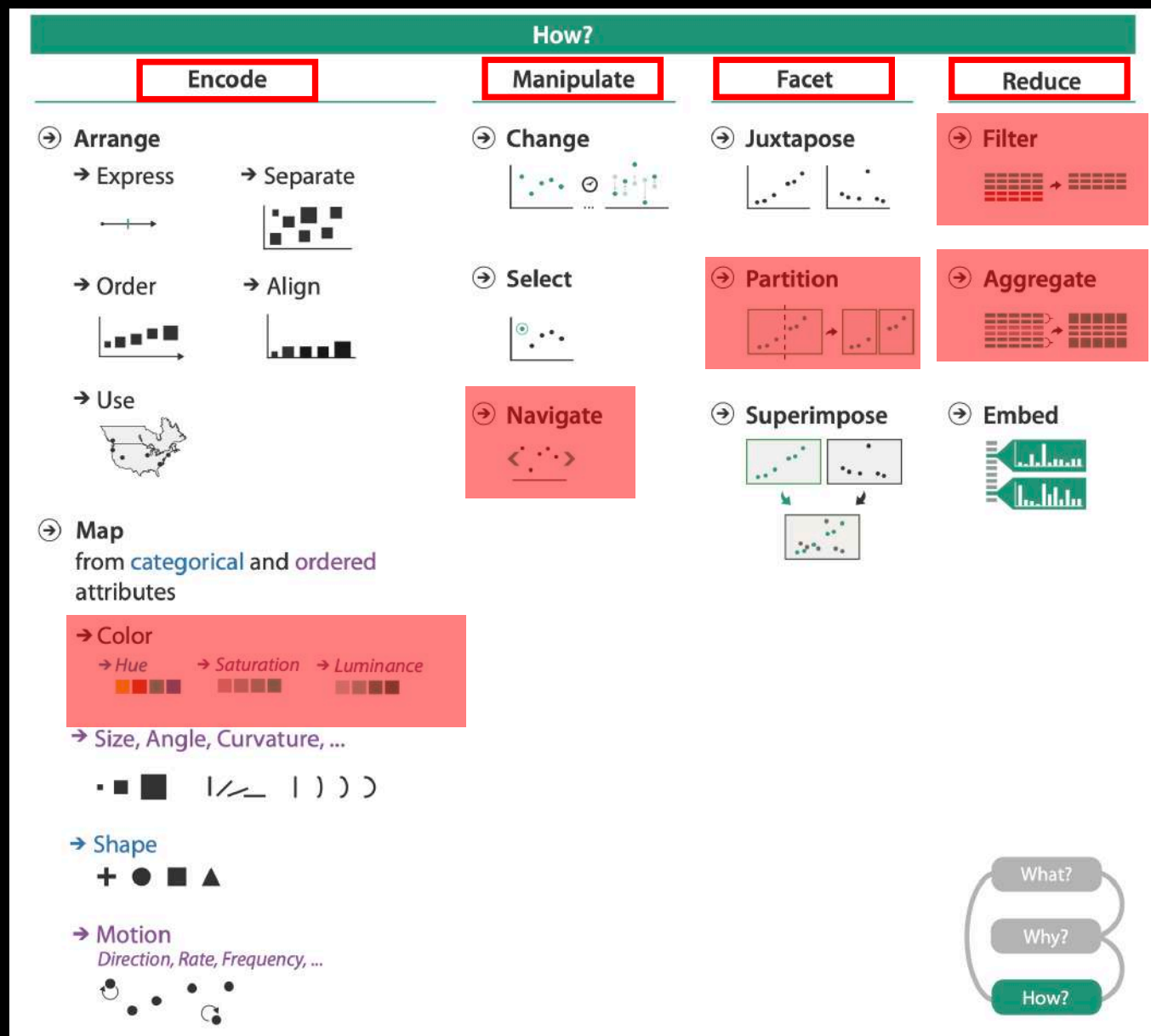
What can we visualize?



Why do people use our visualization?



How to design our visualization?





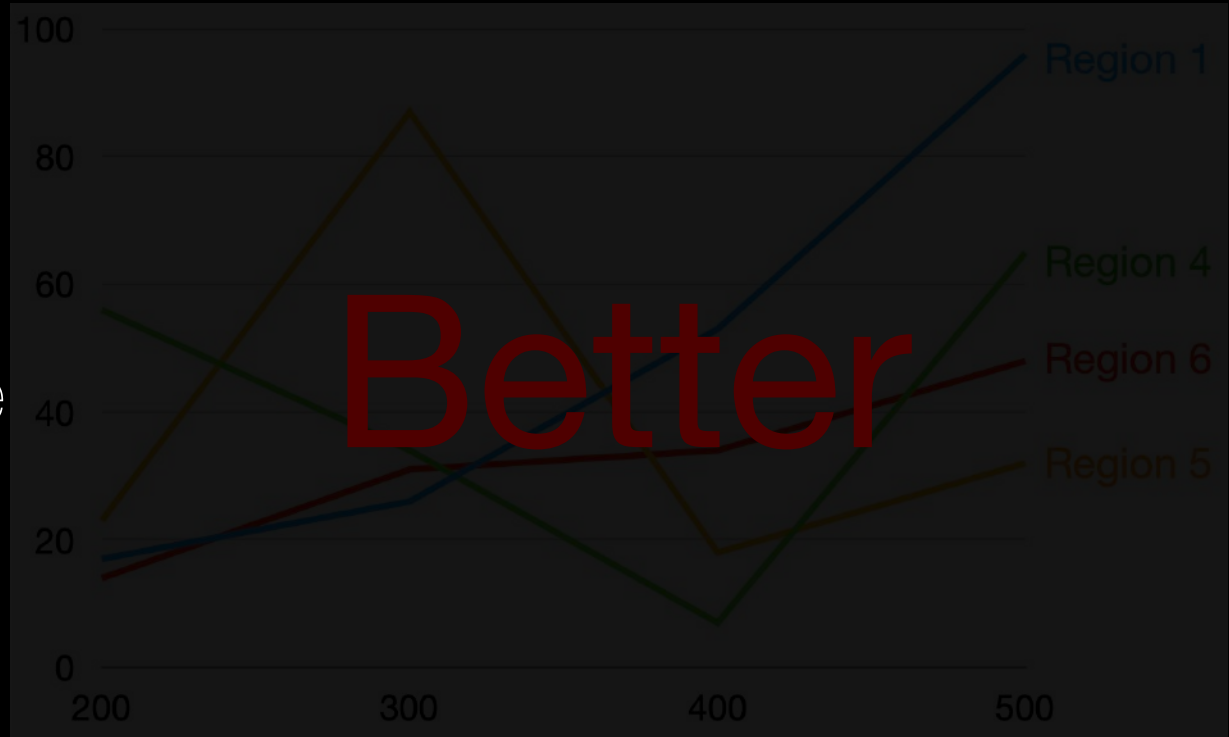
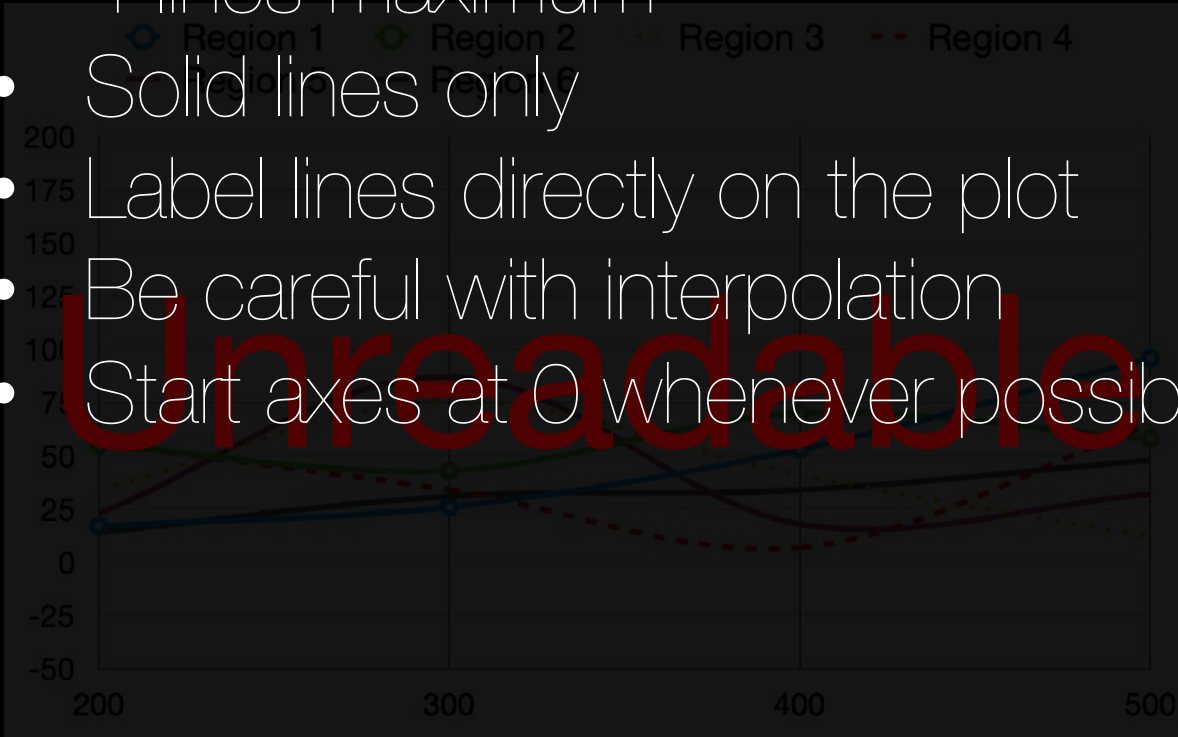
2. DOs and DON'Ts

DOs and DON'Ts

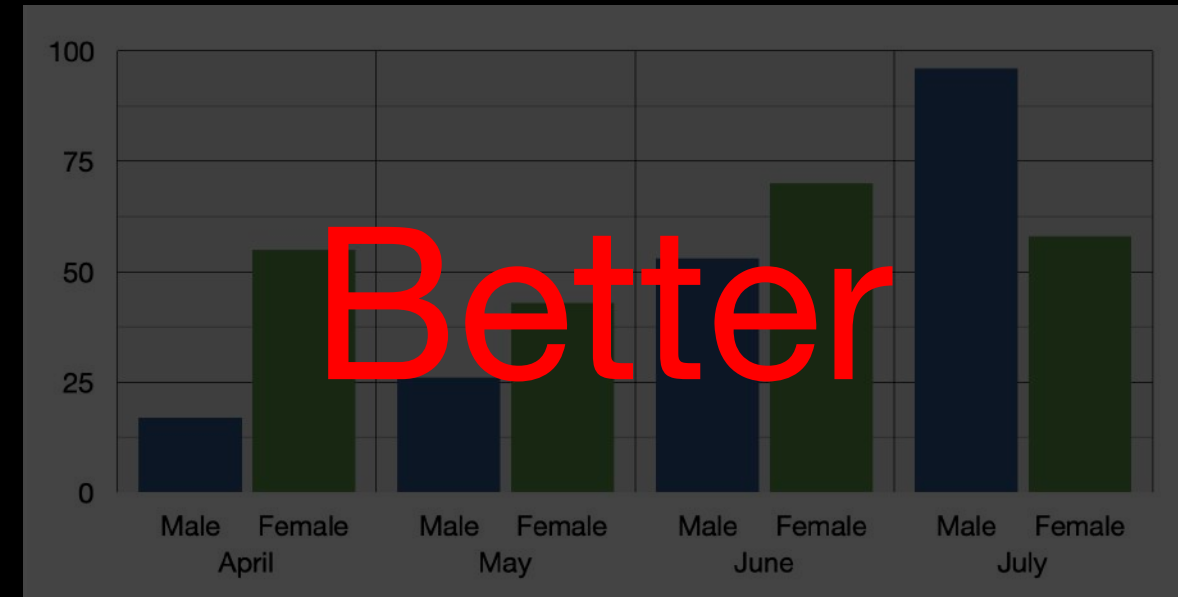
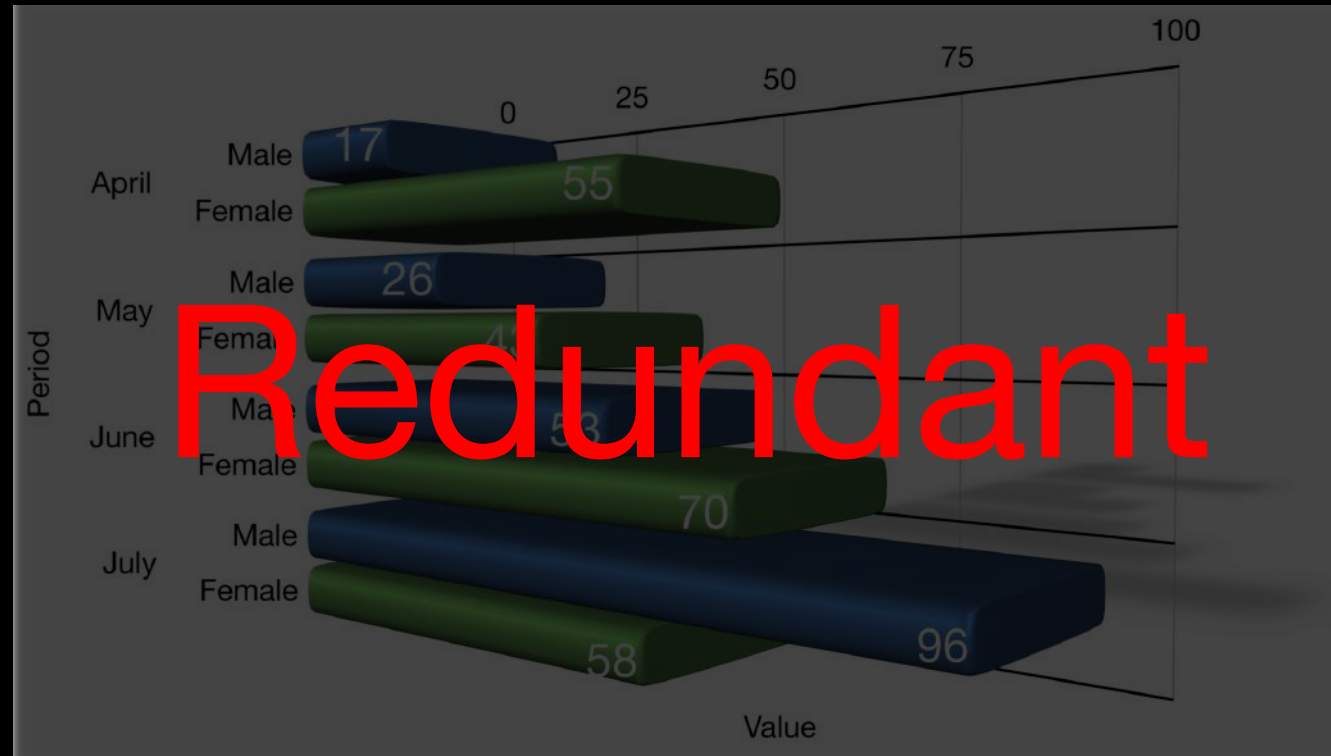
- Best practices, Tips and Tricks
- Misinformation
 - Manipulation
 - Misleading visualization
 - False insights

