# 생명정보학에서 쓰이는
# 컴퓨터 Perl 언어의 기초 교육
## Raw sequences of NGS technologies and Perl program

**2017/9/16**
**Jongsun Park, Ph. D.**

성신여자대학교
SUNGSHIN UNIVERSITY

서울시 강서구 화곡동 359-63 영주빌딩 201호
Tel 02 2698 1188 / Fax 02 6280 8821

www.infoboss.co.kr

**InfoBoss**
Infomations & Systems

- We will reconstruct rice (*Oryza sativa*) genomes based on raw data of 3000 rice genomes.

ⓘ gigadb.org/dataset/200001

Revolutionizing data dissemination, organization, and use

**Search**

Data released on May 27, 2014

**The Rice 3000 Genomes Project Data.**

**The 3000 Rice Genomes Project**, (2014): The Rice 3000 Genomes Project Data. GigaScience Database.
http://dx.doi.org/10.5524/200001 RIS BIBTEX TEXT

**Genomic**

Rice, *Oryza sativa* L., is the staple food for half the world's population. By 2030, rice production must increase by at least 25% to keep pace with population growth. Accelerated genetic gains in rice improvement are needed to mitigate the effects of climate change and loss of arable land and to ensure global food supply.
Here, we include data from an international effort resequencing a core collection of 3,000 rice accessions from 89 countries as a global public good. The 3,000 sequenced rice genomes had an average sequencing depth of 14X, average genome coverage and mapping rates of 94.0% and 92.5%, respectively.
This data provides a foundation for large-scale discovery of novel alleles for important rice phenotypes using various bioinformatics and/or genetic approaches. It also serves to understand at a higher level of detail the genomic diversity within *O. sativa*. With the release of the sequencing data, the project calls for the global rice community to take advantage of this data as a foundation for establishing a global, public rice genetic/genomic database and information platform for advancing rice breeding technology for future rice improvement.
Keywords: Oryza sativa, genetic resources, genome diversity, next generation sequencing

The 3000 rice genomes sequence data are now completely uploaded into the INSDC databases (the Sequence Read Archives (SRA) at EBI, DDBJ and NCBI), and available for easy download from the links noted below, rather than from GigaDB. NB - the mapping of each file as previously stored in GigaDB to the relevant location on the EBI FTP server can be found in the file seq_file_mapping_to_SRA.txt.

Please choose the appropriate geographical location:

Europe: EBI - PRJEB6180

USA: NCBI - PRJEB6180

Asia: DDBJ - ERP005654

- Raw data can be found in NCBI Short Read Archive (SRA) or EMBL European Nucleotide Archive or

  DDBJ DDBJ Read Archive (DRA).

| NCBI (SRA) | EMBL (ENA) | DDBJ (DRA) |

- Usually, we can download data via DDBJ or EMBL because NCBI ask us to use special toolkit to download.

- Three archives are synchronized periodically so that data in three archives are same.

**InfoBoss**
Infomations & Systems

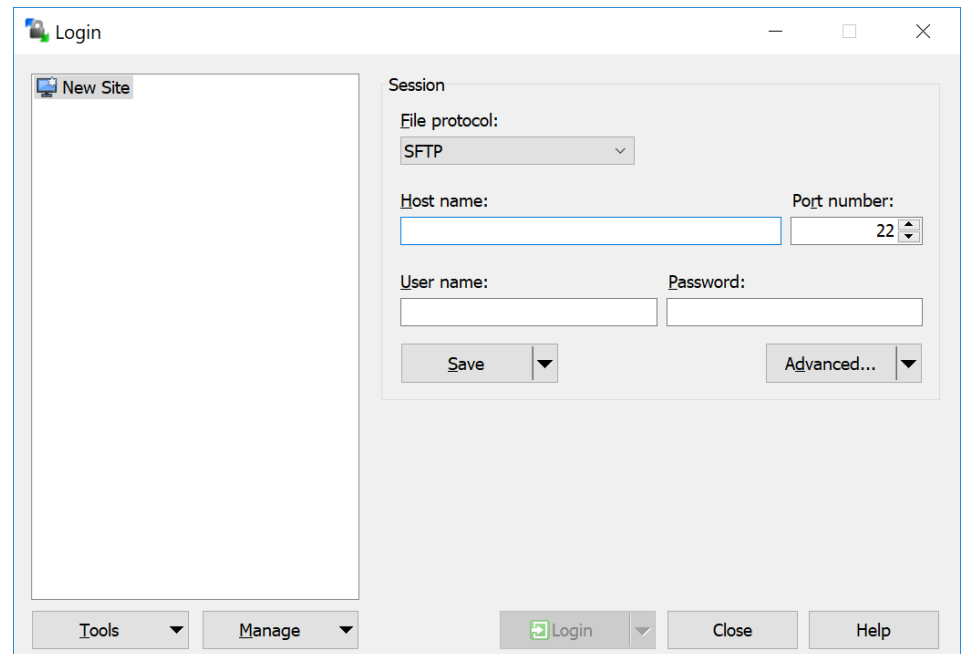- Download fastq files from the below link.

