# Tidy Finance and Accessing Financial Data

R Consortium Webinar

Christoph Scheuch

2024-03-06

# What is Tidy Finance?

A **transparent**, **open-source** approach to research in financial economics, featuring **multiple programming languages**

**tidy-finance.org** offers tools to:

- Learn about empirical applications using tidy principles
- Learn to work with financial data in a tidy manner
- Teach students the importance of reproducible research
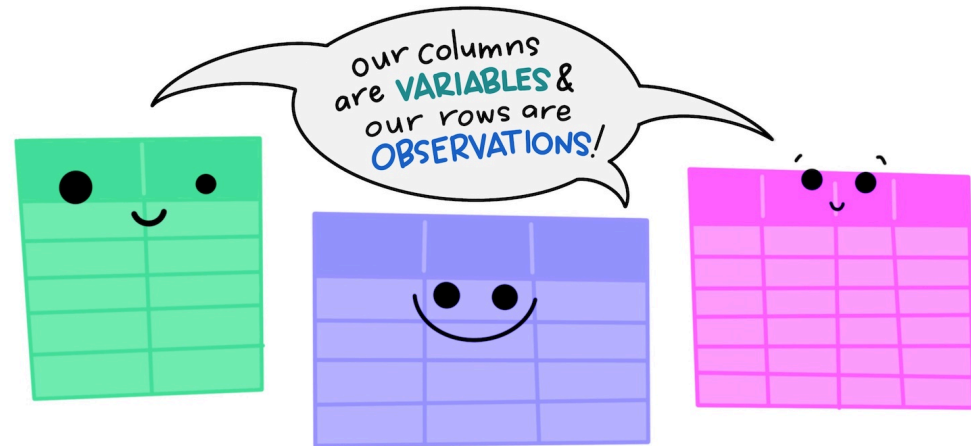- Contribute to reproducible finance research via our blog

# Why Tidy?

Code should not just be correct, but also follow principles:

1. Design so that code is **easy to read** for humans

2. **Compose simple functions** to solve complex problems

3. **Embrace functional programming** for reproducible results

4. **Reuse data structures** across applications
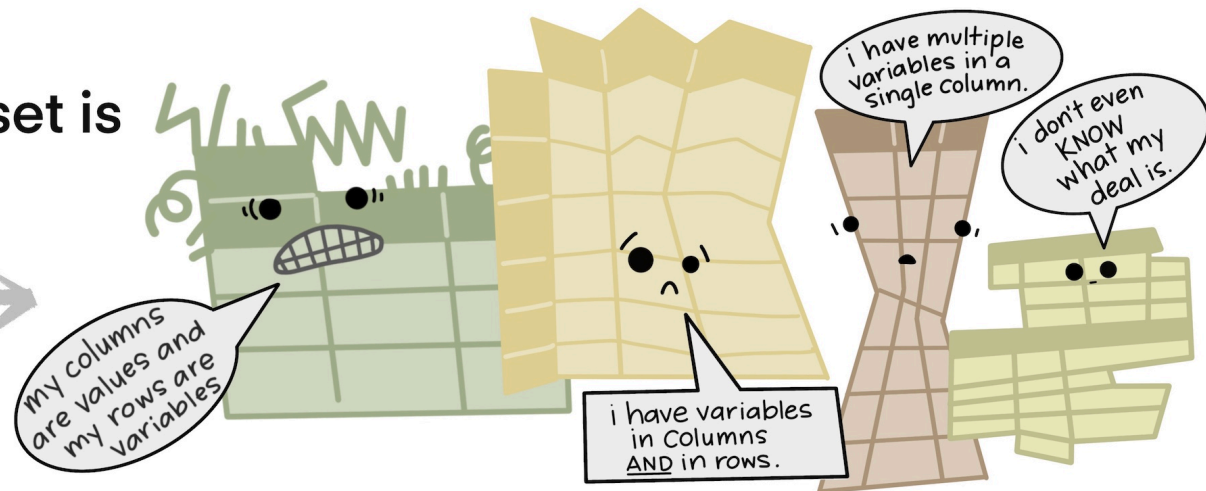
Focus of this talk: **tidy data**

# Recap: what is *tidy data?*

The standard structure of tidy data means that "tidy datasets are all alike..."

"...but every messy dataset is messy in its own way."

