

Making Sense of Genomes

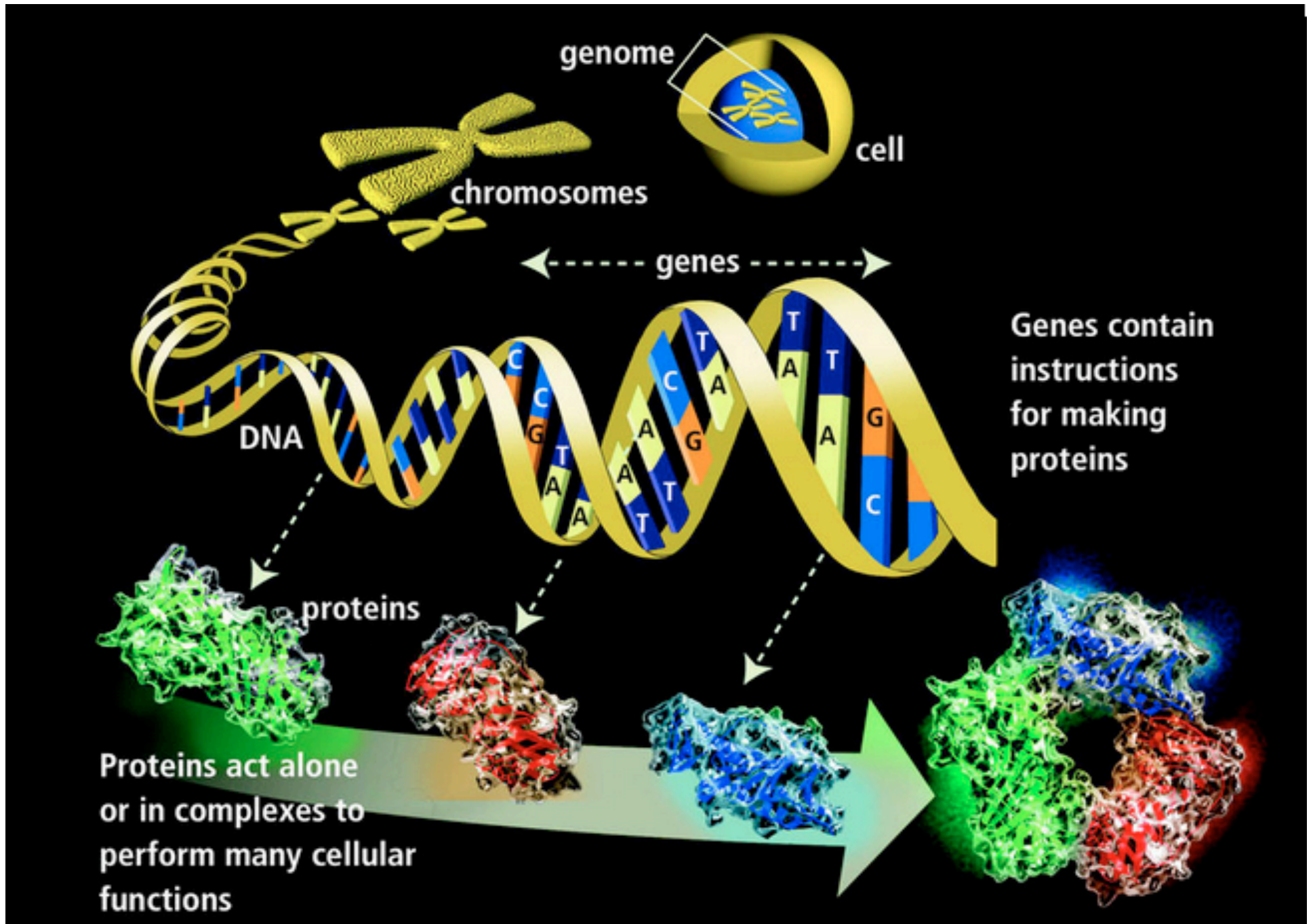
Jane Loveland

Open Door Workshop
11th May 2015
Hinxton

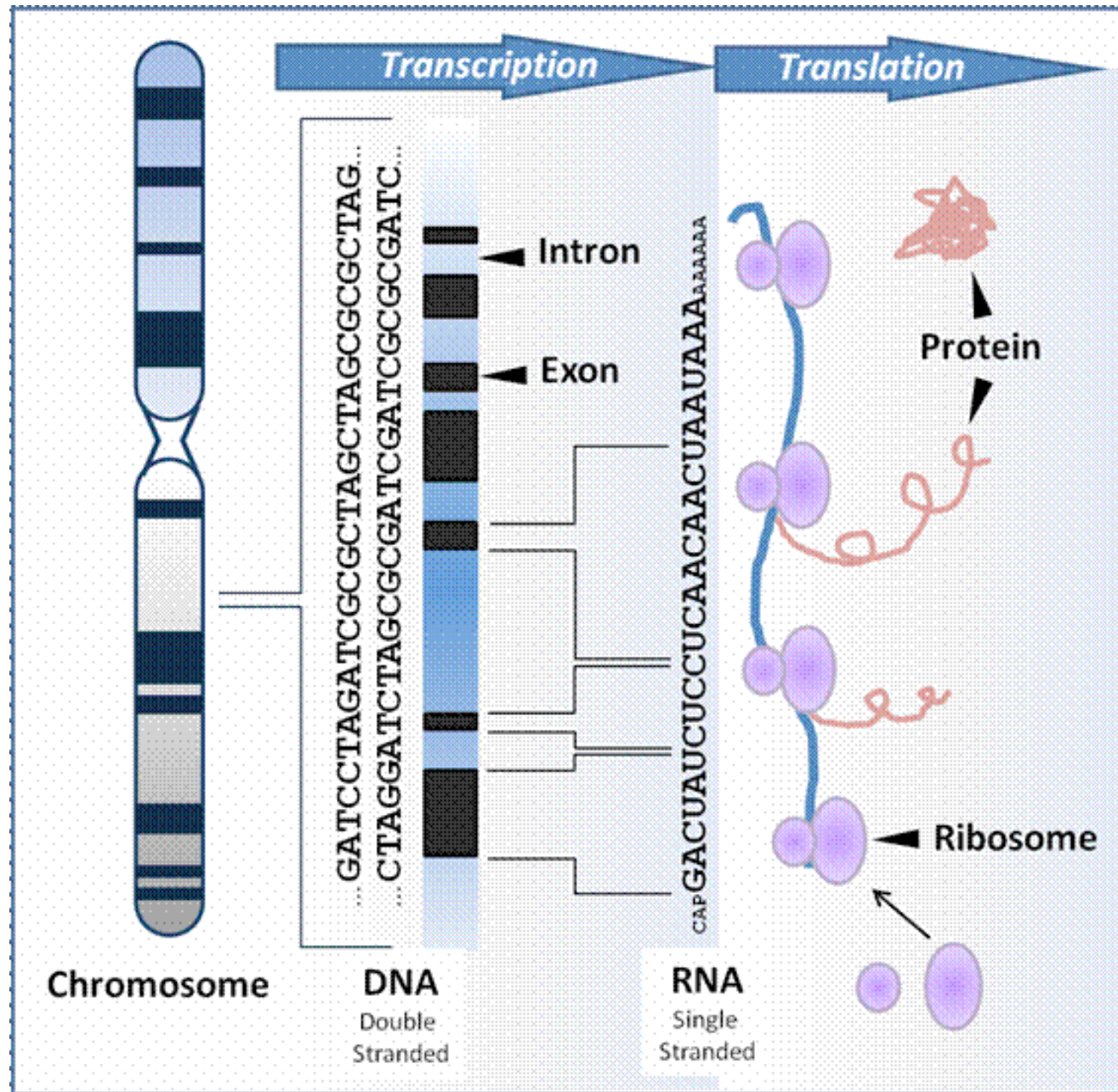
Overview

- A bit of background
- Genomes
- Genes
- Some bioinformatics basics

- **A bit of background**
- Genomes
- Genes
- Some bioinformatics basics



Central Dogma Of Molecular Biology

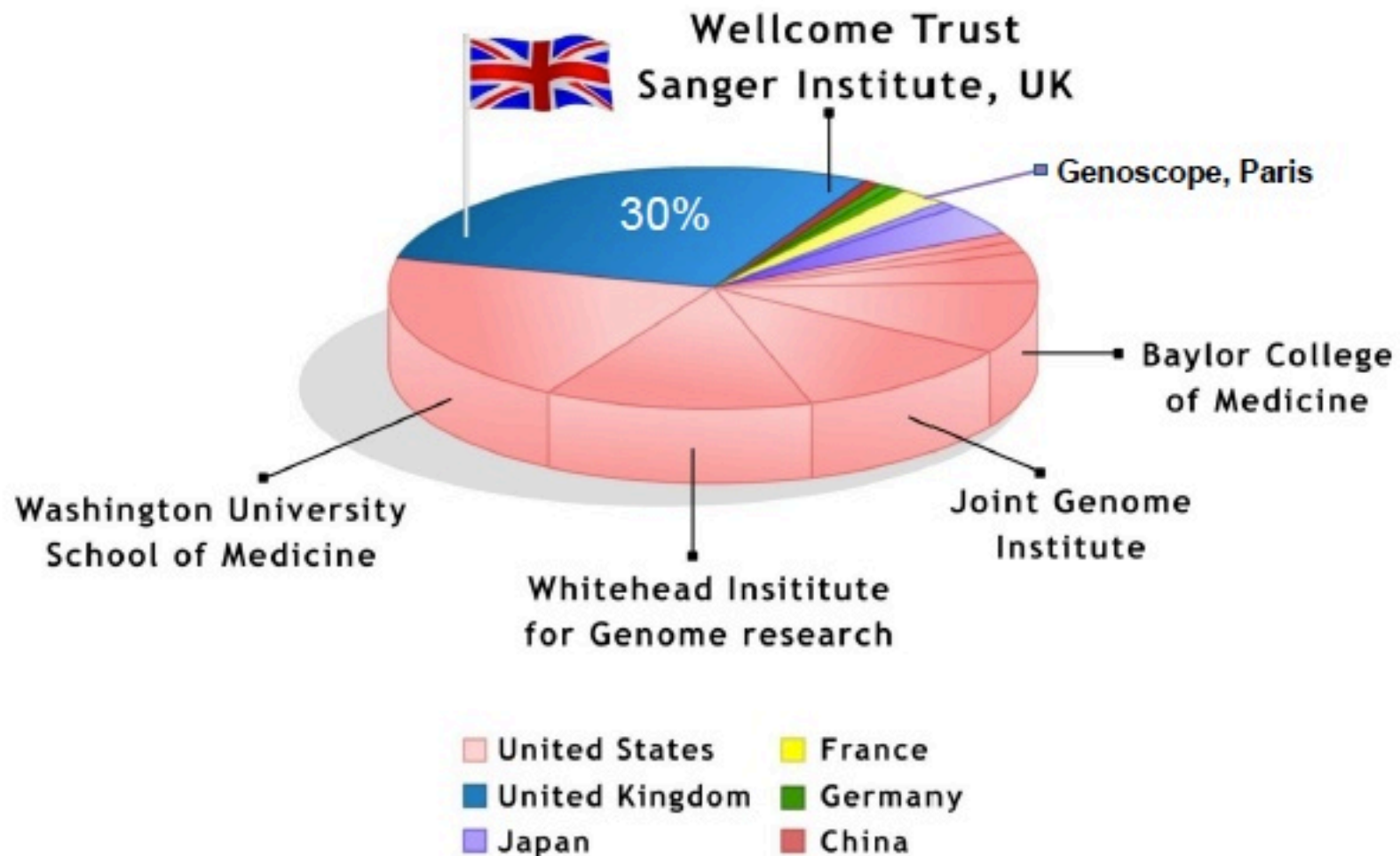


- A bit of background
- **Genomes**
- Genes
- Some bioinformatics basics

The Human Genome Project

- Possibility first raised in the mid 1980s
- 1990 a plan was created to begin work over 15 years
- Automatic release of draft data (1996 Bermuda Statement)
- Immediate submission of finished sequence

