



GÖTEBORGS
UNIVERSITET

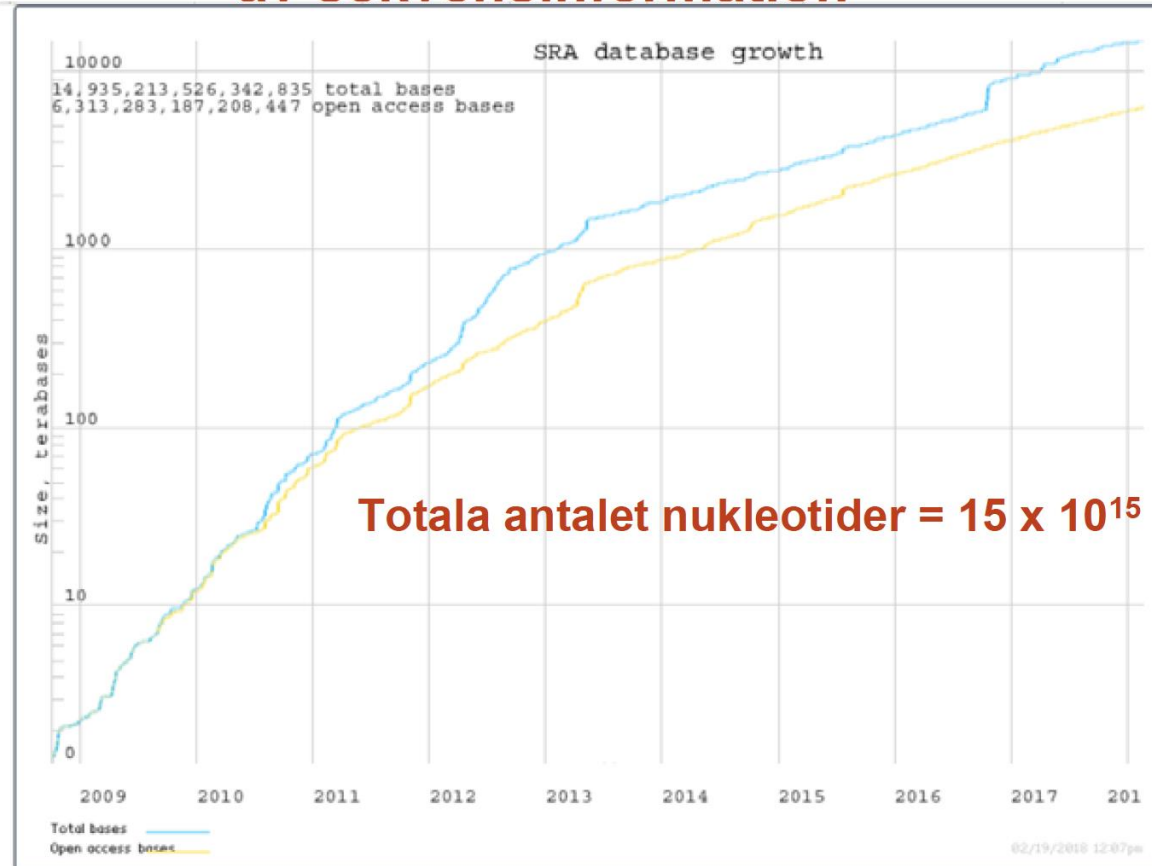
Bioinformatics
Core Facility

BMA Molekylärbiologisk metodik Bioinformatik

Feb/2019

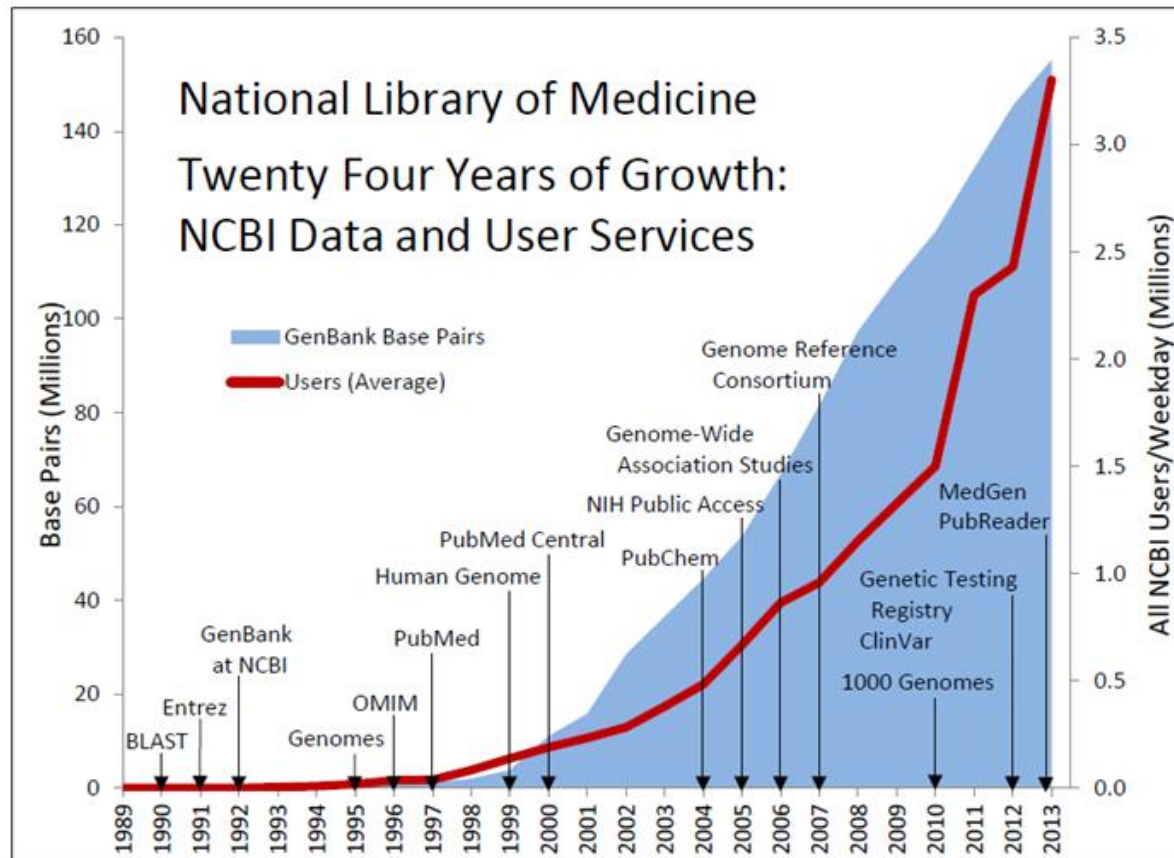


Explosiv tillväxt av sekvensinformation





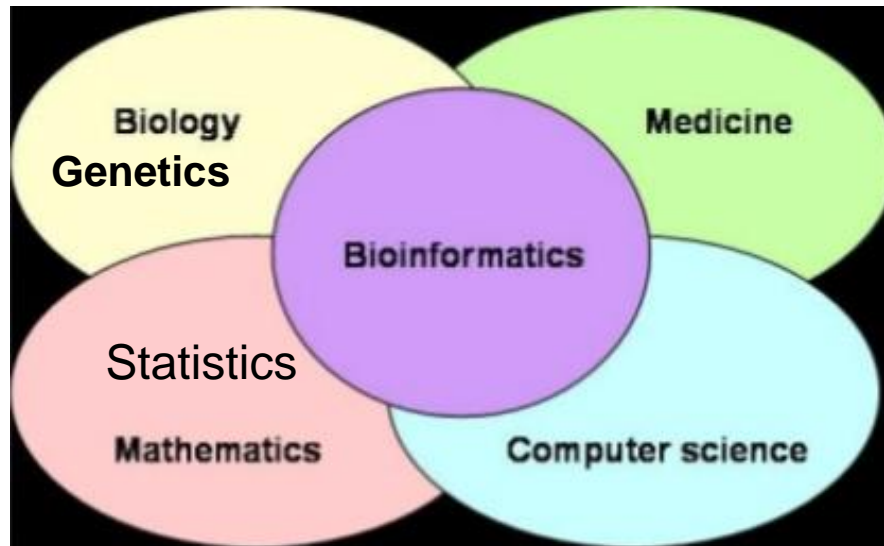
Growth of GenBank



Report from NCBI 2013

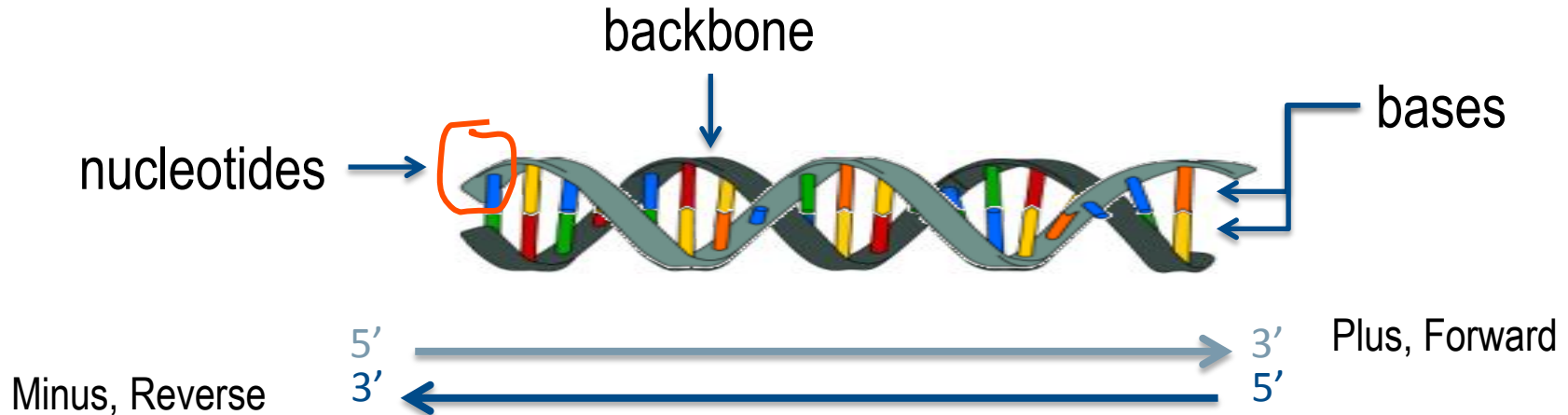
GenBank world's largest collection of DNA sequences in collaboration with UK/EMBL and Japan/DDBJ (DNA data bank of Japan)

Hantering och analys av information
från genteknologisk forskning
Tvärvetenskaplig disciplin





DNA



AGCTGACGATGGGCAGATACACAGTAAC
TCGACTGCTACCCGTCTATGTGTCATTG

Genes can be located on both strands. So the mRNA sequence always corresponds to the 5-3 coding sequence of a gene.

- Double stranded
- Antiparallel
- Complementary

A=T

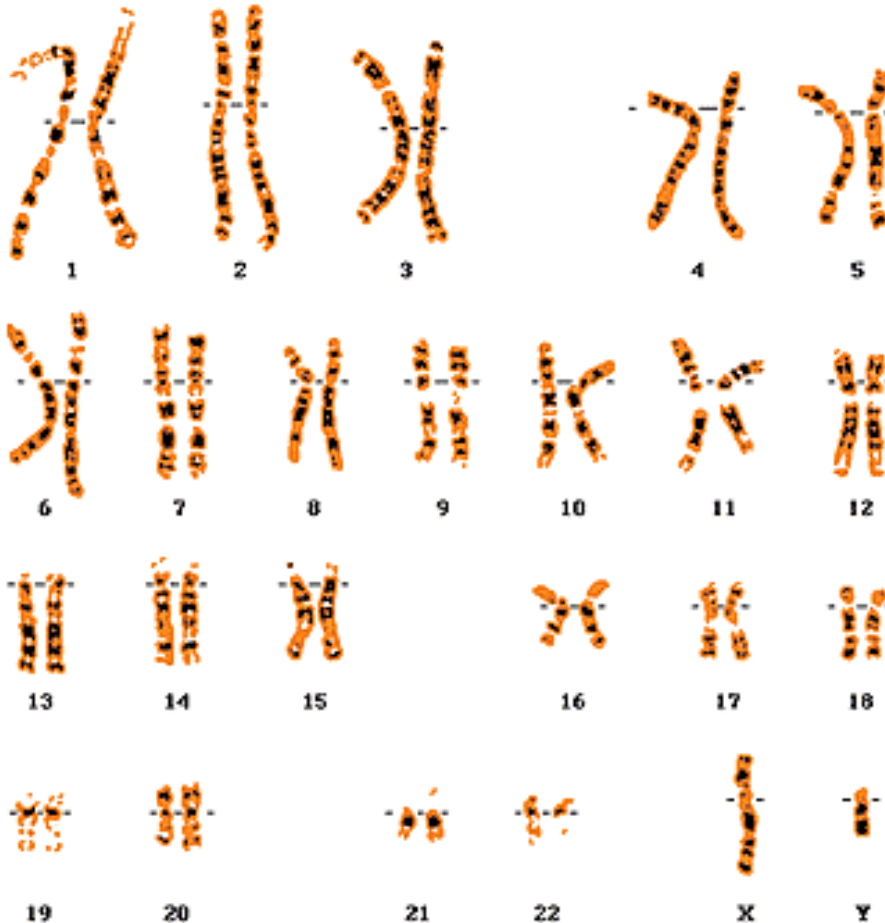
adenine=thymine

C≡G

cytosine≡guanine



Genomes split into chromosomes

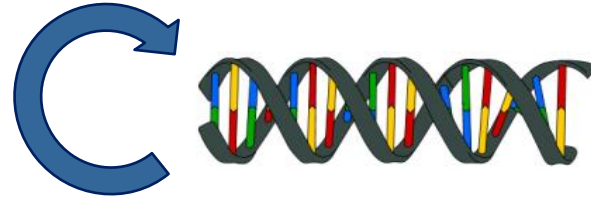


Human Genome
22 autosomal pairs
1 allosome pair
Diploid

Species	# Chr	Ploidy
A.thaliana	10	Diploid
Bread Wheat	42	Hexaploid
Tobacco	48	Tetraploid
Fruit fly	8	Diploid
Earth worm	36	Diploid
Mouse	40	Diploid
Human	46	Diploid
