

Foundations of Data Science for Biologists

Best practices: data visualization

BIO 724D

18-MAR-2024

Instructors: Greg Wray and Paul Magwene

Examples of graphical excellence

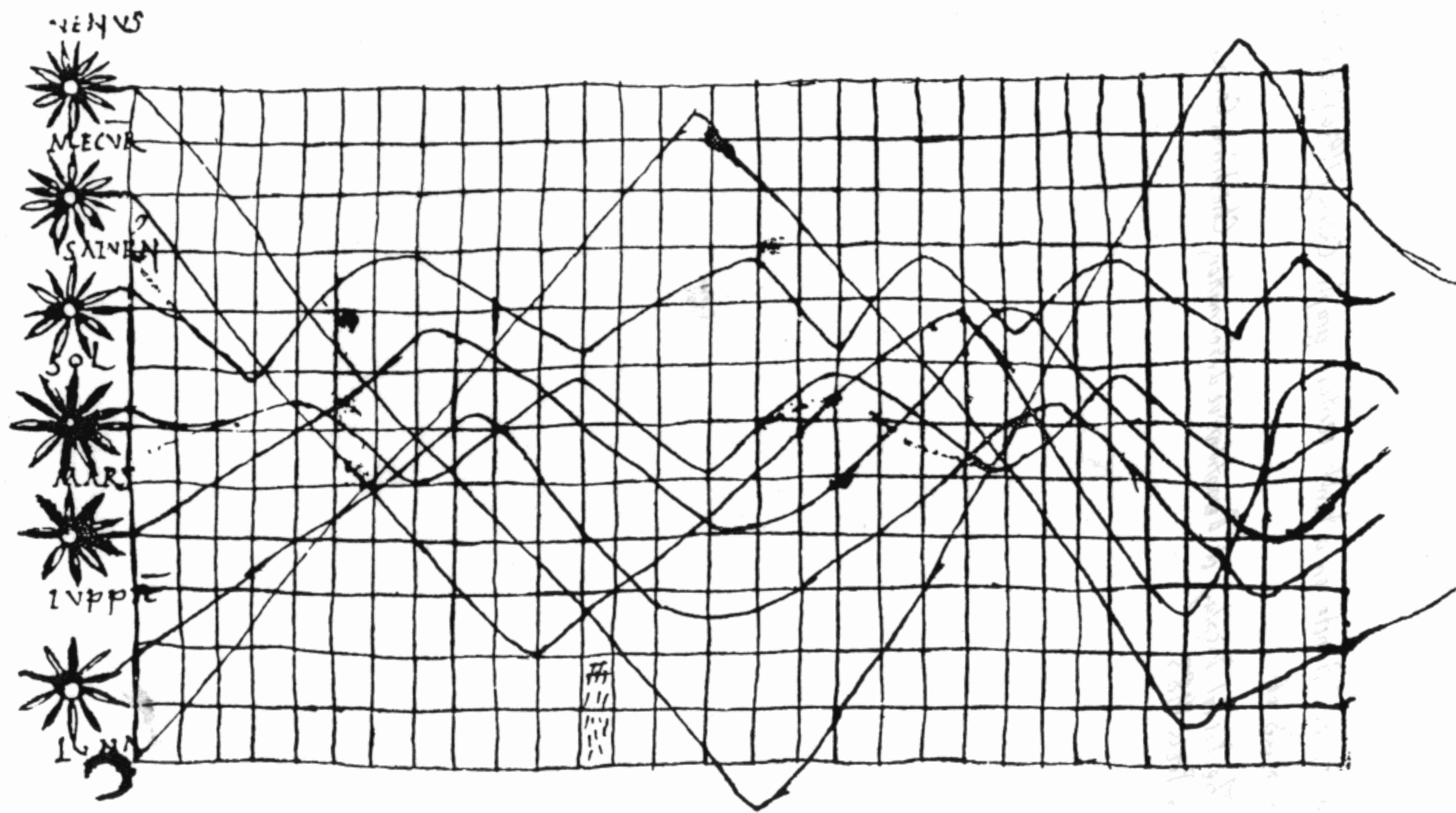


The earliest known abstract graphic? (right)

Chauvet-Pont-d'Arc Cave, France; unknown artists; 32,000 - 30,000 BCE



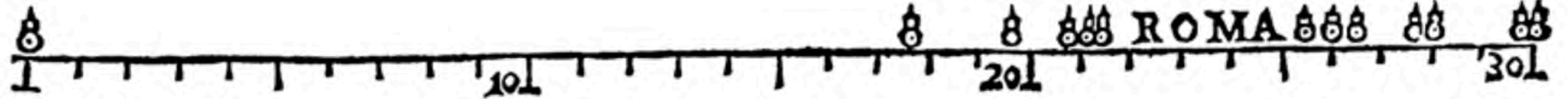
Earliest known “world map”
Babylon; unknown designer; early 6th century BCE



Earliest known “graph” — annotated as *De cursu per zodiacum*
Unknown designer; early 11th century

TOLEDO.

GRADOS DE LA LONGITUD.

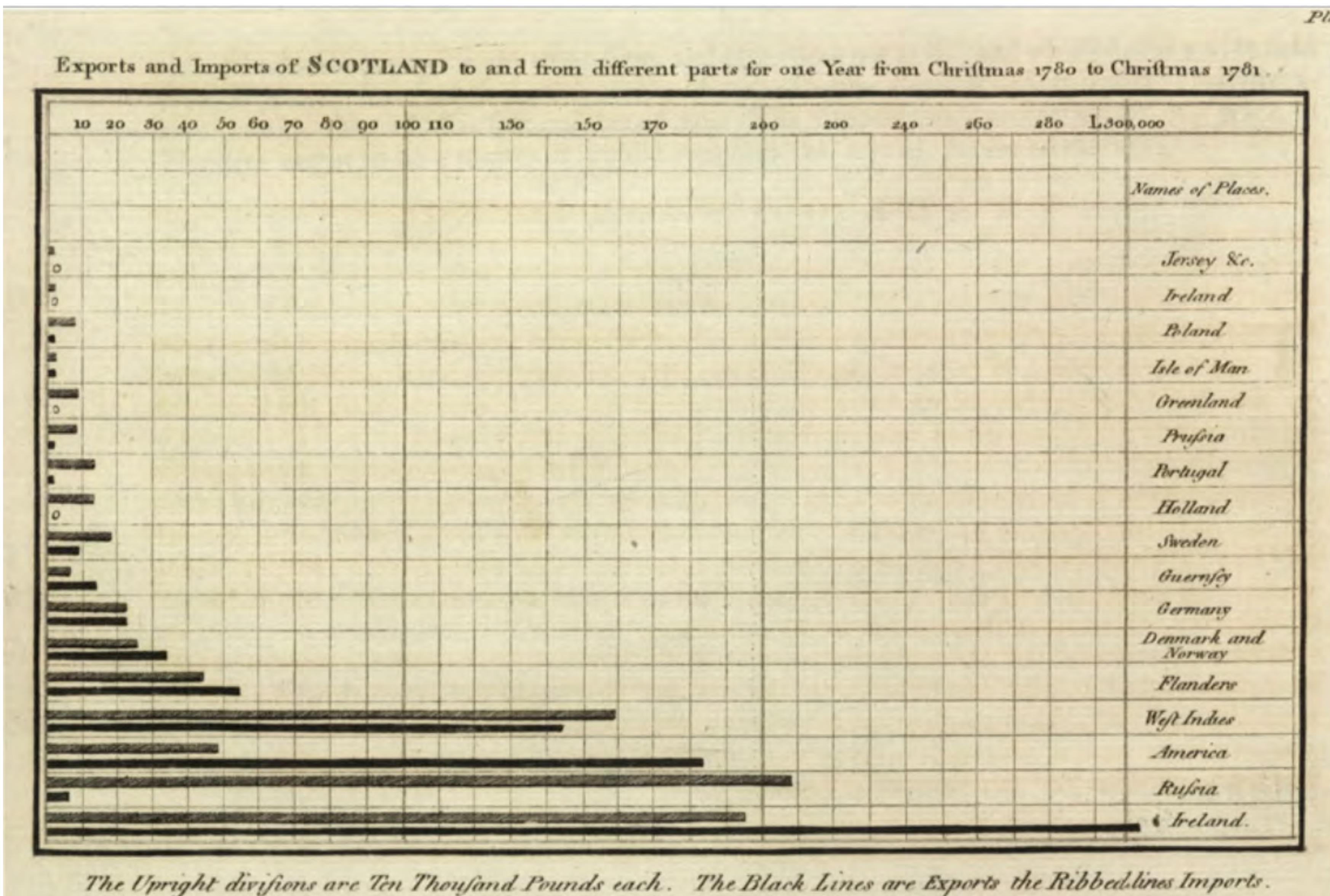
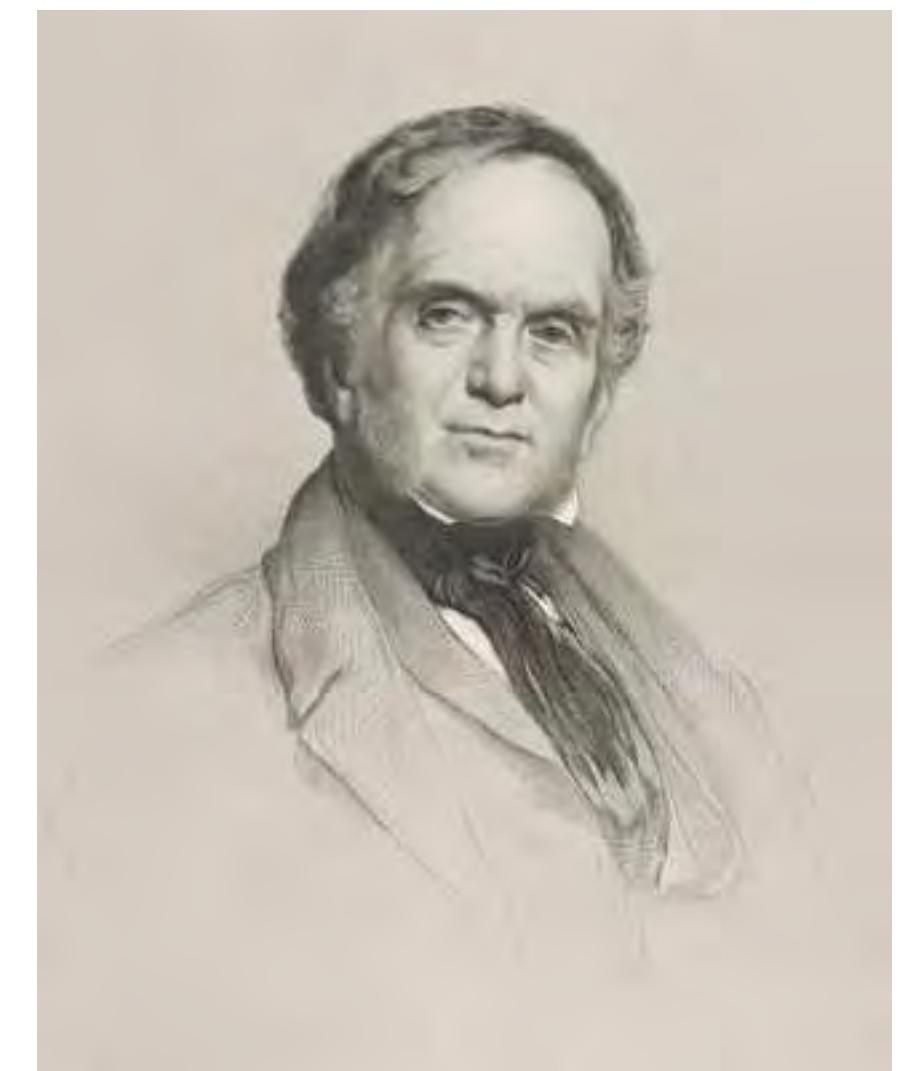


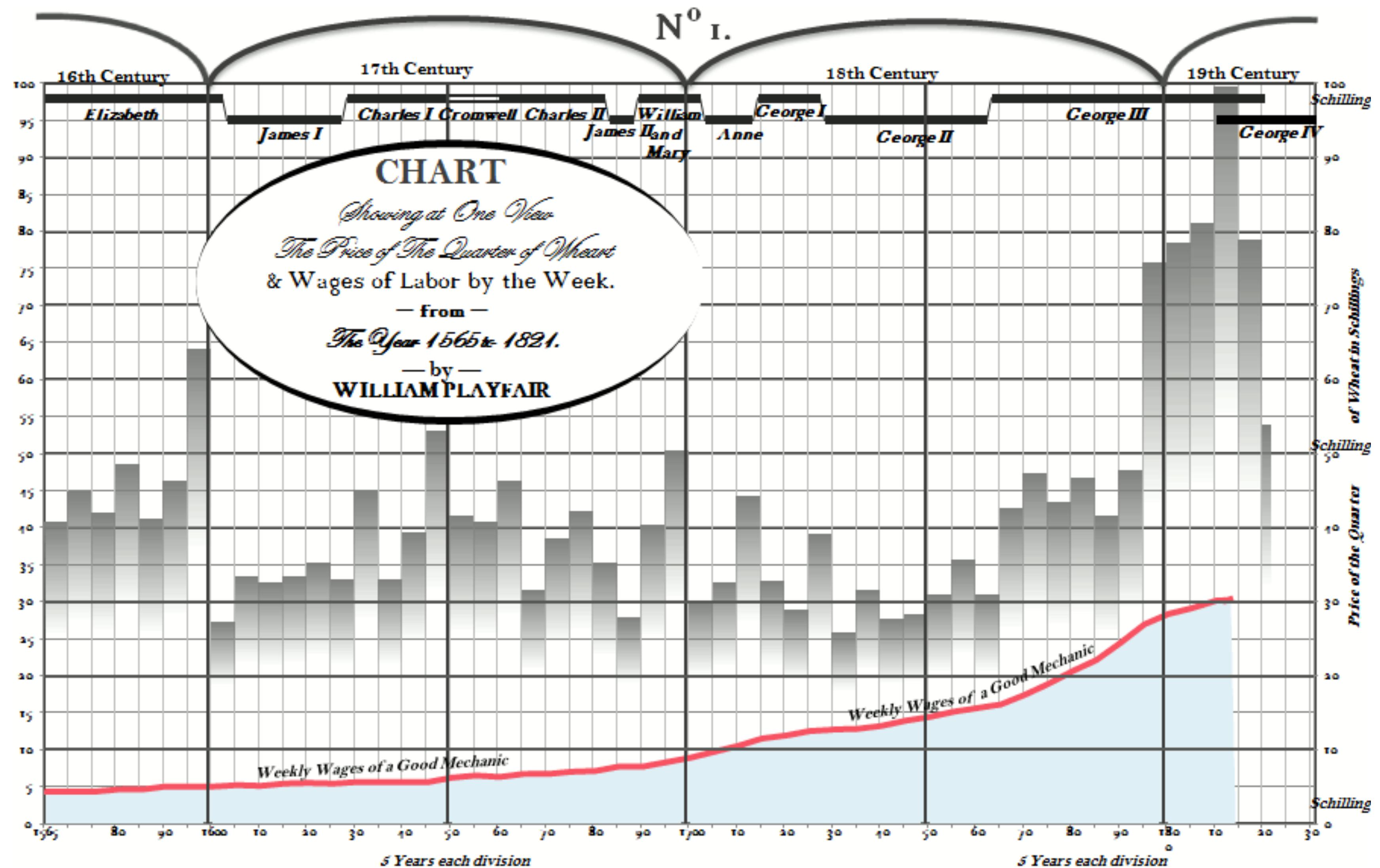
Earliest known statistical graphic

Michael Florent van Langren; 1644 (two earlier manuscript versions are known)

The first bar graph
William Playfair; 1781

(he also invented
the pie chart)



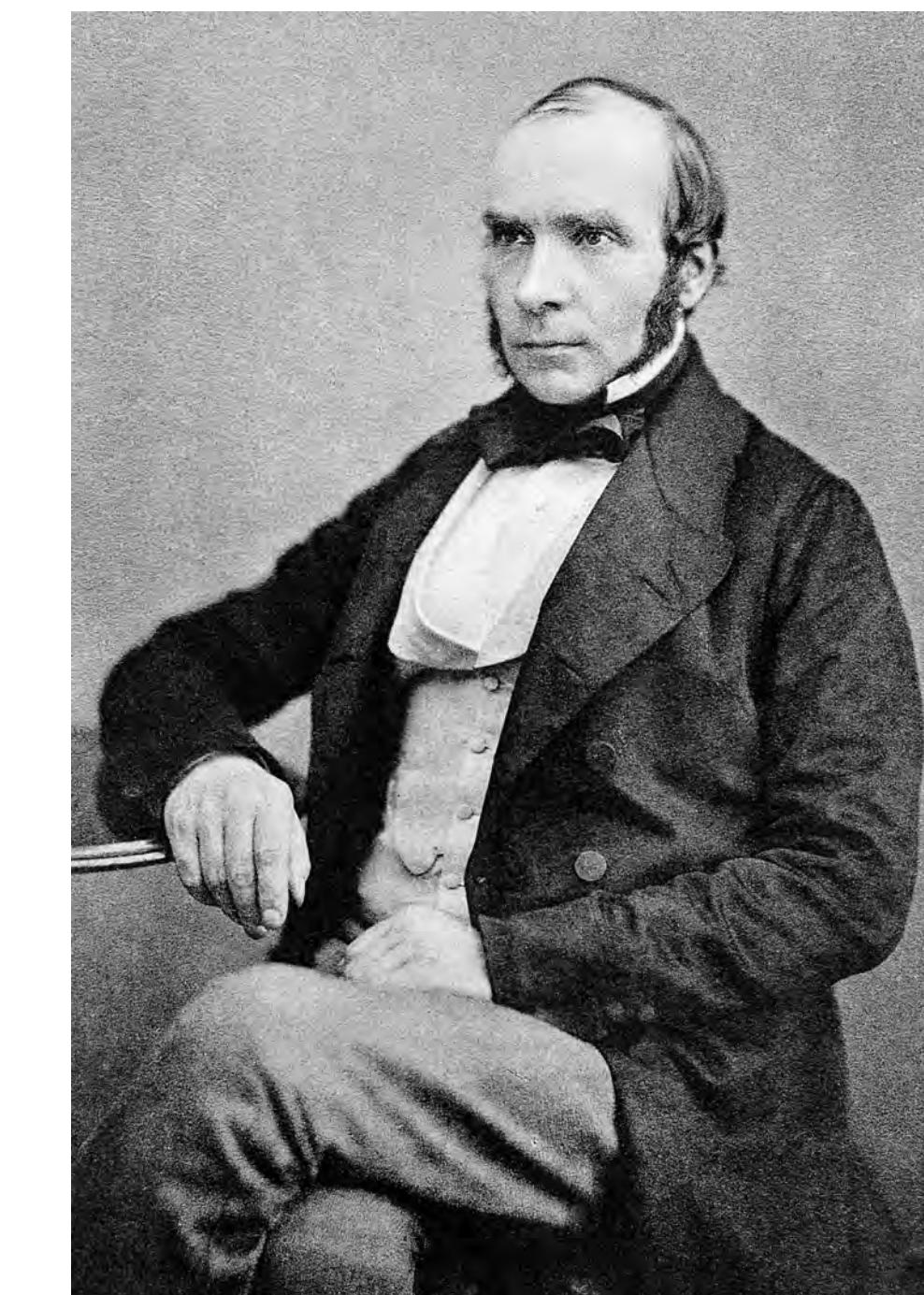


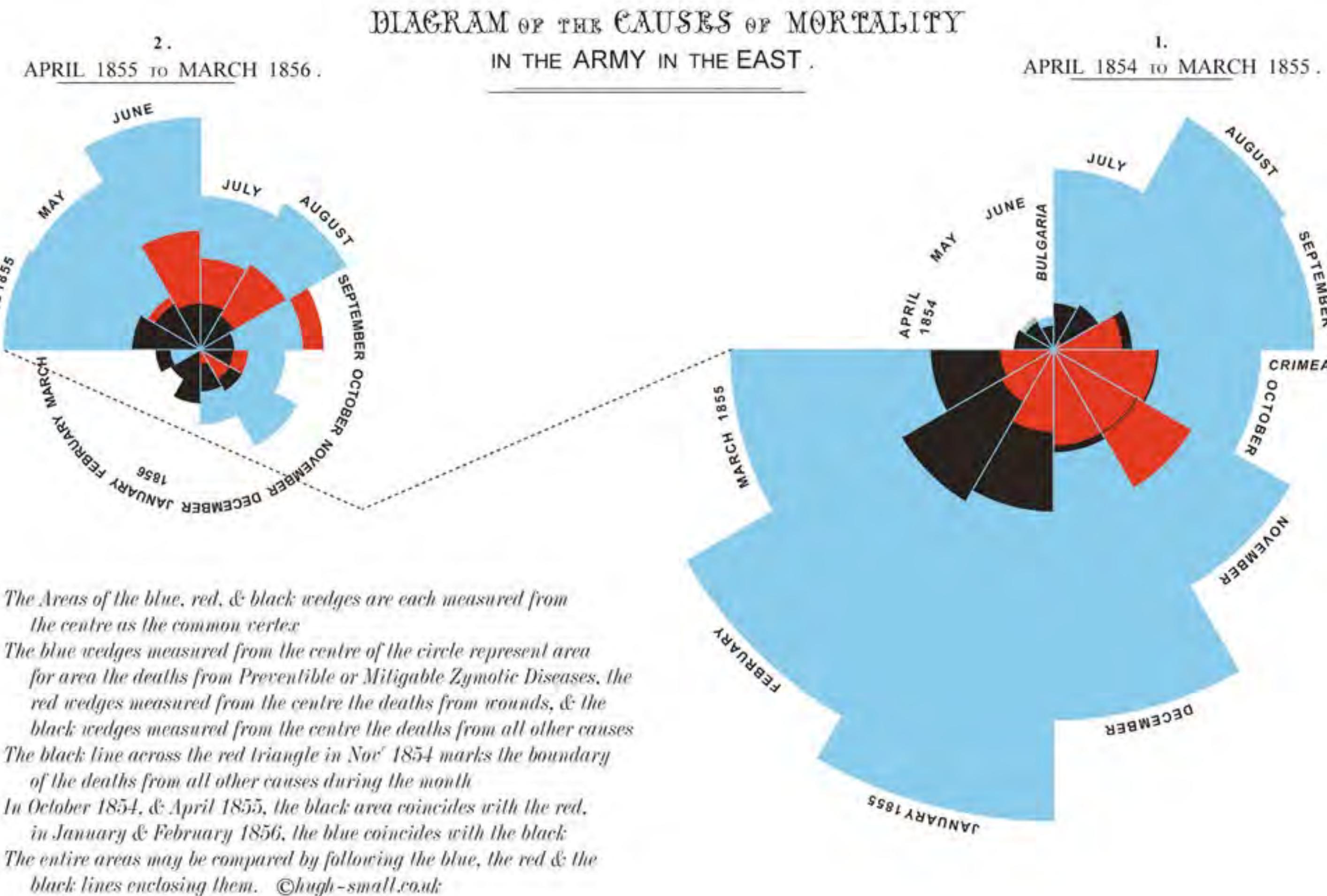
The invention of time-series graphs (also multivariate!) by Joseph Priestly and William Playfair
 William Playfair; 1821

