

Identification of selenoprotein-related genes in a vertebrate genome

- Bioinformatics 2022/23 -

Marco Mariotti, Diego Garrido, Roderic Guigó



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



Selenoproteins as test case

- Selenoproteins have the particular characteristic of including a **UGA codon**, recoded because of the presence of the **SECIS element**.



BIOINFORMATICS PROJECT

Find all **selenoprotein-related genes**
in a **vertebrate genome**

- If you learn how to predict selenoproteins, you will be able to do the same with any protein family.

UPF Human Biology.

Bioinformatics Courses 2007-2022

- 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 – **2022/23**: find all selenoproteins in a given **vertebrate** genome

Projectes de l'assignatura de Bioinformàtica

Facultat de Ciències de la Salut i de la Vida

Universitat Pompeu Fabra

Curs 2021/2022

1. <i>Haplochromis burtoni</i> Supp. Miguel Sánchez miguel.sanchez@upf.edu	2. <i>Pyrgilauda ruficollis</i> Supp. Valèria del Olmo valeria.delolmo@upf.edu	3. <i>Hystrix brachyura</i> Supp. Pablo Carrón pablo.carron@upf.edu
4. <i>Atractosteus spatula</i> Supp. Marina Marín marina.marin@upf.edu	5. <i>Paroedura picta</i> Supp. Guàrdia Tereixa guardia.terexia@upf.edu	6. <i>Nanger dama</i> Supp. Miguel Rodríguez miguel.rodriguez@upf.edu
7. <i>Ophiodon elongatus</i> Supp. Jorge Garrido jorge.garrido@upf.edu	8. <i>Glandirana rugosa</i> Supp. Silvia Pérez silvia.perez@upf.edu	9. <i>Centrocercus urophasianus</i> Supp. Marina Marín marina.marin@upf.edu
10. <i>Otocyon megalotis</i> Supp. Valèria Marín valeria.marin@upf.edu		

<http://bioinformaticaupf.crg.eu>

Project 2022-2023

selenoproteins in vertebrates

- **Web page:** Structure of a scientific paper
- **Wikipedia:** Species description



Selenoproteins in *Miichthys miiuy*

[HOME](#)

[ABSTRACT](#)

[INTRODUCTION](#)

[MATERIALS AND METHODS](#)

[RESULTS](#)

[DISCUSSION](#)

[CONCLUSIONS](#)

[REFERENCES](#)

[ACKNOWLEDGMENTS](#)

[CONTACT](#)

Selenoproteins are a group of proteins characterized by the presence of, at least, one Selenocysteine (Sec) residue in its chain. Since this residue is codified by UGA, which is normally considered as a stop codon, some of this proteins are dismissed in genome databases.

Moreover, the inclusion of Selenocysteine residue depends on the presence of an element called Selenocystein Insertion Sequence (SECIS), which is a secondary mRNA structure that allows the insertion of a selenocysteine instead of a stop codon.

The aim of our study is to predict the selenoproteins of *Miichthys miiuy*, a Japanese benthic fish, performing an homology-based in silico search. In order to assess the characteristics of the *Miichthys miiuy's* selenoproteome, we have compared the genome of this species with *Danio rerio's* and *Homo sapiens's* selenoproteins annotations obtained from SelenoDB. For the prediction, different bioinformatic tools such as BLAST, Exonerate, Genewise, T_coffee, Seblastian and SECISearch3 were needed. Additionally, we have designed an automatic program to speed up the process.

Our results show a high conservation between Zebrafish' and *Miichthys miiuy'* selenoproteome. We have found 33 selenoproteins, 8 Cys-containing homologous proteins, 5 machinery proteins and 11 proteins related to selenium metabolism.

This study contributes with the identification of selenoproteins in new-sequenced organisms.



VIQUIPÈDIA
L'enciclopèdia lliure

Portada

Article a l'atzar

Articles de qualitat

Comunitat

Portal **vi**quipedista

Canvis recents

La taverna

Contacte

Xat

Donatius

Ajuda

Eines

Què hi enllaça

Canvis relacionats

Pàgines especials

Enllaç permanent

Informació de la pàgina

Element a Wikidata

Citau aquest article

Imprimeix/exporta

Crear un llibre

Baixa com a PDF

Sense sessió iniciada [Discussió per aquest IP](#) [Contribucions](#) [Crea un compte](#) [Inicia la sessió](#)

Pàgina **Discussió**

Mostra

Modifica

Modifica el codi

Mostra l'historial

Més ▾

Cerca a **Viquipèdia**



Miichthys miiuy

Miichthys miiuy és una espècie de peix de la família dels esciènids i de l'ordre dels perciformes.