Dog	78	Diploid
Goldfish	100-104	Diploid

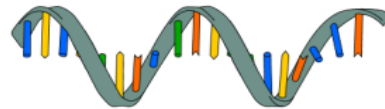
Central dogma - Flow of genetic information

Replication



Transcription

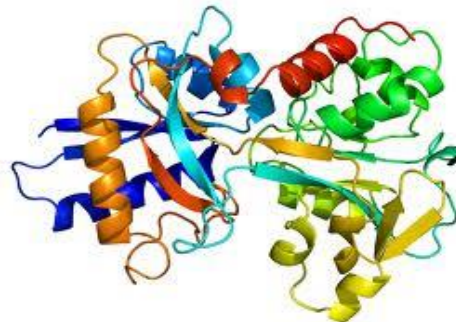
DNA/RNA



mRNA



Translation



Protein

20 amino acids Proteins - long chain of amino-acids

nonpolar polar basic acidic (stop codon)

Standard genetic code									
1st base	2nd base								3rd base
	U		C		A		G		
U	UUU	(Phe/F) Phenylalanine	UCU	(Ser/S) Serine	UAU	(Tyr/Y) Tyrosine	UGU	(Cys/C) Cysteine	U
	UUC		UCC		UAC		UGC		C
	UUA	(Leu/L) Leucine	UCA		UAA	Stop (Ochre)	UGA	Stop (Opal)	A
	UUG		UCG		UAG	Stop (Amber)	UGG	(Trp/W) Tryptophan	G
C	CUU		CCU	(Pro/P) Proline	CAU	(His/H) Histidine	CGU	(Arg/R) Arginine	U
	CUC		CCC		CAC		CGC		C
	CUA		CCA		CAA	(Gln/Q) Glutamine	CGA		A
	CUG		CCG		CAG		CGG		G
A	AUU	(Ile/I) Isoleucine	ACU	(Thr/T) Threonine	AAU	(Asn/N) Asparagine	AGU	(Ser/S) Serine	U
	AUC		ACC		AAC		AGC		C
	AUA		ACA		AAA	(Lys/K) Lysine	AGA	(Arg/R) Arginine	A
	AUG ^A	(Met/M) Methionine	ACG		AAG		AGG		G
G	GUU	(Val/V) Valine	GCU	(Ala/A) Alanine	GAU	(Asp/D) Aspartic acid	GGU	(Gly/G) Glycine	U
	GUC		GCC		GAC		GGC		C
	GUA		GCA		GAA	(Glu/E) Glutamic acid	GGA		A
	GUG		GCG		GAG		GGG		G

