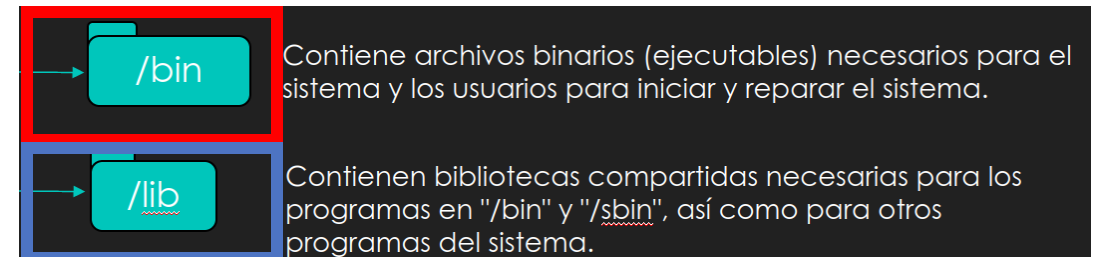
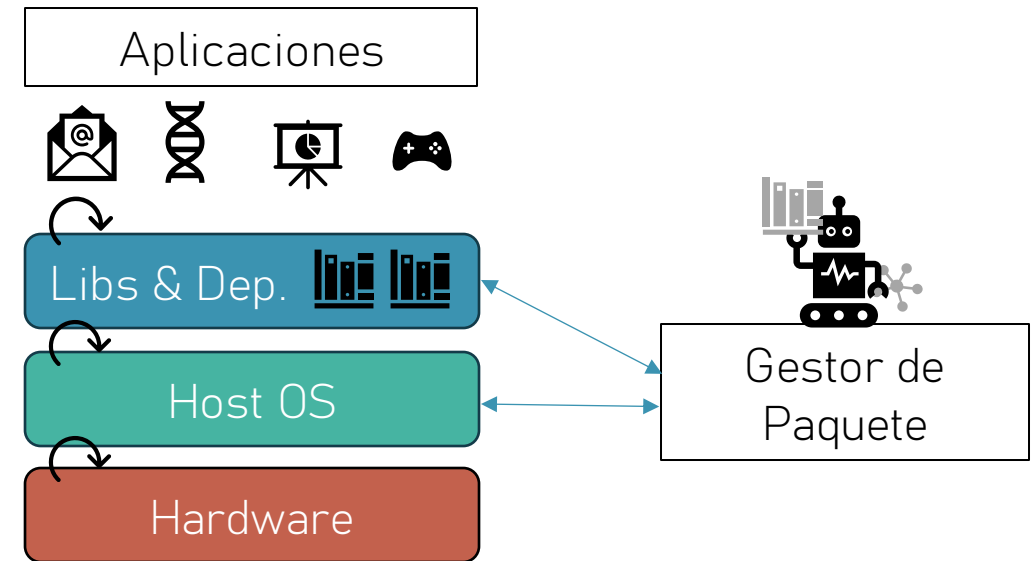


ARCHIVOS Y PROGRAMAS BIOINFORMÁTICOS

Murilo Cassiano, M.Sc.

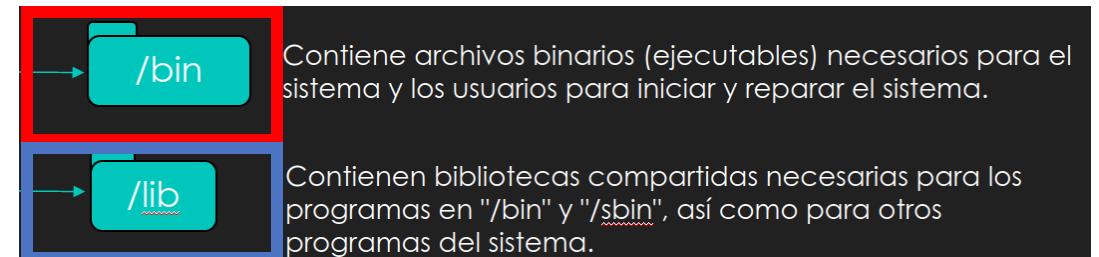
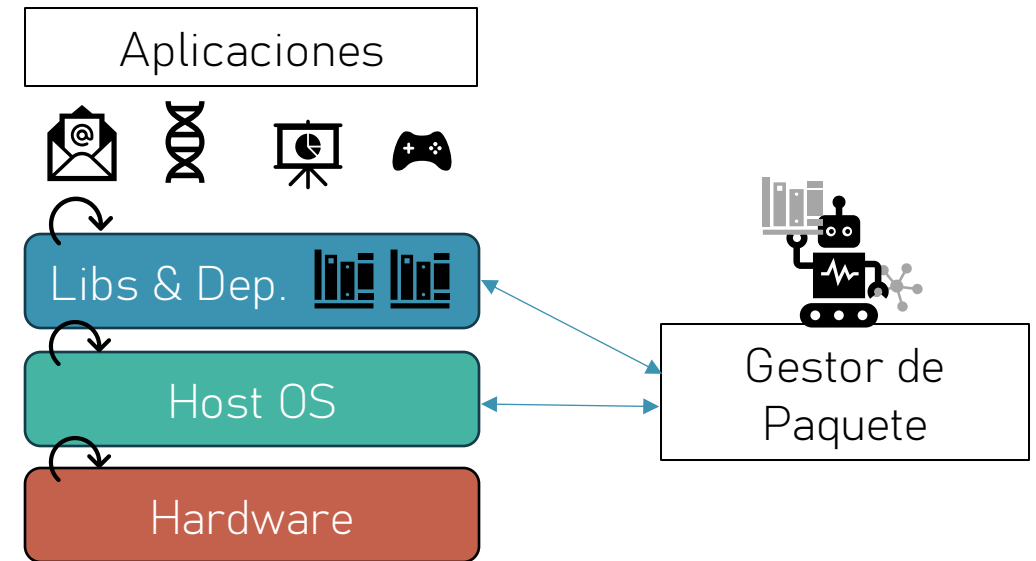
GESTORES DE PAQUETES

- Es una herramienta que ayuda a administrar software en un sistema operativo.
- Facilita la instalación, actualización, configuración y eliminación de programas en una computadora.
- Mantienen un registro de las dependencias del software, asegurando que las **aplicaciones** funcionen correctamente y que todas las **bibliotecas** necesarias estén presentes en el sistema.



GESTOR DE PAQUETE: APT

- APT (Advanced Package Tool) es el sistema de gestión de paquetes utilizado en sistemas basados en Debian
 - Ubuntu
 - Mint
 - PopOs
 - Debian Linux
- APT simplifica la instalación y actualización de software, automatizando muchos procesos.
- A través de la línea de comandos, los usuarios pueden utilizar varios comandos para interactuar con APT.



GESTOR DE PAQUETE:

APT

- Buscar un paquete:

`apt search término_de_búsqueda`

Este comando busca paquetes relacionados con el término de búsqueda especificado.

- Mostrar información sobre un paquete:

`apt show nombre_del_paquete`

Este comando muestra información detallada sobre un paquete específico, incluyendo su descripción y dependencias.

```
murilo@muca10-t14:~$ apt search htop
Sorting... Done
Full Text Search... Done
aha/stable,stable,now 0.5.1-3 amd64 [installed,automatic]
  ANSI color to HTML converter

bpytop/stable,stable 1.0.68-1 all
  Resource monitor that shows usage and stats

htop/stable,stable 1.2.13-1 amd64
  Modern and colorful command line resource monitor that shows usage and stats

htop/stable,stable 3.2.2-2 amd64
  interactive processes viewer

libauthen-oath-perl/stable,stable 2.0.1-2 all
  Perl module for OATH One Time Passwords

pftools/stable,stable 3.2.12-1 amd64
  build and search protein and DNA generalized profiles
```

```
murilo@muca10-t14:~$ apt show htop
Package: htop
Version: 3.2.2-2
Priority: optional
Section: utils
Maintainer: Daniel Lange <DLange@debian.org>
Installed-Size: 387 kB
Depends: libc6 (>= 2.34), libncursesw6 (>= 6), libnl-3-200 (>= 3.2.7), libnl-genl-3-200 (>= 3.2.7), libtinfo6 (>= 6)
Suggests: lm-sensors, lsof, strace
Homepage: https://htop.dev/
Tag: admin::monitoring, implemented-in::c, interface::text-mode,
  role::program, scope::utility, uitoolkit::ncurses, use::monitor,
  works-with::software:running
Download-Size: 152 kB
APT-Sources: http://ftp.au.debian.org/debian bookworm/main amd64 Packages
Description: interactive processes viewer
  Htop is an ncurses-based process viewer similar to top, but it
  allows one to scroll the list vertically and horizontally to see
  all processes and their full command lines.
.
  Tasks related to processes (killing, renicing) can be done without
  entering their PIDs.
```


GESTOR DE PAQUETE:

APT

```
murilo@muca10-t14:~$ apt search mafft
Sorting... Done
Full Text Search... Done
mafft/stable,stable 7.505-1 amd64
  Multiple alignment program for amino acid or nucleotide sequences

muscle/stable,stable 1:5.1.0-1 amd64
  Multiple alignment program of protein sequences

probalign/stable,stable 1.4-10 amd64
  multiple sequence alignment using partition function posterior probabilities

t-coffee/stable,stable 13.45.0.4846264+really13.41.0.28bdc39+dfsg-1 amd64
  Multiple Sequence Alignment
```

Intenta instalar con `apt install htop` y
`apt install muscle`

```
murilo@muca10-t14:~$ apt show muscle
Package: muscle
Version: 1:5.1.0-1
Priority: optional
Section: science
Maintainer: Debian Med Packaging Team <debian-med-packaging@lists.aliases.debian.org>
Installed-Size: 863 kB
Provides: muscle-doc
Depends: libc6 (>= 2.34), libgcc-s1 (>= 3.0), libgomp1 (>= 4.9), libstdc++6 (>= 11)
Conflicts: muscle-doc
Replaces: muscle-doc
Enhances: seaview, t-coffee
Homepage: https://www.drive5.com/muscle/
Tag: biology::format:aln, biology::peptidic, field::biology,
  field::biology:bioinformatics, implemented-in::c++,
  interface::commandline, role::program, scope::utility, use::comparing,
  works-with-format::plaintext, works-with::TODO
Download-Size: 297 kB
APT-Sources: http://ftp.au.debian.org/debian bookworm/main amd64 Packages
Description: Multiple alignment program of protein sequences
 MUSCLE is a multiple alignment program for protein sequences. MUSCLE
 stands for multiple sequence comparison by log-expectation. In the
 authors tests, MUSCLE achieved the highest scores of all tested
 programs on several alignment accuracy benchmarks, and is also one of
 the fastest programs out there.

.
Muscle v5 is a major re-write of MUSCLE based on new algorithms.
.
Users should be aware that command line arguments compared to version
3.x of MUSCLE have changed!
.
Highest accuracy, scalable to thousands of sequences
.
Compared to previous versions, Muscle v5 is much more accurate, is often
faster, and scales to much larger datasets. At the time of writing (late
2021), Muscle v5 has the highest scores on multiple alignment benchmarks
including Balibase, Bralibase, Prefab and Balifam. It can align tens of
thousands of sequences with high accuracy on a low-cost commodity computer
(say, an 8-core Intel CPU with 32 Gb RAM). On large datasets, Muscle v5
is 20-30% more accurate than MAFFT and Clustal-Omega.
.
Alignment ensembles
.
Muscle v5 can generate ensembles of high-accuracy alternative alignments.
All replicates have equal average accuracy on benchmark test, including
the MSA made with default parameters. By comparing results of downstream
analysis (trees, structure prediction...) on different replicates, you can
assess the effects of alignment errors on your study.
```

GESTOR DE PAQUETE: APT

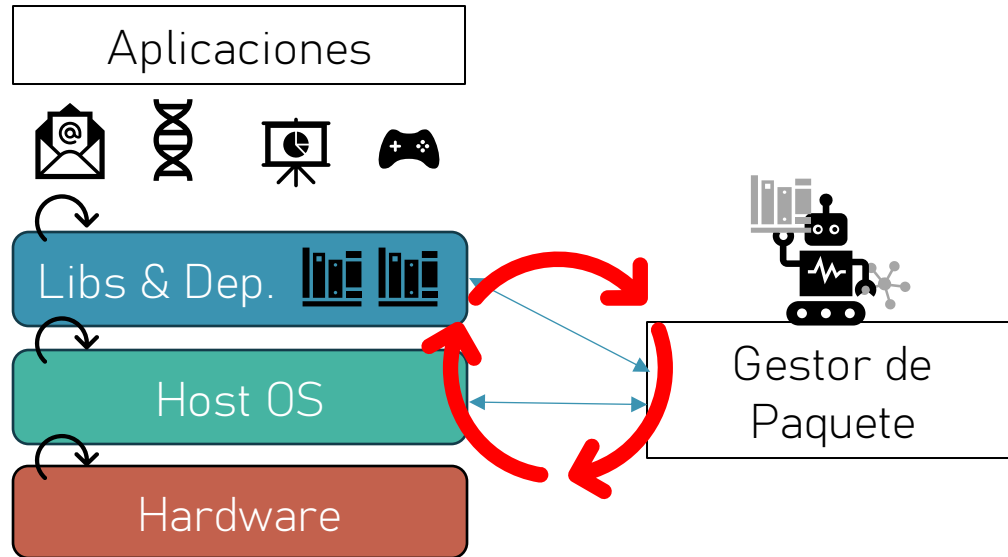
- La instalación de un paquete requiere que se modifiquen las dependencias y otras bibliotecas que están instaladas o que faltan en el sistema operativo.
- ¡El usuario debe ser parte de los administradores! ¡SUDO!
- Siempre es arriesgado porque los programas instalados ya utilizan librerías que están configuradas

Instalar un nuevo paquete:

```
sudo apt install nombre_del_paquete
```

Este comando instala un nuevo paquete en el sistema.

```
sudo apt install muscle mafft
```



Eliminar un paquete:

```
sudo apt remove nombre_del_paquete
```

Este comando elimina un paquete del sistema, pero mantiene sus configuraciones.

