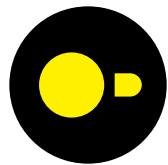


# DuckDB dplyr



Kirill Müller, posit::conf(2025), cynkra.■



# DuckDB

- Embedded analytical database
- Command line interface
- Python/Rust/WASM/... package
- R package

# Large



# Fast

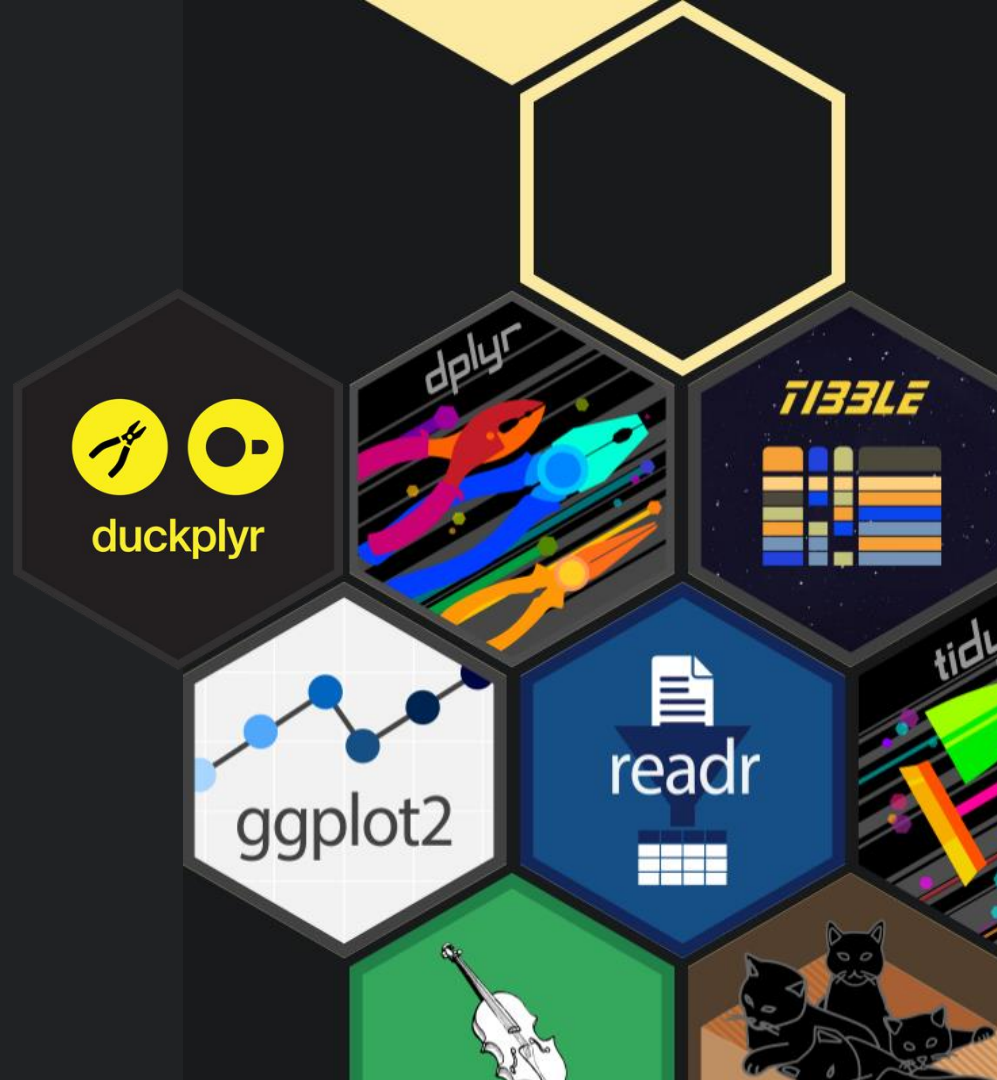


# Featureful









# DuckDB interfaces

dplyr-like



SQL

DBI



DuckDB

```
SELECT COUNT(*)  
FROM 'personas.parquet';
```

```
--  
--  
--  
--  
--  
--  
--
```

count_star() int64
1000000



```
library(DBI)
con <- dbConnect(duckdb::duckdb())

tbl <- duckdb::tbl_file(
  con,
  "personas.parquet")
tbl |> dplyr::count()
#> # Source:   SQL [?? x 1]
#> # Database: DuckDB 1.4.0
#>       n
#>   <dbl>
#> 1 100000

dbGetQuery(
  con,
  "SELECT COUNT(*)
   FROM 'personas.parquet'")
#>   count_star()
#> 1          1e+05
```

```
SELECT COUNT(*)
FROM 'personas.parquet';
```