—HADLEY WICKHAM



Illustrations from the Openscapes blog Tidy Data for reproducibility, efficiency, and collaboration by Julia Lowndes and Allison Horst

# Example chunks with *tidy code*

R  Python

```r
 1  # Load packages
 2  library(tidyverse)
 3  library(tidyquant)
 4
 5  # Download symbols of DOW index
 6  symbols <- tq_index(x = "DOW") |>
 7    filter(company != "US DOLLAR")
 8
 9  # Download prices of DOW index constituents
10  prices <- tq_get(x = symbols, get = "stock.prices",
11                   from = "2000-01-01", to = "2022-12-31")
12
13  # Calculate returns
14  returns <- prices |>
15    group_by(symbol) |>
16    mutate(ret = adjusted / lag(adjusted) - 1) |>
17    select(symbol, date, ret) |>
18    drop_na(ret)
```

# Welcoming contributions on our blog

# Maintainers of tidy-finance.org

**Christoph Scheuch**

Head of Artificial
Intelligence at
wikifolio.com

**Stefan Voigt**

Assistant Professor of
Finance at University
of Copenhagen

**Patrick Weiss**

Assistant Professor of
Finance at Reykjavik
University

**Christoph Frey**

Quantitative
Researcher at
Pinechip Capital

# We also wrote books



Tidy Finance with R — The R Series — Christoph Scheuch, Stefan Voigt, Patrick Weiss — CRC Press, A Chapman & Hall Book



Tidy Finance with Python — Christoph Frey, Christoph Scheuch, Stefan Voigt and Patrick Weiss — Chapman & Hall, CRC Press

# Accessing & Managing Financial Data

# Importance of organizing data efficiently

- **Challenge:** ensure consistency across various data sources

- **Solution:**

  - Use R to import, prepare & store data

  - Use SQLite to organize data in a database

- **R Packages:**

  - Manipulation: `tidyverse`

  - Import: `tidyquant`, `frenchdata`, `readxl`

  - Storage: `RSQLite`

# Fama-French factors & portfolios

Most popular data for asset pricing tests since Fama and French (1993)

```r
1   library(frenchdata)
2
3   factors_ff3_monthly_raw <- download_french_data("Fama/French 3 Factors")
4   factors_ff3_monthly <- factors_ff3_monthly_raw$subsets$data[[1]] |>
5     mutate(
6       month = floor_date(ymd(str_c(date, "01")), "month"),
7       across(c(RF, `Mkt-RF`, SMB, HML), ~as.numeric(.) / 100),
8       .keep = "none"
9     ) |>
10    rename_with(str_to_lower) |>
11    rename(mkt_excess = `mkt-rf`) |>
12    select(month, everything())
13
14  print(factors_ff3_monthly, n = 5)
```

```
# A tibble: 1,170 × 5
  month      mkt_excess     smb      hml      rf
  <date>          <dbl>   <dbl>    <dbl>   <dbl>
1 1926-07-01     0.0296 -0.0256  -0.0243  0.0022
2 1926-08-01     0.0264 -0.0117   0.0382  0.0025
3 1926-09-01     0.0036 -0.014    0.0013  0.0023
4 1926-10-01    -0.0324 -0.0009   0.007   0.0032
5 1926-11-01     0.0253 -0.001   -0.0051  0.0031
# i 1,165 more rows
```

# q-Factors

Alternative to Fama-French data by Hou, Xue, and Zhang (2014)

```r
 1  library(readr)
 2
 3  factors_q_monthly_link <-
 4    "https://global-q.org/uploads/1/2/2/6/122679606/q5_factors_monthly_2022.csv"
 5
 6  factors_q_monthly <- read_csv(factors_q_monthly_link) |>
 7    mutate(month = ymd(str_c(year, month, "01", sep = "-"))) |>
 8    select(-R_F, -R_MKT, -year) |>
 9    rename_with(~ str_remove(., "R_")) |>
10    rename_with(~ str_to_lower(.)) |>
11    mutate(across(-month, ~ . / 100))
12
13  print(factors_q_monthly, n = 5)
```

```
# A tibble: 672 × 5
  month           me       ia      roe       eg
  <date>        <dbl>    <dbl>    <dbl>    <dbl>
1 1967-01-01   0.0683  -0.0297   0.0192  -0.0218
2 1967-02-01   0.0165  -0.00227  0.0354   0.0222
3 1967-03-01   0.0200  -0.0178   0.0184  -0.0104
4 1967-04-01  -0.00690 -0.0288   0.0106  -0.0173
5 1967-05-01   0.0285   0.0252   0.00692  0.00158
# i 667 more rows
```

