

Recerca de selenoproteïnes en el genoma d'organismes eucariotes

- Bioinformàtica 2019/20 -

Aida Ripoll (PhD Student)

Didac Santesmasses (PhD)



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



What are **selenoproteins**?

What are selenoproteins?

❖ Role of selenium (Se)

- One of the **nine** essential trace elements
- **Vital** functions (homeostasis)
 - Se deficiency* → pathophysiological status (Keshan disease)
 - Se excess* → toxic

What are selenoproteins?

❖ Role of selenium (Se)

- One of the **nine** essential trace elements
- **Vital** functions (homeostasis)
 - Se deficiency* → pathophysiological status (Keshan disease)
 - Se excess* → toxic

❖ Selenoproteins (Se-containing proteins)

- All **3 domains of life**: eukarya, archea, and eubacteria
- Human **selenoproteome** is encoded in 25 selenoprotein genes (but the number varies for different taxa)
- Sometimes the **orthologue** of a selenoprotein has **Cys** instead of Sec
- Mostly **redox enzymes** → antioxidant protection capacity

What are selenoproteins?

❖ Role of selenium (Se)

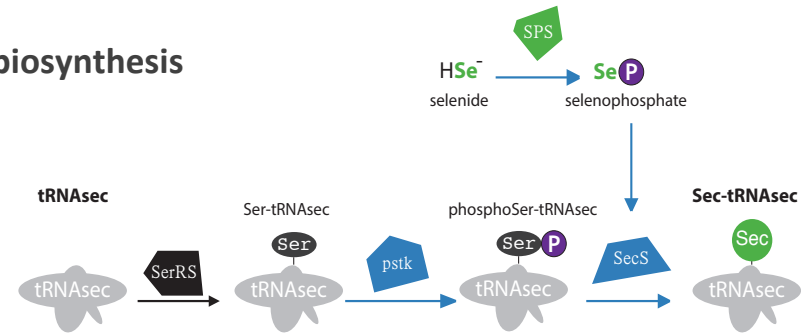
- One of the **nine** essential trace elements
- **Vital functions** (homeostasis)
 - Se deficiency* → pathophysiological status (Keshan disease)
 - Se excess* → toxic

❖ Selenoproteins (Se-containing proteins)

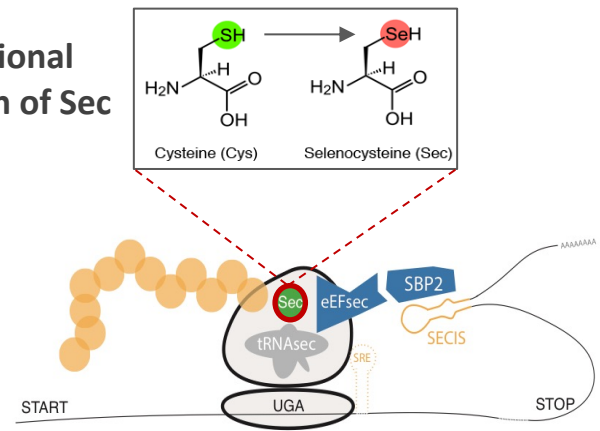
- All **3 domains of life**: eukarya, archea, and eubacteria
- Human **selenoproteome** is encoded in 25 selenoprotein genes (but the number varies for different taxa)
- Sometimes the **orthologue** of a selenoprotein has **Cys** instead of Sec
- Mostly **redox enzymes** → antioxidant protection capacity

