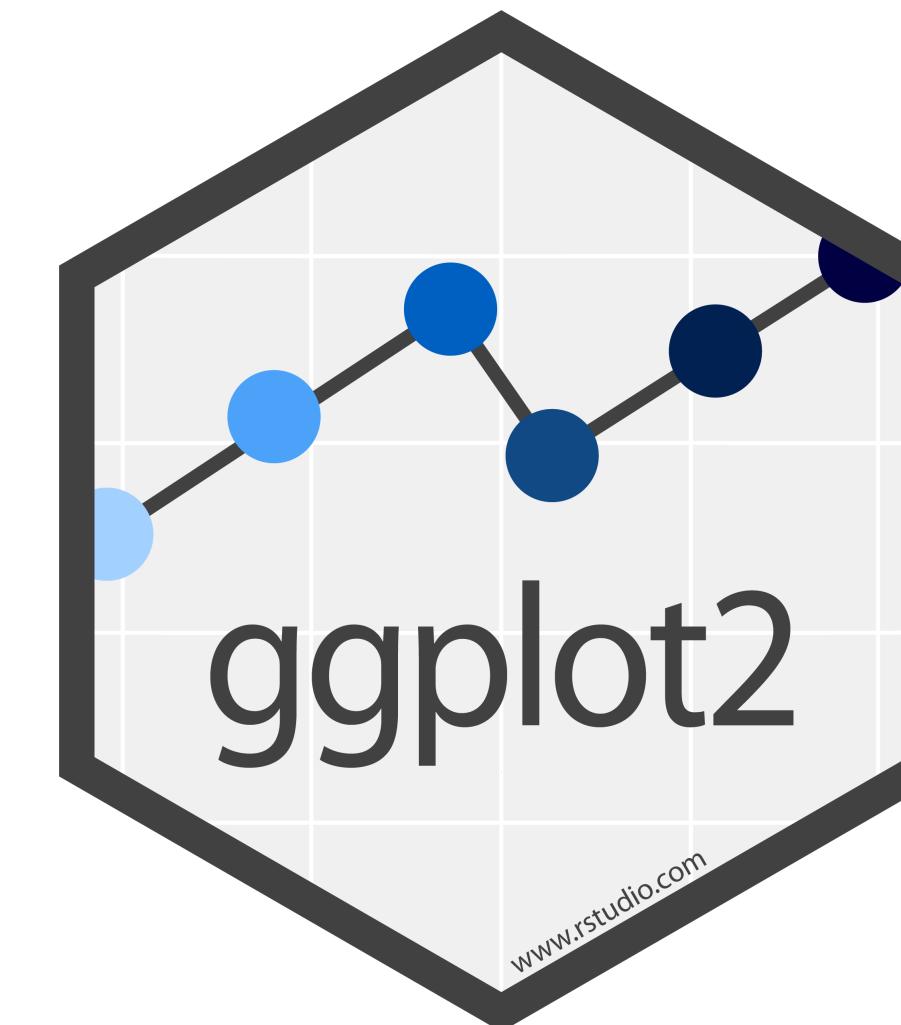


Programme for today

- Introductory lecture
- Hands-on dataviz tutorial
- Individual exercises



At the end of the day you should be familiar with making visualisations in R using the ggplot2 package



The greatest value of a picture
is when it forces us to notice
what we never expected to see.

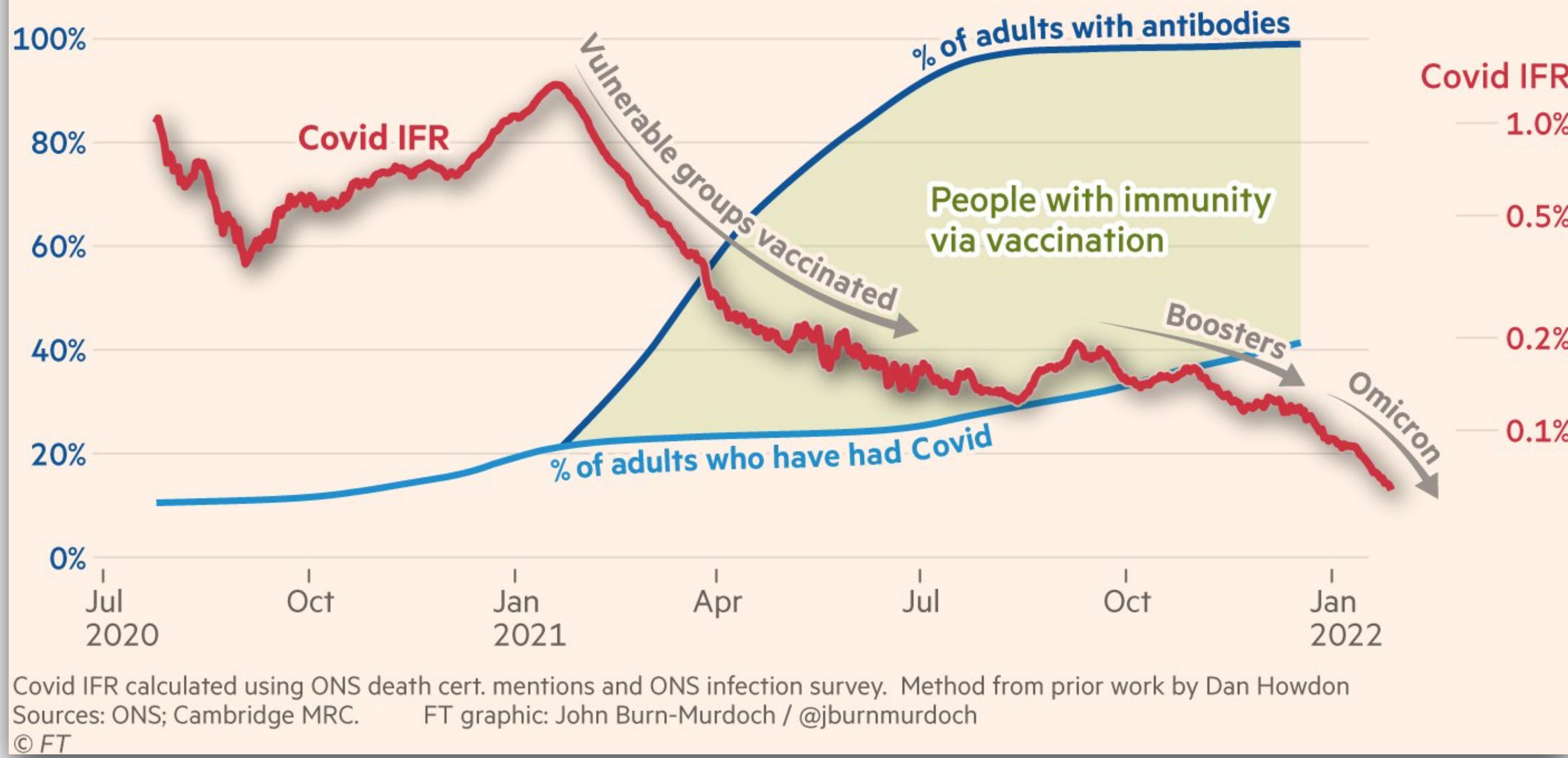
John Tukey

Data visualisation - Telling a story with data

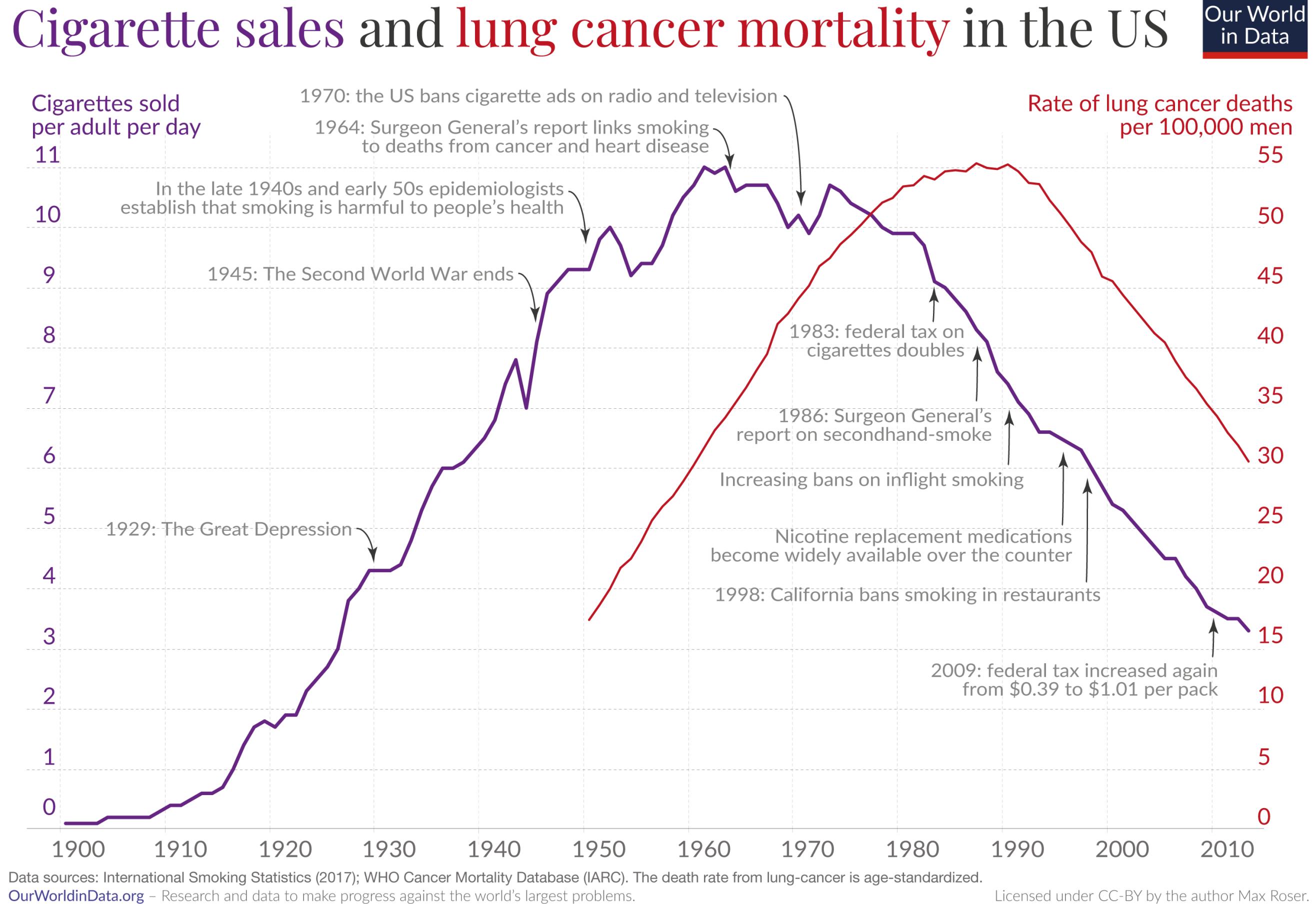
Data visualisation - Telling a story with data

Covid has grown gradually less lethal over the pandemic, mainly due to immunity, the majority of which has come via vaccines

Evolution of Covid's infection fatality ratio in England, overlaid on levels and sources of immunity



Data visualisation - Telling a story with data



Data visualisation - Telling a story with data

Cigarette sales and lung cancer mortality in the US

Our World
in Data

Cigarettes sold
per adult per day

1970: the US bans cigarette ads on radio and television
1964: Surgeon General's report links smoking
to deaths from cancer and heart disease

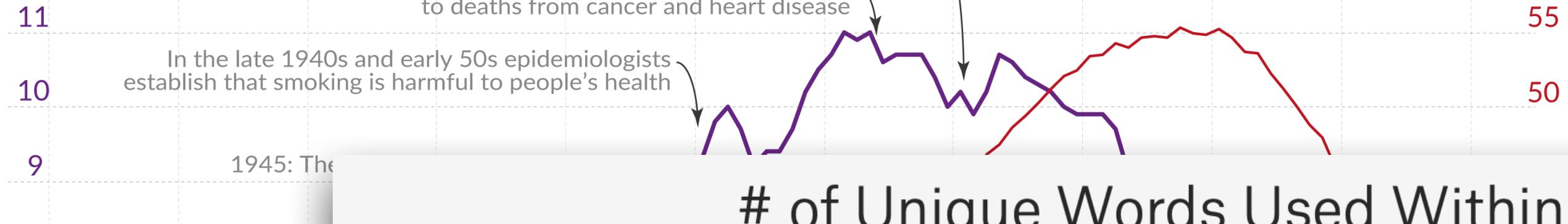
In the late 1940s and early 50s epidemiologists
establish that smoking is harmful to people's health

1945: The

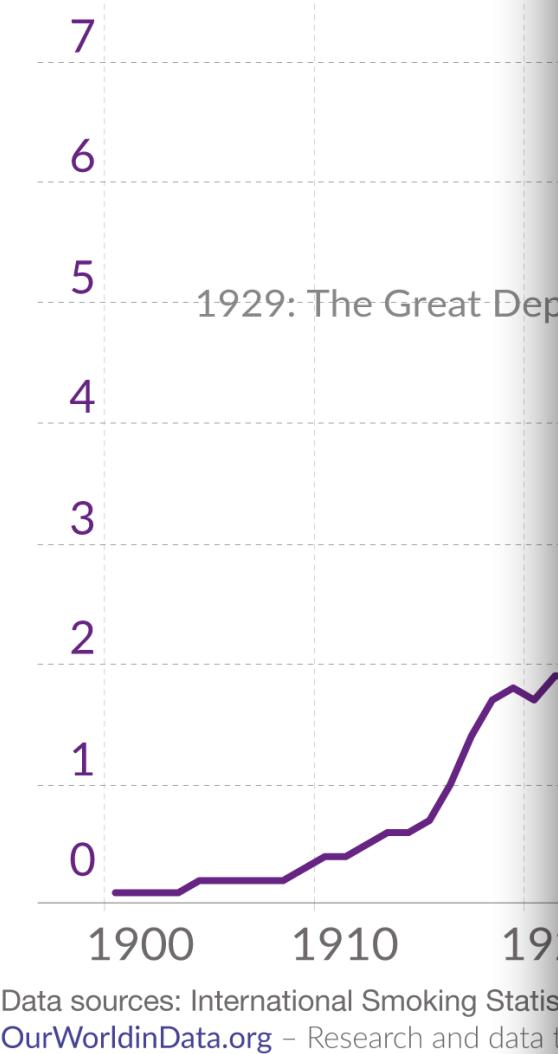
Rate of lung cancer deaths
per 100,000 men

50

55



of Unique Words Used Within Artist's First 35,000 Lyrics



3,000 words

4,000

5,000

6,000 words

All Just Find an Artist

1929: The Great Dep
Data sources: International Smoking Statistics
OurWorldInData.org – Research and data t