50 0 50 100 150 200
Yards
X Pump • Deaths from cholera



The birth of epidemiology
John Snow; 1854





Invented the polar area diagram

Color and proportion



Policy reform and social justice through graphics
Florence Nightingale; 1858

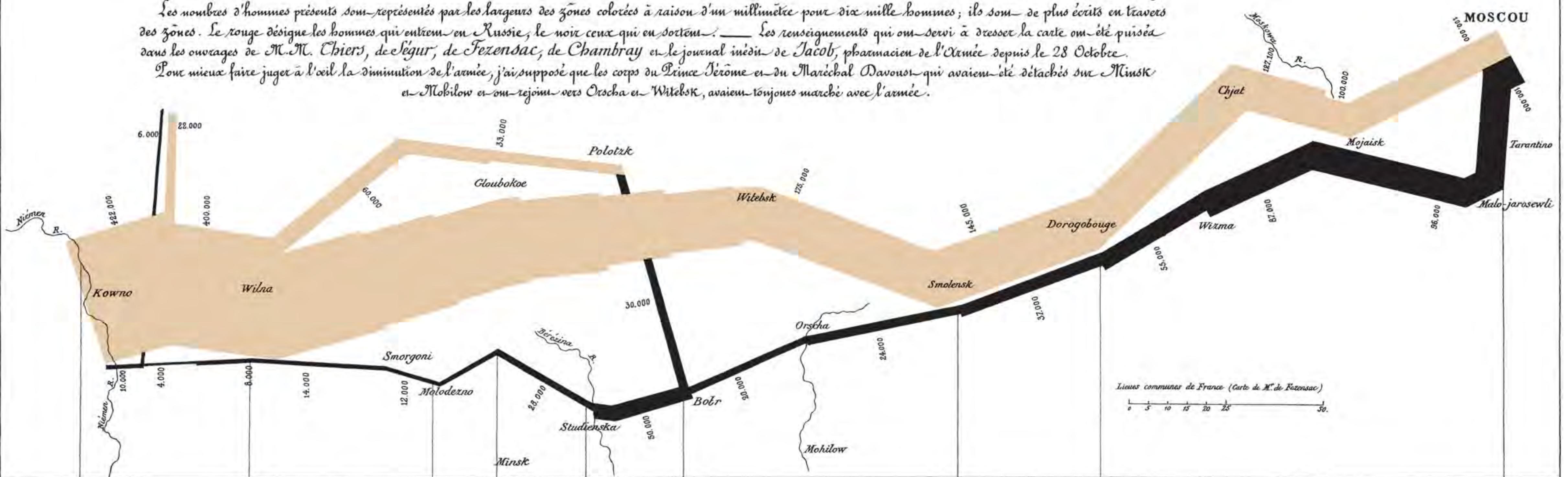
Carte Figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813.

Dressée par M. Minard, Inspecteur Général des Ponts et Chaussees en retraite.

Paris, le 20 Novembre 1869.

Les nombres d'hommes présents sont représentés par les largeurs des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en travers des zones. Le rouge désigne les hommes qui entrent en Russie, le noir ceux qui en sortent. — Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M. M. Chiers, de Séguir, de Fezensac, de Chambray et le journal médical de Jacob, pharmacien de l'Armée depuis le 28 Octobre.

Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps du Prince Jérôme et du Maréchal Davout qui avaient été détachés sur Minsk et Mohilow et qui rejoignirent vers Orscha et Witebsk, avaient toujours marché avec l'armée.



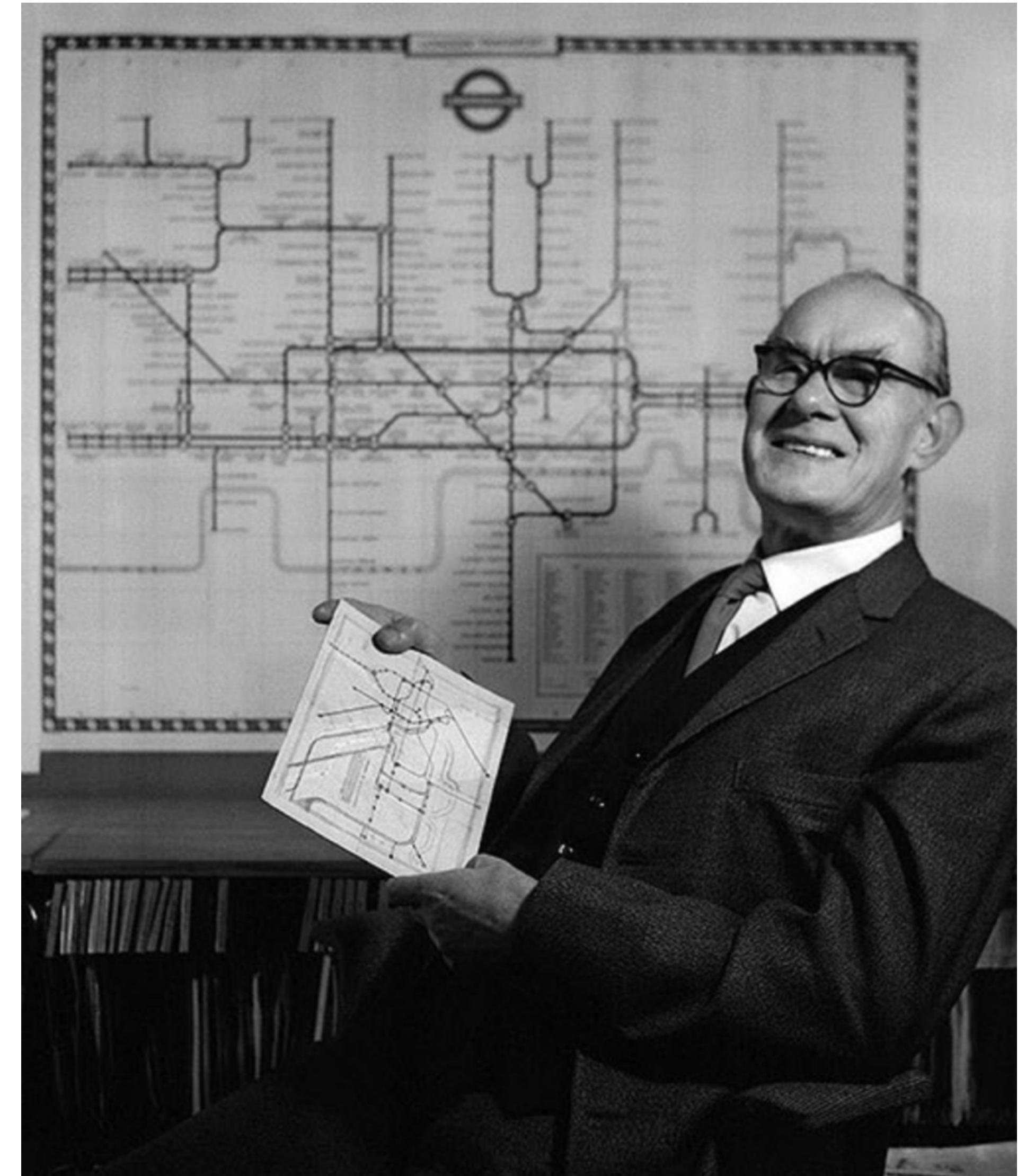


Enriching a univariate measure

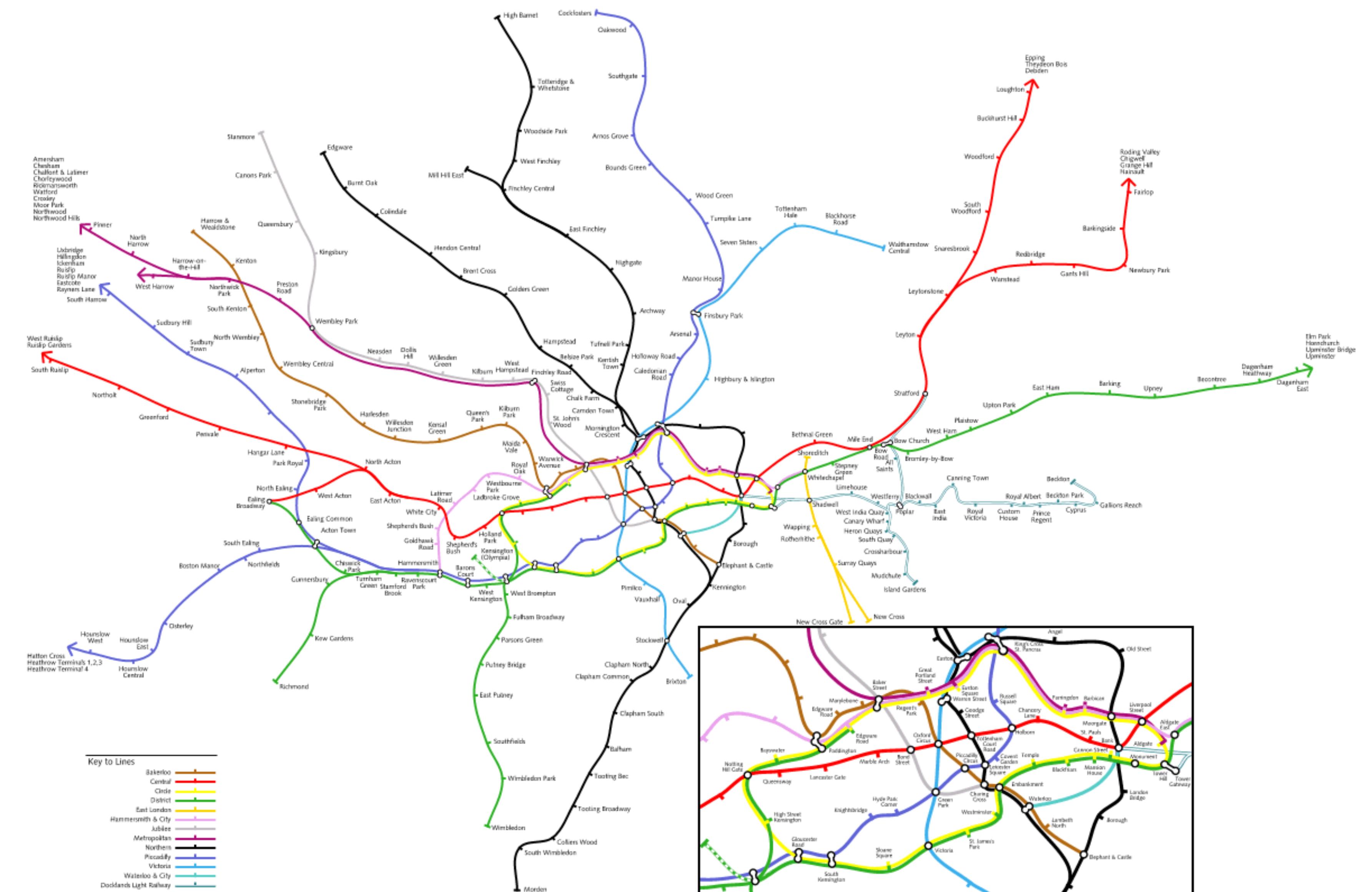
Joseph Hutchins Colton; mid 1800s



Practical, intelligent simplification



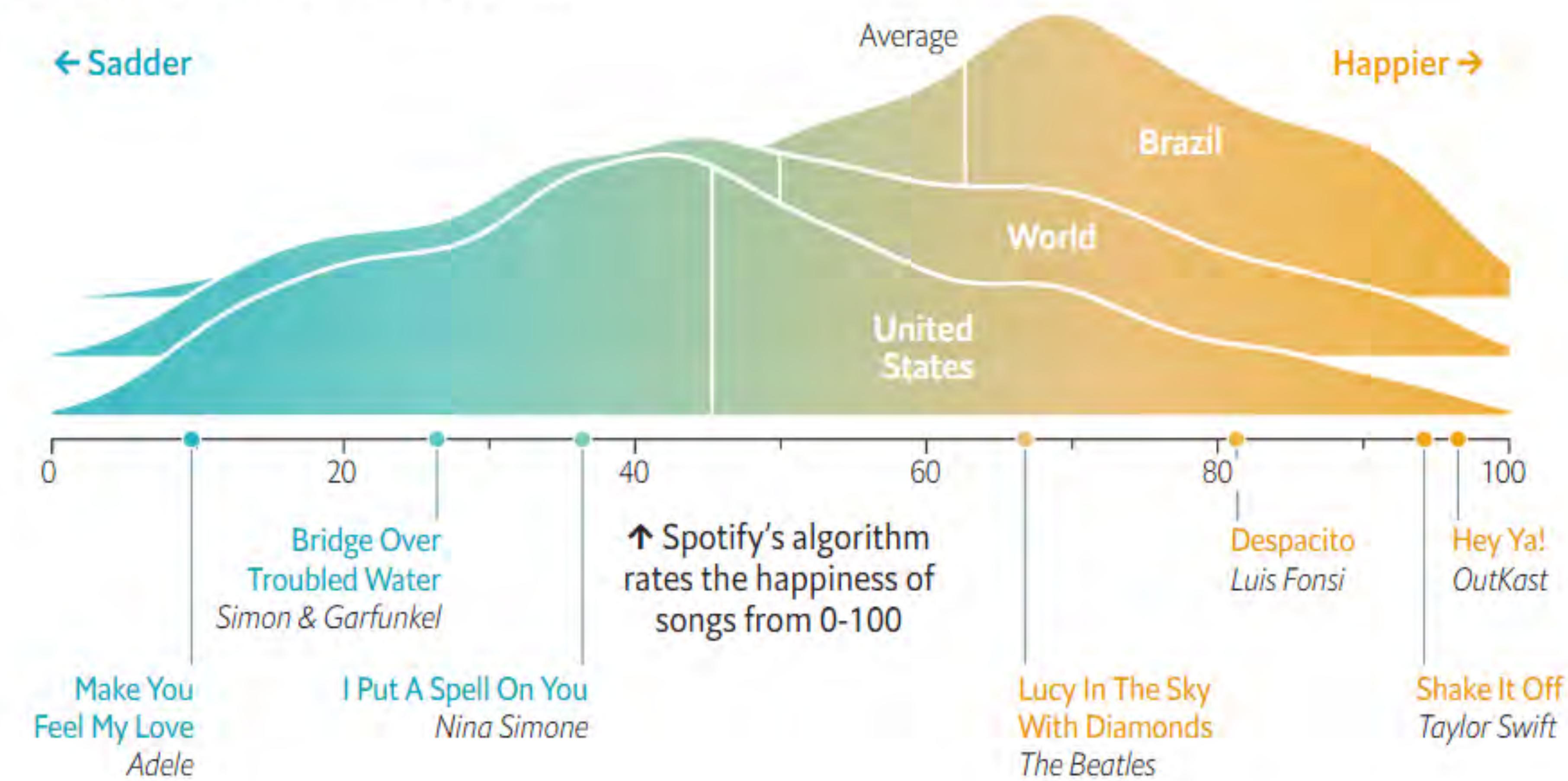
Henry "Harry" Beck; 1933



Geographically accurate London Underground map: not as easy to understand!

