

# UPF Bioinformatics course projects

## Students guide 2017

Didac Santesmasses (PhD) [didac.santesmasses@crg.eu](mailto:didac.santesmasses@crg.eu)

Aida Ripoll (Master student) [aida.ripoll@crg.eu](mailto:aida.ripoll@crg.eu)



Bioinformatics and genomics programme  
Roderic Guigó's group  
Centre for Genomic Regulation, Barcelona



# Protocol overview

## Tools:

- **BLAST** - typically **tblastn**
- **Exonerate** - protein2genome mode
- **Genewise**
- **T-coffee**



**U** ESTUDIS DE CIÈNCIES DE LA SALUT I DE LA VIDA

**S13. Elaboració de pàgines Web**  
Professor: Toni Gabaldón  
grups 1,2: 16 d'octubre, 08:40 (61.303).  
grups 3,4: 17 d'octubre, 08:40 (61.303).

**S14. Anotació de genomes (I)**  
Professor: Toni Gabaldón  
grups 1,2: 17 d'octubre, 13:10 (61.303).  
grups 3,4: 17 d'octubre, 16:10 (61.329-331).

**S15. Anotació de genomes (II)**  
Professor: Toni Gabaldón  
grups 1,2: 18 d'octubre, 13:10 (61.303).  
grups 3,4: 18 d'octubre, 09:40 (61.303).

**S16. Genome Browsers**  
Professor: Toni Gabaldón  
grups 1,2: 18 d'octubre, 16:10 (61.303).  
grups 3,4: 25 d'octubre, 18:10 (61.303).

**S17. El Projecte ENCODE**

<http://bioinformatica.upf.edu/>

- Webserver with **SECISearch3** and **Seblastian**:

<http://seblastian.crg.es/>

# Protocol overview

## 1st step: Get selenoprotein sequences

- **SelenoDB 2.0 (and 1.0)** **SelenoDB**  
<http://www.selenodb.org> (2.0; automatic annotation)  
<http://www1.selenodb.org> (1.0; manually curated, less species)
- **Protein databases**  
<https://www.ncbi.nlm.nih.gov/protein/>  
<http://www.uniprot.org>
- **Past year projects:**  
<http://bioinformatica.upf.edu/>

# Genomes

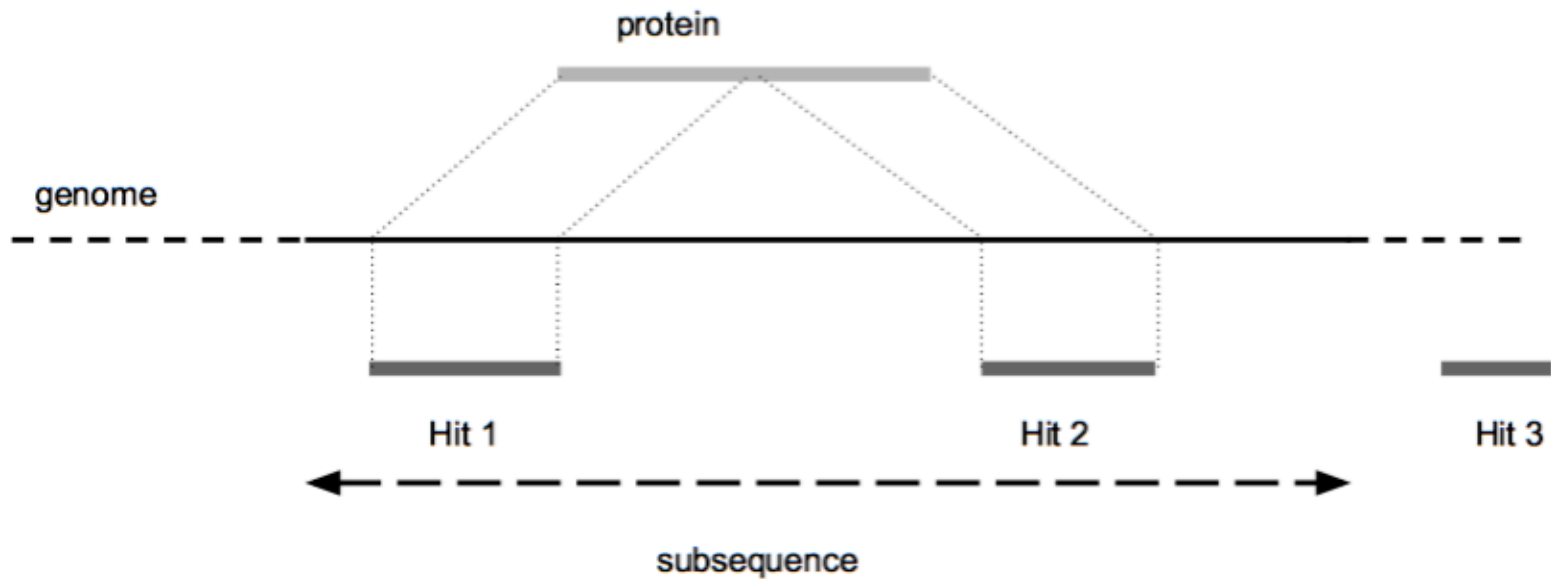
**Genomes** (blast formatted and indexed)