# Macroeconomic predictors

Collection of variables for equity premium prediction (Welch & Goyal, 2008)

```r
 1  library(readxl)
 2
 3  download.file(
 4    url = "https://docs.google.com/spreadsheets/d/1g4LOaRj4TvwJr9RIaA_nwrXXWTOy46bP/expo
 5    destfile = "macro_predictors.xlsx",
 6    mode = "wb"
 7  )
 8
 9  macro_predictors <- read_xlsx("macro_predictors.xlsx", sheet = "Monthly") |>
10    mutate(
11      # Several cleaning steps & variable transformations...
12    )
```

```
# A tibble: 1,152 × 15
  month          rp_div      dp     dy     ep     de      svar     bm    ntis     tbl
  <date>          <dbl>   <dbl>  <dbl>  <dbl>  <dbl>     <dbl>  <dbl>   <dbl>   <dbl>
1 1926-12-01   -0.0220    -2.97  -2.96  -2.39  -0.586 0.000465  0.441  0.0509  0.0307
2 1927-01-01    0.0422    -2.94  -2.96  -2.37  -0.568 0.000470  0.444  0.0508  0.0323
3 1927-02-01    0.00363   -2.98  -2.93  -2.43  -0.549 0.000287  0.429  0.0517  0.0329
4 1927-03-01    0.0142    -2.98  -2.97  -2.45  -0.531 0.000924  0.470  0.0464  0.032
5 1927-04-01    0.0459    -2.98  -2.97  -2.47  -0.513 0.000603  0.457  0.0505  0.0339
# i 1,147 more rows
# i 5 more variables: lty <dbl>, ltr <dbl>, tms <dbl>, dfy <dbl>, infl <dbl>
```

# Other macroeconomic data

10K data sets available via Federal Reserve Economic Data (FRED) database

```r
 1  library(tidyquant)
 2
 3  # Example: consumer price index (CPI)
 4  cpi_monthly <- tq_get("CPIAUCNS", get = "economic.data") |>
 5    mutate(
 6      month = floor_date(date, "month"),
 7      cpi = price / price[month == max(month)],
 8      .keep = "none"
 9    )
10  print(cpi_monthly, n = 5)
```

```
# A tibble: 121 × 2
  month         cpi
  <date>      <dbl>
1 2014-01-01  0.758
2 2014-02-01  0.761
3 2014-03-01  0.766
4 2014-04-01  0.769
5 2014-05-01  0.771
# i 116 more rows
```

# Use SQLite database for storage

```r
1  library(RSQLite)
2  library(dbplyr)
3
4  # Create database
5  tidy_finance <- dbConnect(
6    SQLite(), "tidy_finance_r.sqlite", extended_types = TRUE
7  )
8
9  # Write data to database
10 dbWriteTable(
11   conn = tidy_finance,
12   name = "factors_ff3_monthly",
13   value = factors_ff3_monthly,
14   overwrite = TRUE
15 )
16
17 # Load data from database
18 factors_ff3_monthly <- tbl(tidy_finance, "factors_ff3_monthly") |>
19   collect()
```

# Why SQLite?

**Pros:**

- Lightweight, self-contained, serverless database engine
- Great for education purposes or prototyping

