

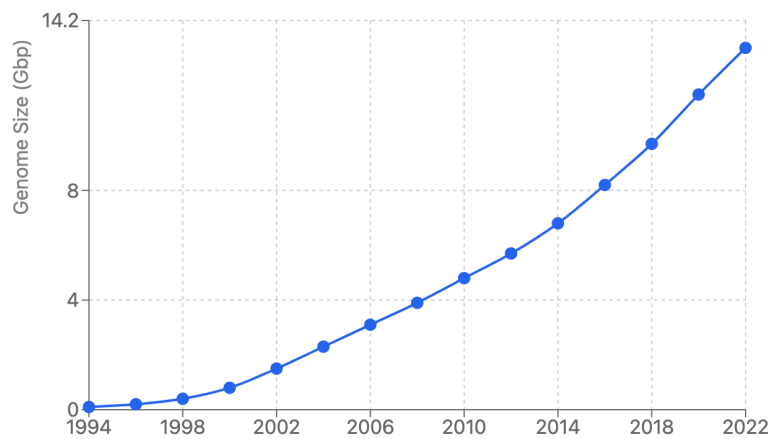
4. BASI DI BIOINFORMATICA PER LO STUDIO DELLA REGOLAZIONE GENICA

1

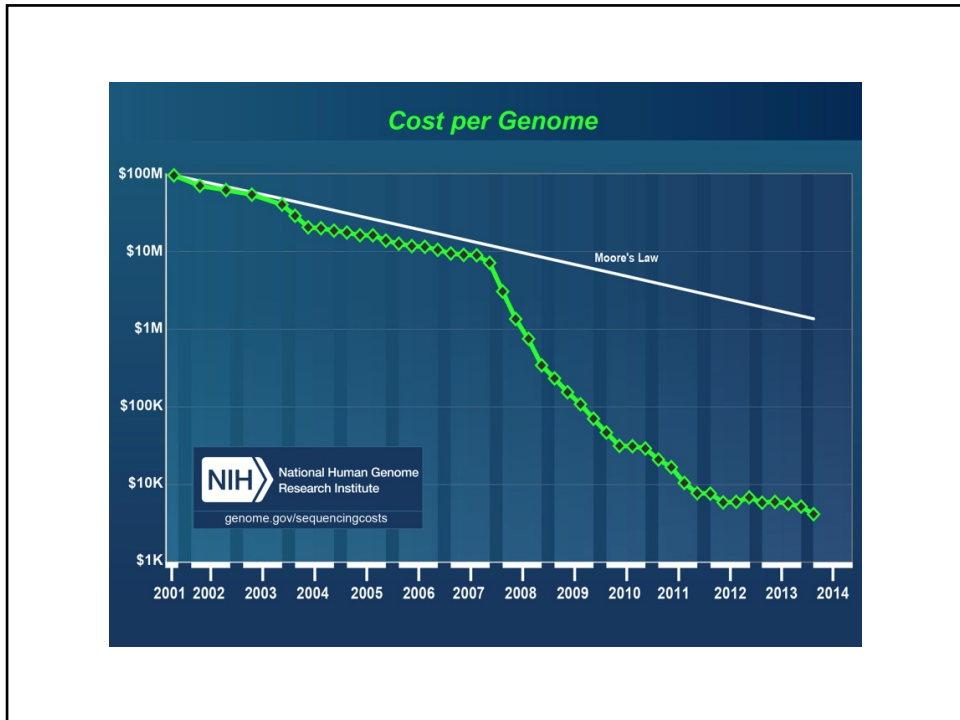
La sfida dei genome browsers

aumentare sempre di più le informazioni disponibili sulle sequenze

Genome Size Evolution (1994-2022)



2



3

Genome Browsers Oggi

- Ensembl Genome browser
<http://www.ensembl.org>
- NCBI Map Viewer
<http://www.ncbi.nlm.nih.gov/mapview/>
- UCSC Genome Browser
<http://genome.ucsc.edu>

4

Differenze tra Ensembl, UCSC ed NCBI?

	UCSC	Ensembl	NCBI
Presentation	Genome in horizontal orientation Main page contains a single graphic displaying annotation ('tracks') Clicking on annotation element presents web page of detailed information and links to other resources	Genome in horizontal orientation Main ContigView page contains three graphics displaying annotations at different resolutions Clicking on annotation element presents box with links to other resources or Views with more detailed information	Genome in vertical orientation Annotations graphically presented in columns ('maps') Clicking on annotation elements or links in columns provides quick access to other, primarily NCBI, resources
Content	13 vertebrate, 15 invertebrate Many cross-species annotations including conservation across eight species ENCODE Project annotations	13 vertebrate, six invertebrate Heavy focus on gene annotations such as Ensembl genes and VEGA HapMap project-related Views	11 vertebrate, five invertebrate, one protozoan, 12 plant, eight fungi Annotations primarily from NCBI resources
Functionality	Text search, BLAT sequence search, isPCR primer search Advanced annotation extraction using Table Browser Ability to upload and view own	Text search, BLAST and SSAHA sequence search, e-PCR primer search Advanced annotation extraction using BioMart Ability to upload and view own annotations	Text search, BLAST sequence search, e-PCR primer search Basic annotation extraction

5

Genome Reference Consortium



Wellcome Sanger Institute



The McDonnell Genome Institute at Washington University



The European Bioinformatics Institute



The National Center for Biotechnology Information



The Zebrafish Model Organism Database



Rat Genome Database

6

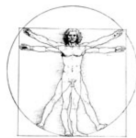
Human genome assemblies

Guard

- GRCh38 (aka hg38)
 - No gaps
 - www.ensembl.org
 - Most up-to-date and supported
- GRCh37 (aka hg19)
 - 250 gaps
 - grch37.ensembl.org
 - Limited data and software updates
- NCBI36 (aka hg18)
 - 150,000 gaps
 - ncbi36.ensembl.org
 - No longer updated



7



Human

The human genome assembly was produced as part of the [Human Genome Project \(HGP\)](#). The previous assembly (NCBI36) was the last one produced by the HGP and was described in 2004 (PMID: 15496913); this was the starting point for the GRC. The assembly is based largely on assembling overlapping clone sequences.

Human assembly information

Current major assembly	GRCh38
Regions with alternate loci	178
Assembly N50	67,794,873 bp
Remaining gaps	875
Patch release version	p14
Patches released	FIX: 164 , NOVEL: 90



Mouse

The GRC has produced an updated assembly (GRCm38). This is an update of the last MGSC assembly (MGSCv37) which was described in 2009 (PMID: 19468303). The primary assembly is based on assembling overlapping BAC clones derived from the C57BL/6J strain and several loci have sequence available from other strains.

