

Data Workshop #5

Visualisation

dataworkshop.eu

Alekseichenko Vladimir

DataWorkshop.eu

Data Workshop

Intro Goal Approach Prerequisite Success metric How to join?

**Talk is cheap. Show me
the data!**

Matters is only ready-made solution with actionable insights. The rest is secondary. Practice and learn.



About me



Vladimir Alekseichenko
Love analyze data



Architect

slon1024

slon1024

vova@vova.me

Disclaimer

Data Workshop [*all time*] focuses on the **intuition** and **practical** tips.

For a formal treatment, see something else^{}.*

^{*} papers or classical machine learning books

Environment

[github.com/dataworkshop/**prerequisite**](https://github.com/dataworkshop/prerequisite)
[github.com/dataworkshop/**environment**](https://github.com/dataworkshop/environment)

[github.com/dataworkshop/**visualization**](https://github.com/dataworkshop/visualization)

Packages

github.com/datworkshop/prerequisite

```
$ python run.py
seaborn-0.7.0 - OK
xgboost-0.4 - OK
matplotlib-1.5.1 - OK
IPython-4.1.2 - OK
numpy-1.11.0 - OK
pandas-0.18.0 - OK
sklearn-0.17.1 - OK
```

```
=====
All right, you are ready to go on Data Workshop!
```

```
$ python run.py
seaborn-0.6 should be upgraded to seaborn-0.7
xgboost-0.4 - OK
matplotlib-1.5.1 - OK
IPython-4.1.2 - OK
numpy-1.11.0 - OK
pandas-0.18.0 - OK
sklearn-0.17.1 - OK
```

```
=====
RECOMENDATION (without upgrade some needed features could be missing)
pip install --upgrade seaborn
```

```
$ python run.py
seaborn-0.7.0 - OK
xgboost - missing
matplotlib-1.5.1 - OK
IPython-4.1.2 - OK
numpy-1.11.0 - OK
pandas-0.18.0 - OK
sklearn-0.17.1 - OK
```

```
=====
REQUIRED
Please install those packages before Data Workshop: xgboost
pip install xgboost
More info how to install xgboost: http://xgboost.readthedocs.org/en/latest/build.html
```

jupyter notebook



```
$ jupyter notebook
[I 22:17:17.650 NotebookApp] The port 8888 is already in use, trying another random port.
[I 22:17:17.650 NotebookApp] The port 8889 is already in use, trying another random port.
[I 22:17:17.651 NotebookApp] The port 8890 is already in use, trying another random port.
[I 22:17:17.651 NotebookApp] The port 8891 is already in use, trying another random port.
[I 22:17:17.657 NotebookApp] Serving notebooks from local directory: /Users/vova/src/github/dataworkshop/titanic/vladimir/tmp
[I 22:17:17.657 NotebookApp] 0 active kernels
[I 22:17:17.657 NotebookApp] The IPython Notebook is running at: http://localhost:8892/
[I 22:17:17.657 NotebookApp] Use Control-C to stop this server and shut down all kernels (twice to skip confirmation).
```

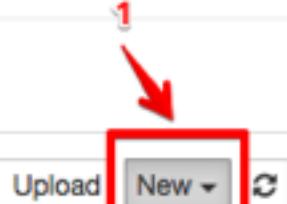


Files Running Clusters

Select items to perform actions on them.



Notebook list empty.



Text File

Folder

Terminal

Notebooks

Haskell

Julia 0.3.8

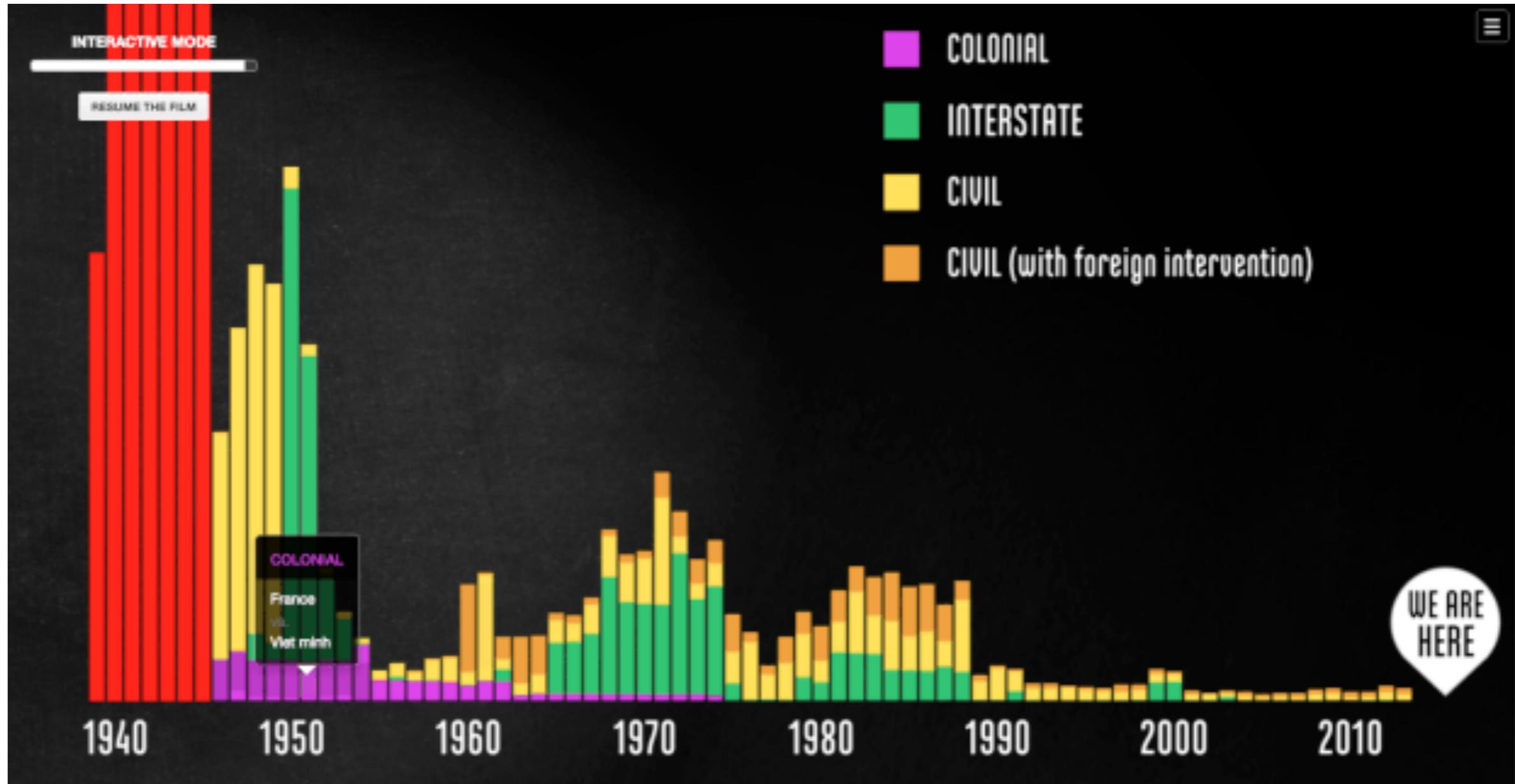
Python 2



Motivation

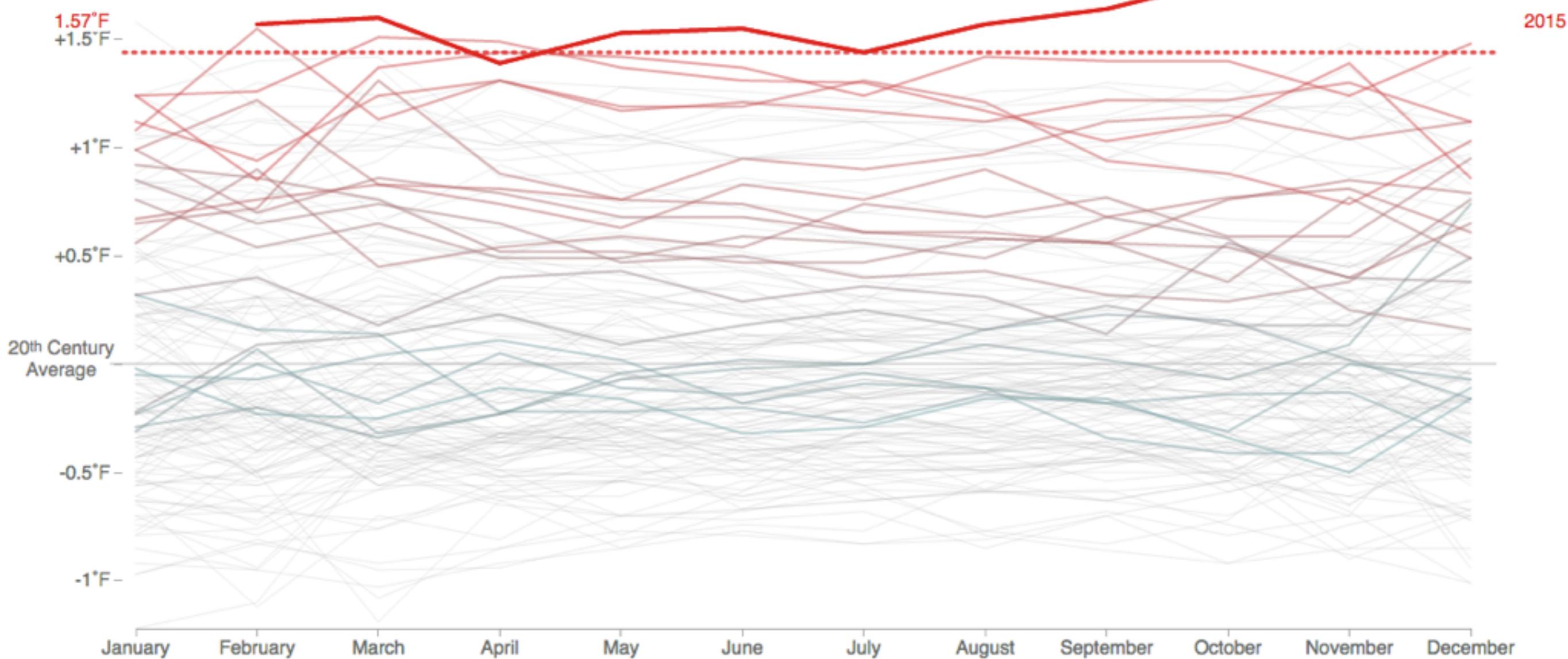
Millions of people died during World War II, but it's difficult to grasp what all the big numbers associated with the war mean. Neil Halloran explains in *The Fallen of World War II*, a hybrid between interactive visualization and documentary.

