

Winter School on Imaging Genetics

Whole Transcriptome Data Analysis

Francesca CORDERO

Computer Science Department
University of Turin

francesca.cordero@unito.it

Marco BECCUTI

Computer Science Department
University of Turin

marco.beccuti@unito.it

Verona, 26-29 November, 2019

Winter School on Imaging Genetics

Whole Transcriptome Data Analysis

Reproducibility and docker

R - basic functions

Theoretical and practical session for RNA-seq analysis

Verona, 26-29 November, 2019

Reproducible research for NGS analysis

The package **docker4seq** was developed in **Reproducible Bioinformatics Project** to facilitate the use of computing demanding applications in the field of NGS data analysis.



Reproducible Bioinformatics Project

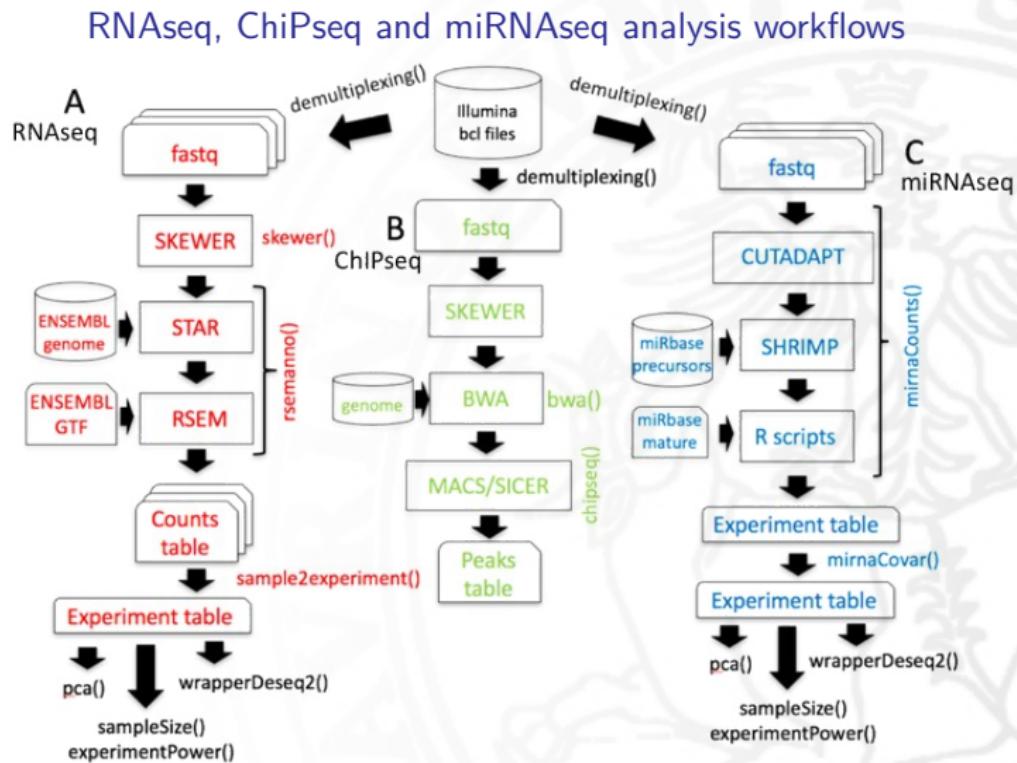
A project to provide reproducible results in Bioinformatics using Docker images

- Reproducible Bioinformatics Project (RBP) is community open to anyone interested to shared workflows under the umbrella of reproducibility;
- To enable an easy access NGS data analysis pipeline to users without advanced computer science skills;
- To provide robust, reproducible, portable bioinformatics workflows.