Contingut [amaga]

- Morfologia
- Hàbitat
- Distribució geogràfica
- Ús comercial
- Observacions
- Referències
- Bibliografia
- Enllaços externs

Morfologia [modifica | modifica el codi]

Els mascles poden assolir 70 cm de longitud total.^{[5][6]} Com la resta de peixos de la família Sciaenidae, *M. miiuy* és conegut per tenir uns otòlits excepcionalment grans que els doten d'un sistema auditiu molt desenvolupat.^[7] Aquests peixos s'anomenen sovint peixos tambors o corballs a causa dels sons que produeixen amb les seves bufetes natatòries.

Hàbitat [modifica | modifica el codi]

És un peix de clima temperat i demersal que viu entre 15-100 m de fondària.^{[5][6]} Eviten les aigües clares, prefereixen viure en estuaris, badies i riberes de rius fangosos. Són organismes carnívors bentònics.^[7]



Miichthys miiuy

Taxonomia

Super-regne	Eukaryota
Regne	Animalia
Fílum	Chordata
Classe	Actinopterygii
Ordre	Perciformes
Família	Sciaenidae
Gènere	<i>Miichthys</i>
Espècie	<i>Miichthys miiuy</i> (Basilewsky, 1855) ^{[1][2][3]}

Nomenclatura

Sinònims

- Argyrosomus miiuy* (Basilewsky, 1855)
- Miichthys imbricatus* (Matsubara, 1937)
- Nibea imbricata* (Matsubara, 1937)
- Otolithus fauvelii* (Peters, 1881)
- Sciaena miiuy* (Basilewsky, 1855)^[4]

Bioinformatics methods for selenoprotein prediction

- ***De novo***: Selenogeneid (Castellano et al. 2001)
- **Homology-based approaches:**
 - UGA/Sec or UGA/Cys alignments (e.g. Kryukov et al. 2003)
 - Selenoprofiles (Mariotti and Guigó 2010)
 - Seblastian (Mariotti et al. 2013)
- **SECIS prediction:**
 - SECISearch (Kryukov et al. 2003)
 - SECISearch3 (Mariotti et al. 2013)
- **tRNA-Sec prediction:**
 - Secmarker

Protocol steps

- **SelenoDB 2.0**

<http://selenodb.crg.es>

- **BLAST** (tblastn)
- **Exonerate** - protein2genome mode
- **Genewise**
- **T-coffee**