The Fallen of World War II



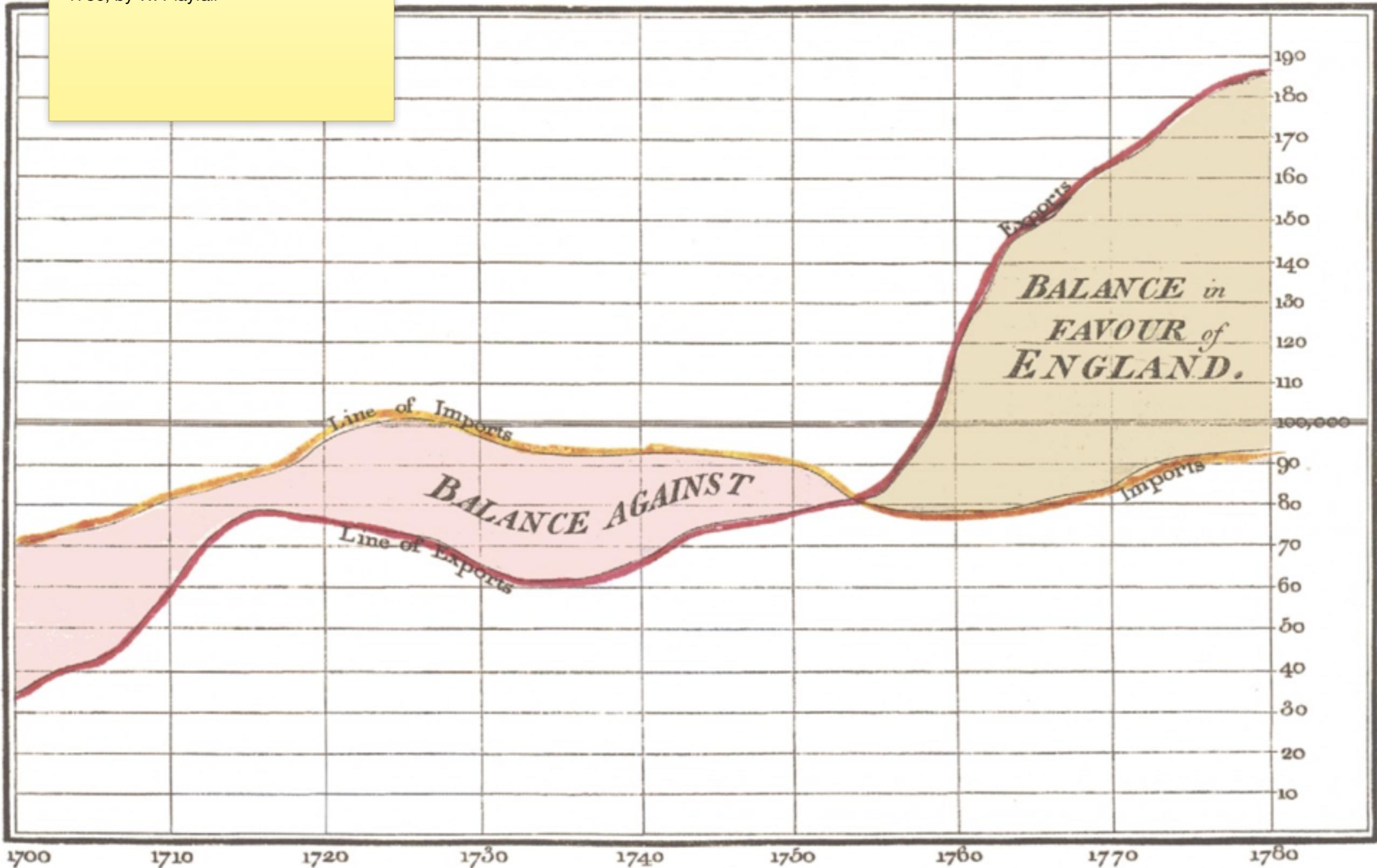
The animated chart by Tom Randall and Blacki Migliozi for Bloomberg shows average monthly temperature. Each line represents a full calendar year, and as you get closer to present-day the line

2015: The Hottest Year



Published as the Act directs, 14th May
1786, by W. Playfair

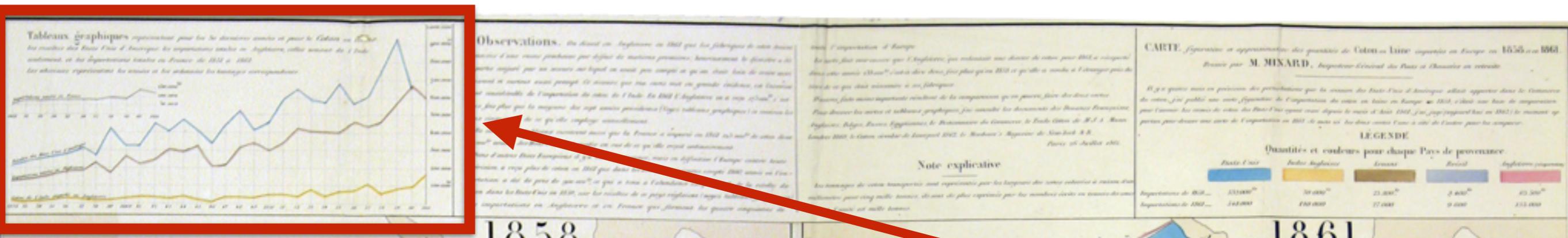
Imports to and from DENMARK & NORWAY from 1700 to 1780.



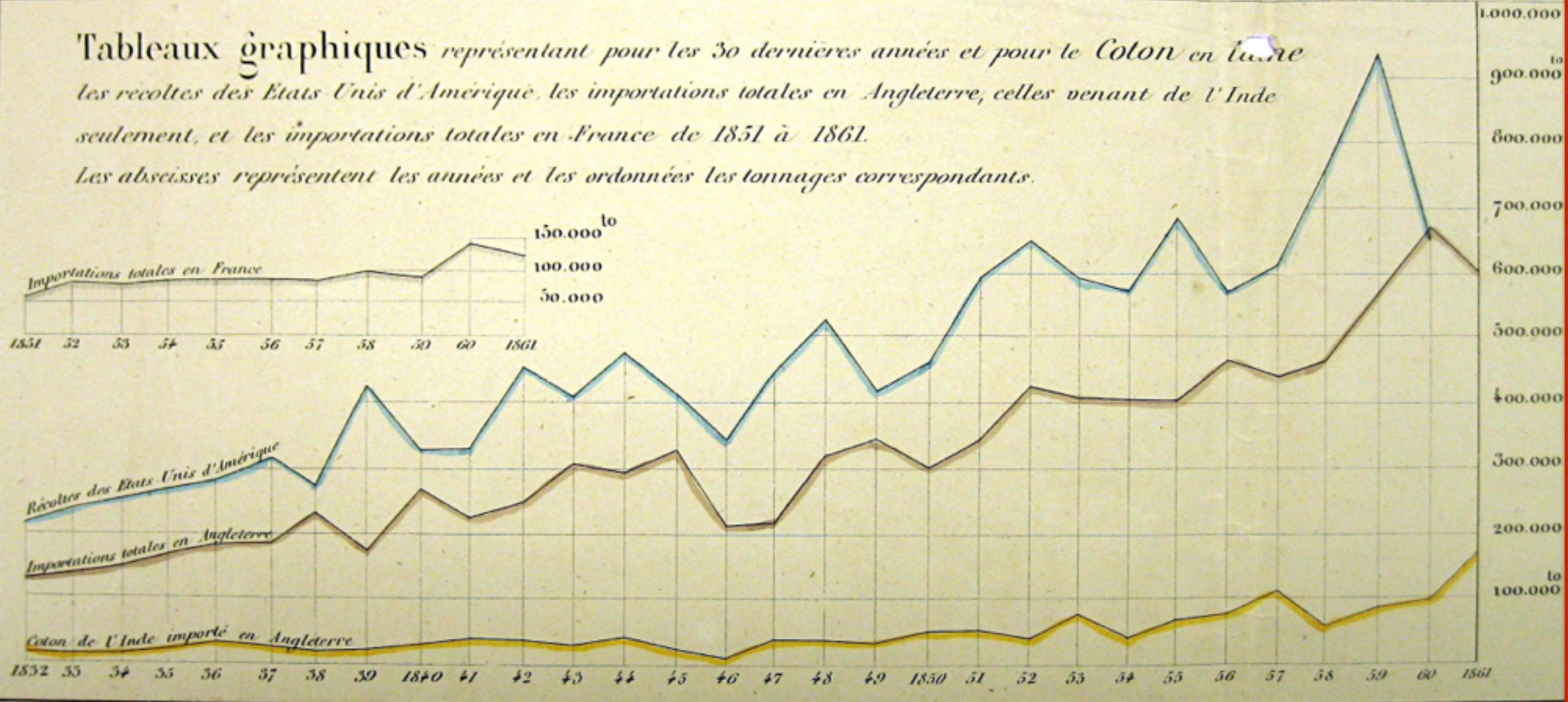
The Bottom line is divided into Years, the Right hand line into £10,000 each.