→ Some countries listen to happier music than others

Distribution of tracks streamed*, by mood

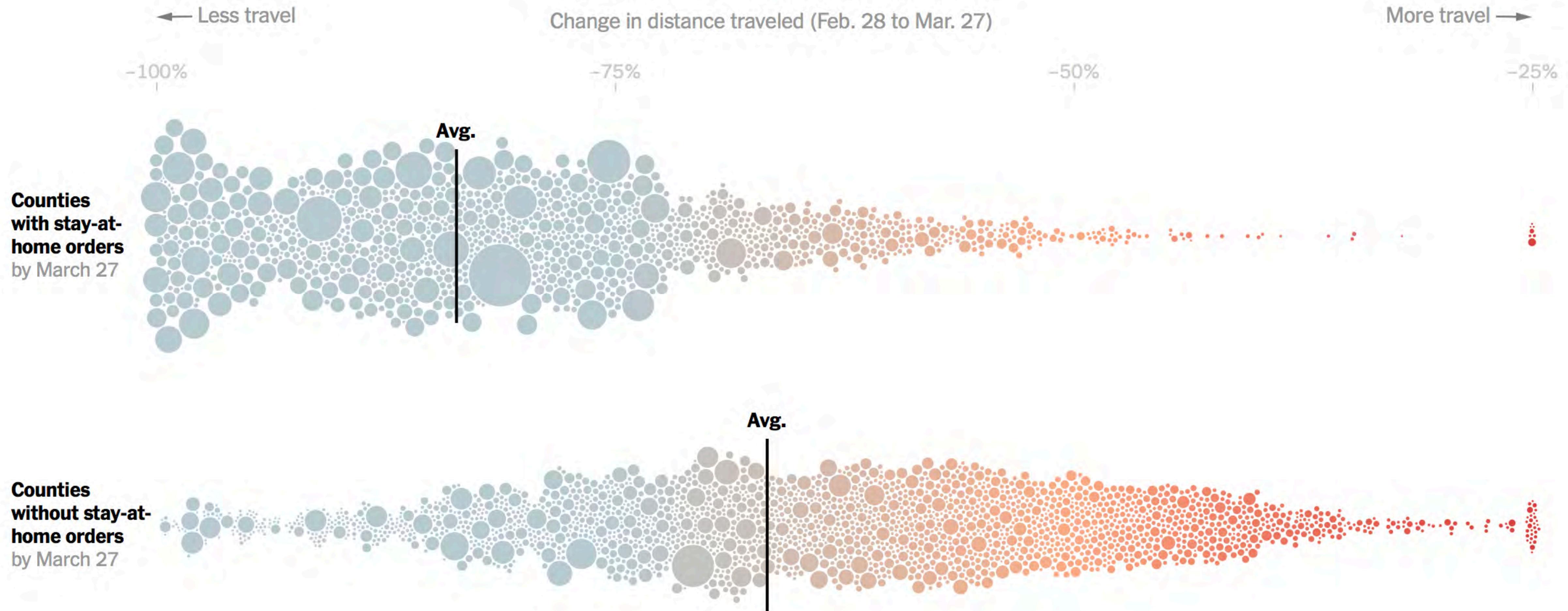


*200 most-streamed songs on each day, January 1st 2017-January 29th 2020

Economist; 2020

Which counties reduced travel the most

Change in distance traveled (Feb. 28 to Mar. 27)



New York Times; 2020

March of the Oaks

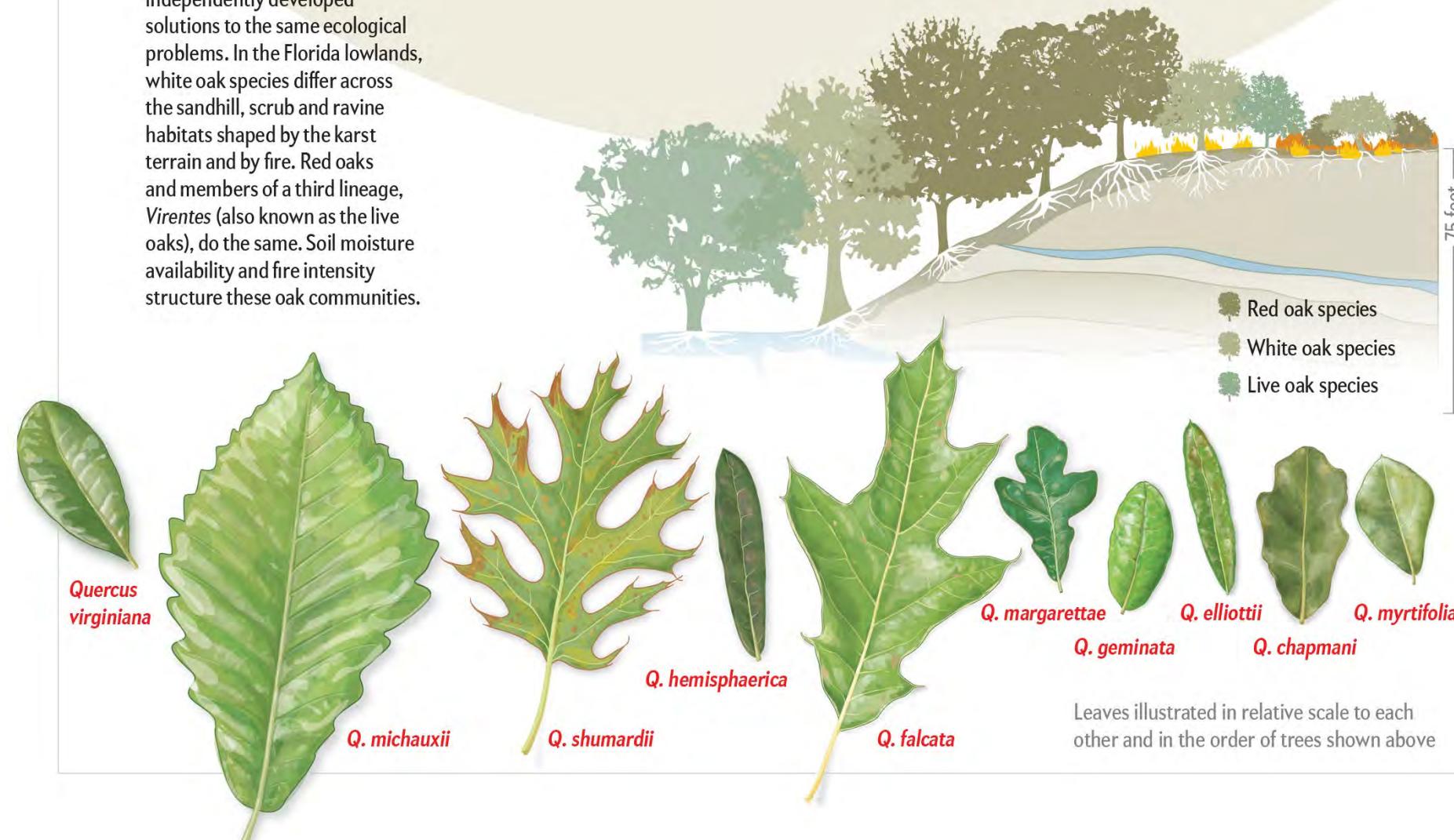
Over some 56 million years oaks have diversified into the 435 species alive today that together span five continents. Genome studies have allowed researchers to reconstruct the history of speciation in oaks. The findings help to explain how oaks came to be so diverse, particularly in the Americas, where some 60 percent of oak species reside.

Oak Classification

All living oak species are members of the genus *Quercus*, which comprises eight major lineages or sections, as they are termed. Two of these have dominated the Americas: the *Lobatae* section (also known as the red oaks) and the *Quercus* section (also known as the white oaks).

Diversity within Communities

Red oaks and white oaks often grow together in the same habitats. These two lineages colonized the same areas and independently developed solutions to the same ecological problems. In the Florida lowlands, white oak species differ across the sandhill, scrub and ravine habitats shaped by the karst terrain and by fire. Red oaks and members of a third lineage, *Virentes* (also known as the live oaks), do the same. Soil moisture availability and fire intensity structure these oak communities.



Leaves illustrated in relative scale to each other and in the order of trees shown above

Millions of years ago: 56

PALOEocene

33.9

EOCENE

23

OLIGOCENE

5.3

MIOCENE

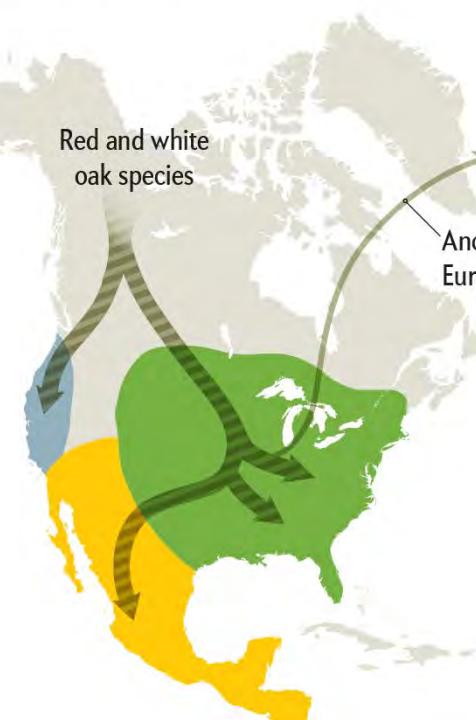
2.6

PLIOCENE

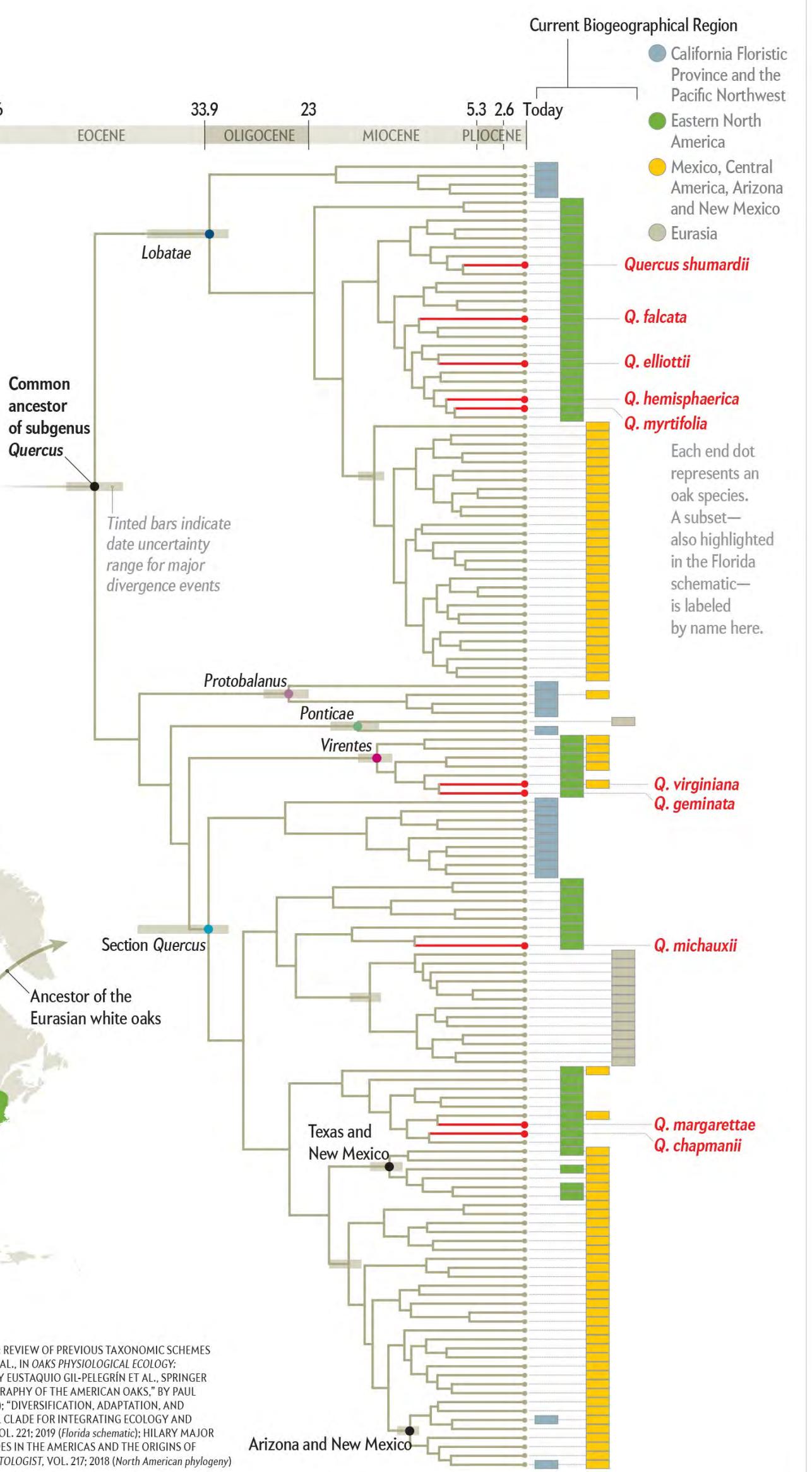
Today

Into the Americas

White oaks and red oaks arose and diversified simultaneously in the Americas. As these two groups moved south, each split into a lineage on the western side of the Rocky Mountains that gave rise to the oaks of California and the Pacific Northwest and into a lineage on the eastern side of the Rockies that gave rise to the oaks of eastern North America. On the eastern side, the red and white oaks each subdivided into northeastern, southeastern and Texan lineages. The red and white oaks then spread from eastern North America into Mexico, where they underwent another burst of diversification.



SOURCES: "AN UPDATED INFRAGENERIC CLASSIFICATION OF THE OAKS: REVIEW OF PREVIOUS TAXONOMIC SCHEMES AND SYNTHESIS OF EVOLUTIONARY PATTERNS," BY THOMAS DENK ET AL., IN *OAKS PHYSIOLOGICAL ECOLOGY: EXPLORING THE FUNCTIONAL DIVERSITY OF GENUS QUERCUS*; EDITED BY EUSTAQUIO GIL-PELEGRÍN ET AL., SPRINGER INTERNATIONAL PUBLISHING, 2017; AND "SYSTEMATICS AND BIOGEOGRAPHY OF THE AMERICAN OAKS," BY PAUL MANOS, IN *INTERNATIONAL OAKS*, VOL. 27: 2016 (ranges and classification); "DIVERSIFICATION, ADAPTATION, AND COMMUNITY ASSEMBLY OF THE AMERICAN OAKS (*QUERCUS*): A MODEL CLADE FOR INTEGRATING ECOLOGY AND EVOLUTION," BY JEANNINE CAVENDER-BARES, IN *NEW PHYTOLOGIST*, VOL. 221: 2019 (Florida schematic); HILARY MAJOR (leaves); "SYMPATRIC PARALLEL DIVERSIFICATION OF MAJOR OAK CLADES IN THE AMERICAS AND THE ORIGINS OF MEXICAN SPECIES DIVERSITY," BY ANDREW L. HIPP ET AL., IN *NEW PHYTOLOGIST*, VOL. 217: 2018 (North American phylogeny)



Infographic based on research done at Duke

Paul Manos; Biology Dept; 2020

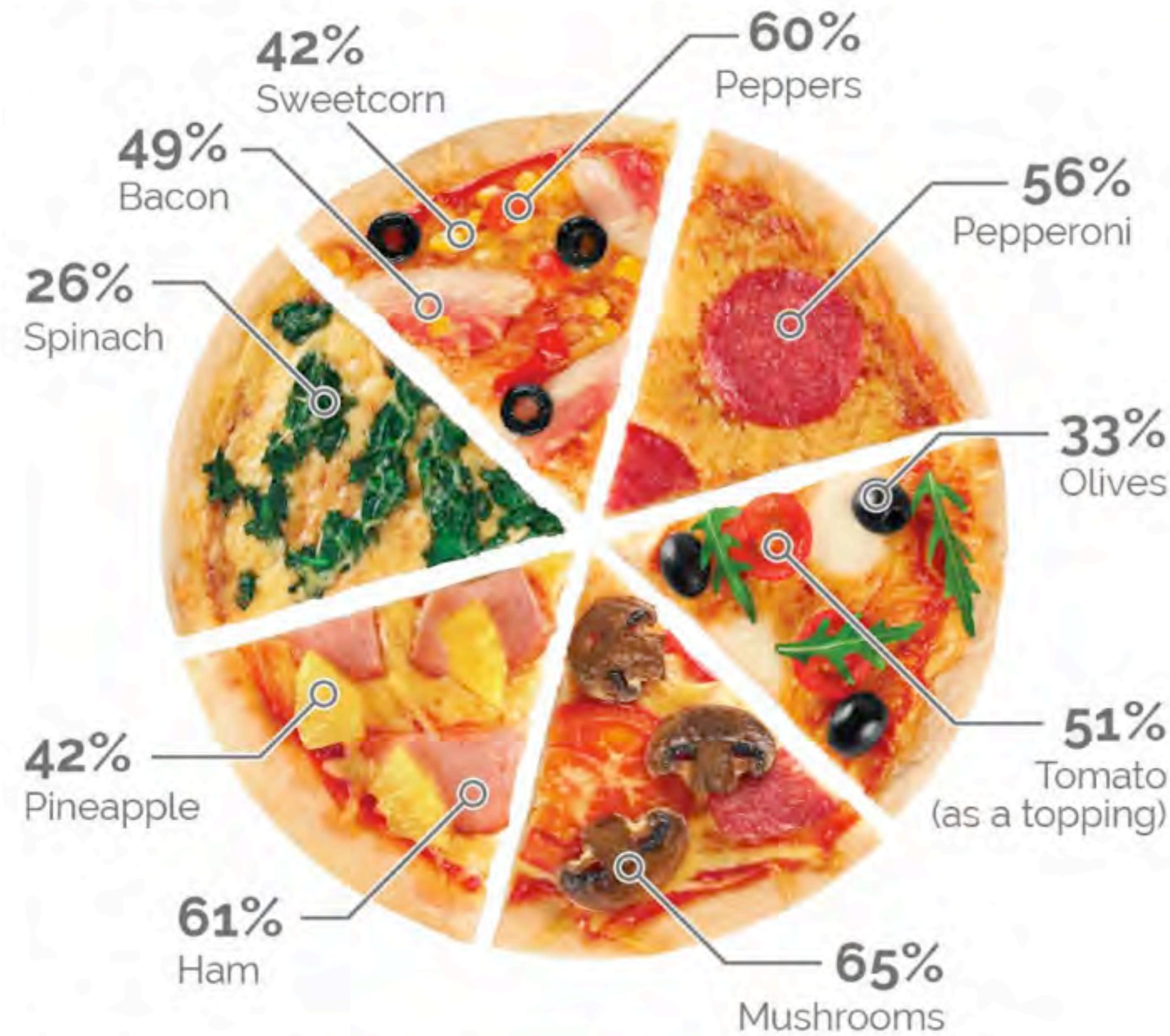
Graphical excellence and counter-examples

