# Fast automatic indexing with data.table

**R/Finance, Chicago**

**30 May 2015**

**Matt Dowle**

# Yesterday

***Thomas to me:***

**"dplyr has completely killed off data.table"**

**I've tilted this presentation to address this.**

# 1964

*U.S. Supreme Court Justice Stewart :*

**"I can't define it but I know it when I see it."** (paraphrased)

**data.table users know they need data.table**

# https://github.com/Rdatatable/data.table/wiki

fast **aggregation** of large data; e.g. 100GB in RAM (see benchmarks on up to two billion rows)

fast **ordered joins**; e.g. rolling forwards, backwards, nearest and limited staleness

fast **overlapping range joins**; e.g. GenomicRanges

fast add/modify/delete of columns **by reference** by group using no copies at all

cells may themselves contain vectors/objects/functions; i.e. **columns of type list**

fast and friendly file reader: **fread**

+ research into production (e.g. daily or intra-day) with no code changes
+ or might in future
+ brief syntax to prevent code bloat; e.g. do anything in j
+ optimization of combined [where, select|update, by]

H$_2$O.ai
Machine Intelligence

```
> DT     # 1.5GB

          id val

1e+00:  BAR     5

2e+00:  FOO     1

3e+00:  REW     4

4e+00:  NUR     5

5e+00:  AMW     3

   ---

1e+08:  QNP     1

1e+08:  HXB     2

1e+08:  FOO     1

1e+08:  CYY     2

1e+08:  VKG     1
```

```
> DT[id=="FOO",]
        id val
  1: OSK    1
  2: OSK    3
 ---
5813: OSK    5
5814: OSK    1

  user   system elapsed
 1.928    0.064   1.991
```

```
> DT[id=="BAR",]

  user   system elapsed
 0.000    0.000   0.001
```

```
> DT[id %in% c("FOO","BAR"),]

  user   system elapsed
 0.000    0.000   0.001
```

H₂O.ai
Machine Intelligence

```
> options(datatable.verbose=TRUE)

> DT[id=="FOO",]

creating new index 'id'

forder took 1.932 sec

Starting bmerge ...done in 0.00 secs


> DT[id=="BAR",]

using existing index 'id'

Starting bmerge ...done in 0.00 secs
```

```
> DF %>% filter(id=="FOO")

   user   system elapsed
  1.952    0.020   1.970
> DF %>% filter(id=="FOO")

   user   system elapsed
  1.940    0.012   1.949


> DF[DF$id=="FOO", ]

   user   system elapsed
  2.244    0.124   2.367
> DF[DF$id=="FOO", ]

   user   system elapsed
  2.260    0.112   2.369
```

```
> DT %>% filter(id=="FOO")   # v0.3.0.2
                             # Oct 2014

using existing index 'id'
Starting bmerge ...done in 0 secs
   user   system elapsed
  0.000    0.000   0.001


> DT %>% filter(id=="FOO")   # v0.4.0
                             # Jan 2015

   user   system elapsed
  1.952    0.020   1.982
```
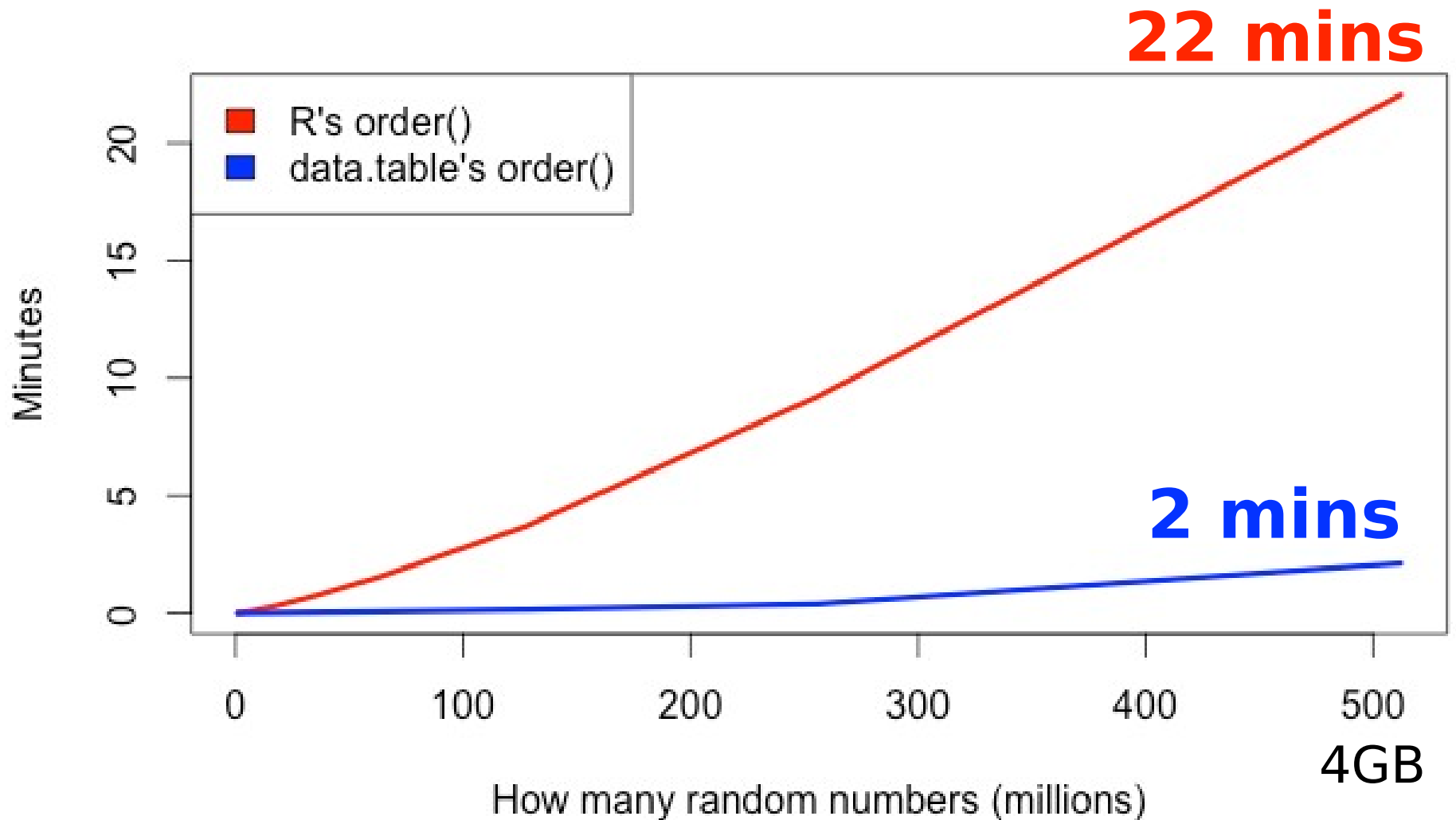
**22 mins**

