

# ADaCGH2: parallelized analysis of (big) CNA data

Ramon Diaz-Uriarte

Department of Biochemistry, Universidad Autónoma de Madrid, Instituto de Investigaciones Biomédicas ‘Alberto Sols’ (UAM-CSIC), 28029 Madrid, Spain

Associate Editor: John Hancock

## ABSTRACT

**Motivation:** Studies of genomic DNA copy number alteration can deal with datasets with several million probes and thousands of subjects. Analyzing these data with currently available software (e.g. as available from BioConductor) can be extremely slow and may not be feasible because of memory requirements.

**Results:** We have developed a BioConductor package, ADaCGH2, that parallelizes the main segmentation algorithms (using forking on multicore computers or parallelization via message passing interface, etc., in clusters of computers) and uses *ff* objects for reading and data storage. We show examples of data with 6 million probes per array; we can analyze data that would otherwise not fit in memory, and compared with the non-parallelized versions we can achieve speed-ups of 25–40 times on a 64-cores machine.

**Availability and implementation:** ADaCGH2 is an R package available from BioConductor. Version 2.3.11 or higher is available from the development branch: <http://www.bioconductor.org/packages/development/html/ADaCGH2.html>.

**Contact:** [ramon.diaz@iib.uam.es](mailto:ramon.diaz@iib.uam.es)

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

Received on July 1, 2013; revised on January 31, 2014; accepted on February 11, 2014

## 1 INTRODUCTION

Current studies of genomic copy number alterations (CNA) are using platforms with several million probes per array and several thousand subjects (e.g. Grozeva *et al.*, 2012), but the canonical implementations of the widely used open-source packages for the analysis of CNA data did not allow for parallelized execution of the analysis. This makes it difficult to use clusters of servers and does not take advantage of the trends in parallel computing toward multicore machines (servers with 16–124 cores or laptops with 2 or 4 cores are nowadays common). Moreover, especially for R/BioConductor software, the available packages will not be able to analyze big datasets if these are larger than about a quarter to a half of the available memory (unless one uses time-consuming, and *ad hoc*, manual partition of the input and subsequent recombination of the output—see discussion in Section ‘Why ADaCGH2 instead of a “manual” solution’ in the vignette of the package).

Here I describe ADaCGH2, a BioConductor package designed to address the above deficiencies. I leverage on two R packages, *parallel*, part of the standard set of R packages,

and *ff* (Adler *et al.*, 2013), and combine them, to provide the following:

- Parallelized analysis. I have parallelized the most widely used segmentation approaches that can be applied to CNA data, including both comparative genomic hybridization (CGH) and SNP arrays (Valsesia *et al.*, 2013), and also covering sequencing data, when these have been appropriately processed, but see Duan *et al.* (2013), Wu *et al.* (2013), Zhao *et al.* (2013) and Zheng *et al.* (2013). The methods available are CBS (Venkatraman and Olshen, 2007), HaarSeg (Ben-Yaacov and Eldar, 2008), HMM (Fridlyand *et al.*, 2004), BioHMM (Marioni *et al.*, 2006), the wavelet-based method from Hsu *et al.* (2005), GLAD (Hupe *et al.*, 2004) and CGHseg (Picard *et al.*, 2005) and two merging algorithms. Some of those methods, however, are not suitable for large datasets—see details in Section 1.2.1 of the ‘benchmarks.pdf’ package vignette.
- I use package *parallel* to provide parallelization using (i) forking, for single multicore computers; (ii) parallelization with message passing interface (MPI), sockets, parallel virtual machine, etc., for clusters built of several computers.
- The ability to analyze data that cannot fit in memory. Using *ff*, we only access the section of the data currently being analyzed, keeping in RAM and moving between processes only a pointer to the rest of the data on disk.
- Parallelization of data input and output and plotting.
- Input from, and output to, other BioConductor packages.

Here I present the main functions of the package, the differences with former version and some benchmarks. A full set of examples, further benchmarks and detailed suggestions for usage are included in the package vignettes.