Published as the Act directs, 14th May 1786. by W^m Playfair

Neale sculpt 352, Strand, London.



Tableaux graphiques représentant pour les 30 dernières années et pour le Coton en laine
les récoltes des Etats-Unis d'Amérique, les importations totales en Angleterre, celles venant de l'Inde
seulement, et les importations totales en France de 1851 à 1861.
Les abscisses représentent les années et les ordonnées les tonnages correspondants.

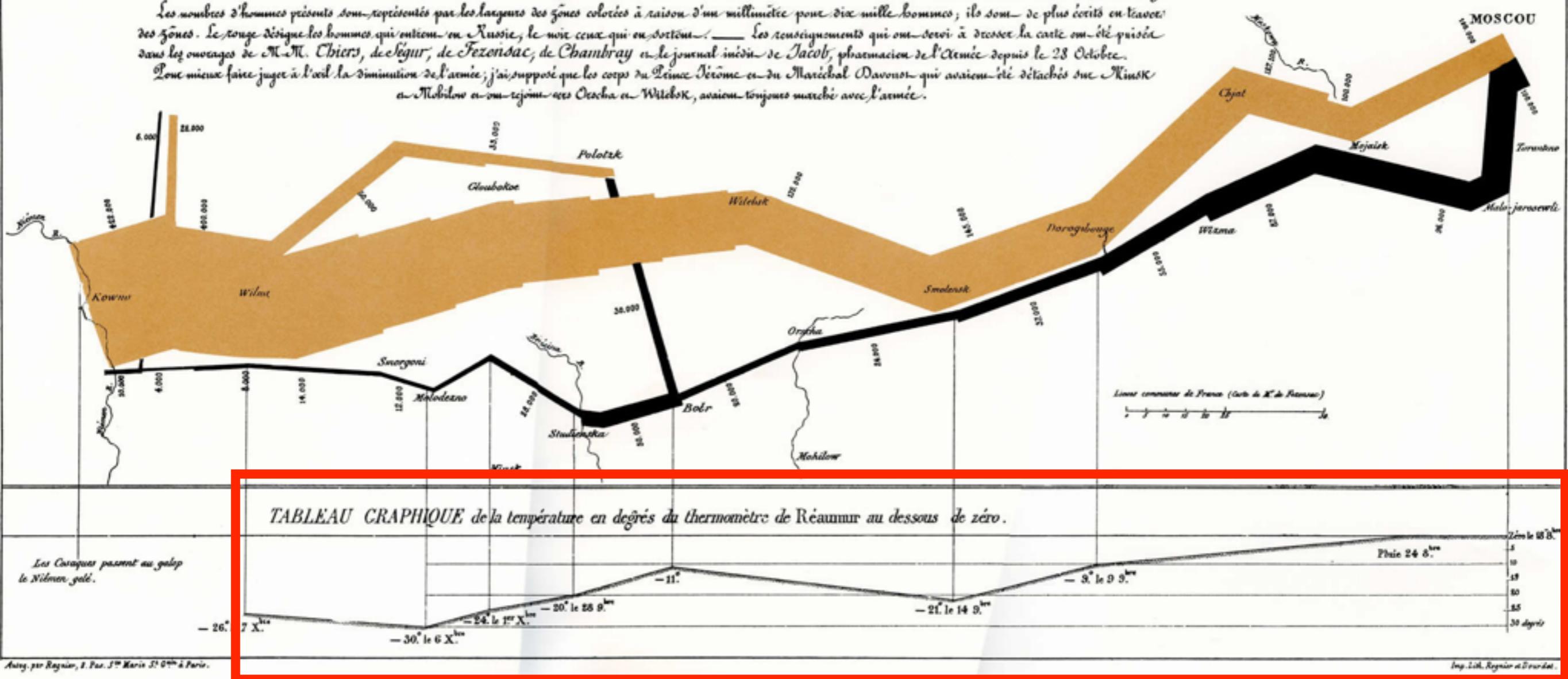


Sankey diagram

Carte Figurative des pertes successives en hommes de l'Armée Française dans la Campagne de Russie 1812-1813.
 Dessinée par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite
 Paris, le 20 Novembre 1869

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

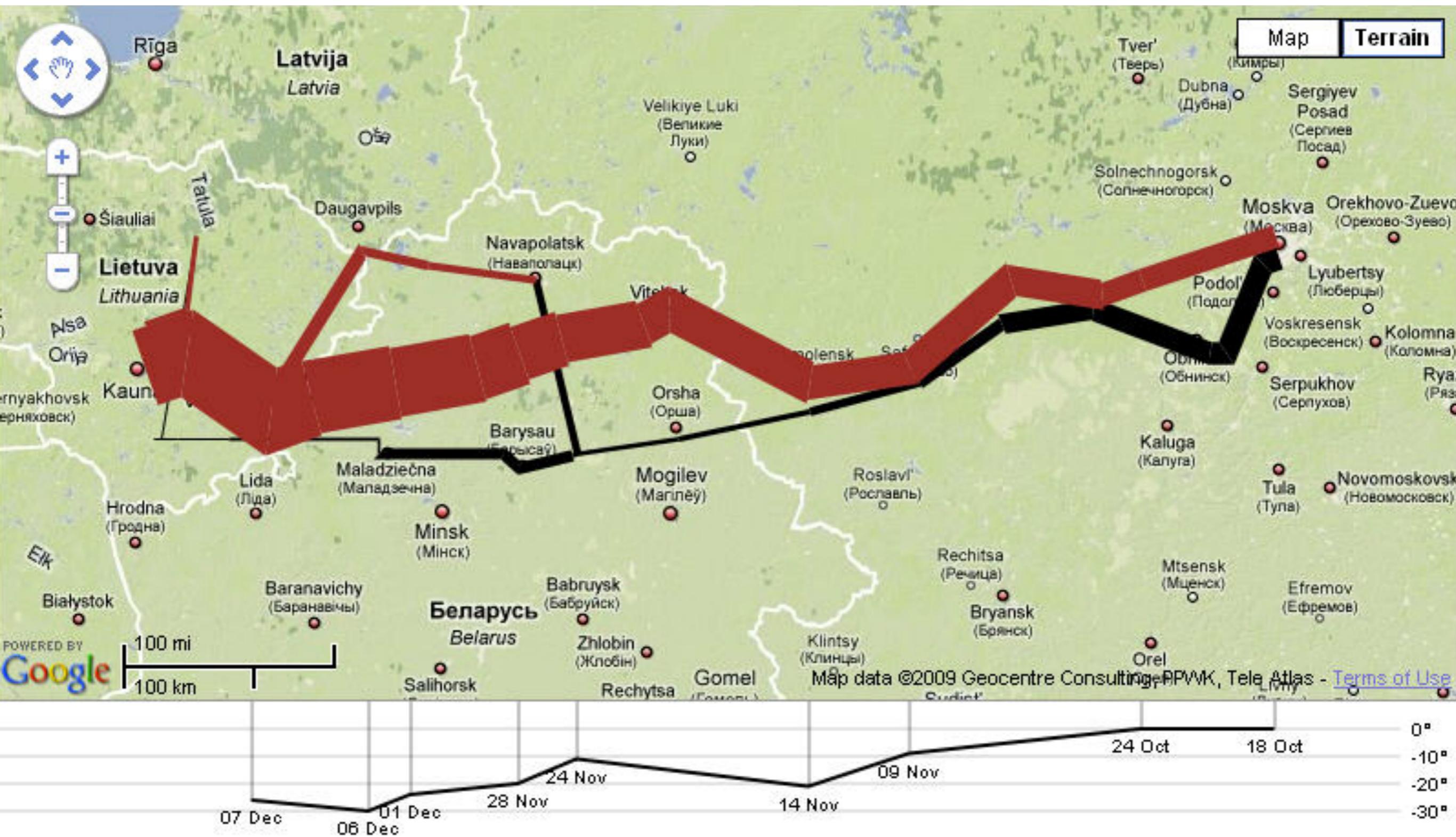
Pour mieux faire juger à l'œil la diminution de l'armée; j'ai supposé que les corps du Prince Jérôme et du Maréchal Davout qui avaient été détachés sur Minsk et Maliblou et qui rejoignirent Ossia et Wileïk, avaient toujours marché avec l'armée.