Mouse assembly information

Current major assembly	GRCm39
Regions with alternate loci	0
Assembly N50	106,145,001 bp
Remaining gaps	347
Patch release version	None



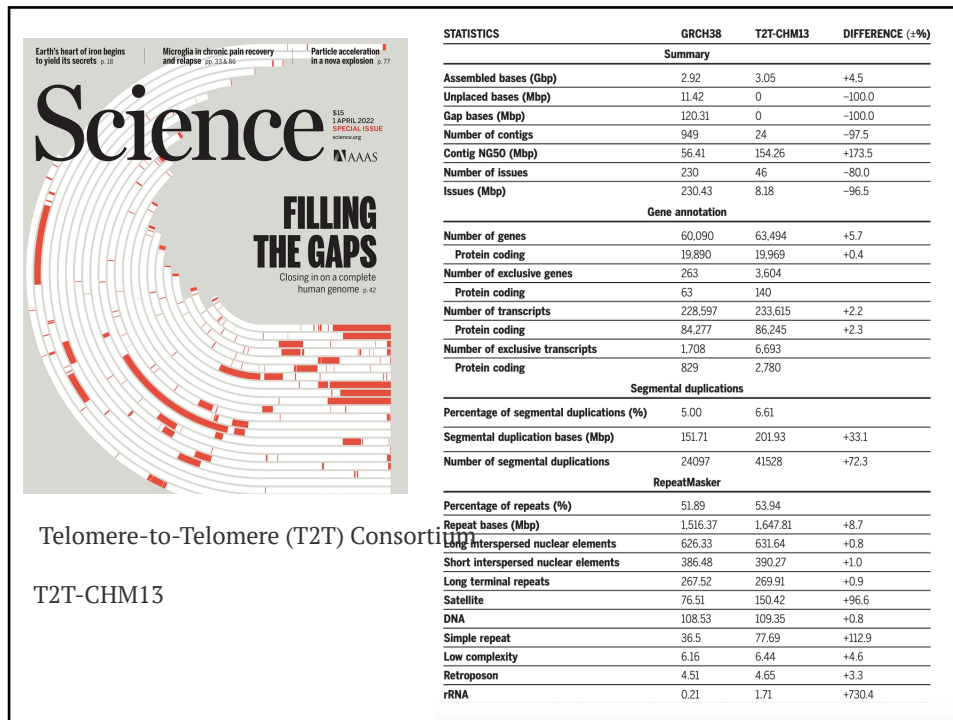
Zebrafish

The zebrafish genome assembly was produced at the [Wellcome Sanger Institute](#). The last assembly produced from the original project was Zv9 and was described in 2013 (PMID: 23594743). This assembly is the starting point for the GRC. The assembly is based on assembling overlapping BAC clones and integrating these sequences with the whole genome shotgun assembly.

Zebrafish assembly information

Current major assembly	GRCz11
Regions with alternate loci	607
Assembly N50	7,379,053 bp
Remaining gaps	18,736

8



9

Cosa impariamo da queste annotazioni?

- **All'interno di un genoma: elementi regolatori, ordine dei geni, struttura della cromatina.....**
- **Facendo studi comparativi: evoluzione, regioni conservate, riarrangiamenti...predizione di geni.**

10

Studiare i genomi con **Ensembl**

The screenshot shows the Ensembl website header with navigation links: BLAST/BLAT, BioMart, VEP, Tools, Downloads, Help & Docs, and Blog. Below the header are three main tool categories: Tools (with a link to All tools), BioMart (for exporting custom datasets), and BLAST/BLAT (for searching genomes). The Variant Effect Predictor (VEP) is also listed for analyzing variants. A search bar is present with a dropdown for 'All species' and a 'Go' button. Below the search bar, there are sections for 'All genomes' (with a species selector) and 'Favourite genomes' (listing Human, Mouse, and Zebrafish).

11

Nomenclatura in Ensembl

- **ENSG###** Ensembl **Gene** ID
- **ENST###** Ensembl **Transcript** ID
- **ENSP###** Ensembl **Peptide** ID
- **ENSE###** Ensembl **Exon** ID

- **Per specie diverse dall'uomo è aggiunto un suffisso**

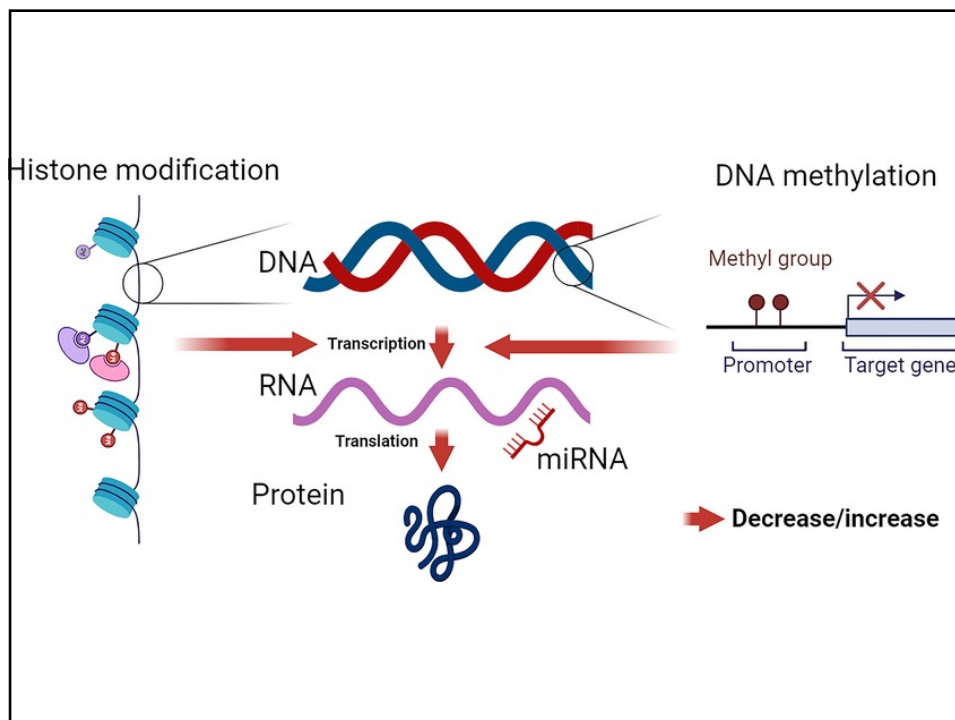
MUS (*Mus musculus*) for mouse: **ENSMUSG###**
DAR (*Danio rerio*) for zebrafish: **ENSDARG###**, etc.

12

Quali annotazioni sono disponibili?

