# Data Science



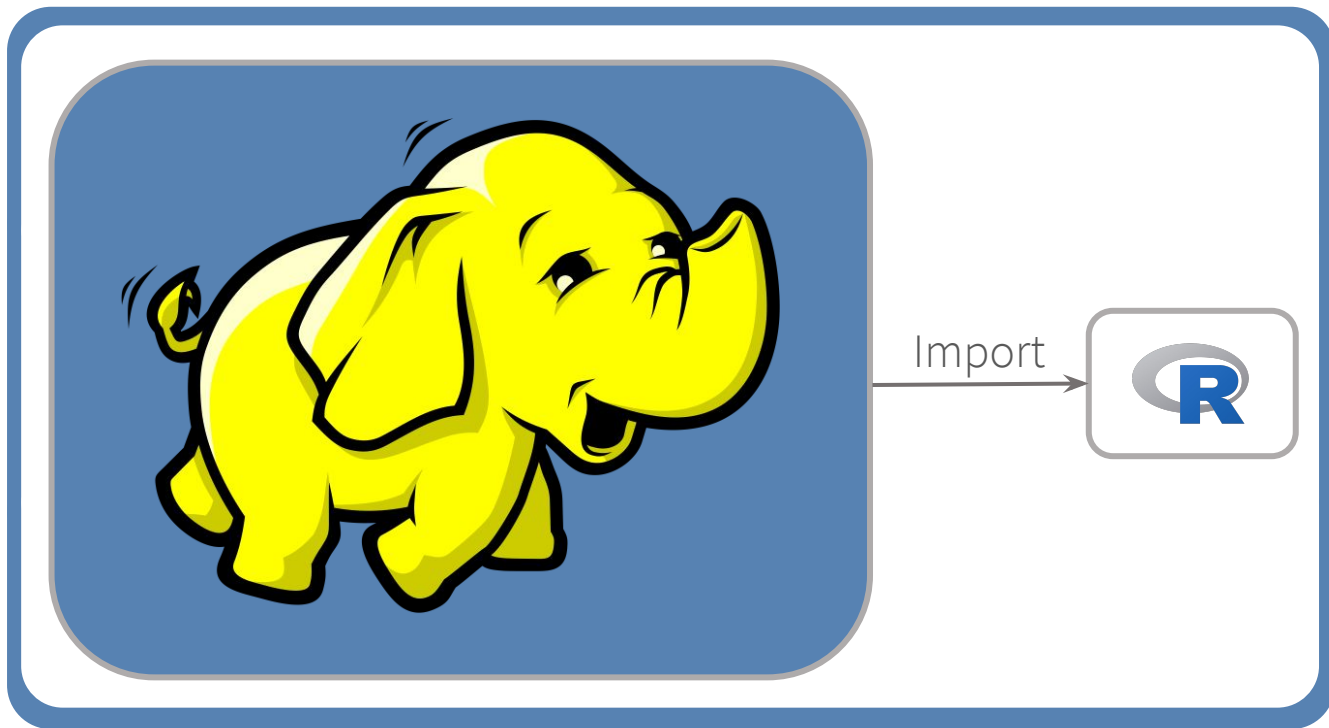R for Data Science, Wickham & Grolemund

# Big Data as a Data Source

Too large to download into memory



Import

R Studio

# R inside Spark

Only appropriate for Embarrassingly Parallel cases

# We keep forgetting…

"Spark is a unified analytics engine for large-scale data processing"

# A better way… R as an *interface* for Spark

ml_linear_regression()

R function wraps Spark Scala code

```scala
val lr = new LinearRegression()
val lrModel = lr.fit(training)
```