Para eliminar también las configuraciones, se puede usar

```
sudo apt purge nombre_del_paquete
```

.

```
sudo apt purge mafft
```

GESTOR DE PAQUETE: APT

- Actualizar la lista de paquetes:

```
sudo apt update
```

Este comando actualiza la lista de paquetes disponibles en los repositorios.

- Actualizar los paquetes instalados:

```
sudo apt upgrade
```

Este comando actualiza todos los paquetes instalados en el sistema a sus últimas versiones.

```
murilo@muca10-t14:~$ sudo apt update
Hit:1 https://brave-browser-apt-release.s3.brave.com stable InRelease
Get:2 http://deb.debian.org/debian-security bookworm-security InRelease [48.0 kB]
Hit:3 http://repository.spotify.com stable InRelease
Hit:4 http://deb.debian.org/debian bookworm InRelease
Get:5 http://deb.debian.org/debian-security bookworm-security/main amd64 Packages [87.2 kB]
Get:6 http://deb.debian.org/debian-security bookworm-security/main Translation-en [49.2 kB]
Hit:7 http://cloud.r-project.org/bin/linux/debian bookworm-cran40/ InRelease
Hit:8 http://ftp.au.debian.org/debian bookworm InRelease
Get:9 http://ftp.au.debian.org/debian bookworm-updates InRelease [52.1 kB]
Fetched 236 kB in 3s (75.6 kB/s)
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
203 packages can be upgraded. Run 'apt list --upgradable' to see them.
```

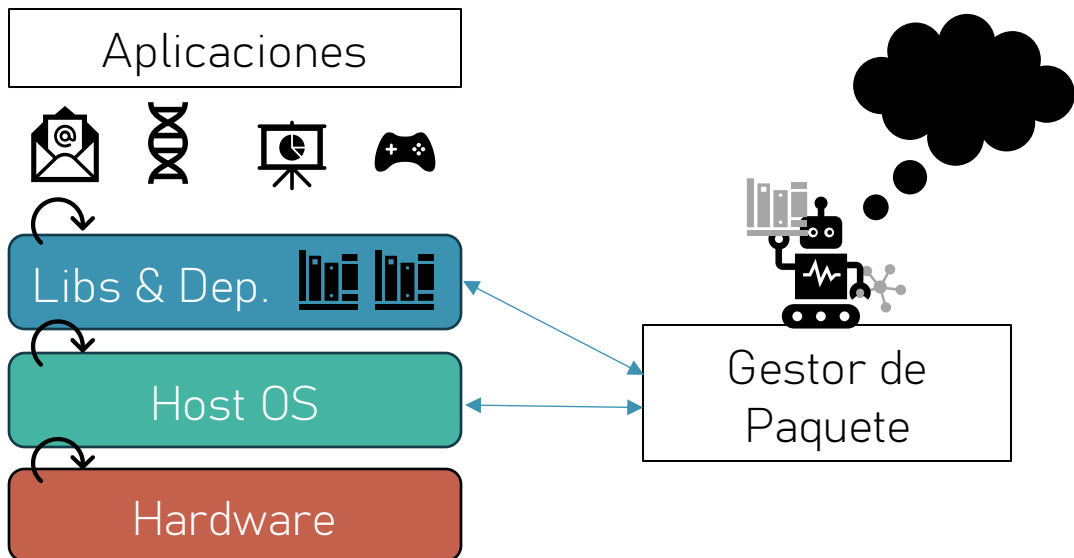
```

nurltoemca10-t14:~$ sudo apt upgrade
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
Calculating upgrade... Done

The following packages were automatically installed and are no longer required:
coincor-libbcb3 coinor-libcgl1 coinor-libcpl1 coinor-libcoinimpv5 coinor-libcoinutils3v5 coinoir-libosiv5
libabw-0.1-1 libbox2d2 libcdr-0.1-1 libe-book-0.1-1 libepubgen-0.1-1 libetonyek-0.1-1 libfreehand-0.1-1
libmspub-0.1-1 libmwaw-0.3-3 liboddfont-0.1-1 libpagemaker-0.0-0 libqxp-0.0-0 liboffice-base-core
libstaroffice-0.0 libvisio-0.1-1 libwpd-0.10-10 libwpg-0.3-3 libwps-0.4-4 libzmf-0.0-0 lp-solve

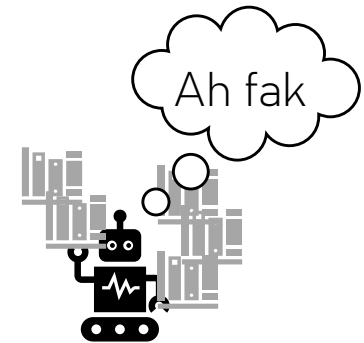
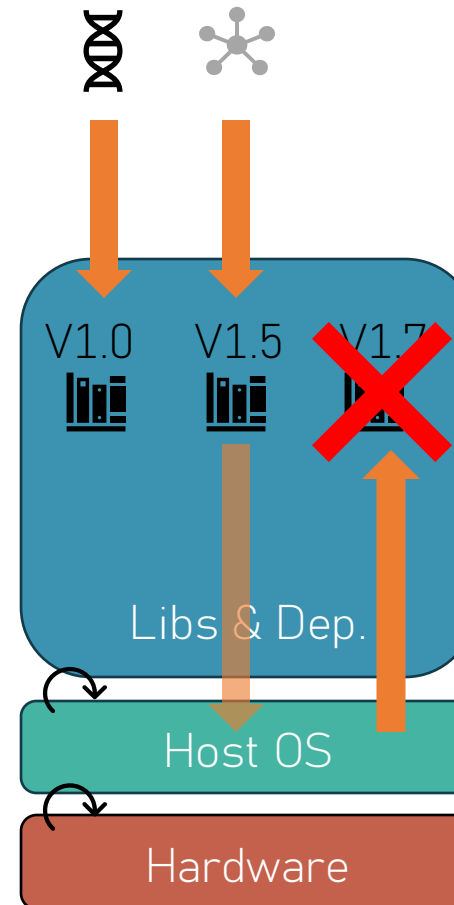
Use 'sudo apt autoremove' to remove them.

The following NEW packages will be installed:
libb3 libb4 libb5 libb6 libb7 libb8 libb9 libb10 libb11 libb12 libb13 libb14 libb15 libb16 libb17 libb18 libb19 libb20 libb21 libb22 libb23 libb24 libb25 libb26 libb27 libb28 libb29 libb30 libb31 libb32 libb33 libb34 libb35 libb36 libb37 libb38 libb39 libb40 libb41 libb42 libb43 libb44 libb45 libb46 libb47 libb48 libb49 libb50 libb51 libb52 libb53 libb54 libb55 libb56 libb57 libb58 libb59 libb60 libb61 libb62 libb63 libb64 libb65 libb66 libb67 libb68 libb69 libb70 libb71 libb72 libb73 libb74 libb75 libb76 libb77 libb78 libb79 libb80 libb81 libb82 libb83 libb84 libb85 libb86 libb87 libb88 libb89 libb90 libb91 libb92 libb93 libb94 libb95 libb96 libb97 libb98 libb99 libb100 libb101 libb102 libb103 libb104 libb105 libb106 libb107 libb108 libb109 libb110 libb111 libb112 libb113 libb114 libb115 libb116 libb117 libb118 libb119 libb120 libb121 libb122 libb123 libb124 libb125 libb126 libb127 libb128 libb129 libb130 libb131 libb132 libb133 libb134 libb135 libb136 libb137 libb138 libb139 libb140 libb141 libb142 libb143 libb144 libb145 libb146 libb147 libb148 libb149 libb150 libb151 libb152 libb153 libb154 libb155 libb156 libb157 libb158 libb159 libb160 libb161 libb162 libb163 libb164 libb165 libb166 libb167 libb168 libb169 libb170 libb171 libb172 libb173 libb174 libb175 libb176 libb177 libb178 libb179 libb180 libb181 libb182 libb183 libb184 libb185 libb186 libb187 libb188 libb189 libb190 libb191 libb192 libb193 libb194 libb195 libb196 libb197 libb198 libb199 libb200 libb201 libb202 libb203 libb204 libb205 libb206 libb207 libb208 libb209 libb210 libb211 libb212 libb213 libb214 libb215 libb216 libb217 libb218 libb219 libb220 libb221 libb222 libb223 libb224 libb225 libb226 libb227 libb228 libb229 libb230 libb231 libb232 libb233 libb234 libb235 libb236 libb237 libb238 libb239 libb240 libb241 libb242 libb243 libb244 libb245 libb246 libb247 libb248 libb249 libb250 libb251 libb252 libb253 libb254 libb255 libb256 libb257 libb258 libb259 libb260 libb261 libb262 libb263 libb264 libb265 libb266 libb267 libb268 libb269 libb270 libb271 libb272 libb273 libb274 libb275 libb276 libb277 libb278 libb279 libb280 libb281 libb282 libb283 libb284 libb285 libb286 libb287 libb288 libb289 libb290 libb291 libb292 libb293 libb294 libb295 libb296 libb297 libb298 libb299 libb300 libb301 libb302 libb303 libb304 libb305 libb306 libb307 libb308 libb309 libb310 libb311 libb312 libb313 libb314 libb315 libb316 libb317 libb318 libb319 libb320 libb321 libb322 libb323 libb324 libb325 libb326 libb327 libb328 libb329 libb330 libb331 libb332 libb333 libb334 libb335 libb336 libb337 libb338 libb339 libb340 libb341 libb342 libb343 libb344 libb345 libb346 libb347 libb348 libb349 libb350 libb351 libb352 libb353 libb354 libb355 libb356 libb357 libb358 libb359 libb360 libb361 libb362 libb363 libb364 libb365 libb366 libb367 libb368 libb369 libb370 libb371 libb372 libb373 libb374 libb375 libb376 libb377 libb378 libb379 libb380 libb381 libb382 libb383 libb384 libb385 libb386 libb387 libb388 libb389 libb390 libb391 libb392 libb393 libb394 libb395 libb396 libb397 libb398 libb399 libb400 libb401 libb402 libb403 libb404 libb405 libb406 libb407 libb408 libb409 libb410 libb411 libb412 libb413 libb414 libb415 libb416 libb417 libb418 libb419 libb420 libb421 libb422 libb423 libb424 libb425 libb426 libb427 libb428 libb429 libb430 libb431 libb432 libb433 libb434 libb435 libb436 libb437 libb438 libb439 libb440 libb441 libb442 libb443 libb444 libb445 libb446 libb447 libb448 libb449 libb450 libb451 libb452 libb453 libb454 libb455 libb456 libb457 libb458 libb459 libb460 libb461 libb462 libb463 libb464 libb465 libb466 libb467 libb468 libb469 libb470 libb471 libb472 libb473 libb474 libb475 libb476 libb477 libb478 libb479 libb480 libb481 libb482 libb483 libb484 libb485 libb486 libb487 libb488 libb489 libb490 libb491 libb492 libb493 libb494 libb495 libb496 libb497 libb498 libb499 libb500 libb501 libb502 libb503 libb504 libb505 libb506 libb507 libb508 libb509 libb510 libb511 libb512 libb513 libb514 libb515 libb516 libb517 libb518 libb519 libb520 libb521 libb522 libb523 libb524 libb525 libb526 libb527 libb528 libb529 libb530 libb531 libb532 libb533 libb534 libb535 libb536 libb537 libb538 libb539 libb540 libb541 libb542 libb543 libb544 libb545 libb546 libb547 libb548 libb549 libb550 libb551 libb552 libb553 libb554 libb555 libb556 libb557 libb558 libb559 libb560 libb561 libb562 libb563 libb564 libb565 libb566 libb567 libb568 libb569 libb570 libb571 libb572 libb573 libb574 libb575 libb576 libb577 libb578 libb579 libb580 libb581 libb582 libb583 libb584 libb585 libb586 libb587 libb588 libb589 libb590 libb591 libb592 libb593 libb594 libb595 libb596 libb597 libb598 libb599 libb600 libb601 libb602 libb603 libb604 libb605 libb606 libb607 libb608 libb609 libb610 libb611 libb612 libb613 libb614 libb615 libb616 libb617 libb618 libb619 libb620 libb621 libb622 libb623 libb624 libb625 libb626 libb627 libb628 libb629 libb630 libb631 libb632 libb633 libb634 libb635 libb636 libb637 libb638 libb639 libb640 libb641 libb642 libb643 libb644 libb645 libb646 libb647 libb648 libb649 libb650 libb651 libb652 libb653 libb654 libb655 libb656 libb657 libb658 libb659 libb660 libb661 libb662 libb663 libb664 libb665 libb666 libb667 libb668 libb669 libb670 libb671 libb672 libb673 libb674 libb675 libb676 libb677 libb678 libb679 libb680 libb681 libb682 libb683 libb684 libb685 libb686 libb687 libb688 libb689 libb690 libb691 libb692 libb693 libb694 libb695 libb696 libb697 libb698 libb699 libb700 libb701 libb702 libb703 libb704 libb705 libb706 libb707 libb708 libb709 libb710 libb711 libb712 libb713 libb714 libb715 libb716 libb717 libb718 libb719 libb720 libb721 libb722 libb723 libb724 libb725 libb726 libb727 libb728 libb729 libb730 libb731 libb732 libb733 libb734 libb735 libb736 libb737 libb738 libb739 libb740 libb741 libb742 libb743 libb744 libb745 libb746 libb747 libb748 libb749 libb750 libb751 libb752 libb753 libb754 libb755 libb756 libb757 libb758 libb759 libb760 libb761 libb762 libb763 libb764 libb765 libb766 libb767 libb768 libb769 libb770 libb771 libb772 libb773 libb774 libb775 libb776 libb777 libb778 libb779 libb780 libb781 libb782 libb783 libb78
```