<http://reproducible-bioinformatics.org/>

Reproducible research for NGS analysis

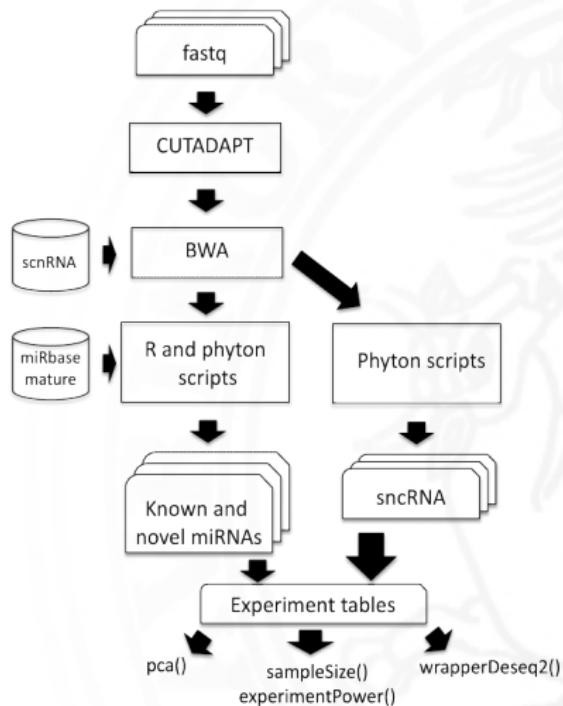
Current available workflows in RBP (docker4seq):



Reproducible research for NGS analysis

Current available workflows in RBP(docker4seq):

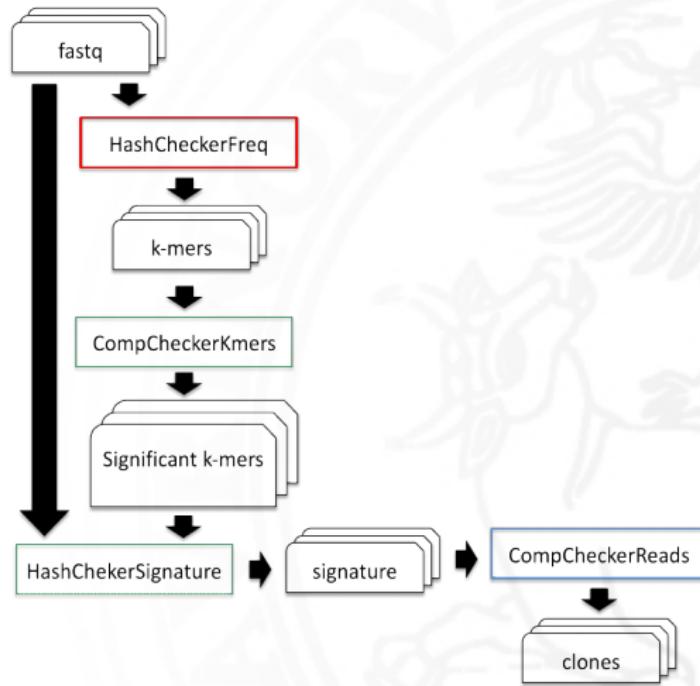
Small non-coding RNA workflow



Reproducible research for NGS analysis

Current available workflows in RBP(docker4seq):