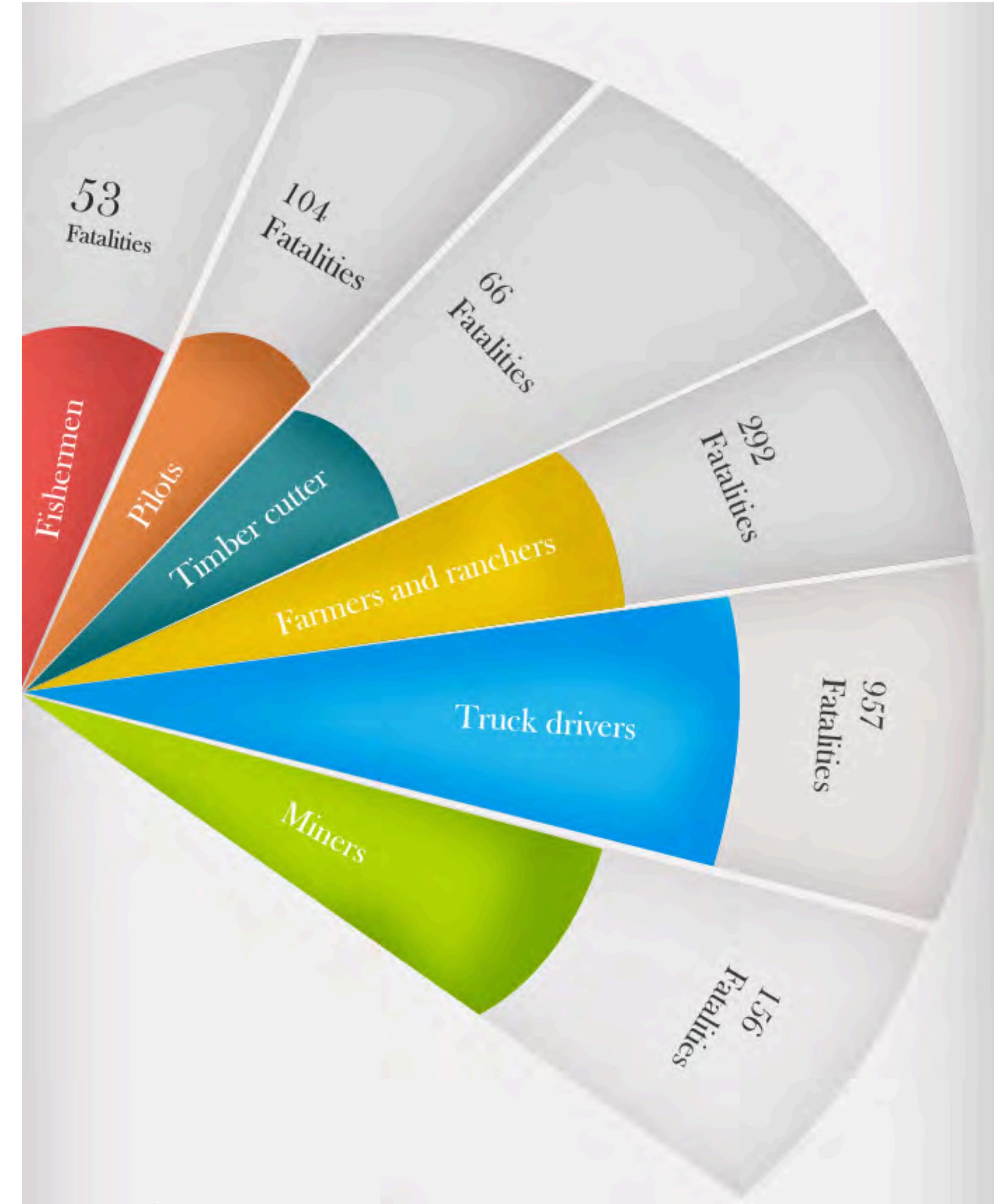
Tufte on graphical excellence

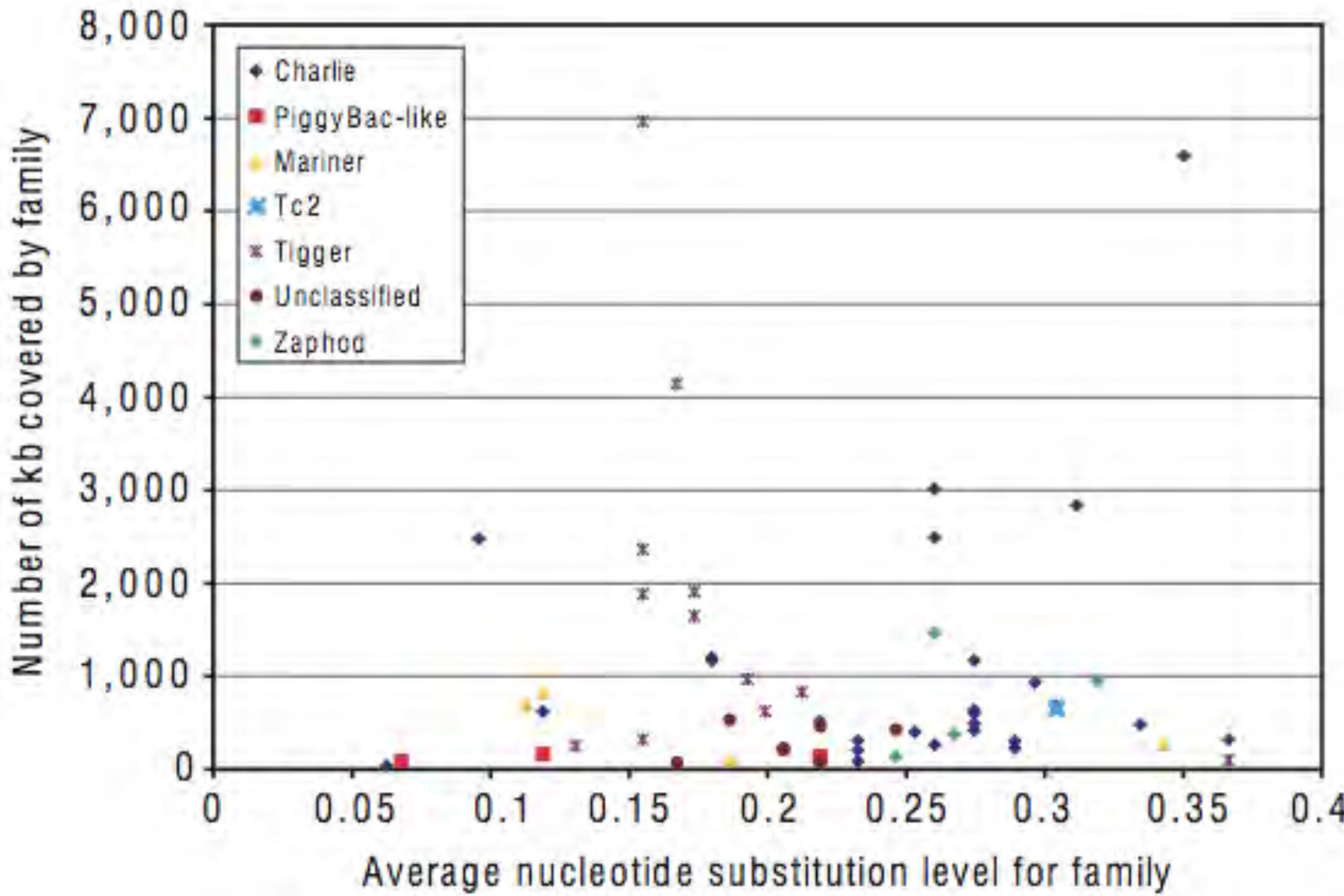
“Excellence in statistical graphics consists of complex ideas communicated with clarity, precision, and efficiency.”

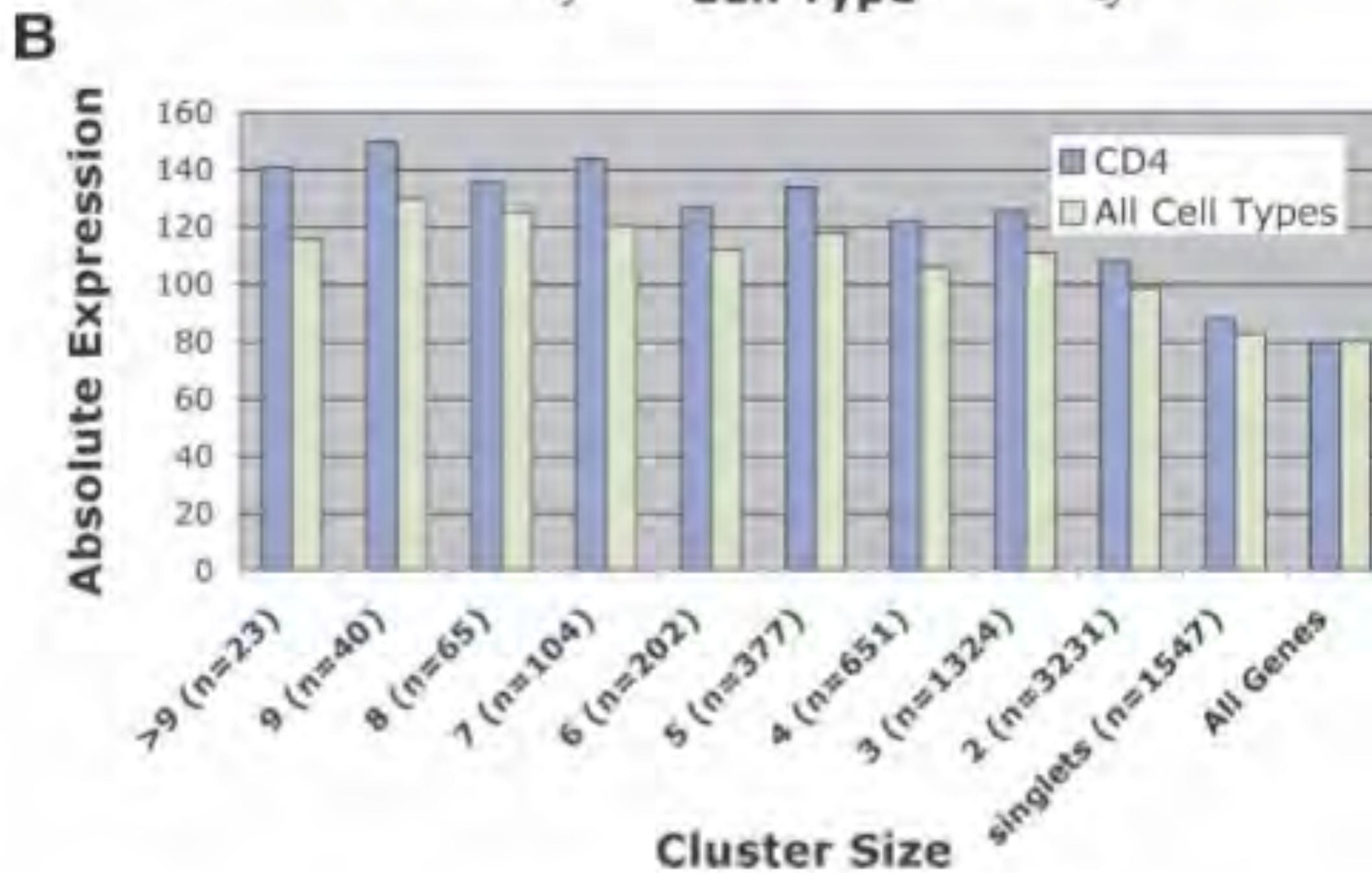
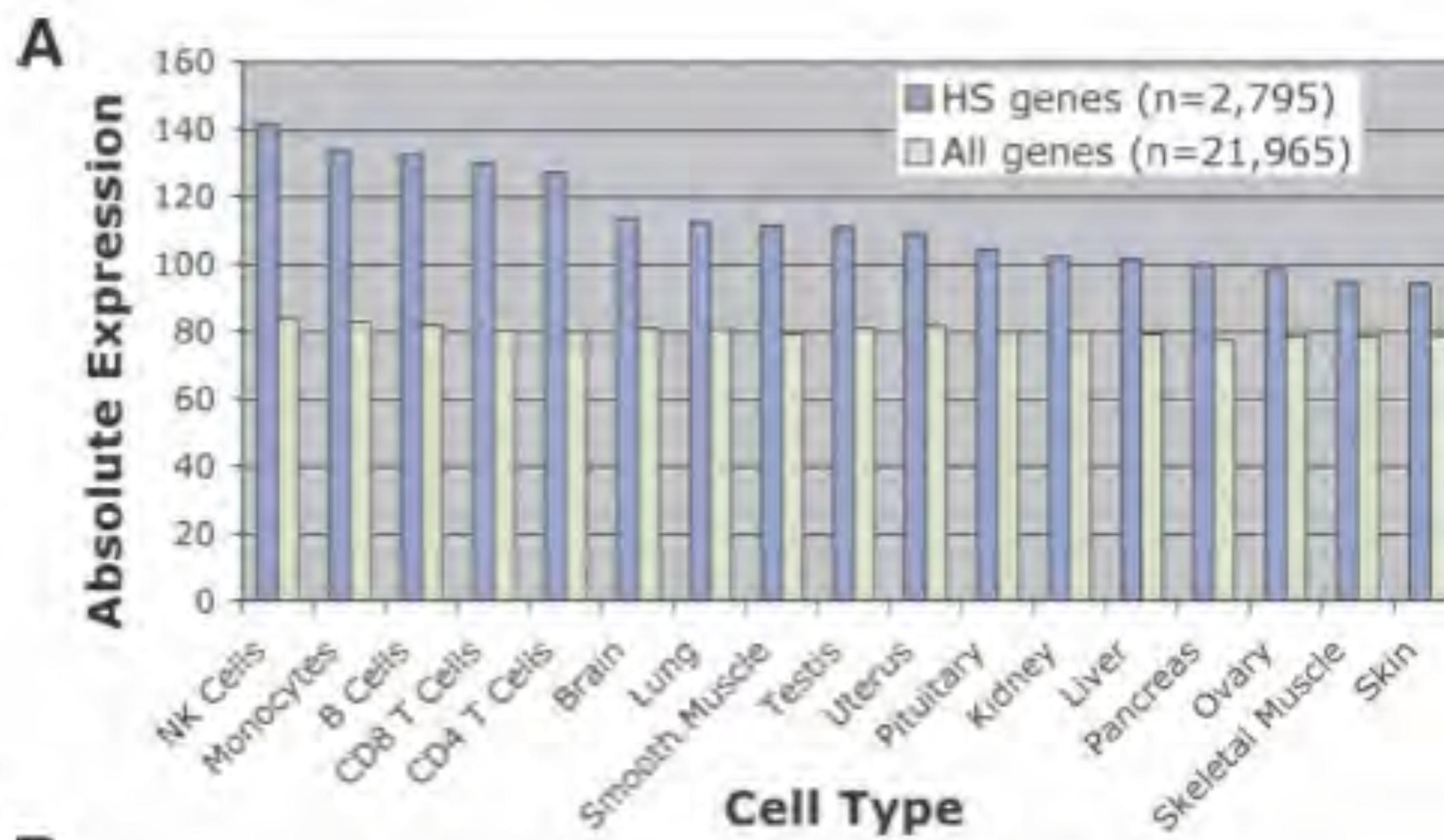
“a matter of *substance*, of *statistics*, and of *design*.”

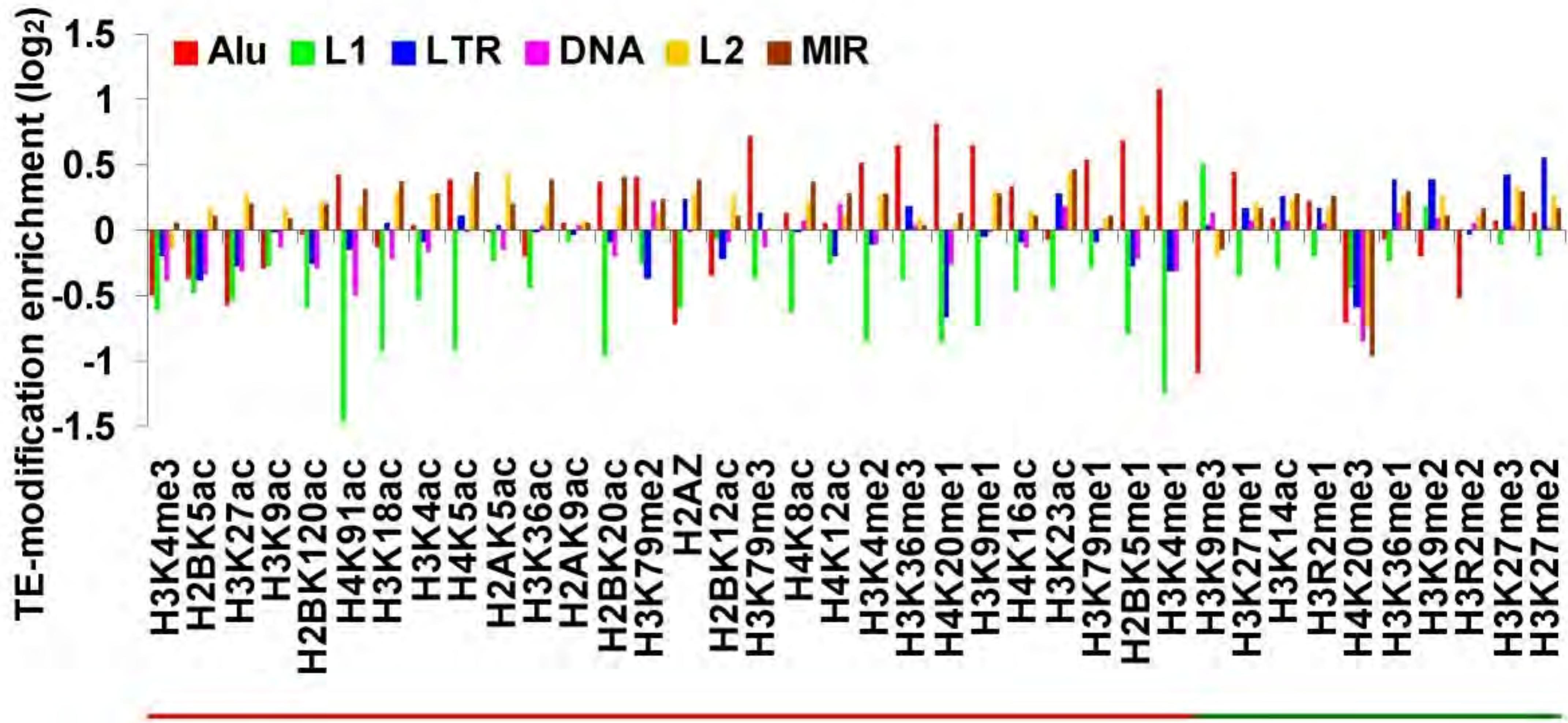
Edward Tufte, “The Visual Display of Quantitative Information”

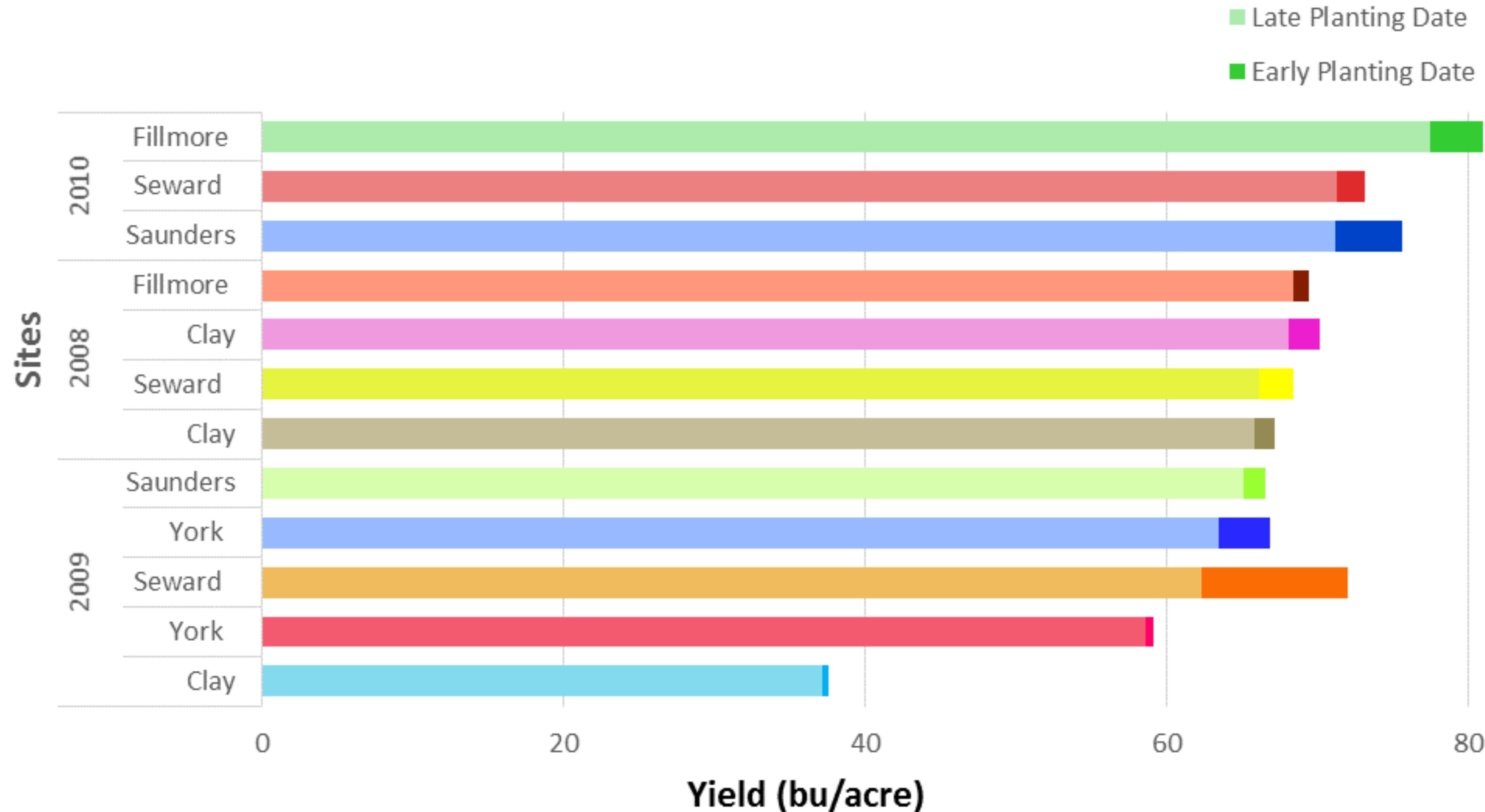


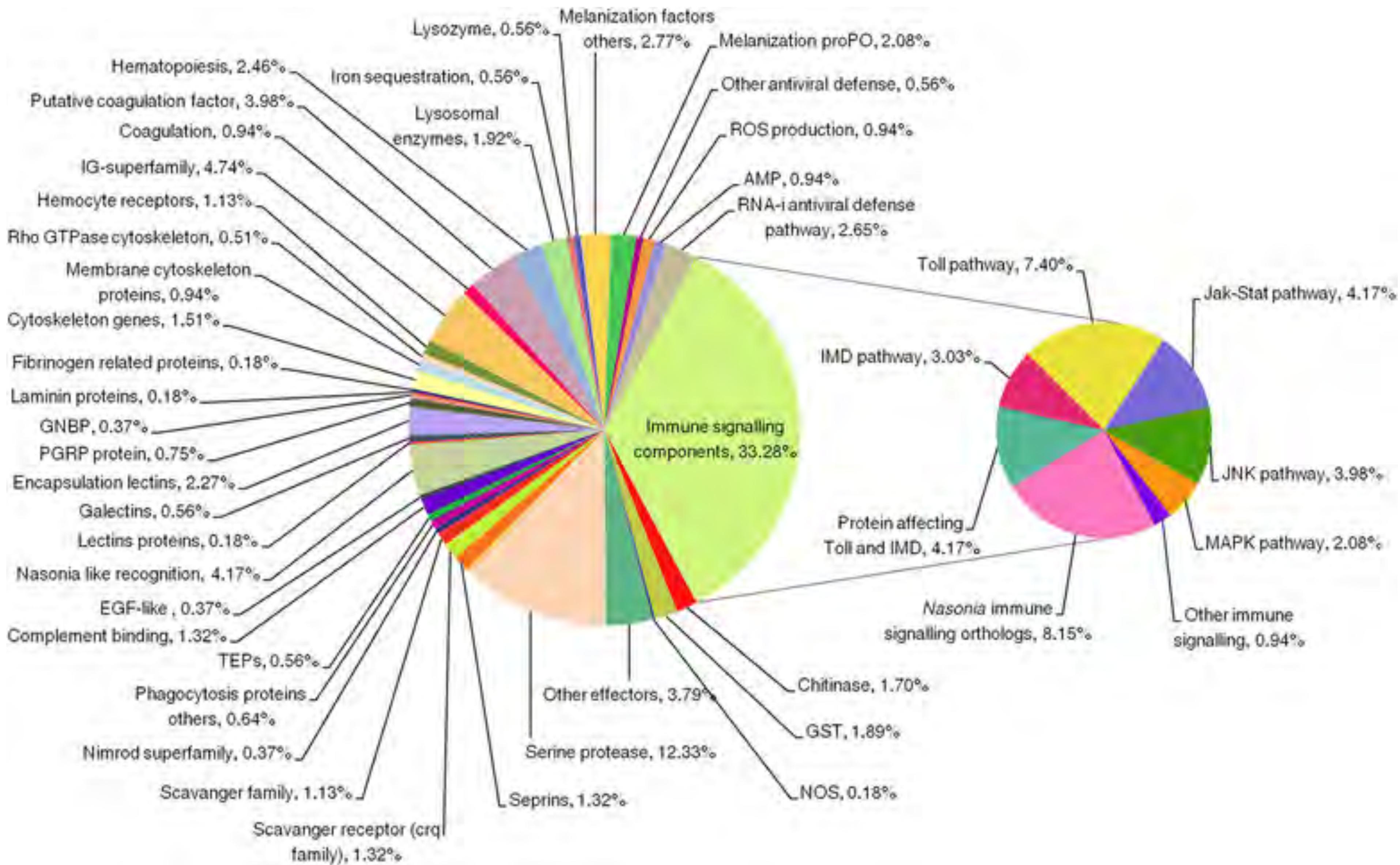












Goals of graphical excellence

A graphic should convey information:

Accurately: no distortions, distractions, or critical omissions in the data presented

Clearly: the focus is on the data

Efficiently: with minimal need for explanation

Revealingly: uncovers trends and organization in data; ideally leads to insight

Engagingly: invites the viewer to think about what results mean

And it should fit into the overall presentation:

Consistently: orientation, scale, color, font, encodings, expected location, etc.

