

'\$' ile başlayanlar bash komutları.  
'#' ile başlayanlar yorum ve açıklamalar.

## 1) GIRIS

# bilgisayarımızda hangi shell var?  
\$ echo \$SHELL

\$ pwd

\$ cd /directory1  
\$ cd directory1/directory2/directory3

\$ cd ..  
# üst konum  
\$ cd ../../..

\$ ls  
\$ ls -t  
\$ ls -l  
\$ ls -1  
\$ ls -l -1  
\$ ls -lt1  
# sıralama önemli değil

\$ man ls  
# köşeli parantez içindekiler opsiyonel

\$ mkdir isim\_soyad  
# dosyaya girerken TAB'e basma alışkanlığı.

#brace expansion kullanmak  
\$ echo sekans\_{ornek1,ornek2,ornek3}  
# ne oldu?

# yukarı ok tuşunun faydası

# mkdir -p zmays\_snps/{data,seqs,scripts,analysis}

# *silmek*  
\$ touch test1 test2  
\$ ls

\$ rm touch1 touch2  
# silerken onay isteyebilir, bu nedenle çok tehlikeli bir komut.  
# rm ile bir şeyi silerseniz geri gelmez!  
# tüm dosyalar silabilirsiniz!

#bin/ folder's icini silebilirsiniz

\$ touch test

\$ rm -i test

# -i parametresi ile sildiğinizde kesinlikle onay ister.

Ctrl + a     basa git

Ctrl + e     sona

Ctrl + w     komuttaki bir önceki kelimeyi sil

# hangi programlar çalışıyor, ne kadar CPU kullanıyorlar, 'top' komutu ile öğrenebilirsiniz.

\$ top

# devam eden işlemler neler:

\$ jobs

# dosya adlandırma mantığı:

# Yanlış: genes\_1.txt, genes\_2.txt,...,genes\_11.txt - numeric sıralamayı bilgisayar algılayamaz bu şekilde numaralandırırsanız.

#doğru adlandırma:

genes\_001.txt, genes\_002.txt,...genes\_011.txt

## 2) UNIX DATA ARAÇLARI'na giriş:

SIMSDI İSİMİZE YARAYACAK ASIL KONUYA GİRİYORUZ, unix-bash ile text dataları (sekans dataları da bir basit text datası zaten) ile çalışmak,

# cat, grep, cut, awk gibi programları/komutları kullanarak data analizinde ilk adımı atacağız!

Mantık: sekans datası üstüne düşüneceğiz (fakat big data analizi için de kullanabilirsiniz bu araçları). Gb boyutunda, satırlarca text datası. Data'yi acmak, manipüle etmek, memory işgal etmek olasılık dışı. Read-only analiz gerekiyor. Linux, bash read-only data analizi için birebir bilgisayar işlemleri ortamı.

CAT:

# cat, memory'ye almadan "stream" olarak bash ekranında dosyayı görüntüler

#tga1-protein.fasta dosyasına bakalım:

\$ cat tga1-protein.fasta

# catkin nasıl çalıştırıldığını öğrenmek için manual'ine bakalım:

\$ man cat

#Birden fazla dosyaya tek komut ile işlem uygulayabilirsiniz:

\$ cat tb1-protein.fasta tga1-protein.fasta

# \* wild card'ini kullanarak deneyelim?

\$ cat \*.fasta

# .fasta ortak ismine sahip tüm dosyalar '\*' wild card'ini kullanarak görüntüledik.

# Input'u output'a yönlendirmek (yeni bir dosya olarak kaydetmek):

'>' veya '>>' kullanarak yönlendirmek:

'>' girdileri yeni bir dosyaya yönlendirip/kaydeder - eğer aynı isimli başka bir dosya var ise, sormadan üzerine yazar, dikkat!

'>>' yine yönlendiriyor fakat eğer o isimde bir dosya varsa, ilgili dosyanın sonuna ekler yeni veriyi.

#CAT programı ile append edebilirsiniz.

