# Predicting genes using RNA sequencing
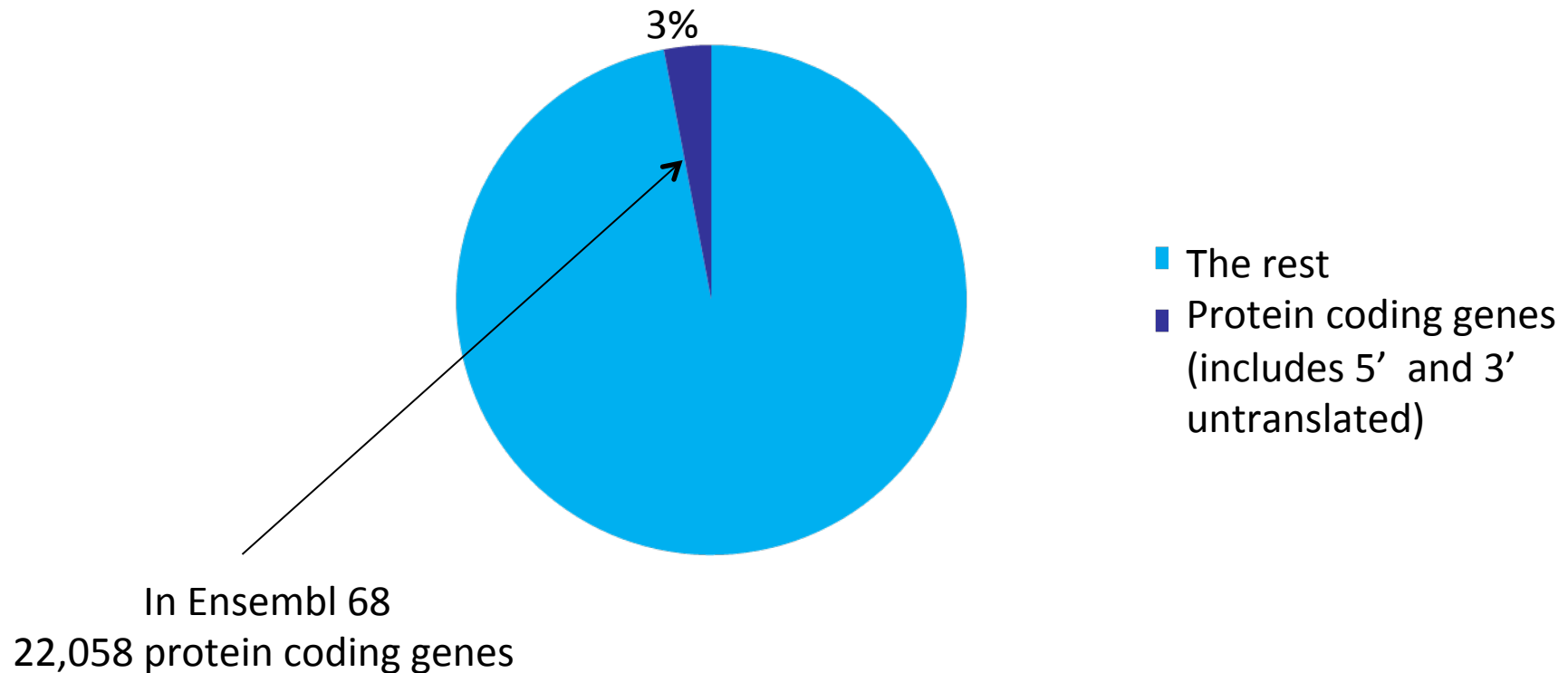
9th October 2012

# Predicting genes – finding the critical 3%

Human Genome of 3,300,000,000 bases

3%



The rest
Protein coding genes
(includes 5' and 3'
untranslated)

In Ensembl 68
22,058 protein coding genes

# 3% - like finding 7or 8 people in this crowd

# Easier if you know who you are looking for…

# Predicting genes using RNA sequencing
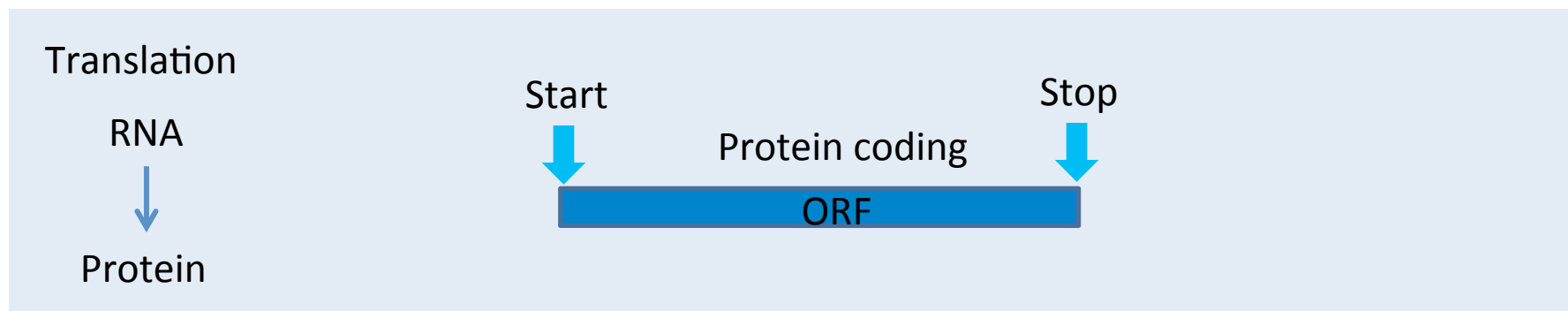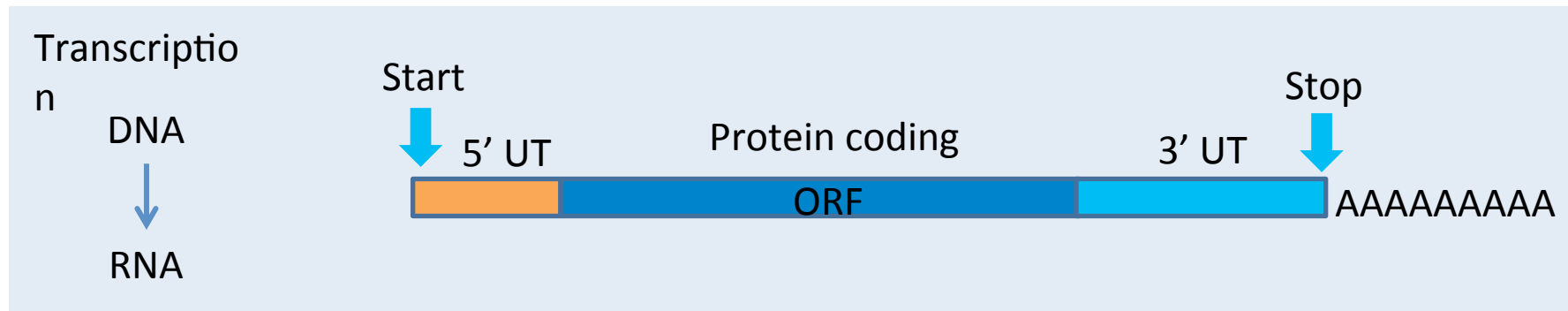
Identify the genes by
sequencing RNA

↓

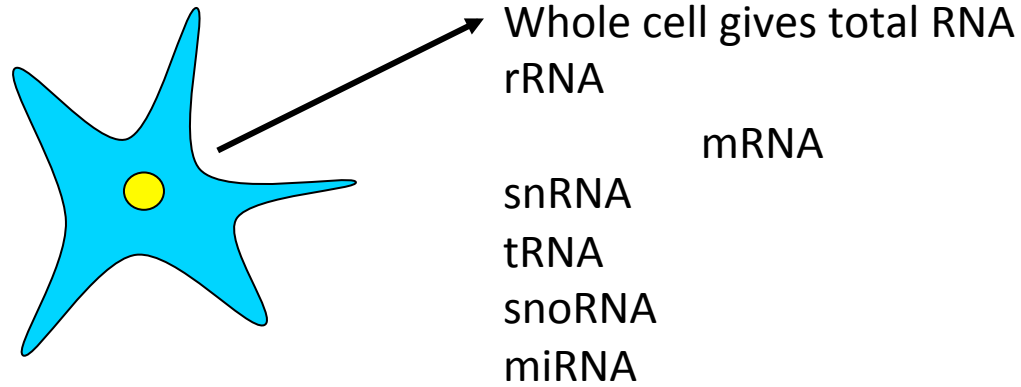Find the location of the
genes in the genome