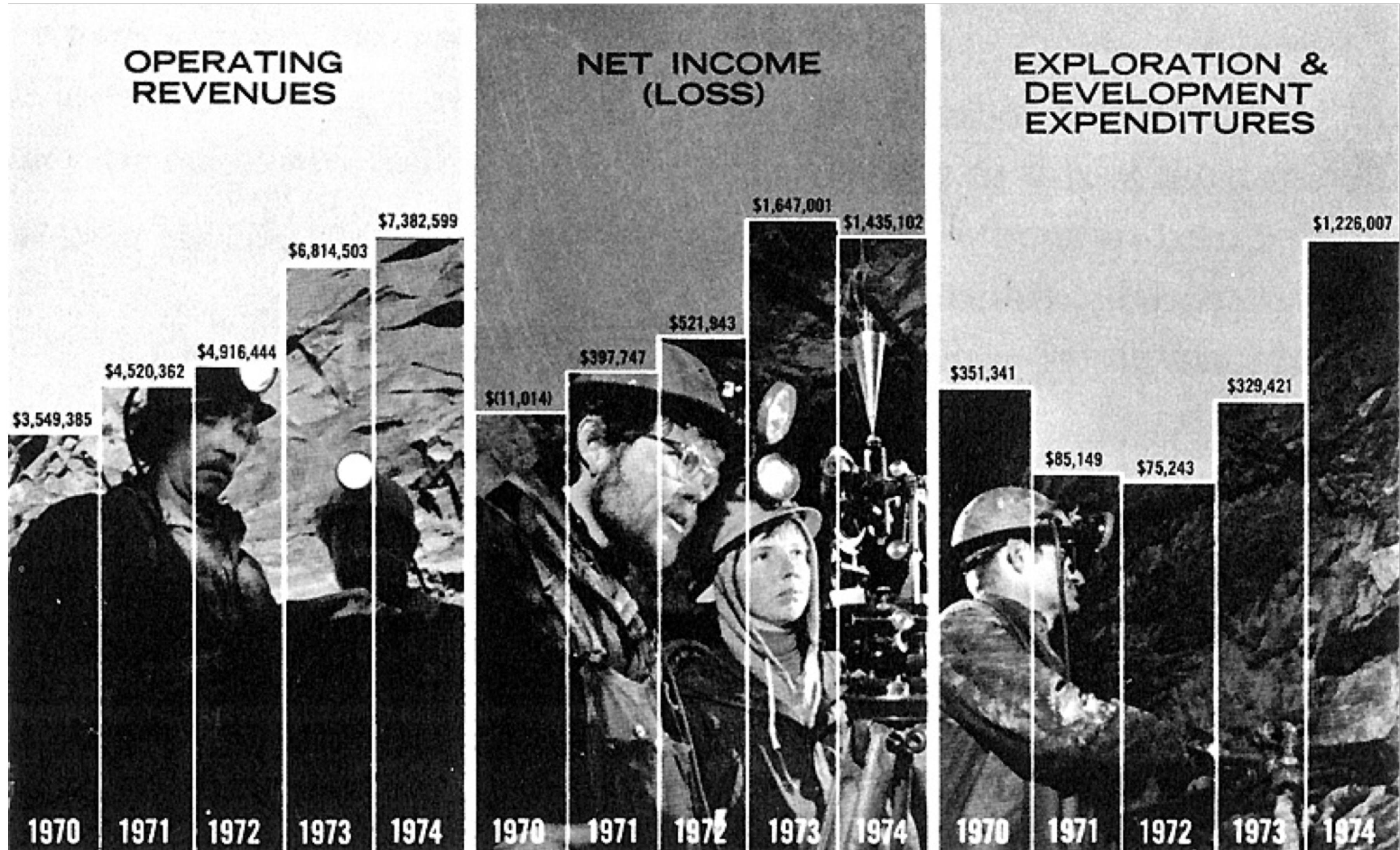
Relationally: connects to the overall narrative and to other visuals

Graphical integrity and counter-examples

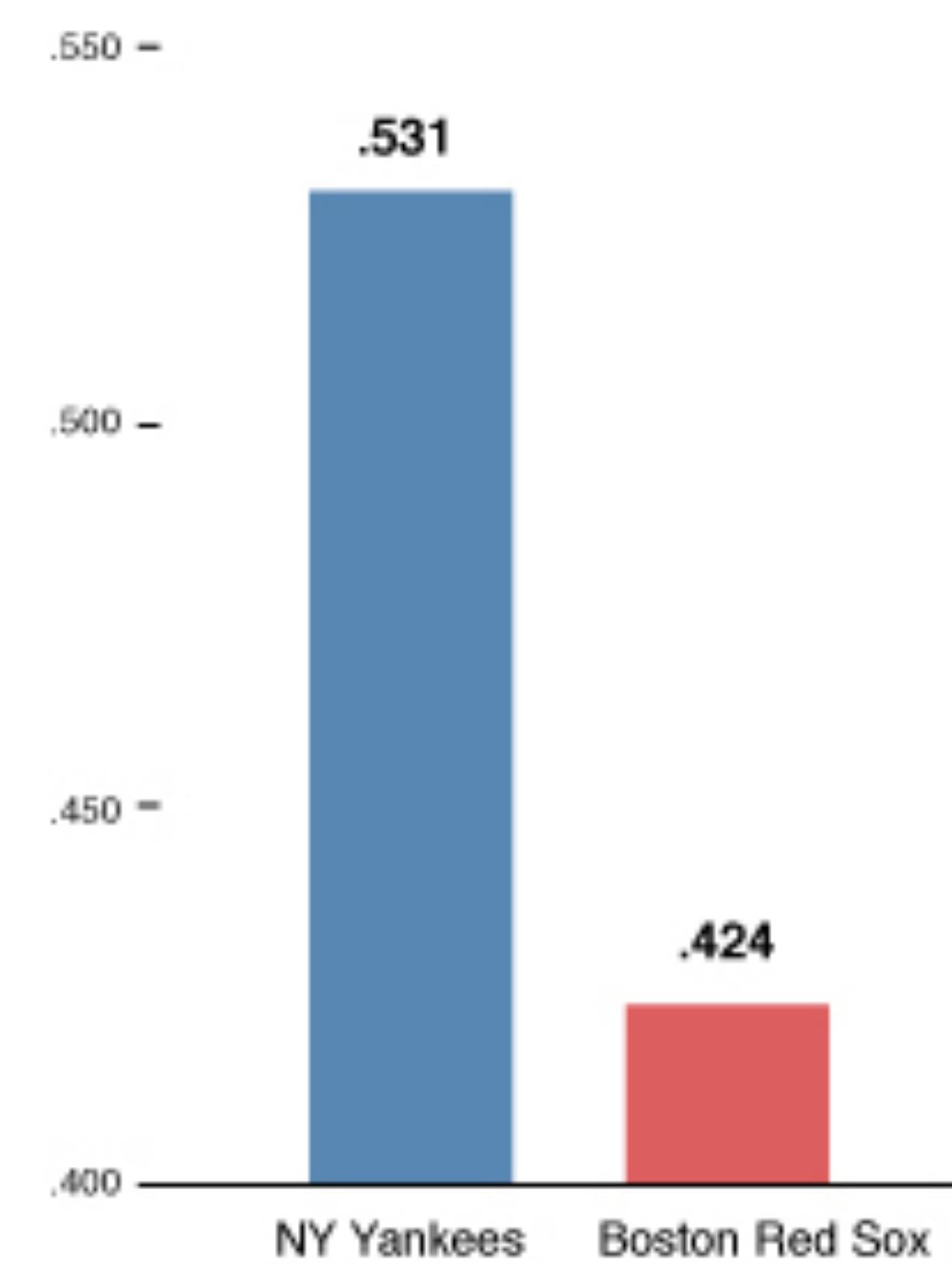
Tufte on graphical integrity

“Graphical integrity begins with telling the truth about the data”

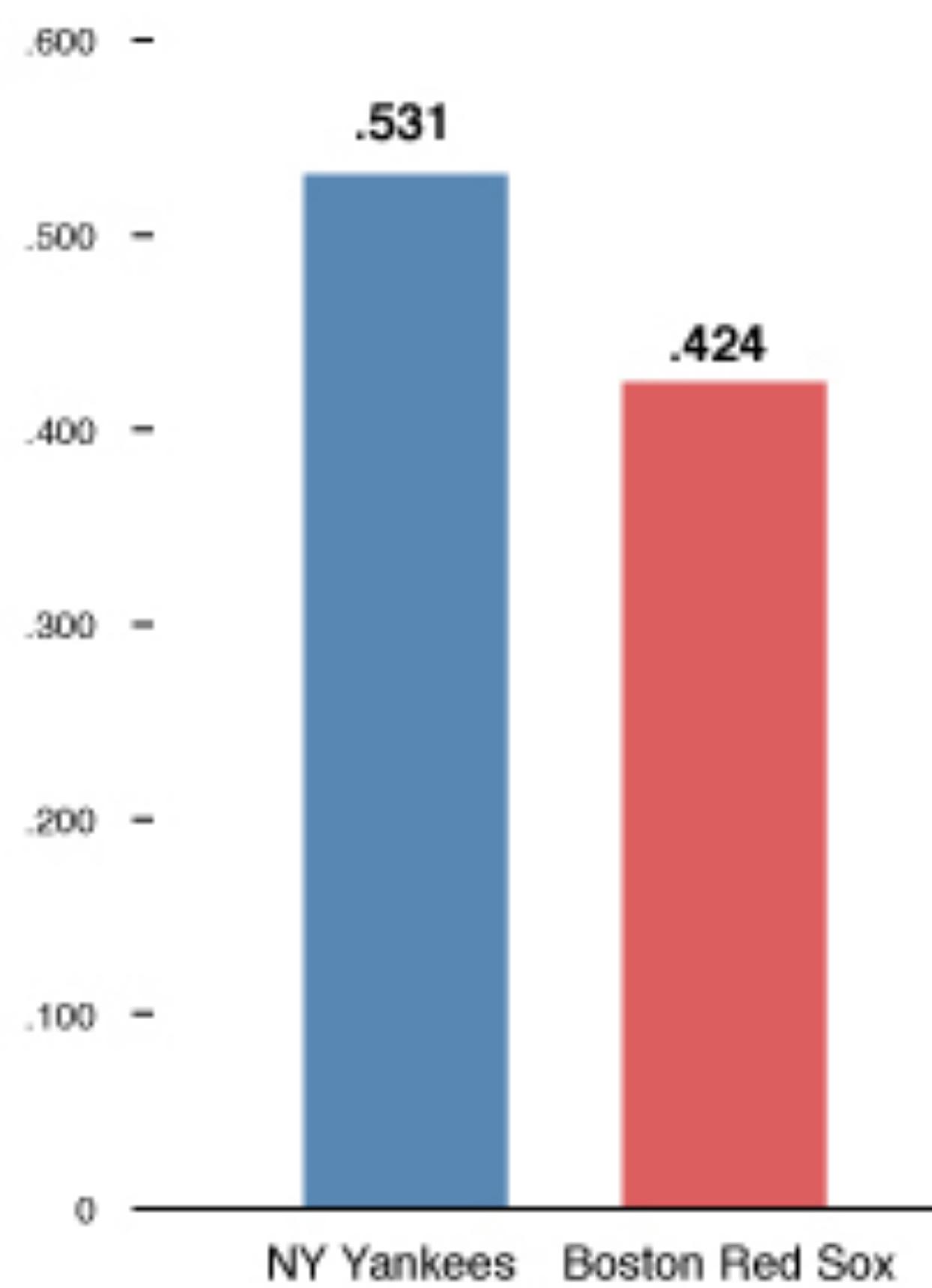
Edward Tufte, “The Visual Display of Quantitative Information”



