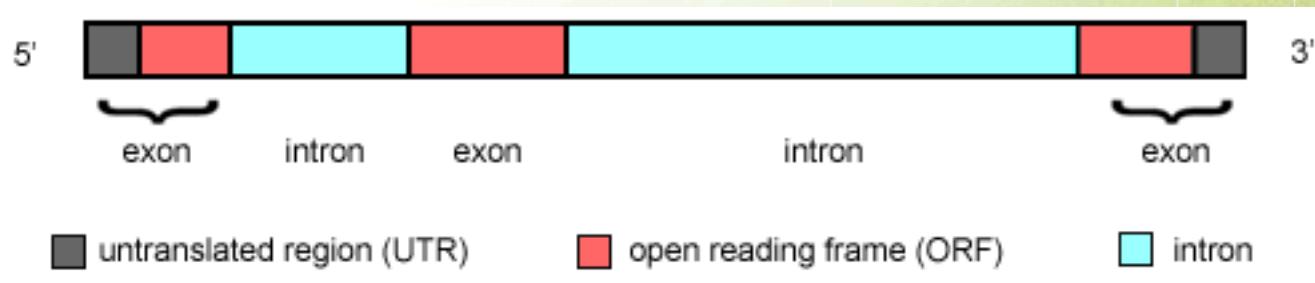


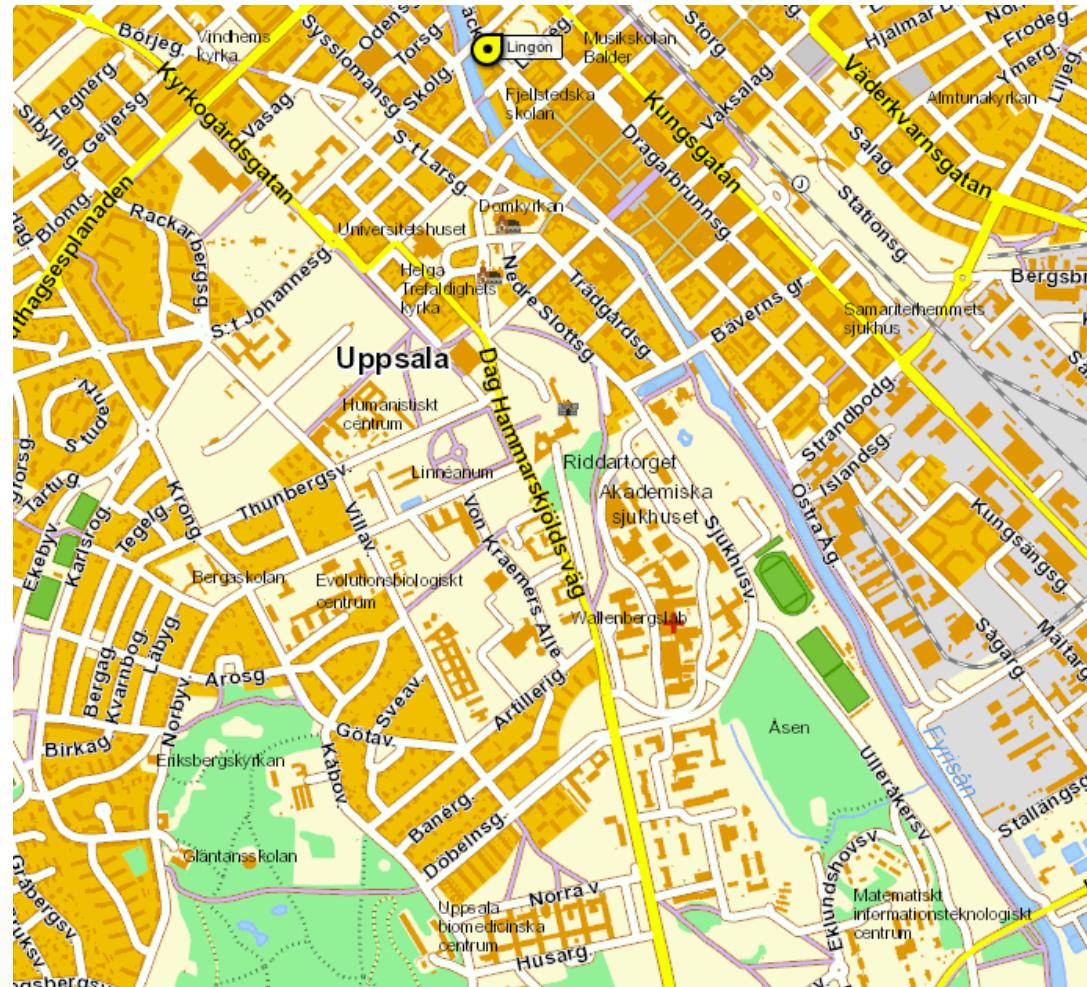
# Introduction to genome annotation - practical information



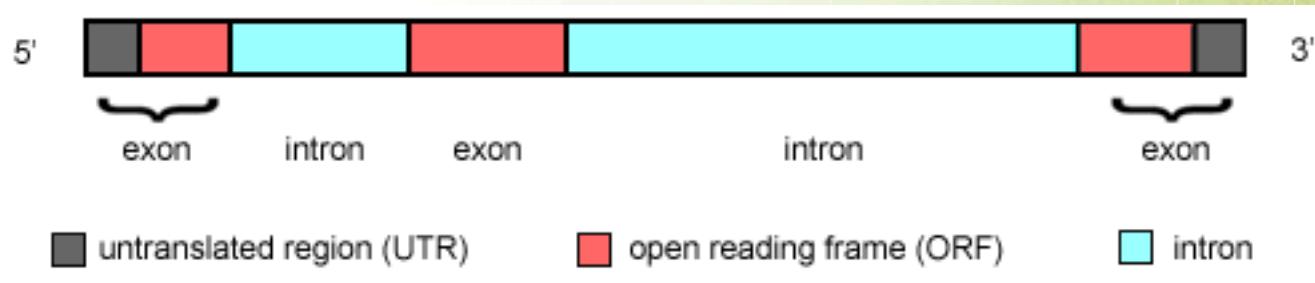
Enabler for Life Sciences

# Practical info

- Coffee breaks
- Lunch
- Dinner at  
Lingon 19.00  
Svartbäcksg. 30



# Understanding annotation



Henrik Lantz, BILS/SciLifeLab

Enabler for Life Sciences

# Lecture synopsis

- What is annotation?
- Structural genome annotation
- Types of data used
- Transcriptome annotation
- Functional annotation

