

# Software & Requirements

---

Javier Perales-Patón  
[jperales@cniio.es](mailto:jperales@cniio.es)



Bioinformatics Unit  
CNIO. Madrid, Spain.

Fátima Al-Shahrour  
[\[falshahrour@cniio.es\]](mailto:falshahrour@cniio.es)  
Elena Piñeiro-Yáñez  
[\[epineiro@cniio.es\]](mailto:epineiro@cniio.es)  
Pedro Fernandes  
[\[pfern@igc.gulbenkian.pt\]](mailto:pfern@igc.gulbenkian.pt)



# Computational Facilities

- The **genome reference** must be **indexed** only once in your computer, but it takes a long time (2-3 days). The sizes of the indexes are huge (~100 Gb). See “Before running a variant analysis for first time” at:  
<http://rubioseq.bioinfo.cnio.es/sites/default/files/PDF/RUbioSeq-book.pdf>
- You need a huge **hard disk capacity**:
  - **Human genome reference** (including indexes) : 150 Gb.
  - **Bundle of files for Variant Calling Analysis** (GATK) : 14 Gb.
  - **For each sample**:
    - Raw data : 7 - 15 Gb.
    - Intermediate files for each sample : 15 - 20 Gb (whole-exome seq).

## WorkStation:

- **Minimal requirements**:
  - 16 Gb RAM
  - 500 Gb hard disk
  - 8 threads ~ 8 cores.
- **Recommended requirements (multi-sample analysis, samples storage)**:
  - 25 Gb RAM
  - 1 Tb hard disk (e.g. Exome-seq analysis from 50 samples → 1.5 Tb).
  - 16 threads.

