

# RStudio IDE :: CHEAT SHEET

## Documents and Apps

Open Shiny, R Markdown, Knitr, Sweave, Latex, Rd files and more in Source Pane

Check spelling  
Render output  
Choose output format  
Jump to previous chunk  
Jump to next chunk  
Access markdown guide at **Help > Markdown Quick Reference**

RStudio recognizes that files named **app.R**, **server.R**, **ui.R**, and **global.R** belong to a shiny app

## Debug Mode

Open with **debug()**, **browse()**, or a breakpoint. RStudio will open the debugger mode when it encounters a breakpoint while executing code.

Click next to line number to add/remove a breakpoint.  
Highlighted line shows where execution has paused.  
Run commands in environment where execution has paused.

Examine variables in executing environment.  
Select function in traceback to debug.

Step through code one line at a time.

Resume execution mode.  
Quit debugging.

## Write Code

Open in new window  
Save  
Find and replace  
Compile as notebook  
Run selected code

Multiple cursors/column selection with **Alt + mouse drag**.  
Code diagnostics that appear in the margin. Hover over diagnostic symbols for details.

Syntax highlighting based on your file's extension.  
Tab completion to finish function names, file paths, arguments, and more.

Multi-language code snippets to quickly use common blocks of code.  
Change file type.

## R Support

Import data with wizard  
History of past commands to run/copy  
Display .Rpres slideshows  
**File > New File > R Presentation**

Choose environment to display from list of parent environments.  
Search inside environment.  
Display objects as list or grid.

Displays saved objects by type with short description.  
View in data viewer  
View function source code

## Version Control with Git or SVN

Turn on at **Tools > Project Options > Git/SVN**

Show file diff  
Commit staged files to remote  
Push/Pull  
View History  
Open shell to type commands  
current branch

## Package Writing

**File > New Project > New Directory > R Package**

Turn project into package.  
Enable oxygen documentation with **Tools > Project Options > Build Tools**  
Rxygen guide at **Help > Rxygen Quick Reference**

## Pro Features

Share Project with Collaborators... collaborators  
Start new R Session in current project

RStudio saves the call history, workspace, and working directory associated with a project. It reloads each when you re-open a project.

RStudio opens plots in a dedicated Plots pane  
Name of current project  
Export plot  
Delete all plots

Install Packages  
Update Packages library for your project  
Click to load package with **library()**. Undo/Click to detach package with **detach()**  
Delete Package from installed library

RStudio opens documentation in a dedicated Help pane  
Home page of helpful links  
Search within help file  
Search for help file

Viewer Pane displays HTML content, such as Shiny apps, RMarkdown reports, and interactive visualizations  
Stop Shiny app  
Publish to shinyapps.io, Refresh rpubs, RSCONnect, ...

**View(-data-)** opens spreadsheet like view of data set  
Filter rows by value  
Sort by values  
Search for value



# Tidy evaluation with rlang :: CHEAT SHEET



## Vocabulary

**Tidy Evaluation (Tidy Eval)** is not a package, but a framework for doing non-standard evaluation (i.e. delayed evaluation) that makes it easier to program with tidyverse functions.

**pi** **Symbol** - a name that represents a value or object stored in R, `is_symbol(expr(pi))`

**Environment** - a list-like object that binds symbols (names) to objects stored in memory. Each env contains a link to a second, **parent env**, which creates a chain, or search path, of environments. `is_environment(current_env())`



**rlang::caller\_env(n = 1)** Returns calling env of the function it is in.

**rlang::child\_env(parent, ...)** Creates new env as child of parent. Also **env**.

**rlang::current\_env()** Returns execution env of the function it is in.

**1** **Constant** - a bare value (i.e. an atomic vector of length 1). `is_bare_atomic(1)`

**abs (1)** **Call object** - a vector of symbols/constants/calls that begins with a function name, possibly followed by arguments. `is_call(expr(abs(1)))`