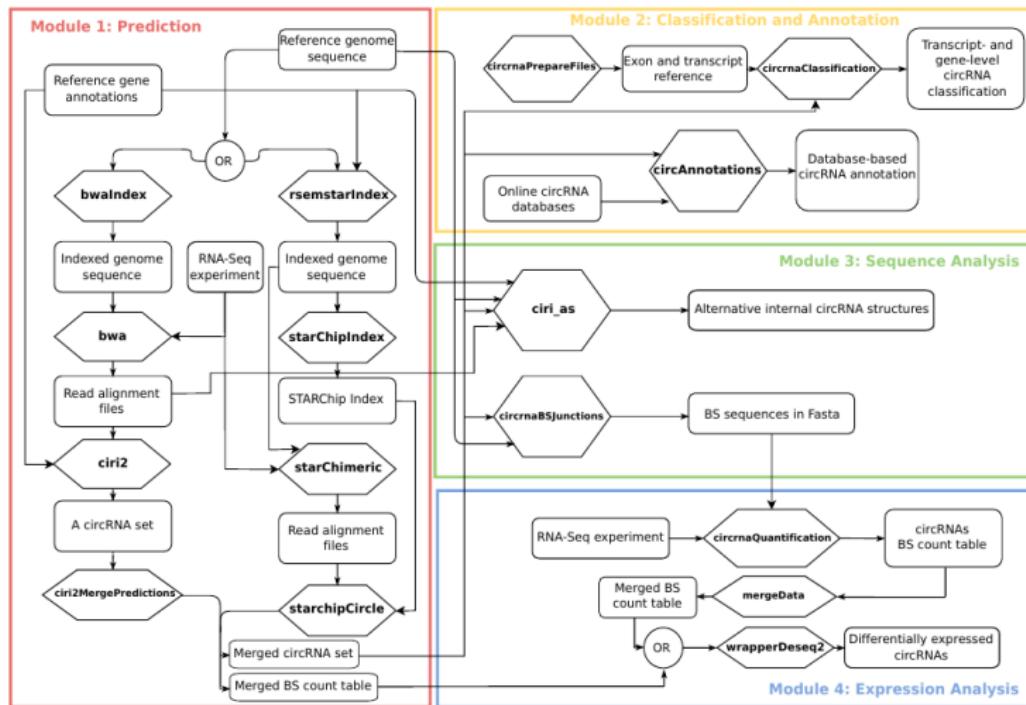
Minimal residual disease in B-cell lymphoma workflow



Reproducible research for NGS analysis

Current available workflows in RBP(docker4seq):

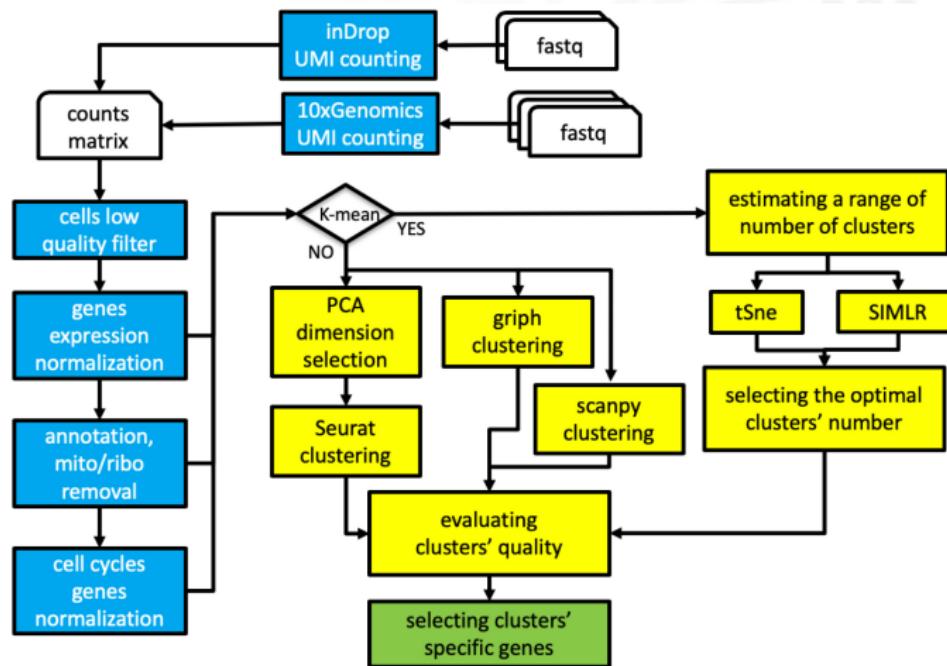
Prediction circular RNAs workflow



Reproducible research for NGS analysis

Current available workflows in RBP(rCasc):

