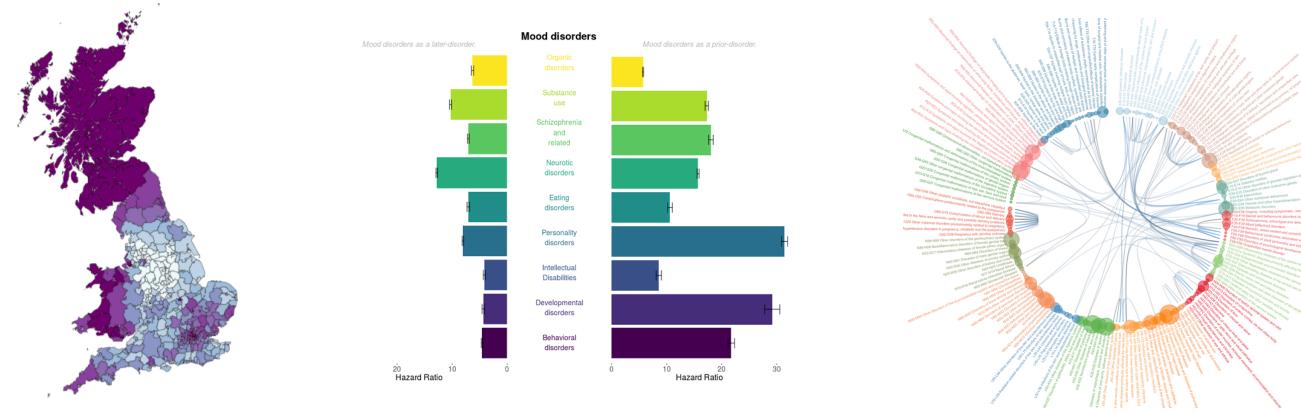
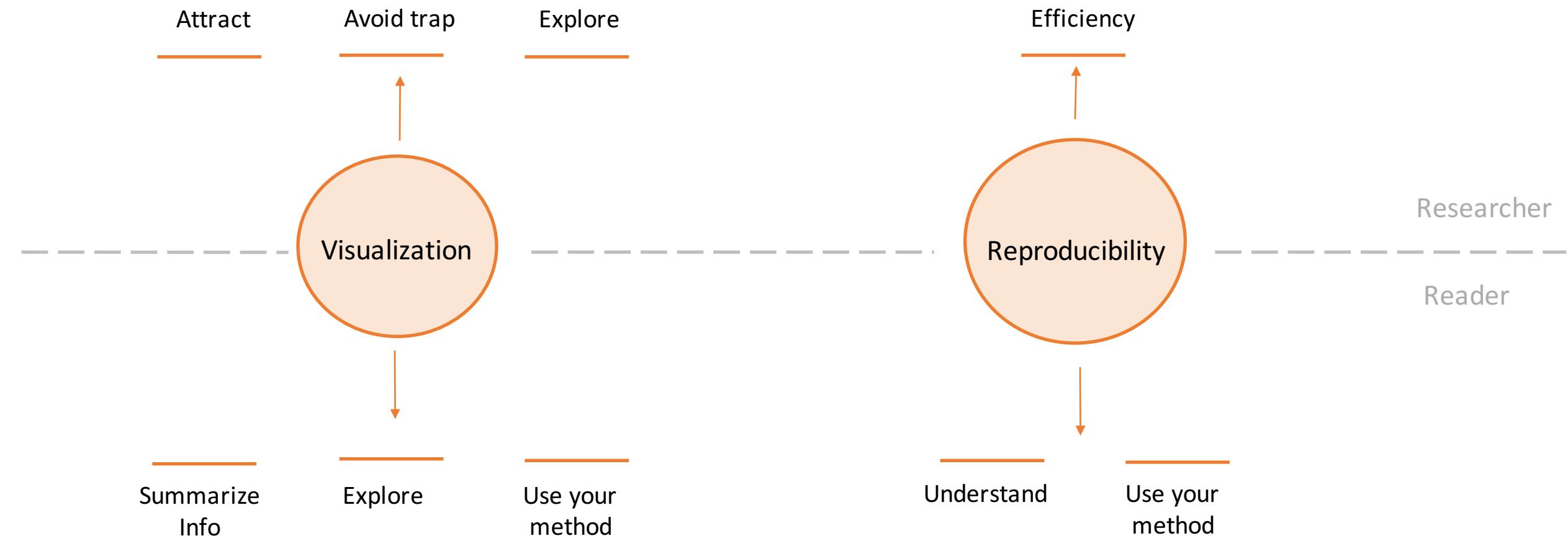


# DATAVIZ, REPRODUCIBILITY AND IMPACTFUL RESEARCH

---

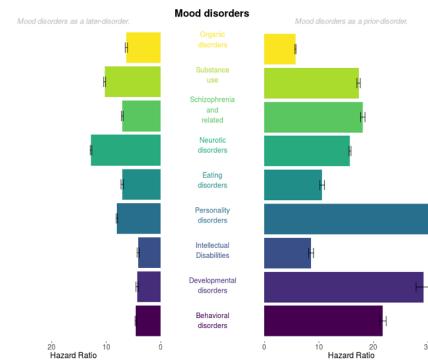




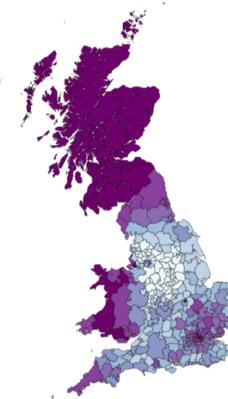
# DATAVIZ



Comorbidity  
in the UK  
Biobank



Comorbidity  
in the Danish  
register

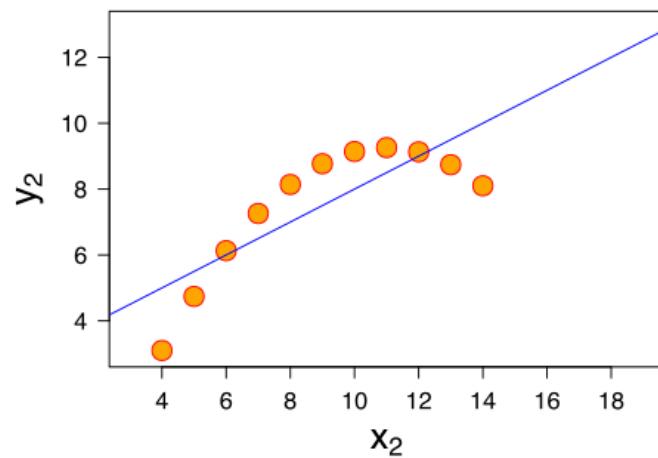
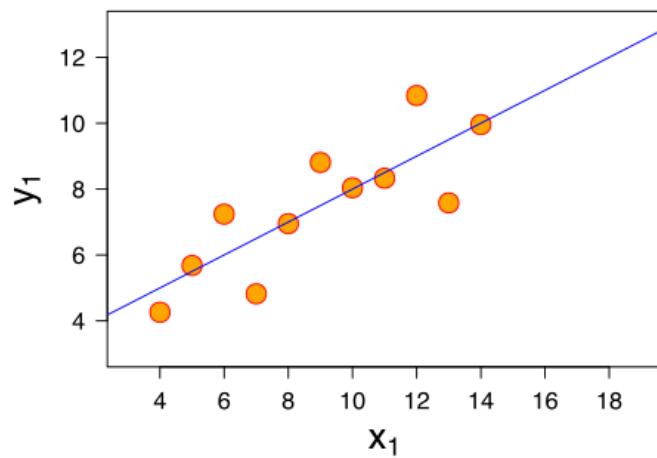


Geography in  
the UK  
Biobank

$$Y = 3 + 0.5x$$
$$\text{Cor} = 0.8$$

$$\text{Mean}(x) = 9$$
$$\text{Var}(x) = 11$$

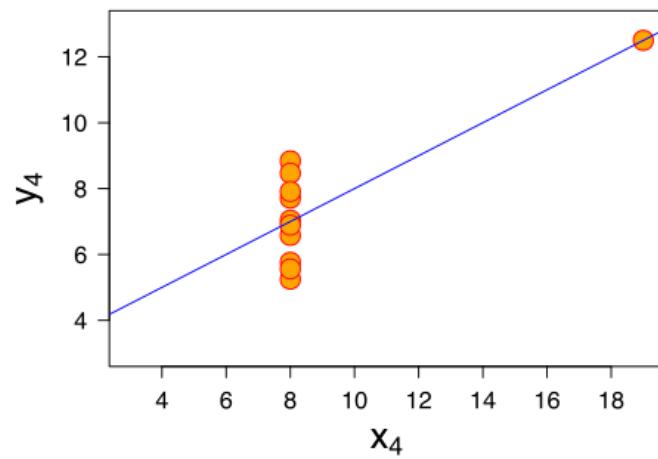
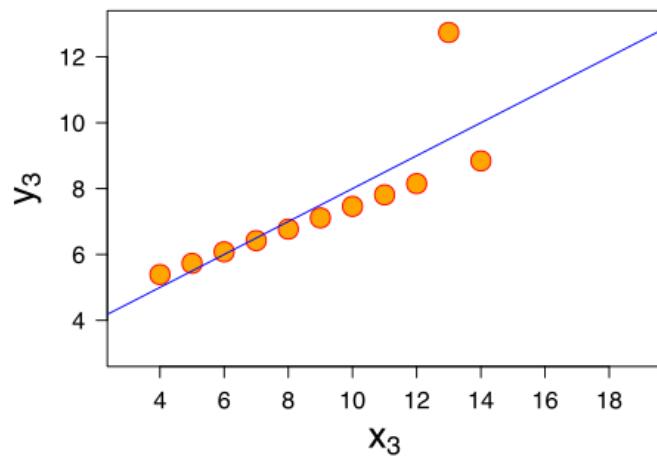
$$\text{Mean}(Y) = 7.5$$
$$\text{Var}(Y) = 4.1$$



