# oRganization

How to make internal R
packages a part of your team

Emily Riederer @emilyriederer  emily.rbind.io

data access
server connection
proxies, ssh, ssl

right problems
tribal knowledge
intuition

team norms
meetings
communication

my first day

data access
server connection
proxies, ssh, ssl

right problems
tribal knowledge
intuition

team norms
meetings
communication

my
first
day

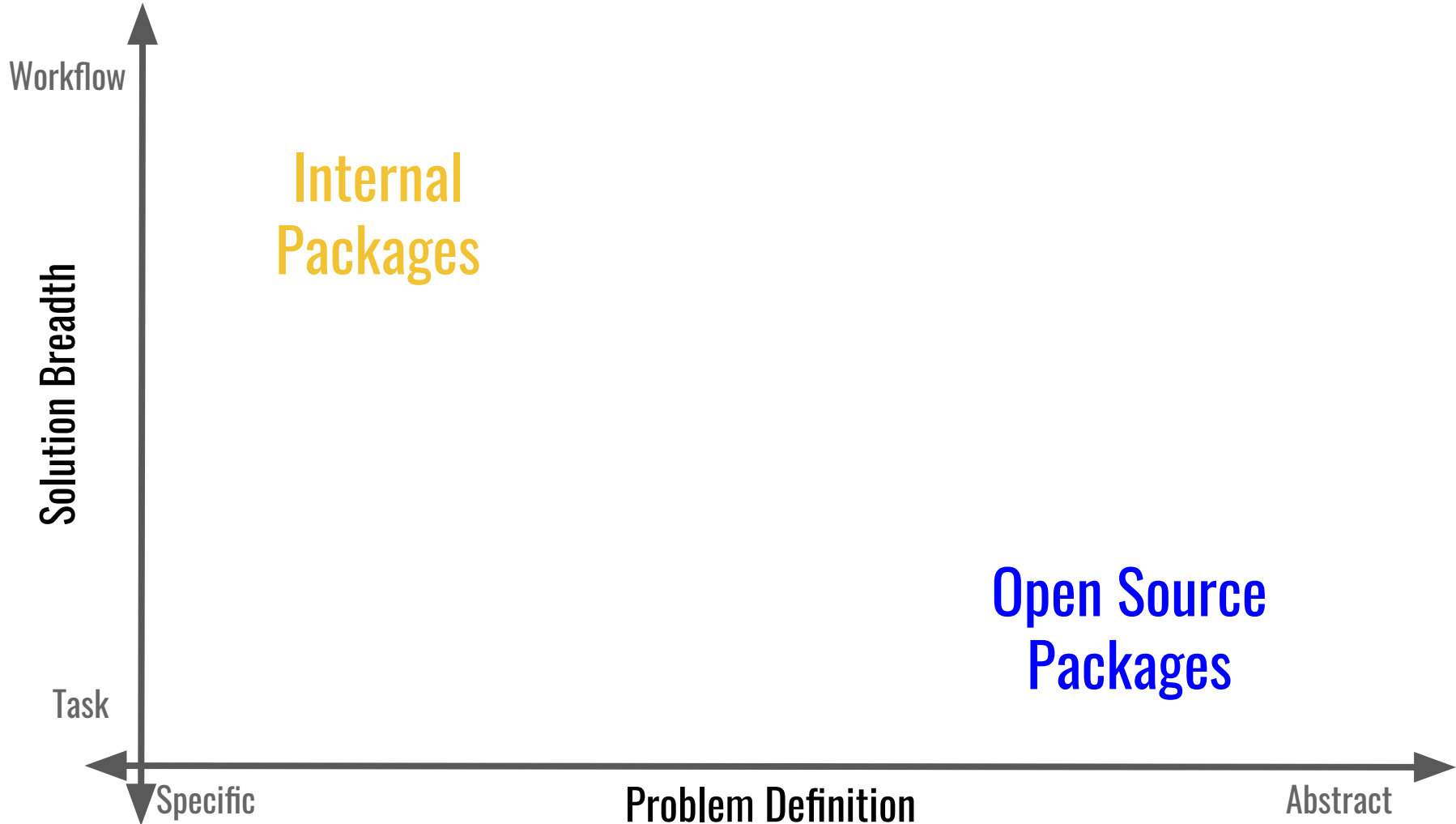## utilities packages

data access
server connection
proxies, ssh, ssl

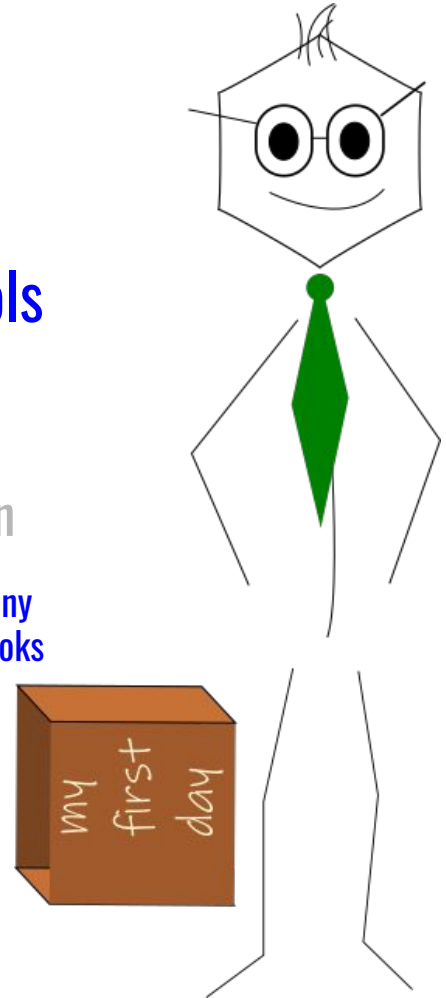e.g. abstraction layer for
infrastructure

## analysis packages

right problems
tribal knowledge
intuition

e.g. curated workflow, tailored
function calls, automated
result generation

## developer tools

team norms
meetings
communication

e.g. color palettes, Shiny
modules, linters, git hooks

# Jobs-to-be-Done

We

**hire a product**

to do a

**job**

that helps us make

**progress**

towards a goal

# Jobs-to-be-Done

We

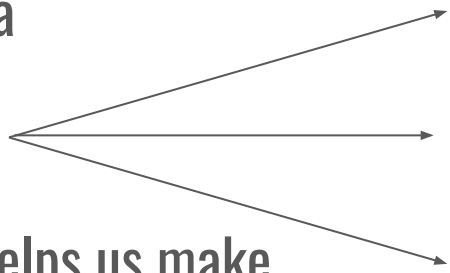**hire a product**

to do a

**job**

that helps us make

**progress**

towards a goal

# Jobs-to-be-Done

We

**hire a product**

to do a          <span style="color:blue">functional</span>

**job**          <span style="color:green">social</span>

that helps us make          <span style="color:gold">emotional</span>

**progress**

towards a goal

# Jobs-to-be-Done

We

**hire a product**

to do a       → functional

**job**       → social

that helps us make      → emotional

**progress**

towards a goal

Let's

**build a team of packages**

to do the

**jobs**

that helps our org

**answer impactful questions**

with efficient workflows

