

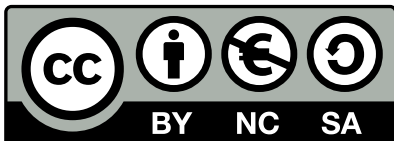
# R, paralelización, datos masivos y aplicaciones web: ejemplos del uso de R en bioinformática

Ramón Díaz-Uriarte

Dept. Bioquímica  
Universidad Autónoma de Madrid  
Madrid, Spain  
rdiaz02@gmail.com  
<http://ligarto.org/rdiaz>

III Jornadas de Usuarios de R  
17-Noviembre-2011

## License and copyright



This work is Copyright, ©, 2011, Ramón Díaz-Uriarte, and is licensed under the **Creative Commons Attribution-NonCommercial-ShareAlike** License. To view a copy of this license, visit

<http://creativecommons.org/licenses/by-nc-sa/3.0/> or send a letter to Creative Commons, 559 Nathan Abbott Way, Stanford, California 94305, USA.

\*\*\*\*\*

Please, **respect the copyright**. This material is provided freely, and if you use it, I only ask that you use it according to the (very permissive) terms of the license: attribution, non-commercial use, and a share alike license. If you have any doubts, ask me.

# Outline

Context

Context

Parallelizing  
code

Parallelizing code

Web applications

Web applications

Large data sets  
and  
parallelization

Large data sets and parallelization

R, C, and  
compression on  
the fly

R, C, and compression on the fly

Conclusions,  
future, et al.

Conclusions, future, et al.

Appendix

Appendix

Context

Biological context

Computational context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

## Context

Biological context

Computational context

Parallelizing code

Web applications

Large data sets and parallelization

R, C, and compression on the fly

Conclusions, future, et al.

Appendix

# Chromosomes

Context

Biological context

Computational context

Parallelizing  
code

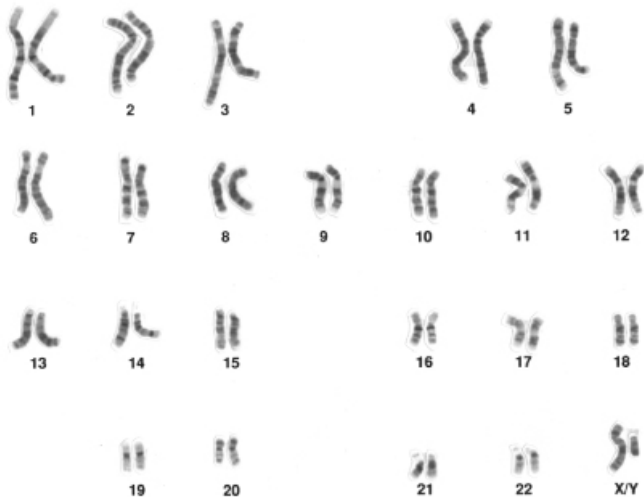
Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

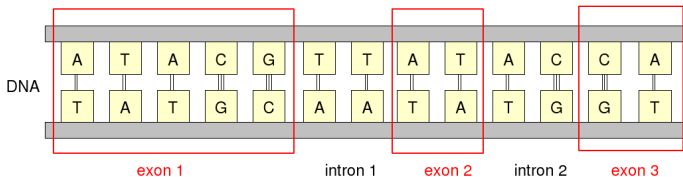
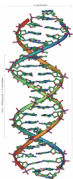
Appendix



From the Wikipedia; original source

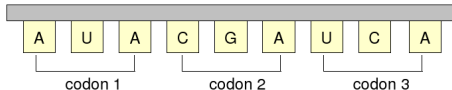
<http://www.genome.gov/Pages/Hyperion//DIR/VIP/Glossary/Illustration/karyotype.shtml>

# DNA → protein



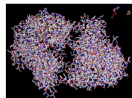
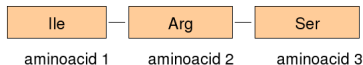
TRANSCRIPTION

RNA



TRANSLATION

PROTEIN



[http://en.wikipedia.org/wiki/Image:Hexokinase\\_ball\\_and\\_stick\\_model%2C\\_with\\_substrates\\_to\\_scale\\_copy.png](http://en.wikipedia.org/wiki/Image:Hexokinase_ball_and_stick_model%2C_with_substrates_to_scale_copy.png)

(From O. Rueda's PhD Thesis)