$Y = 3 + 0.5x$   
Cor = 0.8

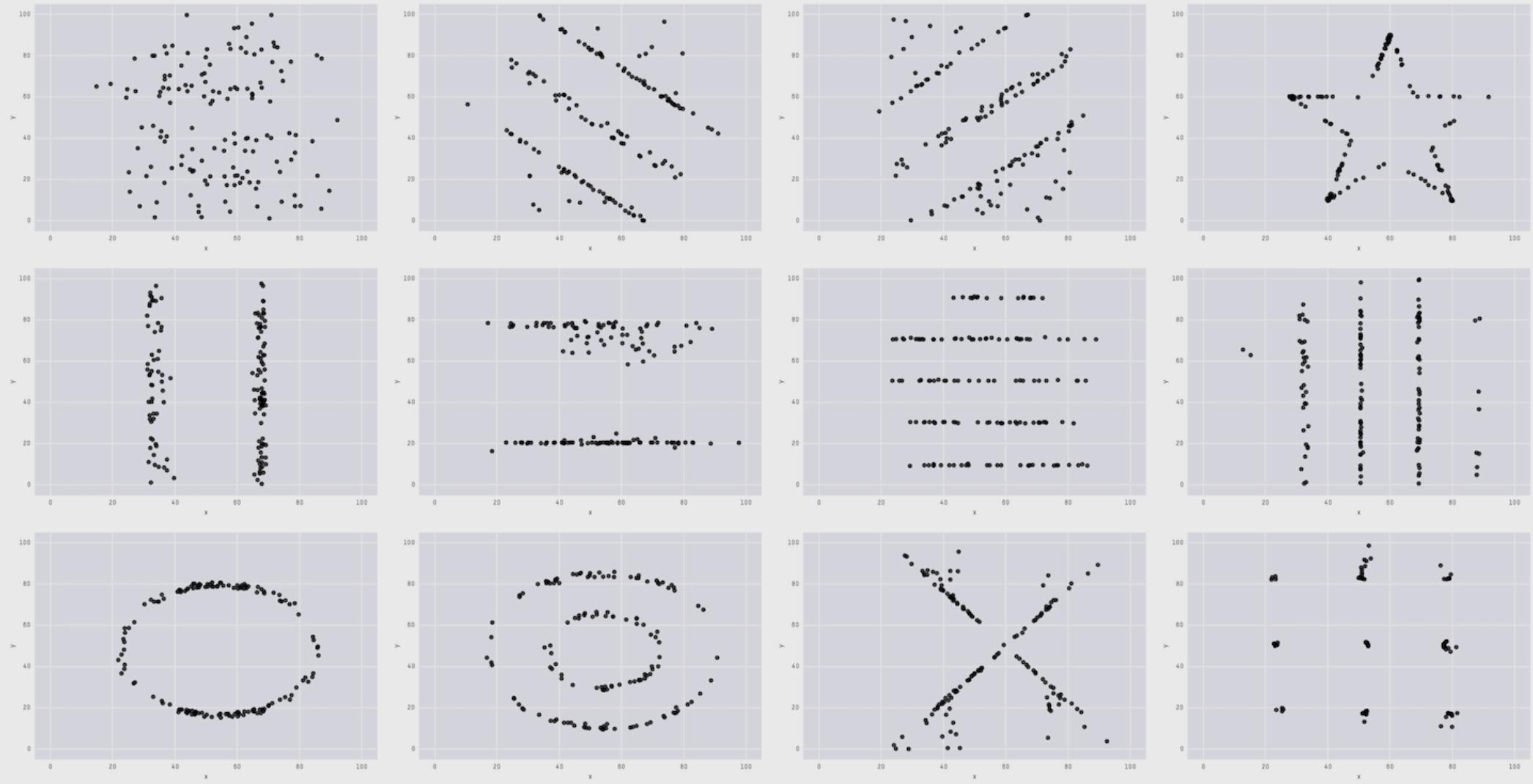
Mean( $x$ ) = 9  
Var( $x$ ) = 11

Mean( $Y$ ) = 7.5  
Var( $Y$ ) = 4.1

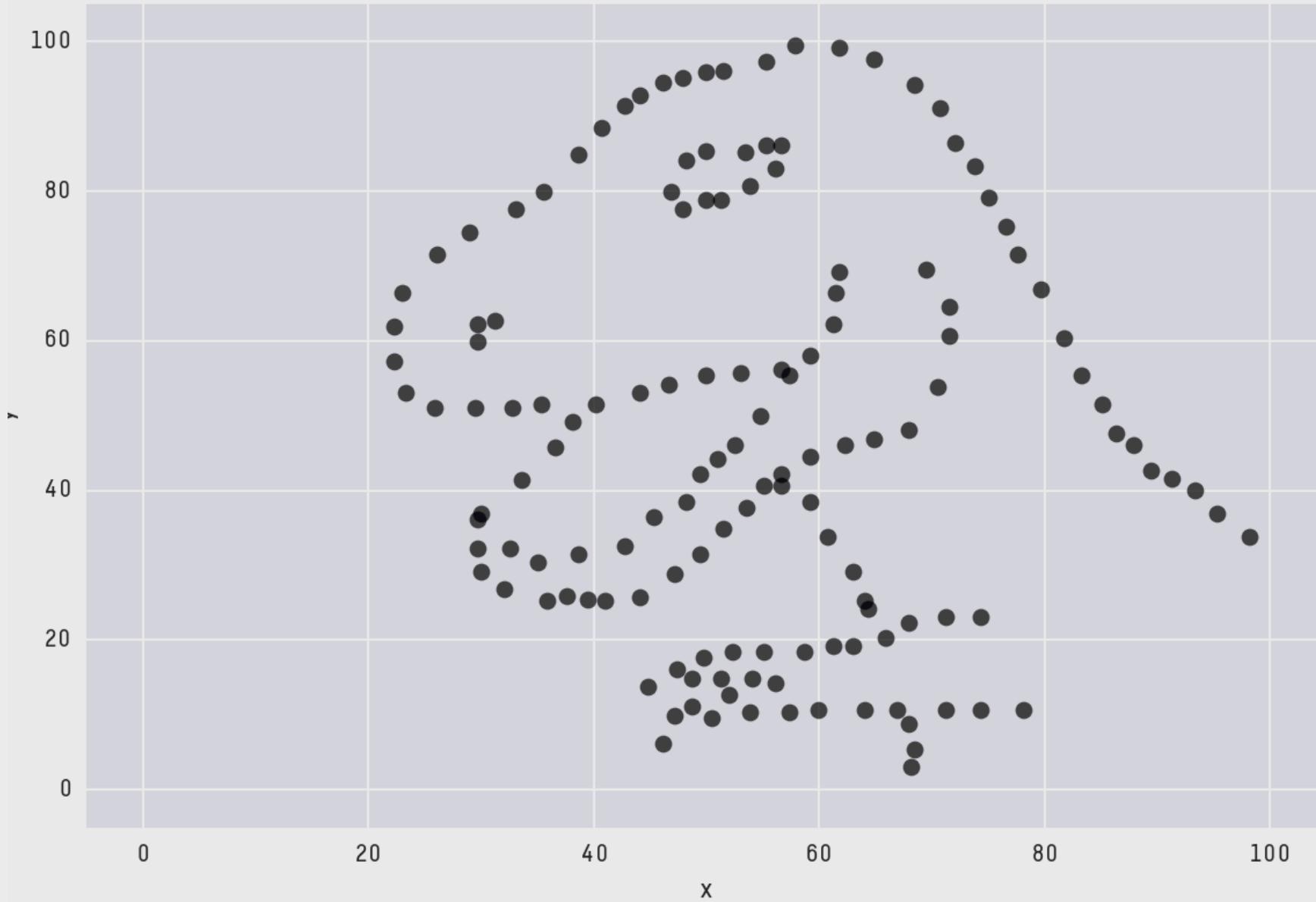


Plot your data

## Anscombe's quartet



## The Datasaurus Dozen

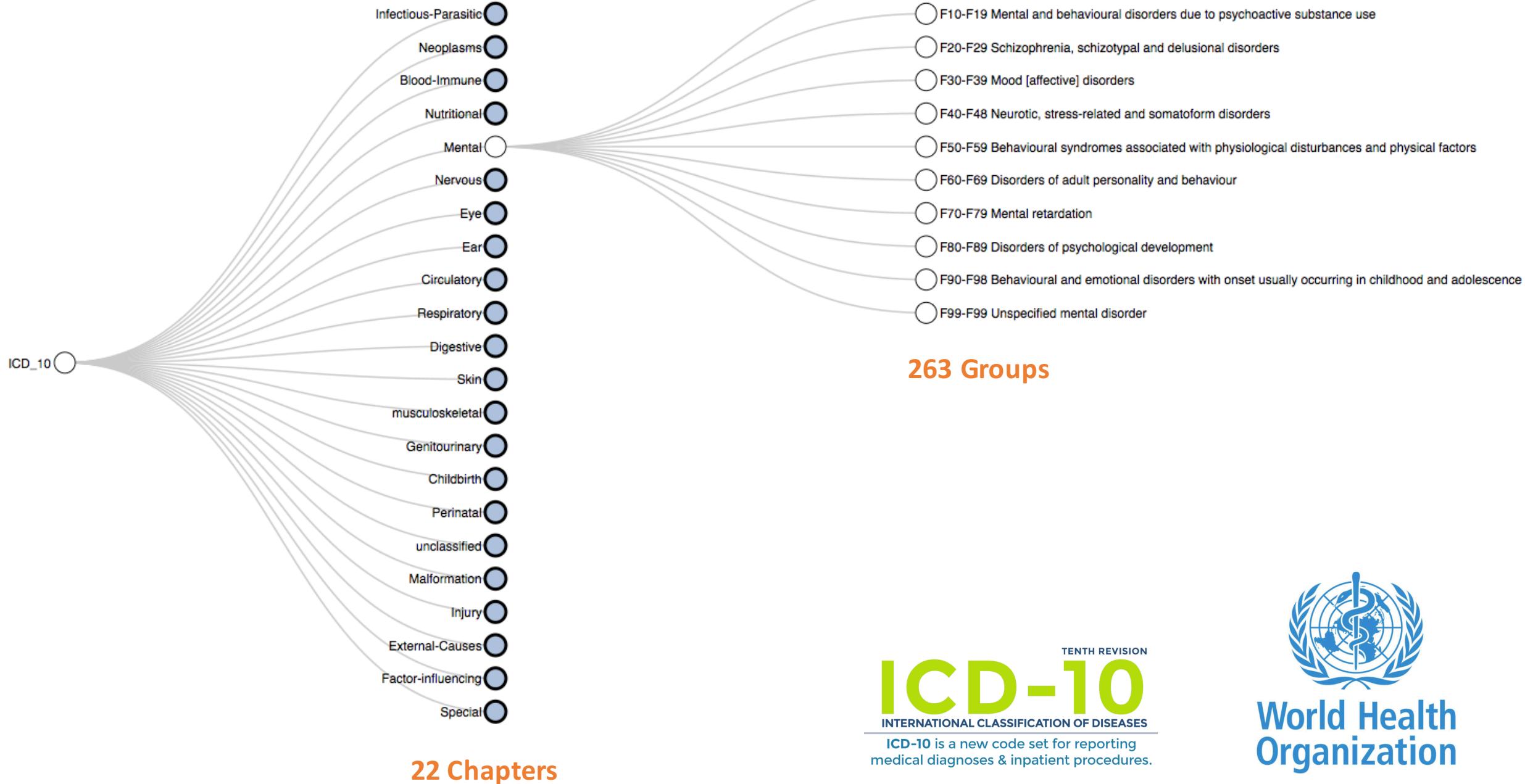


The Datasaurus Dozen

Cairo, 2017

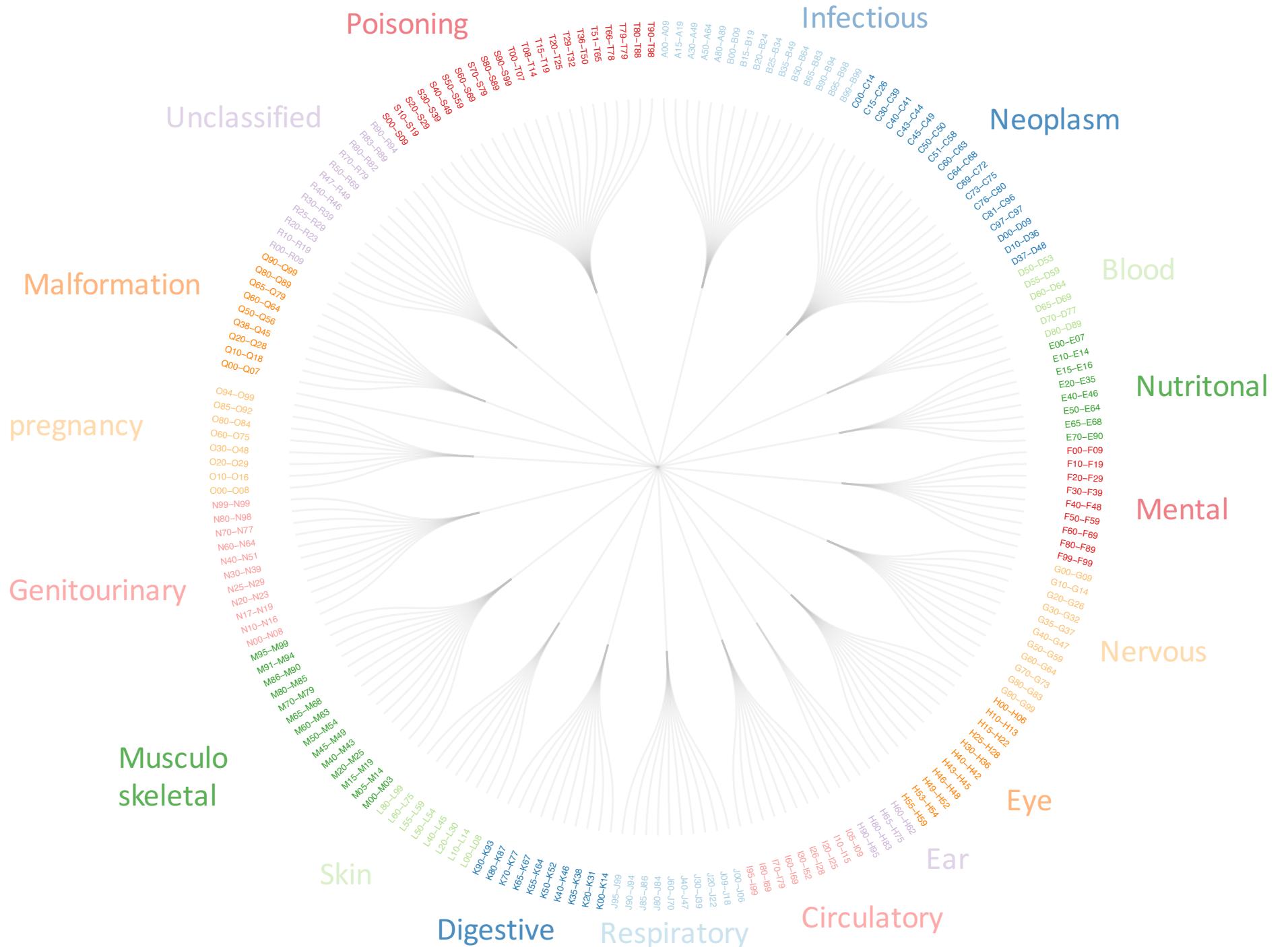