# The IT Guy



**functional** handle quirks of infrastructure

**social** promote or enforce good practices

**emotional** avoid frustration or stress of time lost

# The IT Guy



**functional**   handle quirks of infrastructure

**social**   promote or enforce good practices

**emotional**   avoid frustration or stress of time lost

-> utility functions

-> opinionated design

-> helpful error messages

```r
get_database_conn <- function(username, password) {

conn <-
  DBI::dbConnect(
    drv = odbc::odbc(),
    driver = {driver name},
    server = {server},
    UID = username,
    PWD = password,
    port = {port number}
  )

return(conn)

}
```

```
get_database_conn <- function(~~username, password~~) {

conn <-
  DBI::dbConnect(
    drv = odbc::odbc(),
    driver = {driver name},
    server = {server},
    UID = **Sys.getenv("DB_USER")** ~~username~~,
    PWD = **Sys.getenv("DB_PASS")** ~~password~~,
    port = {port number}
  )

return(conn)

}
```

```r
get_database_conn <- function() {

if (any(Sys.getenv(c("DB_USER", "DB_PASS")) == "")) {
  stop(
    "DB_USER or DB_PASS environment variables are missing.",
    "Please read set-up vignette to configure your system."
  )
}

conn <-
  DBI::dbConnect(
    drv = odbc::odbc(),
    driver = {driver name},
    server = {server},
    UID = Sys.getenv("DB_USER"),
    PWD = Sys.getenv("DB_PASS"),
    port = {port number}
  )

return(conn)

}
```

```r
get_database_conn <- function() {

if (any(Sys.getenv(c(“DB_USER”, “DB_PASS”)) == “”)) {
  stop(
    “DB_USER or DB_PASS environment variables are missing.”,
    “Please read set-up vignette to configure your system.”
  )
}

conn <-
  DBI::dbConnect(
    drv = odbc::odbc(),
    driver = {driver name},
    server = {server},
    UID = Sys.getenv(“DB_USER”),
    PWD = URLencode(Sys.getenv(“DB_PASS”), reserved = TRUE),
    port = {port number}
  )

return(conn)

}
```

