

# **Parallel fread() and data.table update**

**Bay Area R User Group Meetup**

**11 April 2017, Google**

**Matt Dowle**

# Overview

`data.table::fread()` is now parallel and available in dev; please test and report problems.

## Highlights from recent versions

# Try dev 1.10.5

```
install.packages("data.table", type = "source",  
  repos = "http://Rdatatable.github.io/data.table")
```

Windows binary on AppVeyor. See **Installation**.

Rdatatable/data.table

LATEST BUILD HISTORY

1.0.632  
9 days ago by Matt Dowle master 79a961e1 9 days ago in 7 min 14 sec

Jobs

JOB NAME	TESTS	DURATION
Environment: R_VERSION=devel, R_ARCH=x64		3 min 19 sec
Environment: R_VERSION=release, R_ARCH=x64		3 min 41 sec

LATEST BUILD HISTORY

1.0.632  
9 days ago by Matt Dowle master 79a961e1 9 days ago in 7 min 14 sec

Environment: R\_VERSION=release, R\_ARCH=x64

FILE NAME	TYPE	SIZE	DEPLOYMENT NAME
data.table.Rcheck\00check.log	File	3 KB	Logs
data.table.Rcheck\00install.out	File	6 KB	Logs
data.table.Rcheck\data.table-Ex.Rout	File	125 KB	Logs
data.table.Rcheck\tests\autoprnt.Rout	File	4 KB	Logs
data.table.Rcheck\tests\knitr.Rout	File	2 KB	Logs
data.table.Rcheck\tests\main.Rout	File	2 KB	Logs
data.table.Rcheck\tests\testthat.Rout	File	865 bytes	Logs
data.table_1.10.5.tar.gz	File	905 KB	Bits
data.table_1.10.5.zip	Zip archive	1 MB	Bits



**Guillermo Ponce** @guillermo\_ponce · Mar 20

Replying to @guillermo\_ponce @MattDowle

"fread" improvement for a 44 GB csv file: loading time went from 27mins (single-core) to 1.34 minutes (multi-core). 👍

← 1

↻ 5

❤️ 16

```

1  [ 100.0%]
2  [ 100.0%]
3  [ 100.0%]
4  [ 100.0%]
5  [ 100.0%]
6  [ 100.0%]
7  [ 100.0%]
8  [ 100.0%]
9  [ 100.0%]
10 [ 100.0%]
11 [ 100.0%]
12 [ 100.0%]
13 [ 100.0%]
14 [ 100.0%]
15 [ 100.0%]
16 [ 100.0%]
17 [ 100.0%]
18 [ 100.0%]
19 [ 100.0%]
20 [ 100.0%]
21 [ 100.0%]
22 [ 100.0%]
23 [ 100.0%]
24 [ 100.0%]
25 [ 100.0%]
26 [ 100.0%]
27 [ 100.0%]
28 [ 100.0%]
29 [ 100.0%]
30 [ 100.0%]
31 [ 100.0%]
32 [ 100.0%]
33 [ 100.0%]
34 [ 100.0%]
35 [ 100.0%]
36 [ 100.0%]
Mem [ 58.8G/1008G]
Swp [ 14.8M/4.88G]
56 [ 100.0%]
57 [ 100.0%]
58 [ 100.0%]
59 [ 100.0%]
60 [ 100.0%]
61 [ 100.0%]
62 [ 100.0%]
63 [ 100.0%]
64 [ 100.0%]
65 [ 77.9%]
66 [ 100.0%]
67 [ 100.0%]
68 [ 100.0%]
69 [ 100.0%]
70 [ 100.0%]
71 [ 100.0%]
72 [ 100.0%]
92 [ 100.0%]
93 [ 100.0%]
94 [ 100.0%]
95 [ 100.0%]
96 [ 100.0%]
97 [ 100.0%]
98 [ 100.0%]
99 [ 100.0%]
100 [ 100.0%]
101 [ 100.0%]
102 [ 100.0%]
103 [ 100.0%]
104 [ 100.0%]
105 [ 100.0%]
106 [ 100.0%]
107 [ 100.0%]
108 [ 100.0%]
128 [ 100.0%]
129 [ 100.0%]
130 [ 100.0%]
131 [ 100.0%]
132 [ 100.0%]
133 [ 100.0%]
134 [ 100.0%]
135 [ 100.0%]
136 [ 100.0%]
137 [ 100.0%]
138 [ 100.0%]
139 [ 100.0%]
140 [ 100.0%]
141 [ 100.0%]
142 [ 100.0%]
143 [ 100.0%]
144 [ 100.0%]
Tasks: 569, 2490 thr; 145 running
Load average: 0.54 1.65 1.88