# Data from a microarray experiment

DATOS

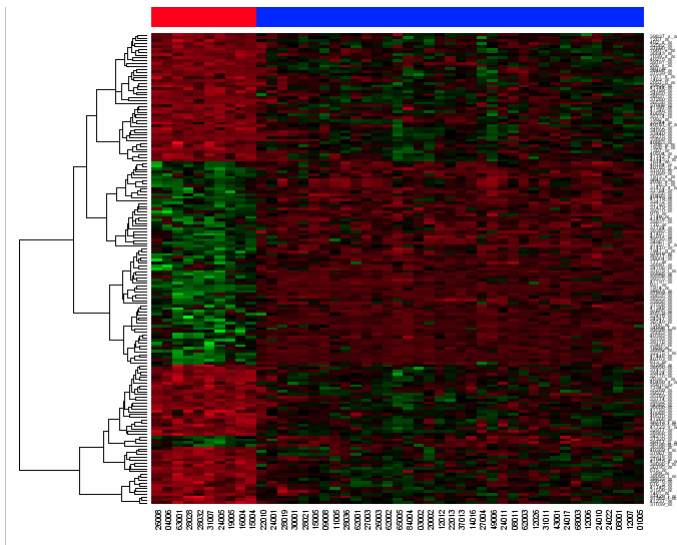
MUESTRAS /CASOS

**GENES**

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1	GENES	CPV													
2	BLZFP	0.06	0.3	0.02	-0.31	-0.29	0.04	0.83	-0.12	1.74	0.12	0.7	0.12	0.05	0.85
3	BLZFP	-0.51	-0.21	0.73	0.24	0.05	0.5	0.9	-0.12	0.11	-0.06	0.3	-0.66	0.15	0.72
4	NP	0.4	-0.62	-0.37	-0.20	-0.47	4.2	-0.39	-0.34	-0.94	-0.27	-0.74	-0.4	-0.6	-0.4
5	ARM	0.27	2.65	0.74	0.59	0.06	0.86	-0.42	0.73	-0.36	1.26	0.23	0.39	0.58	-0.01
6	ARM	0.42	1.95	0.49	0.23	-0.13	0.59	4.4	0.79	-0.8	1.6	-0.28	-0.12	0.03	-0.2
7	ARGALT	-1.33	0.19	3.33	0.14	0.33	0.4	0.14	0.24	-0.51	0.36	0.52	0.36	0.33	0.41
8	AACS	1.42	-0.05	0.04	0.86	0.57	0.14	0.26	0.21	0.24	0.22	0.27	-0.14	1.06	0.15
9	ARAT	0.52	-0.03	-0.46	0.19	0.47	-0.45	0.16	-0.06	0.42	0.3	0.31	-0.46	0.16	0.19
10	ARAT	-1.23	-0.07	0.3	0.07	0.29	0.36	0.37	0	0.38	0.02	0.49	0.25	0.19	0.41
11	ASBPPT	-0.73	-0.39	-0.43	-0.41	-1.09	-0.64	-1.07	0.27	4.30	0.24	0.14	1.18	-0.57	1.14
12	ARCA1	0.46	1.16	1.19	0.49	1.2	1.46	0.36	1.24	0.96	0.71	0.66	1.07	1.33	0.33
13	ARCA1	-0.40	1.15	1.42	0.42	1.14	1.13	0.34	1.11	0.56	0.64	0.7	0.74	1.3	0.31
14	ARCA12	-0.23	-0.19	1.5	-0.29	-0.45	1.3	-0.37	1.13	2.96	0	0.66	-1.05	1.22	-0.7
15	ARCA1	-0.5	1.62	-0.09	0.37	-0.16	0.27	-0.09	-0.18	-0.19	0.47	0.14	0.04	0.56	-0.18
16	ARCA1	-1.13	0.65	0.46	0.16	0.91	-0.03	-0.11	-0.11	0.22	0.03	0.31	0.16	0.63	0.29
17	ARCA6	-0.93	-0.65	-0.46	-0.36	-0.3	-0.3	-0.17	-0.94	0.44	-0.74	0.44	-0.02	-0.04	0.73
18	ARCA7	-0.38	-0.19	-0.4	-0.02	0.26	-0.46	-0.57	-0.51	0.22	-0.17	-0.39	-0.36	-0.16	0.05
19	ARCA1	0.53	-1.31	-0.69	-1.1	-0.57	4.6	-1.29	-0.82	-1.11	-0.67	-0.7	-0.36	-1.03	-1.54
20	ARCA12	0.71	-0.03	1.22	1.64	0.39	0.8	1.03	0.96	1.22	0.89	-0.4	-1.23	0.48	0.58
21	ARCA2	-0.30	-1.47	-1.67	-1.6	-1.61	-1.96	-1.74	-1.13	-1.66	-1.22	-1.67	-1.41	-1.66	-2.26
22	ARCA3	-0.42	-0.14	0.2	-0.16	0.41	0.26	-0.42	-0.59	-0.07	-0.76	-0.42	0	-0.59	0
23	ARCA9	-0.73	2.08	0.29	0.18	0.02	0.54	-0.32	0.57	0.81	0.31	0.37	0.82	0.51	0.1
24	ARCA3	1.56	0.24	-0.21	1.11	1.49	-0.06	2.83	1.2	0.54	0.89	0.06	0.34	0.43	0.89
25	ARCA1	1.09	-1.05	0.09	-0.23	-0.16	0.07	-0.27	-0.91	-0.36	-0.63	-0.39	-0.62	-0.63	0.26
26	ARCA7	0.31	-0.35	-1.5	-1.06	-1.21	-1.42	-1.31	-0.46	-0.79	-0.46	-1.36	-0.26	-1.35	-0.97
27	ARCA1	0.47	0.9	-0.21	0.14	0.62	-0.06	0.61	-0.52	0.7	0.52	0.59	0.69	1.02	-0.19
28	ARCA2	-1.68	-0.09	-0.48	-0.28	-0.39	-0.49	-0.46	-0.45	-0.14	0.14	-0.27	-0.12	0.55	-0.63
29	ARCA2	0.46	-0.23	-0.87	0.42	-0.15	-0.73	0.97	0.3	-0.24	0.46	-0.52	-1.11	-0.95	0.66
30	ARCA2	0.46	1.29	1.61	0.83	1.25	1.79	2.06	1.22	0.32	0.41	1.63	1.63	0	0.64
31	ARBP	-1.52	-0.43	0.37	0.4	-0.11	0.17	0.12	0.06	0.7	-0.4	0.69	0.64	0.18	0.23
32	ARCA1	-0.16	-0.03	0.12	-0.46	-0.18	0.26	-0.91	-0.36	-1	0.24	0.16	0.63	-0.62	-0.1
33	ARCA1	0.68	-0.26	-0.21	-0.26	-0.16	-0.13	-0.53	-0.11	-0.96	0.19	-0.36	-0.66	-0.8	-0.28
34	ARCA1	2.6	0.42	1.37	0.7	-1.37	1.19	1.37	0.52	-1.36	0.97	0.54	0.1	0.69	-0.62
35	ARCA1	-0.69	-0.69	-0.29	-0.3	-0.32	-0.34	-0.20	-1.03	-0.09	-1.36	-0.35	-0.05	-0.08	-0.2
36	ARCA2	-0.61	-0.66	0.06	-0.46	-0.39	0.06	-0.96	-0.36	-0.43	-0.72	-0.38	0.32	-0.38	0.64
37	ARCA1	-2.22	0.4	0.66	0	0.71	0.36	-0.07	0.43	0.36	0.44	0.16	0.33	0.53	0.31
38	ARCA1	1.09	-0.16	-0.27	-0.44	0.10	-0.20	-0.24	-0.26	-0.36	0.31	-0.73	0.00	-0.53	-0.76
39	ARCA1	-0.6	0.07	-0.09	0.26	0.1	-0.09	1.20	0.5	0.06	-0.13	-0.3	-0.62	0.09	0.26
40	ARCA1	1.37	0.38	-0.48	0.44	0.47	-0.46	0.44	0.86	-0.28	0.53	-0.2	-0.05	0.6	-0.17
41	ARCA2	0.6	-0.06	-0.46	-0.78	-0.71	-0.6	-0.41	-0.87	-1.36	0.52	-1.22	-0.17	-1.03	-1.11
42	ARCA1	0.31	1.15	0.2	2.89	0.81	0.46	0.61	0.4	0.21	0.41	0.2	-0.16	-0.09	0.69
43	ARCA1	0.72	0.09	0.45	-0.59	-0.5	0.96	0.16	0.11	0.01	0.1	-0.69	-0.71	-0.69	-0.2
44	ARCA1	0.02	0.03	-1.46	-1.13	-1.03	-1.1	-1.16	-0.52	-1.16	-0.96	-1.44	-1.46	-0.61	-2.31
45	ARCA1	-0.16	0.37	0.27	0.5	0.06	0.11	-0.69	0.36	0.56	0.42	1.23	0.22	0.41	-0.38

Slide from Gema Moreno Bueno, Department of Biochemistry, UAM

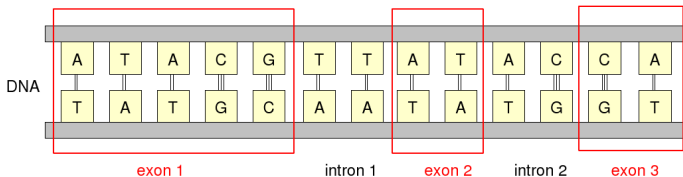
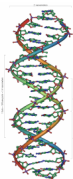
# More microarray data



Modified from [http://www2.warwick.ac.uk/fac/sci/moac/students/peter\\_cock/r/heatmap/scaled\\_color\\_key.png](http://www2.warwick.ac.uk/fac/sci/moac/students/peter_cock/r/heatmap/scaled_color_key.png)

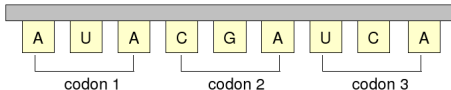


# DNA → protein



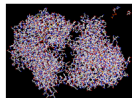
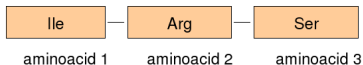
TRANSCRIPTION