- **SECISearch3** and **Seblastian**

<http://seblastian.crg.es>

S13. Elaboració de pàgines Web
Professor: Toni Gabaldón

S14. Anotació de genomes (!)
Professor: Toni Gabaldón

S15. Anotació de genomes (!!)
Professor: Toni Gabaldón

S16. Genome Browsers
Professor: Toni Gabaldón

<http://bioinformaticaupf.crg.eu>

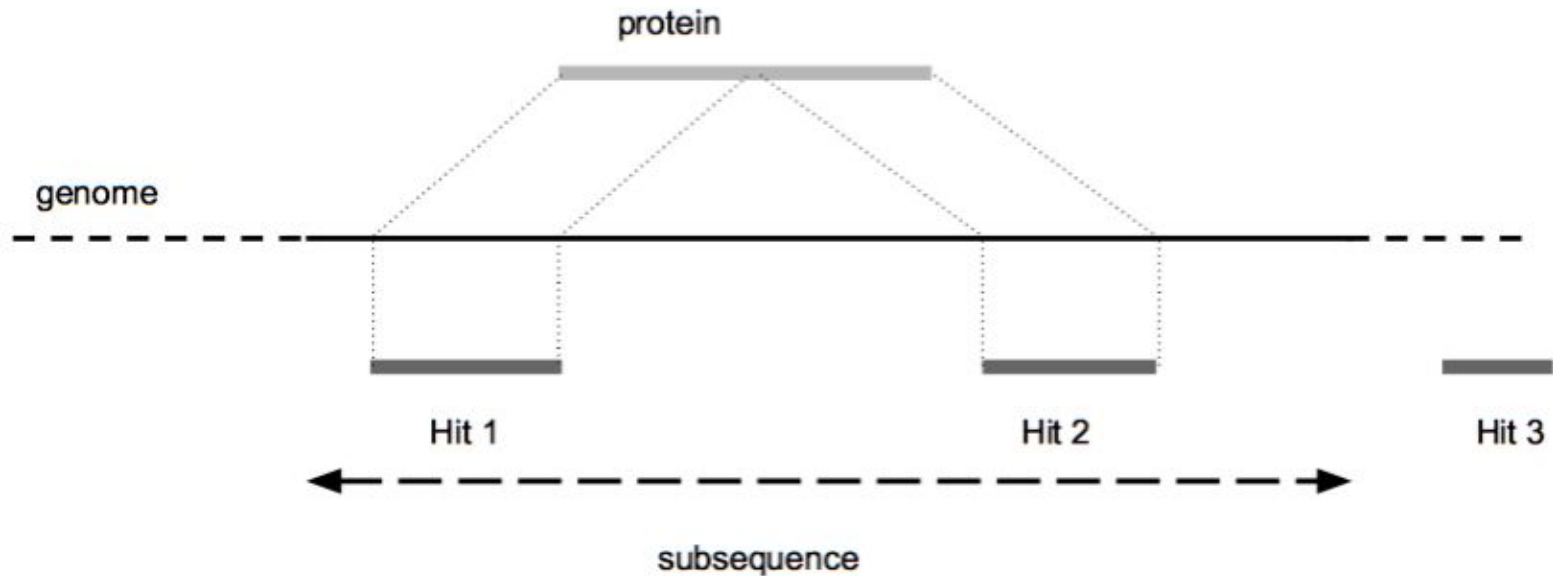
Protocol steps

Get selenoprotein sequences from a related organism

- **SelenoDB 2.0** **SelenoDB**
<http://selenodb.crg.es> (automatic annotation)

Protocol steps

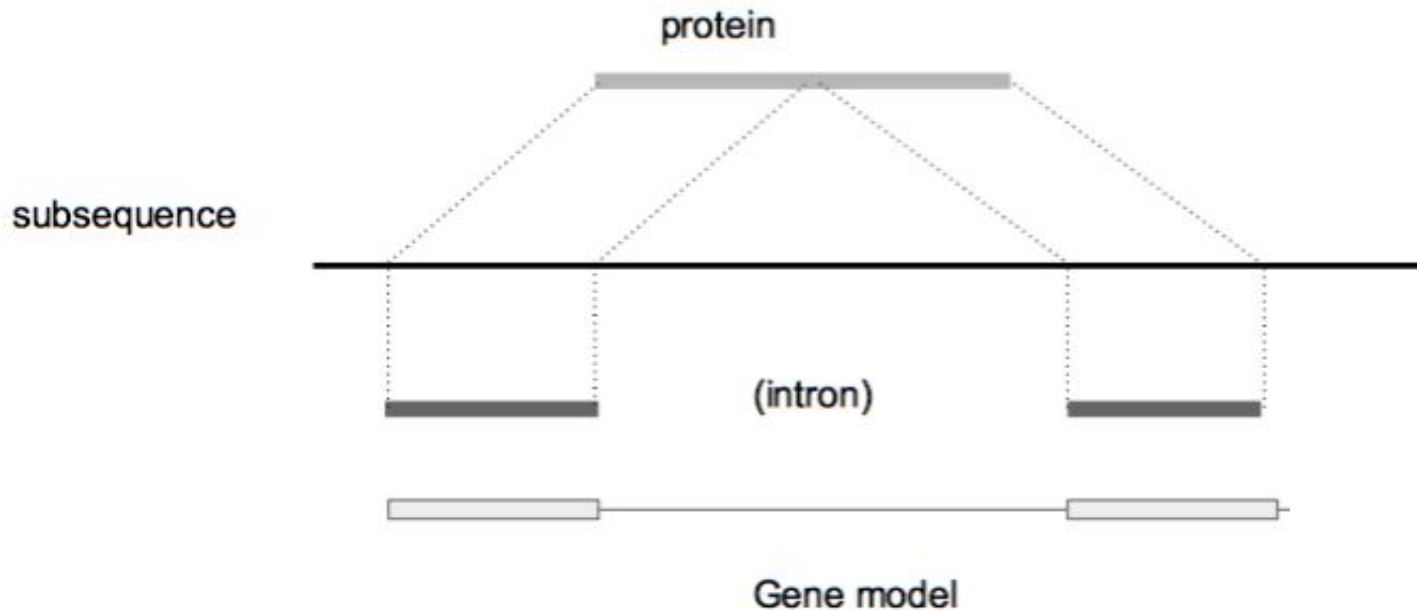
Locate gene exons (independent **tblastn** hits) in the genome of your organism