Percentage of victories

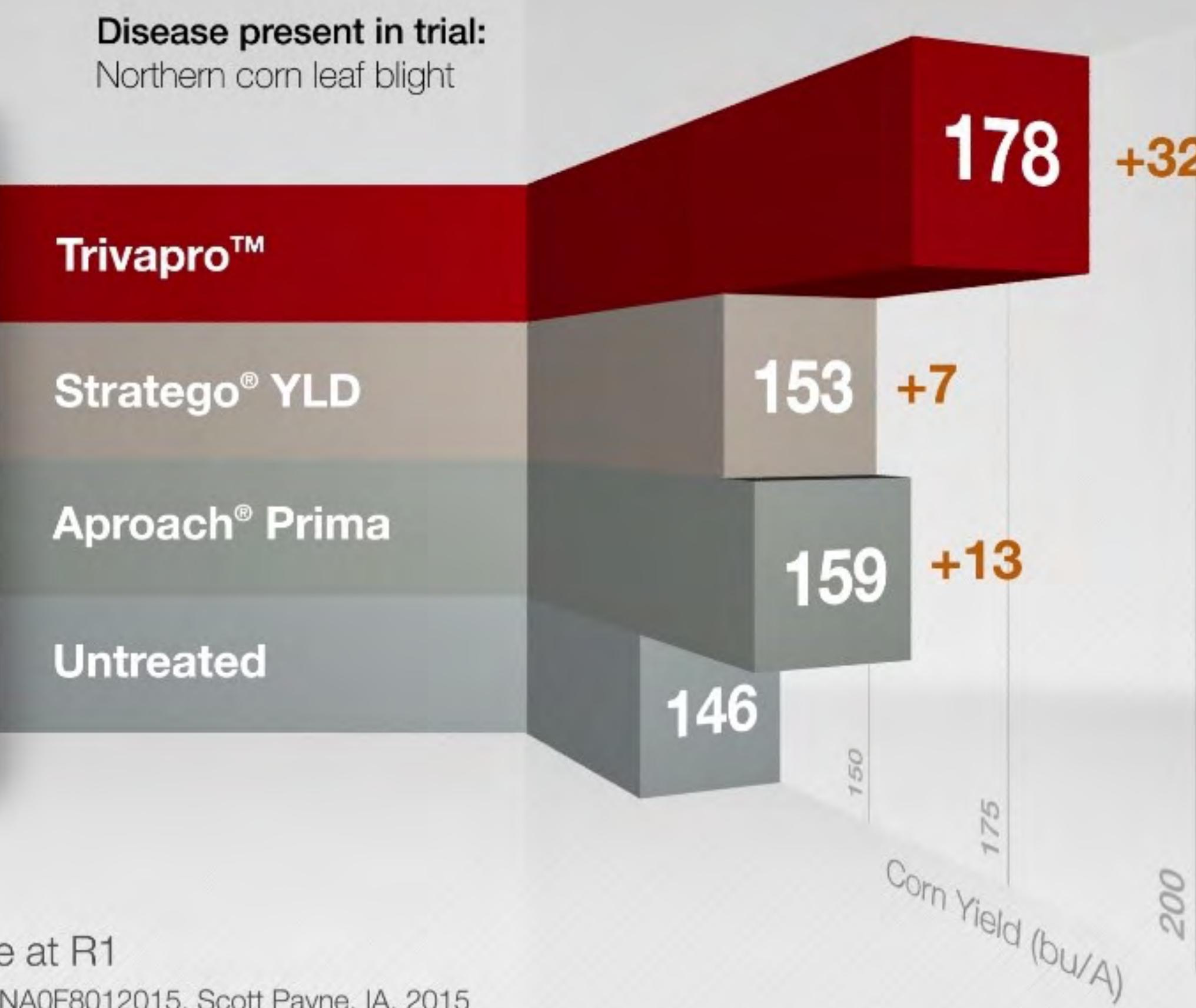


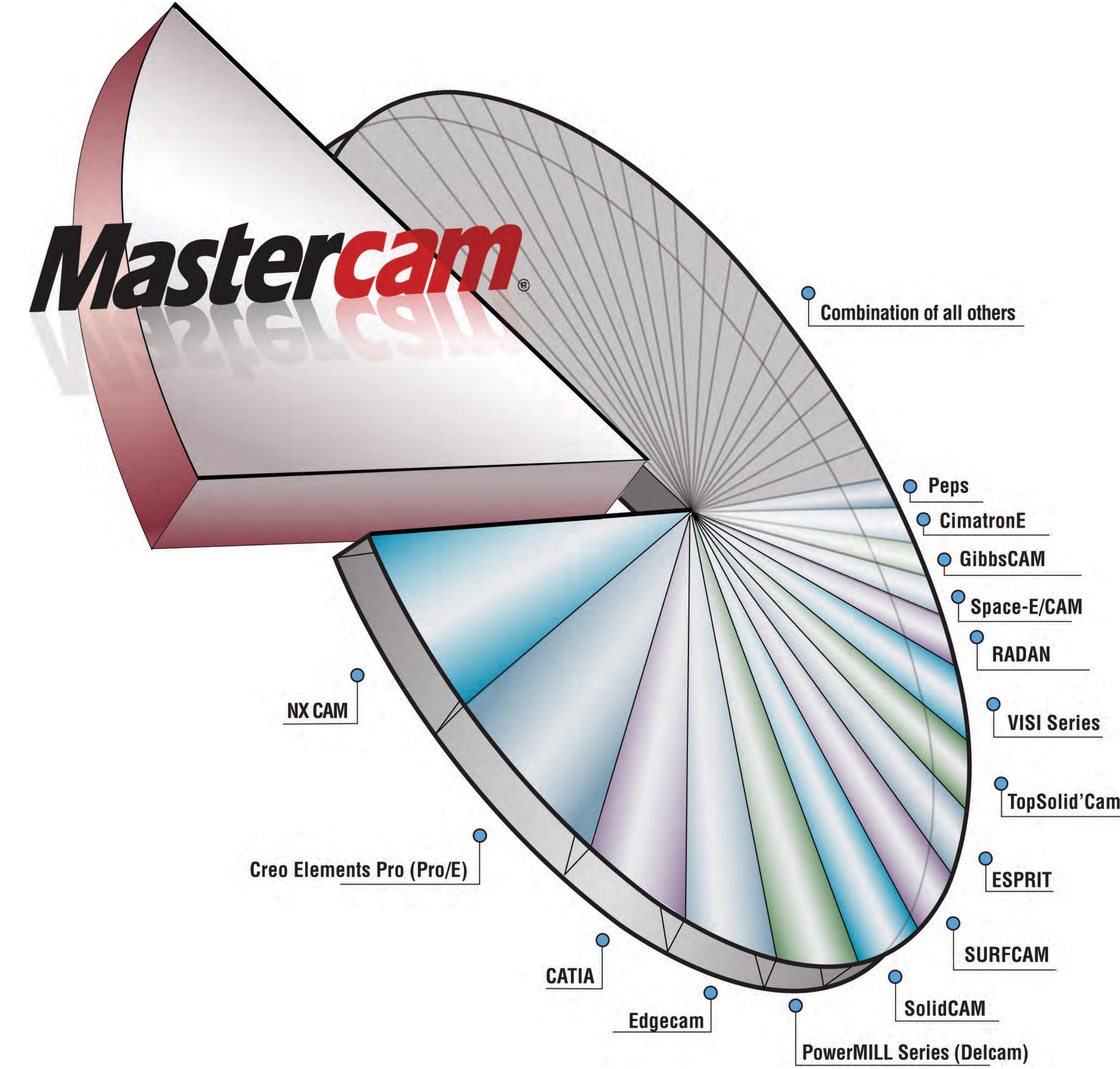
Percentage of victories



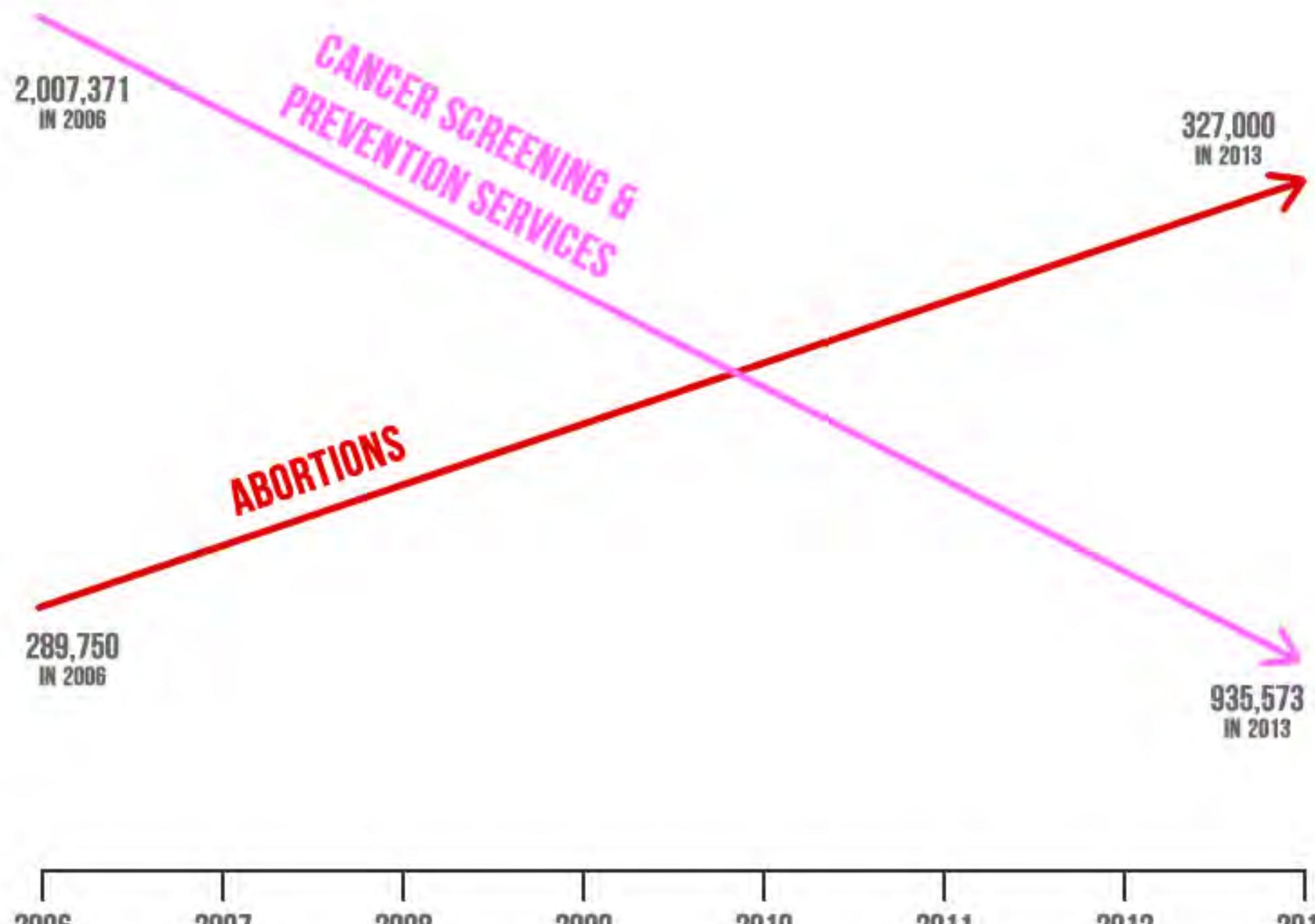
Trivapro corn yield response

in Boone, IA



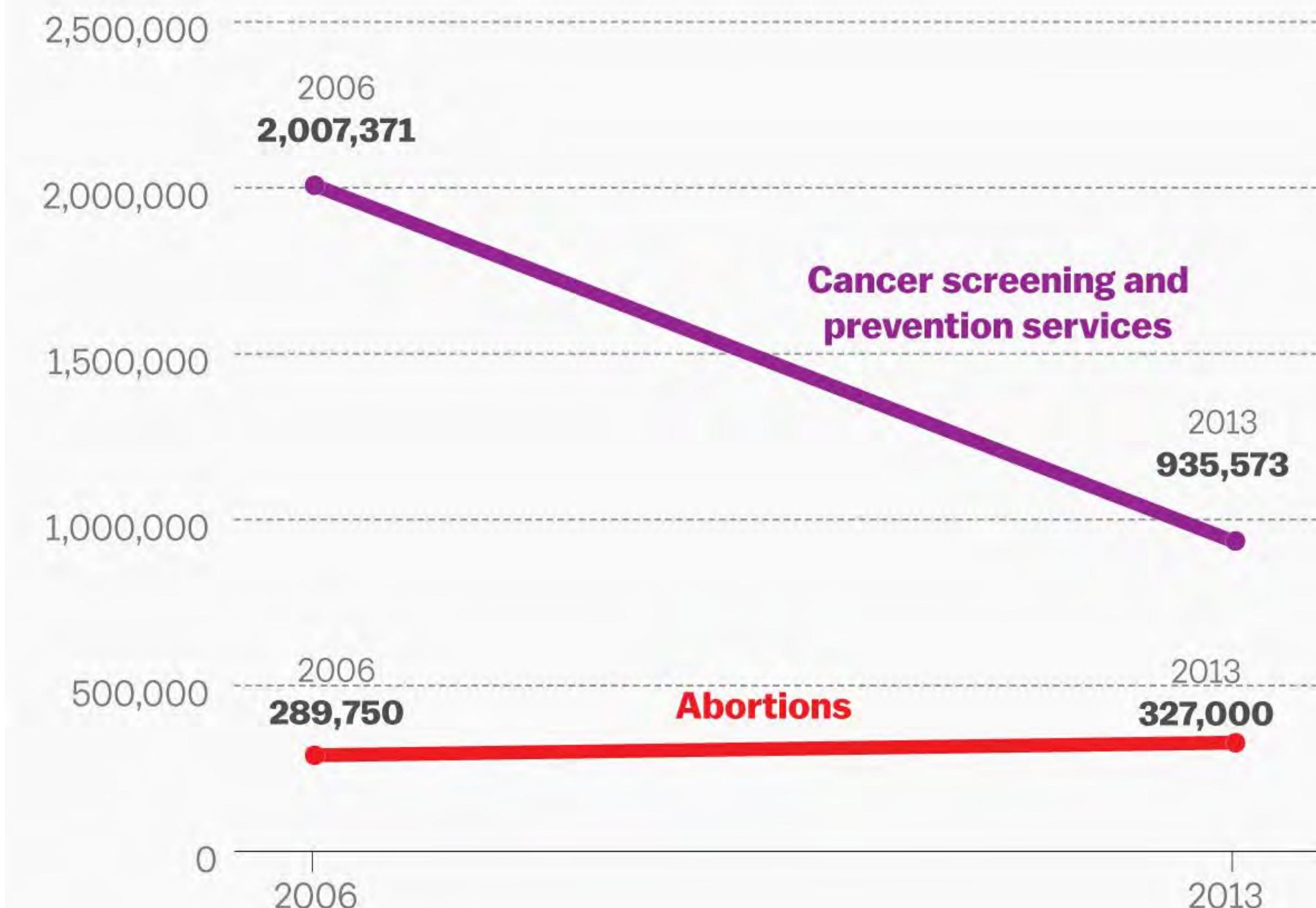


PLANNED PARENTHOOD FEDERATION OF AMERICA: ABORTIONS UP – LIFE-SAVING PROCEDURES DOWN



SOURCE: AMERICANS UNITED FOR LIFE

Services provided by Planned Parenthood

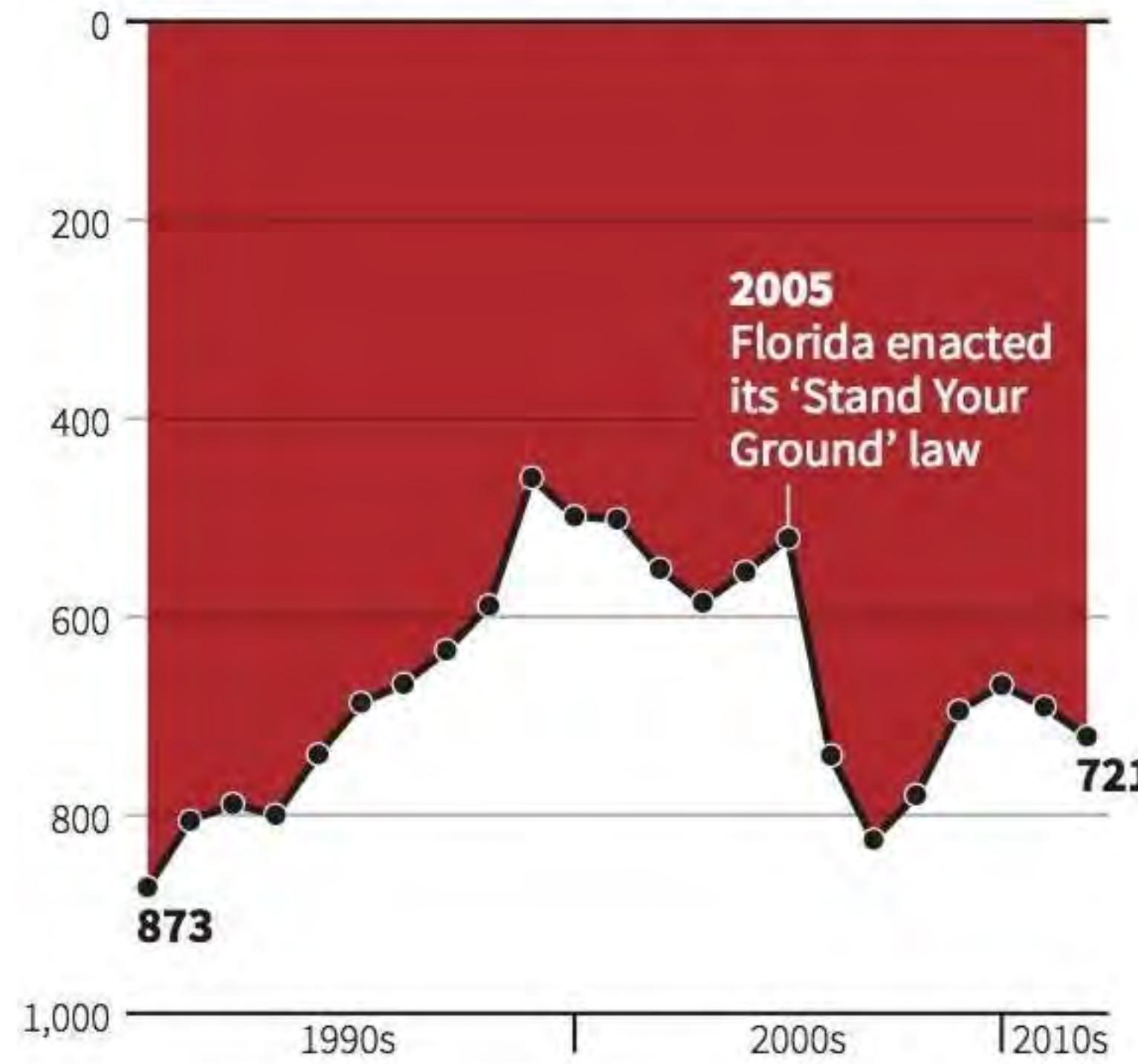


SOURCE: Planned Parenthood

Vox

Gun deaths in Florida

Number of murders committed using firearms



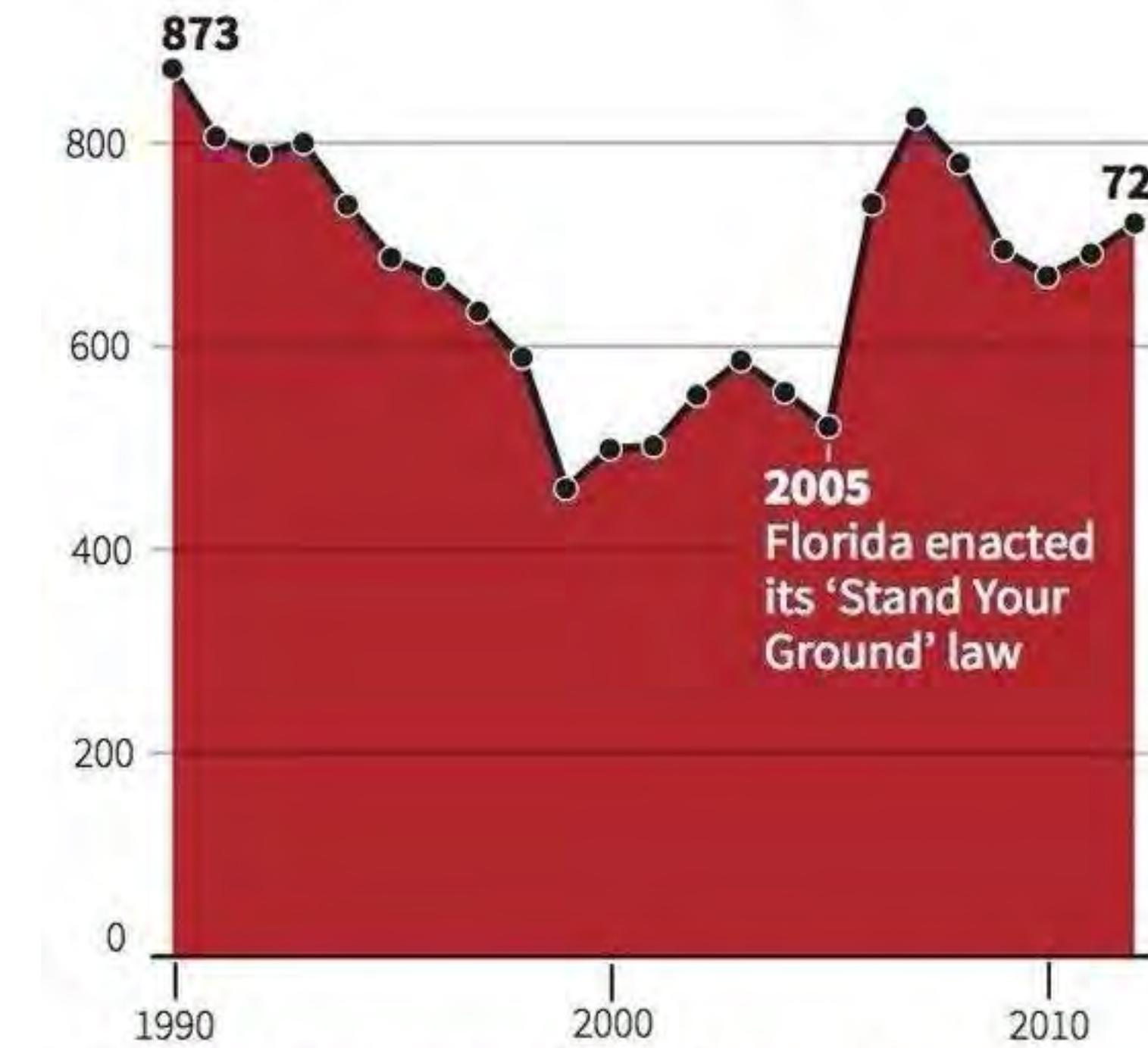
Source: Florida Department of Law Enforcement

C. Chan 16/02/2014

REUTERS

Gun deaths in Florida

Number of murders committed using firearms

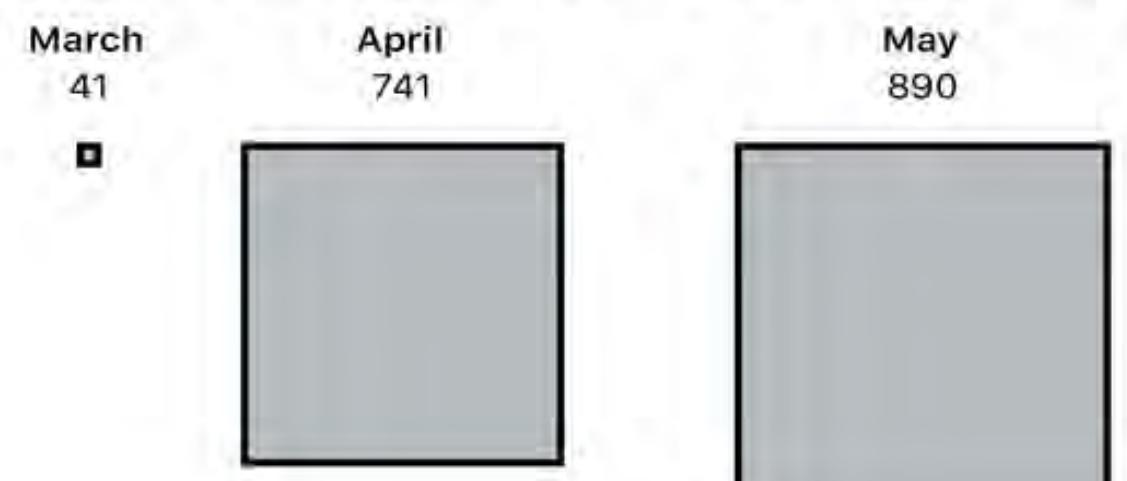


Source: Florida Department of Law Enforcement

C. Chan 16/02/2014

REUTERS

NEW DEATHS IN TEXAS BY MONTH



Note: July data as of July 21.

