

# VRROOM



*Life's Too Short To Drive Slow*

*Read*



*Jim Hester*



 *@jimhester*  *@jimhester\_*

 [SPEAKERDECK.COM/JIMHESTER/VROOM](https://speakerdeck.com/jimhester/vroom)

*Photo by Joe Neric on Unsplash*

 [SPEAKERDECK.COM/JIMHESTER/VROOM](https://SPEAKERDECK.COM/JIMHESTER/VROOM)



*Photo by Joshua Reddekopp on Unsplash*



**.1** *second*  
**1** *second*  
**10** *seconds*



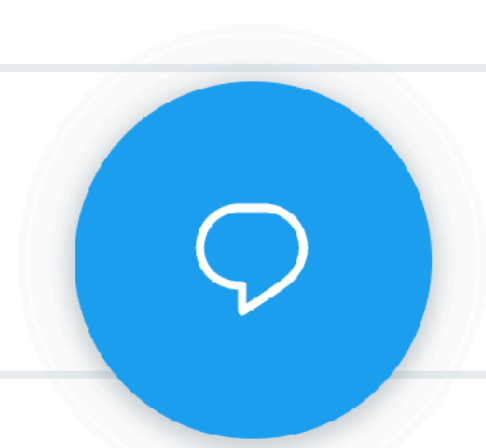
**Amjad Masad** 

@amasad

Interactive programming environments should give you feedback faster than your impulse to open Twitter or HN.

3:34 PM · Jun 18, 2019 · [Twitter Web Client](#)