## 2 DIFFERENCES WITH THE FORMER VERSION

ADaCGH was first developed to provide parallelized analysis of CNA data for web-based applications (Carro *et al.*, 2010; Diaz-Uriarte and Rueda, 2007). The first version parallelized eight segmentation algorithms (using MPI), was available from CRAN and was last updated in 2009, but will no longer run without tweaks because it depends on *papply*, a package that will not install in old versions of R. Next, parallelization was extended so clusters were not limited to MPI clusters, and *ff* objects were used for storage; that version is available as version 1.10 from BioConductor 2.12. For the current version, most of

the code has been rewritten to use forking, data handling and reading of input data has been completely modified so that data much larger than available memory can be read and analyzed and missing value handling has been changed to use all available data per array. The vignette `benchmarks.pdf` provides extensive comparisons between the new ( $\geq 2.3.5$ ) and latest previous running versions (version 1.10), but the main differences between these two versions are as follows:

- **Reading and analysis of large datasets.** The new version can read datasets much larger than the old one and datasets much larger than available memory (see details in Section 3). As a consequence of being able to read much larger datasets, the new version can analyze datasets much larger than the old one.
- **Missing value handling.** The old version used row-wise deletion of missing values when reading data (i.e. a probe would be deleted from the data if it had one missing value in any array/column). The new version deals with missing values array by array, so for each array (or column) all available data (or probes) are used in the segmentation.
- **Forking and clusters.** The new version of ADaCGH2 allows for the usage of forking or an explicit cluster (e.g. MPI, sockets) to parallelize reading and analysis. In POSIX operating systems (including Unix, GNU/Linux and Mac OS), forking can be faster, less memory consuming and much easier to use than a cluster.
- **Flexibility of reading data and compatibility with former version.** The new version of ADaCGH2 has not removed the mechanisms of reading data available in the old version (when data are small or memory is plentiful, reading data from a single RData is an available option) and accepts data read by the former version. However, the old version cannot accept data read by the new version because it assumes that data that have been read contain no missing values.

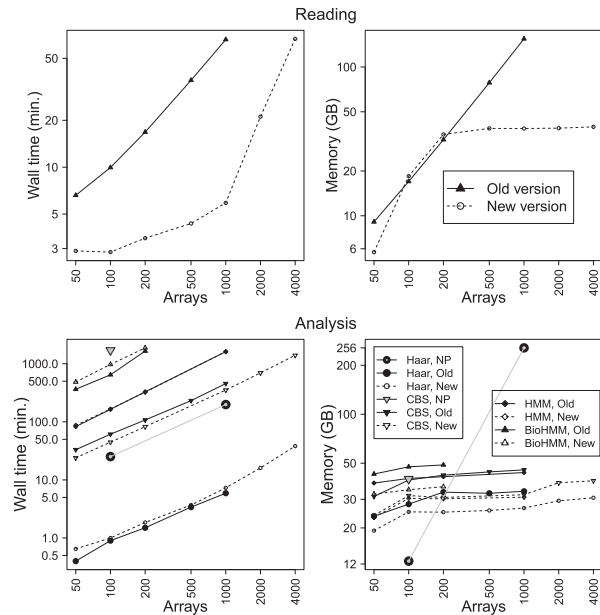
These differences in implementation, however, do not affect the underlying core code for the algorithms, which is the same as in the previous version. There have been, however, changes in some defaults to adapt the package to really large data (e.g. using MAD as merging default or using 'haarseg' as the 'smoothfunc' for daglad, following recommendations in the package vignette for GLAD).

### 3 BENCHMARKS

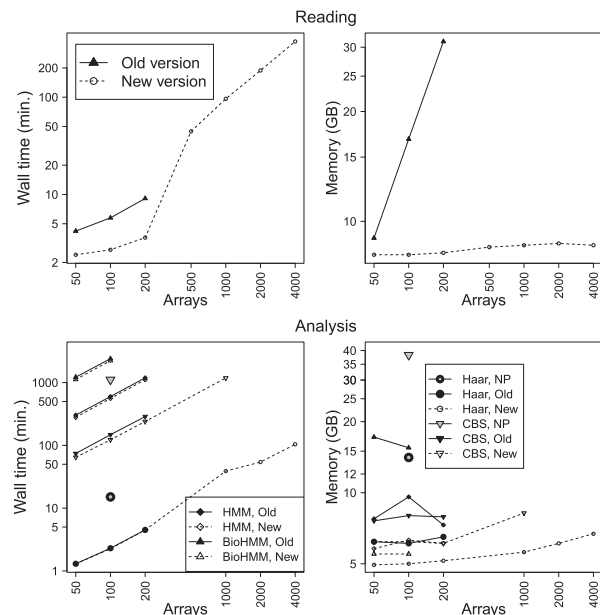
Figure 1 shows benchmarks of reading and analyzing data with 6067433 probes per array/column. Those figures compare memory usage and wall time of the old and new versions and of the non-parallelized (NP) versions in two different machines (data for the figures, as well as benchmarks for a third machine, and with MPI over two machines, are available from the vignette 'benchmarks.pdf'). To give an idea of sizes, the RData file for the 1000 arrays data is of ~41 GB, and the directory with the data for 2000 columns/arrays occupies ~198 GB.