# GATK's bundle from FTP repository (14 Gb)



http://ftp.broadinstitute.org/bundle/2.8/hg19/

## Index of /bundle/2.8/hg19/

Name	Size	Date Modified
[parent directory]		
1000G_omni2.5.hg19.sites.vcf.gz	49.4 MB	12/8/13, 1:00:00 AM
1000G_omni2.5.hg19.sites.vcf.gz.md5	97 B	12/8/13, 1:00:00 AM
1000G_omni2.5.hg19.sites.vcf.idx.gz	464 kB	12/8/13, 1:00:00 AM
1000G_omni2.5.hg19.sites.vcf.idx.gz.md5	101 B	12/8/13, 1:00:00 AM
1000G_phase1.indels.hg19.sites.vcf.gz	42.9 MB	12/8/13, 1:00:00 AM
1000G_phase1.indels.hg19.sites.vcf.gz.md5	103 B	12/8/13, 1:00:00 AM
1000G_phase1.indels.hg19.sites.vcf.idx.gz	326 kB	12/8/13, 1:00:00 AM
1000G_phase1.indels.hg19.sites.vcf.idx.gz.md5	107 B	12/8/13, 1:00:00 AM
1000G_phase1.snps.high_confidence.hg19.sites.vcf.gz	1.7 GB	12/8/13, 1:00:00 AM
1000G_phase1.snps.high_confidence.hg19.sites.vcf.gz.md5	117 B	12/8/13, 1:00:00 AM
1000G_phase1.snps.high_confidence.hg19.sites.vcf.idx.gz	3.4 MB	12/8/13, 1:00:00 AM
1000G_phase1.snps.high_confidence.hg19.sites.vcf.idx.gz.md5	121 B	12/8/13, 1:00:00 AM
CEUTrio.HiSeq.WGS.b37.bestPractices.hg19.vcf.gz	407 MB	12/8/13, 1:00:00 AM
CEUTrio.HiSeq.WGS.b37.bestPractices.hg19.vcf.gz.md5	119 B	12/8/13, 1:00:00 AM
CEUTrio.HiSeq.WGS.b37.bestPractices.hg19.vcf.idx.gz	3.2 MB	12/8/13, 1:00:00 AM
CEUTrio.HiSeq.WGS.b37.bestPractices.hg19.vcf.idx.gz.md5	123 B	12/8/13, 1:00:00 AM
Mills_and_1000G_gold_standard.indels.hg19.sites.vcf.gz	19.1 MB	12/8/13, 1:00:00 AM
Mills_and_1000G_gold_standard.indels.hg19.sites.vcf.gz.md5	120 B	12/8/13, 1:00:00 AM
Mills_and_1000G_gold_standard.indels.hg19.sites.vcf.idx.gz	426 kB	12/8/13, 1:00:00 AM
Mills_and_1000G_gold_standard.indels.hg19.sites.vcf.idx.gz.md5	124 B	12/8/13, 1:00:00 AM
dbsnp_138.hg19.excluding_sites_after_129.vcf.gz	334 MB	12/8/13, 1:00:00 AM
dbsnp_138.hg19.excluding_sites_after_129.vcf.gz.md5	119 B	12/8/13, 1:00:00 AM
dbsnp_138.hg19.excluding_sites_after_129.vcf.idx.gz	3.6 MB	12/8/13, 1:00:00 AM
dbsnp_138.hg19.excluding_sites_after_129.vcf.idx.gz.md5	123 B	12/8/13, 1:00:00 AM
dbsnp_138.hg19.vcf.gz	1.4 GB	12/8/13, 1:00:00 AM
dbsnp_138.hg19.vcf.gz.md5	93 B	12/8/13, 1:00:00 AM
dbsnp_138.hg19.vcf.idx.gz	3.8 MB	12/8/13, 1:00:00 AM
dbsnp_138.hg19.vcf.idx.gz.md5	97 B	12/8/13, 1:00:00 AM
hapmap_3.3.hg19.sites.vcf.gz	58.0 MB	12/8/13, 1:00:00 AM
hapmap_3.3.hg19.sites.vcf.gz.md5	94 B	12/8/13, 1:00:00 AM
hapmap_3.3.hg19.sites.vcf.idx.gz	807 kB	12/8/13, 1:00:00 AM
hapmap_3.3.hg19.sites.vcf.idx.gz.md5	98 B	12/8/13, 1:00:00 AM

You can download the different bundles from GATK's FTP (Broad Institute) visiting this URL with your Internet Browser:

<ftp://gsapubftp-anonymous@ftp.broadinstitute.org/bundle/2.8/>

1. **Genome Reference** (standard 1000 Genomes, fasta).
2. List of **Target beats or intervals** of genomic regions sequenced by the Library protocol.
3. **dbSNP** (VCF file) for a recent dbSNP release (build 138, it includes the 1000 Genomes).
4. HapMap genotypes and sites VCFs
5. **OMNI 2.5 genotypes for 1000 Genomes samples** (VCF).
6. The current best set of **known indels** to be used for local realignment); use both files:
  - 1000G\_phase1.indels.b37.vcf (currently from the 1000 Genomes Phase I indel calls)
  - Mills\_and\_1000G\_gold\_standard.indels.b37.sites.vcf

 Tip for home: The following UNIX command downloads the whole bundle for **hg19** in one step (~hrs) :

```
$ wget -r -nH --cut-dirs=2 --reject-regex "NA12878|CEUTrio" \  
-P /path/to/your_directory/ ftp://gsapubftp-anonymous@ftp.broadinstitute.org/bundle/hg19/
```

