## PAC2. Anàlisi bioinformàtic amb la terminal

Magí Bas

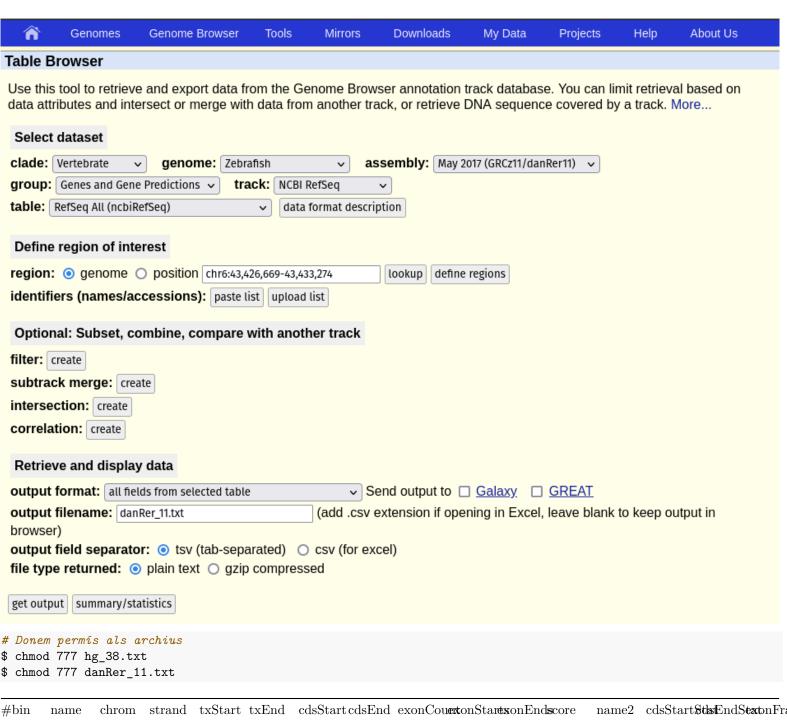
2023-04-23

## Exercici 1 – Descripció dels catàlegs de gens (15%)

Anàlisi comparativa de les col·leccions de gens de diverses espècies. De la mateixa manera que heu vist durant els exercicis pràctics, connecteu-vos al servidor UCSC per accedir als fitxers refSeq.txt de diverses espècies. L'objectiu és que empleneu la següent Taula amb les dades que obtindreu fent servir les comandes apropiades en el vostre terminal sobre els catàlegs de gens per a aquesta versió dels genomes. Us hem anotat dues espècies. Afegiu una petita interpretació biològica dels resultats anotats (1 punt):

Genoma	Chr	Gens	Tr	TrC	TrNC	Tr/Gen	Ex/Tr	nuc/Tr
H. sapiens (hg38) Zebrafish (danRer11)								
S. cerevisiae (sacCer3)	17	6125	6125	5983	123	1,00	1,058	1467,7
D. melanogaster (dm6)	8	17202	3446	3070	374	2,00	5,4687	9940,3
			3	4	6			

Importem els archius rerSeq des de UCSC:



#bin	name chrom strand	txStart txEnd cdsStartcdsEnd exonC	ouentonStantxonEndscore name2 cdsStart&taxtEndStaxtonFra
0	XM_0116341469.2	670921646710907267093004671033825	670921645 <b>679962</b> \$)65 <b>679962</b> \$) <b>1666996828</b> \$(\$ <b>5708962</b> }6;71 <b>0</b> 92)720,-
0	XM 0170011276.2	670921646713122767093004671272409	1 67092164;6 <b>799623</b> ;6 <b>799623;66790623;67103352</b> (6 <b>7101£3</b> 40 <b>;</b> 7
	_		1
0	XM_0116441467.2	670921646713122767093004671272409	670921645 <b>6799623</b> 45 <b>67996231666906235</b> 5 <b>67023345</b> 766 <b>701254</b> 405 <b>6</b>
0	NM_001276352.2	670921646713497067093579671272409	670921645 <b>779962</b> 436 <b>7796323367703B5Z6</b> 7 <b>711635</b> 367 <b>7235</b> 05363
0	NNA 001dFd9F1 0	6H0001 6 HH10 40HMH00000 HH10H04M	1,-1
0	NM_001 <b>276</b> 351.2	670921646713497067093004671272408	67092164, 67095234, 67096251, 67115351, 671257

```
# Modifiquem la taula perquè ens quedi amb les columnes que ens interesen
$ gawk 'BEGIN {FS="\t"; OFS =","} {print $13 "\t" $1 "\t" $2 "\t" $3 "\t" $4 "\t" $5 "\t" $6"\t"}' hg_38.txt > hg3
```

