

Advanced sparklyr features

ADVANCED SPARKLYR FEATURES

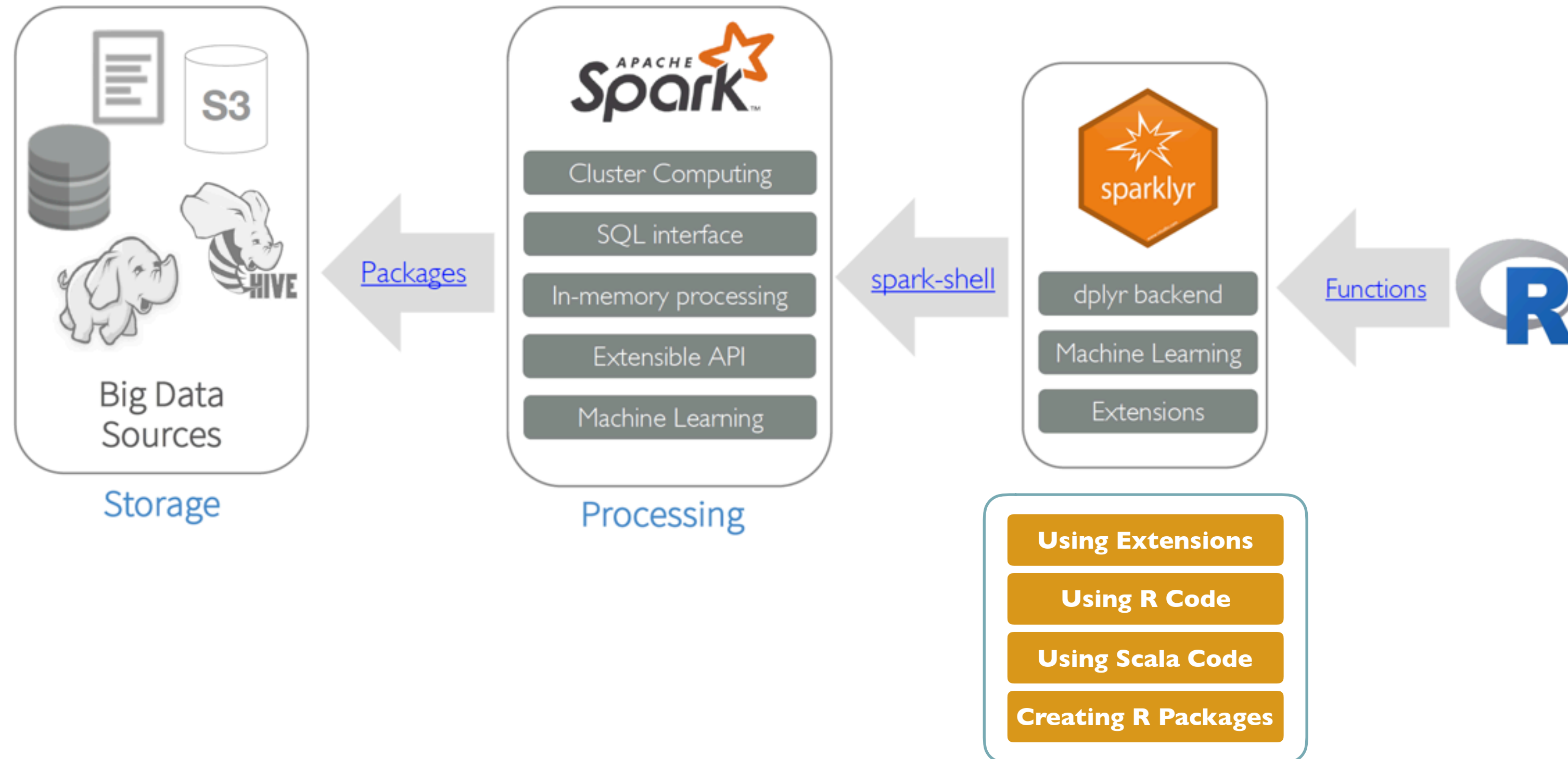
- ✓ Overview
- ✓ Scaling R in Spark with `spark_apply()`
- ✓ Spark from your desktop using Livy
- ✓ Spark Applications with Shiny
- ✓ Questions