# Predicting genes using RNA sequencing

1.  Getting the RNA template

2.  Sequencing the RNA

3.  Predicting gene models

4.  Producing a gene build

5.  Using the gene build
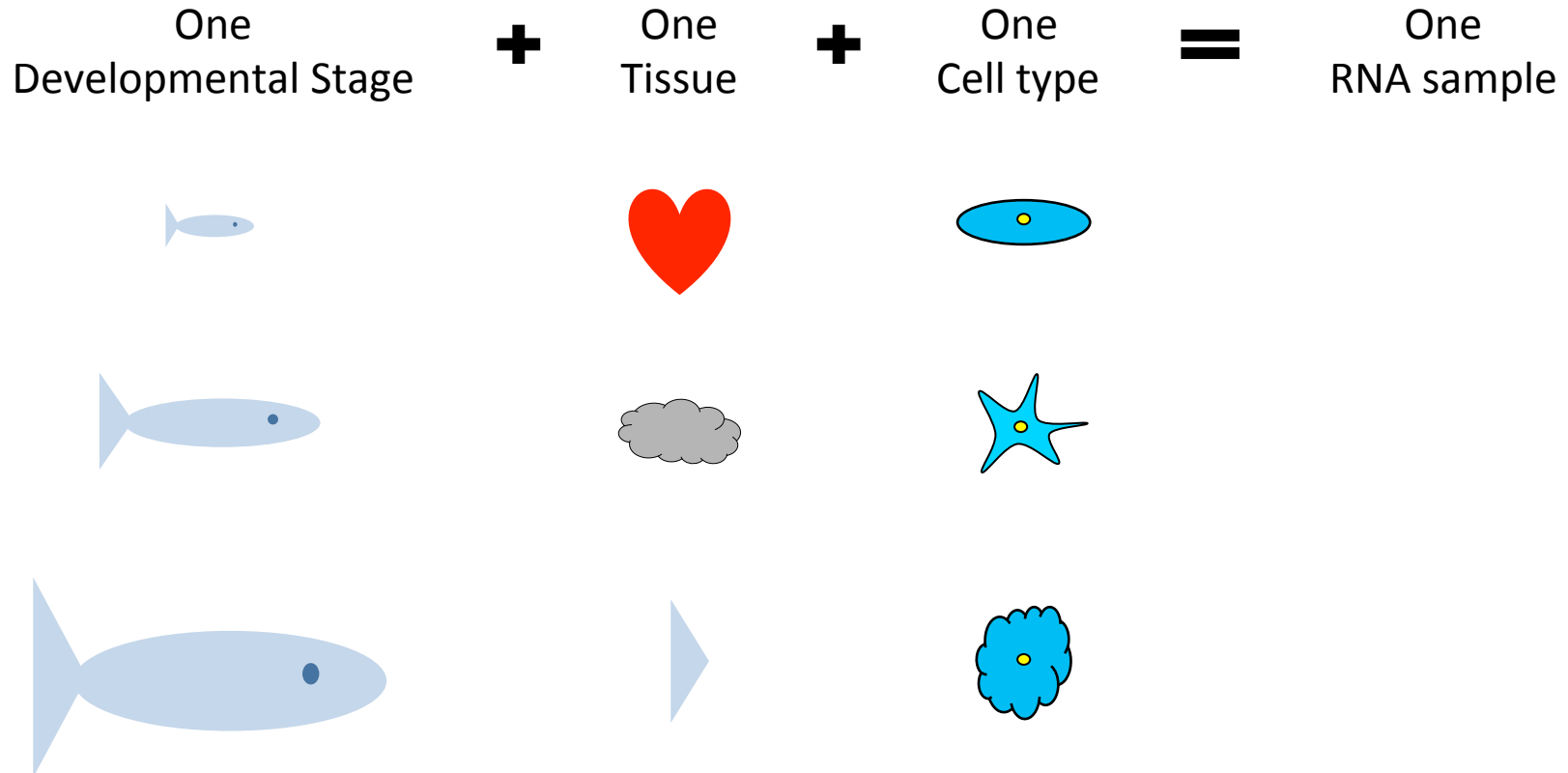
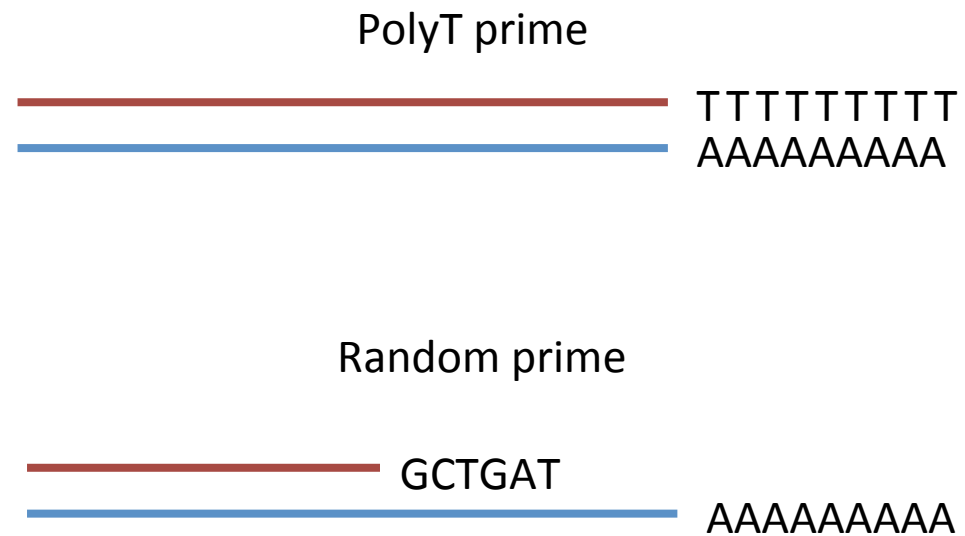# Predicting genes – what are we looking for?

# Getting hold of the RNA

Whole cell gives total RNA
rRNA

mRNA

snRNA

tRNA

snoRNA

miRNA

T T T T T T T T T---

ORF

AAAAAAAAAA

# Getting hold of the RNA

For each sample pick can pick …
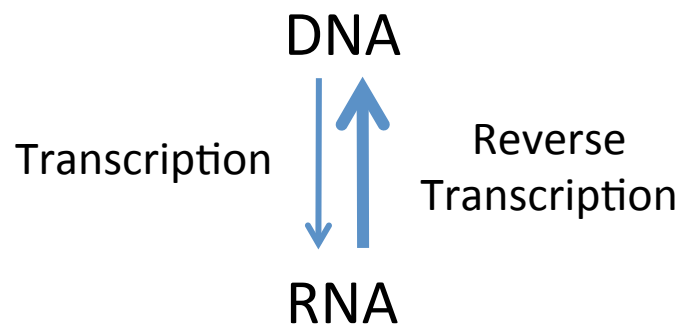
One
Developmental Stage **+** One
Tissue **+** One
Cell type **=** One
RNA sample



Each sample has a different set of mRNAs each representing a genes –
the transcriptome of the sample

# Preparing the sequence template – reverse transcription

Sequencing machines use DNA as a template
so the RNA must be converted back into
complementary DNA or cDNA

DNA

Transcription

Reverse
Transcription

RNA

PolyT prime

TTTTTTTTT
AAAAAAAAA

Random prime

GCTGAT

AAAAAAAAA

# Sequencing methods

**Capillary sequencing**

One sequence from one template about 800 bases long

Template linked directly to sequence

**Illumina sequencing or GAII or HiSeq or MiSeq**

Millions of sequences from millions of templates all at the same time about 100 bases long

No template/sequence information
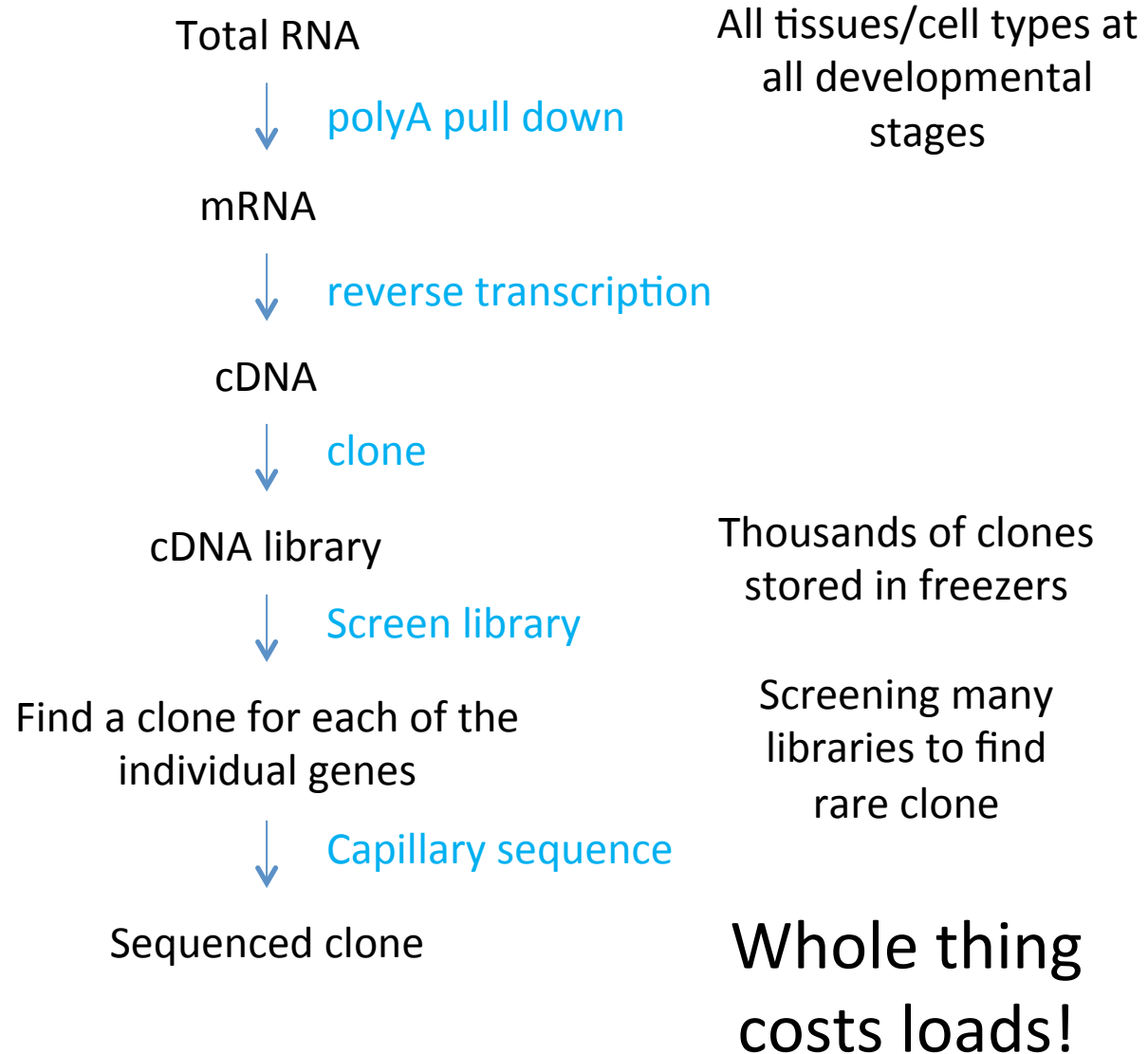
# Sequencing methods