**pi** **code** `3.14` **result**

**Code** - a sequence of symbols/constants/calls that will return a result if evaluated. Code can be:

1. Evaluated immediately (Standard Eval)
2. Quoted to use later (Non-Standard Eval) `is_expression(expr(pi))`

**e** **Expression** - an object that stores quoted code without evaluating it. `is_expression(expr(a + b))`

**a + b** **q** **Quosure** - an object that stores both quoted code (without evaluating it) and the code's environment. `is_quosure(quo(a + b))`

**rlang::quo\_get\_env(quo)** Return the environment of a quosure.

**rlang::quo\_set\_env(quo, expr)** Set the environment of a quosure.

**a + b** **rlang::quo\_get\_expr(quo)** Return the expression of a quosure.

**Expression Vector** - a list of pieces of quoted code created by base R's `expression` and `parse` functions. Not to be confused with `expression`.

## Quoting Code

Quote code in one of two ways (if in doubt use a quosure):

### QUOSURES

**Quosure** - An expression that has been saved with an environment (aka a closure). A quosure can be evaluated later in the stored environment to return a predictable result.



**rlang::quo(expr)** Quote contents as a quosure. Also **quos** to quote multiple expressions. `a <- 1 b <- 2 q <- quo(a + b) qs <- quos(a, b)`

**rlang::enquo(arg)** Call from within a function to quote what the user passed to an argument as a quosure. Also **enquos** for multiple args. `quote, this <- function(x) enquo(x) quote_these <- function(...) enquos(...)`

**rlang::new\_quosure(expr, env = caller\_env())** Build a quosure from a quoted expression and an environment. `new_quosure(expr(a + b), current_env())`

## Parsing and Deparsing

**"a + b"** **parse** **e** **deparse** **"a + b"**

**Parse** - Convert a string to a saved expression.

**Deparse** - Convert a saved expression to a string.

**rlang::parse\_expr(x)** Convert a string to an expression. Also **parse\_exprs**, **sym**, **parse\_quo**, **parse\_quos**. `e <- parse_expr("a+b") expr_text(e)`

## Building Calls

**rlang::call2(fn, ..., ns = NULL)** Create a call from a function and a list of args. Use `exec` to create and then evaluate the call. (See back page for `!!!`) `args <- list(x = 4, base = 2)`

**log(x = 4, base = 2)**

`call2("log", x = 4, base = 2)`  
`exec("log", !!!args)`  
`exec("log", !!!args)`



### EXPRESSION

**Quoted Expression** - An expression that has been saved by itself. A quoted expression can be evaluated later to return a result that will depend on the environment it is evaluated in

**rlang::expr(expr)** Quote contents. Also **exprs** to quote multiple expressions. `a <- 1 b <- 2 e <- expr(a + b) es <- exprs(a, b, a + b)`

**rlang::enexpr(arg)** Call from within a function to quote what the user passed to an argument. Also **enexprs** to quote multiple arguments. `quote_that <- function(x) enexpr(x) quote_those <- function(...) enexprs(...)`

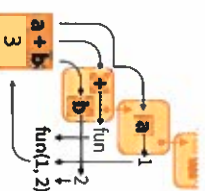
**rlang::ensym(x)** Call from within a function to quote what the user passed to an argument as a symbol, accepts strings. Also **ensyms**. `quote_names <- function(...) ensyms(...)`

## Evaluation

To evaluate an expression, R:

1. Looks up the symbols in the expression in the active environment (or a supplied one), followed by the environment's parents
2. Executes the calls in the expression

*The result of an expression depends on which environment it is evaluated in.*



### QUOTED EXPRESSION

**rlang::eval\_bare(expr, env = parent.frame())** Evaluate `expr` in `env`, using `bare/e`. `env = GlobalEnv()`



**QUOSURES** (and quoted `exprs`)

**rlang::eval\_tidy(expr, data = NULL, env = caller\_env())** Evaluate `expr` in `env`, using data as a **data mask**. Will evaluate quosures in their stored environment. `eval_tidy(q)`

**Data Mask** - If data is non-NULL, `eval_tidy` inserts data into the search path before `env`, matching symbols to names in data.