RNA



TRANSLATION

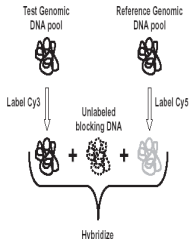
PROTEIN



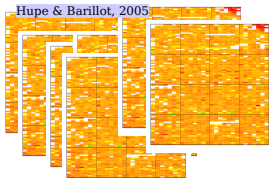
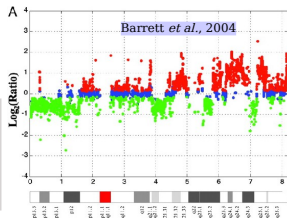
[http://en.wikipedia.org/wiki/Image:Hexokinase\\_ball\\_and\\_stick\\_model%2C\\_with\\_substrates\\_to\\_scale\\_copy.png](http://en.wikipedia.org/wiki/Image:Hexokinase_ball_and_stick_model%2C_with_substrates_to_scale_copy.png)

(From O. Rueda's PhD Thesis)

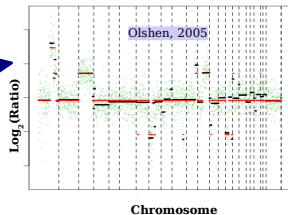
J. Fridlyand et al. / Journal of Multivariate Analysis 90 (2004) 132–153



Calling gains and losses: hypothesis testing



Inferring number of copy gains/losses: estimation



Arrays: a dot is a DNA fragment.  
Each array a sample. Each array all chromosomes.  
(For analysis, location in chromosome matters)

# Data, data, data (in Gigabytes)

## Context

Biological context

Computational context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

Expression arrays (mRNA) > 40,000 probes

Copy number with aCGH > 400,000 common;

some >  $4 \times 10^6$

... ..

# Multicores and computing clusters

- Increases in CPU speed slowed down (< 20% per year since 2002).
- Increase in the number of “cores”: 2, 4, 8. Next 10 years?
- Inexpensive computing clusters with off-the-shelf components.
- Must design our programs from the start: **parallel programming**

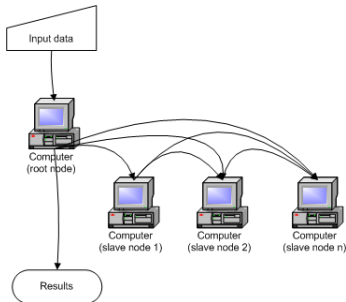


Image from <http://faq.distributed.net/>

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

Context

**Parallelizing code**

Web applications

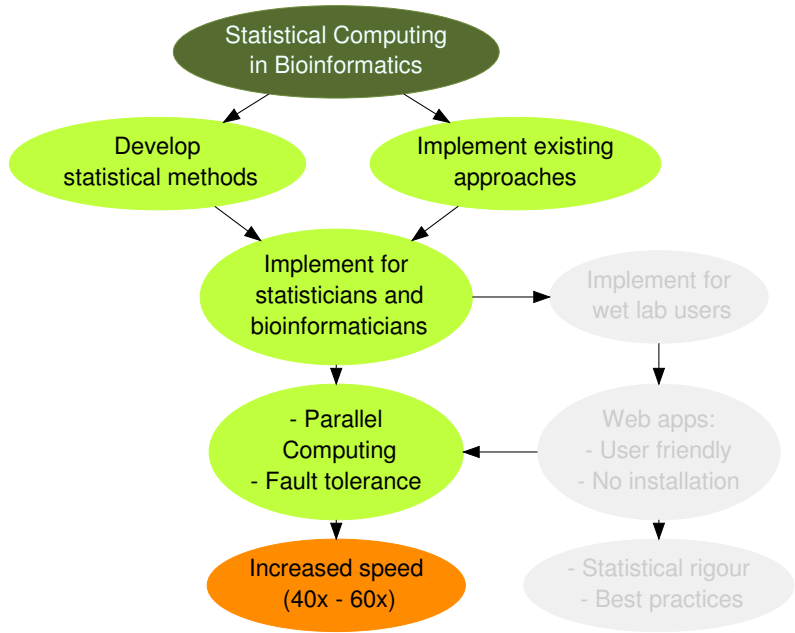
Large data sets and parallelization

R, C, and compression on the fly

Conclusions, future, et al.

Appendix

# Standalone



- Context
- Parallelizing code
- Web applications
- Large data sets and parallelization
- R, C, and compression on the fly
- Conclusions, future, et al.
- Appendix

# R code

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- Code available for many procedures (but a few years ago none parallelized!)
- Many computations embarrassingly parallelizable:
  - ▶ bootstrapping and cross-validation
  - ▶ arrays (or samples)
  - ▶ arrays by chromosomes
  - ▶ parallel chains in MCMC
  
- Figure production can be parallelized

# Parallelizing R code

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- (Implement missing functionality: R/C)
- MPI: R packages Rmpi, papply, snow, snowfall
- Load balanced
- Wrappers over “mid level” functions in package: ease updating
- Parallelize:
  - ▶ Bootstrap samples/Cross-val. runs.
  - ▶ arrays
  - ▶ arrays by chromosomes
  - ▶ (or a combination of both)



# Is it worth it?

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- Are speed improvements really worth the effort?
- Over what range of problems do we see improvements?
- With what hardware can we see improvements?

# What do we gain?

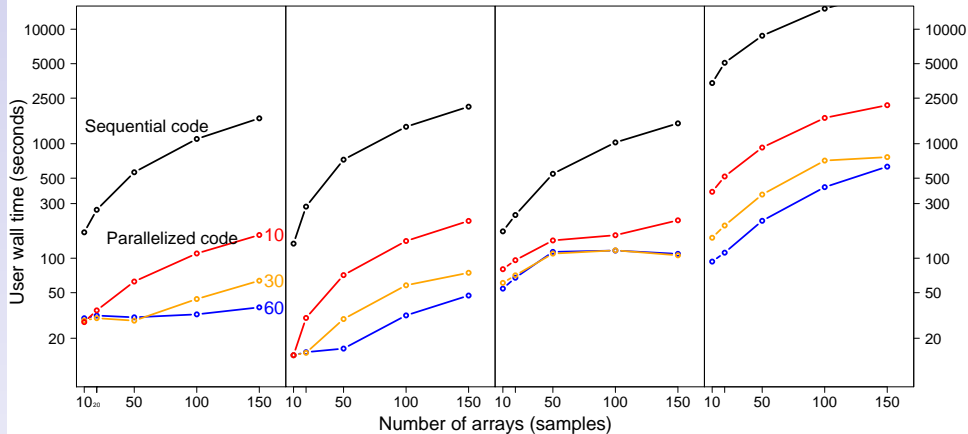
20,000 genes

HMM

GLAD

CBS

BioHMM



# What do we gain?

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

Are speed improvements really worth the effort?

Your effort: “R CMD INSTALL ADaCGH2”.

Over what range of problems do we see improvements?

10 to  $10^3$  arrays/samples;

$10^4$  to  $10^6$  spots/genes.

With what hardware can we see improvements?

2 cores to 120 cores.

- Smaller clusters: more cost effective
- Single node/multi-core: less communication overhead

# Where is this running?

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- *varSelRF* (CRAN)
- *ADaCGH2* (BioConductor)
- *SignS* (launchpad:  
<http://launchpad.net/signs>)

Context

Parallelizing  
code

**Web applications**

Web apps: how

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

Context

Parallelizing code

**Web applications**

**Web apps: how**

Large data sets and parallelization

R, C, and compression on the fly

Conclusions, future, et al.

Appendix

# Applications for wet lab researchers

Context

Parallelizing  
code

**Web applications**

Web apps: how

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- Analyze data in a reasonably short time.
- User friendly access to methods that are statistically rigorous.

