# Annotating Genomes:
## where are all the genes?

## Jane Loveland
### Havana group

Introduction to Making Sense of Genomes
2nd October 2012

# Overview

- A bit of background

- What is manual genome annotation?
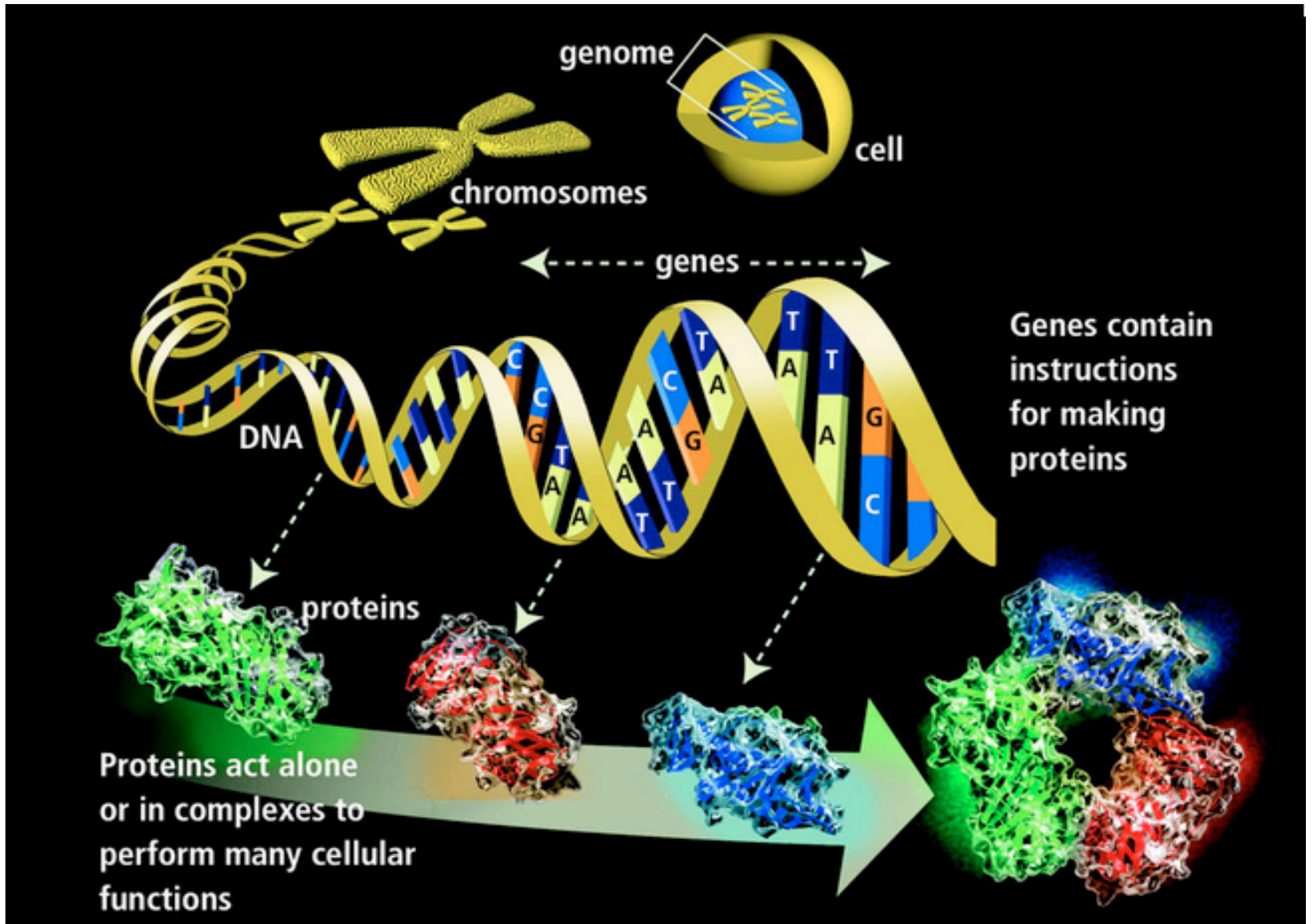
- Why do we need it?