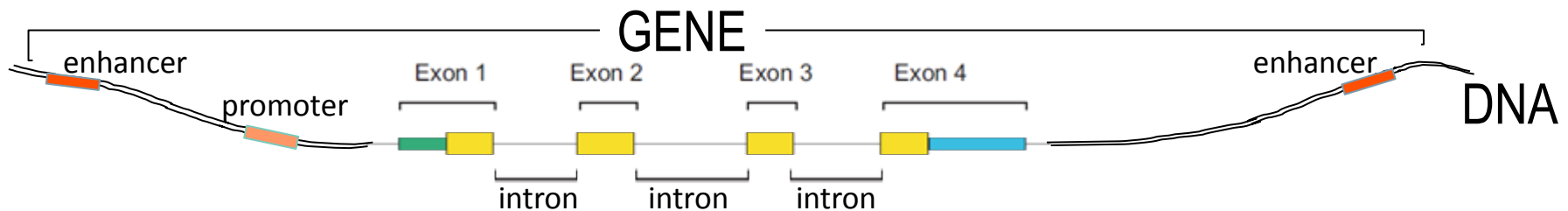
^A The codon AUG both codes for methionine and serves as an initiation site: the first AUG in an mRNA's coding region is where translation into protein begins.^[30]

Several cellular functions: regulation, structure, movement, catalysis, transport, signaling

Genes

A gene is a segment of DNA that encodes for a function.

Below portion of the DNA that contains all the information for production of a protein



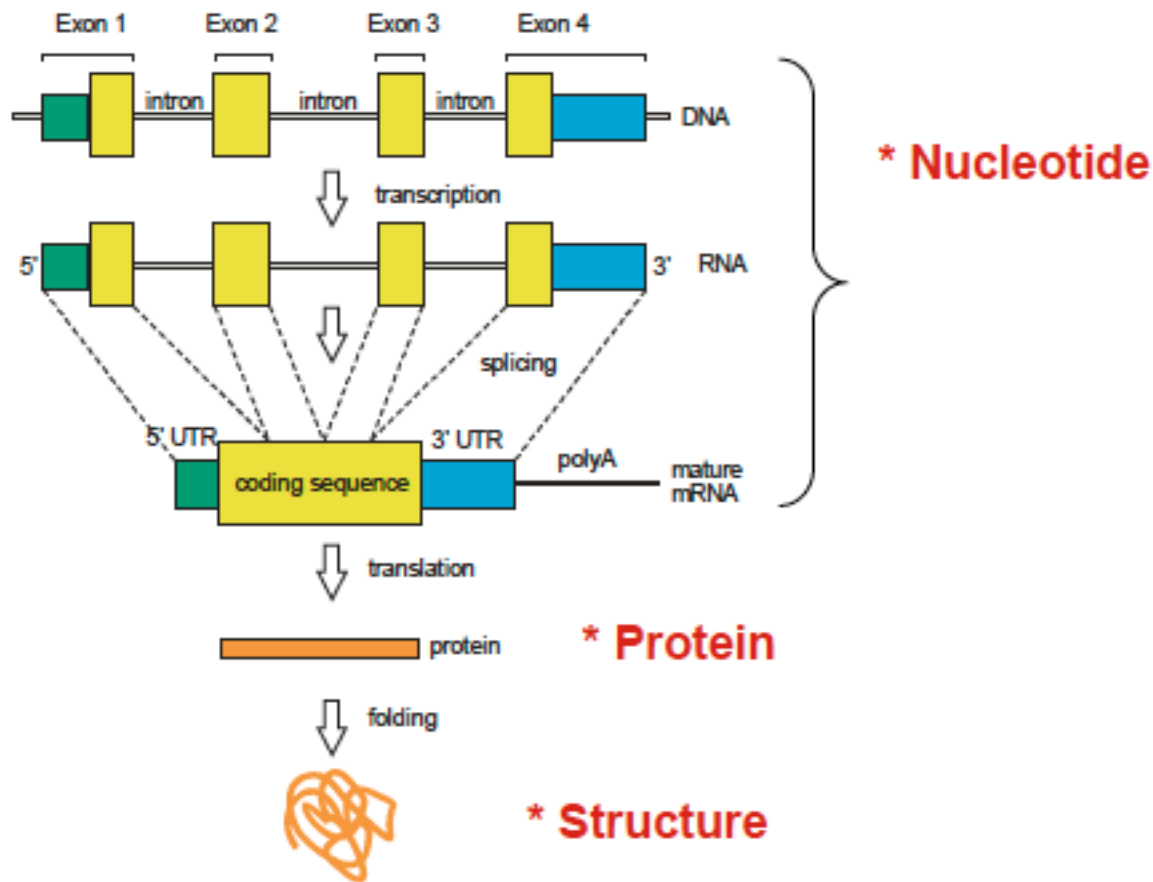
20,000 protein-coding genes (~1.0%)

Unevenly distributed across the chromosomes

HIST1H1A	781 nt	lack introns
DMD	2.2Mb	largest gene
TTN	80Kb (364 exons)	longest coding sequence



Coding sequence

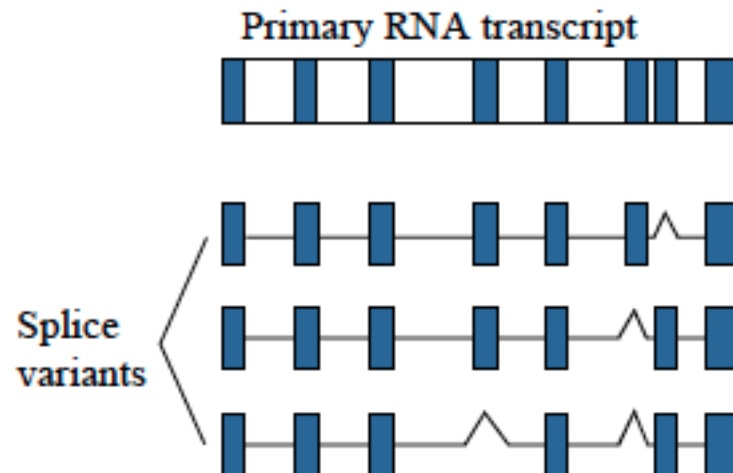




One gene - multiple transcripts and proteins

Alternative :

- splicing
- promoters
- polyA / termination sites





Variation mellan individer

Single nucleotide polymorphism (SNPs)

- Två mänskliga genom är i genomsnitt 99,9% lika
- Genom är stora så innebär ca 3 miljoner skillnader, men de flesta av dessa har ingen funktion
- De flesta varianterna 90% är i form av förändring av en nukleotid SNP





Functional Consequences

Type	Consequence
SNPs in coding area that alter aa sequence	Cause of most monogenic disorders, e.g: Hemochromatosis (HFE) Cystic fibrosis (CFTR) Hemophilia (F8)
SNPs in coding areas that don't alter aa sequence	May affect splicing
SNPs in promoter or regulatory regions	May affect the level, location or timing of gene expression
SNPs in other regions	No direct known impact on phenotype Useful as markers



Practical applications of genomic variation

- Disease diagnosis
- Association studies
- Pharmacogenomics
- Forensic testing
- Population genetics and evolutionary studies

