

# BIOINFORMATICS SESSION 7. PRACTICE

2023-10-16

Evolution: What makes us human?

# Basic Shell Commands

```
$ cd 2023123456_hyunwoo  
$ mkdir week7_7  
$ cd week7_7
```

# Vim editor - Vim: text editor for Linux

```
$ ll  
$ vi print_script.py
```

[In Vim]  
i: insert

```
print("Today is Wednesday!")
```

[In Vim] "ESC"  
:wq save and quit

```
VIM - Vi IMproved  
  
version 8.2.2072  
by Bram Moolenaar et al.  
Modified by <bugzilla@redhat.com>  
Vim is open source and freely distributable  
  
Help poor children in Uganda!  
type :help iccf<Enter>      for information  
  
type :q<Enter>              to exit  
type :help<Enter> or <F1>   for on-line help  
type :help version8<Enter> for version info
```

# Basic Shell Commands

```
$ ll  
$ python print_script.py
```

**[Output]** Today is Wednesday!

```
$ less print_script.py  
$ cat print_script.py
```

**[Output]** print “Today is Wednesday!”

# Basic Shell Commands

```
$ mv print_script.py new_script.py
$ cp new_script.py new_script_copy.py
$ ll
```

**mv** : **move file**

**cp** : **copy** and **paste**

```
$ rm new_script_copy.py
$ ll
$ ln -s /path/to/file/filename
$ unlink filename
```

**rm** : **remove**

do not use “rm” for symbolic-linked files

please use “unlink”

# Comparing FOXP2 in different animals

---

FOXP2_HUMAN	<i>Homo sapiens</i>
FOXP2_GORGO	<i>Gorilla gorilla</i>
FOXP2_MACMU	<i>Macaca mulatta (Rhesus macaque)</i>
FOXP2_PANTR	<i>Pan troglodytes (Chimpanzee)</i>
FOXP2_HYLLA	<i>Hylobates lar (Common gibbon)</i>
FOXP2_PONPY	<i>Pongo pygmaeus (Bornean orangutan)</i>
FOXP2_MOUSE	<i>Mus musculus (Mouse)</i>
FOXP2_XENLA	<i>Xenopus laevis (African clawed frog)</i>

---



```
$ ln -s /home/biguser/tutor/session6/swissprot.* .
```

# Comparing FOXP2 in different animals

```
[biguser@R440 session7]$ cat seqids.txt  
FOXP2_HUMAN  
FOXP2_GORGO  
FOXP2_MACMU  
FOXP2_PANTR  
FOXP2_HYLLA  
FOXP2_PONPY  
FOXP2_MOUSE  
FOXP2_XENLA
```

# Comparing FOXP2 in different animals

```
$ blastdbcmd -entry FOXP2_HYLLA -db swissprot -long_seqids > FOXP2_HYLLA.fa
```

↳ retrieving a single sequence with id containing FOXP2\_HYLLA

```
$ blastdbcmd -entry_batch seqids.txt -db swissprot -long_seqids > foxp2.fa
```

↳ retrieving multiple sequences with a list of identifiers

```
[biguser@R440 session7]$ blastdbcmd -entry_batch seqids.txt -db swissprot -long_seqids > foxp2.fa
[biguser@R440 session7]$ ll
total 24
-rw-rw-r-- 1 biguser biguser 6518 Oct 10 16:18 foxp2.fa
-rw-r--r-- 1 biguser biguser 5278 Oct 9 18:43 foxp2.fasta
-rw-rw-r-- 1 biguser biguser 813 Oct 9 18:47 FOXP2_HYLLA.fa
-rw-r--r-- 1 biguser biguser 96 Oct 9 18:43 seqids.txt
lrwxrwxrwx 1 biguser biguser 21 Oct 9 18:47 swissprot -> ../session6/swissprot
lrwxrwxrwx 1 biguser biguser 25 Oct 9 18:47 swissprot.phr -> ../session6/swissprot.phr
lrwxrwxrwx 1 biguser biguser 25 Oct 9 18:47 swissprot.pin -> ../session6/swissprot.pin
lrwxrwxrwx 1 biguser biguser 25 Oct 9 18:47 swissprot.pog -> ../session6/swissprot.pog
lrwxrwxrwx 1 biguser biguser 25 Oct 9 18:47 swissprot.psd -> ../session6/swissprot.psd
lrwxrwxrwx 1 biguser biguser 25 Oct 9 18:47 swissprot.psi -> ../session6/swissprot.psi
lrwxrwxrwx 1 biguser biguser 25 Oct 9 18:47 swissprot.psq -> ../session6/swissprot.psq
```



