

Reproducible research is still a challenge

Rich FitzJohn  | Matt Pennell | Amy Zanne | Will Cornwell

June 9, 2014

why is this still hard?

complexity

why is this still hard?

tooling

why is this still hard?

incentives + motivation

nobody cares

(Apart from people doing it out of the goodness of their hearts of course)

nobody cares?

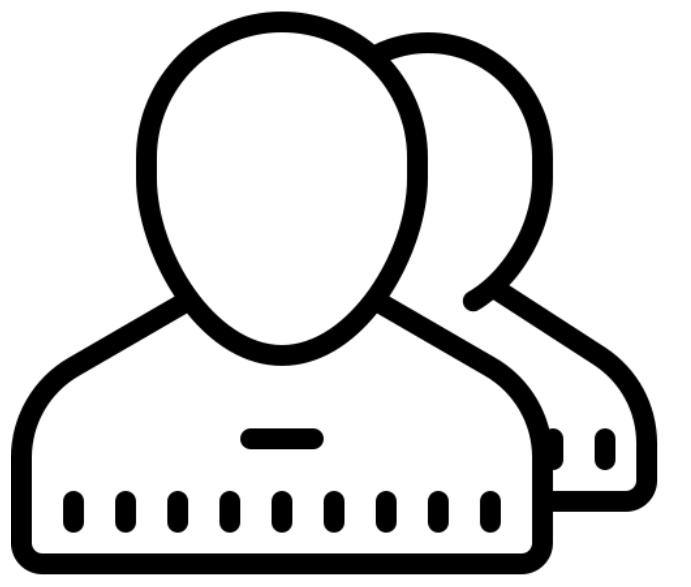
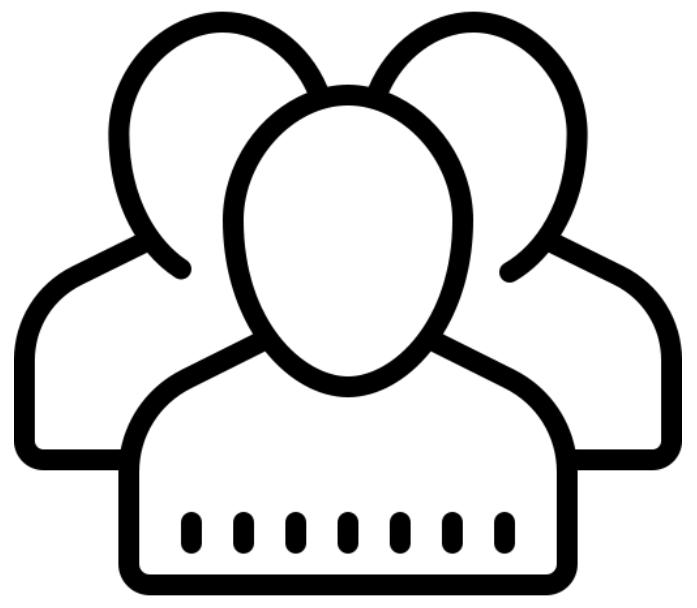
(Apart from people doing it out of the goodness of their hearts of course)



VACCINE IMPACT
MODELLING CONSORTIUM

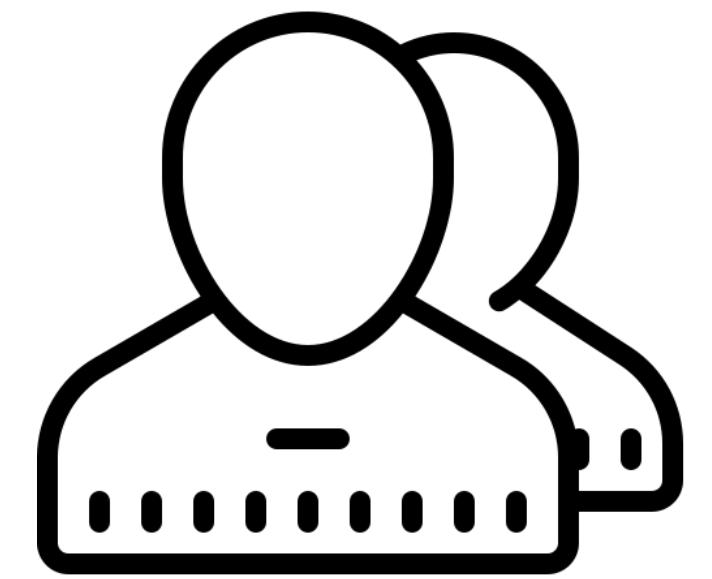
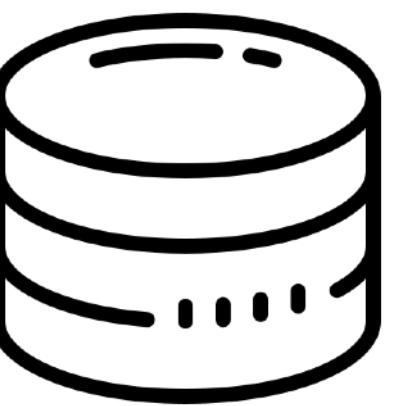
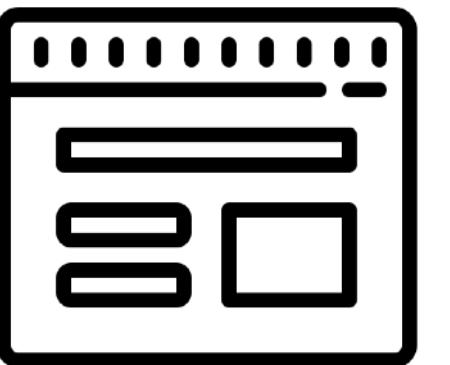
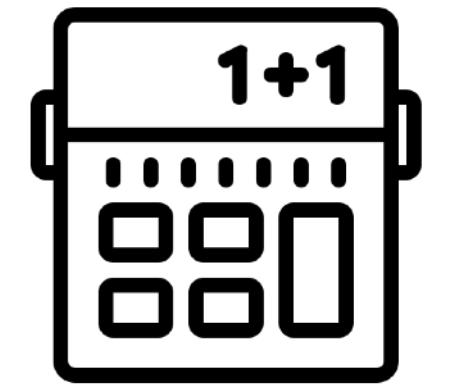
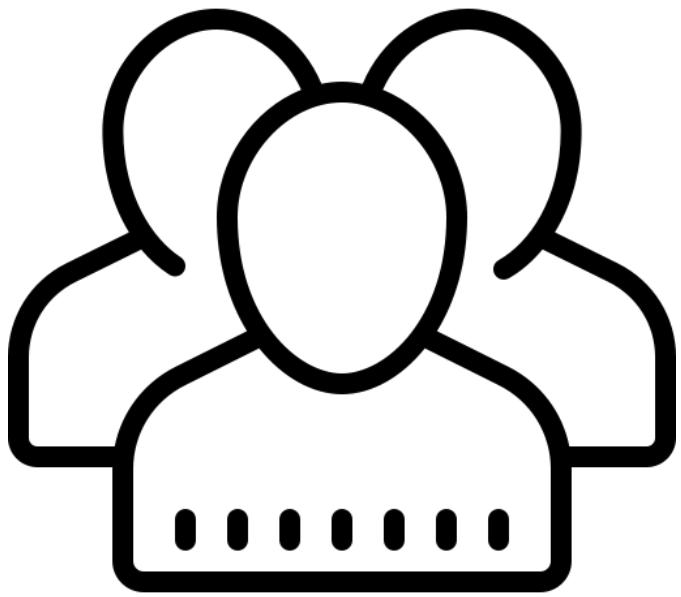
BILL & MELINDA
GATES foundation





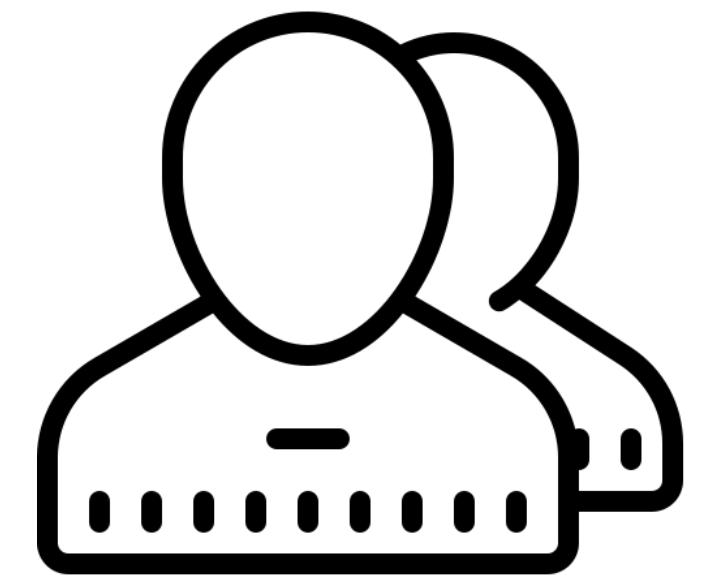
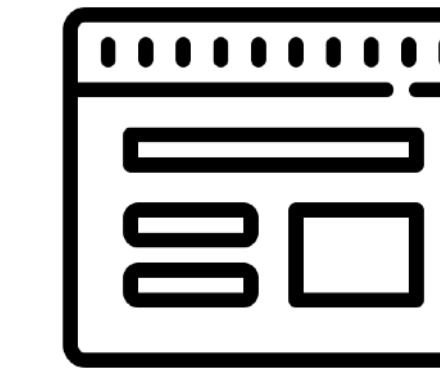
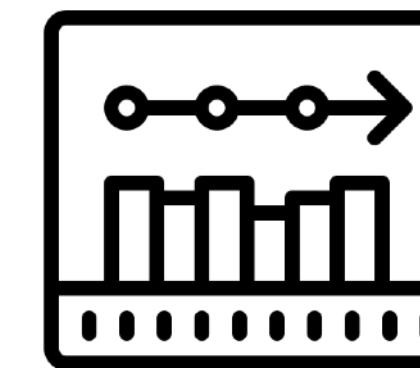
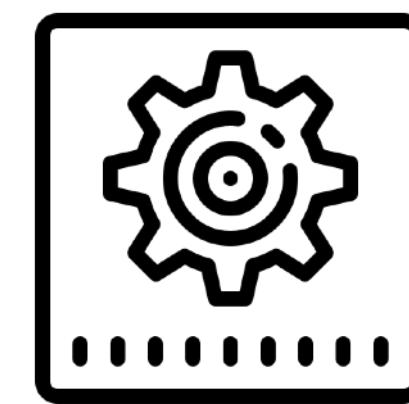
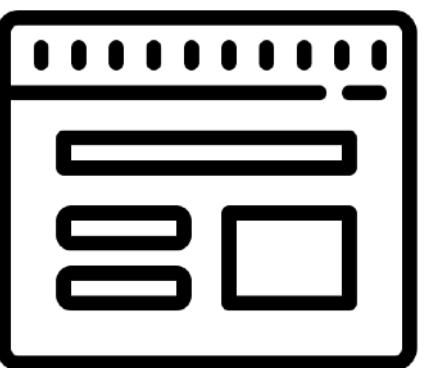
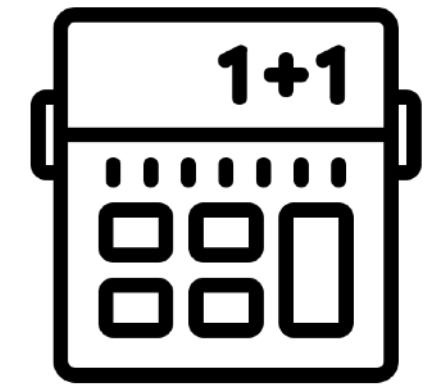
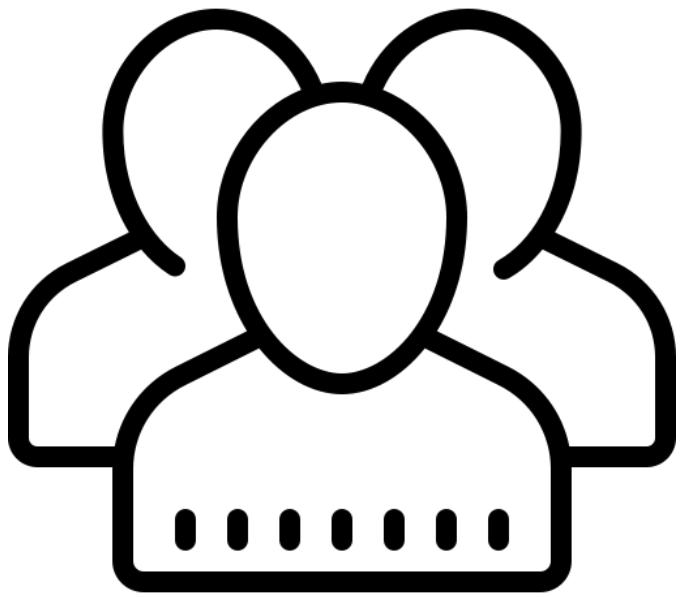
BILL & MELINDA
GATES foundation