# Traditional RNA sequencing

Total RNA

&darr; polyA pull down

mRNA

&darr; reverse transcription

cDNA

&darr; clone

cDNA library

&darr; Screen library

Find a clone for each of the individual genes

&darr; Capillary sequence

Sequenced clone

All tissues/cell types at all developmental stages

Thousands of clones stored in freezers

Screening many libraries to find rare clone

Whole thing costs loads!

# Next generation RNA sequencing

RNAseq

PolyA pull down for mRNA

Synthesis double strand
cDNA

Chemical hydrolysis

Make a library and Sequence

Note each fragment is sequenced from
both ends to make a pair of reads

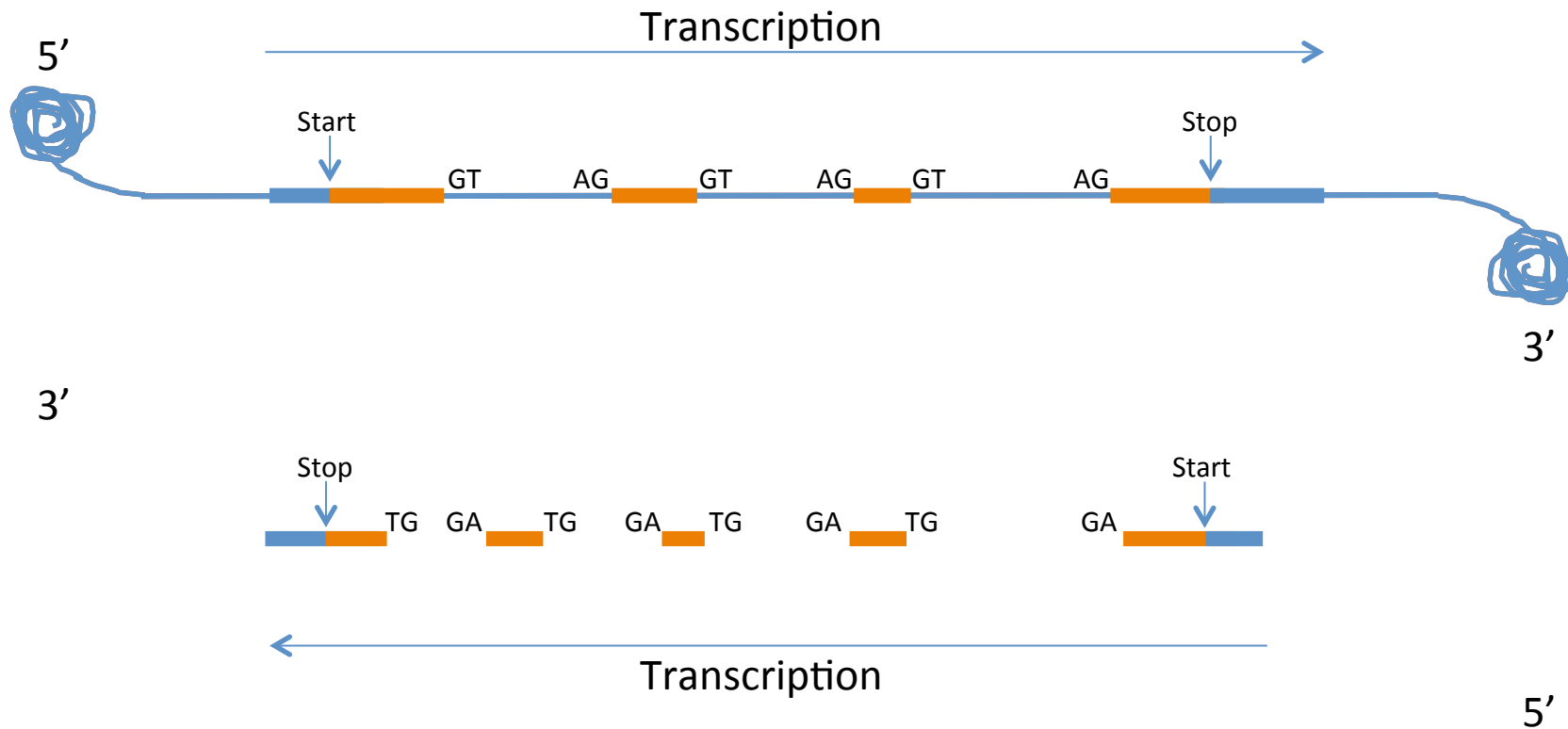# Annotation genes on genome reference sequence

Species specific whole gene sequences – **cDNA library sequencing**

Species specific small fragments of gene sequences - **RNAseq**

All the sequence from other animals in the public databases

Make **GENE MODELS** for a **GENE BUILD**

# Annotation gene models

# Building genes from species-specific cDNAs

>part of human C22 reference sequence

gccgggctcatgctgggcacctccgggcaagctttgtgacttagaggtgtggggccactggtcaccctcggtggctcagaggctgtggctcc
atggctcatgagcgcctcctgtgtcccagacctgagaccatgcgtatgtcccgaggcggctgctccggtctcttgcggtgtggaaagggaac
gacctacgagggcggtgtccgagagcctgccttggccttctggccaggtcatatcgctcccggtcagtccgcaggccctctccttggaaccct
ggccccaccaccccaaccttgatggcgaactgagtgactgaccagcctcctgcccccaggcgtgacccacggagctggccagctccctg
gacctgctgcctaccctggcagccctggctggggccccactgcccaatgtcaccttggatggctttgacctcagcccccctgctgctgggcaca
ggcaaggtagggccggtgaccctgatcccagatccttggcccctgtcctggccttccctgggtgagtgtgggcagtgcctgagagtctgt
gcctcagtgcctcctgcactgagtggcatccaagtggcgccacctctcaggttcctgggtgggcaagaagcggtgcacg
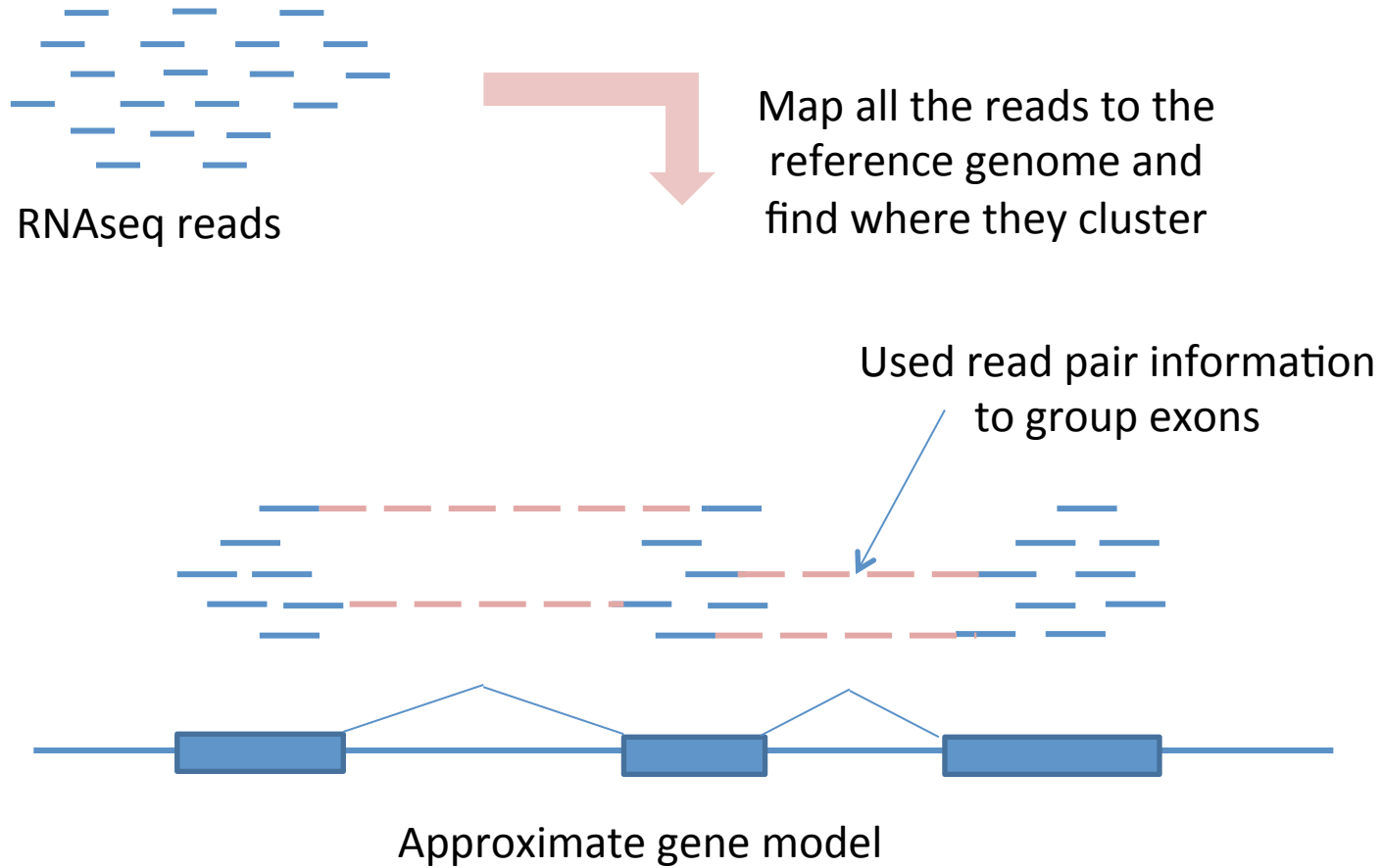
## PLUS