Search Name	Query Type	Database Type	Translation
tblastn	Peptide	Nucleotide	Database

Protocol steps

Build a multi-exonic gene model (**exonerate** or **genewise**)



and translate it into a protein

Protocol steps

Compare the known sequence that you obtained from the database (query protein) with the homologous sequence of the genome of your organism (predicted protein) (**t-coffee**)

```

T-COFFEE, Version_9.01 (2012-01-27 09:40:38)
Cedric Notredame
CPU TIME:0 sec.
SCORE=75
*
BAD AVG GOOD
*
1j46_A : 69
2lef_A : 67
1k99_A : 75
1aab_  : 70
cons   : 75

1j46_A MQ-----DRVKRP---MNAFIVWSRDQRRKMALENPRMRN-
2lef_A MH-----IKKP---LNAFMLYMKEMRANVVAESTLKES-
1k99_A MKKLKKHPDFPKKP---LTPYRFFMEKRKAKYAKLHPMSN-
1aab_  GK-----GDPKKPRGKMSSYAFFVQTSREEHKKKHPDASVN

cons   : ██████████ *.* ██████████ :.:. : * : ██████████ .

1j46_A -SEISKQLGYQWKMLTEAEKWPFQEAQKLOA-----MHR
2lef_A -AAINQILGRRWHALSREEQAKYELARKERQ-----LHM
1k99_A -LDLTKILSKKYKELPEKKKMKYIQDFQREKQEFERNLARFR
1aab_  FSEFSKKCSERWKTMSAKEKGFEDMAKADKA-----RYE

cons   : ██████████ :.:. ██████████ :.:. :. : : : : : : : : : ██████████ .

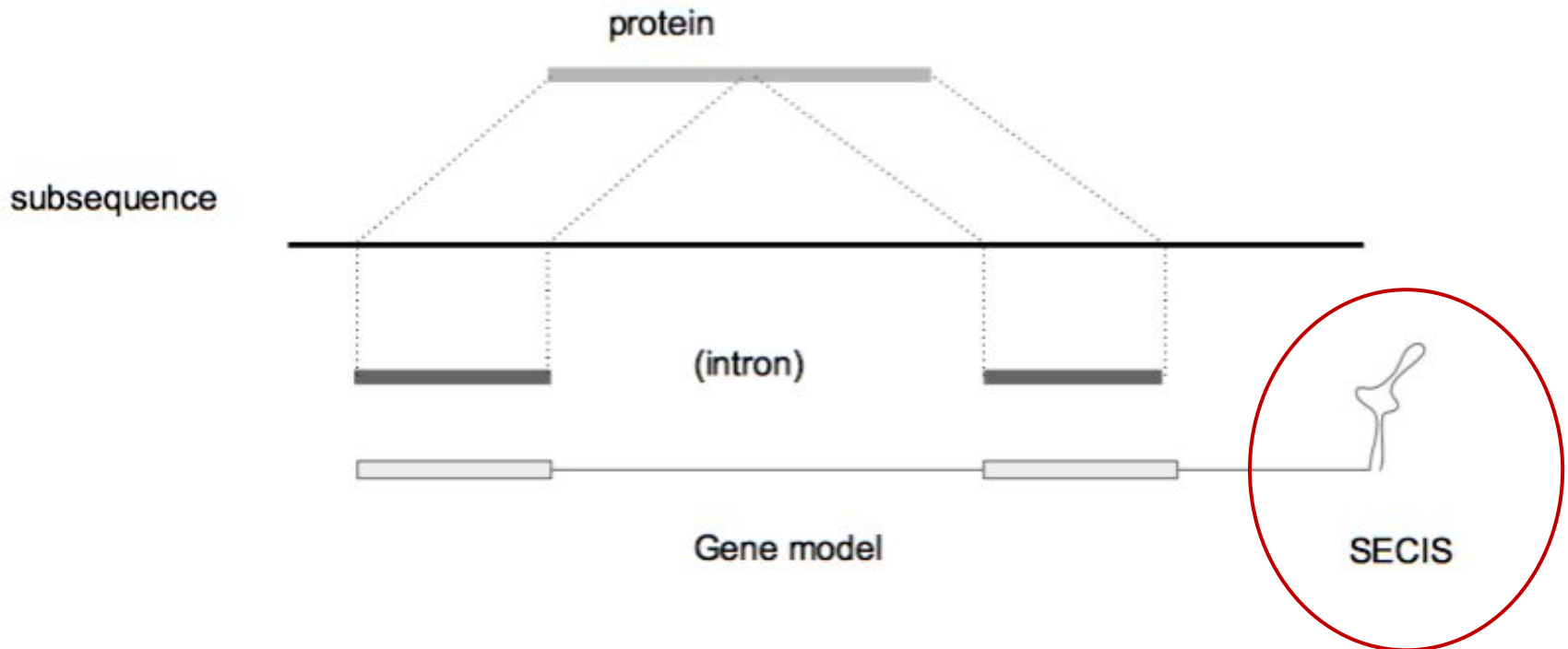
1j46_A EKYPNYKYRP---RRKAKMLPK
2lef_A QLYPGWSARDNYGKKKKRKRK
1k99_A EDHPDLIQNA-----KK
1aab_  REMKTYIPK-----GE

cons   : ██████████ ██████████ ██████████ ██████████ :