Matejka & Fitzmaurice, 2017

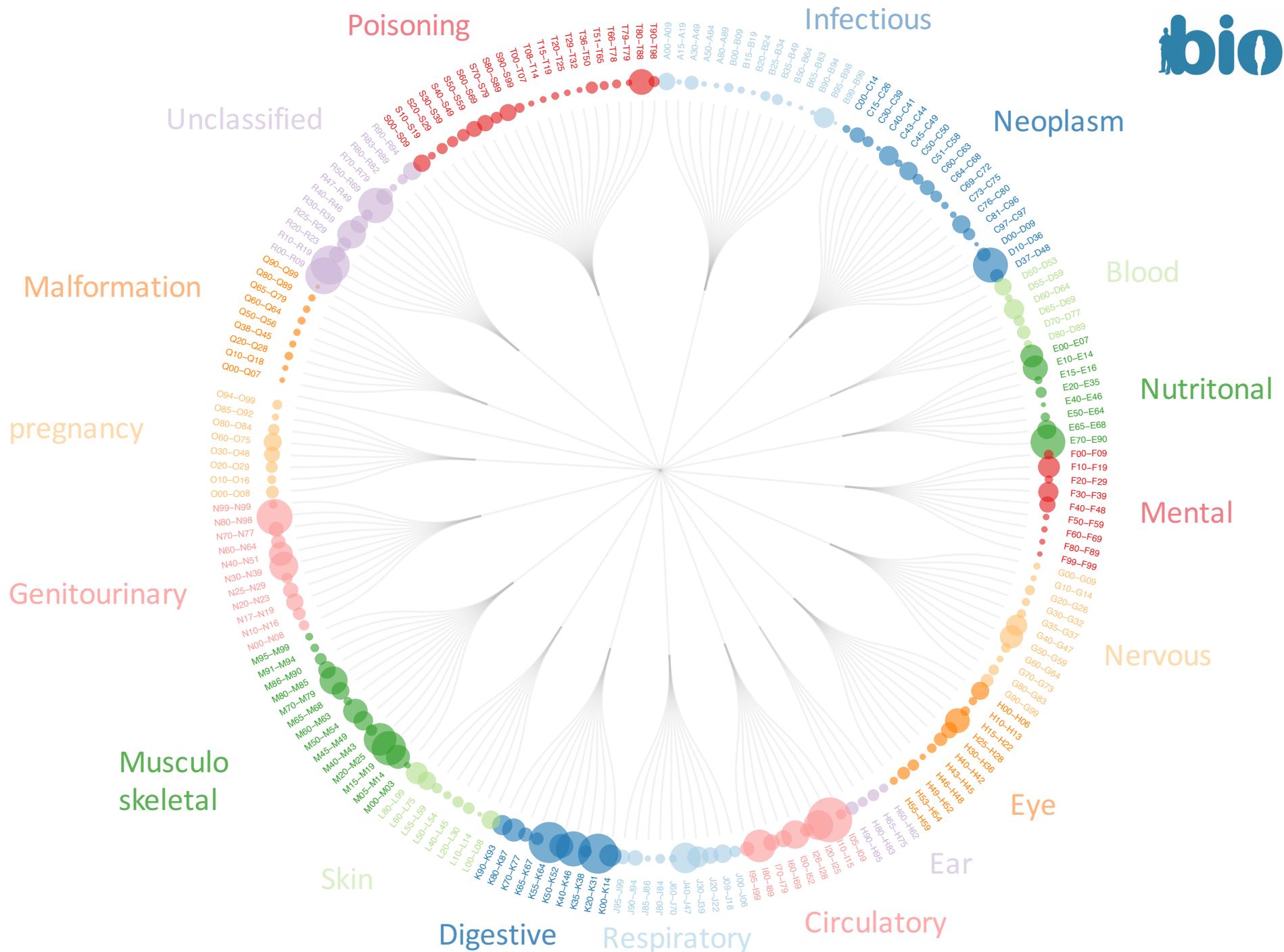
# The ICD 10 classification

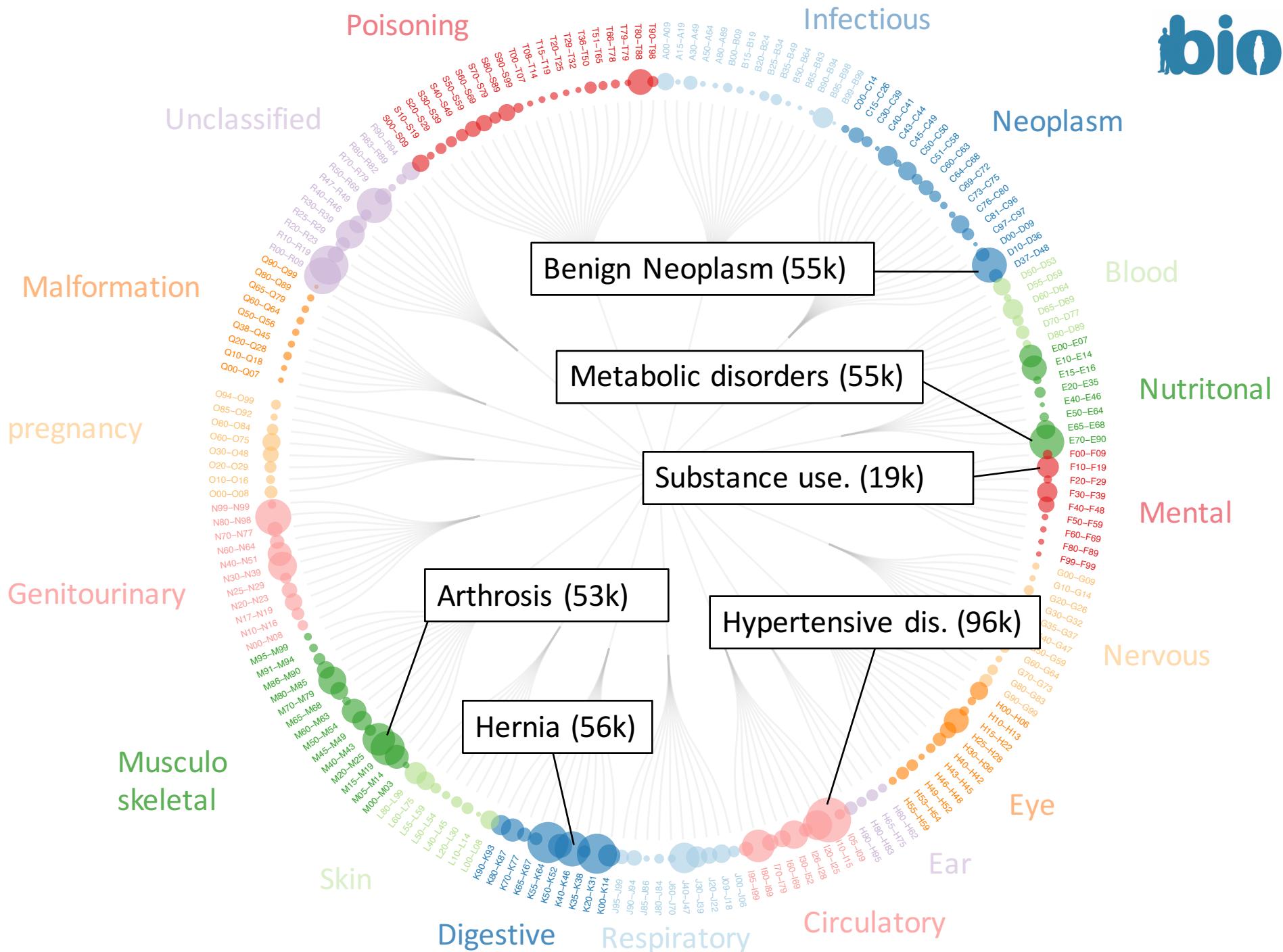


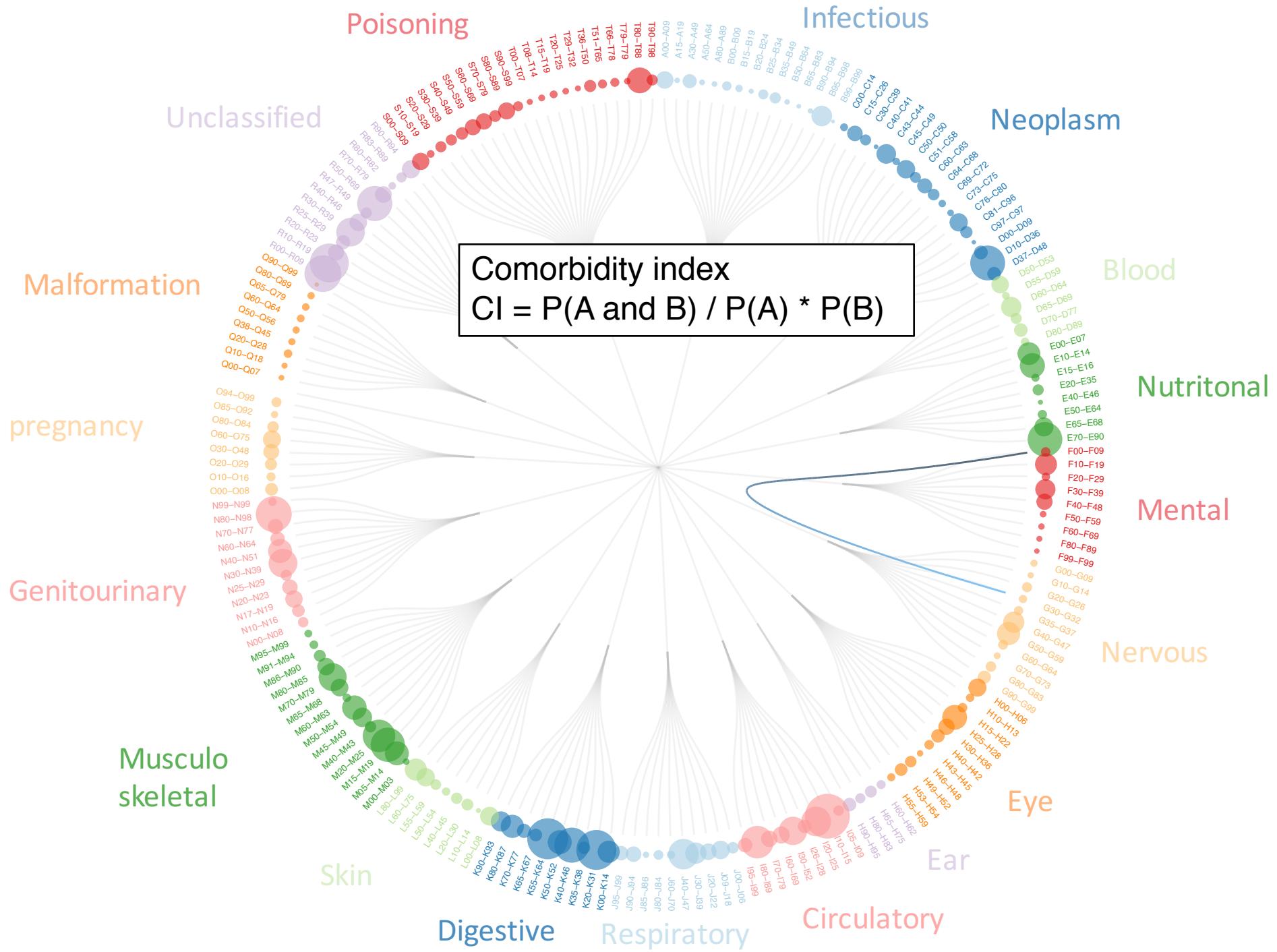
TENTH REVISION  
**ICD-10**  
INTERNATIONAL CLASSIFICATION OF DISEASES  
ICD-10 is a new code set for reporting  
medical diagnoses & inpatient procedures.