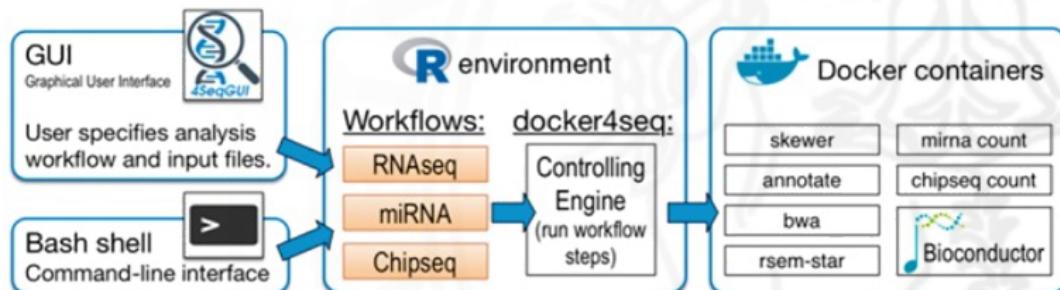
Single cell analysis workflow



Reproducible research for NGS analysis

The package **docker4seq** general schema:

- Any workflow is specified as a set of **R functions** that defines and controls the correct execution of all its tasks;
- any single task is encapsulated into a **container** (i.e. docker images) to guarantee the computation reproducibility and portability.



Any workflow must be supported by an explanatory vignette

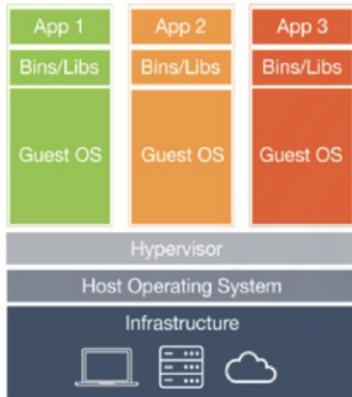
Reproducible research for NGS analysis

How to make easier the use of these workflows for beginner users:

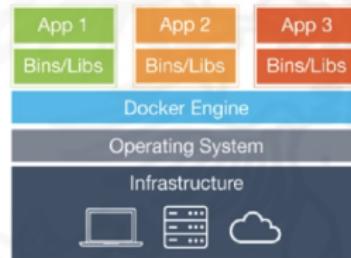
The screenshot shows a software interface for reproducible NGS analysis. On the left, a sidebar lists various analysis tools: RNAseq, miRNA, Chipseq, Tools, Genome indexing STAR-RSEM, Genes, isoforms counting RSEM, Trans. pseudo-reference building (Salmon), Trans. and genes counting (Salmon), From samples to experiment, PCA analysis, Sample size estimator, Experiment statistical Power, DESeq2 analysis, and Count Filter. The main panel is titled "PCA" and contains fields for "FPKM/TPM file:" (with "Browse" and "Cancel" buttons) and "Output folders:" (with "Cancel" and "Browse" buttons). Below these are fields for "Component 1:" and "Component 2:". Under "Data type:", the radio button for "FPKM" is selected. A "Legend position:" dropdown menu is set to "bottomleft". Under "Covariates:", the radio button for "yes" is selected. At the bottom right are buttons for "Execute", "Save conf.", "Reset", and "Close". A "Process status" section at the bottom is currently empty.

Reproducible research for NGS analysis

Container and Virtual machine are two virtualization techniques:



Virtual Machines



Containers



Virtual Machines



Size		
Startup Execution		

In Containers the sharing of same OS Kernel with the real machine reduces the portability, but

Docker container, VM and real server: a comparison

In [1] a comparison among physical server, KVM, and Docker is reported.