--	
--	count_star()
--	int64
--	
--	100000
--	



```
library(DBI)
con <- dbConnect(duckdb::duckdb())

tbl <- duckdb::tbl_file(
  con,
  "personas.parquet")
tbl |> dplyr::count()
#> # Source:   SQL [?? x 1]
#> # Database: DuckDB 1.4.0
#>       n
#>   <dbl>
#> 1 100000

dbGetQuery(
  con,
  "SELECT COUNT(*)
  FROM 'personas.parquet'")
#> count_star()
#> 1          1e+05
```

```
tbl <- duckplyr::read_parquet_duckdb(
  "personas.parquet")
tbl |> dplyr::count()
#> # A duckplyr data frame: 1 variable
#>       n
#>   <int>
#> 1 100000
```

```
SELECT COUNT(*)
FROM 'personas.parquet';
```

count_star()
int64
100000

```
library(duckplyr)

personas <- read_parquet_duckdb(
  "hf://datasets/.../*/*.parquet"
)

personas_count <-
  personas |>
  count(sex, education_level)

personas_count |>
  explain()
```

ORDER\_BY  
sex A, education\_level A  
~14 rows

HASH\_GROUP\_BY  
Groups:  
#0, #1  
Aggregates:  
count\_star()  
~100,000 rows

READ\_PARQUET  
Projections:  
sex, education\_level  
~100,000 rows



# Large data

- Peek with remote access
- Prudence avoids crashes
- Iterate with local file
- **Production: remote**



```
library(tidyverse)
```

```
G <- 200
```

```
N <- 5000
```

```
data <-
```

```
  crossing(g1 = 1:G, g2 = 1:G, id = 1:N) |>  
  mutate(x = rnorm(n()))
```

```
data |>
```

```
  summarize(.by = c(g1, g2), m = mean(x)) |>
```

```
  system.time()
```

```
#>   user  system elapsed
```