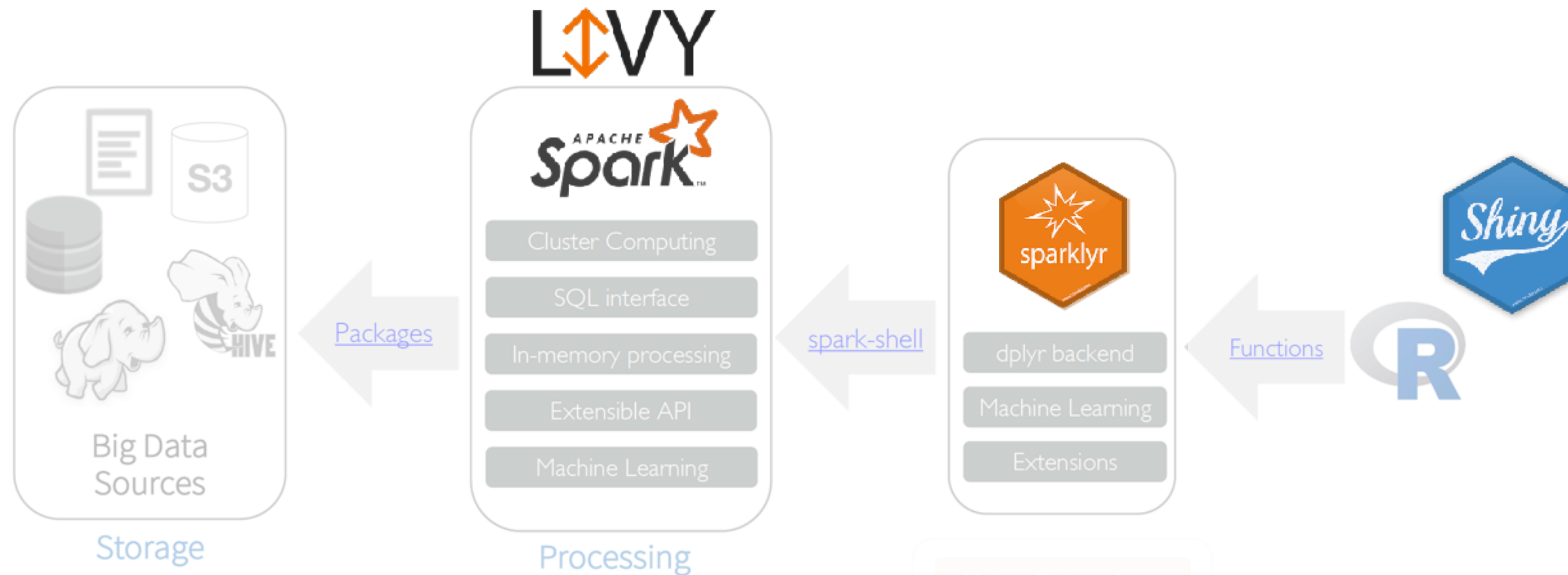
Overview



OVERVIEW



OVERVIEW



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

Scaling R



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   columns = c("term", "estimate", "std.error", "statistic", "p.value"),  
6   group_by = "Species"  
7 )  
8
```



USING SPARK_APPLY()

