# RStudio®

```
R Studio
https://beta.rstudio.org/s/c573a7568bf7fb4b43cbb/

File    Edit    Code    View    Plots    Session    Build    Debug    Tools    Help                                    Radar    R 3.2.0

Go to file/function                          Addins                                                      Sessions (3)

chicagoFood.R                                                        Run    Source    R Script        Environment    History    Git
                                                                                                       Import Dataset
  1  url <- "http://data.cityofchicago.org/api/views/4ijn-s7e5/rows.csv?a       Global Environment
  2  data <- read.csv(url, header = TRUE) # takes a minute...
  3  names(data) <- tolower(names(data))                                        Data
  4  data1 <- subset(data, risk %in% c("Risk 1 (High)","Risk 2 (Medium)",       data      118607 obs. of 17 variables
  5  data1$risk <- droplevels(data1$risk)                                       data1     50 obs. of 17 variables
  6
  7  data1 <- data1[1:50,]                                                      Files    Plots    Packages    Help    Viewer
  8  library(leaflet)                                                           Zoom    Export
  9  leaflet(data1) %>%
 10     addTiles() %>%
 11     addMarkers(lat = ~latitude, lng = ~longitude)

11:33   (Top Level)

Console  ~/Radar/
> data1 <- subset(data, risk %in% c("Risk 1 (High)","Risk 2 (Medium)","Ris
k 3 (Low)"))
> data1$risk <- droplevels(data1$risk)
>
> data1 <- data1[1:50,]
> library(leaflet)
> leaflet(data1) %>%
+    addTiles() %>%
+    addMarkers(lat = ~latitude, lng = ~longitude)
>
```