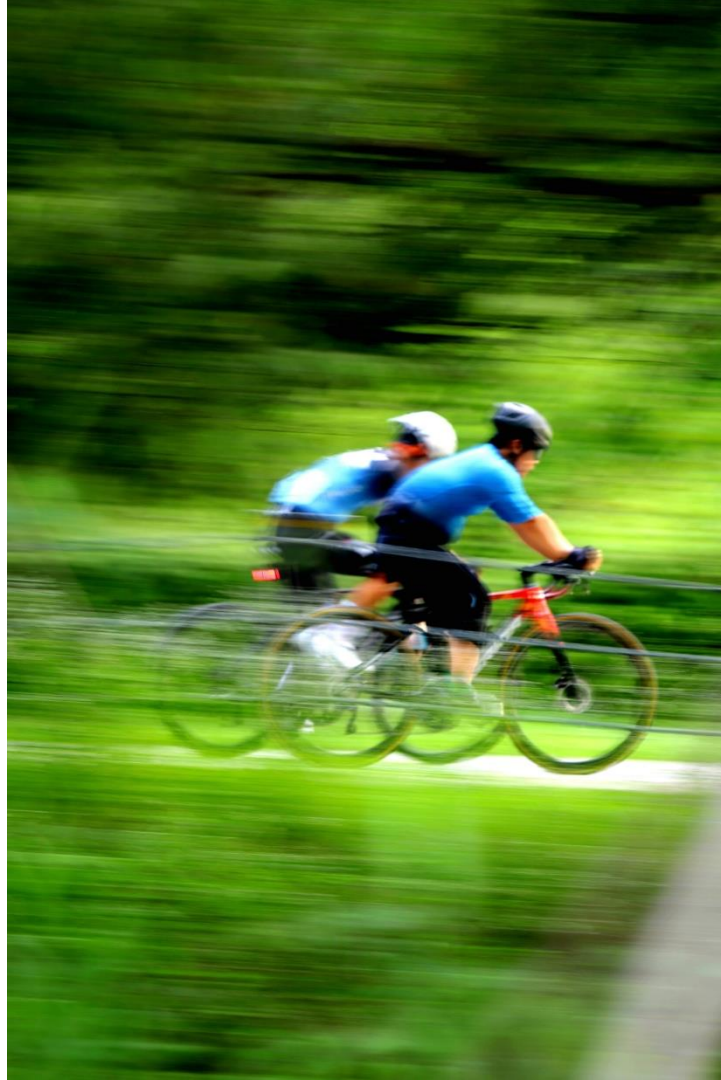
```
#> 2.397   1.720   5.172
```





# Existing pipelines

- All in
- Best effort
- Targeted



```
library(tidyverse)
Sys.setenv(DUCKPLYR_FORCE = TRUE)
```

```
G <- 200
N <- 5000
```

```
data <-
  crossing(g1 = 1:G, g2 = 1:G, id = 1:N) |>
  mutate(x = rnorm(n()))
#> Error in `mutate()`:
#> ! Can't translate function `rnorm()`.
```



```
library(tidyverse)
library(duckplyr)
#> ✓ Overwriting dplyr methods with duckplyr methods.
#> i Turn off with `duckplyr::methods_restore()`.
```

```
G <- 200
N <- 5000
```

```
data <-
  crossing(g1 = 1:G, g2 = 1:G, id = 1:N) |>
  mutate(x = rnorm(n()))
```

```
data |>
  summarize(.by = c(g1, g2), m = mean(x)) |>
```

```
collect() |>
system.time()
#>   user  system elapsed
#> 4.079   0.016   0.600
```



```
library(tidyverse)
library(duckplyr)
#> ✓ Overwriting dplyr methods with duckplyr methods.
#> i Turn off with `duckplyr::methods_restore()`.
methods_restore()
#> i Restoring dplyr methods.
```

```
data <-
  crossing(g1 = 1:G, g2 = 1:G, id = 1:N) |>
  mutate(x = rnorm(n()))
```

```
data |>
  as_duckdb_tibble(prudence = "stingy") |>
  summarize(.by = c(g1, g2), m = mean(x)) |>
  collect() |>
  system.time()
#>      user  system elapsed
#>  4.142    0.017    0.659
```



```
library(tidyverse)
```

```
idx <- duckplyr::duckdb_tibble(id = 1:1e4)
```

```
idx |>
```

```
  mutate(x = dd::random()) |>
```

```
  ggplot(aes(x)) +  
  geom_histogram()
```



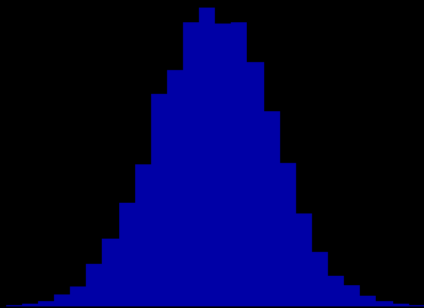


```
library(tidyverse)

idx <- duckplyr::duckdb_tibble(id = 1:1e4)

duckplyr::db_exec("LOAD stochastic")

idx |>
  mutate(x = dd::dist_normal_sample(0, 1)) |>
  ggplot(aes(x)) +
  geom_histogram()
```





Search



R 4.4.3

posit-conf-2025

42-dd.R U ●

43-nested.R U ●



scripts &gt; 42-dd.R

1 ?dd::random()

2

3 ?dd::struct\_pack()

4 dd::unnest(col0 = ANY)

5 ?dd::list()

6 ANY

7 dd::unnest()

8

SESSION

CONNECTIONS

HELP

VIEWER



R: DuckDB function list



list {dd}

R Documentation

## DuckDB function list

### Description

Returns a LIST containing all the values of a column.

### Usage

```
list(arg = T)
```

### Arguments

arg T

### Value

CONSOLE

TERMINAL

PROBLEMS

OUTPUT

PORTS

...



~/git/R/cynkra/posit-conf-2025



```
install.packages("duckplyr")  
library(duckplyr)  
methods_restore()
```

cynkra■

[tinyurl.com/  
duckplyr-25](https://tinyurl.com/duckplyr-25)