Showing results 1 - 1 of 1 results

| Study accession | Sample accession | Secondary sample accession | Experiment accession | Run accession | Tax ID | Scientific name | Instrument model | Library layout | FASTQ files (FTP) | FASTQ files (Galaxy) | Submitted files (FTP) | Submitted files (Galaxy) | NCBI SRA file (FTP) | NCBI SRA file (Galaxy) | CRAM Index files (FTP) | CRAM Index files (Galaxy) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PRJEB6180 | SAMEA2569476 | ERS469754 | ERX562431 | ERR605660 | 4530 | Oryza sativa | Illumina HiSeq 2000 | PAIRED | File 1 File 2 | File 1 File 2 | Fastq file 1 Fastq file 2 | Fastq file 1 Fastq file 2 | File 1 | File 1 | | |

- Downloaded files are compressed files because of large size of file.

- gz format is better uncompressed in server side because we have to transfer file from local computer to server.
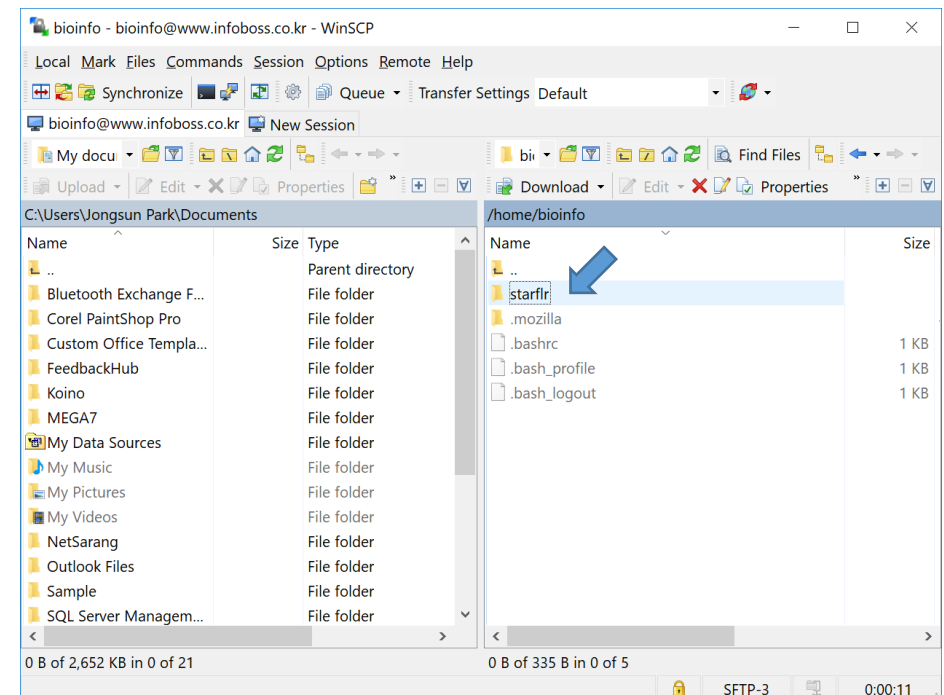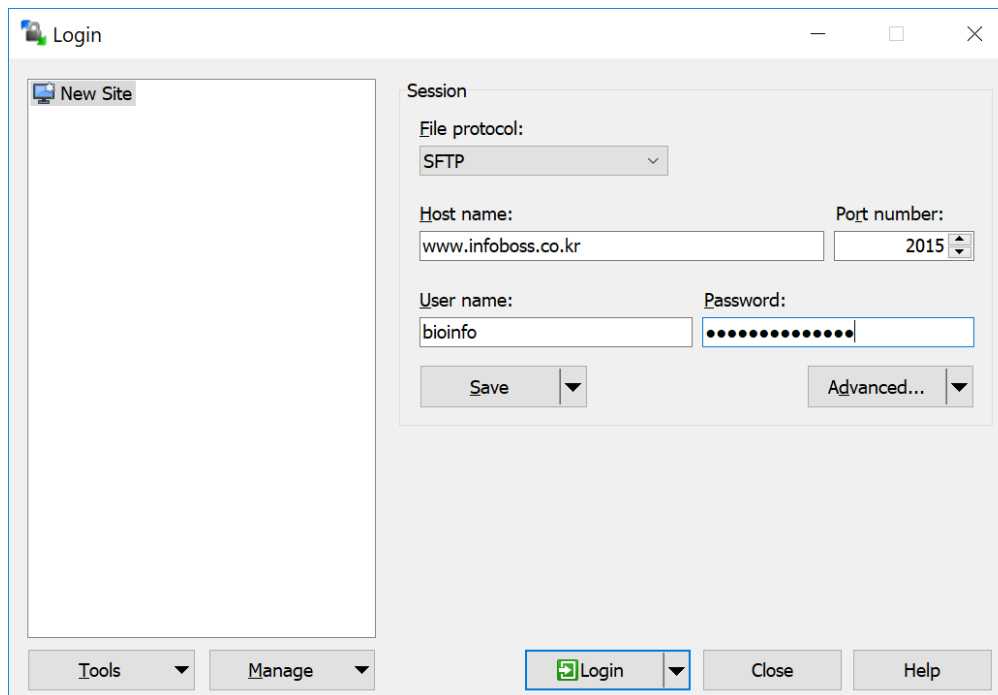
- To transfer files, WinSCP program can be used.