**2 mins**

4GB

MacBook Pro 2.8GHz Intel Core i7 16GB
R 3.1.3   data.table 1.9.4

$H_2O$.ai
Machine Intelligence

# References

**Terdiman, 2000**

http://codercorner.com/RadixSortRevisited.htm

**Herf, 2001**

http://stereopsis.com/radix.html

Arun Srinivasan implemented forder() in data.table entirely in C for integer, character and double

Matt Dowle changed from LSD (backwards) to MSD (forwards)

H<sub>2</sub>O.ai
Machine Intelligence

## Pros

- Index storage is small and fixed:  nrow * 4|8 bytes
- No collisions in hash table (no hash table)
- Building new indexes may be able to reuse existing indexes
- Rolling joins and overlapping range joins

## Cons

- Insert and delete of rows requires memmove
- Binary search vs direct hash table lookup (note though collisions)

# H2O

Machine learning e.g. Deep Learning (GBM)

In-memory, parallel and distributed

1. Data >> 250GB needle-in-haystack e.g. fraud

2. Data < 250GB compute intensive, in parallel on 1000 cores

3. Data < 250GB but where feature engineering > 250GB

Speed for production

Open source on GitHub, liberal Apache license

$H_2O$.ai
Machine Intelligence

# Install H2O

```
$ sudo add-apt-repository -y ppa:webupd8team/java

$ sudo apt-get update

$ sudo apt-get -y install oracle-java8-installer

$ sudo apt-get -y install oracle-java8-set-default

$ java -version


$ R

> install.packages("h2o")
```

# Start H2O

> require(h2o)

> h2oServer <- h2o.init()

H2O is not running yet, starting it now...

Successfully connected to http://127.0.0.1:54321

R is connected to H2O cluster:

    H2O cluster uptime:        1 sec 397 ms

    H2O cluster version:       2.8.4.4

    H2O cluster total nodes:    1

    H2O cluster total memory:   26.67 GB

    H2O cluster total cores:    32

# h2o.importFile

- 23GB .csv, 9 columns, 500e6 rows

  ```
  > system.time(DF <- h2o.importFile(h2oServer, path = "/dev/shm/test.csv"))
  ```

  ```
   user  system elapsed
   0.775   0.058  50.559
  ```

  ```
  > head(DF)
    id1   id2        id3 id4 id5  id6 v1 v2      v3
  1 id076 id035 id0000003459  20  80 8969  4  3 43.1525
  2 id062 id023 id0000002848  99  49 7520  5  2 86.9519
  3 id001 id052 id0000007074  89  16 8183  1  3 19.6696
  ```

h2o.importFile                          # 50 sec


require(data.table)

fread("/dev/shm/test.csv")          # 290 sec


require(readr)

read_csv("/dev/shm/test.csv")      # 720 sec

# Why not 30x?

- Maybe IO bound but this test.csv was on ramdisk /dev/shm

- H2O compresses the data in RAM

- Profiles the data while reading e.g. min and max per column, for later efficiency gains.

# Thank you

https://github.com/Rdatatable/data.table/wiki

http://h2o.ai/product