# Workshop 1:
# The tidyverse and beyond



**Brendan Palmer,**
**Statistics & Data Analysis Unit,**
**Clinical Research Facility - Cork**

# The tidyverse package

A bundle of ~20 individual R packages

The six main ones are loaded when the tidyverse package is called

```
> library(tidyverse)
Loading tidyverse: ggplot2
```
→ Data visualisation

```
Loading tidyverse: tidyr
```
→ Data tidying
```
Loading tidyverse: readr
```
→ Data import

```
Loading tidyverse: dplyr
```
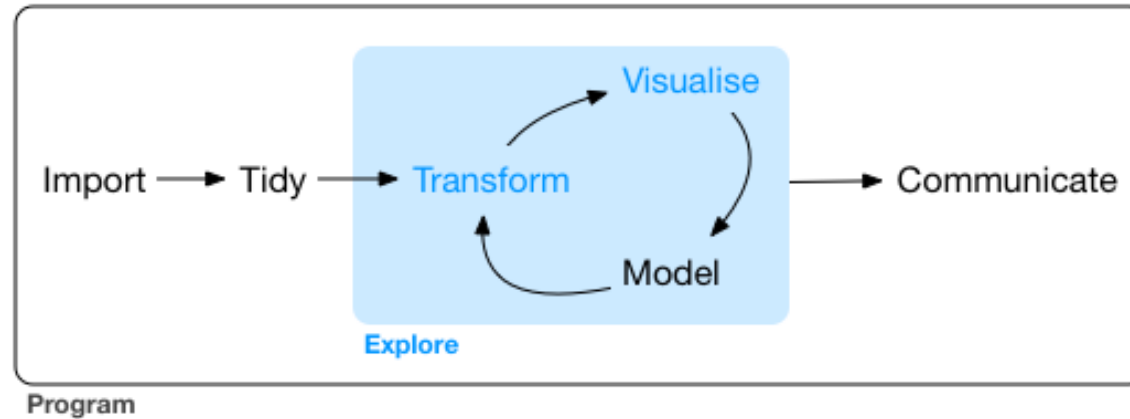→ Data manipulation
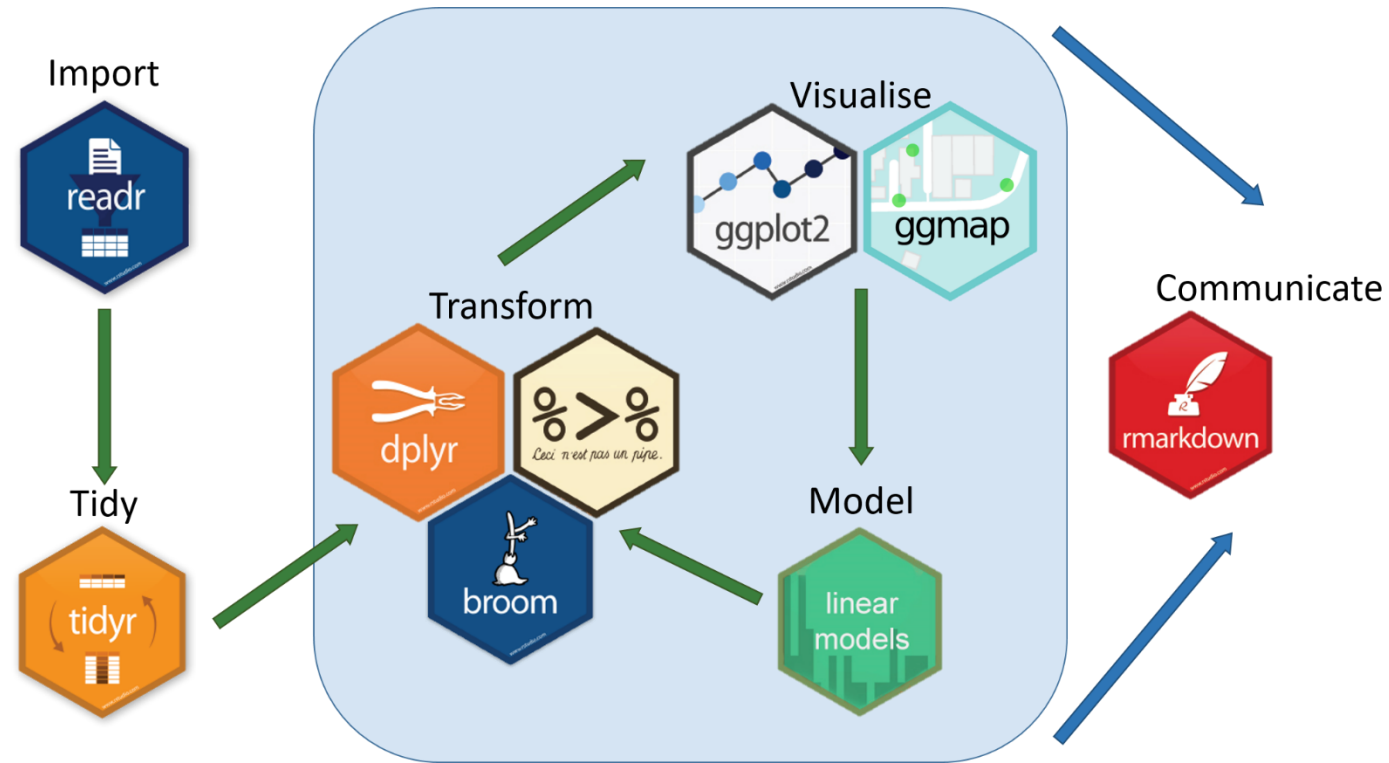
# You could write a book on that!!