# Web-based applications

- User-friendly interface.
- No hardware/software hassles for end users.
- Parallelization is transparent.
- Method selection can be partially transferred (to us).
- Short user wall time: use (hardware/software) resources rarely available to individual biomedical researchers
- **Just type in a URL:**  
`http://www.some-application`

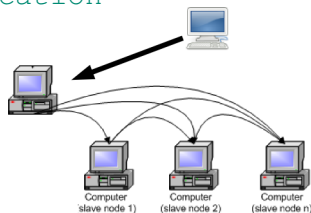


Image modified from <http://faq.distributed.net/>

# Sometimes collaborations feel like ...

Context

Parallelizing  
code

Web applications

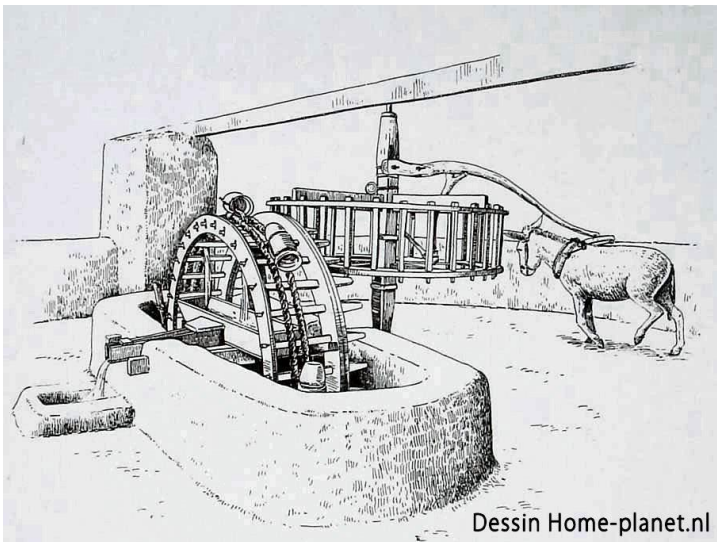
Web apps: how

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

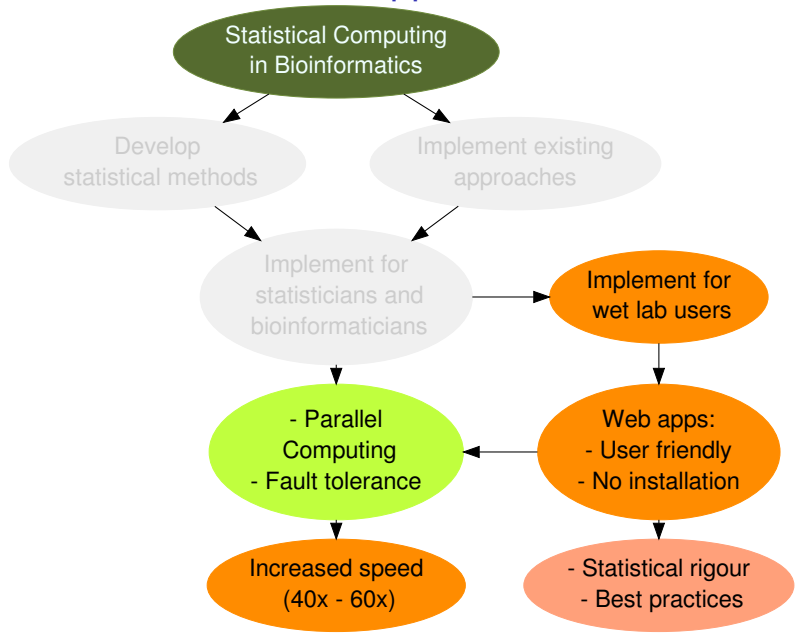
Appendix



(From <http://www.bitacoradegalileo.com/2010/11/16/giordano-bruno-en-la-cara-oculta-de-la-luna/>)

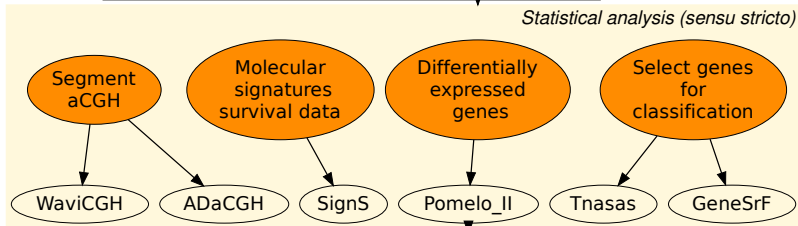
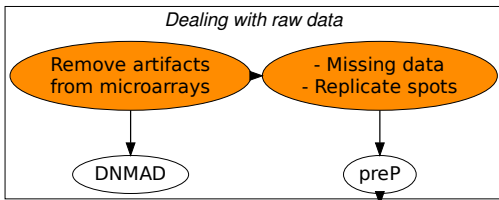


# Parallelization in web applications

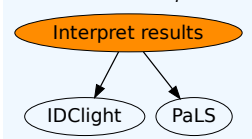


- Context
- Parallelizing code
- Web applications**
  - Web apps: how
- Large data sets and parallelization
- R, C, and compression on the fly
- Conclusions, future, et al.
- Appendix

# Main web-based applications



*Annotation and Interpretation*



Context

Parallelizing  
code

Web applications

Web apps: how

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

# What do we gain?

Context

Parallelizing  
code

Web applications

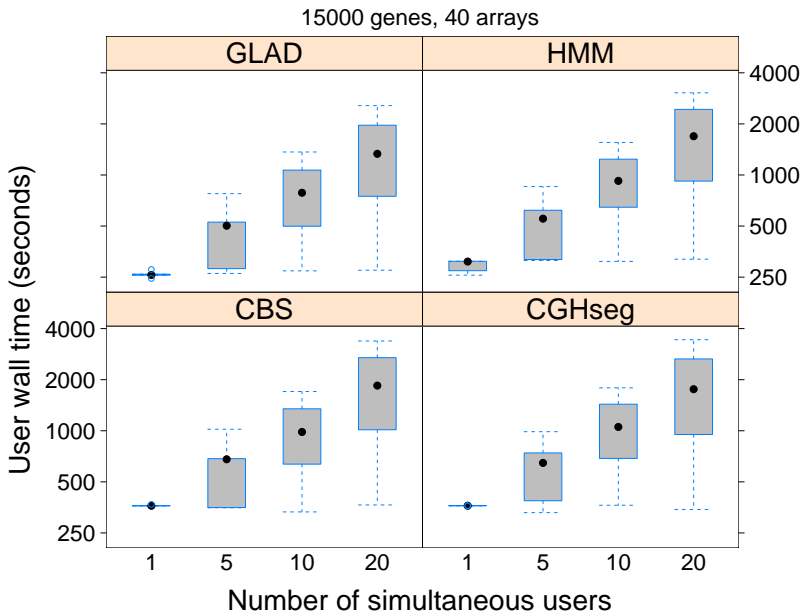
Web apps: how

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix



# How it works: some key ideas

Context

Parallelizing  
code

Web applications  
Web apps: how

Large data sets  
and  
parallelization

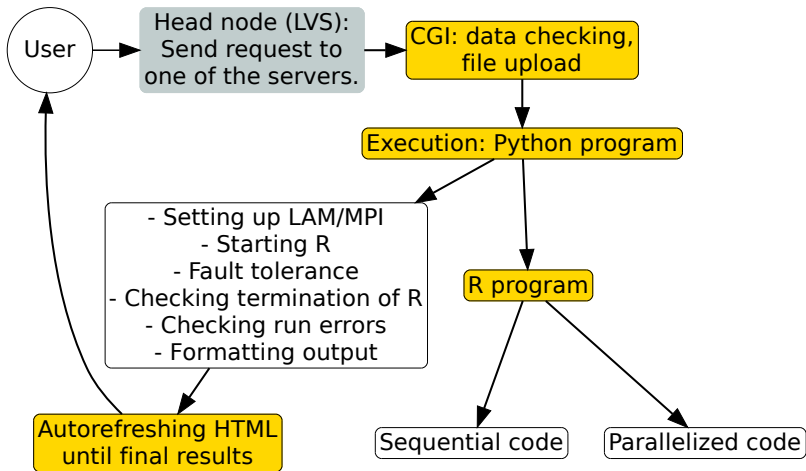
R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