Use the pronoun `data$` to force a symbol to be matched in data, and `!!` (see back) to force a symbol to be matched in the environments.

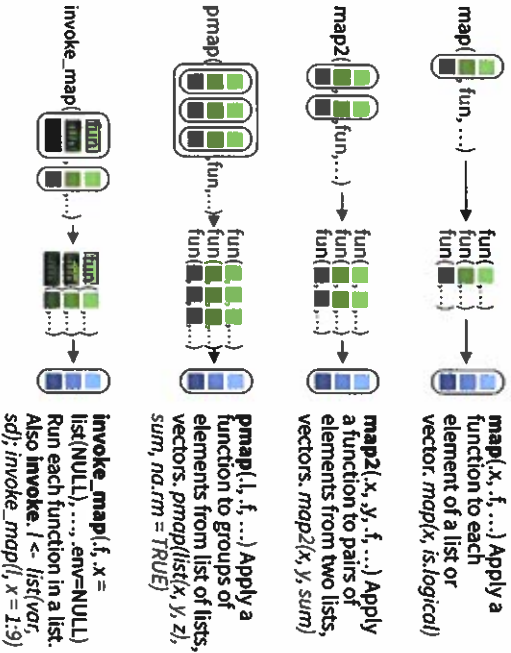
`a <- 1; b <- 2`  
`p <- quo(data$a + !!b)`  
`mask <- tibble(a = 5, b = 6)`  
`eval_tidy(p, data = mask)`

# Apply functions with purrr :: CHEAT SHEET



## Apply Functions

Map functions apply a function iteratively to each element of a list or vector.



**lmap(x, f, ...)** Apply function to each list-element of a list or vector.  
**imap(x, f, ...)** Apply f to each element of a list or vector and its index.

### OUTPUT

**map()**, **map2()**, **pmap()**, **imap()** and **invoke\_map** each return a list. Use a suffixed version to return the results as a specific type of flat vector, e.g. **map2\_chr**, **pmap\_lgl**, etc.

Use **walk**, **walk2**, and **pwalk** to trigger side effects. Each return its input invisibly.

### SHORTCUTS - within a purrr function:

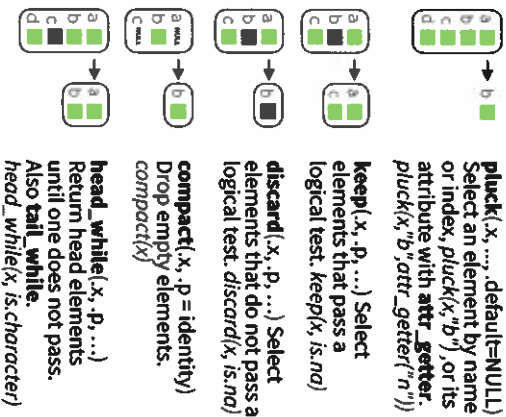
"name" becomes **function(x) x[["name"]]**, e.g. *map(f, "a")* extracts *a* from each element of *l*

~ **x.y** becomes **function(x, y) x.y**, e.g. *map2(f, p, ~ x + y)* becomes *map2(f, p, function(l, p) l + p)*

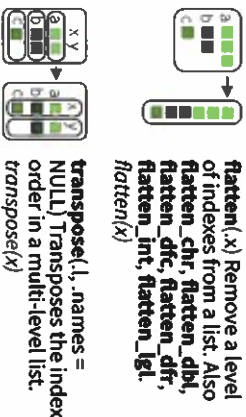
~ **..1..2** becomes **function(..1, ..2, etc) ..1..2**, e.g. *pmap(list(a, b, c), ~.3 + ..1..2)* becomes *function(a, b, c) c + a + b)*