# Data analysis in a nutshell



The only obstacle to this is getting the information inside your brain translated into a machine readable format
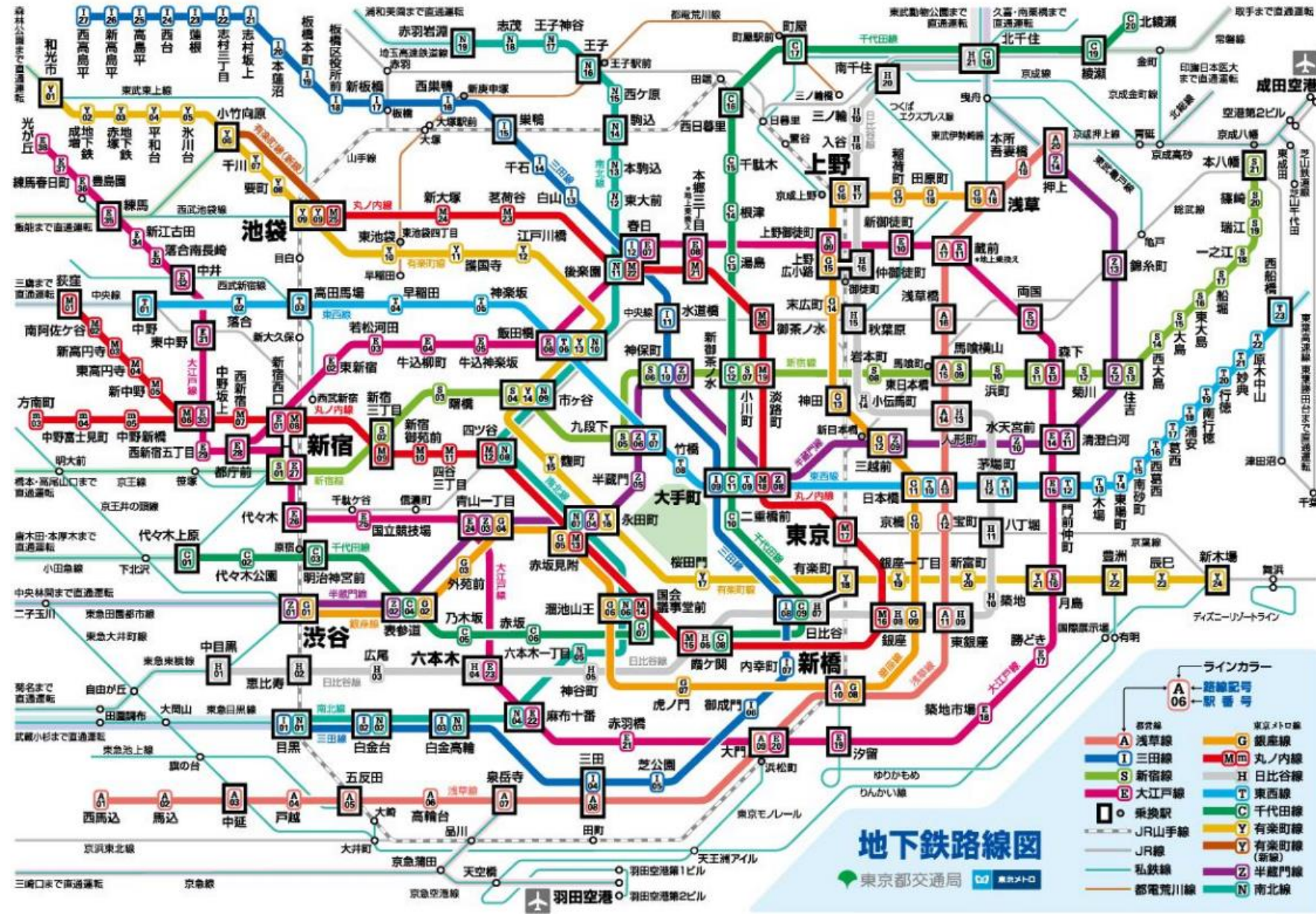
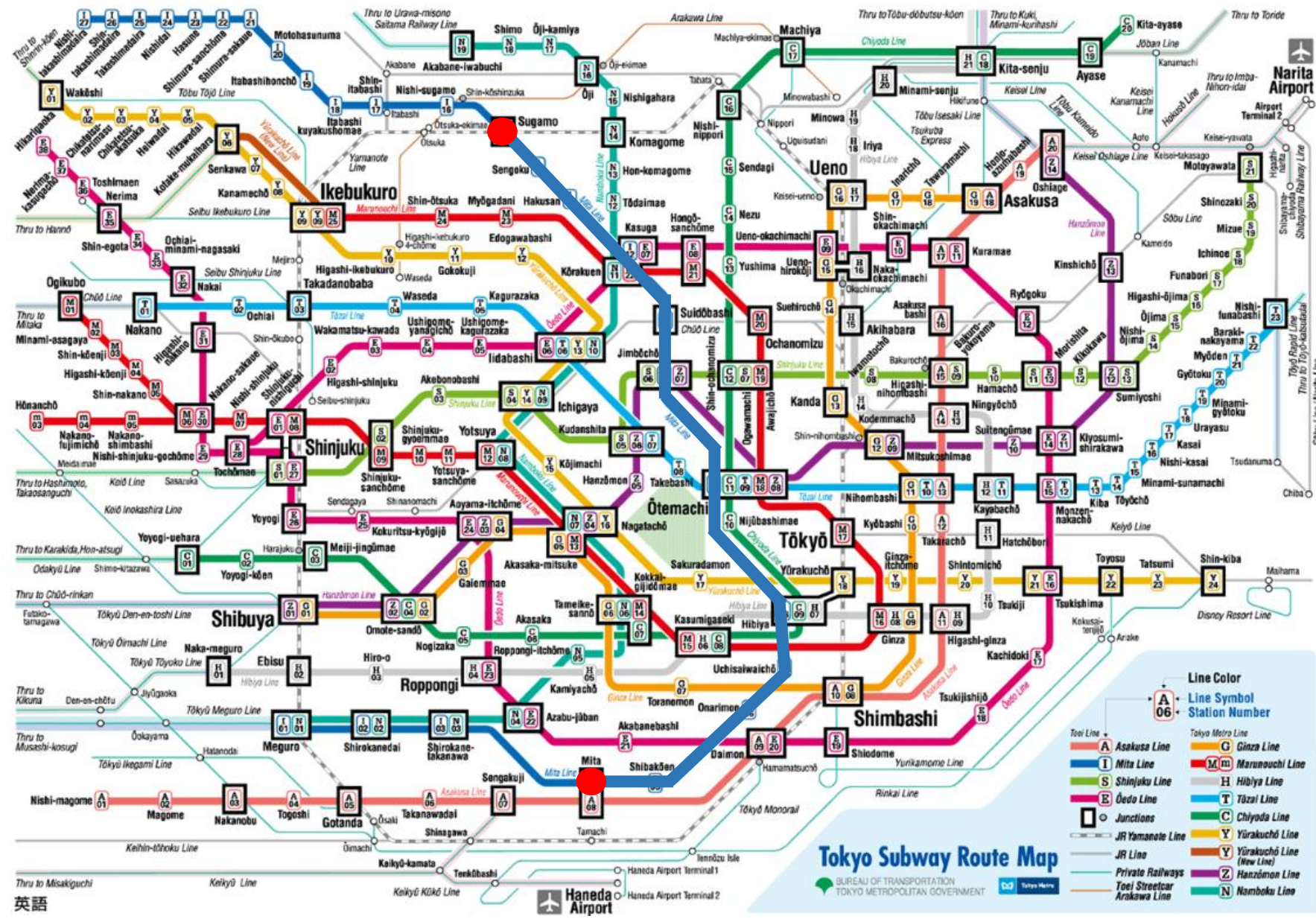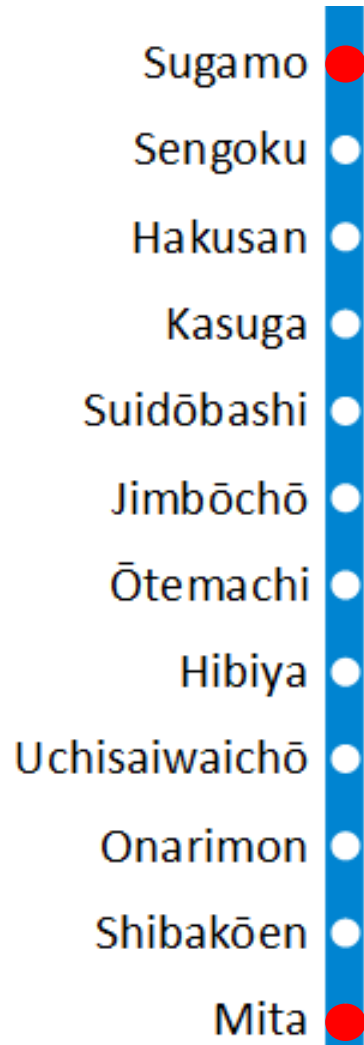# Data analysis in a tidyverse nutshell