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

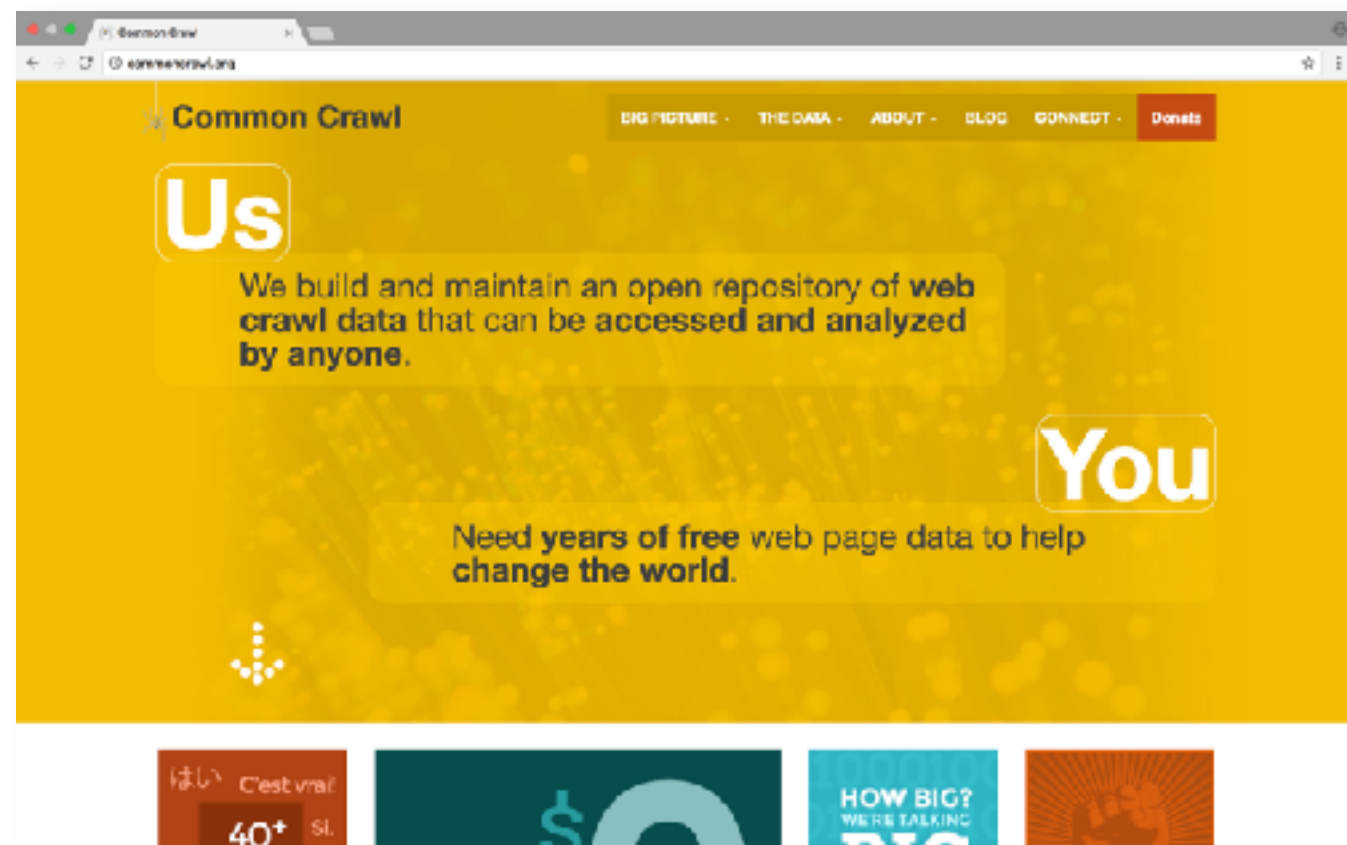


```
1 sapply(e[,1:4], jitter)
```



A large, light blue watermark of the R logo, consisting of a stylized 'R' inside a circle, is positioned on the right side of the slide.

PARSING AND FILTERING THE COMMONCRAWL

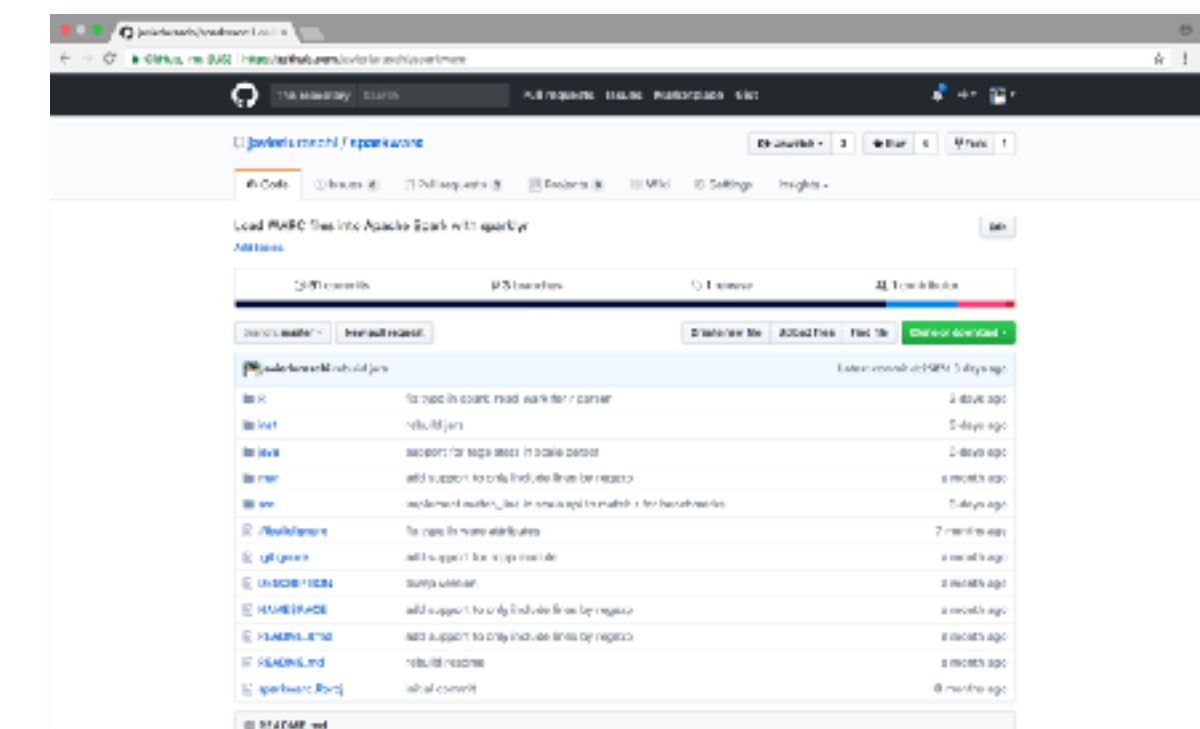


240 TiB
72K Files