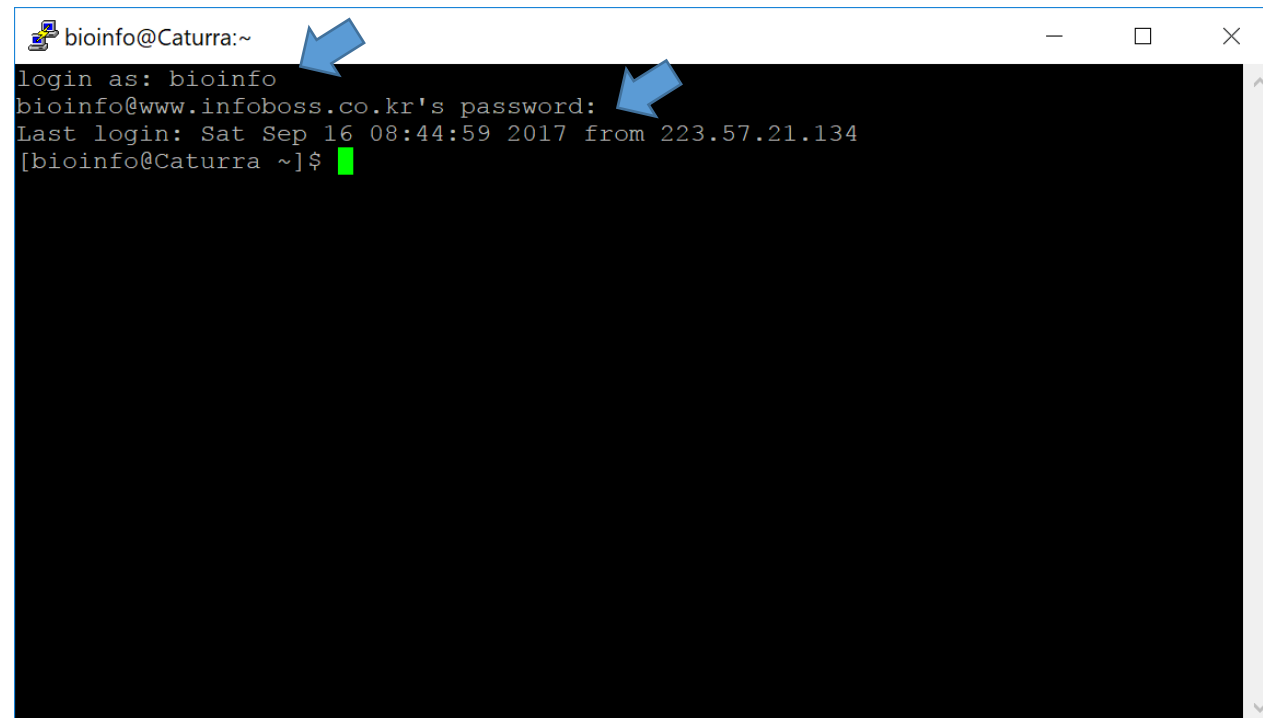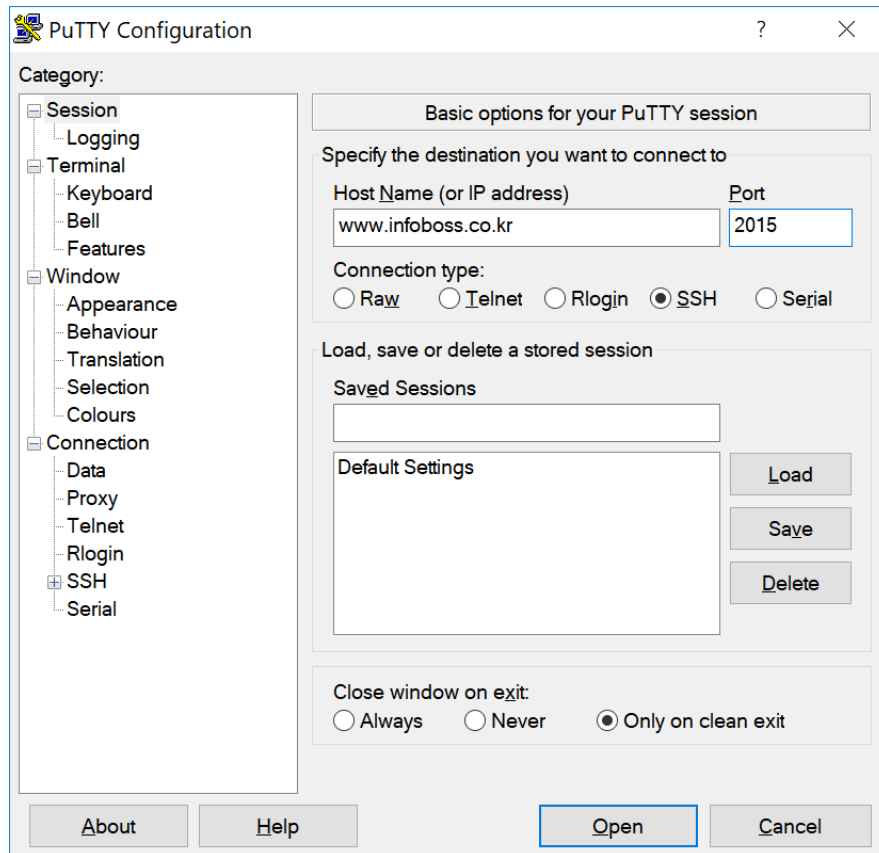- Server IP, port, user account information are required.