Compared with the NP version, in the analysis of data, ADaCGH2 leads to speed increases by factors of 25–40 times

(a) AMD Opteron 6276, 64 cores, 256 GB RAM



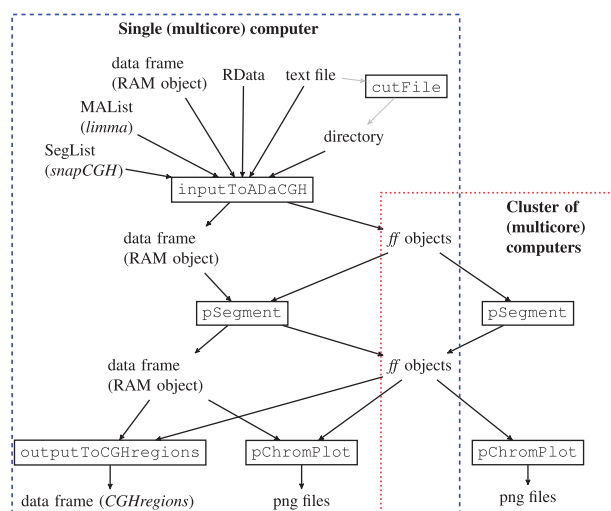
(b) Intel Xeon E5645, 12 cores, 64 GB RAM



**Fig. 1.** Wall time and memory use (summed over all spawned processes) of reading and analysis as a function of number of arrays. Reading: comparison between new and old versions. Analysis: new and old versions with four segmentation methods and NP for two methods. No benchmark allowed to run for >36 h. Without parallelization, in the AMD machine no runs of CBS with 1000 arrays or HaarSeg with 2000 can be done (R runs out of memory); in the Intel machine no runs for 1000 arrays with any method can be done (R runs out of memory)

in the 64 cores machines and 7–10 times in the 12 cores machines, and allows us to analyze data that would not fit in memory.

Compared with the former version, the new version uses less memory for analysis. More importantly, the new version allows us to read and analyze much larger datasets. In the 256 and



**Fig. 2.** Work flow from input data to figures. R functions are shown with courier font inside boxes. In italics are the names of other packages, which can provide input/accept output. Data frames are stored in memory, in contrast to *ff* objects. Data input and conversion to *ff* objects are done with `inputToADaCGH` (maybe after using `cutFile` for parallelized reading of single-column files). Segmentation is performed with `pSegment`, and results can then be plotted (`pChromPlot`) or used by other packages (`outputToCGHregions`). When using *ff* objects, after data input (in a single machine) the remaining analyses can be conducted in a cluster

384 GB machines the old version cannot read datasets with  $\geq 2000$  arrays (R runs out of memory), and in the machine with 64 GB of RAM it cannot read data with  $\geq 500$  arrays (R runs out of memory); as can be seen from the figure, the old version shows a steep linear increase in memory consumption with number of arrays. In sharp contrast, with the new version we can read and analyze 4000 arrays in a machine with only 64 GB of RAM (Fig. 1b), and the scaling of memory usage with number of arrays suggests that much larger datasets could be read and analyzed. In addition, we can obtain speedups by factors of 2–10 times (depending on machine and number of arrays) in the reading step as it is parallelized.

## 4 WORK FLOW

Figure 2 shows the usual sequence of calls with ADaCGH2. `inputToADaCGH` accepts input in different formats, including objects used by `limma` (Smyth, 2005) and `snapCGH` (Smith *et al.*, 2009), and produces R data frames or *ff* objects, after performing several checks and data sanitation. If data are read from a directory with one-column files, reading is parallelized (`cutFile` allows splitting a text file into one-column files). `pSegment` can take as input R data frames and *ff* objects produced by `inputToADaCGH`. `pSegment` can use multiple cores or multiple computers, and it can accept as input data frames or *ff* objects; when running on a cluster, only *ff* objects are used (to avoid passing around large objects and to allow analyzing large datasets). The output from `pSegment` can be converted so it is accepted by the `CGHregions` package (Vosse and van de Wiel, 2009), and creation of figures is also parallelized.