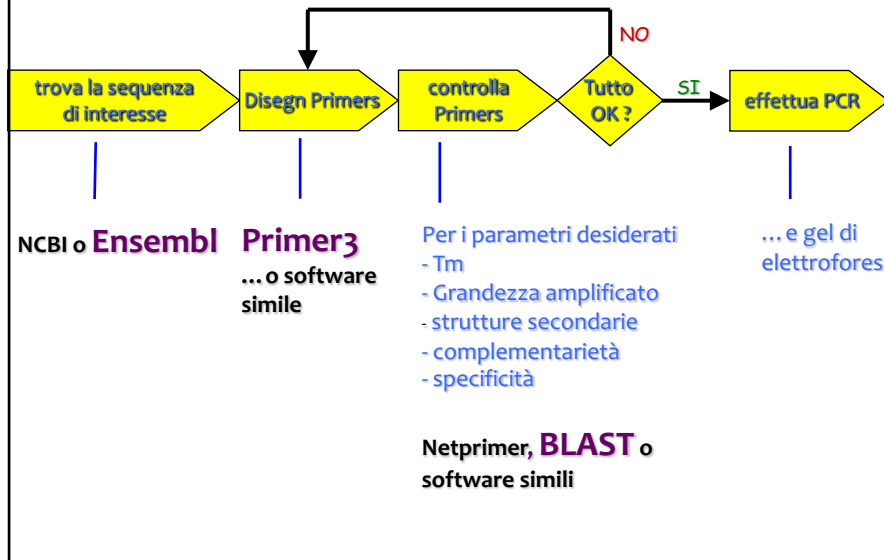
- Gene/transcript/peptide models (coding and noncoding (ncRNAs))
- IDs in other database
- Mapped cDNAs, peptides, micro array probes, BAC clones etc.
- Cytogenetic bands, markers, repeats etc.
- Comparative data:
- orthologues and paralogues, protein families, whole genome alignments, syntenic regions
- Variation data:
- Single Nucleotide Polymorphisms (SNPs)
- Regulatory data:
- “best guess” set of regulatory elements from ENCODE
- Data from external sources (DAS)

13



14

Primer design-come procedere



15

disegno dei primers

Un buon disegno dei primers è la chiave di una PCR di successo !!

16

Proprietà importanti di una buona coppia di primers

Ogni primers dovrà avere

- lunghezza basi compresa tra 18-24
- 40-60% G/C
- Distribuzione bilanciata di basi G/C e A/T
- T_m che permette un annealing tra 55-65° C
- NO strutture secondarie interne (hair-pins)

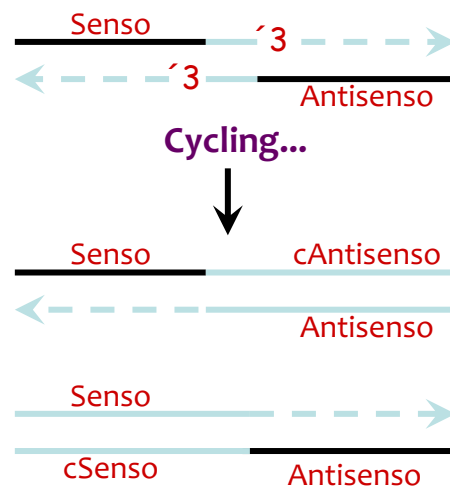
Le coppie di primers inoltre dovrebbero avere:

- T_m simile (max 2-3° C di differenza)
- **NO** complementarità (> 2-3 bp) in particolare al 3'

17

Il problema dei dimeri di primers

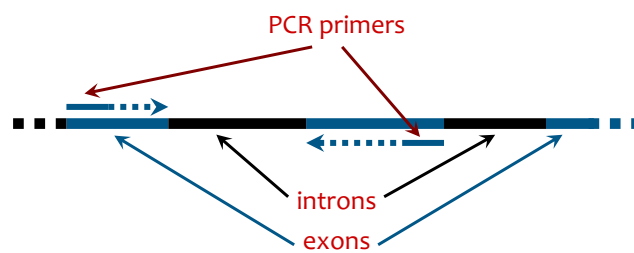
- Primers che interagiscono tra di loro sono **AMPLIFICATI** dalla PCR
- La formazione dei dimeri di primers compete con la PCR e può quindi compromettere l'efficienza della reazione.



18

Considerazioni importanti

- Evitare di avere come bersaglio della qPCR delle **strutture secondarie**
- Evitare contaminazione genomica disegnando primers che **includono esoni diversi e tagliano introni**



19

Links per disegno primers

- <http://www.tataa.com/>
- <http://www.ncbi.nlm.nih.gov/BLAST/>
- www.premierbiosoft.com/netprimer/netplaunch/netplaunch.html
- www.ensembl.org
- http://www-genome.wi.mit.edu/cgi-bin/primer/primer3_www.cgi
- <http://www.bioinfo.rpi.edu/applications/mfold/dna/form1.cgi>
- <http://primer3.ut.ee/> **Primer3**

20

BLAST

Basic Local Alignment Search Tool

21

BLAST Programs

Program	Database (Subject)	Query
BLASTN	Nucleotide	Nucleotide
BLASTP	Protein	Protein
BLASTX	Protein	Nt. → Protein
TBLASTN	Nt. → Protein	Protein
TBLASTX	Nt. → Protein	Nt. → Protein

22

BLAST confronta le sequenze

- **usa** una sequenza *“query”*
- **La confronta con** milioni di sequenze nei database *GenBank*® costruendo local alignments
- **elenca** quelle che sembrano simili alla query
- **dice perchè sono eventualmente omologhe**

NB

- *BLAST suggerisce*
- *Il ricercatore trae le sue conclusioni*

23

23

BLAST output

**descrive in che modo le sequenze
allineate sono simili**

- *Quanto sono lunghi i segmenti allineati?*
- *BLAST ha dovuto introdurre degli spazi per allineare i segmenti?*
- *Quanto sono simili i segmenti allineati?*

24

24

U.S. National Library of Medicine
National Center for Biotechnology Information

COVID-19 Information
Public health information (CDC) | Research information (NIH) | SARS-CoV-2 data (NCBI) | Prevention and treatment information (HHS) | Español

BLAST® Home Recent Results Saved Strategies Help

Basic Local Alignment Search Tool
BLAST finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance. [Learn more](#)

NEWS
A new feature was added to the NCBI IgBLAST webpage
IgBLAST is now able to determine Ig isotypes
Mon, 01 Nov 2021 12:00:00 EST [More BLAST news...](#)

Web BLAST

Nucleotide BLAST
nucleotide ▶ nucleotide

blastx
translated nucleotide ▶ protein

tblastn
protein ▶ translated nucleotide

Protein BLAST
protein ▶ protein

<https://blast.ncbi.nlm.nih.gov/Blast.cgi>

25

Graphic Display

- How good is the match?
 - **Red = excellent!**
 - **Pink = pretty good**
 - **Green = OK, but look at other factors**
 - **Blue = bad**
 - **Black = really bad!**
- How long are the matched segments?
Longer = better

26

26

BLAST fa una lista dei migliori accoppiamenti (hits)

Per ogni accoppiamento fornisce:

- **Max Score:** the highest alignment score calculated from the sum of matched nucleotides and penalties for mismatches and gaps.
- **Tot Score:** the sum of alignment scores of all segments
- **Query Cover:** the % of the query length included in the aligned segments.
- **E[xpect] Value:** the number of alignments expected by chance with the calculated score or better. for significant alignments the E value should be very close to zero.
- **Ident[ity]:** the highest % identity for a set of aligned segments to the same subject sequence.