Contributors to finished human sequence



C. elegans

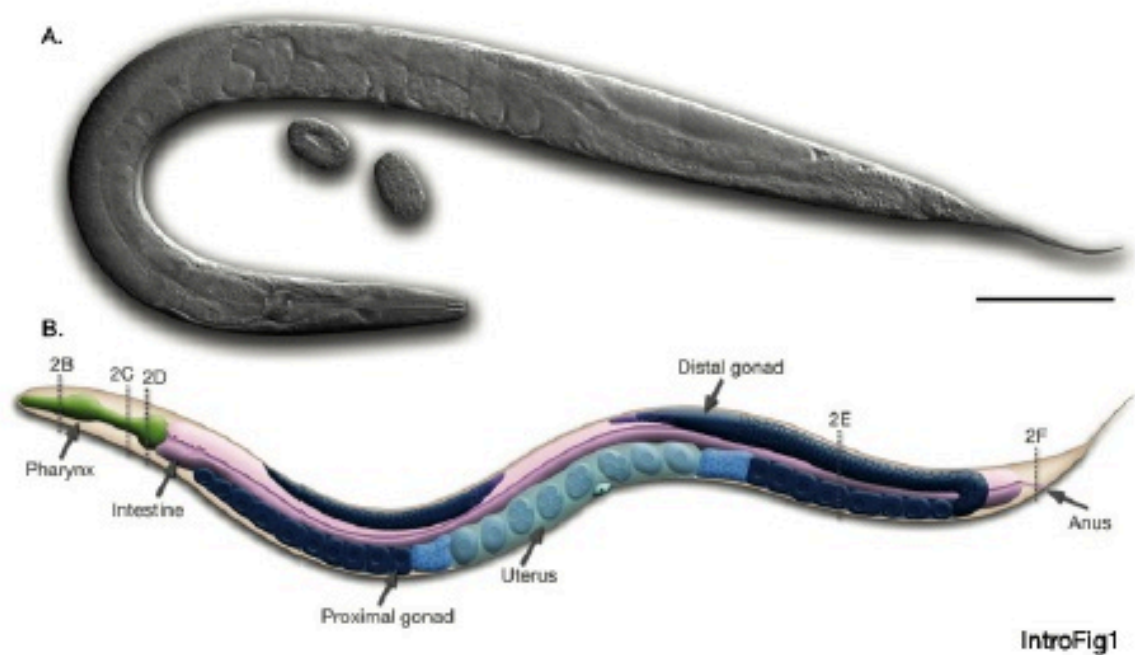
Number of cells – 959

100Mb genome

5 pairs of
chromosomes

XX/XO sex
chromosomes

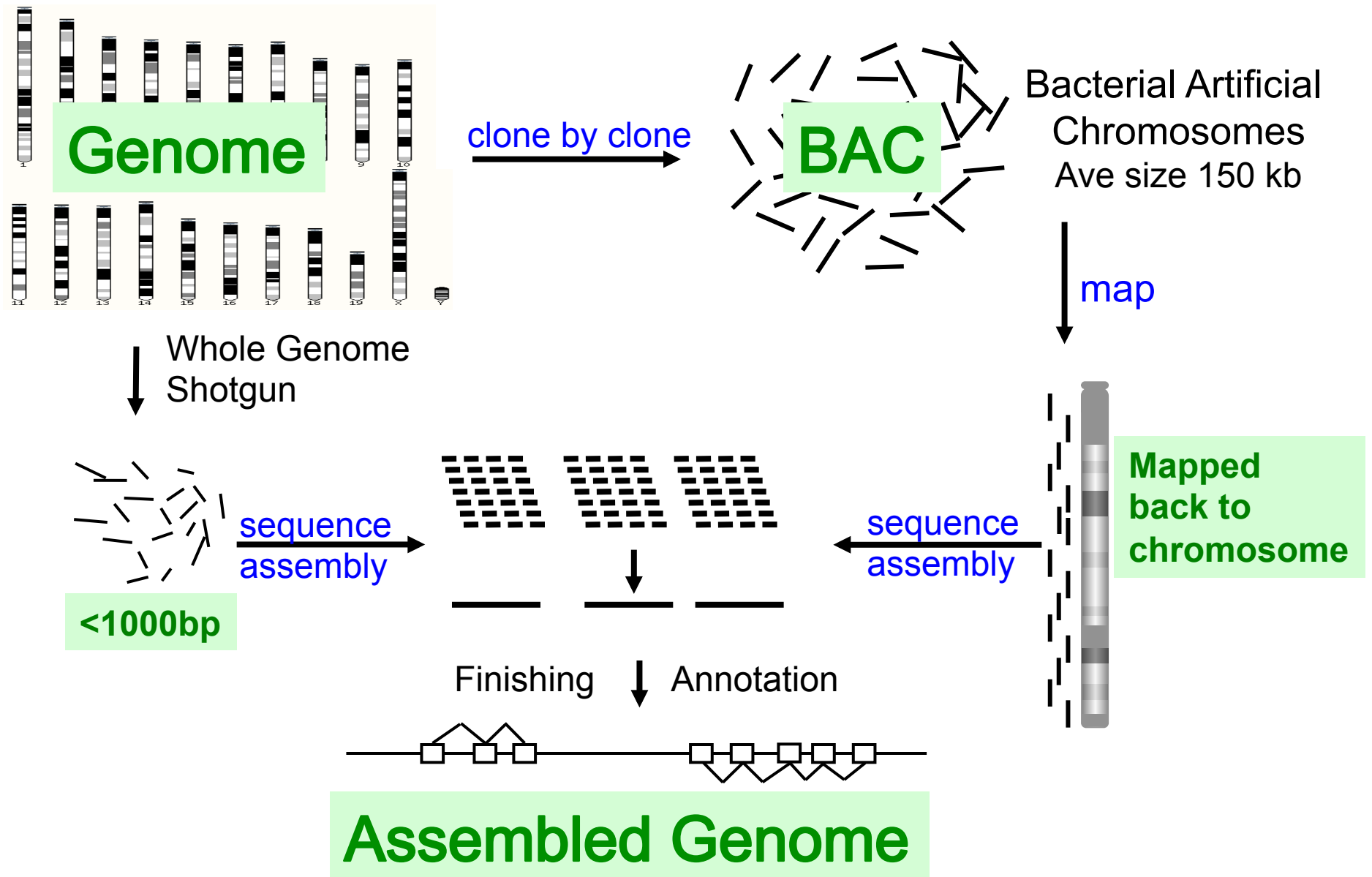
~20,500 genes
Published 1998



The race to the finish

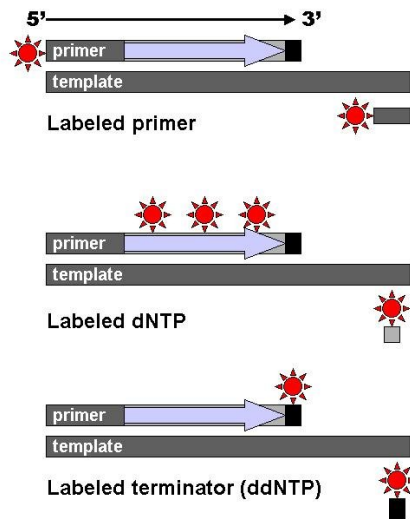
- 1998 Celera Genomics began private sector sequencing, led by J. Craig Venter using WGS
- 2000 President Clinton announced that the genome could not be patented and should be freely available to everyone
- The human genome working draft announced in 2000
- Publicly funded sequence published in 2001 in Nature
- Celera data was published in 2001 in Science
- Essentially complete genome in 2003
- 2006 the last chromosome published (chr 1)

Hybrid Sequencing Strategy



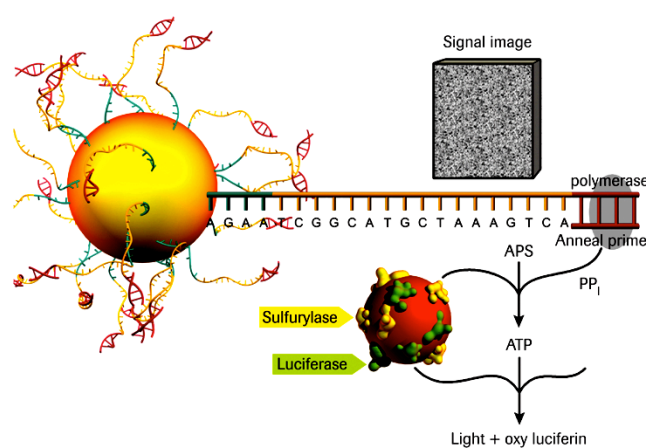
Sequencing chemistry changed in next gen sequencing technologies

Sanger sequencing



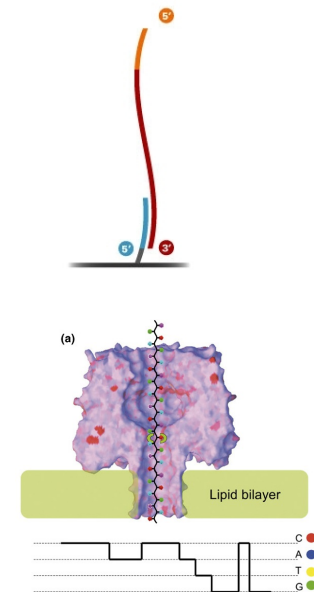
ABI 3730

Sequencing-by-synthesis



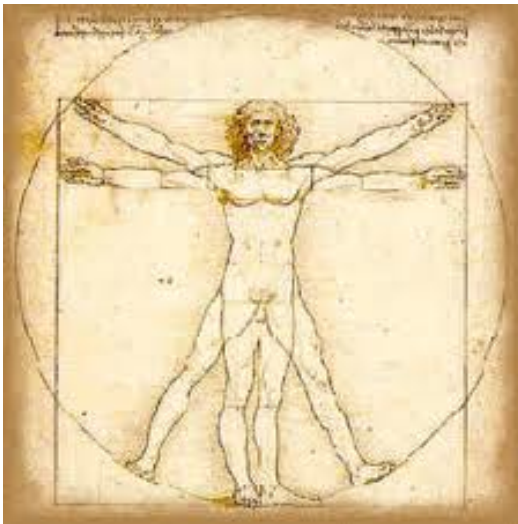
454, Illumina GAIIx and HiSeq2000

Single molecule sequencing



PacBio, Nanopore

Reference genomes:



Human ~3Gb:
22 chromosomes + sex
chromosomes
GRCh38



Mouse ~3 Gb:
19 chromosomes +
sex chromosomes
GRCm38

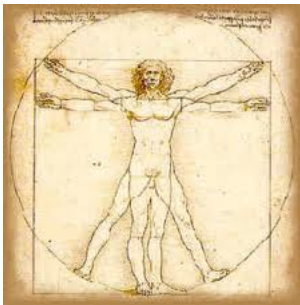


Zebrafish ~1.4 Gb:
25 chromosomes, no
specific sex
chromosomes
Zv10
(GRCz10 coming soon)



Do we know how many genes there are?