Line Charts

- 4 lines maximum
- Solid lines only
- Label lines directly on the plot
- Be careful with interpolation
- Start axes at 0 whenever possible

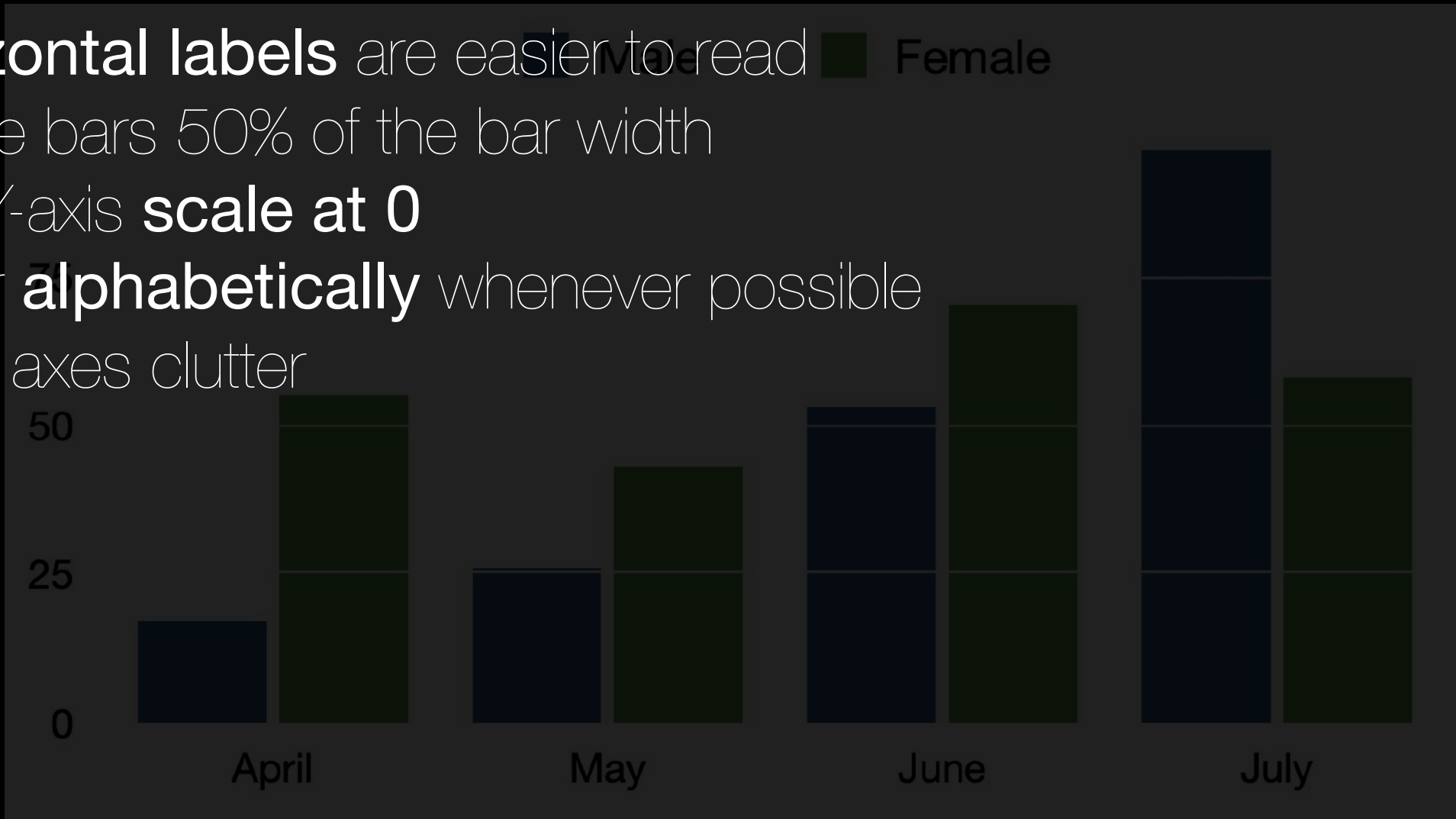


3D or not 3D. Bar charts



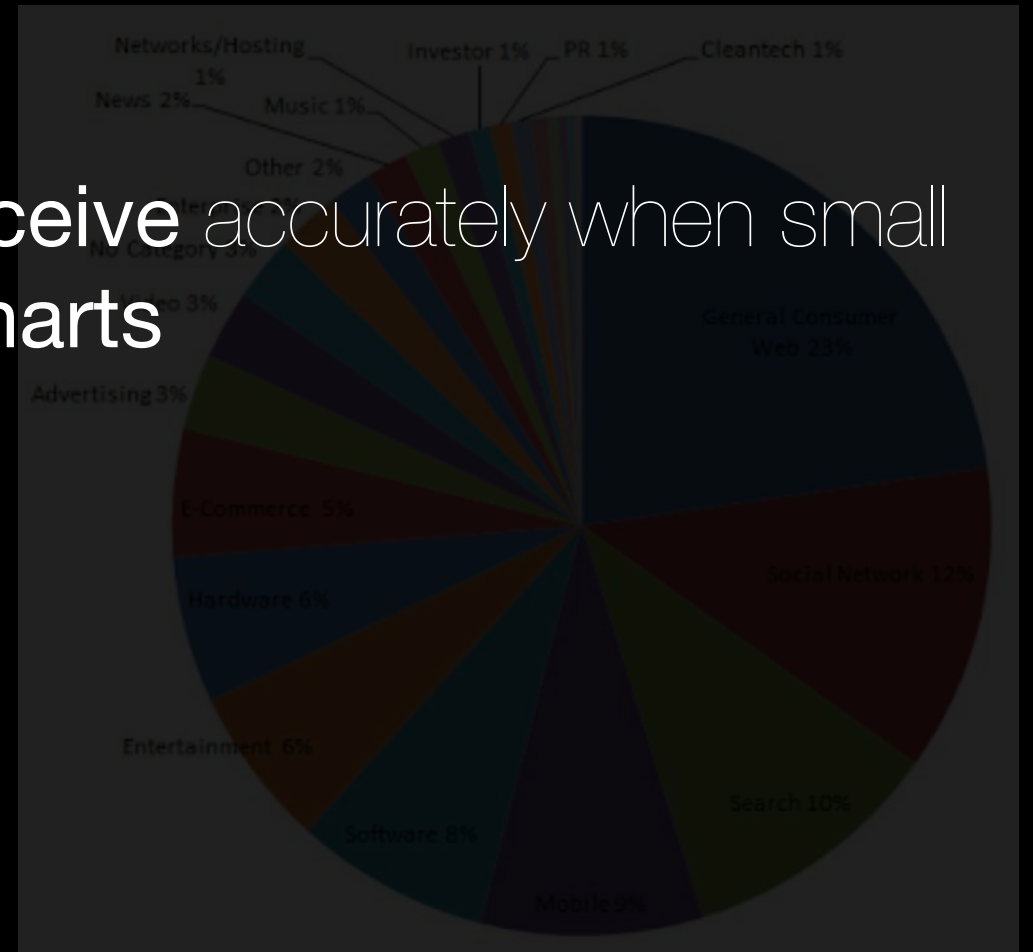
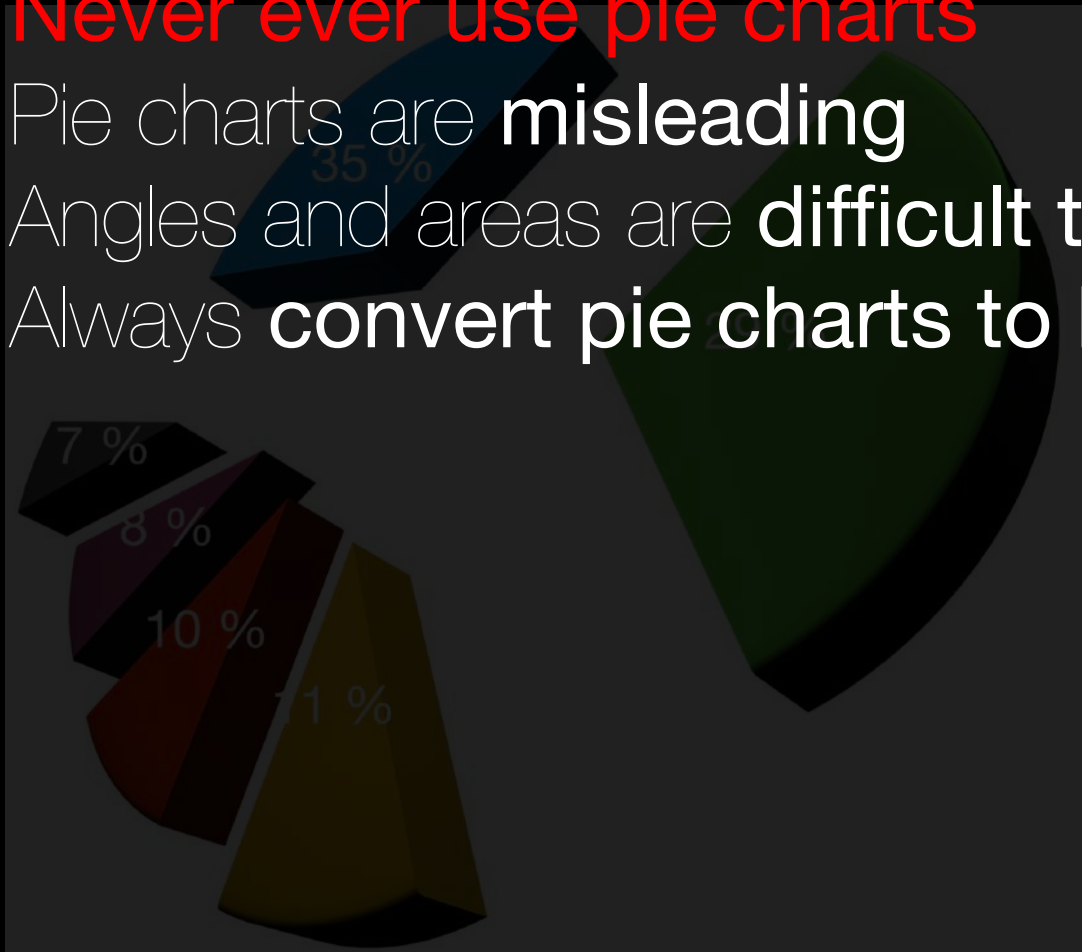
Bar Charts