GESTORES DE PAQUETES: PROBLEMAS

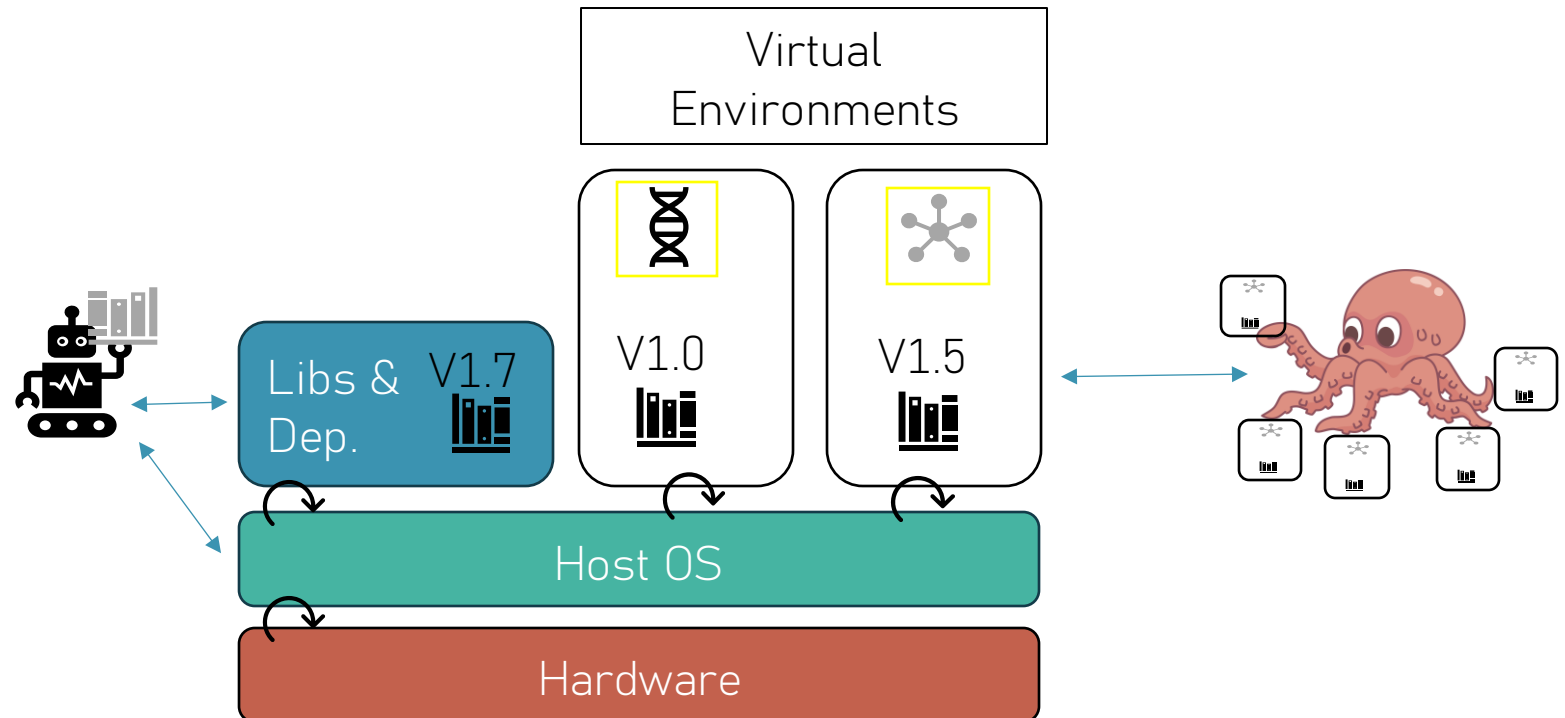
- ¡El usuario debe ser parte de los administradores!
¡SUDO!
- Siempre es arriesgado porque los programas instalados ya utilizan librerías que están configuradas
- En algunas situaciones, como cuando se utilizan HPC, no hay forma de utilizar "sudo"
- Gestionar las dependencias del paquete: a menudo, diferentes proyectos requerirán diferentes versiones del mismo software.
 - Raramente, pero no imposible, un mismo proyecto requiere diferentes versiones de un mismo programa



UTILIZAR ENTORNOS VIRTUALES TE PERMITE INSTALAR LA VERSIÓN QUE NECESITAS PARA TU PROYECTO.

- **Reproducibilidad mejorada:** Usar un entorno virtual permite un seguimiento preciso del software y sus versiones, facilitando la recreación del entorno.
- **Sistema más limpio y ligero:** Evitar múltiples instalaciones previene el caos y mantiene el sistema operativo ágil.
- ¡El usuario no necesita ser parte de los administradores! ¡SUDO!

Software que solo utilizaremos adentro el entorno virtual





Conda: Sistema de gestión de paquetes y entornos para lenguajes variados.

- **Compatibilidad:** Windows, macOS, y Linux.
- **Funcionalidad:** Instalación, ejecución y actualización ágil de paquetes y sus dependencias.
- **Versatilidad:** Originalmente para Python, pero extensible a lenguajes como R, Ruby, Java y más.
- **Gestión de entornos:** Facilita la creación, guardado y cambio entre diferentes entornos locales.
- **Flexibilidad:** Permite manejar diversas versiones de Python sin cambiar de gestor de entornos.



Mamba: Reimplementación del gestor de paquetes conda en C++ (para máxima eficiencia).

- Descarga paralela de datos de repositorio y archivos de paquetes con multi-threading.
- Uso de **libsolv** para resolución de dependencias mucho más rápida
- Al mismo tiempo, mamba utiliza el mismo analizador de línea de comandos, código de instalación y desinstalación de paquetes y rutinas de verificación de transacciones que conda para mantener la máxima compatibilidad.

INSTALACIÓN RÁPIDA DESDE LA LÍNEA DE COMANDOS

Ingresa al sitio web <https://docs.conda.io/projects/miniconda/en/latest/index.html>

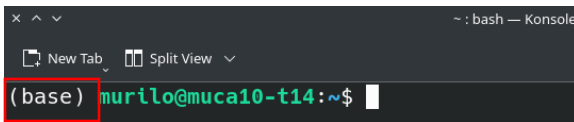
Windows macOS Linux

These four commands quickly and quietly install the latest 64-bit version of the installer and then clean up after themselves. To install a different version or architecture of Miniconda for Linux, change the name of the `.sh` installer in the `wget` command.

```
mkdir -p ~/miniconda3
wget https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh -O ~/miniconda3/miniconda.sh
bash ~/miniconda3/miniconda.sh -b -u -p ~/miniconda3
rm -rf ~/miniconda3/miniconda.sh
```

After installing, initialize your newly-installed Miniconda. The following commands initialize for bash and zsh shells:

```
~/miniconda3/bin/conda init bash
~/miniconda3/bin/conda init zsh
```



En general no es recomendable instalar nada en nuestro entorno "base"

Excepto la biblioteca mamba

De forma predeterminada, cuando abres la terminal, el entorno "base" siempre está activo. Podemos cambiar esto con:

```
conda config --set auto_activate_base false
```

USANDO A CONDA

- Para agregar un "channels" a conda:

```
conda config --add channels conda-forge  
conda config --add channels bioconda
```

- Para agregar el paquete mamba que usaremos para los demás casos:

```
conda install -c conda-forge mamba
```

- -c para es un parámetro con el nombre del channel que usaremos

```
-c conda-forge
```

- conda se instala desde "channels"

The screenshot shows the 'conda-forge' profile page. The 'Profile' section on the left includes the organization name 'conda-forge', its creation date 'Apr 11, 2015', and a description: 'A community-led collection of recipes, build infrastructure, and distributions for the conda package manager.' The 'Packages' section on the right, titled 'View all (23152)', lists several packages with their update times: lammps (a few seconds ago), findent (1 minute and a few seconds ago), pyclesperanto (3 minutes and a few seconds ago), vtkbool (9 minutes and a few seconds ago), cognite-sdk (24 minutes and a few seconds ago), azure-mgmt-containerservice (26 minutes), jupyterlab-git (26 minutes and a few seconds ago), pytest-env (39 minutes and a few seconds ago), vegafusion-python-embed (45 minutes and a few seconds ago), and numbagg (45 minutes and a few seconds ago).

The screenshot shows the 'bioconda' profile page. The 'Profile' section on the left includes the organization name 'bioconda', its creation date 'Sep 11, 2015', and a description: 'Bioconda is a distribution of bioinformatics software realized as a channel for the versatile Conda package manager.' The 'Packages' section on the right, titled 'View all (10268)', lists several packages with their update times: galaxy-objectstore (1 hour and 7 minutes ago), virheat (1 hour and 47 minutes ago), galaxy-files (1 hour and 47 minutes ago), mehari (1 hour and 48 minutes ago), melon (3 hours and 59 minutes ago), plassembler (4 hours and 4 minutes ago), vg (5 hours and 46 minutes ago), bioconda-repodata-patches (11 hours and a few seconds ago), mavis-config (12 hours and 5 minutes ago), and genomertools-genomertools (12 hours and a few seconds ago).

USANDO A CONDA

- Para crear un entorno, utilice mamba create
`mamba create -n {your-env-name}`
- Para activar/desactivar un entorno, utilice conda activate/deactivate
`conda activate {your-env-name}`
`conda deactivate`
- para instalar paquetes en un entorno activo, use mamba install
`mamba install {package-name}`
o especificando el channel
`mamba install -c conda-forge mamba`

```
{YOUR-ENV-NAME} = ENTRENO, ENTRENO2  
{CHANNEL-NAME} = BIOCONDA  
{PACKAGE-NAME} = FASTQC, MAFFT, BWA
```

- Para listar los entornos disponibles
`conda env list`
- Para remover a un entorno
`conda env remove {your-env-name}`
- Para instalar paquetes en un entorno específico no activo, agregue el indicador "-n", seguido del nombre del entorno.
`mamba install -n {YOUR-ENV-NAME} -c {CHANNEL-NAME} {PACKAGE-NAME}`
- Para desinstalar paquetes en un entorno activo
`mamba uninstall {PACKAGE-NAME}`

Practica: crear dos entornos, ver la lista de entornos disponibles, activarlos, remover el entorno 2, remover los paquetes del entorno uno menos a fastqc

USANDO A CONDA

- También podemos controlar la versión de los paquetes que instalamos

```
mamba create -n {YOUR-ENV-NAME} -c {CHANNEL-NAME} {PACKAGE-NAME}={VERSION}
```

ANACONDA.ORG

Search Anaconda.org

About Anaconda Help Download Anaconda Sign In

bioconda / packages / fastqc

A quality control tool for high throughput sequence data.