# The Junior Analyst

**functional**    perform work with reasonable assumptions

**social**    flexible to feedback, trying new things

**emotional**    builds trust so you can focus on other things

# The Junior Analyst

**functional**    perform work with reasonable assumptions

**social**    flexible to feedback, trying new things

**emotional**    builds trust so you can focus on other things

-> default arguments
-> reserved keywords
-> ellipsis

```r
viz_cohort <- function(data, time, metric, group) {

  gg <-
    ggplot(data) +
    aes(x = .data[[time]],
        y = .data[[metric]],
        group = .data[[group]]) +
    geom_line() +
    my_org_theme()

  return(gg)

}
```

```
viz_cohort <- function(data, ~~time~~, metric, group) {

  gg <-
    ggplot(data) +
    aes(x = .data[["MONTHS_SUBSCRIBED"]],
        y = .data[[metric]],
        group = .data[[group]]) +
    geom_line() +
    my_org_theme()

  return(gg)

}
```

```
viz_cohort <- function(data,
                       metric = "IND_ACTIVE",
                       time = "MONTHS_SUBSCRIBED",
                       group = "COHORT") {

  gg <-
    ggplot(data) +
    aes(x = .data[[time]],
        y = .data[[metric]],
        group = .data[[group]]) +
    geom_line() +
    my_org_theme()

  return(gg)

}
```

```r
viz_cohort <- function(data,
                       metric = "IND_ACTIVE",
                       time = "MONTHS_SUBSCRIBED",
                       group = "COHORT") {

  gg <-
    ggplot(data) +
    aes(x = .data[[time]],
        y = .data[[metric]],
        group = .data[[group]]) +
    geom_line() +
    my_org_theme()

  return(gg)

}
```

**Reserved Keywords:**

TIME_SUBSCRIBED
CUSTOMER_COHORT
CUSTOMER_SEGMENT

...

```r
viz_cohort <- function(data,
                       time = "MONTHS_SUBSCRIBED",
                       metric = "IND_ACTIVE",
                       group = "COHORT",
                       ...) {

  gg <-
    ggplot(data) +
    aes(x = .data[[time]],
        y = .data[[metric]],
        group = .data[[group]]) +
    geom_line(aes(...)) +
    my_org_theme()

  return(gg)

}
```

```
viz_cohort <- function(data,
                       time = "MONTHS_SUBSCRIBED",
                       metric = "IND_ACTIVE",
                       group = "COHORT",
                       ...) {

  gg <-
    ggplot(data) +
    aes(x = .data[[time]],
        y = .data[[metric]],
        group = .data[[group]]) +
    geom_line(aes(...)) +
    my_org_theme()

  return(gg)

}
```



```
> viz_cohort(my_data)
```



```
> viz_cohort(my_data,
             color = COHORT,
             linetype = COHORT)
```