# Advanced sparklyr features

# ADVANCED SPARKLYR FEATURES

✓ Overview
✓ Scaling R in Spark with spark_apply()
✓ Spark from your desktop using Livy
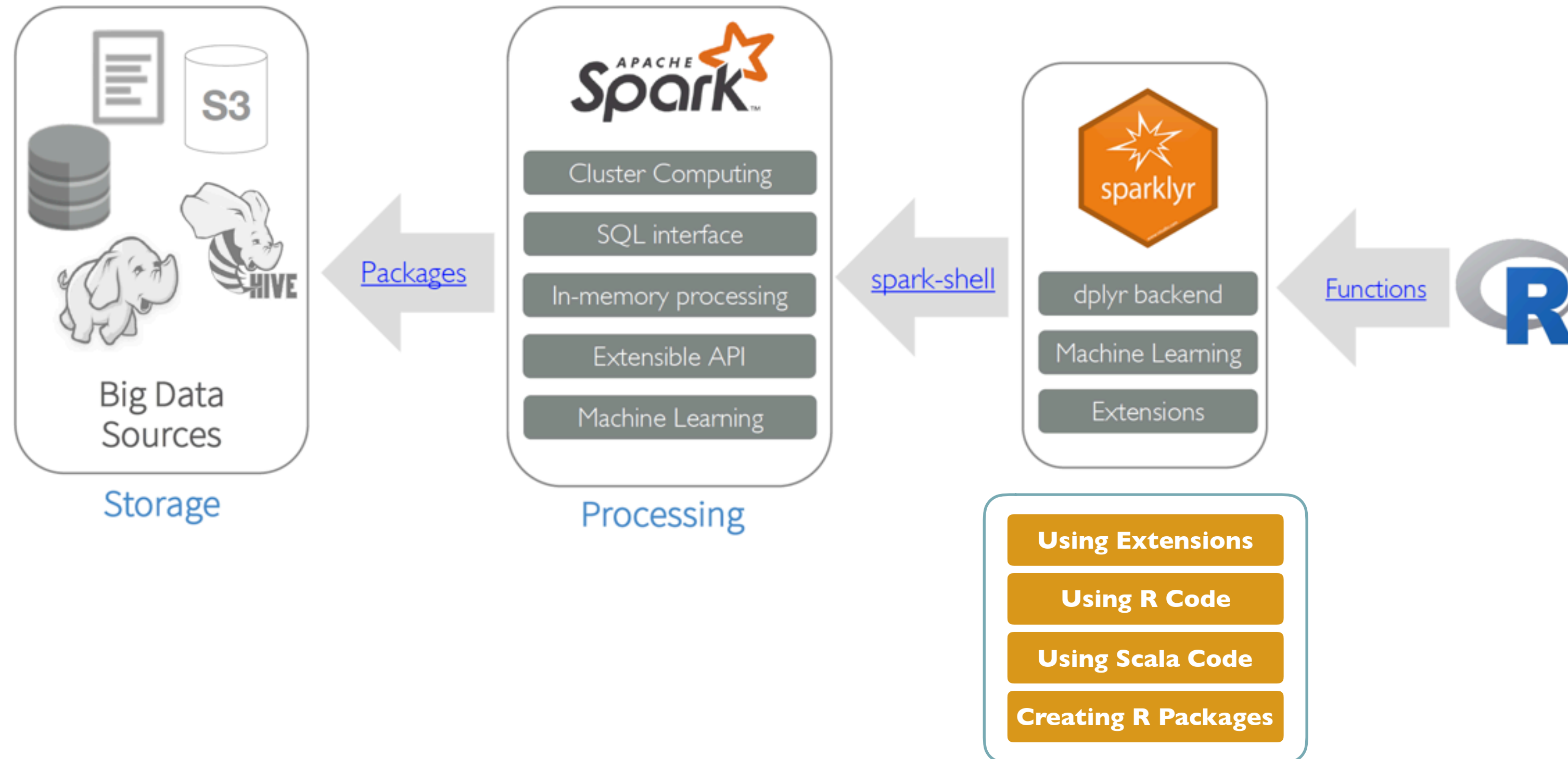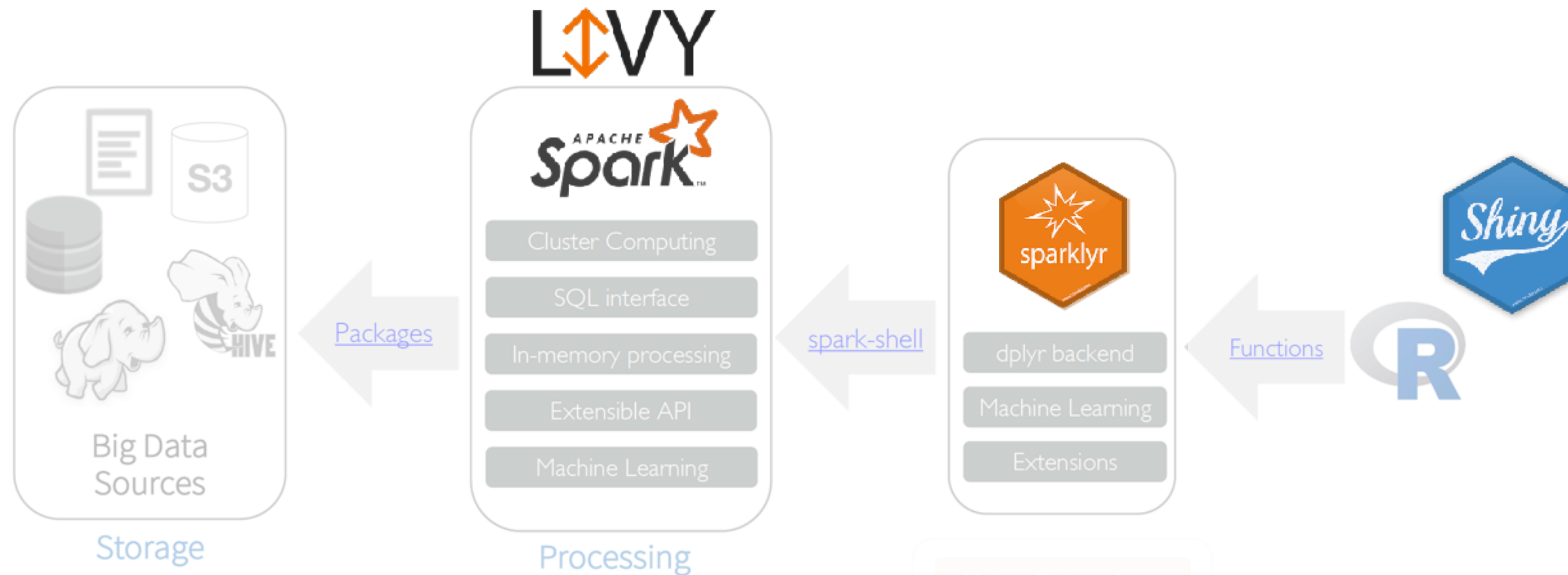✓ Spark Applications with Shiny
✓ Questions

# Overview

# OVERVIEW

# OVERVIEW



**LIVY**

Apache Spark
- Cluster Computing
- SQL interface
- In-memory processing
- Extensible API
- Machine Learning

Packages

spark-shell

sparklyr
- dplyr backend
- Machine Learning
- Extensions

Functions

R

Shiny

S3

Big Data Sources

Storage

Processing

**Common Crawl**

240 TiB
72K Files

spark_apply()

✓ Scaling R in Spark with spark_apply()
✓ Spark from your desktop using Livy
✓ Spark Applications with Shiny