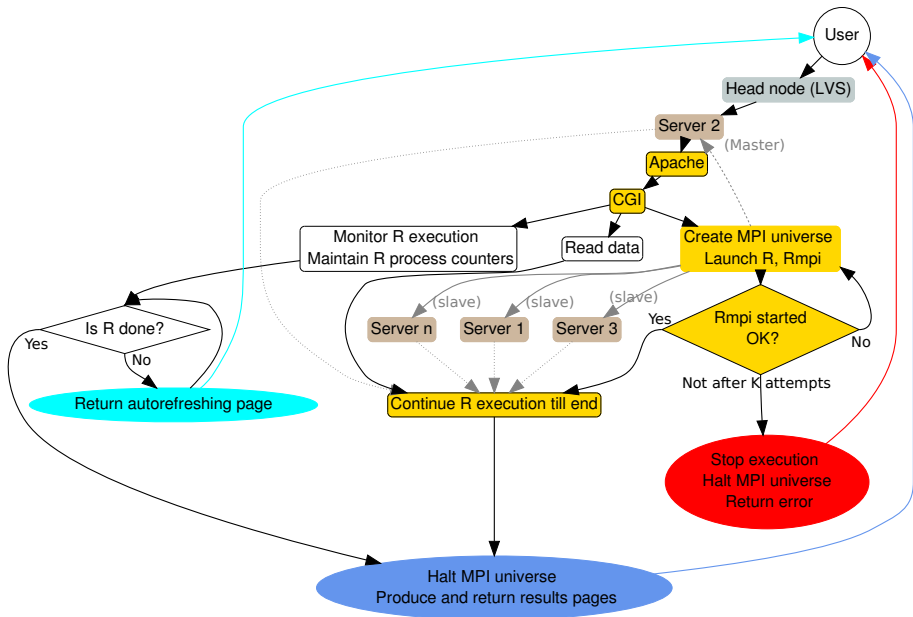
Appendix

- Each run
  - ▶ Parallelization (transparent for users)
  - ▶ Fault-tolerance (network problems, machine crashes, bugs)
    - ▶ Check-pointing
  
- Periodic tasks (keep system running 24h, 365 d)
  - ▶ Automatic monitorization
  - ▶ Automated testing suite

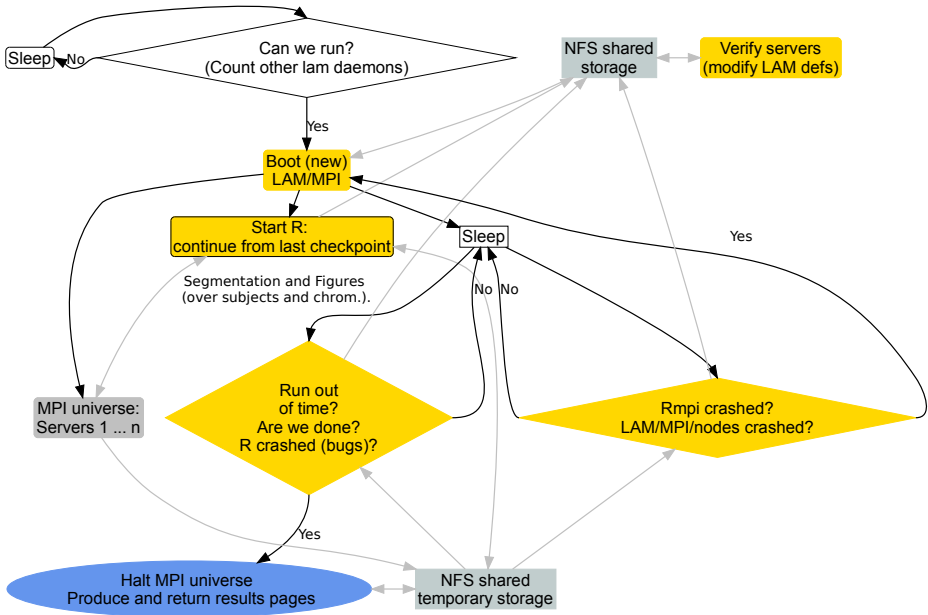
# What happens



# What happens: details



# MPI details



# Where is this running?

Context

Parallelizing  
code

Web applications

Web apps: how

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- <http://signs.bioinfo.cnio.es>
- <http://wavi.bioinfo.cnio.es>
- <http://genesrf.bioinfo.cnio.es>



Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

Context

Parallelizing code

Web applications

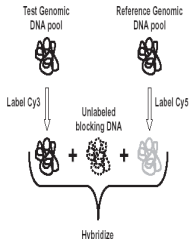
**Large data sets and parallelization**

R, C, and compression on the fly

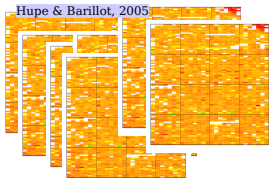
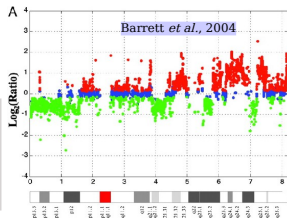
Conclusions, future, et al.

Appendix

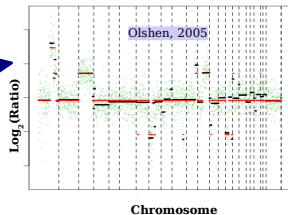
*J. Fridlyand et al. / Journal of Multivariate Analysis 90 (2004) 132–153*



Calling gains and losses: hypothesis testing



Inferring number of copy gains/losses: estimation



Arrays: a dot is a DNA fragment.  
Each array a sample. Each array all chromosomes.  
(For analysis, location in chromosome matters)

# Large data sets

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- Millions of spots
- Hundreds or thousands of subjects.
- No need to hold everything in RAM at once.
- Package *ff*: “memory-efficient storage of large data on disk and fast access functions.”
- Combined with:
  - ▶ parallelization
  - ▶ shared storage

# *ff* and parallelization

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

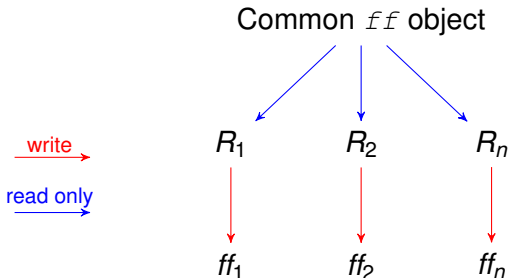
R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- *ff* stores the object on disk.
- Read that object from various R processes.
- Different R processes can write in different *ff* objects