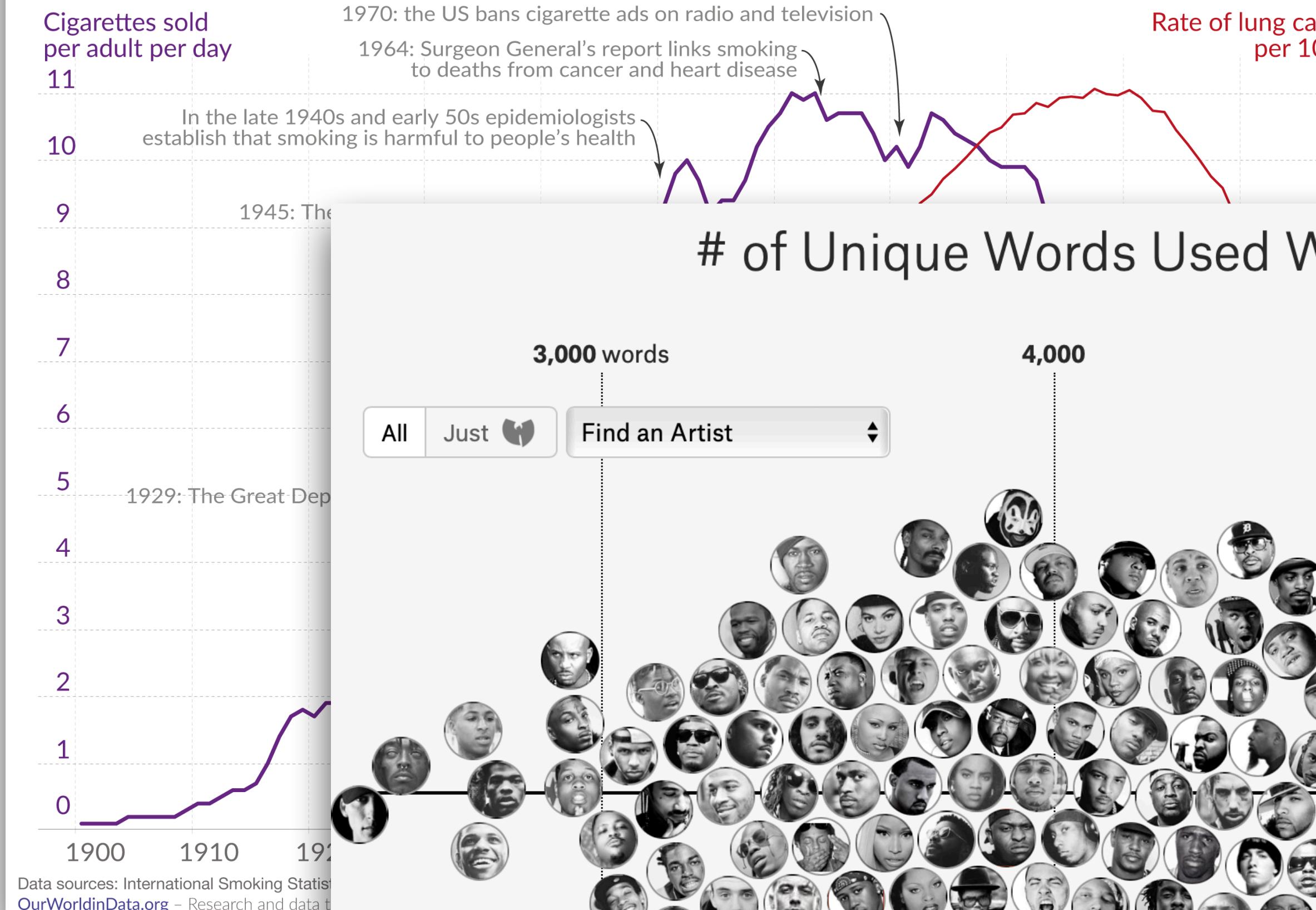
Aesop
Rock:
7,879
unique
words
used

7,300
words

ThePudding

Data visualisation - Telling a story with data

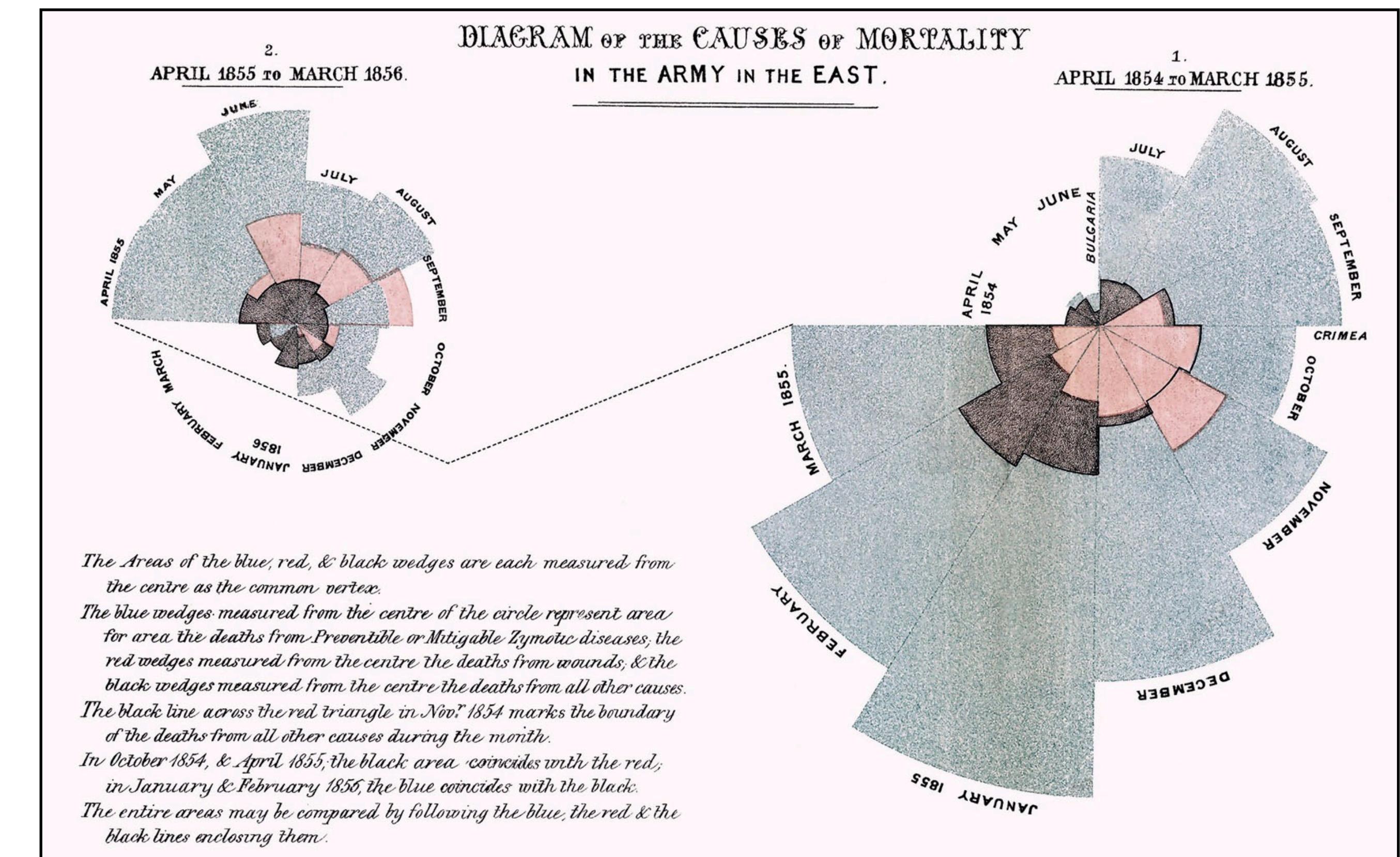
Cigarette sales and lung cancer mortality in the US



Florence Nightingale - “The lady with the lamp”



Florence Nightingale - Data visualisation pioneer

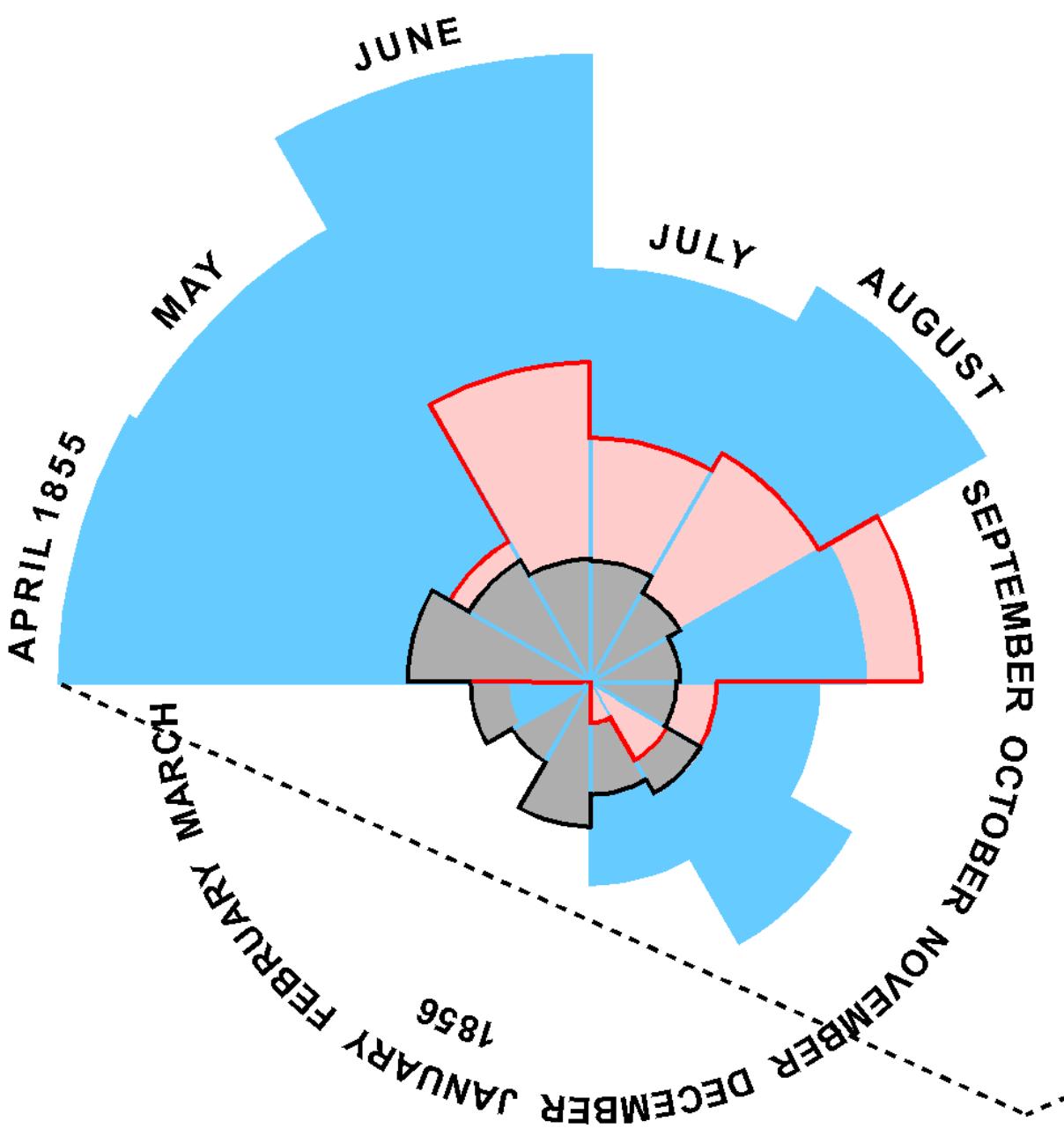


First woman to become a member of the Royal Statistical Society (1858)

DIAGRAM OF THE CAUSES OF MORTALITY IN THE ARMY IN THE EAST.

2.

APRIL 1855 TO MARCH 1856.



After sanitary commission arrival

The Areas of the blue, red, & black wedges are each measured from the centre as the common vertex

The blue wedges measured from the centre of the circle represent area for area the deaths from Preventible or Mitigable Zymotic Diseases, the red wedges measured from the centre the deaths from wounds, & the black wedges measured from the centre the deaths from all other causes

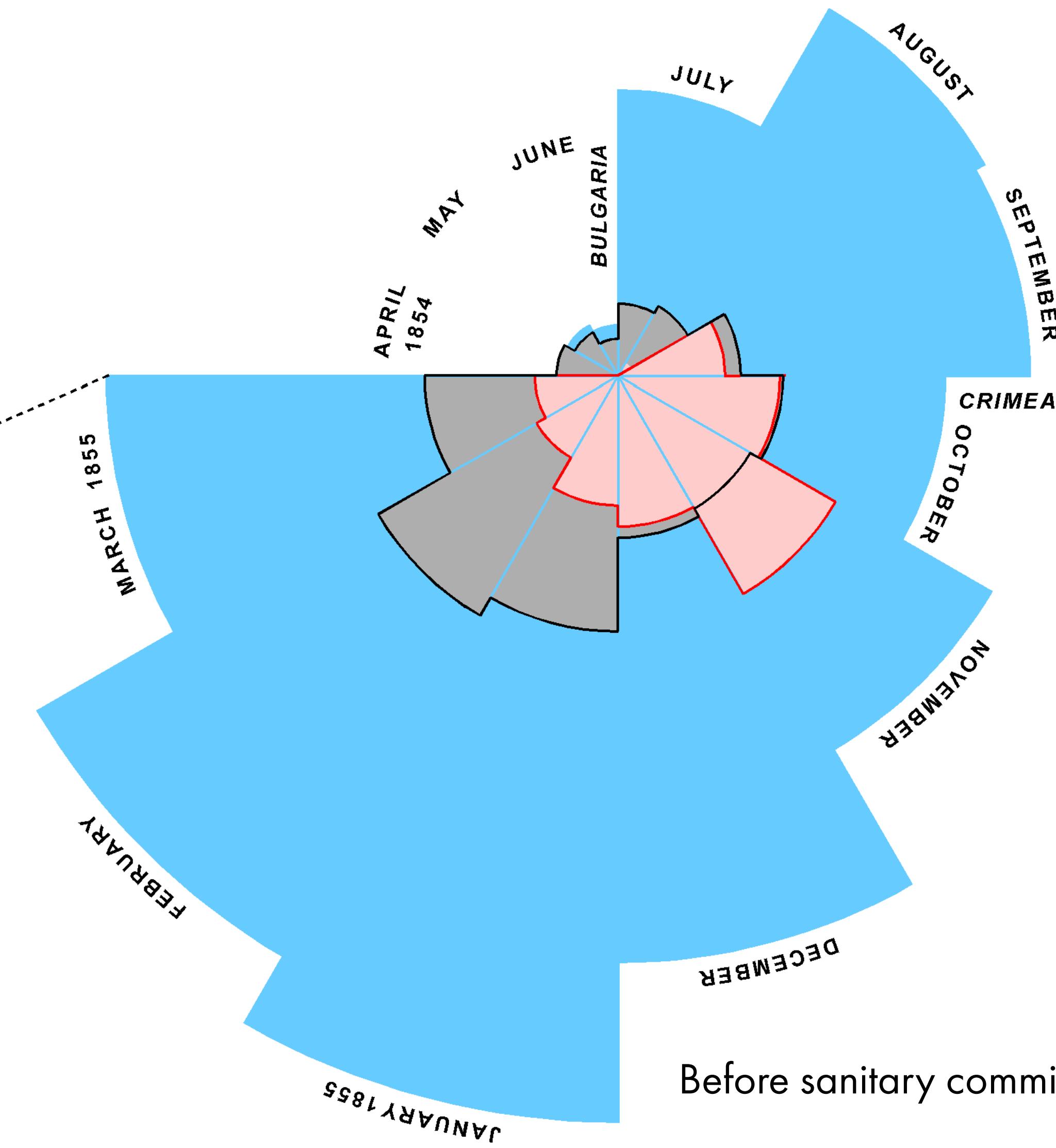
The black line across the red triangle in Nov 1854 marks the boundary of the deaths from all other causes during the month

In October 1854, & April 1855, the black area coincides with the red, in January & February 1856, the blue coincides with the black

The entire areas may be compared by following the blue, the red & the black lines enclosing them

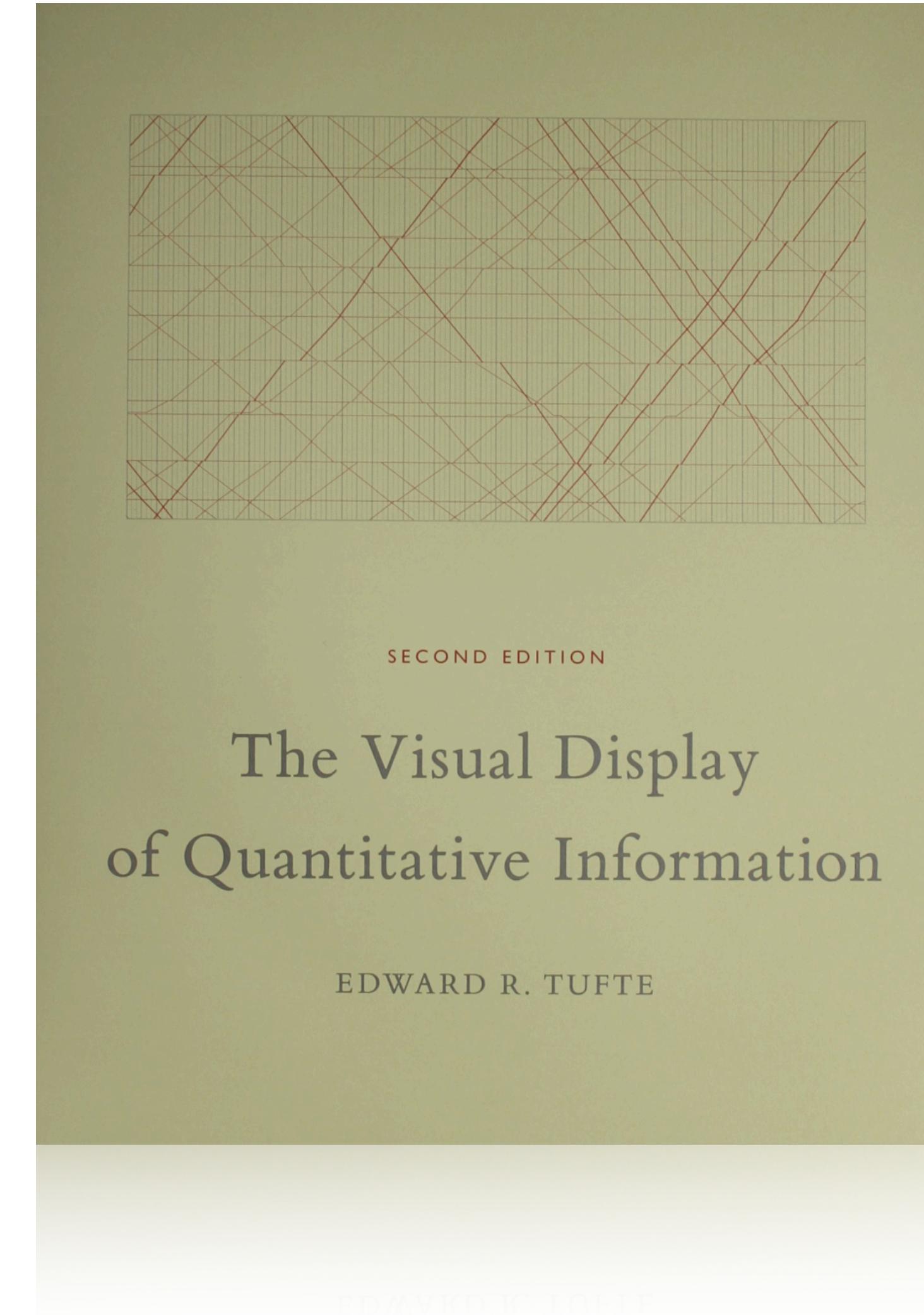
1.

APRIL 1854 TO MARCH 1855.



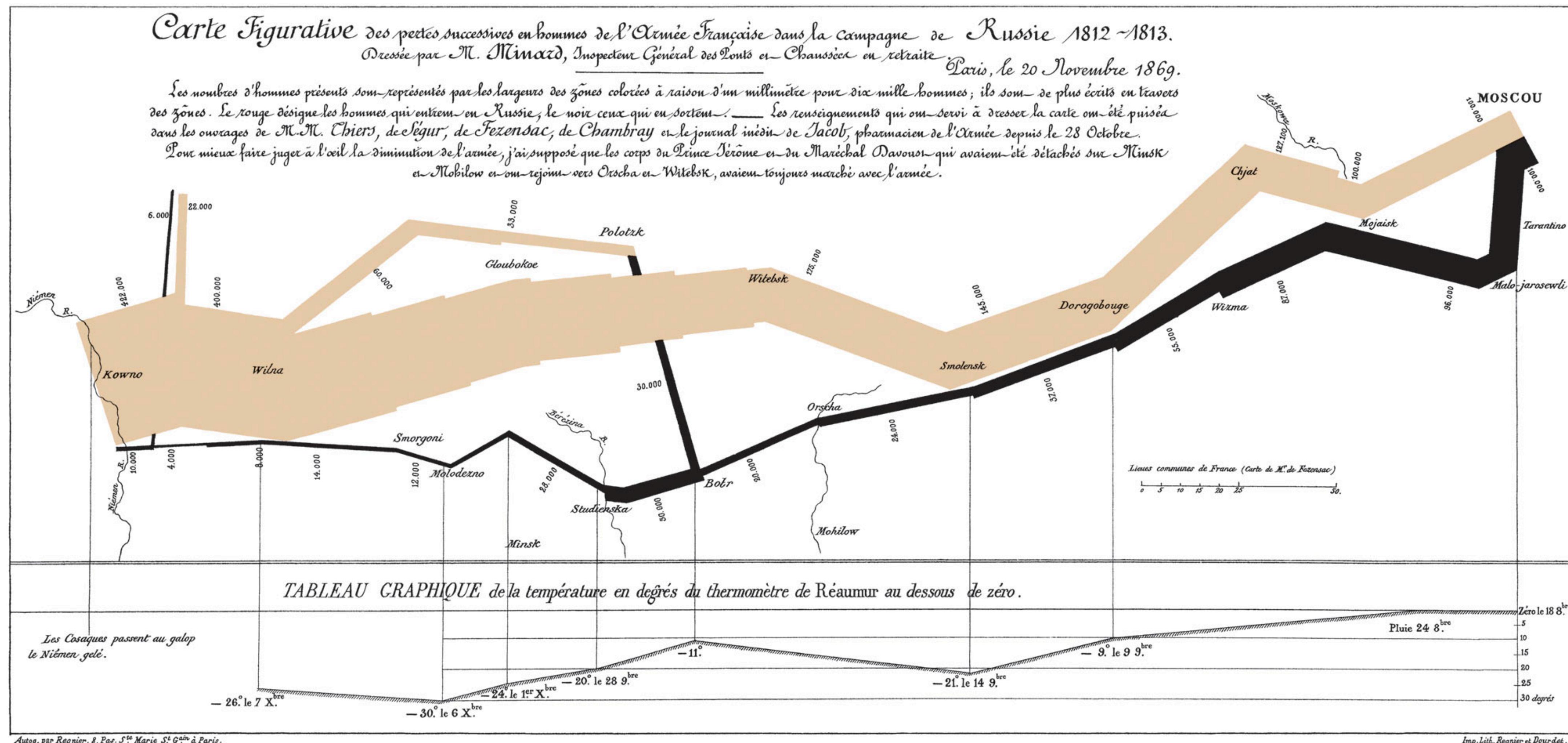
Before sanitary commission arrival

Pioneers of data visualisation - Edward Tufte ("ET")