# USE CASES

✓ Leverage R skills

```
1
2  spark_apply(iris_tbl, function(e) sapply(e[,1:4], jitter))
3
```

✓ Complement R

```
1
2  spark_apply(
3    iris_tbl,
4    function(e) broom::tidy(lm(Petal_Width ~ Petal_Length, e)),
5    colums = c("term", "estimate", "std.error", "statistic", "p.value"),
6    group_by = "Species"
7  )
8
```

# USING SPARK_APPLY()

```
1
2  spark_apply(
3    iris_tbl,
4    function(e) sapply(
5      e[,1:4], jitter
6    )
7  )
8
```

✓ Install R in every node, once per cluster.

✓ spark_apply() distributes packages, once per session.

✓ spark_apply() distributes your code, once per call.

```
1  sapply(e[,1:4], jitter)
```

```
1  sapply(e[,1:4], jitter)
```

```
1  sapply(e[,1:4], jitter)
```

```
1  sapply(e[,1:4], jitter)
```

# Parsing and Filtering the CommonCrawl



We build and maintain an open repository of web crawl data that can be accessed and analyzed by anyone.

Need years of free web page data to help change the world.

240 TiB

72K Files

```r
library(sparkwarc)
spark_read_warc(sc, "warcs", cc_warc(1, 100))
```

**sparkwarc uses**
**Rcpp:: + spark_apply()**

```r
df <- spark_apply(paths_tbl, function(df) {
  entries <- apply(df, 1, function(path) {
    if (grepl("s3n://", path)) {
      path <- sub("s3n://commoncrawl/", "https://commoncrawl.s3.amazonaws.com/", path)
      temp_warc <- tempfile(fileext = ".warc.gz")
      download.file(url = path, destfile = temp_warc)
      path <- temp_warc
    }

    sparkwarc::rcpp_read_warc(path, filter = match_warc, include = match_line)
  })

  if (nrow(df) > 1) do.call("rbind", entries) else data.frame(entries)
}, names = c("tags", "content")) %>% spark_dataframe()
```

github.com/javierluraschi/sparkwarc

# Livy

# USING LIVY

[ http ]

```
1
2  config <- livy_config("<username>",
3                        "<password>")
4
5  sc <- spark_connect(master = "<address>",
6                      method = "livy",
7                      config = config)
8
```
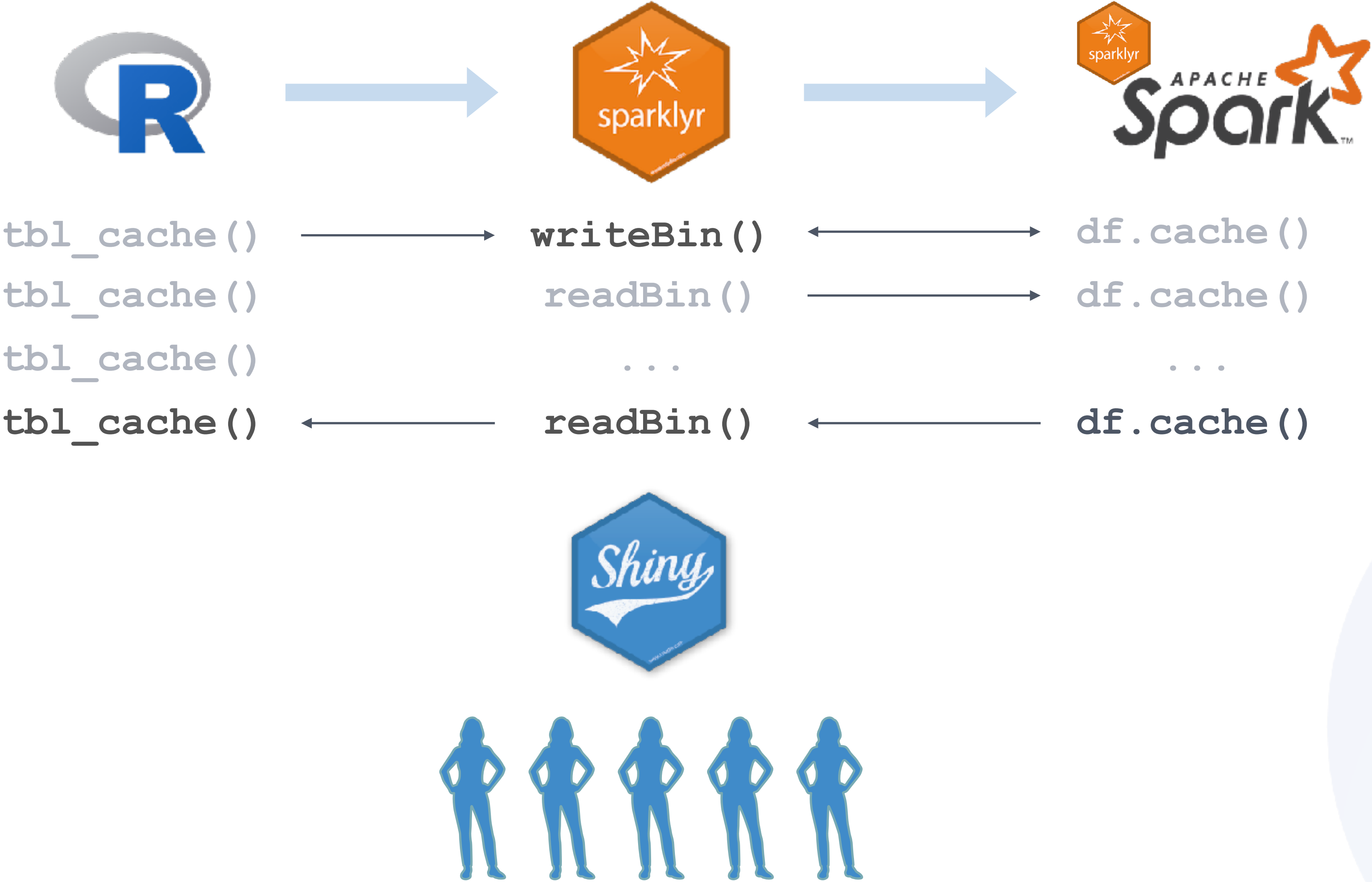
# SPARKLYR CONNECTIONS