# What is annotation?

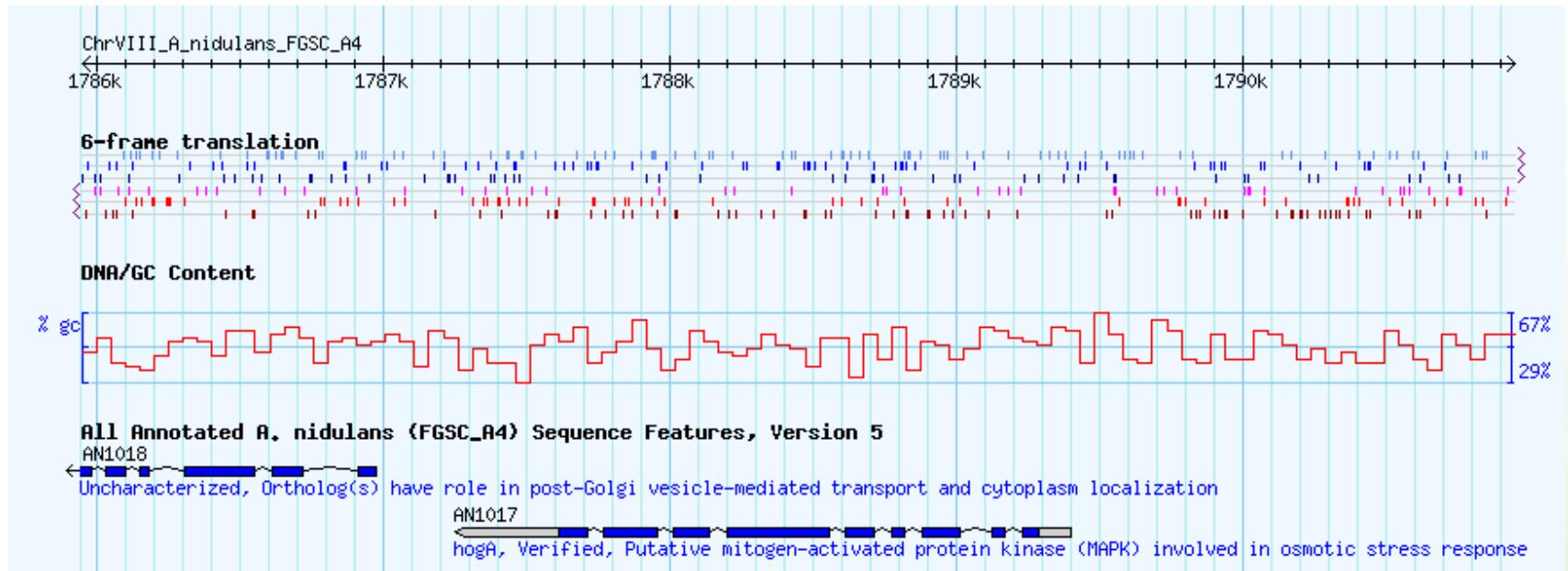
- Identification of regions of interest in sequence data

# From a genome...

```
> scaffold_26
AGTCACACACCCCTTACGCTTACACCCCTGACTGCAGCCCCCTAAACAA
TTCCAGCCAGGAAGATGCTCGCACACGGCTTGTGATGCCGCTCTGAC
GTCGAACGCCGCCGCCGCCGAAAATCTGGCAGCGTGGTACCGCGAGAT
CCGAAGCCGCCCTGGGGACCTCGGAGACACAGGGAGGGGTCAACGAGAC
GCCGAGGGCTGGAGTTTCCCACACGGGCCGGTAAGTTTCTACCA
AAAACCCATAAGAAAGAGATGAAACCTAAAGTTGTAACTCTCTACTT
AACCGTGACCTAATGTGCGGGGCAGGGCAAGCTTGACCCTAAGGGCAC
AGCAACAGGTGGTGGCCAAATAAACAAAGATGTAAGGGCTTG
AATAAAATCTCGGAAGATAATTCTCGAGCCGACACGGCTTGAGGCAGC
GGAACCTACAGAACACCGCAGTCAGTGAAGAGCTAATCTCTCCA
AAGAGAACTCAAGGGAAATGGAGGGTCAAAGAGGTCCTTACAAAGC
GAGAAAGGAAGATGGATGAGACATCTTGATCTCTCTGCTCAA
AGCAAAATGTAAGGATGCCAGACTAAGCCGAGTCTGAGAAAGTACGCCA
GCAGAGACCCCCGCTGCCGTGCCCCAGCACAGTGGCATAAAGGCC
GAGACATAACAAAGGCCCTGTGACACAAAGACGATGACACAAACTACAT
AACACAGACAAACAATAATGACACAGAGAGAAGTGAACACTCTGGGA
AGTAACATTTCTGAACACATCTACCAACAACTCGCTCATATATTTCCA
TTCCACGGGACTTCTGTGTTGTATATGCGTGTAAACGTAATCCCGCT
GTAGCAATACCAACTACGATAATTCTTGGAGGTTGCTGAGT
ATCATCTTATCAGTCTTATTTCCTGGCTCTGGCTCGGCTTCTTT
TTTTCTCTGATAAGATTTCCAGGAATGTGAAGACCCCCGATCT
TCCAAACTGACACCCAAATAAGACATCTATAGCATACATTACAC
AACCTAGGCAAAGTTCTAACATTAAAGGAACATGAAAAAAAGCCACAT
CACAAATATTCTATAAACATTAACTGGAAATGCAAAAGCCAAATCACAG
TACATTATAAACATACTCTCCCTTTCTTAAAGAGATCATATGCT
TGACCCGCCCTCTGCCGGGCCACCGCTGAGTACTGCCGTCGCGAGTC
ACGGAGCAGTCCCCGGGCCACGGCCTCTGCCGGGCCACGGAG
ATCGGCTGCCACTCCGAGCTGCCGTGCCATGCCGGCCCCCCGCG
GGGTCCCCGGCN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN
NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN
NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>NN>

```

# ...to an annotated gene





# GFF3 file format

<b>Seqid</b>	<b>source</b>	<b>type</b>	<b>start</b>	<b>end</b>	<b>score</b>	<b>strand</b>	<b>phase</b>	<b>attributes</b>
Chr1	Snap	gene	234	3657	.	+	.	ID=gene1; Name=Snap1;
Chr1	Snap	mRNA	234	3657	.	+	.	ID=gene1.m1; Parent=gene1;
Chr1	Snap	exon	234	1543	.	+	.	ID=gene1.m1.exon1; Parent=gene1.m1;
Chr1	Snap	CDS	577	1543	.	+	0	ID=gene1.m1.CDS; Parent=gene1.m1;
Chr1	Snap	exon	1822	2674	.	+	.	ID=gene1.m1.exon2; Parent=gene1.m1;
Chr1	Snap	CDS	1822	2674	.	+	2	ID=gene1.m1.CDS; Parent=gene1.m1;
		start_codon						Alias, note, ontology_term ...
		stop_codon						



# GTF file format

<b>Seqid</b>	<b>source</b>	<b>type</b>	<b>start</b>	<b>end</b>	<b>score</b>	<b>strand</b>	<b>phase</b>	<b>attributes</b>
Chr1	Snap	exon	234	1543	.	+	.	gene_id "gene1"; transcript_id "transcript1";
Chr1	Snap	CDS	577	1543	.	+	0	gene_id "gene1"; transcript_id "transcript1";
Chr1	Snap	exon	1822	2674	.	+	.	gene_id "gene1"; transcript_id "transcript1";
Chr1	Snap	CDS	1822	2674	.	+	2	gene_id "gene1"; transcript_id "transcript1";
		start_codon						
		stop_codon						