www.edwardtufte.com

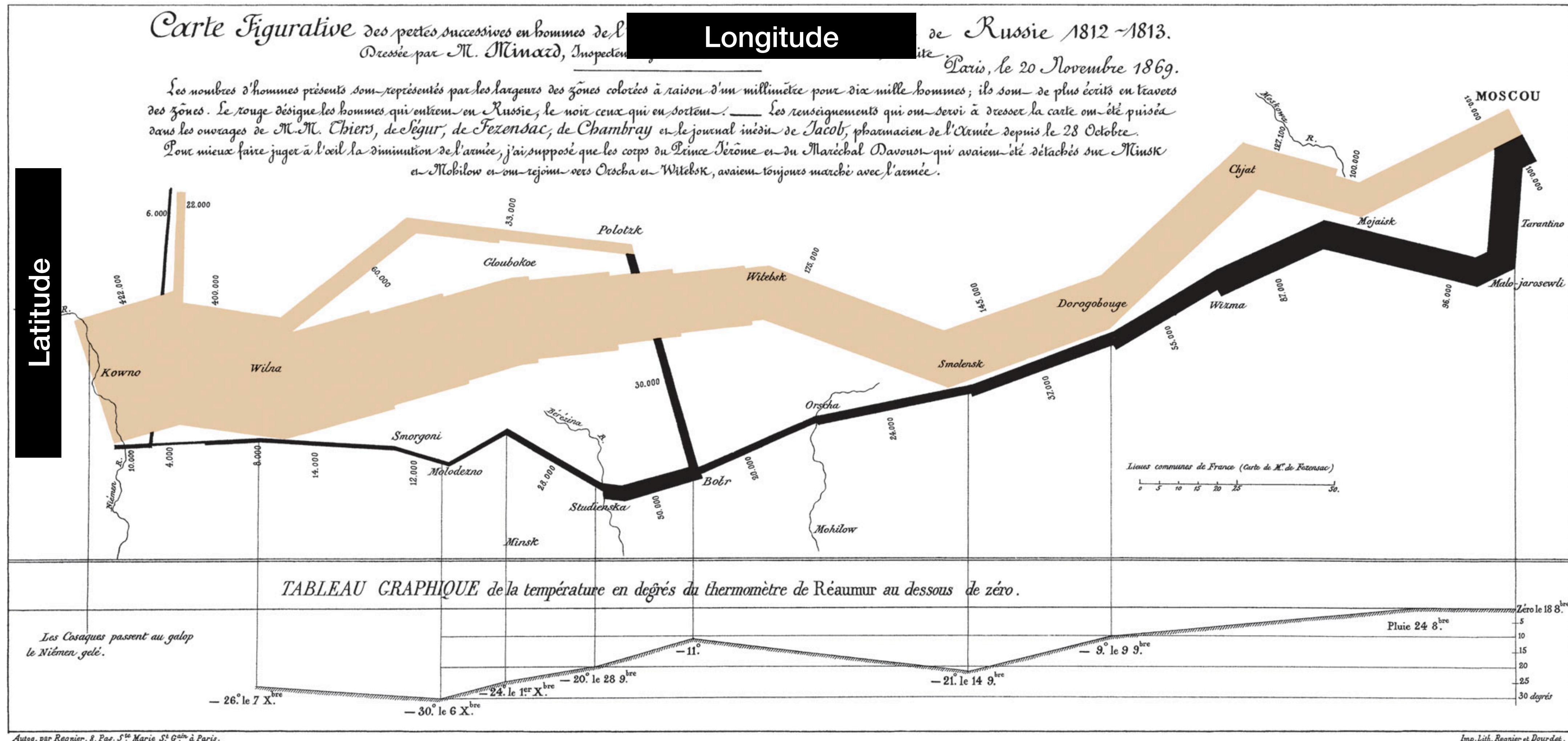
"The best statistical graph ever drawn"



Charles Joseph Minard

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

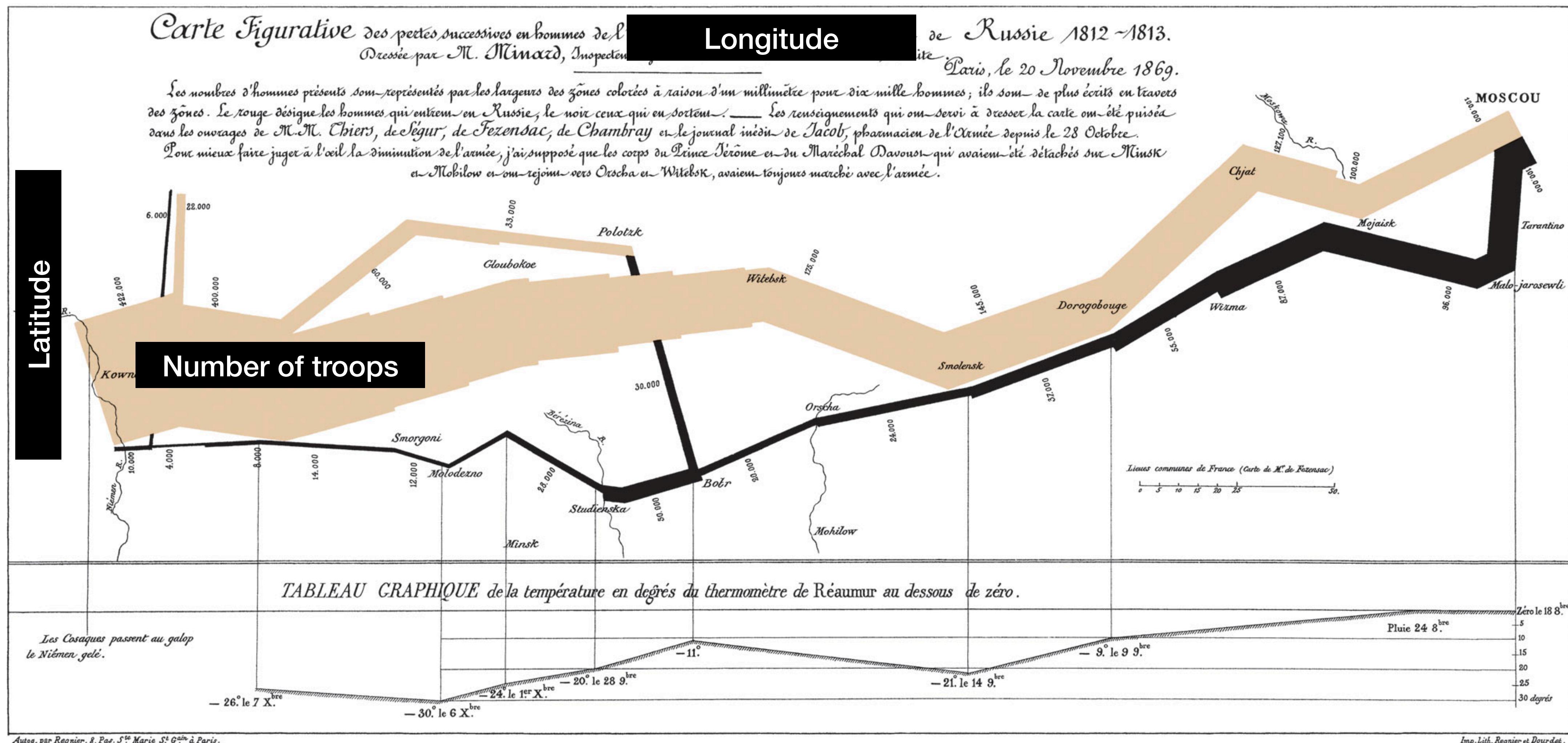
"The best statistical graph ever drawn"



Charles Joseph Minard

"Carte figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812–1813"

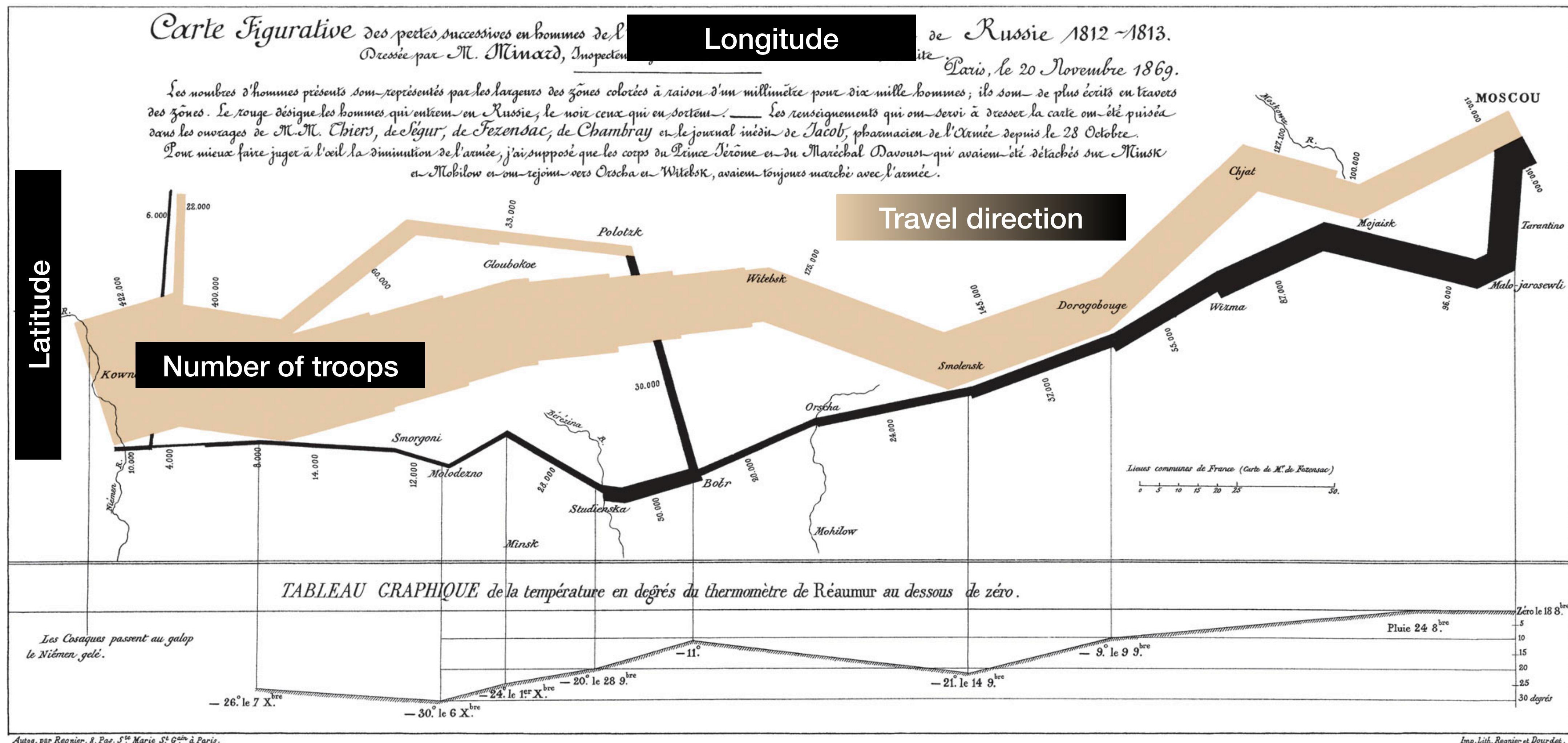
"The best statistical graph ever drawn"



Charles Joseph Minard

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

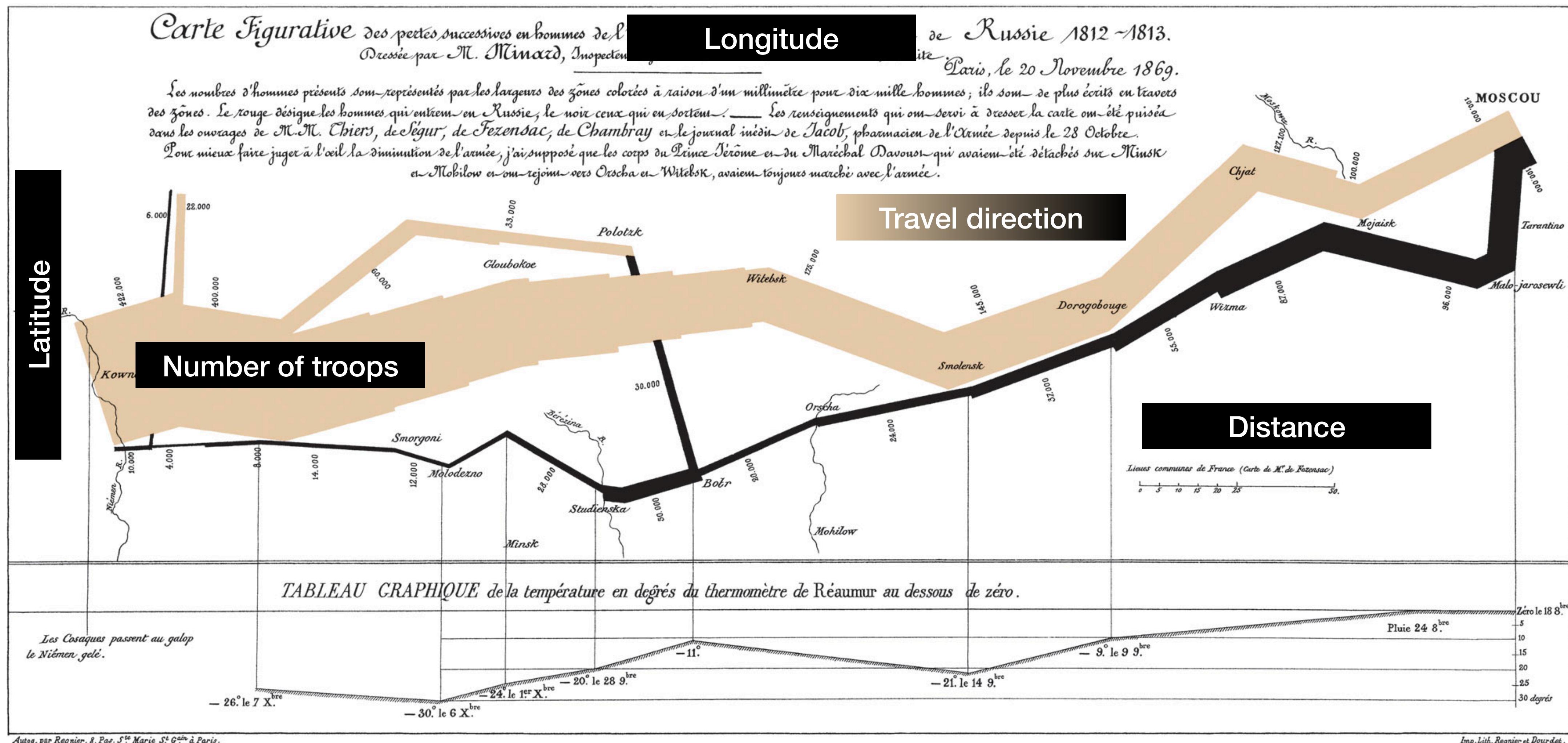
"The best statistical graph ever drawn"



Charles Joseph Minard

"Carte figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812–1813"

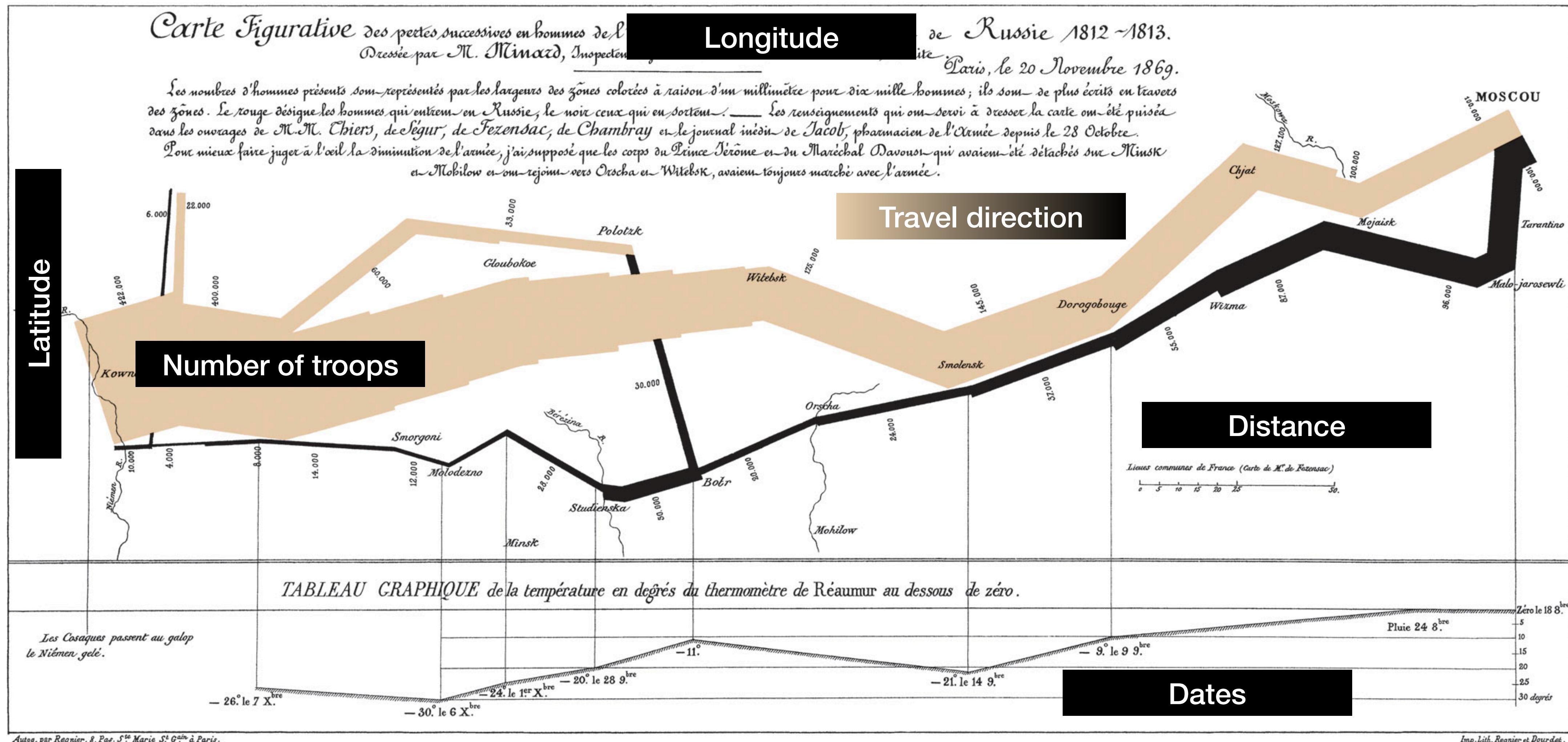
"The best statistical graph ever drawn"



Charles Joseph Minard

"Carte figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812–1813"