## Work with Lists

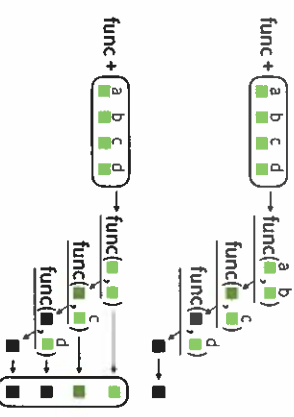
### FILTER LISTS



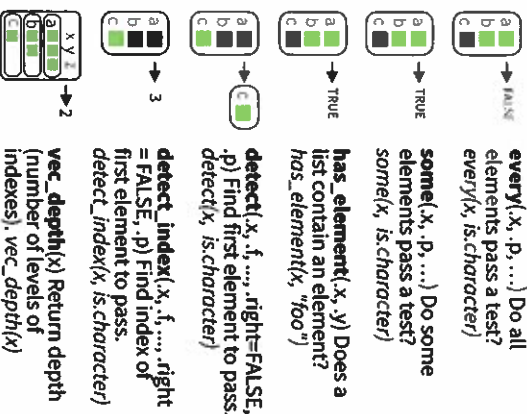
### RESHAPE LISTS



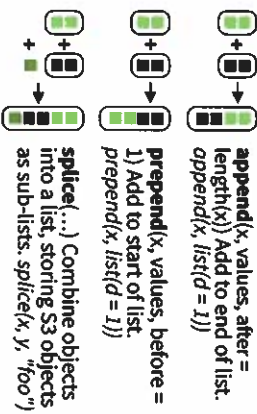
## Reduce Lists



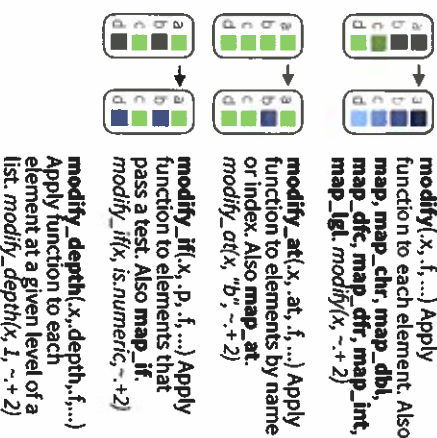
### SUMMARISE LISTS



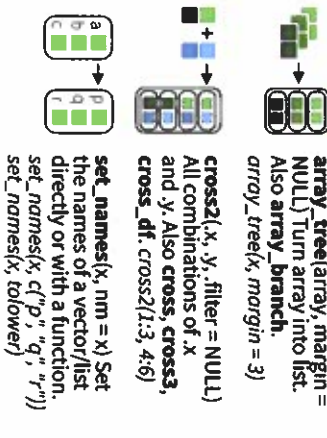
### JOIN (TO) LISTS



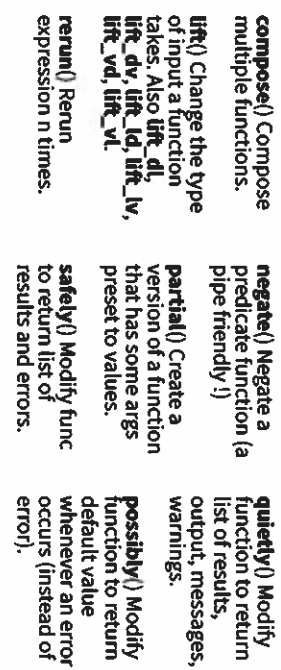
### TRANSFORM LISTS



### WORK WITH LISTS



## Modify function behavior



# Data Import :: CHEAT SHEET

R's **tidyverse** is built around **tidy data** stored in **tibbles**, which are enhanced data frames.



The front side of this sheet shows how to read text files into R with **readr**.



The reverse side shows how to create tibbles with **tibble** and to layout tidy data with **tidyr**.