name2	#bin	name	chrom	strand	txStart	txEnd
Clorf141	0	XM_011541469.2	chr1	-	67092164	67109072
C1orf141	0	$XM\_017001276.2$	$\operatorname{chr}1$	-	67092164	67131227
C1orf141	0	$XM\_011541467.2$	$\mathrm{chr}1$	-	67092164	67131227
C1orf141	0	${\rm NM}\_001276352.2$	chr1	_	67092164	67134970

name2	#bin	name	chrom	strand	txStart	txEnd
C1orf141	0	NM_001276351.2	chr1	-	67092164	67134970
C1orf141	0	$XM\_011541465.3$	chr1	-	67092164	67134970
C1orf141	0	XM_011541466.3	$\operatorname{chr}1$	-	67092164	67141646
C1orf141	0	$NR\_075077.2$	chr1	-	67092164	67134970
C1orf141	0	$XM\_047420474.1$	$\operatorname{chr}1$	-	67096250	67131227
name2	$\#\mathrm{bin}$	name	$\operatorname{chrom}$	strand	txStart	txEnd

```
# Eliminem les capçaleres
$ sed -i "1d" hg38modif1.txt
$ head -5 hg38modif1.txt
C1orf141
            0
                XM_011541469.2 chr1
                                             67092164
                                                         67109072
C1orf141
                XM_017001276.2
                                             67092164
                                 chr1
                                                         67131227
C1orf141
            0
                XM 011541467.2
                                             67092164
                                                         67131227
                                 chr1
                NM_001276352.2
                                                         67134970
C1orf141
            0
                                 chr1
                                             67092164
C1orf141
            0
                NM_001276351.2
                                 chr1
                                             67092164
                                                         67134970
$ gawk '{print $4}' hg38modif1.txt |sort| uniq | grep -v "_" | wc -1
24
# Porque se tiene en cuenta el comosoma mitocondrial
```

Ens surt un total de 24 parells de cromosomes, però nosaltres sabem que en H. sapiens el número de cromosomes es 22 més el sexual (xy)/(xx). El resulat és de 24 perquè estem contant també el cromosoma mitocondrial que també esta en la base de dades.

```
$ gawk '{print $4}' danRer11_modify.txt |sort| uniq | grep -v "_" | wc -l
26
# Cuenta también el cromosoma Me
```

Aquí pasa el mateix. El peix zebra té 24 + 1 parells de cromosomes. I 26 si contem també el mitocondrial.

```
# nº de gens
gawk '{print $1}' hg38modif1.txt |sort| uniq | wc -l
42776

$ gawk '{print $1}' danRer11_modify.txt |sort| uniq | wc -l
30419
```

Veiem com el número de gens es més gran en H. sapiens però sabem que no correlaciona segons la complexitat de l'organisme.

```
# Transcrits
$ gaWk '{print $4}' hg38modif1.txt |sort| uniq | wc -1
521
$ gaWk '{print $4}' danRer11_modify.txt |sort| uniq | wc -1
1451
```

Com podem veure un número més gran de gens no repercuteix sempre en un número més àmpli de transcrits.