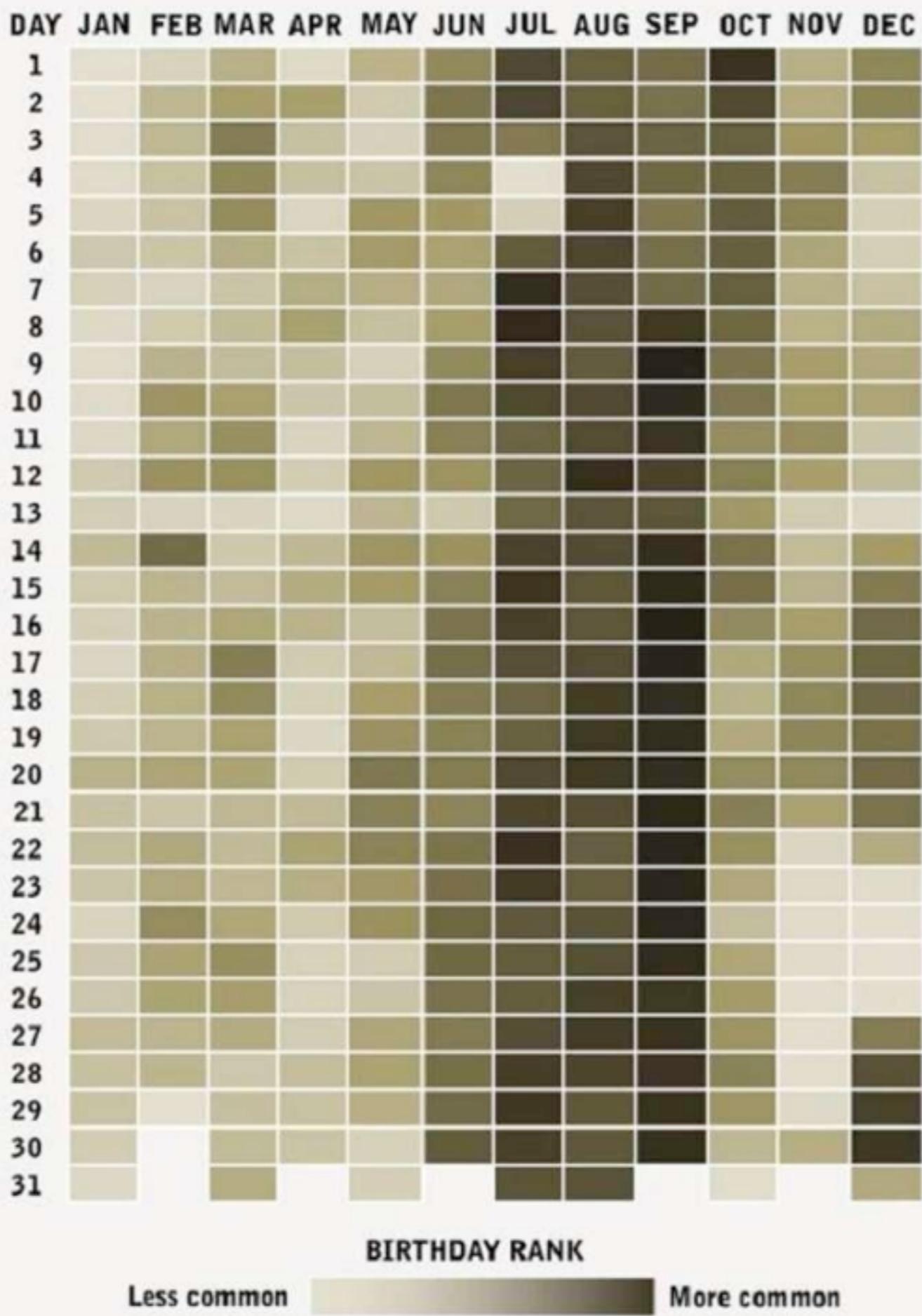
Temperature

https://en.wikipedia.org/wiki/Sankey_diagram

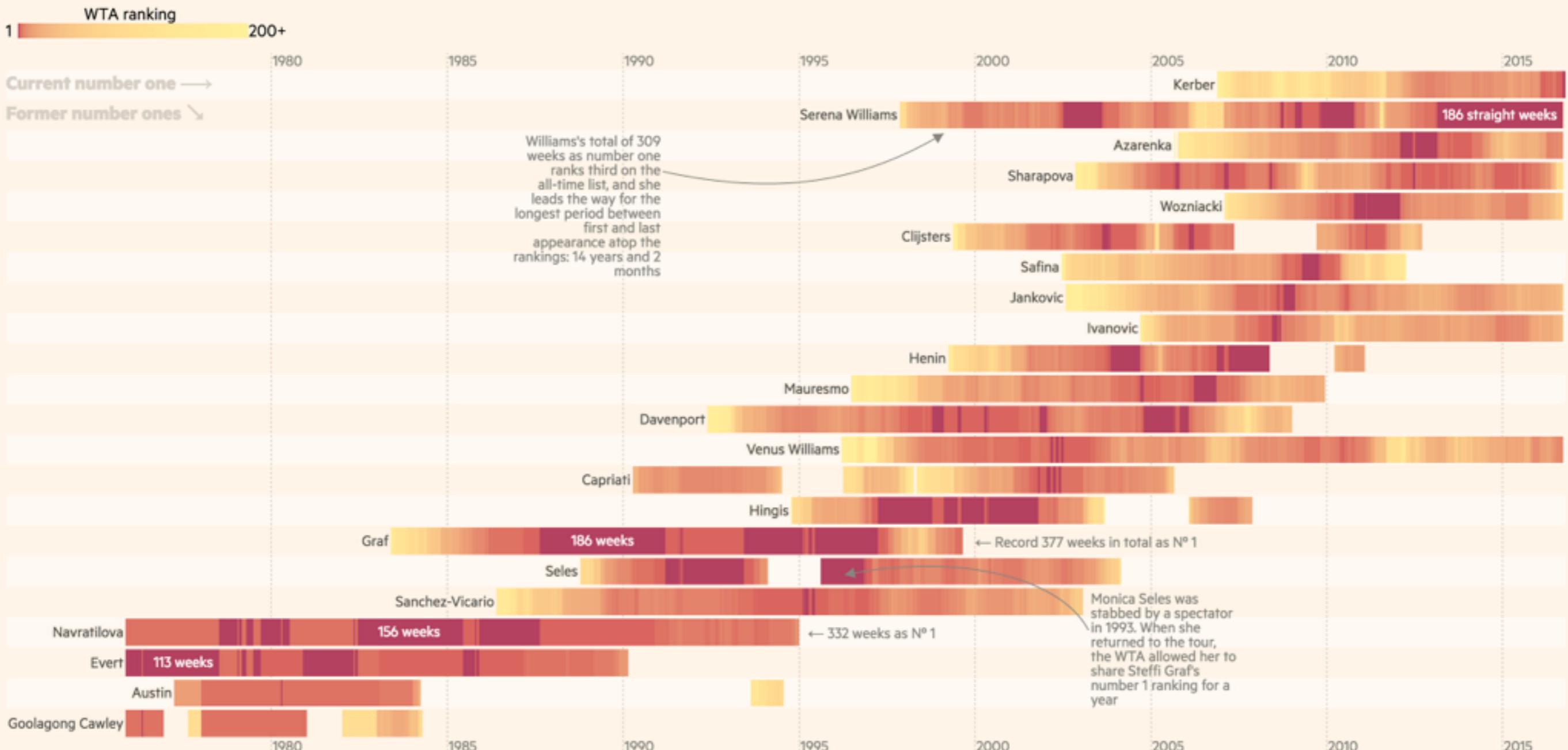
Sankey diagram



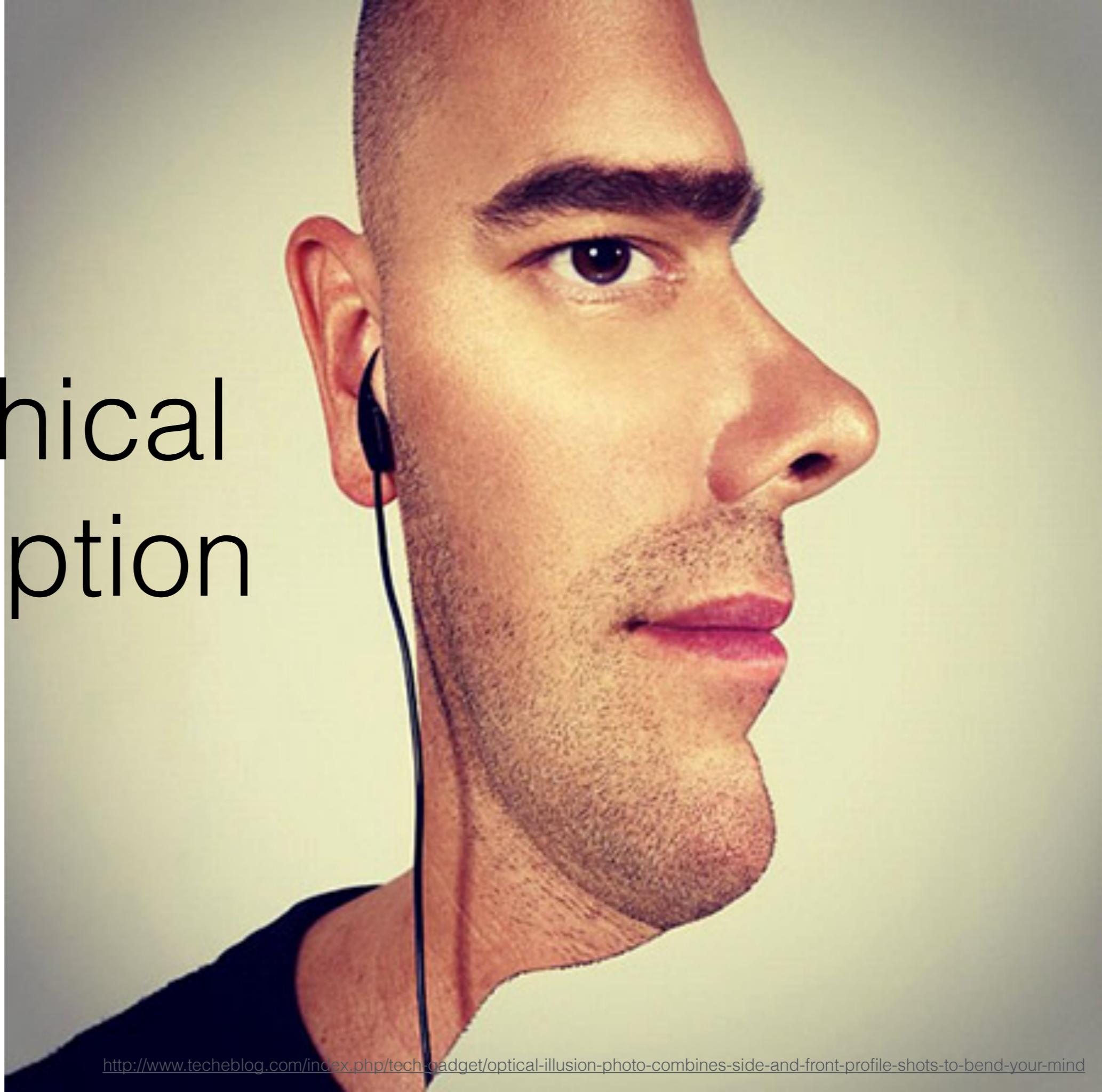
Which Birth Dates Are Most Common?