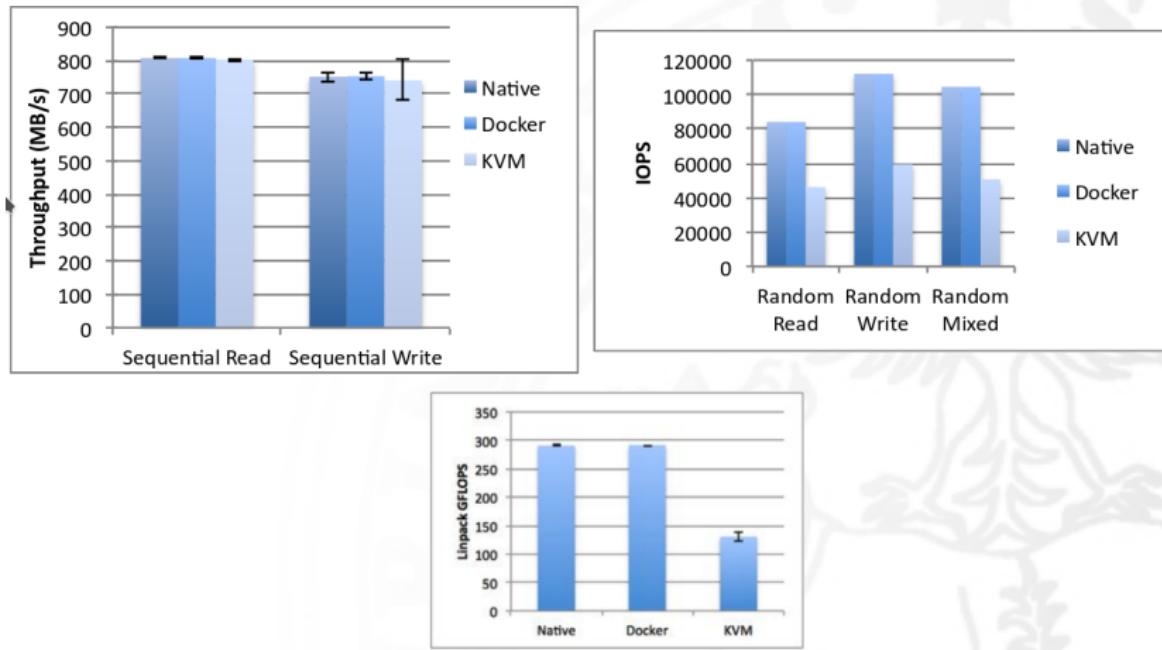
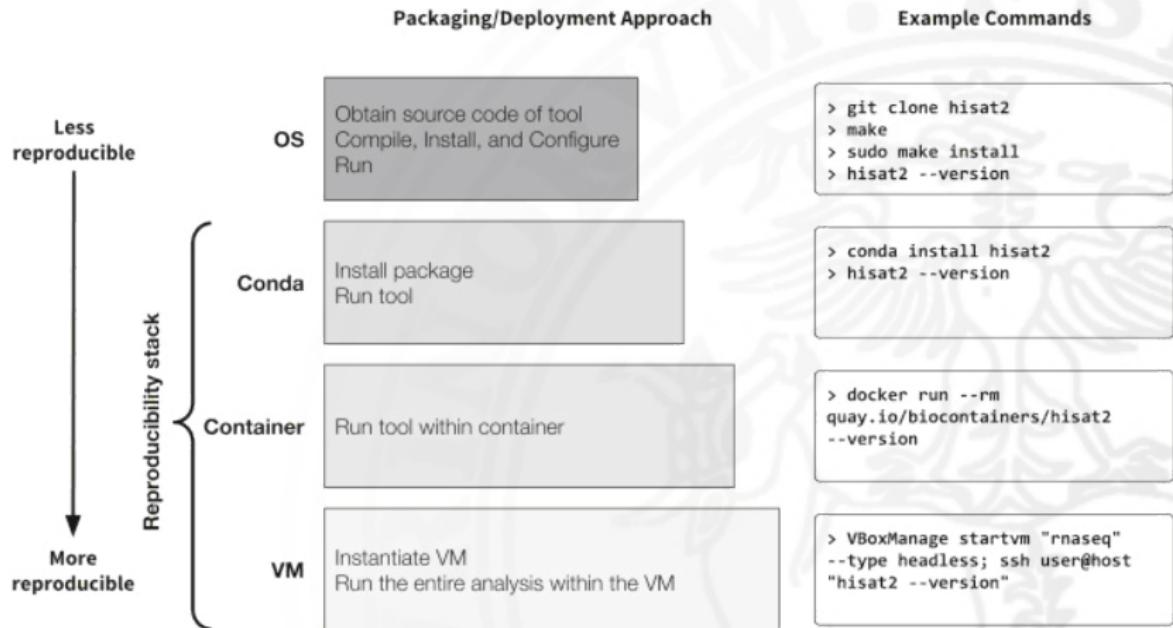