- How do we do it?
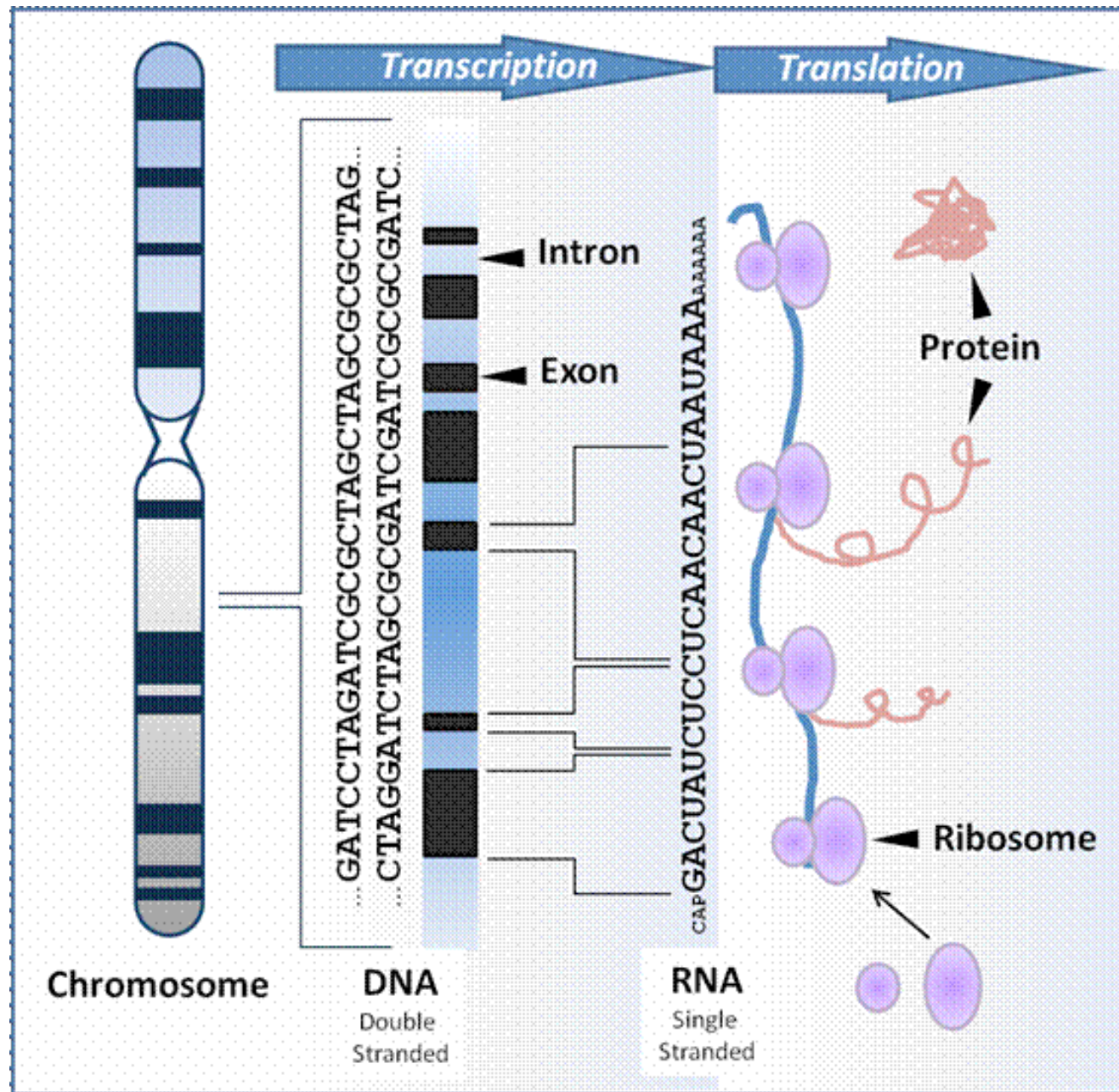
# Havana: Human And Vertebrate Analysis aNd Annotation

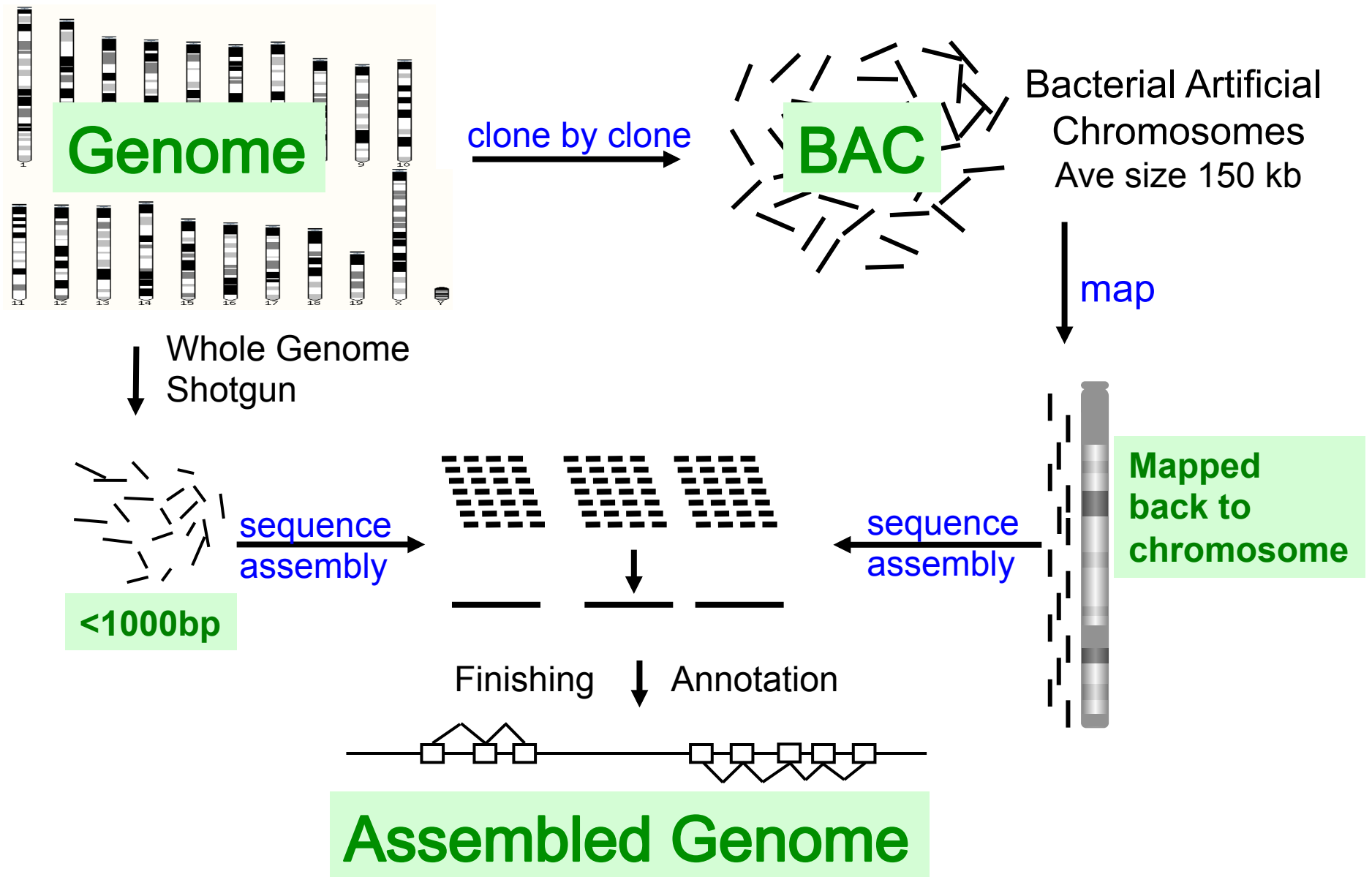Manual annotation of human, mouse and zebrafish whole chromosomes or genomes