# 1. Most Data Science steps can run inside Spark

# 2. Best of both worlds: Spark API + R ecosystem

# 3. Hand off to Production: ML Pipelines built in R

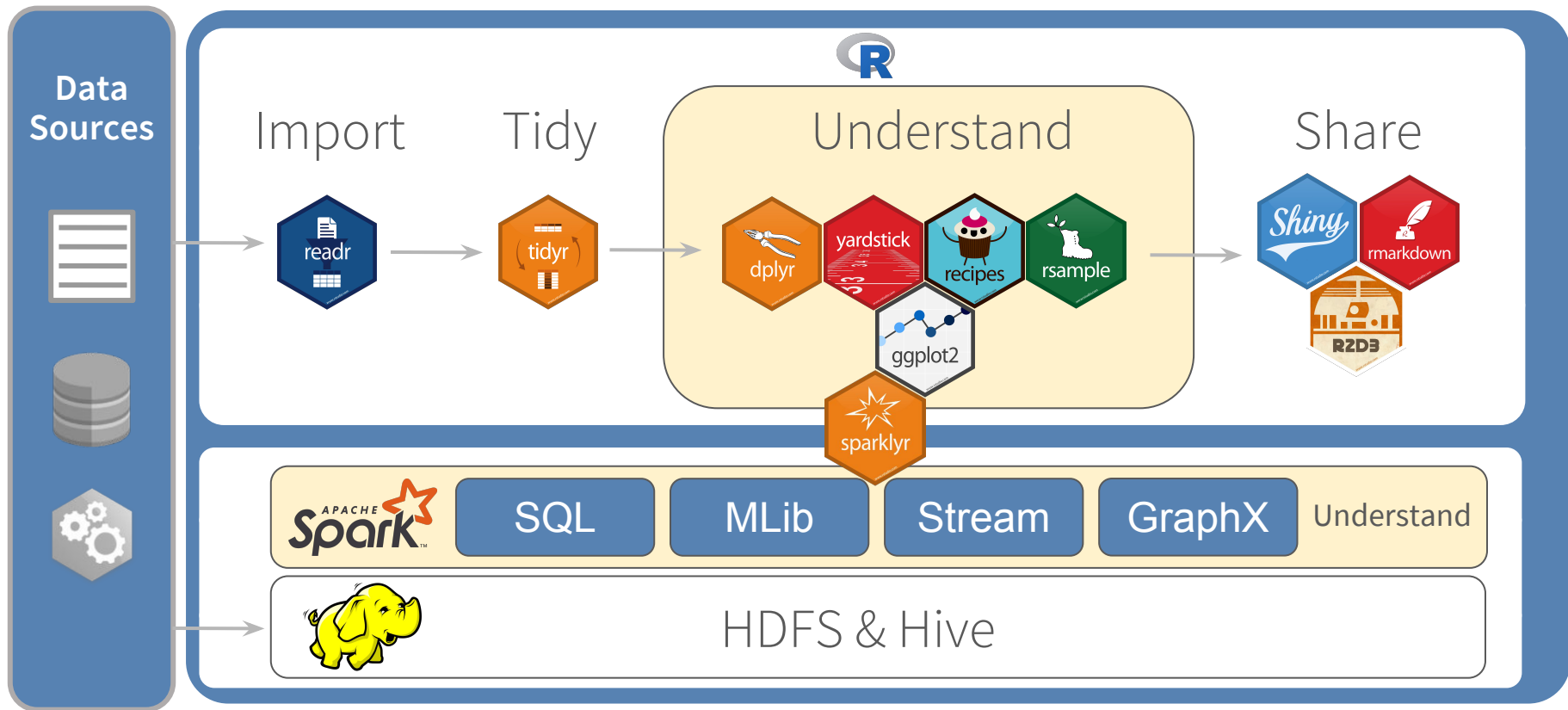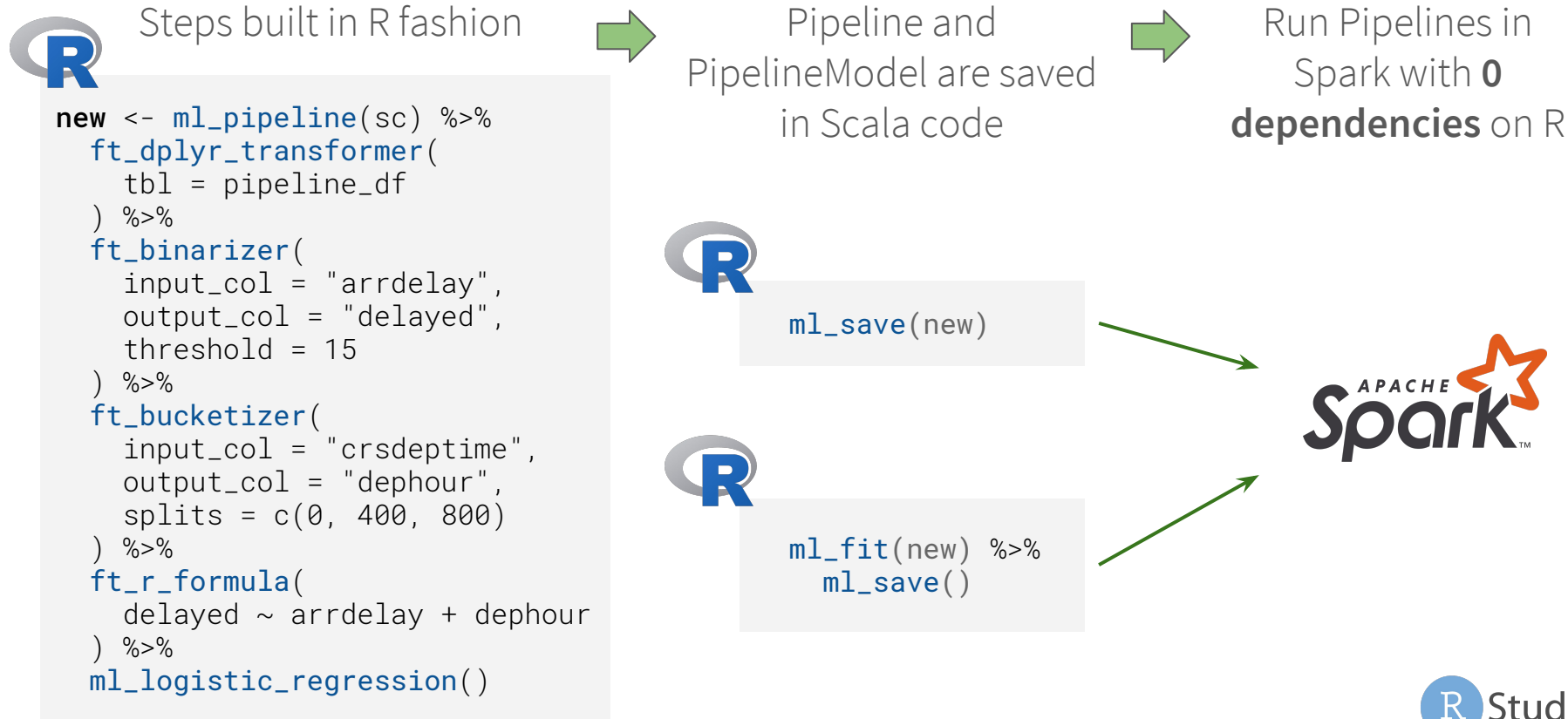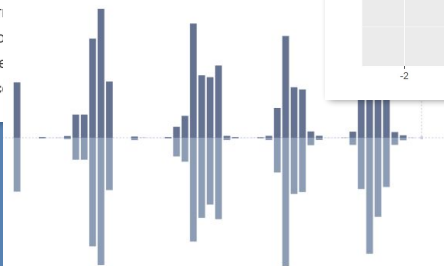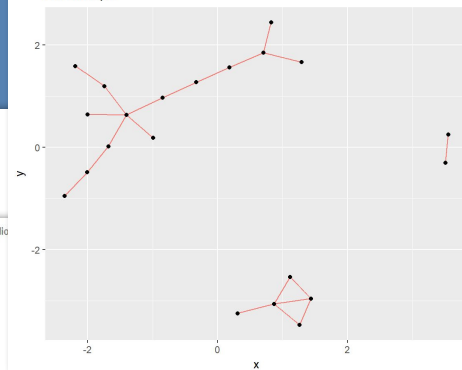Steps built in R fashion → Pipeline and PipelineModel are saved in Scala code → Run Pipelines in Spark with **0 dependencies** on R

```r
new <- ml_pipeline(sc) %>%
  ft_dplyr_transformer(
    tbl = pipeline_df
  ) %>%
  ft_binarizer(
    input_col = "arrdelay",
    output_col = "delayed",
    threshold = 15
  ) %>%
  ft_bucketizer(
    input_col = "crsdeptime",
    output_col = "dephour",
    splits = c(0, 400, 800)
  ) %>%
  ft_r_formula(
    delayed ~ arrdelay + dephour
  ) %>%
  ml_logistic_regression()
```

```r
ml_save(new)
```

```r
ml_fit(new) %>%
  ml_save()
```