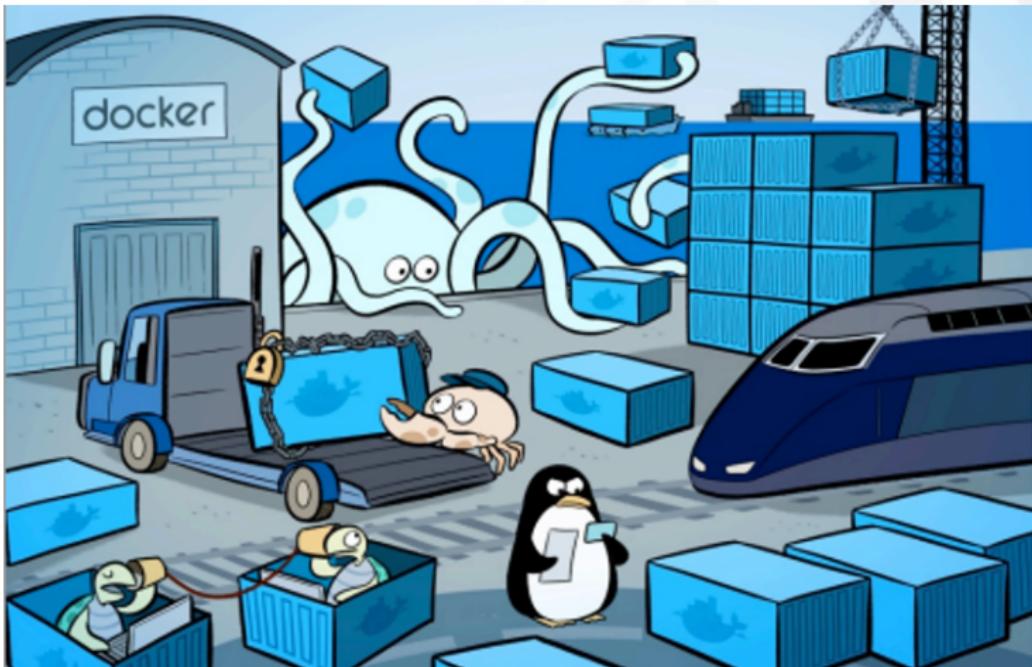
Figure 1. Linpack performance on two sockets (16 cores). Each data point is the arithmetic mean obtained from ten runs. Error bars indicate the standard deviation obtained over all runs.

[1] W. Felter, A. Ferreira, R. Rajamony and J. Rubio, *An updated performance comparison of virtual machines and Linux containers*, 2015 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), Philadelphia, PA, 2015, pp. 171-172.

Computational Reproducibility Stack

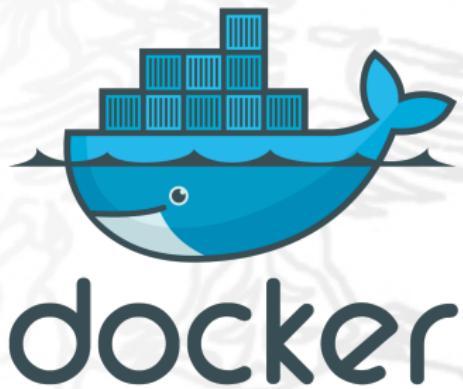


Docker project



Docker project

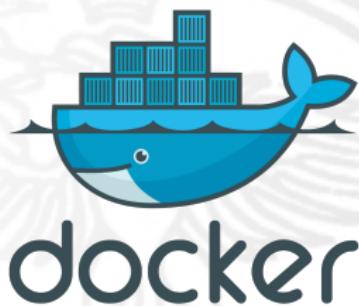
- Docker is an open-source project based on Linux containers.
- Others Linux container technologies include Solaris Zones, BSD jails, and LXC, which have been around for many years.