\$ cat tb1-protein.fasta tga1-protein.fasta > zea\_proteins.fasta

\$ cat zea\_proteins.fasta

HEAD & TAIL

\$ man head

\$ head Mus\_musculus.GRCm38.75\_chr1.bed

\$ man tail

\$ tail Mus\_musculus.GRCm38.75\_chr1.bed

\$ head -17 Mus\_musculus.GRCm38.75\_chr1.bed

# ilk 17 satiri gösteriyor

\$ tail -17 Mus\_musculus.GRCm38.75\_chr1.bed

\$ (head -n 2; tail -n 2) < Mus\_musculus.GRCm38.75\_chr1.bed

\$ head yeast\_chr1\_orfs.fa.txt

\$ tail yeast\_chr1\_orfs.fa.txt

BED & GTF dosyaları farklı kalıplarda dosyalardır:

\$ head -15 Mus\_musculus.GRCm38.75\_chr1.gtf

## LESS

# cat ekrandan stream olarak dosya icerigini okutuyordu. Akip-gitmeden bakmak istersek, an basit kullanabileceğimiz program 'less'. Çok büyük text dosyalarını / genom datası dosyalarını hiç sorun olmadan acabilir (ornegin word ile devasa text dosyalarını acmaya calisirsaniz bilgisayar zorlanır). 'Less'in hızlı olmasının mantigi memory'ye almadan okuyabilmesi.

\$ man less

# yukarı asabi nasıl hareket edebileceğinizi manual'den ogrenin.

\$ less contaminated.fastq

# cikis icin 'q'ya basiniz.

#daha büyük bir dosyaya bakalım:

\$ less celegans.sam

## WC

# The WC utility displays the number of lines, words, and bytes contained in each input file, or standard input (if no file is specified) to the standard output.

\$ wc Mus\_musculus.GRCm38.75\_chr1.bed

# 3 adet bilgi verdi. Neler?

\$ wc mm\_GRCm38.75\_protein\_coding\_genes.gtf

# genelde ilgili datanın/dosyanın kaç satır olduğu ile ilgiliyiz:

\$ wc -l Mus\_musculus.GRCm38.75\_chr1.bed

# uyari: wc empty lines larida sayiliyor.

## CUT - SUTUN VERISI ILE CALISMAK

# cut'in tab-delimited dosya secimi "by default"

# Once cat ile dosyaya bakalım:

\$ cat Mus\_musculus.GRCm38.75\_chr1.bed

# 2. sutunun icerigine ulaşmak:

\$ cut -f 2 Mus\_musculus.GRCm38.75\_chr1.bed

\$ cut -f 2 Mus\_musculus.GRCm38.75\_chr1.bed | head -14

# -f argument birden fazla sutuna ulaşmak icin de kullanilabilir. 3, 4, 5, 6, 7, 8. sutunlarin

hepsine ulaşmak için -f 3-8 argümanı yeterlidir.

# 1'den 3'e:

```
$ cut -f1-3 Mus_musculus.GRCm38.75_chr1.bed | head -14
```

```
$ cut -f2-3 Mus_musculus.GRCm38.75_chr1.bed | head -14
```

## GREP

# grep'in kullanım mantığı: aradığınız bir 'kalıp' (pattern) ister (regular expression olabilir, harf ve rakamları içeren bir kalıp olabilir), bir de kalıbı içinde arayacağınız datayı ister.

```
$ man grep
```

```
$ grep Man heroes_and_villains.txt
```

```
# case sensitive!
```

```
$ grep -i Man heroes_and_villains.txt
```

```
# 'i' seçeneği: case-insensitive arama
```

```
# ilgili pattern'i buluyor, tüm satırı ekrana veriyor.
```

```
$ grep -i -v Man heroes_and_villains.txt
```

```
# 'v' aramanın mantığını zıt yöne çeviriyor. Yani 'Man' i "ignore" ediyor bu komutta.
```

```
$ grep -w Man heroes_and_villains.txt
```

```
# 'w' sadece tek başına "Man" i içeren satırları buluyor.
```

```
#grep ile bir kalıbın kaç kere tekrarlandığını bulabilirsiniz:
```

```
$ grep -c Villain heroes_and_villains.txt
```

```
# grep -c: grep count.
```

```
# peki grep ne işe yarayabilir data analizinde:
```

```
$ grep -n -i CGTATAT yeast_chr1_orfs.fa.txt
```

```
# 7 nükleotidlik kalıbın bulunduğu satırları buldu.
```

```
#Bu doğru cevap olmayabilir! Dikkat! 'grep', eğer sekans iki farklı satıra yayılmışsa, aranan diziyi bulamaz!
```

```
# çalıştığımız bir gene dair özellikleri grep ile bulalım:
```

```
$ grep 'gene_id "ENSMUSG00000025907"' Mus_musculus.GRCm38.75_chr1.gtf | head
```

```
-n 1
```

```
$ grep "Olfr418-ps1" Mus_musculus.GRCm38.75_chr1_genes.txt
```

```
# "partial match" mantığı ile çalışıyor:
```

```
$ grep Olfr Mus_musculus.GRCm38.75_chr1_genes.txt | head -n 5
```