# Why is annotation important?

Example: Differential expression

Mapped reads - condition 1



Genome

Mapped reads - condition 2

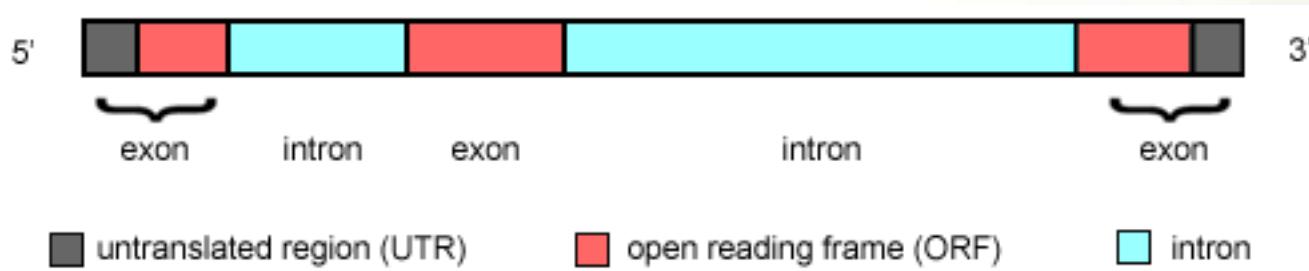
# Why is annotation important?

RNA-seq reads



# There are two major parts of annotation

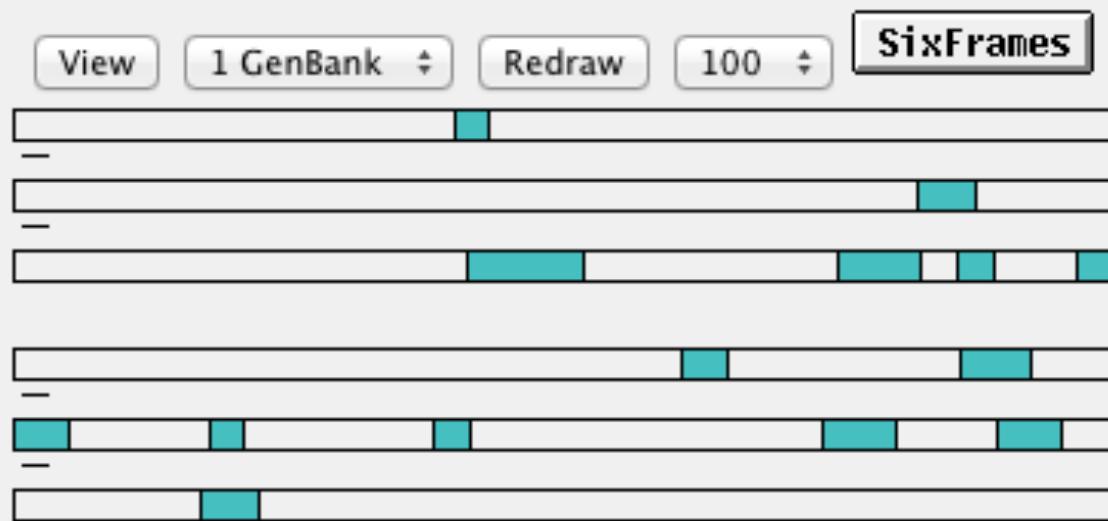
- 1) Structural: Find out where the regions of interest (usually genes) are in the genome and what they look like. How many exons/introns? UTRs? Isoforms?



- 2) Functional: Find out what the regions do. What do they code for?

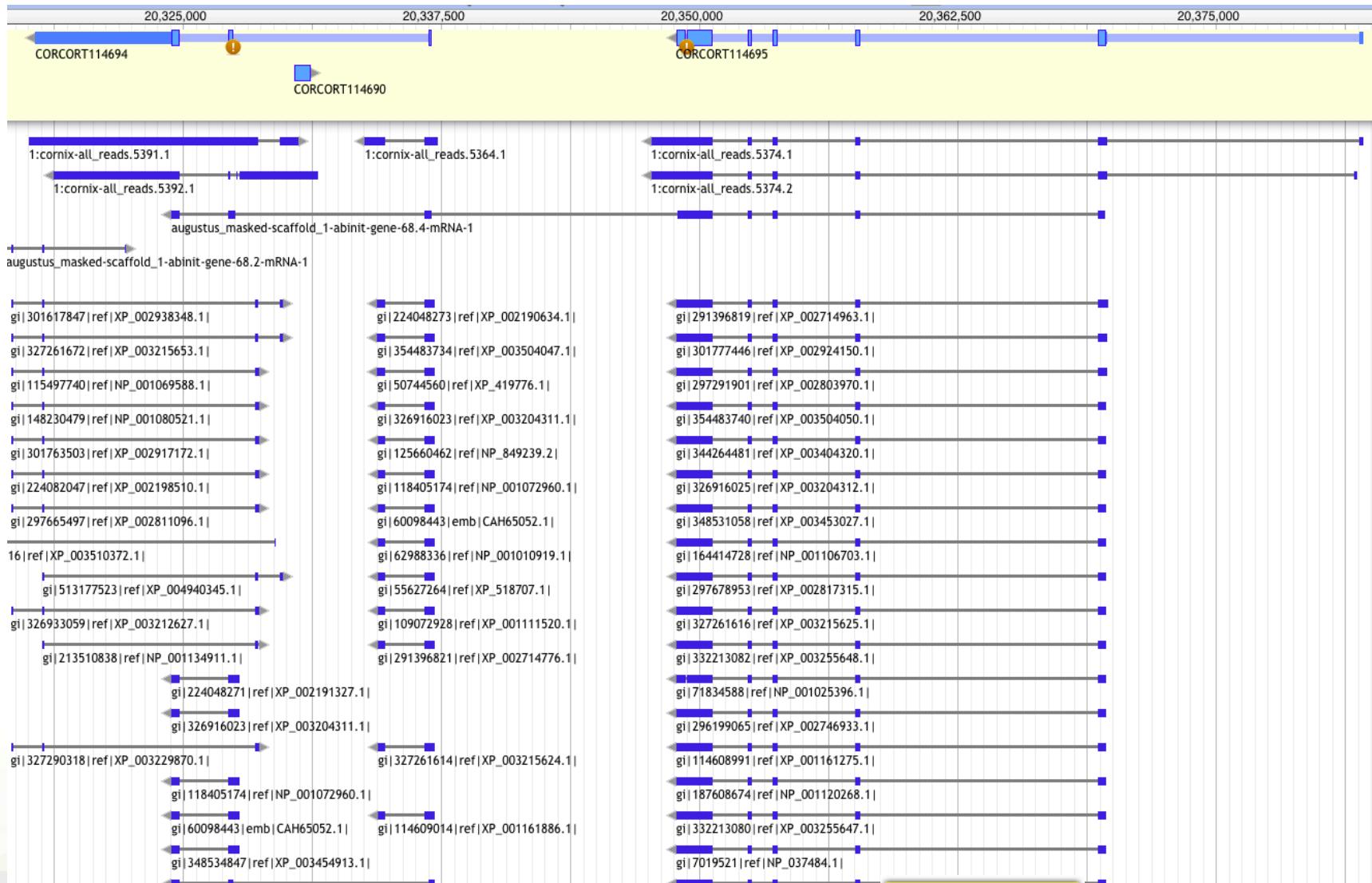
# Open reading frames

## Anonymous



Frame	from	to	Length
+3	1383..1736	354	
+3	2511..2756	246	
-2	2465..2686	222	
-1	2880..3092	213	
-2	2996..3187	192	
-3	574.. 753	180	
+2	2753..2929	177	
-2	2.. 172	171	
-1	2034..2171	138	
+3	3237..3349	114	
+3	2874..2984	111	
-2	1283..1393	111	
+1	1345..1446	102	
-2	599.. 700	102	