# $ff$ and parallelization (I)



Context

Parallelizing  
code

Web applications

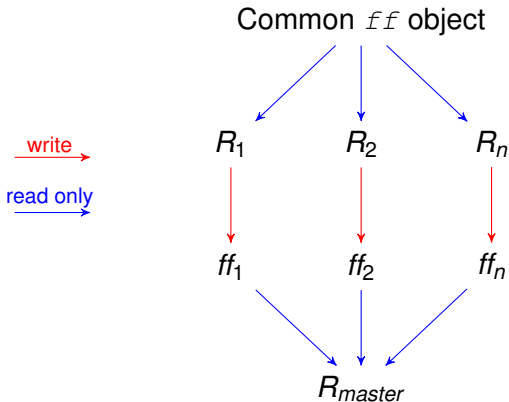
Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

# $ff$ and parallelization (I)



Context

Parallelizing  
code

Web applications

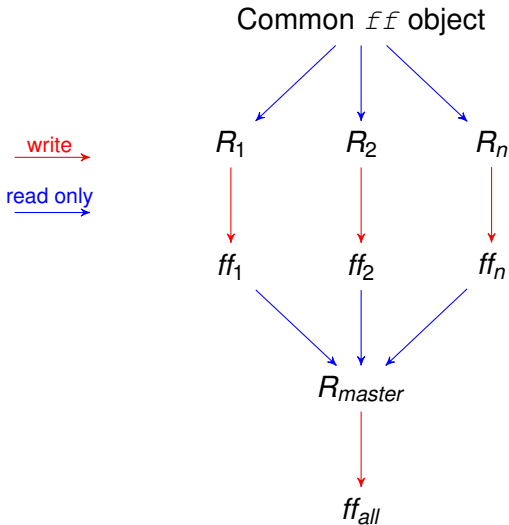
Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

# $ff$ and parallelization (I)



Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

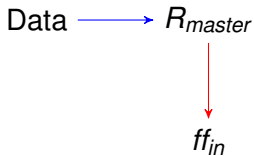
R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

# *ff* and parallelization (II)

write →  
read only →



Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

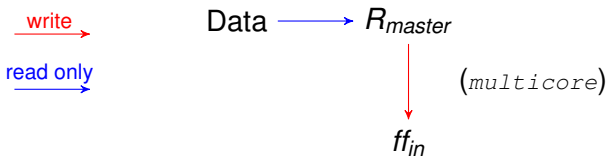
R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix



# *ff* and parallelization (II)



Context

Parallelizing  
code

Web applications

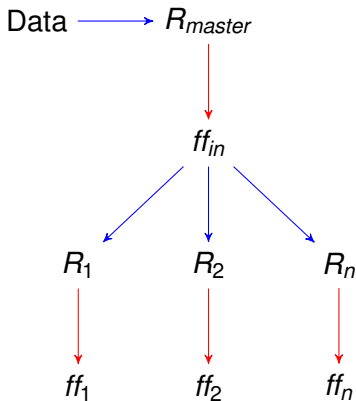
Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

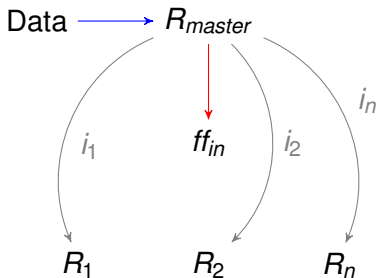
Appendix

# $ff$ and parallelization (II)



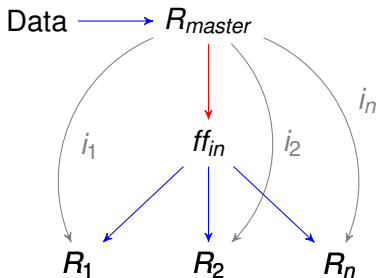
## $ff$ and parallelization (II)

write →  
read only →



## $ff$ and parallelization (II)

write →  
read only →



Context

Parallellizing  
code

Web applications

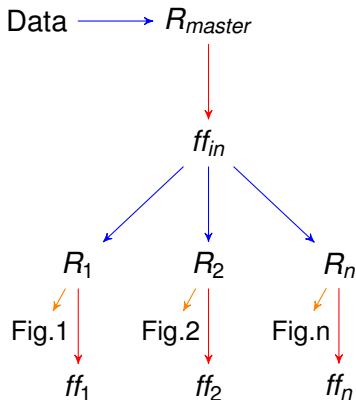
Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

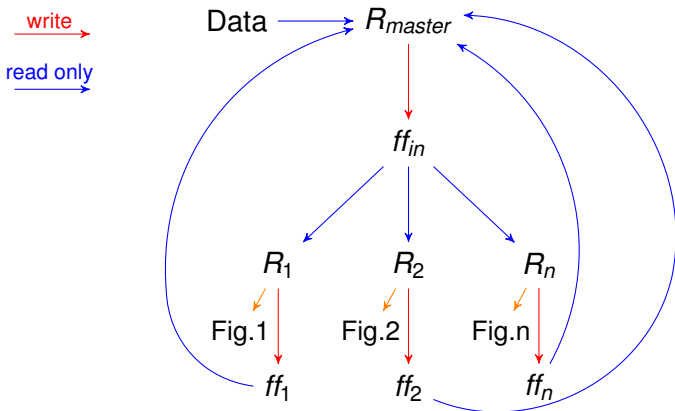
Conclusions,  
future, et al.

Appendix

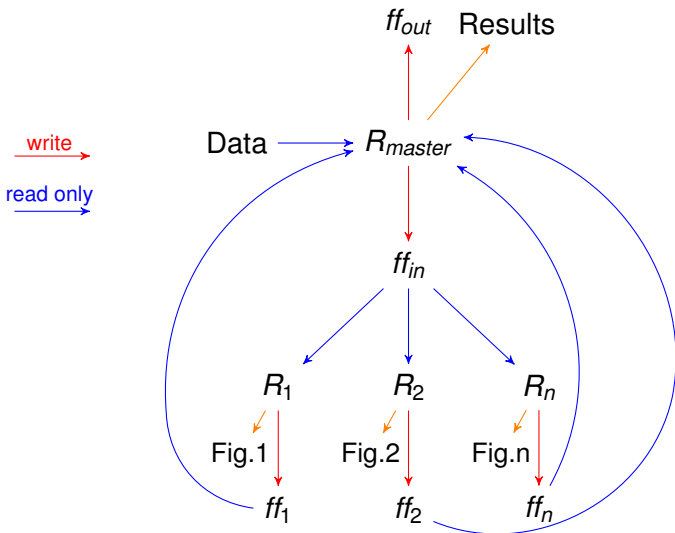
# $ff$ and parallelization (II)



## $ff$ and parallelization (II)



## $ff$ and parallelization (II)



Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

# Where is this running?

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- *ADaCGH2* (BioConductor package)
- **Web-based application**  
<http://wavi.bioinfo.cnio.es>.



Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

**R, C, and  
compression on  
the fly**

Conclusions,  
future, et al.

Appendix

Context

Parallelizing code

Web applications

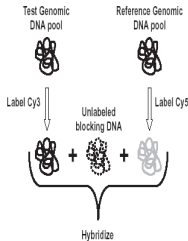
Large data sets and parallelization

**R, C, and compression on the fly**

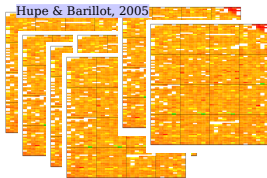
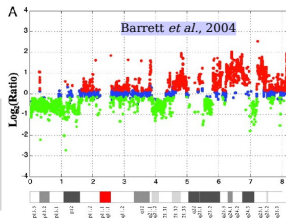
Conclusions, future, et al.