# Vega: VErtebrate Genome Annotation database

genome

cell

chromosomes

genes

Genes contain instructions for making proteins

DNA

proteins

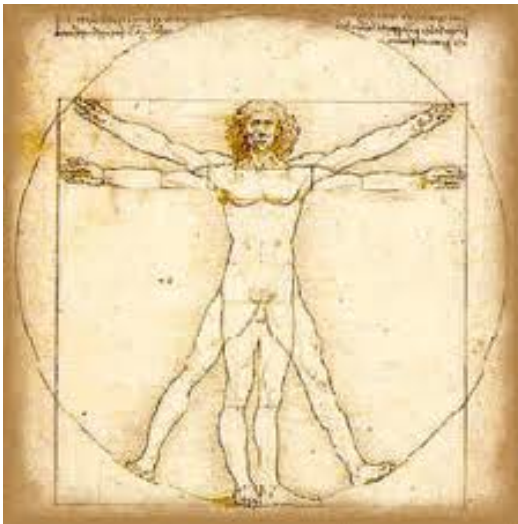Proteins act alone or in complexes to perform many cellular functions

Central Dogma Of Molecular Biology

# Hybrid Sequencing Strategy

# Reference genomes:
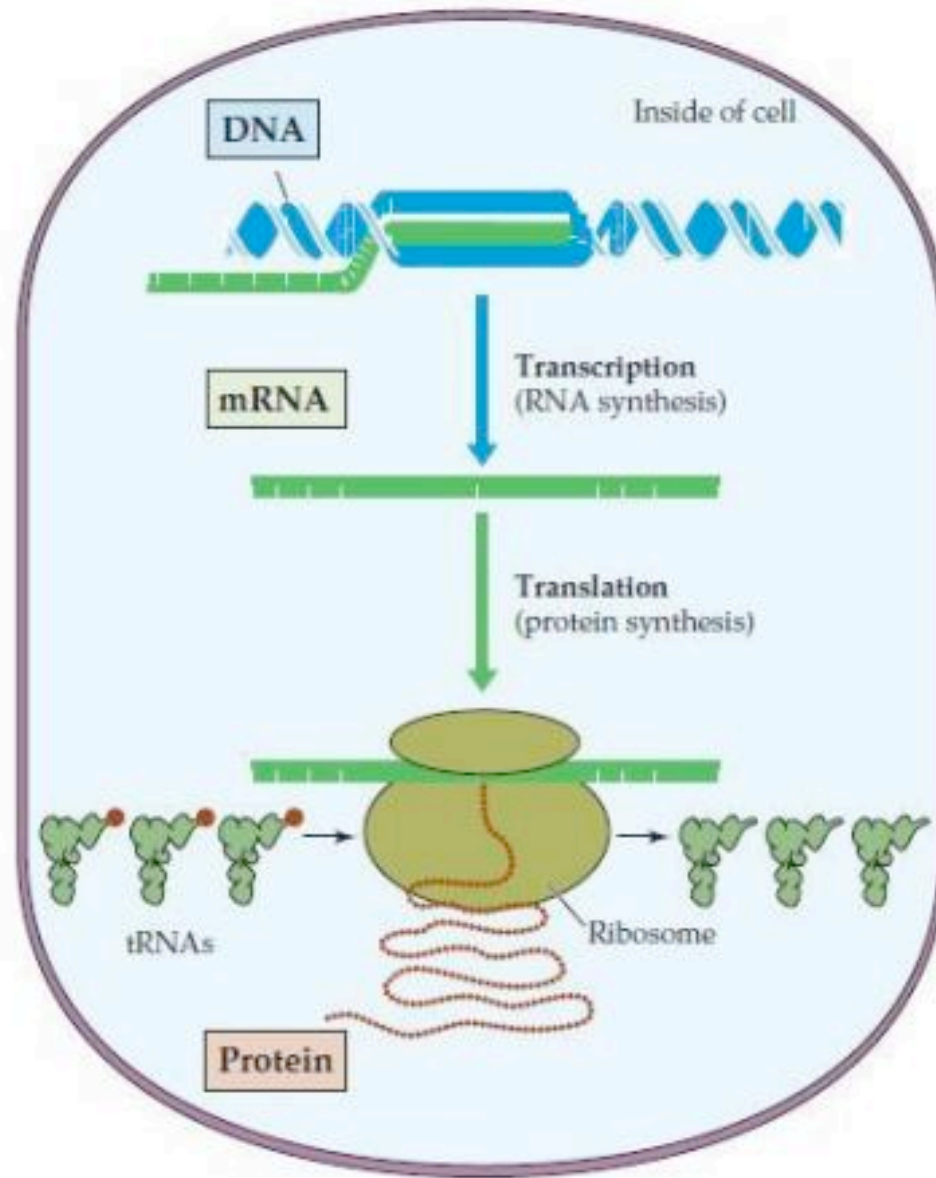


Human ~3Gb:
22 chromosomes + sex chromosomes

Mouse ~3 Gb:
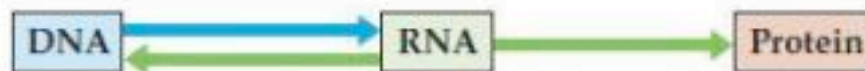19 chromosomes + sex chromosomes