# The Tech Lead

**functional**    coach you through issues & alternatives

**social**    share collected knowledge

**emotional**    inspire you to do your best work

# The Tech Lead



**functional**  help navigate common issues & alternatives

**social**  share collected knowledge

**emotional**  connect to latent community of practice

-> vignettes

-> templates

# Vignettes as a time capsule for knowledge transfer



**Crash course**
(`dplyr`)

**Method Overview**
(`survival`)

# Vignettes as a time capsule for knowledge transfer



**Conceptual Overview**

**Workflow & Key Questions**

**Process Documentation**

# Vignettes as a time capsule for knowledge transfer

**Conceptual Overview**

**Workflow & Key Questions**

**Process Documentation**

**Technical Overview**

**Methods Comparison**

# Vignettes as a time capsule for knowledge transfer



**Conceptual Overview**

**Workflow & Key Questions**

**Process Documentation**

**Technical Overview**

**Methods Comparison**

**Lessons Learned**

**Past Examples**

# Expand your reach with `pkgdown`

```
> pkgdown::build_site()
```

Daily Calendar

7
8
9
10
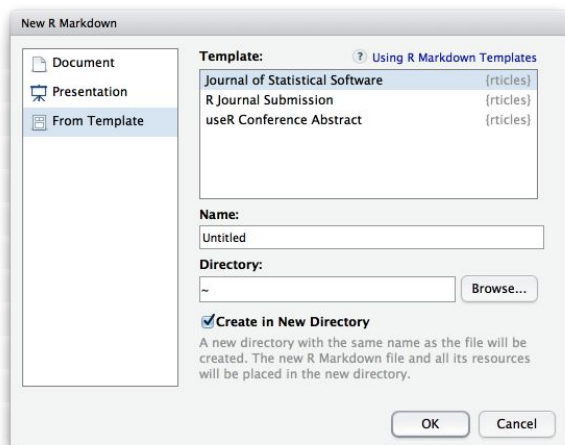11
12
1
2
3
4
5
6
7

# Templates as coach

## Structure
(`flexdashboard`)



```
---
title: "Untitled"
output:
  flexdashboard::flex_dashboard:
    orientation: columns
    vertical_layout: fill
---

```{r setup, include=FALSE}
library(flexdashboard)
```

Column {data-width=650}
-------------------------------

### Chart A

```{r}

```

Column {data-width=350}
-------------------------------

### Chart B

```{r}
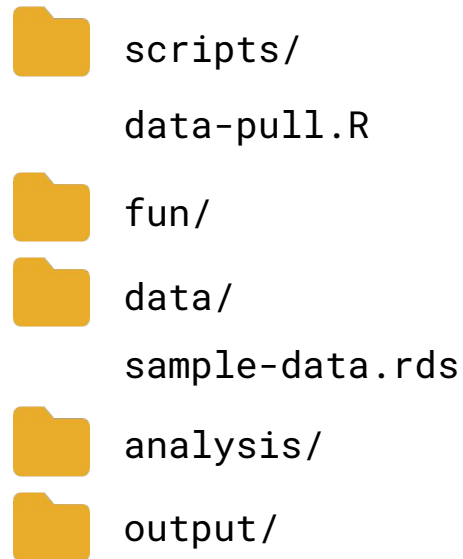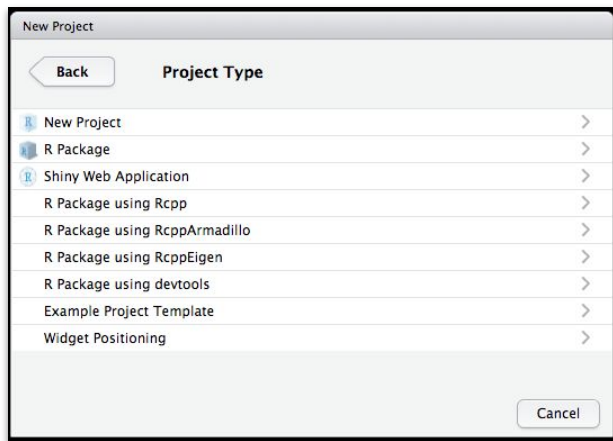
```
```

# Templates as coach

## Process walk-through

```
---
title: "Data Validation"
output: html_document
---

## Censored Data

Run the following code to visualize how many
observations were censored. Depending on what
you find you will want to...

```{r censored}
```

## Analysis outline

```
---
title: "Final Report"
output: html_document
params:
  month: September
---

## Final Report

TODO: UPDATE COMMENTARY SUMMARIZING TRENDS

```{r dashboard}
```

New R Markdown

Template:                    ? Using R Markdown Templates

Document        Journal of Statistical Software    {rticles}
Presentation    R Journal Submission               {rticles}
From Template   useR Conference Abstract           {rticles}

Name:
Untitled

Directory:
~                                          Browse...

☑ Create in New Directory
A new directory with the same name as the file will be
created. The new R Markdown file and all its resources
will be placed in the new directory.