## OTHER TYPES OF DATA

Try one of the following packages to import other types of files

- **haven** - SPSS, Stata, and SAS files
- **readxl** - excel files (.xls and .xlsx)
- **DBI** - databases
- **jsonlite** - json
- **xmll2** - XML
- **httr** - Web APIs
- **rvest** - HTML (Web Scraping)

## Save Data

Save **x**, an R object, to **path**, a file path, as:

### Comma delimited file

`write_csv(x, path, na = "NA", append = FALSE, col_names = !append)`

### File with arbitrary delimiter

`write_delim(x, path, delim = " ", na = "NA", append = FALSE, col_names = !append)`

### CSV for excel

`write_excel_csv(x, path, na = "NA", append = FALSE, col_names = !append)`

### String to file

`write_file(x, path, append = FALSE)`

### String vector to file, one element per line

`write_lines(x, path, na = "NA", append = FALSE)`

### Object to RDS file

`write_rds(x, path, compress = c("none", "gz", "bz2", "xz", ...))`

### Tab delimited files

`write_tsv(x, path, na = "NA", append = FALSE, col_names = !append)`

## Read Tabular Data - These functions share the common arguments:

`read_* (file, col_names = TRUE, col_types = NULL, locale = default_locale(), na = c("", "NA"), quoted_na = TRUE, comment = "#", trim_ws = TRUE, skip = 0, n_max = Inf, guess_max = min(1000, n_max), progress = interactive())`

### Comma Delimited Files

`read_csv("file.csv")`

To make file csv run:

`write_file(x = "a,b,c\n1,2,3\n4,5,NA", path = "file.csv")`

### Semi-colon Delimited Files

`read_csv2("file2.csv")`

`write_file(x = "a,b,c\n1,2;3\n4;5;NA", path = "file2.csv")`

### Files with Any Delimiter

`read_delim("file.txt", delim = "|")`

`write_file(x = "a|b|c\n1|2|3\n4|5|NA", path = "file.txt")`

### Fixed Width Files

`read_fwf("file.fwf", col_positions = c(1, 3, 5))`

`write_file(x = "a b c\n1 2 3\n4 5 NA", path = "file.fwf")`

### Tab Delimited Files

`read_tsv("file.tsv")` Also `read_table()`.

`write_file(x = "a\tb\tc\n1\t2\t3\n4\t5\tNA", path = "file.tsv")`

## USEFUL ARGUMENTS

### Example file

`write_file("a,b,c\n1,2,3\n4,5,NA", "file.csv")`

`f <- "file.csv"`

a,b,c  
1,2,3  
4,5,NA

### No header

`read_csv(f, col_names = FALSE)`

A B C  
1 2 3  
4 5 NA

### Provide header

`read_csv(f, col_names = c("x", "y", "z"))`

x y z  
A B C  
1 2 3  
4 5 NA

1 2 3  
4 5 NA

### Skip lines

`read_csv(f, skip = 1)`

A B C  
1 2 3

### Read in a subset

`read_csv(f, n_max = 1)`

A B C  
NA 2 3  
4 5 NA

### Missing Values

`read_csv(f, na = c("1", ""))`

## Data types

`readr` functions guess the types of each column and convert types when appropriate (but will NOT convert strings to factors automatically).

A message shows the type of each column in the result.

```
## Parsed with column specification:
## cols(
##   age = col_integer(),
##   sex = col_character(),
##   earn = col_double()
## )
```

earn is a double (numeric)

sex is a character

1. Use `problems()` to diagnose problems.  
`x <- read_csv("file.csv"); problems(x)`

2. Use a `col_` function to guide parsing.
  - `col_guess()` the default
  - `col_character()`
  - `col_double()`, `col_euro_double()`
  - `col_datetime(format = "")` Also
  - `col_date(format = "")`, `col_time(format = "")`
  - `col_factor(levels, ordered = FALSE)`
  - `col_integer()`
  - `col_logical()`
  - `col_number()`, `col_numeric()`
  - `col_skip()`

```
x <- read_csv("file.csv", col_types = cols(
  A = col_double(),
  B = col_logical(),
  C = col_factor()))
```

3. Else, read in as character vectors then parse with a `parse_` function.
  - `parse_guess()`
  - `parse_character()`
  - `parse_datetime()` Also `parse_date()` and `parse_time()`
  - `parse_double()`
  - `parse_factor()`
  - `parse_integer()`
  - `parse_logical()`
  - `parse_number()`
  - `x$A <- parse_number(x$A)`