- **Horizontal labels** are easier to read
- Space bars 50% of the bar width
- Start Y-axis **scale at 0**
- Order **alphabetically** whenever possible
- Avoid axes clutter

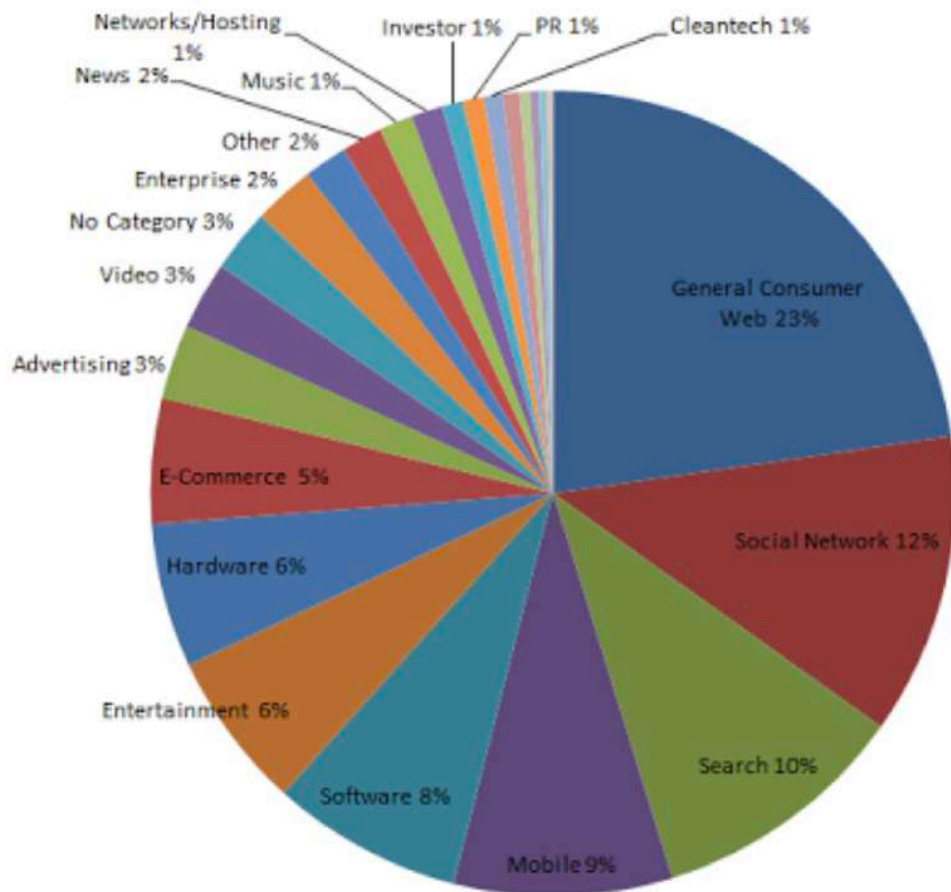


Pie Charts

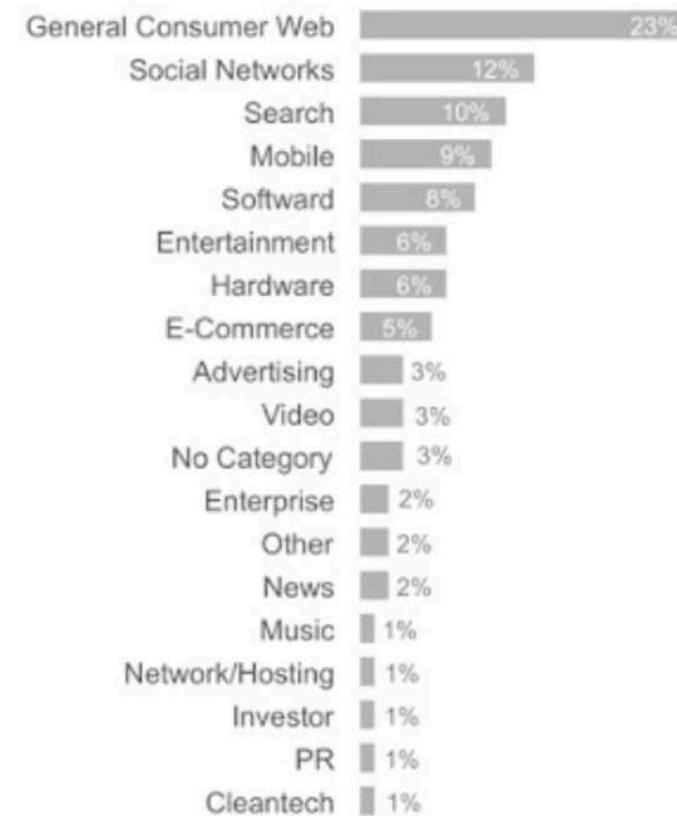
- **Never ever use pie charts**
- Pie charts are **misleading**
- Angles and areas are **difficult to perceive** accurately when small
- Always **convert pie charts to bar charts**



Convert pie to bar charts



TechCrunch Coverage: 2005 - 2011
Bars are best!

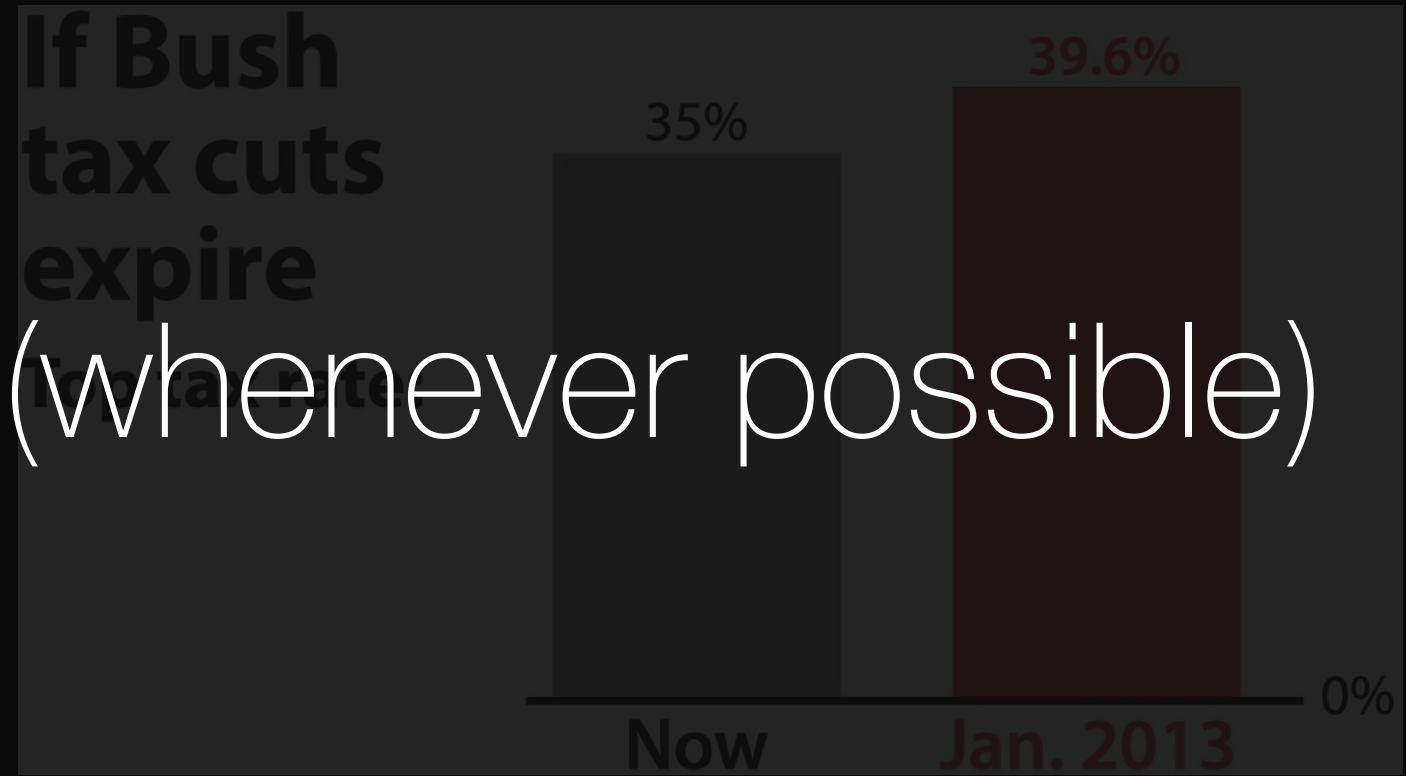
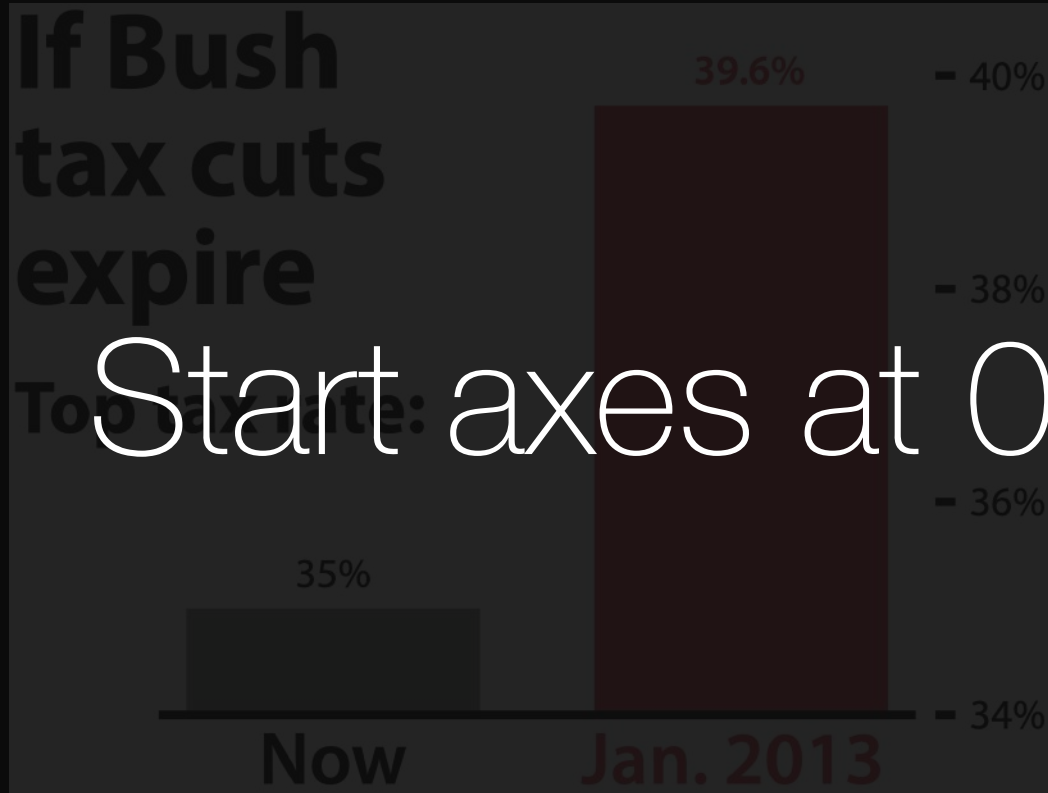


Chartjunk

- Pros
 - Catchy
 - Easier to attract attention
- Cons
 - Hard to interpret
 - Bias
 - Less trustworthy
 - Distracts the viewer from the main content



Misleading scale



Start axes at 0 (whenever possible)

Misleading scale

Image classification

ResNet-50

image classification model (ImageNet-1K)

FULLY-SUPERVISED

76.4%

WEAKLY SUPERVISED

78.2%

SEMI-WEAKLY SUPERVISED

84.4%

75.8%

Top-1 accuracy

WEAKLY SUPERVISED
RESNEXT-101-32X48
TEACHER MODEL

Video classification

R(2+1)D-18

video action classification model(kinetics-400)

FULLY-SUPERVISED

64.8%

WEAKLY SUPERVISED

71.5%

SEMI-WEAKLY SUPERVISED

84.2%

61.8%

Top-1 accuracy

WEAKLY SUPERVISED
R(2+1)D-152
TEACHER MODEL

Start axes at 0 (whenever possible)

Accurate image and video classification is important for a wide range of computer vision applications, from identifying harmful content, to making products more accessible to the visually impaired, to helping people more easily buy and sell things on products like Marketplace. Facebook AI is developing alternative ways to train our AI systems so that we can do more with less labeled training data overall, and also deliver accurate results even when large, high-quality labeled data sets are simply not available. Today, we are sharing details on a versatile new model training technique that delivers state-of-the-art accuracy for image and video classification systems.