                              OK        Cancel
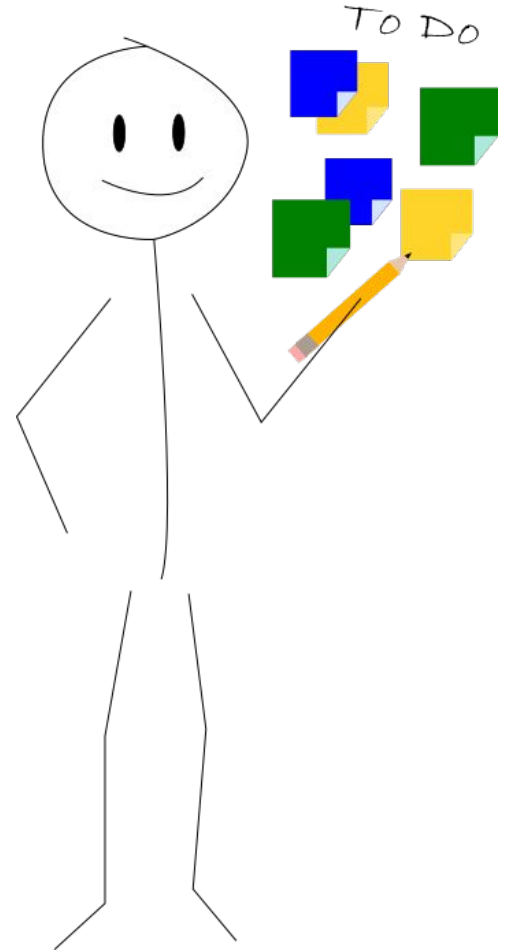
# Templates as code reviewer

# The Project Manager

**functional**   integrates work

**social**   finds common ground

**emotional**   meets you where you are
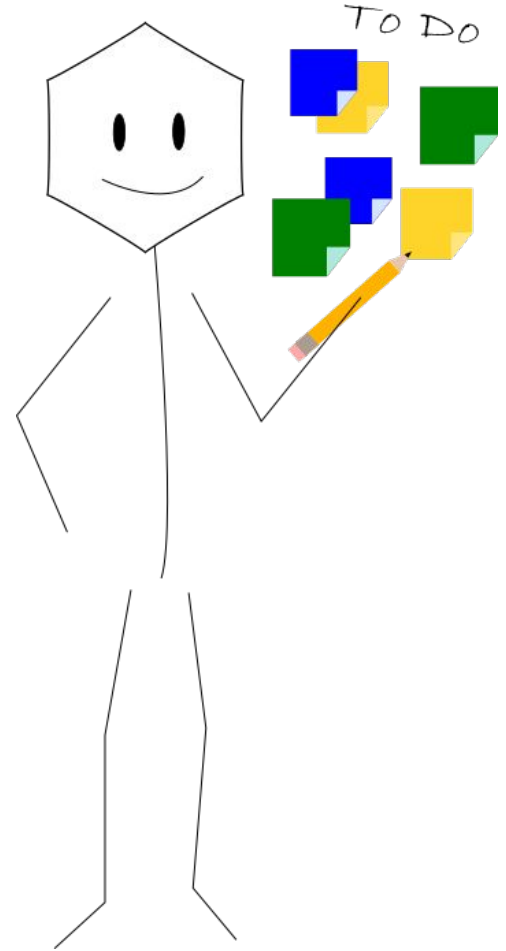
# The Project Manager

**functional**   integrates work

**social**   finds common ground

**emotional**   meets you where you are

-> modularized workflow

-> IDE support

TO DO

# Modularization

```
---
title: "My Document"
output: html_document
---

```{r setup, include=FALSE}
knitr::opts_chunk$set(echo = TRUE)
```

## Section 1

```{r cars}
summary(cars)
```

```{r child = "commentary.md"}
```
```

**Main R Markdown**

**commentary.md**

```
### My observations

This is what we noticed…
```

# IDE Support



**Visual Editor**

**Add-Ins (e.g. `esquisse`)**

# Collaboration



**functional** clear communication

**social** keeps promises

**emotional** confident yet engaged

# Collaboration



**functional** — clear communication

**social** — keeps promises

**emotional** — confident yet engaged

-> naming
-> scope
-> dependencies
-> testing

# Clear communication

# Clear ownership

# Clear ownership



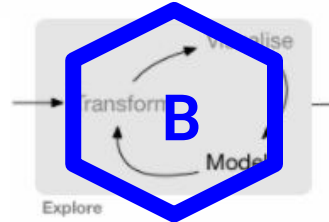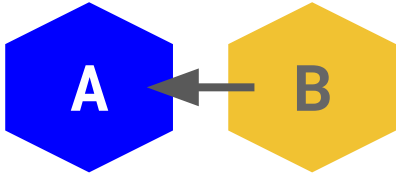Image from *R for Data Science (Wickham & Grolemund)*