Conda Files Labels Badges

Filters

Type: All Version: All Label: All

Type	Size	Name	Version	Uploaded	Downloads	Labels
conda	11.1 MB	noarch/fastqc-0.12.0	0.12.1	7 months and 19 days ago	71164	main
conda	9.7 MB	noarch/fastqc-0.11.9	0.11.9	2 years and 6 months ago	249628	main
conda	9.6 MB	noarch/fastqc-0.11.8	0.11.8	2 years and 6 months ago	4152	main
conda	9.5 MB	noarch/fastqc-0.11.7	0.11.7	2 years and 6 months ago	3203	main
conda	9.7 MB	noarch/fastqc-0.11.6	0.11.6	3 years and 8 months ago	94670	main
conda	9.6 MB	noarch/fastqc-0.11.5	0.11.5	3 years and 11 months ago	3694	main
conda	9.6 MB	noarch/fastqc-0.11.4	0.11.4	3 years and 11 months ago	62589	main
conda	9.5 MB	noarch/fastqc-0.11.3	0.11.3	4 years and 6 days ago	2022	main
conda	9.6 MB	osx-64/fastqc-0.11.3	0.11.3	4 years and 10 months ago	6410	main cf201901
conda	9.6 MB	linux-64/fastqc-0.11.3	0.11.3	4 years and 10 months ago	44824	main cf201901
conda	9.6 MB	linux-64/fastqc-0.11.2	0.11.2	5 years and 19 days ago	16253	main cf201901
conda	9.6 MB	osx-64/fastqc-0.11.2	0.11.2	5 years and 19 days ago	2144	main cf201901
conda	9.6 MB	linux-64/fastqc-0.11.1	0.11.1	5 years and 19 days ago	2144	main cf201901
conda	9.6 MB	linux-64/fastqc-0.11.0	0.11.0	5 years and 1 month ago	12097	main cf201901
conda	9.6 MB	osx-64/fastqc-0.11.7-5.tar.bz2	0.11.7	5 years and 1 month ago	645	main cf201901
conda	507.7 kB	linux-64/fastqc-0.10.1-1.tar.bz2	0.10.1	5 years and 3 months ago	579	main cf201901
conda	9.5 MB	linux-64/fastqc-0.11.3-1.tar.bz2	0.11.3	5 years and 3 months ago	807	main cf201901

You can **install packages** when creating a new environment by specifying the package name(s) afterwards

```
mamba create -n {your-env-name} {package-name-1} {package-name-2}
```

You can install a **specific version** by adding the '=' sign, followed by the version (no spaces)

```
mamba create -n {your-env-name} {package-name}={version}
```

```
mamba install {package-name}={version}
```

You can install from a **specific channel** by using the '-c' flag, followed by the channel

```
mamba create -n {your-env-name} -c {a-channel} {package-name}
```

```
mamba install -c {a-channel} {package-name}
```

Here are a couple of examples:

```
mamba create -n my-proj -c conda-forge python=3.8 pandas
```

```
mamba install -c conda-forge python=3.8 pandas
```

¿QUÉ ES EXPORTAR ENTORNOS CONDA Y POR QUÉ HACERLO?

- Exportar un entorno crea un archivo (yaml) que lista todos los paquetes y versiones que están instalados en ese entorno.
- Este archivo se puede guardar en un repositorio de tu código de proyecto para que otros puedan crear fácilmente un entorno que ejecutará exitosamente tu código.
- También puede ser utilizado por ti para recrear el entorno (en caso de que no hayas trabajado en el proyecto durante meses).

```
{YOUR-ENV-NAME} = ENTRENO3  
{CHANNEL-NAME} = BIOCONDA  
{PACKAGE-NAME} = FASTQC, MAFFT, BWA
```

Practica: crear un entornoA, instalar paquetes, exportar este entorno y crear un entornoB con el yaml

- Para exportar un entorno activo, utilice `mamba env export > {mi-entorno}.yaml`
- Para recrear un entorno de un archivo yaml, utilice `mamba env create -f {mi-entorno}.yaml`

```
name: snakemake  
channels:  
  - bioconda  
  - conda-forge  
  - defaults  
dependencies:  
  - _libgcc_mutex=0.1=conda_forge  
  - _openmp_mutex=4.5=2_gnu  
  - aioeasywebdav=2.4.0=pyha770c72_0  
  - aiohttp=3.8.5=py311h459d7ec_0  
  - aiosignal=1.3.1=pyhd8ed1ab_0  
  - amply=0.1.6=pyhd8ed1ab_0  
  - appdirs=1.4.4=pyh9f0ad1d_0  
  - async-timeout=4.0.3=pyhd8ed1ab_0  
  - attmap=0.13.2=pyhd8ed1ab_0  
  - attrs=23.1.0=pyh71513ae_1  
  - backports=1.0=pyhd8ed1ab_3  
  - backports.functools_lru_cache=1.6.5=pyhd8ed1ab_0  
  - bcrypt=4.0.1=py311h46250e7_1  
  - boto3=1.28.55=pyhd8ed1ab_0  
  - botocore=1.31.55=pyhd8ed1ab_0
```




PARTE 2

Comprendiendo la estructura de algunos archivos y utilizando programas para los primeros pasos

MODELO DE ARCHIVO: FASTA

- Una secuencia comienza con un carácter mayor que (">") seguido o no de una descripción de la secuencia (todo en una sola línea).
- Las siguientes líneas inmediatamente después de la línea de descripción son la representación de la secuencia.
- Una letra por aminoácido o ácido nucleico y, por lo general, no tienen más de 80 caracteres de longitud.

```
>SEQUENCE_1
MTEITAAMVKELRESTGAGMMDCKNALSETNGDFDKAVQLLREKGLGKAAKKADRLAAEG
LVSVKVSDDFTIAAMRPSYLSYEDLDMTFVENEYKALVAELEKENEERRRLKDPNKPEHK
IPQFASRKQLSDAILKEAEEKIKEELKAQ GKPEKIWDNIIPGKMNSFIADNSQLDSKLT L
MGQFYVMDDKKTVEQVIAEKEKEFGGKIKIVEFICFEVGEGLEKKTEDFAAEVAAQL
```

```
>P01013 GENE X PROTEIN (OVALBUMIN-RELATED)
QIKDLLVSSSTDLDTTLLVLVNAIYFKGMWKTAFAEDTREMPFHVTKQESKPVQMMCMNNSFNVATLP AE
KMKILELPFASGDL SMLVLLPDEVSDLERIEKTINFEKLTEWTNPNTMEKRRVKVYLPQMKIEEKYNLTS
VLMALGMTDLFIPSANLTGISSAESLKISQAVHGAFMELSEDGIEMAGSTGVIEDIKHSPESQFRADHP
FLFLIKHNPTNTIVYFGRYWSP
```


MODELO DE ARCHIVO: (MULTI)FASTA

- Multifasta: un único archivo que contiene varias secuencias biológicas.
- Tenga en cuenta la falta de patrón en los nombres y descripciones de cada secuencia.
- Lo ideal sería no utilizar espacios y elegir algún carácter como separador.
 - facilita la construcción de scripts en lenguaje de programación para manipular estos datos

```
;LCBO - Prolactin precursor - Bovine
; a sample sequence in FASTA format
MDSKGSSQKGSRLLLLLVVSNLLCQGVVSTPVCNPGPGNCQVSLRDLFDRAVMVSHYIHDLS
EMFNEFDKRYAQKGKGFITMALNSCHTSSLPTPEDKEQAQQTHHEVLMSLILGLLRSWNPPLYHL
VTEVRGMKGAPDAILSR AIEIEEENKRLLEGMEMIFGQVIPGAKETEPYPVWSGLPSLQTKDED
ARYSAFYNLLHCLRRDSSKIDTYLKLLNCRIIYNNNC*

>MCHU - Calmodulin - Human, rabbit, bovine, rat, and chicken
MADQLTEEQIAEFKEAFSLFDKDGDTITTKELGTVMRSLGQNPTEAELQDMINEVDADGNGTID
FPEFLTMMARKMKD TDSEEEIREAFRVFDKDGNGYISAAELRHVMTNLGEKLTDEEVDEMIREA
DIDGDGQVNYEEFVQMMTAK*

>gi|5524211|gb|AAD44166.1| cytochrome b [Elephas maximus maximus]
LCLYTHIGRNIYYGSYLYSETWNTGIMLLLITMATAFMGYVLPWGQMSFWGATVITNLFSAIPYIGTNLV
EWIWGGFSVDKATLNRFFAFHFILPFTMVALAGVHLTFLHETGSNNPLGLTSDSDKIPFHPYYTIKDFLG
LLILILLLLLLALLSPDMLGDPDNHMPADPLNTPLHIKPEWYFLFAYAILRSVPNKLGGVLALFLSIVIL
GLMPFLHTSKHRSMMLRPLSQALFWTLTMDLLTLTWIGSQPVEYPYTIIGQMASILYFSIILAFPLIAGX
IENY
```

MODELO DE ARCHIVO: FASTA

código de letras

convenciones para facilitar
la identificación

A	adenosine	C	cytidine	G	guanine
T	thymidine	N	A/G/C/T (any)	U	uridine
K	G/T (keto)	S	G/C (strong)	Y	T/C (pyrimidine)
M	A/C (amino)	W	A/T (weak)	R	G/A (purine)
B	G/T/C	D	G/A/T	H	A/C/T
V	G/C/A	-	gap of indeterminate length		

A	alanine	P	proline
B	aspartate/asparagine	Q	glutamine
C	cystine	R	arginine
D	aspartate	S	serine
E	glutamate	T	threonine
F	phenylalanine	U	selenocysteine
G	glycine	V	valine
H	histidine	W	tryptophan
I	isoleucine	Y	tyrosine
K	lysine	Z	glutamate/glutamine
L	leucine	X	any
M	methionine	*	translation stop
N	asparagine	-	gap of indeterminate length

Extension	Meaning	Notes
fasta, fa	generic FASTA	Any generic fasta file. See below for other common FASTA file extensions
fna	FASTA nucleic acid	Used generically to specify nucleic acids.
ffn	FASTA nucleotide of gene regions	Contains coding regions for a genome.
faa	FASTA amino acid	Contains amino acid sequences. A multiple protein fasta file can have the more specific extension mpfa.
frn	FASTA <u>non-coding RNA</u>	Contains non-coding RNA regions for a genome, in DNA alphabet e.g. tRNA, rRNA

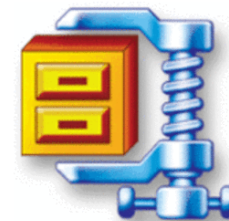
MODELO DE ARCHIVO: FASTQ

Un archivo FASTQ tiene cuatro líneas con cuatro campos separados por secuencia:

- El campo 1 comienza con un carácter '@' y va seguido de un identificador de secuencia.
- El campo 2 son las letras de secuencia.
- El campo 3 contiene un carácter '+'.
+
! ' ' * (((* * * +)) % % % + +) (% % % %) . 1 * * * - + * ' ')) * * 5 5 C C F > > > > > C C C C C C C C 6 5
- El campo 4 codifica los valores de calidad para la secuencia en el campo 2 y debe contener la misma cantidad de símbolos que letras en la secuencia.

```
@SEQ_ID
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
! ' ' * ( ( ( * * * + ) ) % % % + + ) ( % % % % ) . 1 * * * - + * ' ' ) ) * * 5 5 C C F > > > > > C C C C C C C C 6 5
```

En general, los archivos FASTQ pueden contener desde miles hasta millones de lecturas.



*fastq.gz

MODELO DE ARCHIVO: FASTQ

Element	Requirements	Description
@	@	Each sequence identifier line starts with @
<instrument>	Characters allowed: a-z, A-Z, 0-9 and underscore	Instrument ID
<run number>	Numerical	Run number on instrument
<flowcell ID>	Characters allowed: a-z, A-Z, 0-9	Flowcell ID
<lane>	Numerical	Lane number
<tile>	Numerical	Tile number
<x_pos>	Numerical	X coordinate of cluster
<y_pos>	Numerical	Y coordinate of cluster
<read>	Numerical	Read number. 1 can be single read or Read 2 of paired-end
<is filtered>	Y or N	Y if the read is filtered (did not pass), N otherwise
<control number>	Numerical	0 when none of the control bits are on, otherwise it is an even number. On HiSeq X systems, control specification is not performed and this number is always 0.
<sample number>	Numerical	Sample number from sample sheet

identificador de una lectura
Illumina

@<instrument>:<run number>:<flowcell ID>:<lane>:<tile>:<x-pos>:<y-pos> <read>:<is filtered>:<control number>:<sample number>

Identifier — | @HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Sequence — | TTAATTGGTAAATAAAATCTCCTAATAGCTTAGATNTTACCTTNNNNNNNNNTAGTTTCTTGAGA
+ sign & identifier — | +HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Quality scores — | efcfffffcfeefffcfffffdff`feed]`_]_Ba^__[YBBBBBBBBBRTT\]][] dddd`
Base T
phred Quality] = 29

¡Esta información es
utilizada posteriormente por
otros programas!