This approach, which we call semi-weak supervision, is a new way to combine the merits of two different training methods: [semi-supervised learning](#) and [weakly supervised learning](#). It opens the door to creating more accurate, efficient production classification models by using a teacher-student model training paradigm and billion-scale weakly supervised data sets. If the weakly supervised data sets (such as the hashtags associated with publicly available photos) are not available for the target classification task, our method can also make use of unlabeled data sets to produce highly accurate semi-supervised models.

Related posts

[Creating a data set and a challenge for deepfakes](#)

September 05, 2019

[Mapping roads through deep learning and weakly supervised training](#)

July 23, 2019

[RoBERTa: An optimized method for](#)
[robustness self-supervised BERT pre-training](#)

Misleading attributes

Votes for Donald Trump

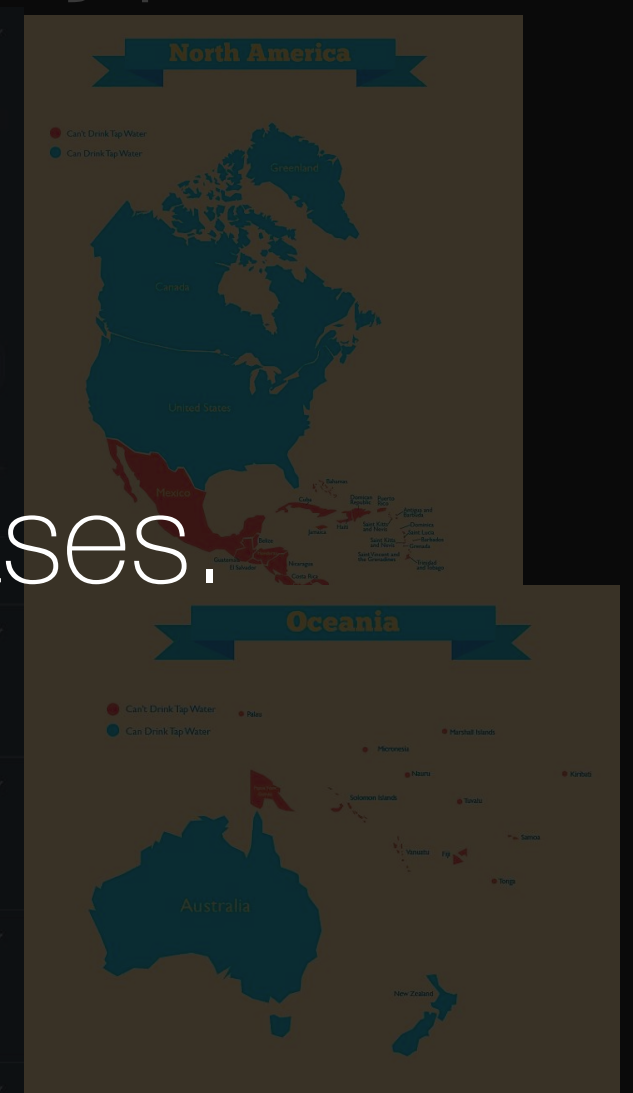
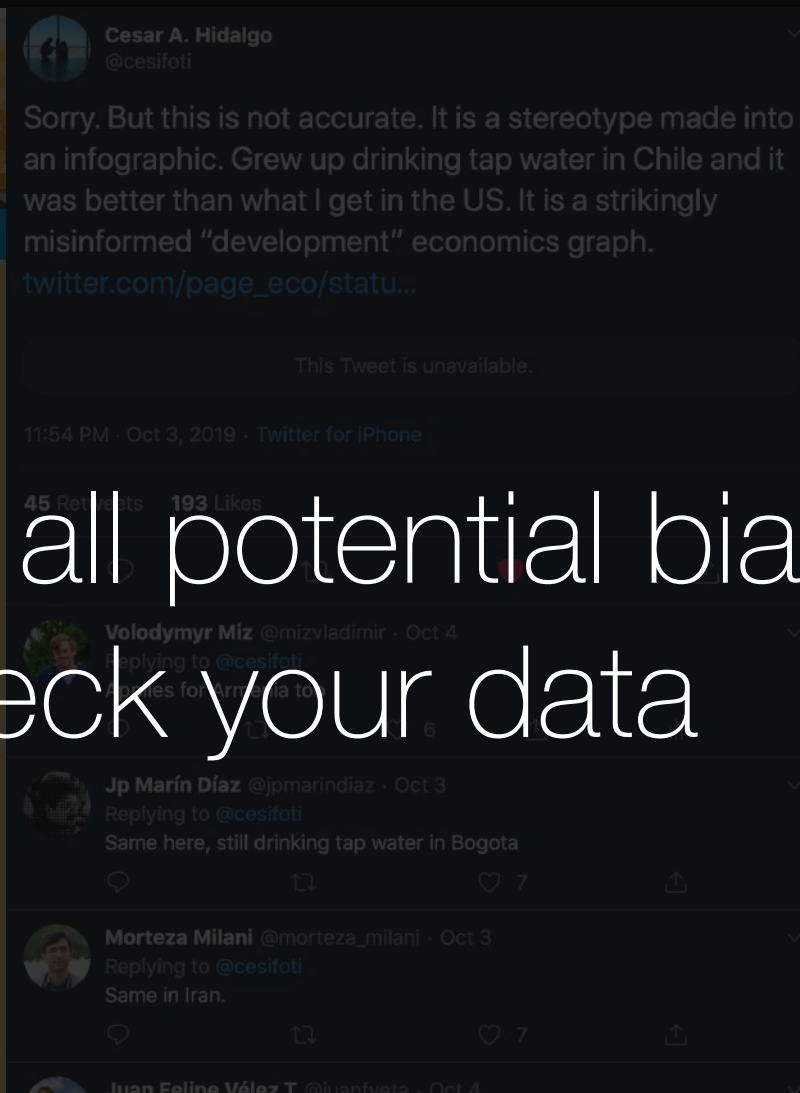
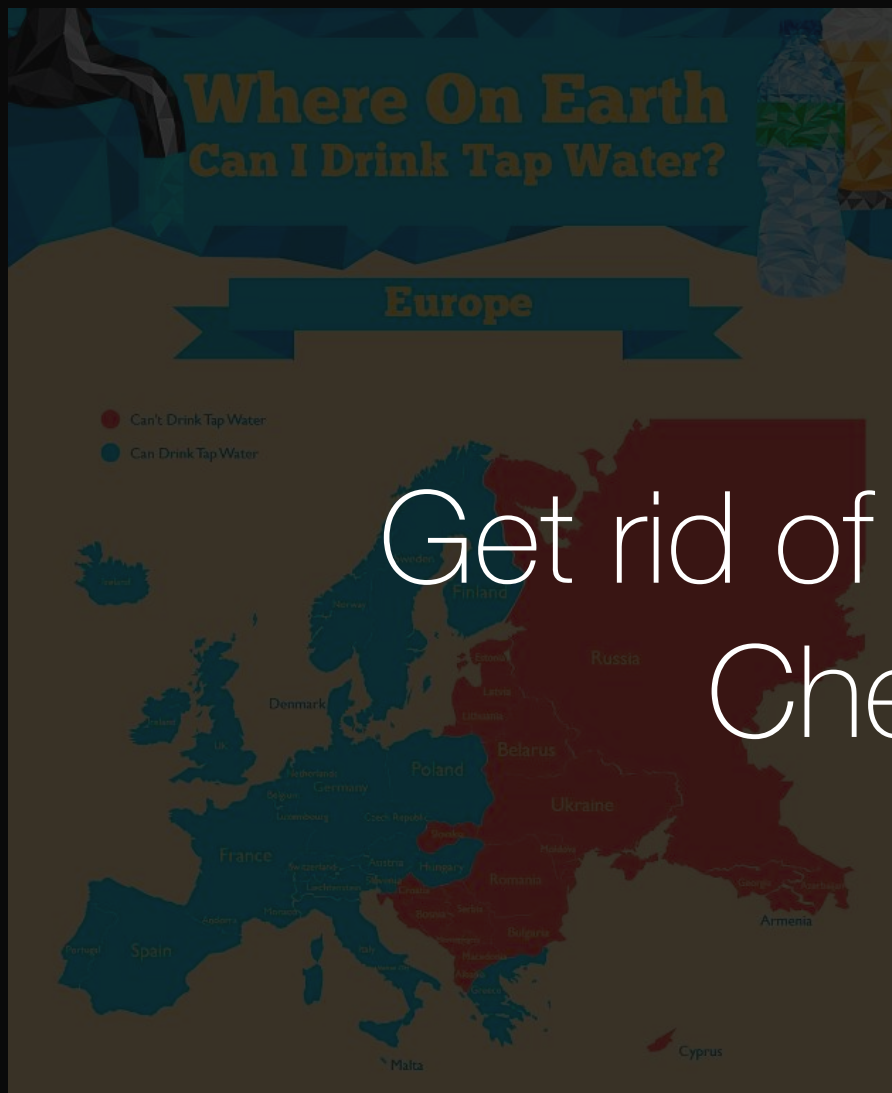
Votes for Hillary Clinton

Attributes can be misleading.