# Comparing FOXP2 in different animals

```
$ clustalw2 foxp2.fa
```

```
[biguser@R440 session7]$ clustalw2 foxp2.fa
```

```
CLUSTAL 2.1 Multiple Sequence Alignments
```

```
Sequence format is Pearson
```

```
Sequence 1: sp|015409|FOXP2_HUMAN 715 aa  
Sequence 2: sp|Q8MJ99|FOXP2_GORGO 713 aa  
Sequence 3: sp|Q8MJ97|FOXP2_MACMU 714 aa  
Sequence 4: sp|Q8MJA0|FOXP2_PANTR 716 aa  
Sequence 5: sp|Q5QL03|FOXP2_HYLLA 713 aa  
Sequence 6: sp|Q8MJ98|FOXP2_PONPY 713 aa  
Sequence 7: sp|P58463|FOXP2_MOUSE 714 aa  
Sequence 8: sp|Q4VYS1|FOXP2_XENLA 706 aa
```

```
Start of Pairwise alignments  
Aligning...
```

```
Sequences (1:2) Aligned. Score: 99  
Sequences (1:3) Aligned. Score: 99  
Sequences (1:4) Aligned. Score: 99  
Sequences (1:5) Aligned. Score: 99  
Sequences (1:6) Aligned. Score: 99  
Sequences (1:7) Aligned. Score: 99  
Sequences (1:8) Aligned. Score: 95  
Sequences (2:3) Aligned. Score: 100  
Sequences (2:4) Aligned. Score: 100  
Sequences (2:5) Aligned. Score: 98  
Sequences (2:6) Aligned. Score: 98  
Sequences (2:7) Aligned. Score: 99  
Sequences (2:8) Aligned. Score: 95  
Sequences (3:4) Aligned. Score: 100  
Sequences (3:5) Aligned. Score: 100  
Sequences (3:6) Aligned. Score: 99  
Sequences (3:7) Aligned. Score: 98  
Sequences (3:8) Aligned. Score: 95  
Sequences (4:5) Aligned. Score: 100  
Sequences (4:6) Aligned. Score: 99  
Sequences (4:7) Aligned. Score: 99  
Sequences (4:8) Aligned. Score: 95  
Sequences (5:6) Aligned. Score: 99  
Sequences (5:7) Aligned. Score: 99  
Sequences (5:8) Aligned. Score: 95  
Sequences (6:7) Aligned. Score: 99  
Sequences (6:8) Aligned. Score: 95  
Sequences (7:8) Aligned. Score: 95  
Guide tree file created: [foxp2.dnd]
```

```
There are 7 groups  
Start of Multiple Alignment
```

```
Aligning...  
Group 1: Sequences: 2 Score: 15411  
Group 2: Sequences: 3 Score: 15416  
Group 3: Sequences: 4 Score: 15408  
Group 4: Sequences: 2 Score: 15421  
Group 5: Sequences: 3 Score: 15400  
Group 6: Sequences: 7 Score: 15379  
Group 7: Sequences: 8 Score: 14978  
Alignment Score 118874
```

```
CLUSTAL-Alignment file created [foxp2.a1n]
```

```
-rw-rw-r-- 1 biguser biguser 10280 Oct 10 17:03 foxp2.a1n  
-rw-rw-r-- 1 biguser biguser 313 Oct 10 17:03 foxp2.dnd
```