**35** Retweets   **183** Likes





# *Taxi Trip Fare*

*Observations: 14,776,615*

*Variables: 11*

*File Size: 1.55G*

*chr [4]: medallion, hack\_license, v*

*dbl [6]: fare\_amount, surcharge,*

*dtm [1]: pickup\_datetime*



# *Taxi Trip Fare*

*Observations: 14,776,615*

*Variables: 11*

*File Size: 1.55G*

*chr [4]: medallion, hack\_license, v*

*dbl [6]: fare\_amount, surcharge,*

*dtm [1]: pickup\_datetime*



# *Taxi Trip Fare*

*Observations: 14,776,615*

*Variables: 11*

*File Size: 1.55G*

*chr [4]: medallion, hack\_license, v*

*dbl [6]: fare\_amount, surcharge,*

*dtm [1]: pickup\_datetime*



# *Taxi Trip Fare*

*Observations: 14,776,615*

*Variables: 11*

*File Size: 1.55G*

***chr** [4]: medallion, hack\_license, v*

*dbl [6]: fare\_amount, surcharge,*

*dtm [1]: pickup\_datetime*





# *Taxi Trip Fare*

*Observations: 14,776,615*

*Variables: 11*

*File Size: 1.55G*

*chr [4]: medallion, hack\_license, v*

***dbl [6]: fare\_amount, surcharge,***

*dtm [1]: pickup\_datetime*



# *Taxi Trip Fare*

*Observations: 14,776,615*

*Variables: 11*

*File Size: 1.55G*

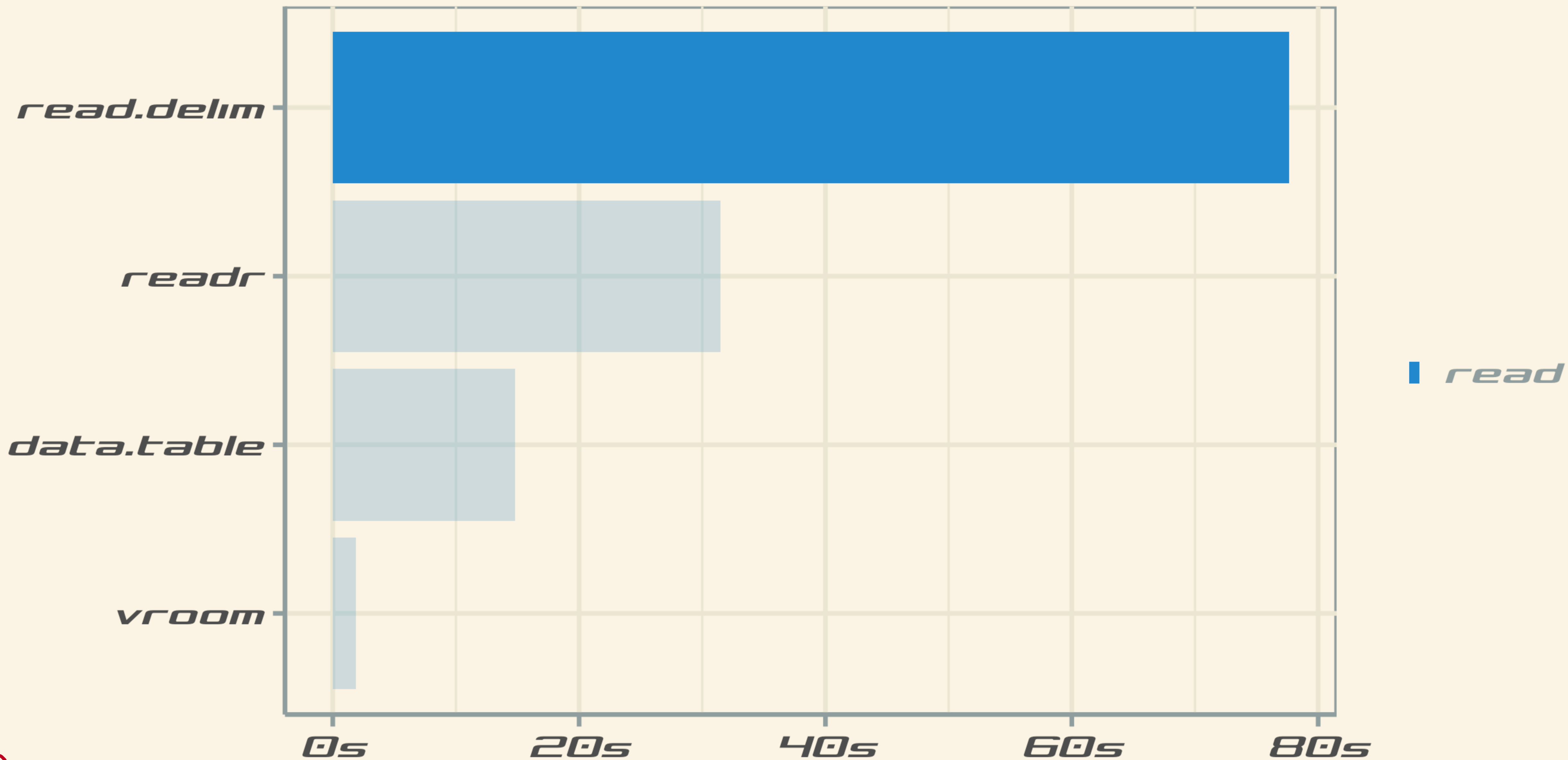
*chr [4]: medallion, hack\_license, v*

*dbl [6]: fare\_amount, surcharge,*

***dtm [1]: pickup\_datetime***

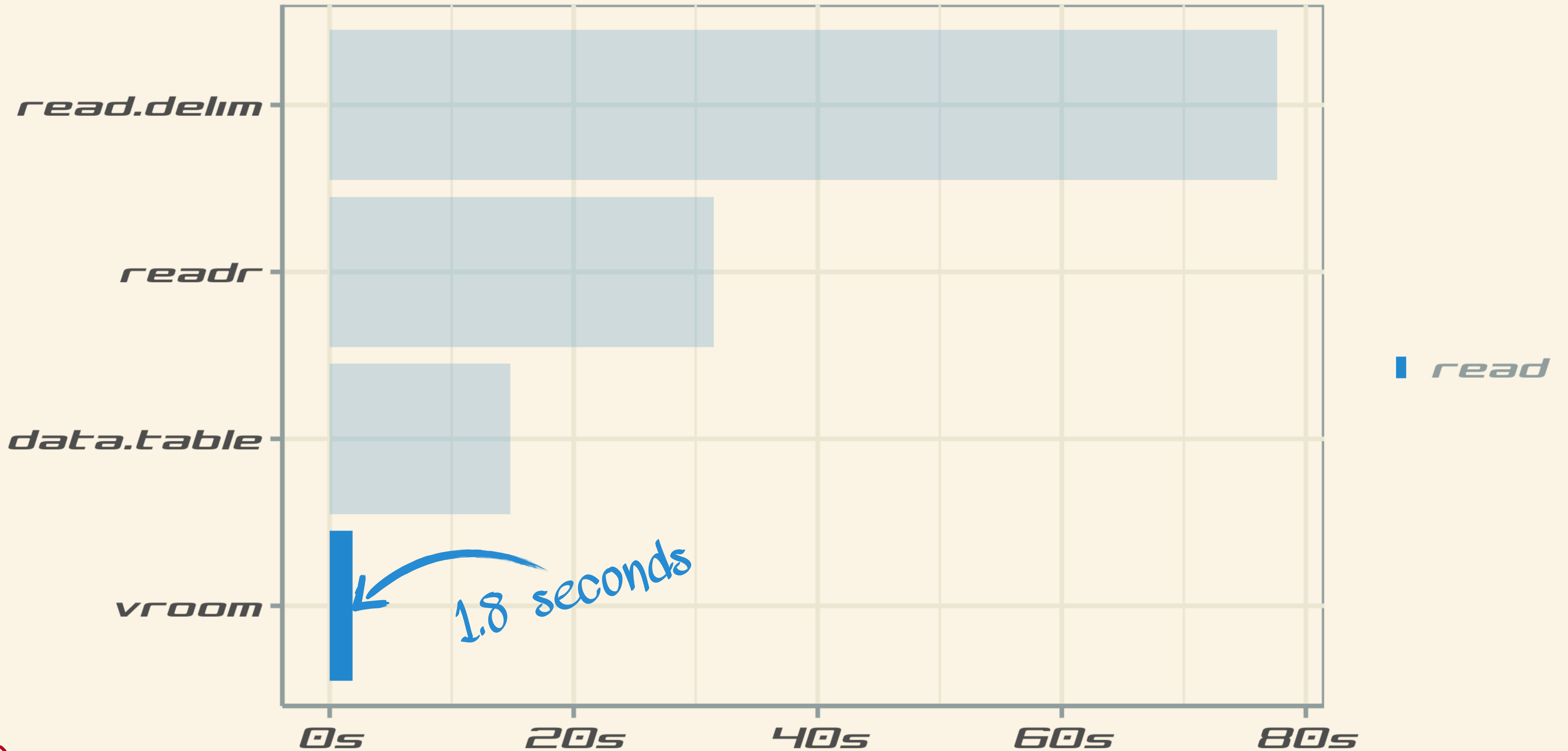
# *Taxi trip fare*

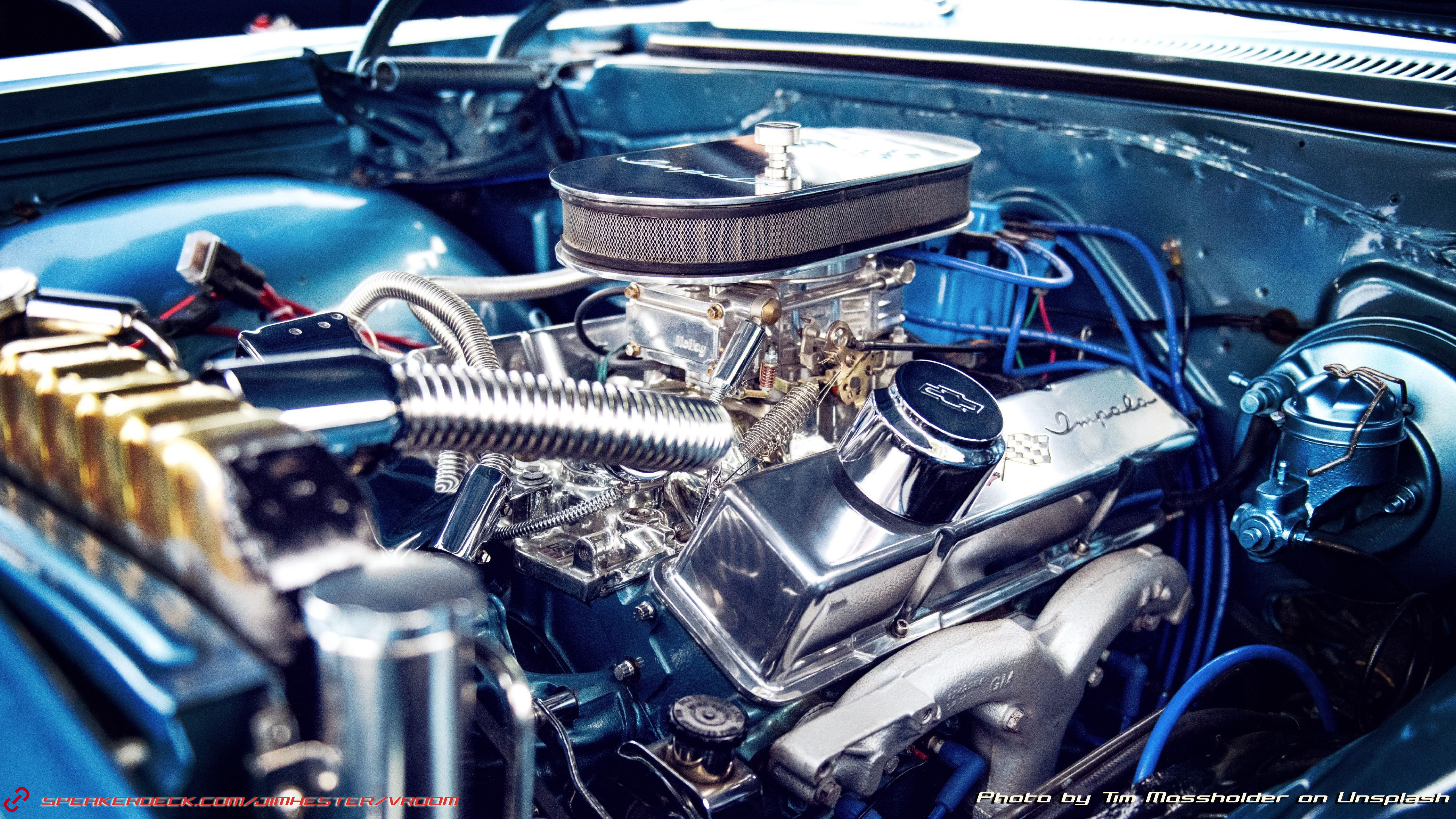
*14,776,615 x 11 - 1.55GB*



# Taxi trip fare

14,776,615 x 11 - 1.55GB





# *Memory mapped*

*Importance*



*Memory mapped*

*Multi-threaded*

Importance



*Memory mapped*

*Multi-threaded*

*strcpy()*

Importance





*Memory mapped*

*Multi-threaded*

*strcpy()*

*Altrep*

Importance



*ALTRIP*



*ALTRIP*

*Alternative representation*



# ALTRIP

*Alternative representation*