# SORU: "Olfr1413" dışında tüm "Olfr" genlerine ihtiyacınız varsa ne yapacaksınız?  
# pipeline oluşturarak, ilk kısımda bir işlem yapıyoruz, output'u 'l' işaretinden sonra ikinci bir programa input olarak veriyoruz:

```
$ grep Olfr Mus_musculus.GRCm38.75_chr1_genes.txt | grep -v Olfr1413
```

#Emin miyiz? Çalıştı mı? Yazdığınız komut ve kodların çalışıp çalışmadığının SAGLAMASI her seferinde yapılmalı. Çoğu zaman şikayet etmeden komut bir output verebilir. "Garbage in garbage out" deyişi, bu riskli durumu anlatır. Kötü bir komut/kod yazıp, yanlış sonuçlar almanız sürekli mümkündür. Doğru işlem mi yaptırdınız, sürekli test etmek lazım. Sağlamayı nasıl wc -l ile yaptırabilirsiniz? Düşün...

# Çözüm:

```
$ grep Olfr Mus_musculus.GRCm38.75_chr1_genes.txt | grep -v Olfr1413 | wc -l  
$ grep Olfr Mus_musculus.GRCm38.75_chr1_genes.txt | wc -l
```

# Önemli uyarı: bu tip arama işlemlerinizi yaparken kuracağınız mantık minimum hedeflenmeyen kalıp bulmaya göre önceden düşünülerek yapılmalı.

# grep'in diğer işe yarayan bir özelliği, ilgili kalıbın (pattern) önünde (-B), sonrasında (-A), veya hem önünde hem arkasında bulunan daha fazla kısmını da gösterebilir.

```
$ grep -B1 "AGATCGG" contam.fastq | head -n 6  
# "Print one line of context before (-B) the matching line."
```

# GREP çok kuvvetli bir unix data aracıdır. Çok hızlı. Birçok regular expression arayıp bulan programdan daha hızlı. Sebebi de ilgili kalıbı satırda bulunca, satırın gerisini aramadan (daha fazla zaman kaybetmeden) ilgili satırı çıktıya verir.

```
$ grep -A2 "AGATCGG" contam.fastq | head -n 6
```

# "Olfr218" veya "Olfr1416." genlerini aramak için:

```
$ grep -E "(Olfr1413|Olfr1416)" Mus_musculus.GRCm38.75_chr1_genes.txt  
$ man grep  
# -E opsiyonunun ne işe yaradığını okuyun.
```

```
$ grep "Olfr141[13]" Mus_musculus.GRCm38.75_chr1_genes.txt
```

#count seçeneği:

```
$ grep -c "\tOlfr" Mus_musculus.GRCm38.75_chr1_genes.txt  
#İlgili datada Olfr gen ismi kalıbından kaç adet olduğunu bulduk.
```

# Mus\_musculus kromozom 1'in özellikleri bilgisini içeren .gtf dosyasının ilk başında bulunan yorum (# ile başlayan) satırları eleterek, asıl ilgilendiklerimiz bilginin olduğu ilk 10 satıra bakalım:

```
$ grep -v "^#" Mus_musculus.GRCm38.75_chr1.gtf | cut -f1-8 | head
```