Explore all attributes and know unique aspects of the data

Bubble size is proportional to the number of votes per county

Misinformation. Stereotypes



Get rid of all potential biases.
Check your data



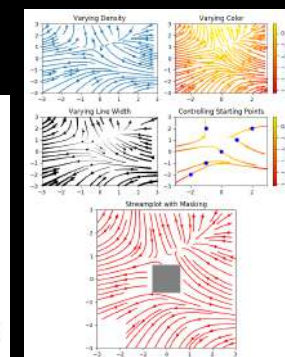
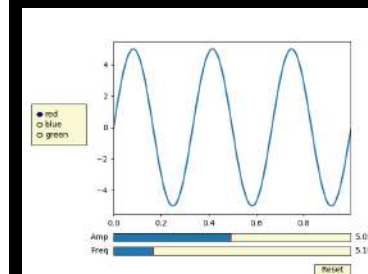
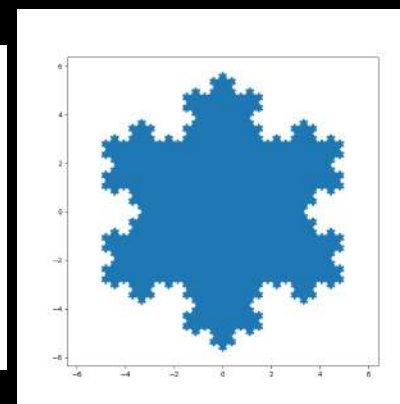
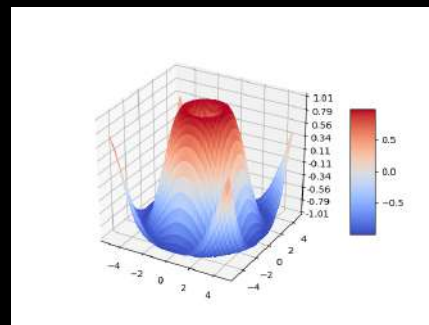
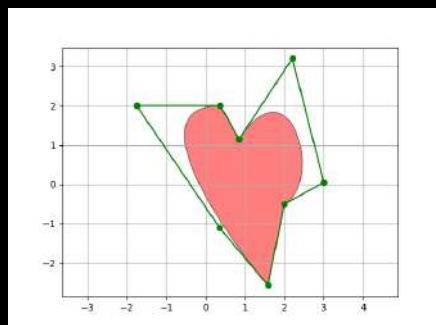
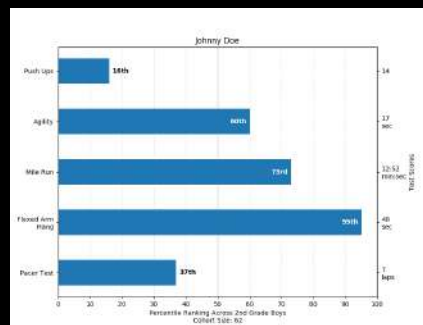
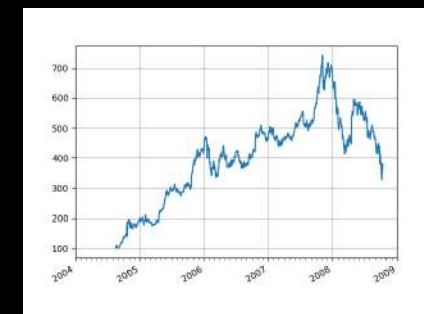
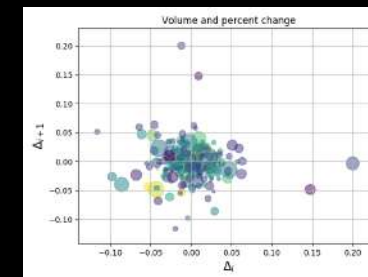
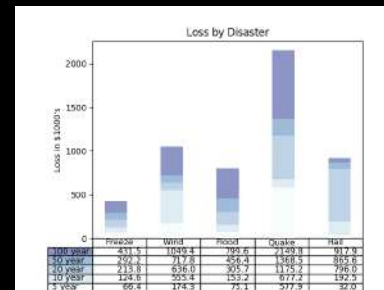
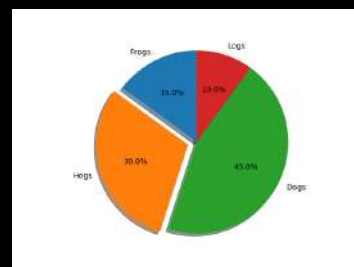
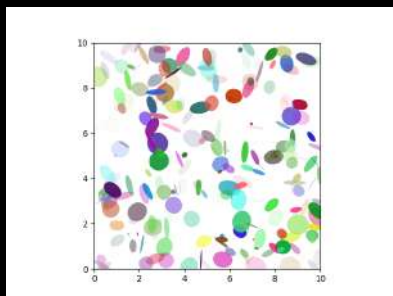
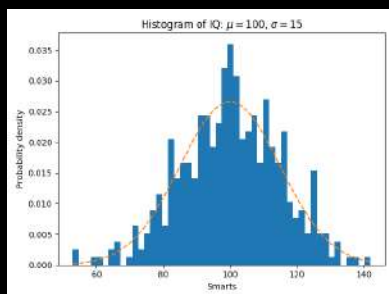
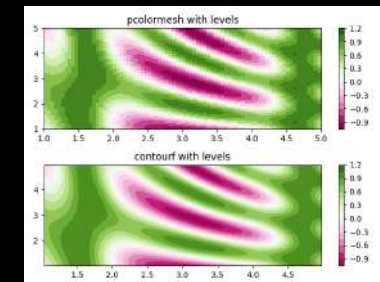
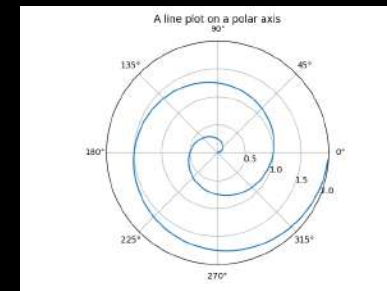
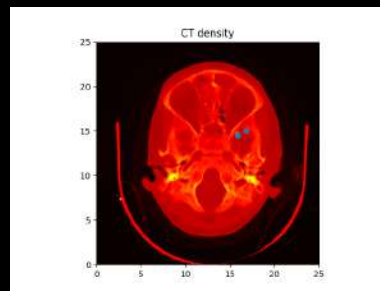
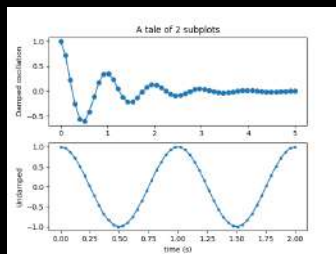
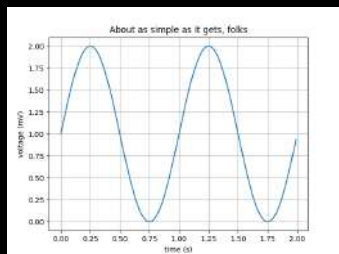
3. Tools

Tools

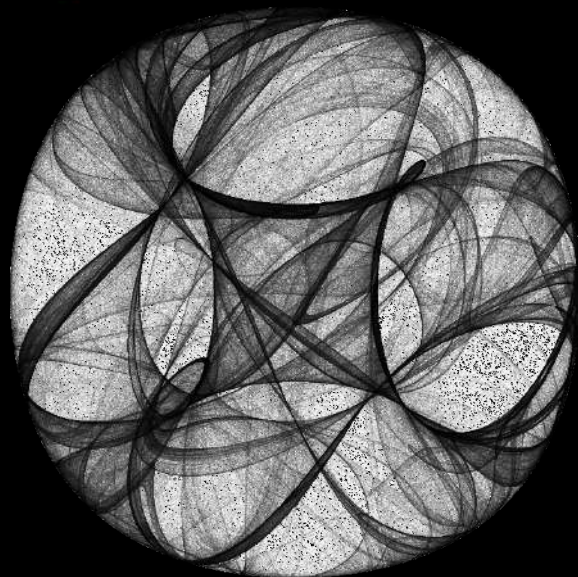
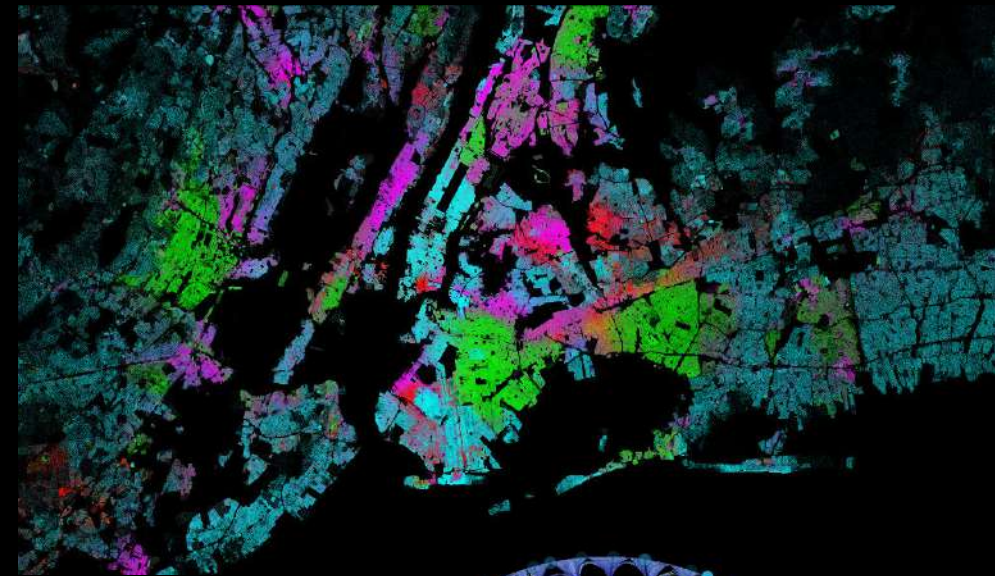
- **Open Source** and commercial products
- Python, **Scientific visualization**
 - Matplotlib
 - Bokeh
 - Datashader
- JavaScript, Publishing **interactive visualization online**
 - D3JS
- **No code** tools
 - RAWgraphs



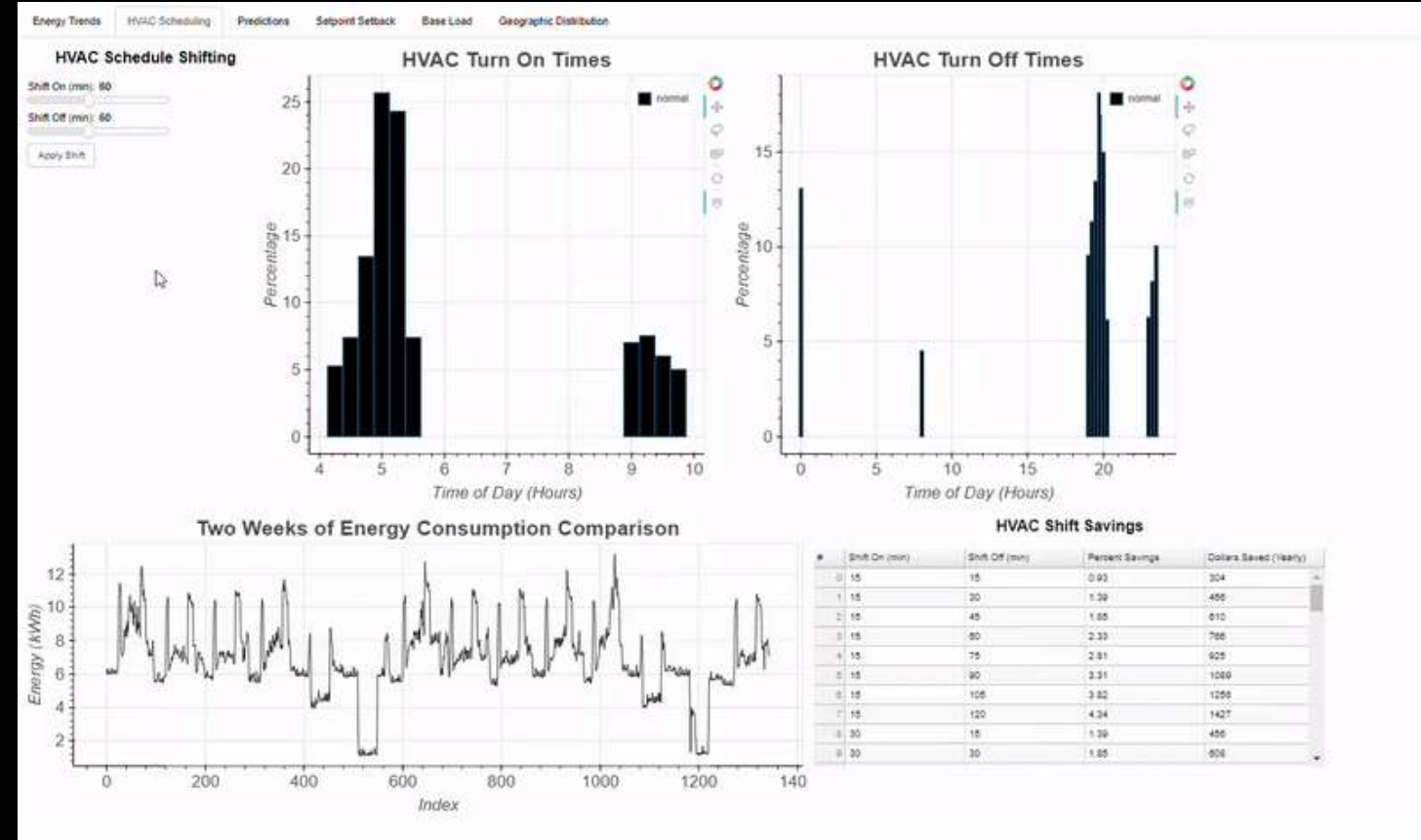
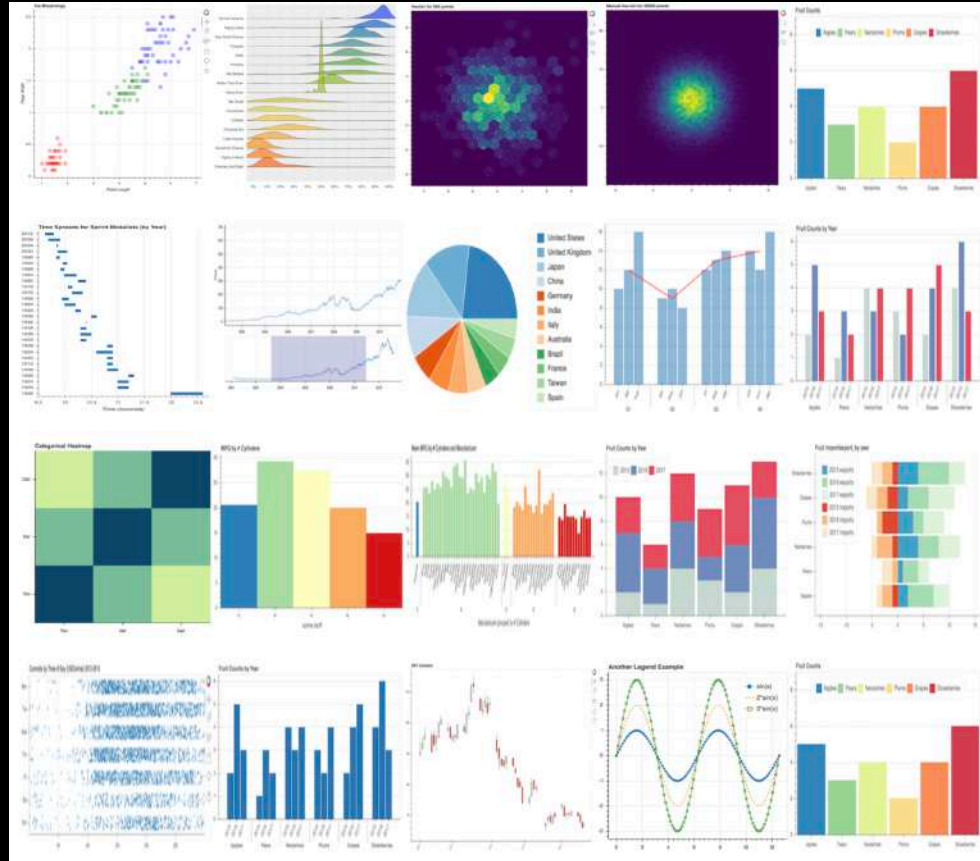
Matplotlib