Appendix

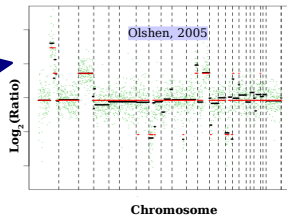
*J. Fridlyand et al. / Journal of Multivariate Analysis 90 (2004) 132–153*



Calling gains and losses: hypothesis testing



Inferring number of copy gains/losses: estimation



Arrays: a dot is a DNA fragment.  
Each array a sample. Each array all chromosomes.  
(For analysis, location in chromosome matters)

# Store and access (large) pre-computed results

- HMM for aCGH data with Reversible Jump: Viterbi
- Common regions: “count” on the Viterbi paths.
  
- Fitting HMM/common regions: distinct operations.
  
- C: number-crunching.
- R: wrapper and figures/tables.
- C: creates large amounts of data.
  
- In package *RJaCGH* (CRAN).

Context

Parallelizing  
code

Web applications

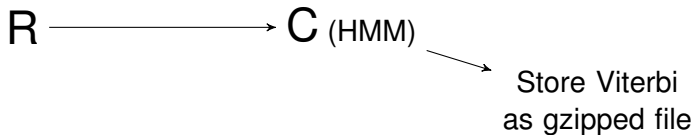
Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

## Fit HMM



Context

Parallelizing  
code

Web applications

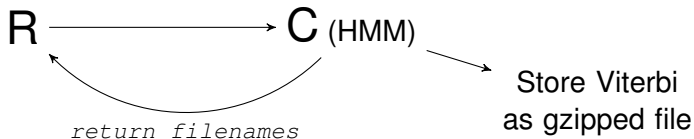
Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

## Fit HMM



Context

Parallelizing  
code

Web applications

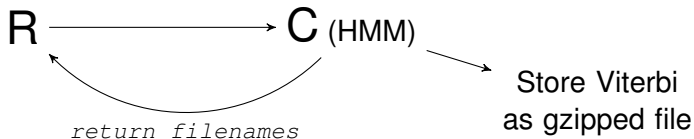
Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

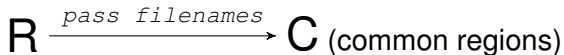
Conclusions,  
future, et al.

Appendix

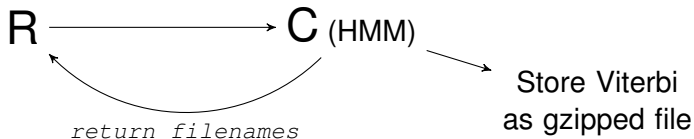
## Fit HMM



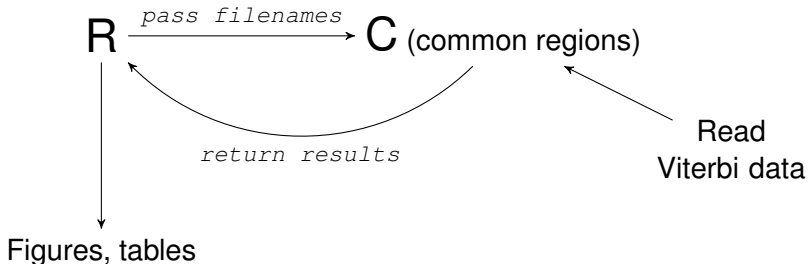
## Find common regions



## Fit HMM



## Find common regions



Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

Context

Parallelizing code

Web applications

Large data sets and parallelization

R, C, and compression on the fly

**Conclusions, future, et al.**

Appendix



# Web-based: A few things we've learned

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

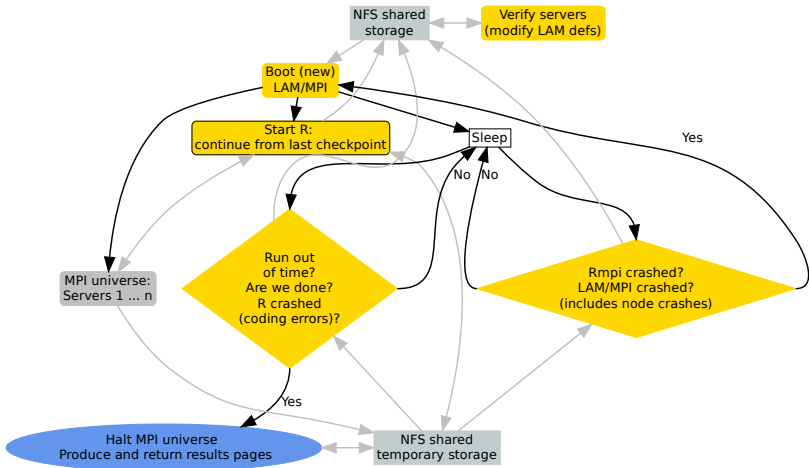
Conclusions,  
future, et al.

Appendix

- Configuration sucks (if you need to modify  $> 1$  file)
- Too many languages
- Adding test cases to the testing suites: web, R
- Documentation: in the code, web pages,  $\text{\LaTeX}$  ...
- Too much R code to catch errors
- User interfaces: who designs them?

# Fault tolerance and communication

- Manual check for errors (R ain't Erlang)
- Too much network traffic



# Solutions?

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- Literate programming and org-mode
- Alternatives to MPI and/or use Erlang
- . . .
- Keep things as they are (only a few painful events a year)

# Single machine applications

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- Run applications in a single, multicore (e.g., 12) machine
- Just verify if the machine is up

# Rethinking web-based applications

- Users can get into trouble.

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

# Rethinking web-based applications

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- Users can get into trouble.

Sure, but we can do a good job . . .

- ▶ Provide state-of-the art statistical and computational approaches (R ;-)
- ▶ Pedagogical examples and pipelines
- ▶ Minimize the chance of users getting into trouble

# Rethinking web-based applications

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- Users can get into trouble.

Sure, but we can do a good job . . .

- ▶ Provide state-of-the art statistical and computational approaches (R ;-)
- ▶ Pedagogical examples and pipelines
- ▶ Minimize the chance of users getting into trouble

- Web-based applications are here to stay

## ... SO ...

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- Forget about them: just write R (plus C, etc) code
- Go for it
  - ▶ R will do its job (which is only part of the job)
  - ▶ HPC available
  - ▶ No problem with large data sets
  - ▶ But other tools and work necessary



# Regardless of web-based applications ...

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

- Parallel computing can be used routinely
  - ▶ (`library(parallel)` in R  $\geq$  2.14.0)
- Large data sets with *ff* + parallelization.

# Acknowledgements

- O. M. Rueda, A. Alibés, A. Cañada, E. R. Morrissey, M. L. Neves, D. Rico.
- Funding: Fundación de Investigación Médica Mutua Madrileña, Project TIC2003-09331-C02-02 of the Spanish MEC and BIO2009-12458 of the Spanish MICINN. Ramón y Cajal Programme of the Spanish Ministry of Education and Science.
- CNIO (Spanish National Cancer Research Center).
- The R users and developers for a vibrant statistical computing community and amazing platform.
- The organizing and scientific committees of the III Jornadas de Usuarios de R.