❖ Sec biosynthesis & insertion mechanism

I. Sec biosynthesis



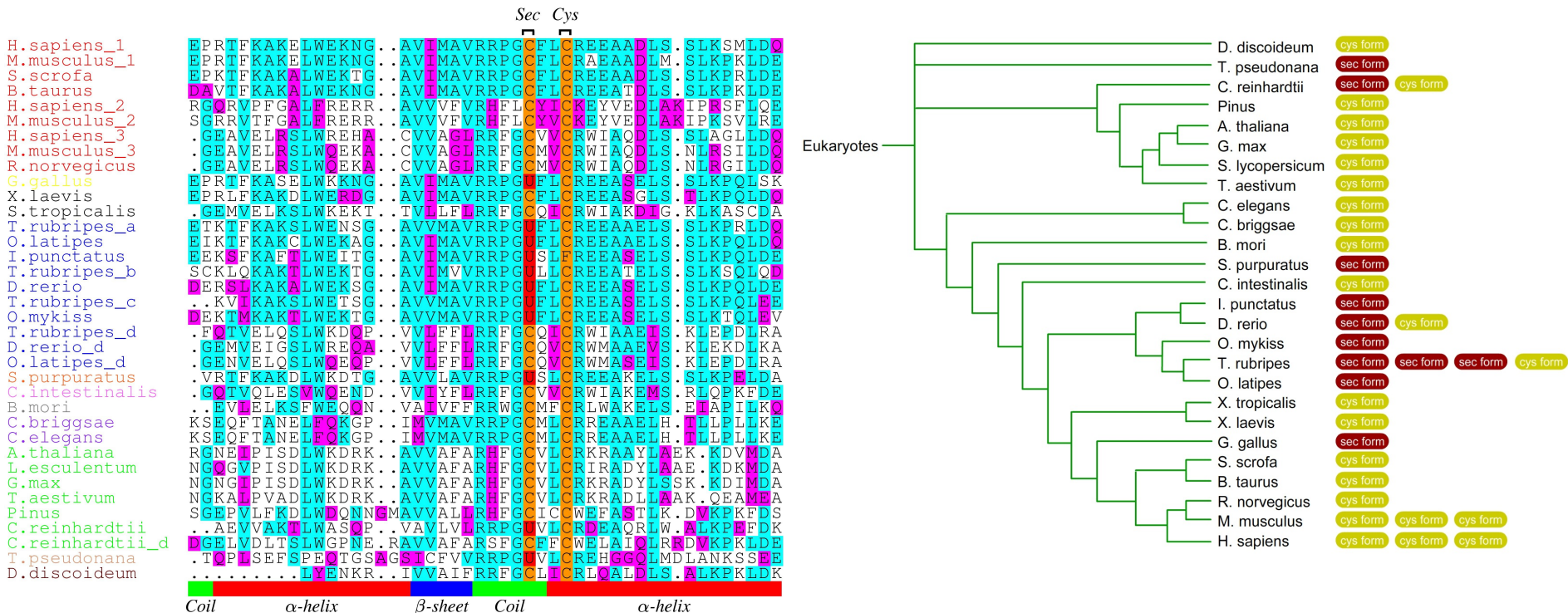
II. Cotranslational incorporation of Sec



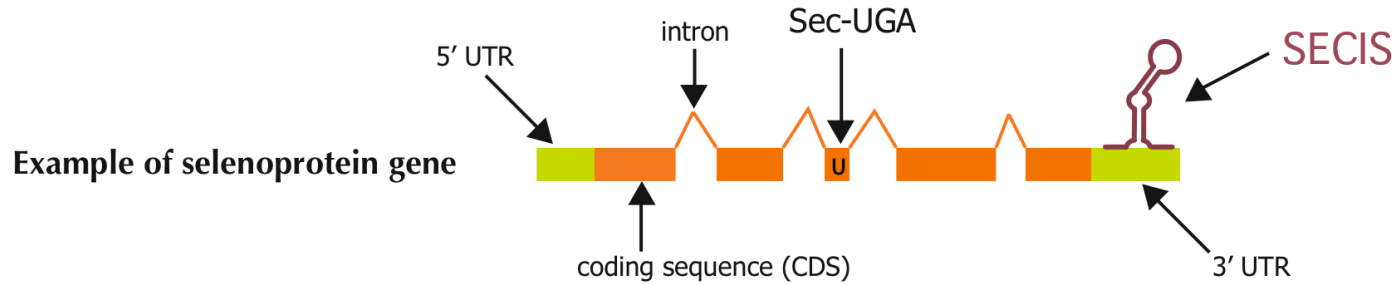
Extracted from:
 Mariotti M et al. Mol Biol Evol. 2016;33(9):2441-2453

Selenoprotein families include:

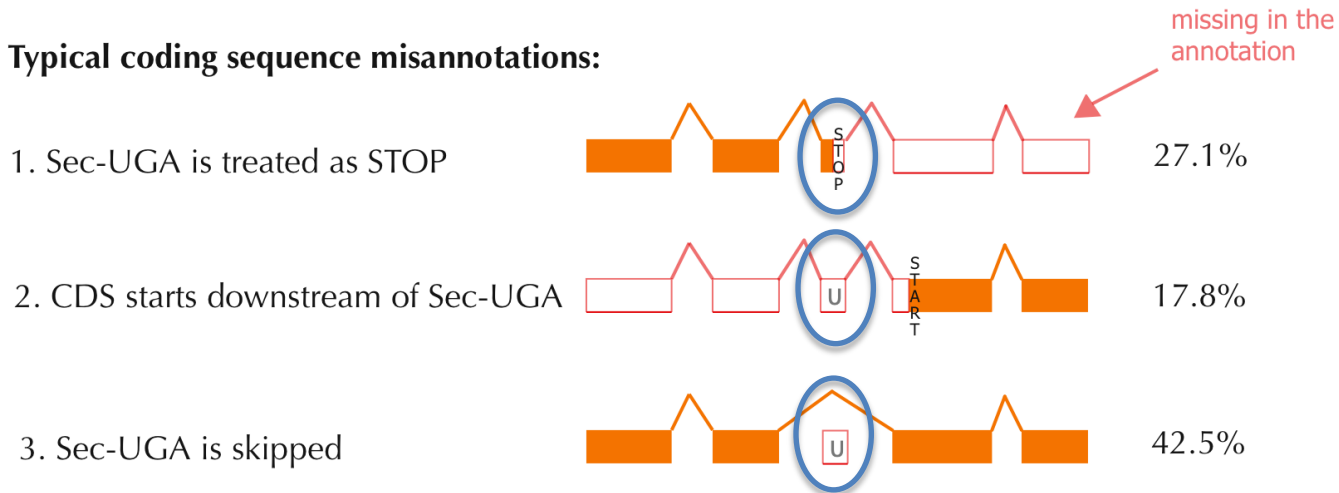
- Selenoproteins (Sec-containing proteins)
- Cysteine homologues (Cys-containing proteins)
 - Orthologues (*speciation* event)
 - Paralogues (*duplication* event)