```

Protocol steps

SECIS and selenoprotein prediction (**SECISearch3** and **Seblastian**)



Protocol steps

Seblastian: Predict SECIS in the 3'UTR (using SECISearch3), and then searches upstream for selenoprotein coding sequences.

Selenoprotein prediction server

Mouse over the forms to display help information

SECIS prediction
SECISearch3

search also complementary strand
 filter improbable structures
 generate SECIS images (dpi: 15)
 predict SECIS type

SECISearch3 method:

Infernal
score threshold: 10
 Covels
 Original SECISearch

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

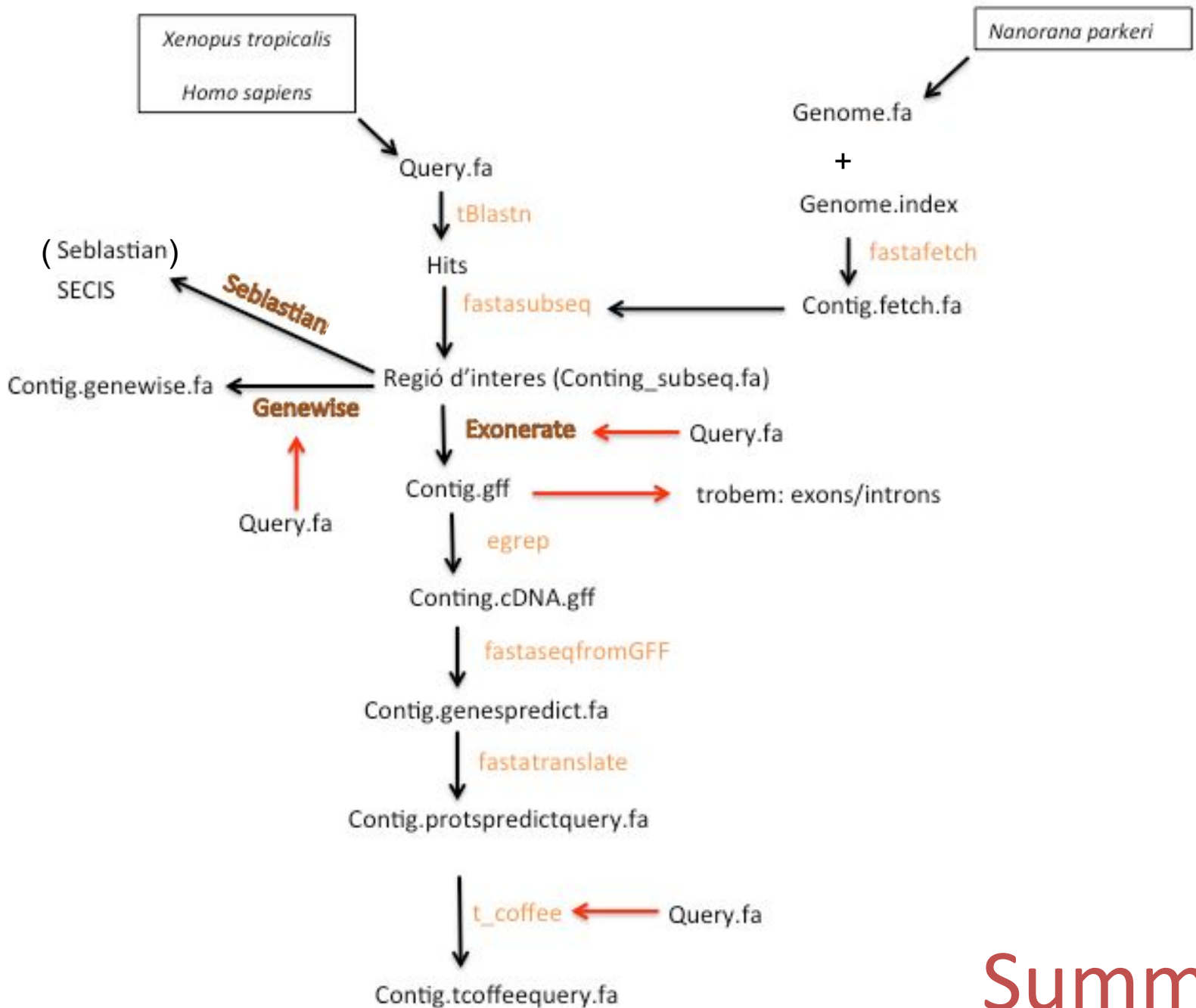
Selenoprotein prediction
Seblastian

Search for: known selenoproteins
upstream sequence length: 5000
blastx evaluate 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

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



Summary

Protocol steps

Gene finding tools:

- **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.pl:** 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" as multiple sequence alignment programs just ignore "*"

Notes for the project

- Results must be presented in a **web page** with the **structure of a scientific paper**
 - ✓ Aminoacid sequences + SECIS sequences
 - ✓ Genes in GFF format (absolute coordinates)
- All **genes** should be **as complete as possible**: starting with a AUG, ending with a STOP codon, and with an identified SECIS element downstream.
- **Ignore alternative isoforms** (if any), just choose one
- Report also the **genes** of:
 - ✓ **selenoprotein machinery**: SEPSECS, EEFSEC, PSTK, SBP2, SECP43, SEPHS1, SEPHS2, etc.
 - ✓ **Cys-containing homologs**
- **Other helpful resources** to biologically interpret and visualize the results (**phylogenetic trees**):
 - ✓ Phylogeny.fr: http://www.phylogeny.fr/simple_phylogeny.cgi (.mfa)
 - ✓ phyloT: <https://phylot.biobyte.de/> (from NCBI taxonomy .nw)
 - ✓ iTOL: <https://itol.embl.de/> (.nw)
 - ✓ Etetoolkit: <http://etetoolkit.org/treeview/> (.nw or .msa)
- In some cases, the predicted protein can be **located in more than one contigs/scaffolds**. You will notice this if you try to predict the protein in both of them, and you pay attention at both sequence alignments performed by T-coffee.
- **HTML language**
 - ✓ <https://www.w3schools.com/html/default.asp>
 - ✓ <https://getbootstrap.com/>