# bash ekranında biraz daha düzenli görmek isterseniz 'column -t' programından yardım alabilirsiniz:

```
$ grep -v "^#" Mus_musculus.GRCm38.75_chr1.gtf | cut -f 1-8 | column -t | head -n 3
```

CUT + GREP (birden çok programı, pipeline ile çalıştırmak):

# comma separated bir dosyaya (.csv uzantılı) yine column -t ile daha düzenli bakmak mümkün, fakat -s"," parametresi ile separator'un ',' (virgül) olduğunu programa bildirmemiz gerekiyor:

```
# once column -t kullanmadan bakalım.
```

```
$ head Mus_musculus.GRCm38.75_chr1_bed.csv
```

```
# şimdi column -t kullanarak bakalım.
```

```
$ column -s"," -t Mus_musculus.GRCm38.75_chr1_bed.csv | head -n 3
```

— —

Egzersiz: CUT'i kullanarak Mus musculus kromozom 1 özelliklerini içeren .gtf dosyasını 3 sütunun bilgisini içeren, tab-ile-ayrılmış (tab-delimited) özelliklerin chromosome, start, & end position'i veren daha sade bir dosyaya dönüştürelim:

#problemi adım adım düşünürseniz kolayca çözebilirsiniz:

- adım 1: en basta '#' ile başlayan metadata'dan kurtulalım.

- (chromosome, start, end) bilgileri 1., 4., 5. sütunlarda. CUT kullanarak ayıklayalım.

Çözüm:

# öncelikle # içeren satırlardan kurtulalım demistik. # içeren satırlar derken:

```
$ head "^#" Mus_musculus.GRCm38.75_chr1.gtf | column -t
```

```
$ grep -v "^#" Mus_musculus.GRCm38.75_chr1.gtf | cut -f1,4,5 | head -n 3
```

# -v'nin (**-invert-match, Selected lines are those not matching any of the specified patterns.**) istediğimiz gibi çalıştığına emin olduktan sonra yeni dosyaya kaydedebiliriz:

```
$ grep -v "^#" Mus_musculus.GRCm38.75_chr1.gtf | cut -f1,4,5 > test.txt
```

```
$ ls -ltr
```

```
$ cat test.txt
```

# Mus musculus gif özellikler dosyasında small nuclear RNA sayısını bulmak istiyoruz. snRNA özellikleri nasıl kayıtlı ilgili .gtf dosyasında: gene\_biotype "snRNA" - buradan yola çıkıp bu özellikleri sayalım.

```
$ grep -c 'gene_biotype "snRNA"' Mus_musculus.GRCm38.75_chr1.gtf
```

UYARI: ' ' ve " " ayrimina dikkat!

# pseudogene bayisini bulalim:

```
$ grep -c 'pseudogene' Mus_musculus.GRCm38.75_chr1.gtf
```

# grep hızlı çünkü satırda match bulunca satirin gerisine bakmadan tum satiri veriyor.  
Sadece istediğiniz kısmi vermesi için '-o' seçeneği kullanılır:

```
$ grep -o "Olfr.*" Mus_musculus.GRCm38.75_chr1_genes.txt | head -n 3
```

**SORT**

#sort alfa-numerik calisir.

```
$ cat example.bed
```

# sutunlarin nasıl sıralandigini gordunuz.

# sort yaparken farklı isteklerimiz var:

- belli bir sutuna gore sort etmek
- sort'a belli sutunlarin "sayısal" oldugunu söylemek.
- birden cok sutunu dikkate alarak sort etmek.