*R 3.5+*



# ALTRREP

*Alternative representation*

*R 3.5+*

*Custom memory storage*



# *ALTRREP*

*Alternative representation*

*R 3.5+*

*Custom memory storage*

*Transparent to R & C/C++*



# ALTRREP

*Alternative representation*

*R 3.5+*

*Custom memory storage*

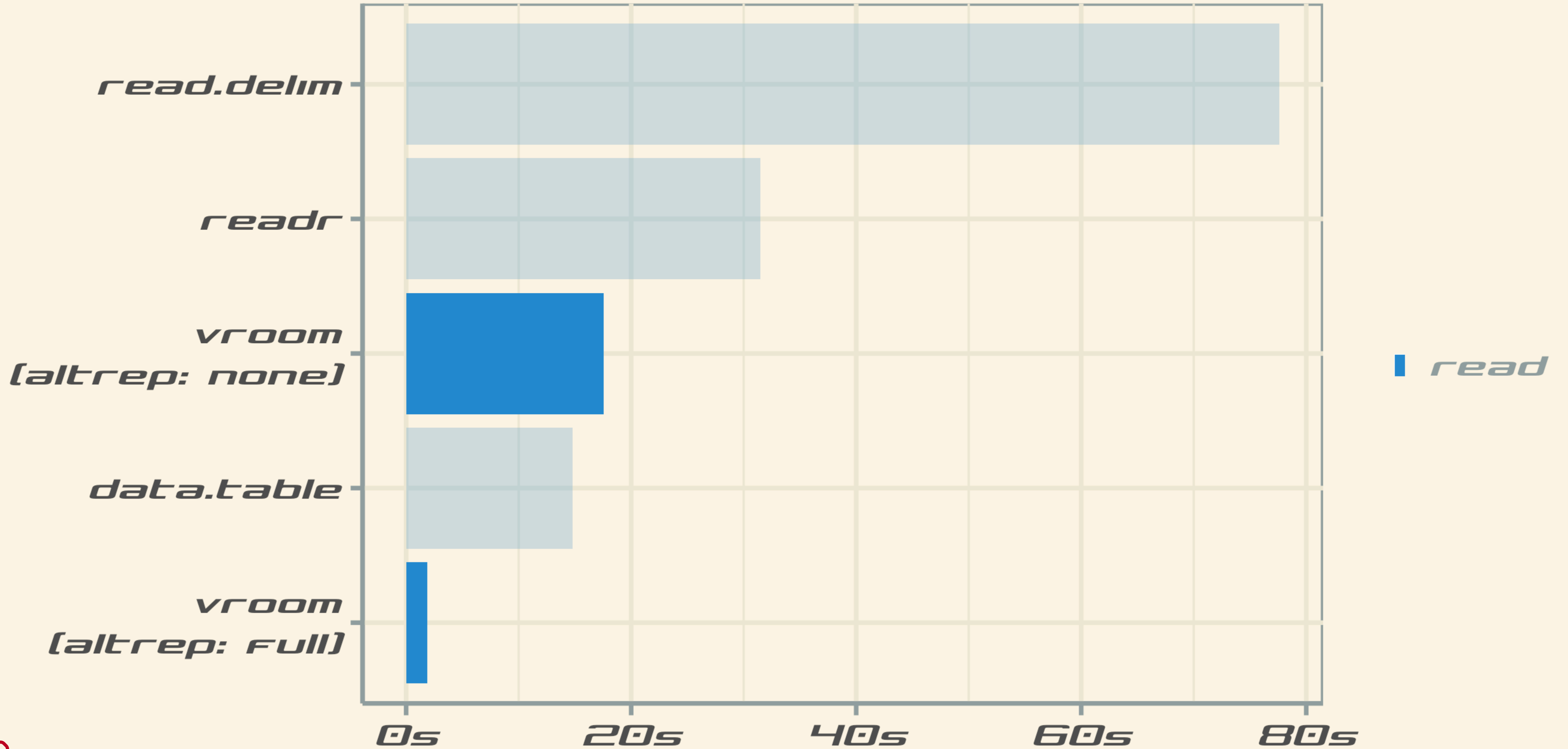
*Transparent to R & C/C++*

*On-demand parsing*

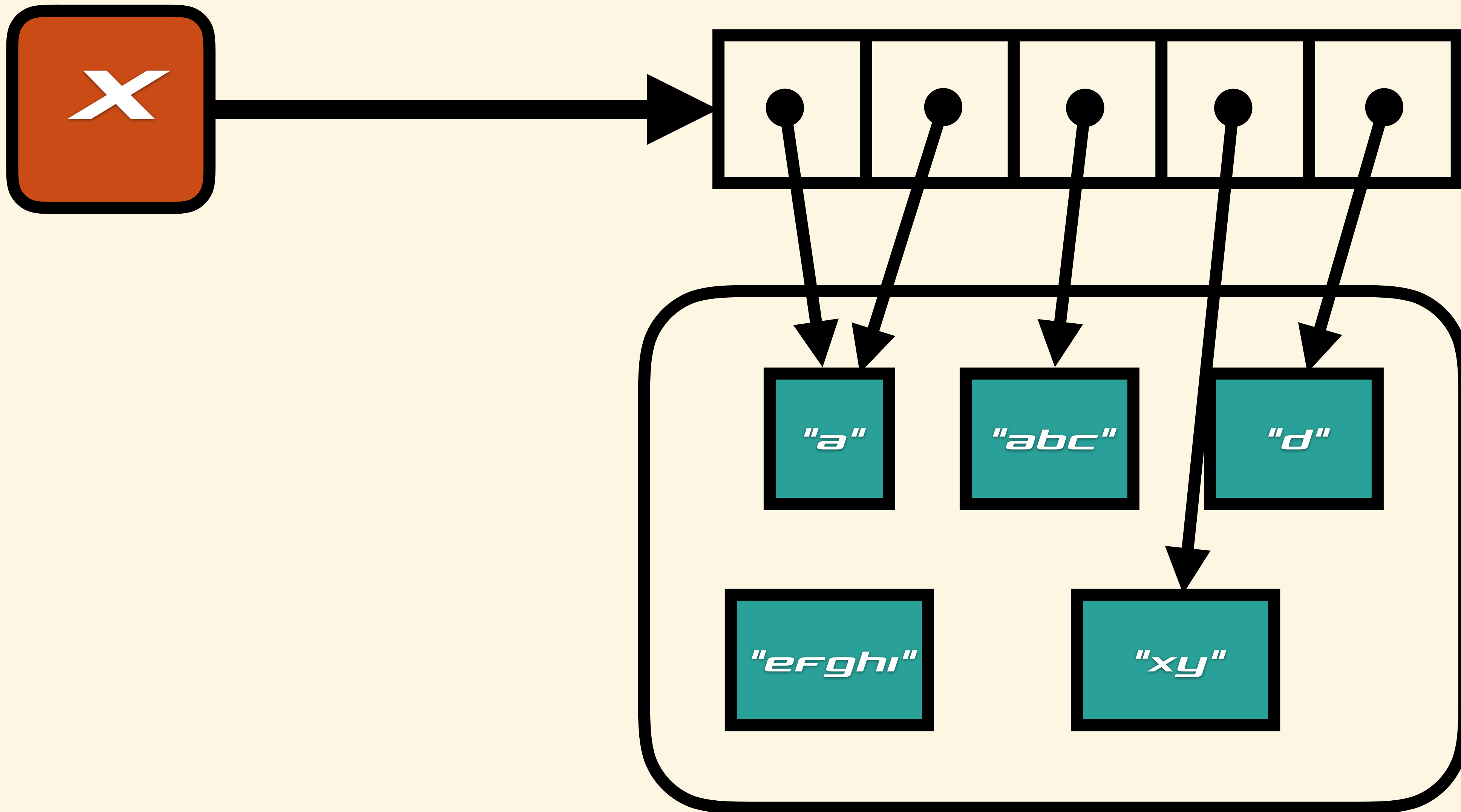


# Taxi trip fare

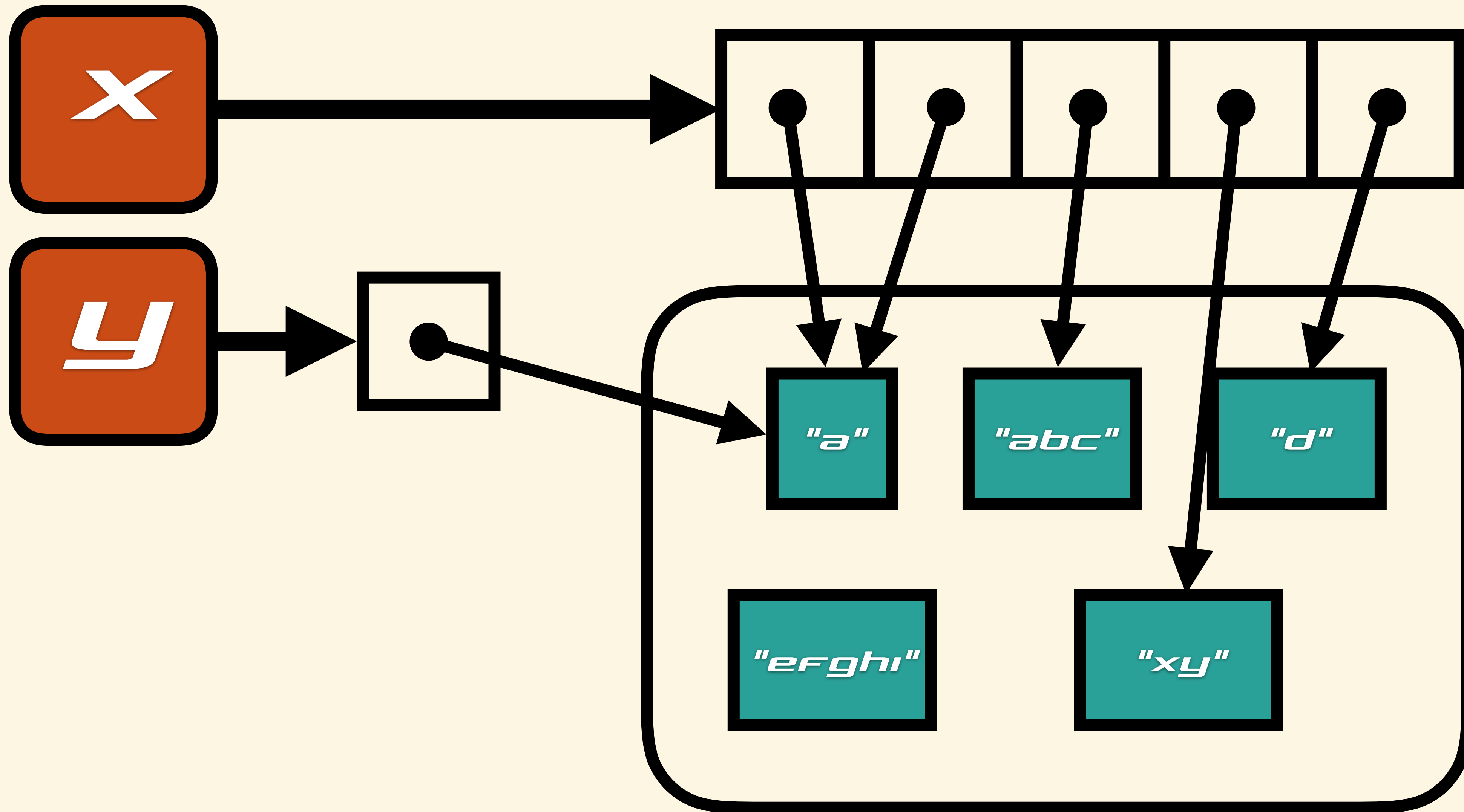
14,776,615 x 11 - 1.55GB



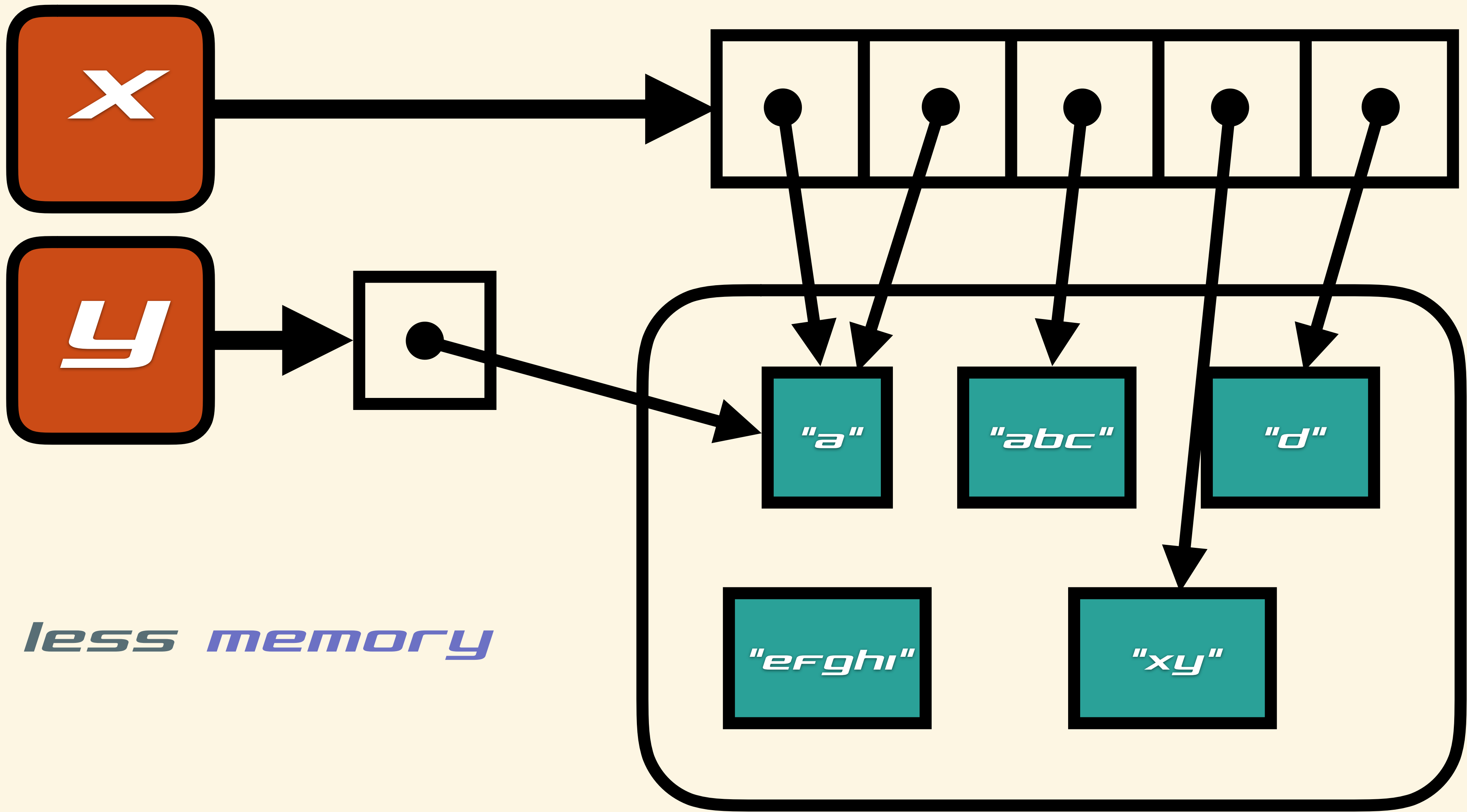




*Global string pool*

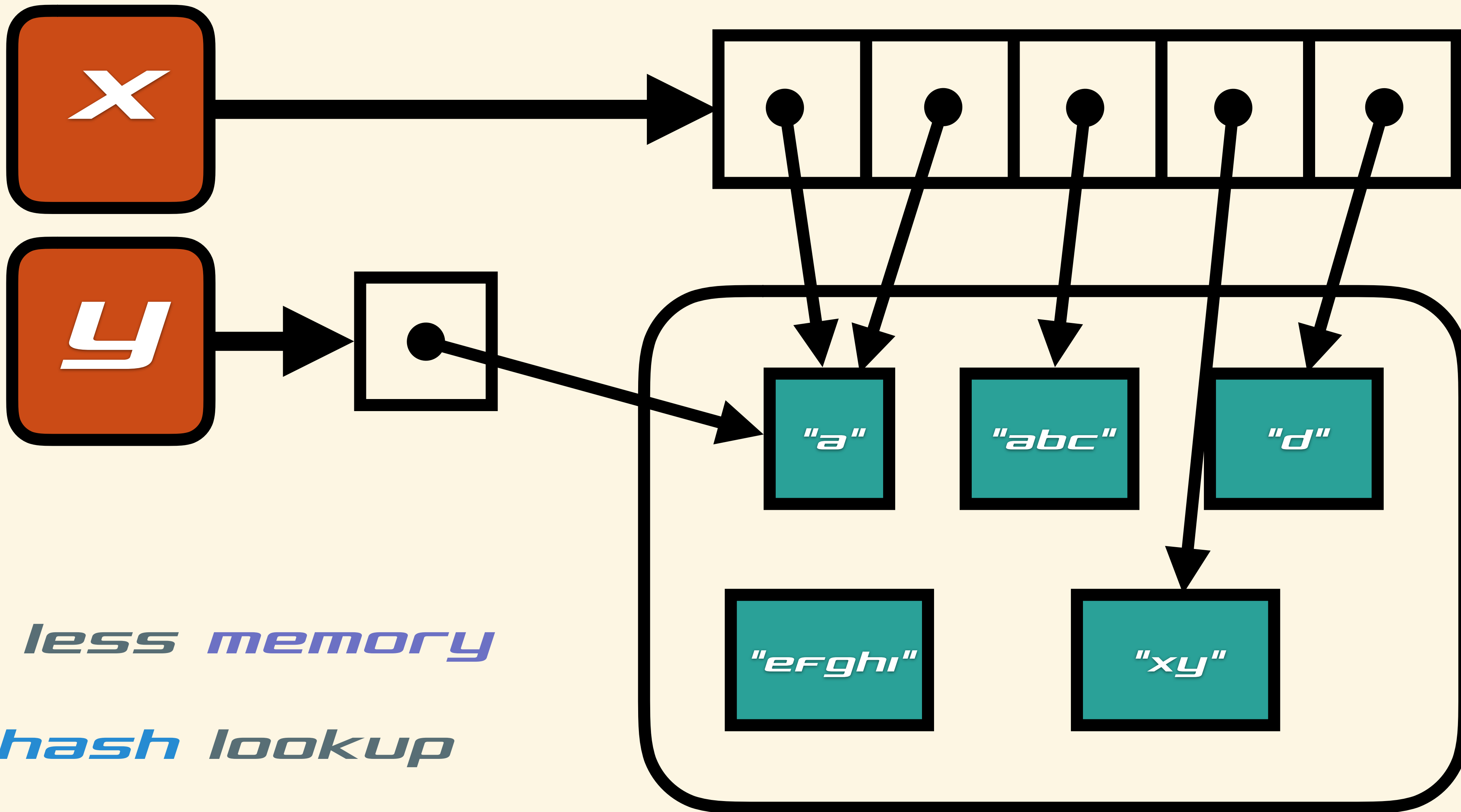


*Global string pool*



*+++ less memory*

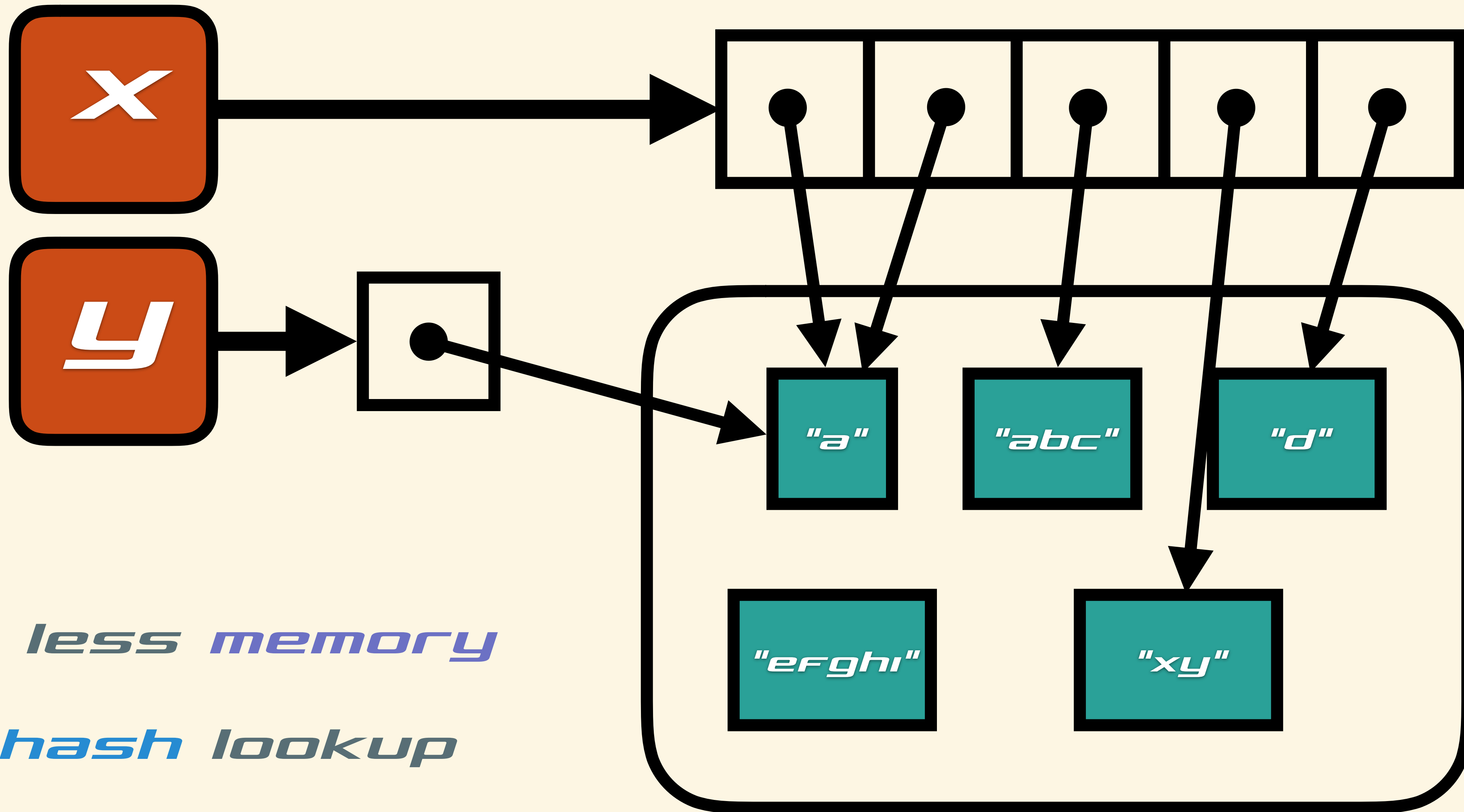
*Global string pool*



*+++ less memory*

*--- hash lookup*

## *Global string pool*



*+++ less memory*

*--- hash lookup*

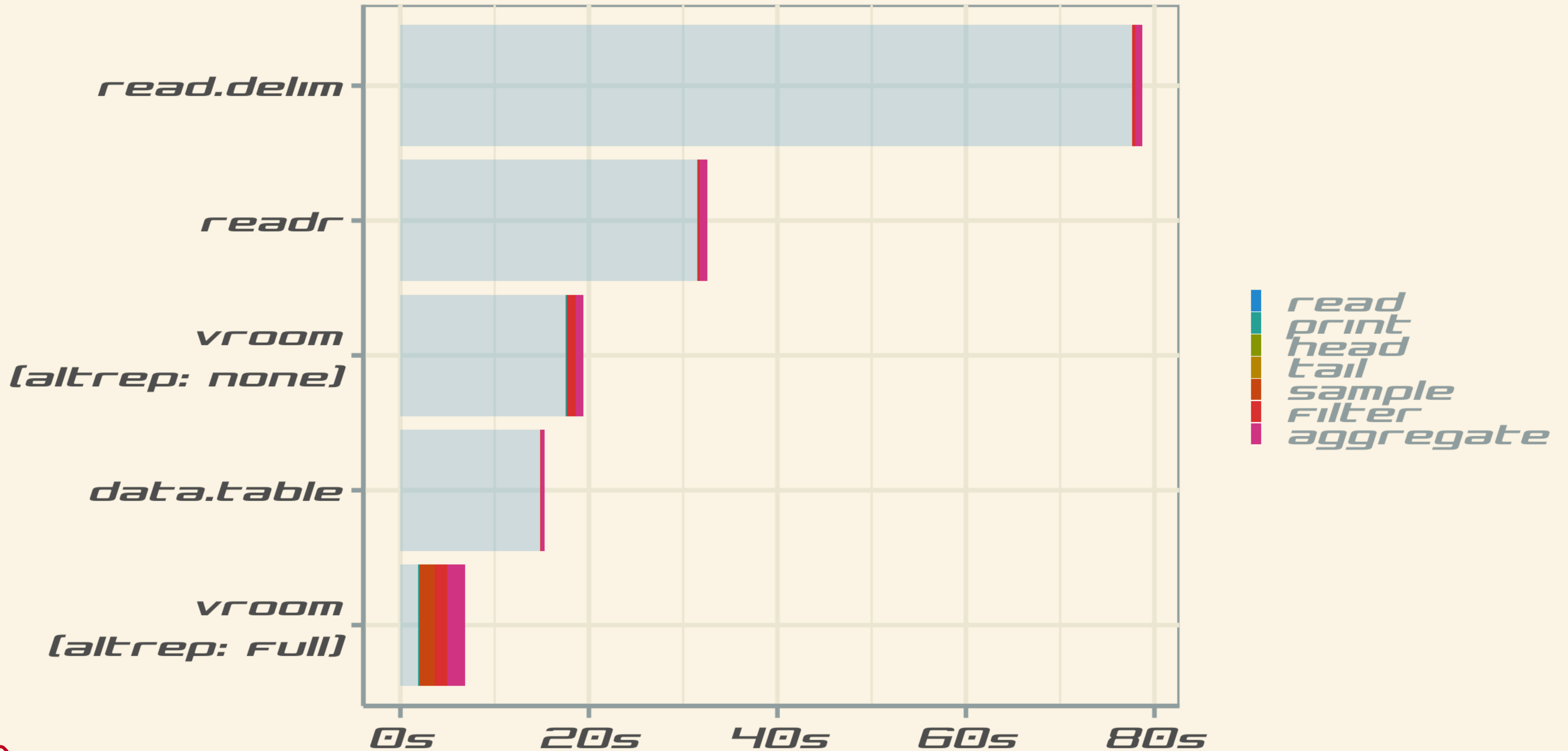
*--- single threaded*