```
13 // [[Rcpp::export]]
14 Dataframe rcpp_read_warc(std::string path, std::string filter, std::string include) {
15     FILE *fp = fopen(path.c_str(), "rb");
16     if (!fp) Rcpp::stop("Failed to open WARC file.");
17     gzFile gzf = gzopen(fp, "rb");
18     if (!gzf) Rcpp::stop("Failed to open WARC as a compressed file.");
19
20     const int buffer_size = 4 * 1024;
21     char buffer[buffer_size] = {"\0"};
22     std::list<std::string> warc_entries;
23     const int warc_max_size = 40 * 1024;
24     std::string warc_entry;
25     warc_entry.reserve(warc_max_size);
26     bool one_matched = false;
27     const std::string warc_separator = "WARC/1.0";
28
29     long stats_tags_total = 0;
30     std::list<long> warc_stats;
31
32     while(gzgets(gzf, buffer, buffer_size) != Z_NULL) {
33         std::string line(buffer);
34
35         if (!filter.empty() && !one_matched) {
36             one_matched = line.find(filter) != std::string::npos;
37         }
38
39         if (std::string(line).substr(0, warc_separator.size()) == warc_separator && warc_entry.size() > 0) {
40             if (!filter.empty() || one_matched) {
41                 warc_entries.push_back(warc_entry);
42                 warc_stats.push_back(stats_tags_total);
43                 stats_tags_total = 0;
44             }
45
46             one_matched = false;
47             warc_entry.clear();
48         }
49
50         if (!include.empty() && line.find(include) != std::string::npos)
51             warc_entry.append(line);
52
53         std::size_t tag_start = rcpp_find_tag(line, 0);
54         while(tag_start != std::string::npos) {
55             stats_tags_total += 1;
56             tag_start = rcpp_find_tag(line, tag_start + 1);
57         }
58     }
59
60     if (gzf) gzclose(gzf);
61     if (fp) fclose(fp);
62
63     long idxEntry = 0;
64     CharacterVector results(warc_entries.size());
65     std::for_each(warc_entries.begin(), warc_entries.end(), [&results, &idxEntry](std::string &entry) {
66         results[idxEntry++] = entry;
67     });
68
69     long idxStat = 0;
70     NumericVector stats(warc_stats.size());
71     std::for_each(warc_stats.begin(), warc_stats.end(), [&stats, &idxStat](long &stat) {
72         stats[idxStat++] = stat;
73     });
74
75     return Dataframe::create(Named("tags") = stats, _["content"] = results);
76 }
```

Rcpp::

or use sparkwarc::