- Server information is like below:

> IP : www.infoboss.co.kr
> Port : 2015
> ID : bioinfo
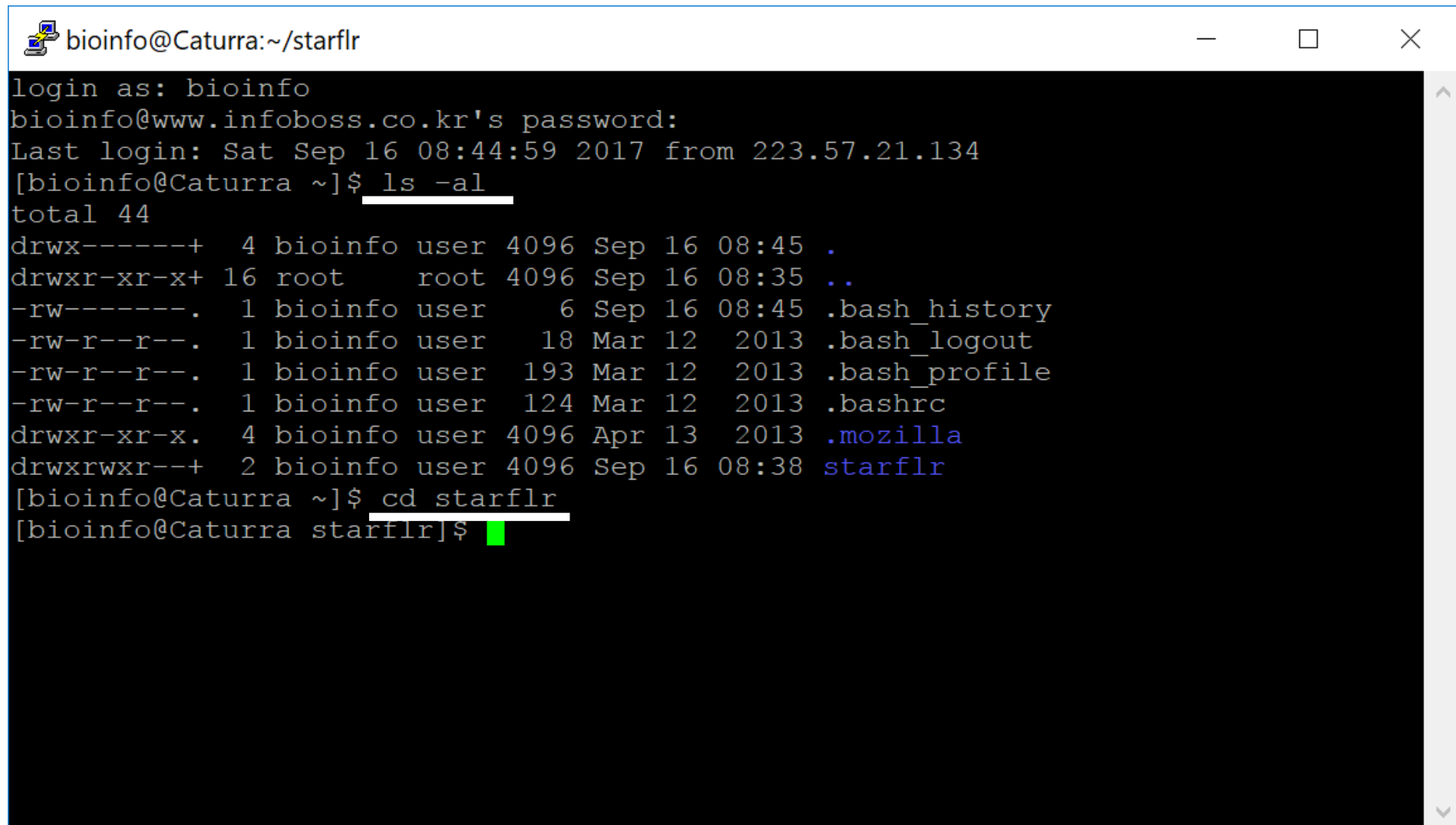> Password : todanfwjdqhgkr

- Using WinSCP, let's upload files after connecting servers.

- Please use personal folder after login, if not, you may not lose your data for the remaining lecture!