*Global string pool*

*COST OF  
LAZINESS*

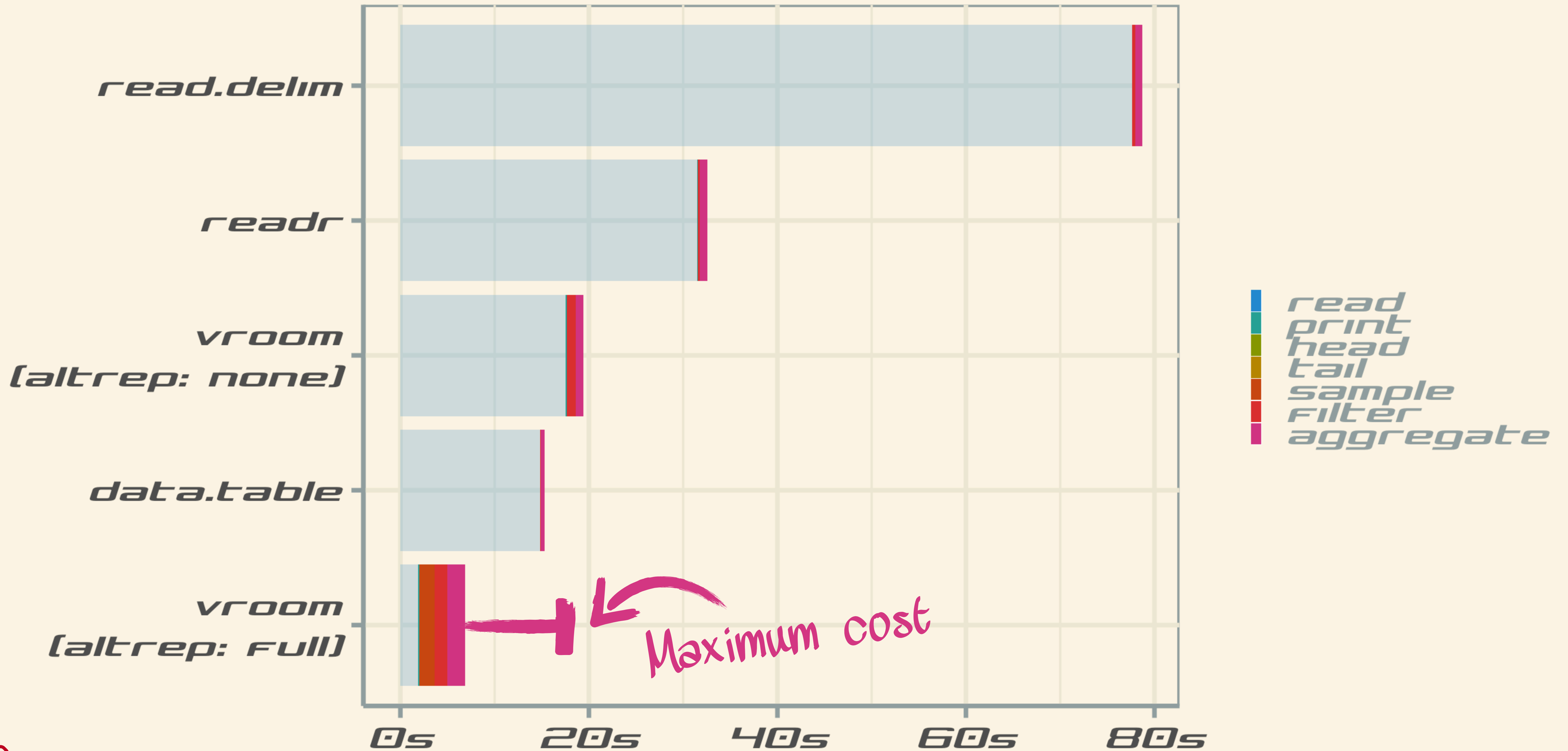
# Taxi trip fare

14,776,615 x 11 - 1.55GB



# Taxi trip fare

14,776,615 x 11 - 1.55GB

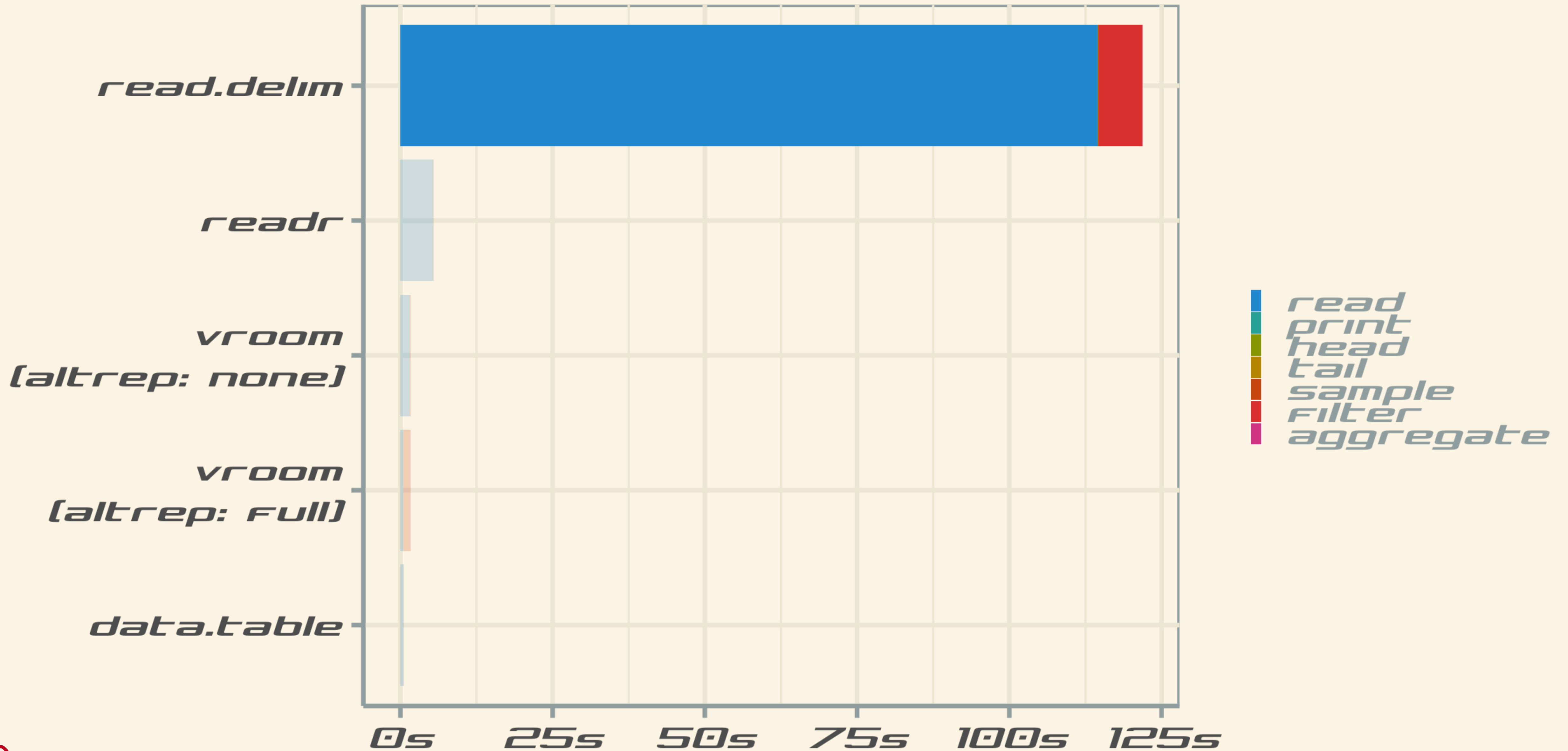




*ALL DOUBLES*

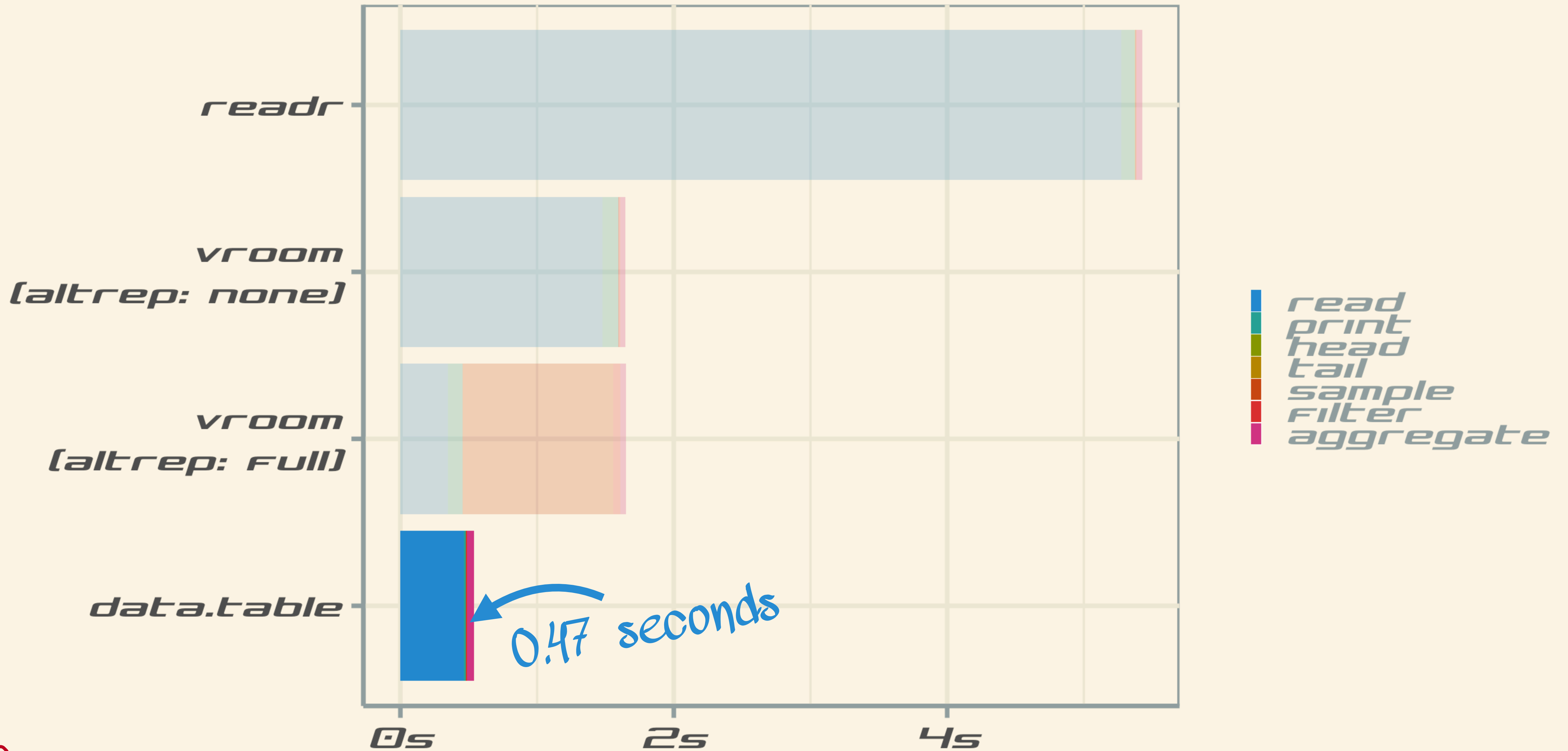
# All double

1,000,000 x 25 - 468MB



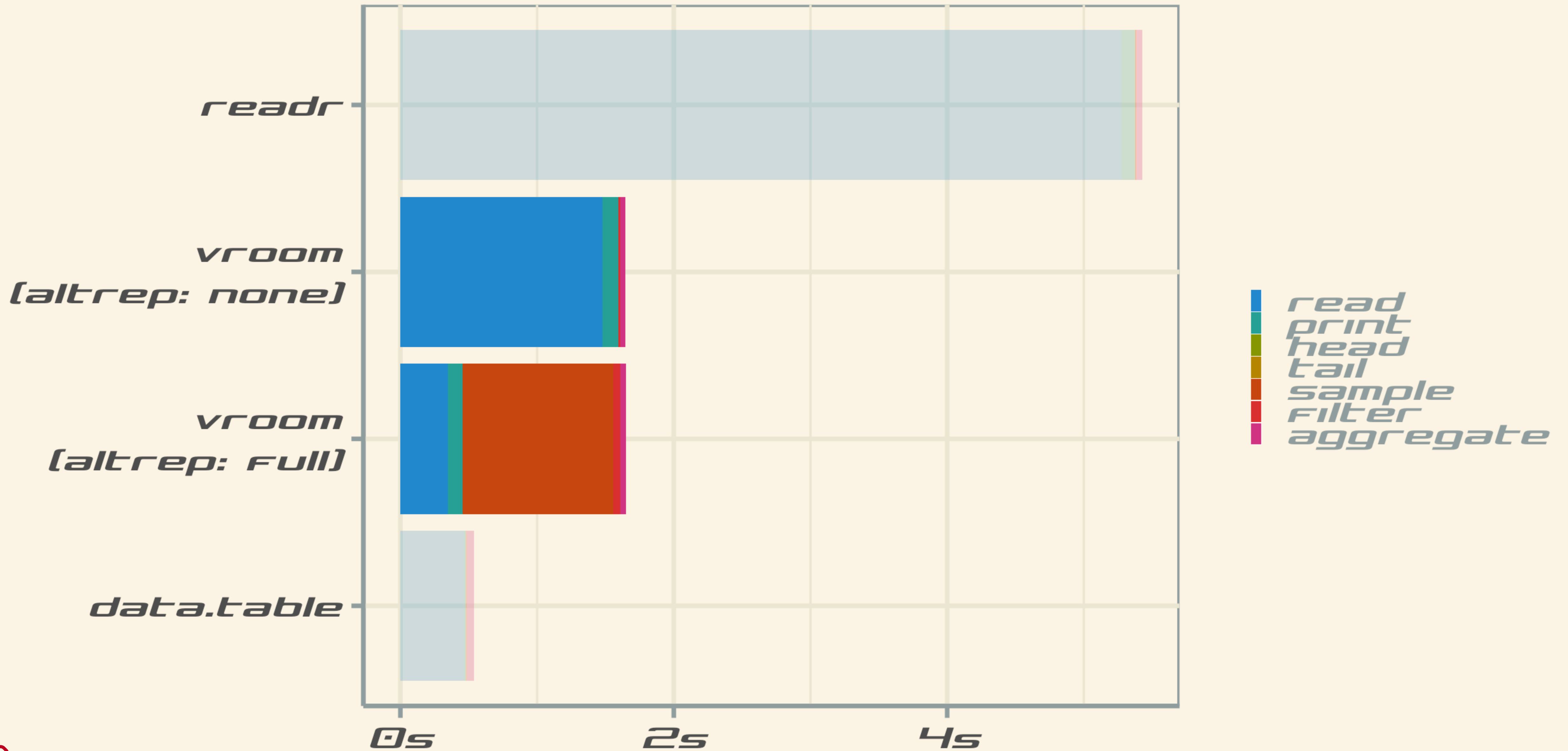
# All double

1,000,000 x 25 - 468MB



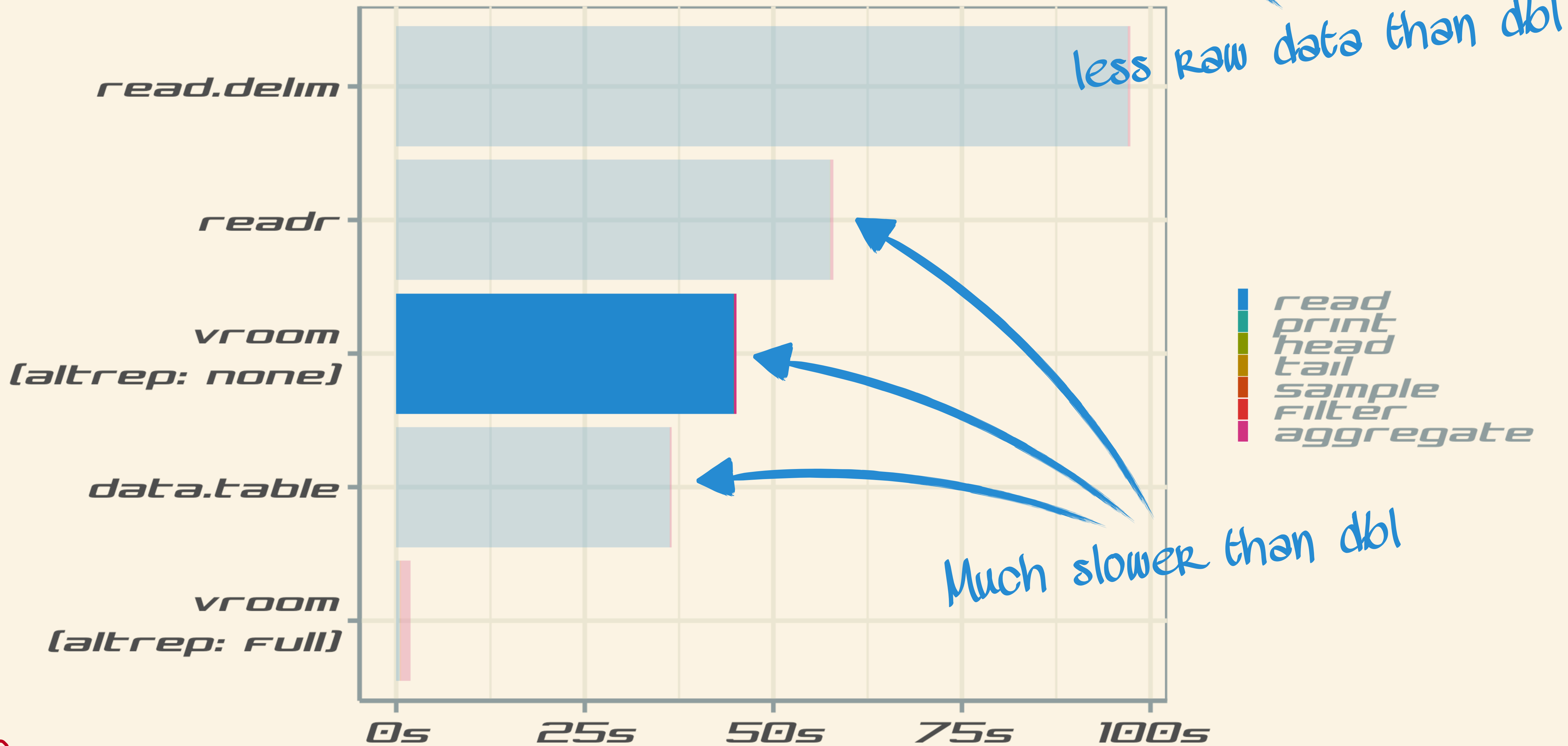
# All double

1,000,000 x 25 - 468MB



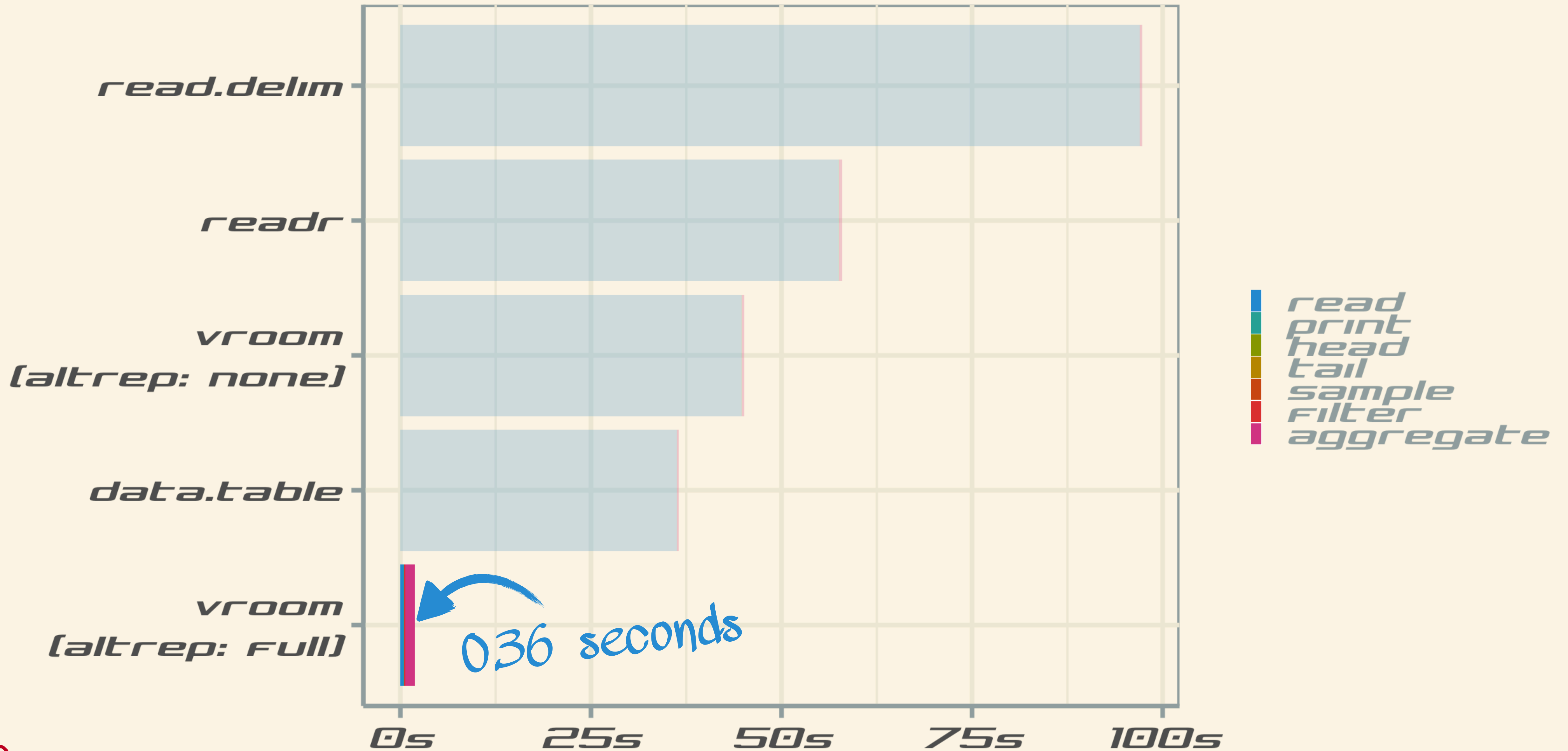
*ALL CHARACTERS*

*All character*  
*1,000,000 x 25 - 382MB*



# All character

1,000,000 x 25 - 382MB



*OTHER FEATURES*



# SELECTION

```
vroom("taxi.csv") %>%  
  dplyr::select(medallion, pickup_datetime, ends_with("amount"))  
  
# select  
vroom("taxi.csv",  
  col_select = list(medallion, pickup_datetime, ends_with("amount")))  
  
# remove  
vroom("taxi.csv", col_select = -hack_license)  
  
# rename  
vroom("taxi.csv", col_select = list(taxi = medallion, everything()))
```

# SELECTION

```
vroom("taxi.csv") %>%  
  dplyr::select(medallion, pickup_datetime, ends_with("amount"))  
  
# select  
vroom("taxi.csv",  
  col_select = list(medallion, pickup_datetime, ends_with("amount")))  
  
# remove  
vroom("taxi.csv", col_select = -hack_license)  
  
# rename  
vroom("taxi.csv", col_select = list(taxi = medallion, everything()))
```

# SELECTION

```
vroom("taxi.csv") %>%  
  dplyr::select(medallion, pickup_datetime, ends_with("amount"))  
  
# select  
vroom("taxi.csv",  
  col_select = list(medallion, pickup_datetime, ends_with("amount")))  
  
# remove  
vroom("taxi.csv", col_select = -hack_license)  
  
# rename  
vroom("taxi.csv", col_select = list(taxi = medallion, everything()))
```