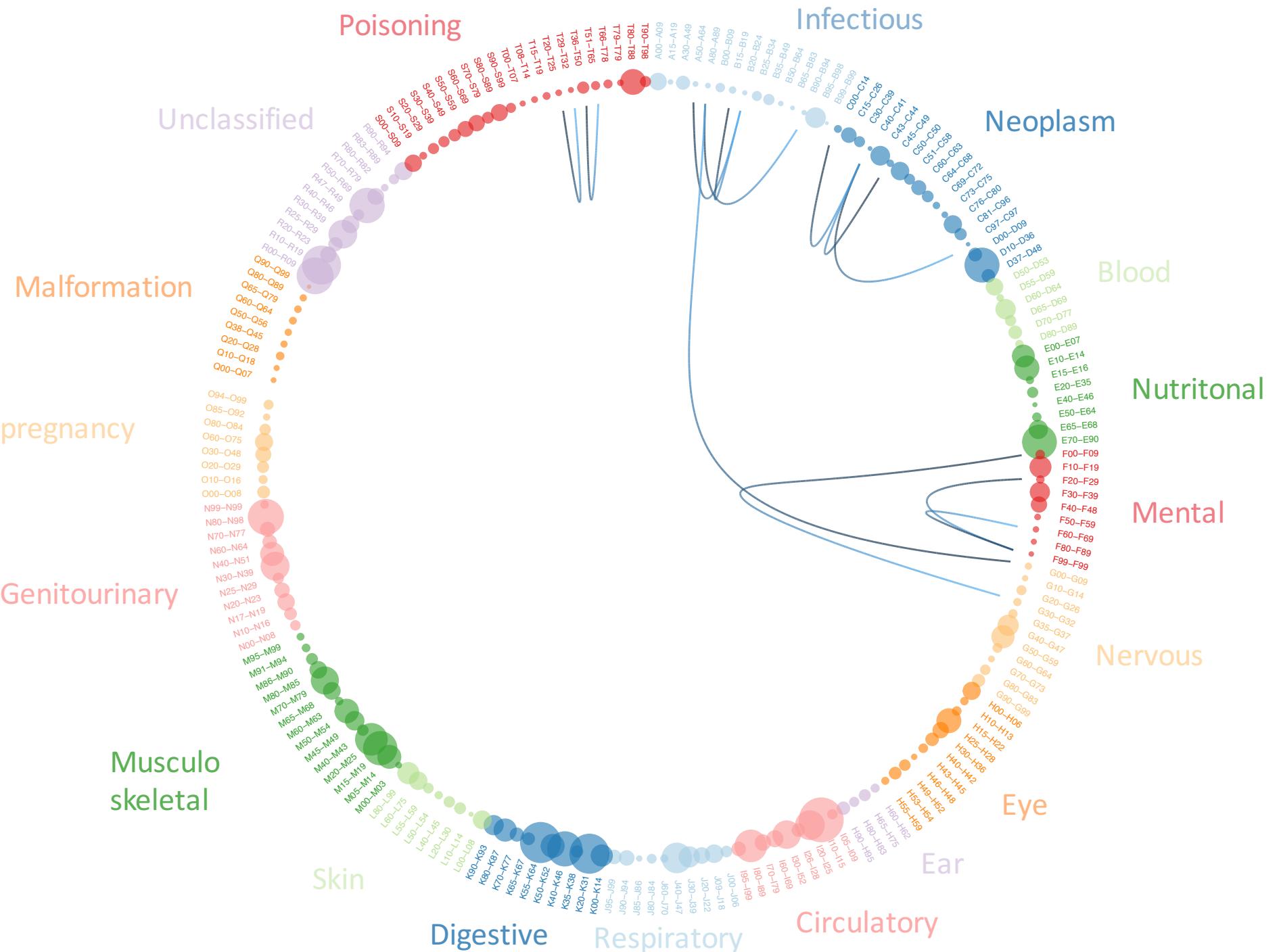




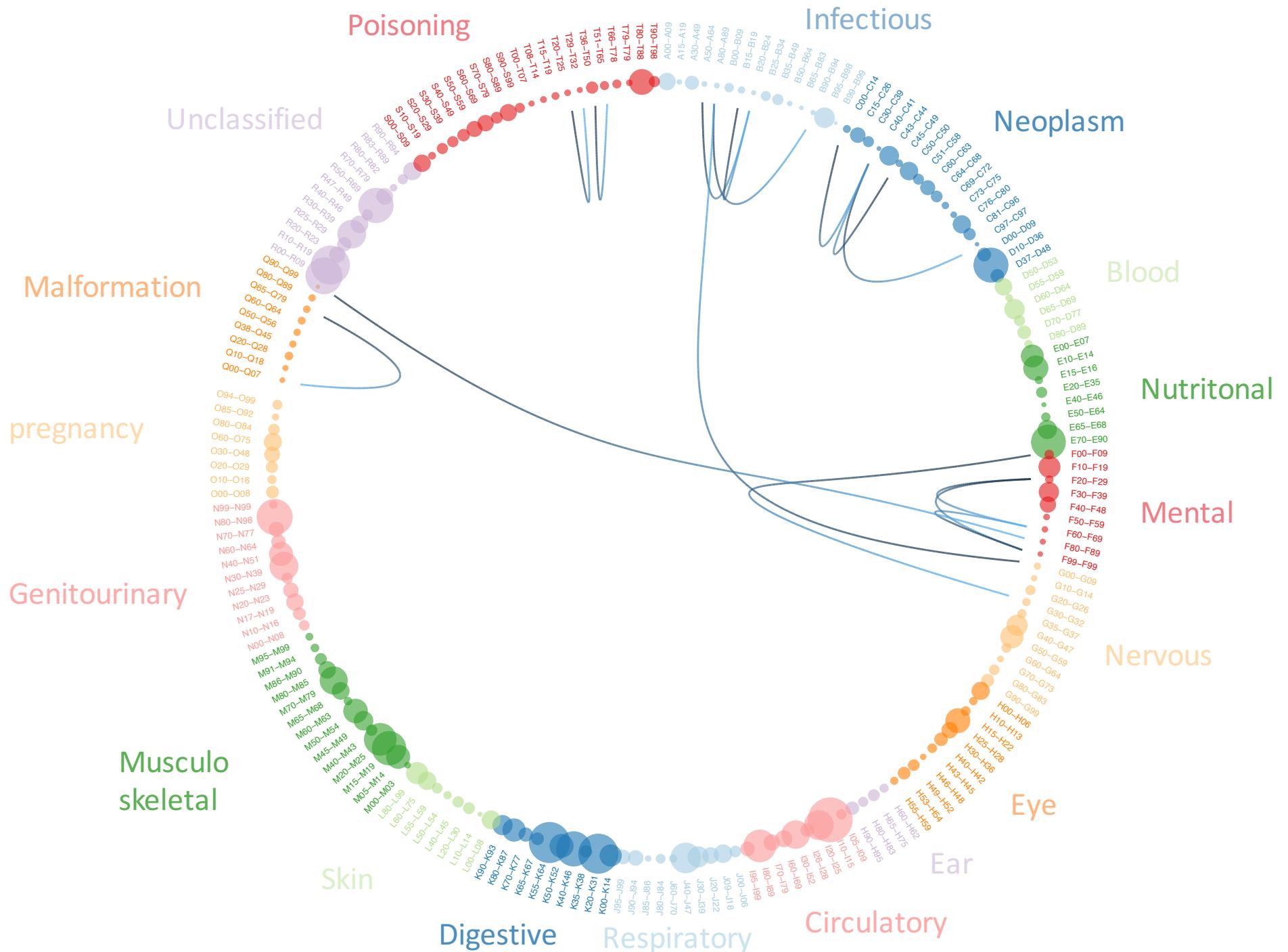




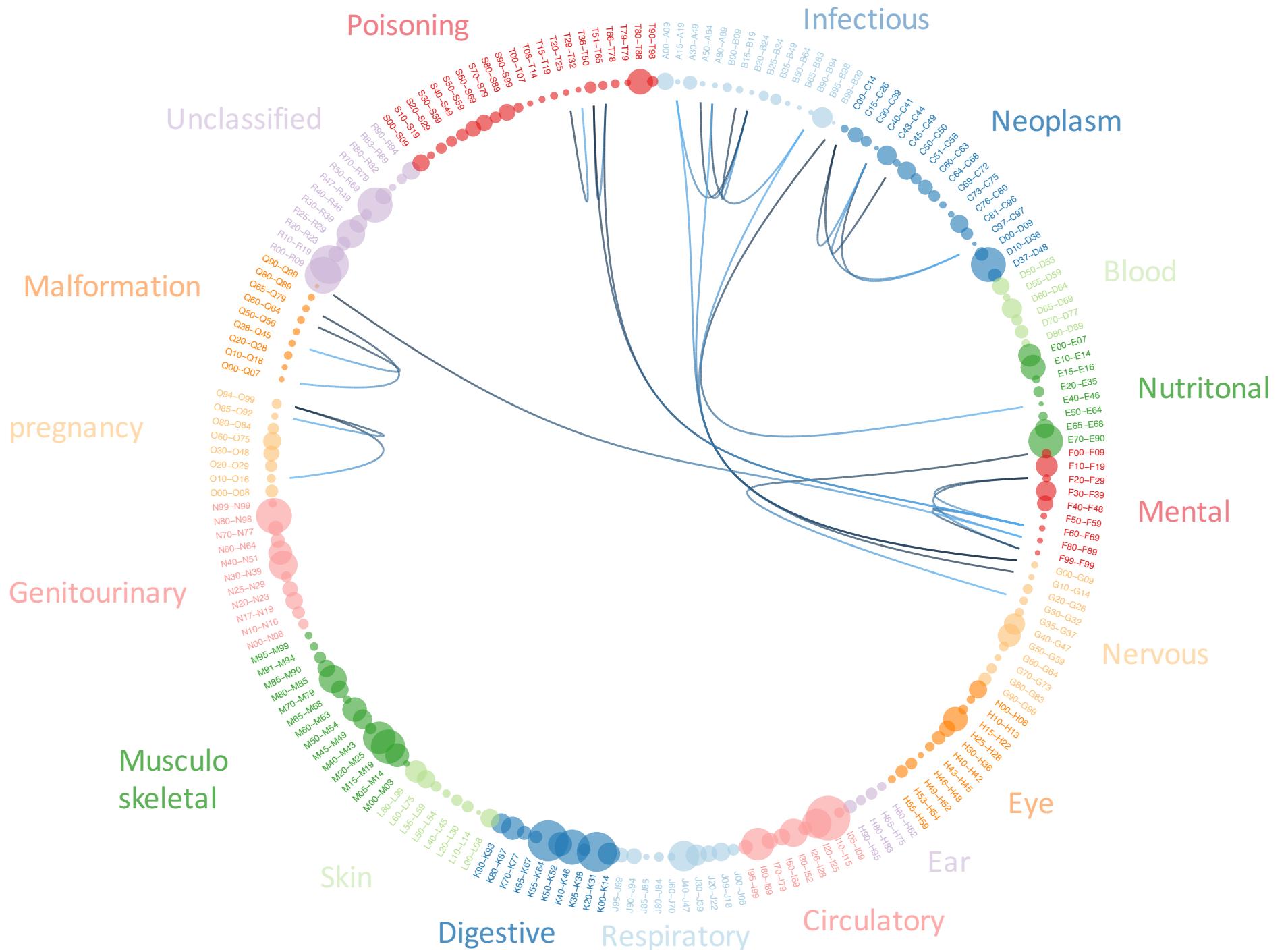
CI > 70



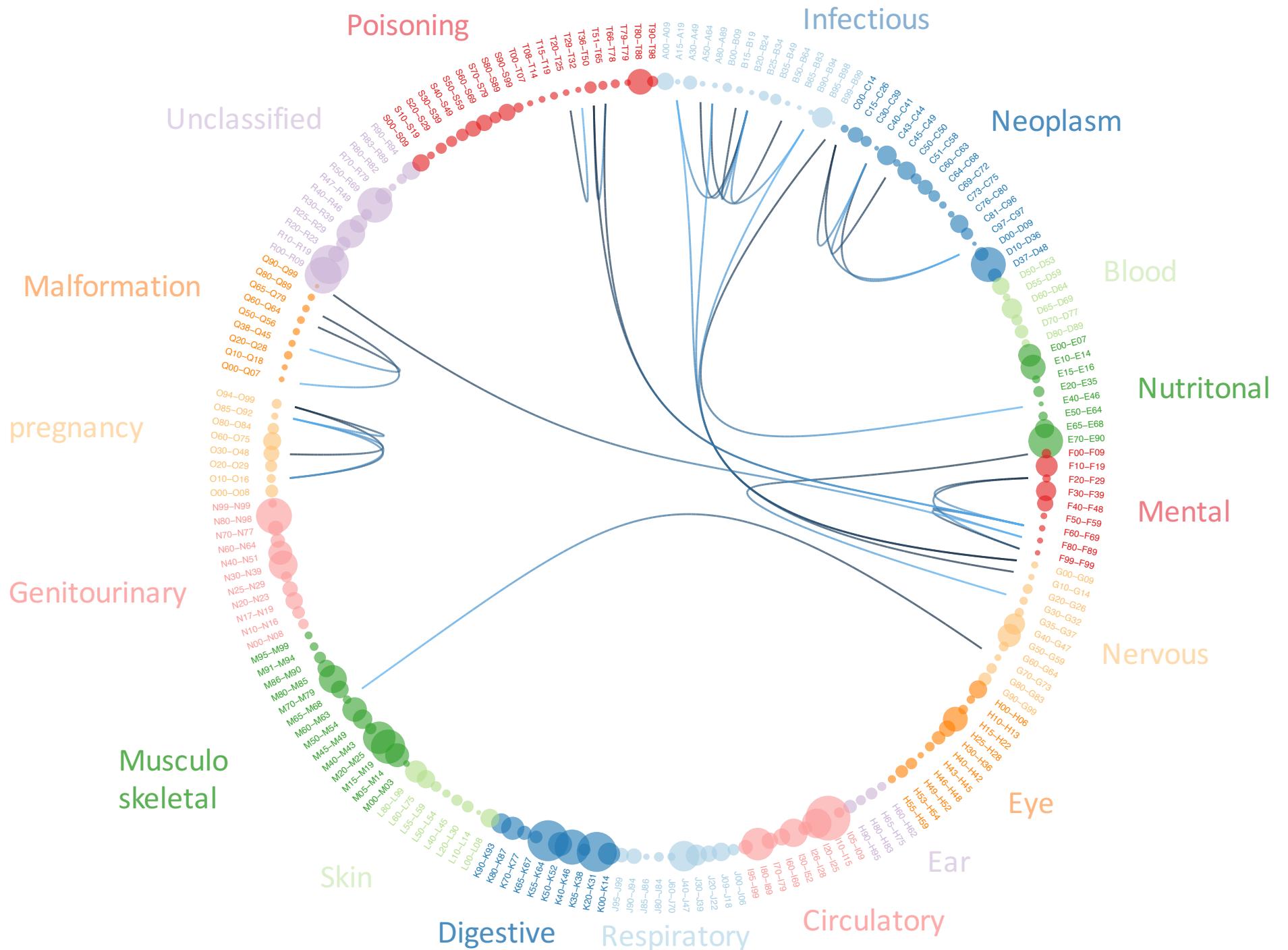
CI > 60



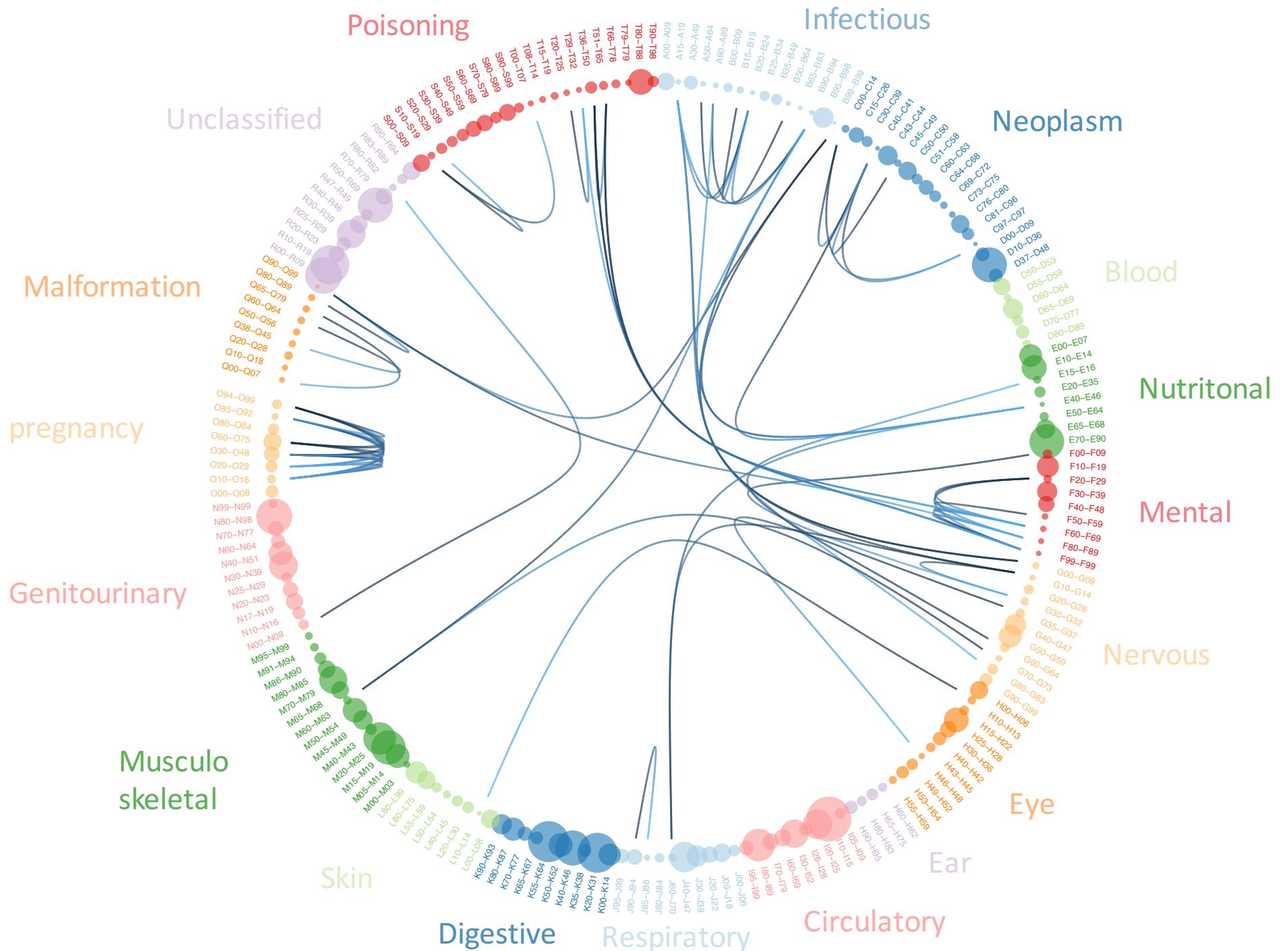
CI > 50



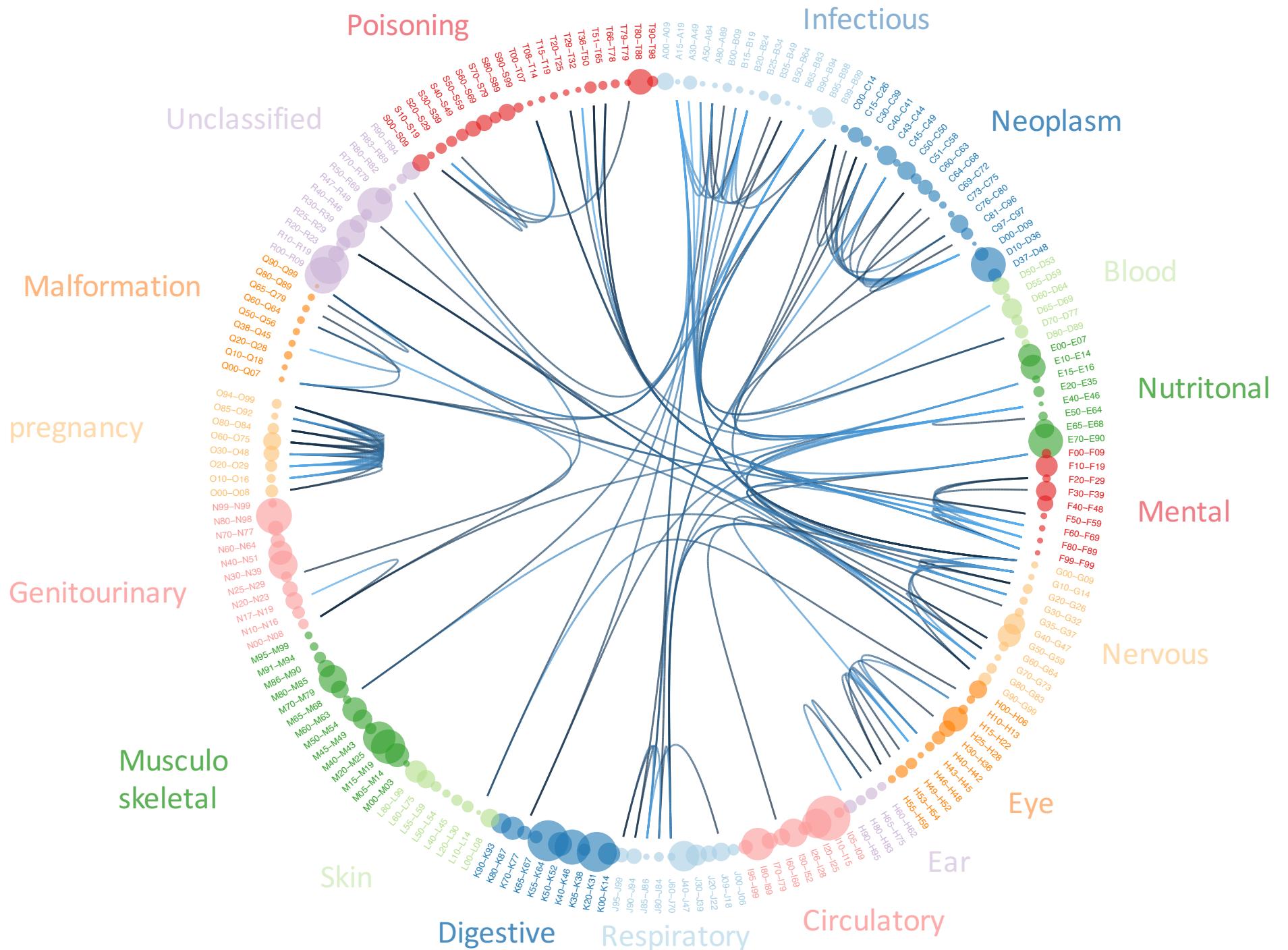
CI > 40



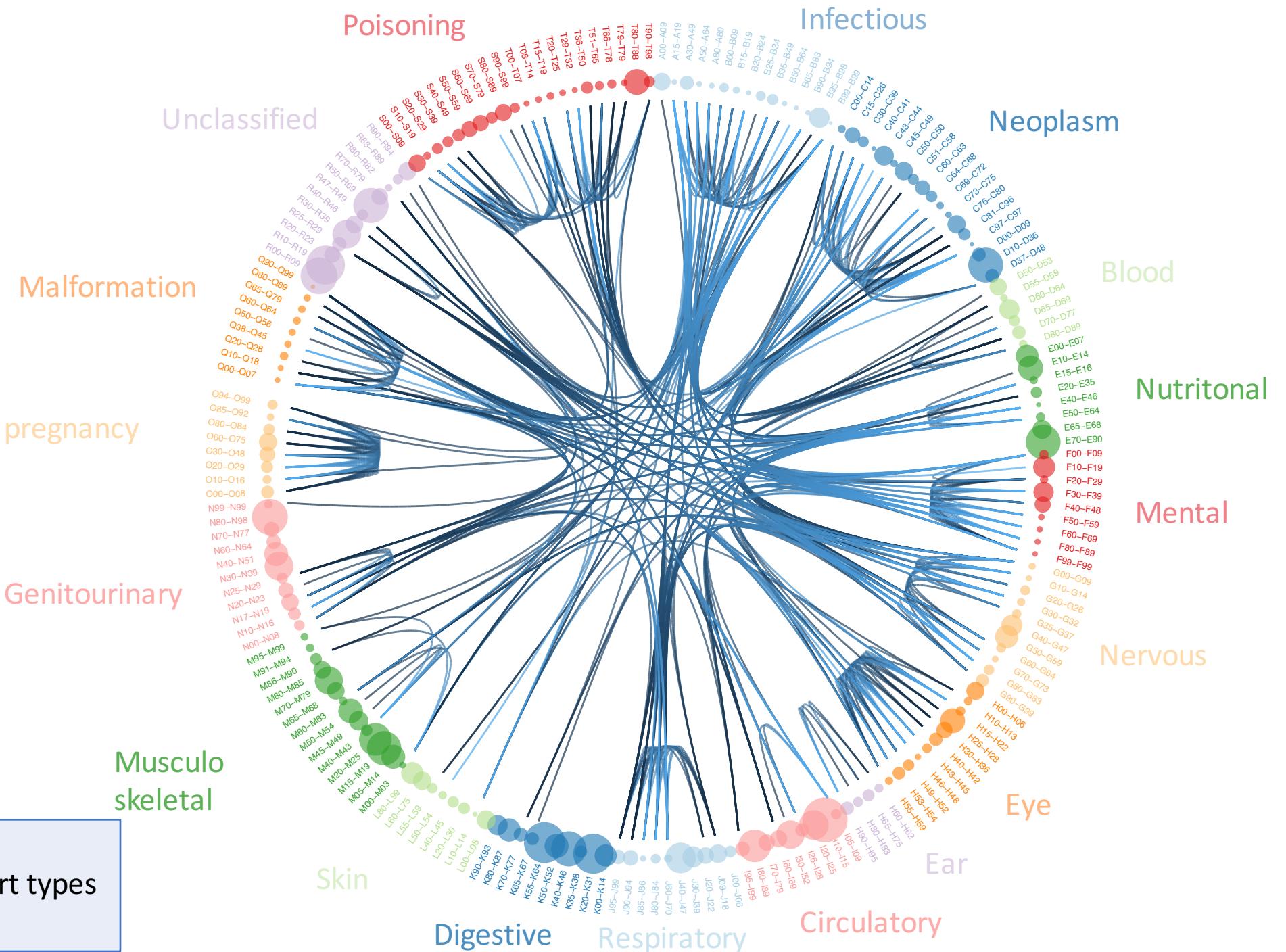
CI > 30



CI > 20



**CI > 10**



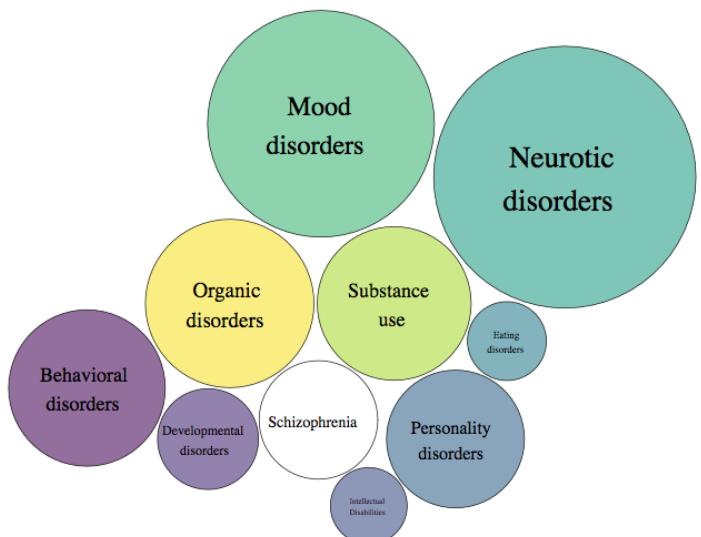