# Difficult in practice



# Combine data - use Maker!

- External data - proteins, rna-seq (incl. ESTs)
- Ab-initio gene finders
- (Lift-overs from closely related genomes)



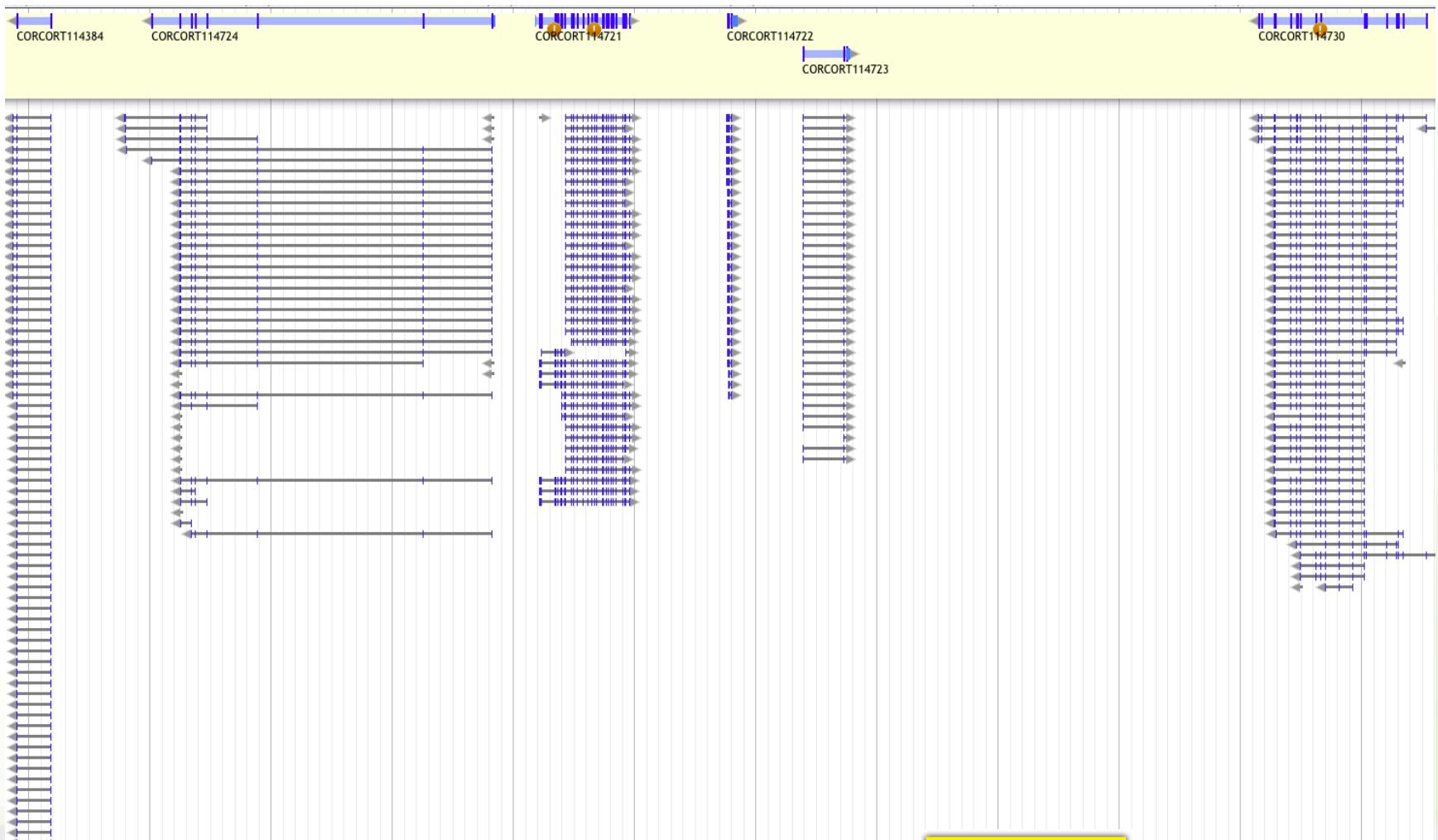
Combined annotation

# Transcriptomes are different but have their own challenges

- No introns, but where are the start and stop codons?
- Still needs functional annotation

```
>asmb1_2719
AGCACCTAGACGAGATGGGAGGTCTCTCTTGCTGTGCAGAGGCAGATCTCTTTCC
AACACCTAGCAGTATGAACACTAGTGAGCTCTGACTGTTTCCAGTGGTAATGAGGTGTGA
CCCGCTGCACTGCAACACTGAATTCTTCAGTCCCCGAGGCCAGCCCAGCTGTGGGC
AATGCTTGTGTTGTGCTGTGACCATTC
>asmb1_2702
GTCCTGACTGGGAATGCCCTGGAGCAGAACCTTGCATGGATAAGGACACTACATT
CTGGTGTAAAGGTGAATATAACCTTCAAGGTTAAAGGTGACATTAAATTCAATTACAGCT
TGCTCTTGTAAAGCTAACAGTTAACTAACAAAGCTATACTGTGACTACACCCCTTAGATCA
ATAGCTGGGAAACATACCTCCCCAAATACTCCACCTTAACTGCACTCTTGAAAG
AAGTACAGGCCAGGTTAGCTGATCCATCTGTGCTAATGCTGCTTACAAGCTG
CAATATTTTAAACACCAGAACATTGGTAGAGGTTAACATCAGCAAGCTTCAATT
TACAGCAGGTTAACCTCTGAACACTGTGATCACTGATATTTGGCTAGTCAGATGT
CTTGTAGTGTCTT
>asmb1_2701
ACAAACAAAACAAAAAAACAAAGGAACAAAGCAAAAAAAACATCATACAATCCCAGT
TGTCAGAGCTTACTGTGAAATCAATCTGGACTAAACAAATGAAAAGCTTCCAGA
TTCTGTATTCCAGGCTGAGACAAGTTGTTAAACTCTCCAGAAATTGCAACAGCTG
CAGGGTAACATCTTAATGCAACCTCTGTACAGAAATGCGAGACCTTAACTCTT
CAGGCCCTCCCCAGTCAACACCCAGTATAAACTCTGCCCTTCACTTGTGGAATATCTA
TCATAAGGGAAAGCATTTTTAGGCTGAGAAATACAATCCACCTTGTGAGGCCGGTCA
GCATATCATGGCTATGCTGTGATAGGTTGACCAAGCACTCTGTGAGAAATA
CTTAGAGGTGACCTAACAGGTCATTTGCAACTAACCTGGTGGCAGTATCTTTA
TTCCAACCTCCCACCTTCCCCAGAGAACAAAGCTGTATTGCGAGTAGCAGTGTGTTG
AAGGTAAACCTGCACTGTACTAGTAGCTGCGAGGCACAACCTTCAACACTAGGCC
CTAGTCTAAAGTAACCTCTTGCACAGGAAGAACATGGAACACACAGGCCACACTTGCAG
AGGATCTGAGGCTGAGCTGCTTCCAGGAGGCCATGGGTGAGGCCAGTACAGAAGG
GCACTGAGGCTGCTTCAACCAAGCTGGCACAGAGGGCTGATCAGGAAA
ACCCACCATCAGCACAAAAGGAGCCCTGCAAACTCAGCCAGTGTAGGTTACTGGGTG
GAAATCAACTCTGCTGTGAAAGCTCTGTATACCCACATTACCTCATCCAGTGC
CAGAACACAAAGAGAGGAAATGGGAGGAGCAGGAATGTGCACTGAGGAGCAGGG
CATCTTGTGCTCAGGCTGCTGCACTGAGCTTCACTGGCAAGGCAGTCTGAGTCC
TACCGGTGAGGACATCTGTCATGTAATGCCCCCTCTGTGACAGCAATCTCAAG
TGGTCTTAAATGGCTCTCACTTGGGAGCTCACTGGCACAGCTACTGCCA