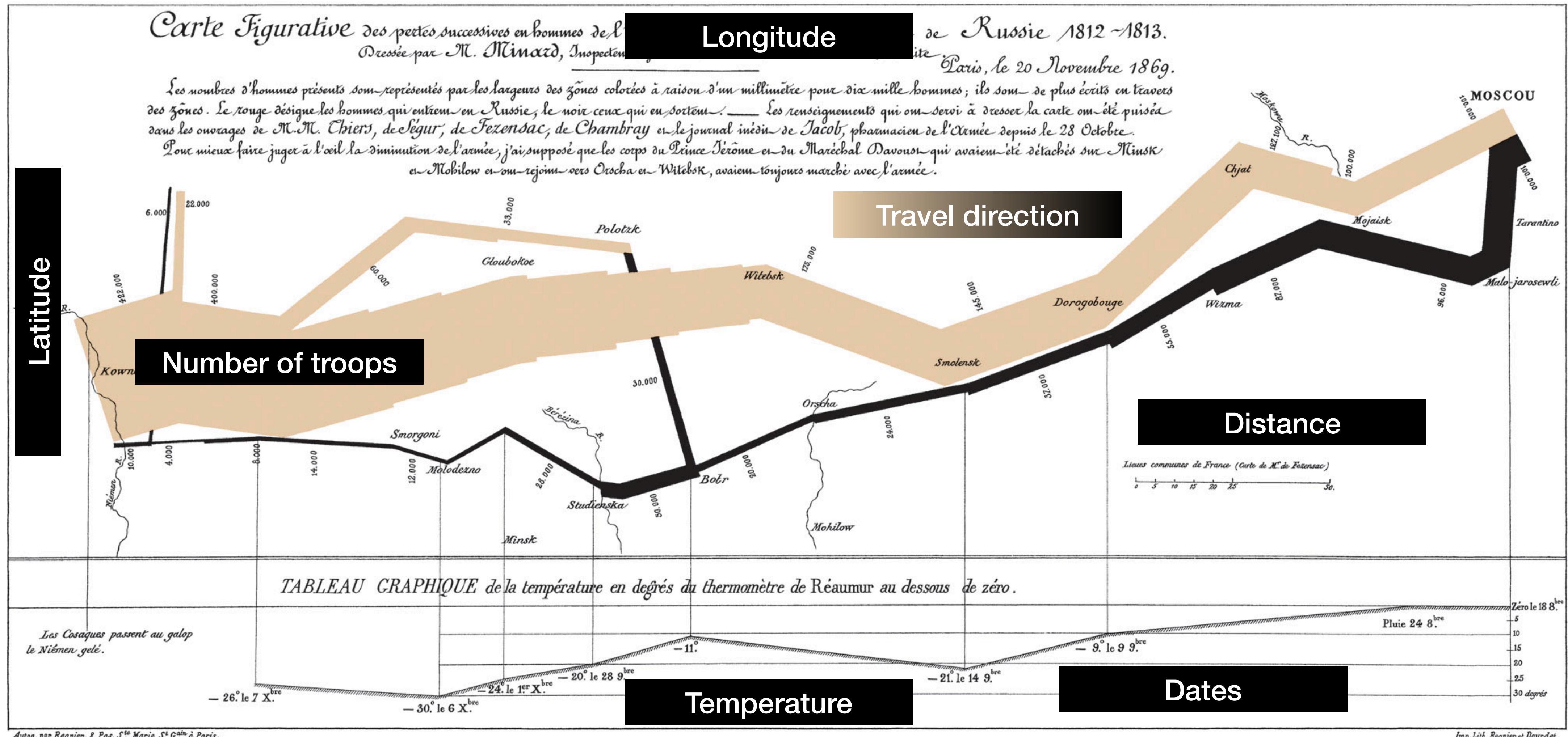
"The best statistical graph ever drawn"



Charles Joseph Minard

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

"The best statistical graph ever drawn"

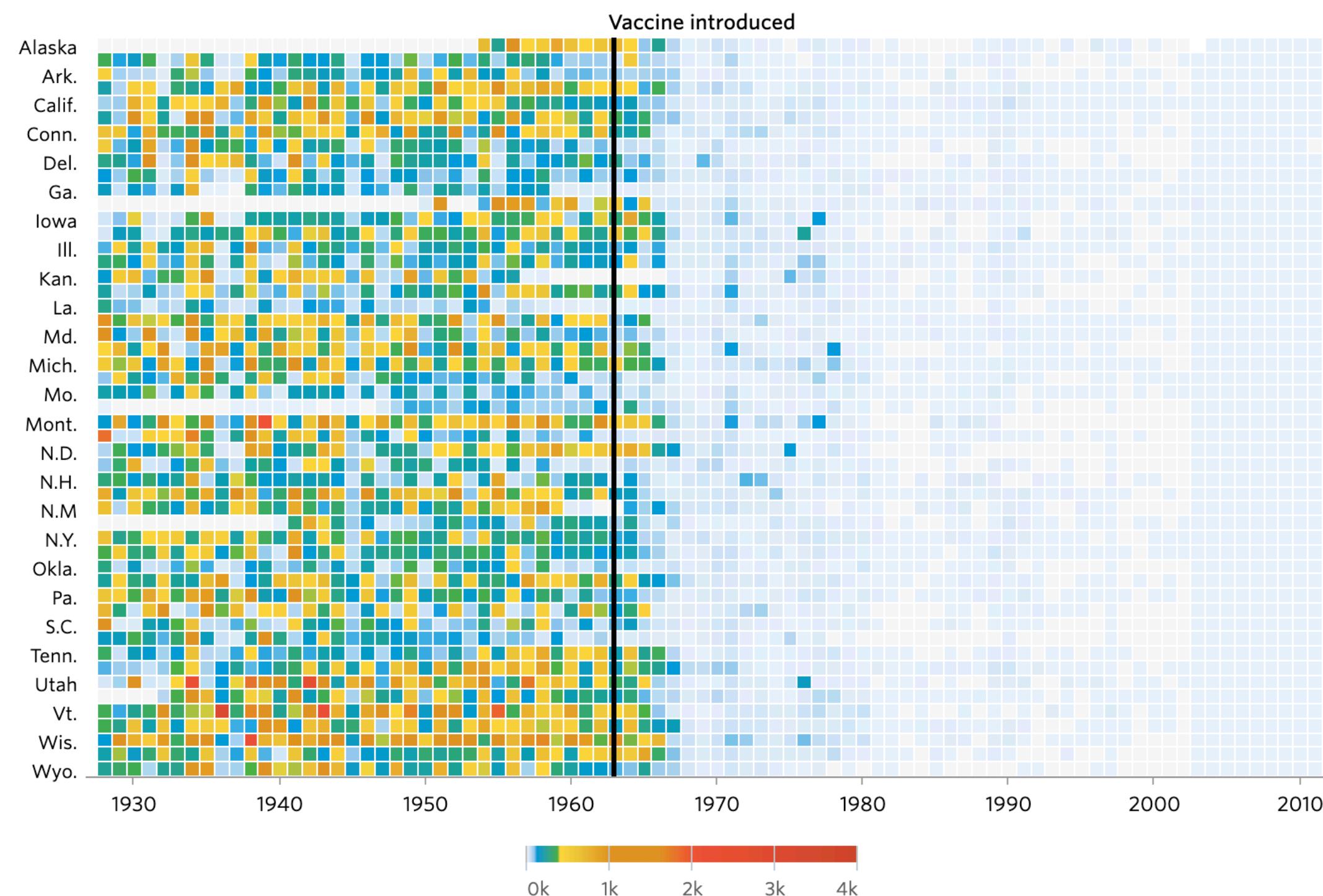


Charles Joseph Minard

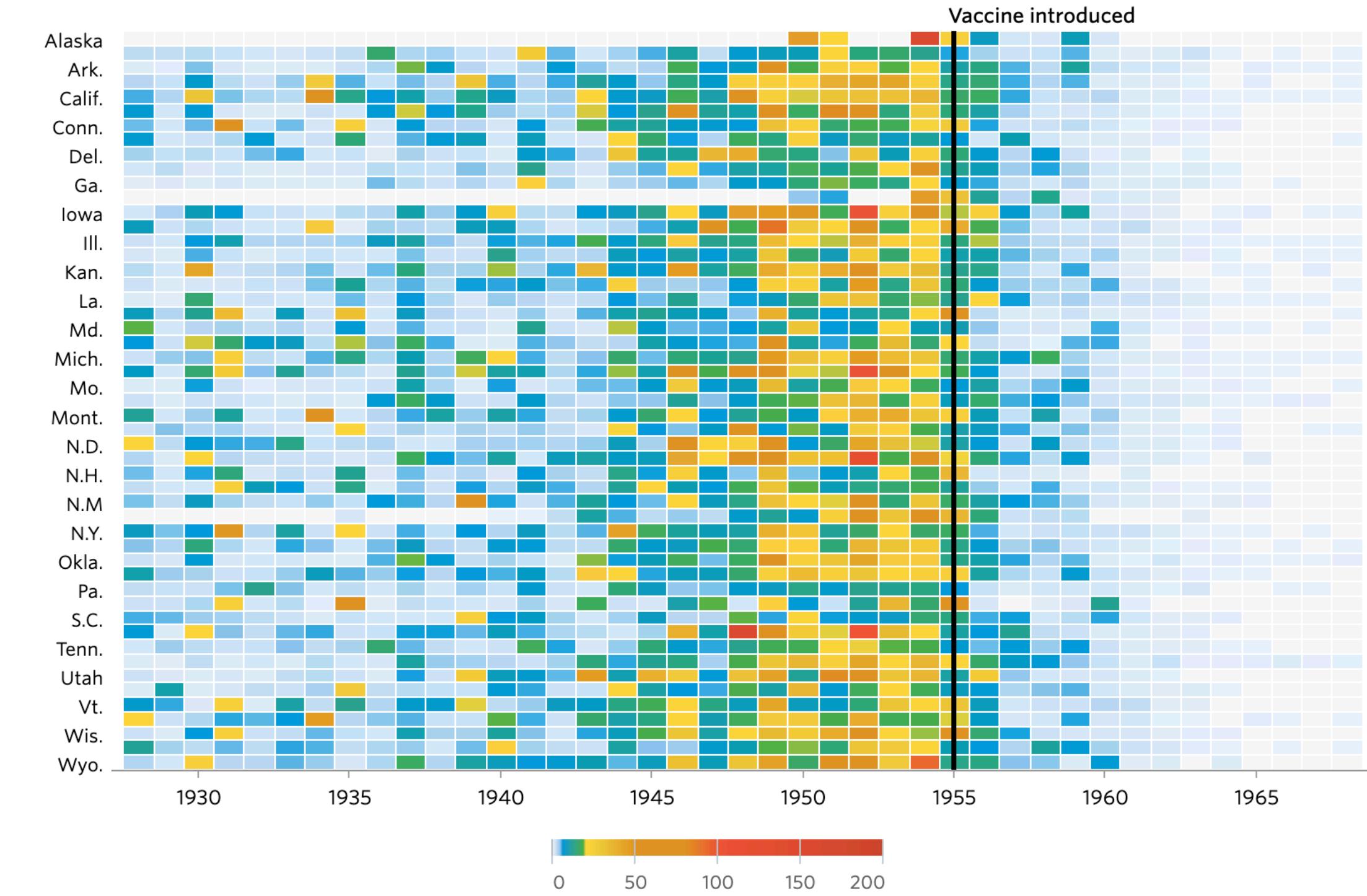
"Carte figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813"

The impact of vaccines on infectious diseases

Measles



Polio



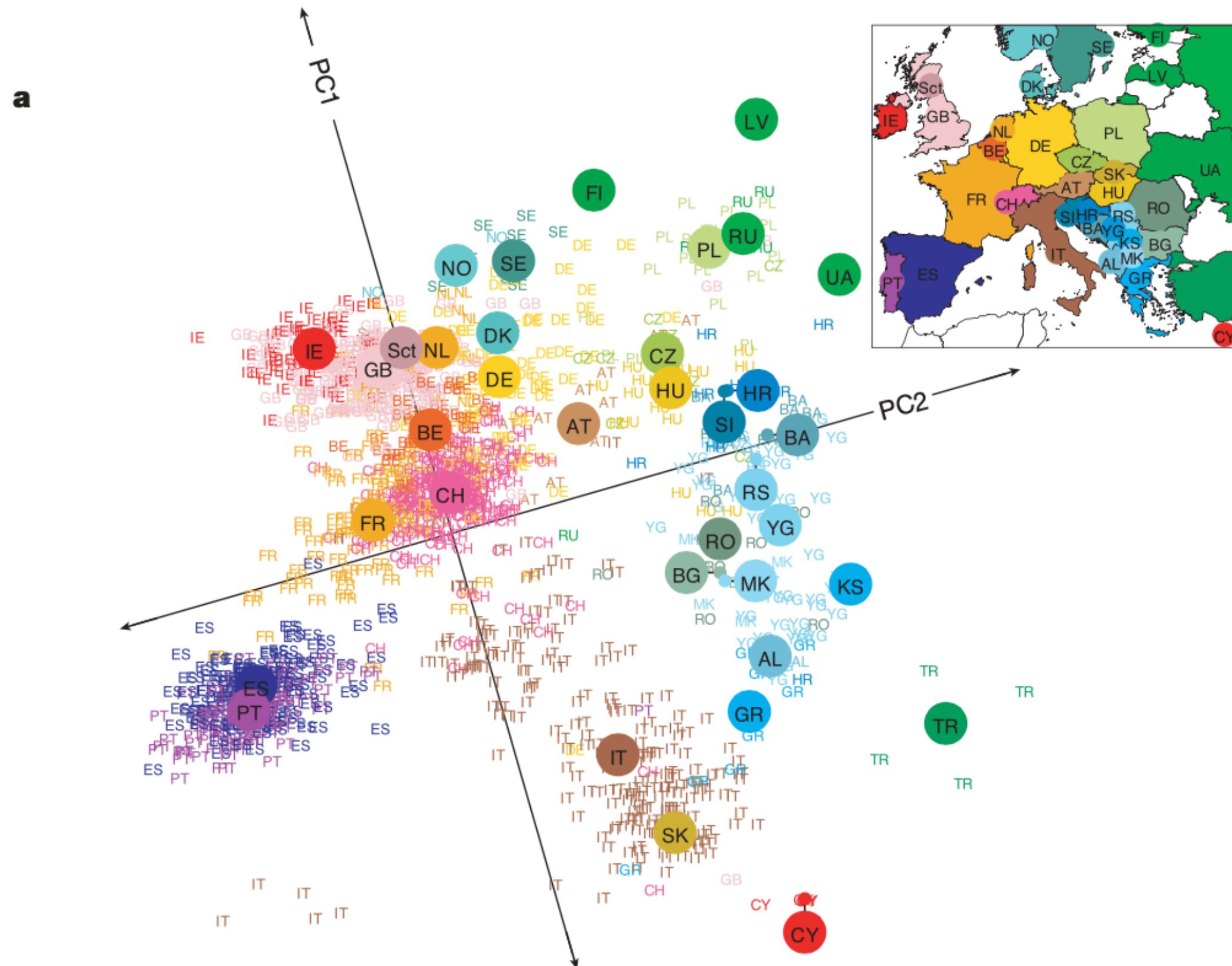
Battling Infectious Diseases in the 20th Century: The Impact of Vaccines

By [Tynan DeBold](#) and [Dov Friedman](#)

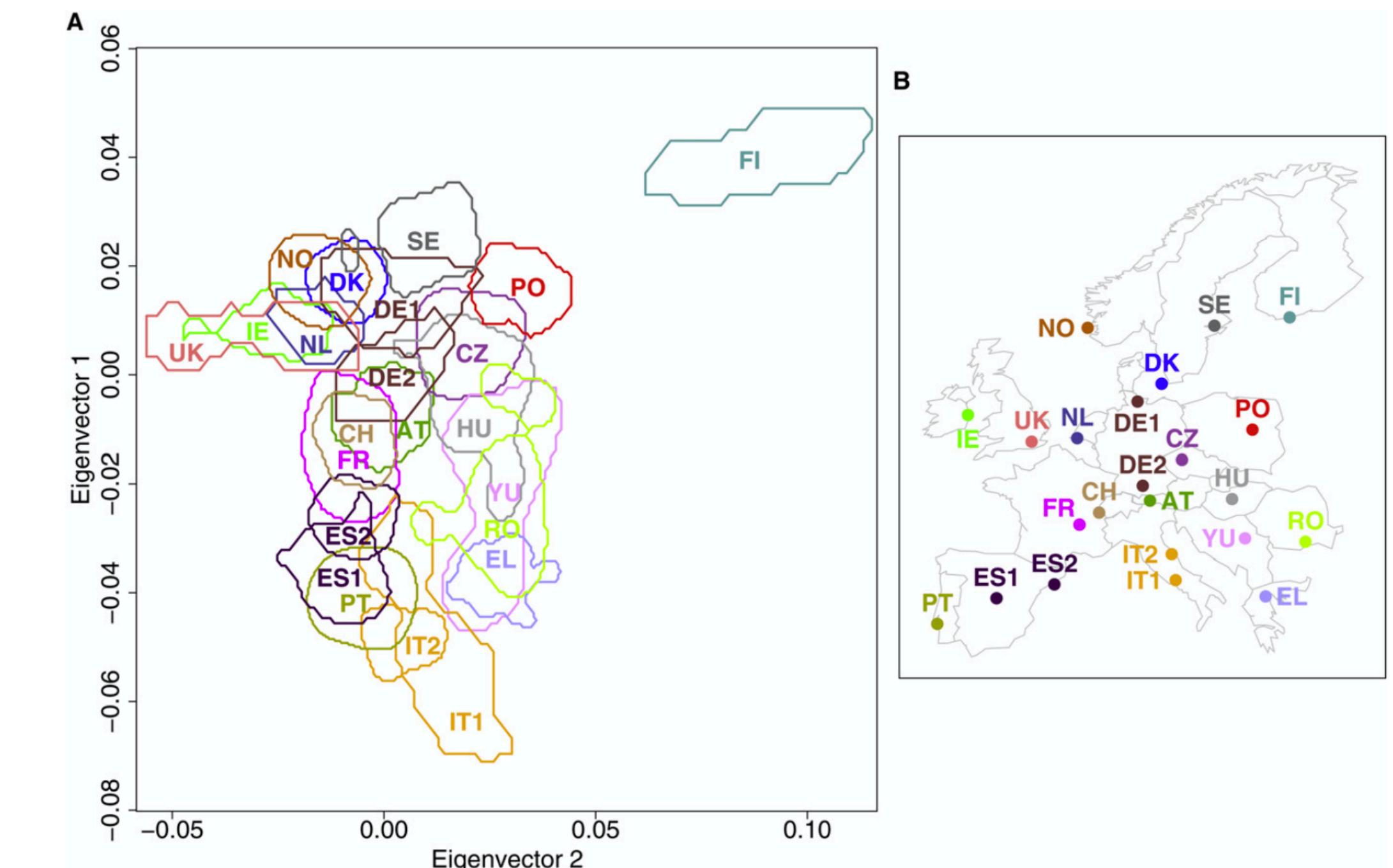
Published Feb. 11, 2015 at 3:45 p.m. ET

A tale of two visualisations

Novembre et al 2008



Lao et al 2008

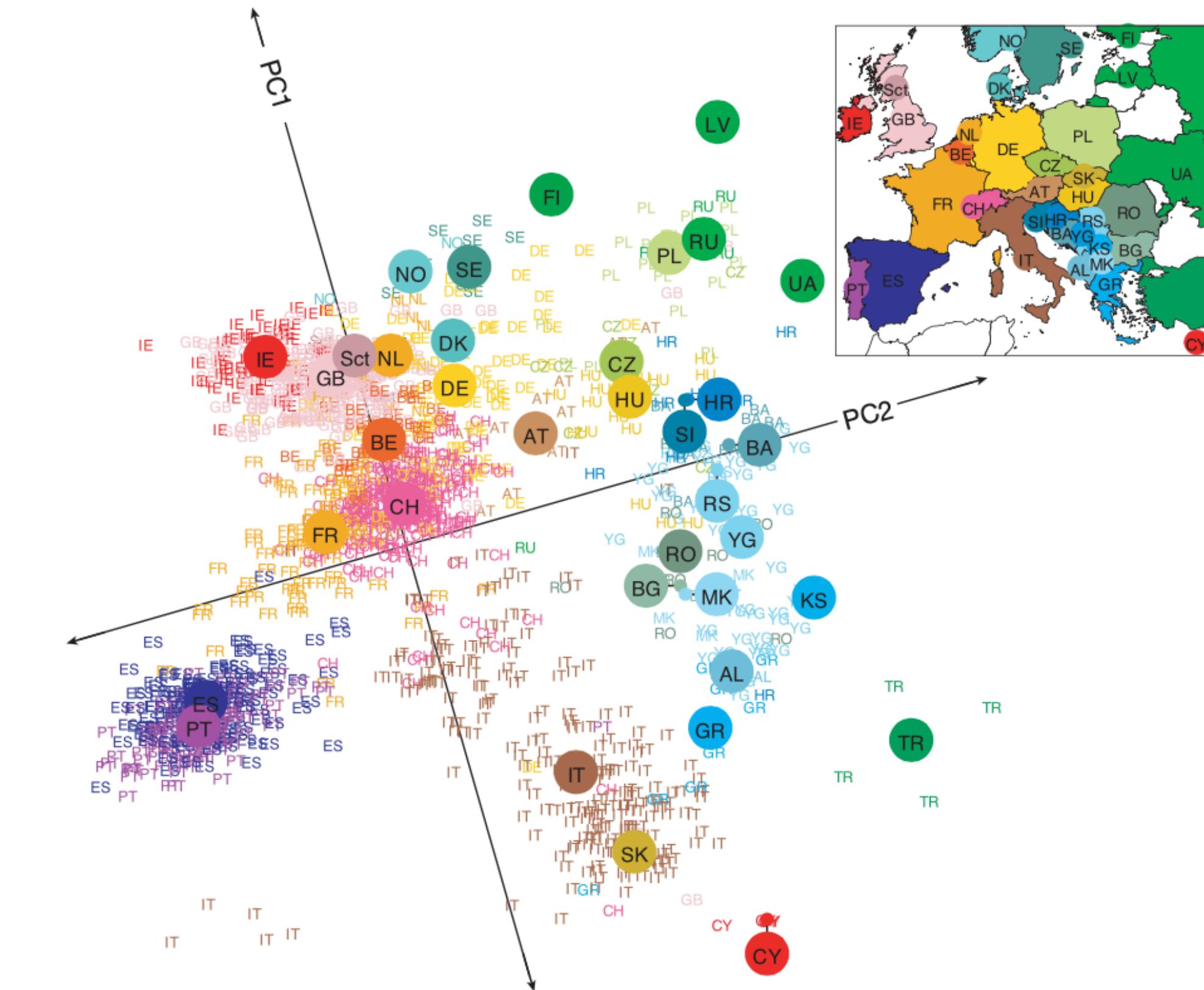


Which visualisation do you prefer?