Zebrafish ~1.4 Gb:
25 chromosomes, no specific sex chromosomes
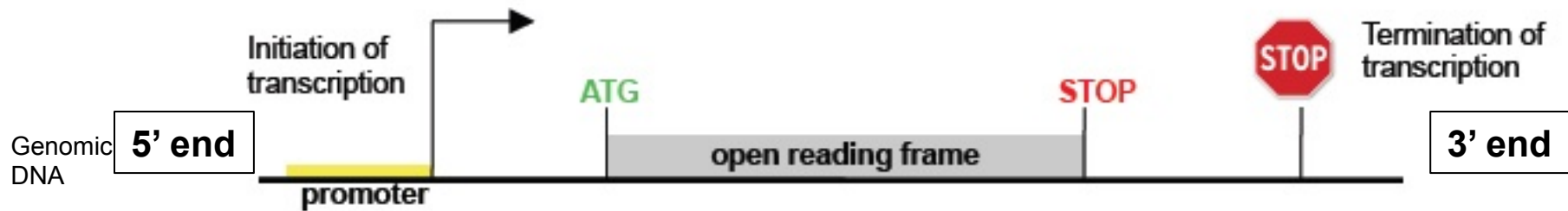
RNA polymerase
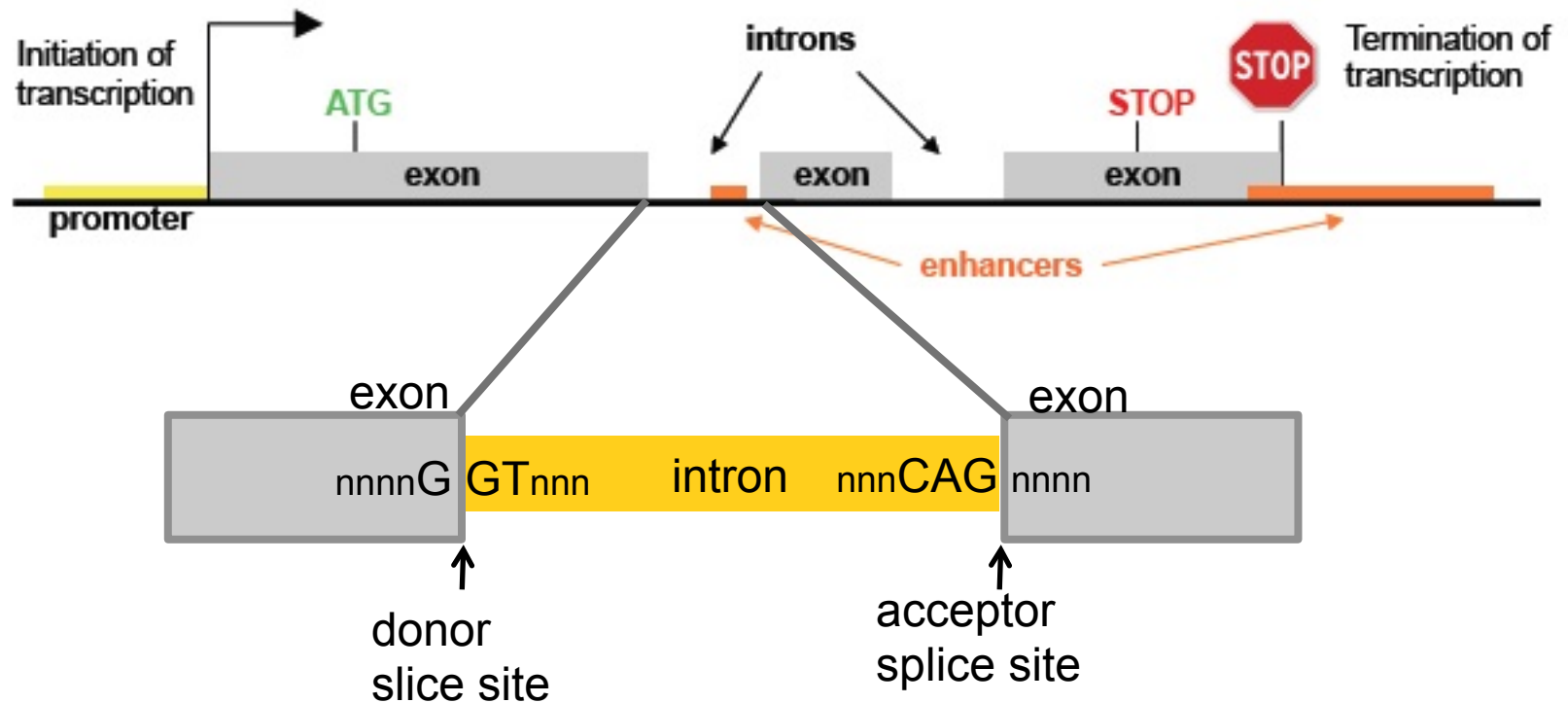
Ribosome binds to RNA generates amino acids

DNA

Inside of cell

Transcription (RNA synthesis)
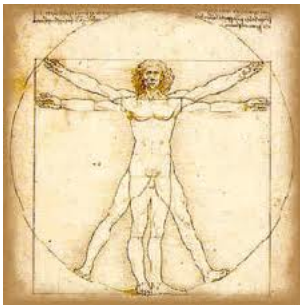
mRNA

Translation (protein synthesis)

tRNAs

Ribosome

Protein

DNA → RNA → Protein

# Prokaryotes:
## Simple protein-coding gene



Initiation of transcription

ATG

STOP

STOP Termination of transcription

Genomic DNA  **5' end**  open reading frame  **3' end**

promoter

# Eukaryotes:
## More complex: Introns and Exons



Initiation of transcription

ATG

introns

STOP STOP Termination of transcription

promoter  exon  exon  exon

enhancers

exon  exon

nnnnG GTnnn  intron  nnnCAG nnnn

donor slice site

acceptor splice site

# Do we know how many genes there are?

## Protein coding genes



| | |
|---|---|
| 1980's | 100, 000 |
| 2000 | 40, 000 |
| Today | ~ 21, 000 |

 21,638

 25,120

# Evidence for genes: DNA (*in vitro*)

**5' end**

mRNA extracted from tissue sample

**3' end**

**mRNA**  AUGCCTG~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~AATAAA~~~~~~~~AAAAAAAA

Reverse transcriptase transcribes the mRNA into a complimentary strand of DNA (cDNA)