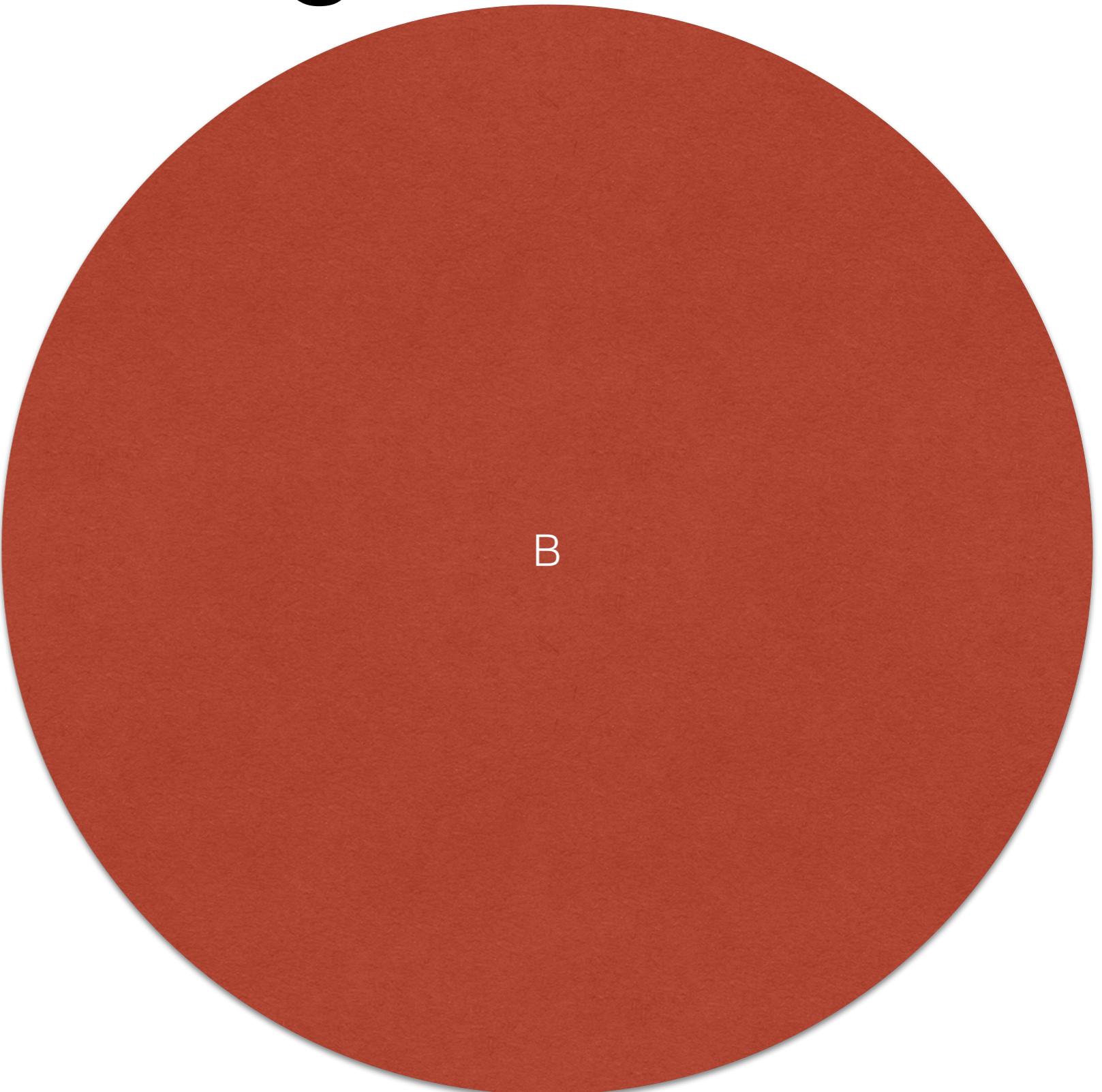
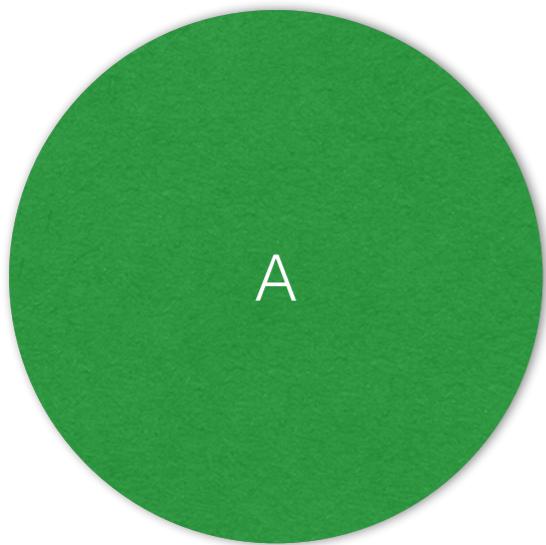
A visual history of women's tennis



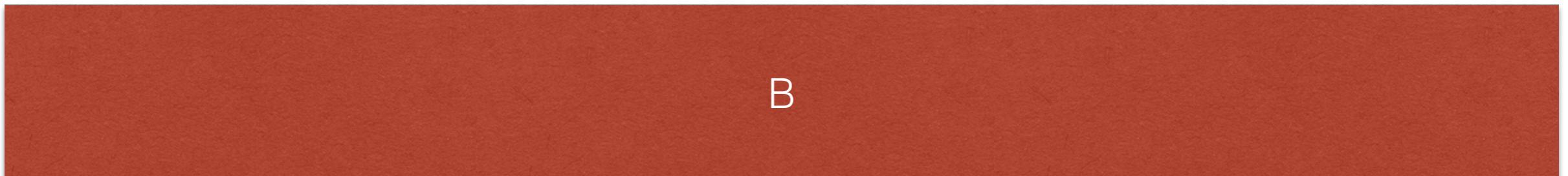
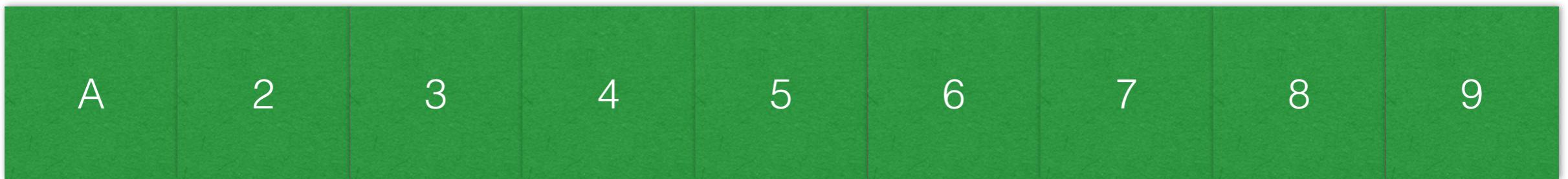
Graphical Perception



How much larger is area?