APACHE Spark™

R Studio®

# Interface highlights
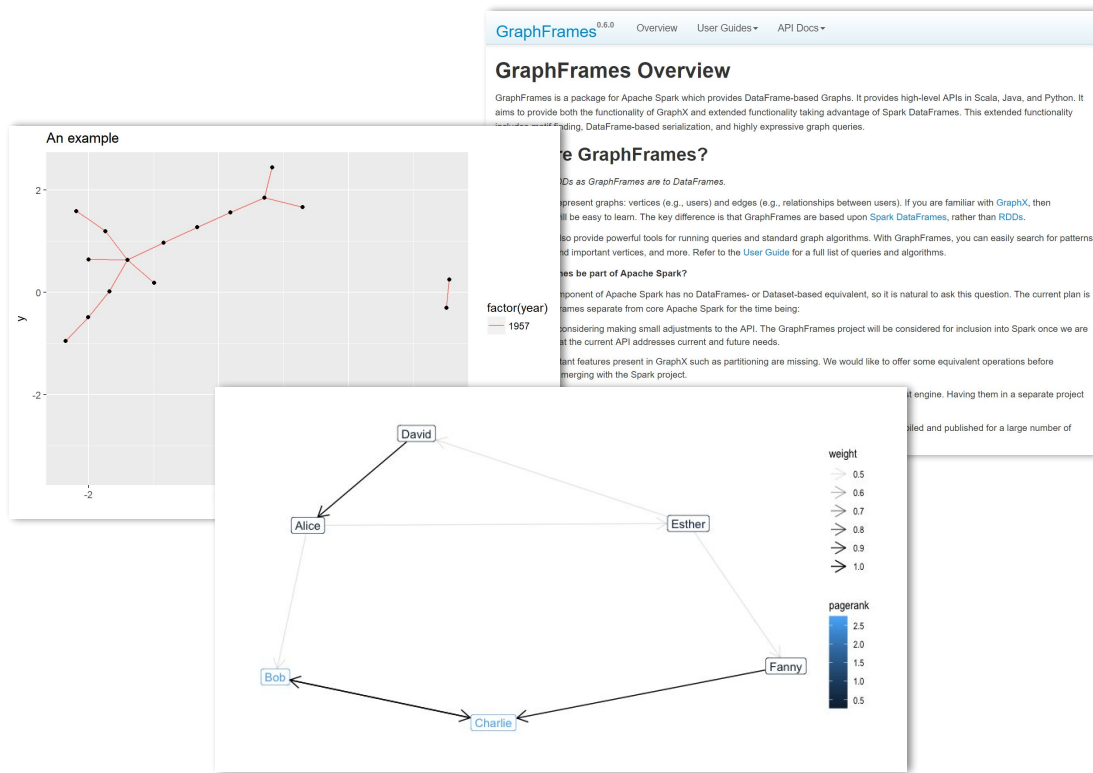
# Streaming

- Ability to run `dplyr`, SQL, and PipelineModels against a stream
- Read & write stream results to Spark memory and files
- An out-of-the box graph visualization to monitor the stream
- `reactiveSpark()` function allows Shiny apps to poll the contents of the stream



https://spark.rstudio.com/guides/streaming/

# Graph analysis

- Support for GraphFrames which aims to provide the functionality of GraphX.

- Perform graph algorithms such as: PageRank, ShortestPaths and many others

- Designed to work with sparklyr and the sparklyr extensions



https://spark.rstudio.com/graphframes/

# Production pipelines with Mleap

`mleap` is a `sparklyr` extension that provides an interface to MLeap, which allows us to take Spark pipelines to production.