DESCARGAR LOS ARCHIVOS DE SECUENCIACIÓN EJEMPLO

- `wget -q -O - "https://trace.ncbi.nlm.nih.gov/Traces/sra-reads-be/fastq?acc=SRR24053947" >
SRR24053947.fastq.gz`
- `wget -q -O - "https://trace.ncbi.nlm.nih.gov/Traces/sra-reads-be/fastq?acc=SRR24054051" >
SRR24054051.fastq.gz`
- `wget -q -O - "https://trace.ncbi.nlm.nih.gov/Traces/sra-reads-be/fastq?acc=SRR24053926" >
SRR24053926.fastq.gz`

USAREMOS EN
DOS CLASES

ALGUNAS COMBINACIONES DE COMANDOS INTERESANTES PARA EVALUAR RÁPIDAMENTE ARCHIVOS FASTA Y FASTQ

Visualizar contenido de archivos

a. Usar cat para visualizar el contenido de ejemplo.fasta

cat ejemplo.fasta

b. Usar zcat para visualizar el contenido de ejemplo.fastq.gz

zcat ejemplo.fastq.gz

¿Qué pasa si intentamos lo siguiente?

head -n 4 ejemplo.fastq.gz

Visualizar las primeras y últimas líneas de un archivo

a. Usar head para ver las primeras 4 líneas de ejemplo.fastq

head -n 4 ejemplo.fasta

Intenta:

zcat ejemplo.fastq.gz | head -n 4

b. Usar tail para ver las últimas 4 líneas de ejemplo.fastq

tail -n 4 ejemplo.fasta

ALGUNAS COMBINACIONES DE COMANDOS INTERESANTES PARA EVALUAR RÁPIDAMENTE ARCHIVOS FASTA Y FASTQ

Contar cantidad de secuencias en un archivo fasta

Podemos contar cuántas hay:

```
grep -c '^>' ejemplo.fasta
```

O podemos visualizar todas las líneas de encabezamiento

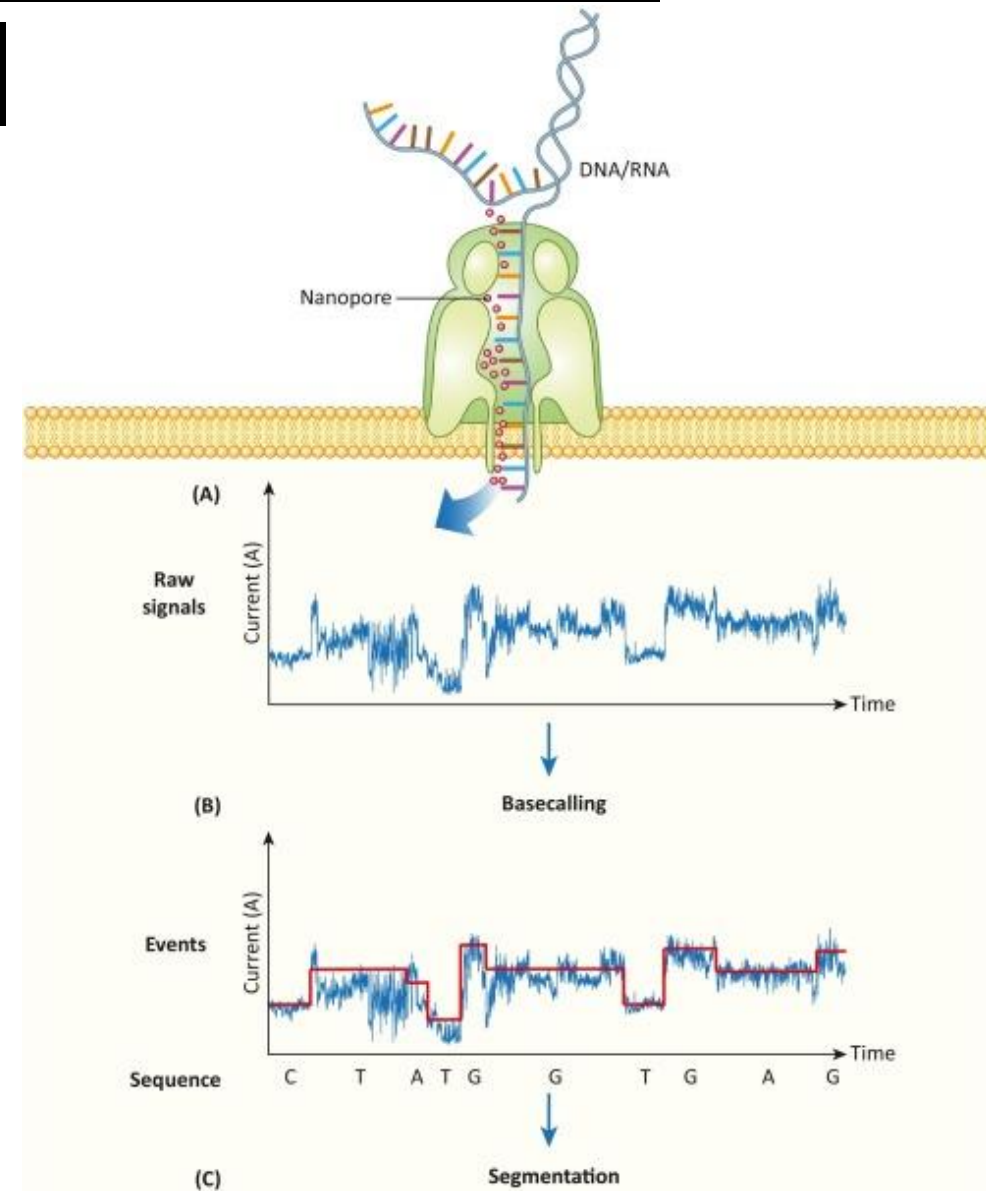
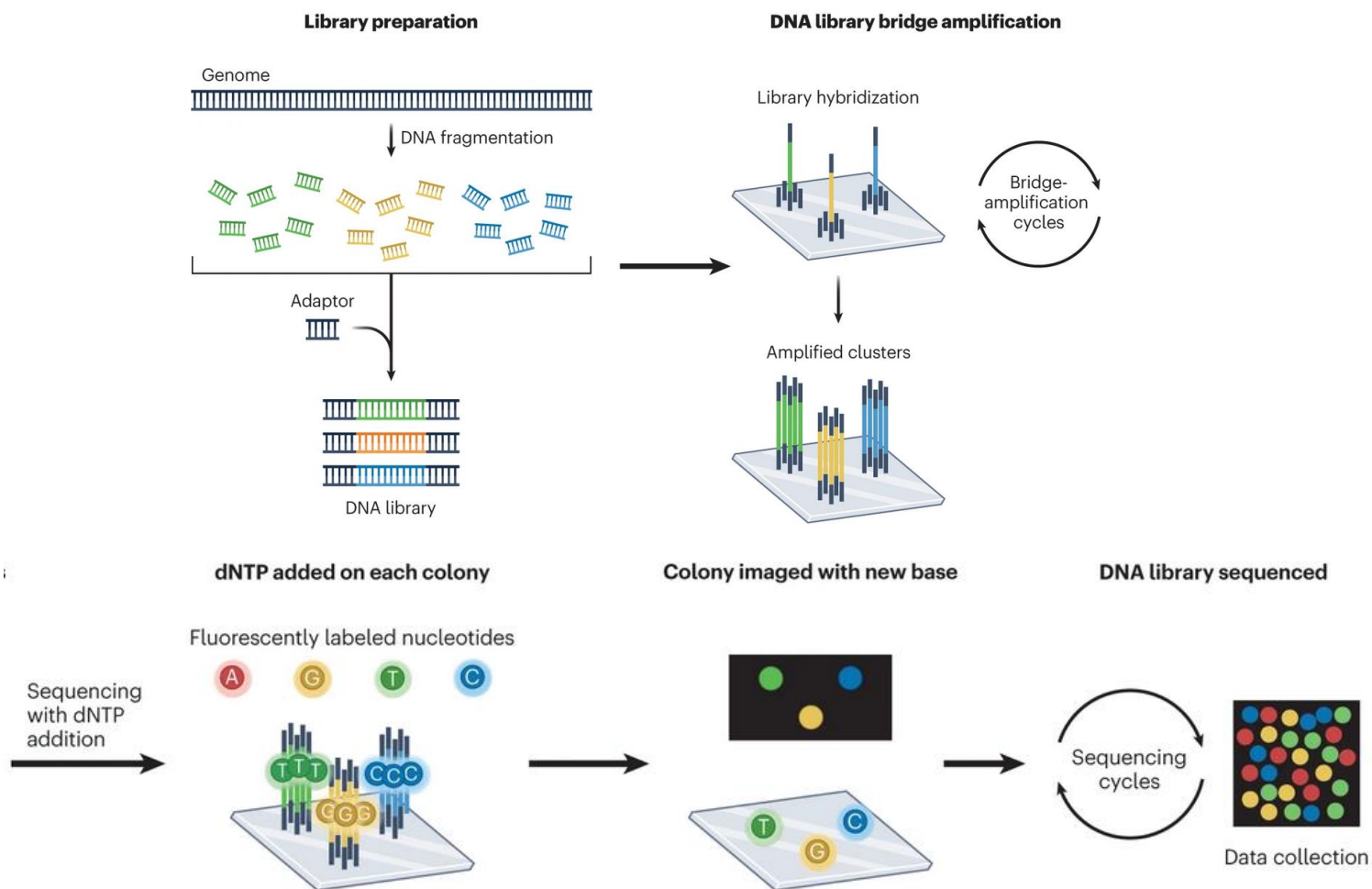
```
grep '^>' ejemplo.fasta
```

De manera análoga

a. Usar head para ver las primeras 4 líneas de ejemplo.fastq

```
zcat ejemplo.fastq.gz | grep -c '^@'
```

¿DE DÓNDE PROVIENE LA INFORMACIÓN DE CALIDAD DE LAS 'LECTURAS'?

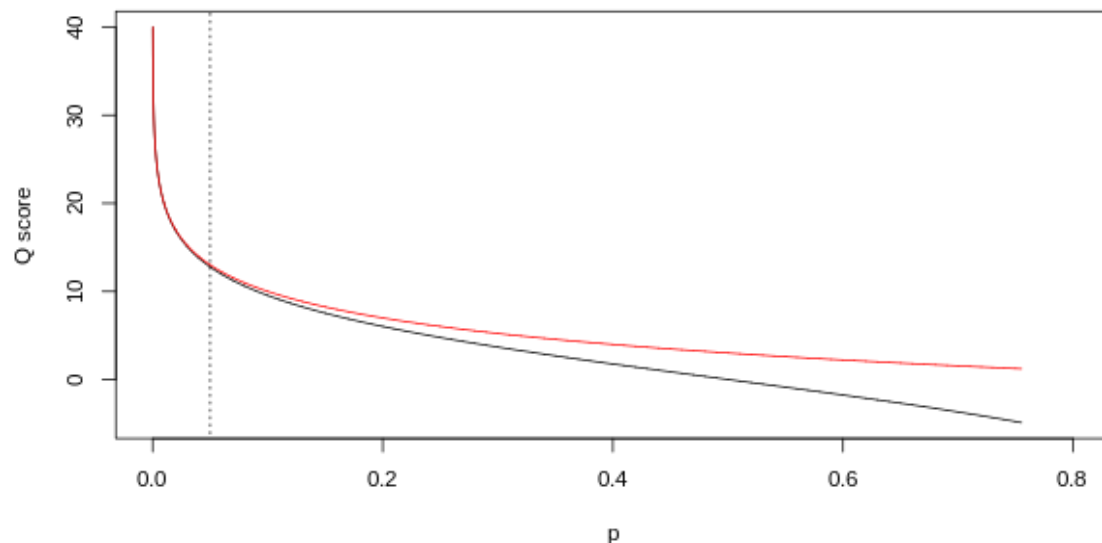


LIDIAR CON LA INCERTIDUMBRE

Dada una afirmación, A, el puntaje de calidad, Q(A), expresa la probabilidad de que A no sea cierta, P(~A), según la relación:

$$Q(A) = -10 \log_{10}(P(\sim A))$$

donde P(~A) es la probabilidad estimada de que una afirmación A sea incorrecta.



Quality score, Q(A)	Error probability, P(~A)
10	0.1
20	0.01
30	0.001

Symbol	ASCII Code	Q-Score			
!	33	0	?	63	30
"	34	1	@	64	31
#	35	2	A	65	32
\$	36	3	B	66	33
%	37	4	C	67	34
&	38	5	D	68	35
'	39	6	E	69	36
(40	7	F	70	37
)	41	8	G	71	38
*	42	9	H	72	39
+	43	10	I	73	40

FASTQC

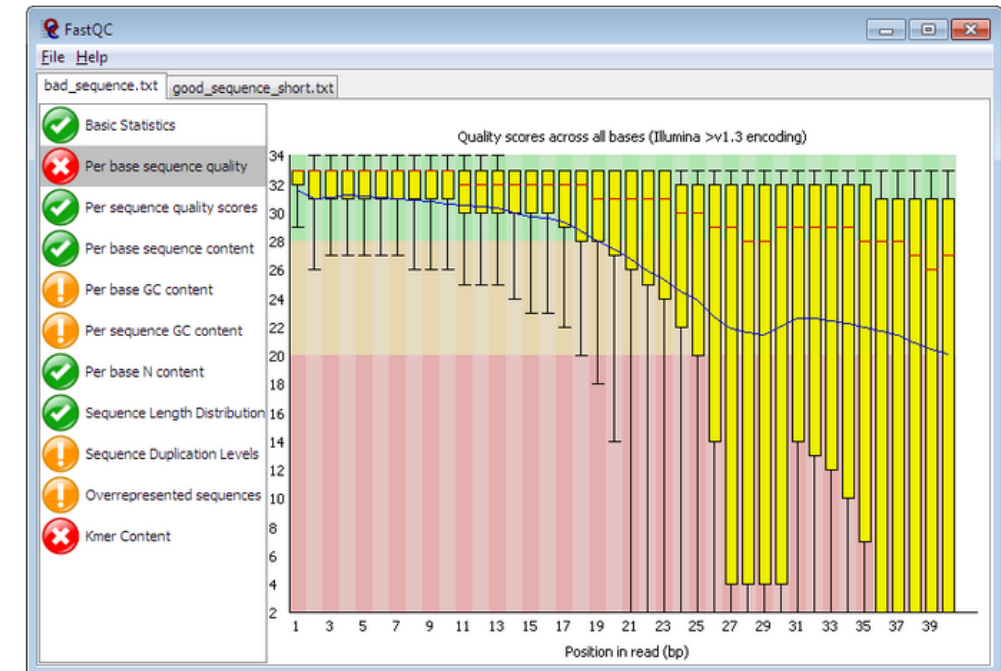
FastQC

- FastQC tiene como objetivo proporcionar una manera sencilla de realizar algunas verificaciones de control de calidad en datos de secuencia crudos provenientes de pipelines de secuenciación de alto rendimiento.
- Ofrece un conjunto modular de análisis que puedes usar para obtener una rápida impresión de si tus datos tienen algún problema del cual deberías estar al tanto antes de realizar cualquier análisis adicional.