Protein coding genes



1980's 100,000

2000 40,000

Today ~ 20,000



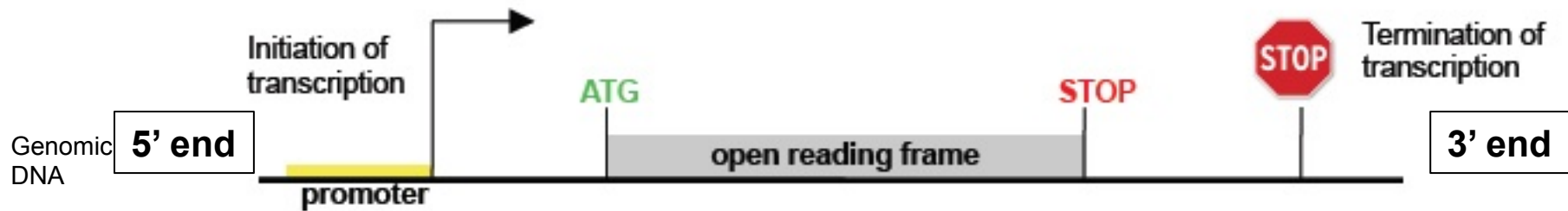
~22,000



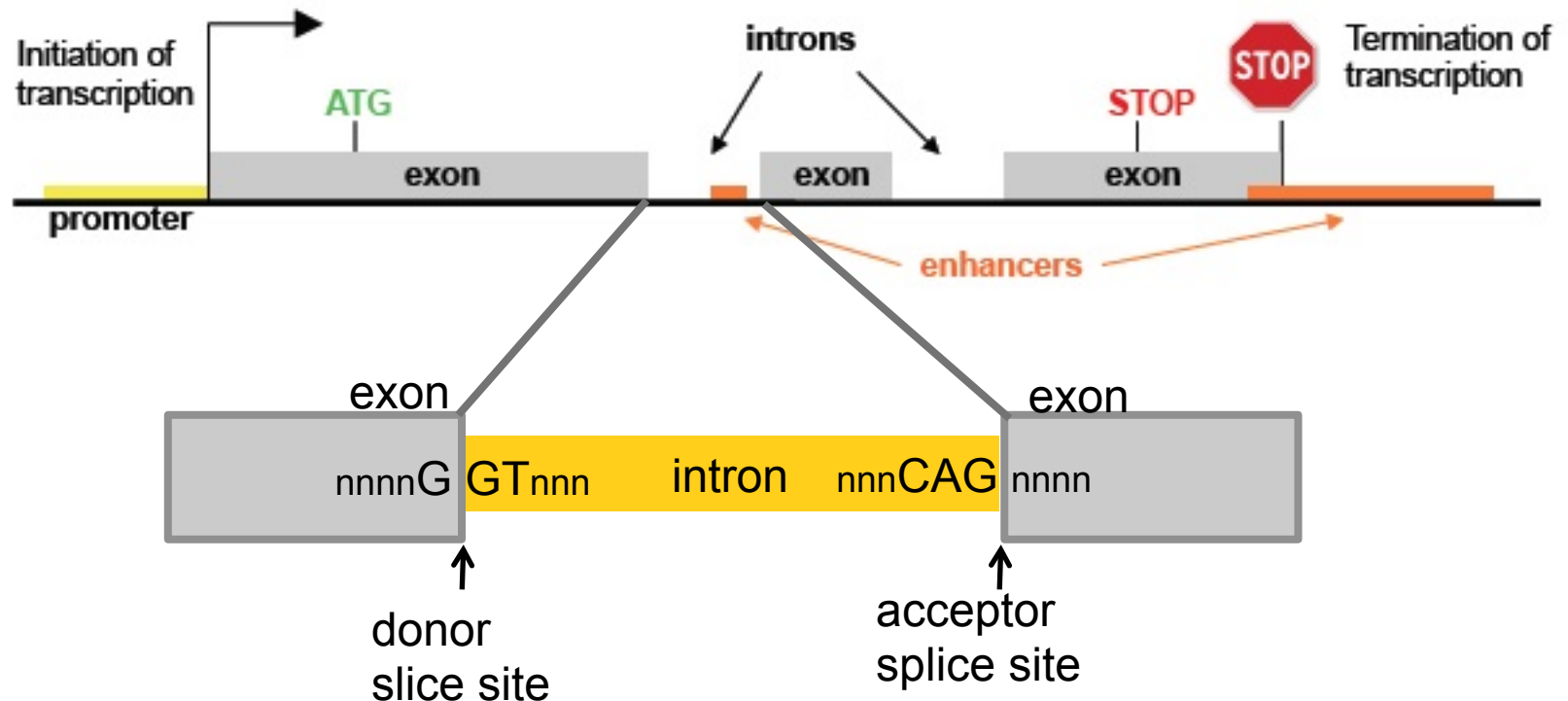
~26,000

- A bit of background
- Genomes
- **Genes**
- Some bioinformatics basics

Prokaryotes: Simple protein-coding gene

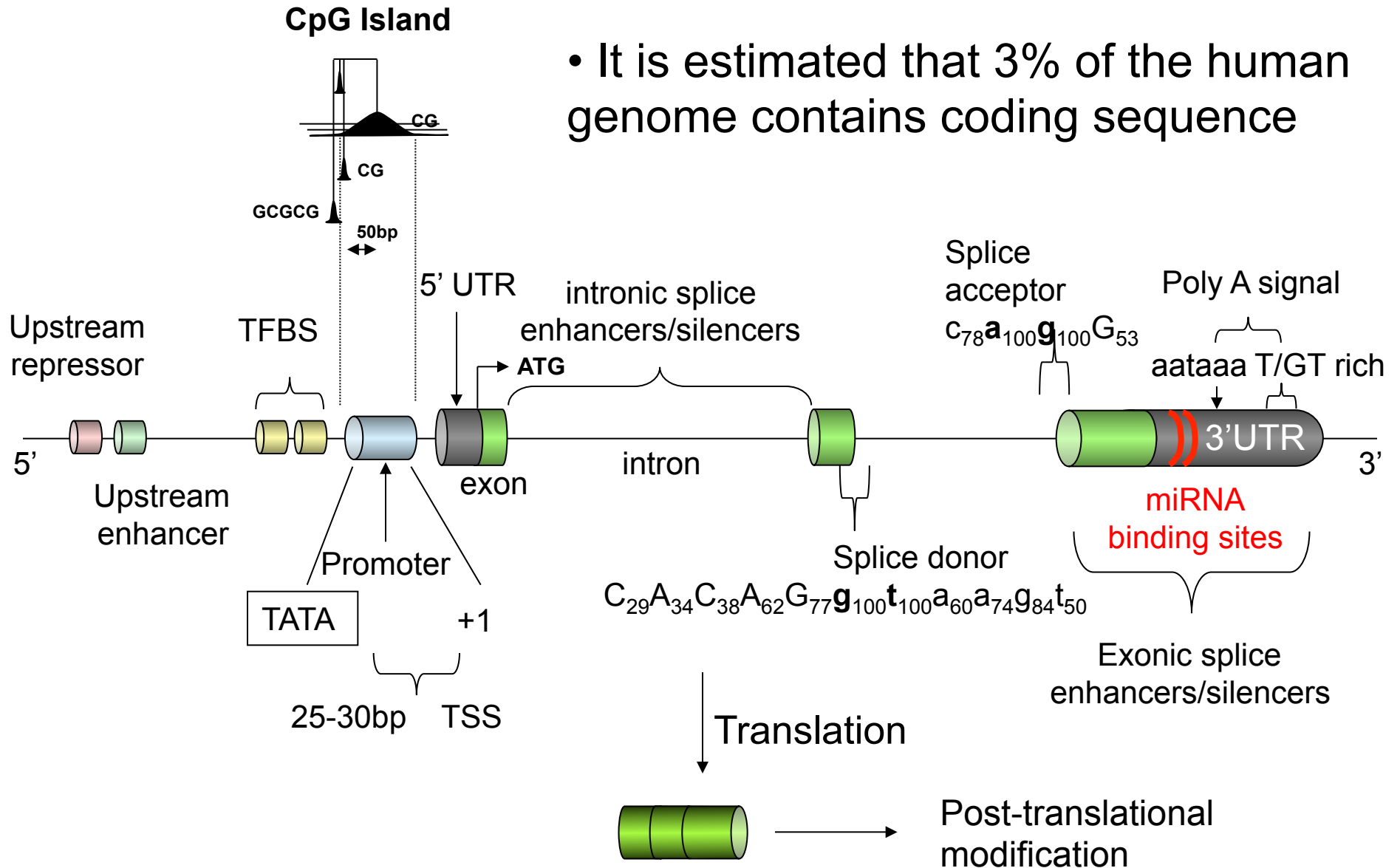


Eukaryotes: More complex: Introns and Exons

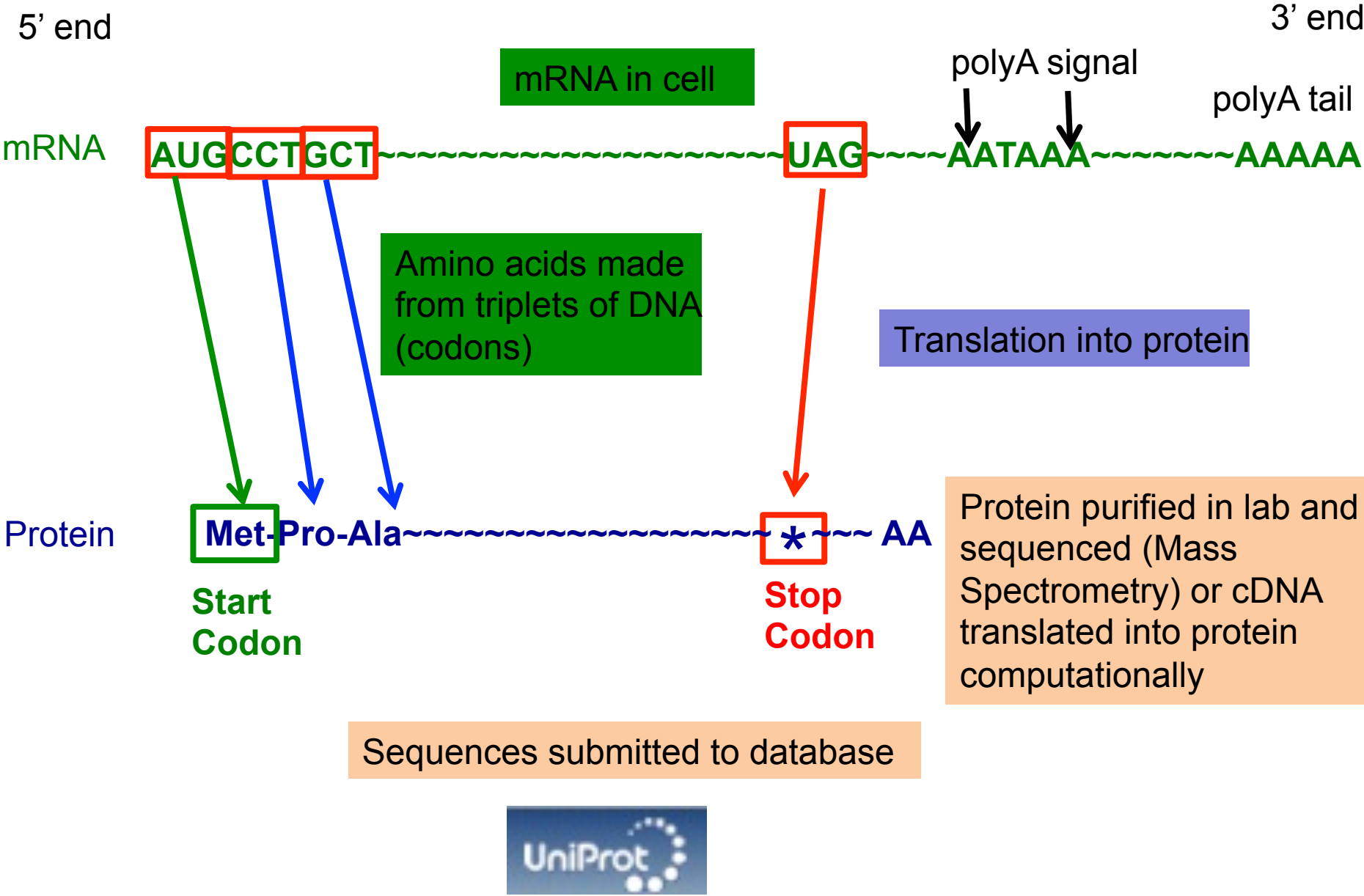


What is in a gene?