**cDNA**  TACGGAC~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~TTATTT~~~~~~~~TTTTTTTTT
ATGCCTG~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~AATAAA~~~~~~~~AAAAAAAA

mRNA degraded and a second strand of cDNA made

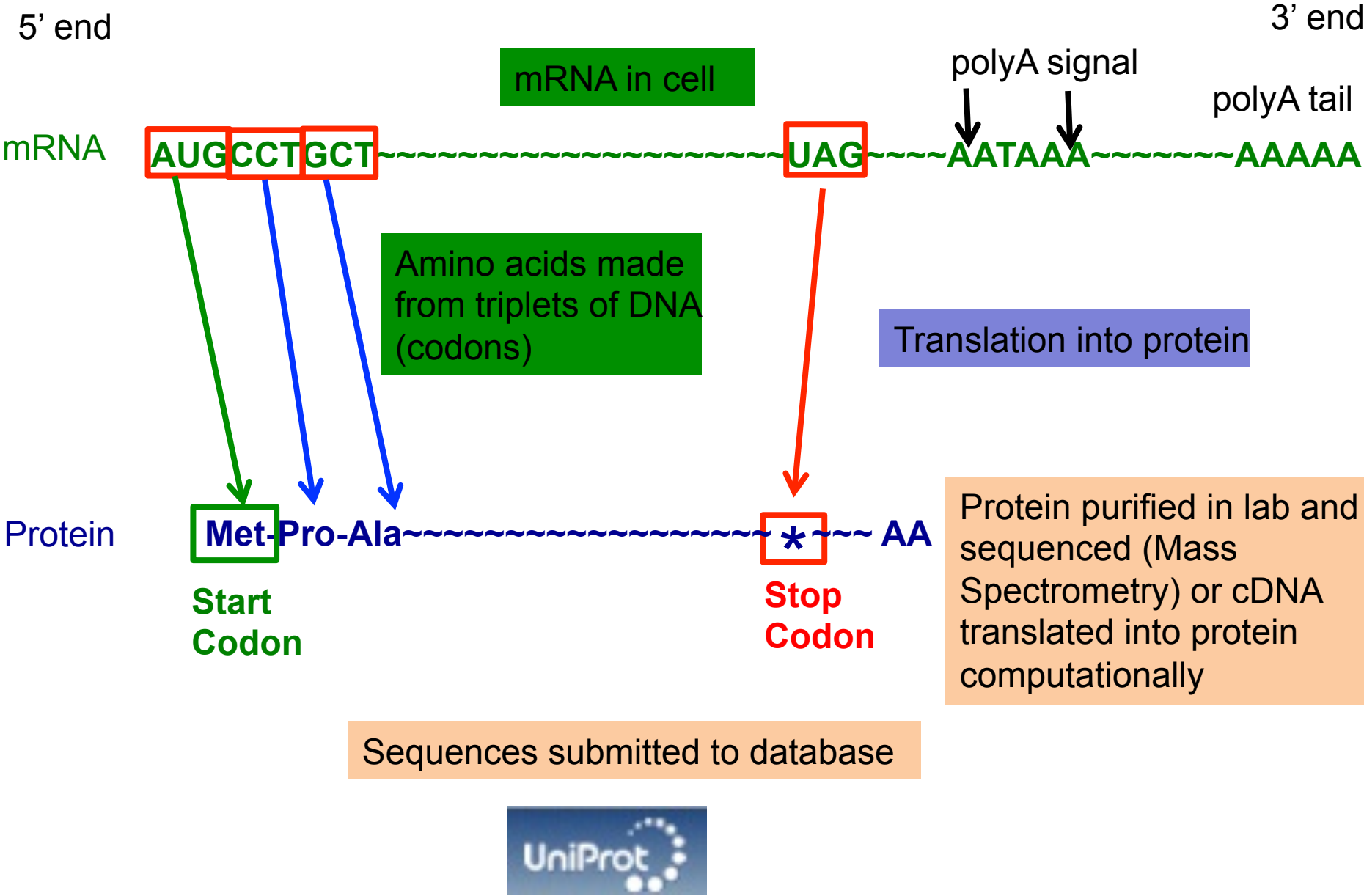**EST**  ATGCCTG~~~~~~~~~~~~~

~TTATTT~~~~~~~~TTTTTTTTT

**5' –> 3'**

Expressed Sequence Tags (ESTs) shorter stretches of cDNA (`800bp) sequences from both ends

**3' <– 5'**

Sequences submitted to databases

**ENA**
European Nucleotide Archive

**NCBI**
GenBank

**DDBJ**
DNA Data Bank of Japan

# Evidence for genes: Protein

5' end

3' end

polyA signal

polyA tail

mRNA in cell

mRNA  **AUG CCT GCT** ~~~~~~~~~~~~~~~~~~~ **UAG** ~~~ **AATAAA** ~~~~~~~~ **AAAAA**

Amino acids made from triplets of DNA (codons)

Translation into protein

Protein  **Met-Pro-Ala** ~~~~~~~~~~~~~~~~~ **∗** ~~~ **AA**

Protein purified in lab and sequenced (Mass Spectrometry) or cDNA translated into protein computationally

**Start Codon**

**Stop Codon**

Sequences submitted to database

UniProt

# 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| **U E G M D** 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

**BLAST results and alignment**

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

Query  67   CTGCTGCGCCTCGGGTGTCTTTTGCGGCGGTGGGTCGCCGCCGGGAGAAGCGTGAGGGGA  126
            ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct  61   CTGCTGCGCCTCGGGTGTCTTTTGCGGCGGTGGGTCGCCGCCGGGAGAAGCGTGAGGGGA  120

Query  127  CAGATTTGTGACCGGCGCGGTTTTTGTCAGCTTACTCCGGCCAAAAAAGAACTGCACCTC  186
            ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct  121  CAGATTTGTGACCGGCGCGGTTTTTGTCAGCTTACTCCGGCCAAAAAAGAACTGCACCTC  180