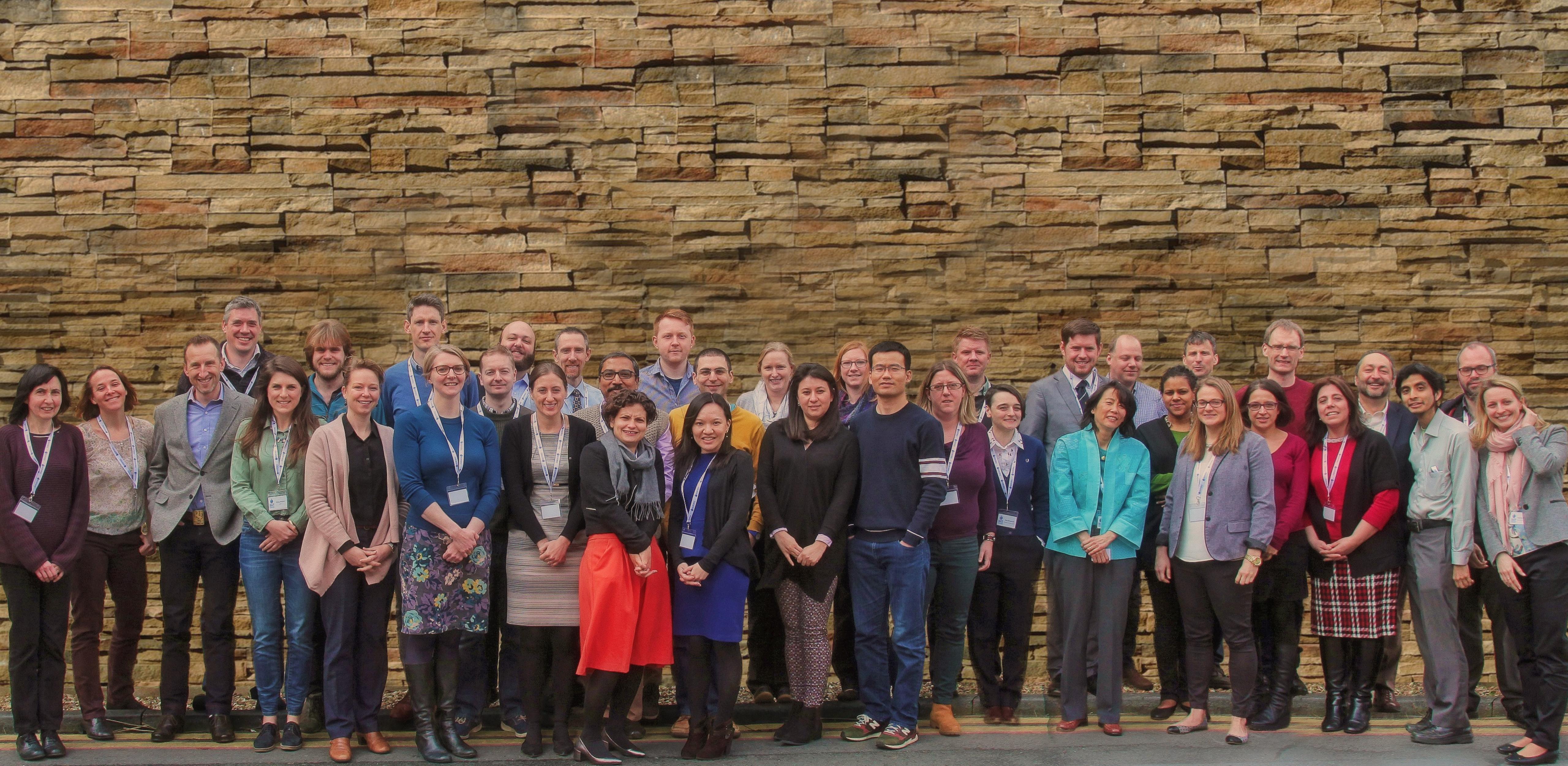
BILL & MELINDA
GATES foundation

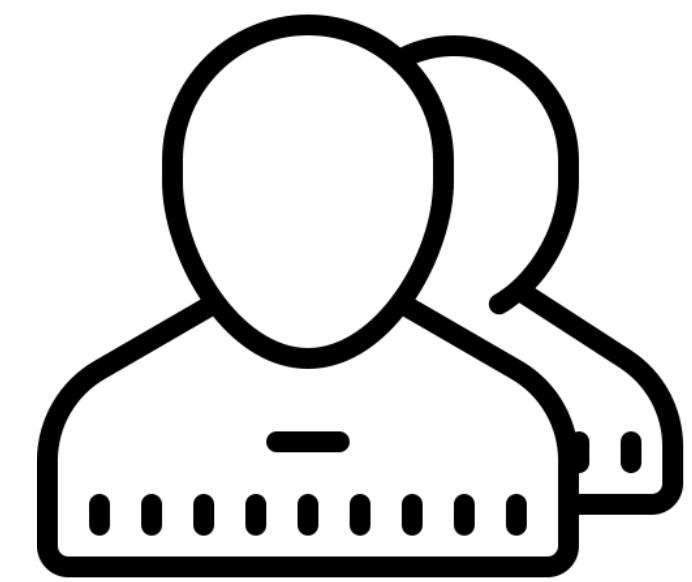
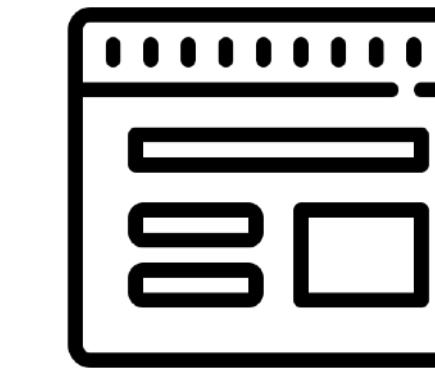
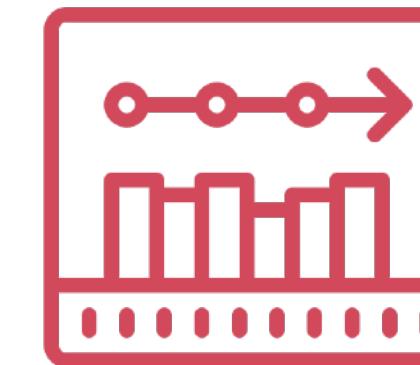
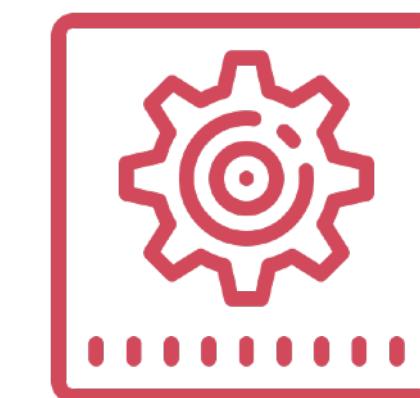
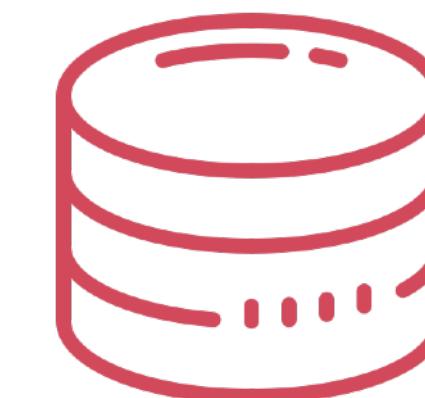
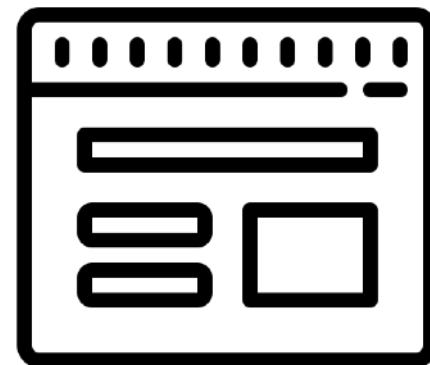
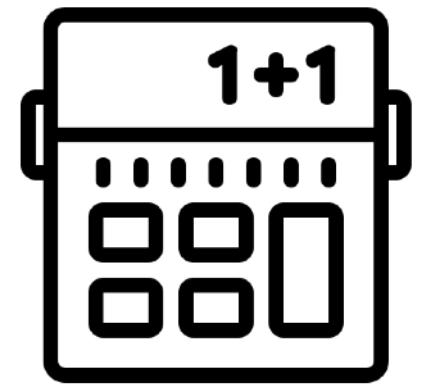
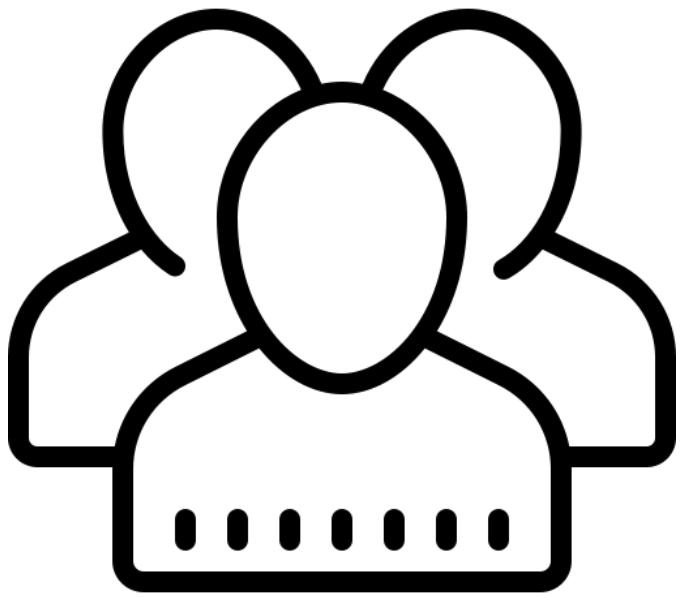