A tale of two visualisations

Novembre et al, 31/8/2008

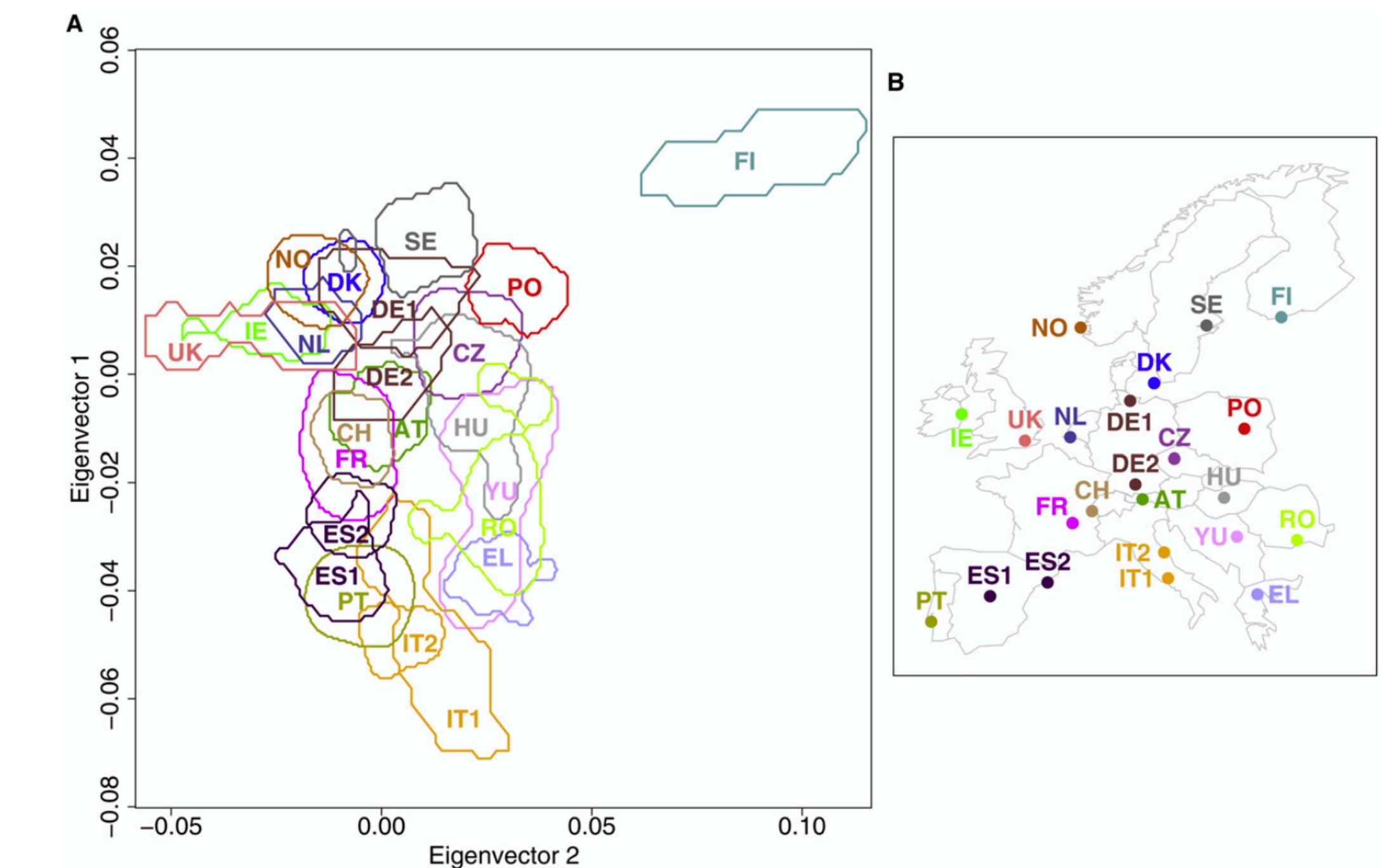
"Genes mirror geography in Europe"



1,822 citations (Nature)

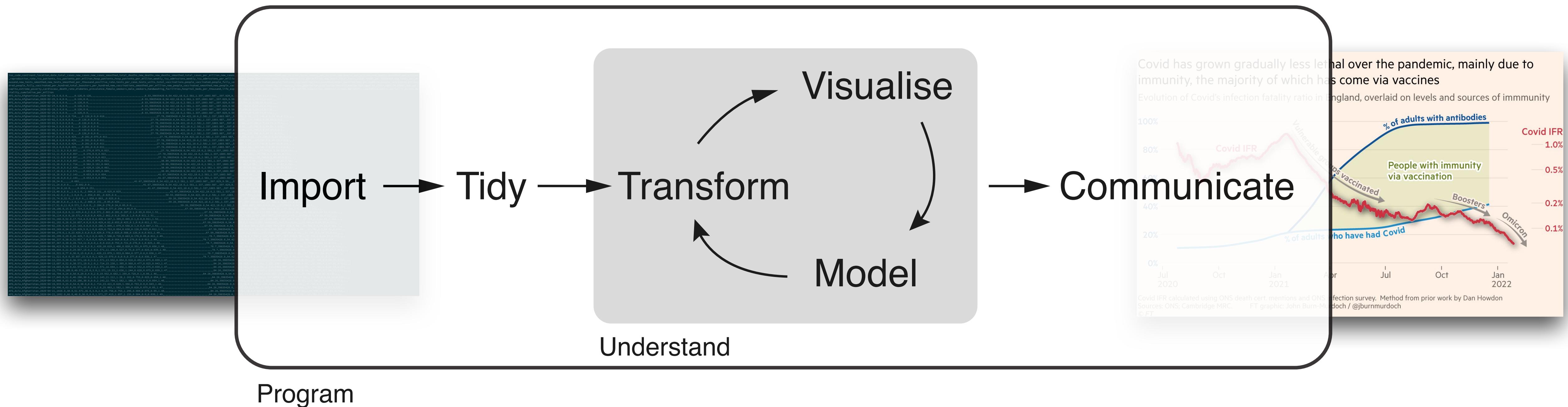
Lao et al, 26/8/2008

"Correlation between Genetic and Geographic Structure in Europe"

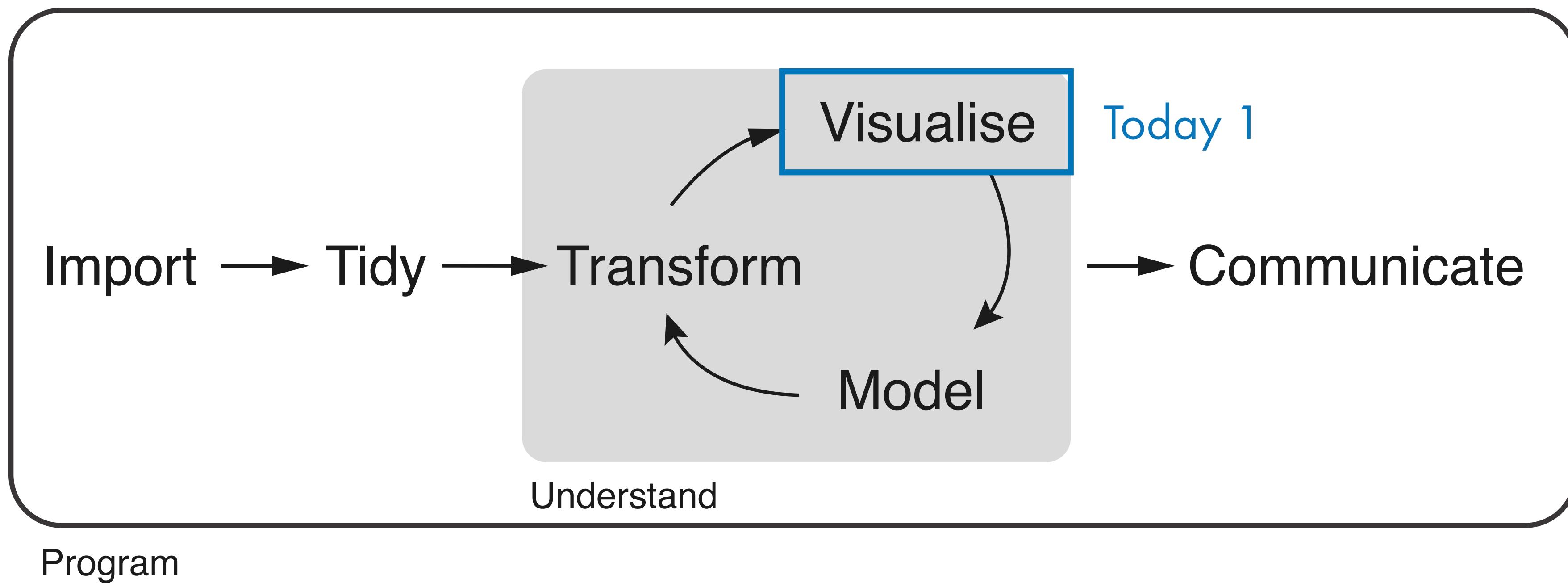


610 citations (Current Biology)