Selenoproteins are generally misannotated



Typical coding sequence misannotations:



Bioinformatics methods for selenoproteins

- **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

Composition and Evolution of the Vertebrate and Mammalian Selenoproteomes

Marco Mariotti^{1,2,3}, Perry G. Ridge^{3,5}, Yan Zhang^{1,4,5}, Alexei V. Lobanov¹, Thomas H. Pringle⁵, Roderic Guigo², Dolph L. Hatfield⁶, Vadim N. Gladyshev^{1*}

1 Brigham and Women's Hospital and Harvard Medical School, Boston, Massachusetts, United States of America, **2** Center for Genomic Regulation and Universitat Pompeu Fabra, Barcelona, Spain, **3** Department of Biochemistry and Redox Biology Center, University of Nebraska, Lincoln, Nebraska, United States of America, **4** Key Laboratory of Systems Biology, Shanghai Institutes for Biological Sciences, Chinese Academy of Sciences, Shanghai, China, **5** Sperlberg Foundation, Eugene, Oregon, United States of America, **6** Laboratory of Cancer Prevention, National Cancer Institute, National Institutes of Health, Bethesda, Maryland, United States of America

Abstract

Background: Selenium is an essential trace element in mammals due to its presence in proteins in the form of selenocysteine (Sec). Human genome codes for 25 Sec-containing protein genes, and mouse and rat genomes for 24.

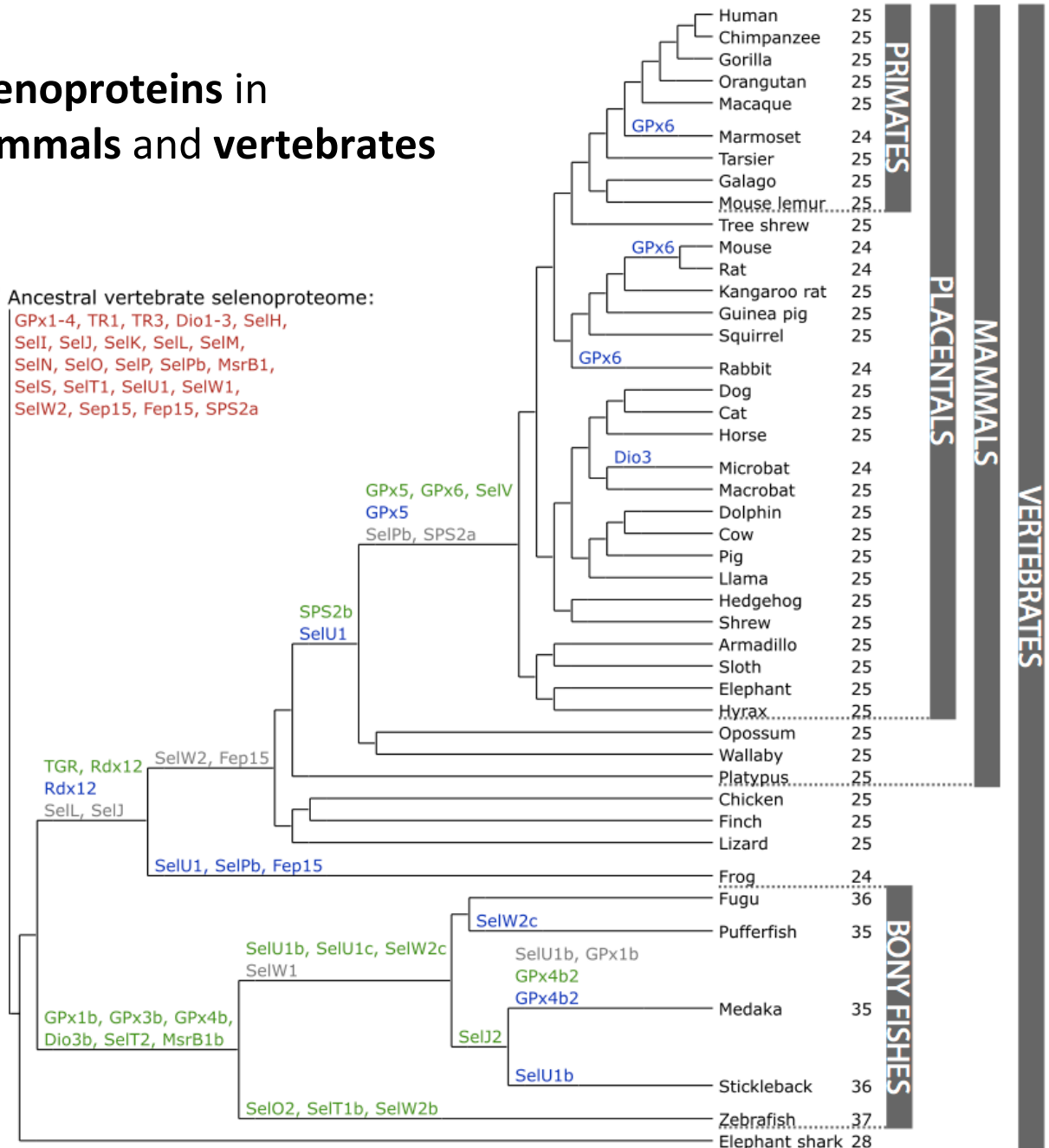
Methodology/Principal Findings: We characterized the selenoproteomes of 44 sequenced vertebrates by applying gene prediction and phylogenetic reconstruction methods, supplemented with the analyses of gene structures, alternative splicing isoforms, untranslated regions, SECIS elements, and pseudogenes. In total, we detected 45 selenoprotein subfamilies. 28 of them were found in mammals, and 41 in bony fishes. We define the ancestral vertebrate (28 proteins) and mammalian (25 proteins) selenoproteomes, and describe how they evolved along lineages through gene duplication (20 events), gene loss (10 events) and replacement of Sec with cysteine (12 events). We show that an intronless selenophosphate synthetase 2 gene evolved in early mammals and replaced functionally the original multiexon gene in placental mammals, whereas both genes remain in marsupials. Mammalian thioredoxin reductase 1 and thioredoxin-glutathione reductase evolved from an ancestral glutaredoxin-domain containing enzyme, still present in fish. Selenoprotein V and GPx6 evolved specifically in placental mammals from duplications of SelW and GPx3, respectively, and GPx6 lost Sec several times independently. Bony fishes were characterized by duplications of several selenoprotein families (GPx1, GPx3, GPx4, Dio3, MsrB1, SelJ, SelO, SelT, SelU1, and SelW2). Finally, we report identification of new isoforms for several selenoproteins and describe unusually conserved selenoprotein pseudogenes.