## 5 CONCLUSIONS

ADaCGH2 should be of immediate use for researchers involved in the analysis of CNA data. Parallelization allows it to speed up data processing, and it can handle data that will not fit in memory with excellent scaling of memory usage with number of arrays. These behaviors are needed for the analyses of platforms with increasing number of probes and multicenter studies with thousands of subjects.

## ACKNOWLEDGEMENTS

The author thanks O. M. Rueda and D. Rico for collaboration in previous versions of ADaCGH and for help debugging under Mac OS, O. M. Rueda for the initial benchmarking data, C. Lázaro-Perea, O. M. Rueda, I. López and three anonymous reviewers for comments on the manuscript.

*Funding:* Spanish MINECO (Project BIO2009-12458).

*Conflict of Interest:* none declared.

## REFERENCES

- Adler, D. *et al.* (2013) *ff: Memory-Efficient Storage of Large Data on Disk and Fast Access Functions*. R package version 2, pp. 2–11.
- Ben-Yaacov, E. and Eldar, Y.C. (2008) A fast and flexible method for the segmentation of aCGH data. *Bioinformatics*, **24**, i139–i145.
- Carro, A. *et al.* (2010) waviCGH: a web application for the analysis and visualization of genomic copy number alterations. *Nucleic Acids Res.*, **38**, W182–W187.
- Diaz-Uriarte, R. and Rueda, O.M. (2007) ADaCGH: a parallelized web-based application and R package for the analysis of aCGH data. *PLoS One*, **2**, e737.
- Duan, J. *et al.* (2013) Comparative studies of copy number variation detection methods for next-generation sequencing technologies. *PLoS One*, **8**, e59128.
- Fridlyand, J. *et al.* (2004) Hidden Markov models approach to the analysis of array CGH data. *J. Multivar. Anal.*, **90**, 132–153.
- Grozeva, D. *et al.* (2012) Independent estimation of the frequency of rare CNVs in the UK population confirms their role in schizophrenia. *Schizophr. Res.*, **135**, 1–7.
- Hsu, L. *et al.* (2005) Denoising array-based comparative genomic hybridization data using wavelets. *Biostatistics*, **6**, 211–226.
- Hu, P. *et al.* (2004) Analysis of array CGH data: from signal ratio to gain and loss of DNA regions. *Bioinformatics*, **20**, 3413–3422.
- Marioni, J.C. *et al.* (2006) BioHMM: a heterogeneous hidden Markov model for segmenting array CGH data. *Bioinformatics*, **22**, 1144–1146.
- Picard, F. *et al.* (2005) A statistical approach for array CGH data analysis. *BMC Bioinformatics*, **6**, 27.
- Smith, M.L. *et al.* (2009) *snapCGH: Segmentation, Normalisation and Processing of aCGH Data*. R package version 1.31.0.
- Smyth, G.K. (2005) *Limma: linear models for microarray data*. In: Gentleman, R. *et al.* (ed.) *Bioinformatics and Computational Biology Solutions Using R and Bioconductor*. Springer, New York, pp. 397–420.
- Valsesia, A. *et al.* (2013) The growing importance of CNVs: new insights for detection and clinical interpretation. *Front. Genet.*, **4**, 1–19.
- Venkatraman, E.S. and Olshen, A.B. (2007) A faster circular binary segmentation algorithm for the analysis of array CGH data. *Bioinformatics*, **23**, 657–663.
- Vosse, S. and van de Wiel, M. (2009) *CGHregions: Dimension Reduction for Array CGH Data with Minimal Information Loss*. R package version 1.7.1.
- Wu, Y. *et al.* (2013) MATCHCLIP: locate precise breakpoints for copy number variation using CIGAR string by matching soft clipped reads. *Front. Genet.*, **4**, 1–7.
- Zhao, M. *et al.* (2013) Computational tools for copy number variation (CNV) detection using next-generation sequencing data: features and perspectives. *BMC Bioinformatics*, **14**, S1.
- Zheng, C. *et al.* (2013) Determination of genomic copy number alteration emphasizing a restriction site-based strategy of genome re-sequencing. *Bioinformatics*, **29**, 2813–2821.