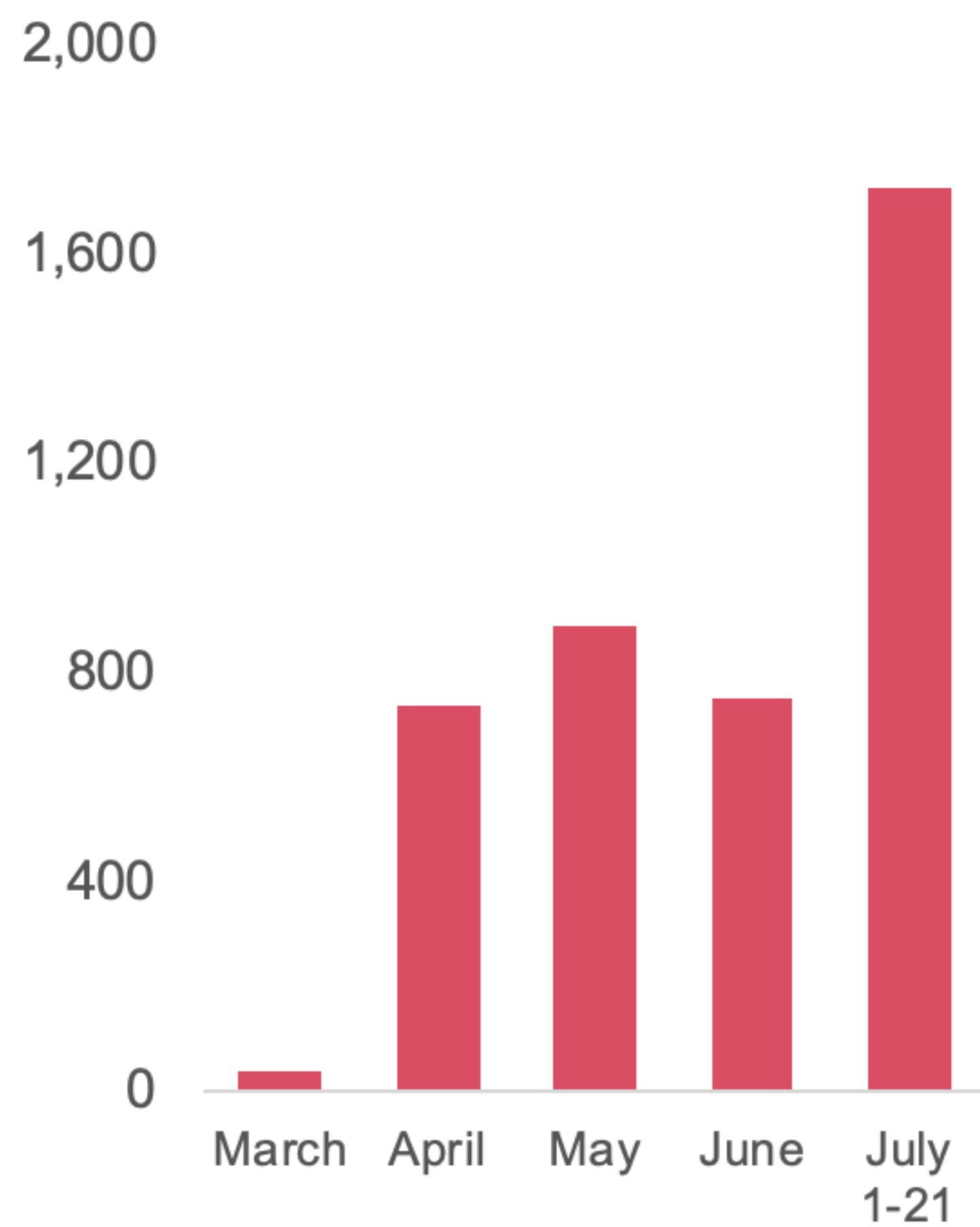
Source: The COVID Tracking Project

By Julia Redden | The COVID Tracking Project

Published July 21



Texas Deaths Due to Covid-19



Adapted from COVID Tracking Project
Kaiser Fung / JunkCharts



No. of Meetings

0

1

2

3

5

7

Goals of graphical excellence

Visual representations should be consistent with numerical values

- Avoid non-linear representations (e.g., perspective, truncated, or log plots)

- Use consistent scales in multicomponent plots

- Avoid broken scales (often, there are better ways)

Use clear, explicit, and appropriate labeling

Use the appropriate number of variable dimensions

- Never exceed the number of data dimensions

Do not quote data out of context

- Make comparisons clear: controls, baselines, neutral models, etc.

- For time series, provide enough points to show trends

Do not selectively present data to alter conclusions

Component-based design

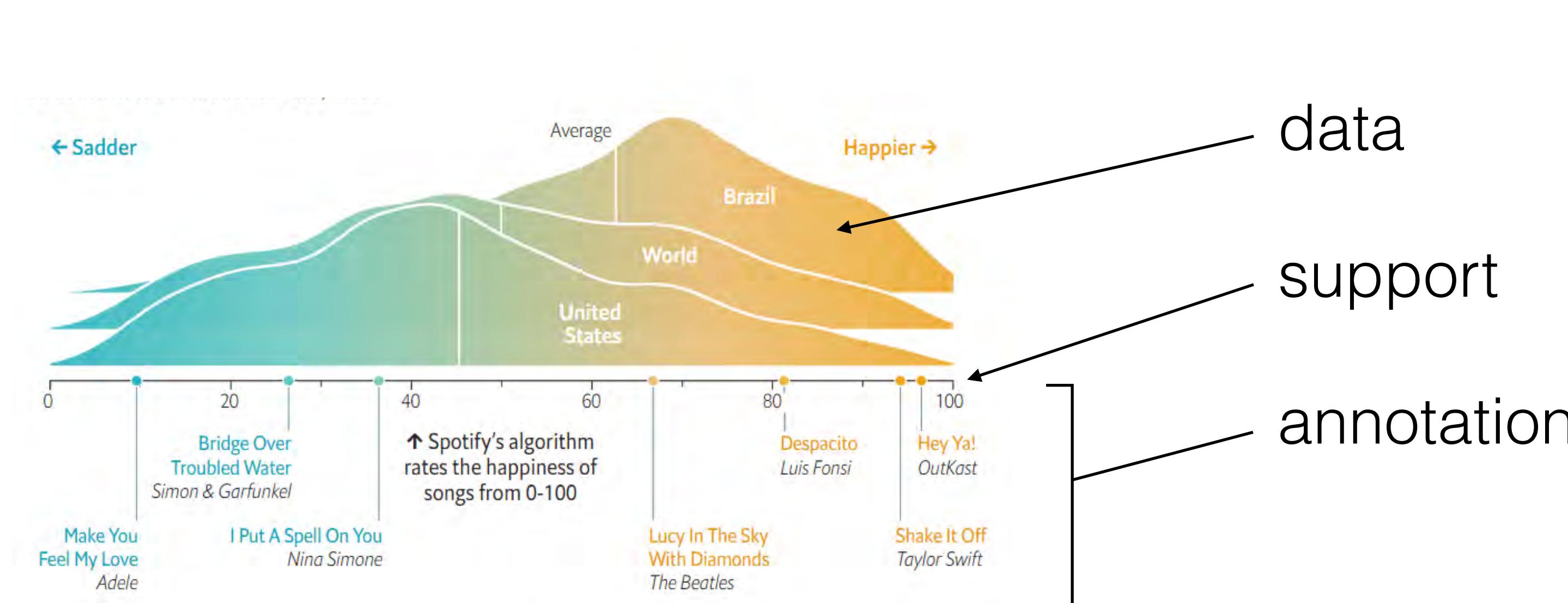
Component-based design

Quantitative graphics contain three kinds of components:

Data: geometric representations of numerical values (points, lines, bars, areas, etc.)

Data support: axes, axis labels, scales, grids, tick marks, encoding keys, etc.

Annotation: groupings, trend lines, arrows, highlights, text boxes, callouts, etc.



Data components

Primary objective:

Data values and distributions should be easy to perceive

Overall trends and differences/similarities among groups should be clear

Points to consider

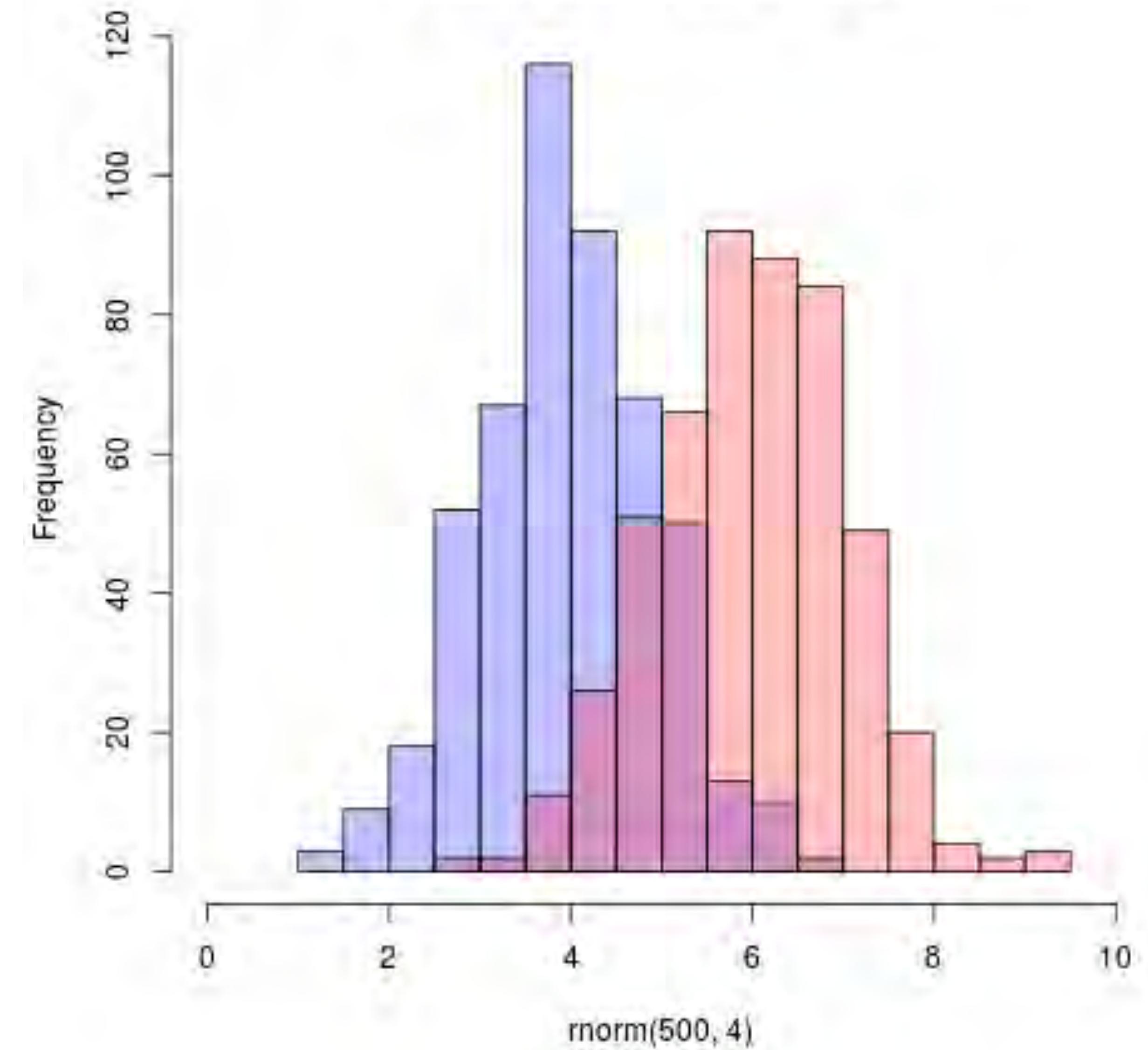
Data representation: graph type, data point/line/area, size, color/shape/transparency

How much data: all or summary or both?

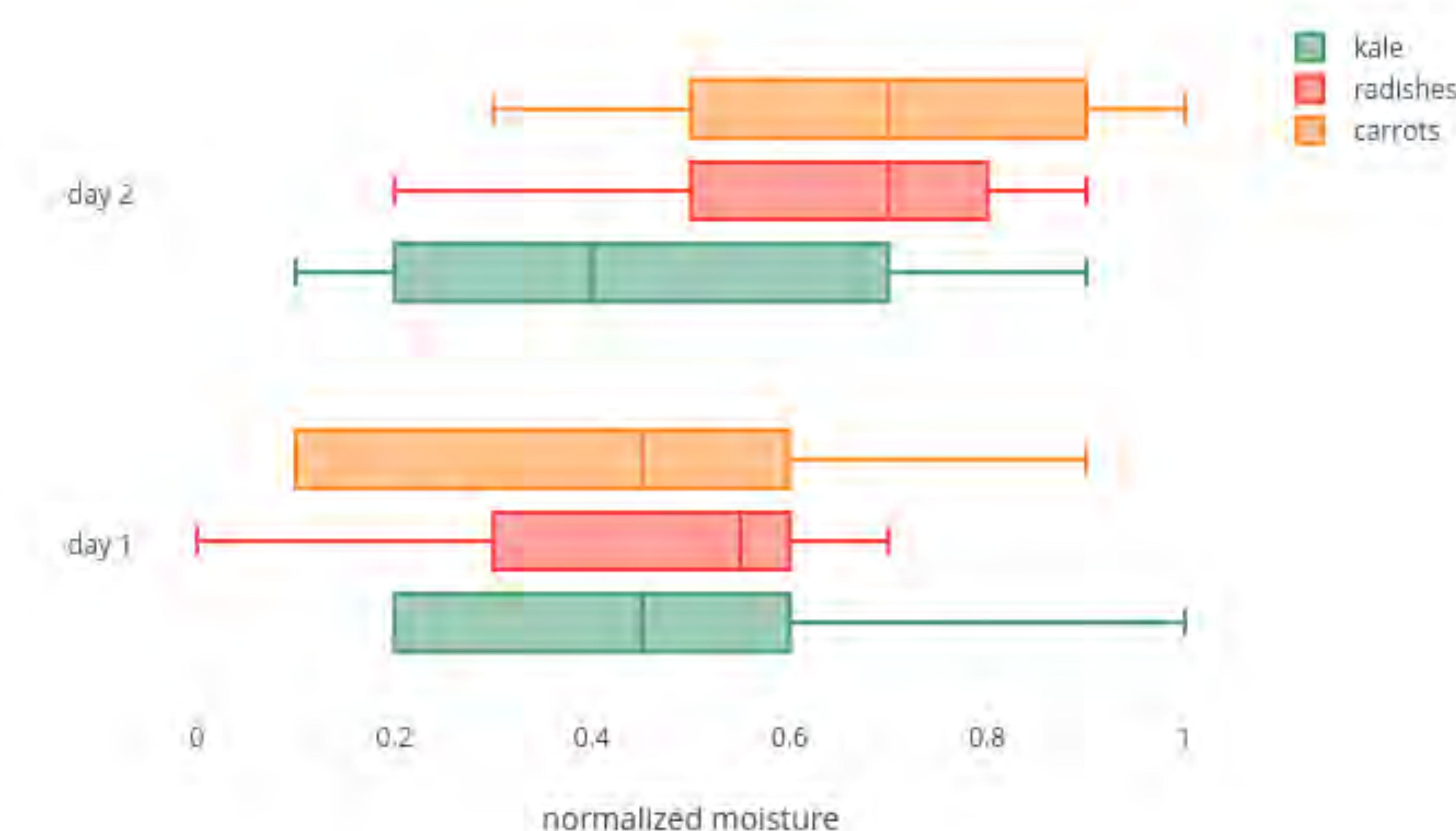
Encodings: aim to be intuitive and consistent