# Clear ownership



Image from *R for Data Science (Wickham & Grolemund)*

# Clear ownership
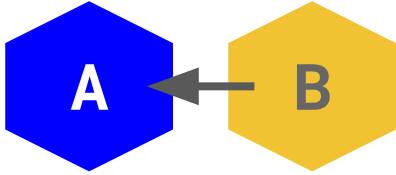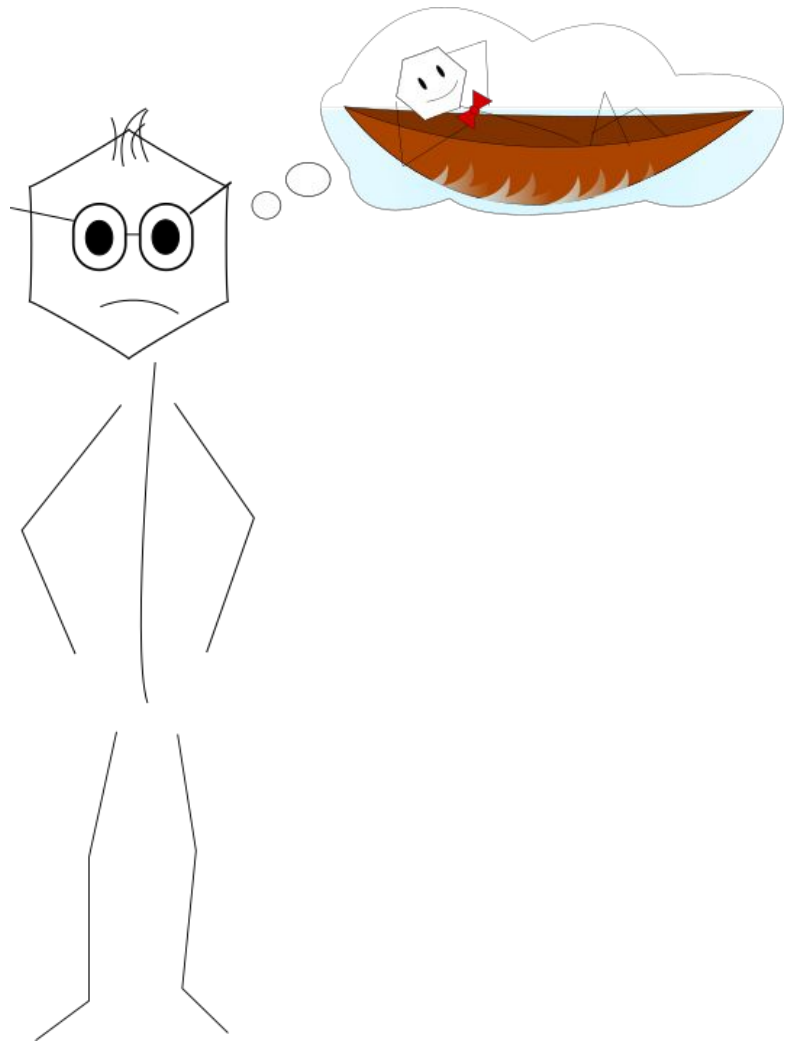
# Dependency structures



```
a_fx <- function() {...}
```

```
b_fx <- function() {
    ...
    a_fx()
    ...
}
```
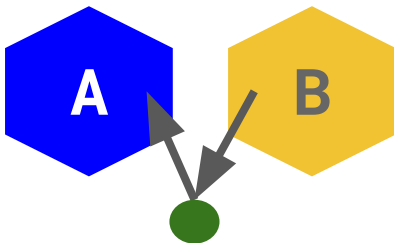
**Direct Dependency**

# Dependency structures
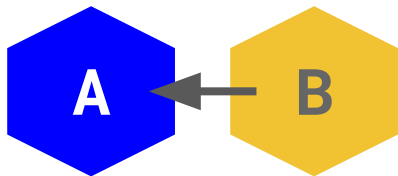


```
a_fx <- function() {...}
```

```
b_fx <- function() {
    ...
    a_fx()
    ...
}
```

**Direct Dependency**

# Dependency structures



```
a_fx <- function() {...}
```
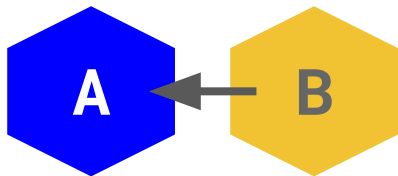
```
b_fx <- function() {
    ...
    a_fx()
    ...
}
```

```
a_fx <- function() {...}
```

```
b_fx <- function(a_input) {
    ...
    do_something(a_input)
    ...
}
```