Datashader



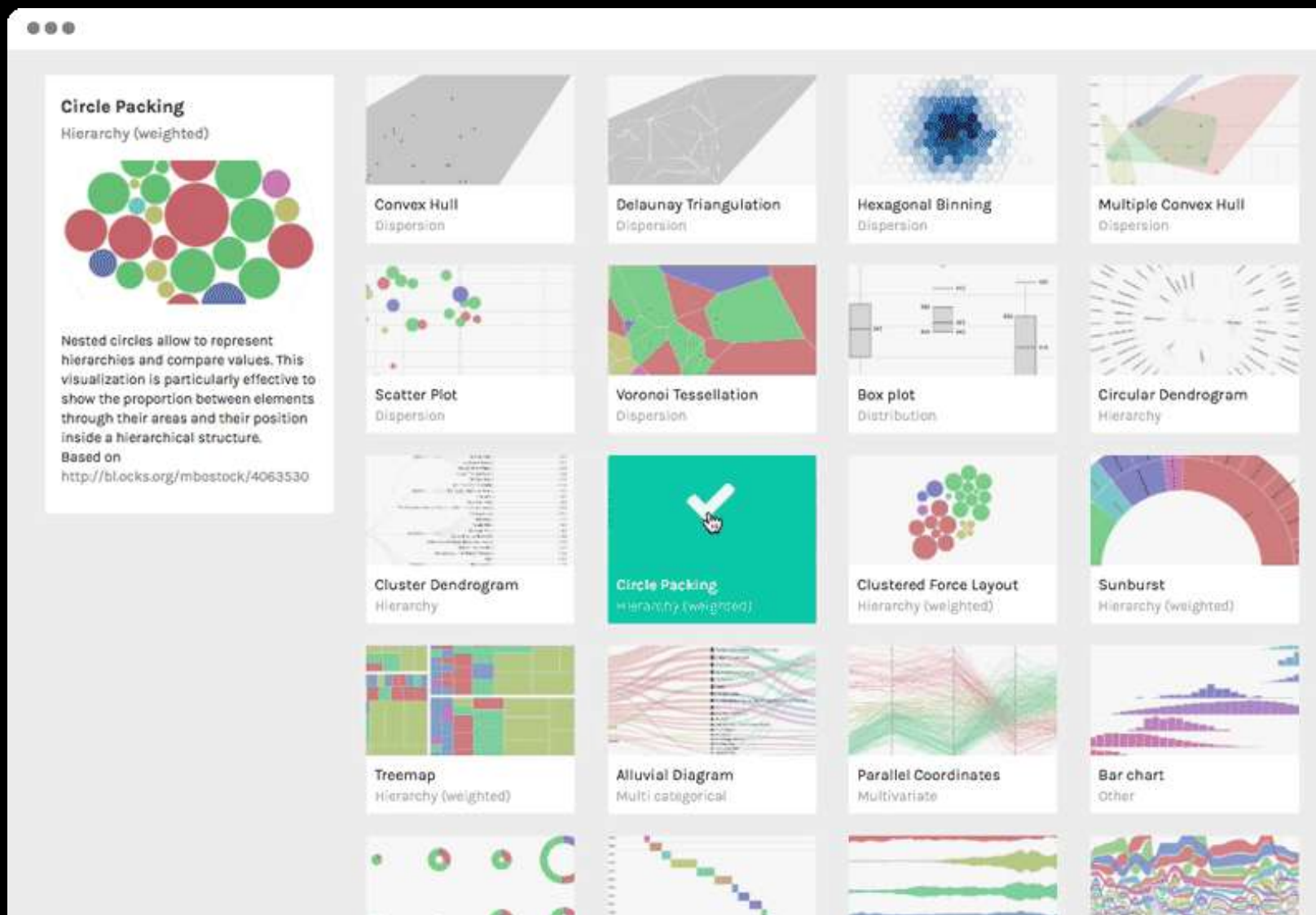
Bokeh



JavaScript. D3JS



RAWgraphs



The background of the slide is a dense, intricate network graph. It consists of numerous small, green circular nodes connected by thin, green lines representing edges. The connections are complex and non-uniform, creating a web-like structure that fills the entire frame. The overall aesthetic is technical and data-driven.

4. Graph Visualization

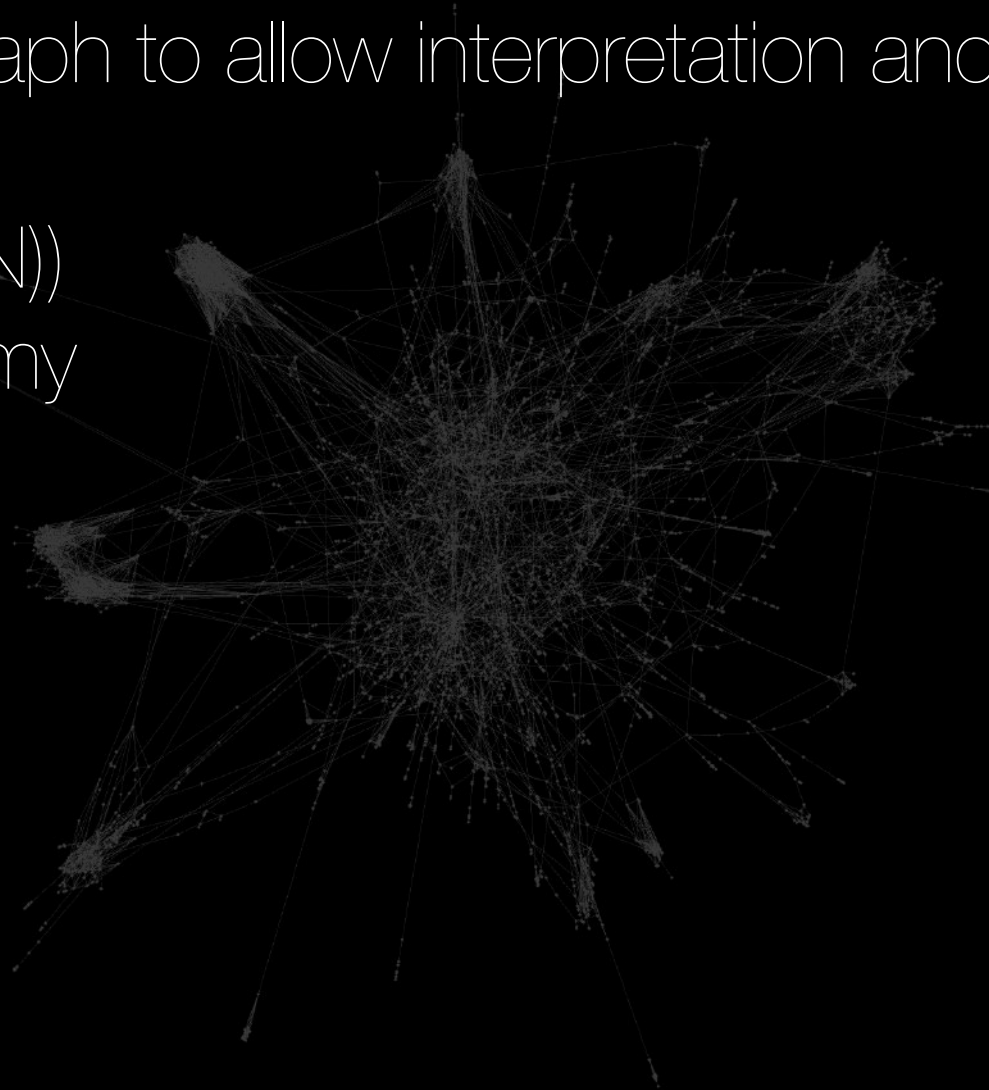
Some information appearing in this section is taken from Gephi Tutorial on Layouts by Gephi contributors
<https://gephi.org/about/people/>

Graph Visualization

- **Layouts**. Getting **visual insights** about **the structure**
 - Force-directed
 - Circular (radial)
 - Splitters
- **Attributes**. Highlighting graph **properties, attributes and clusters**
 - Community detection
 - Graph metrics
 - Dynamics

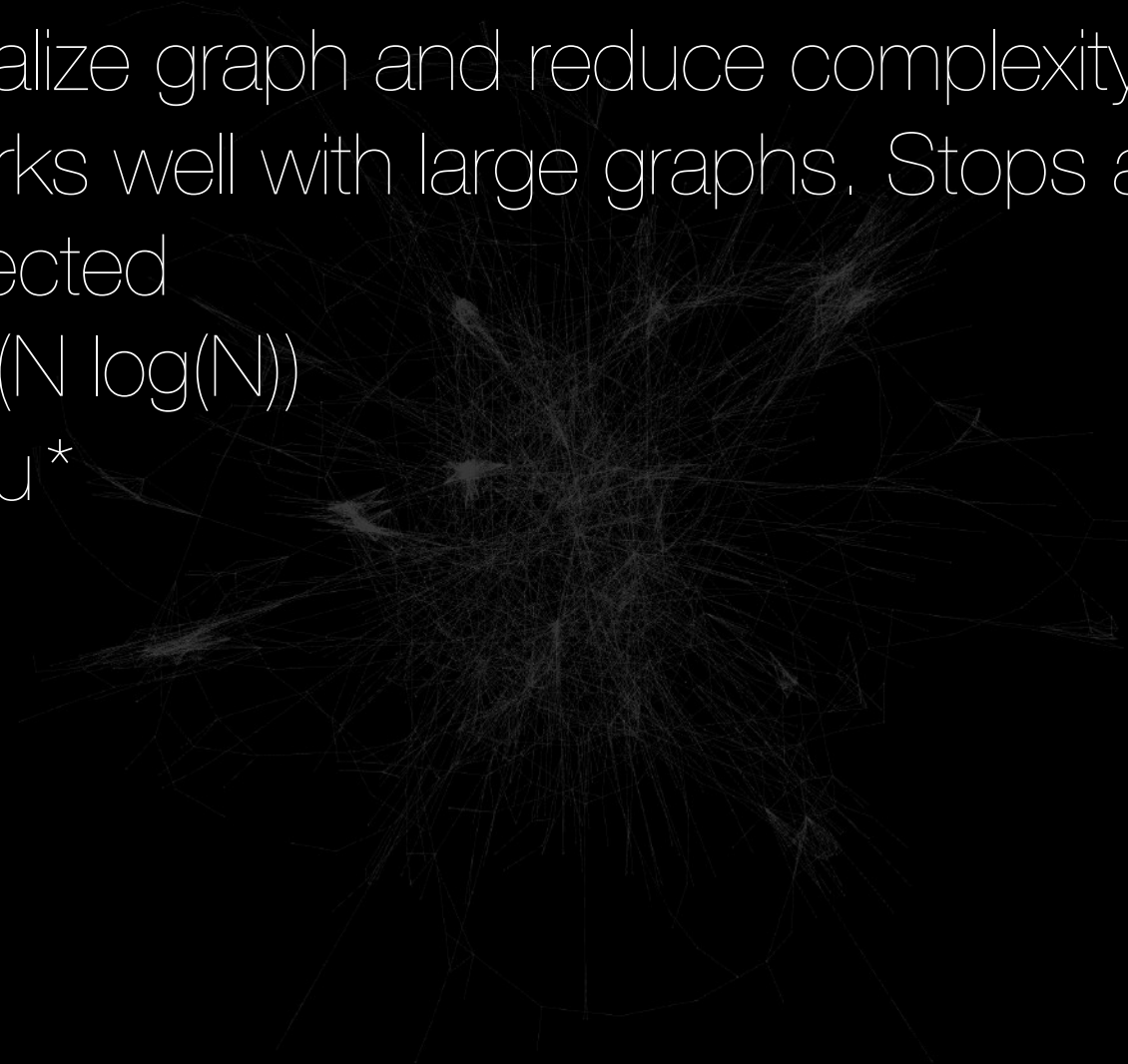
Layouts. Force Atlas (and variations)

- **Purpose:** Spatialize graph to allow interpretation and good readability
- **Kind:** Force-directed
- **Complexity:** $O(N \log(N))$
- **Author:** Mathieu Jacomy