# Comparing FOXP2 in different animals

foxp2\_sequences.a1n

```

CLUSTAL 2.1 multiple sequence alignment

sp|Q8MJ99|FOXP2_GORGO      MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
sp|Q8MJ97|FOXP2_MACMU     MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
sp|Q8MJA0|FOXP2_PANTR     MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
sp|015409|FOXP2_HUMAN     MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
sp|Q5QL03|FOXP2_HYLLA    MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
sp|Q8MJ98|FOXP2_PONPY     MMQESVTEITISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
sp|P58463|FOXP2_MOUSE    MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
sp|Q4VYS1|FOXP2_XENLA    MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSDTSSEVSTVELL
                          *****

sp|Q8MJ99|FOXP2_GORGO     HLQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
sp|Q8MJ97|FOXP2_MACMU     HLQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
sp|Q8MJA0|FOXP2_PANTR     HLQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
sp|015409|FOXP2_HUMAN     HLQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
sp|Q5QL03|FOXP2_HYLLA    HLQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
sp|Q8MJ98|FOXP2_PONPY     HLQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
sp|P58463|FOXP2_MOUSE     HLQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
sp|Q4VYS1|FOXP2_XENLA     HLQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
                          *****

sp|Q8MJ99|FOXP2_GORGO     PQMQQILQQQVLSPPQLQALLQQQAVMLQQQQLQEFYKKQEQQLHLQL
sp|Q8MJ97|FOXP2_MACMU     PQMQQILQQQVLSPPQLQALLQQQAVMLQQQQLQEFYKKQEQQLHLQL
sp|Q8MJA0|FOXP2_PANTR     PQMQQILQQQVLSPPQLQALLQQQAVMLQQQQLQEFYKKQEQQLHLQL
sp|015409|FOXP2_HUMAN     PQMQQILQQQVLSPPQLQALLQQQAVMLQQQQLQEFYKKQEQQLHLQL
sp|Q5QL03|FOXP2_HYLLA    PQMQQILQQQVLSPPQLQALLQQQAVMLQQQQLQEFYKKQEQQLHLQL
sp|Q8MJ98|FOXP2_PONPY     PQMQQILQQQVLSPPQLQALLQQQAVMLQQQQLQEFYKKQEQQLHLQL
sp|P58463|FOXP2_MOUSE     PQMQQILQQQVLSPPQLQALLQQQAVMLQQQQLQEFYKKQEQQLHLQL
sp|Q4VYS1|FOXP2_XENLA     PQMQQILQQQVLSPPQLQALLQQQAVMLQQQQLQEFYKKQEQQLHLQL
                          *****
    
```

#dendrogram

foxp2\_sequences.dnd

```

(
(
(
sp|015409|FOXP2_HUMAN:0.00122,
(
(
sp|Q8MJ99|FOXP2_GORGO:0.00117,
sp|Q8MJ97|FOXP2_MACMU:-0.00117)
:0.00196,
sp|Q8MJA0|FOXP2_PANTR:-0.00196)
:0.00157)
:0.00059,
sp|Q4VYS1|FOXP2_XENLA:0.04051)
:0.00083,
(
sp|Q5QL03|FOXP2_HYLLA:0.00000,
sp|Q8MJ98|FOXP2_PONPY:0.00140)
:0.00066,
sp|P58463|FOXP2_MOUSE:0.00074);
    
```

- \* -- all residues or nucleotides in that column are identical
- : -- conserved substitutions have been observed
- . -- semi-conserved substitutions have been observed
- no match.

# Identification of mutation specific to human

```
$ clustalw2 foxp2.fa -output=fasta  
$ less foxp2.fasta
```

```
>sp|Q8MJ99|FOX P2_GORGO  
MMQESATETISNSSMÑNGMSTLSSQLDAGSRDRSSGDTSSSEVSTVELL  
HLQQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT  
PQQMQQILQQQVLSPPQLQALLQQQAVMLQQQQLQEFYKKQQEQLHLQL  
LQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQ - - HPGKQAKE  
QQQQQQQQQLAAQQLVFQQQLLQMQLQQQHLLSLQRQGLISIPPGQA  
ALPVQSLPQAGLSPAIEIQLWKEVTGVHSMEDNGIKHGGLDLTTNNSST  
TSSTTSKASPPITHHSIVNGQSSVLNARRDSSSHEETGASHTLYGHGVCK  
WPGCESICEDFGQFLKHLNNEHALDDRSTAQCRVQMQQVQLEIQLSKER  
ERLQAMMTHLHMRPSEPKPSKPLNLVSSVTMSKNMLETSPQSLPQTPTT  
PTAPVTPITQGPSVITPASVPNVGAIARRRHS DKYNIPMSSEIAPNYEFYK  
NADV RPPFTYATLIRQAIMESSDRQLTLNEIYSWFTRTFAYFRRNAATWK  
NAVRHNL SLHKCFVRVENVKGAVWTVDEVEYQKRRSQKITGSP TLVKNIP  
TSLGYGAALNASLQAALAESSLPLLSNPGLINNASSGLLQAVHEDLNGSL  
DHIDSNGNSSPGCSPQPHIHSIHVKEEPVIAEDEDCPMSLVTTANHSPEL  
EDDREIEEEPLSEDL E
```

만약에 linux 상에서 작업할 수 없거나 결과를 만들지 못했으면  
`ln -s /home/biguser/tutor/session7/foxp2.fasta .`  
foxp2.fasta 파일 링크가 형성됐는지 확인하세요!