Comparative Genomics



D. sp	UCGUAUACGAA	AGGACAGAAACUUGGCGA	UCCCCCUUGACAAGGAC	GGAAU	GAGC	UUUUUCUC	aaacU	UA	CA	GC	A	CCU	CGU	G	UGUUUGa	AGGC	ACCGGCUU	U	AGGACCAUUA	Aggca	auuugC	CACCAAAUUUUUUU	
R. oryzae	UGGUAUAUGAA	AGGACAGAAACUUGGCGA	AUCCCCUUGACAAGGAC	GGAAU	GAGC	UUUUUCUC		Aa	CA	UA	AC	CAC	CGU	G	UGUUUGa	AGGC	AGGCUUU	U	CGGACCAUUCU	ggcau	UAU	CCCAUAUUUUU	
P. blakesleezanus	AAACUAUAUGAA	AGGACAGAAACUUGGCGA	AUCCCCUUGACAAGGAC	GGAAU	GAGC	UUUUUCUC		A	CA	CU	AC	ACA	CCU	C	GGUUUGa	AGGC	AUUGCUU	A	CGGACCGUUU	ggcau	UAU	CCCAAAUUUUU	
H. sapiens	UGUUUGUAUGAA	AGGACAGAAACUUUACCA	UCCCCCUUGACAAGCAU	GGAAGA	GGCC	UUGGGGCG		U	CA	CA	AC	CCC	AUA	C	GGUUA	AGGC	AUUGCCACCUUACU	CGUGCA	U		CUAACAC	UGUUUUU	
G. gallus	UGUUUGUAUGAA	AGGACAGAAACUUUACCA	UCCCCCUUGAUAAGCA	GugGAACAGCGCC	UAGUGGCC		Ua	UA	CA	AC	A	CCC	AUA	C	GGUUA	AGGC	ACUGCCACCUUACU	CGUGCA	U		CUAACCA	UGUUUUU	
D. rerio	UGUUUGUAUAAU	AGGACAGAAAGUUUACCA	UCCCCCUUGACAAGAGU	GGAAAGCGUCC	UAGUGGCC		U	UA	CA	AC	A	CCC	AUA	C	GGUUA	AGGC	AUUGCCACCUUACU	CGUGCAUCU			UAACCAA	UUUUUUU	
T. nigroviridis	UGUUUGUAUGAA	AGGACAGAAAGUUUACCA	UCCCCCUUGACAAGAGU	GGAAAGCGCC	UAGUGGCC		U	UA	CA	AC	A	CAC	AUA	C	GGUUA	AGGC	aaUUUGCCACCUUACU	UGUGGCAUUC			UAACC	AAUUUUUUU	
F. rubripes	-UGUUUGUAUGAA	AGGACAGAAAGUUUACCA	UCCCCCUUGACAAGAGU	GGAAAGCGCC	UUGUGGCC		U	UA	CA	AC	A	CAC	AUA	C	GGUUA	AGGC	AUUGCCACCUUACU	CGUGGCAUUC			UAACC	AAUUUUUUU	
B. floridae	UGUUUGUAUGAA	AGGACAGAAAGUUUACCA	UCCCCCUUGAUAAGAGC	GGAAAGA	CCCC	UUGGGGGC		U	UA	CA	AC	CAC	AUA	C	GGUUA	AGGC	AUUGCGGGCUUACU	CGGGCA	U		CUAAC	AAUUUUUUU	
X. laevis	UGUUUGUAUGAA	AGGACAGAAAGUUUACCA	UCCCCCUUGACAAGAGU	GGAAAGCGUCC	CAGUGGCC		U	UA	CA	AC	A	CAC	AUA	C	GGUUA	AGGC	AUUGCCACCUUACU	CGUGGCAUUC			UAACC	AAUUUUUUU	
C. intestinalis	UGUUUGAAUAAA	AGGACAGCUGUUUUACCA	UCCCCCUUGACAAGAGC	GGACAUGUUUA	UGUGUUCA	Cccu	UA	CA	AU	A	CAC	AUA	A	UGGAUA	AGGC	ACUGUGCAACAU		UGGGCA	U		UAUCCAU	UUUUUUU	
S. purpuratus	UGUUUGUAUGAA	AGGACAGAAAGUUUACCA	UCCCCCUUGAUAAGGAC	GGAAU	U	CACC	UUGGGGUG		A	CA	AC	CAC	AUA	U	GGUUA	AGGC	ACUGCCACCUUACU	UGUGGCA	U		CUAACAC	GUUUUUU	
T. spiralis	UGCGGAUAUGA	AGGACAGCUAGUUUACCA	UCCCCCUUGACAAGAGC	GGCAAAU			U	CU	Cag	AC	A	AGC	ACA	A	CAGCAA	AGC		AUCCGCAU	CCU			UUGCUCCCCAAUUUUU	
A. aegypti	UGUCUUUUUAA	AGGACAGCAGGUUUACCA	UCCCCCUUGACAAGGAC	GGAA				C	UA	AAC	C	GGCG					GGAGACGGCACAA	AAUUUCC				GGCCC	UAGUCCAAU
A. gambiae	UGUCAUUUCAA	AGGACAGCAGGUUUACCA	UCCCCCUUGACAAGGAC	GGAA				C	CAU	ACC	C	UGCAU	CACU	GG			GCACCCA	GACAAAACUG				AUUCUGGGCCAUUUU	
A. mellifera	UGUGUAUUAAG	AGGACAGACCUUUUACCA	UCCCCCUUGACAAGGAC	GGAAUAC				AC	A	UAU	AU	U	GGUAA		GGC		ACUGUGCU	CCUUU	GGCGACA	C		UUACC	AAAUUUUUU
B. mori	UGUCUAUUCAA	AGGACAGCAGGUUUACCA	UCCCCCUUGACAAGGAC	GGAA				C	A	CAC	AAU	U	UGGUUA		AGGC	AC		U	A			UAACCAAACAUUUU	
D. melanogaster	UGUUGUUUGCA	AGGACAGCAAGUUUACCA	UCCCCCUUGACAAGGAC	GGAA				C	CAU	AAU	C	GGUCG					CGUAGGCACACA	CAAAAGC				CGUCCACAAUUUUU	
D. pseudoobscura	UGUUGUUUGCA	AGGACAGCAAGUUUACCA	UCCCCCUUGACAAGGAC	GGAA				C	CAU	AAU	C	GAGCG					CGUGCG	CAACA	AGAAAGC			CGUUUACCAUUUUU	
T. castaneum	UGUGUAUUGAA	AGGACAGACGUUUUACCA	UCCCCCUUGACAAGGAC	GGAA				AC	A	CAA	AUA	C	GGUUUA				GC	CU	C	CGACAC		UCAAC	CCAUUUUUU
D. pulex	UGUGUAUCAA	AGGACAGCUUUUACCA	UCCCCCUUGACAAGGAC	GGAA				C	CAU	AAU	U	GAGUCA	CGGG				ACAAACUUU	UGCG	A	GAAGU	G	UA	UCAUAAUUUUU
T. adhaerens	UGUUUGUAUGAC	CCGACAGAAAGUUUACCA	UUCUUCCUGCAUAAAGGAC	GAGAAA	ACUC	UUGGGACU	U	U	UC	AAc	auaC	A	CAC	AAU	U	GGUUA	AGGC	AUUGCCACCUUACU	UGUGGCA	U		CUGACACUUUUUAA	
H. magnipapillata	UGUUUGUAUGAC	CCGACAGCCGUUUACCA	UUCUUCCUGCAUAAAGGAC	GAGAAUAGCAC	CUGA	GUC	U	AU	CA	AC													

BLAST (basic local alignment search tool) är en algoritm* som används för att jämföra sekvens information.

Klistra in en sekvens och hitta liknande sekvenser

BLASTN programs search nucleotide databases us

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#) [Clear](#)

Query subrange [?](#)

From

To

Or, upload file no file selected [?](#)

Job Title

Enter a descriptive title for your BLAST search [?](#)

☐ Align two or more sequences [?](#)

Choose Search Set

Database

☐ Human genomic + transcript ☐ Mouse genomic + transcript ☒ Others (nr etc.):

Nucleotide collection (nr/nt) [?](#)

Enter organism name or id--completions will be suggested ☐ exclude [+](#)

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown [?](#)

☐ Models (XM/XP) ☐ Uncultured/environmental sample sequences

☐ Sequences from type material

[YouTube](#) [Create custom database](#)

Enter an Entrez query to limit search [?](#)

Organism
Optional

Exclude
Optional

Limit to
Optional

Entrez Query
Optional

*Algoritm är en mängd väl definierade instruktioner för att lösa en uppgift

Sequences producing significant alignments:

Select: [All](#) [None](#) Selected:0

AT AT		Alignments	Download	GenBank	Graphics	Distance tree of results						
		Description	Max score	Total score	Query cover	E value	Ident	Accession				
<input type="checkbox"/>		PREDICTED: Nomascus leucogenys rhodopsin (RHO), mRNA	187	187	100%	4e-44	100.00%	XM_003265030.3				
<input type="checkbox"/>		PREDICTED: Pan troglodytes rhodopsin (RHO), mRNA	182	182	100%	2e-42	99.01%	XM_516740.7				
<input type="checkbox"/>		PREDICTED: Pongo abelii rhodopsin (RHO), mRNA	182	182	100%	2e-42	99.01%	XM_002813145.4				
<input type="checkbox"/>		PREDICTED: Pan paniscus rhodopsin (RHO), mRNA	182	182	100%	2e-42	99.01%	XM_003829435.3				
<input type="checkbox"/>		Homo sapiens rhodopsin (RHO) gene, RHO-930CG allele, complete cds	182	182	100%	2e-42	99.01%	KP718610.1				
<input type="checkbox"/>		Homo sapiens rhodopsin (RHO) gene, RHO-284TC allele, complete cds	182	182	100%	2e-42	99.01%	KP734176.1				