## Read Non-Tabular Data

### Read a file into a single string

`read_file(file, locale = default_locale())`

### Read each line into its own string

`read_lines(file, skip = 0, n_max = -1L, na = character(), locale = default_locale(), progress = interactive())`

### Read Apache style log files

`read_log(file, col_names = FALSE, col_types = NULL, skip = 0, n_max = -1, progress = interactive())`

### Read a file into a raw vector

`read_file_raw(file)`

### Read each line into a raw vector

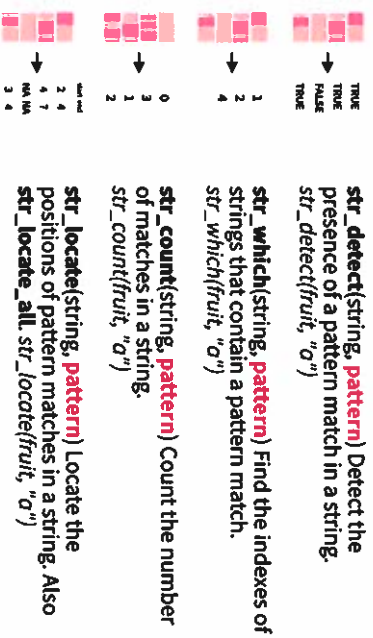
`read_lines_raw(file, skip = 0, n_max = -1L, progress = interactive())`

# String manipulation with stringr :: CHEAT SHEET

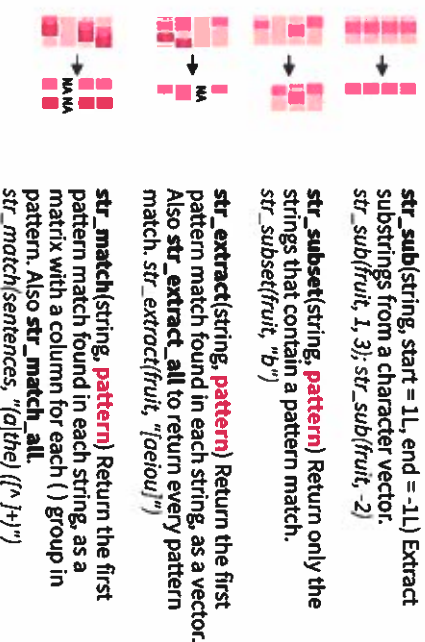


The `stringr` package provides a set of internally consistent tools for working with character strings, i.e. sequences of characters surrounded by quotation marks.