Conclusions/Significance: This analysis represents the first comprehensive survey of the vertebrate and mammal selenoproteomes, and depicts their evolution along lineages. It also provides a wealth of information on these selenoproteins and their forms.

Selenoproteins in mammals and vertebrates

28

Ancestral vertebrate selenoproteome:
 GPx1-4, TR1, TR3, Dio1-3, SelH,
 SelI, SelJ, SelK, SelL, SelM,
 SelN, SelO, SelP, SelPb, MsrB1,
 SelS, SelT1, SelU1, SelW1,
 SelW2, Sep15, Fep15, SPS2a



20 duplications

9 gene losses

13 Sec → Cys

Protocol overview

Tools:

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

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/>
<http://bioinformaticaupf.crg.eu>

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

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

Useful resources

Your assigned genomes will be available in the **UPF computers** when you will start the project

- Ensembl:** Collection of genomes (and annotations)

- NCBI nucleotide:** Collection of all sequences (genomes, ESTs, etc)

The screenshot shows the Ensembl website interface. At the top, there is a navigation bar with links for BLAST/BLAT, BioMart, Tools, Downloads, Help & Documentation, Blog, and Mirrors. Below this is a search bar with the text "Search all species...". The main content area is titled "Find a Species" and includes a "Species tree" section. A grid of species icons and names is displayed, including Alpaca, Anole lizard, Armadillo, Baboon, Budgerigar, Bushbaby, Clona intestinalis, Clona savignyi, Gibbon, Gorilla, Guinea Pig, Hedgehog, Horse, Human, Hyrax, Kangaroo rat, Platyfish, Platypus, Rabbit, Rat, Saccharomyces cerevisiae, Sheep, Shrew, and Sloth. Each species entry includes a small icon and the species name followed by its Ensembl ID.

The screenshot shows the NCBI Nucleotide Advanced Search Builder interface. The search query is "ailuropoda melanoleuca"[Organism]. The search results are displayed in a list, with the first result, "ailuropoda melanoleuca (182905)", highlighted. The list includes other related species such as ailuropoda melanoleuca david, 1869 (182905), ailuropoda melanoleuca qinlingensis (1), ailuropoda melanoleuca qinlingensis wan, wu and fang, 2005 (1), ailuropoda melanoleura (182905), ailuropus (7), ailuropus ursinus (7), ailurus (311), and ailurus fulgens (311). The interface also includes a search bar, a "Builder" section, and a "Hide index list" button.