[R-graph-gallery.com](http://R-graph-gallery.com)

[Python-graph-gallery.com](http://Python-graph-gallery.com)



Danish register  
(n ~ 5M)



10 Mental Disorder  
Groups (F chapter  
of ICD 10)

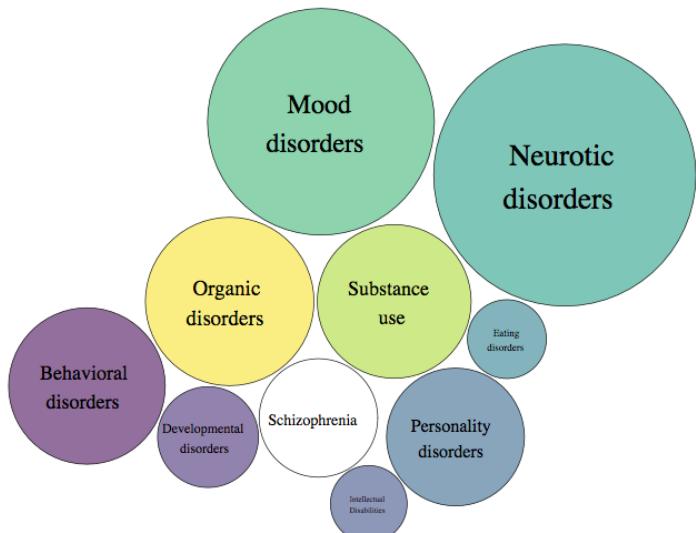


#### Compute Hazard Ratios:

- 90 pairs
- Males and Female
- Several Models
- Many angles of analysis



Danish register  
(n ~ 5M)