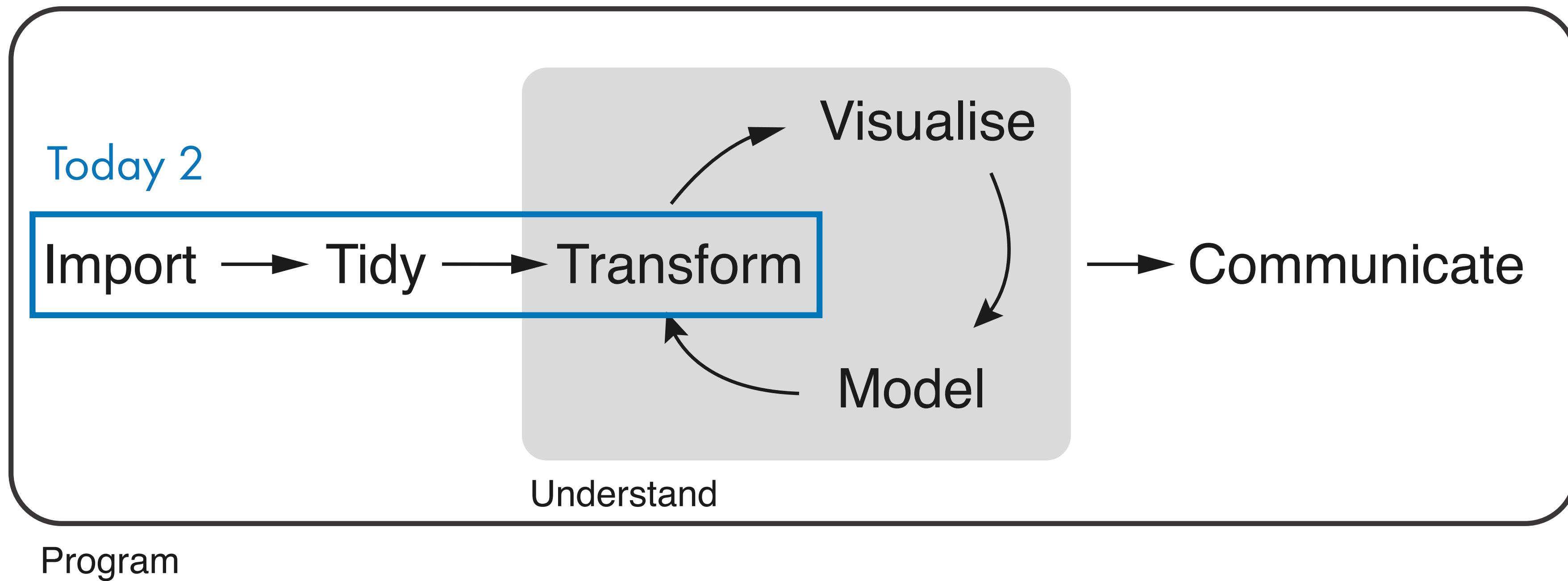
A data science project workflow



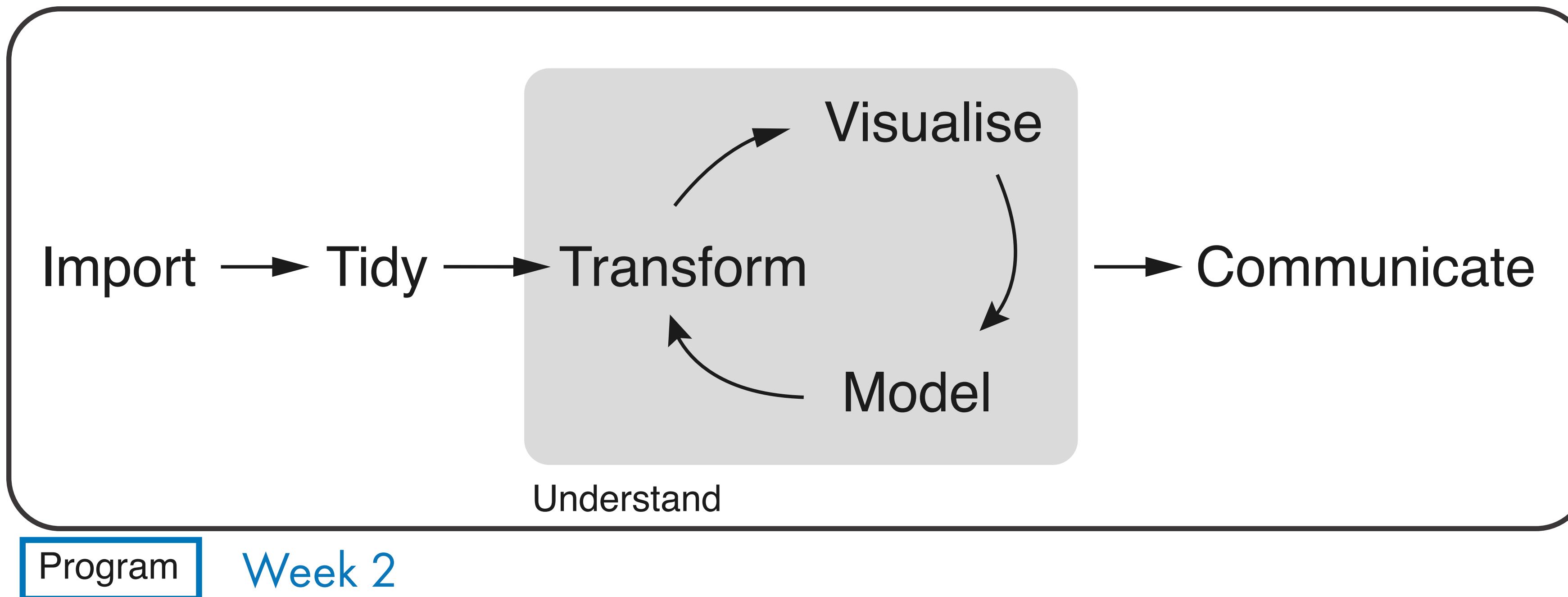
A data science project workflow



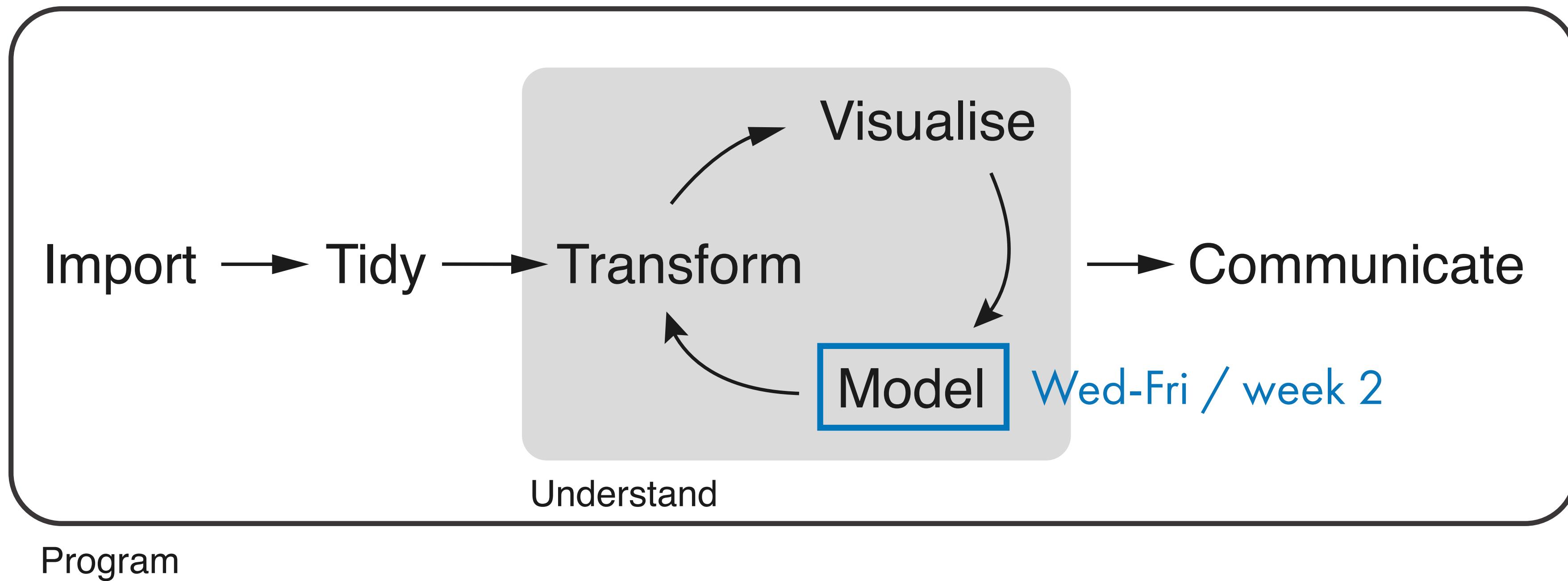
A data science project workflow



A data science project workflow



A data science project workflow



The tidyverse



Tidyverse

Packages Blog Learn Help Contribute

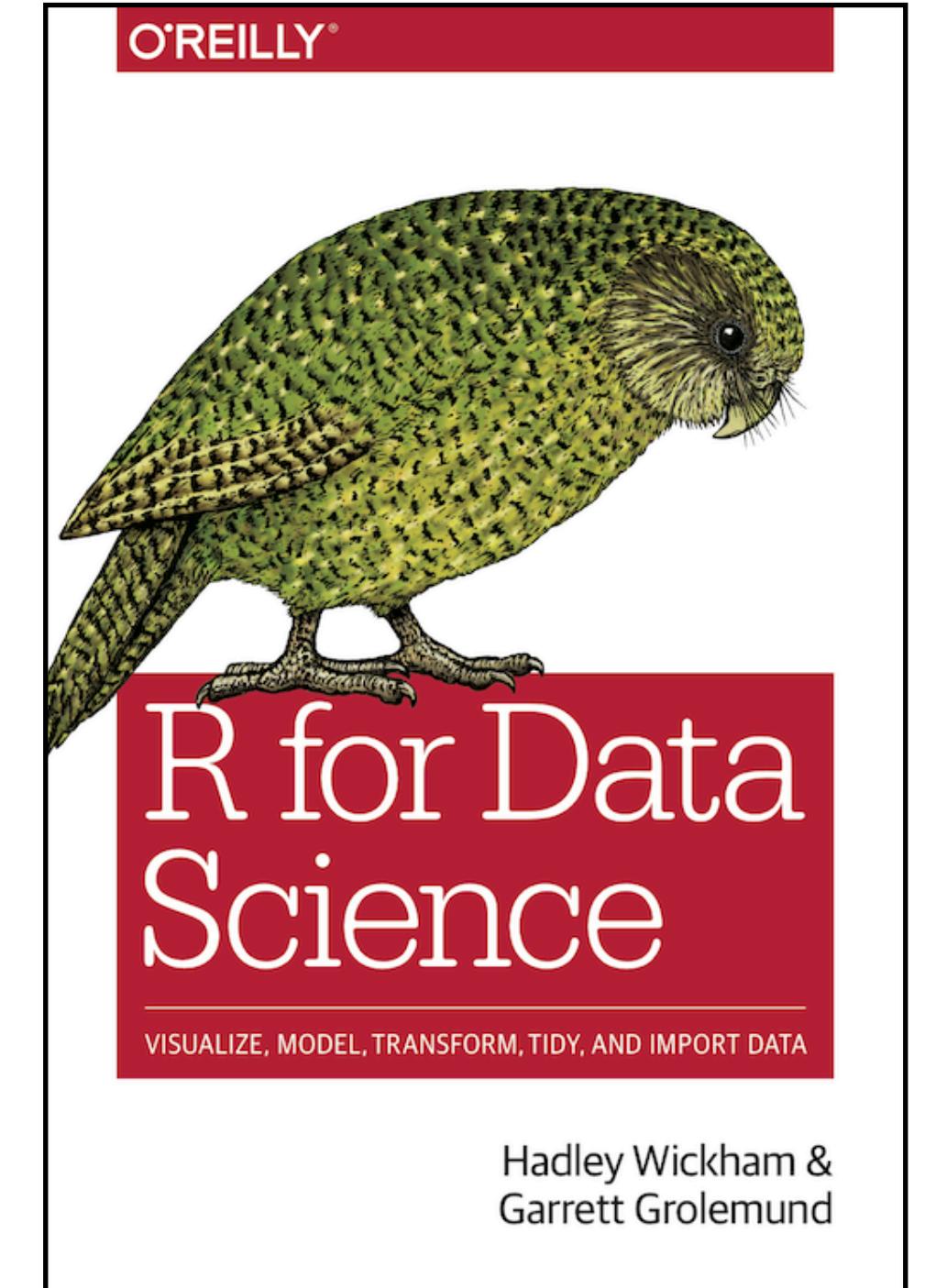
R packages for data science

The tidyverse is an opinionated **collection of R packages** designed for data science. All packages share an underlying design philosophy, grammar, and data structures.

Install the complete tidyverse with:

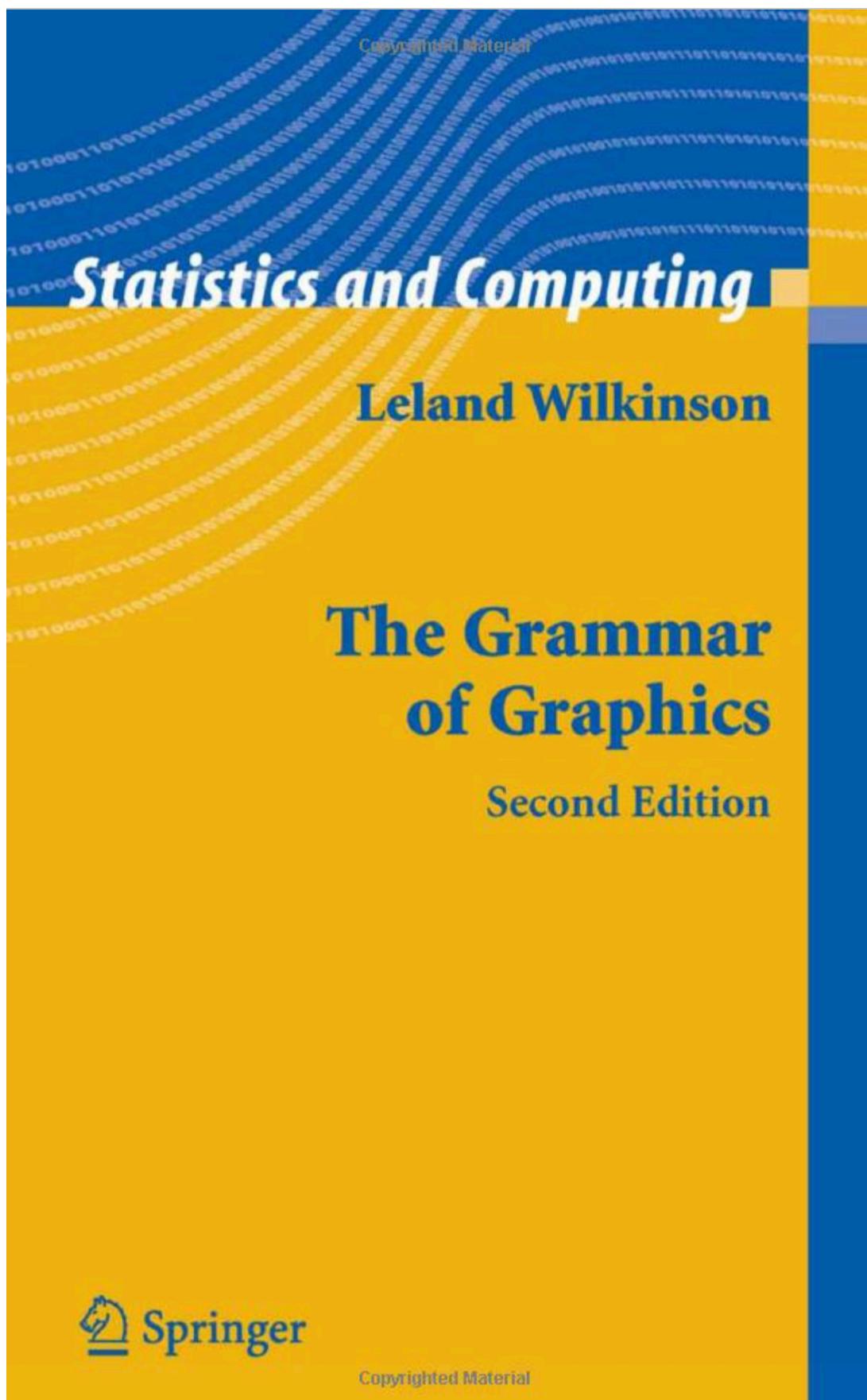
```
install.packages("tidyverse")
```

www.tidyverse.org

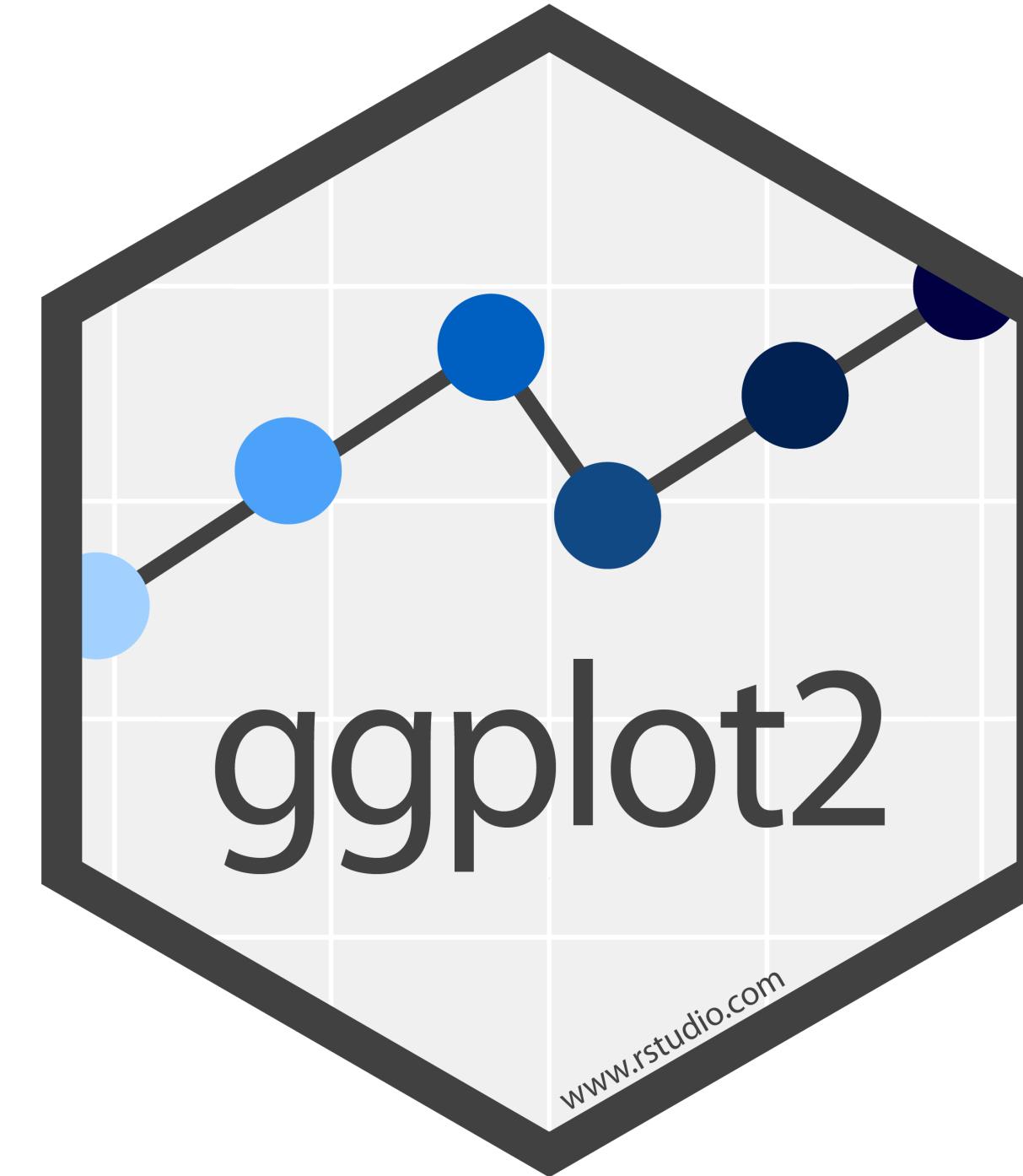


r4ds.had.co.nz

A layered grammar of graphics



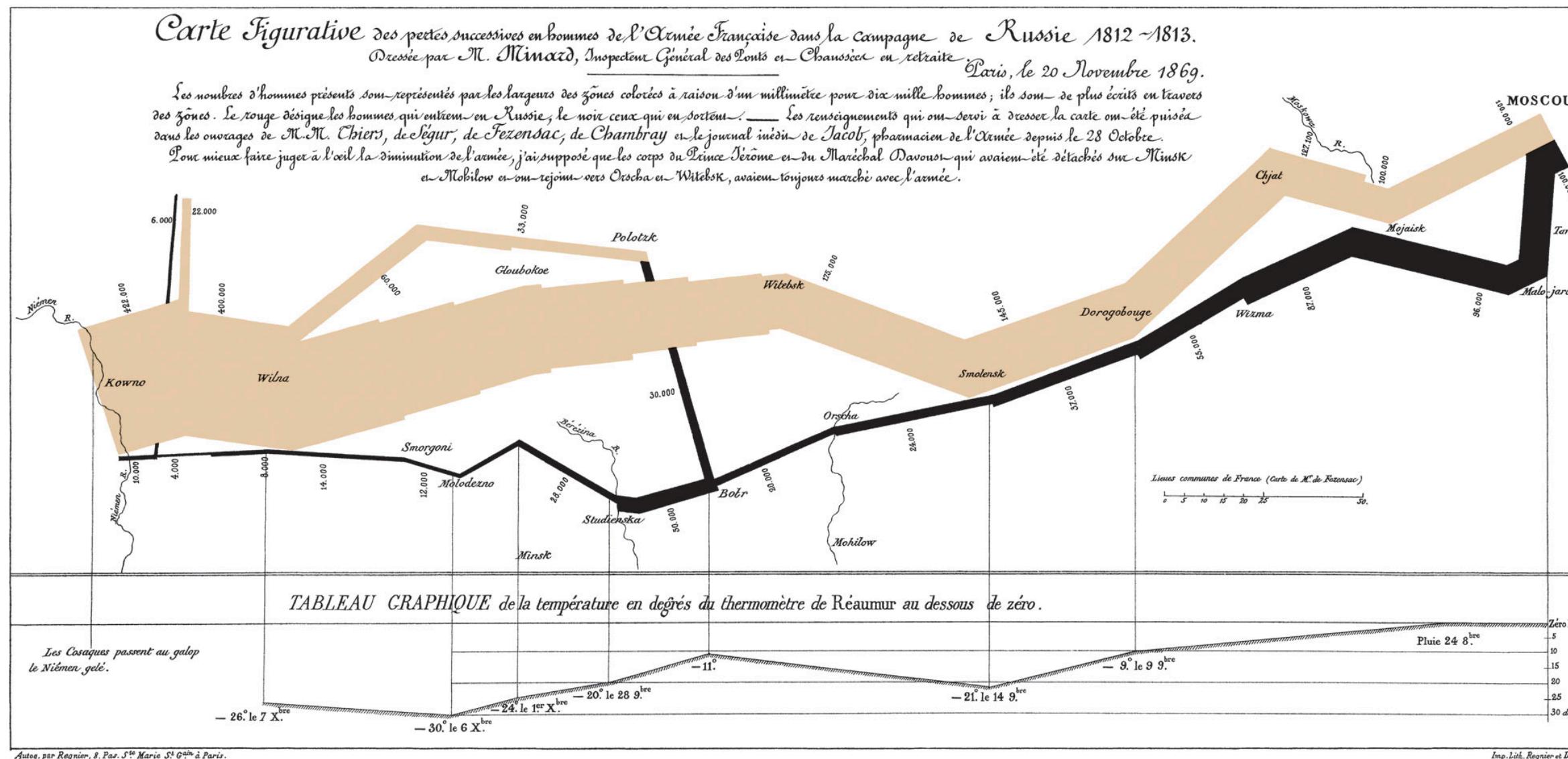
A description of the fundamental features
that underlie all statistical graphics



A toolkit implementing the grammar of
graphics in R

A layered grammar of graphics

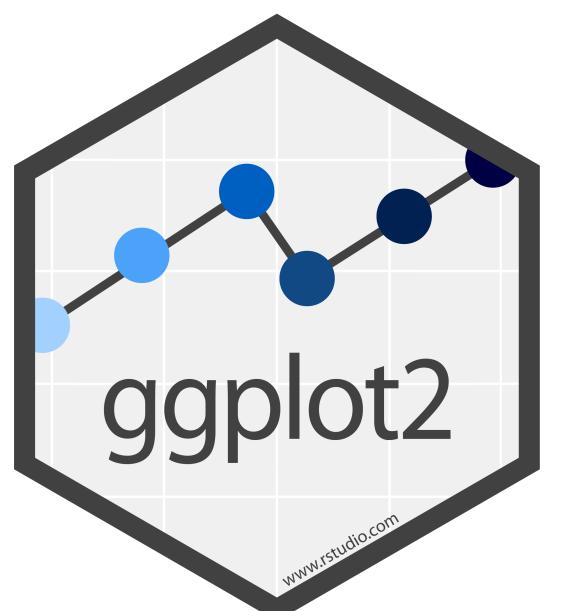
Original



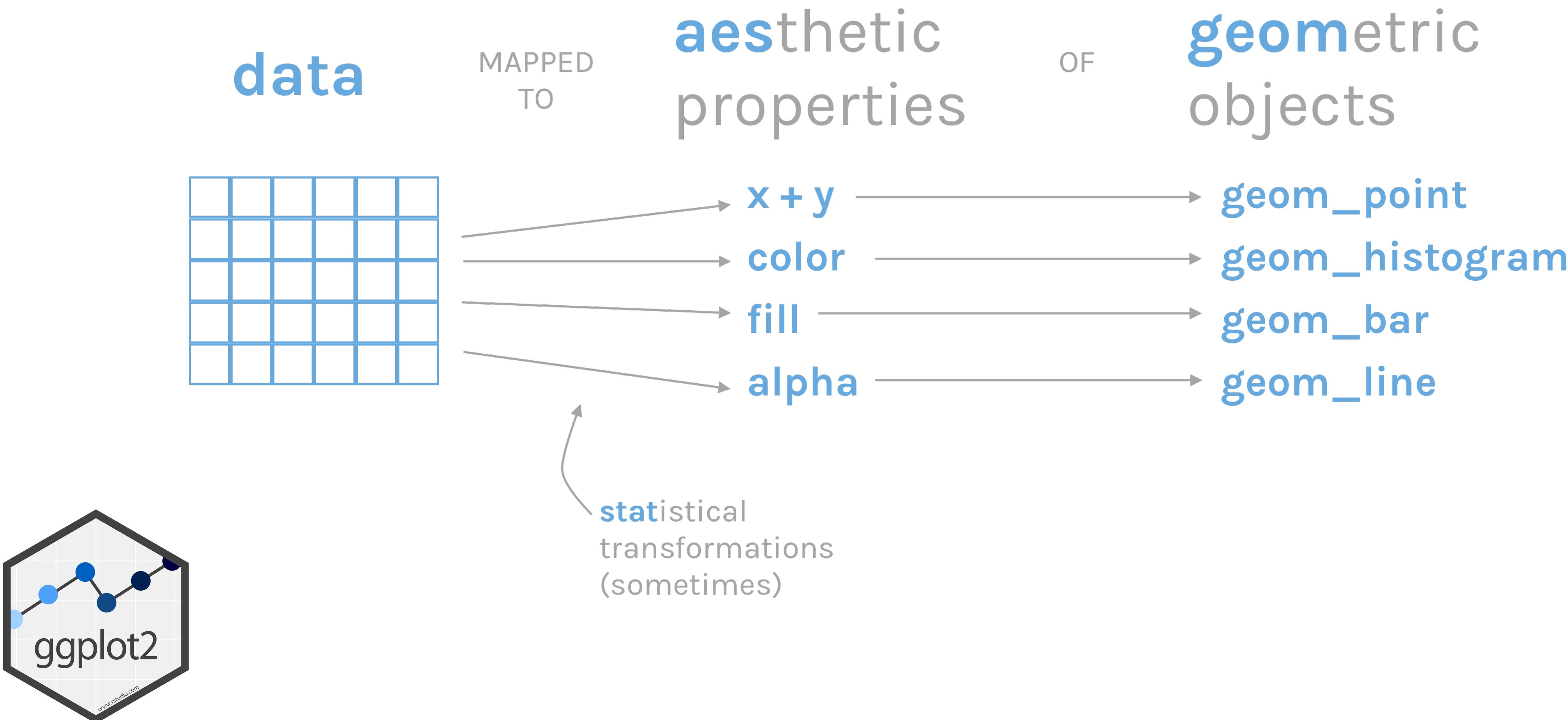
ggplot2 version



```
ggplot() +  
  geom_path(data = troops, aes(x = long, y = lat, group = group,  
                               color = direction, size = survivors),  
            lineend = "round") +  
  geom_point(data = cities, aes(x = long, y = lat),  
             color = "#DC5B44") +  
  geom_text_repel(data = cities, aes(x = long, y = lat, label = city),  
                 color = "#DC5B44", family = "Open Sans Condensed Bold") +  
  scale_size(range = c(0.5, 15)) +  
  scale_colour_manual(values = c("#DFC17E", "#252523")) +  
  labs(x = NULL, y = NULL) +  
  guides(color = FALSE, size = FALSE)
```