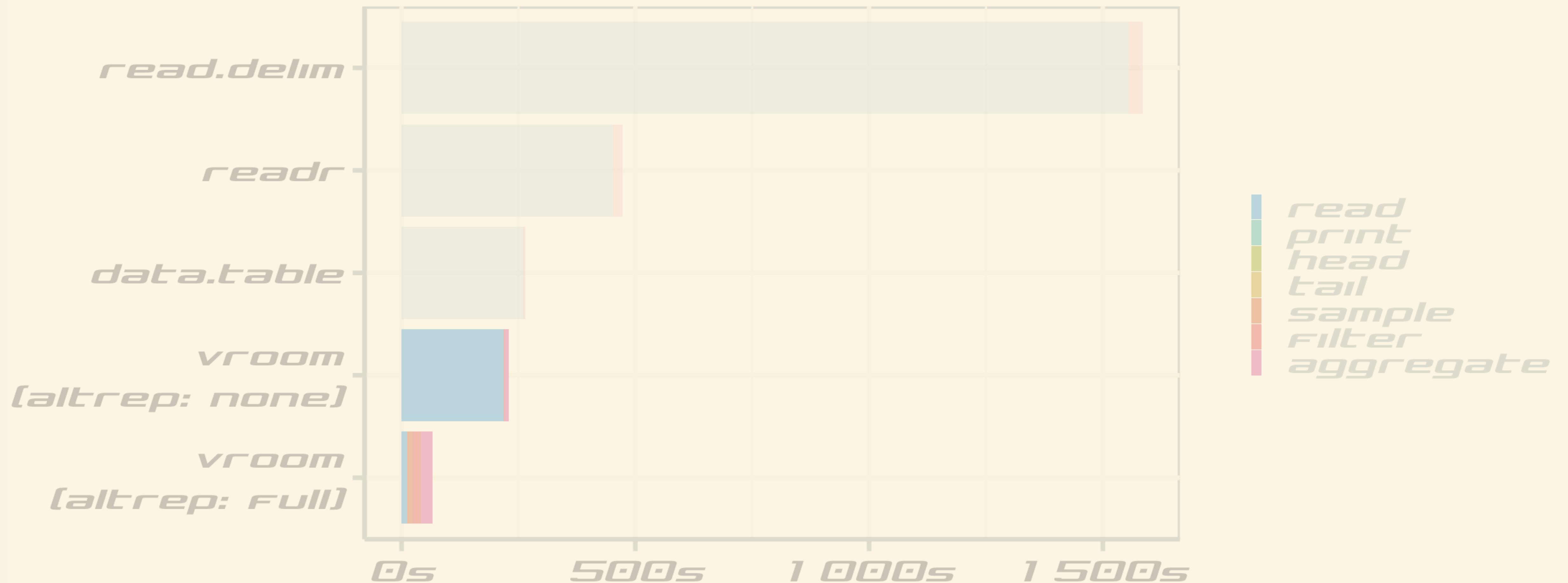
# SELECTION

```
vroom("taxi.csv") %>%  
  dplyr::select(medallion, pickup_datetime, ends_with("amount"))  
  
# select  
vroom("taxi.csv",  
  col_select = list(medallion, pickup_datetime, ends_with("amount")))  
  
# remove  
vroom("taxi.csv", col_select = -hack_license)  
  
# rename  
vroom("taxi.csv", col_select = list(taxi = medallion, everything()))
```

# MULTIPLE FILES

```
vroom(c("taxi_1.csv", "taxi_2.csv"), id = "path")
```

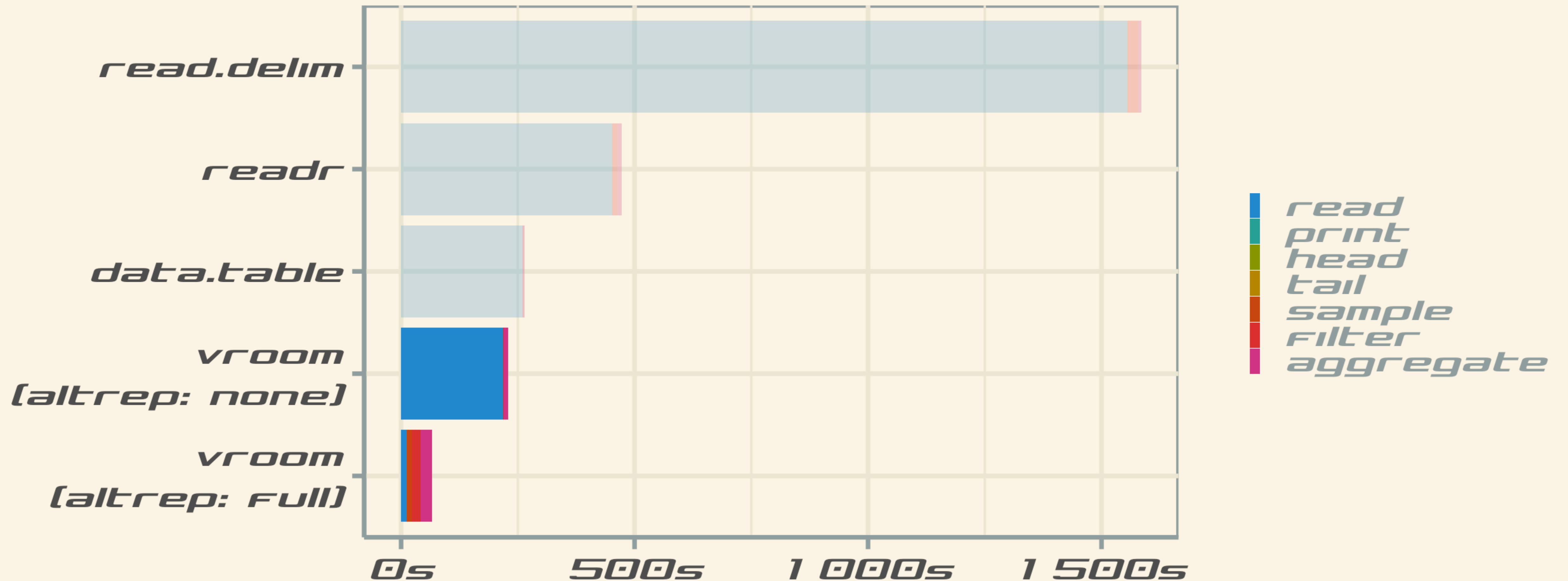
*Taxi trip fare (12 1.5Gb files)  
173,179,759 x 12 - 18.4GB*



# MULTIPLE FILES

```
vroom(c("taxi_1.csv", "taxi_2.csv"), id = "path")
```

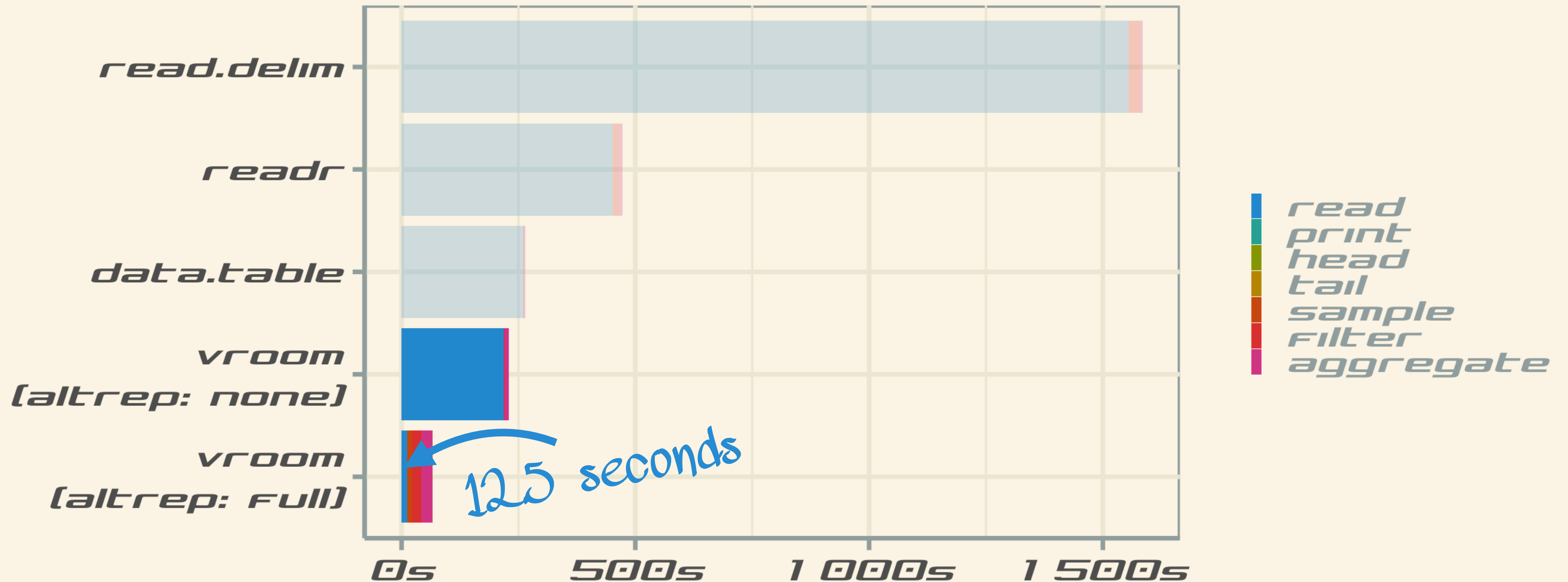
*Taxi trip fare (12 1.5Gb files)  
173,179,759 x 12 - 18.4GB*



# MULTIPLE FILES

```
vroom(c("taxi_1.csv", "taxi_2.csv"), id = "path")
```

*Taxi trip fare (12 1.5Gb files)  
173,179,759 x 12 - 18.4GB*



# *FIXED WIDTH FILES*

<code>vroom_fwf(altrep = TRUE)</code>	740 ms
<code>vroom_fwf(altrep = FALSE)</code>	10.1sec
<code>readr::read_fwf()</code>	24.3sec
<code>read.fwf()</code>	17.8min