BILL & MELINDA
GATES foundation







BILL & MELINDA
GATES foundation



Our goal

Code written by domain scientists
runs on multiple machines, first time

Trace back to origins of changes
in any report



knitr

www.rstudio.com



rmarkdown

www.rstudio.com

RStudio

example.Rmd x

ABC Knit Insert Run

```
1 ---  
2 title: "Example"  
3 output: html_document  
4 ---  
5 ```{r setup, include=FALSE}  
6 knitr::opts_chunk$set(echo = TRUE)  
7 ```  
8  
9  
10 ## R Markdown  
11  
12 This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.  
13  
14 When you click the Knit button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:  
15 ```{r cars}  
16 summary(cars)  
17 ```  
18  
19 ## Including Plots  
20  
21 You can also embed plots, for example:  
22  
23 ```{r pressure, echo=FALSE}  
24 plot(pressure)  
25 ```  
26  
27 Note that the `echo = FALSE` parameter was added to the code  
28 chunk to prevent printing of the R code that generated the  
29 plot.
```

2:16 # Example R Markdown

Console

example.html Open in Browser Find

~/example.html

Example

R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

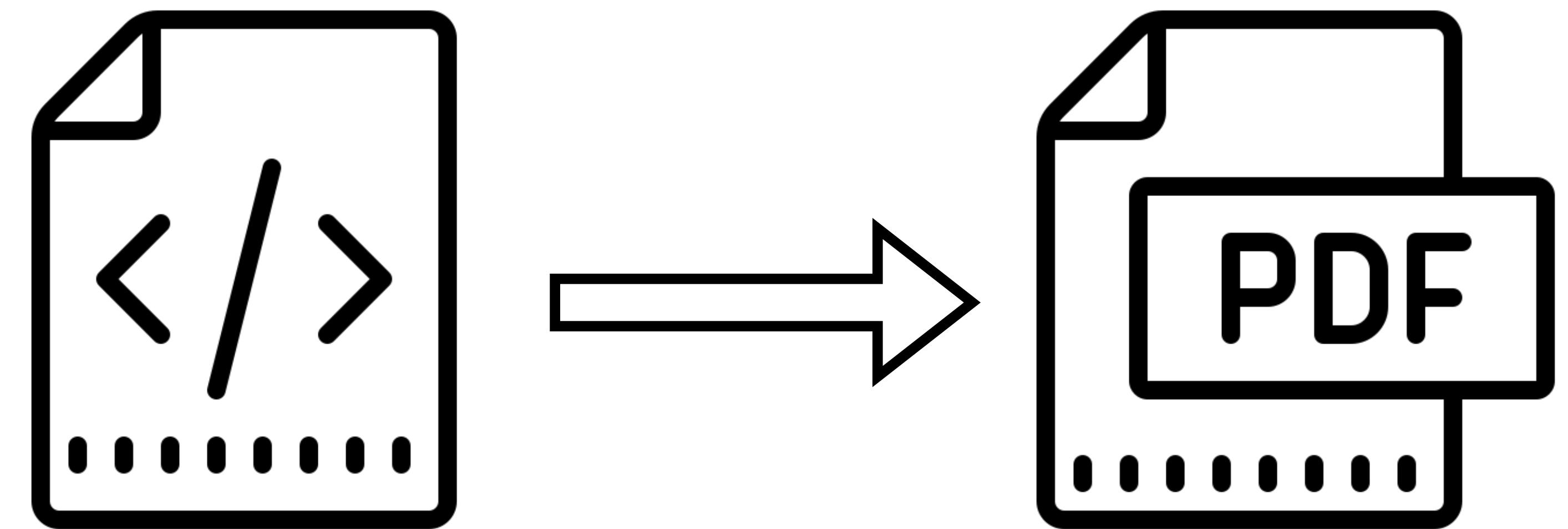
```
summary(cars)
```

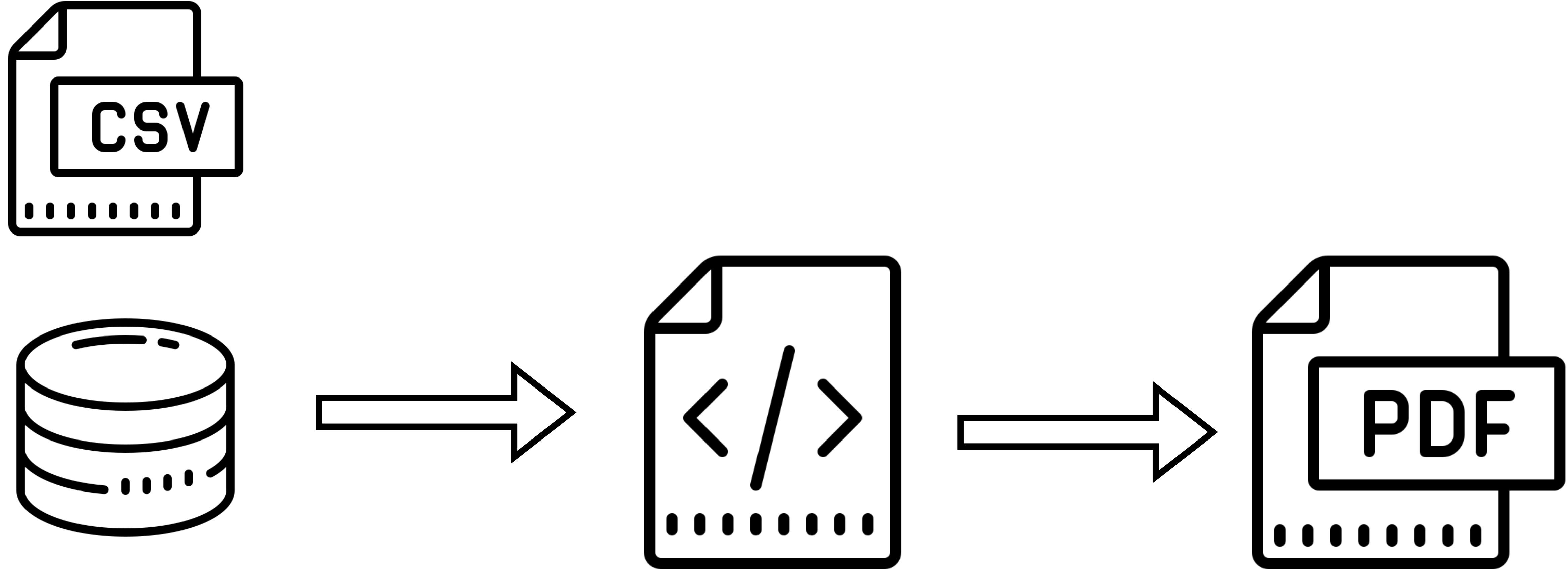
```
##      speed          dist  
##  Min.   : 4.0   Min.   : 2.00  
##  1st Qu.:12.0   1st Qu.: 26.00  
##  Median :15.0   Median : 36.00  
##  Mean   :15.4   Mean   : 42.98  
##  3rd Qu.:19.0   3rd Qu.: 56.00  
##  Max.   :25.0   Max.   :120.00
```

Including Plots

You can also embed plots, for example:

A scatter plot showing the relationship between pressure and distance. The x-axis is labeled "pressure" and ranges from approximately 400 to 800. The y-axis ranges from 400 to 800. There are three data points plotted at approximately (400, 400), (600, 600), and (800, 800).







drake

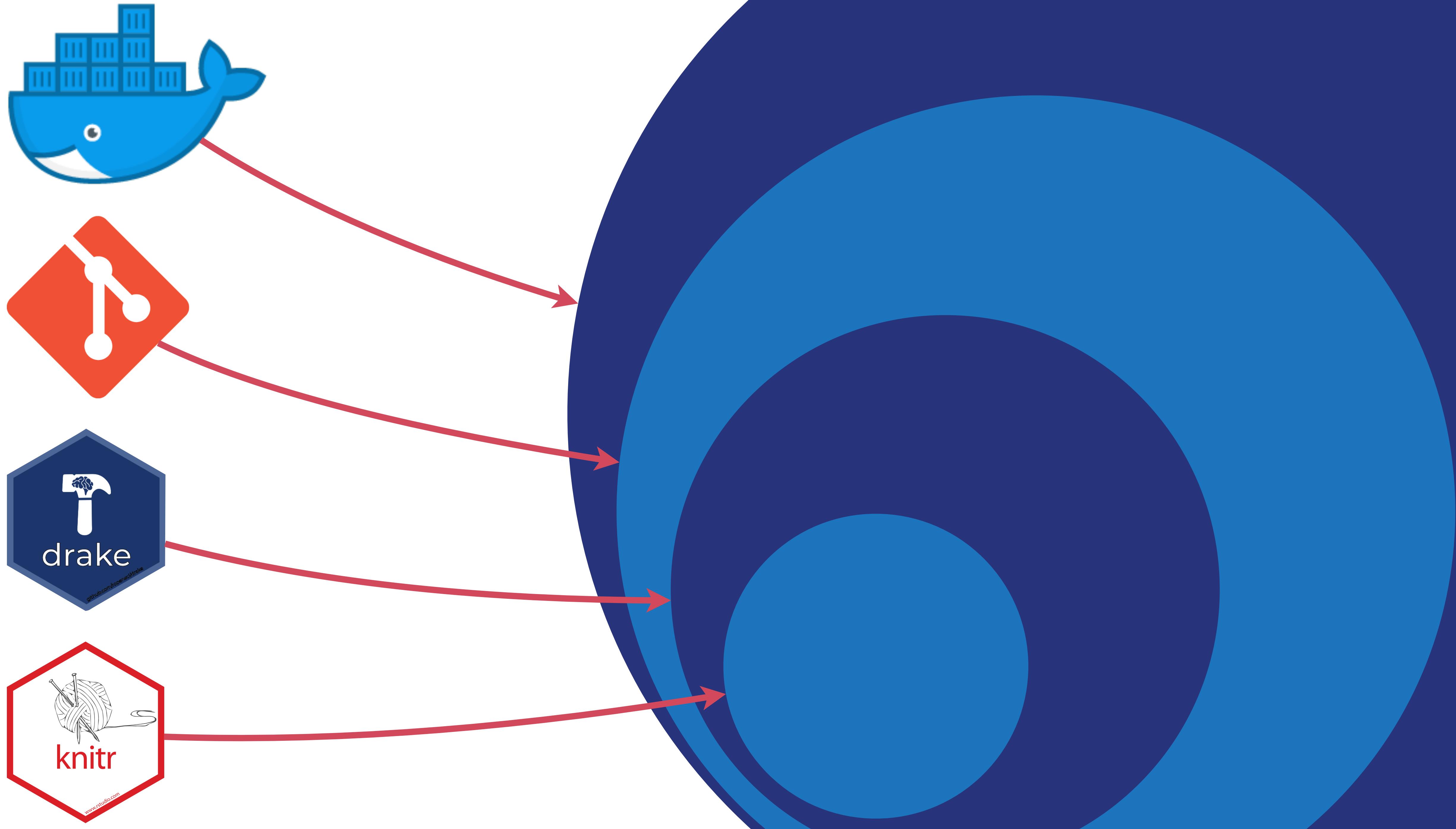
github.com/openscience/drake



git



docker





Go To Statement Considered Harmful

Key Words and Phrases: go to statement, jump instruction, branch instruction, conditional clause, alternative clause, repetitive clause, program intelligibility, program sequencing

