# Studio®

## You Love Big Data!
## So Does R.

Alex K Gold
@alexkgold
#rstats

https://github.com/akgold/bdl_2019

@alexkgold

# In-Memory

- doFuture
- RStudio Server Pro Launcher

st Enough (profvis)

Rcpp

@alexkgold

R Studio®

# In-Memory



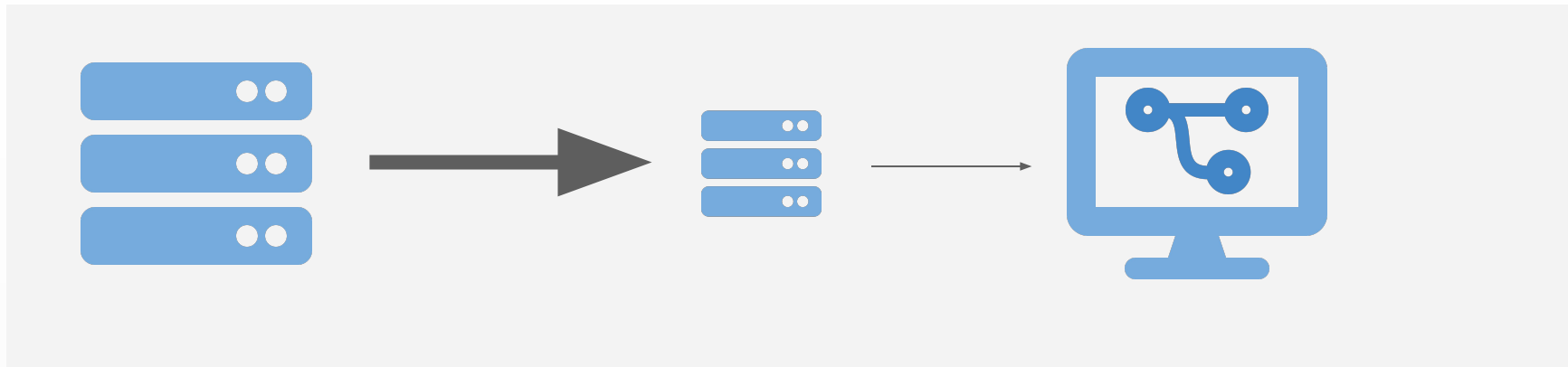To Scale

Database

Dev Machine

@alexkgold

R Studio®

# Big Data Strategies for R

@alexkgold

# Strategy 1: Sample and Model



😃 Use favorite R modeling package `(Caret/Parsnip/rsample)`.

😃 Really good for iterating/prototyping.

☹️ Requires care for sampling and scaling.

☹️ Not good for BI tasks.

@alexkgold

R Studio®

# Demo!



@alexkgold

R Studio

# Strategy 2: Chunk and Pull



😃 Great when discrete chunks exist.
😃 Facilitates parallelization.
☹️ Can't have interactions between chunks.
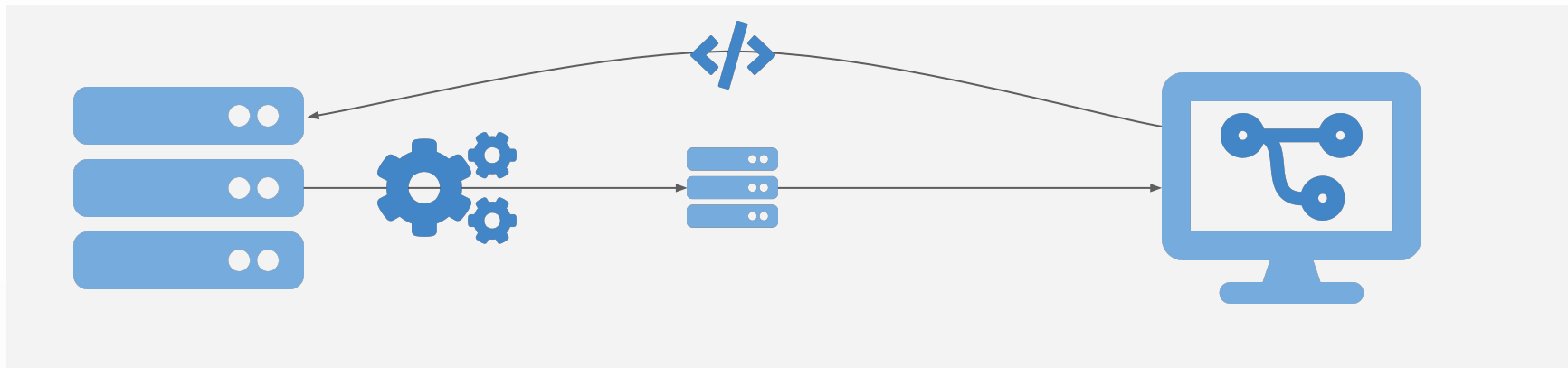☹️ Eventually pull in all data.