Context

Parallelizing  
code

Web applications

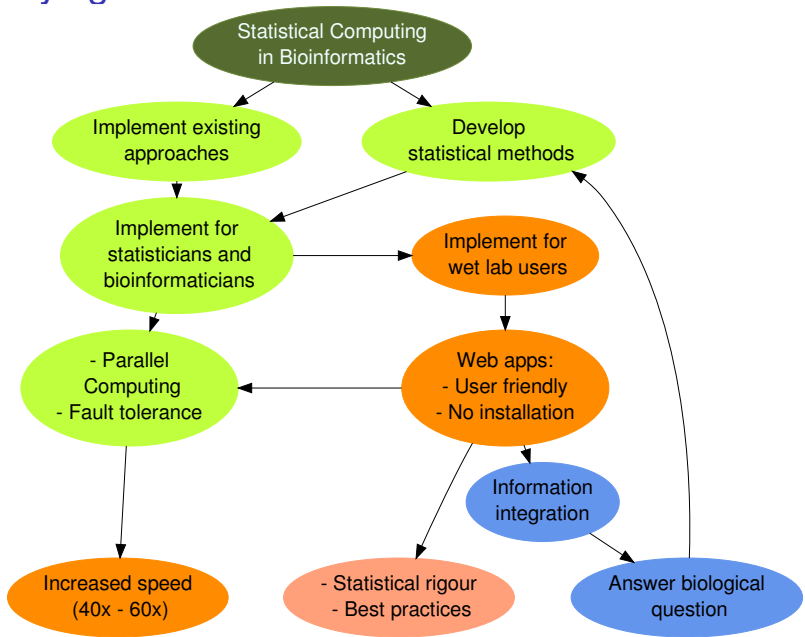
Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

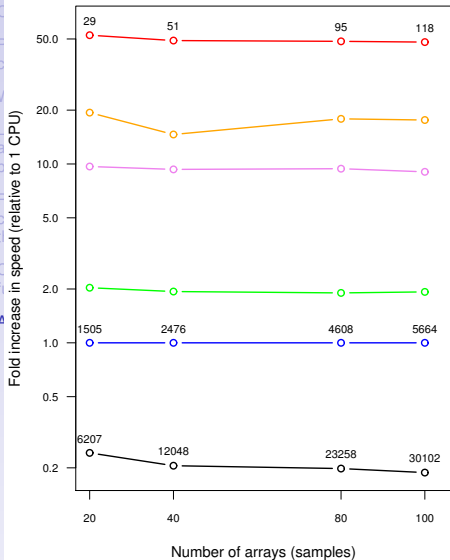
# Trying to fit it all



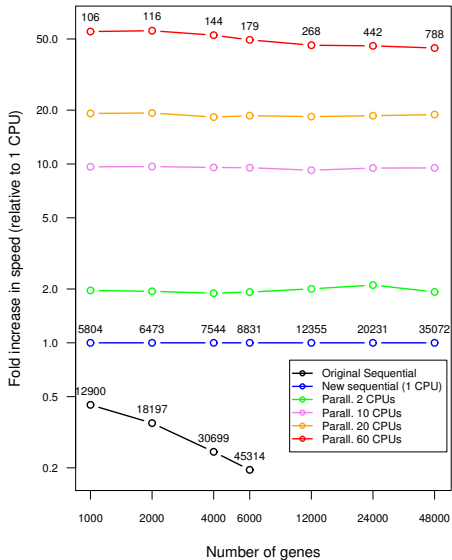
Context  
Parallelizing code  
Web applications  
Large data sets and parallelization  
R, C, and compression on the fly  
Conclusions, future, et al.  
Appendix

# What do we gain?

Effect of number of arrays (number of genes = 7399)

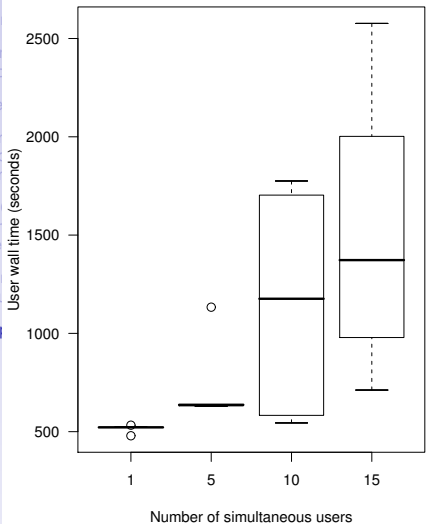


Effect of number of genes (number of arrays = 160)

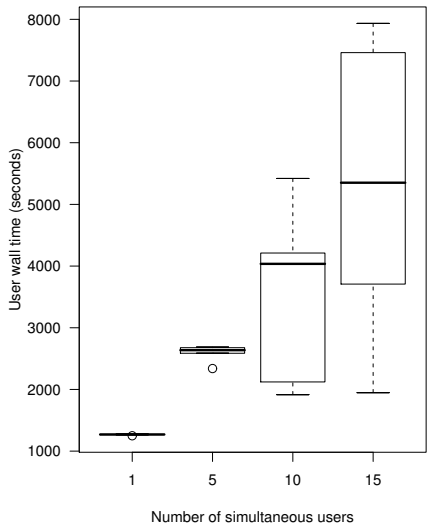


# What do we gain?

### Breast data set (78 arrays x 4751 genes)



### DLBCL data set (160 arrays x 7399 genes)



# Tools...

Context

Parallelizing  
code

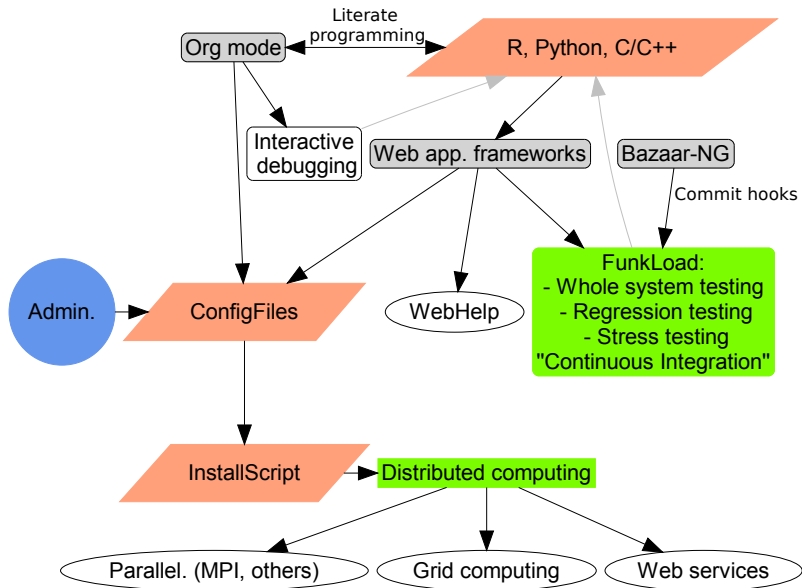
Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix



# Too many languages

Impedance mismatch problem:

“Building Web-based applications requires the mastering of a number of languages/technologies (e.g. HTML, CSS, CGI, ASP, PHP, XML, etc..). Such languages and technologies were created to address different aspects on a by-need evolutionary manner. The result is a plethora of tools that are fitted together in an ad hoc fashion.” El-Ansary, Grolaux, Van Roy, Rafea (2005) *“Overcoming the Multiplicity of Languages and Technologies for Web-Based Development Using a Multi-paradigm Approach”*.

- R and C
- HTML and Python: CGI, data entry, display
- Python (and others): control and monitor MPI
- Javascript: AJAX and figures

# Other solutions?

Context

Parallelizing  
code

Web applications

Large data sets  
and  
parallelization

R, C, and  
compression on  
the fly

Conclusions,  
future, et al.

Appendix

**Too many languages** Use languages designed to overcome this problem: Hop, Links, QHTML.

**Fault tolerance and too much traffic** Alternatives to MPI?

- Linda and tuple spaces (also between-language funct.)
- Roll-our-own based on Rserve
- Have Erlang control R processes?