Why to use Docker??

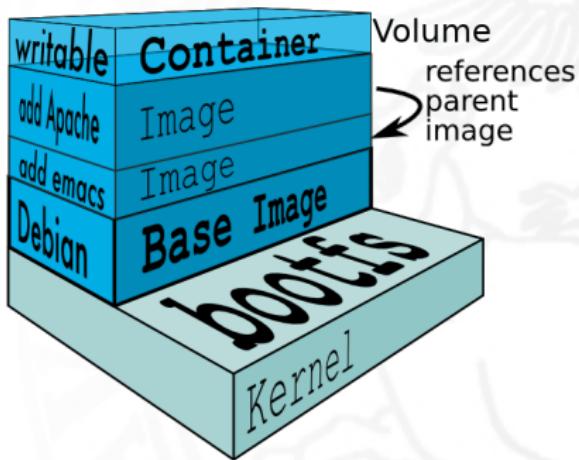
Why to use Docker?

- **Ease of use:** Docker has made it much easier for anyone to take advantage of containers in order to quickly build and test portable applications;
- **Speed:** Docker containers are very lightweight and fast.
- **Docker Hub:** Docker users also benefit from the increasingly rich repository of Docker Hub, which you can think of as an "app store for Docker images";
- **Modularity and Scalability:** Docker makes it easy to break out your application's functionality into individual containers.



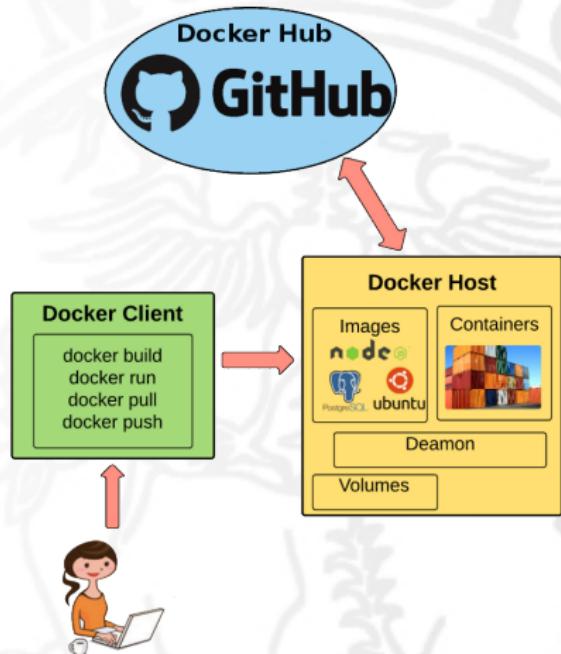
Docker basics

- **Image** is an executable package that includes everything needed to run an application (i.e. its code, libraries, environment variables, and configuration files);
- **Container** is a run-time instance of an image;
- **Volume** is used to share files from the host machine to containers.



Docker schema

- **Docker client:** provides an interface for users;
- **Docker Host:** executes the commands sent to the Docker Client;
- **Docker Hub:** remote repository storing docker's images.



How to install docker4seq packages



How to install docker4seq packages

- You need to install on your machine:
 - ▶ A Linux-based operating system (i.e. Fedora, Ubuntu, Idots);
 - ▶ R environment (see <http://cran.mirror.garr.it/mirrors/CRAN/>)
 - ▶ Docker environment (see <https://www.docker.com/>)
- The minimal hardware requirements are an x86_64 compatible processor, 32Gb RAM, an SSD 250GB.