10 Mental Disorder  
Groups (F chapter  
of ICD 10)

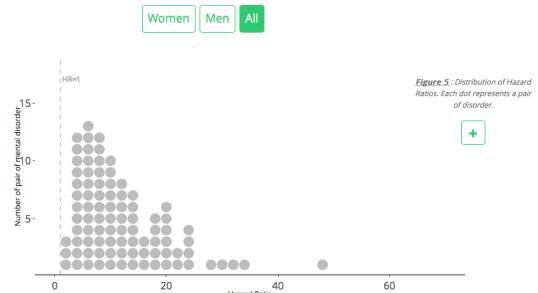
The NB-COMO project

[Methodology](#) [Results](#) [Meet the team](#)

#### 1. Pairwise, temporally-ordered hazard ratios

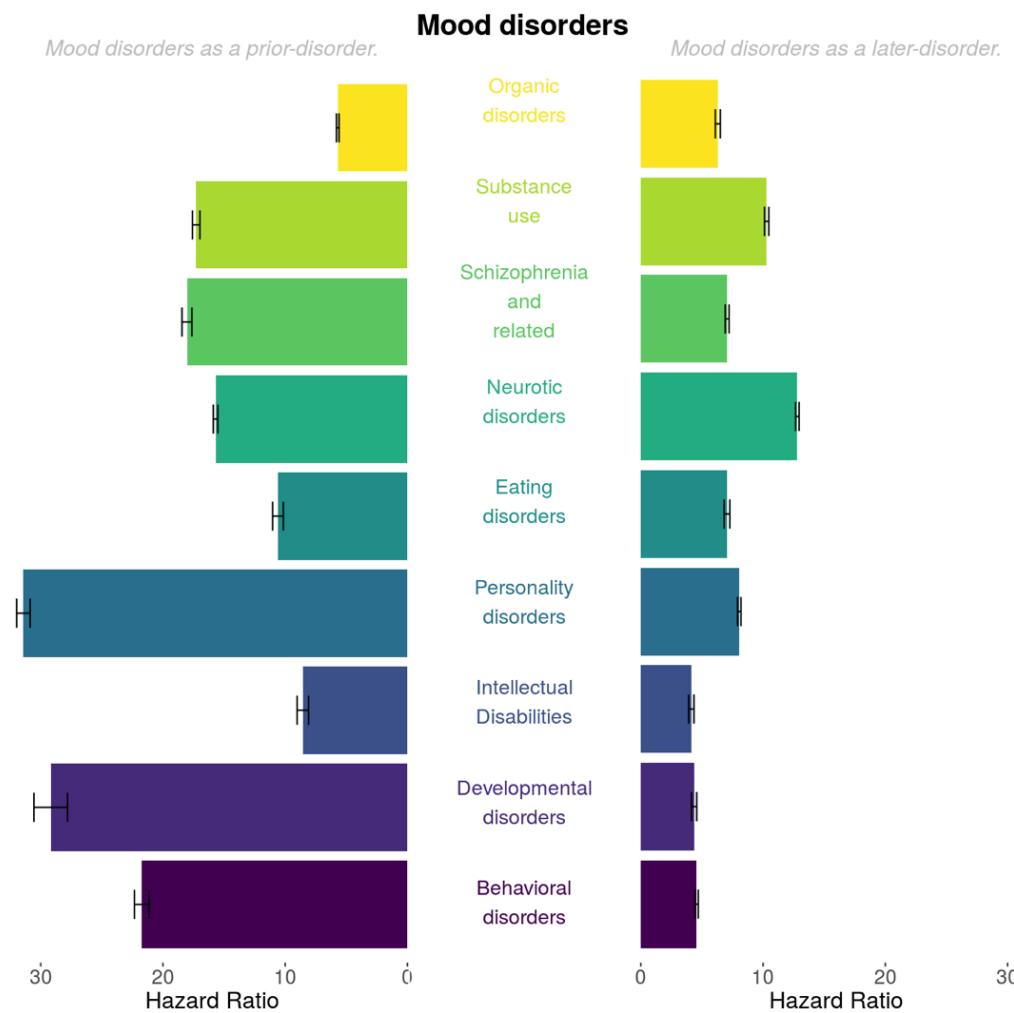
We calculated hazard ratios between each pair of Disorder Groups. The vast majority of these disorder estimates are over 1, meaning that individuals with a prior mental disorder (prior-disorder), are at greater risk in developing subsequent mental disorder (later-disorder).

The plot below shows the distribution of hazard ratios for all COMO pairs. Note that all hazard ratios are greater than 1.



#### Compute Hazard Ratios:

- 90 pairs
- Males and Female
- Several Models
- Many angles of analysis

[Organic](#)[Substance](#)[Schizophrenia](#)[Mood](#)[Neurotic](#)[Eating](#)[Personality](#)[Intellectual Dis.](#)[Developmental](#)[Behaviour](#)

**Figure 8 :** Description of the symmetry between Disorder Groups. Hazard ratio between all possible disorders and your selected disorder are represented on the left panel and vice-versa.



Allow exploration

[Men](#)[Women](#)[All](#)

# Dataviz Pipeline

---



Data prep



Charts



Interactivity



Button  
Website

**biobank<sup>uk</sup>**

(N ~ 500k)



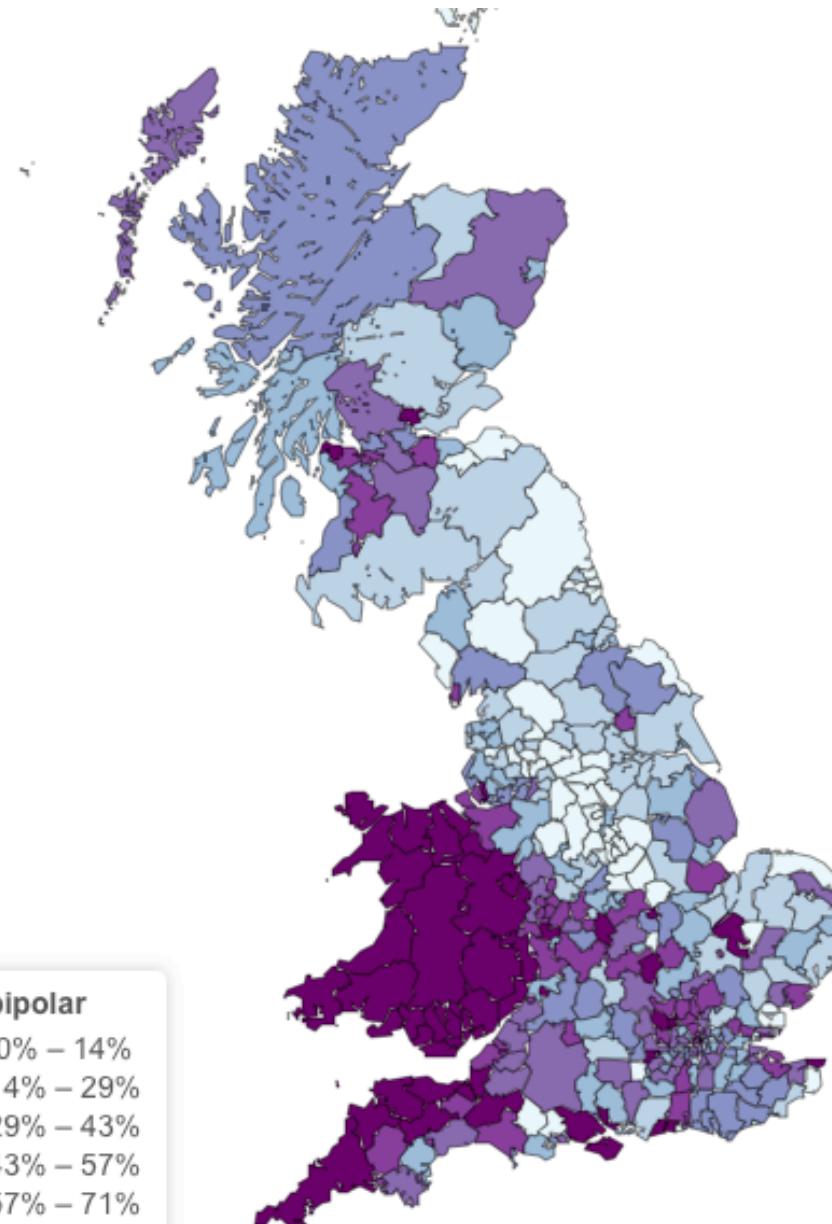
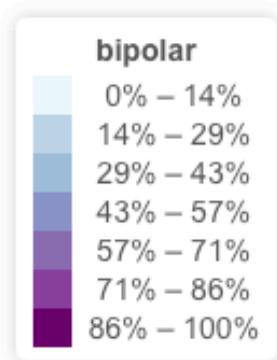
Compute PCs