Sequences producing significant alignments:

Select: All None Selected 0

Alignments Download GenBank Graphics Distance tree of results

Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/> PREDICTED: Homo sapiens prodynorphin (PDYN), transcript variant X8, mRNA	1585	1585	100%	0.0	100%	XM_011529250.2
<input type="checkbox"/> PREDICTED: Homo sapiens prodynorphin (PDYN), transcript variant X7, mRNA	1585	1585	100%	0.0	100%	XM_011529249.2
<input type="checkbox"/> PREDICTED: Homo sapiens prodynorphin (PDYN), transcript variant X6, mRNA	1585	1585	100%	0.0	100%	XM_017027879.1
<input type="checkbox"/> PREDICTED: Homo sapiens prodynorphin (PDYN), transcript variant X4, mRNA	1585	1585	100%	0.0	100%	XM_011529246.2
<input type="checkbox"/> PREDICTED: Homo sapiens prodynorphin (PDYN), transcript variant X5, mRNA	1585	1585	100%	0.0	100%	XM_011529248.1
<input type="checkbox"/> PREDICTED: Homo sapiens prodynorphin (PDYN), transcript variant X3, mRNA	1585	1585	100%	0.0	100%	XM_011529247.1
<input type="checkbox"/> PREDICTED: Homo sapiens prodynorphin (PDYN), transcript variant X2, mRNA	1585	1585	100%	0.0	100%	XM_011529245.1
<input type="checkbox"/> PREDICTED: Homo sapiens prodynorphin (PDYN), transcript variant X1, mRNA	1585	1585	100%	0.0	100%	XM_011529244.1
<input type="checkbox"/> Homo sapiens preproenkephalin (enkeB) gene, partial cds	1585	1585	100%	0.0	100%	AF4002816.3
<input type="checkbox"/> Homo sapiens prodynorphin (PDYN), transcript variant 1, mRNA	1585	1585	100%	0.0	100%	NM_024411.4
<input type="checkbox"/> Homo sapiens prodynorphin (PDYN), transcript variant 4, mRNA	1585	1585	100%	0.0	100%	NM_001190899.2
<input type="checkbox"/> Homo sapiens prodynorphin (PDYN), transcript variant 2, mRNA	1585	1585	100%	0.0	100%	NM_001190898.2

27

Che cos'è l'E-value?

- **E-value**

La probabilità che quel determinato accoppiamento non è casuale

+ basso è l'E-value, + significativo è l'accoppiamento

- **E = 10^{-4}** è considerato il **cutoff point**
- **E = 0** significa che le due sequenze sono statisticamente **identiche**

28

Software to predict CpG islands?

MethPrimer is an online platform which provides a number of tools and databases to facilitate the study of DNA methylation and epigenetics, including tools for designing prim

<http://www.urogene.org/cgi-bin/methprimer/methprimer.cgi>

29

Software to predict miRNA binding ??



<http://mirdb.org/index.html>

30

growing evidence suggest the importance of both environmental and genetic factors in the influence of DNA methylation.

DNA methylation can be influenced by cis-acting DNA sequence variation located on the same chromosome.

- Mill et al., *Am J Hum Genet* 2008
- Zhang et al., *Am J Hum Genet* 2010
- Milani et al., *Genome Res* 2009
- Docherty et al., *Behav and Brain Func* 2012
- Ball et al., *Genome Biol.* 2011

new models have to be developed to integrate genetic variants and DNA methylation.

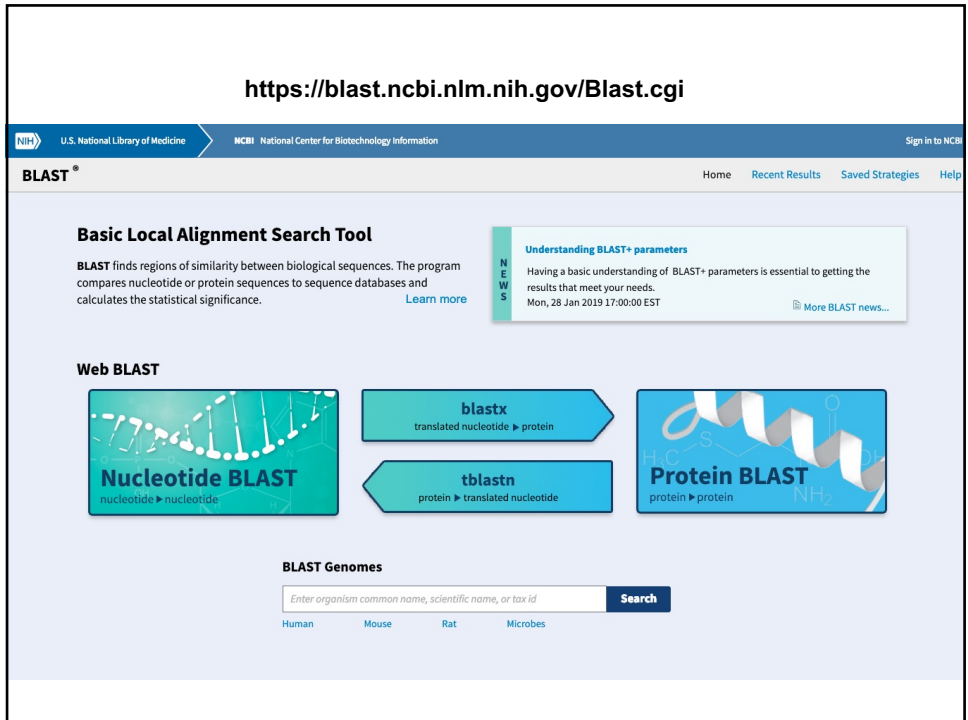
31

miRdSNP

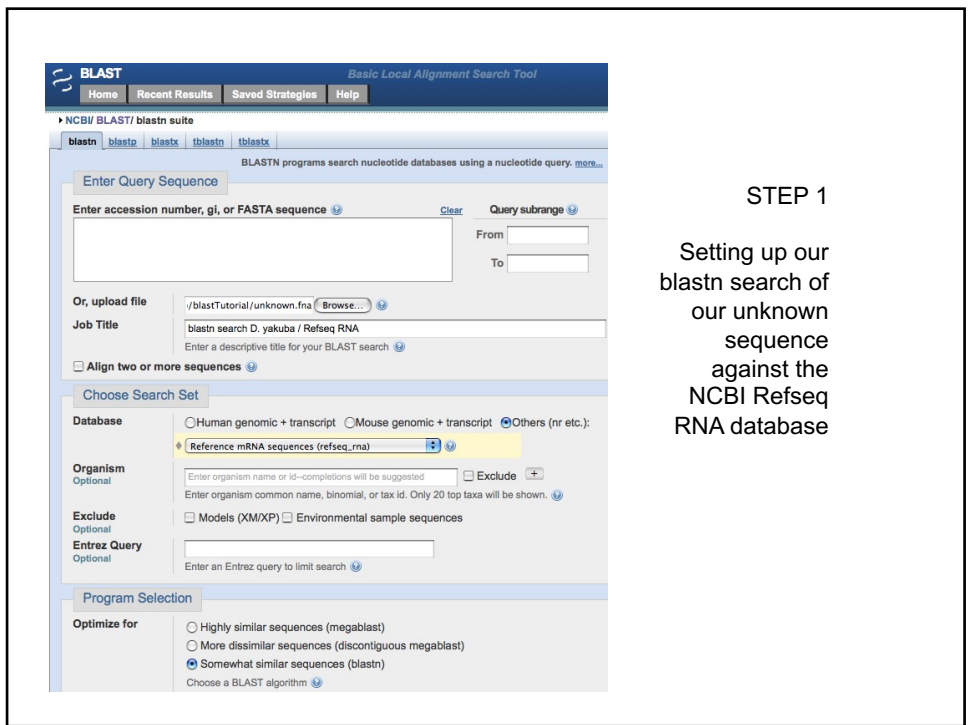
a database of disease-associated SNPs and microRNA target sites on 3'UTRs of human genes