Query  187  TGGAGCGG  194
            ||||||||
Sbjct  181  TGGAGCGG  188
```
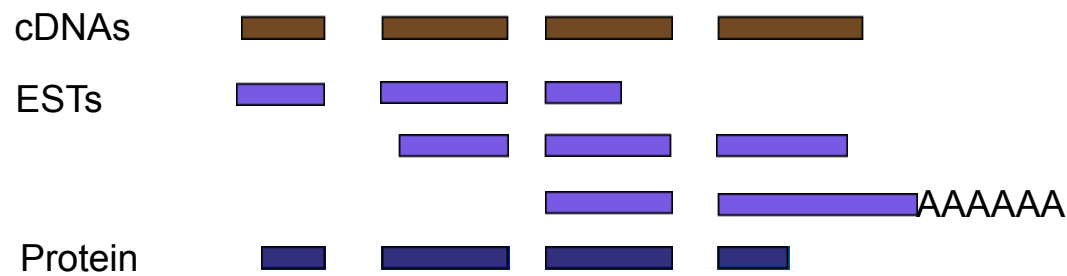
# Conservation:
# GAPDH gene coding exons

# Conservation:
## TCEB2 gene

3' end          Coding exons          5' end



Conserved
regions

# Making the transcript from evidence:
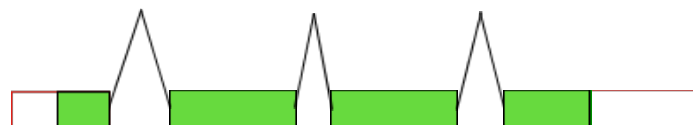
Genomic sequence

Analysis pipeline

cDNAs

ESTs

Sequences from databases

AAAAAA

Protein

Annotation
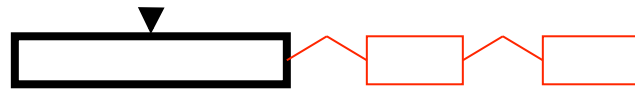
Gene structure

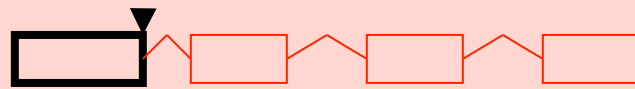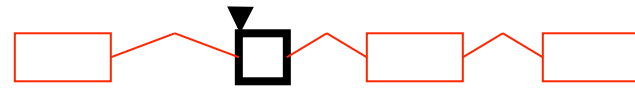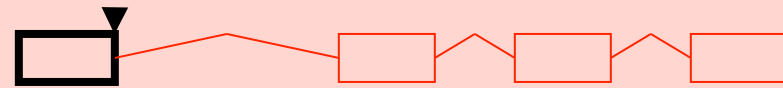# Alternative Splicing



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

# Do we know how many genes there are?

|  | Coding | Non-coding | Pseudogene |
|---|---|---|---|
|  | 21,224 | 15,962 | 14,427 |
|  | 21,638 | 6,875 | 5,510 |
|  | 25,120 | 4,556 | 224 |

# 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) | UGUAAACAUCCCCGACUGGAAGCUGUAAGCCAC....AGCCAAGCUUUCAGUCAGAUGUUUGCU |
| Spermophilus tridecemlineatus (th | UGUAAACAUCCCCGACUGGAAGCUGUAGGACAC....AGCUGAGCUUUCAGUCAGAUGUUUGCU |
| Macaca mulatta (Rhesus monkey) | UGUAAACAUCCUA..CACUCAGCUGUAAUACAU....GGAUUGGCUGGGAGGUGGAUGUUUACU |
| Macaca nemestrina (pig-tailed mac | UGUAAACAUCCUACACUCUCAGCUGUGGAAAGU....AAGAAAGCUGGGAGAAGGCUGUUUACU |