*480,174 x 156 - 1.05 Gb File*



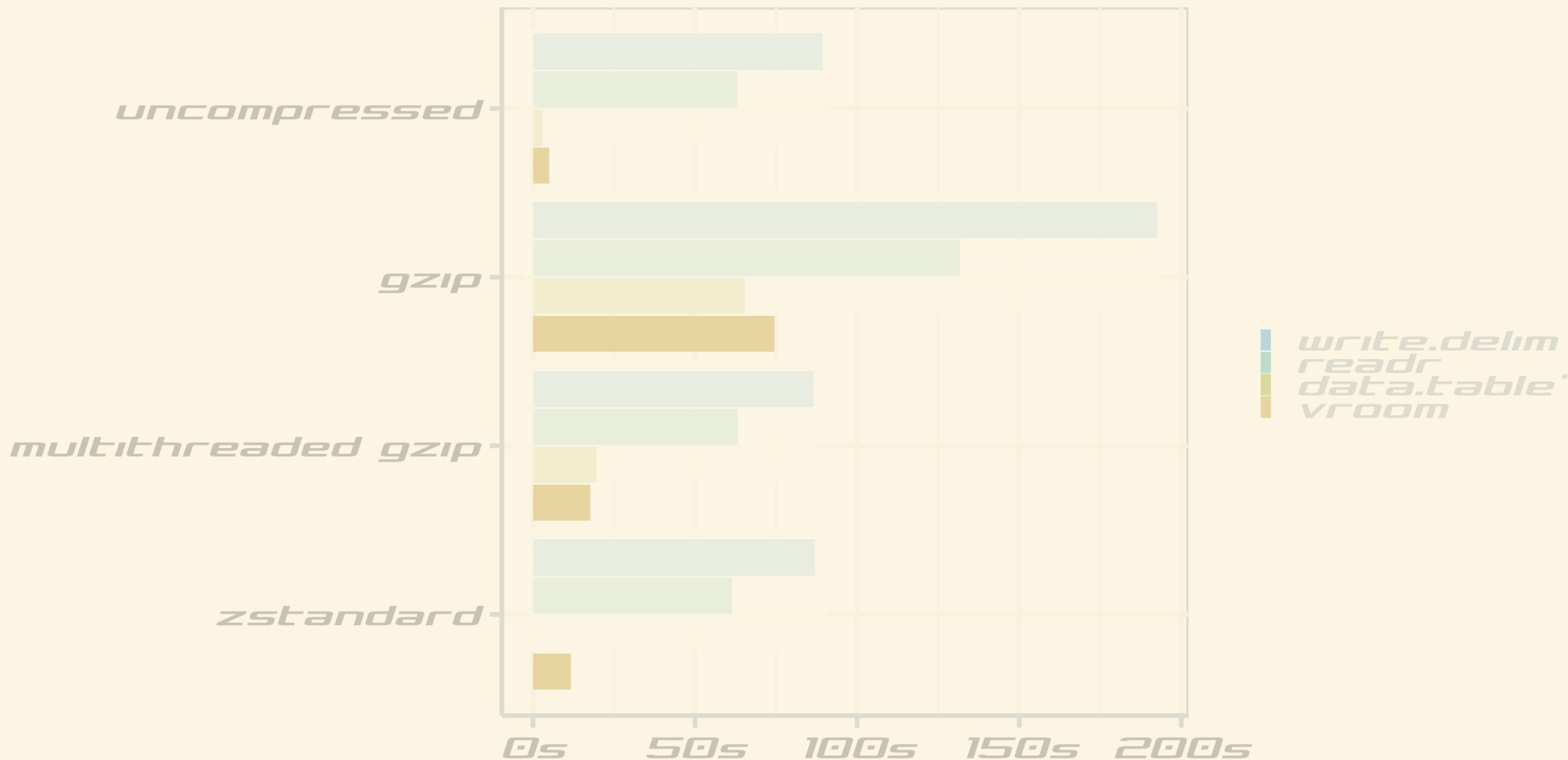
# COUNTING LINES

<code>length(vroom_lines(file))</code>	501 ms
<code>length(data.table::fread(file, sep = "\n", header = FALSE)[1])</code>	9.19sec
<code>length(readr::read_lines(file))</code>	11.82sec
<code>length(readLines(file))</code>	20.72sec

**14,776,616 lines - 1.55 Gb file**

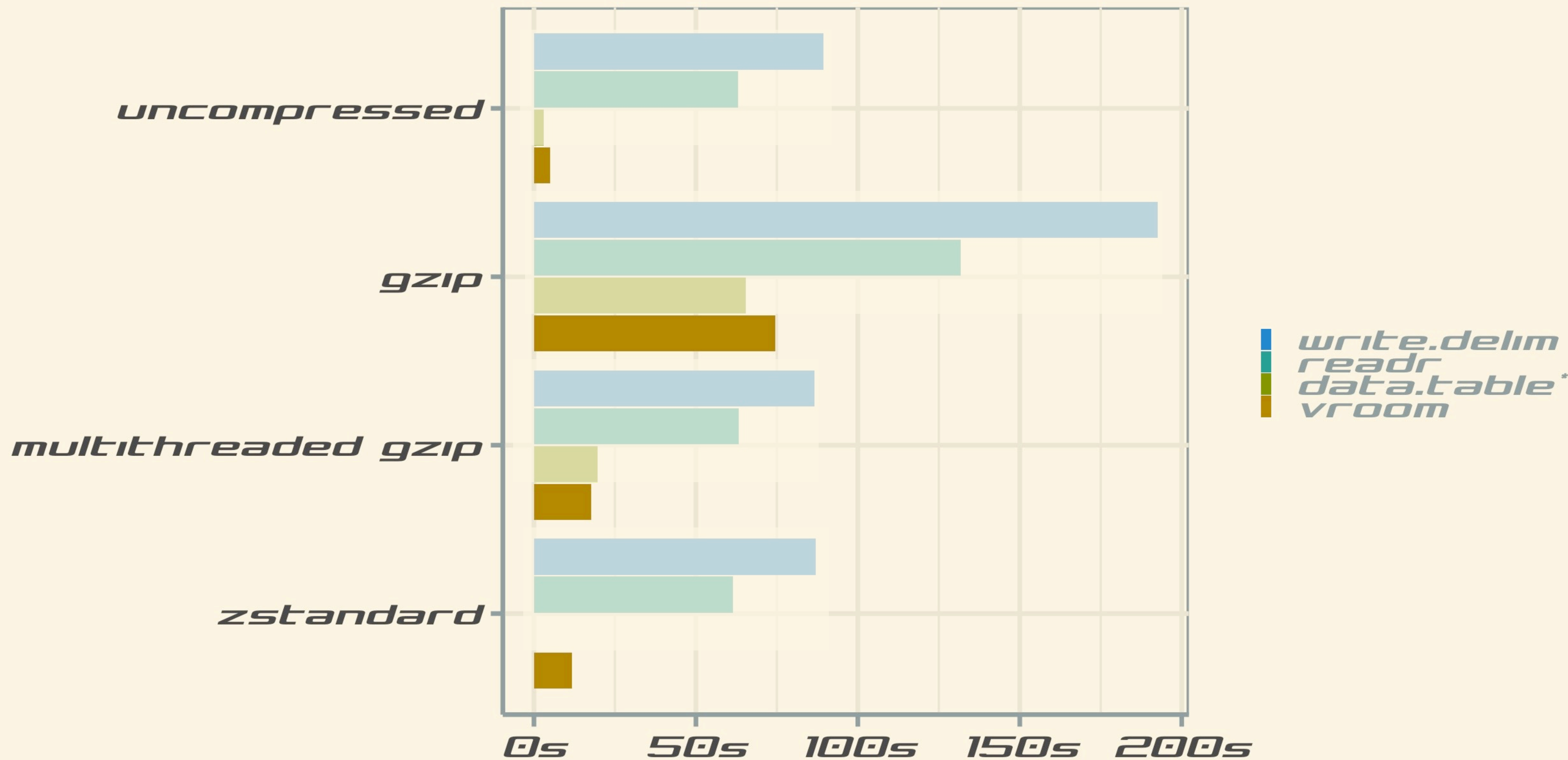
```
vroom_write(df, "taxi.tsv.gz", delim = "\t")
```

## Writing taxi trip fare 14,776,615 x 11 - 1.55GB



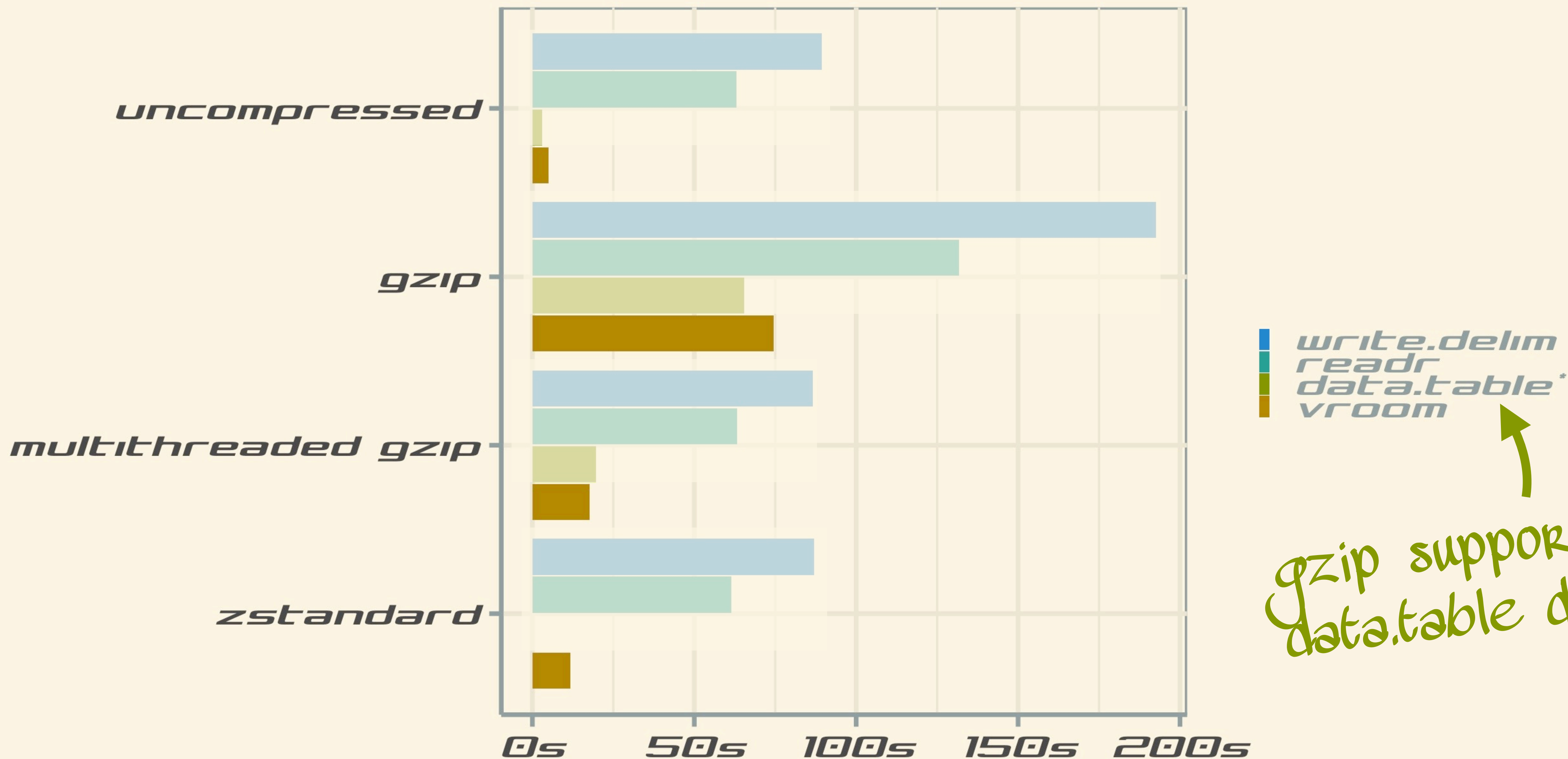
```
vroom_write(df, "taxi.tsv.gz", delim = "\t")
```

## Writing taxi trip fare 14,776,615 x 11 - 1.55GB

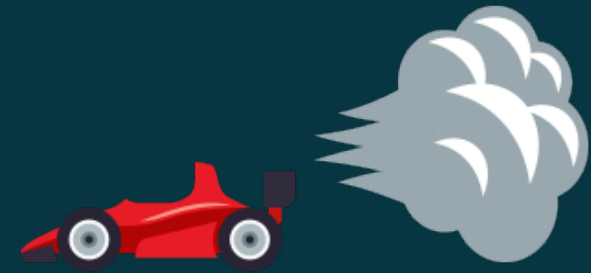


```
vroom_write(df, "taxi.tsv.gz", delim = "\t")
```

## Writing taxi trip fare 14,776,615 x 11 - 1.55GB



gzip support in data.table devel



# VROOM



`install.packages("vroom")`



`vroom.r-lib.org`



`bit.ly/vroom-yt`



`github.com/r-lib/vroom/issues`



@jimhester  
@jimhester\_