R Studio®

# Demo!

R Studio®

# Strategy 3: Push Compute to Data



😃 Take advantage of database strengths.
😃 Get whole dataset, but move less data.
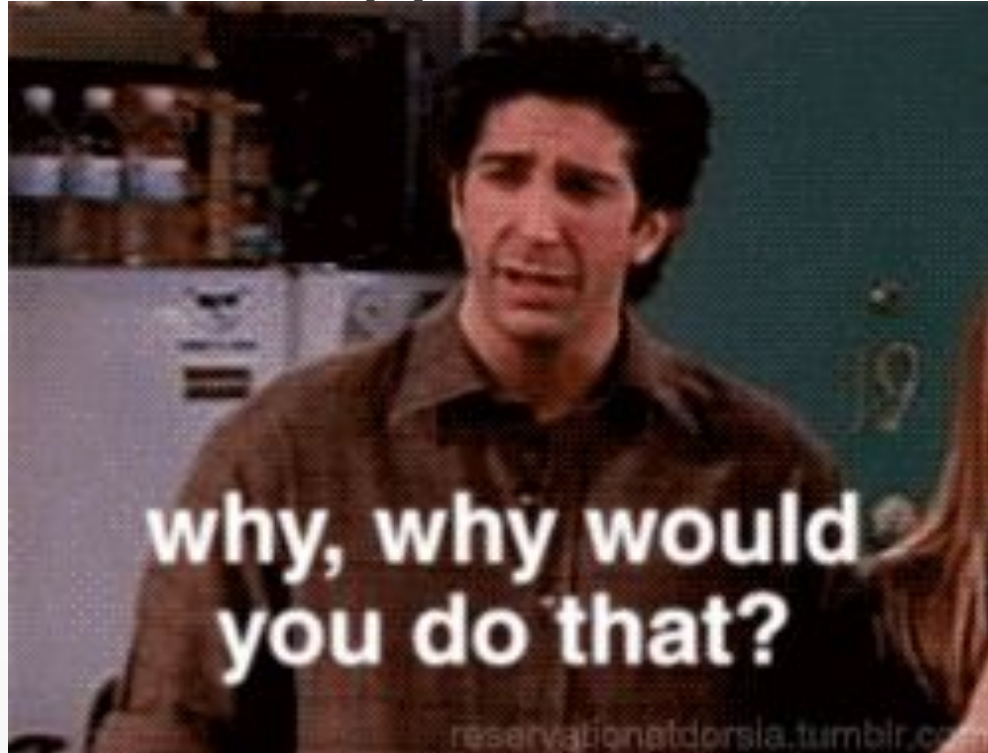☹ Operations might not be permitted in database.
☹ Maybe your database is slow?

@alexkgold

# Demo!

# 3 Big Data Strategies for R

Import �le Tidy

Communicate or Automate

@LateNightSeth

@alexkgold

# What about deployment?


WHAT ABOUT IT?

Open-Source (Free!)

- Build-your-own
- Shiny Server

Enterprise Products

- RStudio Connect
- Evals

R Studio®

# Recommendation Summary

| Problem | Solution |
|---|---|
| Single-Threading | <ul><li>Many R packages<ul><li>My favorite: `doFuture`</li></ul></li><li>RStudio Server Pro Job Launcher</li></ul> |
| R is Slow | <ul><li>Profile with `profvis`</li><li>Write in a faster language, call from R `(Rcpp)`</li></ul> |
| In-Memory Data | <ul><li>Adopt a big data paradigm for R<ol><li>Sample and Model</li><li>Chunk and Pull</li><li>Push Compute to Data</li></ol></li></ul> |

**db.rstudio.com**

**spark.rstudio.com**
**therinspark.com**

@alexkgold
https://github.com/akgold/bdl_2019

R Studio