>exons of the gene ARSA cDNA sequence

acctgagaccatgcgtatgtcccgaggcggctgctccggtctcttgcggtgtggaaagggaacgacctacgagggcggtgtccgagagcctgccttggcctt
ctggccaggtcatatcgctcccggcgtgacccgacgagctggccagctccctggacctgctgcctaccctggcagccctggctggggccccactgcccaatg
tcaccttggatggctttgacctcagccccctgctgctgggcacaggcaag

## EQUALS
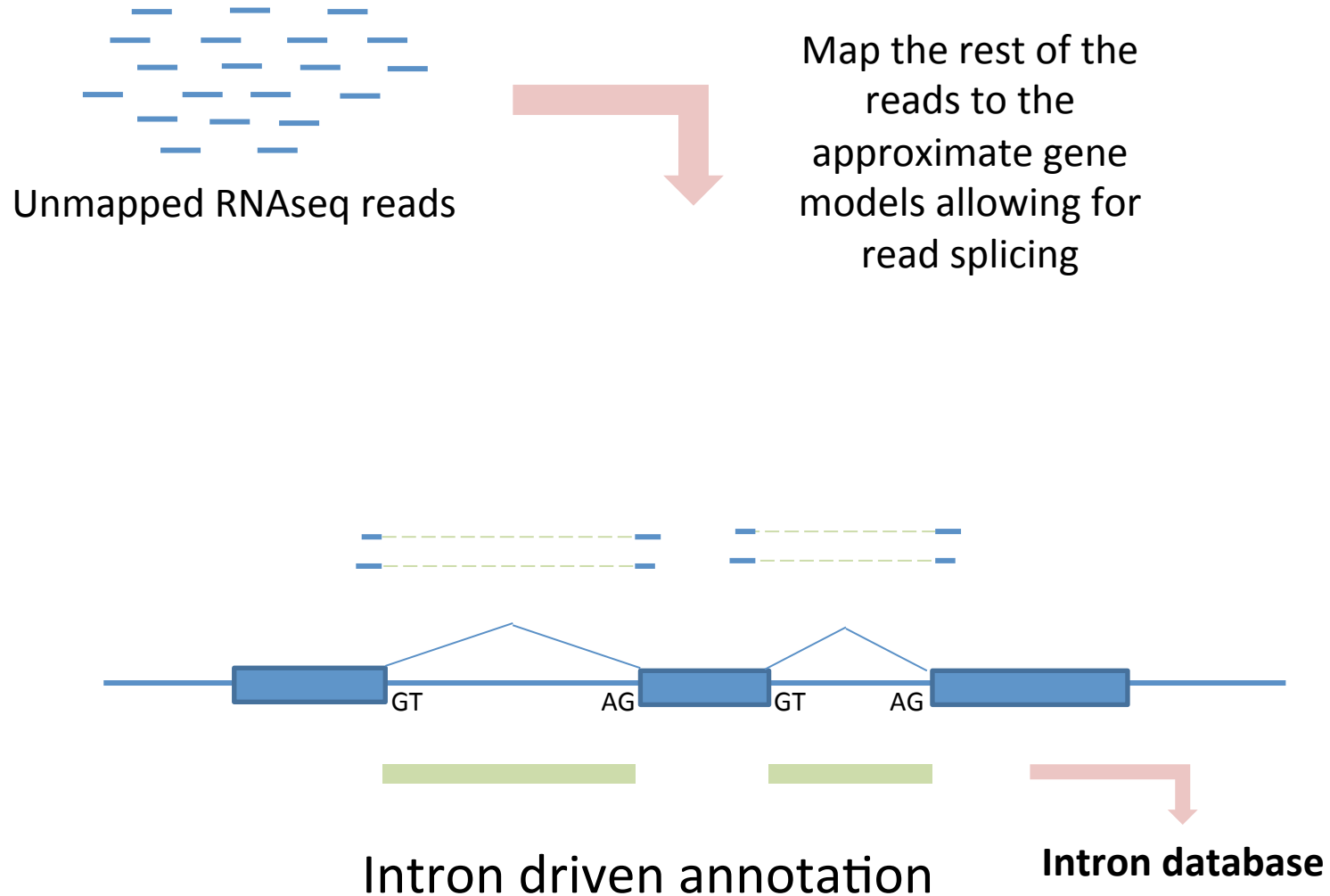
>part of human C22 sequence annotated with two exons from the ARSA cDNA sequence

gccgggctcatgctgggcacctccgggcaagctttgtgacttagaggtgtggggccactggtcaccctcggtggctcagaggctgtggctcc
atggctcatgagcgcctcctgtgtccc**ag**acctgagaccatgcgtatgtcccgaggcggctgctccggtctcttgcggtgtggaaagggaac
gacctacgagggcggtgtccgagagcctgccttggccttctggccaggtcatatcgctcccg**gt**cagtccgcaggccctctccttggaaccct
ggccccaccaccccaaccttgatggcgaactgagtgactgaccagcctcctgccccc**ag**gcgtgacccacggagctggccagctccctg
gacctgctgcctaccctggcagccctggctggggccccactgcccaatgtcaccttggatggctttgacctcagccccctgctgctgggcaca
ggcaag**gt**agggccggtgaccctgatcccagatccttggcccctgtcctggccttccctgggtgagtgtgggcagtgcctgagagtctgt
gcctcagtgcctcctgcactgagtggcatccaagtggcgccacctctcaggttcctgggtgggcaagaagcggtgcacg
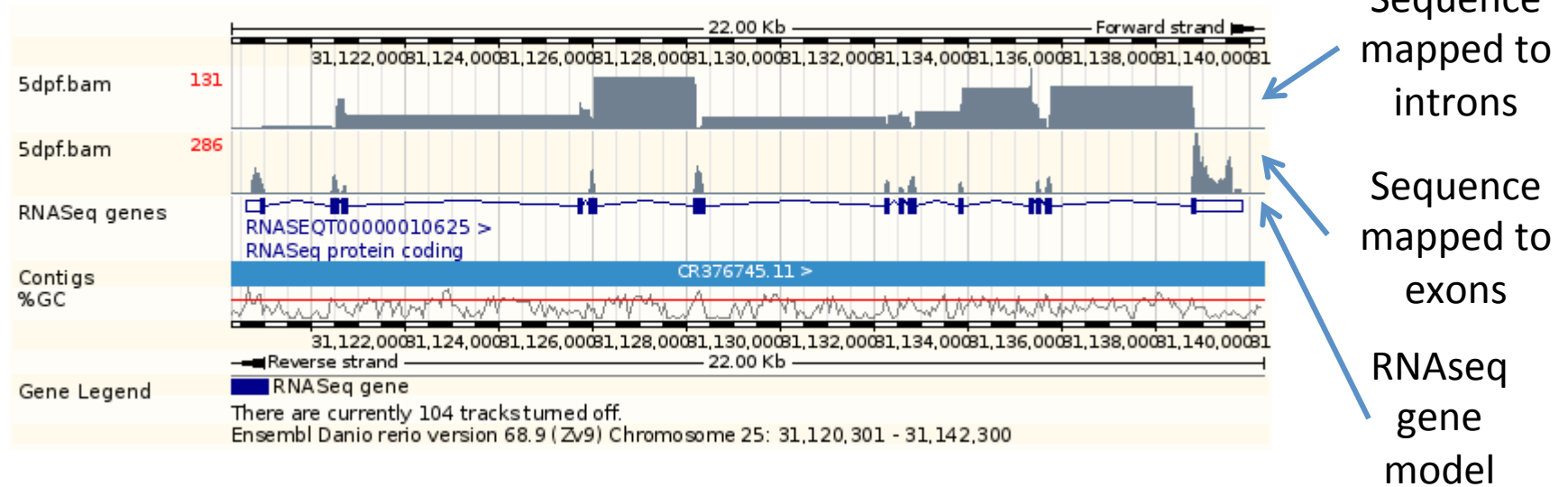
# Building genes from RNAseq