# Basic Linux Commands (1) / Basic commands to operate servers
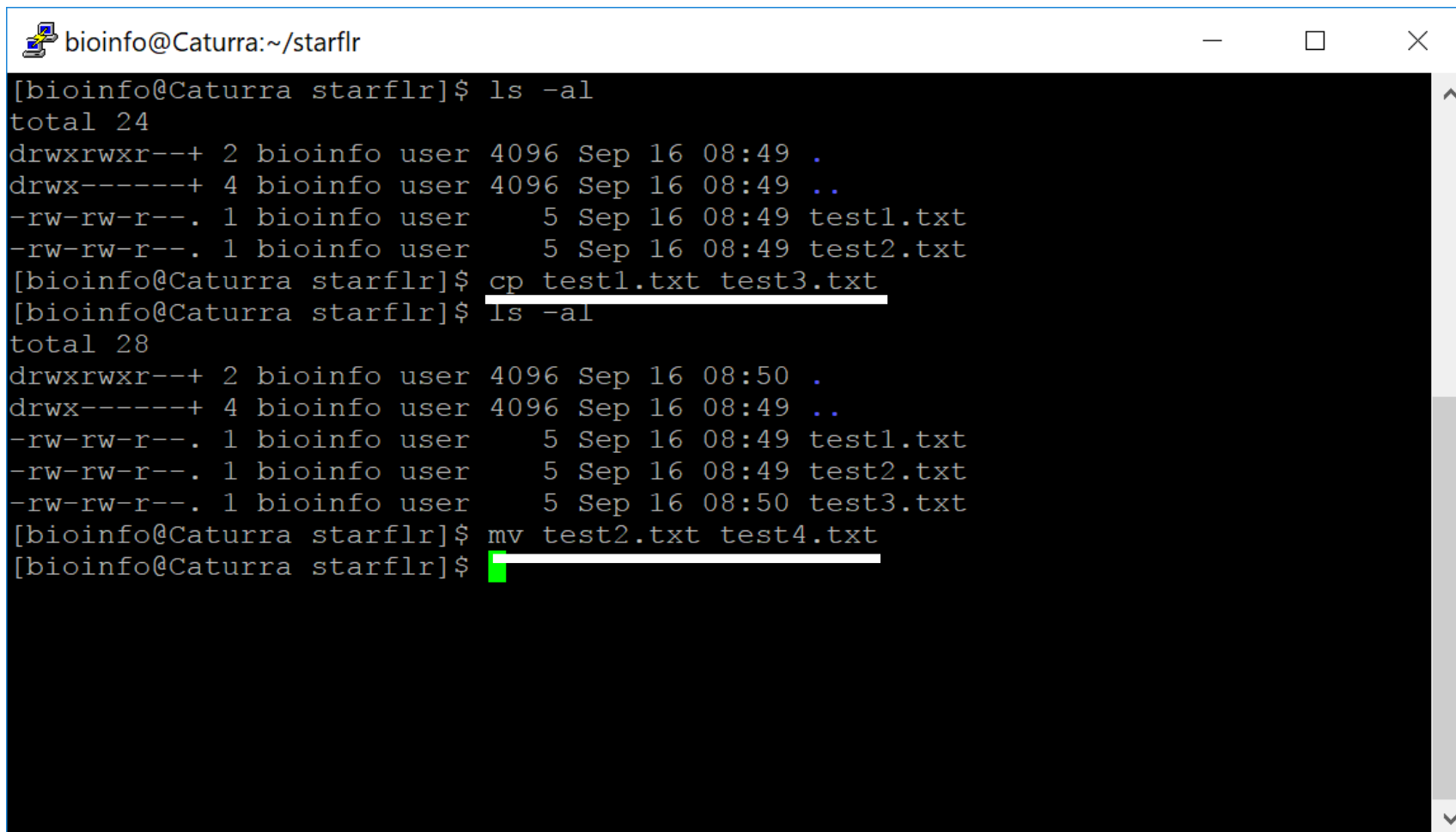


- Let's login to server using Putty.