GGAAACAAAGGCTCAACAGGGTGGGAAACAAACTATTGTCAGTGTGAGGAAATG
TGGATAGAAACAGGGTGTGTAACCTGACTGATAAGAAGAGGAGGAGTGCAGATAGAG
CTGAGAACAGTACTCAGGTTGGAAACAAGCTGTATAAAAGCTTAAGAGGGTGAATGT
GAAAAAAATATGCCGAGCAGAATAAAAGGACTCTATTCCATCCCACCTGGAATCTGA
ACCCAGTTCAGAGTAATGAAGGGCTTTGTGTTGTGCTGTGAGGAGATCACCAGT
AAGCAATGCTCAGGCTCAGCTGCCAGTCTCAGGAAAGGACTCACCCCTGAGA
GGTTGATGTGCTGCAACAGGCCACAGCTGGCATTTGGAGGTGTTAGGTGCTT
GGGTATGCTAAACCAACGTTGAATGAAAGTGTCTGTCACTCACTGAGTGAAGGGAGA
GAACATTTCAGGAGCTGGAAACTCTGGAGGGACTGACTTAATTGTTGAGTCAA
GGTCTGTGCTGGAAATACTTCAGTGGCAATGCTCTGGAGGTGTTGGGGGGAGGCGT
GCACTGAAGCCCTGAGCTGGGAGAGTACACAAAGATATGCAATGCTAGCACAAGC
CAGAACACTTGGCTGAGGCCAGGGTCTACTCTGGCTGGGAGCTGCTCCCTGTG
GTACAGTTGCACTCCCCCTGTGAGCTGCCAGCTGGAGGAGTGTCTCTGTG
CTGGAGGAGACTGCGCACACCTCTGGAGATTTGGAGCTGTAATGTTGAGGCTCTG
TGGATGGAGACTGGCTCTGTGTTGGGGAGCTTCTCTGGCCAGCAGACTGTTAG
CTGAGGCCAGGCCACTCAACTCAGCATAGTGTCCAGAAAGAAGTAAAGTCCCAT
CTGCTTCCCGTGTACCTAACCTGTCAGGCCATTAGTGTCCAGAAAGTAAAGTCCCAT
CTGCTTCCCGTGTACCTAACCTGTCAGGCCATTAGTGTCCAGAAAGTAAAGTCCCAT
crow_gonads.assemblies.fasta
```

# Data used - Proteins



# Data used - Proteins

- Conserved in sequence => conserved annotation with little noise
- Proteins from model organisms often used => bias?
- Proteins can be incomplete => problems as many annotation procedures are heavily dependent on protein alignments

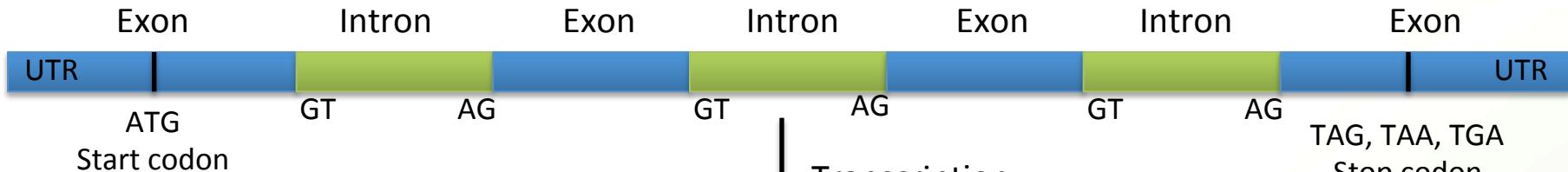
```
>ENSTGUP00000017616 pep:novel chromosome:taeGut3.2.4:8_random:2849599:2959678:-1 gene:ENSTGUG00000017338 transcript:ENSTGUT0000001801  
RSPNATEYNWHHLRYPKIPERLNPPAAAGPALSTAEGWMLPWNGQHPLLARAPGKGRER  
DGKELIKPKTFKFTLKKKKKKKKTFK  
>ENSTGUP00000017615 pep:novel chromosome:taeGut3.2.4:23_random:205321:209117:1 gene:ENSTGUG00000017337 transcript:ENSTGUT00000018017  
PDLRELVLMFEHLHRVRNGGFRNSEVKWPDRSPPPYHSFTPQAQKSFSLAGCSGESTKMG  
IKERMRLSSQRQGSRGRRQQHLGPPPLHRSPSPEDVAEATSPTKVQKSWSFNDRTRFRASL  
RLKPRIPAEGDCPPEDSGEERSSPCDLTFEDIMPRAVKTIRAVRILKFVAKRKFKETLR  
PYDVKDVIQEYQYSAGHLDMLGRIKSLQTRVEQIVGRDRALPADKKVREKGEKPALEAELVD  
ELSMMGRVVVKVERQVQSIEHKDLLLGLYSRCLRKGANSVLAAVRVPPGEPDVTSDYQ  
SPVEHEDISTSAQSLISRLASTNMD
```