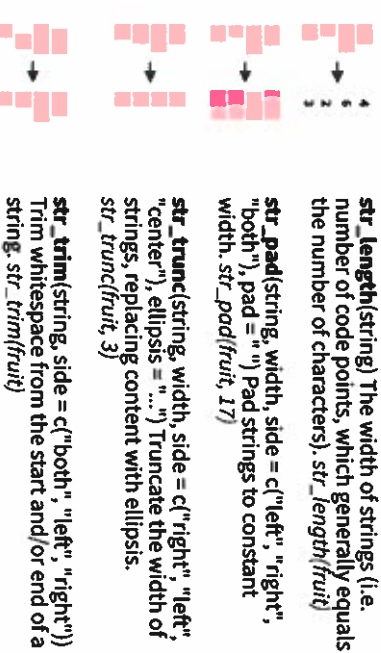
## Detect Matches



## Subset Strings



## Manage Lengths



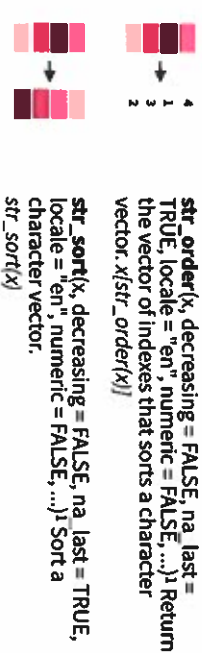
## Mutate Strings



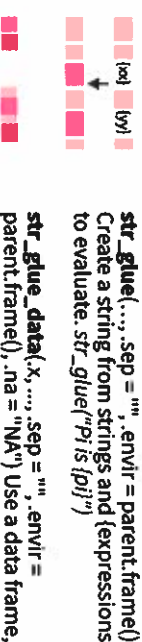
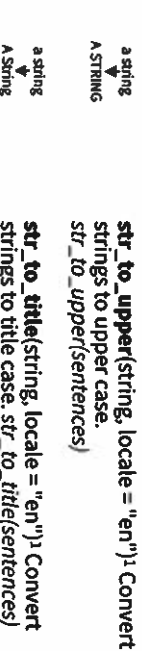
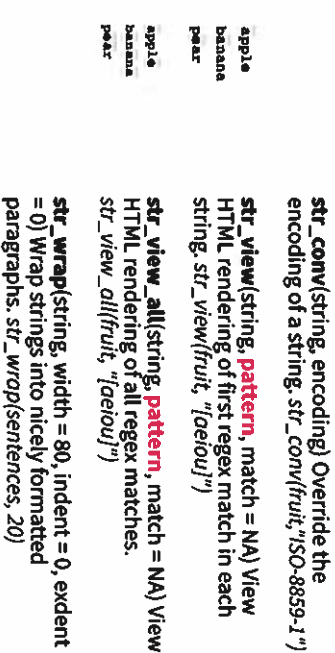
## Join and Split



## Order Strings



## Helpers

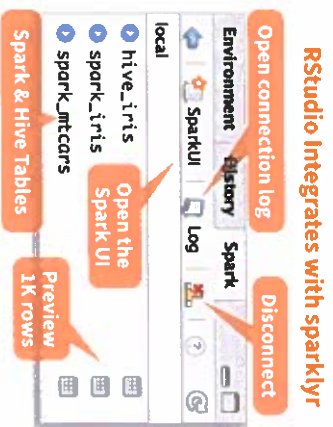


# Data Science in Spark with Sparklyr :: CHEAT SHEET

## Intro

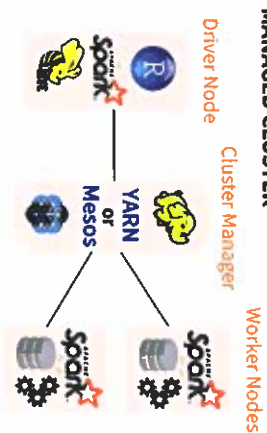
**sparklyr** is an R interface for Apache Spark™, it provides a complete **dplyr** backend and the option to query directly using **Spark SQL** statement. With sparklyr, you can orchestrate distributed machine learning using either **Spark's MLlib** or **H2O Sparkling Water**.

Starting with **version 1.044**, **RStudio Desktop**, **Server and Pro** include **integrated support** for the **sparklyr** package. You can create and manage connections to Spark clusters and local Spark instances from inside the IDE.

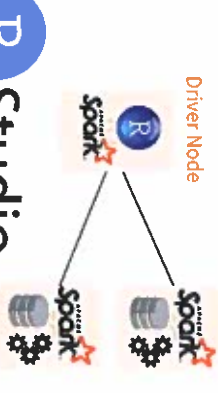


## Cluster Deployment

### MANAGED CLUSTER

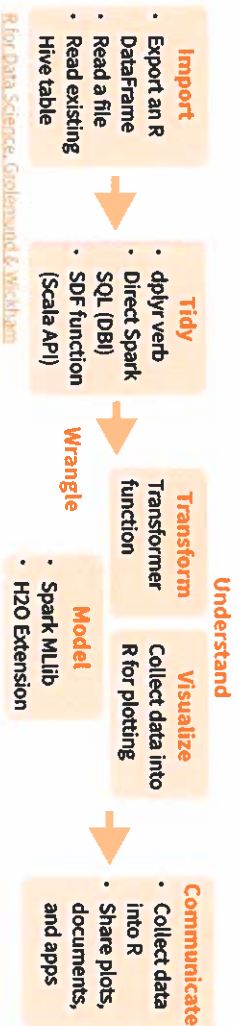


### STAND ALONE CLUSTER



## R Studio

## Data Science Toolchain with Spark + sparklyr



## Getting Started

### LOCAL MODE (No cluster required)

1. Install a local version of Spark:  
`spark_install("2.0.1")`
2. Open a connection  
`sc <- spark_connect(master = "local")`