# Basic Linux Commands (2) / Basic commands to operate servers

- List of files and directories : ls -al

- Change directory: cd

# Basic Linux Commands (3)  / Basic commands to operate servers

**InfoBoss**
Infomations & Systems

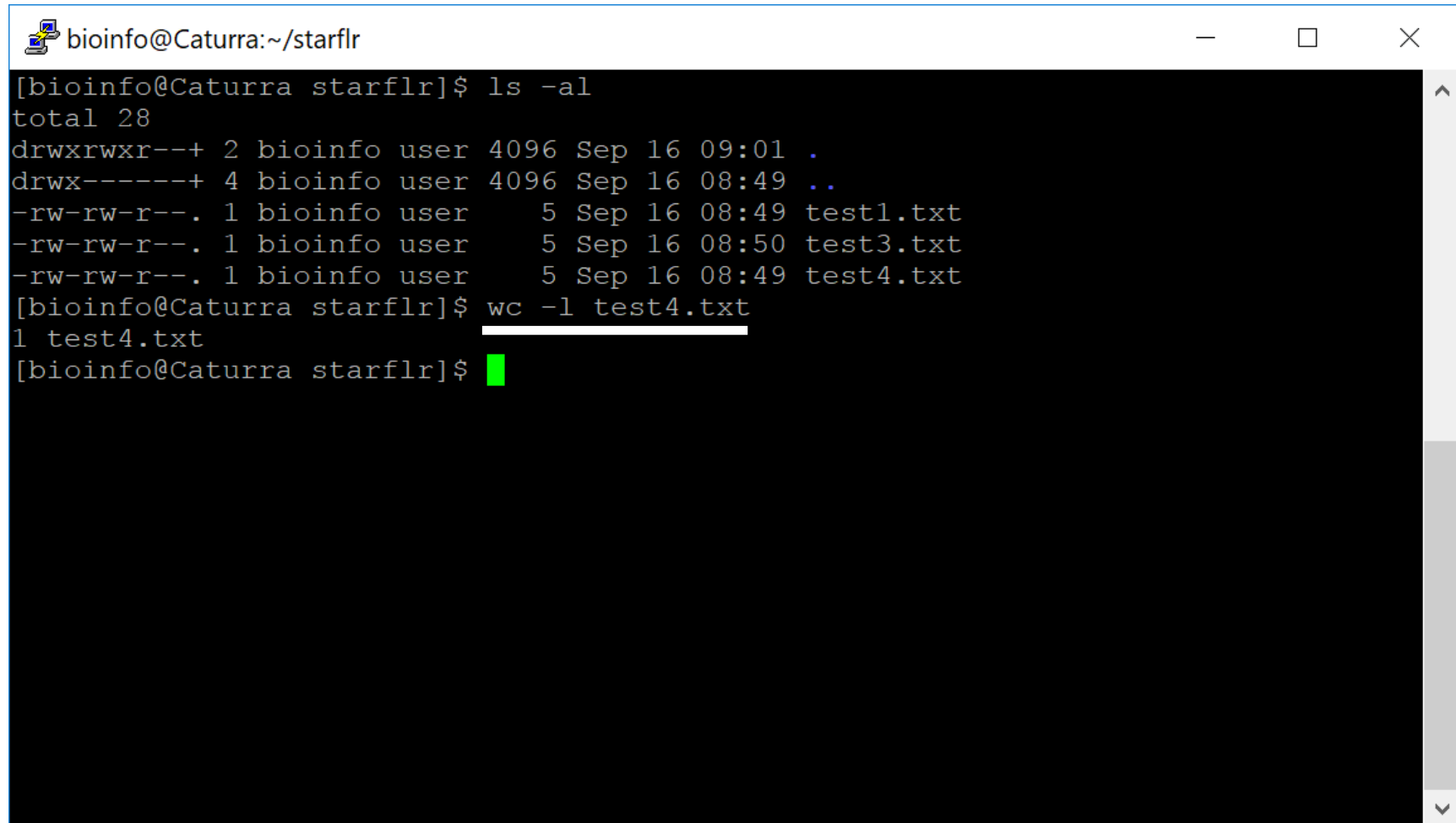- Copy files: cp

- Move files : mv

- Rename filename?

```
bioinfo@Caturra:~/starflr

[bioinfo@Caturra starflr]$ ls -al
total 24
drwxrwxr--+ 2 bioinfo user 4096 Sep 16 08:49 .
drwx------+ 4 bioinfo user 4096 Sep 16 08:49 ..
-rw-rw-r--. 1 bioinfo user    5 Sep 16 08:49 test1.txt
-rw-rw-r--. 1 bioinfo user    5 Sep 16 08:49 test2.txt
[bioinfo@Caturra starflr]$ cp test1.txt test3.txt
[bioinfo@Caturra starflr]$ ls -al
total 28
drwxrwxr--+ 2 bioinfo user 4096 Sep 16 08:50 .
drwx------+ 4 bioinfo user 4096 Sep 16 08:49 ..
-rw-rw-r--. 1 bioinfo user    5 Sep 16 08:49 test1.txt
-rw-rw-r--. 1 bioinfo user    5 Sep 16 08:49 test2.txt
-rw-rw-r--. 1 bioinfo user    5 Sep 16 08:50 test3.txt
[bioinfo@Caturra starflr]$ mv test2.txt test4.txt
[bioinfo@Caturra starflr]$
```

- Compress file using gzip : gzip [filename]

- Uncompress file using gzip : gzip –d [filename]

- Counting number of lines in file: wc –l [filename]

- To run perl program in Linux, please use the below header:

```
#!/usr/bin/perl –w

use strict;
```

- Let's make program to cut 100,000 lines from the downloaded files because that file is too big to test

    during lecture!

PC31.pl

| Write the program in PC | ➡ | Test program in PC |

⬇

| Running perl program in server | ⬅ | Upload program to Server |

# Thank you for your attention!

*If you have question,*
*please ask! =)*

**starflr@infoboss.co.kr**

InfoBoss
Infomations & Systems