# Identification of mutation specific to human

## □ foxp2.ipynb

```
1 import re
2
3 inFile = open('/home/biguser/your_directory/session7/foxp2.fasta', 'r')
4
5 nonhuman = dict()
6
7 id = ''
8 seq = ''
9
10 for line in inFile.readlines():
11     line = line.strip()
12     #print line
13     if re.search('^\>', line):
14         if id != '':
15             if not "HUMAN" in id:
16                 nonhuman[id[1:]] = seq
17             else:
18                 id_human = id[1:]
19                 seq_human = seq
20                 id = line
21                 seq = ''
22         else:
23             id = line
24     else:
25         seq += line
26
27 if not "HUMAN" in id:
28     nonhuman[id[1:]] = seq
29 else:
30     id_human = id[1:]
31     seq_human = seq
32
33 inFile.close()
```

	POSITION 1	POSITION 2
HUMAN	A	T
GORGO	A	A
MACMU	A	C
PANTR	A	C
HYLLA	A	A
PONPY	A	G

```
35 for i in range(0, len(seq_human)):
36     unique = True
37     for id in nonhuman.keys():
38         human_seq_posi = seq_human[i]
39         nonhuman_seq_posi = nonhuman[id][i]
40         if human_seq_posi == nonhuman_seq_posi:
41             unique = False
42     if unique:
43         pos = i + 1
44         print("At position ", pos)
45         aa = human_seq_posi
46         print(id_human, '#t', aa)
47         for ID in nonhuman.keys():
48             print(ID, '#t', nonhuman[ID][i])
```

# Identification of mutation specific to human

## □ results

At position 304

sp 015409 FOXP2_HUMAN	N
sp Q8MJ99 FOXP2_GORGO	T
sp Q8MJ97 FOXP2_MACMU	T
sp Q8MJA0 FOXP2_PANTR	T
sp Q5QL03 FOXP2_HYLLA	T
sp Q8MJ98 FOXP2_PONPY	T
sp P58463 FOXP2_MOUSE	T
sp Q4VYS1 FOXP2_XENLA	T

Human having Asparagine(N)

whereas all non-human sequences have threonine(T)

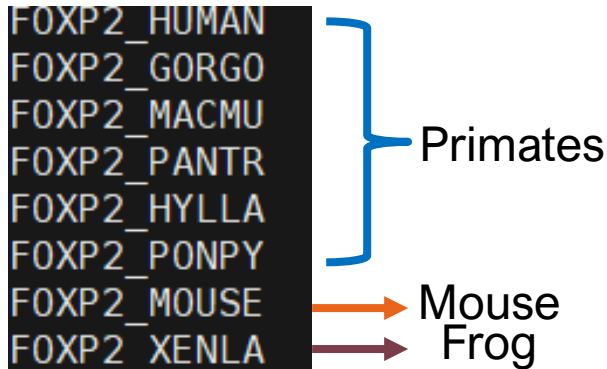
not equal because every sequence has different length  
thus, gaps are inserted

## Human-specific T303N

- T303N and N325S are positively selected for language and speech
- No N325S in our res. due to *Xenopus*

# Exercise

- Modify `foxp2.ipynb` such that you may identify
  1. positions where the human sequence is different from all other primates and mouse;
  2. positions where the human and primate sequences are identical but different to both mouse and frog.



# Exercise -1

```
At position 304
sp|O15409|FOXP2_HUMAN      N
sp|Q8MJ99|FOXP2_GORGO     T
sp|Q8MJ97|FOXP2_MACMU     T
sp|Q8MJA0|FOXP2_PANTR     T
sp|Q5QL03|FOXP2_HYLLA     T
sp|Q8MJ98|FOXP2_PONPY     T
sp|P58463|FOXP2_MOUSE     T
At position 326
sp|O15409|FOXP2_HUMAN      S
sp|Q8MJ99|FOXP2_GORGO     N
sp|Q8MJ97|FOXP2_MACMU     N
sp|Q8MJA0|FOXP2_PANTR     N
sp|Q5QL03|FOXP2_HYLLA     N
sp|Q8MJ98|FOXP2_PONPY     N
sp|P58463|FOXP2_MOUSE     N
```

the first part of the code is same as foxp2.ipynb

```
for i in range(0, len(seq_human)):
    unique = True
    for id in nonhuman.keys():
        if not "XENLA" in id:
            human_seq_posi = seq_human[i]
            nonhuman_seq_posi = nonhuman[id][i]
            if human_seq_posi == nonhuman_seq_posi:
                unique = False
    if unique:
        pos = i + 1
        print("At position ", pos)
        aa = human_seq_posi
        print(id_human, '\t\t', aa)
        for ID in nonhuman.keys():
            if not "XENLA" in ID:
                print(ID, '\t\t', nonhuman[ID][i])
```