\$ awk '{print \$1;}' hg38modif1.txt |sort| uniq -c | gawk 'BEGIN{t=0}{t=t+\$1}END{print t/NR}'

```
# nº transcripts codificants
$ gawk '{print $3}' danRer11_modify.txt |sort| uniq | grep "NM" |wc -1
15393
$ gawk '{print $3}' hg38modif1.txt |sort| uniq | grep "NM" |wc -1
66826
```

No tots els transcrits codifiquen per a una proteïna,

```
# nº transcripts no-codificants
$ gawk '{print $3}' hg38modif1.txt |sort| uniq | grep "NR" |wc -1
20584
$ gawk '{print $3}' danRer11_modify.txt |sort| uniq | grep "NR" |wc -1
476

# existen campos sin "NR" ni "NM"
$ gawk '{print $3}' hg38modif1.txt |sort| uniq | grep -v "NM" |wc -1
117663
```

```
$ awk '{print $1;}' danRer11_modify.txt |sort| uniq -c | gawk 'BEGIN{t=0}{t=t+$1}END{print t/NR}'
2.14402
# nº exones/transcrito
# nº nucleotidos/transcrito
$ gawk 'BEGIN {FS=0FS="\t"; total=0; count=0} {len=$7-$6+1; total+=len; count++} END {print "Número promedio de nu Número promedio de nucleótidos por transcrito: 76371.7
$ gawk 'BEGIN {FS=0FS="\t"; total=0; count=0} {len=$7-$6+1; total+=len; count++} END {print "Número promedio de nu Número promedio de nucleótidos por transcrito: 40588.3
```

## Exercici 2 – Extracció de dades del catàleg OMIM (20%)

OMIM (Online Mendelian Inheritance in Man) és un catàleg de gens, trastorns i trets genètics, amb especial atenció en la relació gen-fenotip.

En aquest exercici es subministra tres dels fitxers que conformen aquesta Basede Dades amb els quals podreu contestar les preguntes que es realitzen.

Imagineu que voleu esbrinar:

4.56564

El nombre de gens associats a la malaltia d'Alzheimer que es troben en els cromosomes sexuals. (5 punts)

El nomb	ore de	e gens associats	s a la mal	altia d'Alz	sheimer qu	e es trob	en en els ci	romosom	es sexuals	(5 punts)		
<pre>\$ head</pre>	-10	genemap2.txt	;									
# Chron	nosom	le Genomia	. Positio	on Start	Genomic	Positio	on End	Cyto Lo	cation	Computed	Cyto Location	MIM Number
chr1	0	27600000	1p36	6074	413 AD70	CNTP Alz	heimer di	sease no	euronal d	thread pro	otein	
chr1	0	27600000	1p36	612	367 ALPO	JTL2 Alk	aline pho	sphatase	e, plasma	a level of	f, QTL 2	100196914
chr1	0	123400000	1p	606788	ANON1	Anorexi	ia nervosa	, susce	ptibility	y to, 1	171514	{Anore:
chr1	0	27600000	1p36	6054	462 BCC1	l Bas	al cell c	arcinoma	a, susce	ptibility	to, 1 10	0307118
chr1	0	27600000	1p36	6069	928 BMNI	)3 Bon	e mineral	density	y QTL 3	246259	<pre>?another 1</pre>	ocus at 3p21
chr1	0	2300000 1p3	36.33	618815	C1DUPp36	3.33, DU	P1p36.33	Chromoso	ome 1p36	.33 duplic	cation syndrome	e, ATAD3 gene
#	Ger	nomic Genomi	c	Comput	$\overline{\mathrm{ed}}$							Mouse
Chro-	Pos	si- Posi-	Cyto	Cyto	MIM	Gene		Approv	redEntrez	Ensembl		Gene
mo-	tion	n tion	Loca-	Loca-	Num-	Sym-	Gene	Sym-	Gene	Gene		Sym-
some	Sta	rt End	tion	tion	ber	bols	Name	bol	ID	ID	CommentPheno	otyp <b>e</b> xol/ID
chr1	0	2760000	0 1p36		612367	ALPO7	TLA lkalina		1001969	01/linkage	∫ Alkaline	