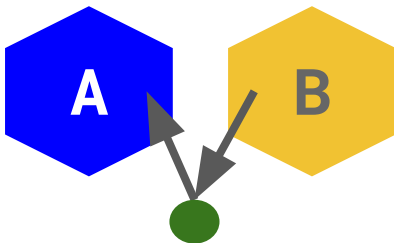**Direct Dependency**                    **Clean Hand Off**

# Dependency structures



```
a_fx <- function() {...}
```
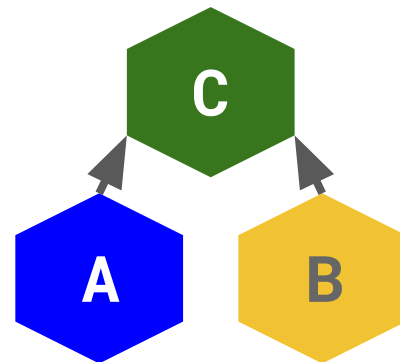
```
b_fx <- function() {
    ...
    a_fx()
    ...
}
```

**Direct Dependency**

```
a_fx <- function() {...}
```

```
b_fx <- function(a_input) {
    ...
    do_something(a_input)
    ...
}
```
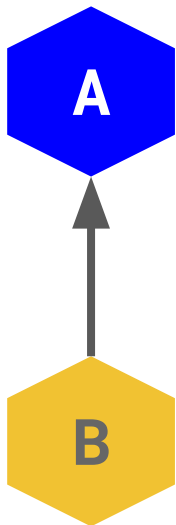
**Clean Hand Off**

```
b_fx <- function() {
    ...
    c_fx()
    ...
}
```

```
b_fx <- function() {
    ...
    c_fx()
    ...
}
```
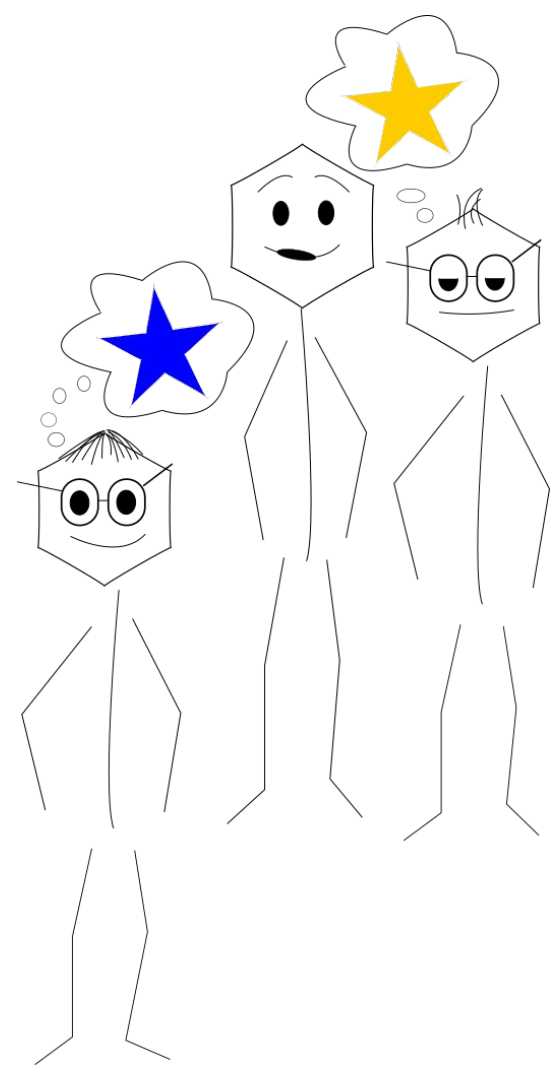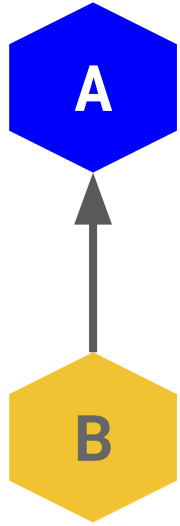
**Common Parent**