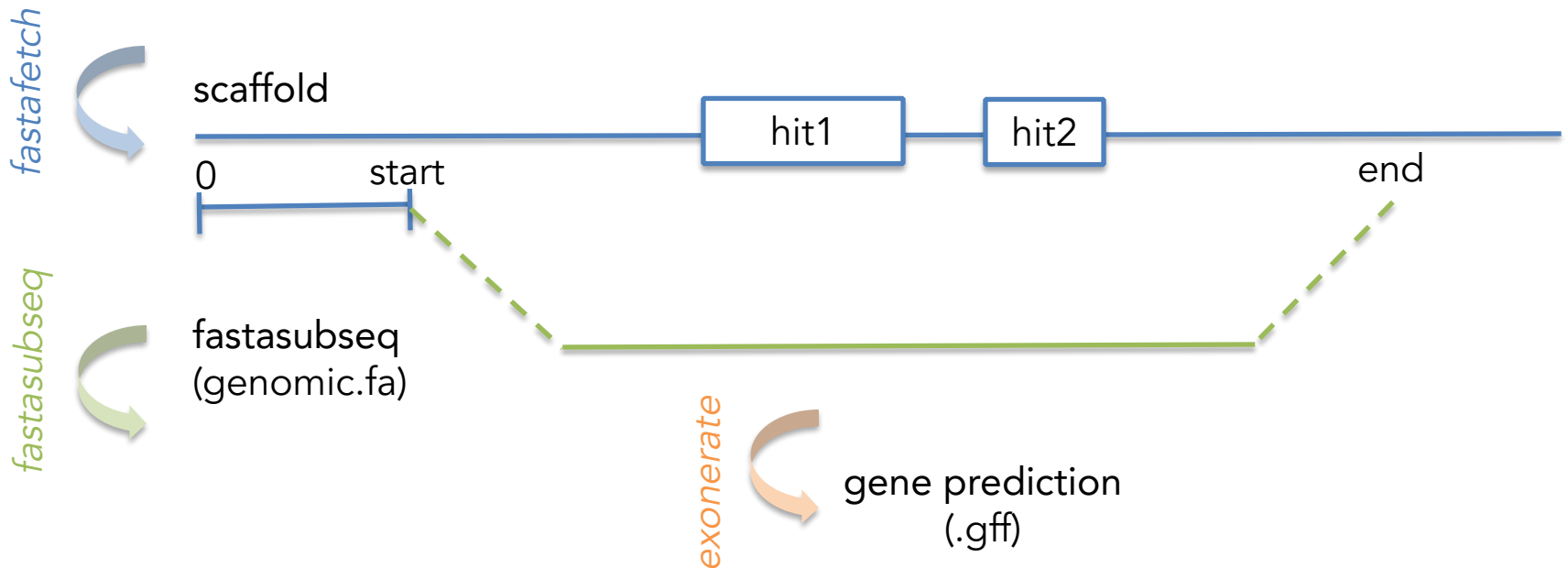
Notes for the project

- We **already provide** you, together with the genome:
 /mnt/NFS_UPF/soft/genomes/2022/**Genus_specie**
 - ✓ BLAST database for the genome
 - ✓ Indexed genome
- **Scaffolds/Contigs lengths** can be found in the *genomes.lengths* file
- **fastatranslate** (option -F 1) to consider only the 1st ORF.
- Before performing the **sequence alignment with T-Coffee**, substitute the “*” with “X” or “U” as multiple sequence alignment programs just ignore “*”
- **Sebastian and SECISearch3 web servers:**
 - Input: Nucleotide sequence (*fastasubseq* file)
 - * DO NOT take into account other nucleotide bases different than A, C, G, T, a, c, g, t, or N. Then, in case you have one of the other symbols from ambiguity code, one solution could be substituting them by an N.

Johnson A.D. An extended IUPAC nomenclature code for polymorphic nucleic acids.