# Installing RUBioSeq+ in your computer



<http://rubioseq.bioinfo.cnio.es/>

## Download and Run RUBioSeq+

RUBioSeq+ LiveDVD

RUBioSeq+ sources

RUBioSeq+ Docker



LINUX



LINUX

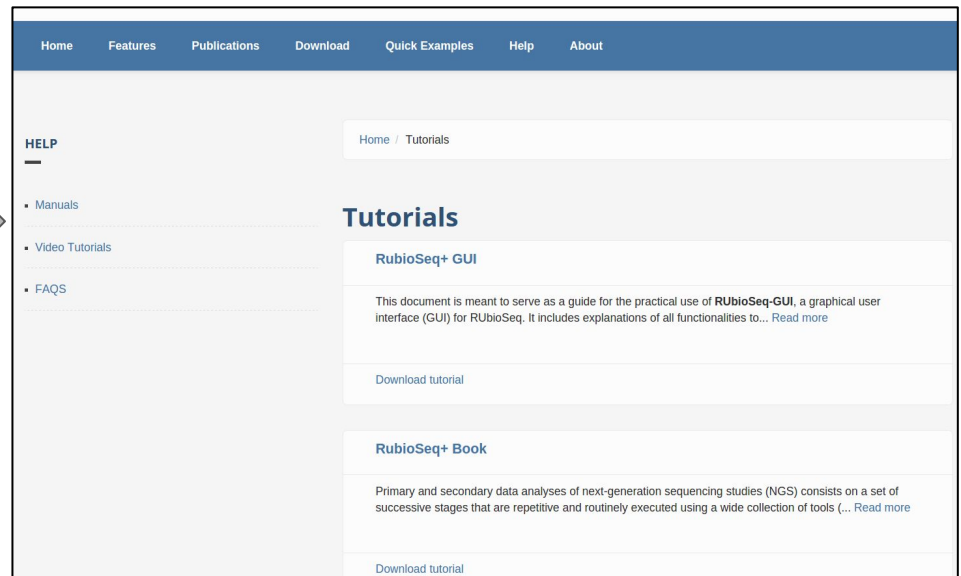
# Ask for some help

**RUBioSeq** is under **development** by the Bioinformatics Unit at CNIO (Madrid, Spain):

- Miriam Rubio-Camarillo ([mrubioc@cnio.es](mailto:mrubioc@cnio.es)).
- José María Fernández ([jmfernandez@cnio.es](mailto:jmfernandez@cnio.es)).
- Gonzalo Gómez-López ([ggomez@cnio.es](mailto:ggomez@cnio.es)).

<http://rubioseq.bioinfo.cnio.es/tutorials>

- Video-tutorials.
- Example exercises.



[RUBioSeq: a suite of parallelized pipelines to automate exome variation and bisulfite-seq analyses](#)

Miriam Rubio-Camarillo\*, Gonzalo Gómez-López, José M. Fernández, Alfonso Valencia and David G. Pisano.

2013, Bioinformatics, 29 (13): 1687-1689

# (Virtual machine, Ubuntu system)

1. First, install Docker client (depending on your Operating System):

Follow the corresponding guidelines:

- Ubuntu (<https://docs.docker.com/installation/ubuntu/linux/>)
- Mac OS X (<https://docs.docker.com/installation/mac/>)
- Windows (<https://docs.docker.com/installation/windows/>)

2. Second, launch the Docker Image (it will download it from the Internet):

- Windows:  
(Start Menu): Program Files > Boot2Docker

- Unix (in terminal):

```
docker run -ti -p 0.0.0.0:8080:8080 --name RUBioSeq ubio/rubioseq:latest /bin/bash
```

3. Try it out:

```
perl /home/RUBioSeq/RUBioSeq3.7/RUBioSeq.pl -h
```

More info: [http://rubioseq.bioinfo.cnio.es/rubioseq\\_docker](http://rubioseq.bioinfo.cnio.es/rubioseq_docker)



<https://cniobu.github.io/pm17/>

## pm17 Precision Medicine

Instituto Gulbenkian de Ciência  
14-17 November 2017

Software Installation

Software Installation Quick Guide (root user)


<https://cniobu.github.io/pm17>


Software Installation

Software Installation Quick Guide (root)



## PM17 - Software Dependencies

Javier Perales-Paton

[jperales@cniobu.es](mailto:jperales@cniobu.es)
Quick manual installation (root user in Ubuntu 16)