- It is estimated that 3% of the human genome contains coding sequence



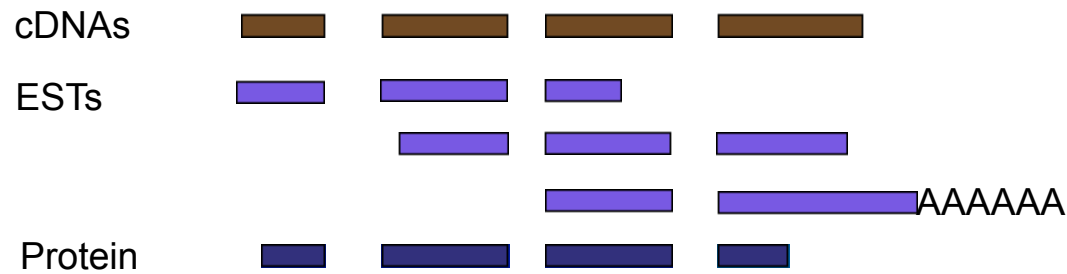
Evidence for genes: Protein



Making the transcript from evidence:

Genomic sequence 

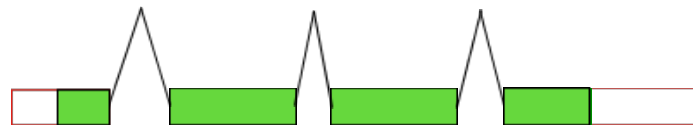
Analysis pipeline



Sequences from databases

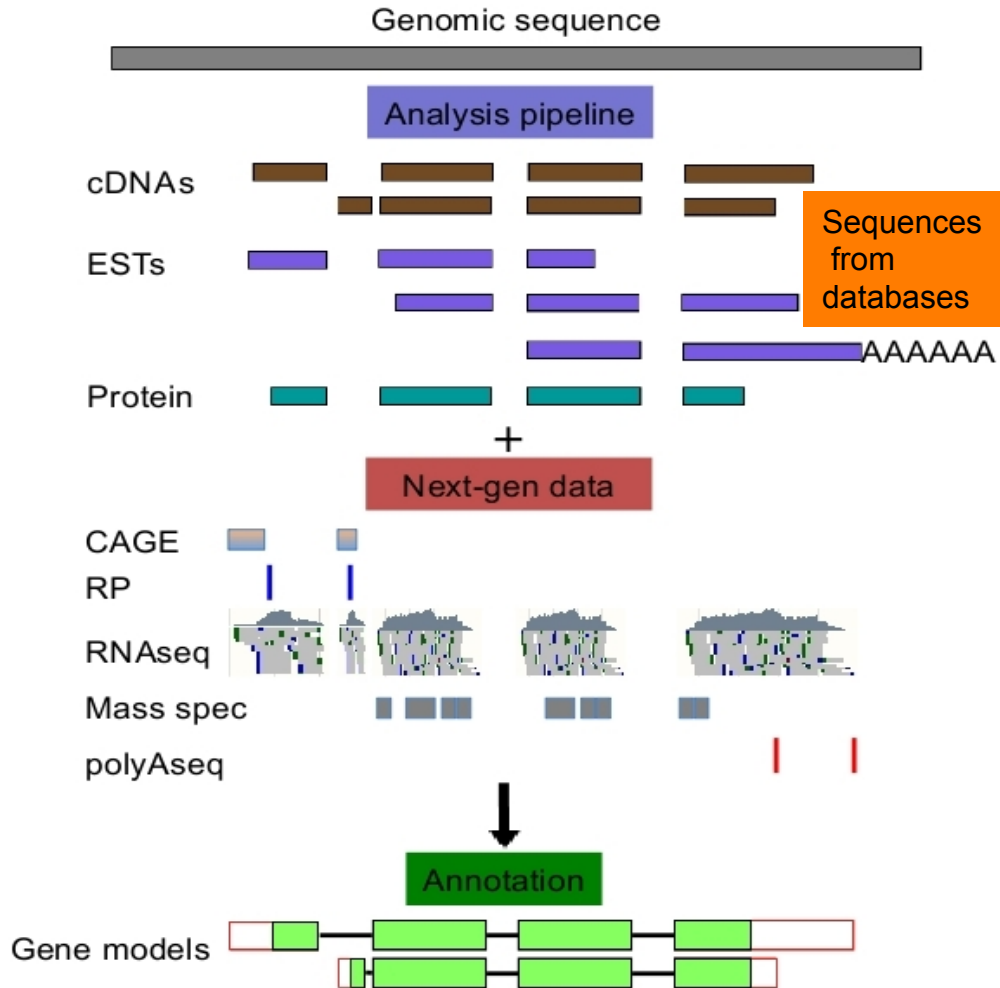
↓ Annotation

Gene structure



Manual Annotation and Biotypes:

Annotation: based on transcriptional evidence



Biotypes

Protein Coding

- Known_CDS
- Novel_CDS
- Putative_CDS
- Nonsense_mediated_decay

Transcript

- retained intron
- putative

Non-coding

- lincRNA
- Antisense
- Sense_intronic
- Sense_overlapping
- 3'_overlapping_ncRNA

Pseudogene

- Processed
- Unprocessed
- Transcribed
- Translated
- Unitary
- Polymorphic

Immunoglobulin

- IG_pseudogene
- IG_Gene
- TR_Gene

Alternative Splicing

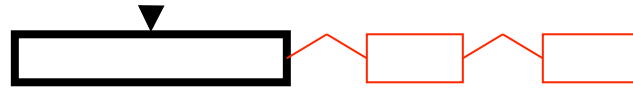
Reference model



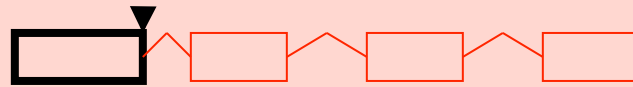
Skipped exon



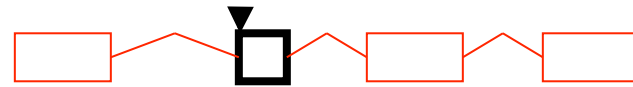
Retained intron



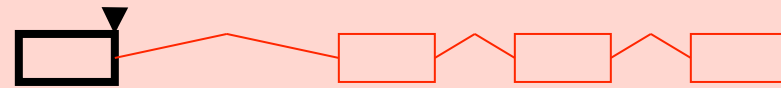
Alternative splice donor



Alternative splice acceptor



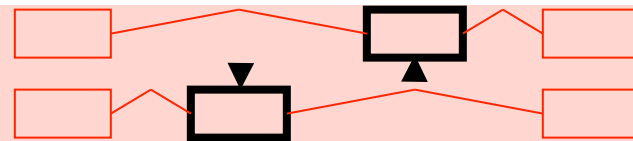
Alternative first exon



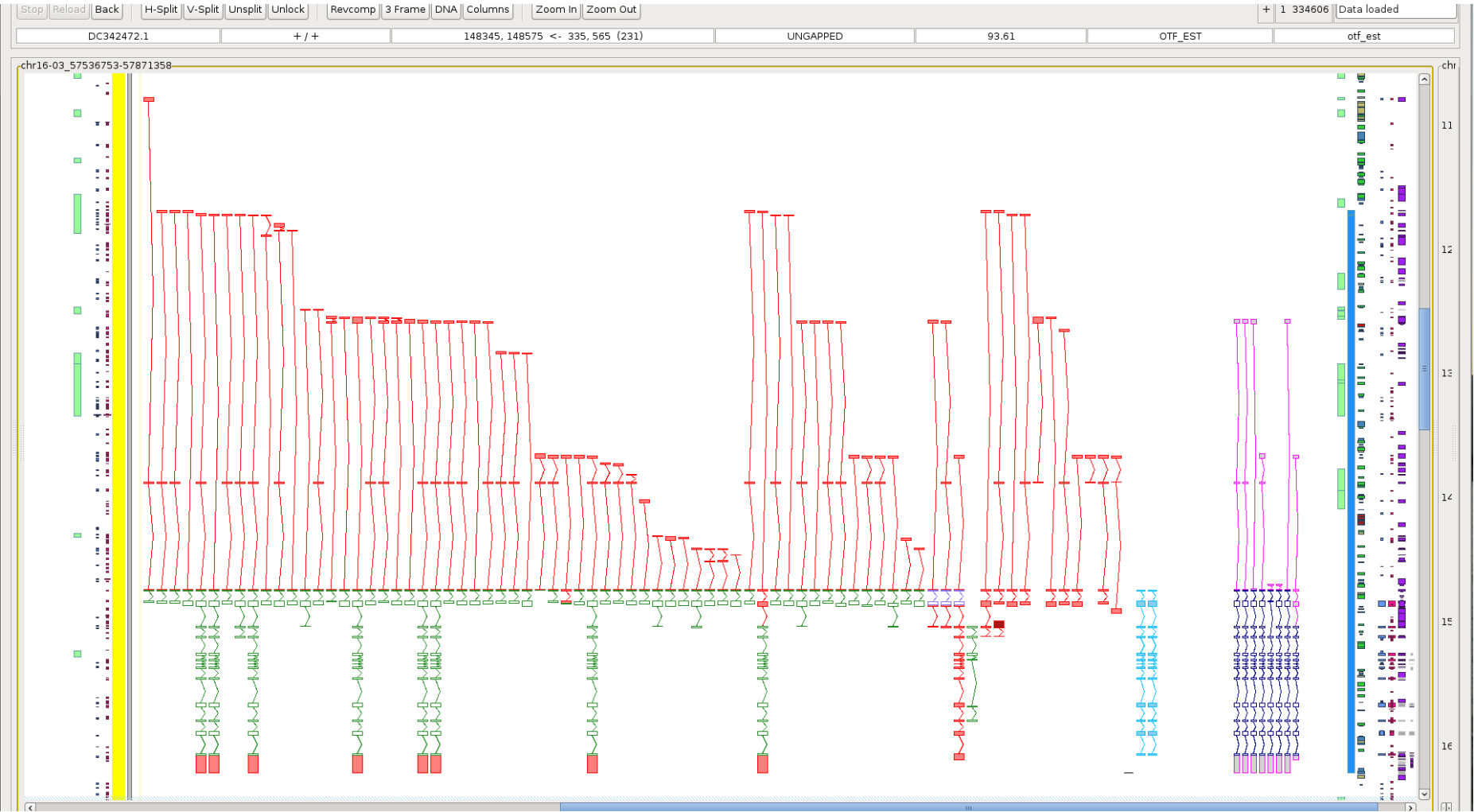
Alternative final exon



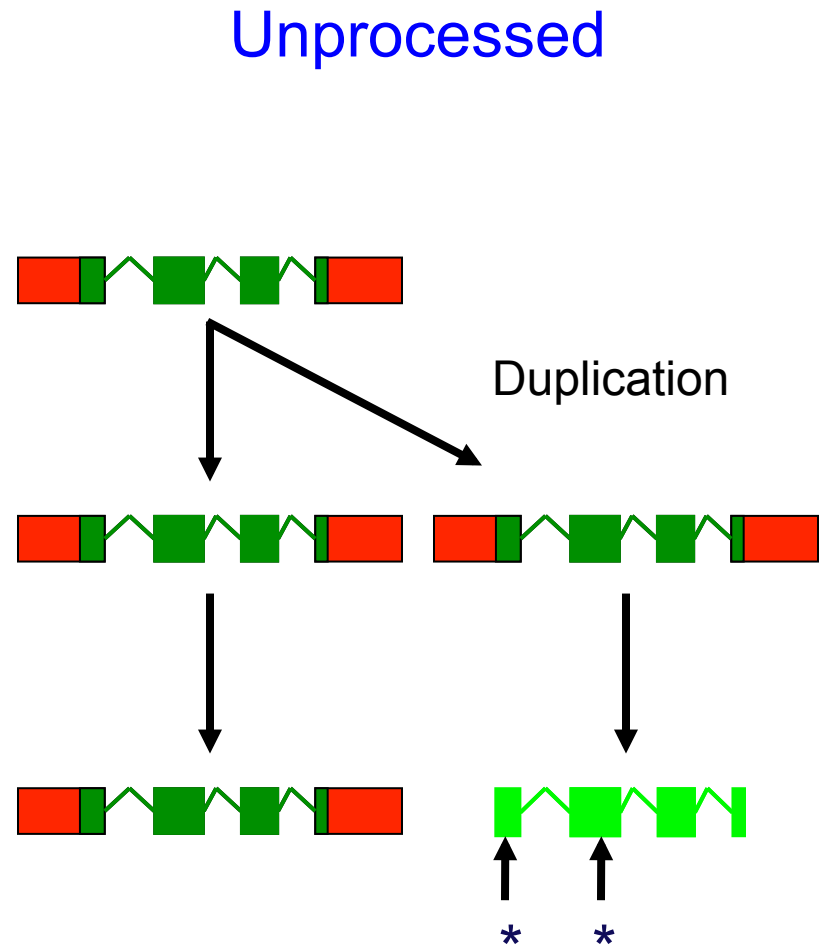
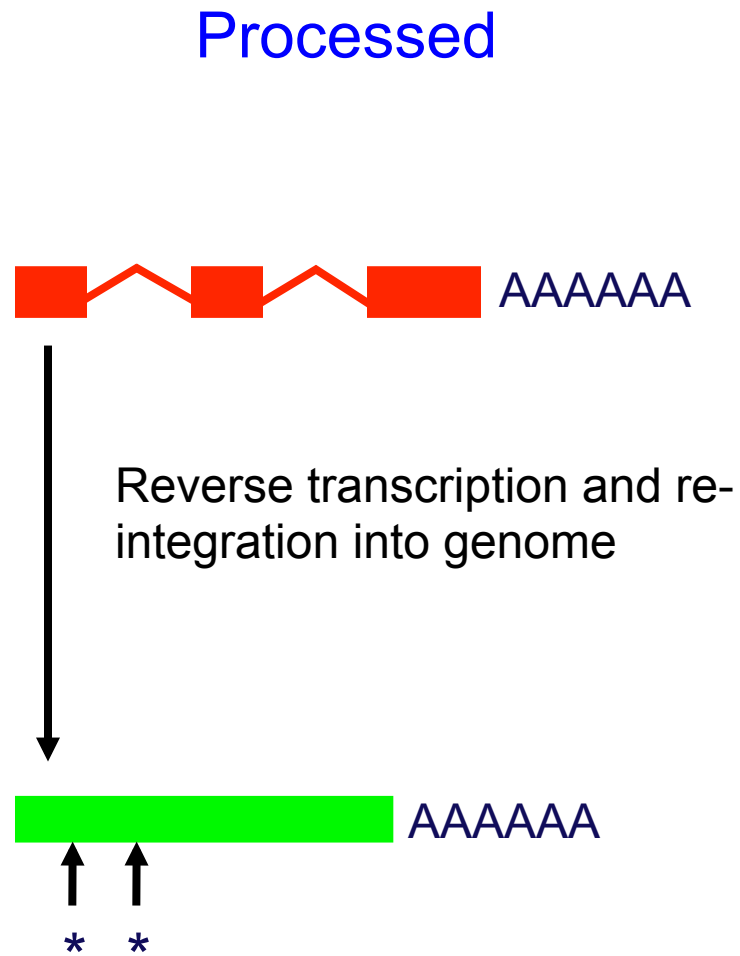
Mutually exclusive