| Macaca nemestrina (pig-tailed mac | UGUAAACAUCCUA..CACUCAGCUGUAAUACAU....GGAUUGGCUGGGAGGUGGAUGUUUACU |
| Gorilla gorilla (western gorilla) | UGUAAACAUCCUA..CACUCAGCUGUAAUACAU....GGAUUGGCUGGGAGGUGGAUGUUUACU |
| Homo sapiens (human) | UGUAAACAUCCCCGACUGGAAGCUGUAAGACAC....AGCUAAGCUUUCAGUCAGAUGUUUGCU |

**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)
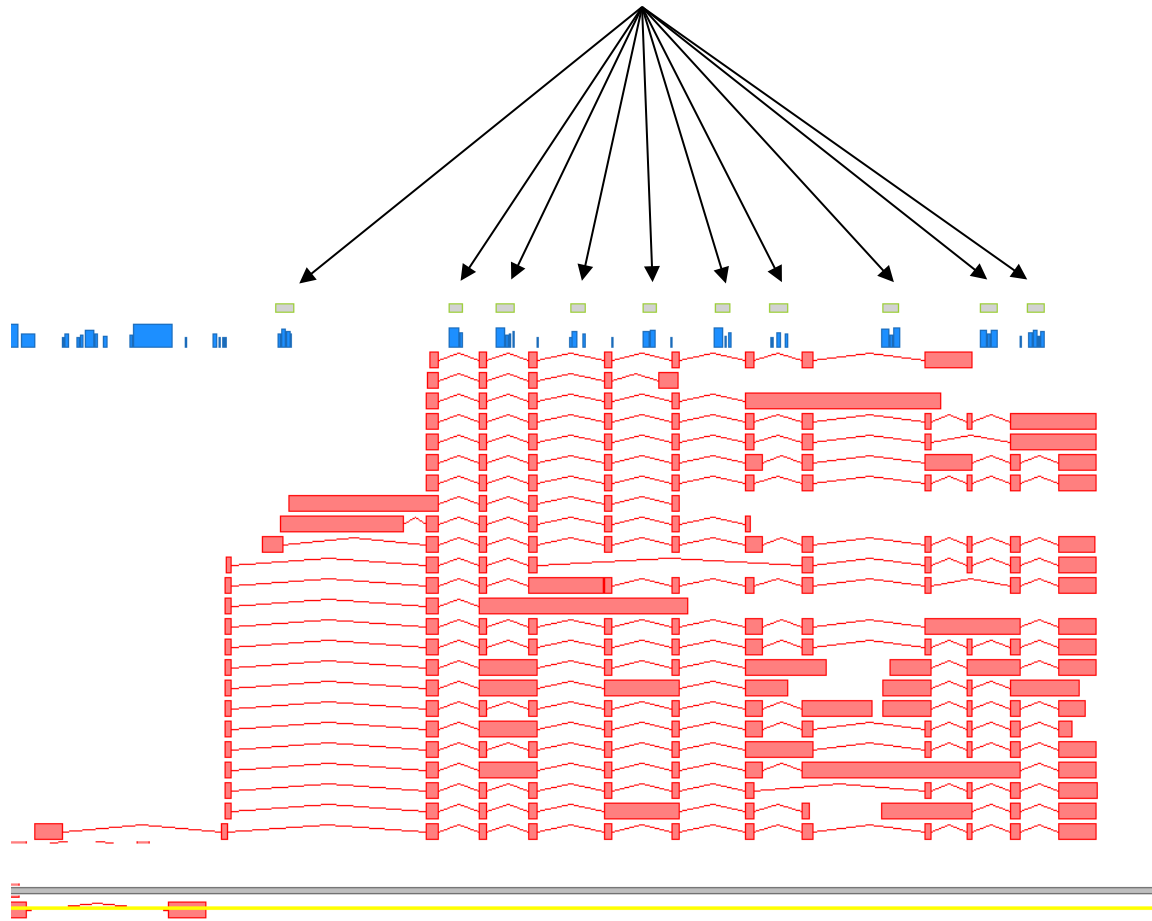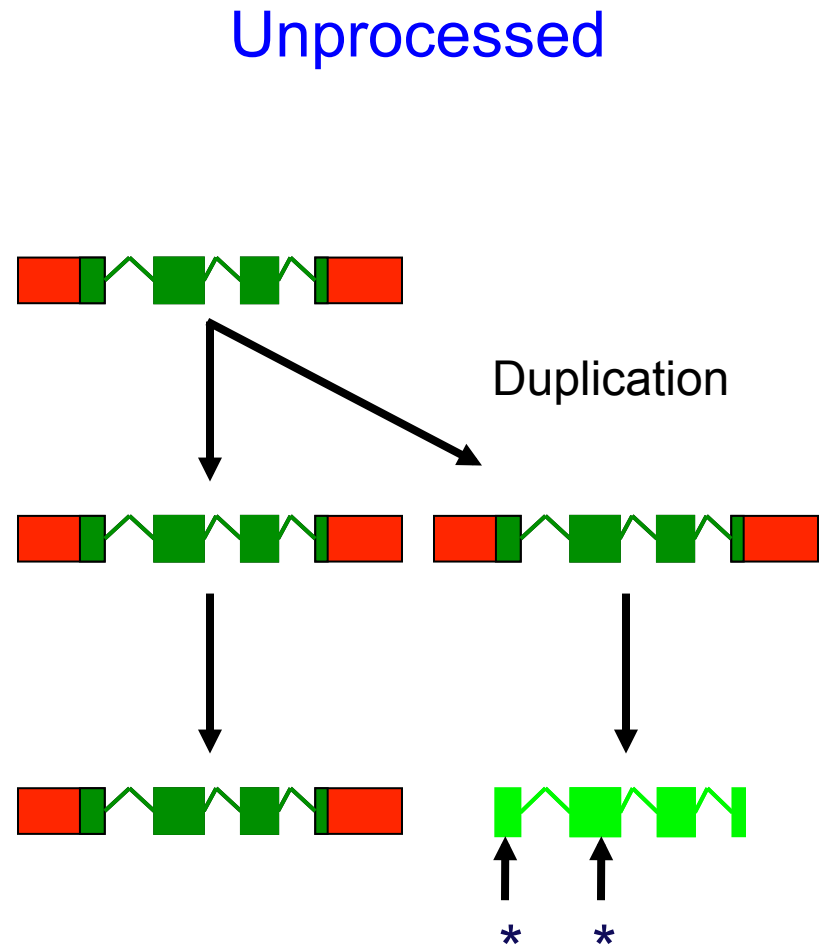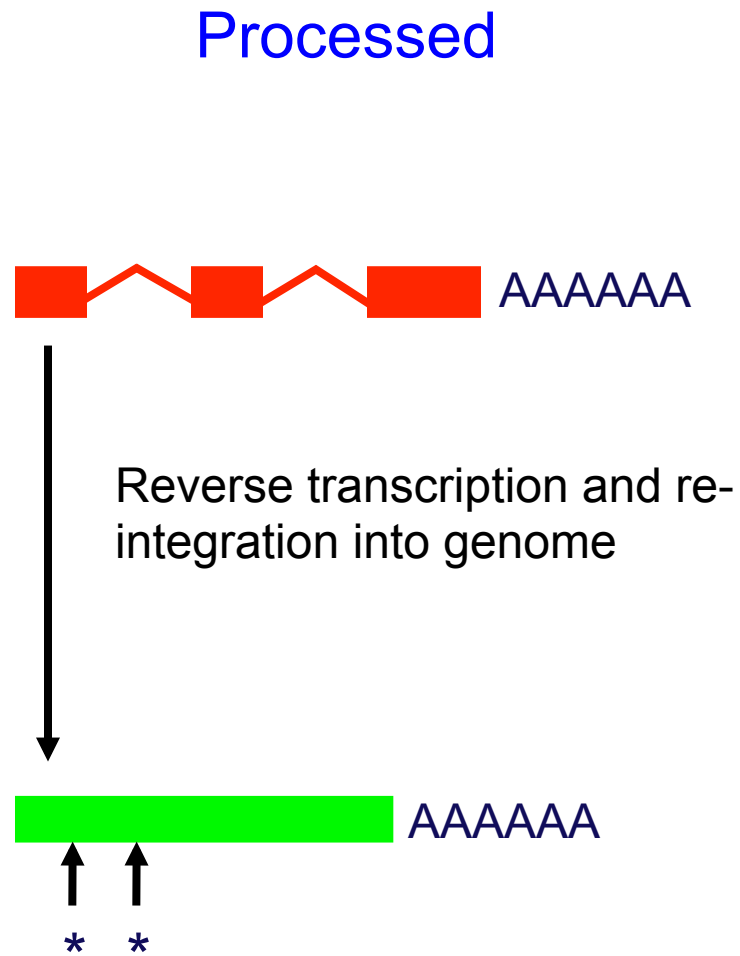Many others not yet characterised

# lncRNAs

sense strand 5′ →

lincRNA

GENE1-IT1

GENE2-OT1

GENE2-UA1

GENE 1 (Coding)

GENE 2 (Coding)

GENE1-AS1

GENE2-AS1

GENE2-AS3

GENE2-GENE1-AS1

GENE2-AS2

5′ antisense strand

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

## intronic snoRNAs

# Havana pseudogenes



Processed

Unprocessed

AAAAAA

Reverse transcription and re-integration into genome

Duplication

AAAAAA

* *

* *

Disruption to coding sequence and stop codons

The Vertebrate Genome Annotation (VEGA) database is a central repository for high quality manual annotation of vertebrate finished genome sequence. Human, mouse and zebrafish are in the process of being completely annotated, whereas for other species the annotation is only of specific genomic regions of particular biological interest. The majority of the annotation is from the HAVANA group at the Welcome Trust Sanger Institute