Mest lik sekvens i toppen av listan. Hög score och E-värde nära 0 betyder att det är liten risk att sekvensen av ren slump är lik din frågesekvens

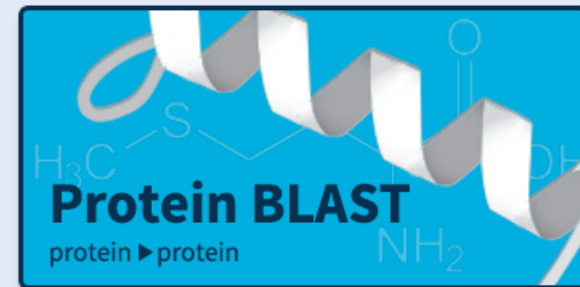
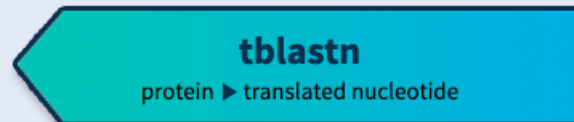
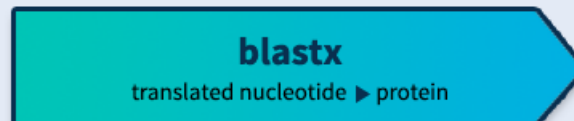
Olika typer av BLAST

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](#)

Web BLAST



Man kan söka med både DNA, RNA och protein sekvens

Sequencing & identification of SNPs



Sequencing



Variant, med allele G eller T



SNP/Indels



Next Generation Sequencing (NGS) Workflow

Ut från sekvenseringsmaskin - Rådata miljontals “reads”
t.ex. exom 30 miljoner reads, 150 baspar långa
fil-format FASTQ



Mappa mot referens
t.ex. BWA (Burrows-Wheeler Aligner)
Output bamfiler binärt format



“Calling” output: identifierade varianter, filer i s.k. VCF format



Annotera, filtrera ut kandidat-varianter som kan förklara fenotyp

FASTQ format

```
@HWI-D00457:83:C6EAMANXX:8:1101:6854:2245 1:N:0:CGCTCATTGGCTCTGA  
GAAGCAGGTGCCATGCCTGCACGTTTGTGGCTTAATGACCAAGGAGGGCCGATGAGCAGCCATGGT  
GGTGATCACTGC+  
B<<BBFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF  
FF
```

Varje “read” består av fyra rader. Första raden börjar med @ och har ett unikt namn.

Andra raden visar nukleotid sekvensen

Tredje raden har ett + tecken

Fjärde raden visar kvalitén för varje läst nukleotid



VCF format

```
##fileformat=VCFv4.0
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=1000GenomesPilot-NCBI36
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=.,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
```

Header

Body

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NA000001
20	14370	rs6054257	G	A	29	PASS	NS=3;DP=14;AF=0.5;DB;H2	GT:GQ:DP	0 0:48:1
20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0.017	GT:GQ:DP	0 0:49:3
20	1110696	rs6040355	A	G,T	67	PASS	NS=2;DP=10;AF=0.333,0.667;AA=T;DB	GT:GQ:DP	1 2:21:6
20	1230237	.	T	.	47	PASS	NS=3;DP=13;AA=T	GT:GQ:DP	0 0:54:7
20	1234567	microsat1	GTCT	G	50	PASS	NS=3;DP=9;AA=G	GT:GQ:DP	0/1:35:4

VCF format

Alternate allele

dbSNP ID

Defined in the
header

Sample
values

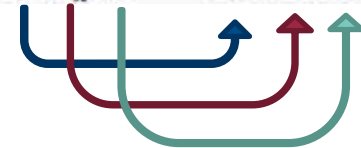
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NA000001
20	14370	rs6054257	G	A	29	PASS	NS=3;DP=14;AF=0.5;DB;H2	GT:GQ:DP	0 0:48:1
20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0.017	GT:GQ:DP	0 0:49:3
20	1110696	rs6040355	A	G,T	67	PASS	NS=2;DP=10;AF=0.333,0.667;AA=T;DB	GT:GQ:DP	1 2:21:6
20	1230237	.	T	.	47	PASS	NS=3;DP=13;AA=T	GT:GQ:DP	0 0:54:7
20	1234567	microsat1	GTCT	G	50	PASS	NS=3;DP=9;AA=G	GT:GQ:DP	0/1:35:4

Reference allele

How reliable is
The SNP

GT: Genotype - 0/0 homozygous ref
- 0/1 heterozygous
- 1/1 homozygous alt

DP: Depth



Variant annotering

- Är det en känd SNP?
- Hur vanlig är den i befolkningen?
- Ligger varianten i en gen?
- Vilken påverkan kan varianten ha på genens funktion?
- Vad vet vi om den genens funktion idag?

Genome Browsers

- A genome browser is a graphical interface to display an integrated picture of data from several databases including e.g.
 - Genes
 - Proteins
 - Expression
 - Regulation
 - Variation
 - Comparative analyses

Big Genome Browser:

Ensembl Genome browser

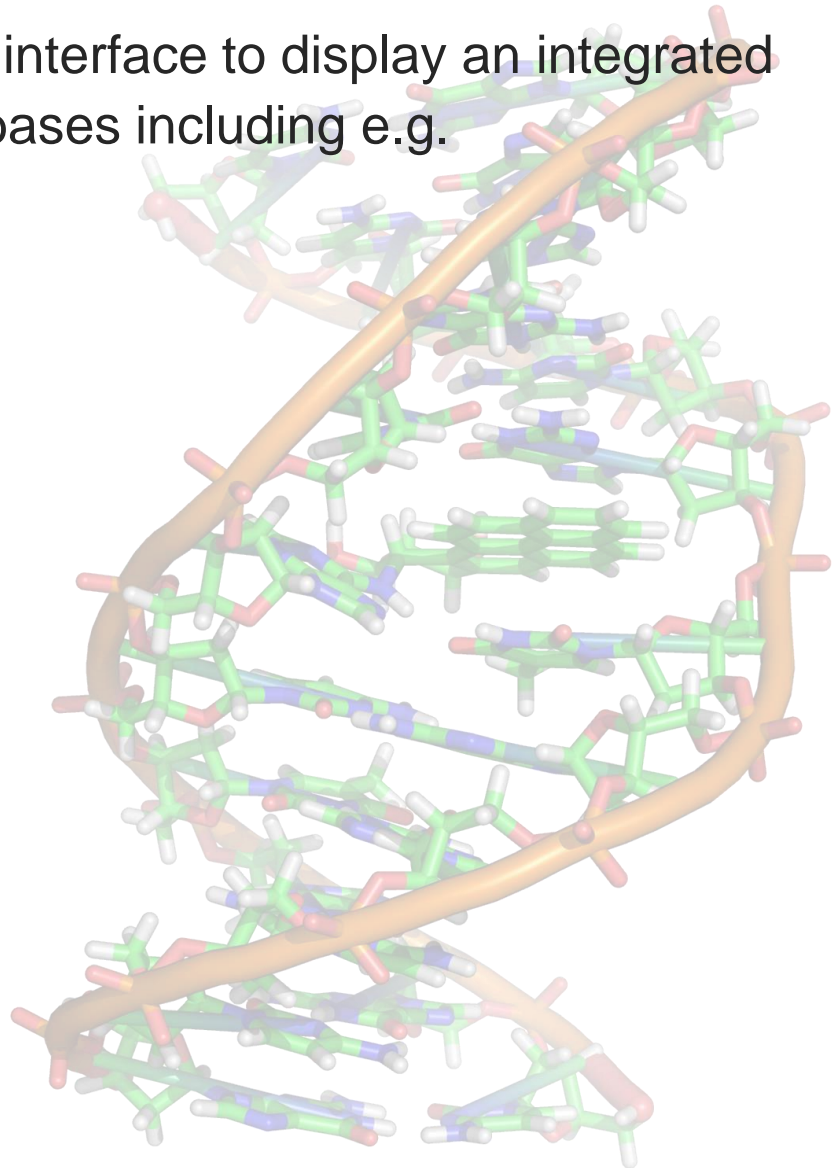
<http://www.ensembl.org>

NCBI

<https://www.ncbi.nlm.nih.gov>

UCSC Genome Browser

<https://genome.ucsc.edu>



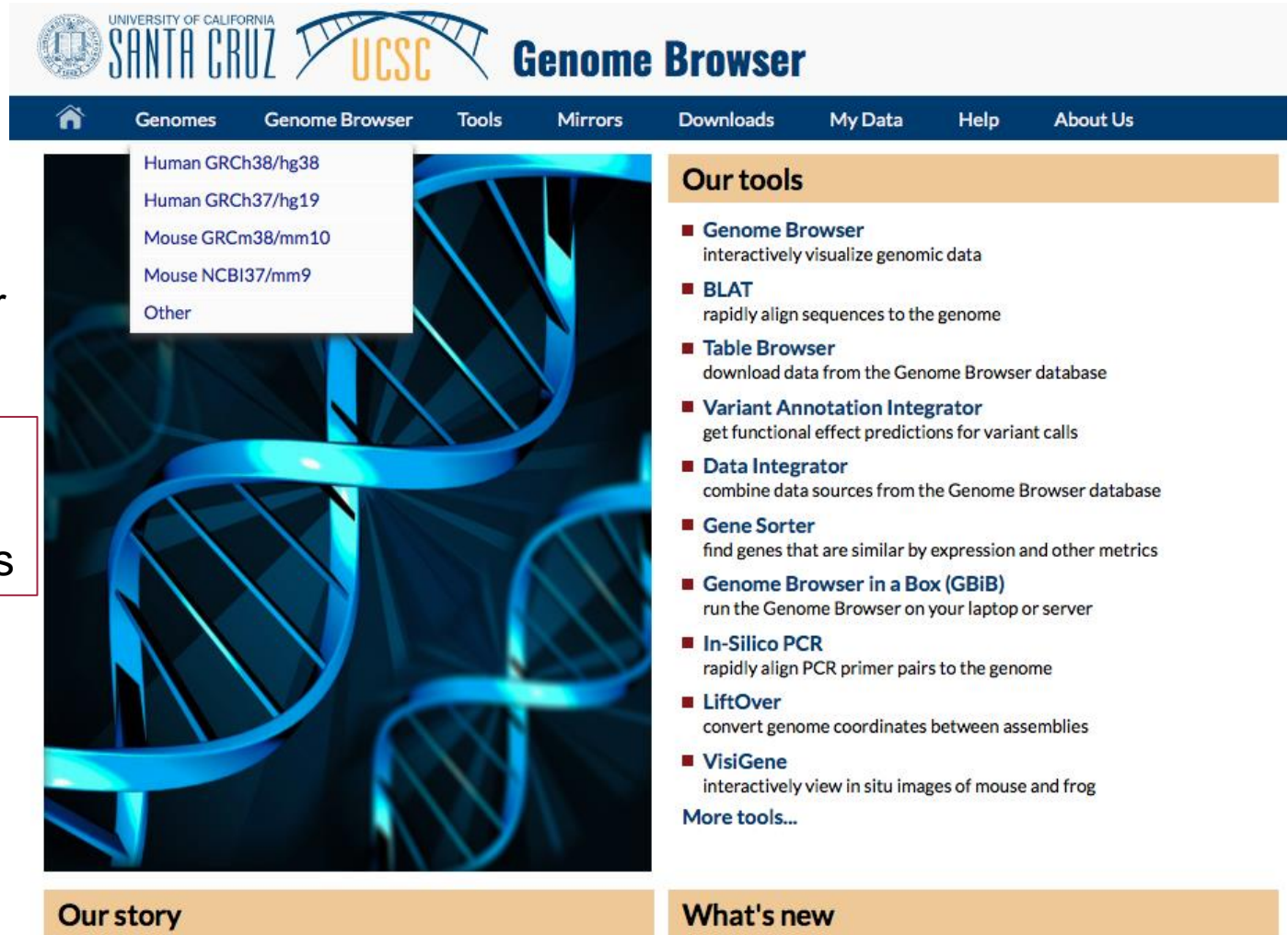
UCSC Browser <https://genome.ucsc.edu>

On-line genome browser with access to genome sequence data from a variety of organisms integrated with a large collection of aligned annotations.



Choose species,
assembly, and go
to Genome Browser

NOTE! Genomic
positions differ
between assemblies



UNIVERSITY OF CALIFORNIA
SANTA CRUZ

UCSC Genome Browser

Home Genomes Genome Browser Tools Mirrors Downloads My Data Help About Us

Human GRCh38/hg38
Human GRCh37/hg19
Mouse GRCm38/mm10
Mouse NCBI37/mm9
Other

Our tools

- Genome Browser**
interactively visualize genomic data
- BLAT**
rapidly align sequences to the genome
- Table Browser**
download data from the Genome Browser database
- Variant Annotation Integrator**
get functional effect predictions for variant calls
- Data Integrator**
combine data sources from the Genome Browser database
- Gene Sorter**
find genes that are similar by expression and other metrics
- Genome Browser in a Box (GBiB)**
run the Genome Browser on your laptop or server
- In-Silico PCR**
rapidly align PCR primer pairs to the genome
- LiftOver**
convert genome coordinates between assemblies
- VisiGene**
interactively view in situ images of mouse and frog

[More tools...](#)

Our story What's new



UCSC Genome Browser

Genomes Genome Browser Tools Mirrors Downloads My Data View Help About Us

UCSC Genome Browser on Human Feb. 2009 (GRCh37/hg19) Assembly

move <<< << < > >> >>> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x 100x

chr7:140,433,813-140,624,564 190,752 bp. enter position, gene symbol or search terms go

chr7 (q34) p21.3 14.3 14.1 q21.11 22.1 q31.1 7q33 q34 q35

Click on a feature for details. Click or drag in the base position track to zoom in. Click side bars for track options. Drag side bars or labels up or down to reorder tracks. Drag tracks left or right to new position. Press "?" for keyboard shortcuts.

move start move end

< 2.0 >

track search default tracks default order hide all manage custom tracks track hubs configure multi-region reverse resize refresh