GPR56: Human G protein-coupled receptor 56 gene



Havana pseudogenes



Disruption to coding sequence and stop codons

Non-coding genes:



microRNAs: Under 200 residues
Highly conserved signature
Imported from Rfam database

Family: *mir-30* (RF00131)

Description: *mir-30* microRNA precursor

```
Mus musculus (house mouse)  UGUAACAUCCCCGACUGGAAGCUGUAAGCCAC...AGCCAAGCUUCAGUCAGAUGUUUGCU
Spermophilus tridecemlineatus (th  UGUAACAUCCCCGACUGGAAGCUGUAGGACAC...AGCUGAGCUUCAGUCAGAUGUUUGCU
Macaca mulatta (Rhesus monkey)  UGUAACAUCCUA..CACUCAGCUGUAUACAU...GGAUUGGCUGGGAGGUGGAUGUUUACU
Macaca nemestrina (pig-tailed mac  UGUAACAUCCUA..CACUCAGCUGUAUACAU...GGAUUGGCUGGGAGGUGGAUGUUUACU
Macaca nemestrina (pig-tailed mac  UGUAACAUCCUA..CACUCAGCUGUAUACAU...GGAUUGGCUGGGAGGUGGAUGUUUACU
Gorilla gorilla (western gorilla)  UGUAACAUCCUA..CACUCAGCUGUAUACAU...GGAUUGGCUGGGAGGUGGAUGUUUACU
Homo sapiens (human)  UGUAACAUCCCCGACUGGAAGCUGUAAGACAC...AGCUAAGCUUCAGUCAGAUGUUUGCU
```

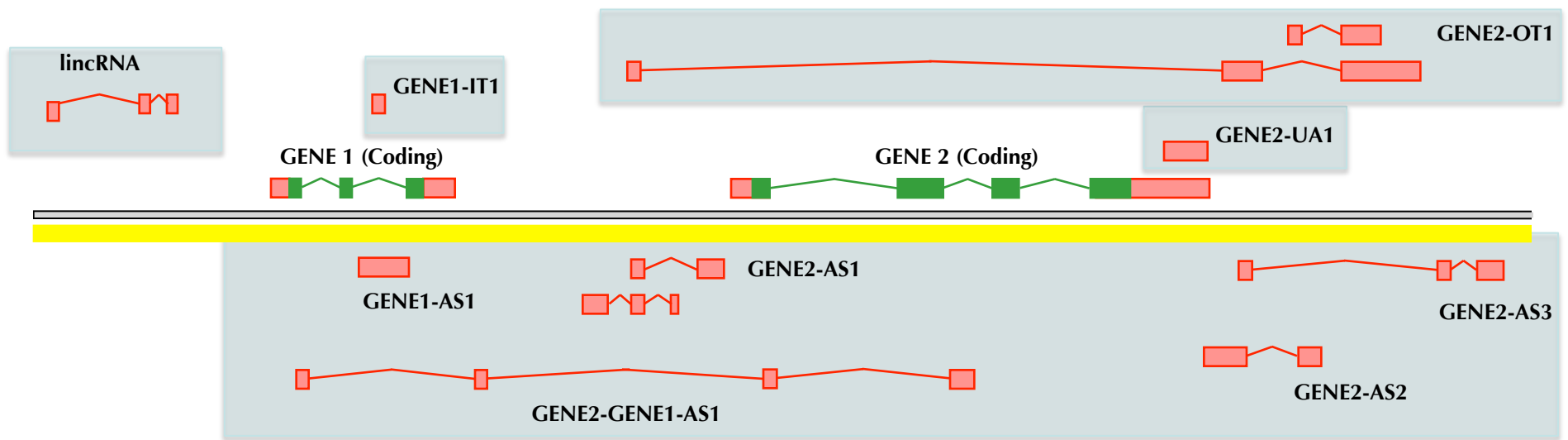
Long non-coding RNAs: (lncRNAs)

Over ~ 200 residues
Not highly conserved between species
Manually annotated
Some very well-known e.g. Hotair (HOX antisense intergenic RNA)
Many others not yet characterised



lncRNAs

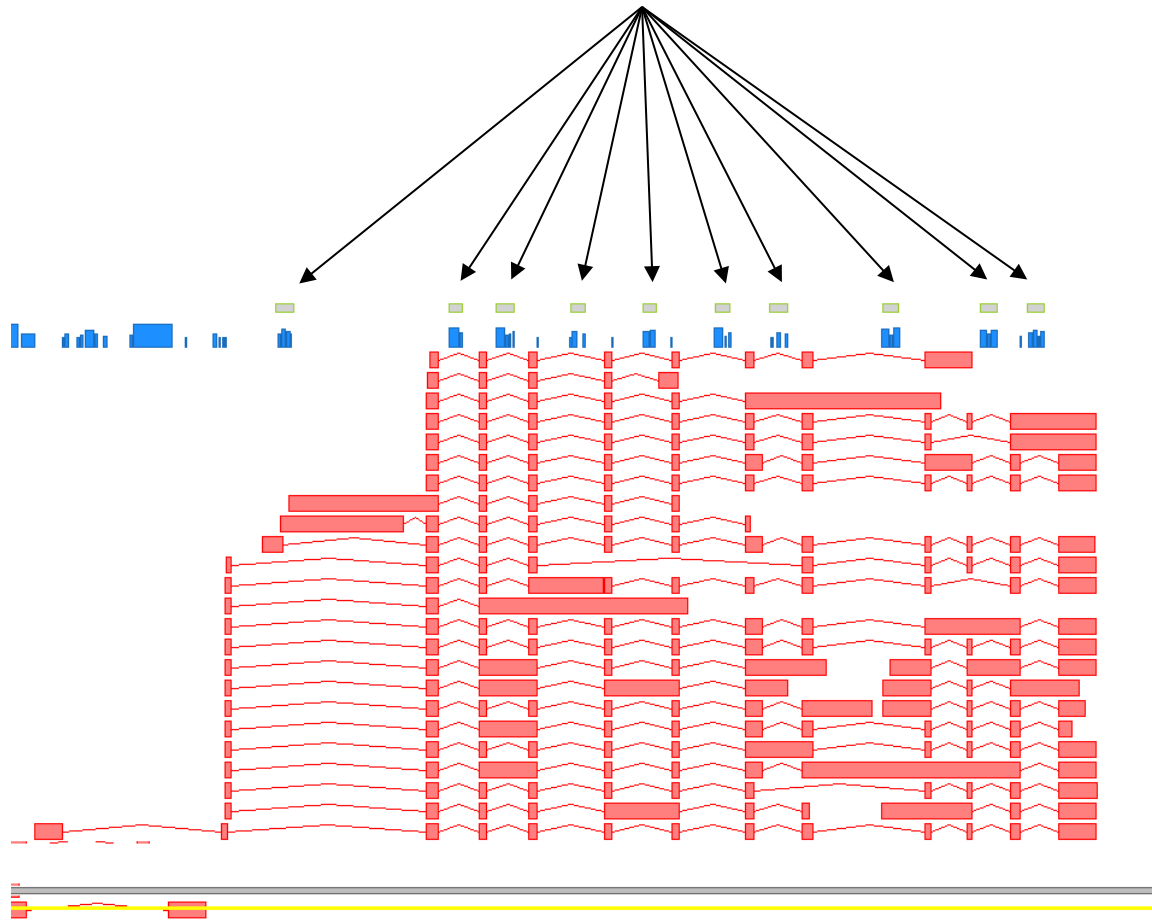
sense strand 5' →



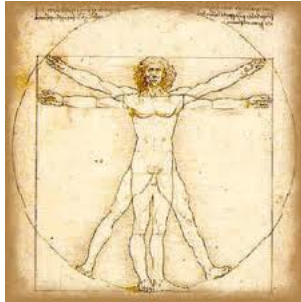
← 5' antisense strand

GAS5: growth arrest-specific 5 (non-protein coding)

intronic snoRNAs



Do we know how many genes there are?



Coding	Non-coding	Pseudogene
20,300	24,885	14,424

22,606	11,662	8,015
--------	--------	-------



26,459	7,014	264
--------	-------	-----

Who provides gene sets?



<http://www.ncbi.nlm.nih.gov/RefSeq/>

The logo for the University of California, Santa Cruz (UCSC) Genome Browser, with the text 'UCSC' in a gold serif font on a yellow background.