RNAseq reads

Map all the reads to the reference genome and find where they cluster

Used read pair information to group exons

Approximate gene model

# Building genes from RNAseq



Unmapped RNAseq reads

Map the rest of the reads to the approximate gene models allowing for read splicing

GT          AG   GT          AG

Intron driven annotation          **Intron database**

# Building genes from RNAseq

Example of an RNAseq gene built from the short reads for apopotosis inhibitor 5 (api5)



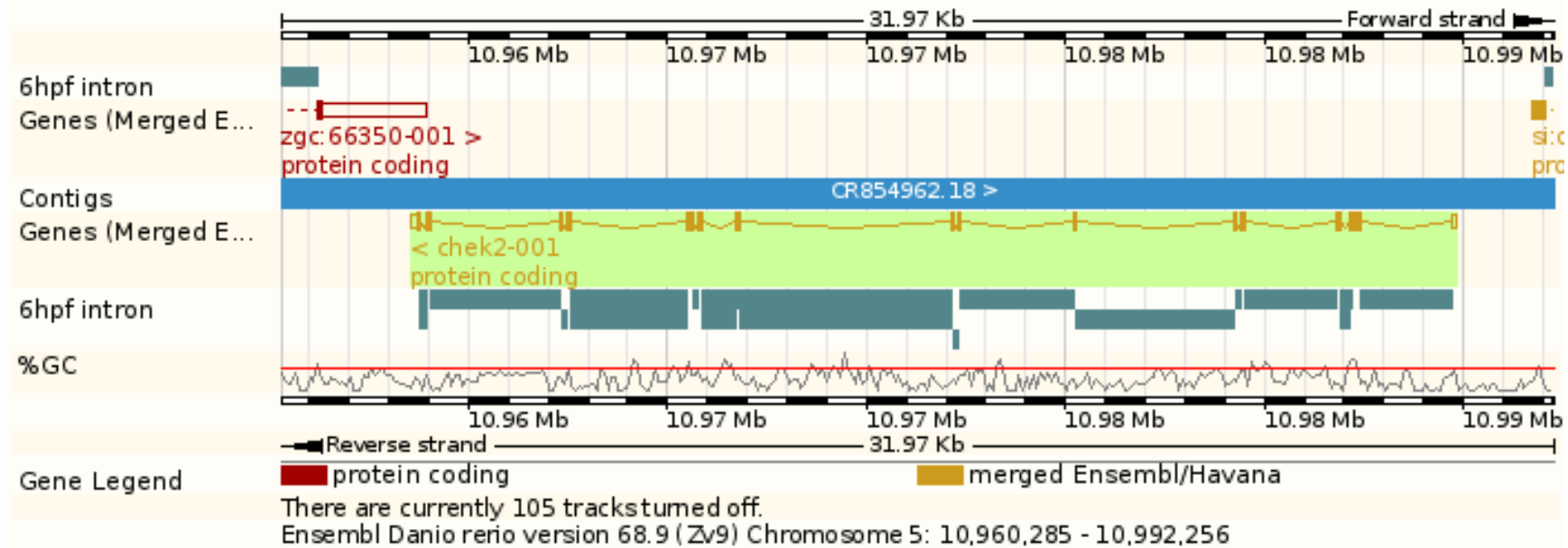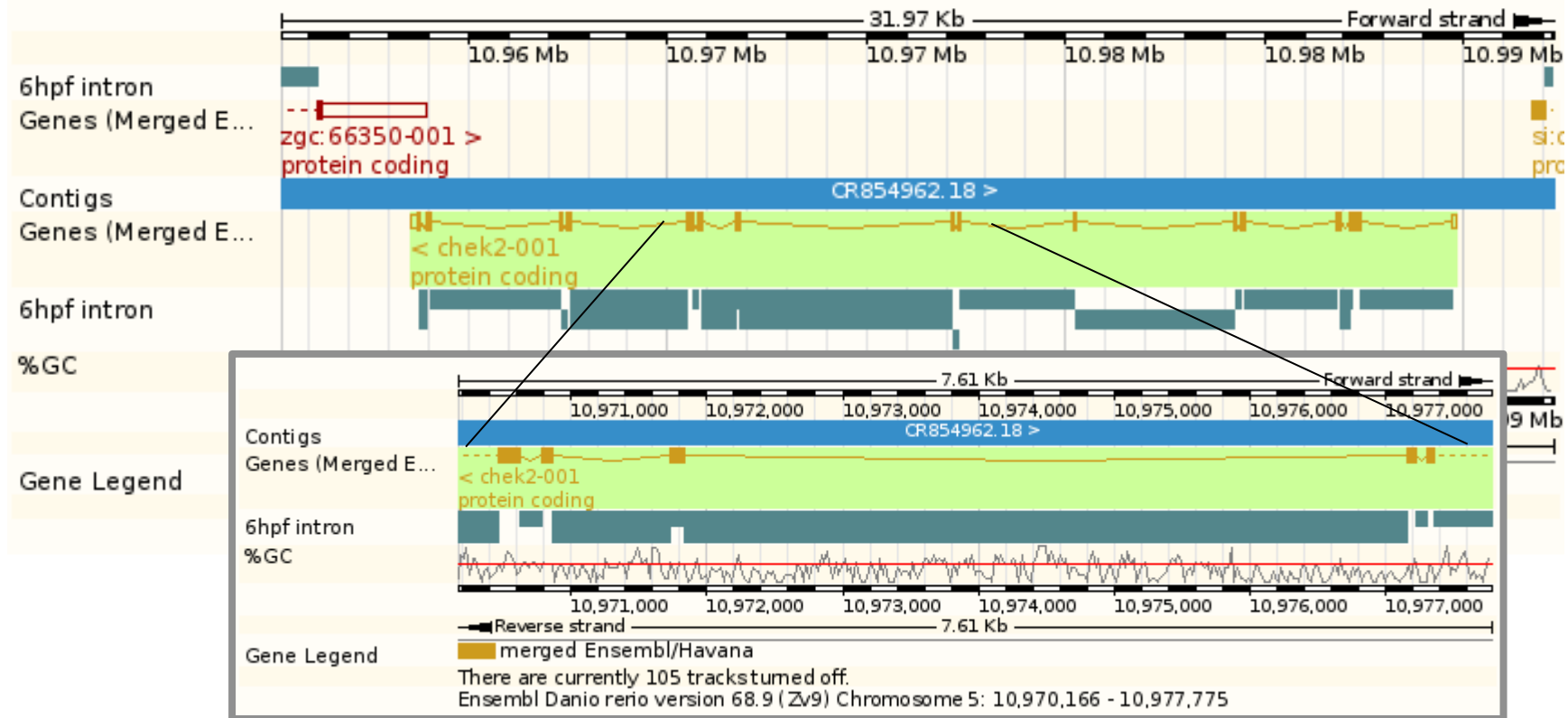Zebrafish gene build   http://www.ensembl.org/index.html release 68

# Building genes from RNAseq

Example of an RNAseq gene built from the short reads for apopotosis inhibitor 5 (api5)