# example.bed'i kromozom (sutun 1) & kromozom uzerinde baslangic konumuna (sutun 2) gore sıralamak istediğimizi düşünelim :

```
$ sort -k1,1 -k2,2n example.bed
```

# n numerical sort

# -k1,1 : start,start

```
$ sort -k1,1 -k2,2n example.bed > example_sorted.bed
```

# sort edilmiş dosyaya yeni bir dosyaya kaydedelim.

Ornek: daginik bir .gtf dosyasini once kromozom sonra da pozisyona gore sıralayalım:

```
$ sort -k1,1 -k4,4n Mus_musculus.GRCm38.75_chr1_random.gtf >
```

```
Mus_musculus.GRCm38.75_chr1_sorted.gtf
```

```
$ sort -k1,1 -k2,2n example2.bed
```

# sort computer tabanla sıralama yapıyor, chr2 vs chr11

**UNIQ -c**

```
Uniq -c Precede each output line with the count of the  
number of times the line occurred in the input, followed by a
```



single space.

```
# lac ozgun ozellik var, letters.txt dosyasinda
$ uniq -c letters.txt
# neden olmadı? Sort etmediğimiz için.
$ sort letters.txt | uniq -c
```

```
# mükemmel örnek. Mus musculus chr1'de bulunan genetik özelliklerin sayısı:
$ grep -v "^#" Mus_musculus.GRCm38.75_chr1.gtf | cut -f3 | sort | uniq -c
# aynı isi çoktan aza doğrular sıralarsak:
$ grep -v "^#" Mus_musculus.GRCm38.75_chr1.gtf | cut -f3 | sort | uniq -c | \
sort -rn
```

```
# bir gene ait olan özellikler:
$ grep "ENSMUSG00000033793" Mus_musculus.GRCm38.75_chr1.gtf | cut -f3 | sort \
| uniq -c
```

```
# bir başka gene ait özelliklere bakmak:
$ grep "Lypla1" Mus_musculus.GRCm38.75_chr1.gtf | cut -f 3 | sort | uniq -c
# cut, sort,uniq sırasının mantığı önemli.
```

## AWK

# awk data analiz ettiğimiz basit bir komut aracı değil aynı zamanda kompleks bir programlama dilidir. Biz sadece komut aracı olarak kullanacağız.

```
# komut mantığı:
$ awk pattern { action }
# ilk bakışta garip geliyor. pattern dışarıda, action {} içerisinde.
# pattern TRUE ise, statement'in içindeki action işlenir.
# eğer action vermezseniz, pattern'a uyan tüm pattern print edilir.
```

```
$ awk '{ print $0 }' example.bed
# $0 tüm sütunları print eder.
```

```
$ awk '{ print $2 }' example.bed
# ikinci satırı print eder (cut -f2 gibi)
```

```
$ awk '{ print $2 "\t" $3 }' example.bed
# cut f2,3 ye eşdeğer.
# "\t" tab-delimited olarak print ediyor.
```

Bir data dosyasında kaç adet sütun var, öğrenmek istiyoruz:

```
$ tail -n +6 Mus_musculus.GRCm38.75_chr1.gtf | head
$ tail -n +6 Mus_musculus.GRCm38.75_chr1.gtf | awk -F "\t" '{print NF; exit}'
```

#NF: "number of fields" için kullanılan bir awk özelliği.  
# -F dosyanın tab-aralıklı veya virgül-aralıklı olduğuna göre önemli bir parametre.

# awk ile aritmetik işlemler: +, -, \*, /, %, ^  
\$ awk '\$3 - \$2 > 18' example.bed  
# belli bir nükleotid uzunluğundan kısa/uzun genleri bulmak vb.

# tek bir pattern değil, birçoklarıyla pattern'i birbirine ekleyerek mantıksal operasyonlar kullanarak daha spesifik aramalar/işlemler yapabiliriz. Logical operators: && (AND), || (OR), and ! (NOT).

kromozom 1'de satır uzunluğu 10'u aşan satırları bulmak:  
\$ awk '\$1 ~ /chr1/ && \$3 - \$2 > 10' example.bed  
# Bu komutta ilk sütunu (\$1) ve (&&) chr1'i seçiyoruz.  
# regular expression'lar '/' arasına yazılıyor.  
# ~ eşleş (match) anlamına geliyor

# sadece kromozom 2 ve 3 için özelliklerin uzunluğunu bulmak:  
\$ awk '\$1 ~ /chr2|chr3/ { print \$0 "\t" \$3 - \$2 }' example.bed  
# Ciktıya bakıp işlemi nereye verdigine dikkat edin.