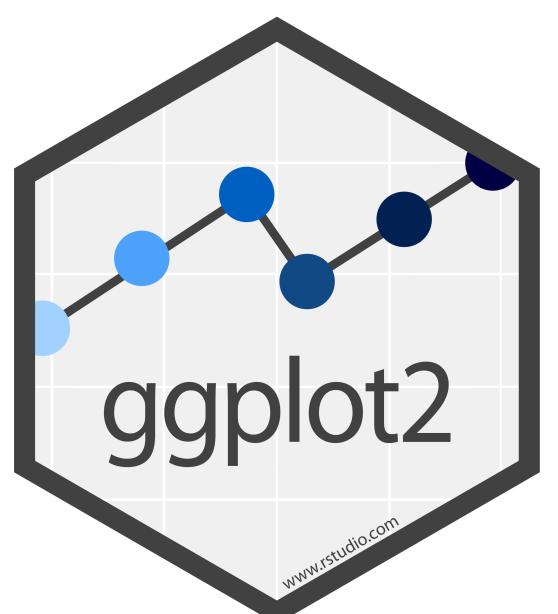
A layered grammar of graphics



A layered grammar of graphics

```
ggplot (data = <DATA>) +  
  <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>),
```

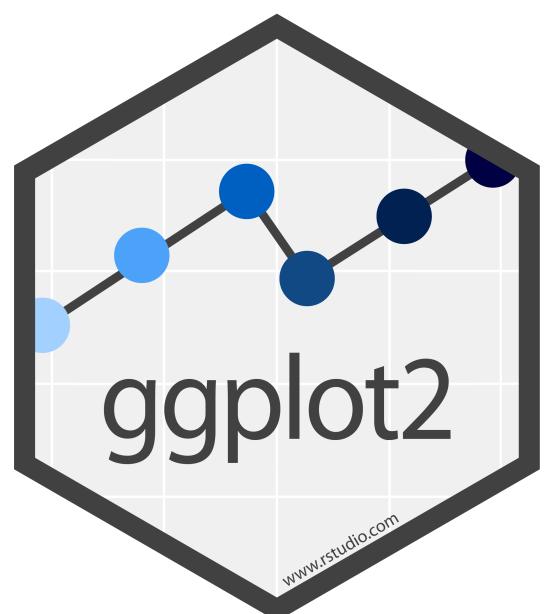
required



A layered grammar of graphics

```
ggplot (data = <DATA>) +  
  <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>),  
    stat = <STAT>, position = <POSITION>) +  
  <COORDINATE_FUNCTION> +  
  <FACET_FUNCTION> +  
  <SCALE_FUNCTION> +  
  <THEME_FUNCTION>
```

] required
] Not required, sensible defaults supplied

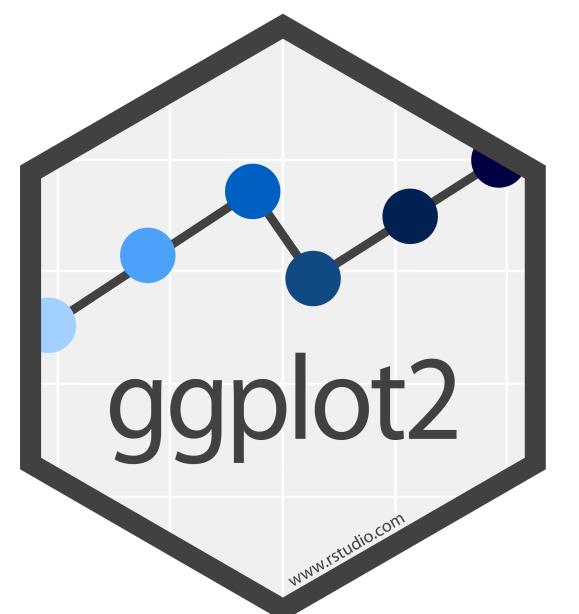


Aesthetics

Geometric shapes

Scales

Themes

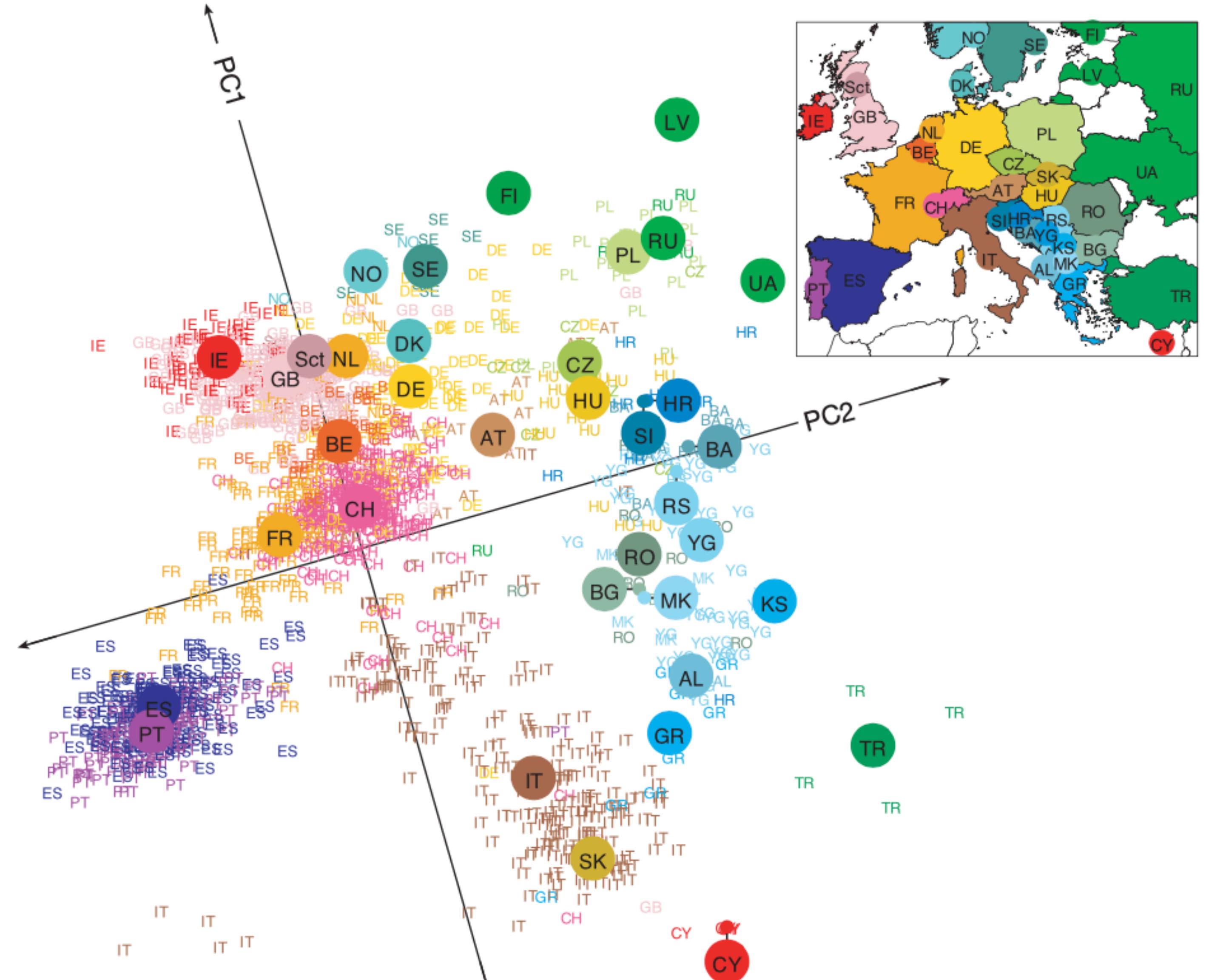
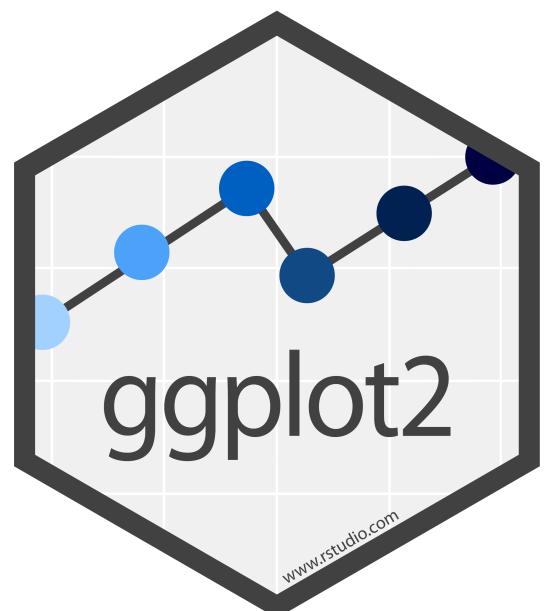


Aesthetics

Geometric shapes

Scales

Themes



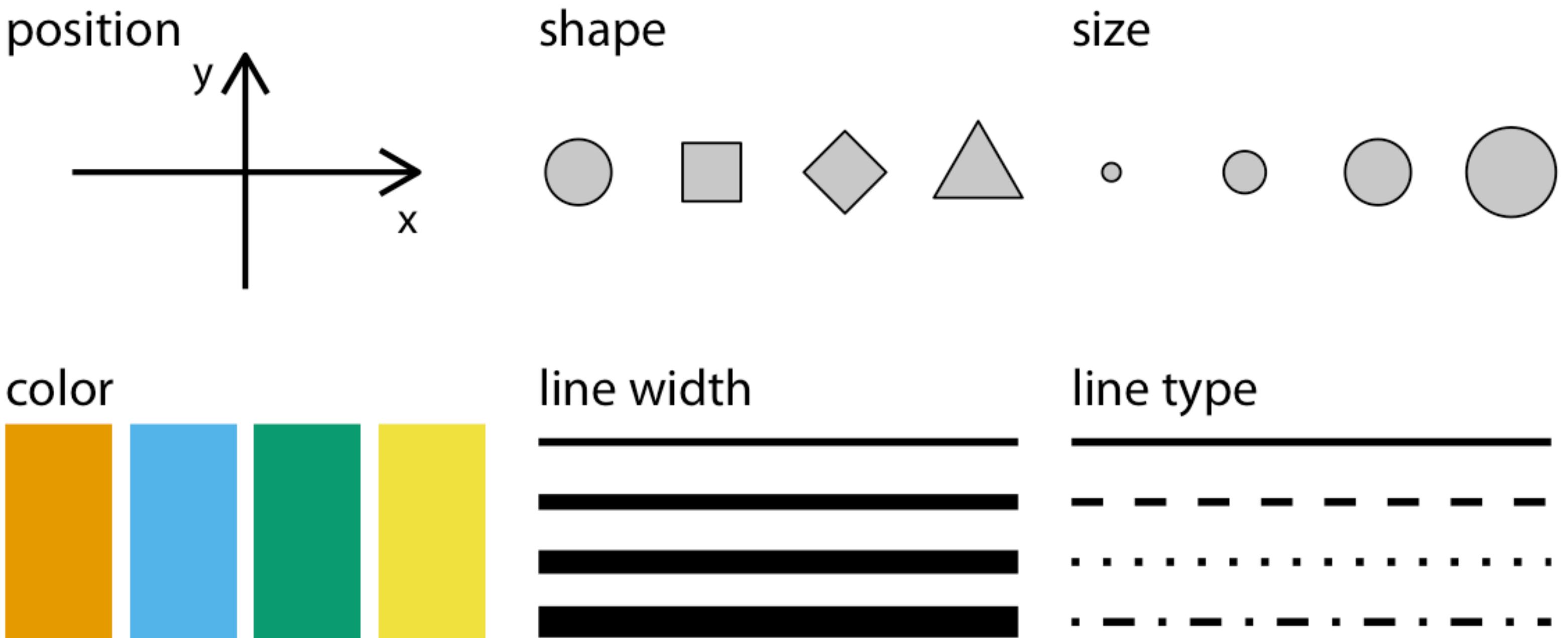
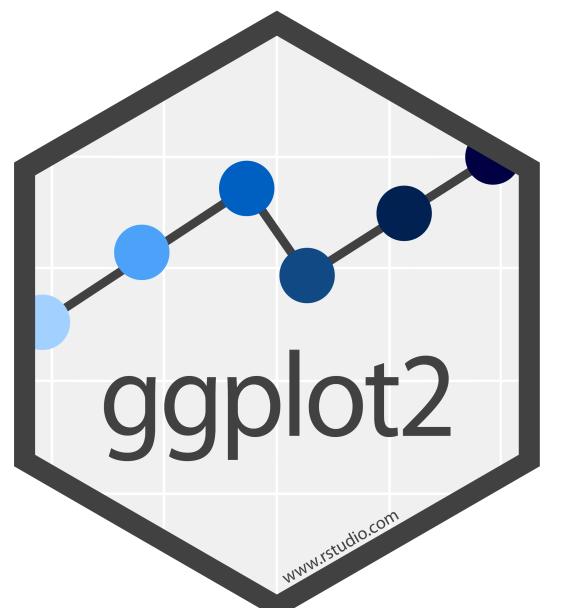
Which aesthetic attributes are used here to visualise attributes of the dataset?

Aesthetics

Geometric shapes

Scales

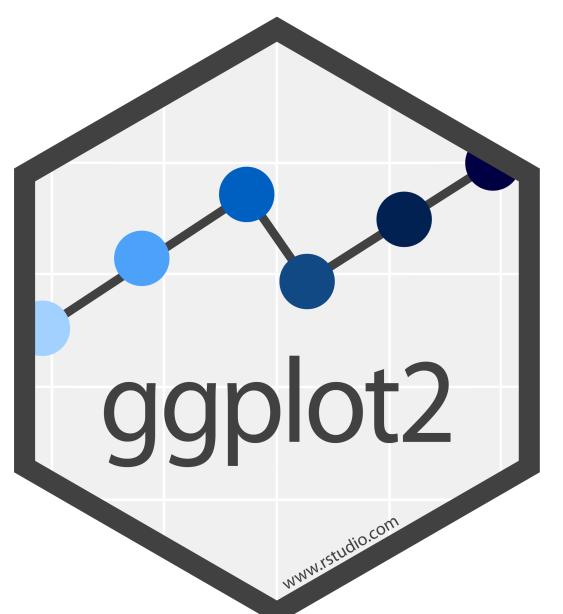
Themes



Commonly used aesthetics

Aesthetics

Geometric shapes



Aes Common aesthetic values.

color and **fill** - string ("red", "#RRGGBB")

linetype - integer or string (0 = "blank", 1 = "solid",
2 = "dashed", 3 = "dotted", 4 = "dotdash", 5 = "longdash",
6 = "twodash")

lineend - string ("round", "butt", or "square")

linejoin - string ("round", "mitre", or "bevel")

size - integer (line width in mm)

shape - integer/shape name or
a single character ("a")

0 1 2 3 4 5 6 7 8 9 10 11 12
□ ○ △ + × ◊ ▽ □ × * ◇ ⊕ ⧺ ▨
13 14 15 16 17 18 19 20 21 22 23 24 25
⊗ □ ○ △ ◊ ○ ○ ○ □ ◇ △ ▽

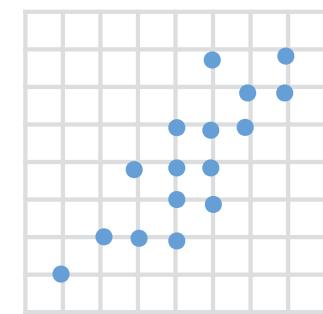
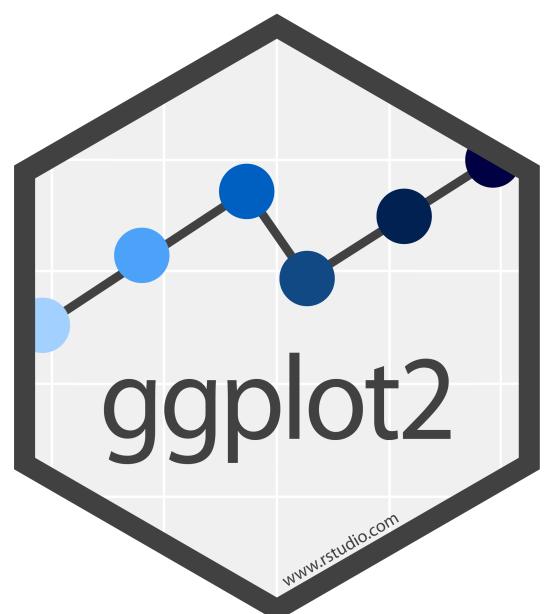
Values for commonly used aesthetics