Function	A quality control tool for high throughput sequence data.
Language	Java
Requirements	A suitable Java Runtime Environment The Picard BAM/SAM Libraries (included in download)
Code Maturity	Stable. Mature code, but feedback is appreciated.
Code Released	Yes, under GPL v3 or later .
Initial Contact	Simon Andrews

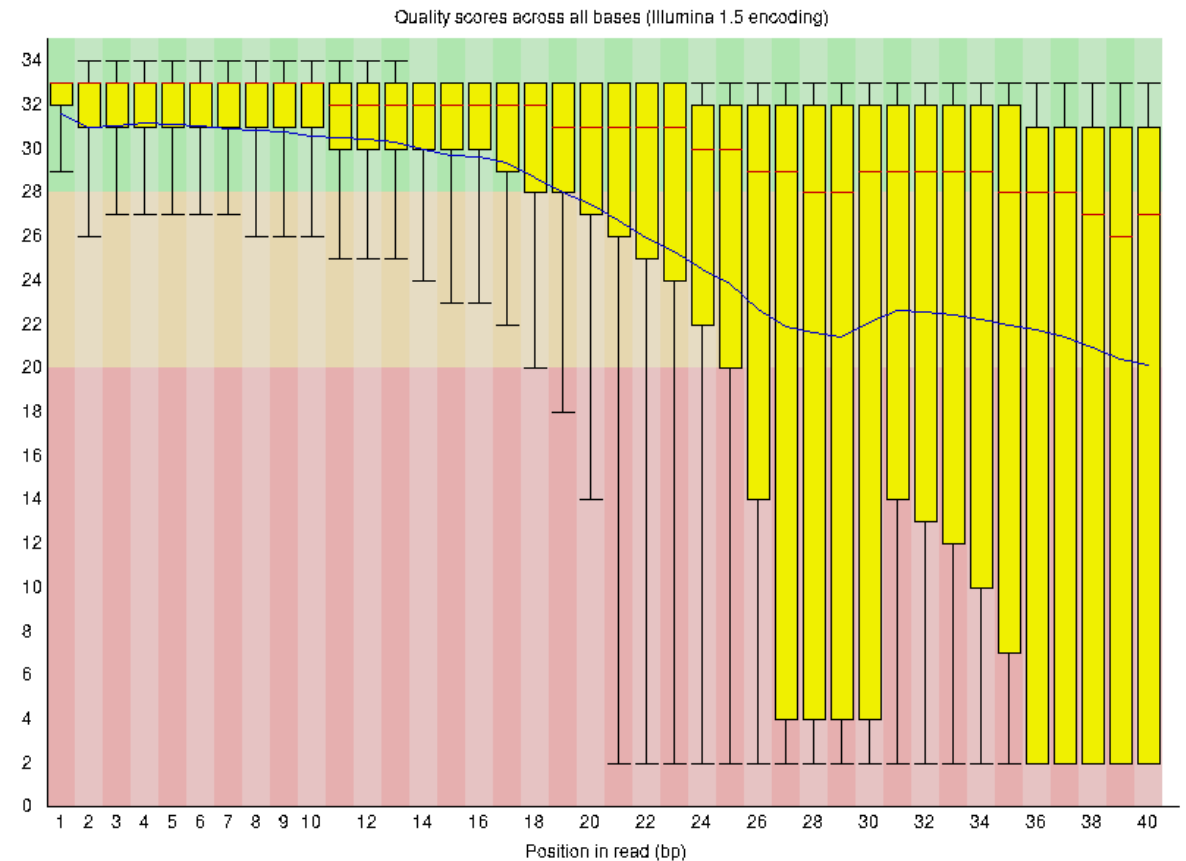
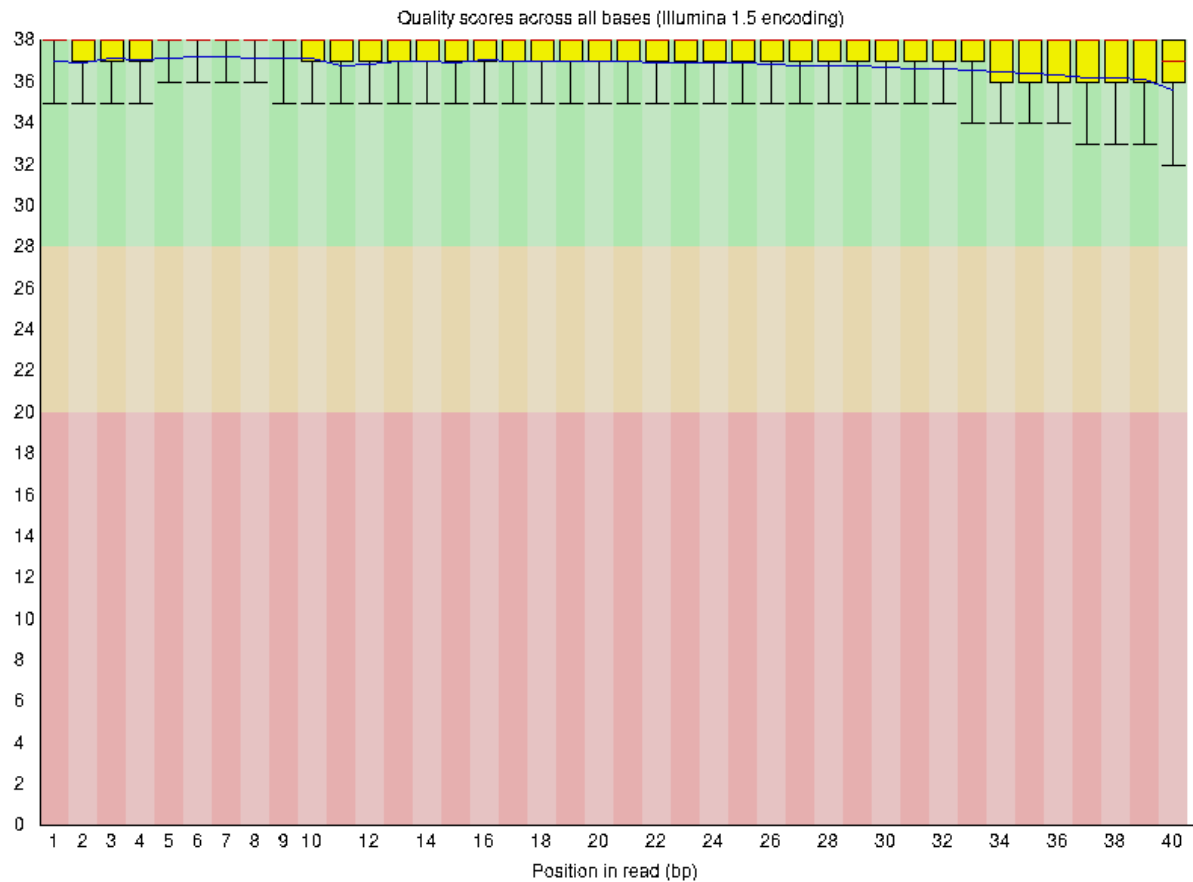
[Download Now](#)

¡12 conjuntos de estadísticas en informes gráficos!



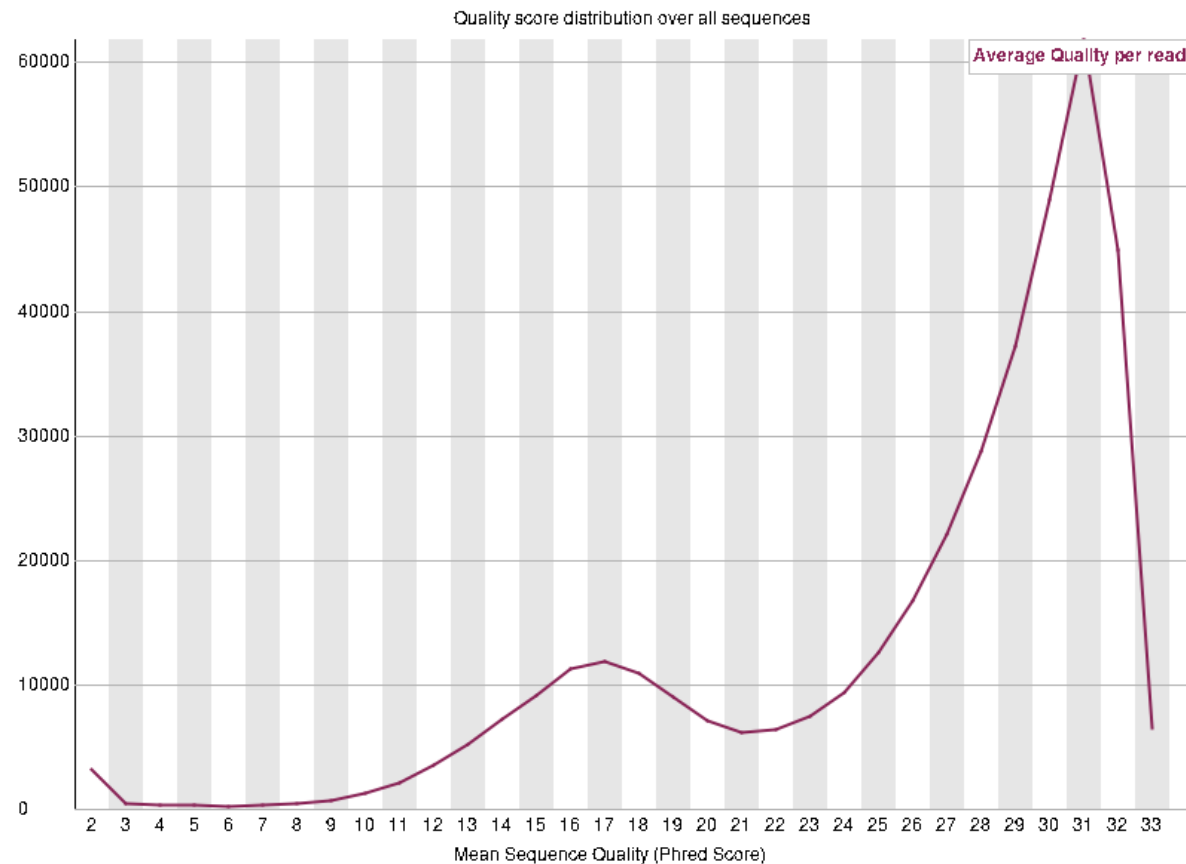
FASTQC

Calidad de la secuencia por base: Muestra la distribución de scores de calidad a través de todas las posiciones en las lecturas. Se espera que la calidad disminuya hacia el final de las lecturas, pero grandes desviaciones pueden indicar problemas.