# Typical unit test with dependency
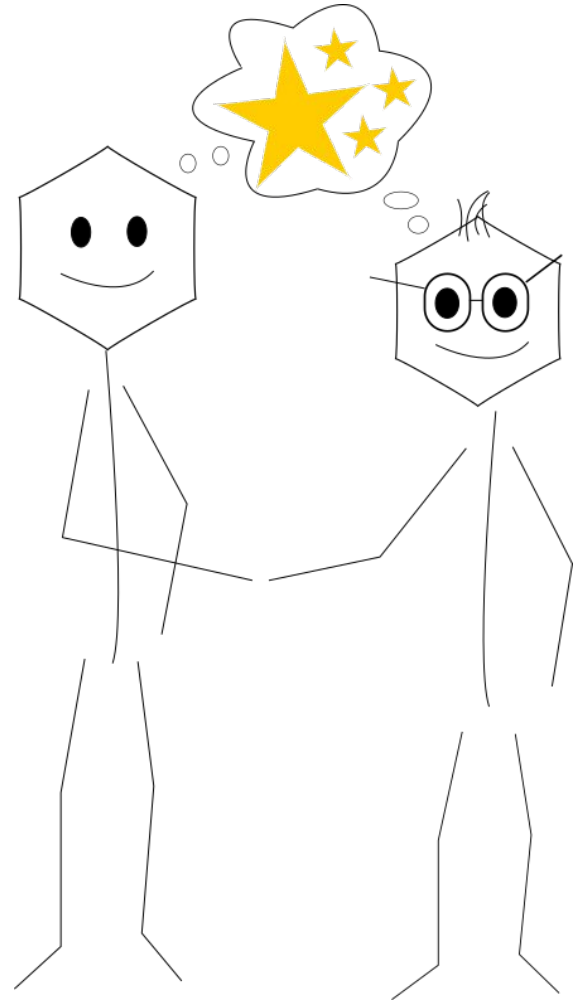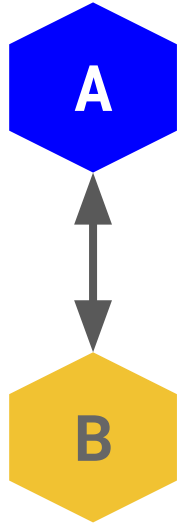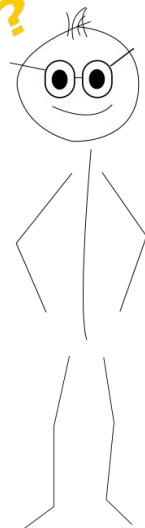


**b**/tests/testthat/test-**pkga**.R

```
test_that(
  "Receives input correctly from a",
  {
    expect_error(fxb(fxa(1)), NA)
  }
)
```

# Typical unit test with dependency

# Typical unit test with dependency

# Integration tests

**a**/tests/testthat/test-**pkgb**.R

```
test_that(
  "Preps input correctly for b",
  {
   expect_error(fxb(fxa(1)), NA)
  }
)
```
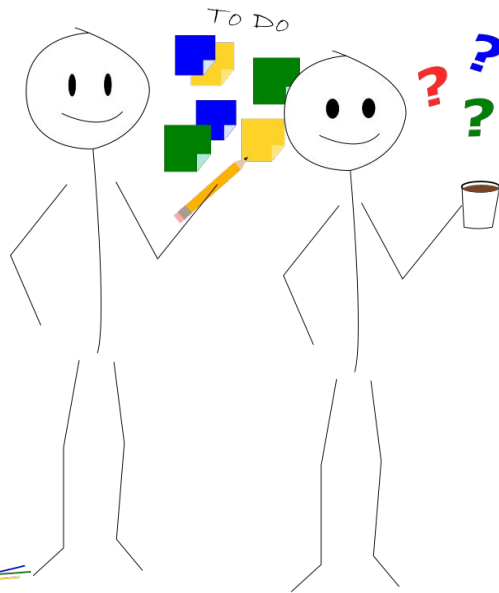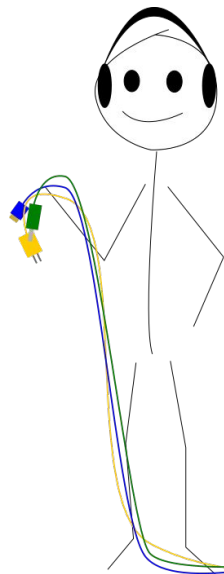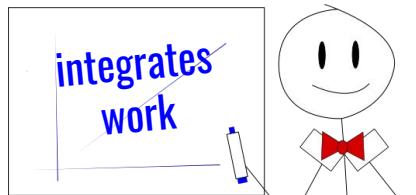
**b**/tests/testthat/test-**pkga**.R

```
test_that(
  "Receives input correctly from a",
  {
   expect_error(fxb(fxa(1)), NA)
  }
)
```