Aesthetics

Geometric shapes

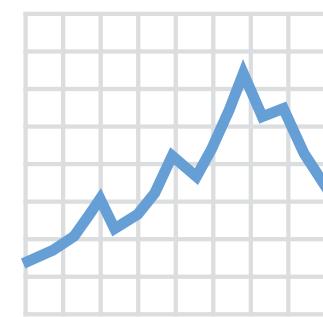
Scales

Themes



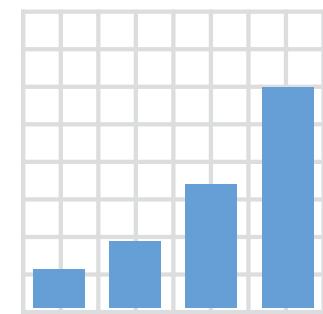
e + geom_point()

x, y, alpha, color, fill, shape, size, stroke



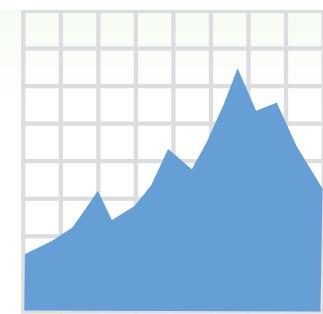
i + geom_line()

x, y, alpha, color, group, linetype, size



f + geom_col()

x, y, alpha, color, fill, group, linetype, size



c + geom_area(stat = "bin")

x, y, alpha, color, fill, linetype, size

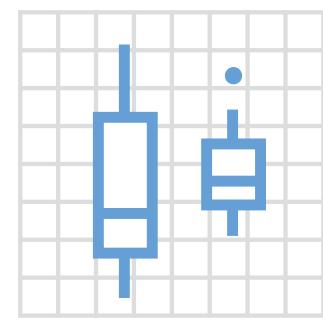
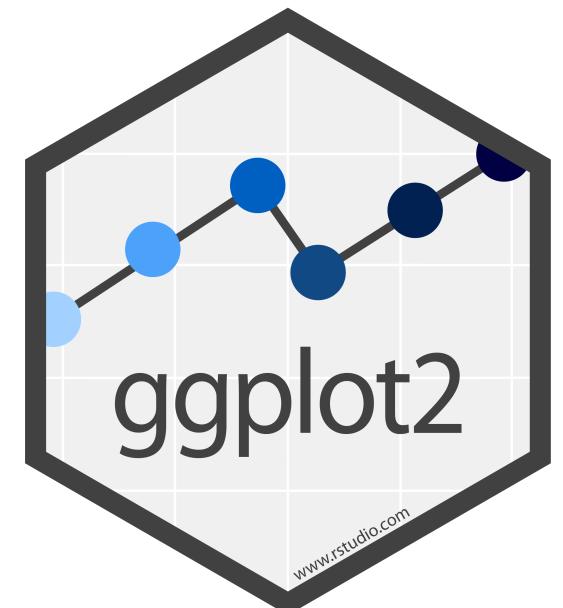
Individual geoms - one distinct object per data point

Aesthetics

Geometric shapes

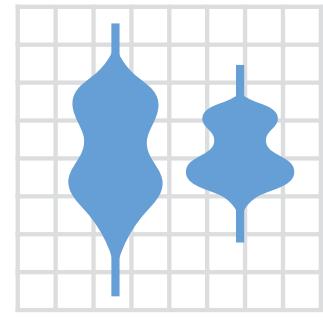
Scales

Themes



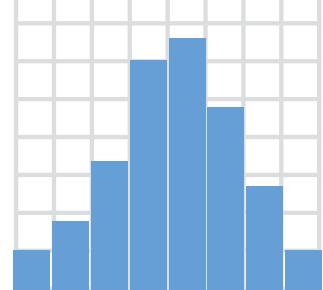
f + geom_boxplot()

x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight



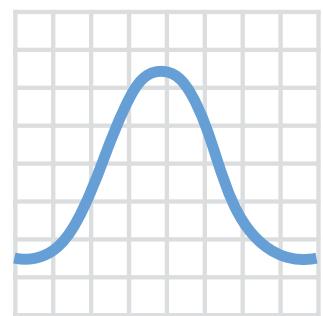
f + geom_violin(scale = "area")

x, y, alpha, color, fill, group, linetype, size, weight



c + geom_histogram(binwidth = 5)

x, y, alpha, color, fill, linetype, size, weight



c + geom_density(kernel = "gaussian")

x, y, alpha, color, fill, group, linetype, size, weight

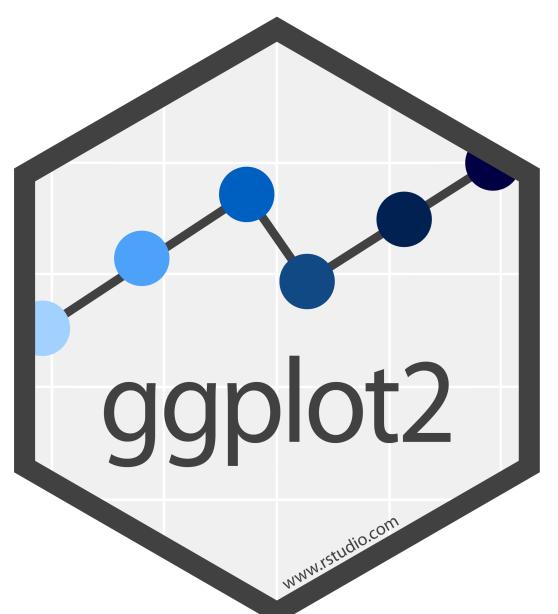
Collective geoms - one distinct object for multiple data points

Aesthetics

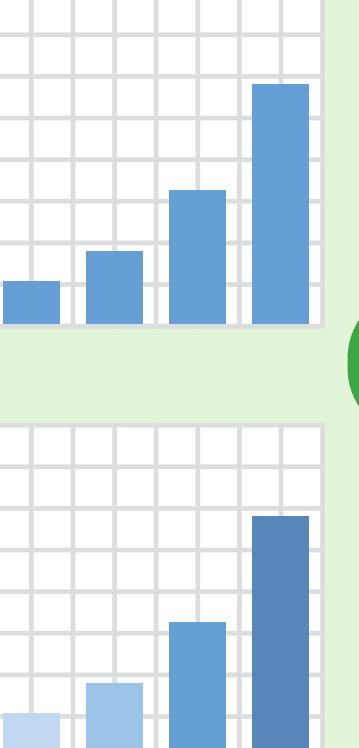
Geometric shapes

Scales

Themes



Scales map data values to the visual values of an aesthetic. To change a mapping, add a new scale.



```
n <- d + geom_bar(aes(fill = fl))  
n + scale_fill_manual(  
  values = c("skyblue", "royalblue", "blue", "navy"),  
  limits = c("d", "e", "p", "r"), breaks =c("d", "e", "p", "r"),  
  name = "fuel", labels = c("D", "E", "P", "R"))
```

scale_ aesthetic to adjust prepackaged scale to use scale-specific arguments

range of values to include in mapping title to use in legend/axis labels to use in legend/axis breaks to use in legend/axis

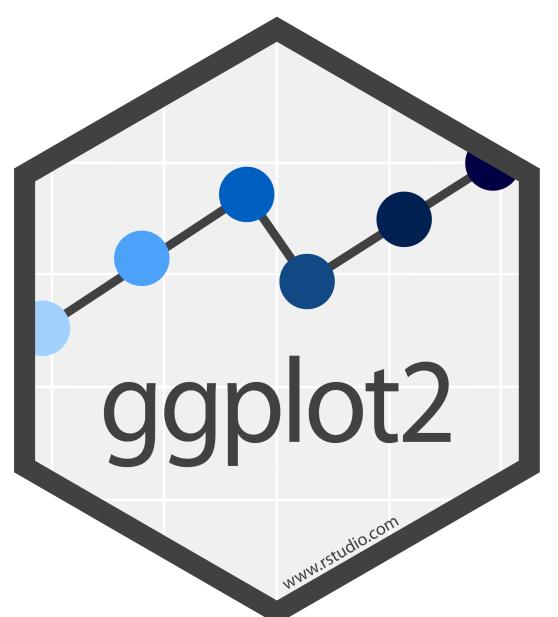
Example of manually setting a scale for aesthetic `fill`

Aesthetics

Geometric shapes

Scales

Themes



scale_*_continuous() - Map cont' values to visual ones.
scale_*_discrete() - Map discrete values to visual ones.
scale_*_binned() - Map continuous values to discrete bins.
scale_*_identity() - Use data values as visual ones.
scale_*_manual(values = c()) - Map discrete values to manually chosen visual ones.

scale_x_log10() - Plot x on log10 scale.
scale_x_reverse() - Reverse the direction of the x axis.
scale_x_sqrt() - Plot x on square root scale.

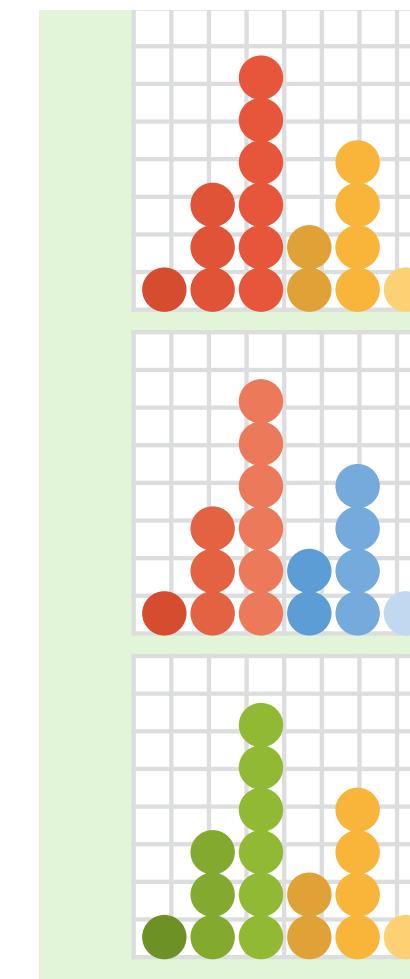
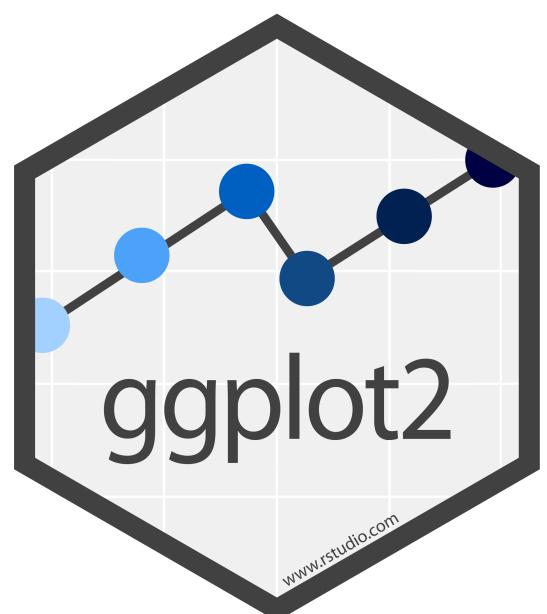
Examples of different scales

Aesthetics

Geometric shapes

Scales

Themes

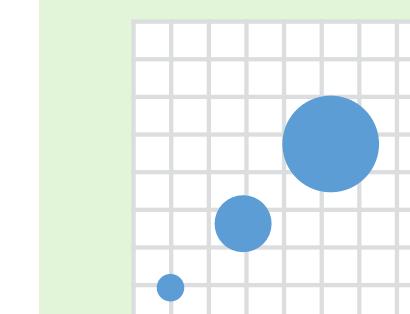
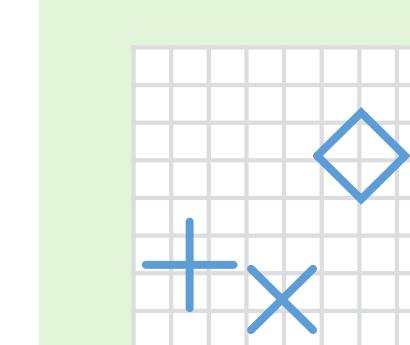


`o + scale_fill_gradient(low="red", high="yellow")`

`o + scale_fill_gradient2(low = "red", high = "blue", mid = "white", midpoint = 25)`

`o + scale_fill_gradientn(colors = topo.colors(6))`

Also: `rainbow()`, `heat.colors()`, `terrain.colors()`,
`cm.colors()`, `RColorBrewer::brewer.pal()`



`p + scale_shape() + scale_size()`

`p + scale_shape_manual(values = c(3:7))`

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25

□○△+×◊▽⊗✳⊕⊛⊛田⊗⊗□○△◊○○●□◆△▽

`p + scale_radius(range = c(1,6))`

`p + scale_size_area(max_size = 6)`

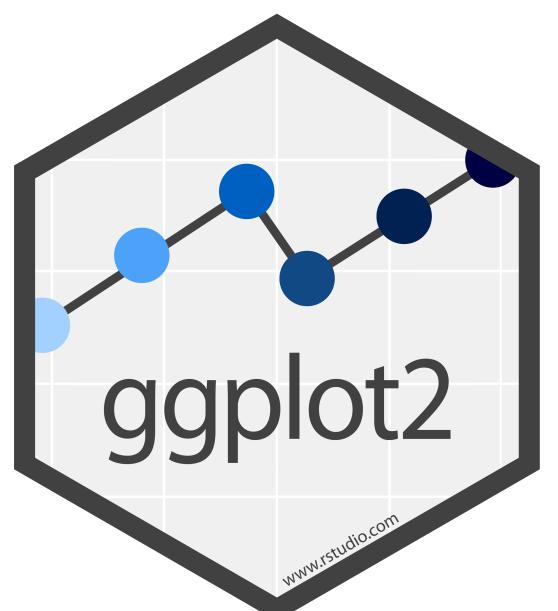
Examples of different scales

Aesthetics

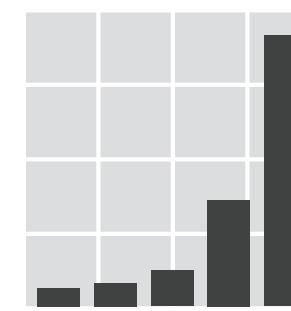
Geometric shapes

Scales

Themes



r + theme() Customize aspects of the theme such as axis, legend, panel, and facet properties.



r + theme_gray()
Grey background
(default theme).

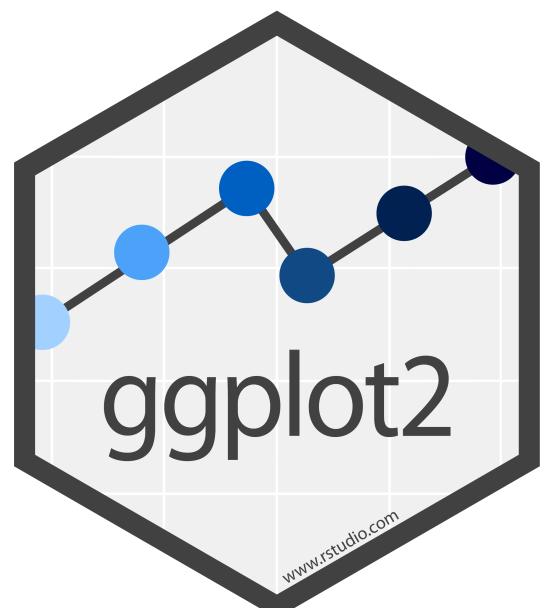
Themes control the non-data elements of a plot

Aesthetics

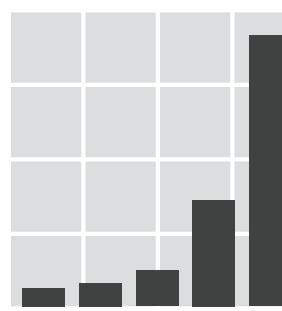
Geometric shapes

Scales

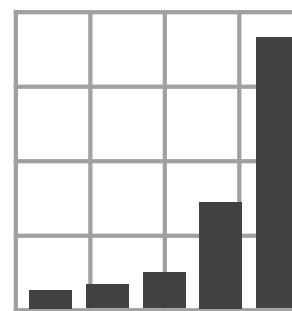
Themes



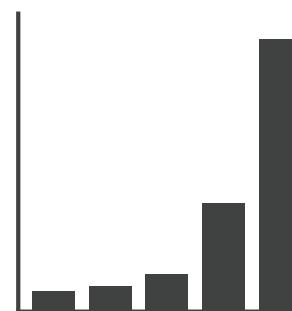
r + theme() Customize aspects of the theme such as axis, legend, panel, and facet properties.



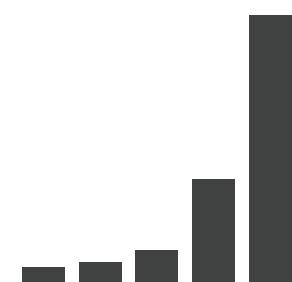
r + theme_gray()
Grey background
(default theme).



r + theme_bw()
White background
with grid lines.



r + theme_classic()

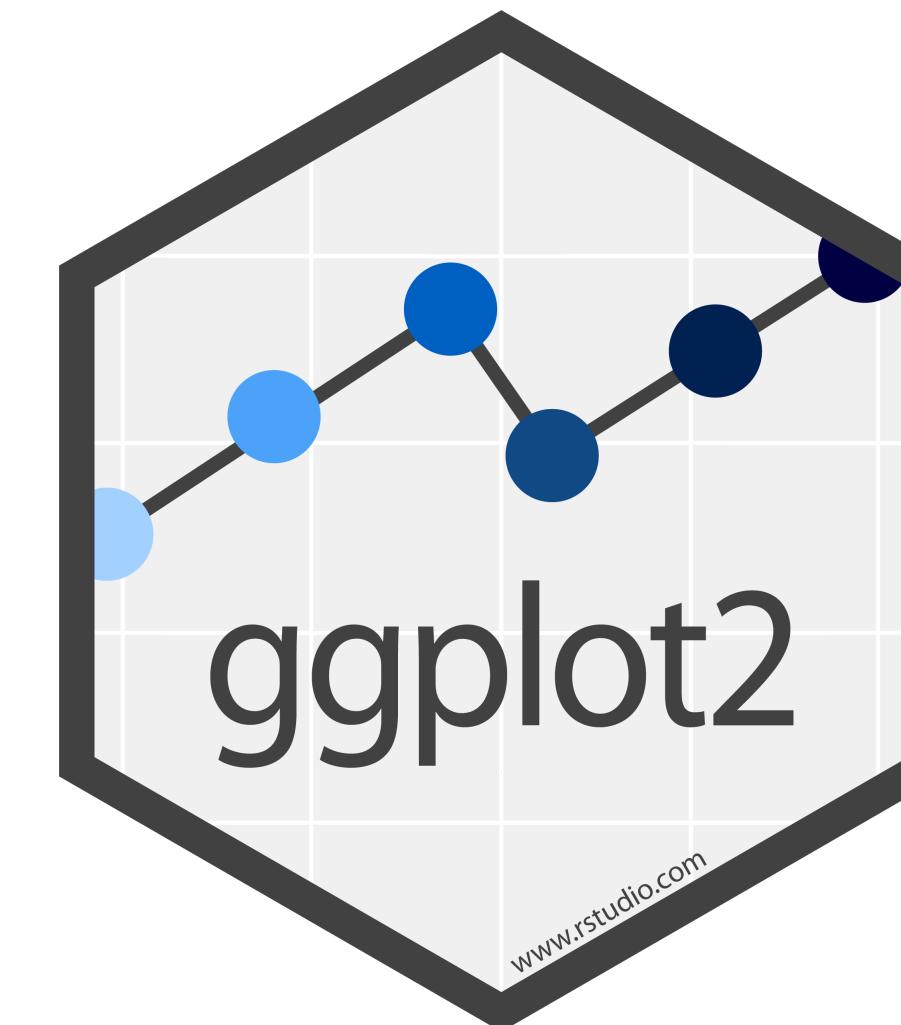


r + theme_void()
Empty theme.

Examples of available pre-defined themes

Programme for today

- Introductory lecture
- Hands-on dataviz tutorial
- Individual exercises



At the end of the day you should be familiar with making visualisations in R using the ggplot2 package