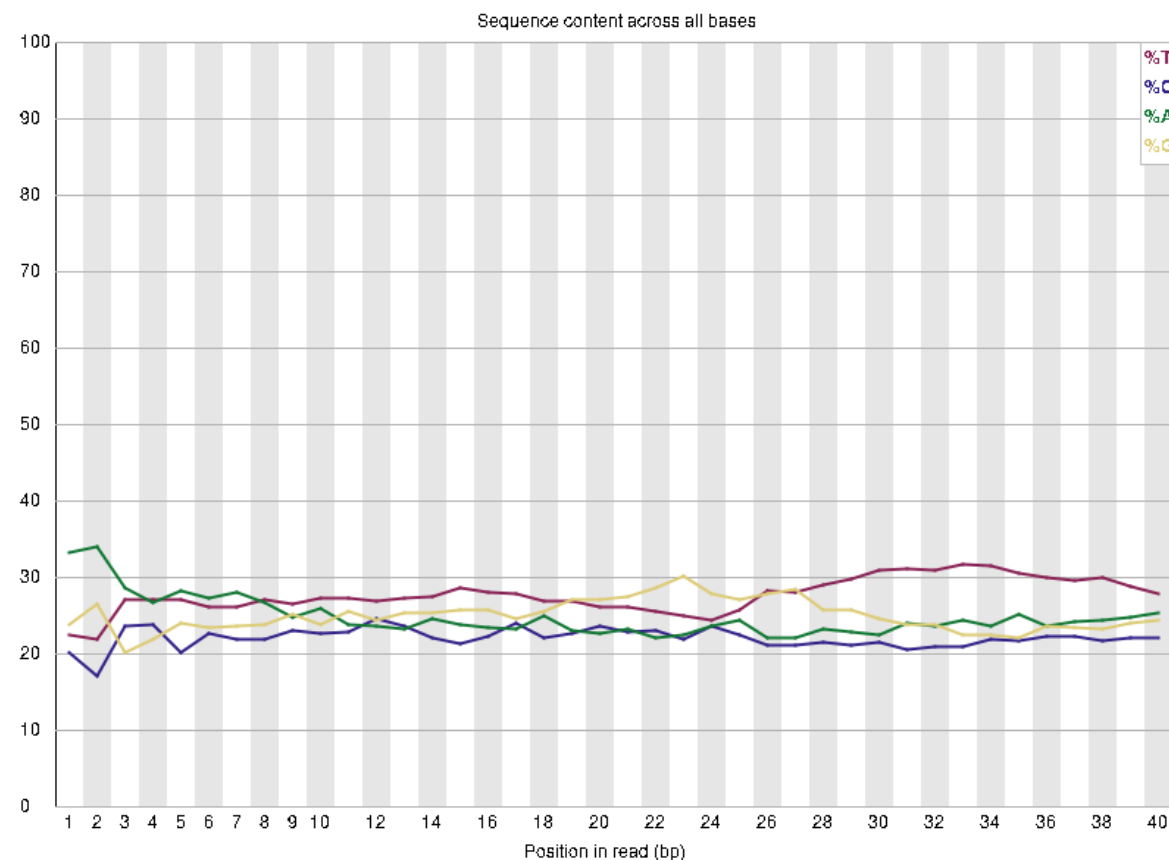
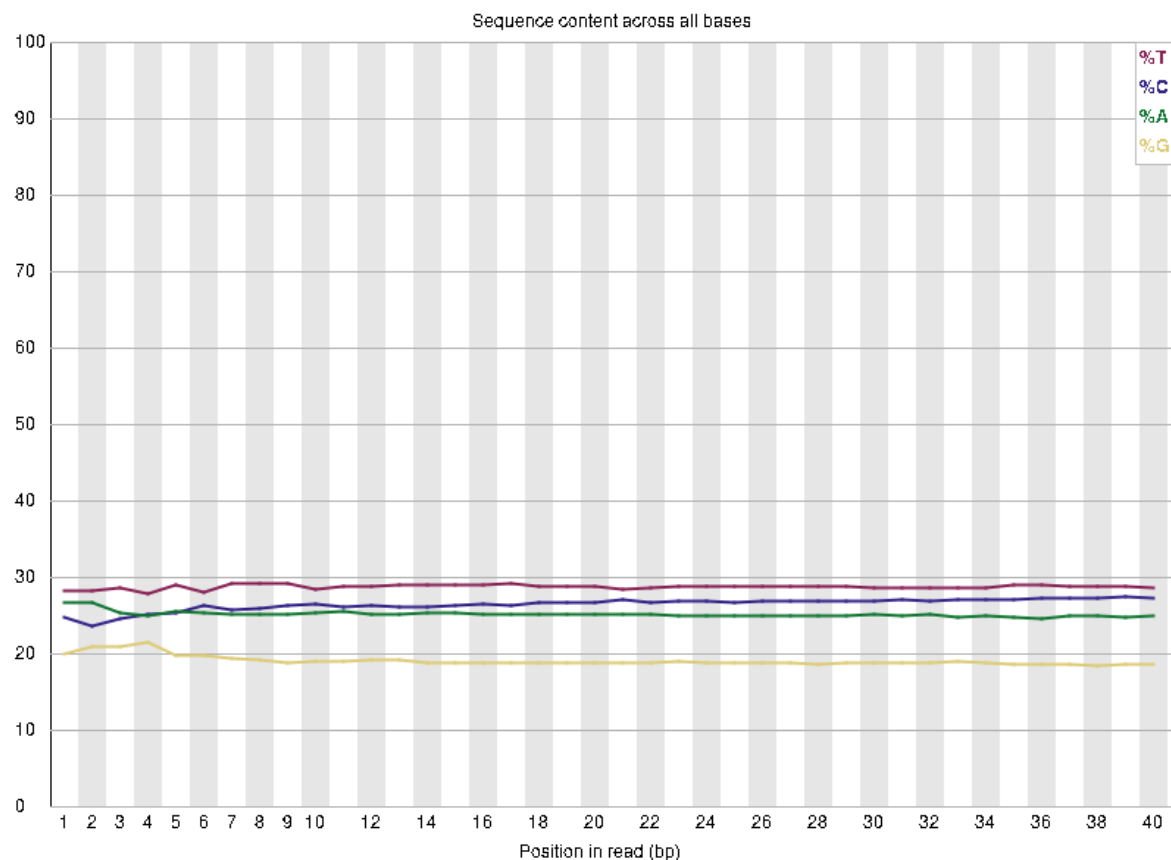
FASTQC

Distribución de calidad de la secuencia: Proporciona una visión general de la calidad de todo el archivo, no solo por base. Se esperaría ver la mayoría de las lecturas con una calidad alta.



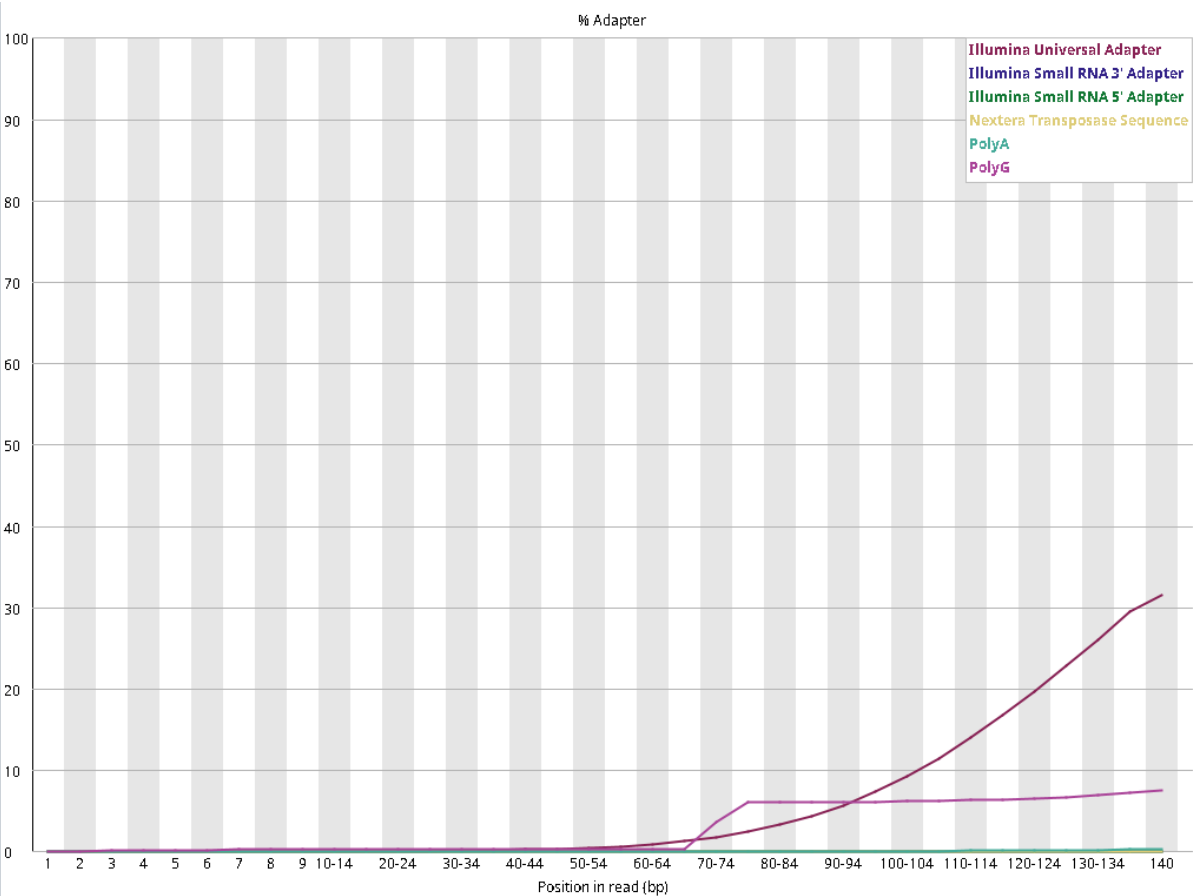
FASTQC

Contenido de base por base: En una biblioteca sin sesgo, el contenido de base debería ser constante. Variaciones pueden deberse a contaminación con adaptadores o a otros tipos de sesgo.



FASTQC

Contaminación por adaptadores: Si no se han eliminado los adaptadores, FastQC identificará las secuencias de adaptadores más comunes.



Contaminación por sobrerepresentación de secuencias: Las secuencias que aparecen con mucha frecuencia pueden ser indicativas de contaminación o secuencias altamente conservadas.

! Overrepresented sequences

Sequence	Count	Percentage	Possible Source
AGAGTTTTCGCTTCCATGACGCAGAAGTTAACACTTTC	2065	0.5224039181558763	No Hit
GATTGGCGTATCCAACCTGCAGAGTTTTATCGCTTCCATG	2047	0.5178502762542754	No Hit
ATTGGCGTATCCAACCTGCAGAGTTTTATCGCTTCCATGA	2014	0.5095019327680071	No Hit
CGATAAAATGATTGGCGTATCCAACCTGCAGAGTTTTAT	1913	0.4839509420979134	No Hit
GTATCCAACCTGCAGAGTTTTATCGCTTCCATGACGCAGA	1879	0.4753496185060066	No Hit
AAAAATGATTGGCGTATCCAACCTGCAGAGTTTTATCGCT	1846	0.4670012750197325	No Hit
TGATTGGCGTATCCAACCTGCAGAGTTTTATCGCTTCCAT	1841	0.46573637449150995	No Hit
AACCTGCAGAGTTTTATCGCTTCCATGACGCAGAAGTTAA	1836	0.46447147396328753	No Hit
GATAAAATGATTGGCGTATCCAACCTGCAGAGTTTTATC	1831	0.4632065734350651	No Hit
ATGATTGGCGTATCCAACCTGCAGAGTTTTATCGCTTCCA	1779	0.45005160794155147	No Hit
AAATGATTGGCGTATCCAACCTGCAGAGTTTTATCGCTTC	1779	0.45005160794155147	No Hit
AATGATTGGCGTATCCAACCTGCAGAGTTTTATCGCTTCC	1760	0.4452449859343061	No Hit
AAAATGATTGGCGTATCCAACCTGCAGAGTTTTATCGCTT	1729	0.4374026026593269	No Hit
CGTATCCAACCTGCAGAGTTTTATCGCTTCCATGACGCAG	1713	0.43335492096901496	No Hit
CGGTTTCAGCAGGAATGCCGAGATCGGAAGAGCGGTTTCAGC	599	0.15153508328105078	Illumina Paired End PCR Primer 2 (96% over 25bp)
TCTGCAGGTTGGATACGCCAATCATTTTTATCGAAGCGCG	585	0.1479933618020279	No Hit
CGCTTAAAGCTACCAAGTTATATGGCTGGGGGTTTTTTTTT	552	0.13964501831575965	No Hit
CTCTGCAGGTTGGATACGCCAATCATTTTTATCGAAGCGCG	532	0.1345854162028698	No Hit
CTGCGTCATGGAAGCGATAAACTCTGCAGGTTGGATACG	515	0.13028475440691342	No Hit
CTGCAGGTTGGATACGCCAATCATTTTTATCGAAGCGCGC	505	0.12775495335046852	No Hit
GCTTAAAGCTACCAAGTTATATGGCTGGGGGTTTTTTTTT	411	0.10397482341988626	No Hit

FASTQC

- Basado en GUI
Abra la terminal, active el entorno virtual que creamos con fastqc y escriba fastqc para iniciar el programa. En "archivo" -> "abrir..."
- Basado en línea de comando
Vaya al directorio donde se encuentra un archivo fastqc (use cd, ls y pwd para navegar si es necesario).
Escriba fastqc --help y vea las opciones para ejecutar el programa a través de la línea de comando.

Intentar:

```
fastqc --threads 8 "archivo.fastq.gz"
```

MODELO DE ARCHIVO: SAM/BAM

- El **S**equence **A**lignment/**M**ap (**SAM**) es un formato de archivo para guardar información de alineación de lecturas cortas mapeadas contra secuencias de referencia.
- Normalmente comienza con una sección de encabezado seguida de información de alineación en líneas separadas por tabuladores para cada lectura.