### ON A YARN MANAGED CLUSTER

1. Install RStudio Server or RStudio Pro on one of the existing nodes, preferably an edge node
2. Locate path to the cluster's Spark Home Directory, it normally is "/usr/lib/spark"
3. Open a connection  
`spark_connect(master="yarn-client", version = "1.6.2", spark_home = [Cluster's Spark path])`

### ON A MESOS MANAGED CLUSTER

1. Install RStudio Server or Pro on one of the existing nodes
2. Locate path to the cluster's Spark directory
3. Open a connection  
`spark_connect(master="[mesos URL]", version = "1.6.2", spark_home = [Cluster's Spark path])`

### ON A SPARK STANDALONE CLUSTER

1. Install RStudio Server or RStudio Pro on one of the existing nodes or a server in the same LAN
2. Install a local version of Spark:  
`spark_install (version = "2.0.1")`
3. Open a connection  
`spark_connect(master="spark://host:port", version = "2.0.1", spark_home = spark_home_dir())`

### USING LIVY (Experimental)

1. The Livy REST application should be running on the cluster
2. Connect to the cluster  
`sc <- spark_connect(method = "livy", master = "http://host:port")`

## Tuning Spark

### EXAMPLE CONFIGURATION

```
config <- spark_config()
config$spark.executor.cores <- 2
config$spark.executor.memory <- "4G"
sc <- spark_connect(master="yarn-client",
config = config, version = "2.0.1")
```

### IMPORTANT TUNING PARAMETERS with defaults

```
• spark.yarn.am.cores
• spark.yarn.am.memory 512m
• spark.network.timeout 120s
• spark.executor.memory 1g
• spark.executor.cores 1
• spark.executor.instances
• spark.executor.extraJavaOptions
• spark.network.timeout 120s
• spark.executor.heartbeatInterval 10s
• sparklyr.shell.executor-memory
• sparklyr.shell.driver-memory
```

## Using sparklyr

A brief example of a data analysis using Apache Spark, R and sparklyr in local mode

```
library(sparklyr); library(dplyr); library(ggplot2);
library(tidy);
set.seed(100)
Install Spark locally
```

```
spark_install("2.0.1")
Connect to local version
```

```
sc <- spark_connect(master = "local")
```

```
import_iris <- copy_to(sc, iris, "spark_iris",
overwrite = TRUE)
```

```
Copy data to Spark memory
partition_iris <- sdf_partition(
import_iris, training=0.5, testing=0.5)
Partition data
```

```
sdf_register(partition_iris,
c("spark_iris_training", "spark_iris_test"))
```

Create a hive metadata for each partition

```
tidy_iris <- tbl(sc, "spark_iris_training") %>%
select(Species, Petal.Length, Petal.Width)
```

```
Spark ML Decision Tree Model
model_iris <- tidy_iris %>%
ml_decision_tree(response="Species",
features=c("Petal.Length","Petal.Width"))
```

```
test_iris <- tbl(sc, "spark_iris_test")
```

```
Create reference to Spark table
pred_iris <- sdf_predict(
model_iris, test_iris) %>%
collect
```

```
Bring data back into R memory for plotting
pred_iris %>%
inner_join(data.frame(prediction=0.2,
lab=model_iris$model$parameters$labels) %>%
ggplot(aes(Petal.Length, Petal.Width, col=lab)) +
geom_point())
```

```
spark_disconnect(sc)
Disconnect
```