CR Categories: 4.22, 5.23, 5.24

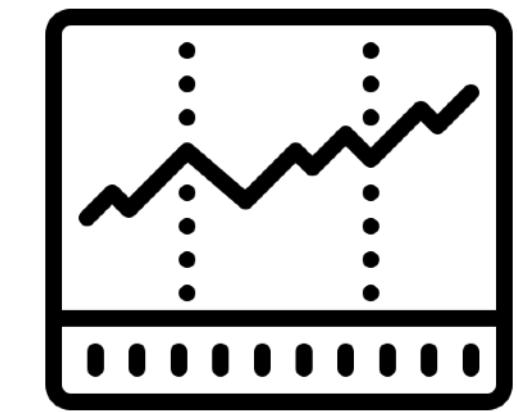
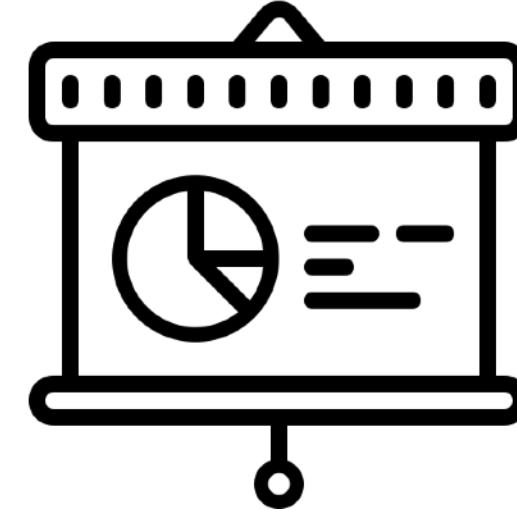
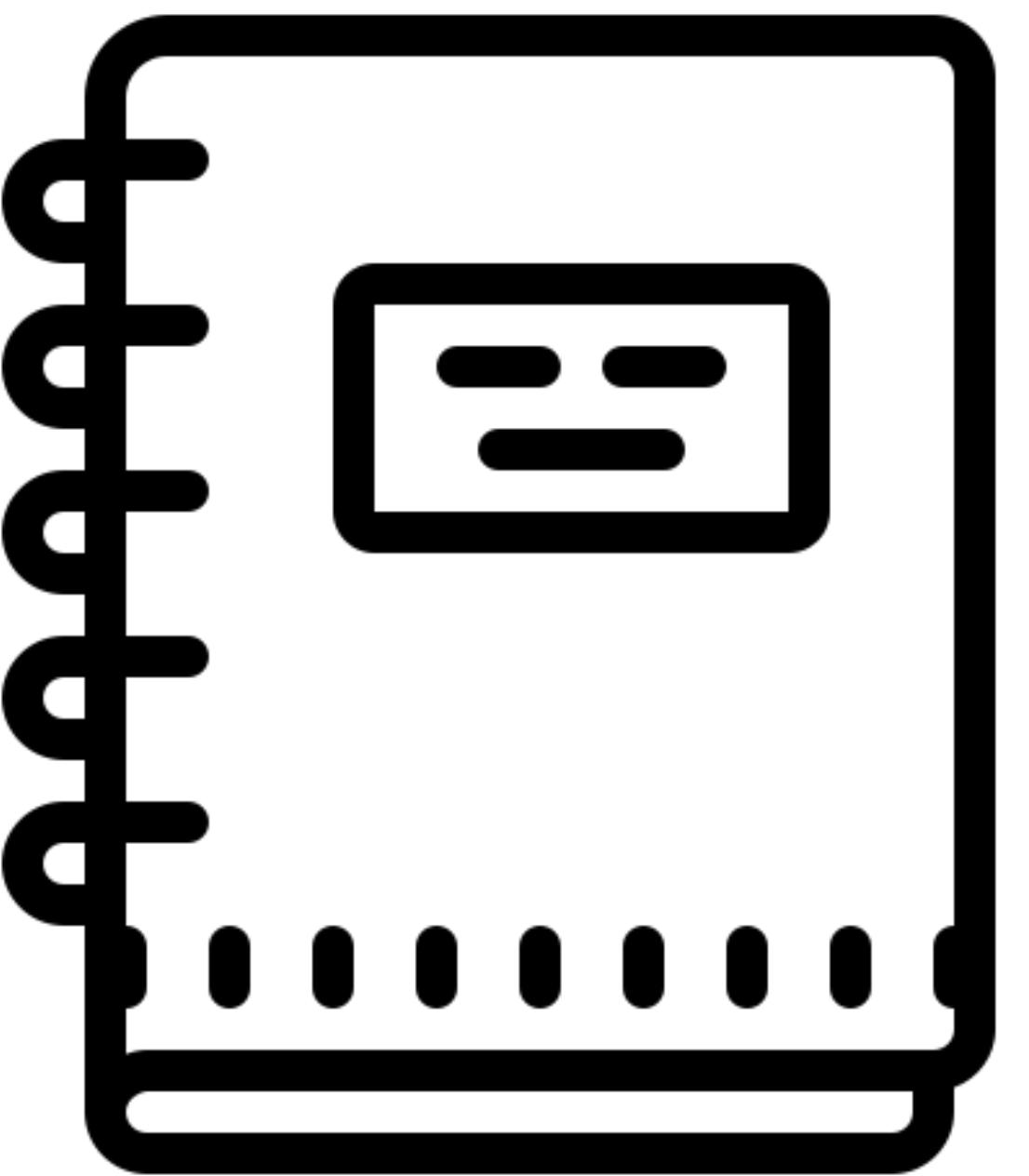
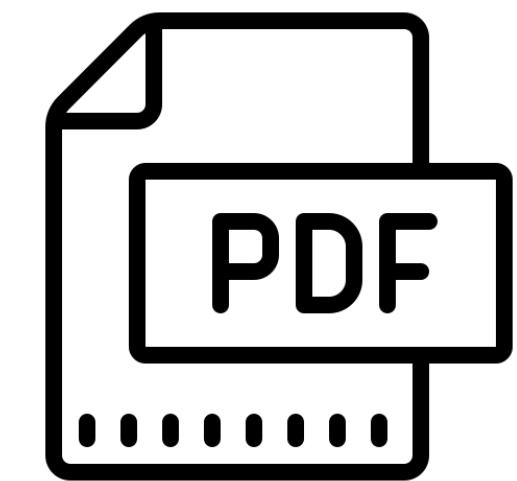
EDITOR:

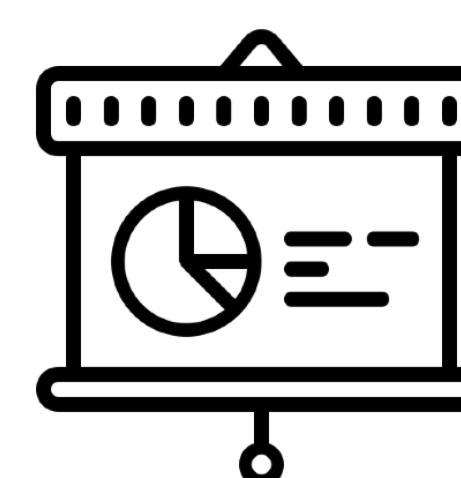
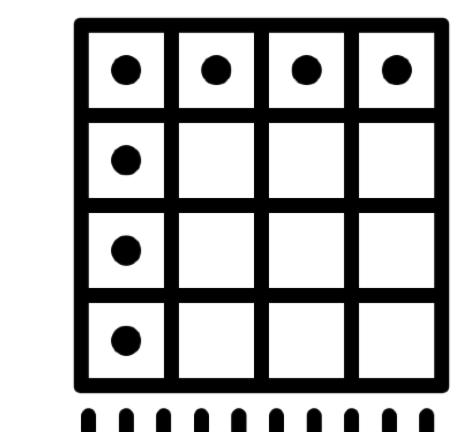
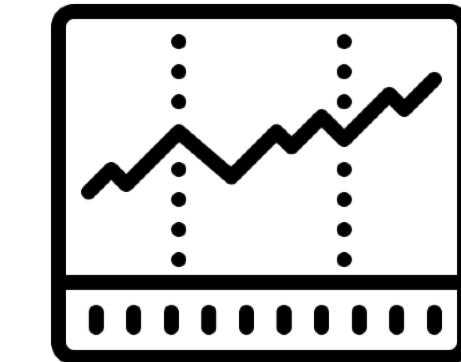
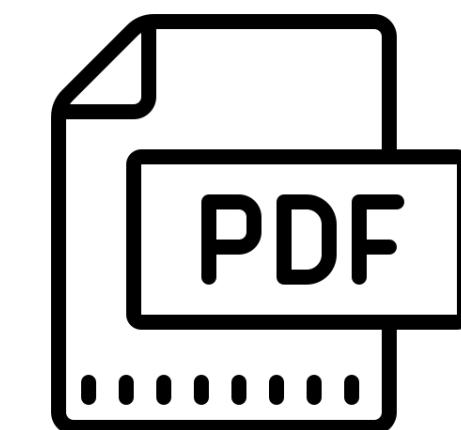
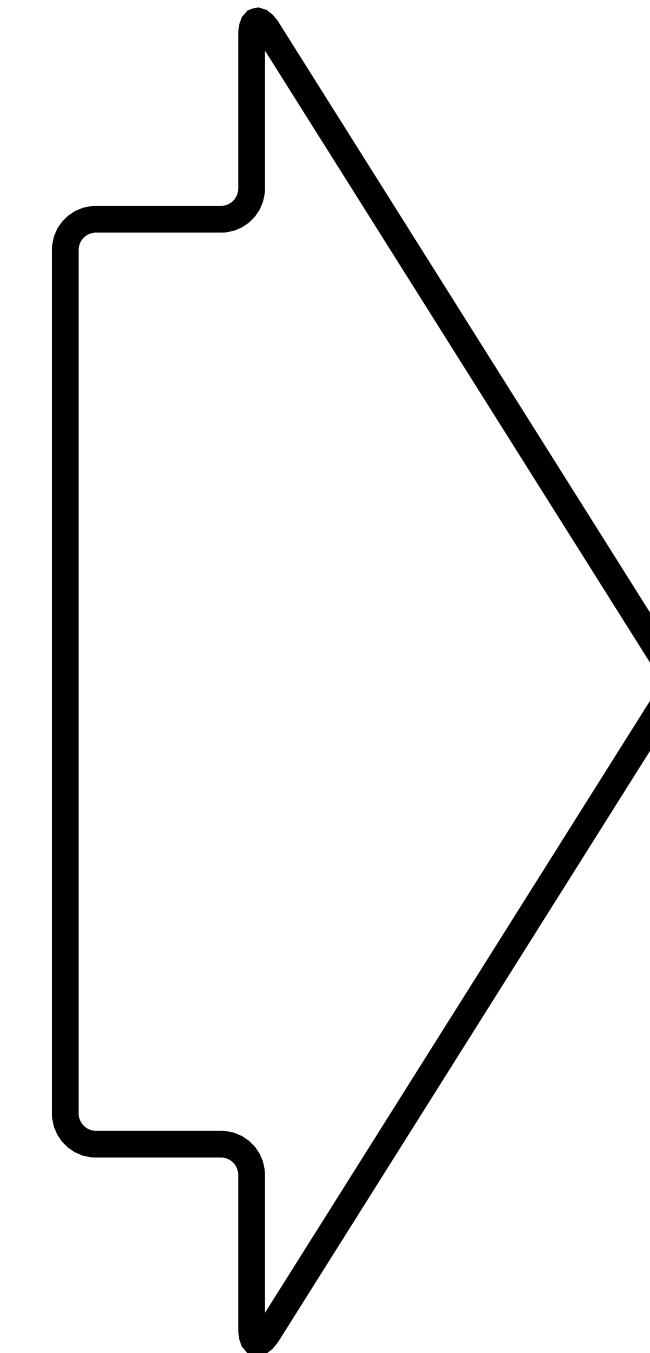
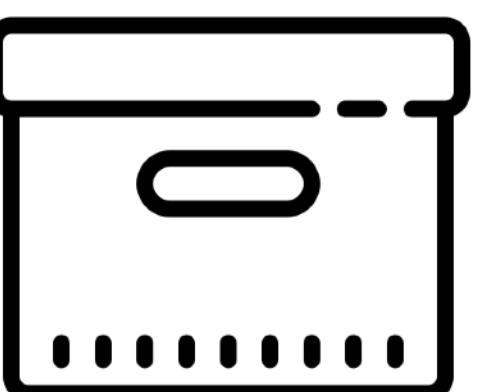
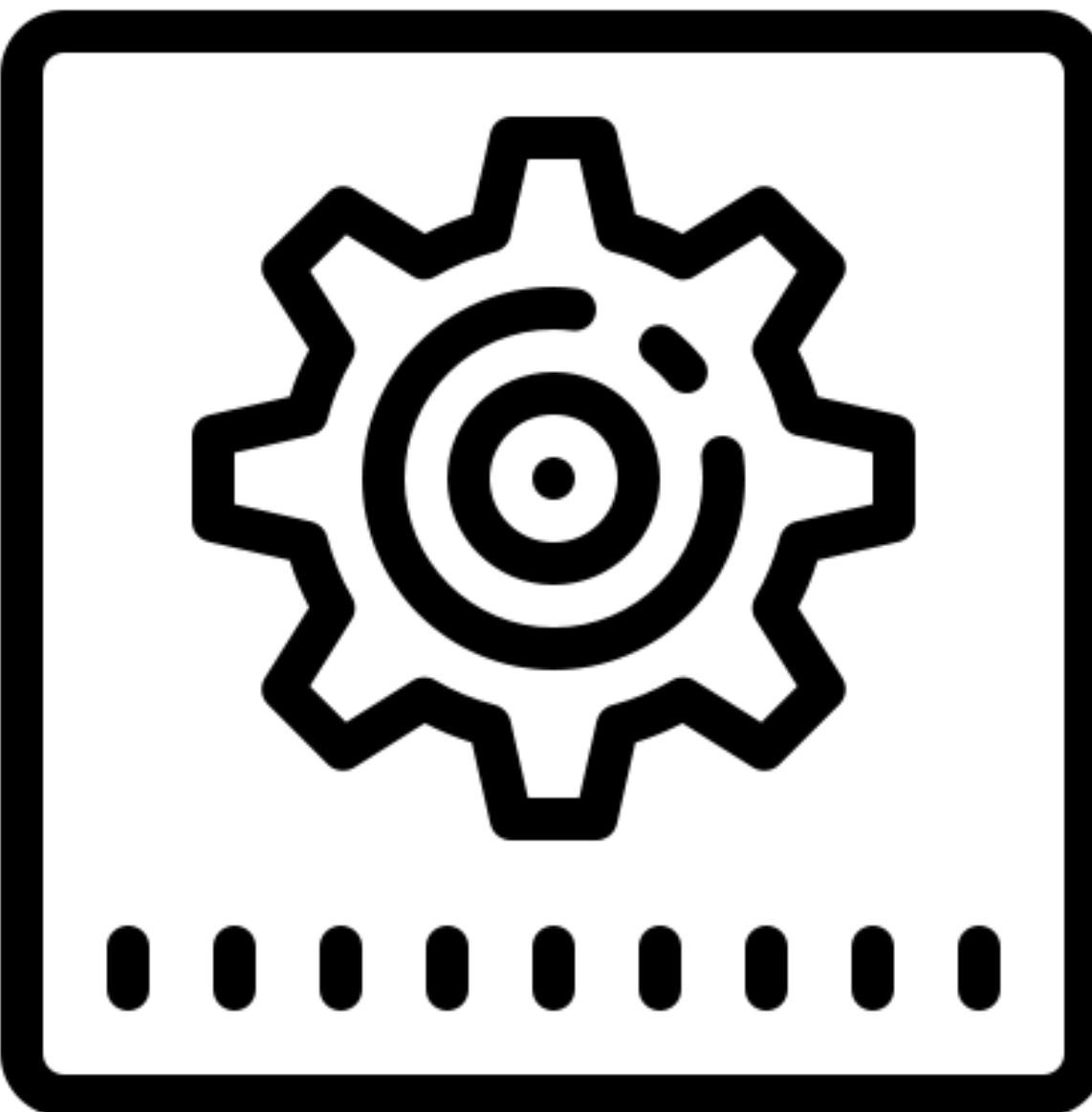
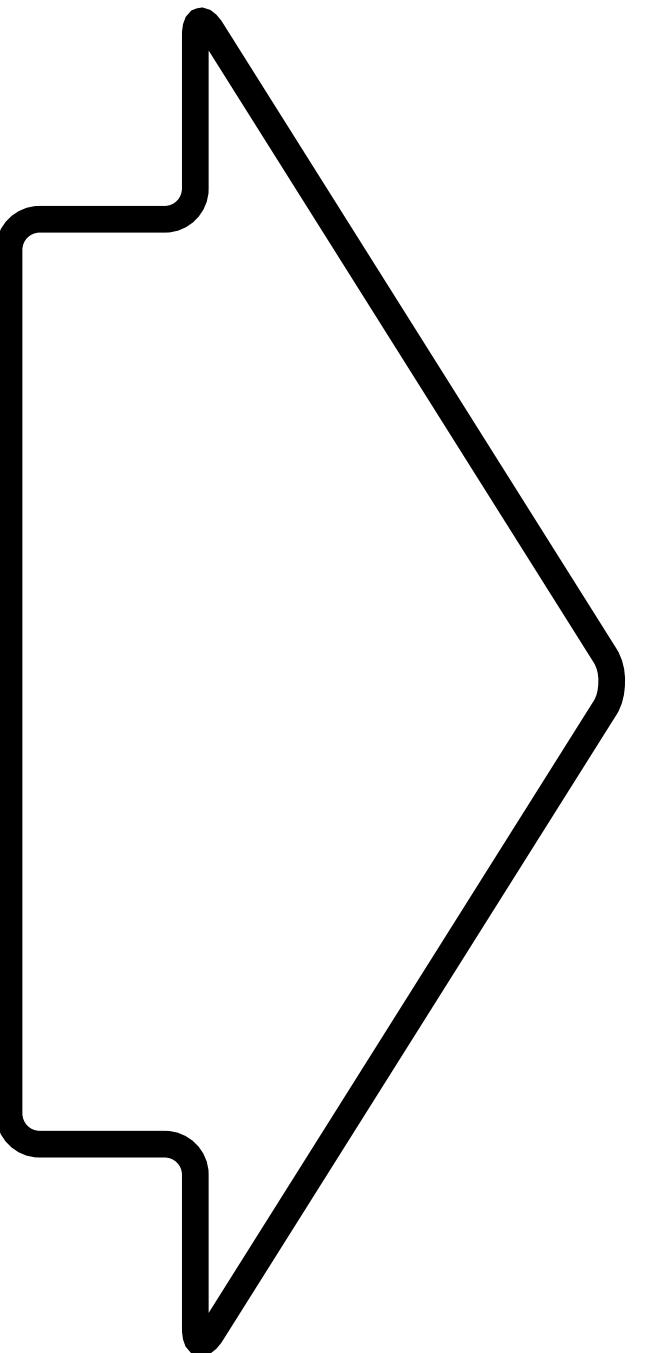
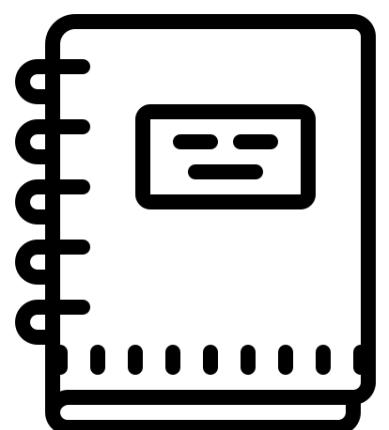
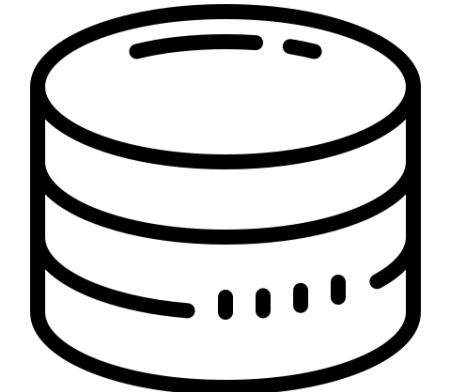
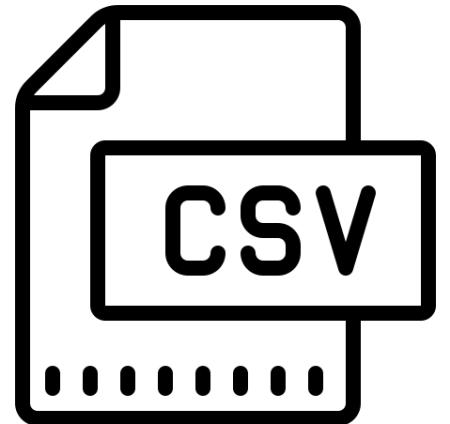
For a number of years I have been familiar with the observation that the quality of programmers is a decreasing function of the density of **go to** statements in the programs they produce. More recently I discovered why the use of the **go to** statement has such disastrous effects, and I became convinced that the **go to** statement should be abolished from all "higher level" programming languages (i.e. everything except, perhaps, plain machine code).

```
add <- function(a, b) {  
  a + b  
}
```

```
expect_equal(  
    add(1, 3),  
    4)
```

```
#' @param a,b inputs
#' @return the sum
#' @export
```