Sequence mapped to introns

Sequence mapped to exons

RNAseq gene model

Zebrafish gene build   http://www.ensembl.org/index.html release 68

# Introns show alternative transcripts



chek2 – checkpoint kinase 2

Zebrafish gene build   http://www.ensembl.org/index.html release 68

# Introns show alternative transcripts



Zebrafish gene build   http://www.ensembl.org/index.html release 68

# Introns show alternative transcripts



Zebrafish gene build   http://www.ensembl.org/index.html release 68
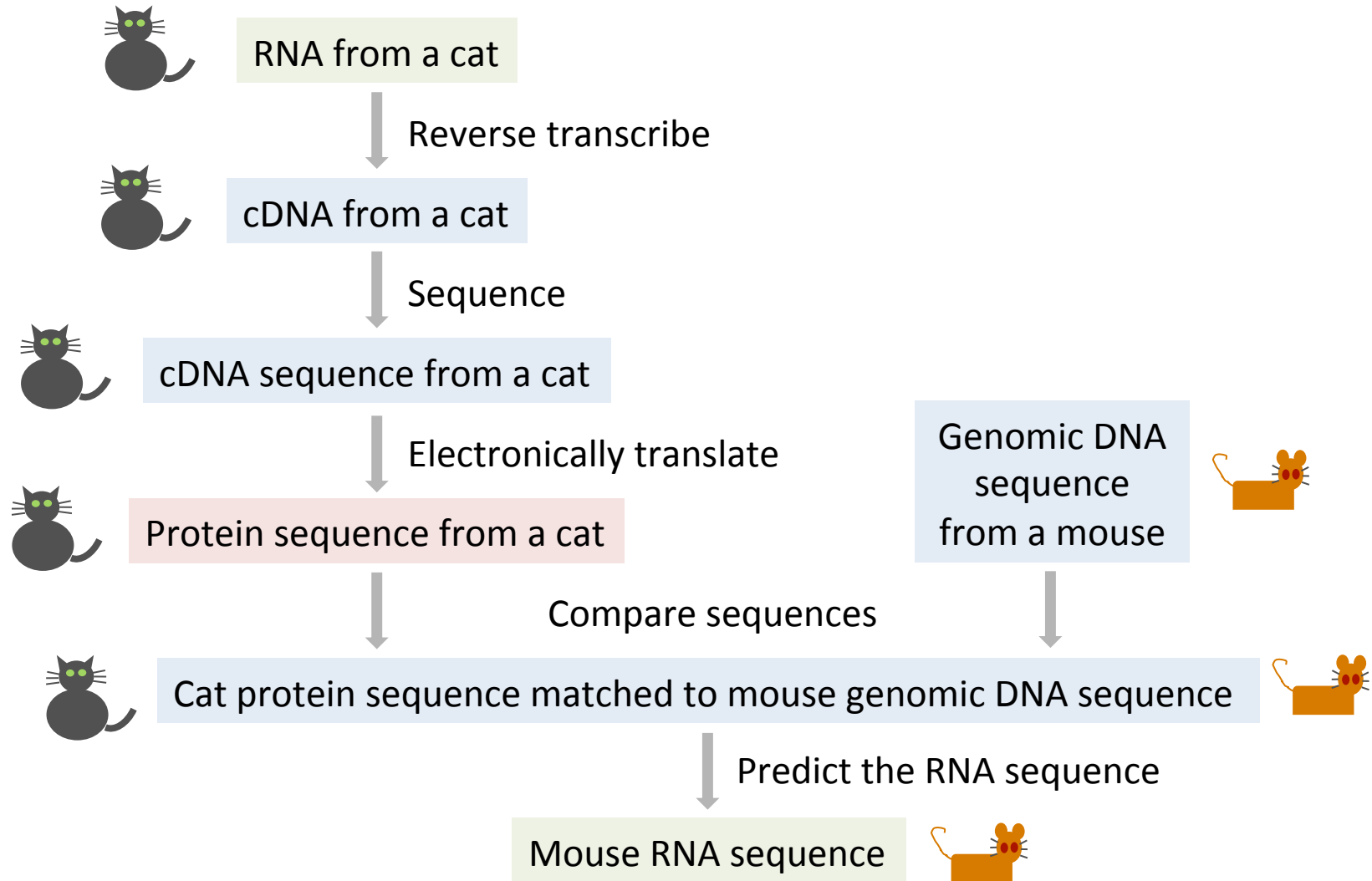
# Building genes from homologues

**paralogue**
Duplicate copies of a gene in the same species

**orthologue**
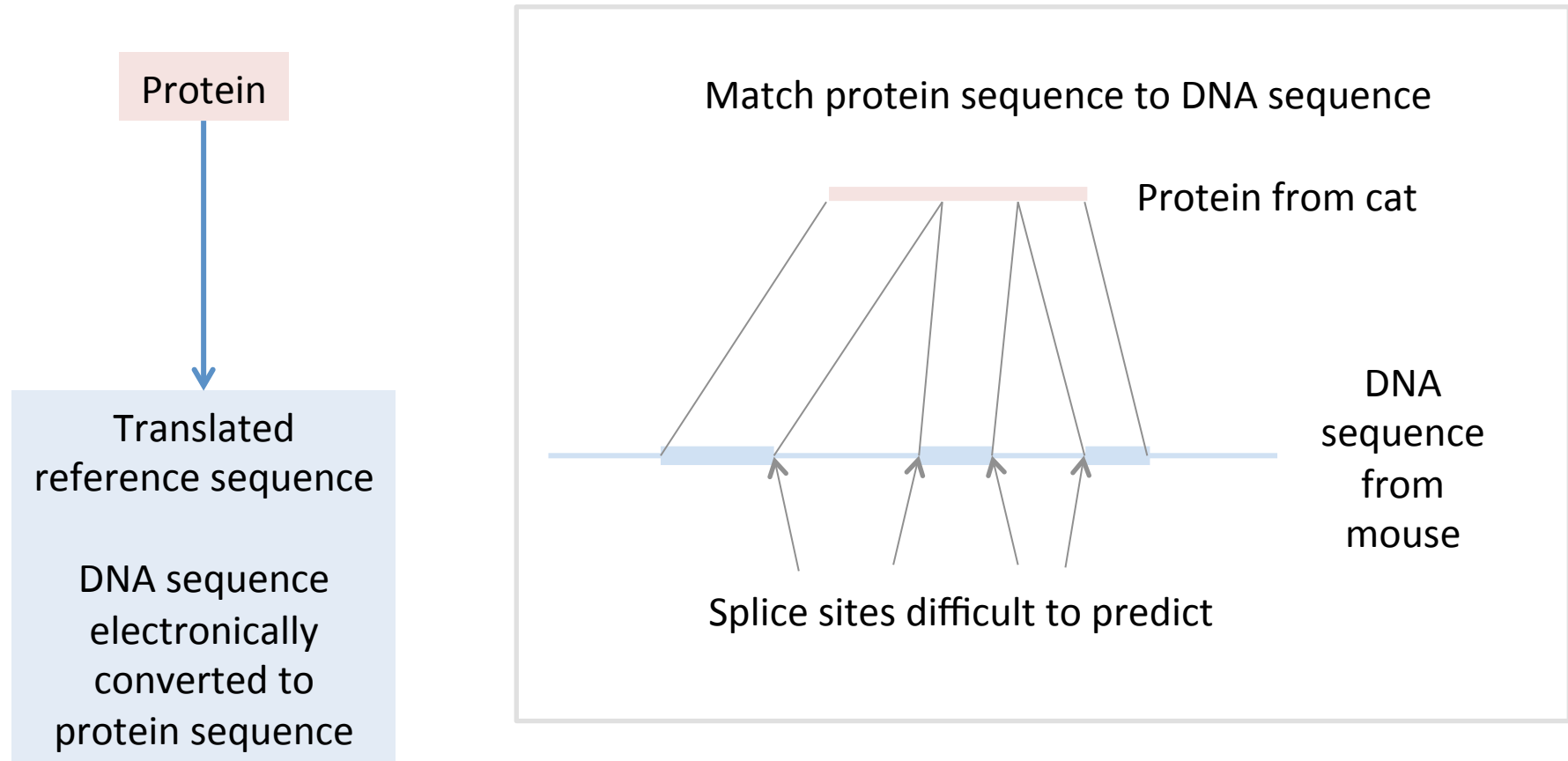Divergent copies of the same gene in different species

If we sequence a cat gene and identified the protein then we can predict the mouse protein using the mouse genome
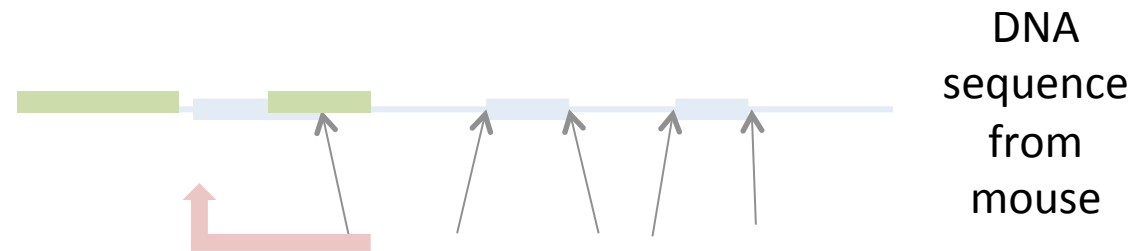
# Building genes from homologues