# Exercise-2

```
At position 80
sp|O15409|FOXP2_HUMAN      D
sp|Q8MJ99|FOXP2_GORGO     D
sp|Q8MJ97|FOXP2_MACMU     D
sp|Q8MJA0|FOXP2_PANTR     D
sp|Q5QL03|FOXP2_HYLLA     D
sp|Q8MJ98|FOXP2_PONPY     D
sp|P58463|FOXP2_MOUSE     E
sp|Q4VYS1|FOXP2_XENLA     E
```

```
1 import re
2
3
4 for i in range(0, len(seq_human)):
5     unique = True
6     for id in nonhuman.keys():
7         human_seq_posi = seq_human[i]
8         if "XENLA" in id or "MOUSE" in id:
9             if human_seq_posi == nonhuman[id][i]:
10                unique=False
11
12            else:
13                if human_seq_posi != nonhuman[id][i]:
14                    unique=False
15
16        if unique:
17            pos = i + 1
18            print("At position ", pos)
19            aa = human_seq_posi
20            print(id_human, '#t', aa)
21            for ID in nonhuman.keys():
22                print(ID, '#t', nonhuman[ID][i])
```



# Assignment

- In the default output from the ClustalW program (a file named as foxp2.aln) there are asterisks (\*) that indicate positions where the sequence is same in all sequences. Make a Python script to count the total length and number of such positions in all alignment.
- Output
- clustalw의 결과 중 foxp2.aln을 이용해서 algin된 총 길이와 sequeunce들이 모두 동일함을 뜻하는 \* 표시가 총 몇 개인지 아래와 같이 print.

```
Total length :
```

```
Number of asterisk (match) :
```

- \* 과제 제출 기한: 10/22 Sunday 23:59 @ LMS
- \* 작성한 코드와 해당 코드의 결과를 캡처한 뒤 워드에 첨부(코드만 굵어와서 붙여넣지 말기), 코드에 대한 설명 간략히 작성  
워드 파일명은 n주차\_학번\_이름 형식으로 제출(e.g. 7주차\_2023123456\_김현우)

# Assignment

In `-s/home/biguser/tutor/Week7/foxp2.aln` .  
를 linux환경에서 여러분의 directory에서 실행시키고 해당 파일 link가 형성 됐는지 확인. 해당 파일을 읽어서 작업할 것

```
CLUSTAL 2.1 multiple sequence alignment

gi|51701430|sp|Q8MJ99.1|FOXP2_      MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
gi|51701429|sp|Q8MJ97.1|FOXP2_      MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
gi|38503046|sp|Q8MJA0.1|FOXP2_      MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
gi|17432967|sp|O15409.2|FOXP2_      MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
gi|62286912|sp|Q5QL03.1|FOXP2_      MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
gi|146345420|sp|Q8MJ98.3|FOXP2      MMQESVTEITISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
gi|51704301|sp|P58463.2|FOXP2_      MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSGDTSSSEVSTVELL
gi|82227296|sp|Q4VYS1.1|FOXP2_      MMQESATETISNSSMNQNGMSTLSSQLDAGSRDGRSSSDTSSEVSTVELL
*****

gi|51701430|sp|Q8MJ99.1|FOXP2_      HLQQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
gi|51701429|sp|Q8MJ97.1|FOXP2_      HLQQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
gi|38503046|sp|Q8MJA0.1|FOXP2_      HLQQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
gi|17432967|sp|O15409.2|FOXP2_      HLQQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
gi|62286912|sp|Q5QL03.1|FOXP2_      HLQQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
gi|146345420|sp|Q8MJ98.3|FOXP2      HLQQQQALQAARQLLLQQQTSGLKSPKSSDKQRPLQVPVSVAMMTPQVIT
gi|51704301|sp|P58463.2|FOXP2_      HLQQQQALQAARQLLLQQQTSGLKSPKSEKORPLQVPVSVAMMTPQVIT
gi|82227296|sp|Q4VYS1.1|FOXP2_      HLQQQQALQAARQLLLQQQTSGLKSPKNNEKORPLQVPVSMAMMTPQVIT
*****
```