`/cursos/20428/BI/genomes/2017/Genus_species/genome.fa`

`/cursos/20428/BI/genomes/2017/Genus_species/genome.index`

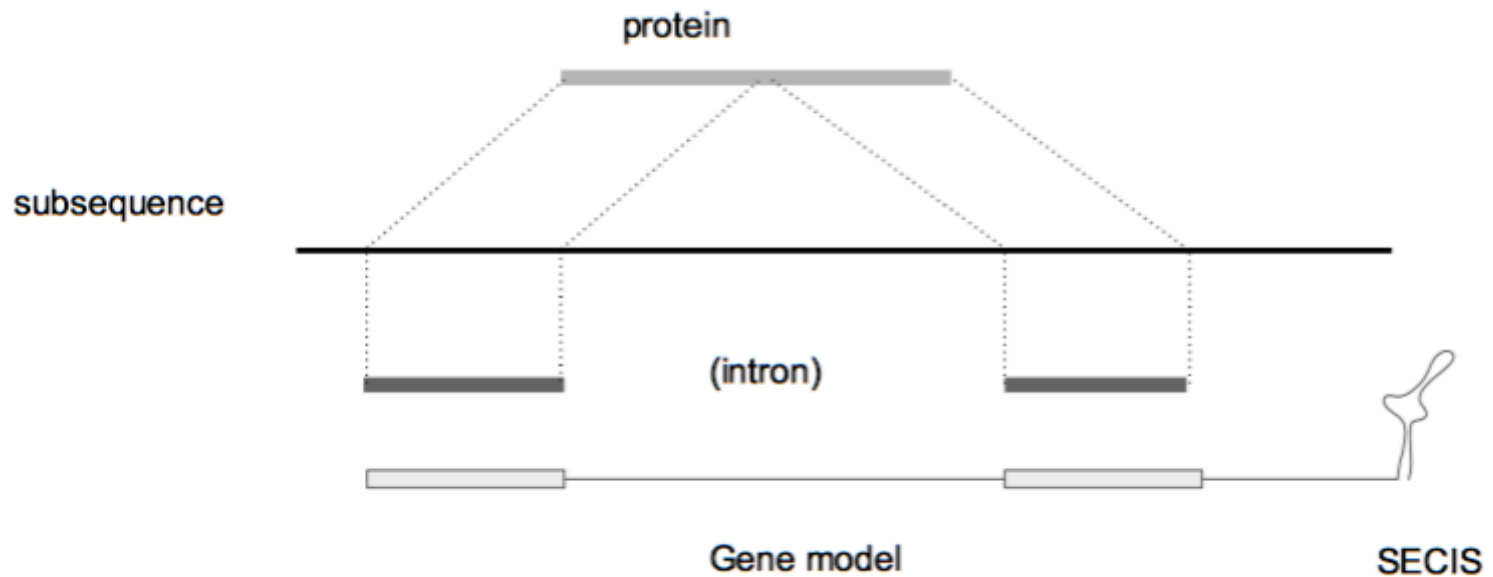
# Protocol overview

- **Tblastn**: locate gene exons (independent blast hits)



# Protocol overview

- **Exonerate** or **genewise**: multi-exonic gene model
- **Seblastian**: SECIS prediction