RNA from a cat

↓ Reverse transcribe

cDNA from a cat

↓ Sequence

cDNA sequence from a cat

↓ Electronically translate

Protein sequence from a cat

Genomic DNA sequence from a mouse

Compare sequences

Cat protein sequence matched to mouse genomic DNA sequence

↓ Predict the RNA sequence

Mouse RNA sequence

# Building genes from homologues

Protein

Translated reference sequence

DNA sequence electronically converted to protein sequence

Match protein sequence to DNA sequence

Protein from cat

DNA sequence from mouse

Splice sites difficult to predict

# Building genes from homologues
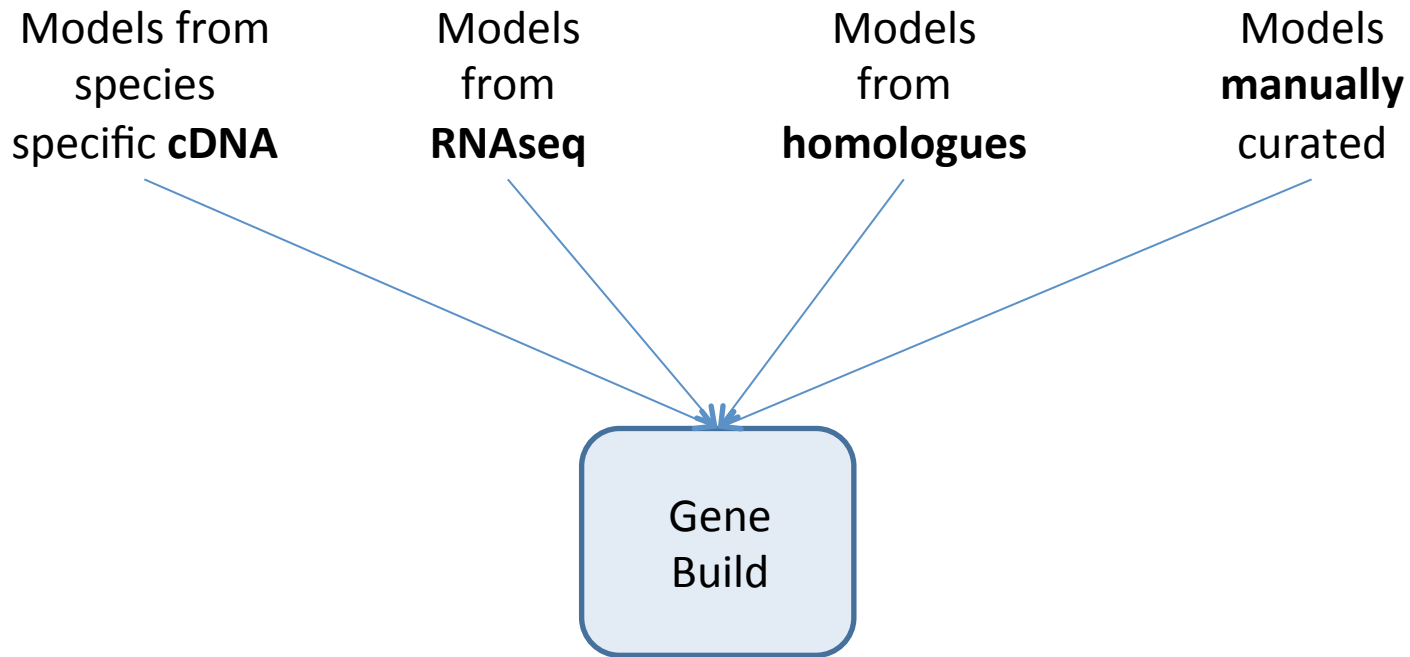
Intron driven annotation



DNA
sequence
from
mouse

Define the exact splice sites
using the mouse RNAseq
**intron database**

# Building genes manually

VEGA gene annotation - Manually curated

- Look at all the evidence for each transcript
- Annotate all the alternative splices for each gene
- Solve problems which aren't possible in an automated system
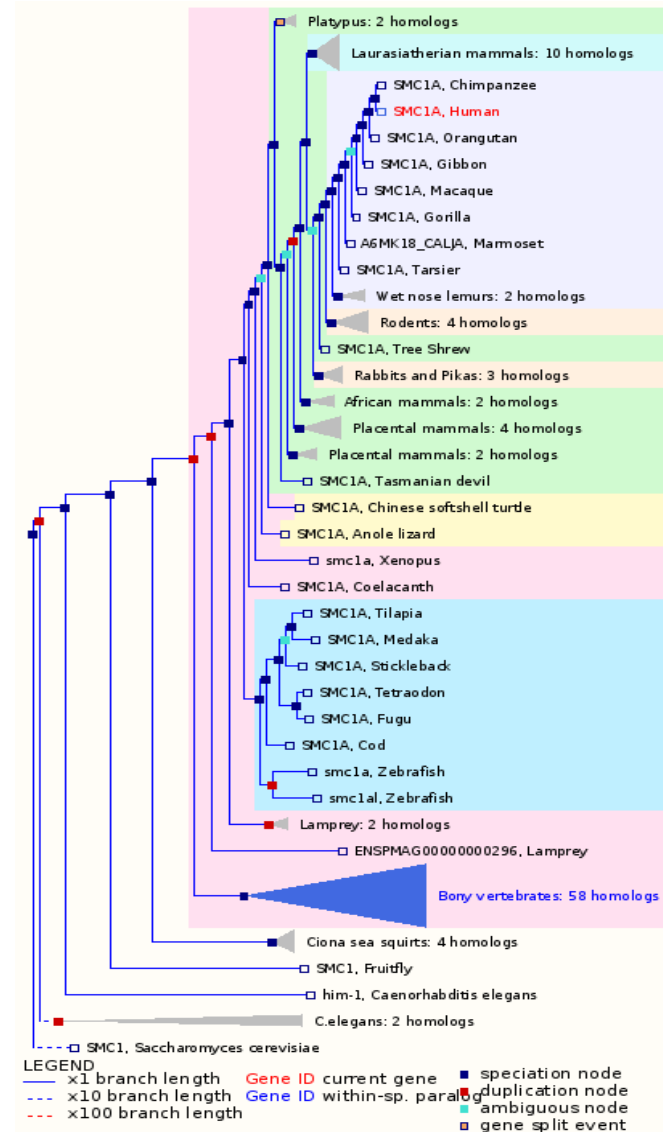
# Bring all the gene model together

Models from species specific **cDNA**

Models from **RNAseq**

Models from **homologues**

Models **manually** curated

Gene Build

# Displaying the gene build



Human gene build  http://www.ensembl.org/index.html release 68

# Comparing gene builds from other organisms

SMC1A is a protein important in cell division.