<http://genome.ucsc.edu/>



<http://www.ensembl.org>



<http://vega.sanger.ac.uk>



<http://www.ncbi.nlm.nih.gov/projects/CCDS/CcidsBrowse.cgi> (CDS only)

Differences in the gene sets

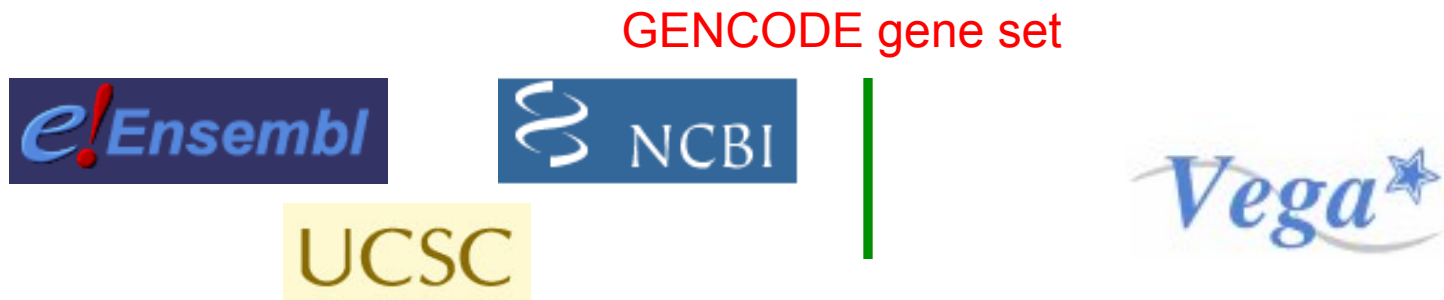
Automatic annotation

- Fast
- Unfinished sequence or shotgun sequence
- Consistent
- Under/Over-prediction
- Limited functional annotation
- Predicts ~75% loci

Manual annotation

- Slow
- Prefer finished sequence
- Flexible – can deal with inconsistencies
- Consult publications
- Extensive biotypes
- Excellent functional annotation

Automated annotation alone is not sufficient for researchers needs





NCBI - RefSeq

- Non-redundant gene set
- Accessed via browsers or Entrez Gene
- Accessions for genomic DNA, transcripts and proteins
- Primarily protein-coding
- Semi-curated

	Automated	Curated
Genomic	NC_12345	
mRNA	XM_12345	NM_12345
ncRNA	XR_12345	NR_12345
Protein	XP_12345	NP_12345

UCSC

UCSC gene set

- Non-redundant gene set
- Automatic annotation based on BLAT alignments
- Transcripts require Genbank accession plus one other supporting feature (eg. Uniprot)
- Includes RefSeq models (require no additional support)
- Both protein-coding and non-coding
- Data hub for ENCODE data, displays GENCODE geneset (for human and mouse)

<http://genome.ucsc.edu/>



Ensembl gene set

- Multiple biotypes (Known, Novel, ESTgenes, Pseudogenes)
- Automatic annotation based on pair-wise alignment using exonerate
- Transcripts require supporting mRNA and protein evidence
- Both protein-coding and non-coding transcripts
- Includes merged data from [VEGA](#) and [CCDS](#) and is called the GENCODE geneset



VEGA gene set

- Manually annotated using Otterlace/ZMAP annotation software
- Based on direct pairwise alignment of mRNA, EST and protein evidence (including cross-species)
- Multiple biotypes, reflect confidence levels
- Includes additional data sources as DAS tracks (eg. CAGE tags, RNAseq)



- Consensus CoDing Sequence project

The Consensus CDS (CCDS) project is a collaborative effort to identify a core set of human and mouse protein coding regions that are consistently annotated and of high quality. The long term goal is to support convergence towards a standard set of gene annotations.

- Havana, Ensembl, RefSeq, HGNC and MGI
- Produce reference CDS: set ATG-STOP on human and mouse genome - must agree. No UTRs.

CCDS website: GATA3 gene ATG->STOP

NCBI Consensus CDS protein set **CCDS Database** EBI • NCBI • UCSC • WTSI

PubMed Entrez Gene BLAST OMIM

Search for in and

Report for CCDS ID CCDS15674.1

CCDS	Status	Species	Chrom.	Gene	NCBI Builds	Links
15674.1	Public	<i>Mus musculus</i>	2	Gata3	36.1 - 37.1	H G G

Sequence IDs included in CCDS 15674.1

Original	Current	Source	Nucleotide ID	Protein ID	Status in CCDS	Seq. Status	Links
✓	✓	EBI,WTSI	ENSMUST00000102976	ENSMUSP00000100041	Accepted	alive	N P N P
✓	✓	EBI,WTSI	OTTMUST00000026063	OTTMUSP00000011932	Accepted	alive	N P N P
✓		NCBI	NM_008091.2	NP_032117.1	Updated	not alive	N P N P
	✓	NCBI	NM_008091.3	NP_032117.1	Accepted	alive	N P N P B

Chromosomal Locations for CCDS 15674.1

On '-' strand of Chromosome 2 (NC_000068.6)

Genome Browser links: [N](#) [U](#) [E](#) [V](#)

Chromosome	Start	Stop	Links
2	9779997	9780281	N U E V
2	9784722	9784847	N U E V
2	9790388	9790533	N U E V
2	9796016	9796552	N U E V
2	9798979	9799216	N U E V

CCDS
Home
FTP
Process
Statistics
AUG-guidelines

Collaborators
EBI
NCBI
UCSC
WTSI

Contact Us
GenComp eMail

Genome Displays
[E](#) Ensembl
[U](#) Genome Browser
[N](#) Map Viewer
[V](#) VEGA

Related Resources
Entrez Gene
HomoloGene
RefSeq
UniGene

CCDS Sequence Data

Blue highlighting indicates alternate exons.

Red highlighting indicates amino acids encoded across a splice junction.

Mouse over the nucleotide or protein sequence below and click on the highlighted codon or residue to select the pair.

Nucleotide Sequence (1332 nt):

ATGGAGGTGACTGCGGACCAGCCGGCTGGGTGAGCCACCATCACCCCGGGTCTCAACGGTCAGCACC
CAGACACGCACCCACCCGGCCCTCGGCCATTCTGACATGGAAGCTCAGTATCCGCTGACGGAAGAGGTGGA
CGTACTTTTAAACATCGATGGTCAAGGCAACCAGTCCCGTCTACTACGGAACCTCCGTGAGGGTACG
GTGCAGAGGTATCTCCGACCCACCACGGGAGCCAGGTATGCCGCCCGCTCTGCTGCACGGATCTCTGC
CCTGGCTGGATGGCGGCAAGCCCTGAGCAGCCACCACCCGCTCGCCCTGGGAACCTCAGCCCTTCTC
CAAGACGTCCATCCACCACGGCTCTCCGGGCTCTGTCGGTTTACCCTCCGGCTTCATCTCTTCTCTG
GCGGCCGCCACTCCAGTCTCATCTTTCACCTTCCCGCCACCCCGCGAAAGACGTCTCCCGAGACC
CGTCGCTGCCACCCCGGATCCCGGGTCCGCGCAGGCAAGATGAGAAAGAGTGCCTCAAGTATCAGT
GCAGCTGCCAGATAGCATGAAGCTGGAGACGTCTACTCTCGAGGAGCATGACCACCTGGGTGGGGCC
TCATCTCAGCCACCCACCTTACCACCTATCCGCCCTATGTGCCGAGTACAGCTCTGGACTCTTCC
CACCCAGCAGCTGTGGGAGGATCCCTACCGGTTCCGATGTAAGTCGAGGCCCAAGGCACGATCCAG
CACAGAAGCCAGGAGTGTGTGAATCGCGGCAACCTTACCCACTGTGGCGGAGATGGTACCAGG
CACTACCTTTGCAATGCCTGCGGACTCTACCATAAAATGAATGGGCAGAACCCGCCCTTATCAAGCCCA
AGCGAAGGCTGTCCGGCAGCAAGGAGACAGGACATCCTGCGGAACTGTGAGACCACCACCACCCT
CTGGAGGAGAACGCTAATGGGACCCCGTCTGCAATGCCTGTGGCTGTACTACAAGCTTCATAATATT
AACAGACCCCTGACTATGAAGAAGAAGGCATCCAGACCCGAAACCGGAAGATGTCTAGCAAAATCGAAAA
AGTGCAAAAAGGTGCATGACGCGCTGGAGGACTTCCCAAGAGCAGCTCTTCAACCCCGCCGCTCTCTC
CAGACACATGTATCCCTGAGCCACATCTCTCCCTCAGCCACTCCAGCCACATGCTGACCACCCGACG
CCCATGCATCCGCCCTCCGGCTCTCTTCCGACCTCACACCCTTCCAGCATGGTACCAGCCATGGGTT
AG

Translation (443 aa):

MEVTADQPRWVSHHHPAVLNGQHPDTHHPGLGHSYMEAQYPLTEEVDVLFNIDGQGNHVPYSYGNVSRAT
VQRYPPTHHGSQVCRPPLLHGSLPWLDDGKALSSHHTASPNWLSPPSKTSHHGSFGLPSVYPPASSSSL
AAGHSSPHLFTFPPTPPKDVSPDPSLSTPGSAGSARQDEKECLKYQVLPDSMKLETSHSRGSMTTLGGAA
SSSAHHPITTYPPYVPEYSSGLFPPSSLLGGSPFTGFCRSRPKARSSTEGRECVNCGATSTPLWRDRDGT
HYLCNACGLYHKMNGQNRPLIKPKRRLSAARRAGTSCANCQTTTTTLWRRNANGDPVNCAGGLYKLNHI
NRPLTMKKEGIQTRNRKMSKSKKCKVHDALEDLPKSSSFNPAALSRHMSLSHISPFSSHSHMLTPTT
PMHPPSGLSFGPHHPSMVTAMG

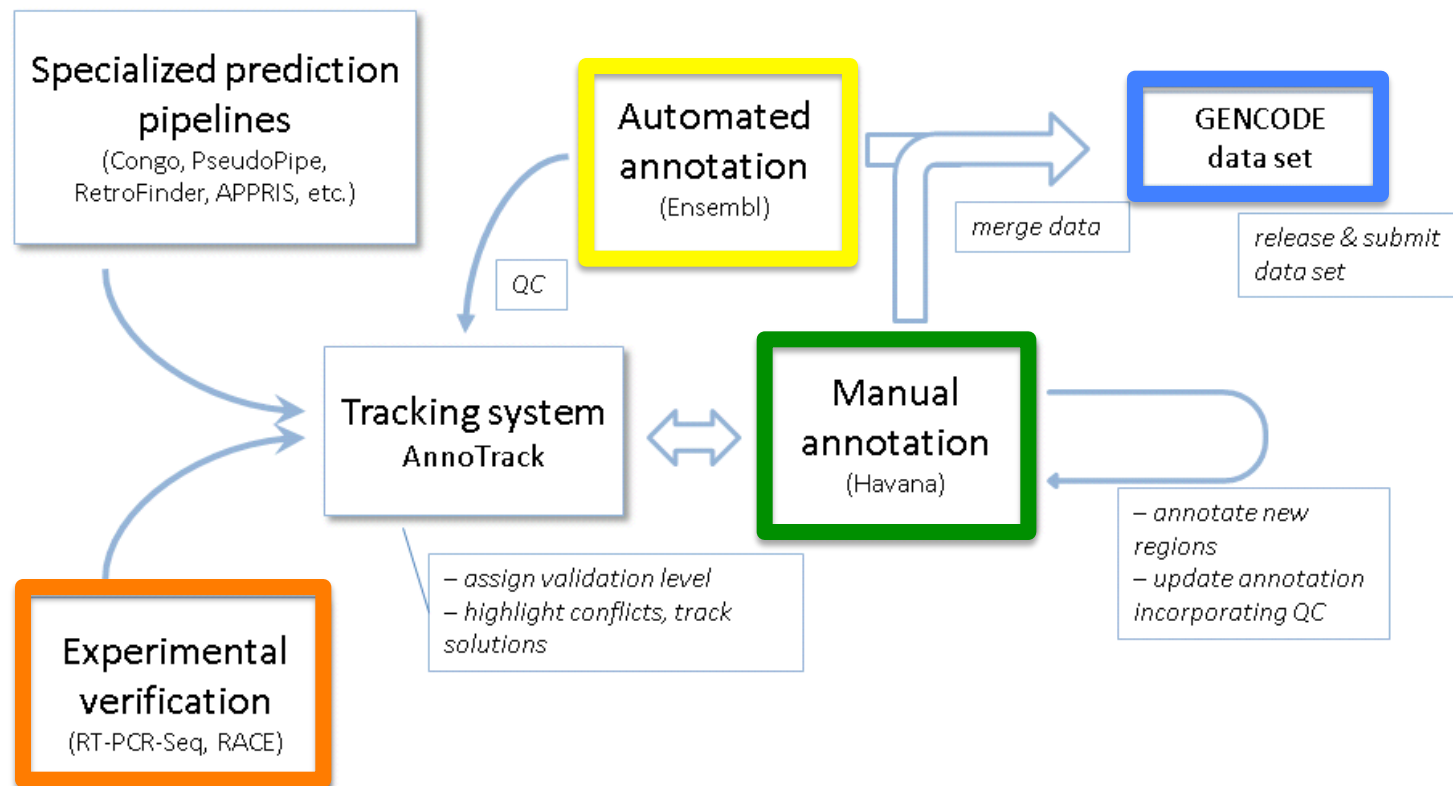
Comparing Gene Sets

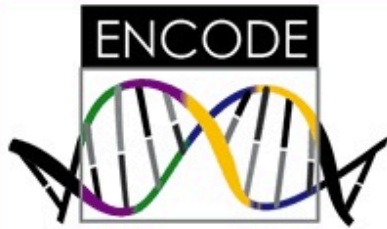
How many genes are there?