# Communication with computers: C++

The challenge:
To find the best route from Mita to Sugamo station on the Toyko metro

# Communication with computers: base R

# Communication with computers: tidyversec

Sugamo
Sengoku
Hakusan
Kasuga
Suidōbashi
Jimbōchō
Ōtemachi
Hibiya
Uchisaiwaichō
Onarimon
Shibakōen
Mita

In general, there are many ways to achieve the same result in programming

With the tidyverse, there is one preferred way to achieve an action!

The tidyverse mantra is that each function does one thing really well

# Basics of R code

| Symbol | What it does | Example 1 | Example 2 |
|---|---|---|---|
| <- | Creates objects | > x <- 5<br>> x<br>[1] 5 | > y <- "This"<br>> y<br>[1] "This" |
| c() | Helps create objects with more than one element | > v <- c(5,6,7,8)<br>> v<br>[1] 5 6 7 8 | > w <- c("This", "is", "easy! ")<br>> w<br>[1] "This" "is" "easy!" |
| # | Computer ignores what is written. Used for adding notes to code | > #print("hello")<br>> | > print("hello")<br>[1] "hello" |
| %>% | Literally translates as "then do this" | > data %<%<br>  do.something.to(data) | |
| %in% | returns a logical vector indicating if there is a match | > "x" %in% c("x", "y", "z")<br>[1] TRUE | > c("x", "y", "z") %in% "x"<br>[1]  TRUE FALSE FALSE |
| ? | Access information | > ?mean() | > ?geom_point() |

**FYI: R is case sensitive!!  Name.of.data ≠ name.of.data**

# Recall:

# R functions

# R packages



Base R:
Comes
pre-
loaded

Other packages:
Install once
Update regularly
Load each session

core
tidyverse

# Worksheet 1
# Part A

# Tidy data should satisfy the following:

Each variable forms a column

Each observation forms a row

**Bauer et al., 2008:**
Column headers are values not variable names

Multiple variable are stored in one column

e.g. column "NAME" contains values such as;
SFB2 || ER to Golgi transport || molecular function unknown || YNL049C || 1082129

These need to be split up

G0.05:U0.03 letter is limiting nutrient and the number is the growth rate

# Try to limit "uninformative" data

"GWEIGHT" contains the same information in every cell
- This isn't going to add to our analysis

"GID" and "YORF" appear to be study specific IDs

"NAME" column contains a lot of information

Going back to the previous example;
SFB2 || ER to Golgi transport || molecular function unknown || YNL049C || 1082129

SFB2: Gene names, but not present in all cases
ER to Golgi transport: Biological process
molecular function unknown: Molecular function
YNL049C: Gene ID listed on public repositories
1082129: Another identifier that does not appear to be useful

# Worksheet 1
# Part B

# Code structure

```
separated_gene_df <- separate(raw_gene_df, NAME, c("name", "BP", "MF", "systematic_name", "number"), sep = "\\|\\|")
```

separated_gene_df          -the new data frame you will create

<-                         -the assign operator

separate                   -the function you are calling on

(raw_gene_df,              -the data frame to be used

NAME,                      -the column to be altered

c("name", "BP", "MF", "systematic_name", "number"),
                   -new columns IDs for the new columns

sep = "\\|\\|")            -identify the separator to be used

# Worksheet 1
# Part C

# How to plot in ggplot

Template:

```
ggplot(data = <DATA>) +

    <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>))    +

      linear model  +

      axes formatting  +

      legend formatting  +

      title    + etc. etc.
```

# Worksheet 1
# Part D