Protocol steps

1st step: Get selenoprotein sequences

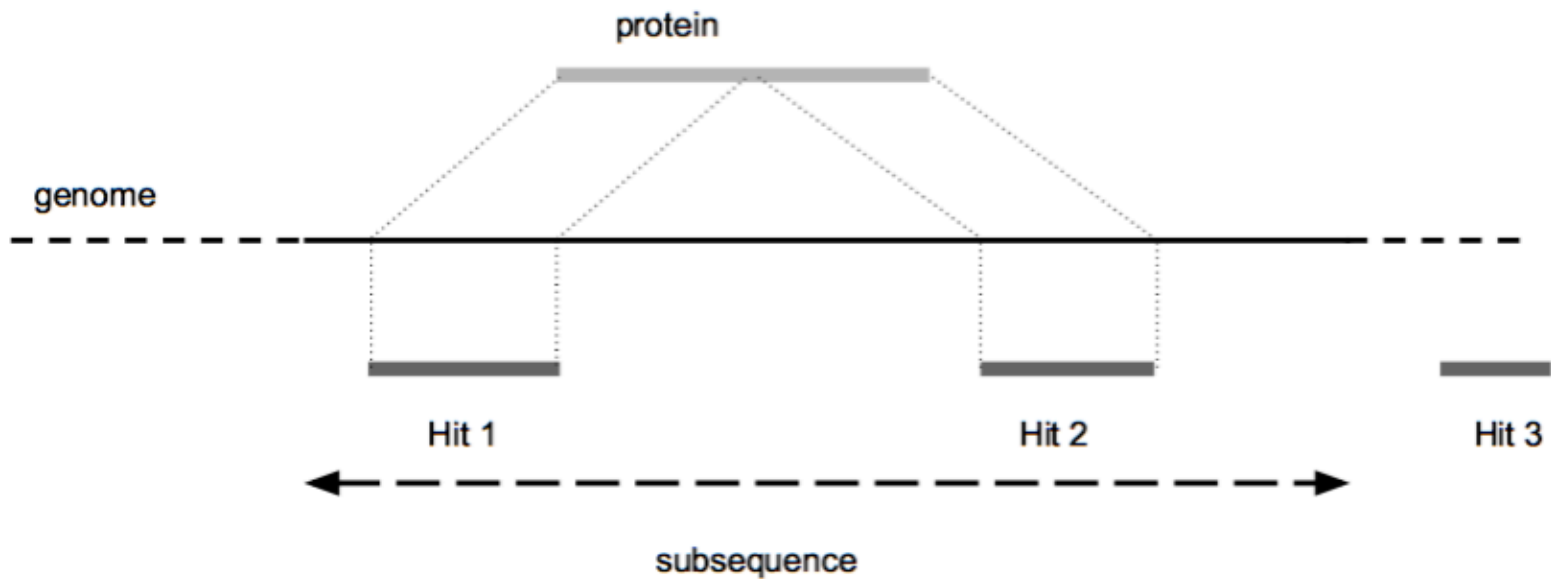
- **SelenoDB 2.0 (and 1.0)** **SelenoDB**
<http://www.selenodb.org> (2.0; automatic annotation)
<http://www1.selenodb.org> (1.0; manually curated, less species)
- **Protein databases**
<https://www.ncbi.nlm.nih.gov/protein/>
<http://www.uniprot.org>
- **Past year projects:**
<http://bioinformatica.upf.edu/>

NCBI BLAST programs... reminder

Search Name	Query Type	Database Type	Translation
blastn	Nucleotide	Nucleotide	None
tblastn	Peptide	Nucleotide	Database
blastx	Nucleotide	Peptide	Query
blastp	Peptide	Peptide	None
tblastx	Nucleotide	Nucleotide	Query and Database