A black and white line-art icon of a 5x5 grid. Some of the cells in the grid are filled with small black dots. The bottom of the icon features a decorative border with small circles.



orderly

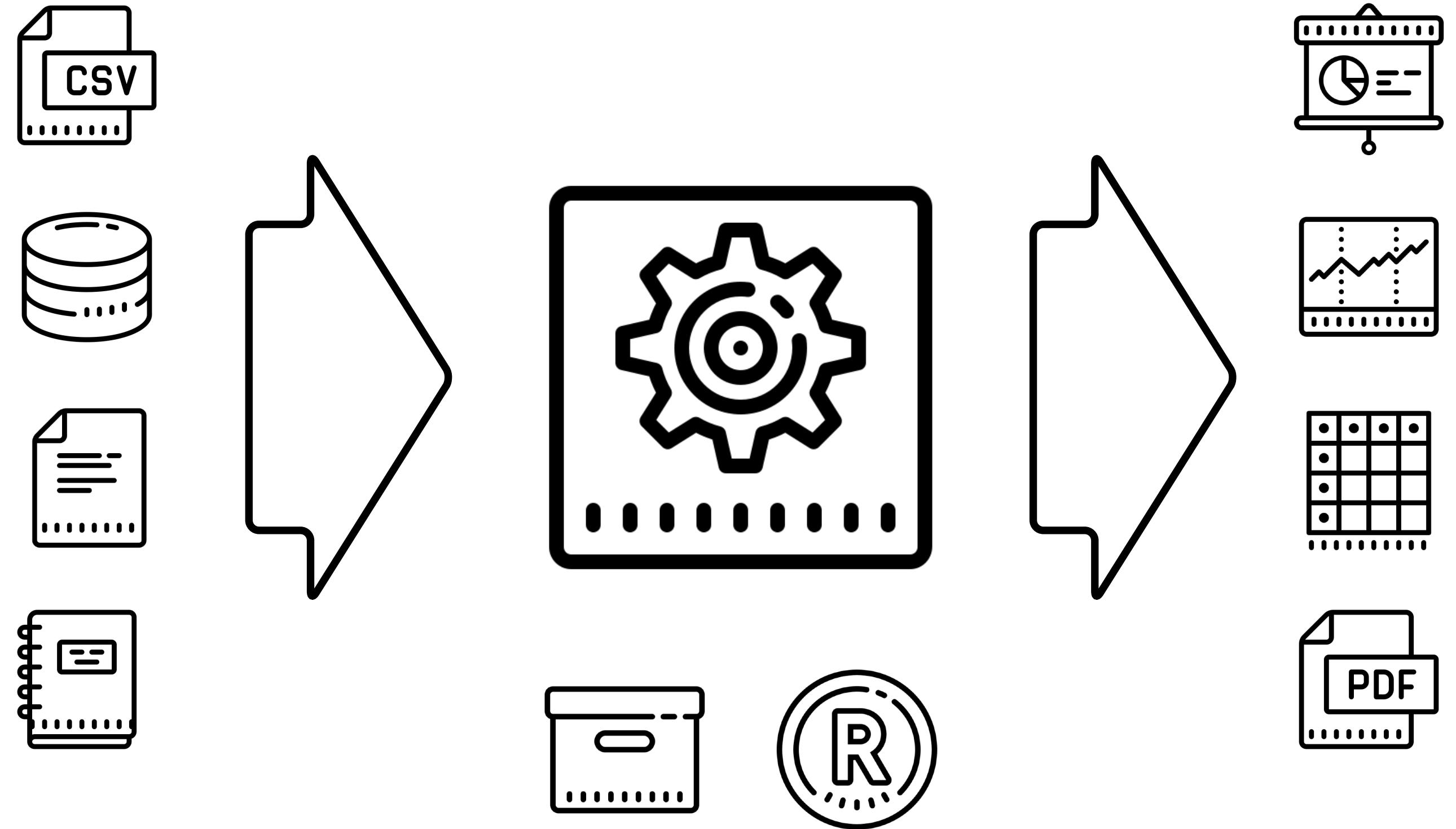
```
data:  
  summary: SELECT * FROM ...
```

```
resources:  
  - support.R  
  - metadata.csv
```

```
packages:  
  - ggplot2  
  - knitr
```

```
script: script.R
```

```
artefacts:  
  - report:  
      description: Summary of results  
      filenames: summary.pdf  
  - data:  
      description: Processed data for further use  
      filenames: data.csv
```



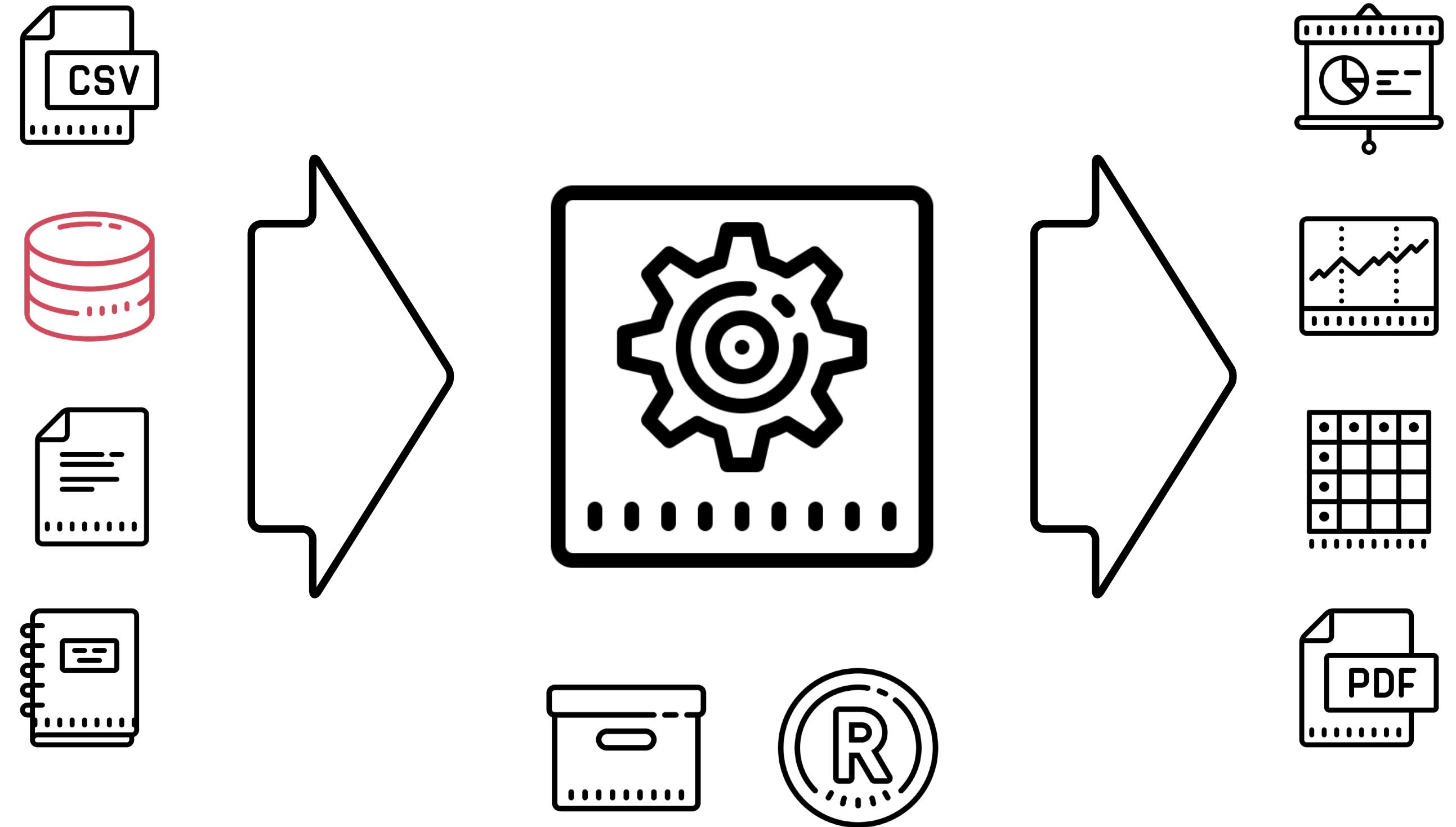
```
data:  
  summary: SELECT * FROM ...
```

```
resources:  
  - support.R  
  - metadata.csv
```

```
packages:  
  - ggplot2  
  - knitr
```

```
script: script.R
```

```
artefacts:  
  - report:  
      description: Summary of results  
      filenames: summary.pdf  
  - data:  
      description: Processed data for further use  
      filenames: data.csv
```



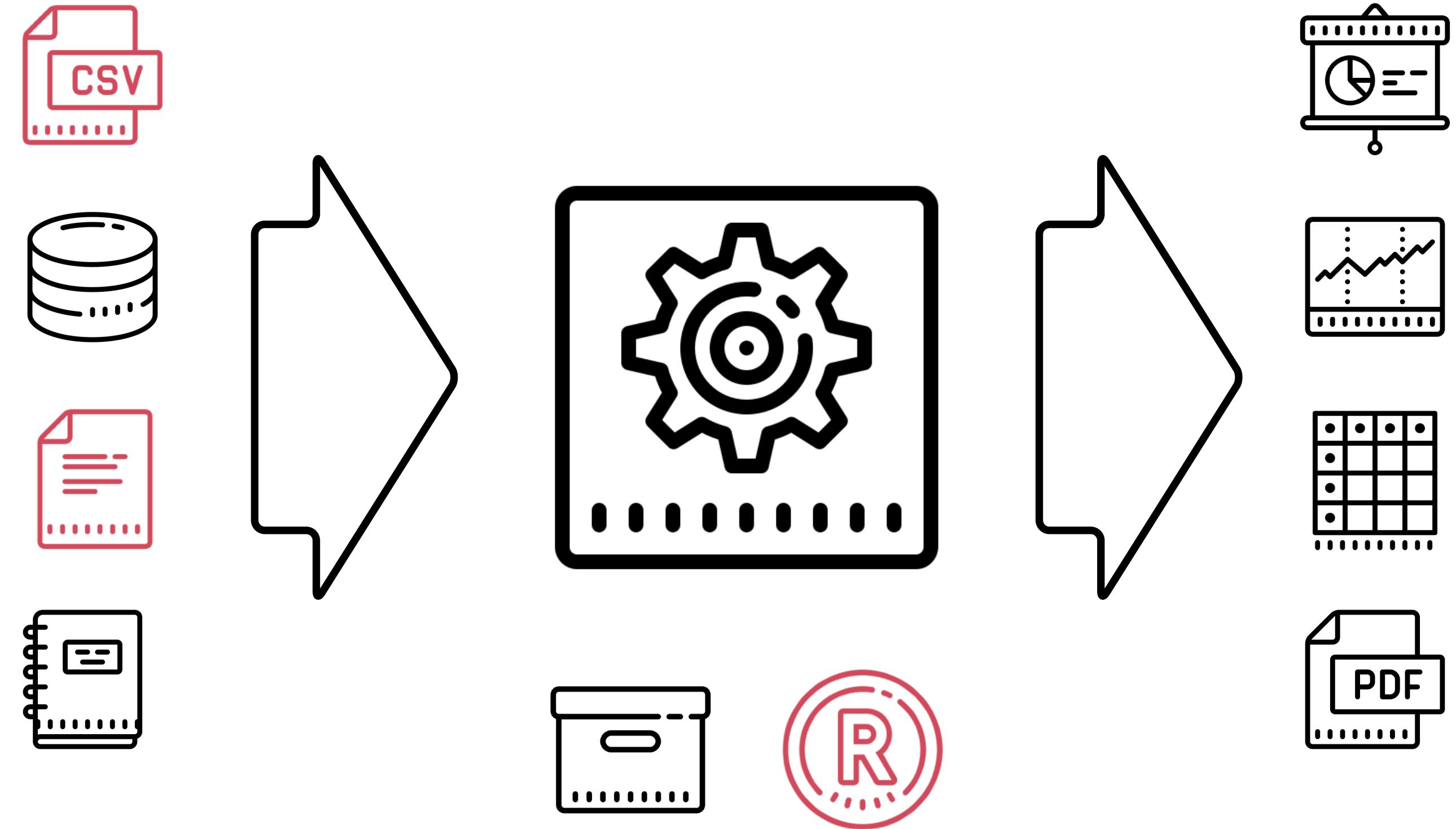
```
data:  
  summary: SELECT * FROM ...
```

```
resources:  
  - support.R  
  - metadata.csv
```

```
packages:  
  - ggplot2  
  - knitr
```

```
script: script.R
```

```
artefacts:  
  - report:  
      description: Summary of results  
      filenames: summary.pdf  
  - data:  
      description: Processed data for further use  
      filenames: data.csv
```