+

```
62 df <- spark_apply(paths_tbl, function(df) {
63     entries <- apply(df, 1, function(path) {
64         if (grepl("s3n://", path)) {
65             path <- sub("s3n://commoncrawl/", "https://commoncrawl.s3.amazonaws.com/", path)
66             temp_warc <- tempfile(fileext = ".warc.gz")
67             download.file(url = path, destfile = temp_warc)
68             path <- temp_warc
69         }
70
71         sparkwarc::rcpp_read_warc(path, filter = match_warc, include = match_line)
72     })
73
74     if (nrow(df) > 1) do.call("rbind", entries) else data.frame(entries)
75 }, names = c("tags", "content")) %>% spark_dataframe()
```

spark_apply

Livy



USING LIVY

✓ LIVY installed in one node



[http]



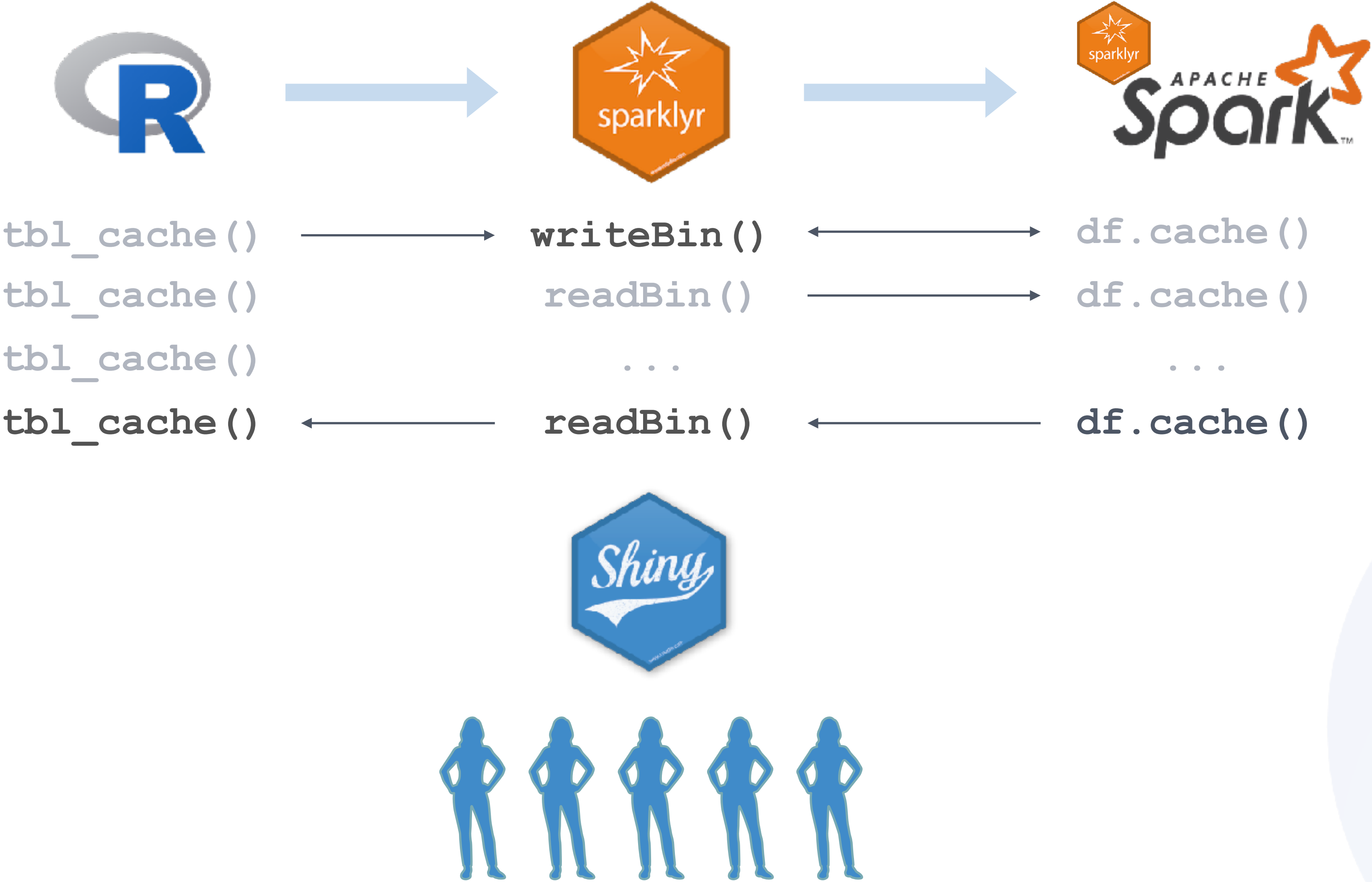
```
1  
2 config <- livy_config_auth("<username>",  
3                             "<password>")  
4  
5 sc <- spark_connect(master = "<address>",  
6                      method = "livy",  
7                      config = config)  
8
```



Shiny



SPARKLYR CONNECTIONS



SHINY AND SPARKLYR

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)
```

```
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
```


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>