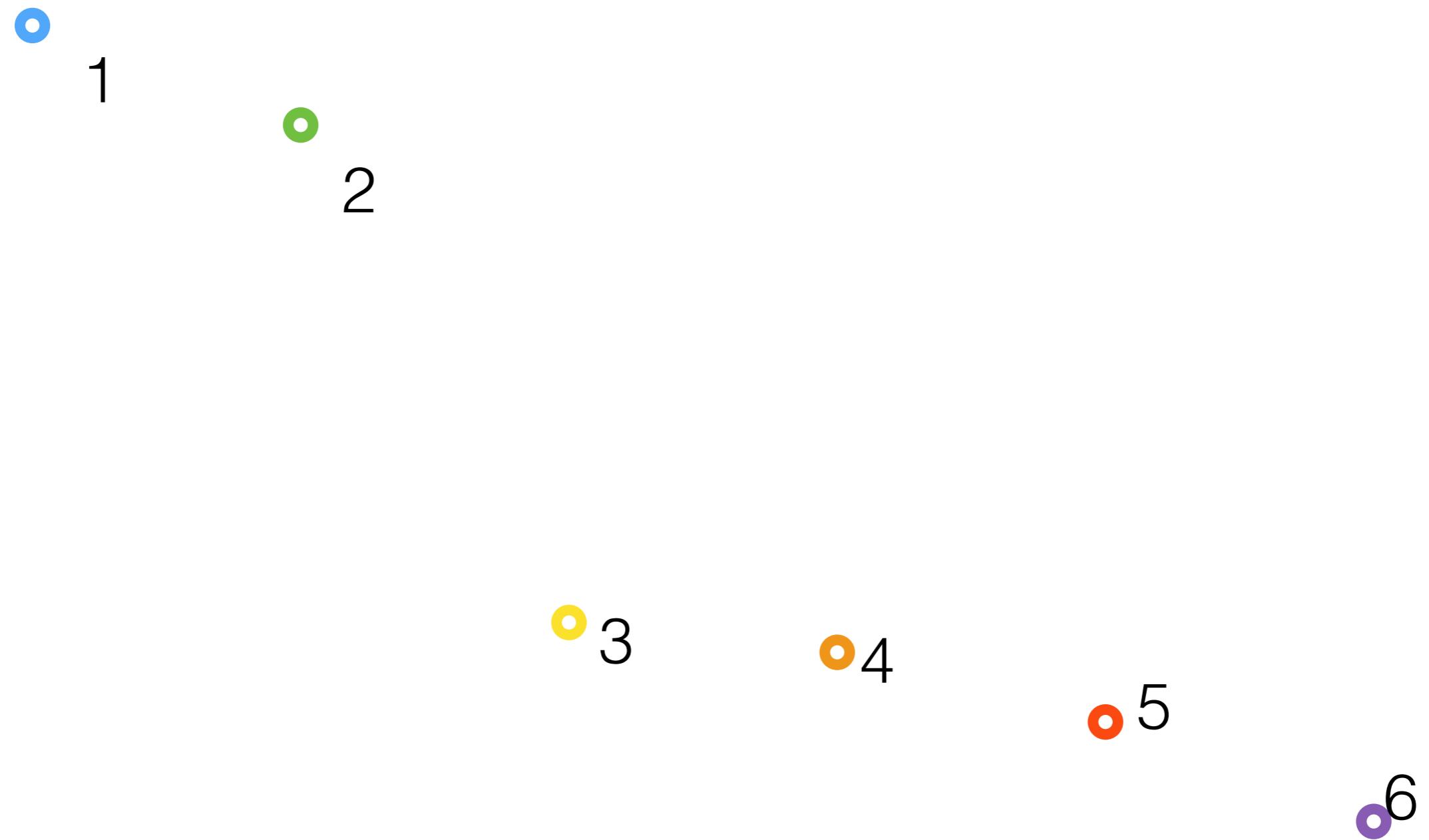
How much larger is area?



Pie Chart



Line Chart



April

May

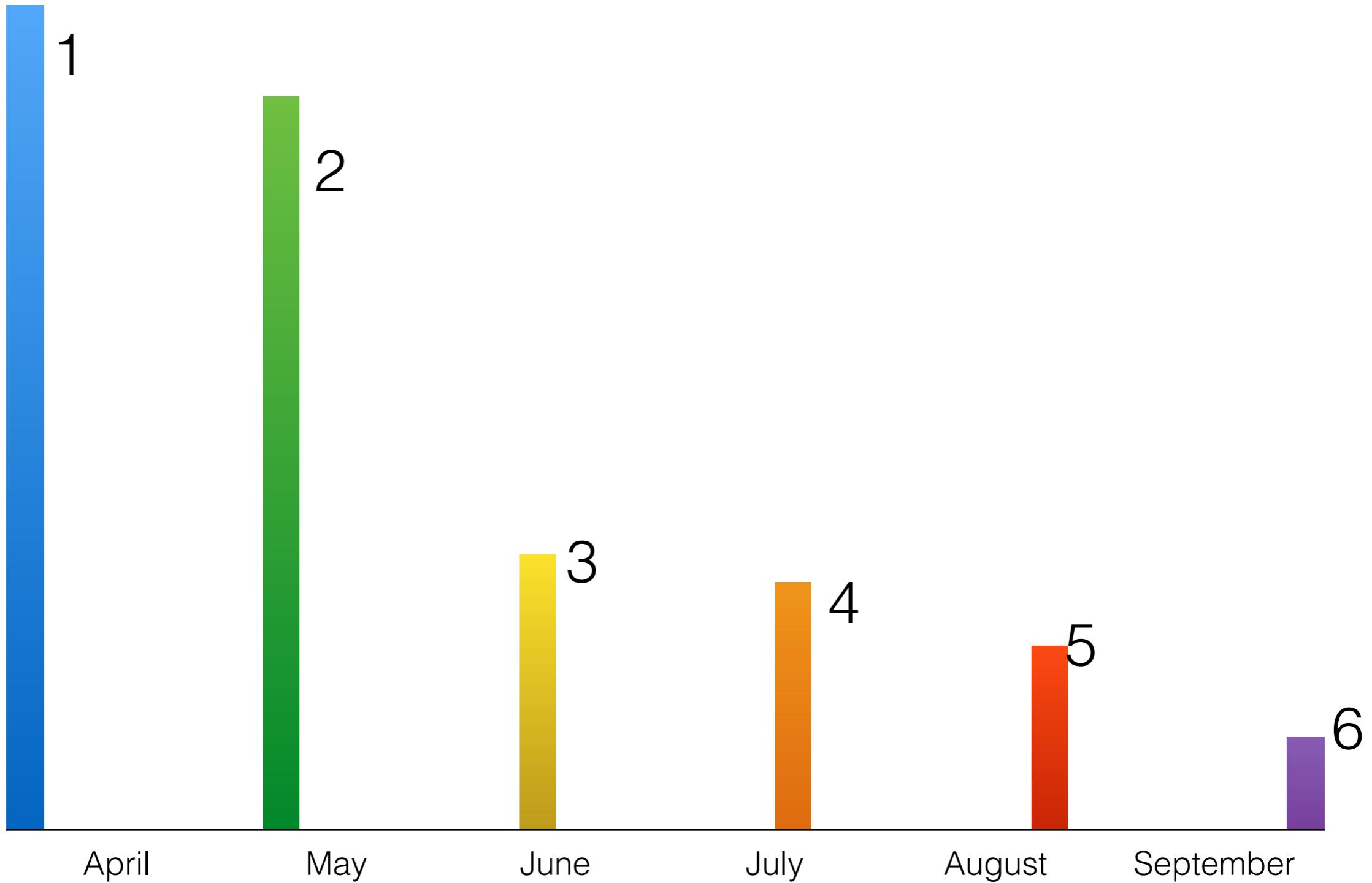
June

July

August

September

Bar Chart

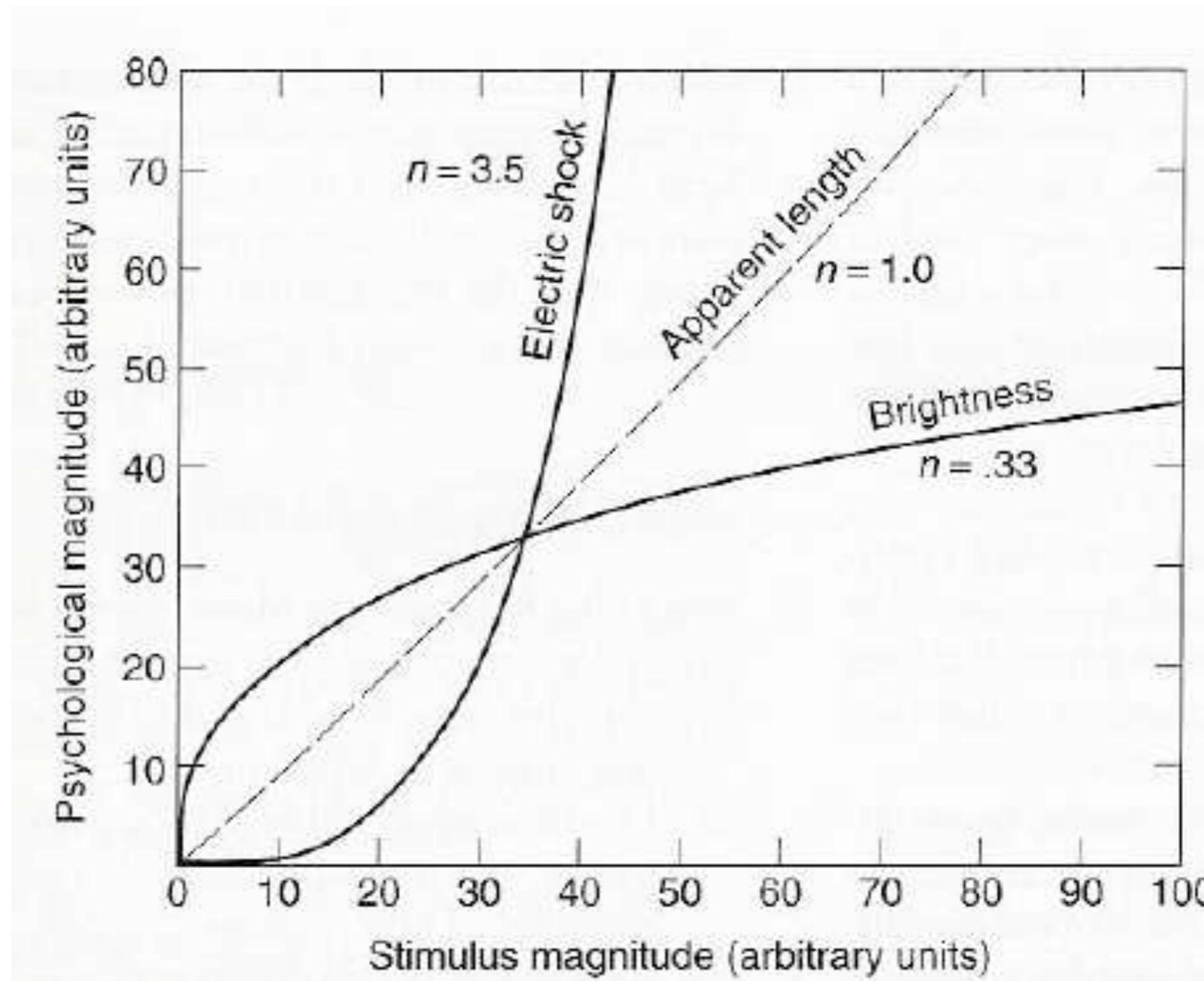


Steven's Power Law

magnitude estimation

$$S = I^n$$

physical intensity



WEBER-FECHNER LAW



“The just-noticeable difference between two stimuli is proportional to the magnitude of the stimuli”

LEVELS OF COMPENSATION.