# daha önce, awk ve cut ile belirli sütunlar üzerinde çalıştık. Peki sadece belirli satır aralıklarına ulaşmak istersek: örn. 3-5 (3, 4, 5) satırları istiyorsan:  
\$ awk 'NR >= 3 && NR <= 5' example.bed

AWK için son örneklerin hepsi yeast\_genome.gff data dosyası üzerinden yapılacak.  
yeast\_genome.gff dosyasını suradan indirebilirsiniz:  
[https://downloads.yeastgenome.org/curation/chromosomal\\_feature/saccharomyces\\_cerevisiae.gff](https://downloads.yeastgenome.org/curation/chromosomal_feature/saccharomyces_cerevisiae.gff)

# dataya bakın  
\$ less yeast\_genome.gff

# input dosyası içerisine awk ile bakmak:  
\$ awk '{print}' yeast\_genome.gff

# bir dosyanın ilk sütununa bakmak:  
\$ **awk '{print \$1}' yeast\_genome.gff**

# önce 5 sonra 1 no'lu sütunlara bakmak:  
\$ awk '{print \$5, \$1}' yeast\_genome.gff

# genomik özelliklerin ilk-son koordinatlarını yazdırın - 10000 nükleotid'den sonrasının:  
\$ awk '{if (\$4 > 10000) print \$4, \$5}' yeast\_genome.gff

# genomik ozelliklerin ilk-son koordinatlarini ve uzunluklarini yazdirin:

```
$ awk '{print $4, $5, $5 - $4 + 1}' yeast_genome.gff
```

# aralik birakmadan yapabilirsiniz:

```
$ awk '{print $4,$5,$5-$4+1}' yeast_genome.gff
```

# 'intron' iceren satirlari yazdirin:

```
$ awk '/intron/ {print}' yeast_genome.gff
```

# 'print' konumu ortuk olarak kullanilabilir:

```
$ awk '/intron/' yeast_genome.gff
```

# ozelliklerin hangi satirlarda oldugunu 'NR'i kullanarak yazdirin (NR = number of records)

```
$ awk '{print NR, $3}' yeast_genome.gff
```

# 'END' komutu ile sadece toplami yazdirin:

```
$ awk 'END {print NR}' yeast_genome.gff
```

# UTR icinde bulunan tum intronlarin uzunluklarinin toplamini yazdirin:

```
$ awk '/UTR_intron/ {lengths += ($4-$3+1)} END {print "total length of UTR  
introns = " lengths}' yeast_genome.gff
```

-----  
EXERCISE (zaman kalirsa):

COK FAYDALI, PRATIK BIR ORNEK:

Elinizde en son sekansladiginiz evrim deneyi sonucu elde ettiginiz mayalarin (S. cerevisiae) kromozom IV'un sonunda bir bolge hic ORF icermiyor. Bu olagan bir durum mu yoksa bir tur noncoding sekans yayilmasi (expansion) mi var? Database'de bulunan bir maya genomunu indirdiniz: 'yeast\_genome.gff'. 16 kromozom var. Sadece kromozom IV'un sonunda genomik ozelliklere bakip karsilastiracaksiniz. Ornegin son 20 genomik ozellik deneyinizde evrilen mayada oldugu gibi ORF iceriginden yoksun mu?

Her seferde bir command ile grep (1), sort (2), head -n 20 (3), cut (4), sort -tekrar- (5), uniq (6) adimda 6 farkli input & output kullanarak yapabilirsiniz elbette. Ornegin 400 .gff file

uzerinde  
yaptiginizi dunusunun...

Pipeline ile tum girdi-cikti kaydindan, hem de teker teker type etmekten kurtuluyorsunuz.

```
$ grep -E "^chrIV" yeast_genome.gff | sort -n -k 4 | head -n 20 | cut -f 3 | sort | uniq
```

iki uc farkli icerige mi sahip?

```
$ grep -E "^chrIV" yeast_genome.gff | sort -r -n -k 4 | head -n 20 | cut -f 3 | sort | uniq
```