By matching gene builds across species we find candidate orthologues in other organisms. The position on the Gene Tree shows the degree of relationship
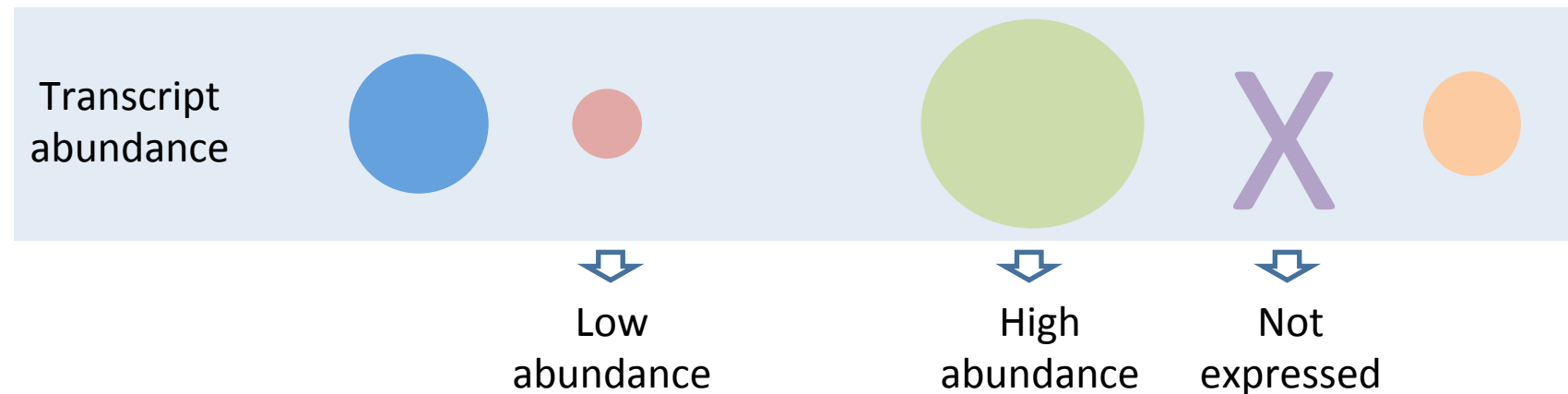
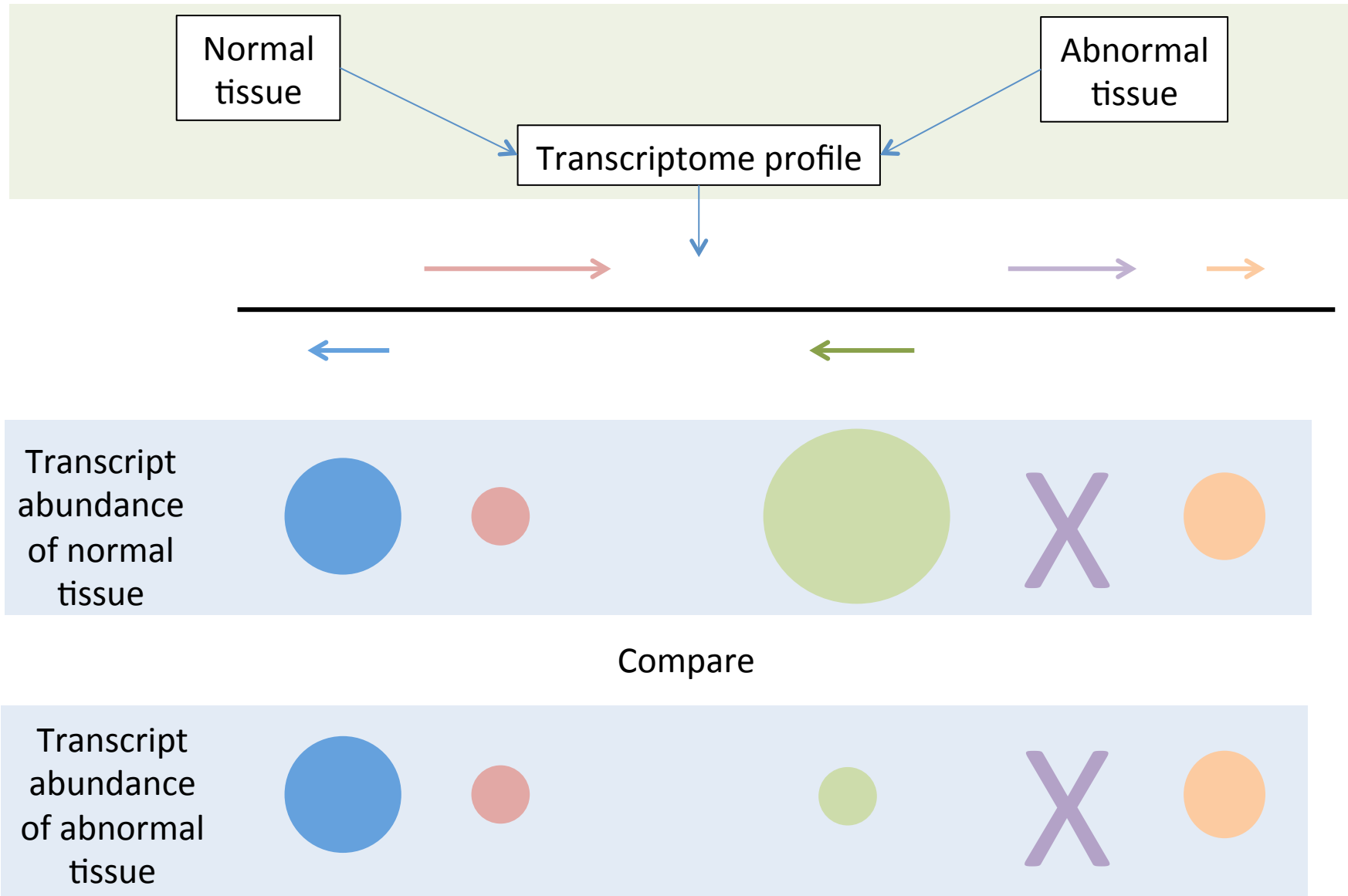Gene Tree in Ensembl 68
ENSGT00580000081569
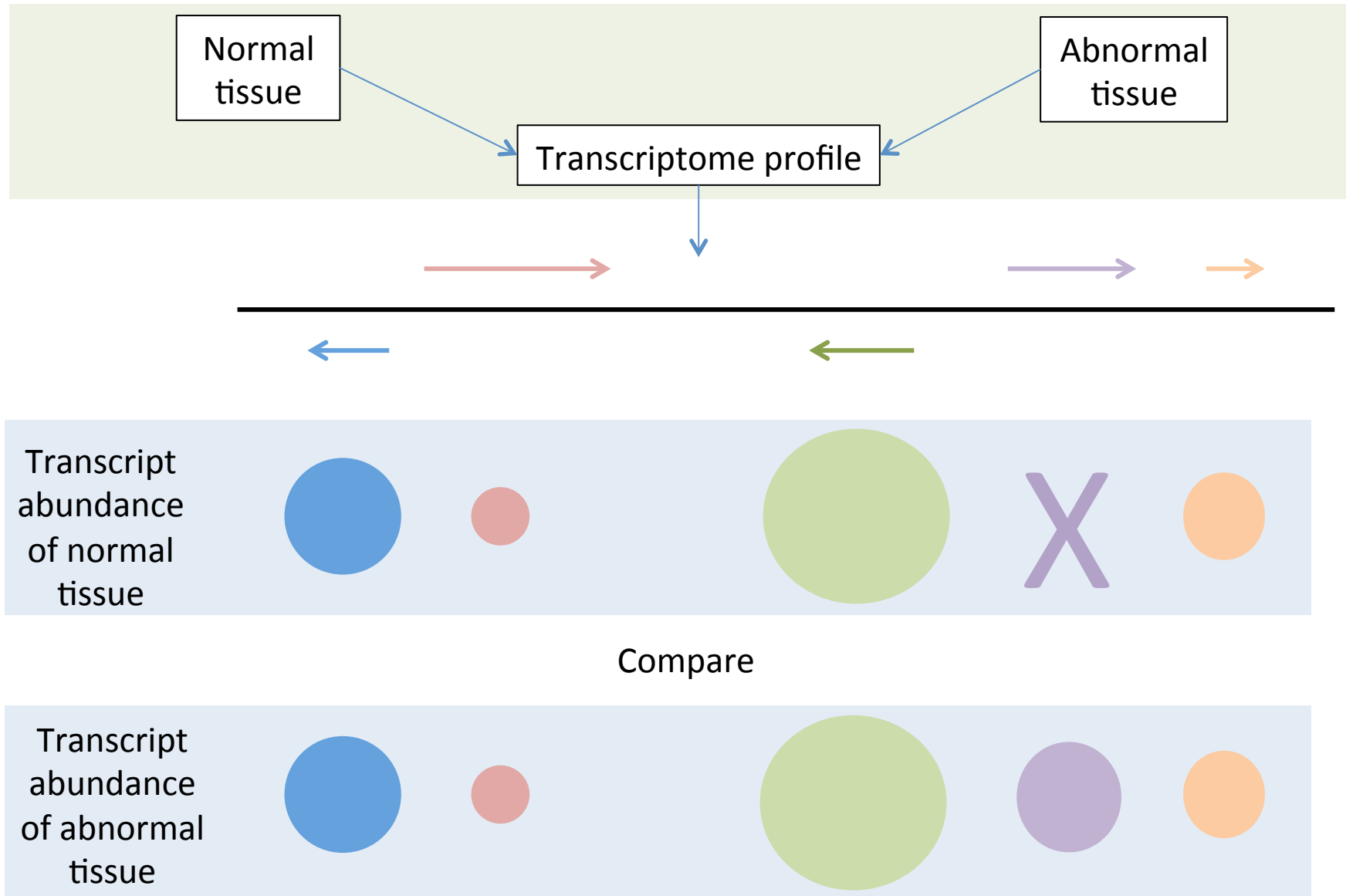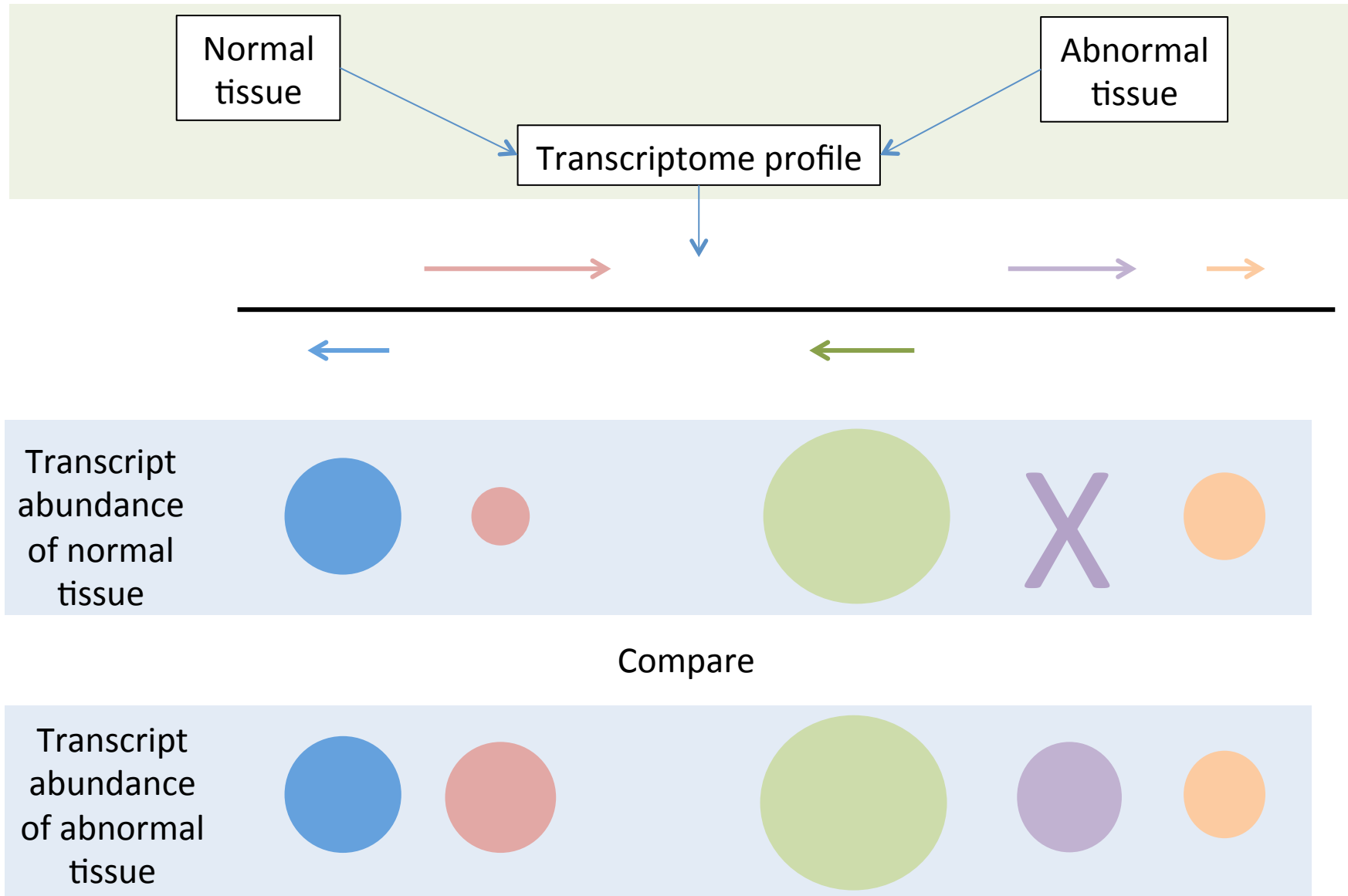
# Using the annotation – transcriptome profiling

# Using the annotation – transcriptome comparison

# Using the annotation – transcriptome comparison

# Summary

Identifying the 3% of the genome which contains the protein coding genes