Scales: linear/semi-log/log, zero-based or not, equal among all panels

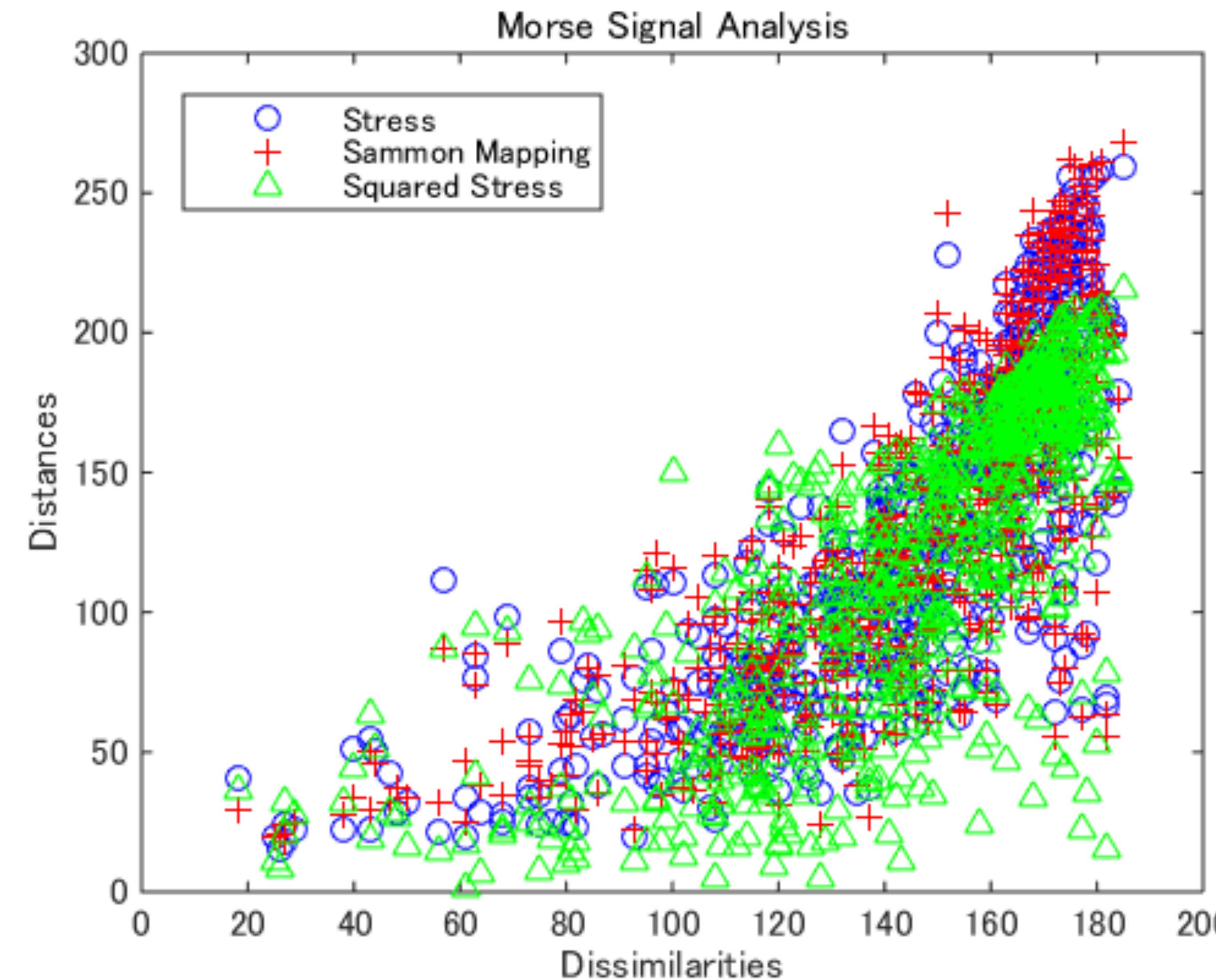
Example of intuitive encoding



Example of intuitive encoding



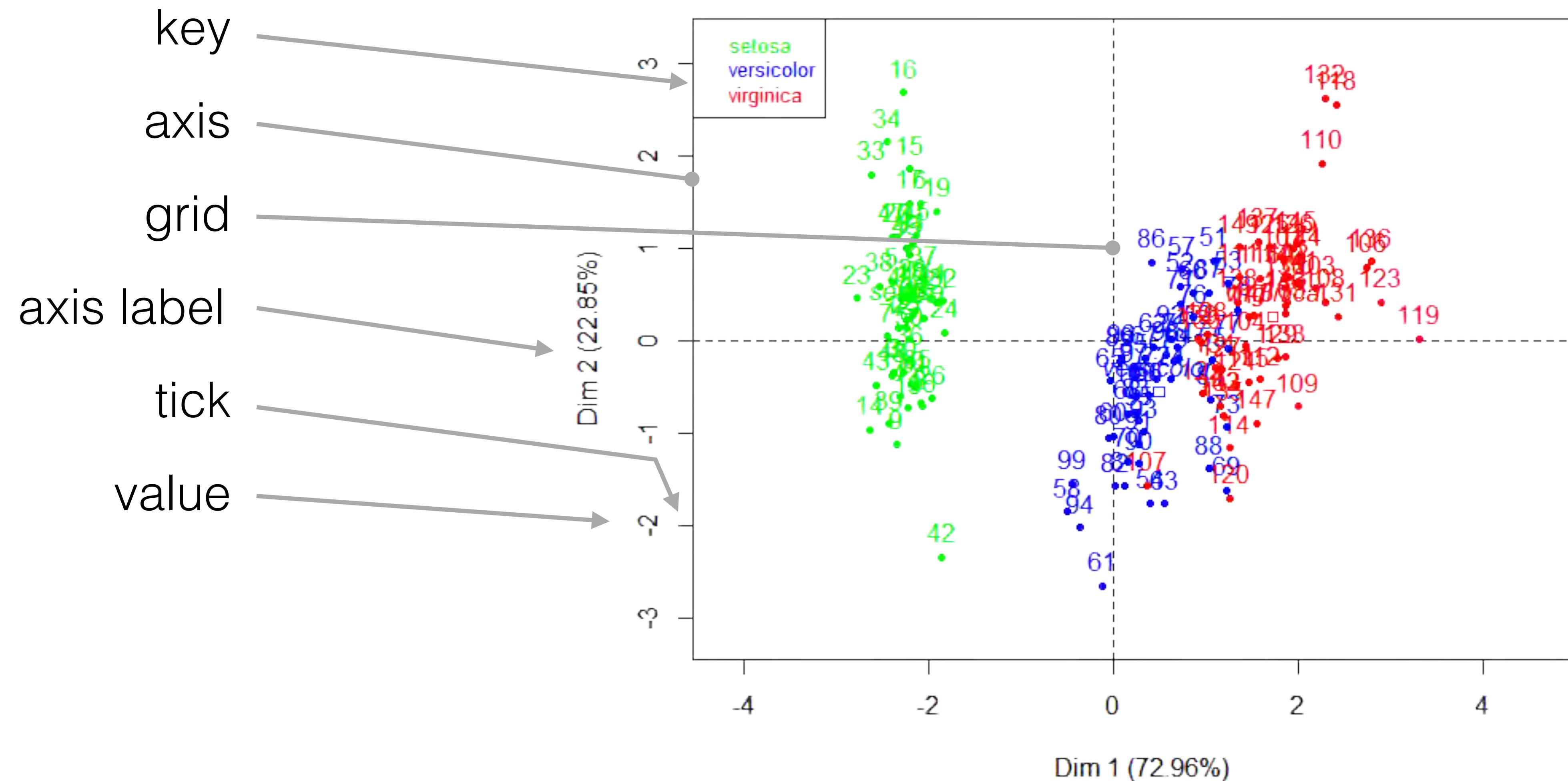
Example of overplotting



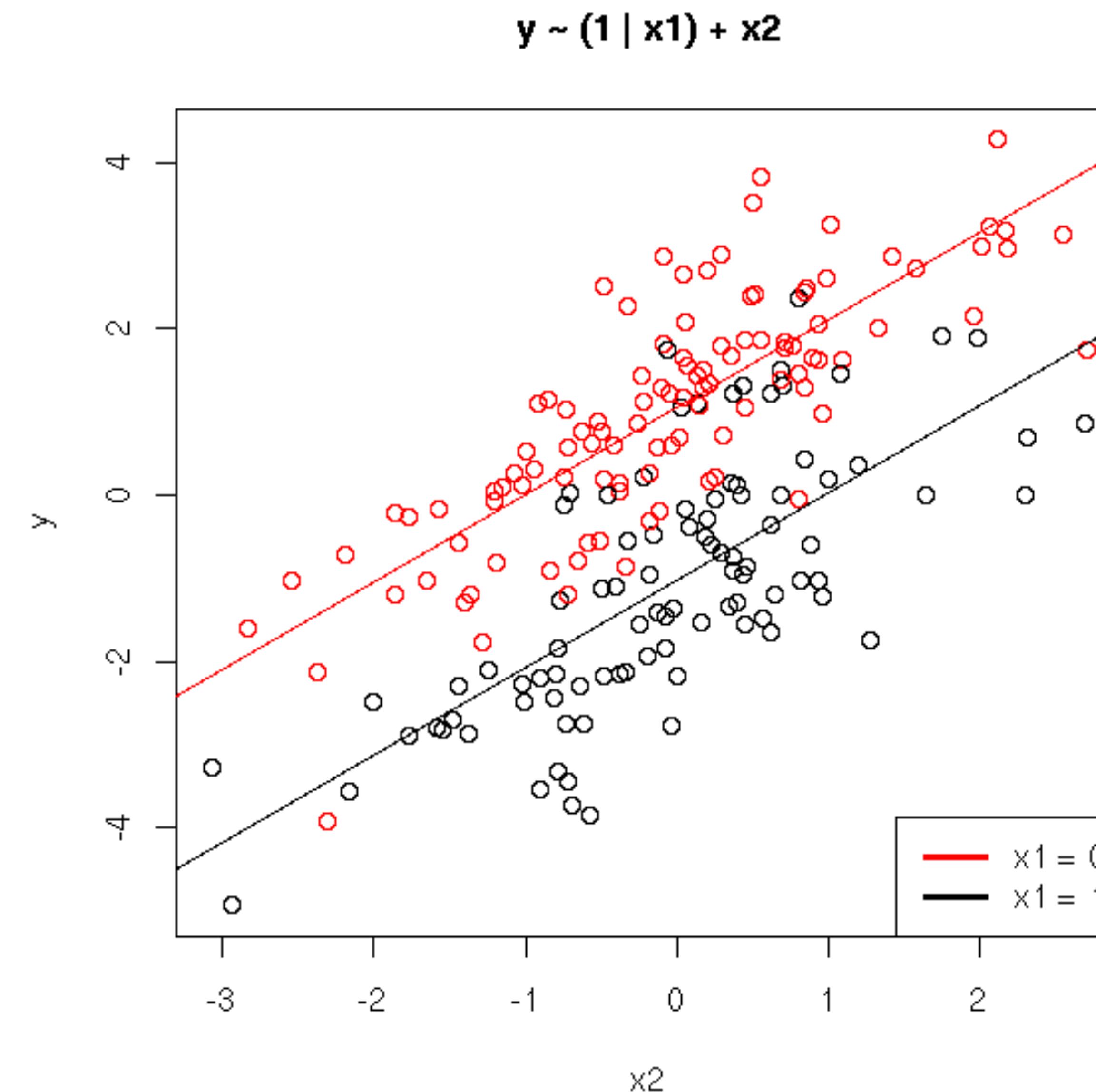
Support components

Primary objective: provide viewers the ability to understand the values of the data

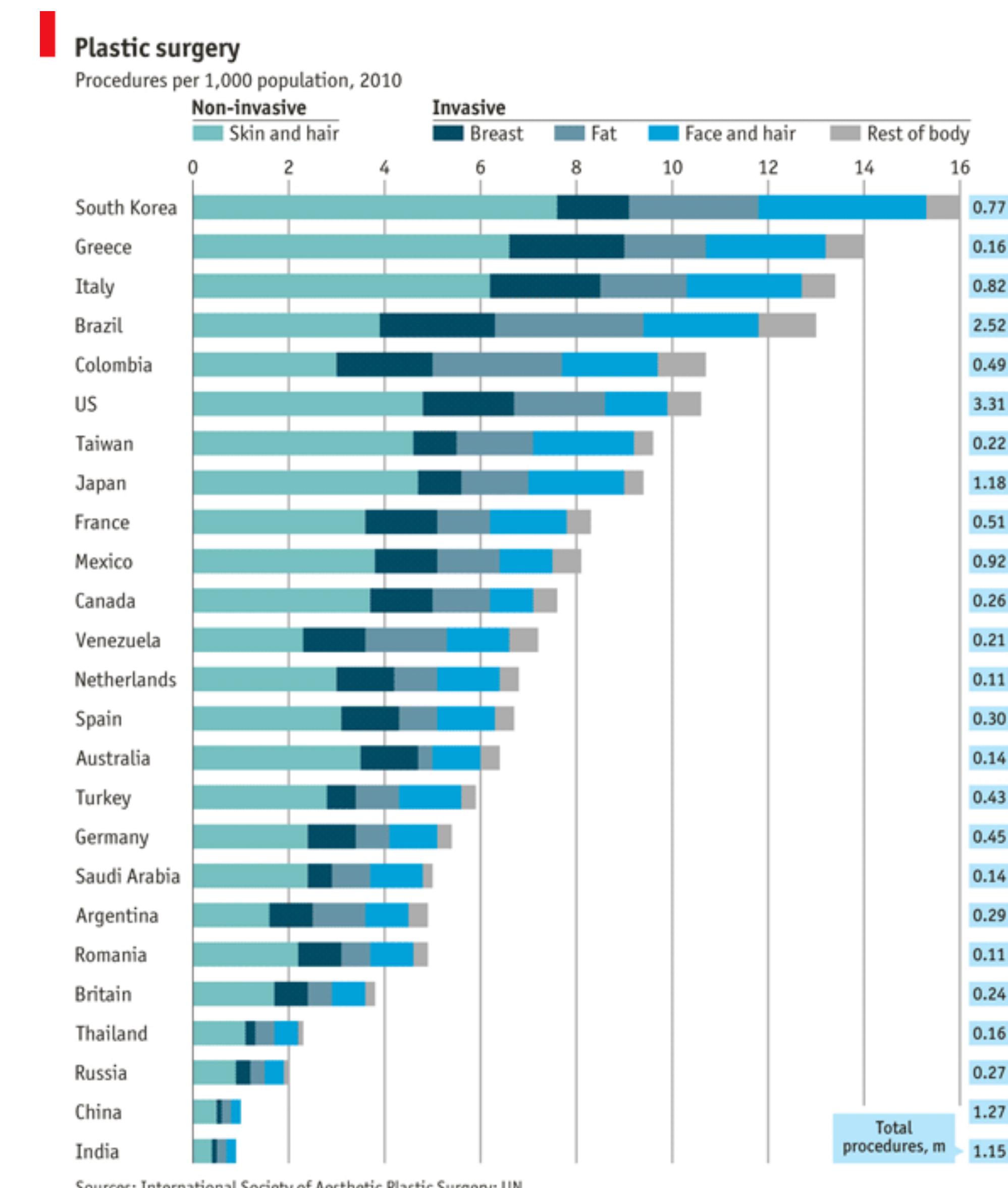
Components to consider:



Example support components: trend lines



Example support components: key and alternate normalization



Annotation components

Primary objective: to tell a story

Components to consider:

Title

Values or labels for individual points

Trend lines

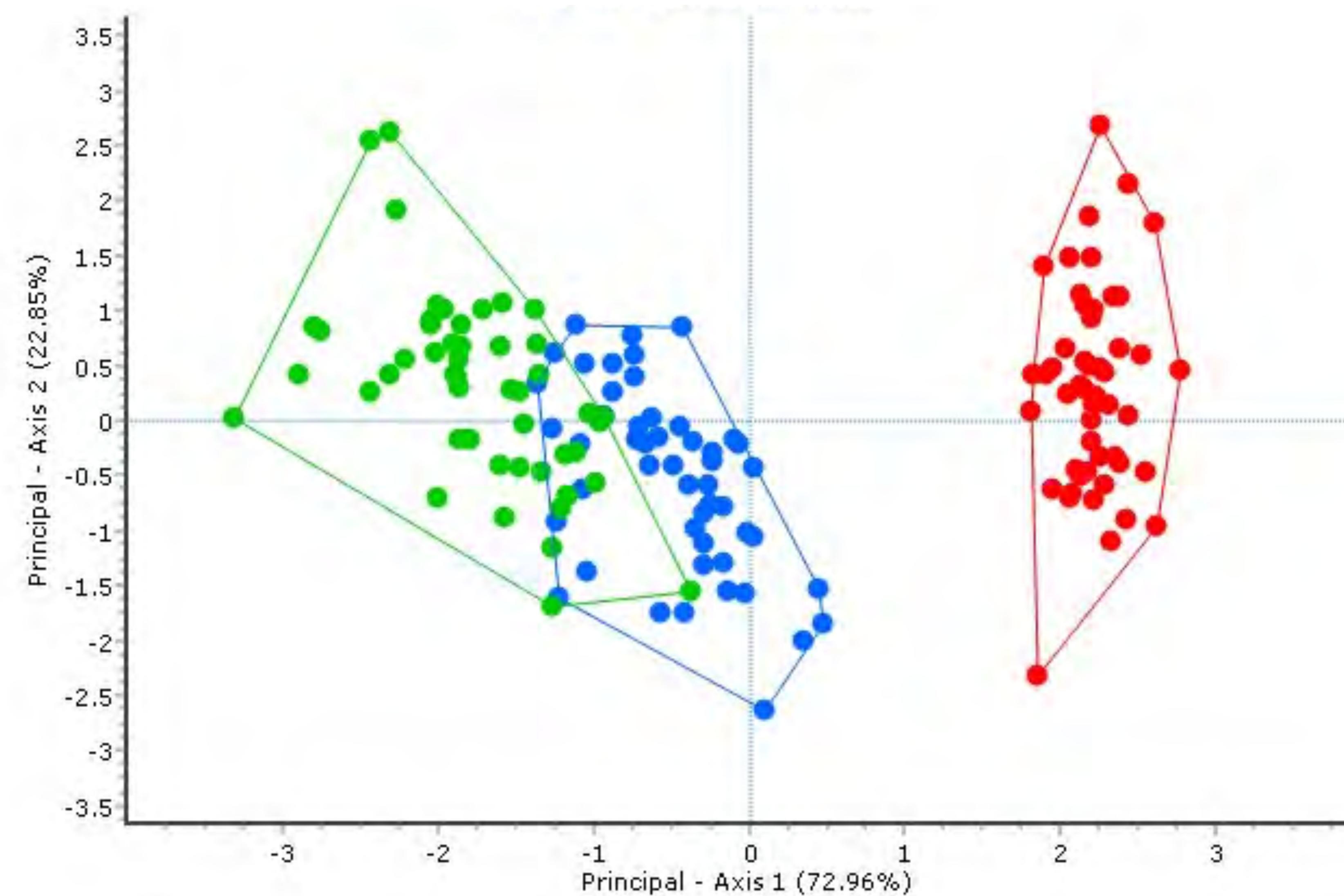
Groupings

Highlighting, outlining, or transparency of data subset

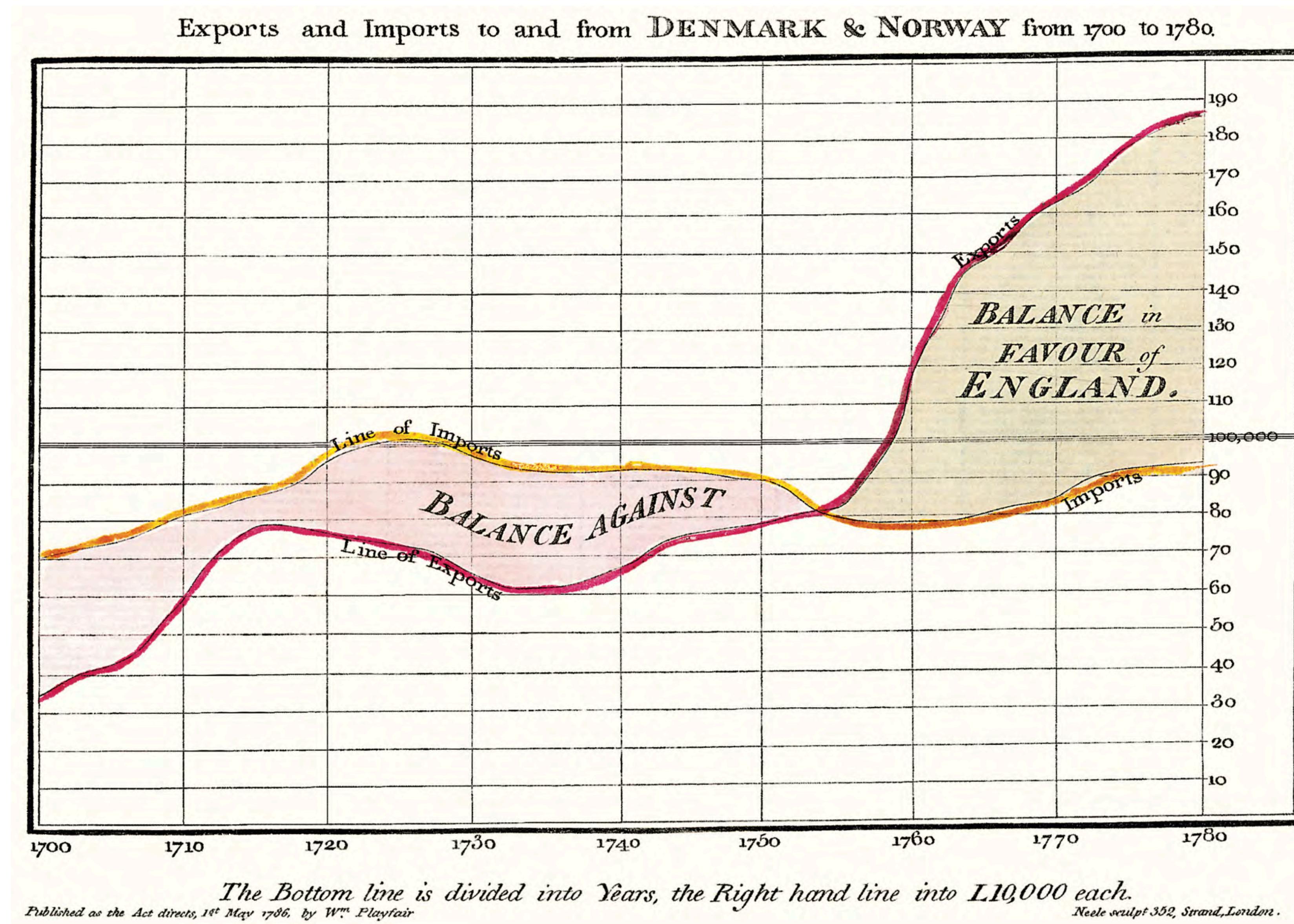
Statistical significance (asterisks, groupings, Tukey groups)

Text labels referring to data components

Example annotation components: grouping



Example annotation components: text



Resources

Recommended resources

Books:

- Tufte, E. 2001. *The Visual Display of Quantitative Information*. Graphics Press. (2nd ed.)
- Wilke, C. 2019. *Fundamentals of Data Visualization*. O'Reilly Media.

Websites:

From Data to Viz: <https://www.data-to-viz.com/>

FriendsDontLetFriends: <https://github.com/cxli233/FriendsDontLetFriends>

R Graph Gallery: <http://rgraphgallery.blogspot.com/>

MatPlotLib Examples: https://matplotlib.org/stable/plot_types/index.html

Seaborn Gallery: <https://seaborn.pydata.org/examples/index.html>

Syntax for adding constraints

```
CREATE TABLE t1 (
    uid INTEGER PRIMARY KEY,
    genus VARCHAR NOT NULL UNIQUE,
    species VARCHAR NOT NULL,
    b_date DATE UNIQUE,
    age INTEGER DEFAULT = 42,
    entry_date DATE DEFAULT CURRENT_DATE,
    p1 INTEGER CHECK (c6 > 10 AND c6 != 0),
    p2 INTEGER CHECK (6 IN ('red', 'green', 'blue')),
);
```

Note: some versions of SQL use a different syntax for **CHECK** constraints and **CURRENT_DATE**

SQL is limited to a small set of operations

Only four kinds of operations are permitted with relational databases

Create, Read, Update, Delete (CRUD)

`CREATE` to create a new database, table, or relation; `INSERT` to add rows to a table

`SELECT` to query a database ✓

`UPDATE` to change the information within a table

`DELETE` to remove rows from a table; `DROP` to remove a database, table, or relation

SQL has some additional keywords, but most work with the above set

Clauses within statements: `WHERE`, `JOIN`, `GROUP BY`, `LIMIT`, etc.

Operators within statements: `=`, `>`, `IN`, `NOT`, `LIKE`, `BETWEEN`, etc.

Functions within statements: `MIN`, `MEAN`, `COUNT`, `DISTINCT`, etc.

`Pandas` and `dplyr` implement many SQL-like operations (and were inspired by it)

CREATE syntax

table to create

```
CREATE TABLE table (
    column_1 DATA_TYPE PRIMARY KEY,
    column_2 DATA_TYPE,
    column_3 DATA_TYPE
);
```

optional but highly recommended

Data types typically include: **INTEGER**, **FLOAT**, **BOOLEAN**, **CHAR(n)**, **VARCHAR**, **DATE**, etc.

Optional qualifiers: **PRIMARY KEY**, **NOT NULL**, **UNIQUE**, etc.