Compute PRS,  
corrected or not

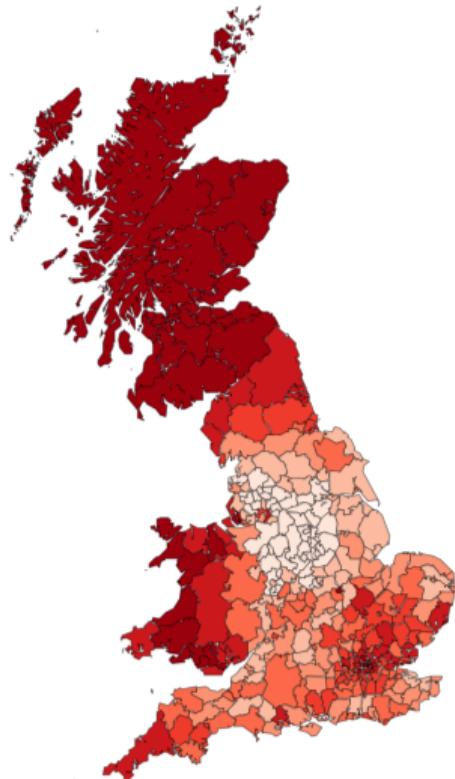


Map

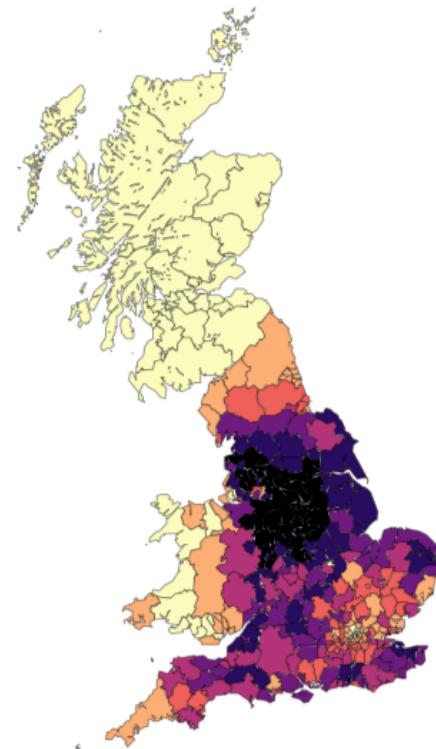


# PC2 representation in the UK Biobank

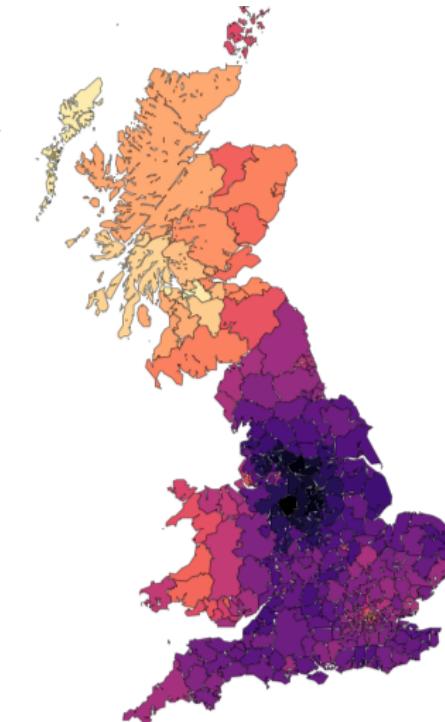
---



Red



BuPu – 6 bins



BuPu – Numeric



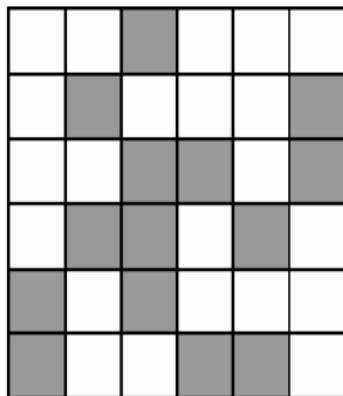
Hexbin



Cartogram

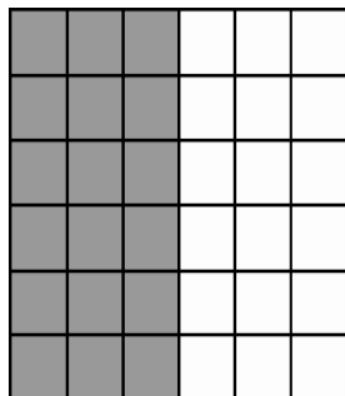
## Compute Moran's Value

Spatially Random



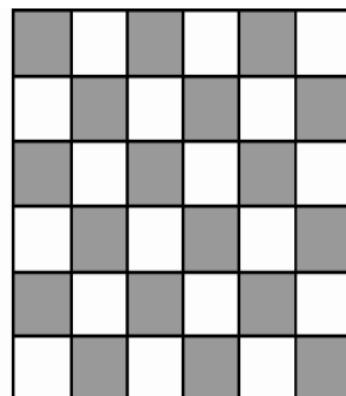
0

Clustered

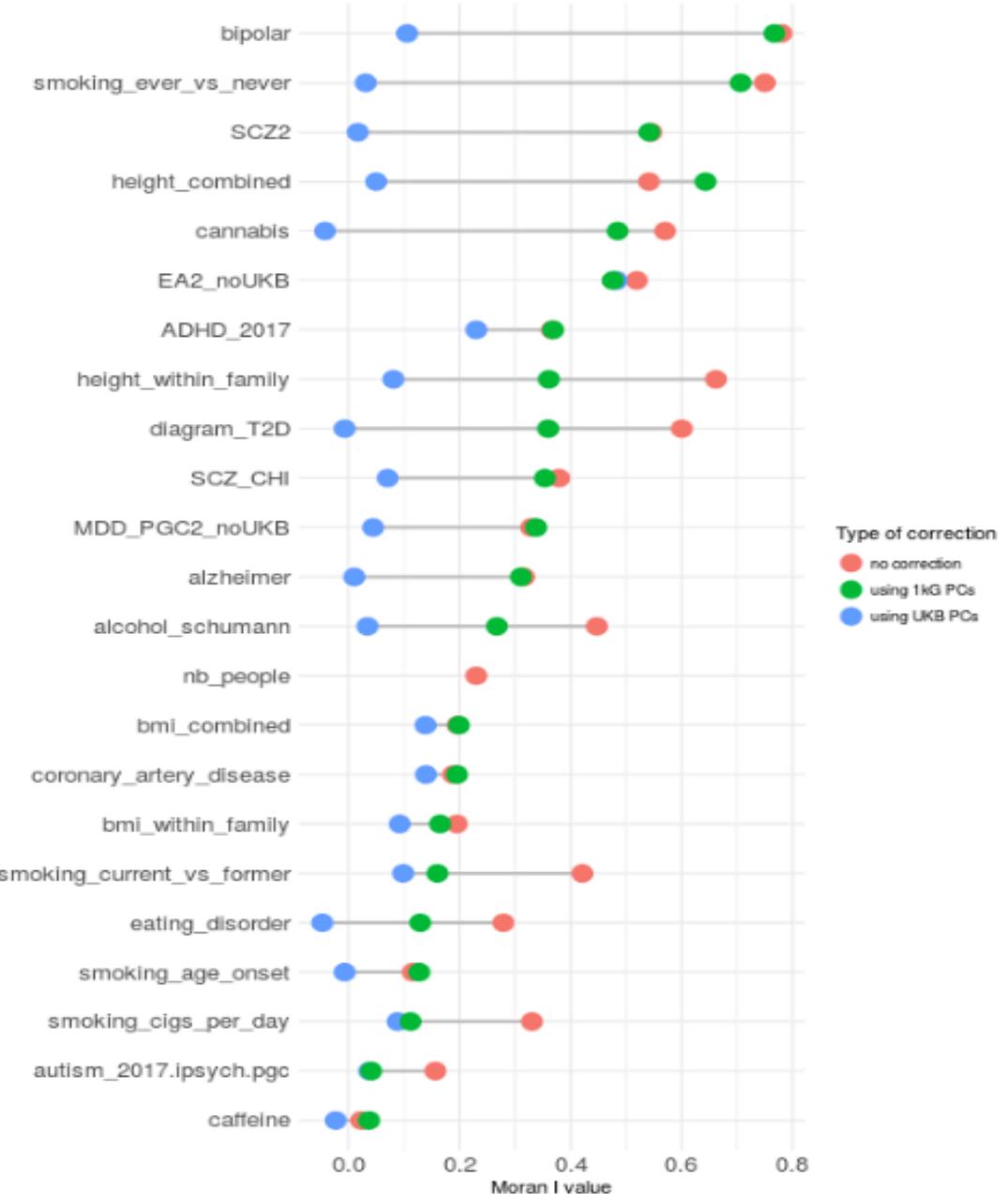


1

Dispersed



-1



Type of correction

- no correction
- using 1KG PCs
- using UKB PCs

# Genes & Geography in great britain

Methods

Explore

Compare

Load your data



Share what you see



Load your data

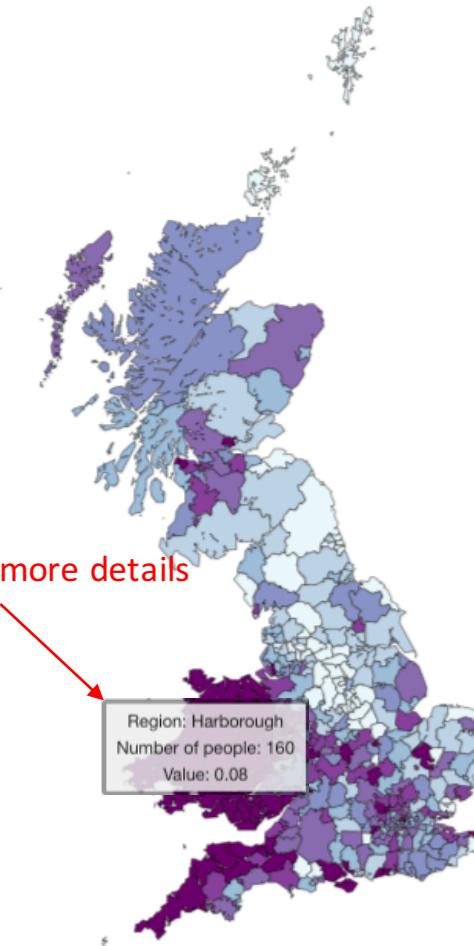
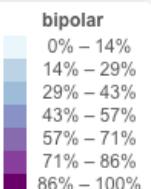
Custom the map

## Welcome

This application describes the geographical distribution of several variable of the [UK Biobank dataset](#) (n=502630).

More than 100 variables are available for visualization. You can observe them using different geographical units. It is possible to custom and export this map following the surrounding buttons. Use the compare tab above if you want to study the relationship between several variables.

Interactive map:  
Hover, zoom for more details



Allow re-utilization

Pick up your traits

bipolar

## Variable

We propose to represent the geographical distribution of 126 variables. These variable are split in several groups: Principal Components (PCs), Polygenic Risk Scores (PRS). To understand how these variable have been computed, visit the method section.

## Moran's I value

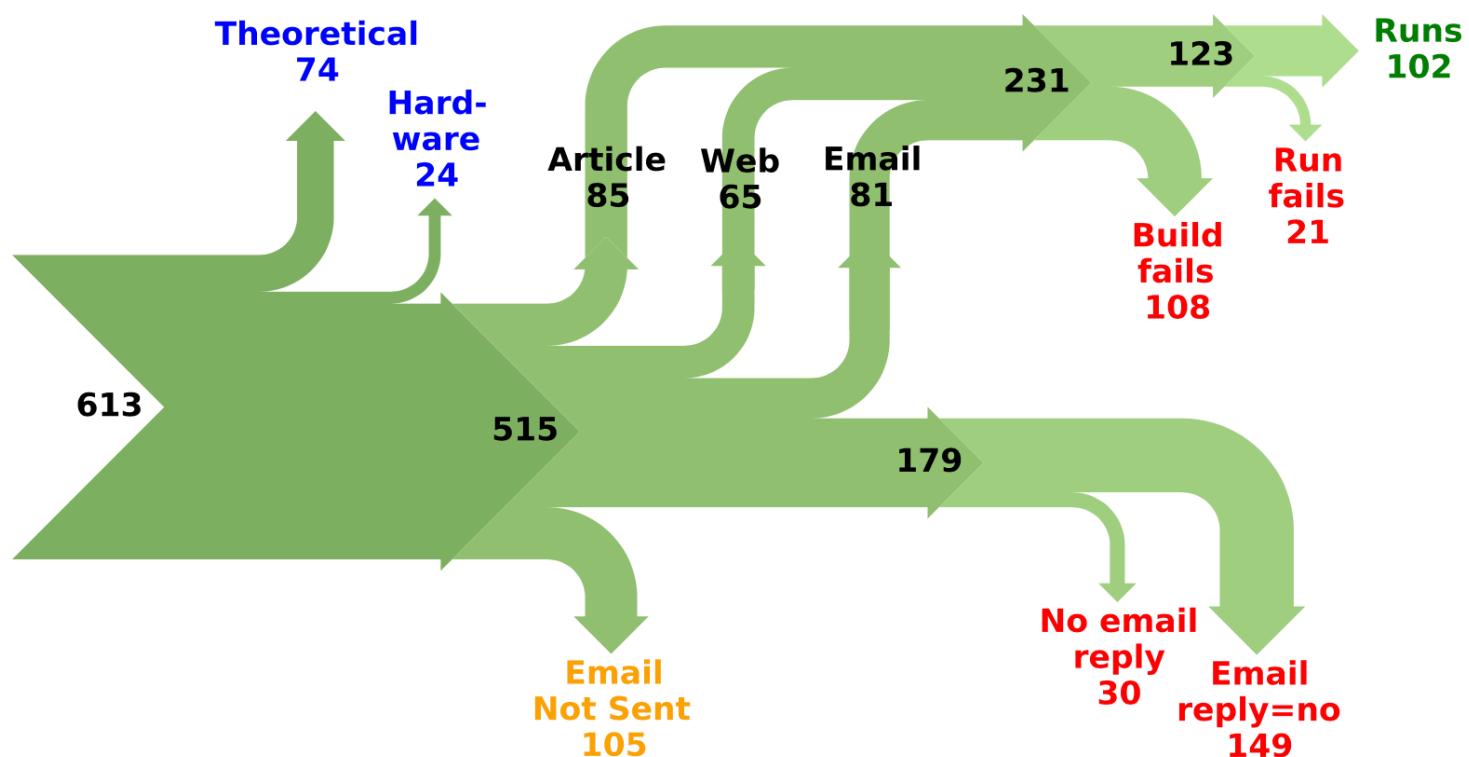
Moran's I value is a measure of spatial autocorrelation, i.e. a correlation in a signal among nearby locations in space. Click on the plus button below to see what variables are the most clustered



Figure 1: Geographical distribution of bipolar in the UK

# REPRODUCIBILITY IN DATA ANALYSIS

---



# An empirical analysis of journal policy effectiveness for computational reproducibility



Victoria Stodden, Jennifer Seiler and Zhaokun Ma

PNAS March 13, 2018; 115 (11) 2584-2589; published ahead of print March 1, 2018  
<https://doi.org/10.1073/pnas.1708290115>

Edited by David B. Allison, Indiana University Bloomington, Bloomington, Indiana, and approved by Susan T. Fiske January 9, 2018 (received for review July 11, 2017)

Article

Figures & SI

Authors & Info

**Table 1. Responses to emailed requests ( $n = 180$ )**

Type of response	Count	Percent, %
Did not share data or code:		
Contact another person	20	11
Asked for reasons	20	11
Refusal to share	12	7
Directed back to supplement	6	3
Unfulfilled promise to follow up	5	3
Impossible to share	3	2
Shared data and code	65	36
Email bounced	3	2
No response	46	26

# An empirical analysis of journal policy effectiveness for computational reproducibility



Victoria Stodden, Jennifer Seiler and Zhaokun Ma

PNAS March 13, 2018; 115 (11) 2584-2589; published ahead of print March 1, 2018  
<https://doi.org/10.1073/pnas.1708290115>

Edited by David B. Allison, Indiana University Bloomington, Bloomington, Indiana, and Susan T. Fiske, January 9, 2018 (received for review July 11, 2017)

Article Figures & SI Authors & Info

**Table 1. Responses to emailed requests ( $n = 180$ )**

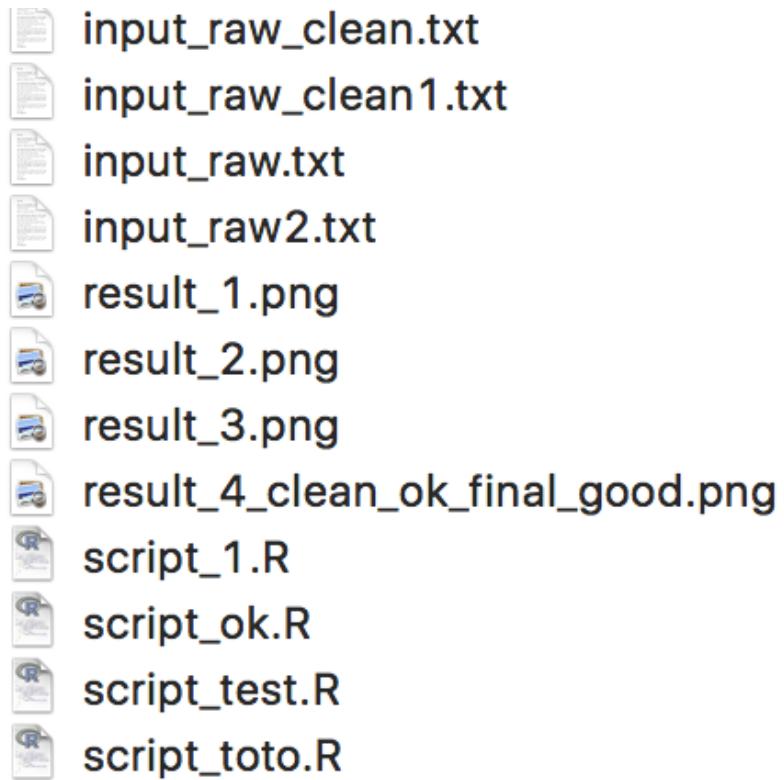
Type of response	Count	Percent, %
Did not share data or code:		
Contact another person	20	11
Asked for reasons	20	11
Refusal to share	12	7
Directed back to supplement	6	3
Unfulfilled promise to follow up	5	3
Impossible to share	3	2

When you approach a PI for the source codes and raw data, you better explain who you are, whom you work for, why you need the data and what you are going to do with it.