**Cons:**

- Limitations with respect to very large data & concurrency
- Transfer to other languages cumbersome (e.g. Python)
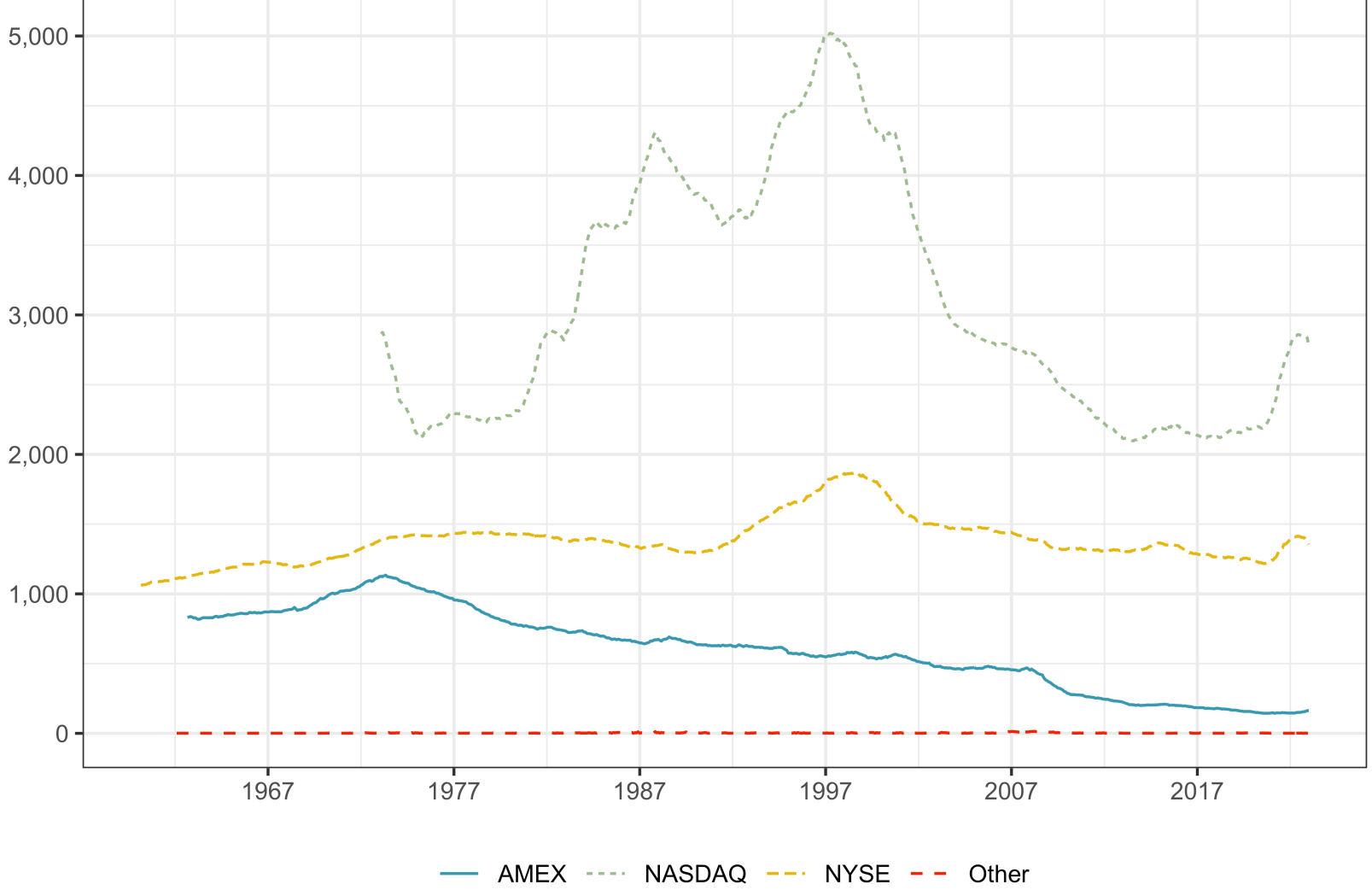
# WRDS & Other Data Providers