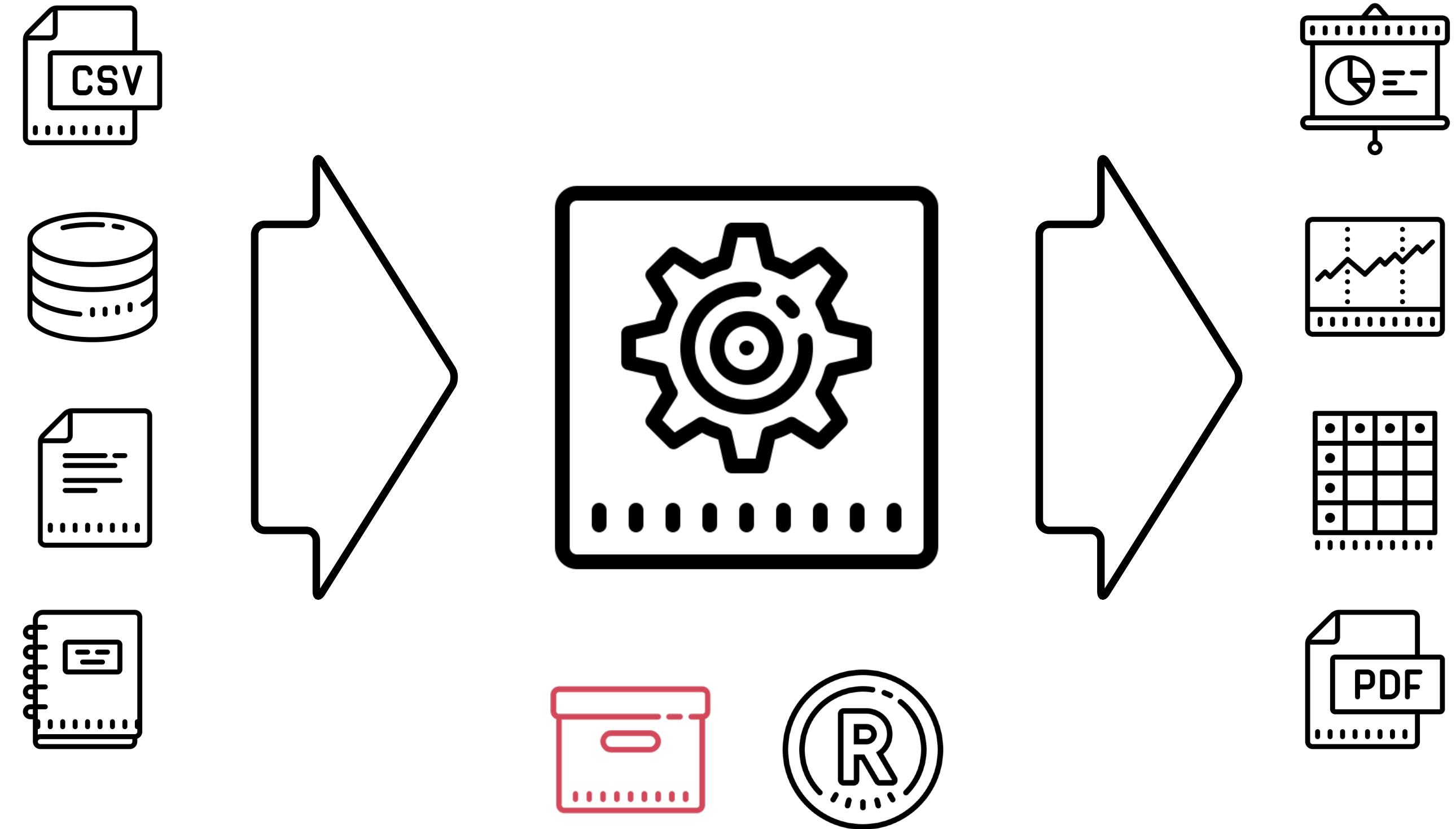
```
data:  
  summary: SELECT * FROM ...
```

```
resources:  
  - support.R  
  - metadata.csv
```

```
packages:  
  - ggplot2  
  - knitr
```

```
script: script.R
```

```
artefacts:  
  - report:  
      description: Summary of results  
      filenames: summary.pdf  
  - data:  
      description: Processed data for further use  
      filenames: data.csv
```



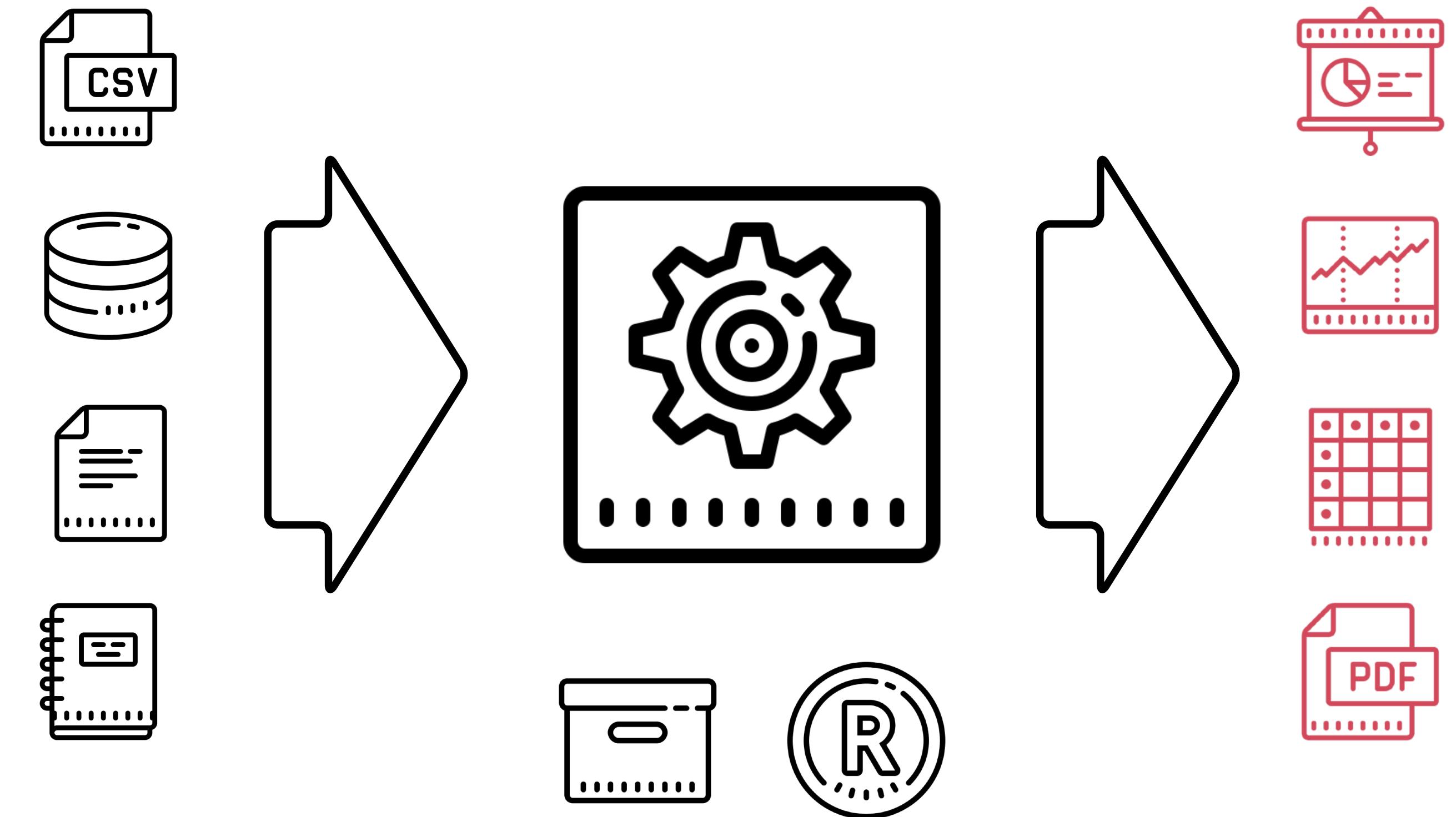
```
data:  
  summary: SELECT * FROM ...
```

```
resources:  
  - support.R  
  - metadata.csv
```

```
packages:  
  - ggplot2  
  - knitr
```

```
script: script.R
```

```
artefacts:  
  - report:  
      description: Summary of results  
      filenames: summary.pdf  
  - data:  
      description: Processed data for further use  
      filenames: data.csv
```



orderly_config.yml

src/

myreport/

 orderly.yml

 script.R

 support.R

 metadata.csv

```
orderly_config.yml  
src/  
myreport/  
    orderly.yml  
    script.R  
    support.R  
    metadata.csv  
archive/  
myreport/  
20190204-143204-f5aa3bc9/  
    orderly.yml  
    script.R  
    support.R  
    metadata.csv  
    summary.pdf  
    data.csv  
    orderly_run.rds
```

```
orderly run myreport
```