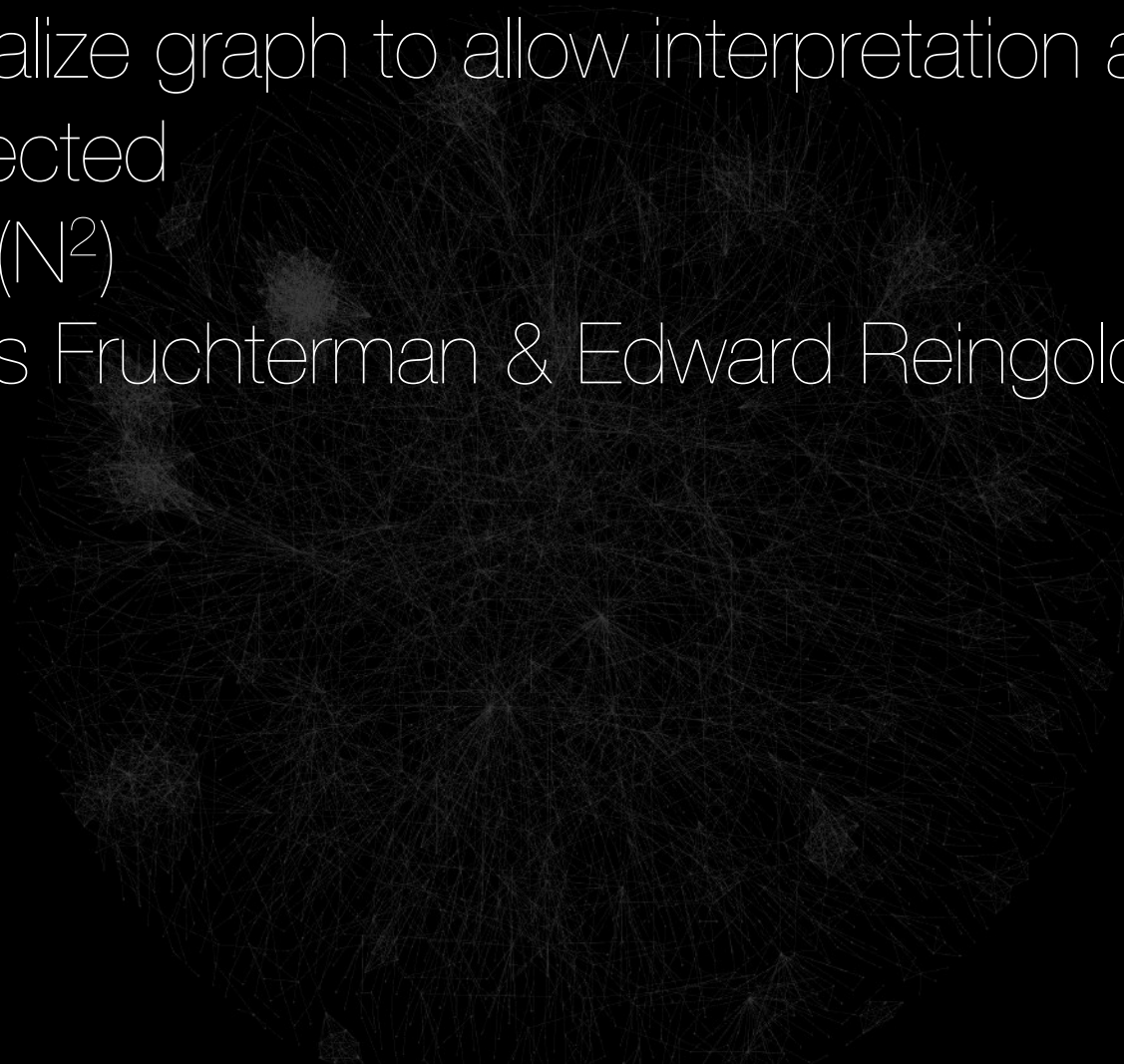
Layouts. Yifan Hu

- **Purpose:** Spatialize graph and reduce complexity with graph coarsening. Works well with large graphs. Stops automatically
- **Kind:** Force-directed
- **Complexity:** $O(N \log(N))$
- **Author:** Yifan Hu*



Layouts. Fruchterman-Reingold

- **Purpose:** Spatialize graph to allow interpretation and good readability
- **Kind:** Force-directed
- **Complexity:** $O(N^2)$
- **Author:** Thomas Fruchterman & Edward Reingold *



Layouts. OpenOrd

- **Purpose:** Distinguish clusters. Cut long edges to separate clusters better
- **Kind:** Force-directed
- **Complexity:** $O(N \log(N))$
- **Author:** S. Martin, W. M. Brown, R. Klavans, and K. Boyack *



Attributes. Community detection

- **Purpose:** Detect strongly-connected communities.
- **Kind:** Community detection algorithm
- **Complexity:** $O(N \log N)$
- **Author:** Louvain * (Vincent D. Blondel, Jean-Loup Guillaume, Renaud Lambiotte, Etienne Lefebvre), Leiden ** (V.A. Traag, L. Waltman, and N.J. van Eck)

* Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, Etienne Lefebvre, Fast unfolding of communities in large networks, in Journal of Statistical Mechanics: Theory and Experiment 2008 (10), P1000

** V.A. Traag, L. Waltman, and N.J. van Eck, From Louvain to Leiden: guaranteeing well-connected communities, Nature Scientific Reports, 2019

Attributes. Node size

- Degree
- Centrality metrics
 - Betweenness centrality
 - Bridging centrality
- Clustering coefficient
- Page rank

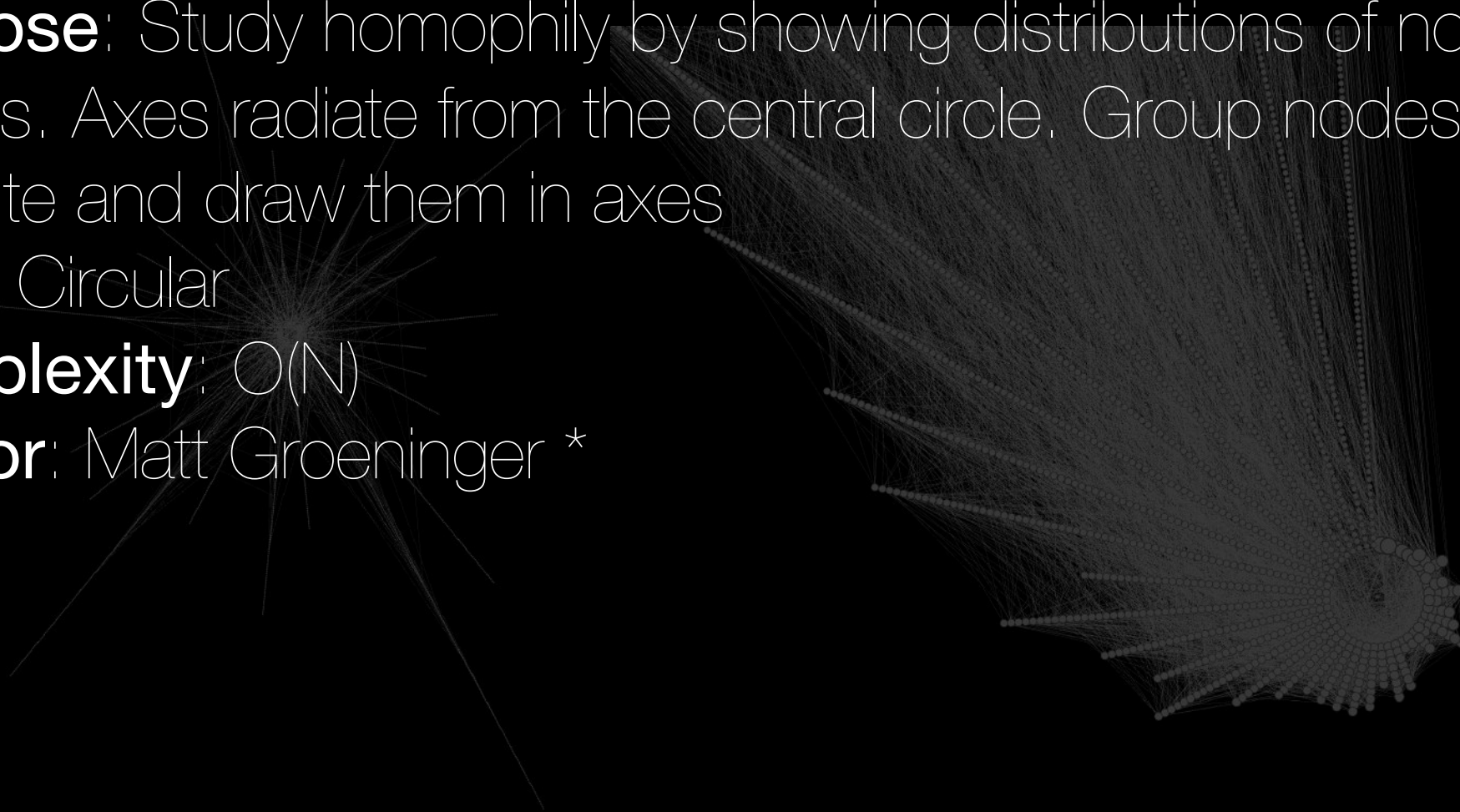


Layouts. Circular

- **Purpose:** Show distribution of nodes and their links. Order nodes by an attribute and draw them on a circle
- **Kind:** Circular
- **Complexity:** $O(N)$
- **Author:** Matt Groeninger *

Layouts. Radial Axis

- **Purpose:** Study homophily by showing distributions of nodes inside groups. Axes radiate from the central circle. Group nodes by an attribute and draw them in axes
- **Kind:** Circular
- **Complexity:** $O(N)$
- **Author:** Matt Groeninger *



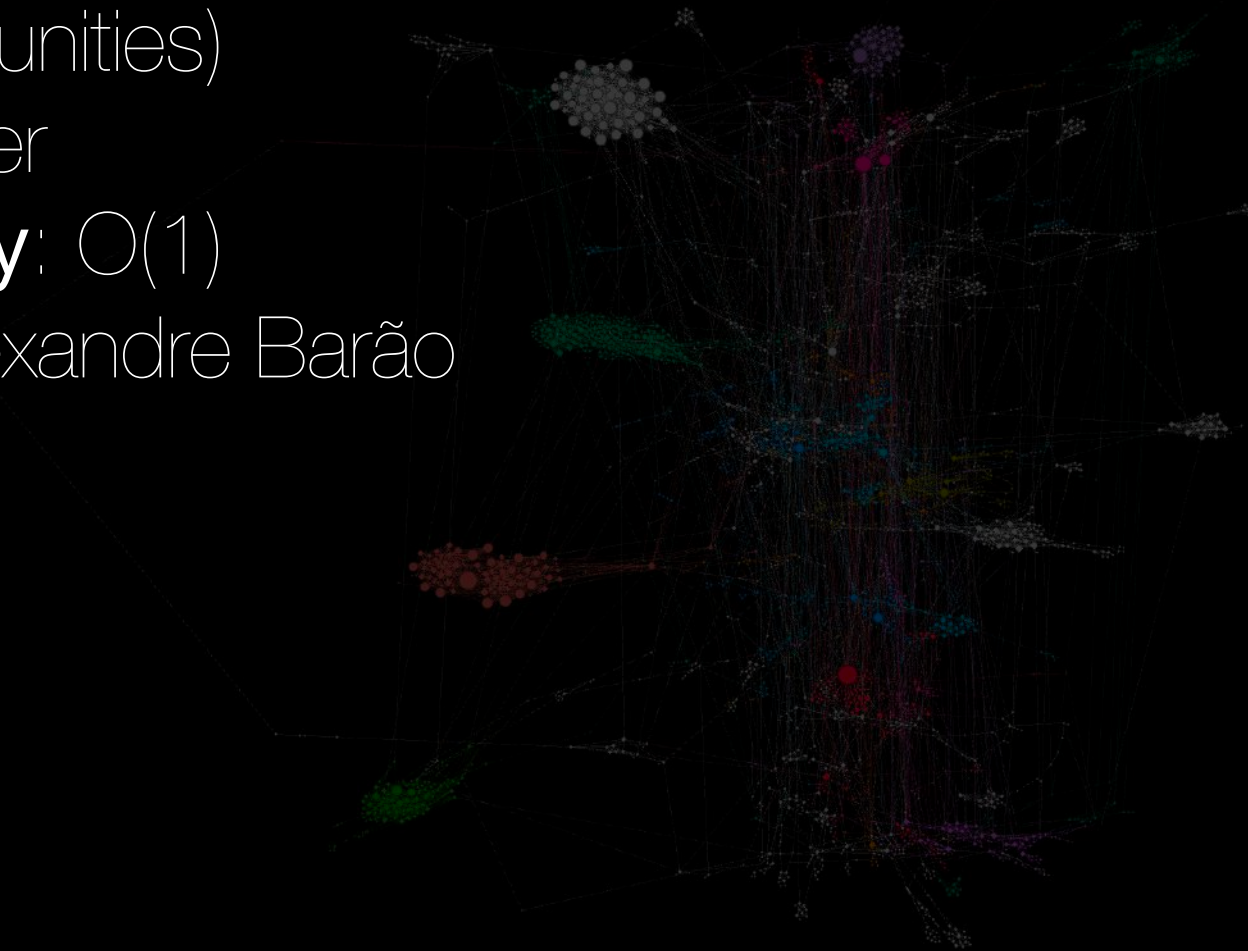
Layouts. Circle pack

- **Purpose:** Spatialize graph using attributes (e.g. community)
- **Kind:** Splitter
- **Complexity:** $O(1)$
- **Author:** Mike Bostock, Patrick Murray, Ken Goulding