# Wharton Research Data Services (WRDS)

- Popular provider of financial & economic data

- Focus on academic audience & research applications

- Access via `RPostgres` package

- Main data used in Tidy Finance

  - **CRSP:** historical monthly & daily returns for US stocks

  - **Compustat:** historical accounting data for US companies

  - **Mergent FISD:** characteristics of US corporate bonds

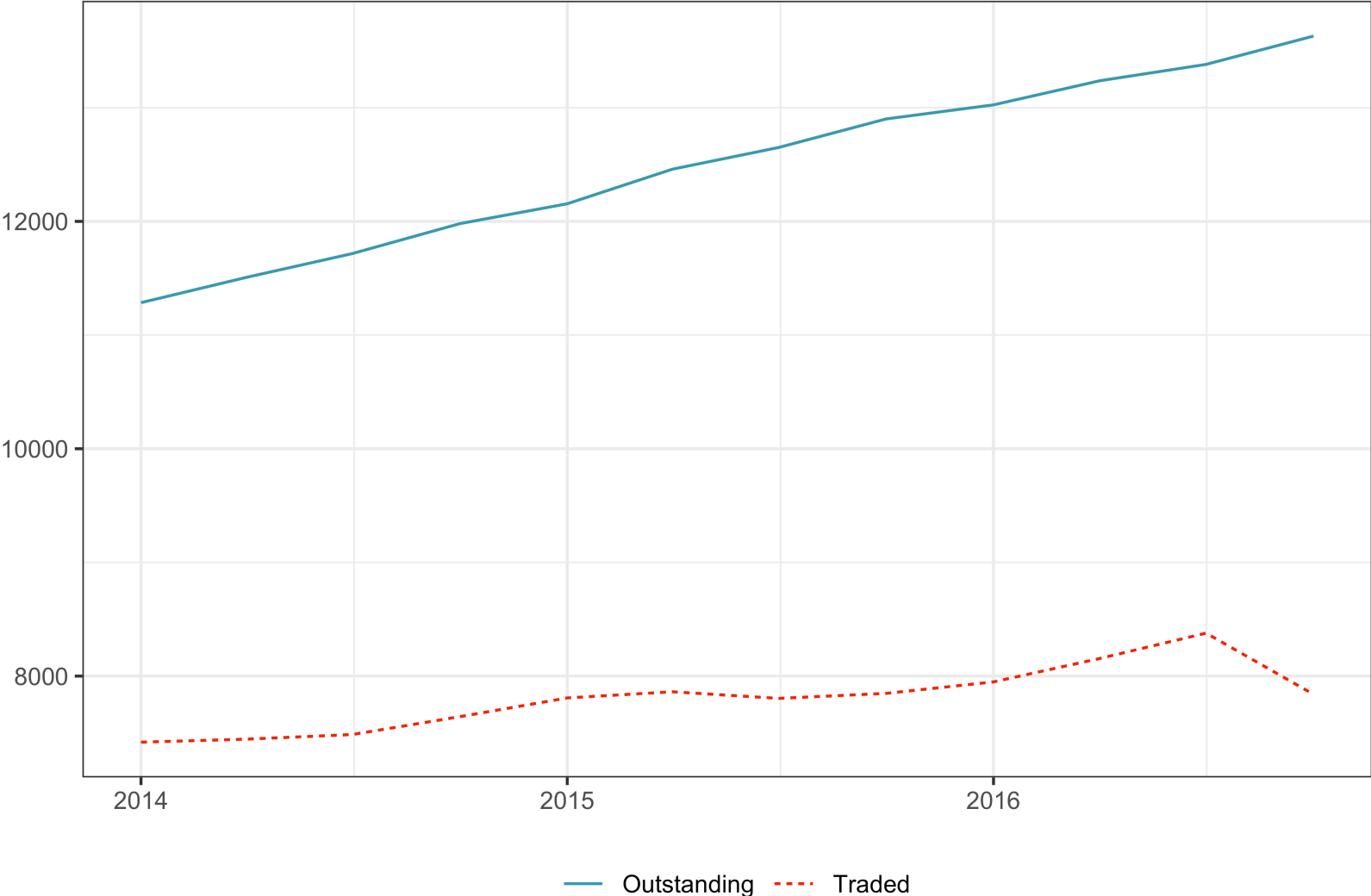  - **TRACE:** detailed US corporate bond transactions

# Glimpse at historical stock data

Monthly number of securities by listing exchange

# Glimpse at historical bond data

Number of bonds outstanding and traded each quarter

# Other data providers

Large ecosystem of alternative data providers

- Extensive list of R packages on tidy-finance.org

- Examples: `fredr`, `ecb`, `Rblpapi`, `Quandl`, `edgarWebR`, etc.

Are we missing an important package?

- please reach out via **contact@tidy-finance.org**

# Wrap-up

# Tidy approach to teaching & research

`tidyfinance` R package to access financial data in a tidy way:

```
1  install.packages("tidyfinance")
2
3  tidyfinance::download_data(
4    type = "wrds_crsp_monthly",
5    start_date = "1960-01-01", end_date = "2020-12-31"
6  )
```

- Detailed open source material at **tidy-finance.org**

- Get in touch for **teaching materials** & to **contribute to blog**

- Follow me for news linkedin.com/in/christophscheuch