<http://mirdsnp.ccr.buffalo.edu/>

32



33



34

► [NCBI/BLAST/Formatting Results - 56RFPSX1012](#) [\[Formatting options\]](#)

Job Title: blastn search D.yakuba / Refseq RNA search

WAITING

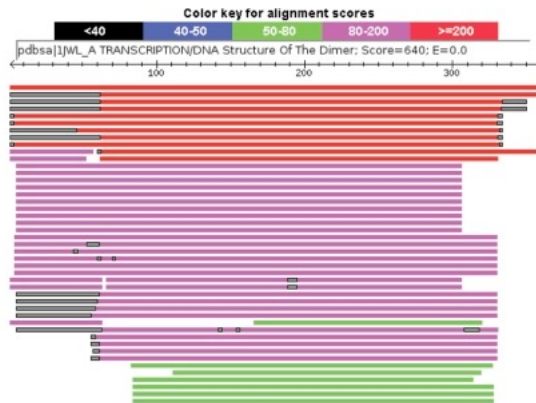
Request ID **56RFPSX1012**
 Status Searching
 Submitted at Tue May 22 17:17:42 2007
 Current time Tue May 22 17:17:45 2007
 Time since submission 00:00:03

This page will be automatically updated in 13 seconds until search is done

When the NCBI web server is busy, the search may take 5 minutes or more

35

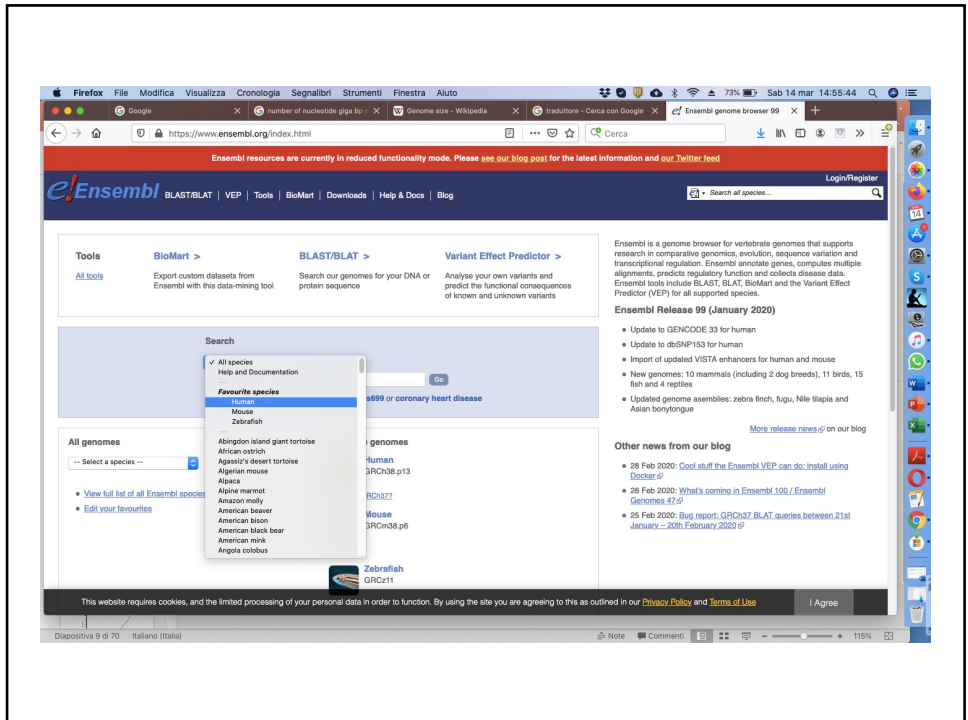
BLAST report



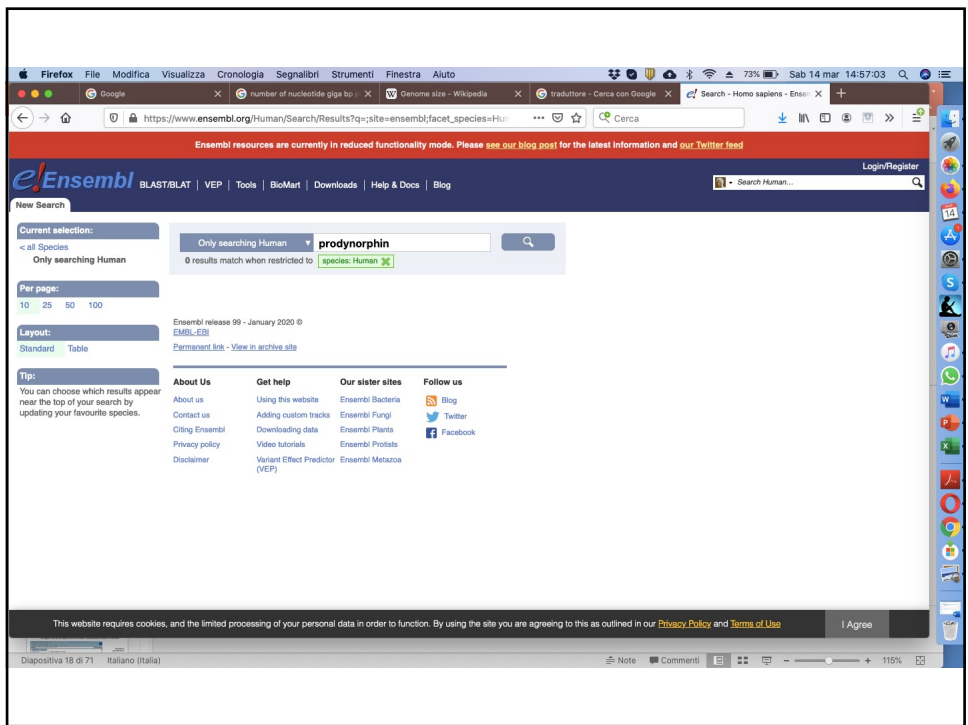
Sequences producing significant alignments:

	Score (bits)	E value	Source DB	NCBI Entrez	Cath Prot/chain
pdbsa 11EG_A	688	0	SDB	NCBI	CATH Prot
pdbsa 11EH_A	688	0	SDB	NCBI	CATH Prot
pdbsa 11EF_A	669	0	SDB	NCBI	CATH Prot
pdbsa 11VE_A	666	0	SDB	NCBI	CATH Prot
pdbsa 11VL_A	640	0	SDB	NCBI	CATH Prot
pdbsa 11FA_A	640	0	SDB	NCBI	CATH Prot
pdbsa 11FA_C	640	0	SDB	NCBI	CATH Prot
pdbsa 11VL_C	640	0	SDB	NCBI	CATH Prot
pdbsa 11FA_B	640	0	SDB	NCBI	CATH Prot
pdbsa 11LF_A	575	1e-164	SDB	NCBI	CATH Prot