Chro- mo- some	Posi- tion Start	Posi- tion End	Cyto Loca- tion	Cyto Loca- tion	MIM Num- ber	Gene Sym- bols	Gene Name	Approv Sym- bol	ed Entrez Gene ID	Ensembl Gene ID	Gene Sym- CommentPhenotypesol/ID
chr1	0	276000	00 1p36		612367	ALPQT	L'Alkaline phos- phatase, plasma level of, QTL 2		1001969	14inkage with rs178032	{Alkaline phos- 4phatase, plasma level of, QTL 2}, 612367 (2)
chr1	0	123400	0001p		606788	ANON1	Anorexia ner- vosa, suscep- tibility to, 1	a	171514		{Anorexia nervosa, susceptibility to, 1}, 606788 (2)

#		c Genomic		Comput		-						Mouse
Chro-	Posi- tion	Posi-	Cyto	Cyto Loca-	MIM Num-	Gene Sym-	C	Approve		Ensembl Gene		Gene
mo- some	Start	$_{ m End}$	Loca- tion	tion	ber	bols	Gene Name	Sym- bol	$_{ m ID}$	ID	Comment Phenoty	Sym- pesol/ID
chr1	0	27600000	1p36		605462	BCC1	Basal		1003071	1&ssociate	d{Basal	· ·
			-				cell			with	cell	
							carci-			rs753887	6carci-	
							noma,				noma,	
							suscep-				suscep-	
							tibility				tibility	
							to, 1				to, $1$ },	
											605462	
											(2)	
chr1	0	27600000	1p36		606928	BMND3		246259			[Bone	
							min-				min-	
							eral				eral	
							den-				den-	
							sity				sity	
							QTL 3				QTL	
											3], 606928	
											(2)	
chr1	0	2300000	1n26 22		618815	C1DHD <sub>r</sub>	<b>36133</b> ,mos	zomo		Chromos		
CIII I	U	2300000	1po0.55		010010	DUP1p3	,	some		1p36.33	ome	
						D01 1p3	dupli-			dupli-		
							cation			cation		
							syn-			syn-		
							drome,			drome,		
							ATAD3			ATAD3		
							gene			gene		
							cluster			cluster,		
										$618815^{'}$		
										(4),		
										Auto-		
										somal		
										domi-		
										nant		
#		c Genomic	Cyto	Comput		Gene	Gene		$\operatorname{edEntrez}$		Comment Phenoty	
Chro-	Posi-	Posi-	Loca-	Cyto	Num-	Sym-	Name	Sym-	Gene	Gene		Gene
mo-	tion	tion	tion	Loca-	ber	bols		bol	ID	ID		Sym-
some	Start	End		tion								bol/ID

```
$ gawk '/Alzheimer/ && ($1 == "chrX" || $1 == "chrY") {print $1}' genemap2.txt | uniq -c
1 chrX
```

El nombre de gens associats a una herència autosòmica dominant i en el fenotip està definida com a síndrome i el gen es troba etiquetat amb el format HGNC. (5 punts)

```
$ gawk '{print $1}' genemap2.txt | cat genemap2.txt | grep "dominant" -A 0 | grep "syndrome" -A 0 | grep -v "^--$
589
```

Determinar quins són els gens que estan associats a la BRCA2 (5 punts)

```
$ awk -F'\t' '$8 ~ /BRCA2/ || $9 ~ /BRCA2/ {print $7}' genemap2.txt
BCCIP, TOK1
BRCA2, FANCD1, BROVCA2, GLM3, PNCA2
PALB2, FANCN, PNCA3
CNTROB, LIP8
BRCC3, BRCC36
```

Els gens associats amb BRCA2 que apareixen a la columna 7 són:

- BCCIP
- BRCA2

- FANCD1
- BROVCA2
- GLM3
- PNCA2
- PALB2
- FANCN
- PNCA3
- CNTROB
- LIP8
- BRCC3
- BRCC36

Quants gens associats al càncer de pit es troben etiquetats a OMIM (5 punts)

```
$ gawk '{print $6}' genemap2.txt | cat genemap2.txt |
grep -i "breast" -A 0 | grep -i "cancer" | grep -v "^--$" | wc -1
36
```

Trobem 36 gens que estan relacionats amb el càncer de pit a OMIM