1. Install Java (version 1.7 or 1.8) in ubuntu.

When the installation is finished, check that Java is correctly installed:

```
participant@machine:~$ java -version
openjdk version "1.8.0_151"
OpenJDK Runtime Environment (build 1.8.0_151-8u151-b12-0ubuntu0.16.04.2-b12)
OpenJDK 64-Bit Server VM (build 25.151-b12, mixed mode)
```

2. Install python

```
sudo apt-get install python-dev
```

3. Get the library libmysqlclient

```
sudo apt-get install libmysqlclient-dev
```

4. Perl Modules ::

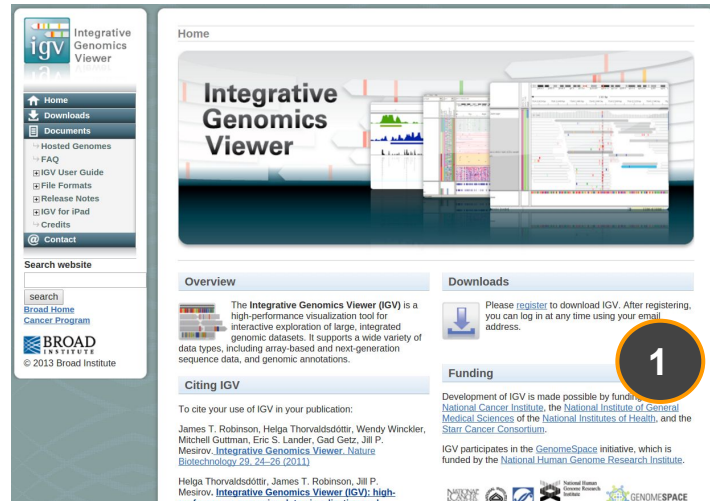
Rubioseq built on perl modules which must be correctly installed. If one of these is missing, the pipeline will get a crash.

- DBI
- DBD::mysql
- XML::LibXML
- Carp
- FindBin
- File::Basename
- File::Spec
- File::Copy
- Getopt::Long
- Class::Inspector

```
sudo cpan DBI DBD::mysql
```



# Integrative Genomics Viewer (IGV)



Open Firefox.

<https://www.broadinstitute.org/igv/>

1. Download section: Register and Fill out the form.

2. Download the Binary distribution: file.zip  
Alternatively, you can use this link:

[http://data.broadinstitute.org/igv/projects/downloads/IGV\\_2.3.66.zip](http://data.broadinstitute.org/igv/projects/downloads/IGV_2.3.66.zip)

## 3. Binary Distribution

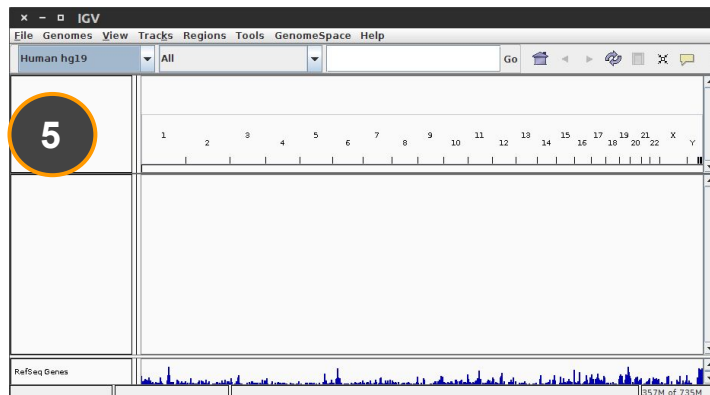
Download and unzip the binary distribution archive in a folder of your choosing. IGV is launched from a command prompt -- follow instructions in the "readme" file. To launch igv on Mac or Linux platforms use the shell script "igv.sh". On Windows use "igv.bat".

Download  
Binary Distribution

2

3. Extract the Zip file.

4. Go to the new Directory. Click on igv.jar



IGV\_2.3.66.zip



IGV\_2.3.66



4