36



37



38

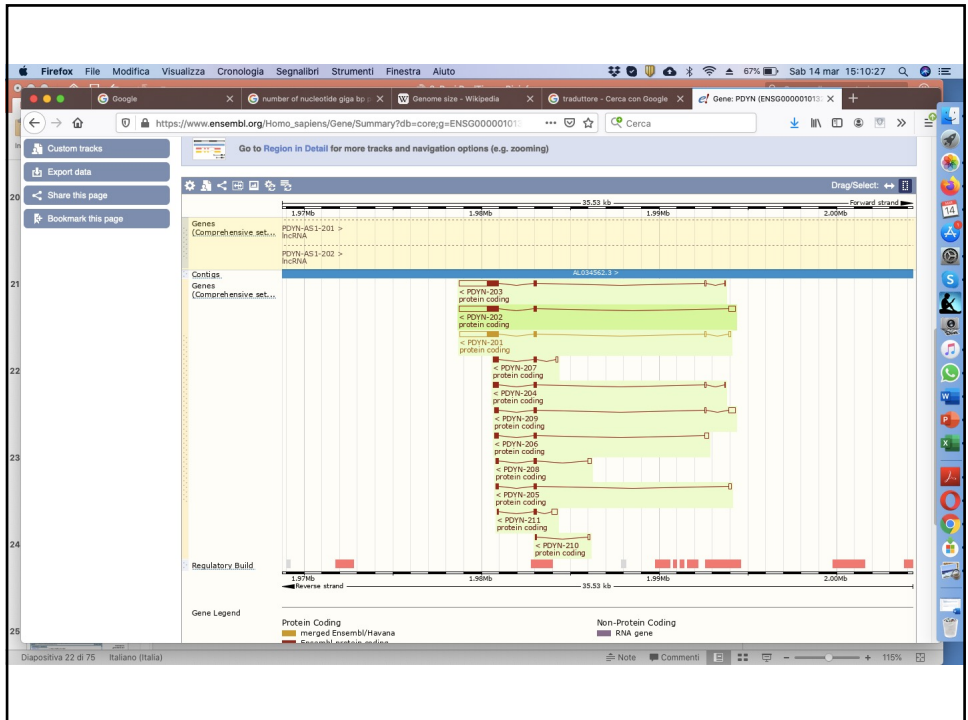
The screenshot shows the Ensembl search results for 'prodynorphin'. The search bar at the top contains 'prodynorphin' and indicates '545 results match prodynorphin'. On the left, there are filters for 'Restrict category to:' (Gene: 244, Transcript: 301) and 'Restrict species to:' (Human: 12, Mouse: 60, Zebrafish: 3, etc.). A 'Tip:' section at the bottom left explains that users can choose which results appear near the top of their search by updating their favourite species. The main content area lists several gene records, including PDYN (Human Gene) and various transcript isoforms like PDYN-210, PDYN-209, PDYN-208, PDYN-211, PDYN-204, and PDYN-205. Each record includes the Ensembl ID, coordinates, and a brief description.

39

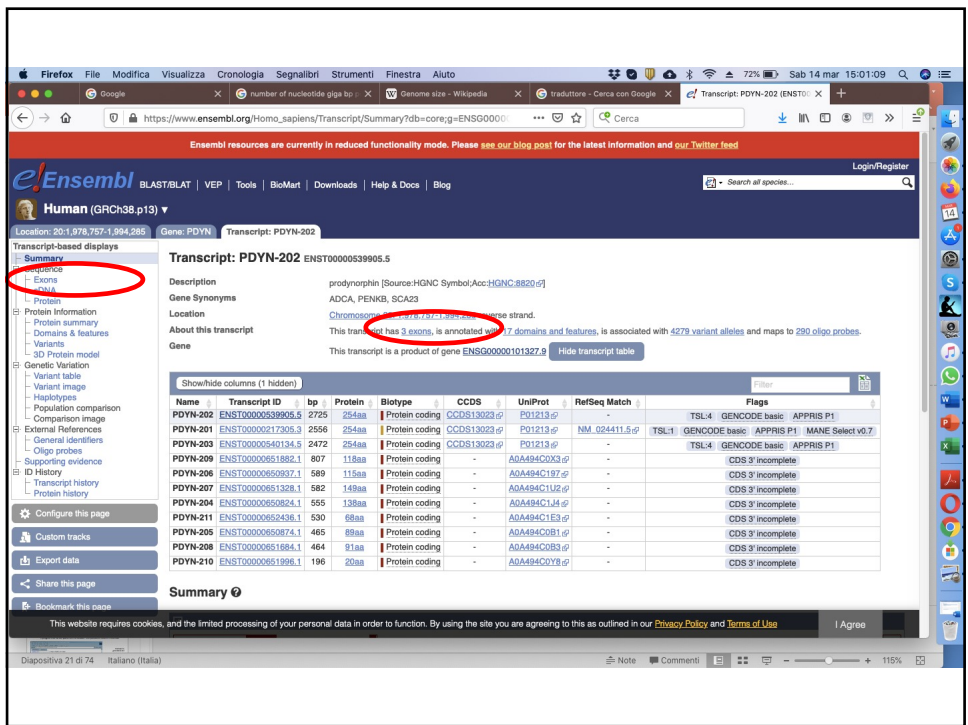
The screenshot shows the Ensembl gene summary page for 'PDYN' (ENSG00000101327). The page is titled 'Human (GRCh38.p13)' and shows the gene's location on Chromosome 20. The 'About this gene' section states: 'This gene has 11 transcripts (splice variants), 22 orthologues, 2 paralogues, is a member of 1 Ensembl protein family and is associated with 2 phenotypes.' The 'Transcripts' section is expanded, showing a table of transcript isoforms. A red circle highlights the text '11 transcripts (splice variants)'. Below the table is a 'Summary' section.

Name	Transcript ID	bp	Protein	Biotype	CCDS	UniProt	RefSeq Match	Flags
PDYN-202	ENST00000309055	2725	254aa	Protein coding	CCDS13023.0	P01213	-	TSL4 GENCODE basic APPRIS P1
PDYN-201	ENST00000217305.3	2556	254aa	Protein coding	CCDS13023.0	P01213	NM_024411.5	TSL1 GENCODE basic APPRIS P1 IMANE Select v0.7
PDYN-203	ENST00000540134.5	2472	254aa	Protein coding	CCDS13023.0	P01213	-	TSL4 GENCODE basic APPRIS P1
PDYN-209	ENST00000651882.1	807	118aa	Protein coding	-	A0A494C0X3	-	CDS 3' incomplete
PDYN-206	ENST00000650937.1	589	115aa	Protein coding	-	A0A494C197	-	CDS 3' incomplete
PDYN-207	ENST00000651328.1	582	148aa	Protein coding	-	A0A494C1U2	-	CDS 3' incomplete
PDYN-204	ENST00000650824.1	555	138aa	Protein coding	-	A0A494C1J4	-	CDS 3' incomplete
PDYN-211	ENST00000652436.1	530	68aa	Protein coding	-	A0A494C1E3	-	CDS 3' incomplete
PDYN-205	ENST00000650874.1	465	89aa	Protein coding	-	A0A494C0B1	-	CDS 3' incomplete
PDYN-208	ENST00000651884.1	464	81aa	Protein coding	-	A0A494C0B3	-	CDS 3' incomplete
PDYN-210	ENST00000651996.1	196	20aa	Protein coding	-	A0A494C0Y8	-	CDS 3' incomplete