Bioinformatics. 2010; 26(10): 1386-1389.

Notes for the project

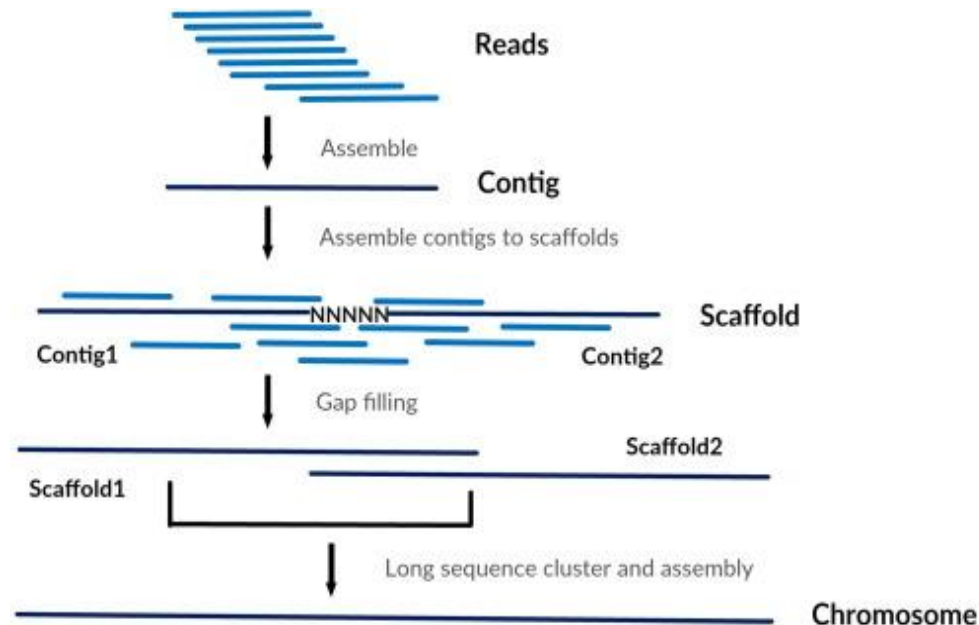
- **Genes prediction (GFF format):** Conversion of **relative** to **absolute coordinates**
 - Apart from obtaining the protein sequence predictions, you should obtain the gene predictions in .gff format considering the absolute coordinates.
 - Remember that, as you made your prediction using the *fastasubseq* file, you will be predicting the genes (.gff file from exonerate) with the relative coordinates instead of the absolute coordinates. Then, to generate the .gff files with absolute coordinates, you will have to convert the your .gff files with relatives coordinates (.gff file from exonerate) considering the **start** you decided to give to the **fastasubseq program**. [In this case, **start**: start_hit1 nt - 50.000 nt]