# Protocol overview

## Gene finding tools: fastasuite (exonerate)

- **Fastafetch:** extracting a single sequence from a multifasta (requires previous run of fastaindex)
- **Fastasubseq:** getting a subsequence of a single sequence, careful with indexes, 0-based! Transform gene positions to absolute coordinates.
- **Exonerate/Genewise:** predict the gene and align it with the sequence of the selenoprotein that encodes, and also recognizes the exons.
- **FastaSeqFromGFF:** obtain the cDNA sequence that encodes the final protein. We get it from the subsequence and the file that contains the exons.
- **Fastatranslate:** translate coding sequences careful with the selenocysteine codon character! It is a good idea to substitute the “\*” with “X” or “U” multiple sequence alignment programs just ignore “\*”

# Protocol overview

- **Tcoffe:** compare two sequences, in this case we compare the known sequence (query protein) with the homologous sequence of the the genome (predicted protein).



# Protocol overview

**Seblastian:** Predict SECIS in the 3'UTR, then upstream for selenoprotein coding sequences

Vadim Gladyshev's lab

**Selenoprotein prediction server**

Roderic Guigo's lab

Mouse over the forms to display help information

SECIS prediction  
SECISearch3

search also complementary strand  
 filter improbable structures  
 generate SECIS images (dpi: 150)  
 predict SECIS type

SECISearch3 method:

Infernal  
score threshold: 10

Covels

Original SECISearch

Upload your sequence file:  
 no file selected  
or paste it here:

Selenoprotein prediction  
Seblastian

Search for: known selenoproteins

upstream sequence length: 5000

blastx evalue threshold: 1e-3

maximum SECIS distance: 3000

output all SECIS elements

*Note: as SECISearch3 is run as a first step, all options on the left are also considered for Seblastian.*

[About](#) | [Contact us](#)

# UPF Human Biology.

## Bioinformatics Courses 2007-18

- 2007/08 – 2008/09: find all selenoproteins in a given protist genome  
2009/10 – 2011/12: find a given selenoprotein family in all protist genomes  
2012/13 – **2017/18**: find all selenoproteins in a given **vertebrate** genome

<http://bioinformatica.upf.edu/>

Projectes de l'assignatura de Bioinformàtica

Facultat de Ciències de la Salut i de la Vida

Universitat Pompeu Fabra

Curs 2012/2013

<b>1A: Ailuropoda melanoleuca</b> <i>AM. Barrios, A. Bellot, S. Castany, M. De Manuel</i>	<b>1B: Cricetulus griseus</b> <i>J. Fernandez, J. Gomez, FD. Jurquiza, A. Lopez</i>	<b>1C: Mustela putorius furo</b> <i>M. Perez, L. Taberner, G. Vilajosana, I. Villate</i>
<b>2A: Nomascus leucogenys</b> <i>M. Alemany, H. Costa, A. Escrig, I. Gafarot</i>	<b>2B: Saimiri boliviensis</b> <i>P. Garcia, J. Latorre, R. Martinez, H. Palma</i>	<b>2C: Sarcophilus harrisii</b> <i>G. Rodriguez, E. Ros, AM. Saludes, H. Xicoy</i>
<b>3A: Chrysemys picta bellii</b> <i>C. Bitlloch, G. Clua, J. Domingo, P. Gelabert</i>	<b>3B: Meleagris gallopavo</b> <i>J. Jancyte, L. Mateo, A. Olle, M. Perera, C. Perez</i>	
<b>4A: Pelodiscus sinensis</b> <i>SU. Abad, A. Almeyda, A. Azagra, R. Bartomeus</i>	<b>4B: Gadus morhua</b> <i>O. Bover, N. Cortell, B. Grau, E. March</i>	<b>4C: Latimeria chalumnae</b> <i>A. Martinez, A. Perlas, T. Robert, S. Walsh</i>