# Data used - Proteins

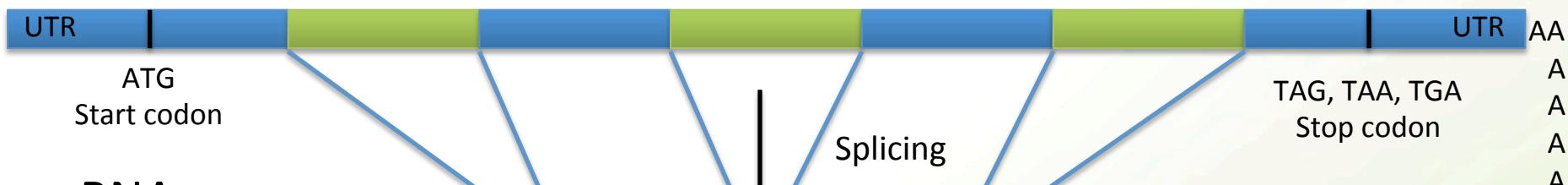
- Maker will align proteins for you: Blast -> Exonerate
- Blast is not structure aware, Exonerate is (splice sites, start/stop codons)
- Preferred file-format: fasta

# RNA-seq

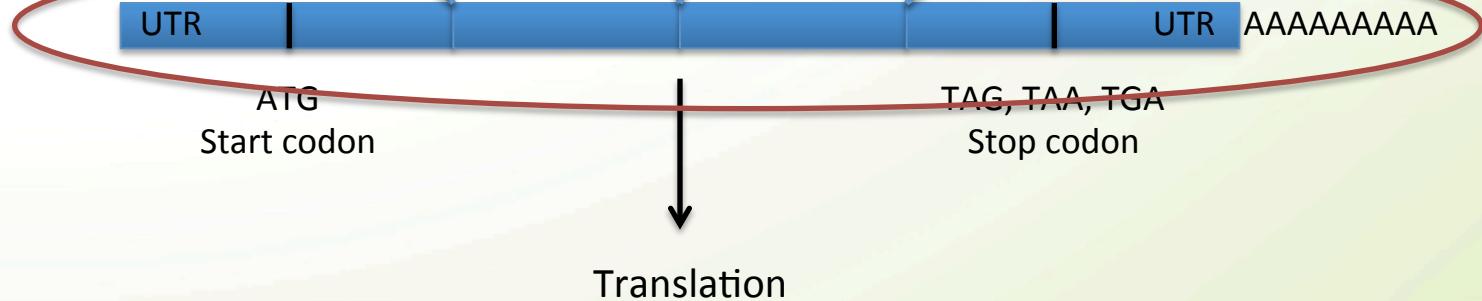
## DNA



## Pre-mRNA



## mRNA



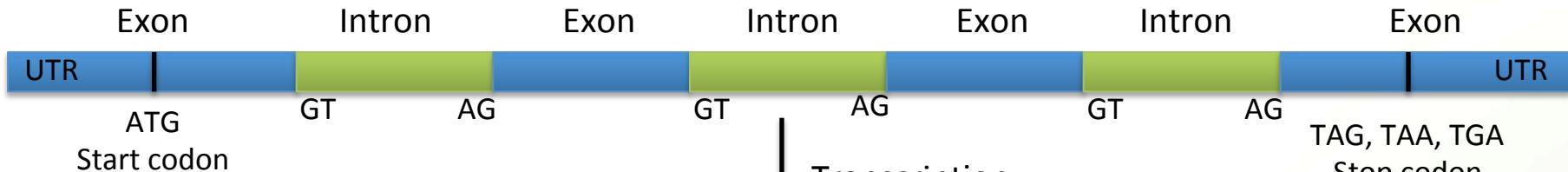
Translation

## Data used - RNA-seq

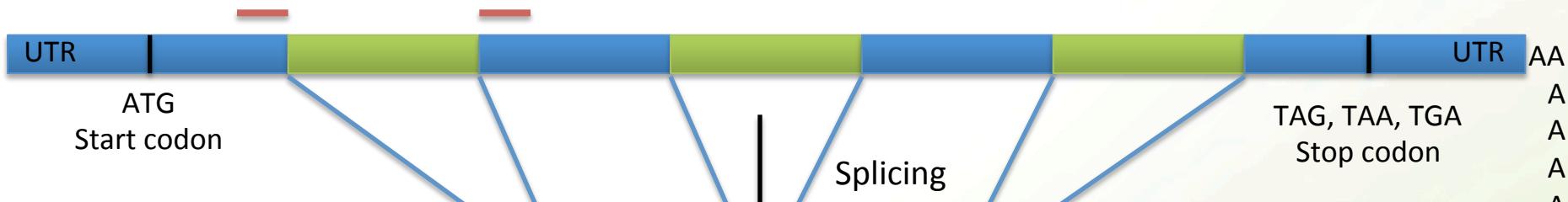
- Should always be included in an annotation project
- From the same organism as the genomic data  
=> unbiased
- Can be very noisy (tissue/species dependent), can include pre-mRNA
- PASA, or some other filtering method, often needed