The website is built upon code from the Ensembl project.

## Browse a Genome

**Human** [25-09-2012]
Ensembl

**Zebrafish** [25-09-2012]
Ensembl

**Mouse** [26-06-2012]
Ensembl

**Pig** [25-09-2012]
Ensembl

**Chimpanzee** [12-01-2012]
Ensembl

**Gorilla** [30-03-2009]
Ensembl

**Wallaby** [30-03-2009]
Ensembl

**Dog** [14-02-2005]
Ensembl

# Genome browsers include Vega genes:

## GENCODE geneset:

# Genome browsers include Vega genes:

## GENCODE geneset:

# Genome browsers include Vega genes:
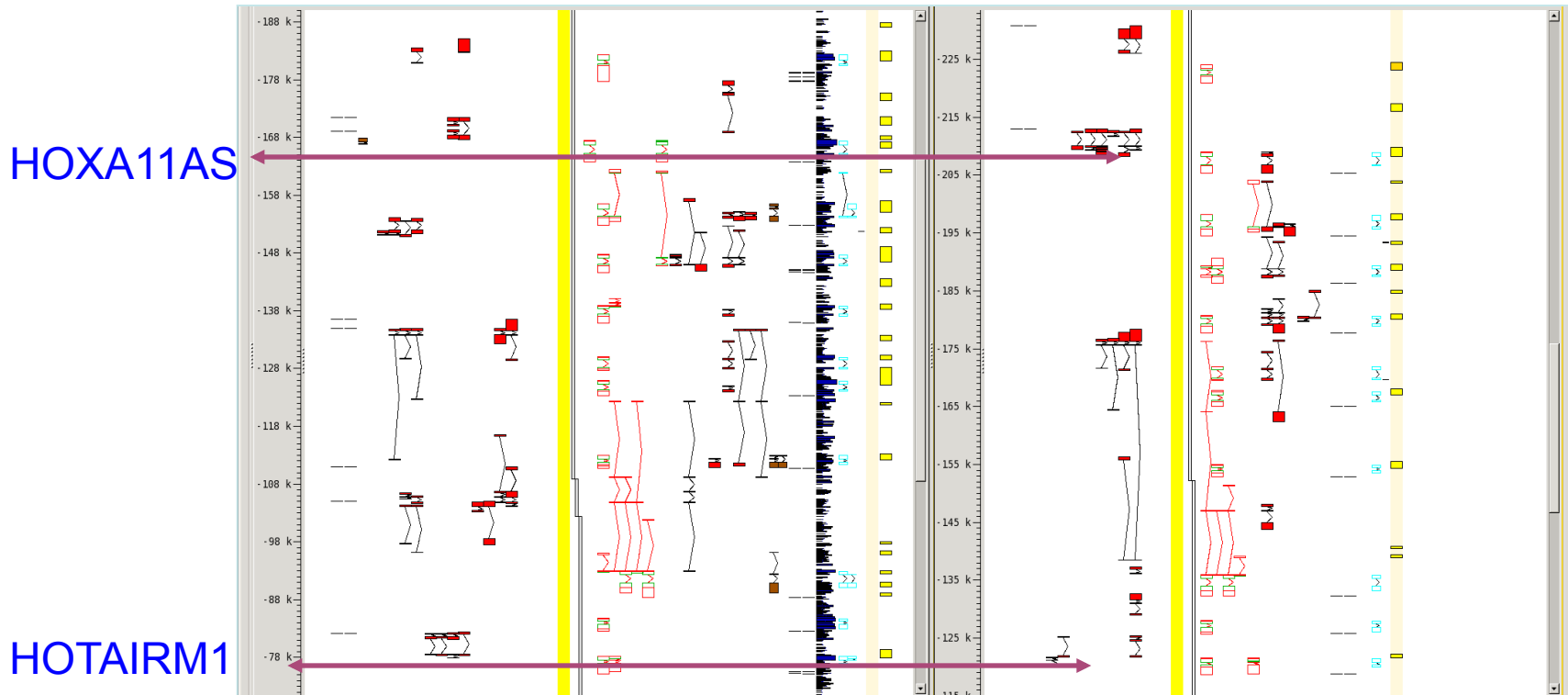
## GENCODE geneset:

# Linked loci

# HOXA gene cluster



Human chr 7p15.2

Mouse chr 6qB3

HOXA11AS

HOTAIRM1

Long non-coding transcripts are conserved across species and regulate expression of HOX genes

# Acknowledgements

**Havana:**
Jen Harrow
If Barnes
Ruth Bennett
Alex Bignell
Veronika Boychenko
Gloria Despacio-Reyes
Sarah Donaldson
Adam Frankish
Matt Hardy
Toby Hunt
Mike Kay
Gavin Laird
David Lloyd
Jane Loveland
Deepa Manthravadi
Gaurab Mukherjee
Jonathan Mudge
Jeena Rajan
Gary Saunders
Catherine Snow
Charles Steward
Marie-Marthe Suner
Mark Thomas
Laurens Wilming

**Encode:**
Jose Gonzalez
Sarah Grubb
Electra Tapanari

**Anacode:**
James Gilbert
Matthew Astley
Michael Gray
Jeremy Henty

**Vega:**
Stephen Trevanion
Dan Sheppard

**Zmap:**
Ed Griffiths
Gemma Barson
Malcolm Hinsley

**Project Coordinator:**
Tim Hubbard

http://vega.sanger.ac.uk