40



41



42

No.	Exon / Intron	Start	End	Start Phase	End Phase	Length	Sequence
1	ENSE00000220217 5' upstream sequence	1,994,285	1,993,911	-	-	375catcaaaaactgaagtggccggtttccaagtgcaaaacagcactacacc AGAAATCCCTCCCTCCACAGAGGCTGCCCTTCTCCACATCTCTGCAGAGAGA AGCCTATTGTCAGGCCCGAGGGCTCCAGTFCGAGGGCTGGGGCTCTCTCTGC TGTCTCAGCCACTCCCAATGGCTCCACAGCCTCTGCTCAGAGAGGCTGAGCAG AGGGAGGCTCTCTCATAAAGGGGGGAGAGGCCCAAGCTGCCATTTTAAAGGGC TTTGTGTGTTCAGAGCTGCTCTTGGACCTCTCCACAGCGAGTCAAGAGGCC CTAGAGCTTGGACAGCCTGACCTCCGACTCCGACTCCGCTCGCGAGAGCTCCAGGGA GAGAGGCTGAG
2	ENSE00000655739 Intron 1-2	1,993,910	1,982,956	-	0	1,087	gtaactccagagppggaactgpa.....tattctctcttctctcccccag CAGGAATTCCTGAGCAGGATGCGCTGCGAGGGCTGCTGGCTGCTGCTCCAT GTTCCTCCACAGAGGCTGCTGCTGCGGGTCTCTGTGTGCTGTAAAGACCA GGATGGTCCAAACTATCAATCCCTG
3	ENSE00000253018 Intron 2-3	1,982,955	1,980,959	-	-	1,997	gtatgtttctccagaggttcctca.....tgggtcttttgggttttccag ATTTGCCTCCGCAAGGCGAGCTGCTGCTGCTGCTGCTGAGAGATGGAGAGTCCAG AGCTTTCTGTCTTTTTCACTCCCTCACTGGCTCAATGACAGAGAGACTTGGG AGCAAGTGGTGGAGAGGGCTCAGAGAGCTGCTGCGAGCTCTGGGCTCTGCTG AGAGGCTGGAGAAAGAGTTCCTCCCAATTTCCAGAAAGAGGAGACTTGAAC AGAGCTGGAGAGAGCTCAGGGCTCTCTCAGCGGTTTGGAGGGAGAGAGTCT GAGCTGATGGAGTCCAGCTGACAGATGTCACAGAGAGGCTGAGCACTATCTC CTGAGAGAGCCCAAGAGAGAGTCAACCTATGAGGGCTTTTGGCAAAATCCD AGAGAGCTCAGAGTGGCTGGAGGGAGACGGAGTACATGGCCATGAGAGCTG TAGAAAGCTATGGGGCTCTTGGGGCTCTGGCCAGCTCAAGCTGAGGGAGACAG AGCGCTATGGGGTTTCTCGGGCCAGTCAAGGTTGAGCTGCTCAAGAGAT CCGAATCTACTCTGGAGCTTTTGGTCAAGAGCTCTTCTATGATATGATG CGAGAAACCTCTGACCTTTCAGTGGTGTGATGATCACTCTTATATG CCCTTCCCAAGCTCAGCTCAGATGATGACAAATCAAGCCAGCTACTCTC TCTGCTGGAGATAGATATGTTATCTCTGGGTCTGTAGGGAGAGTGGATGCT CCCACTATAGCTTATCTTGGCTCAGAGCTTGAATCAAACTATCAGAGCCG AGCAGAGCAGAGAGAGCTTTGATCTCTCTCAAGCTCTCAGGTGACTCTCA ATTCAGGGAGCAGAGAGATGTTCTTGTACCTGTAGGCTATATGCTCAACTA ACAGTCAATGCTCCCTTTAGAGAAATAGAGCTCTCCCTCAGCAGATATAC ACATATAACAAAGAGTAAAGCTTTAAAGAGGTAATAATATCAGCAGAGATCTA ACATATATCCCAACTCAAGAGCTCTGCGACTCTGAGAGCTGTT

43

Configure Page | Personal Data

Display options | Manage configurations | Reset configuration

Select from available configurations: Default

Display options

Flanking sequence at either end of transcript: 50

Number of base pairs per row: 60 bps

Intron base pairs to show at splice sites: 25

Show full intronic sequence:

Show exons only:

Line numbering: None

Show variants: No

Hide variants longer than 10bp:

Hide variants by frequency (MAF): Don't hide

Filter variants by consequence type: No filter

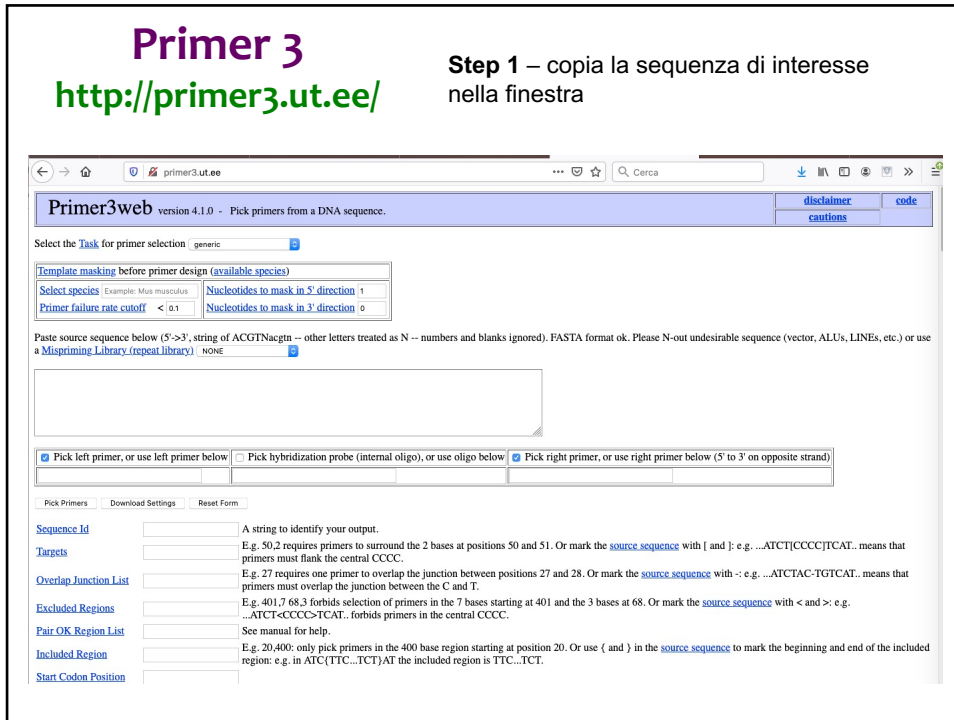
- 3 prime UTR variant
- 5 prime UTR variant
- NMD transcript variant
- coding sequence variant