# Spliced reads

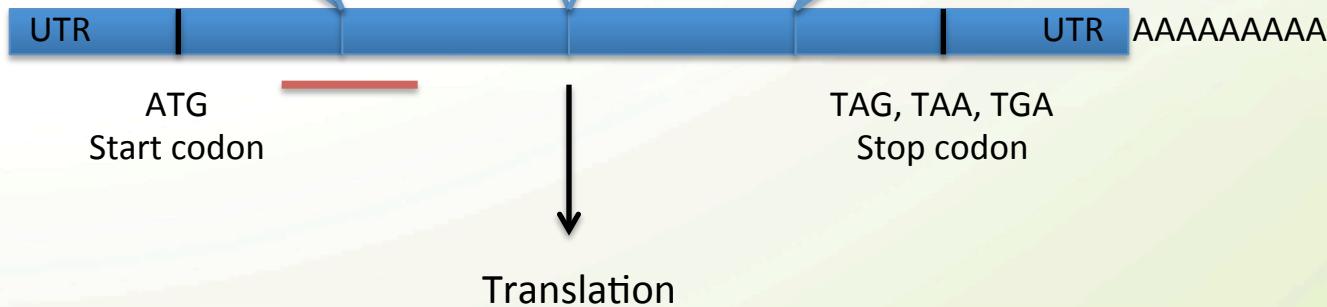
DNA



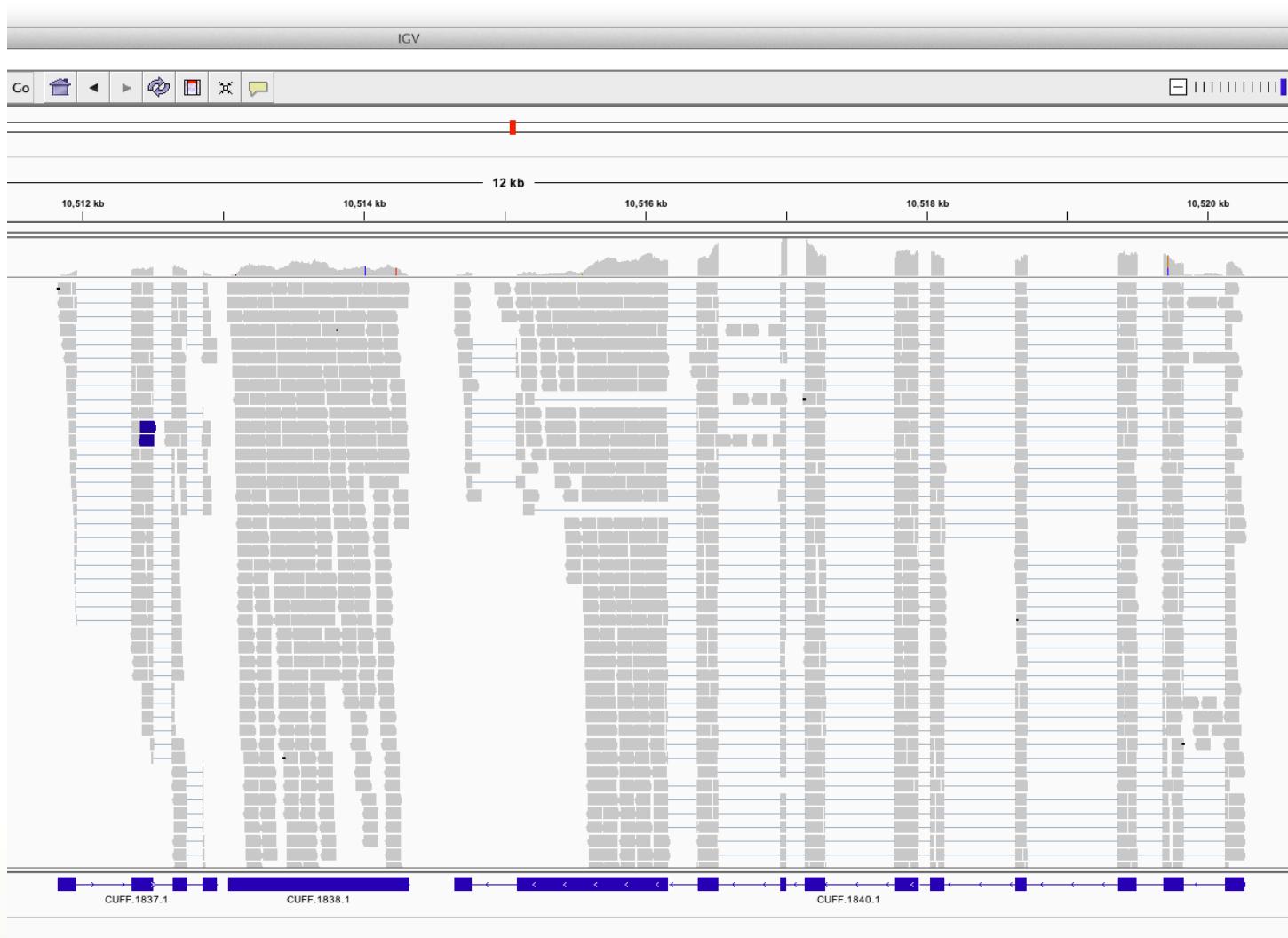
Pre-mRNA



mRNA



# RNA-seq - Spliced reads

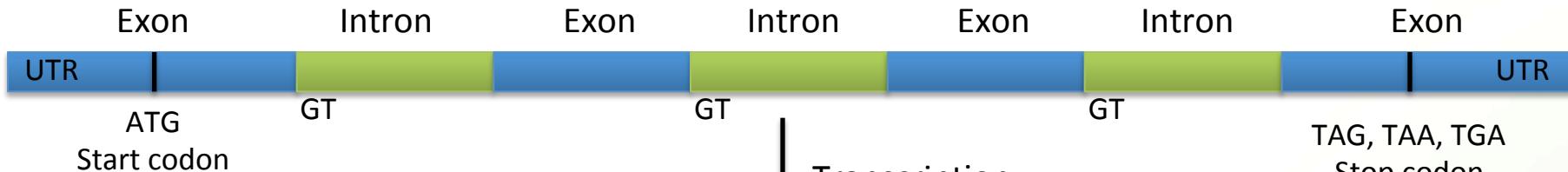