'GOLDEN PARACHUTE'



'PARACHUTE'



'CHUTE'



Relative Magnitude Estimation

Most accurate



Position (common)
Position (non-aligned)

Length
Slope
Angle

Least accurate

Area
Volume
Color (hue,saturation,value)

Variables

Type of Variables

- Nominal (labels or categories)
- Ordered
- Quantitative
 - Interval
 - Ratio

Type of Variables

Name	Operations	Example
Nominal	<code><>, ==</code>	Apple, Orange, Banana Europe, Asia, America
Ordered	<code><>, ==, >, <</code>	Low, Medium, High Child, Adult, Old
Quantitative-Interval	<code><>, ==, >, <, -</code> <i>measure distances</i>	<i>date:</i> 1 Jan 2016, 10 Jan 2016 <i>location:</i> Moscow, London
Quantitative-Ratio (zero fixed)	<code><>, ==, >, <, -, +, %</code> <i>measure ratios</i>	<i>physical measurements:</i> 10cm, 30kg, 25C

Matplotlib

<http://matplotlib.org>

Type of Plots

- **plt.plot** – creates a line plot
- **plt.bar** – creates a bar chart
- **plt.boxplot** – makes a box and whisker plot
- **plt.hist** – makes a histogram
- **plt.scatter** – makes a scatter plot

Title and Labels

- **plt.title** – adds a title to the plot.
- **plt.xlabel** – adds an x-axis label.
- **plt.ylabel** – adds a y-axis label.

Task

Analyze the result of load tests

File Length : 1034494 bytes
File Creation Time : 10/6/2016 3:42:35 PM
File Modified Time : 8/17/2016 8:04:34 PM
File Accessed Time : 10/6/2016 3:42:35 PM

10/7/2016 2:08:39 PM : BEGIN COPYING FILES TO STAGING DIRECTORY
10/7/2016 2:08:39 PM : END COPYING FILES TO STAGING DIRECTORY

10/7/2016 2:08:39 PM : BEGIN PUBLISHING PDF

Server run id : df14a150-571a-46b5-988b-a4aaa62c2442

10/7/2016 2:08:50 PM : Server run status : RunAccepted

10/7/2016 2:10:13 PM : Server run status : RunInProgress

10/7/2016 2:10:23 PM : Server run status : RunQueued

10/7/2016 2:14:02 PM : Server run status : RunInProgress

10/7/2016 2:14:12 PM : Server run status : RunQueued

10/7/2016 2:14:33 PM : Server run status : RunInProgress

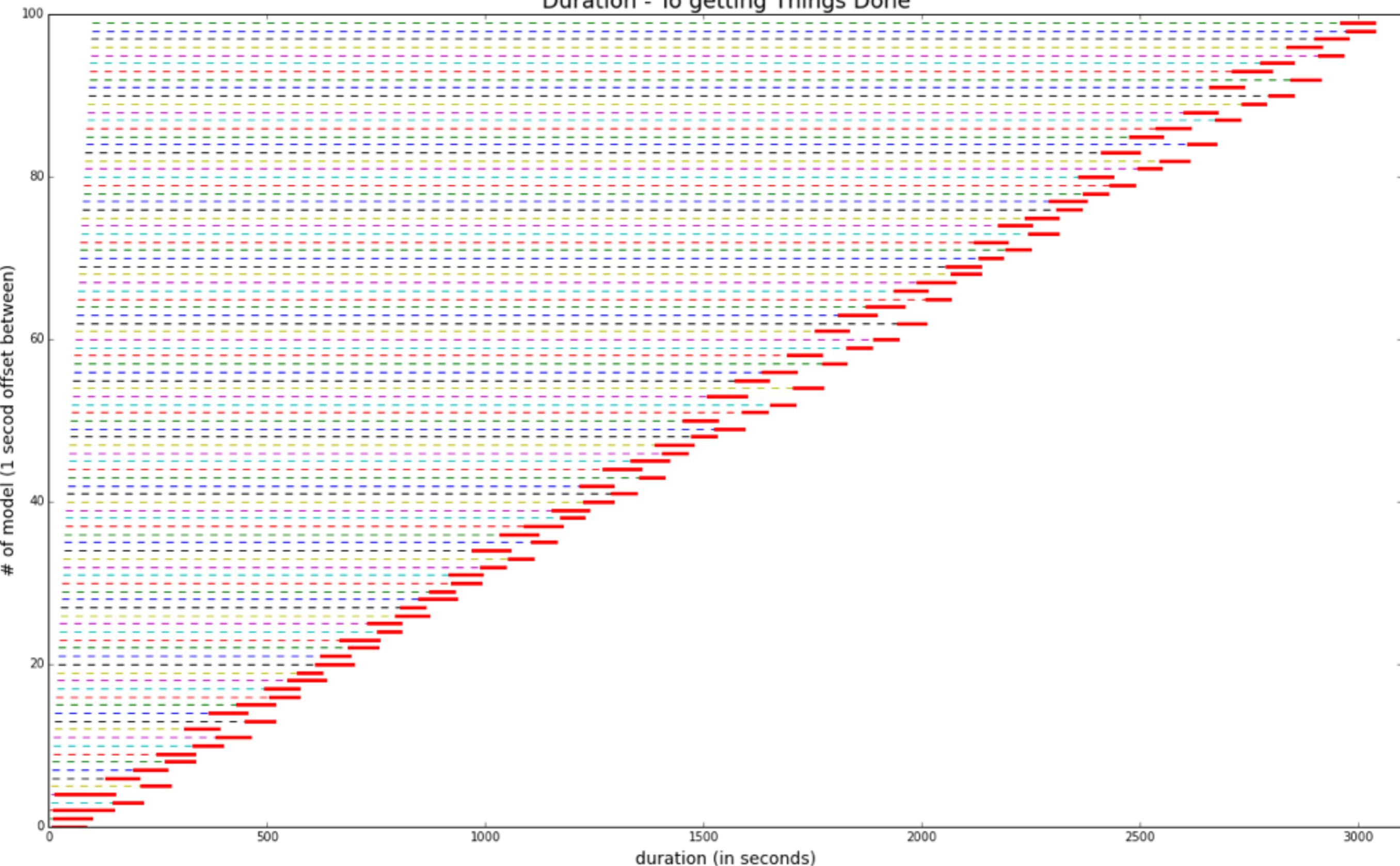
10/7/2016 2:15:57 PM : END PUBLISHING PDF