Share your data

# Can you run your analysis again ?

---



## A few best practices:

- Don't modify input files manually
- Don't copy and paste code / results
- Clean code (comment, use functions)
- Share your code
- Use version control

# R Markdown

from R Studio

```
53  
54 # An interactive manhattan plot | Title  
55 ***  
56  
57 Using `HTML` outputs you can embed some interactive graphics. For example, the  
plotly library can transform any of your ggplot2 graphic in an interactive  
chart:  
58  
59 ```{r, message=FALSE, warning=FALSE, echo=FALSE} | Text  
60 # Libraries  
61 library(plotly)  
62 library(tidyverse)  
63  
64 # Prepare the dataset  
65 don <- gwasResults %>%  
66 ...  
67
```

Code

# R Markdown

from R Studio

```
53  
54 # An interactive manhattan plot  
55 ***  
56  
57 Using `HTML` outputs you can embed some interactive graphics. For example, the  
plotly library can transform any of your ggplot2 graphic in an interactive  
chart:  
58  
59 ```{r, message=FALSE, warning=FALSE, echo=FALSE}  
60 # Libraries  
61 library(plotly)  
62 library(tidyverse)  
63  
64 # Prepare the dataset  
65 don <- gwasResults %>%  
66  
67
```

Title

Code

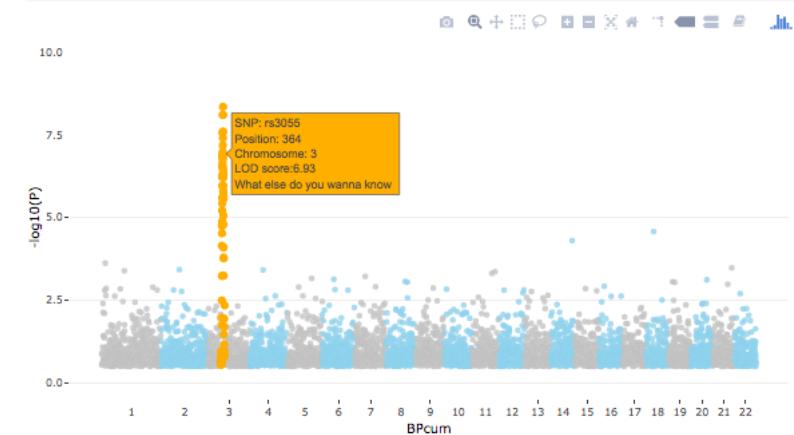
Text



## 3 An interactive manhattan plot

Using `HTML` outputs you can embed some interactive graphics. For example, the `plotly` library can transform any of your `ggplot2` graphic in an interactive chart:

```
# Make the plot  
p <- ggplot(don, aes(x=BPcum, y=-log10(P), text=text)) +  
  
  # Show all points  
  geom_point( aes(color=as.factor(CHR)), alpha=0.8, size=1.3) +  
  scale_color_manual(values = rep(c("grey", "skyblue"), 22 )) +  
  
  # custom X axis:  
  scale_x_continuous( label = axisdf$CHR, breaks= axisdf$center ) +  
  scale_y_continuous(expand = c(0, 0) ) +      # remove space between plot area and x axis  
  
  # Add highlighted points  
  geom_point(data=subset(don, is_highlight=="yes"), color="orange", size=2) +  
  
  # Custom the theme:  
  theme_bw() +  
  theme(  
    legend.position="none",  
    panel.border = element_blank(),  
    panel.grid.major.x = element_blank(),  
    panel.grid.minor.x = element_blank()  
  ) +  
  ylim(0, 10)  
  
ggplotly(p, tooltip="text")
```

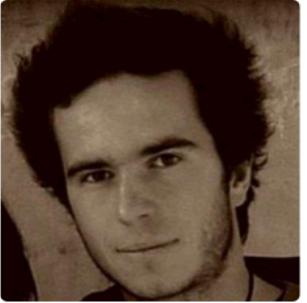


Is your code reusable ?

Interactive report



# GitHub



**Holtz Yan**  
holtzy

Data Analyst with a + in data visualization

[Edit bio](#)

**Queensland Brain Institute**  
Brisbane, Australia  
[yan.holtz.data@gmail.com](mailto:yan.holtz.data@gmail.com)  
<https://holtzyan.wordpress.com/>

**452 contributions in the last year**



Learn how we count contributions.

**Overview**   [Repositories 28](#)   [Stars 22](#)   [Followers 88](#)

**Pinned repositories**

- The-Python-Graph-Gallery**  
A website displaying hundreds of charts made with Python  
Python ★ 122 ⚡ 15
- R-Graph-Gallery**  
A website displaying hundreds of charts made with R  
R ★ 122 ⚡ 15
- Pimp-my-rmd**  
A few tips about R markdown  
HTML ★ 11 ⚡ 2
- epu**  
A clear and effective presentation of your research  
HTML ★ 11 ⚡ 2
- GenMap-Comparator**  
An application to compare genetic maps with D3 & Shiny  
R ★ 8 ⚡ 5
- Pulse**  
Reproductive resistance to climate change  
HTML ★ 8 ⚡ 5

**452 contributions in the last year**



Learn how we count contributions.

**holtzy / the-NB-COMO-Project**

[Code](#)   [Issues 0](#)   [Pull requests 0](#)   [Projects 0](#)   [Wiki](#)   [Insights](#)   [Settings](#)

A Shiny app describing comorbidity in the Danish Health Register

[Edit](#)

**Add topics**

20 commits   1 branch   0 releases   1 contributor

Branch: master   [New pull request](#)   [Create new file](#)   [Upload files](#)   [Find file](#)   [Clone or download](#)

holtzy Change 2 words	Latest commit e1a433e 21 days ago
DATA	Change 2 words
rsconnect/shinyapps.io/holtzyan	Change 2 words
www	John Correction
.DS_Store	John Correction
.Rhistory	First commit
README.md	Last version before sending to collaborators
global.R	Change 2 words
server.R	Change 2 words
ui.R	Change 2 words

**README.md**

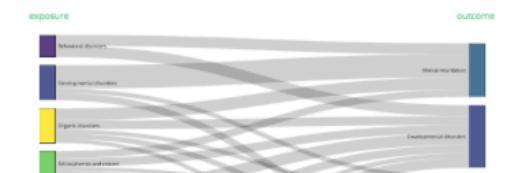
## The NB-COMO Project

### Overview

The NB-BOMO project aims to explore the patterns of comorbidity within treated mental disorders. It explores different ways to capture the complex patterns of comorbidity, notably through data visualization techniques. The first part of this project explores COMO within the [Danish National Patient Registry \(DNPR\)](#), one of the world's oldest nationwide hospital registries.

This repository gives the code of a web application that allows to interactively explore our results. It goes along with our peer reviewed publication (work in progress).

Here is a screenshot of one of the multiple visualizations proposed in the website:



## Acknowledgment

---

John McGrath

Naomi Wray  
Peter Visscher  
Jian Yang

Oleguer Plana-Ripoll  
Abdel Abdellaoui

## Contact

---



@R\_Graph\_Gallery



[github.com/holtzy/Talk](https://github.com/holtzy/Talk)



[Yan.holtz.data@gmail.com](mailto:Yan.holtz.data@gmail.com)



R-graph-gallery.com





## The Scientific Paper Is Obsolete

Here's what's next.

[theatlantic.com](http://theatlantic.com)

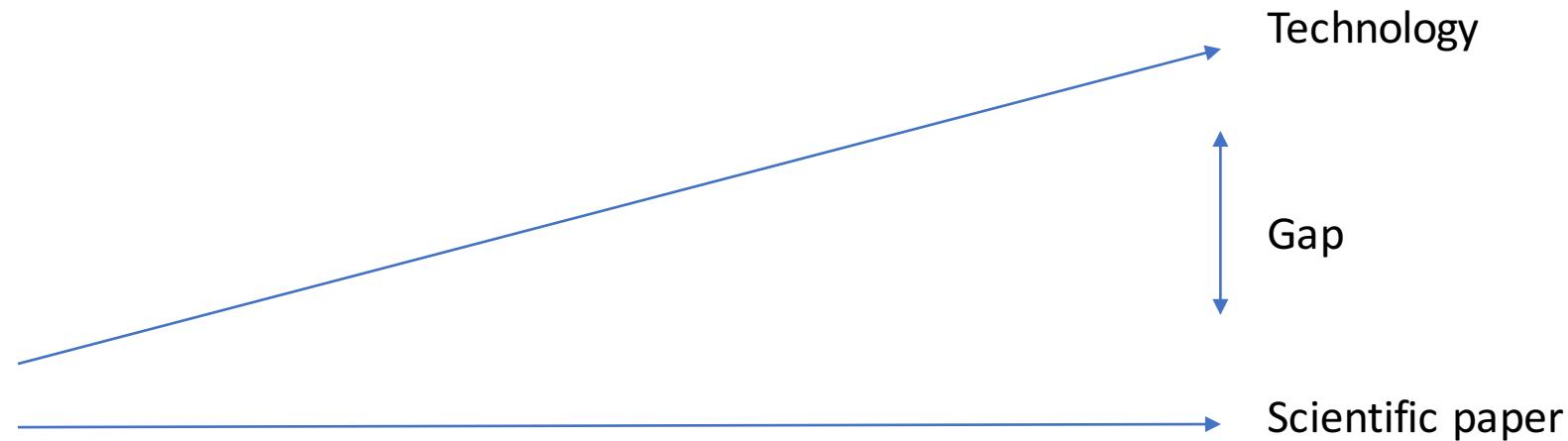
# Observable is a better way to code.

Discover insights faster and communicate more effectively with interactive notebooks for data analysis, visualization, and exploration.

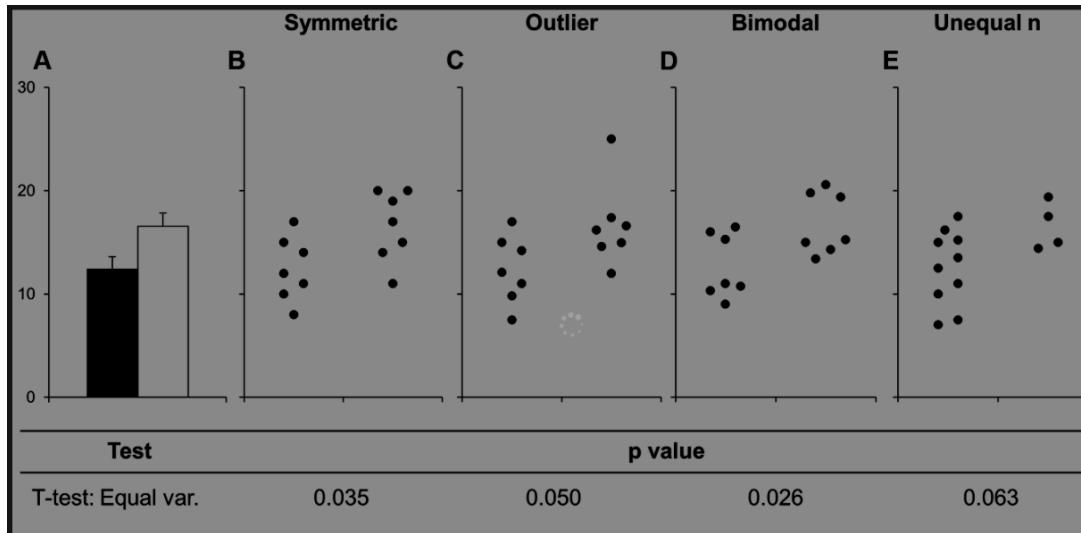
[Sign up for free](#)

[Try the scratchpad](#)

<https://beta.observablehq.com>



# Data viz: Avoid mistake and misleading figures



Show violin plot average odd ratio per group

**PLOS** | BIOLOGY  
FIFTEENTH ANNIVERSARY

Browse | Publish | About | Search | advanced search

OPEN ACCESS

PERSPECTIVE

## Beyond Bar and Line Graphs: Time for a New Data Presentation Paradigm

Tracey L. Weissgerber Natasa M. Milic, Stacey J. Winham, Vesna D. Garovic

Published: April 22, 2015 • <https://doi.org/10.1371/journal.pbio.1002128>

Article	Authors	Metrics	Comments	Related Content
▼				

### Abstract

#### Introduction

Are Your Figures Worth a Thousand Words?

Summary Statistics Are Only Meaningful When There Are Enough Data to Summarize

Recommendations for a New Data Presentation Paradigm

#### Conclusions

#### Supporting Information

### Abstract

Figures in scientific publications are critically important because they often show the data supporting key findings. Our systematic review of research articles published in top physiology journals ( $n = 703$ ) suggests that, as scientists, we urgently need to change our practices for presenting continuous data in small sample size studies. Papers rarely included scatterplots, box plots, and histograms that allow readers to critically evaluate continuous data. Most papers presented continuous data in bar and line graphs. This is problematic, as many different data distributions can lead to the same bar or line graph. The full data may suggest different conclusions from the summary statistics. We recommend training investigators in data presentation, encouraging a more complete presentation of data, and changing journal editorial policies. Investigators can quickly make univariate scatterplots for small sample size studies using our Excel templates.

2,454 Save	106 Citation
281,413 View	7,110 Share

Download PDF

Print

Share

Check for updates

### Included in the Following Collections

Open Data  
Meta-Research: Reporting

ADVERTISEMENT

Subject Areas