tbl_cache()  ⟶  **writeBin()**  ⟷  df.cache()

tbl_cache()  readBin()  ⟶  df.cache()

tbl_cache()  ...  ...

**tbl_cache()**  ⟵  **readBin()**  ⟵  **df.cache()**

# SHINY AND SPARKLYR

```
1
2    library(shiny)
3    library(dplyr)
4    library(sparklyr)
5
6    sc <- spark_connect(master = "local")
7    faithful_tbl <- copy_to(sc, faithful, overwrite = TRUE)
8
9    ui <- fluidPage(
10     titlePanel("Old Faithful Geyser Data"),
11     sidebarLayout(
12       sidebarPanel(
13         sliderInput("bins",
14                     "Number of bins:",
15                     min = 1,
16                     max = 50,
17                     value = 30)
18       ),
19       mainPanel(
20         plotOutput("distPlot")
21       )
22     )
23   )
24
25   server <- function(input, output) {
26     output$distPlot <- renderPlot({
27       x    <- faithful_tbl %>% pull(waiting)
28       bins <- seq(min(x), max(x), length.out = input$bins + 1)
29
30       hist(x, breaks = bins, col = 'darkgray', border = 'white')
31     })
32   }
33
34   shinyApp(ui = ui, server = server)
35
```

## Setup

```
3    library(dplyr)
4    library(sparklyr)
5
6    sc <- spark_connect(master = "local")
7    faithful_tbl <- copy_to(sc, faithful, overwrite = TRUE)
```

## Server

```
27          x      <- faithful_tbl %>% pull(waiting)
```

# USEFUL WEBSITES

RStudio and Spark:
[spark.rstudio.com](spark.rstudio.com)

Shiny:
[shiny.rstudio.com](shiny.rstudio.com)

Apache Livy:
[livy.incubator.apache.org/](livy.incubator.apache.org/)

Sparklyr Github:
[github.com/rstudio/sparklyr](github.com/rstudio/sparklyr)

# Questions

**Open Source & Free**
Desktop: http://www.rstudio.com/products/rstudio/download/
RStudio Server: http://www.rstudio.com/products/rstudio/download-server/
Shiny Server: http://www.rstudio.com/products/shiny/download-server/
shinyapps.io beta: https://www.shinyapps.io/admin/#/signup

**45 Day Evaluation of Pro Products**

RStudio Server Pro: http://www.rstudio.com/products/rstudio-server-pro/evaluation/
Shiny Server Pro: http://www.rstudio.com/products/shiny-server-pro/evaluation/

# PLEASE STAY IN TOUCH

Blog - http://rviews.rstudio.com/

Blog - http://blog.rstudio.org/

Twitter - @rstudio #rstats http://twitter.com/rstudio/

GitHub - https://github.com/rstudio/

LinkedIn - https://linkedin.com/company/rstudio-inc

Facebook - https://www.facebook.com/pages/RStudio-inc

Google+ - https://plus.google.com/110704473211154995841/posts