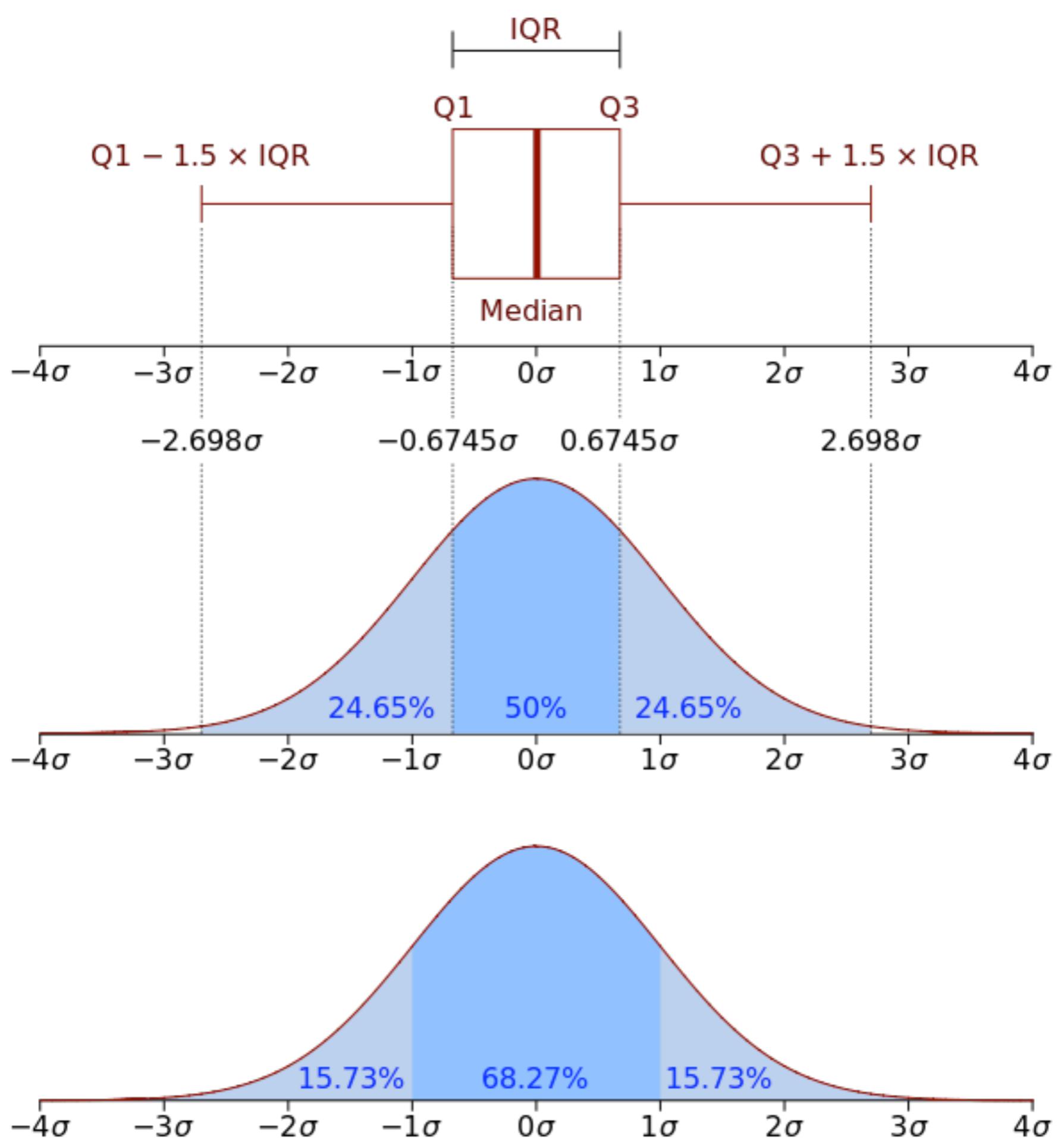
PDF for part.prt.1 created successfully

10/7/2016 2:15:57 PM : BEGIN COPY PDF TO LOCAL FOLDER
10/7/2016 2:15:57 PM : END COPY PDF TO LOCAL FOLDER

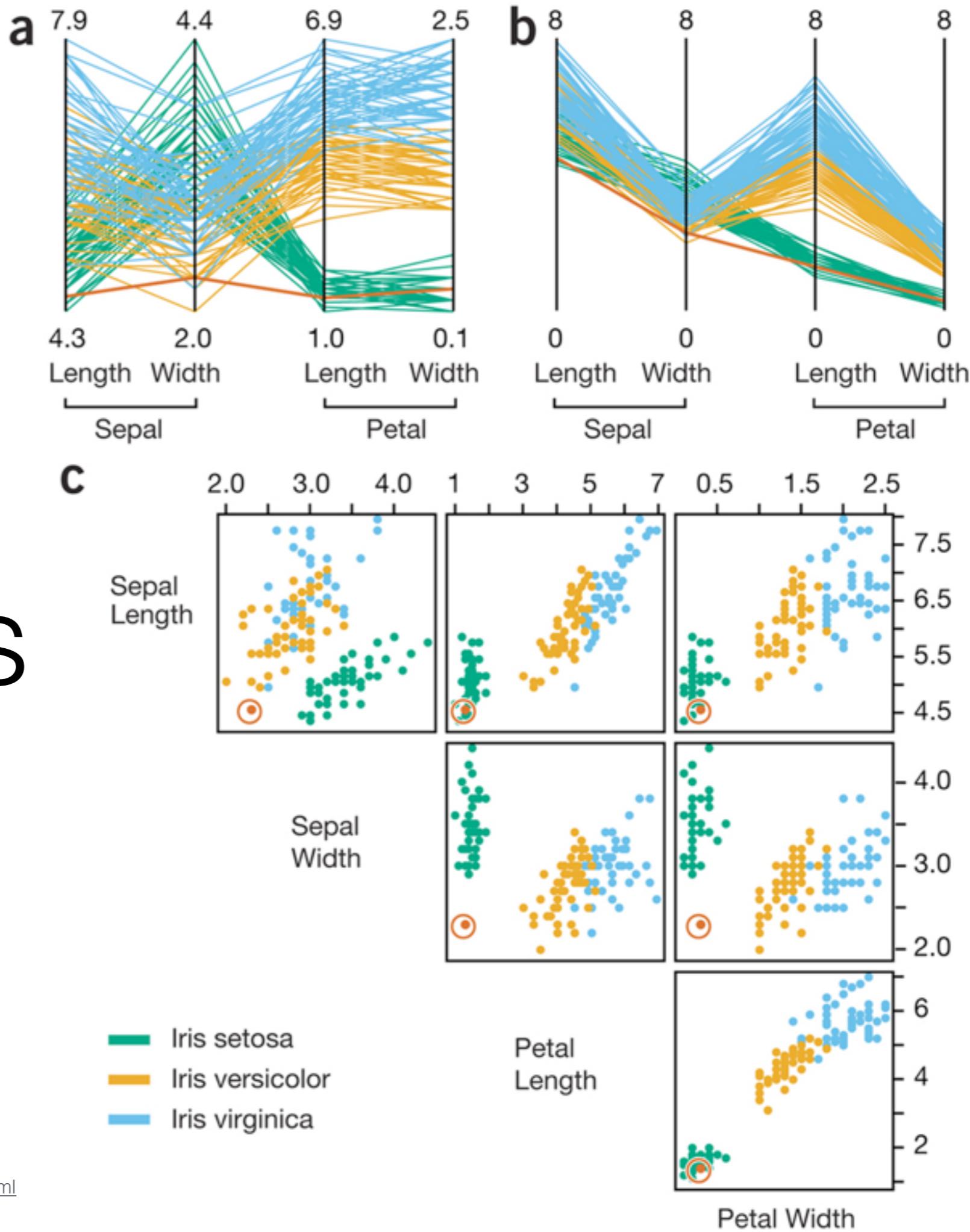
10/7/2016 2:15:57 PM : =====END PROCESSING FILE=====

Duration - To getting Things Done

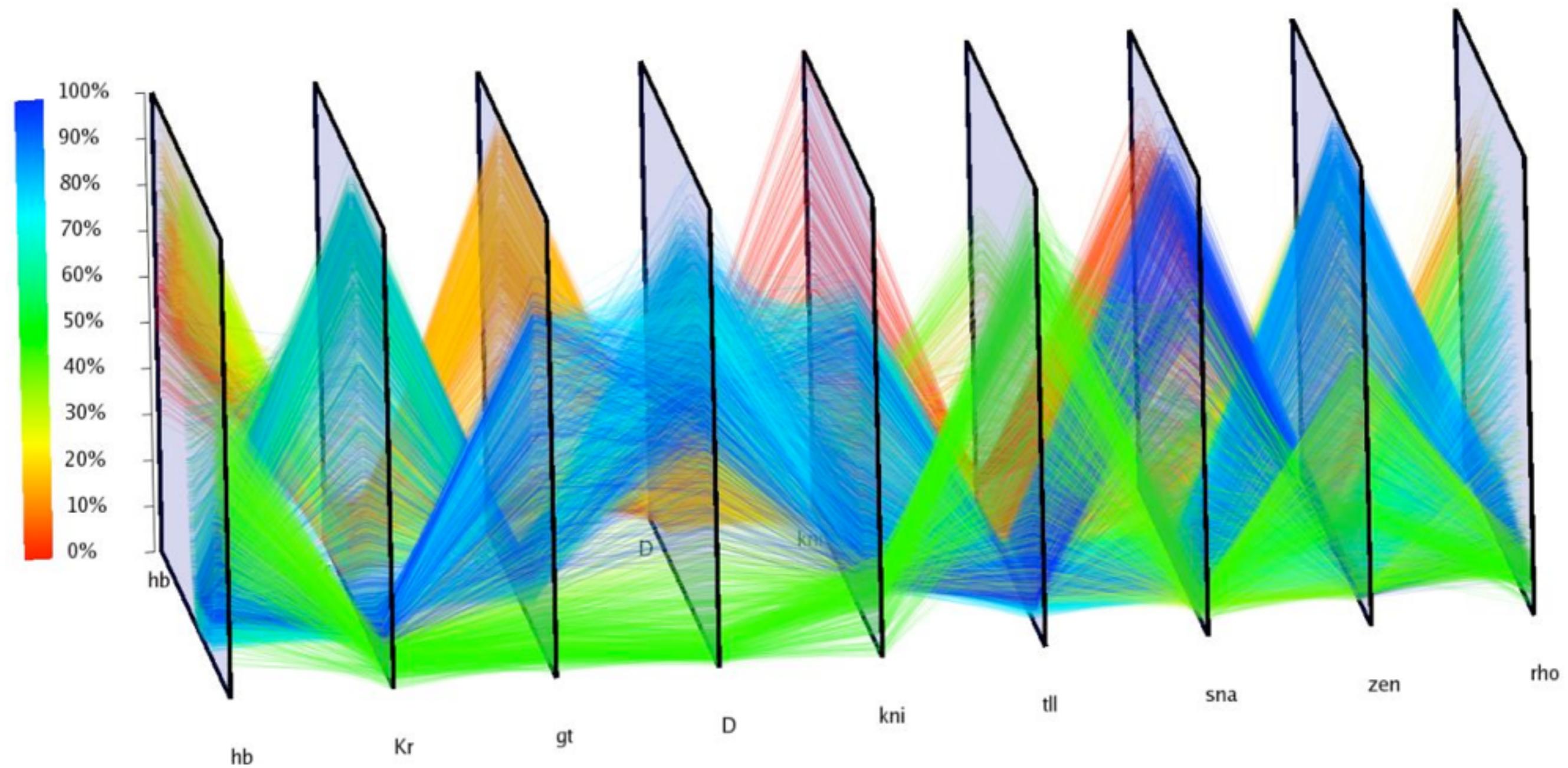




Parallel Coordinates Plot



3D Parallel Coordinates Plot



Summary

Three things

if you can remember only three...

- Visualization help discover data
- Avoid pie charts (use histogram, bartplot ...)
- Show more than 1 or 2 variables

Thank you