```



**Sctt Stnflid** @seesharp · Apr 8

@MattDowle I'm not sure what you did between 1.10.4 and 1.10.5, but fread() on 3.6M row CSV went from 2m 42s down to 24s :)

← 1

↻ 4

❤️ 9

Not new. Prior art.

h2o.importFile 3 years

spark-csv 2 years

Python's Paratext 1 year





**Matt Dowle**  
@MattDowle

## New package "fst: Lightning Fast Serialization of Data Frames" by Mark Klik [fstpackage.org](http://fstpackage.org) #rstats #python

Method	Format	Time (s)	Size (MB)	Speed (MB/s)	N
read.fst	bin	0.0611263	133	3271.9	1397
write.fst	bin	0.1427031	133	1401.5	1397
read_feather	bin	0.2228486	163	897.5	1270
readRDS	bin	0.3089175	200	647.4	1397
saveRDS	bin	0.3406979			
write_feather	bin	0.7631548			
fwrite	csv	0.9813563	392	203.8	1270
fread	csv	17.1259319			

data.table::fwrite so fast that it was deemed on par with binary and included

Should now be faster

RETWEETS  
**86**

LIKES  
**184**



12:43 PM - 23 Jan 2017

# Beware of cache when benchmarking

First timing longest (OS reads from device)

`free -g` (OS cached file in RAM)

`sudo sh -c 'echo 3 >/proc/sys/vm/drop_caches'`

Aside: HD has cache too (burst vs sustained)

`sudo lshw -class disk`

`sudo hdparm -t /dev/sda`

HD 150MB/s ; SSD 800MB/s ; NVMe 3GB/s



**R CMD INSTALL ~/data.table\_1.10.4.tar.gz # CRAN**

**perfbar &**

**htop**

**system.time(fread("~/X3e8\_2c.csv",verbose=TRUE)) # 5.6GB file**

**# 27s first time, 23s second time.**

**# Awful! CPU not IO bound.**

**R CMD INSTALL ~/data.table\_1.10.5.tar.gz # dev**

**system.time(fread("~/X3e8\_2c.csv",verbose=TRUE))**

**# 7s first time (5.6GB file size / 800MB/s SSD speed == 7s)**

**# 3.5s second time**

**free -g**

**system("sudo sh -c 'echo 3 >/proc/sys/vm/drop\_caches'")**

**free -g**

**system.time(fread("~/X3e8\_2c.csv",verbose=TRUE))**

**# 7s first time, 3.5s second time**

```
$ lscpu
Architecture:          x86_64
CPU op-mode(s):        32-bit, 64-bit
Byte Order:             Little Endian
CPU(s):                 8
On-line CPU(s) list:   0-7
Thread(s) per core:    2
Core(s) per socket:    4
Socket(s):              1
NUMA node(s):          1
Vendor ID:              GenuineIntel
CPU family:             6
Model:                  70
Model name:             Intel(R) Core(TM)
                       i7-4980HQ CPU @ 2.80GHz
Stepping:               1
CPU MHz:                2377.156
CPU max MHz:            4000.0000
CPU min MHz:            800.0000
BogoMIPS:               5587.25
Virtualization:        VT-x
L1d cache:              32K
L1i cache:              32K
L2 cache:               256K
L3 cache:               6144K
L4 cache:               131072K
NUMA node0 CPU(s):     0-7

$ sudo hdparm -t /dev/sda
795.15 MB/sec
```

# Arun's update, Jan 2017 Amsterdam

## THE #RDATA TABLE PACKAGE

+ new developments in v1.10.0

### Arun Srinivasan

JAN 20'17, AMSTERDAM



@ARUN\_SRINIV

# Highlights

Already on CRAN :

No longer need `with=FALSE`

`setkey()` partially parallel

`keyby=` much faster than `by=`