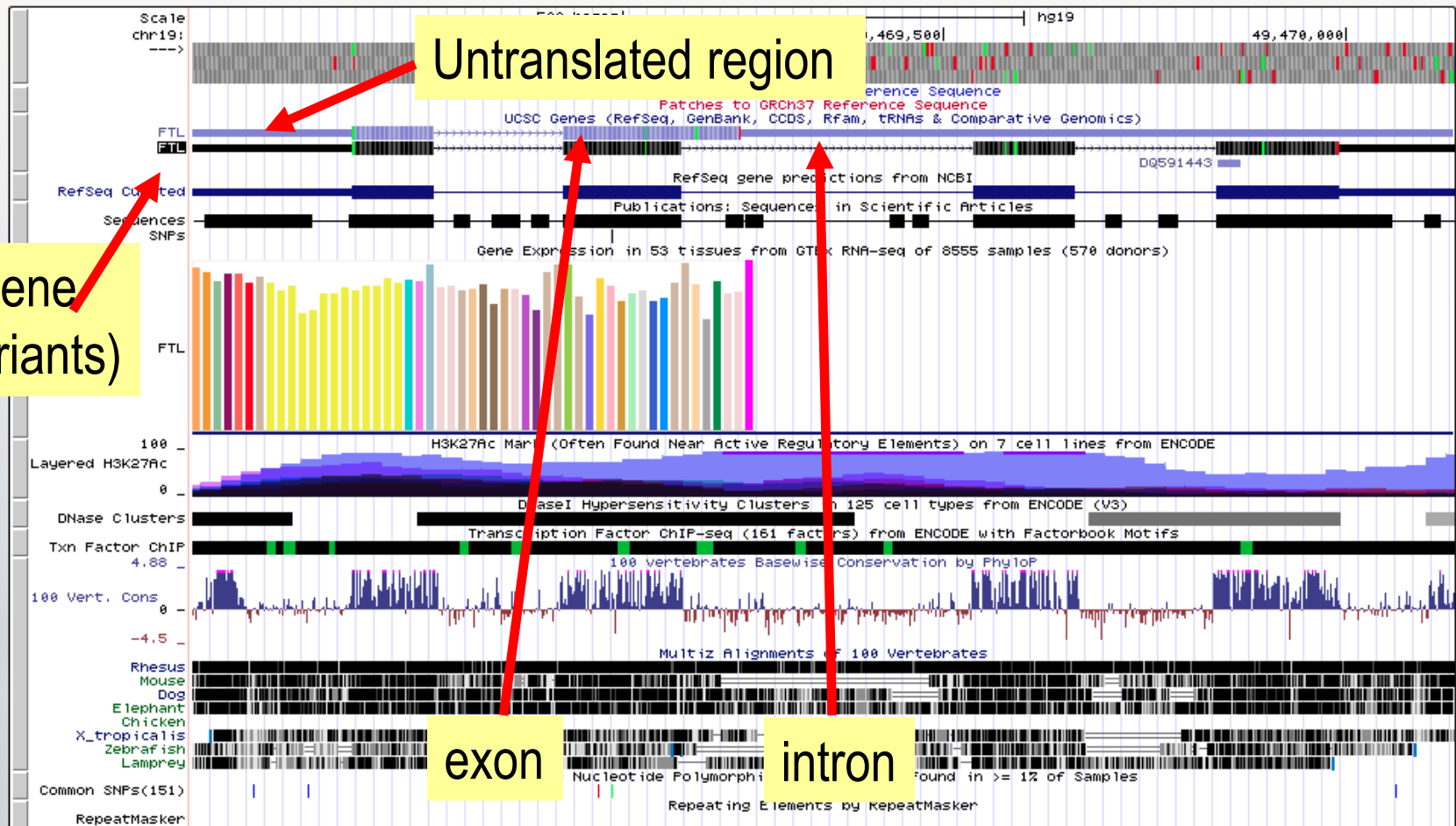
UCSC Genome Browser

UCSC Genome Browser on Human Feb. 2009 (GRCh37/hg19) Assembly

move <<< << < > >> >>> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x 100x

chr19:49,468,566-49,470,136 1,571 bp. go

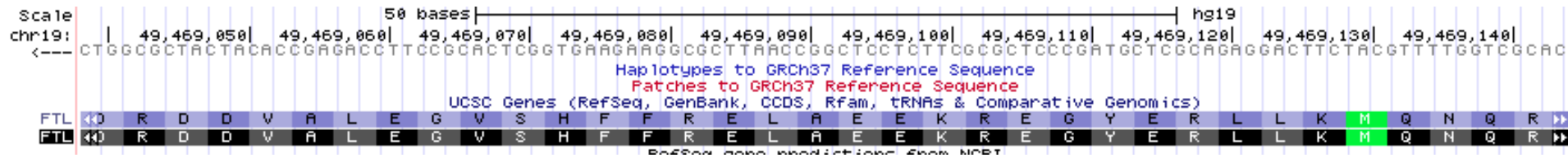
chr19 (q13.33) 19p13.3 19p13.2 p13.11 19p12 19q12 13.11 13.12 19q13.2 13.32 q13.33 13.42 13.43



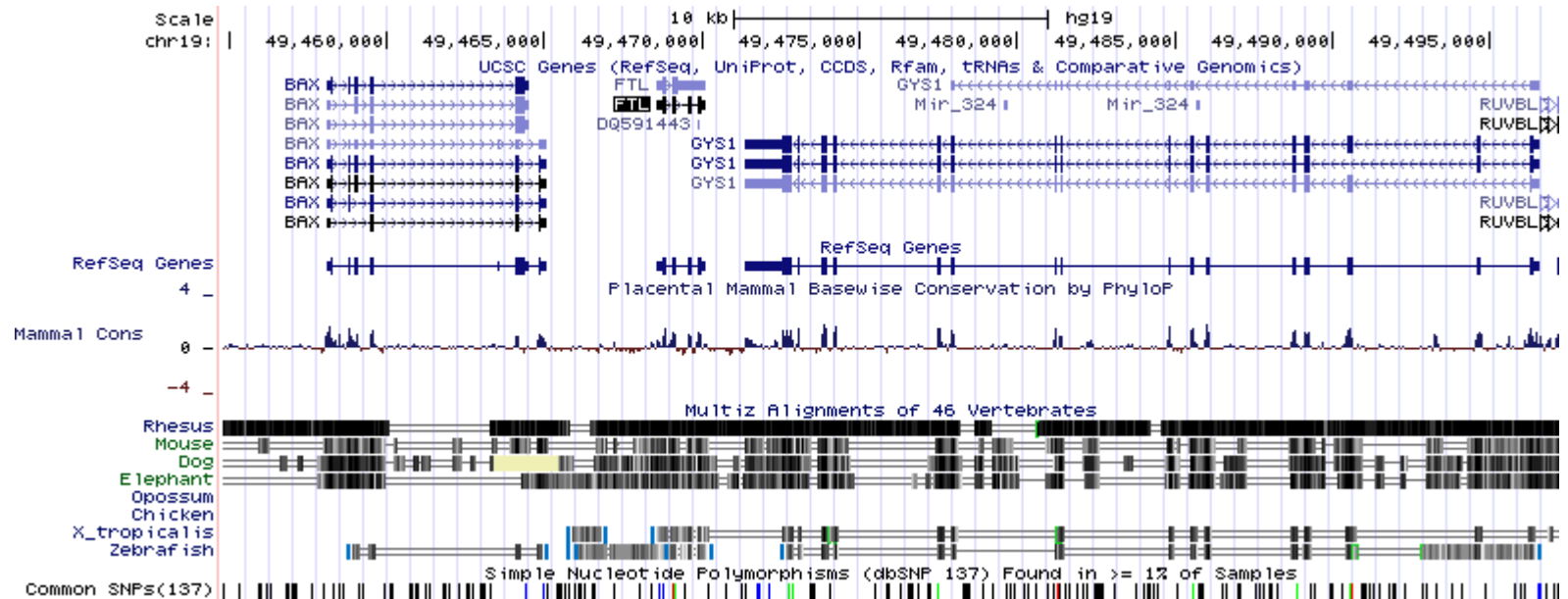


UCSC Genome Browser

Zoom in:

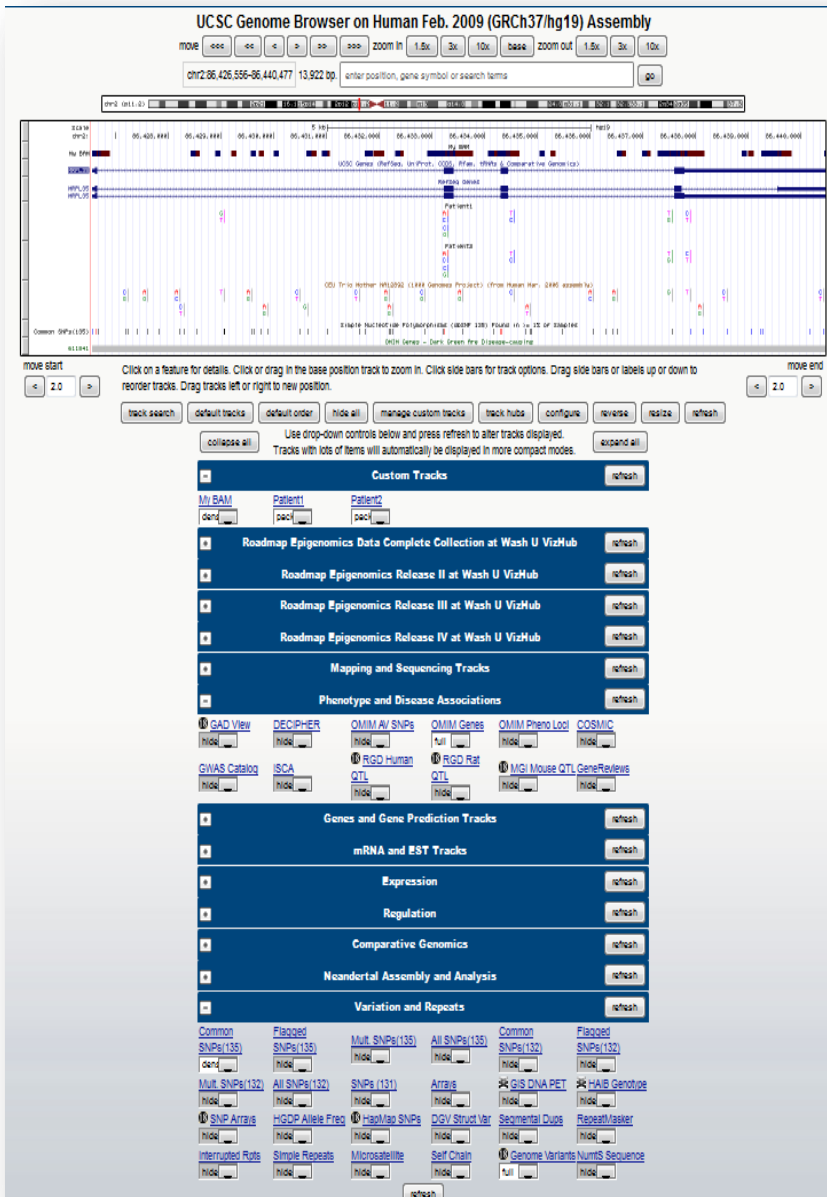


Zoom out:





UCSC Genome Browser



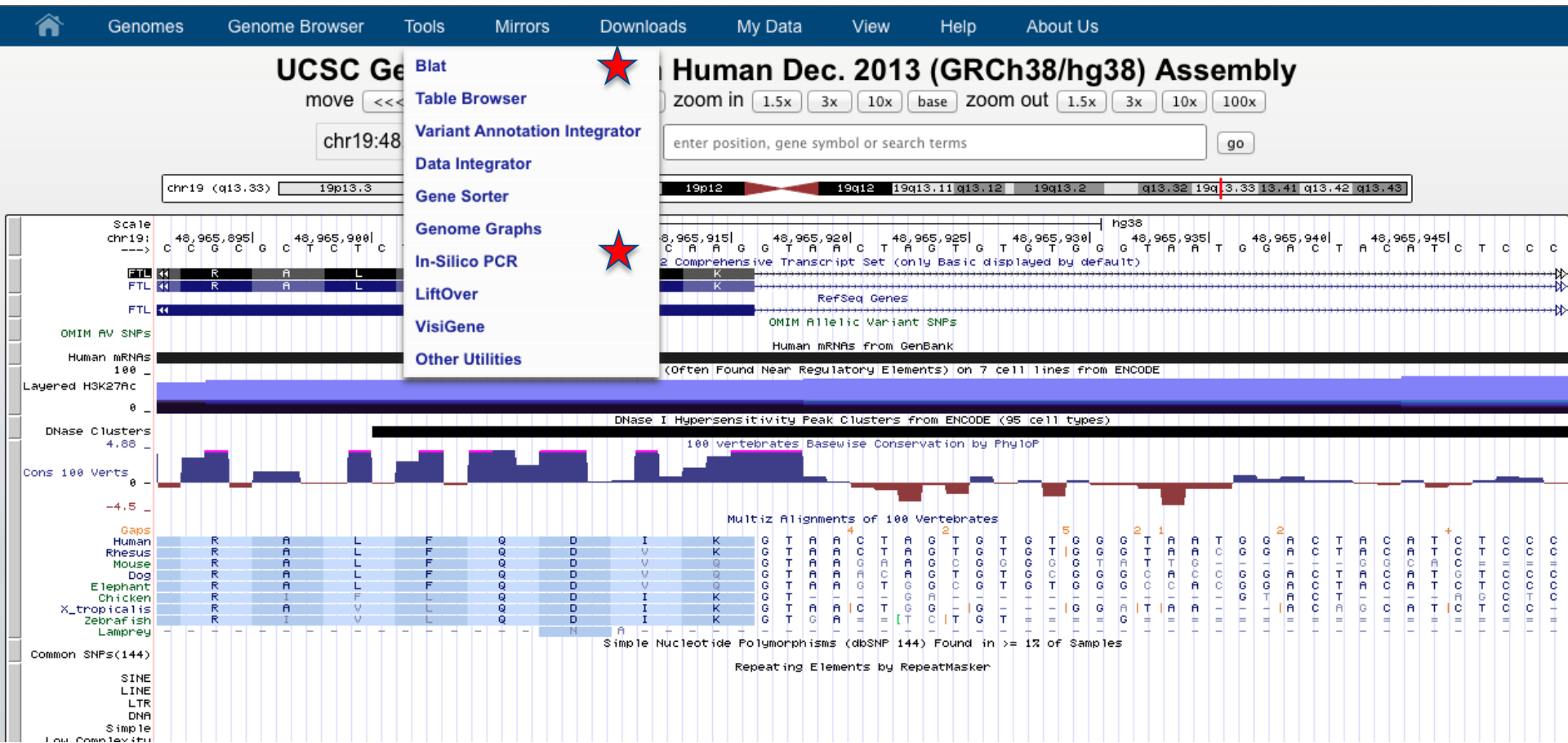
Genome viewer

Custom Annotation tracks
(groups of data)

- Mapping and sequencing
- Gene predictions
- Variation and Repeats
- Comparative genomics
- Custom tracks (BAM, VCF)

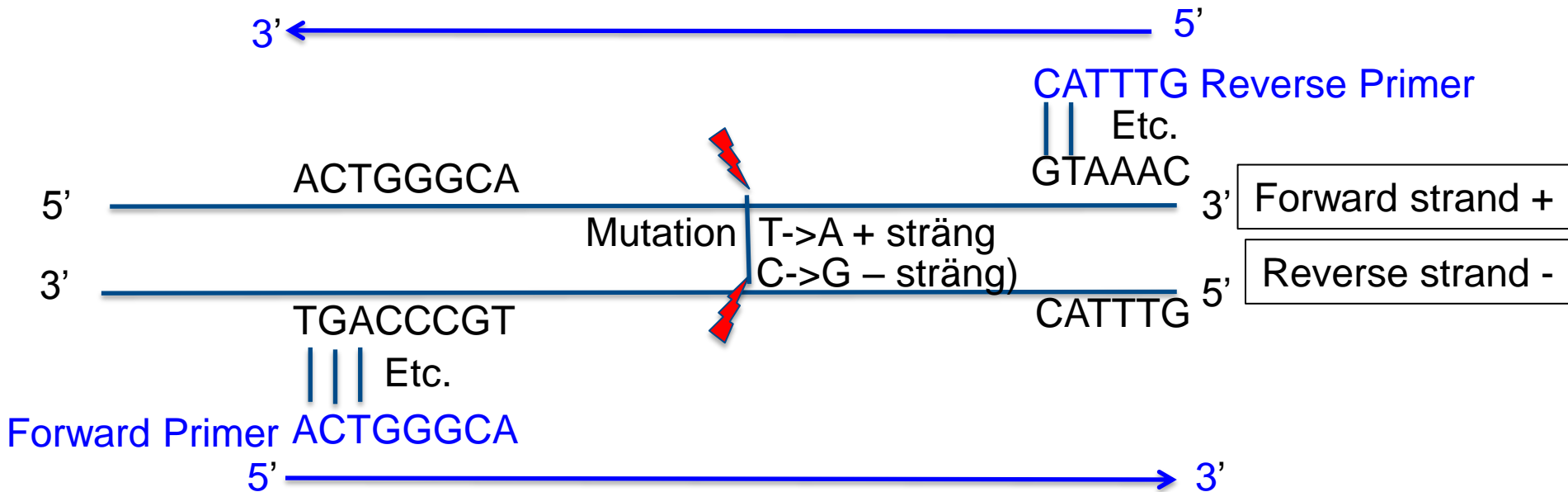


UCSC Genome Browser



BLAT har liknande funktion som BLAST, men behöver en nästan exakt match för att hitta rätt. BLAT passar inte om man vill hitta avlägset relaterade sekvenser.

Polymerase Chain Reaction (PCR)



Polymeras läser DNA mall i 3'-> 5' men bildar ny sträng med riktning 5'->3'

Männsikan läser alltid i riktning 5' -> 3'

Alltså läser vi forward primer: ACTGGGCA

Men reverse primer: GTTTAC

OBS! Det är viktigt att man hittar primers som binder unikt till ett ställe i genomet

Sammanfattning

- Bioinformatik är ett tvärvetenskapligt fält i rörelse
- Central dogma:
 - Replikation, Transkription, Translation
- Online verktyget Blast
- Identifiering av SNPs
- Genome Browsers
- PCR för att snabbt identifiera kända mutationer i patienter