44

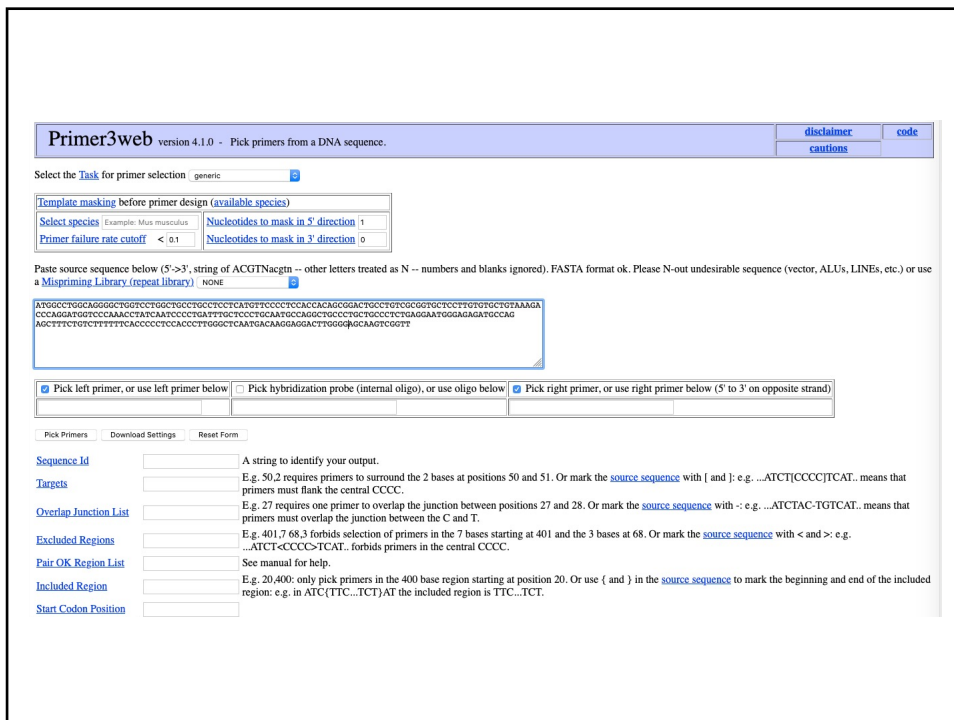
Primer 3

<http://primer3.ut.ee/>

Step 1 – copia la sequenza di interesse nella finestra



45



46

Exons/ Introns	Translated sequence	Flanking sequence	Intron sequence	UTR
Markup loaded				
	5' upstream sequence			
1	ENSE00002202717	1,994,265	1,993,911	- -
		1,993,910	1,983,104	
Intron 1-2				
2	ENSE00000655739	1,983,103	1,982,858	- 0
Intron 2-3				
3	ENSE00002253018	1,980,958	1,978,757	0 -
		1,982,955	1,980,859	1,997

Length	Sequence
375ctcaaaactgaagtggcggttccaagtgacaacagcactacacc AGAAATCCGCCTCCCAACAGCGGGCTGCTTCCCTCCCAACCACCCCGTCACAGAAAG AGCCTATTGTCCAGGCCAAGGAGTTCCAGTTCGAGAGGCCCTGGGGCTGTGCTCTAGC TGCTCAGCCACTTCCCATTTGGCTCCCAAGCTGCTGCTCAGCAAGGGCTGAGGAC AGGGAGGCTCTCTCCAFAMAGGGGGAAAGGACACAACTCCCACTTTGAGAGGGC TTGGTGGTGTCCAGCTGCTCTCTTGGACACTCTCCCAAGCCGGATGAGAGGGCC CTGAGCTTGGACCAGCCACTGCCACTCCAGCTCCGAGCTCCGACAGAGCTCCAGGGA CAAAGCAATCCAG
10,807	gttataatctcaggpppcaactgaaa.....tataatctcttttatctatcctcccccag caggaaattctgagcagagatggccttggcagggggctgctctggctccttccctcctcag gttccctccacacagcggactgccttggcgggcctctcttggtgctgctgfaagacccta ggatggctccaaacctatcaatcccctgt
2,202	gtagggttcaggcaagtctctca.....tgtaggtctttaggtctttag ATTTGCTCCCTGCMAATGCCAGGCTGCCCTGCTGCTGAGGAATGGGAGATGCCAG AGCTTCTGTCTTTTTCAGCCCTCCACCTTGGGCTCAAGCAGAGAGACTTGGG AGCAAGTCGGTTGGGAAAGGCCCTACAGTGAGCTGGCCAGCTCTGGGTCATCTCG AAGGACTGGAGAAAGCAAGTTTCTCCAGATATCTCAAAAGAGAGACACTCTGAGC AAGACCTTGAGAGAACTCAGGGCTCTCTGAGGAGTTTGGAGGGAGCAGAGTCT GAGCTGATGAGGATGCCAGCTGAAAGATGAGTGGCATGGAGACTGGACACTATCTC GCTGAGAGAGCCCAAGGAGCAGGTCAAACCTATGGGGGCTTTTGGCDAATACCCC AAGAGGCTCAGAGGTGGTGGGAGGGAGAGGGAGATGACATGGGCATGAGAGCTG TACAAGCTATGGGGCTTCTTGGGGCAATCTGCCCAAGCTCAGTGGGACACAG AAGCCCTATGGCGTTTCTCCGGGCCAGTCCAGGTGGTGGTCCGCTCAGAGAGAT CCGAATCTTCTGAGAGCTTTTGTAGTATAGAGCCCTTTTCCGAGATAGT CAGGAAACCCCTGACACTTTTCAAGTTGGAGTCACTTCACTCCCTTATATGTG CCCTTCCCCTGCTCAGCTCAGCATTTGGTACAAATCCAGGCCAGCTATCTCTC TTCTGCTGGGATATGATTTCTTCTGGGCTCTGTGAGGAGGGGAGGTGGATCCCTT CCCAATAGGCTTATGCTGGCTGAGCACTAGACTTAAAATATCAGAGGCGG AGCAGAGCAGCAGCAGGTTTGTACTGTCTTCCAACTTTCACGTGACTCTCA ATCCAGGAGCACAGGCGTGTGTCTTGGCTGTTGCTATATATGTCAACTCA ACAGATACATGCCCTCTTGAAGAAATATGAGCATCTCCCTCATGAGATATAC ACTATATAACAAGGATGAGCTTTAAAGAGGTAAATATATCATACAGAAATCTTA ACATTTATATCCAAATCTCAAAAGCTTCTTGCCTCTTCTTGGAGAGGCTTT AGTAAAGCTTAAACATTTGCTTATATTGATCAGAAACCTTACAGTAAAGGCTCA GCTTTAAGGGTTAATTTTACAGGACTCTGAGCTCAATTTAGAAAGG ACTCTTGAAGAACTGCTCTAACTCCAGTATTTGGACCTGGAGAGGCTTA ACTCTTACAGAAAGGAGTGAACCTCTTCCGAAATGATGGAGAGCCCTCTGAAAT GCTTAAATGATCAGGGGGAGAGGCAACAAATTTCTTGGACAACTCAAA TGGGACTATTTCCAGGCCATATGAACCAATGAACGAGGTGATCAGTCTTCA