	Human	Mouse
RefSeq (2013)		
Coding genes	19,119	20,553
Total transcripts	35,539	27,113
UCSC		
Coding genes	21,520	21,181
Total transcripts	82,960	55,121
Ensembl (Gencode)		
Known coding	20,805 (e75)	23,871
Total transcripts	196,501	94,647
VEGA (unfinished)		
Coding genes	19,520 (v59)	16,359
Total transcripts	181,669	73,869
CCDS	18,800 (v106) (30,499)	20,080 (v104) (23,880)

GENCODE pipeline

Aim of GENCODE :annotate all evidence based gene features in the human genome





Project
Phase 2 GENCODE Goals
Data
Statistics - Human
Statistics - Mouse
Genome Browser - Human
Genome Browser - Mouse
Participants
Publications
lncRNA microarray
RGASP 1/2
RGASP 3
Blog
GENCODE workshops
Contact us

The GENCODE Project:

Encyclopædia of genes and gene variants

Current GENCODE version

The current version in **Human** is **Gencode 21**, released on the 2nd October 2014.

For more information about the human releases please see the [README.txt](#)  file.

The current version in **Mouse** is **Gencode M4**, released on the 3rd December 2014.

Introduction

The National Human Genome Research Institute (NHGRI) launched a public research consortium named **ENCODE** , the Encyclopedia Of DNA Elements, in September 2003, to carry out a project to identify all functional elements in the human genome sequence. After a successful pilot phase on 1% of the genome, the scale-up to the entire genome is now underway. The Wellcome Trust Sanger Institute was **awarded a grant**  to carry out a scale-up of the GENCODE project for integrated annotation of gene features.


Having been involved in successfully delivering the definitive annotation of functional elements in the human genome, the GENCODE group were **awarded a second grant** in 2013 in order to continue their human genome annotation work and expand GENCODE to include annotation of the mouse genome.

The **international team** working in the GENCODE project is headed by **Tim Hubbard**  at the **Wellcome Trust Sanger Institute** , and includes members from **Centre de Regulació Genòmica** , **Spanish National Cancer Research Centre** , **The University of Lausanne** , **Massachusetts Institute of Technology** , **Yale University**  and **The University of California, Santa Cruz** .

The GENCODE gene sets are used by the entire ENCODE consortium and by many other projects (eg. 1000 Genomes) as reference gene sets.

Acknowledgements

The GENCODE project is funded through an NHGRI ENCODE grant with additional funding from the Wellcome Trust.

When referencing, please use "Harrow J, et al. (2012) GENCODE: The reference human genome annotation for The ENCODE Project" ([PubMed](#) ).

Genome Reference Consortium

Goal:

- Correct regions in the genome that are currently misrepresented
- To close as many gaps as possible
- To produce alternative assemblies of structurally variant loci where necessary
- Scientific community can report loci in need of review
- Human, mouse and zebrafish

<http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/>



Search: for

gEVAL: Genome Evaluation Browser

The **gEVAL Browser** allows the evaluation of genome assemblies through its tools and pre-computed analyses.

The strength of this browser is the ability to navigate an up to date assembly and identify problematic regions and assist in strategizing potential solutions for these issues.

This facilitates the improvement of overall assemblies to a "gold" standard for release as reference genomes.

Commonly viewed genomes



Zebrafish_GRCz10
CURRENT



Mouse
GRCm38.p3



Human_GRCh38
CURRENT

Browse a Genome

Browse Human Genome Assemblies



Browse Mouse Genome Assemblies



Browse Zebrafish Genome Assemblies



Browse Pig Assemblies



Browse Rat Assemblies



Browse Parasitic Helminth Assemblies



Browse Chicken Assemblies



About the Project

gEVAL utilizes the Ensembl framework and is maintained by the *Genome Reference Informatics Team* at the **Wellcome Trust Sanger Institute**.

The team is part of the *Genome Reference Consortium (GRC)*, a multi-centre collaboration tasked with providing improved reference assemblies that better represent complex diversity.

- A bit of background
- Genomes
- Genes
- **Some bioinformatics basics**

Bioinformatics services

The
Part of

We maintain the world's most comprehensive range of **freely available** and up-to-date molecular databases. Developed in collaboration with our colleagues worldwide, our services let you share data, perform complex queries and analyse the results in different ways. You can work locally by downloading our data and software, or use our web services to access our resources programmatically.

Bioinfo



XXX DNA & RNA

genes, genomes & variation

Gene expression

RNA, protein & metabolite expression

Proteins

sequences, families & motifs

Structures

Molecular & cellular structures

Systems

reactions, interactions & pathways

Chemical biology

chemogenomics & metabolomics

Ontologies

taxonomies & controlled vocabularies

Literature

Scientific publications & patents

Other software

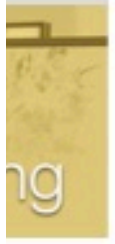
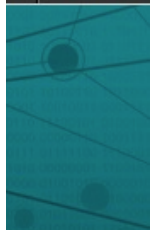
cross-domain tools & resources

Indust



Our m

- o To pro
- o To cor
- o To pro
- o To hel
- o To coc



S

Sequence: BN000065.1

TPA: Homo sapiens SMP1 gene, RHD gene and RHCE gene

[Send Feedback](#) 

View: [TEXT](#) [FASTA](#) [XML](#)

Download: [XML](#) [FASTA](#) [TEXT](#)

Organism Homo sapiens	Molecule type genomic DNA	Topology linear	Data class STD	Taxonomic Division HUM
Sequence length 315,242	Sequence Version 1	First public 23-APR-2002	Last updated 14-NOV-2006	Show Version History BN000065

Keywords

RHCE gene, RhCE protein, RHD gene, RhD protein, small membrane protein 1, SMP1 gene, Third Party Data, TPA, TPA:inferential.

Lineage

[Eukaryota](#), [Metazoa](#), [Chordata](#), [Craniata](#), [Vertebrata](#), [Euteleostomi](#), [Mammalia](#), [Eutheria](#), [Euarchontoglires](#), [Primates](#), [Haplorrhini](#), [Catarrhini](#), [Hominidae](#), [Homo](#)

Navigation

Overview

Source Feature(s)

Sequence

Assembly

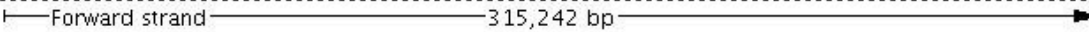
Other Feature(s)

Publications

Base range: -

Overview 

 BN000065.1

Features 

Assembly 





UniProtKB ▾

Advanced ▾



BLAST Align Retrieve/ID Mapping

Help Contact

The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information.

UniProtKB

Swiss-Prot
(547,357)

Manually annotated and reviewed.

TrEMBL
(89,451,166)

Automatically annotated and not reviewed.

UniRef

Sequence clusters



UniParc

Sequence archive



Proteomes



Supporting data

Literature citations



Taxonomy



Subcellular locations



Cross-ref. databases



Diseases

XXX

Keywords



News



Thalidomide, the pharmacological version of yin and yang | Cross-references to DEPOD, MoonProt and Proteomes

[UniProt release 2015_01](#)

Higher and higher | New mouse and zebrafish variation files | Structuring of 'cofactor' annotations

[UniProt release 2014_11](#)

[News archive](#)

Q9HD64 - XAGE1_HUMAN

Protein
Gene
Organism
Status

X antigen family member 1

XAGE1A [more](#)

Homo sapiens (Human)

Reviewed - - Experimental evidence at transcript levelⁱ

Display

None

[BLAST](#) [Align](#) [Format](#) [Add to basket](#) [History](#)

- FUNCTION
- NAMES & TAXONOMY
- SUBCELLULAR LOCATION
- PATHOLOGY & BIOTECH
- PTM / PROCESSING
- EXPRESSION
- INTERACTION
- STRUCTURE
- FAMILY & DOMAINS
- SEQUENCES (2)
- CROSS-REFERENCES
- PUBLICATIONS
- ENTRY INFORMATION
- MISCELLANEOUS
- SIMILAR PROTEINS

[▲ Top](#)

Names & Taxonomyⁱ

Protein names ⁱ	<p><i>Recommended name:</i></p> <p>X antigen family member 1</p> <ul style="list-style-type: none"> ▪ <i>Short name:</i> XAGE-1 <p><i>Alternative name(s):</i></p> <ul style="list-style-type: none"> • Cancer/testis antigen 12.1 <ul style="list-style-type: none"> ▪ <i>Short name:</i> CT12.1 • G antigen family D member 2
Gene names ⁱ	<p><i>Name:</i> XAGE1A Synonyms: GAGED2, XAGE1 AND <i>Name:</i> XAGE1B AND <i>Name:</i> XAGE1C AND <i>Name:</i> XAGE1D AND <i>Name:</i> XAGE1E</p>
Organism ⁱ	<i>Homo sapiens</i> (Human)
Taxonomic identifier ⁱ	9606 [NCBI]
Taxonomic lineage ⁱ	Eukaryota > Metazoa > Chordata > Craniata > Vertebrata > Euteleostomi > Mammalia > Eutheria > Euarchontoglires > Primates > Hominidae > Homo
Proteinomes ⁱ	UP000005640: Chromosome X



All D

How To

Search

Submissions

[BioProject Submission](#)

An online form that provide genomic and genetic data

[ClinVar Submissions](#)

Guidelines and instructions: level/aggregate data); sup

[Database of Genotype and](#)

Guidelines and requiremer

[Database of Major Histoc](#)

Guidelines and template fc

[GenBank: BankIt](#)

A web-based sequence su

[GenBank: Barcode](#)

Tool for submission to the

[GenBank: Sequin](#)

A stand-alone software to c simple submissions that cc sequences from phylogene

[GenBank: tbl2asn](#)

A command-line program t submission of complete ge

[Gene Expression Omnibus](#)

Submit expression data. si

- [Save text searches and set up automated searches with E-mailed results](#)
- [Find bioassays in which a given drug is active](#)
- [Find bioassays that test a particular disease or protein target](#)
- [Submit data to NCBI](#)
- [Download NCBI Software](#)
- [Submit sequence data to NCBI](#)
- [Retrieve all sequences for an organism or taxon](#)
- [Find the function of a gene or gene product](#)
- [View all SNPs associated with a gene](#)
- [Find genes associated with a phenotype or disease](#)
- [Find expression patterns](#)
- [Obtain genomic sequence for/near a gene, marker, transcript or protein](#)
- [Find human variations associated with a phenotype or disease \(clinical association\)](#)
- [Convert feature coordinates between genomic assemblies](#)
- [View/download features around an object or between two objects on a chromosome](#)
- [Compare protein homologs between two microbial genomes](#)
- [Find sequenced genomes, including those in progress, for a taxonomic group](#)
- [Download the complete genome for an organism](#)
- [Display genomic annotation graphically](#)
- [Determine conserved synteny between the genomes of two organisms](#)
- [Find a homolog for a gene in another organism](#)
- [Find articles about a topic similar to that in a given article](#)
- [Obtain the full text of an article](#)
- [View the 3D structure of a protein](#)
- [Find a curated version of a sequence record \(NCBI Reference Sequence\)](#)
- [Align two or more 3D structures to a given structure](#)
- [Find published information on a gene or sequence](#)
- [Find transcript sequences for a gene](#)
- [Link from an object on a map to another resource](#)
- [Run BLAST software on a local computer](#)

t for the submission of

data about a variant (variant

and easy.

n.

i. It is capable of handling ad sets of DNA, as well as

in. It is used primarily for

Display Settings: Full Report

Send to:

Hide sidebar >>

XAGE1B X antigen family, member 1B [*Homo sapiens* (human)]