Header section

```
@HD    VN:1.3    SO:coordinate
@SQ    SN:contigA    LN:443
@SQ    SN:contigB    LN:1493
@SQ    SN:contigC    LN:328
```

Tab-delimited read alignment information lines

read

```
readID43GYAX15:7:1:1202:19894/1 256 contig43 613960 1 65M * 0 0
CCAGCGCGAACGAAATCCGCATGCGTCTGGTCGTTGCACGGAACGGCGCGGTGTGATGCACGGC
EDDEEDEE=EE?DE??DDDBADEBEFFFDDBEFFEBCBC=?BEEEE@=:?:?:7?:8-6?7?@??# AS:i:0 XS:i:0 XN:i:0 XM:i:0 XO:i:0 XG:i:0 NM:i:0
MD:Z:65 YT:Z:UU
```

read

```
readID43GYAX15:7:1:1202:19894/1 272 contig32 21001 1 65M * 0 0
GCCGGACGTACACGGCCGCCGGCTCTACGACCAGACGCATGCGGATTTCTAGAGCCGG
#??@?7?6-8:???:?:?=@EEEEB?=CBCBEFFEBDFFEDEDABDDDD?ED?EE=EEDEEDDE AS:i:-5 XS:i:0 XN:i:0 XM:i:1 XO:i:0 XG:i:0
NM:i:1 MD:Z:42T22 YT:Z:UU
```

read

```
readID43GYAX15:7:1:1202:19894/1 256 contig87 540849 1 65M * 0 0
CCTGCACGAACGAAATCCGCATGCGTCTGGTCGTTGTACGGAACGGCGTTGTGTGACGAACGGC
EDDEEDEE=EE?DE??DDDBADEBEFFFDDBEFFEBCBC=?BEEEE@=:?:?:7?:8-6?7?@??# AS:i:0 XS:i:0 XN:i:0 XM:i:0 XO:i:0 XG:i:0 NM:i:0
MD:Z:65 YT:Z:UU
```

MODELO DE ARCHIVO: SAM

```
Coor      12345678901234  5678901234567890123456789012345
ref       AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

+r001/1      TTAGATAAAGGATA*CTG
+r002        aaaAGATAA*GGATA
+r003        gcctaAGCTAA
+r004                        ATAGCT.....TCAGC
-r003                        ttagctTAGGC
-r001/2                        CAGCGGCAT
```

@HD VN:1.6 SO:coordinate

@SQ SN:ref LN:45

```
r001    99 ref  7 30 8M2I4M1D3M = 37  39 TTAGATAAAGGATACTG *
r002     0 ref  9 30 3S6M1P1I4M * 0   0 AAAAGATAAGGATA   *
r003     0 ref  9 30 5S6M          * 0   0 GCCTAAGCTAA     * SA:Z:ref,29,-,6H5M,17,0;
r004     0 ref 16 30 6M14N5M       * 0   0 ATAGCTTCAGC     *
r003 2064 ref 29 17 6H5M          * 0   0 TAGGC           * SA:Z:ref,9,+,5S6M,30,1;
r001  147 ref 37 30 9M            =  7 -39 CAGCGGCAT      * NM:i:1
```


MODELO DE ARCHIVO: SAM



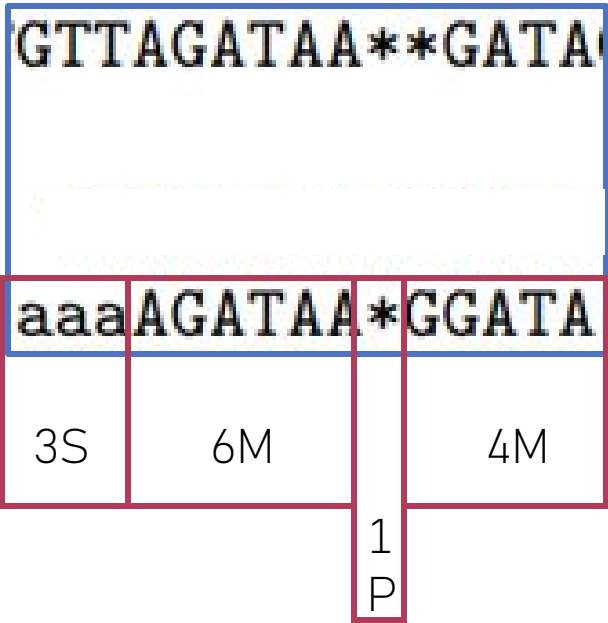
Coord
ref
+r001/1
+r002
+r003
+r004
-r003
-r001/2

12345678901234 5678901234567890123456789012345
AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

TTAGATAAAGGATA*CTG
aaaAGATAA*GGATA
gcctaAGCTAA

ATAGCT.....TCAGC
ttagctTAGGC

CAGCGGCAT



@HD VN:1.6 SO:coordinate
@SQ SN:ref LN:45

Read	Start	End	Score	Ref	Seq	Flags
r001	99	ref	7 30	8M2I4M1D3M	= 37 39 TTAGATAAAGGATACTG	*
r002	0	ref	9 30	3S6M1P1I4M	* 0 0 AAAAGATAAGGATA	*
r003	0	ref	9 30	5S6M	* 0 0 GCCTAAGCTAA	* SA:Z:ref,29,-,6H5M,17,0;
r004	0	ref	16 30	6M14N5M	* 0 0 ATAGCTTCAGC	*
r003	2064	ref	29 17	6H5M	* 0 0 TAGGC	* SA:Z:ref,9,+,5S6M,30,1;
r001	147	ref	37 30	9M	= 7 -39 CAGCGGCAT	* NM:i:1

MODELO DE ARCHIVO: SAM

Las líneas de **encabezado** siempre comenzarán con un símbolo "@" seguido de un identificador que indica el tipo y subtipo de la línea de encabezado. Algunos de los ejemplos más comunes pueden verse de la siguiente manera:



@SQ lleva cada una de las secuencias del conjunto de referencia, en el que mapeamos las lecturas

```
@SQ SN:chr14 LN:107349540
```

```
@PG ID:bwa PN:bwa VN:0.7.7-r441 CL:bwa mem ref/seq.fa r1.fastq r2.fastq
```



@PG describe el programa utilizado para generar el archivo SAM (y por tanto, realizar la alineación). Si se combinan varios archivos SAM, es posible que aparezcan varias líneas @PG, aunque es común tener solo una.

MODELO DE ARCHIVO: SAM

Información de alineación en líneas separadas por tabuladores para cada lectura.

1	2	3	4	5	6	7	8	9	10	11
SRR067577.2766	99	chr14	73240003	60	101M	=	73240004	102	GCTA...	FHG@...

Col	Field	Type	Regexp/Range	Brief description
1	QNAME	String	[!~?A~]{1,254}	Query template NAME
2	FLAG	Int	[0, 2 ¹⁶ - 1]	bitwise FLAG
3	RNAME	String	* [:rname:^*=] [:rname:]*	Reference sequence NAME ¹¹
4	POS	Int	[0, 2 ³¹ - 1]	1-based leftmost mapping POSition
5	MAPQ	Int	[0, 2 ⁸ - 1]	MAPping Quality
6	CIGAR	String	* ([0-9]+[MIDNSHPX=])+	CIGAR string
7	RNEXT	String	* = [:rname:^*=] [:rname:]*	Reference name of the mate/next read
8	PNEXT	Int	[0, 2 ³¹ - 1]	Position of the mate/next read
9	TLEN	Int	[-2 ³¹ + 1, 2 ³¹ - 1]	observed Template LENgth
10	SEQ	String	* [A-Za-z=.]+	segment SEquence
11	QUAL	String	[!~]+	ASCII of Phred-scaled base QUALity+33

Bit	Description
1	0x1 template having multiple segments in sequencing
2	0x2 each segment properly aligned according to the aligner
4	0x4 segment unmapped
8	0x8 next segment in the template unmapped
16	0x10 SEQ being reverse complemented
32	0x20 SEQ of the next segment in the template being reverse complemented
64	0x40 the first segment in the template
128	0x80 the last segment in the template
256	0x100 secondary alignment
512	0x200 not passing filters, such as platform/vendor quality controls
1024	0x400 PCR or optical duplicate
2048	0x800 supplementary alignment

Op	BAM	Description
M	0	alignment match (can be a sequence match or mismatch)
I	1	insertion to the reference
D	2	deletion from the reference
N	3	skipped region from the reference
S	4	soft clipping (clipped sequences present in SEQ)
H	5	hard clipping (clipped sequences NOT present in SEQ)
P	6	padding (silent deletion from padded reference)
=	7	sequence match
X	8	sequence mismatch

MODELO DE ARCHIVO: NEXUS

```
#NEXUS
begin < blockname >;
    < command > < argument > [additional argument];
    [ < another command with args >; ]
end;
[ < another block with commands > ]
```

- Bloque TAXA
Contiene información sobre taxones.
- Bloque DATA
Contiene la matriz de datos (por ejemplo, alineamiento de secuencias).
- Bloque TREES
Contiene árboles filogenéticos descritos usando el formato Newick, por ejemplo, ((A,B),C);:

- Un archivo NEXUS está compuesto por un encabezado fijo **#NEXUS** seguido de múltiples bloques.
- Cada bloque comienza con **BEGIN block_name;** y termina con **END;**.
- Las palabras clave no distinguen entre mayúsculas y minúsculas.
- Los comentarios están encerrados entre corchetes [...]

```
#NEXUS
Begin TAXA;
    Dimensions ntax=4;
    TaxLabels SpaceDog SpaceCat SpaceOrc SpaceElf;
End;

Begin data;
    Dimensions nchar=15;
    Format datatype=dna missing=? gap=- matchchar=.;
    Matrix
    [ "matchchar" means that it is the same as the first entry at the same position. ]
    SpaceDog   atgctagctagctcg
    SpaceCat   .....??...-a.
    SpaceOrc   ...t.....-g. [ same as atgtagctag-tgg ]
    SpaceElf   ...t.....-a.
    ;
End;

BEGIN TREES;
    Tree tree1 = (((SpaceDog,SpaceCat),SpaceOrc,SpaceElf));
END;
```

MODELO DE ARCHIVO: NEXUS

Puedes tener más o menos bloques, con tipos de datos diferentes, dependiendo de cómo una determinada aplicación requiera la entrada de datos.

```
#NEXUS
Begin TAXA;
  Dimensions ntax=4;
  TaxLabels SpaceDog SpaceCat SpaceOrc SpaceElf;
End;

Begin data;
  Dimensions nchar=15;
  Format datatype=dna missing=? gap=- matchchar=.;
  Matrix
    [ "matchchar" means that it is the same as the first entry at the same position. ]
    SpaceDog   atgctagctagctcg
    SpaceCat   .....??...-.a.
    SpaceOrc   ...t.....-.g. [ same as atgttagctag-tgg ]
    SpaceElf   ...t.....-.a.
  ;
End;

BEGIN TREES;
  Tree tree1 = (((SpaceDog,SpaceCat),SpaceOrc,SpaceElf));
END;
```

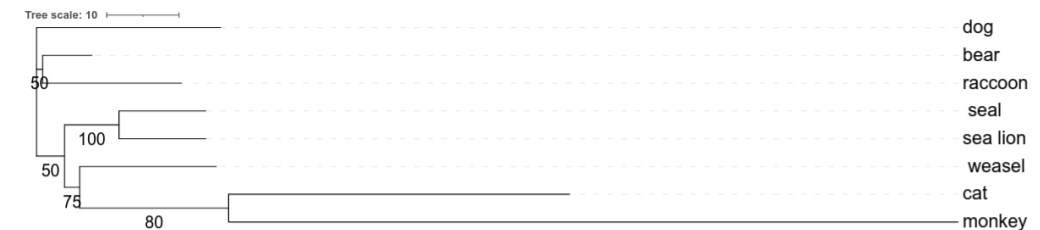
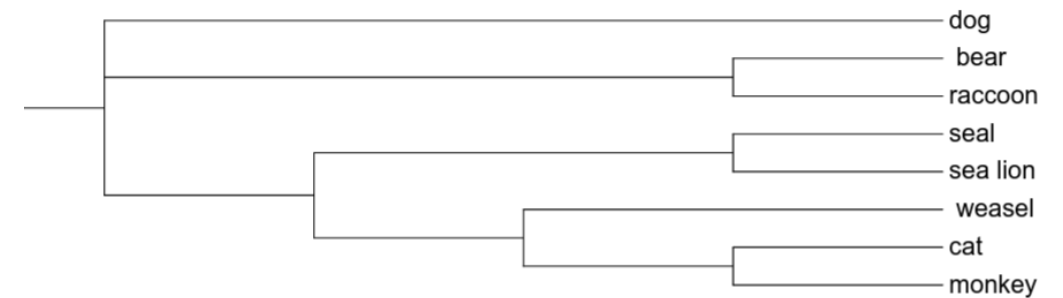
MODELO DE ARCHIVO: NEWICK

- NEWICK es un formato simple que se utiliza para escribir árboles en un archivo de texto.
- Este es un formato difícil de leer para los humanos, pero muy útil para intercambiar árboles entre diferentes tipos de software.
- Siempre se necesita punto y coma para finalizar el árbol).

```
((raccoon, bear), ((sea_lion, seal), ((monkey, cat), weasel)), dog);
```

```
((raccoon:19.19959,bear:6.80041):0.84600, ((sea_lion:11.99700, seal:12.00300):7.52973, ((monkey:100.85930, cat:47.14069):20.59201, weasel:18.87953):2.09460):3.87382, dog:25.46154);
```

```
((raccoon:19.19959,bear:6.80041)50:0.84600, ((sea_lion:11.99700, seal:12.00300)100:7.52973, ((monkey:100.85930, cat:47.14069)80:20.59201, weasel:18.87953)75:2.09460)50:3.87382, dog:25.46154);
```



¡VAMOS A PRACTICAR!