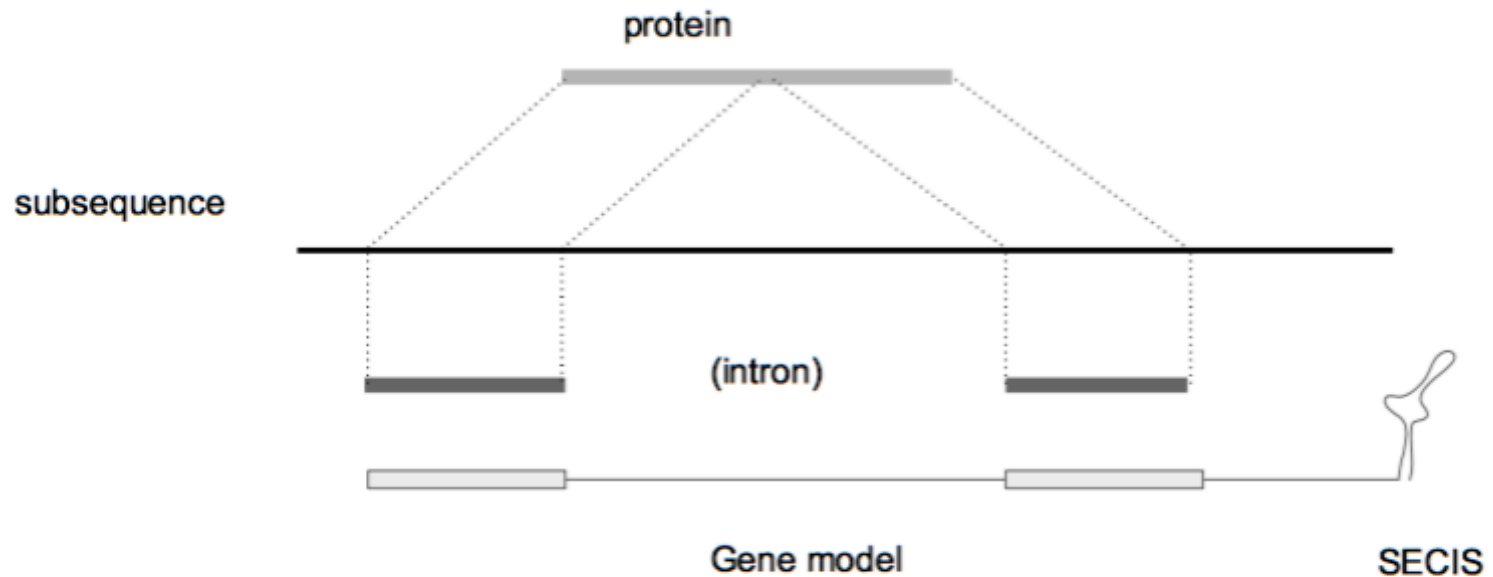
Protocol steps

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



Protocol steps

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



Protocol steps

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.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 "*"

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:

Submit

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.

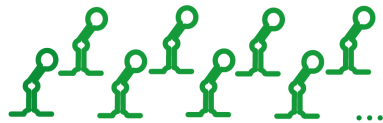
About | Contact us

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

Seblastian

TARGET SEQUENCE

SECISearch3



Assumptions:

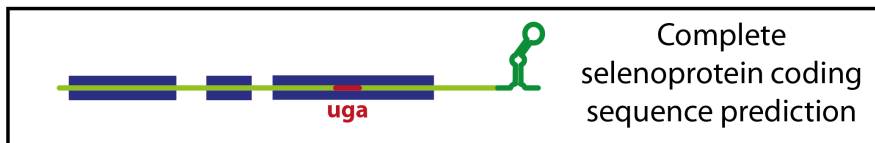
the presence of a detectable SECIS within acceptable genomic distance from the Sec-UGA

annotated homologue(s) (Sec/Cys) in the reference protein database

For each potential SECIS:

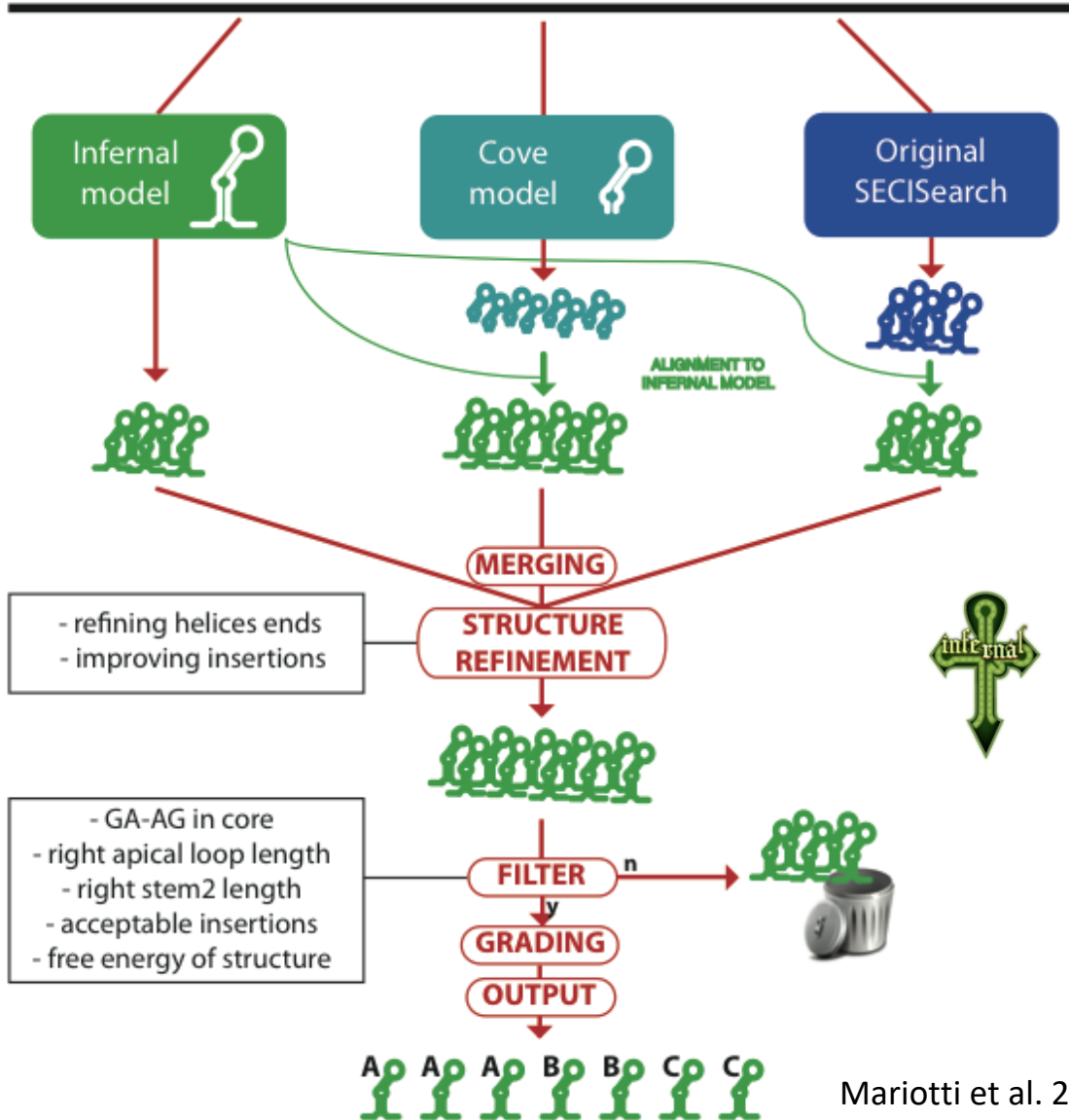


for these candidates



SECISearch 3

TARGET SEQUENCE



Based on a manually curated 2ndary structure alignment

Combines up to 3 methods to ensure maximum sensitivity

Filter and grading procedure based on manual inspection of hundreds of SECIS elements

Infernal: inference of RNA alignments

[infernal home](#) | [rfam database](#) | [eddy lab](#) | [janelia farm](#)

Selenoproteins as **test case**

- Selenoproteins have the peculiar characteristic of possessing a **UGA codon**, recoded because of the presence of the **SECIS element**.
- If you learn how to predict selenoproteins, you are able to do the same with any “*standard*” protein family.



BIOINFORMATICS PROJECT

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

UPF Human Biology.

Bioinformatics Courses 2007-2019

- 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 – **2019/20**: 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

1A: Ailuropoda melanoleuca

*AM. Barrios, A. Bellot,
S. Castany, M. De Manuel*

1B: Cricetulus griseus

*J. Fernandez, J. Gomez,
FD. Jurquiza, A. Lopez*

1C: Mustela putorius furo

*M. Perez, L. Taberner,
G. Vilajosana, I. Villate*

2A: Nomascus leucogenys

*M. Alemany, H. Costa,
A. Escrig, I. Gafarot*

2B: Saimiri boliviensis

*P. Garcia, J. Latorre,
R. Martinez, H. Palma*

2C: Sarcophilus harrisii

*G. Rodriguez, E. Ros,
AM. Saludes, H. Xicoy*

3A: Chrysemys picta bellii

*C. Bitlloch, G. Clua,
J. Domingo, P. Gelabert*

3B: Meleagris gallopavo

*J. Jancyte, L. Mateo,
A. Olle, M. Perera, C. Perez*

4A: Pelodiscus sinensis

*SU. Abad, A. Almeyda,
A. Azagra, R. Bartomeus*

4B: Gadus morhua

*O. Bover, N. Cortell,
B. Grau, E. March*

4C: Latimeria chalumnae

*A. Martinez, A. Perlas,
T. Robert, S. Walsh*

Projects 2019-2020

selenoproteins in vertebrates

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

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

https://ca.wikipedia.org/wiki/Viquiprojecte:Curs_Bioinform%C3%A0tica_UPF_2018

- **SelenoDB:** Insert your selenoprotein genes predictions into a real world database. Available to the scientific community.



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 these 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 a 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 **viqui**pedista

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'història](#)

Més ▾

Cerca a **Viqui**pèdia



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 corbals 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]