Gene ID: 653220, updated on 7-Dec-2014

Summary

Official Symbol XAGE1B provided by HGNC
Official Full Name X antigen family, member 1B provided by HGNC
Primary source HGNC:HGNC:25400
Locus tag RP11-485B17.2
See related Ensembl:ENSG00000204379; MIM:300742; MIM:300743; Vega:OTTHUMG00000021557; Vega:OTTHUMG00000163577
Gene type protein coding
RefSeq status REVIEWED
Organism *Homo sapiens*
Lineage Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorhini; Catarhini; Hominoidea; Homo
Also known as CTP9; XAGE1; CT12.1; GAGED2; XAGE-1; XAGE1A; CT12.1A; CT12.1B
Summary This gene is a member of the XAGE subfamily, which belongs to the GAGE family. The GAGE genes are expressed in a variety of tumors and in some fetal and reproductive tissues. This gene is strongly expressed in Ewing's sarcoma, alveolar rhabdomyosarcoma and normal testis. The protein encoded by this gene contains a nuclear localization signal and shares a sequence similarity with other GAGE/PAGE proteins. Because of the expression pattern and the sequence similarity, this protein also belongs to a family of CT (cancer-testis) antigens. Alternative splicing of this gene, in addition to alternative transcription start sites, results in multiple transcript variants. [provided by RefSeq, Jan 2010]

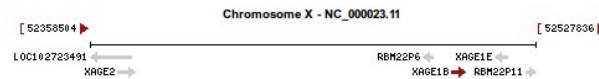
Genomic context

Location: Xp11.22

See XAGE1B in Epigenomics, MapViewer

Exon count: 4

Annotation release	Status	Assembly	Chr	Location
106	current	GRCh38 (GCF_000001405.26)	X	NC_000023.11 (52495668..52500812)
105	previous assembly	GRCh37.p13 (GCF_000001405.25)	X	NC_000023.10 (52238810..52243954)



Genomic regions, transcripts, and products

Genomic Sequence: NC_000023.11 chromosome X reference GRCh38 Primary Assembly

Go to reference sequence details

Go to nucleotide: Graphics FASTA GenBank

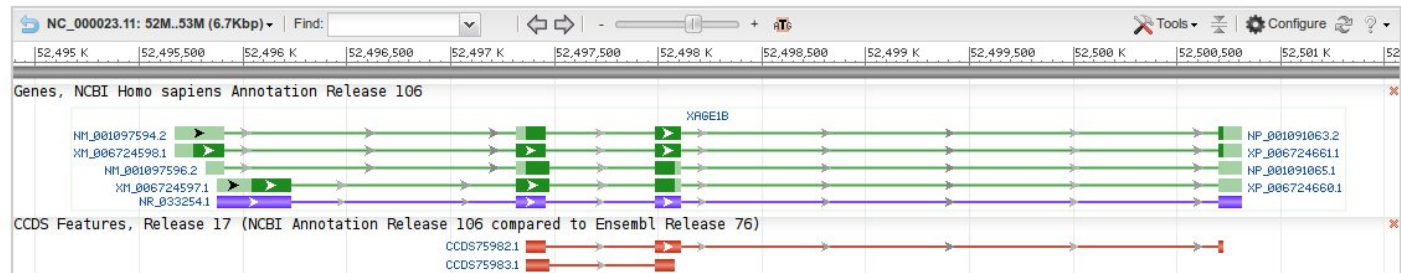


Table of contents

- Summary
- Genomic context
- Genomic regions, transcripts, and products
- Bibliography
- Phenotypes
- Variation
- General protein information
- NCBI Reference Sequences (RefSeq)
- Related sequences
- Additional links

Related information

- BioProjects
- CCDS
- ClinVar
- Conserved Domains
- dbVar
- Full text in PMC
- Full text in PMC_nucleotide
- Gene names
- Genome
- GEO Profiles
- HomoloGene
- Map Viewer
- Nucleotide
- OMIM
- Probe
- Protein
- PubChem Compound
- PubChem Substance
- PubMed
- PubMed (GeneRIF)
- PubMed (OMIM)
- PubMed(nucleotide/PMC)
- RefSeq Proteins
- RefSeq RNAs
- SNP
- SNP: GeneView
- Taxonomy
- UniGene
- Links to other resources
- HGNC
- Ensembl
- Vega

Searching the databases: How to find gene location

BLAST

(Basic Local Alignment Search Tool)



Genomic DNA as query
against the databases

```
ATGCTAGGATCCGATTGCAAG
CCTGAATCCGGCCTAATTTAC
G Pattern matching to CC
A millions of sequences AG
A in the databases AG
ATAGCAGATAGACAGTAAGAC
ATGATAGACGATAGATACAGA
```



```
>ref|NM\_000059.3| UEGMD Homo sapiens breast cancer 2, early onset (BRCA2), mRNA
Length=11386
```

```
GENE ID: 675 BRCA2 | breast cancer 2, early onset [Homo sapiens]
(Over 100 PubMed links)
```

```
Score = 348 bits (188), Expect = 5e-93
Identities = 188/188 (100%), Gaps = 0/188 (0%)
Strand=Plus/Plus
```

```
Query 7 GTGGCGCGAGCTTCTGAAACTAGGCGGCAGAGGCGGAGCCGCTGTGGCACTGCTGCGCCT 66
      |||
Sbjct 1 GTGGCGCGAGCTTCTGAAACTAGGCGGCAGAGGCGGAGCCGCTGTGGCACTGCTGCGCCT 60

Query 67 CTGCTGCGCCTCGGGTGTCTTTTGC CGGCGGTGGGTGCGCCGCGGAGAAAGCGTGAGGGGA 126
      |||
Sbjct 61 CTGCTGCGCCTCGGGTGTCTTTTGC CGGCGGTGGGTGCGCCGCGGAGAAAGCGTGAGGGGA 120

Query 127 CAGATTTGTGACCGGCGCGGTTTTTGT CAGCTTACTCCGGCCAAAAAAGAACTGCACCTC 186
      |||
Sbjct 121 CAGATTTGTGACCGGCGCGGTTTTTGT CAGCTTACTCCGGCCAAAAAAGAACTGCACCTC 180

Query 187 TGGAGCGG 194
      |||
Sbjct 181 TGGAGCGG 188
```

BLAST results
and alignment

NCBI/ BLAST Home

BLAST finds regions of similarity between biological sequences. [more...](#)

New DELTA-BLAST, a more sensitive protein-protein search

Human BLAT Search

BLAST Assembled Genomes

Find Genomic BLAST pages:

Enter organism name or accession number

- [Zebrafish](#)
- [Clawed frog](#)
- [Arabidopsis](#)
- [Rice](#)

BLAT Search Genome

Genome:

Human



HMMER

biosequence analysis using profile hidden Markov models

Home Search Results Software Help About
phmmer hmmscan hmmsearch jackhmmer

protein alignment/profile-HMM vs protein sequence database

[Paste a Sequence](#) | [Upload a File](#) | [Accession Search](#)

Paste in your alignment/hmm or use the [example](#)

Basic BLAST

Choose a BLAST program to use:

- [nucleotide blast](#) Search e-values with *Algo*
- [protein blast](#) Search p-values with *Algo*
- [blastx](#) Search protein sequences against a nucleotide database
- [tblastn](#) Search a nucleotide sequence against a protein database
- [tblastx](#) Search a nucleotide sequence against a protein database

Specialized BLAST

Choose a type of specialized search:

- Make specific protein-protein alignments
- Cluster multiple sequence alignments
- Find conserved domains
- Find sequences similar to a protein
- Search sequences similar to a protein
- Search immunoglobulin sequences
- Screen sequences against a protein database
- Align two (or more) protein sequences
- Search protein sequences
- Search SRA by protein
- Constraint Base
- Needleman-Wunsch [Global Sequence Alignment Tool](#)
- Search [RefSeqGene](#)
- Search [trace archives](#)

Paste in a query sequence to find matches separated by lines starting with >

File Upload: Rather than pasting a sequence, upload a file.

Upload sequence:

Only DNA sequences of 25,000 letters will be processed. Up to multiple sequence submissions

For locating PCR primers, use [Primer3](#)

Sequence Database

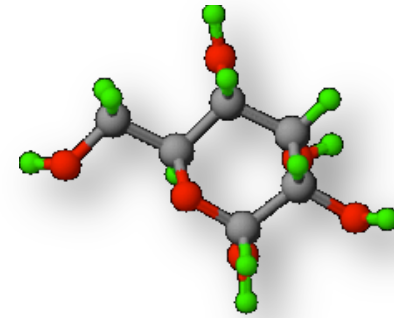
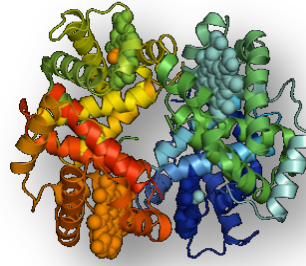
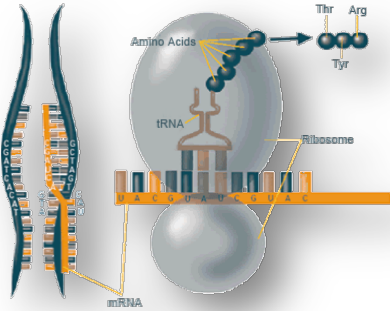
Large collections
 NR RefSeq UniProtKB Pfamseq

Annotated and Structures
 SwissProt PDB

Metagenomics and Environmental
 UniMes env NR

Representative Sets (UniProt)
 rp75 rp55 rp35 rp15 Reference Proteomes

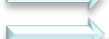
[Restrict by Taxonomy](#)



DNA
 $N \sim 3 \times 10^9$



RNA
 $N \sim 8 \times 10^4$



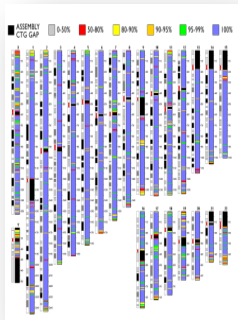
Protein
 $N \sim 10^6$



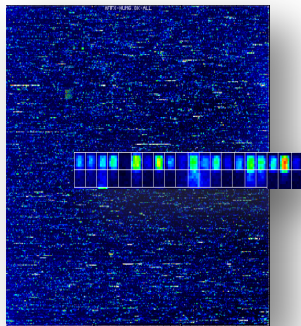
Modified
Proteins
 $N \sim 10^7$

Enzymes

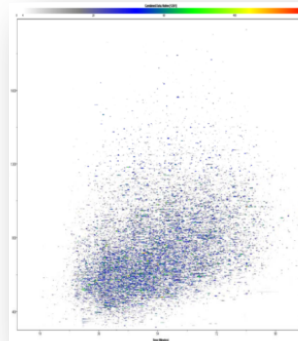
Metabolites
 $N \sim 6 \times 10^3$



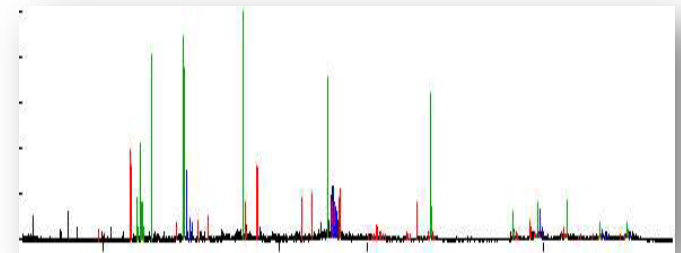
Genome



Transcriptome














Proteome



Metabolome

Introduction

- > [EBI homepage](#) 
- > [EBI bioinformatics tools](#) 
- > [UniProt](#) 
- > [NCBI homepage](#) 
- > [Gquery](#) 
- > [NCBIGene](#) 
- > [NCBI BLAST server](#) 
- > [ORF finder](#) 
- > [Splign](#) 
- > [MUSCLE](#) 
- > [ClustalOmega](#) 

Module 1

- > [Ensembl](#) 

Module 2

- > [UCSC genome browser](#) 
- > [VEGA](#) 
- > [NCBI Map Viewer](#) 

Module 3

- > [ECR browser](#) 
- > [VISTA Enhancer Browser](#) 
- > [BLINK](#) 
- > [Galaxy](#) 
- > [Homologene](#) 
- > [NCBI Trace Archives](#) 
- > [PHYLP](#) 
- > [PhyML](#) 
- > [Pipmaker](#) 
- > [Promo](#) 
- > [RepeatMasker](#) 
- > [rVISTA](#) 
- > [TreeBeST](#) 