How to install docker4seq packages

- You need to execute in R environment the following commands:

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

```
> library(devtools)
```

```
> install_github("kendomaniac/docker4seq", ref = "master")
```

- To use and download docker images (about 50Gb):

```
> library(docker4seq)
```

```
> downloadContainers(group = "docker")
```

During this course ...

- We use Rstudio web server running on two remote machines in the HPC4IA cluster (see <https://hpc4ai.it/>)

The screenshot shows the RStudio web interface running on a remote machine. The top bar indicates the session is on port 8787. The left sidebar shows the project structure with files like 'scripted.R' and 'scriptTest.R'. The main area has tabs for Environment, History, and Connections. The Environment tab shows variables like 'fastes', 'folders', 'home', 'i', and 'my.covariates'. The History tab shows a log of R commands run. The Terminal tab shows Docker command history. The bottom navigation bar includes links to various PDFs and a 'Show all' button.

```
## #folders must contain R1 fastq
71 folders <- paste(getwd(), c("c14", "c15", "c16", "n2", "n4", "n5", "n11"), sep="/")
72 my.covariates <- c(rep("Cov.1", 1), rep("Cov.2", 4))
73 sampleExperiment.table(folders, covariates, my.covariates, batch = NULL,
74   bio.type = "protein coding", "unary pseudogene",
75   "unprocessed pseudogene", "processed pseudogene",
76   "transcribed unprocessed pseudogene", "processed transcript",
77   "antisense", "sense unprocessed pseudogene", "sense pseudogene",
78   "lincRNA", "sense intronic", "transcribed processed pseudogene",
79   "sense overlapping", "IG V pseudogene", "pseudogene", "TR V gene",
80   "sense overlapping", "IG V gene", "bidirectional promoter lncRNA",
81   "snorna", "microRNA", "lncRNA", "lncRNA C gene", "lncRNA",
82   "TR J gene", "TR C gene", "TR V pseudogene", "TR J pseudogene",
83   "IG D gene", "ribosome", "IG C pseudogene", "IG D gene", "TEC",
84   "IG J pseudogene", "scRNA", "scRNA", "svltnRNA", "sRNA", "macro lncRNA",
85   "macro coding", "IG pseudogene", output.prefix = ".")
86 pca.experiment.table <- log2PCA.table(folders, covariates.inNames=TRUE,
87   sampleName=TRUE, principal.components = c(1, 2),
88   legend.position="topright", pdf=TRUE, output.folder = getwd())
89 file.rename(from="pca.pdf", to="pca.wc.c17.pdf")
90 wrapperDeseq2(output.folder=getwd(), group="docker",
91   experiment.table="counts.txt", log2fc=1, fdr=0.1,
92   experiment.table="counts.txt", log2fc=1, fdr=0.1,
93   ref.cover="Cov.1", type="gene", batch=FALSE)
94 |
95 |
96.1 (Top Level) 2 R Script 1
```

```
## + legend.position="topright", pdf=TRUE, output.folder = getwd()
## + file.rename(from="pca.pdf", to="pca.wc.c17.pdf")
[1] TRUE
> wrapperDeseq2(output.folder=getwd(), group="docker",
+   experiment.table="counts.txt", log2fc=1, fdr=0.1,
+   ref.cover="Cov.1", type="gene", batch=FALSE)
```

```
In your system the following version of Docker is installed:
Docker version 19.03.4, build 9013bf583a62902df96203b2aa3033cb9234a2d6b195c4658bbdf2e1bf3e2ee5b956a0c8f
```

```
Docker ID is:
62902df96203
....
```

```
Docker exit status: 0
```

```
DESeq2 analysis is finished
>
```

```
SUPP-D-19-0...pdf ^ SUPP-D-19-0...pdf ^ Suppl..._BMC...pdf ^ Fundamenta l...zip ^ PDP2020_pa...pdf ^ tidyverse_12.1.zip ^ file_show.gz ^ Show all x
```