orderly_config.yml

src/

myreport/

 orderly.yml

 script.R

 support.R

 metadata.csv

archive/

myreport/

 20190204-143204-f5aa3bc9/

 20190204-192249-3bc9f5aa/

 orderly.yml

 script.R

 support.R

 metadata.csv

 summary.pdf

 data.csv

 orderly_run.rds

orderly run myreport

Automate the boring bits

This is embarrassingly
simple

Interface for Stakeholders

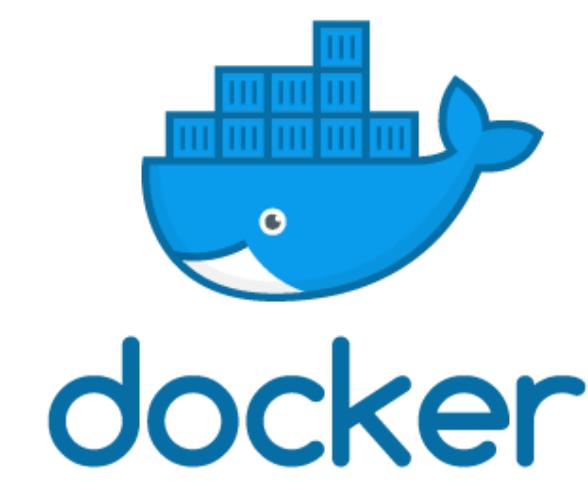
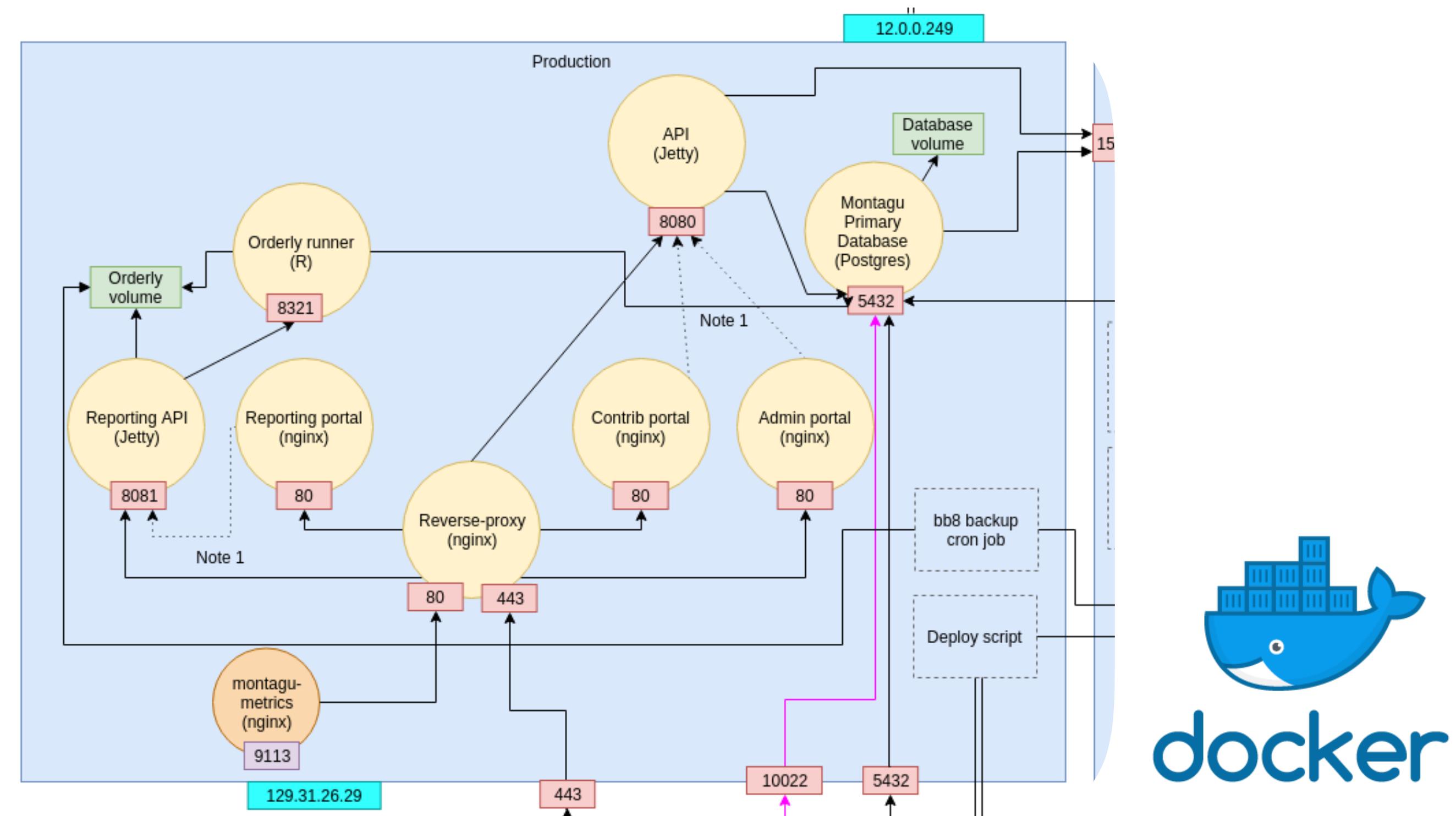
Find a report

Click on a column heading to sort by that field. Hold shift to multi-sort.

[Collapse all reports](#) / [Expand all reports](#)

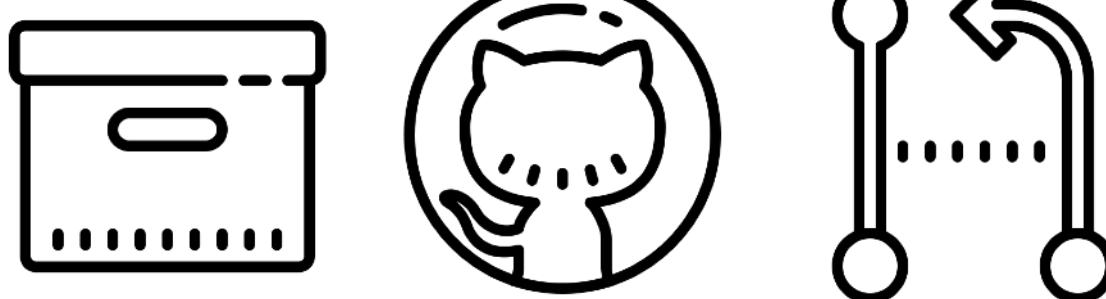
Name	Version	Status	Author	Requester
Type to filter...	<input type="text"/> Mar 1, to <input type="text"/> Feb 6, All ▾	Type to filter...	Type to filter...	
▼ native-diagnostics-burden-report-drafts 5 versions: view latest			Science Team	VIMC
	Tue Feb 05 2019 latest published (20190205-151702-1ba5e47a)		Science Team	VIMC
	Thu Jan 31 2019 out-dated published (20190131-162935-86a83d30)		Science Team	VIMC
	Thu Jan 31 2019 out-dated published (20190131-123847-53fe189e)		Science Team	VIMC
	Mon Jan 28 2019 out-dated published (20190128-151914-a2c1edce)		Science Team	VIMC
	Mon Jan 28 2019 out-dated internal (20190128-125432-31d88275)		Science Team	VIMC

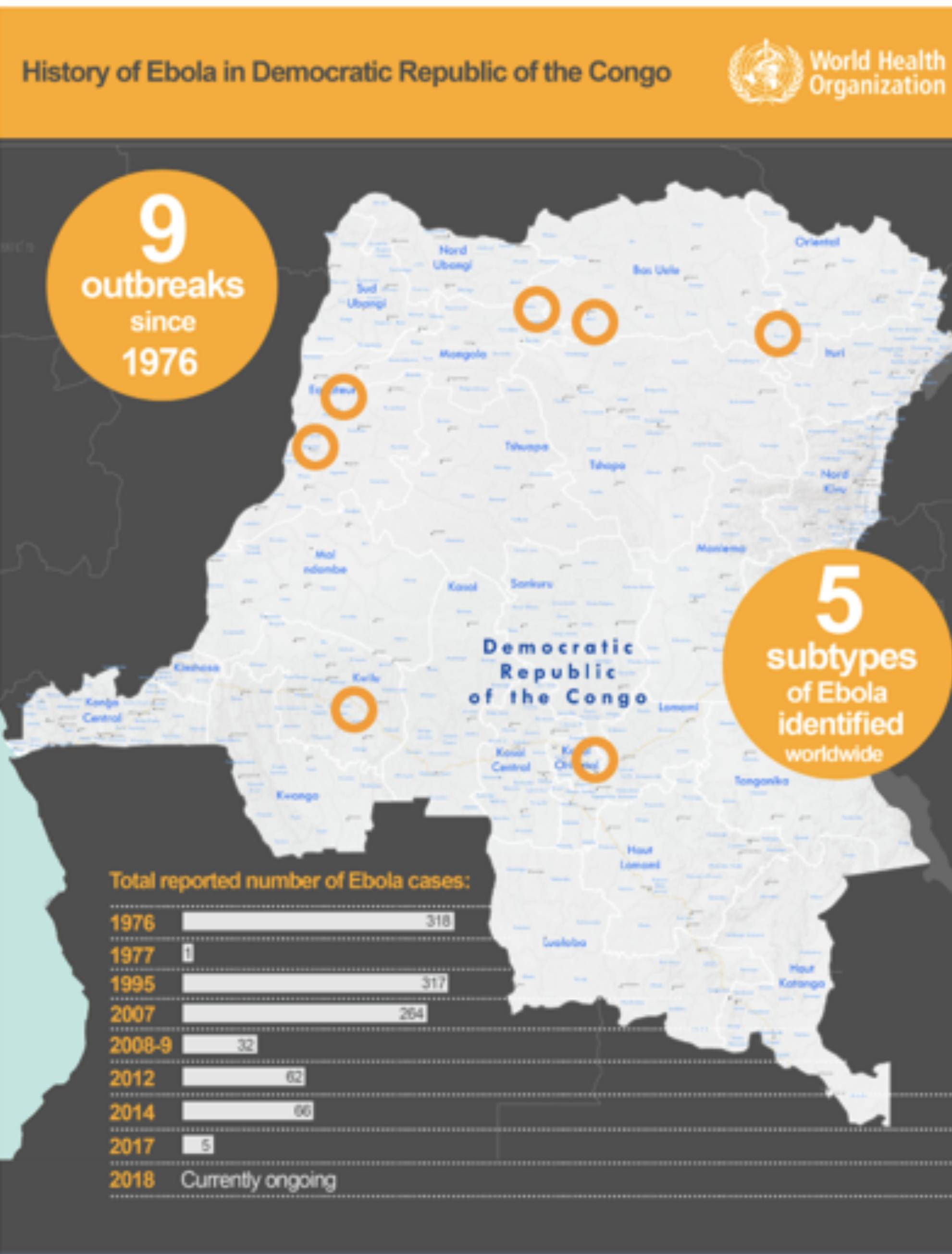
Interface for Engineering team



Interface for Science team

```
install.packages("orderly")
orderly::orderly_new("myreport")
```





Interface for Other groups

Lessons learnt

you can blackbox too much

code reuse is really hard

reproducibility can be easy

problems are as social as technical

Work the way people
want to work





Our team



Vaccine Impact Modelling Consortium
MRC Centre for Global Infectious Disease Analysis
Department of Infectious Disease Epidemiology, Imperial College London