Notes for the project

- **Contigs and Scaffolds**

- ✓ **Contig**: a contiguous stretch of nucleotides resulting from the assembly of several reads
- ✓ **Scaffold**: several contigs stitched together with NNNs in between



Notes for the project

Automation: BASH scripting

- Basics of Bash scripting
 - ✓ [slides](#)
- Bash documentation
 - ✓ <https://www.gnu.org/software/bash/manual/bash.pdf>
- Bash cheatsheets
 - ✓ <https://devhints.io/bash>
 - ✓ <https://github.com/LeCoupa/awesome-cheatsheets/blob/master/languages/bash.sh>

Technical info

- Access to the shared folders in *fs-aules* and VPN usage -

https://bioinformaticaupf.crg.eu/accedir_carpeta_compartida_fs-aules.pdf

Evaluation

The projects will be **evaluated** based on:

- ✓ **Methods:** scripting is encouraged (different levels of automation)
- ✓ **Results:** you are expected to find all selenoprotein-related genes in your assembly
- ✓ **Discussion:** interpret your results logically
- ✓ **Presentation:** the web page should present the work as clearly as possible (including the Wikipedia entry)

Evaluation

Abstract		0.25
Introduction	Selenium, selenoprotein biosynthesis, evolution/phylogeny, families (including selenoproteins, cys-homologues and machinery), links to wikipedia entry, schemes and diagrams	0.50
Materials and methods	Queries selection, description of each step of the annotation pipeline, automatization, SECISearch3/Sebastian server	3.00
Results	Summary table/plot with protein query (+specie/species), Sec/Cys-homologue, tblastn, exonerate/genewise, scaffold, gene prediction, protein prediction, sequence alignment, SECIS and Sebastian predictions. Description of each prediction: protein location (scaffold and coordinates + strand), exons, SECIS elements and Sebastian predictions	2.00
Discussion	Discuss each family, highlight interesting cases (duplications, losses and conversions from Sec-to-Cys events), contrast hypotheses with the literature, use specific phylogenetic terms	2.50
Conclusions	Wrap up the results, further directions, pros and cons of the project, limitations	0.50
References	Citations along the text, proper numeration	0.25
Web page format	Structure of scientific paper, creativity, links working properly	0.50
Wikipedia entry	Extend the content in different languages	0.50

Groups, supervisors and species

Group	Subgroup	Supervisor	E-mail	Species
Grup 101	1	Giovanni Asole	giovanni.asole@crg.eu	<i>Bettongia penicillata</i>
	2	Sadoia Manzano	manzano.saioa@gmail.com	<i>Sceloporus occidentalis</i>
	3	Guifré Torruella	guifftc@gmail.com	<i>Clinocottus analis</i>
Grup 102	4	Hannah Benisty	hannah.benisty@crg.eu	<i>Hipposideros pendleburyi</i>
	5	Sergio Sánchez	sergio.sanchez.moragues@gmail.com	<i>Artemisiospiza belli</i>
	6	Miquel Schikora	miquelangel.schikora@irbbarcelona.org	<i>Bovichtus variegatus</i>
	7	Tamara Perteghella	tamara.perteghella@crg.eu	<i>Symphalangus syndactylus</i>
Group 103	8	Nadia Makarova	na96mk@gmail.com	<i>Hyperoodon ampullatus</i>
	9	Xavier Hernández	xavier.hernandez@crg.eu	<i>Actenoides hombroni</i>
	10	Martina Cardinale	martina.cardinali@irbbarcelona.org	<i>Darevskia valentini</i>
	11	Jaume Reig	jaume.reig@crg.eu	<i>Tragulus kanchil</i>