# Timings of common tasks using the **data.table** package in R

Matthew Dowle

May 3, 2010

* WORK IN PROGRESS *

This document contains a series of tests, followed by a summary table of various timings and comparisons. Please go straight to the summary table first <here> in which each row has a link back to the test.

This document is reproducible. Simply run the .Rnw file yourself in your environment to confirm the results. Also see ?vignette, which says that edit(vignette("datatable-timings")) will extract the code from this document so you can easily work with it.

The .Rnw included in the package has N=10,000,000. This is a small number so that 'R CMD build' completes in a reasonable time (about 5 minutes). We don't want the nightly builds on R-Forge and CRAN to slow down just to run long timing comparisons. We have increased this to N=100,000,000 ourselves, and included the output on the datatable homepage (<link>).

## Contents

## 1 Timing tests

### 1.1 Extraction

This is a repeat of the test in section 1 of the Introduction vignette. The syntax is explained there. This demonstrates the large difference in speed between vector scans and binary search. Therefore, please avoid using `==` in the `i` expression.

```
> n = ceiling(1e7/26^2)   # 10 million rows
> DF = data.frame(x=rep(LETTERS,each=26*n),
+                 y=rep(letters,each=n),
+                 v=rnorm(n*26^2))
> DT = data.table(DF,key="x,y")
> tables()

     NAME        NROW  MB COLS  KEY
[1,] DT    10,000,068 153 x,y,v x,y
Total: 153MB

> tt=system.time(ans1 <- DF[DF$x=="R" & DF$y=="h",]); tt

   user  system elapsed
  4.164   1.084   5.556
```

```
> ss=system.time(ans2 <- DT[J("R","h"),mult="all"]); ss

   user  system elapsed
  0.012   0.000   0.012

> mapply(identical,ans1,ans2)

    x    y    v
 TRUE TRUE TRUE
```

## 1.2 Grouping

This is a repeat of the test in section 2 of the Introduction vignette. The syntax is explained there.

```
> ttt=system.time(ans1 <- tapply(DF$v,DF$x,sum)); ttt

   user  system elapsed
  9.517   0.960  11.636

> sss=system.time(ans2 <- DT[,sum(v),by=x]); sss

   user  system elapsed
  0.548   0.244   0.832

> identical(as.vector(ans1), ans2$V1)

[1] TRUE
```

## 1.3 Test 3

## 1.4 Test 4

## 1.5 Test 5

# 2 Summary table

```
> ans

         base data.table times faster
==      5.556      0.012          462
tapply 11.636      0.832           13

> toLatex(sessionInfo())
```

- R version 2.11.0 (2010-04-22), i486-pc-linux-gnu

- Locale: LC_CTYPE=en_GB.UTF-8, LC_NUMERIC=C, LC_TIME=en_GB.UTF-8,
  LC_COLLATE=en_GB.UTF-8, LC_MONETARY=C, LC_MESSAGES=en_GB.UTF-8,
  LC_PAPER=en_GB.UTF-8, LC_NAME=C, LC_ADDRESS=C, LC_TELEPHONE=C,
  LC_MEASUREMENT=en_GB.UTF-8, LC_IDENTIFICATION=C

- Base packages: base, datasets, graphics, grDevices, methods, stats, tools, utils

- Other packages: data.table~1.4.1