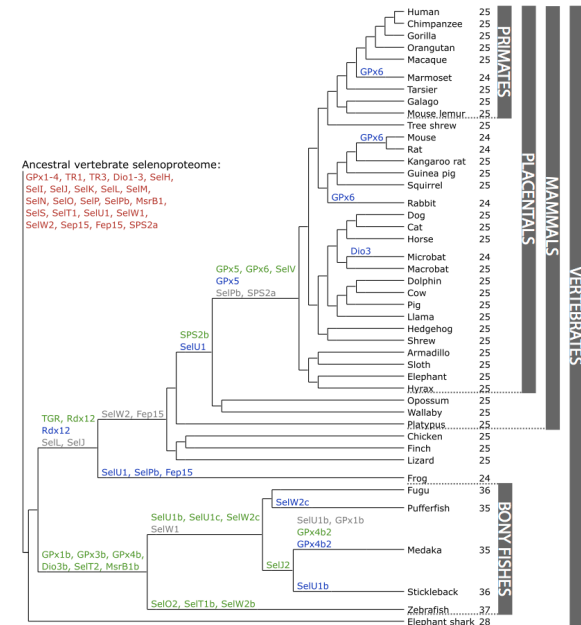
Author	<input type="text" value="email address"/> *
Specie	-- select a specie -- <input type="button" value="+ Add New Specie"/> *
Gene name	-- select a gene family -- <input type="radio"/> Forward <input type="radio"/> Reverse *
Promoter	<input type="text" value="Promoter start"/> - <input type="text" value="Promoter end"/>
Exons	<input type="text" value="Exon start"/> - <input type="text" value="Exon end"/> <input type="button" value="+ Add"/> *
Protein	<input type="text" value="Sec / U"/> <input type="text" value="Protein start"/> - <input type="text" value="Protein start"/> *
Secis	<input type="text" value="Secis start"/> - <input type="text" value="Secis end"/> <input type="button" value="+ Add"/>
Residue	<input type="text" value="Sec / U"/> <input type="text" value="Residue start"/> - <input type="text" value="Residue end"/> <input type="button" value="+ Add"/>
Sequence	<input type="text" value="CHR"/> - <input type="text" value="Sequence start on the c"/> - <input type="text" value="Sequence end on the cl"/> *
	<input type="text" value="ENS00075"/> - <input type="text" value="Ensembl"/> - <input type="text" value="http://www.ensembl.org"/> - <input type="text" value="http://www.ensembl.org"/> *
	<input type="button" value="Choose File"/> No file chosen <input type="text" value="gff format"/> * <input type="button" value="Fetch Sequence"/>
	<div style="border: 1px solid gray; height: 60px; width: 100%;"></div> <div style="border: 1px solid gray; height: 60px; width: 100%;"></div>
<input type="button" value="Submit"/>	

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
- 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.
 - ✓ **Cys-containing homologs**
- **Other helpful resources** to biologically interpret and visualize the results (**phylogenetic trees**):
 - phyloT: <https://phylot.biobyte.de/> (from NCBI taxonomy → .nw)
 - iTOL: <https://itol.embl.de/> (.nw)
 - Etetoolkit: <http://etetoolkit.org/treeview/> (.nw or .msa)
 - Phylogeny.fr: http://www.phylogeny.fr/simple_phylogeny.cgi (.mfa)

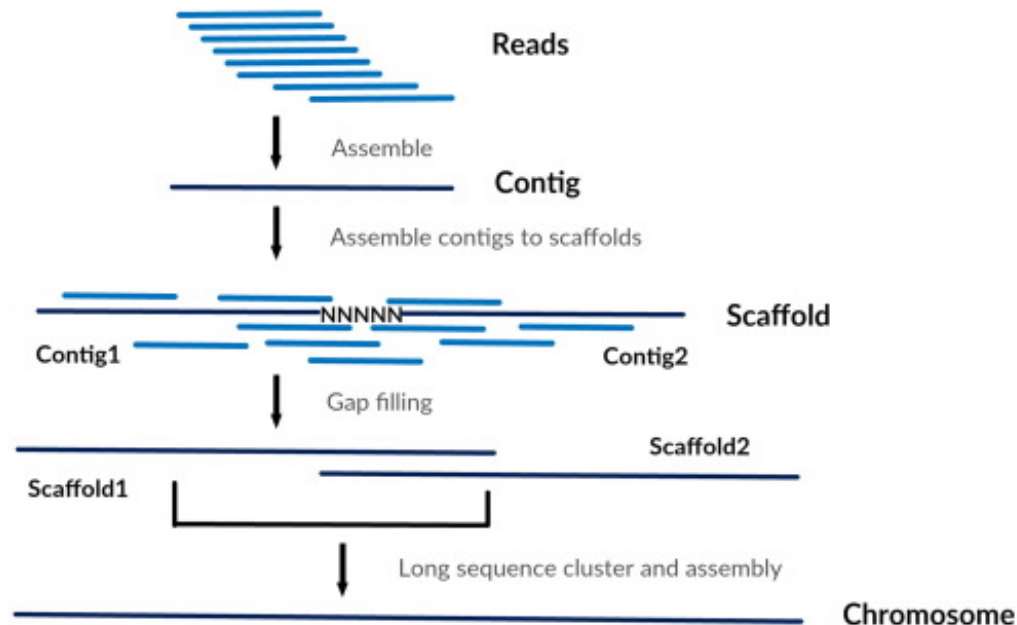
Common pitfalls

- Know what to **expect**
- **Zero, one or many genes?**
(!) careful with superfamilies and gene duplications
- **Genomic context**



Common pitfalls

- **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



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 Wikipedia entry)