# Groups and species

	Treball	Genoma	Contacte	Supervisor
Grup 101	1	<i>Sporophila hypoxantha</i>	laia.carrete@crg.eu	Laia Carrate
	2	<i>Ammotragus lervia</i>	fernando.cid@crg.eu	Fernando Cid
	3	<i>Pogona vitticeps</i>	miki.s.t@hotmail.com / marina.lleal01@estudiant.upf.edu	Miquel Angel Schikora / Marina Lleal
Grup 102	4	<i>Oncorhynchus mykiss</i>	Marina.Marcet-Houben@crg.eu	Marina Marcet-Houben
	5	<i>Monopterus albus</i>	diego.garrido@crg.eu	Diego Garrido
	6	<i>Astyanax mexicanus</i>	aida.ripoll@crg.eu / marta.badia.mb@gmail.com	Aida Ripoll / Marta Badia
Grup 103	7	<i>Meriones unguiculatus</i>	miguelangel.naranjo@crg.eu	Miguel A. Naranjo
	8	<i>Nannopterum harrisi</i>	irene.julca@crg.eu	Irene Julca
	9	<i>Taeniopygia guttata</i>	tonidedios94@gmail.com / ramon.massoni01@estudiant.upf.edu	Toni de Dios / Ramon Massoni
Grup 104	10	<i>Spermophilus dauricus</i>	andres.lanzos@crg.eu	Andres Lanzos
	11	<i>Odocoileus virginianus texanus</i>	cecilia.klein@crg.eu	Cecilia Klein
	12	<i>Anas zonorhyncha</i>	rosamaria.fernandez@crg.eu	Rosa M. Fernandez

## Supervisors Sessions

[https://docs.google.com/spreadsheets/d/1D1jkz1pQD8NyF\\_OPpUesdEEw1KA27qQz9iaetCBgmJE/edit#gid=0](https://docs.google.com/spreadsheets/d/1D1jkz1pQD8NyF_OPpUesdEEw1KA27qQz9iaetCBgmJE/edit#gid=0)

# Project 2017-18

## selenoproteins in vertebrates

<http://bioinformatica.upf.edu/>

- **Web page:** Structure of a scientific paper
- **Wikipedia:** Species description (not english entry)
- **SelenoDB:** Insert your selenoprotein genes predictions into a real world database. Available to the scientific community.

# Technical issues

## Genomes (blast formatted and indexed)

/cursos/20428/BI/genomes/2017/Genus\_species/genome.fa

## Access to the cluster (connect through ssh)

- From informatic rooms (ip):

ssh [UXXXXXX@10.200.127.240](mailto:UXXXXXX@10.200.127.240) -X

(alumnes) password: dd/mm/yyyy (birth date)

- Outside UPF (**VPN**):

ssh UXXXXXX@[sitdoc.s.upf.edu](mailto:UXXXXXX@sitdoc.s.upf.edu) -X

(alumnes) password: dd/mm/yyyy (birth date)

## Modules

```
module load modulepath/goolf-1.7.20
```

```
module load BLAST+/2.2.30-goolf-1.7.20
```

```
module load Exonerate/2.2.0-goolf-1.7.20
```

```
module load T-Coffee/11.00.8cbe486-goolf-1.7.20
```

```
export PATH=/cursos/20428/BI/bin:$PATH
```

```
export PATH=/cursos/20428/BI/soft/genewise/x86_64/bin:$PATH
```

```
export WISECONFIGDIR=/cursos/20428/BI/soft/genewise/x86_64/wise2.2.0/wisecfg/
```

# Cluster

- **VPN connection** <https://www.upf.edu/bibtic/en/guiesiajudes/recinfo/vpn/>  
ssh UXXXXX@sitdoc.s.upf.edu (campus global password)
- File **blastJob.sh**

```
#!/bin/bash
#$ -o tblastn.stdout
#$ -e tblastn.stderr
#$ -q all.q
#$ -N blastJob
#$ -cwd
module load modulepath/goolf-1.7.20
module load BLAST+/2.2.30-goolf-1.7.20
tblastn -query fitxerquery.fa -db nombbddBLAST -out fitxerdesortida
```

# Cluster

- **Queue system**

```
$ qsub blastJob.sh
```

```
$ qstat
```

```
job-ID prior name      user      state submit/start at   queue  
  slots ja-task-ID
```

-----

-----

```
    231 0.55500 blastJob UXXXXX   r   02/10/2010 11:42:09 llicen.q@luke  
1
```

- **To kill a running job**

```
$ qdel 231
```

# Adding genes to SelenoDB

**Before** adding genes to selenodb, you need to **fill the information:**

[http://selenodb.crg.cat/selenodb\\_barcelona/add\\_author.html](http://selenodb.crg.cat/selenodb_barcelona/add_author.html)

**Then** you can use your **email as author** in:

[http://selenodb.crg.cat/selenodb\\_barcelona/add\\_gene.html.mako](http://selenodb.crg.cat/selenodb_barcelona/add_gene.html.mako)