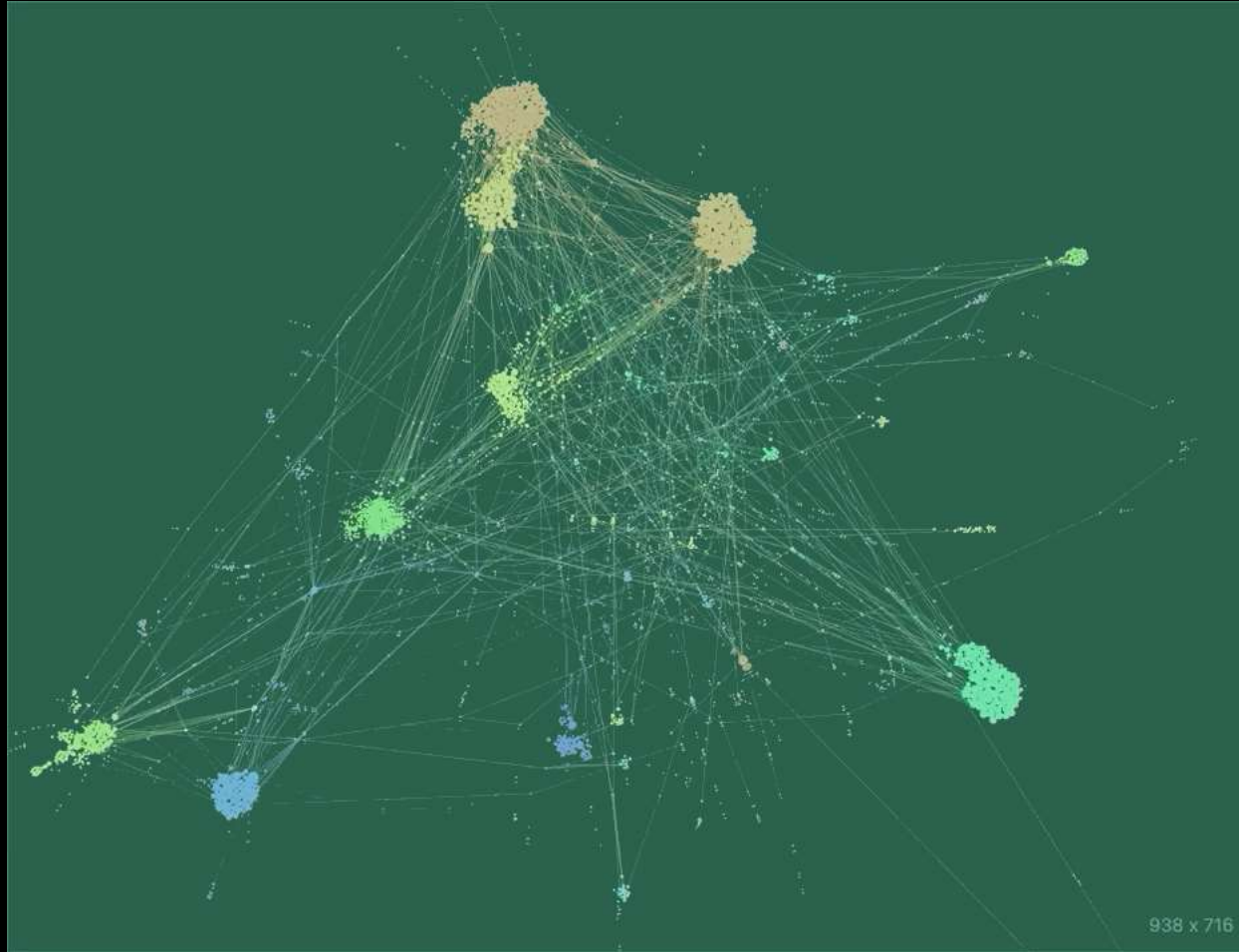


Layouts. 3D Network Splitter

- **Purpose:** Unfold graph in layers based on user-defined attributes (e.g. communities)
- **Kind:** Splitter
- **Complexity:** $O(1)$
- **Author:** Alexandre Barão



Attributes. Dynamics





5. Tutorials

Graph Visualization Tutorials

- Overview of graph visualization tools
 - GraphTool
 - Cytoscape
 - Gephi
- Gephi Tutorials

Tools. GraphTool



graph-tool

Efficient network analysis

Download

Documentation

Mailing List

Git

Issues

What is graph-tool?

Graph-tool is an efficient [Python](#) module for manipulation and statistical analysis of [graphs](#) (a.k.a. [networks](#)). Contrary to most other python modules with similar functionality, the core data structures and algorithms are implemented in [C++](#), making extensive use of template [metaprogramming](#), based heavily on the [Boost Graph Library](#). This confers it a level of [performance](#) that is comparable (both in memory usage and computation time) to that of a pure C/C++ library.

Download version 2.29 

[See Instructions](#) | [See Changelog](#)

▶▶ It is *Fast*!

Despite its nice, soft outer appearance of a regular python module, the core algorithms and data structures of graph-tool are written in C++, with performance in mind. Most of the time, you can expect the algorithms to run just as fast as if graph-tool were a pure C/C++ library. See a [performance comparison](#).

🌐 OpenMP Support

Many algorithms are implemented in parallel using [OpenMP](#), which provides excellent performance on multi-core architectures, without degrading it on single-core machines.

📖 Extensive Features

An extensive array of features is included, such as support for arbitrary vertex, edge or graph [properties](#), efficient "on the fly" [filtering](#) of vertices and edges, powerful graph I/O using the [GraphML](#), [GML](#) and [dot](#) file formats, graph [pickling](#), [graph statistics](#) (degree/property histogram, vertex correlations, average shortest distance, etc.), [centrality measures](#), standard [topological algorithms](#) (isomorphism, minimum spanning tree, connected components, dominator tree, [maximum flow](#), etc.), [generation of random graphs](#) with arbitrary degrees and correlations, [detection of modules and communities](#) via statistical inference, and much more.

👁 Powerful Visualization

Conveniently [draw](#) your graphs, using a variety of algorithms and output formats (including to the screen). Graph-tool has its own layout algorithms and versatile, interactive drawing routines based on [cairo](#) and [GTK+](#), but it can also work as a very comfortable interface to the excellent [graphviz](#) package.

📖 Fully Documented

Every single function in the module is documented in the docstrings and in the online [documentation](#), which is full of examples.

Tools. Cytoscape

A screenshot of the Cytoscape website. The header includes navigation links: Intro, Download, Apps, Documentation, Community, Report a Bug, Help, and a Google Custom Search bar. The main content area features the Cytoscape logo (an orange network graph) and the word 'Cytoscape' in large white text. Below this is the tagline 'Network Data Integration, Analysis, and Visualization in a Box'. Two prominent buttons are displayed: 'Introduction' (a white button with a black border) and 'Download 3.7.2' (an orange button). The background of the website is a dark blue network graph with numerous nodes and edges, some of which are labeled with gene symbols like BRAD1, MDM2, and BRAP.

Tools. Gephi



[Download](#) [Blog](#) [Wiki](#) [Forum](#) [Support](#) [Bug tracker](#)

[Home](#) [Features](#) [Learn](#) [Develop](#) [Plugins](#) [Services](#) [Consortium](#)

The Open Graph Viz Platform

Gephi is the leading visualization and exploration software for all kinds of graphs and networks. Gephi is open-source and free.

Runs on Windows, Mac OS X and Linux.

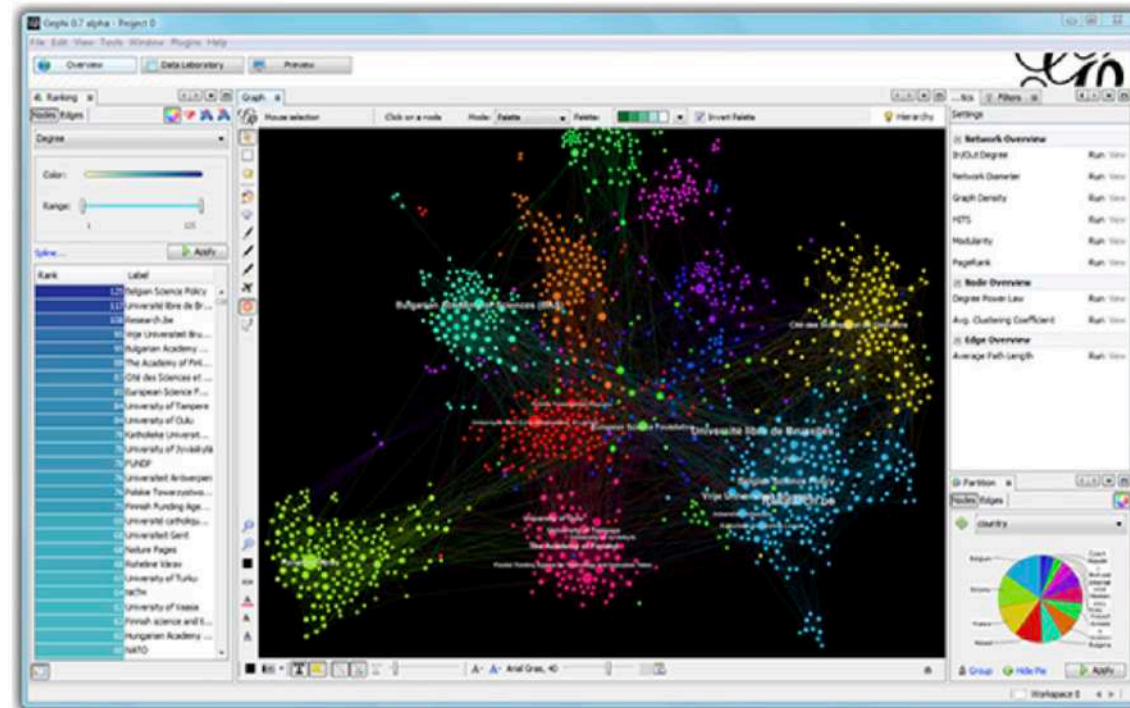
[Learn More on Gephi Platform »](#)



[Release Notes](#) | [System Requirements](#)

► [Features](#)
► [Quick start](#)

► [Screenshots](#)
► [Videos](#)



Gephi pipeline

Load graph

gexf, graphml,
gdf, csv etc.

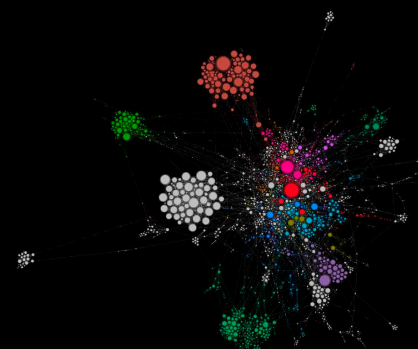
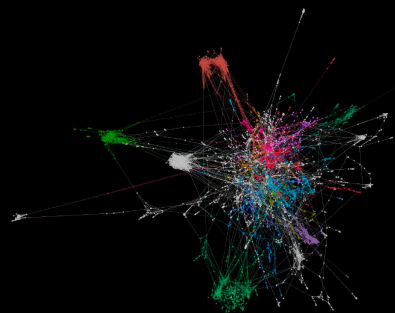
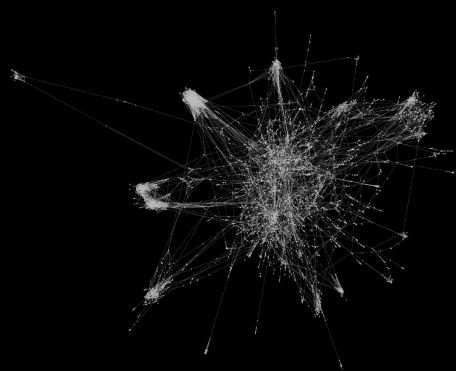
Compute
layout

Compute
attributes

Design
visualization

Publish
online

GitHub pages



GitHub Pages

Tutorials: <https://github.com/mizvol/gephi-tutorials>

Questions?