# Pre-mRNA



# Pre-mRNA

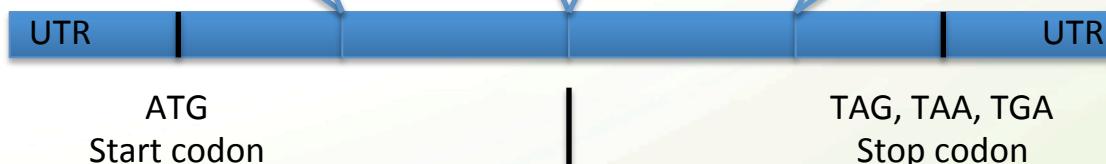
## DNA



## Pre-mRNA

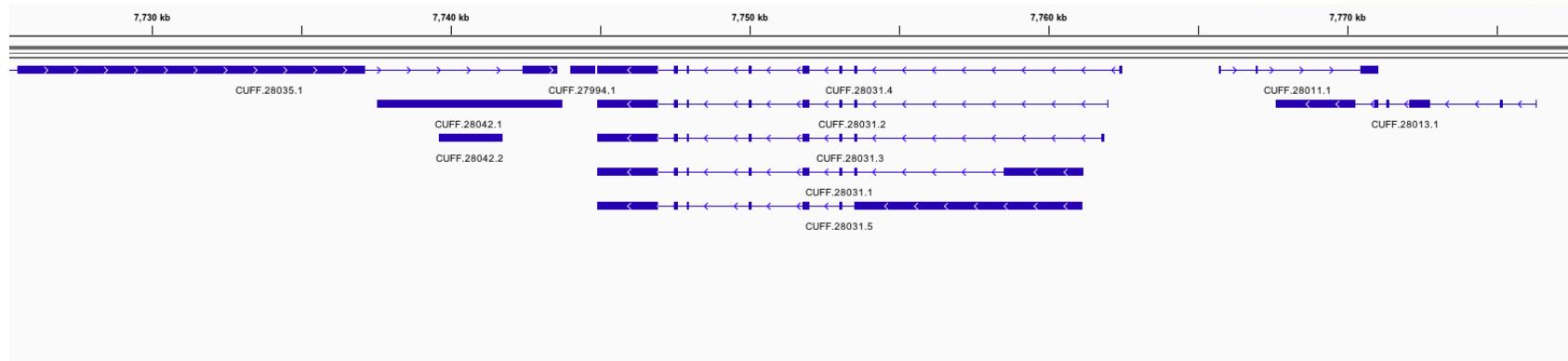


## mRNA

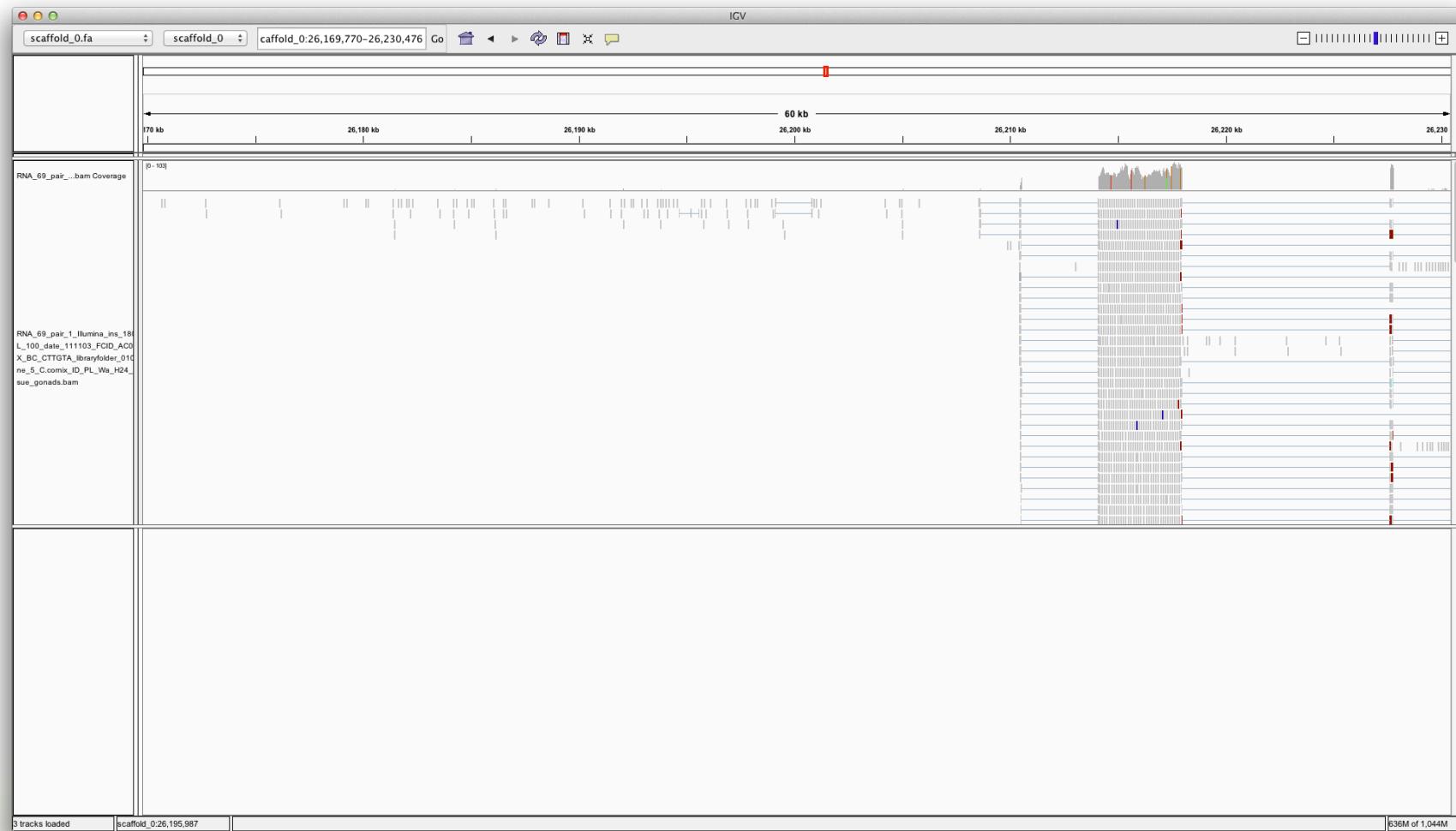


Translation

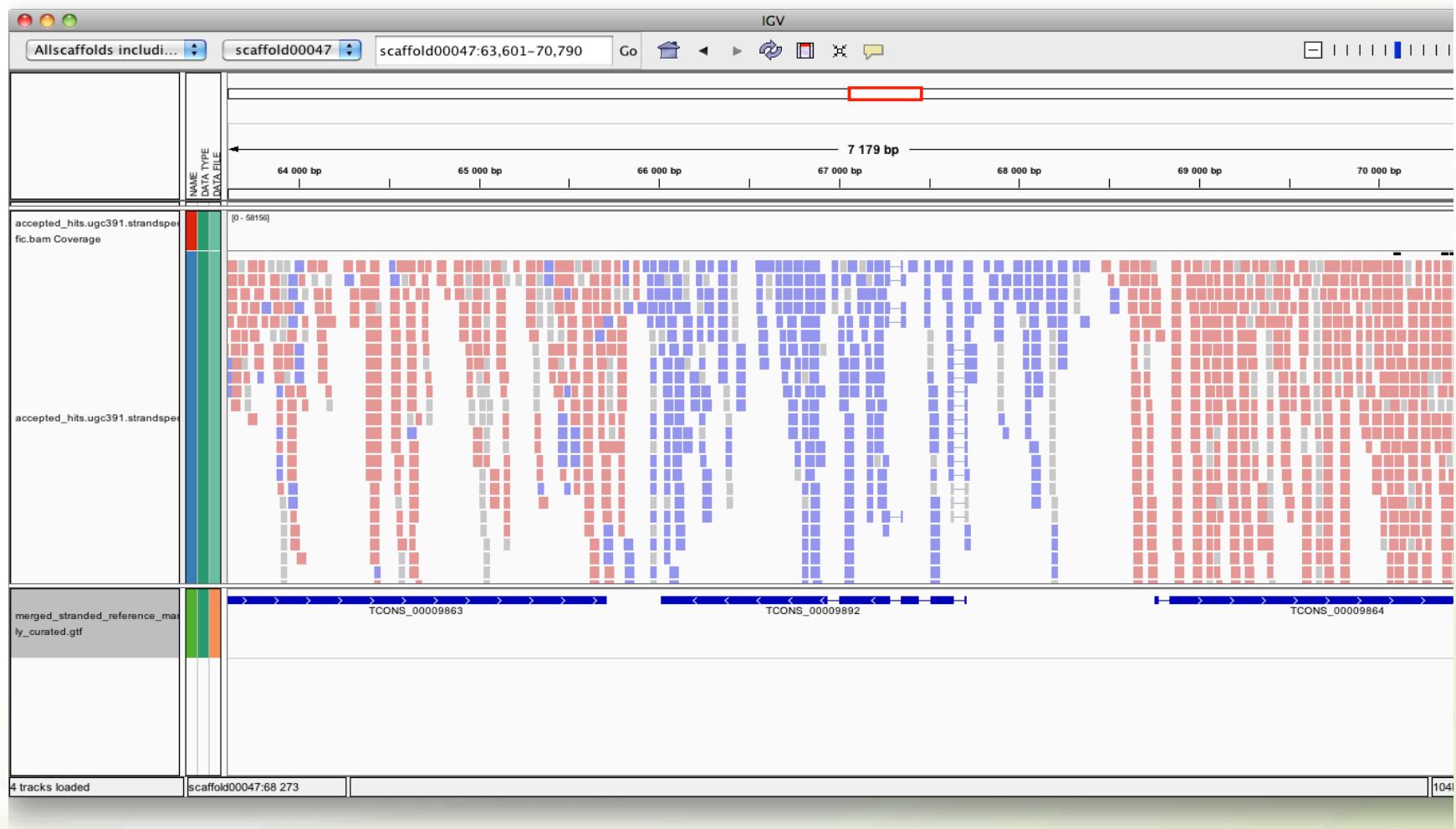
# Pre-mRNA



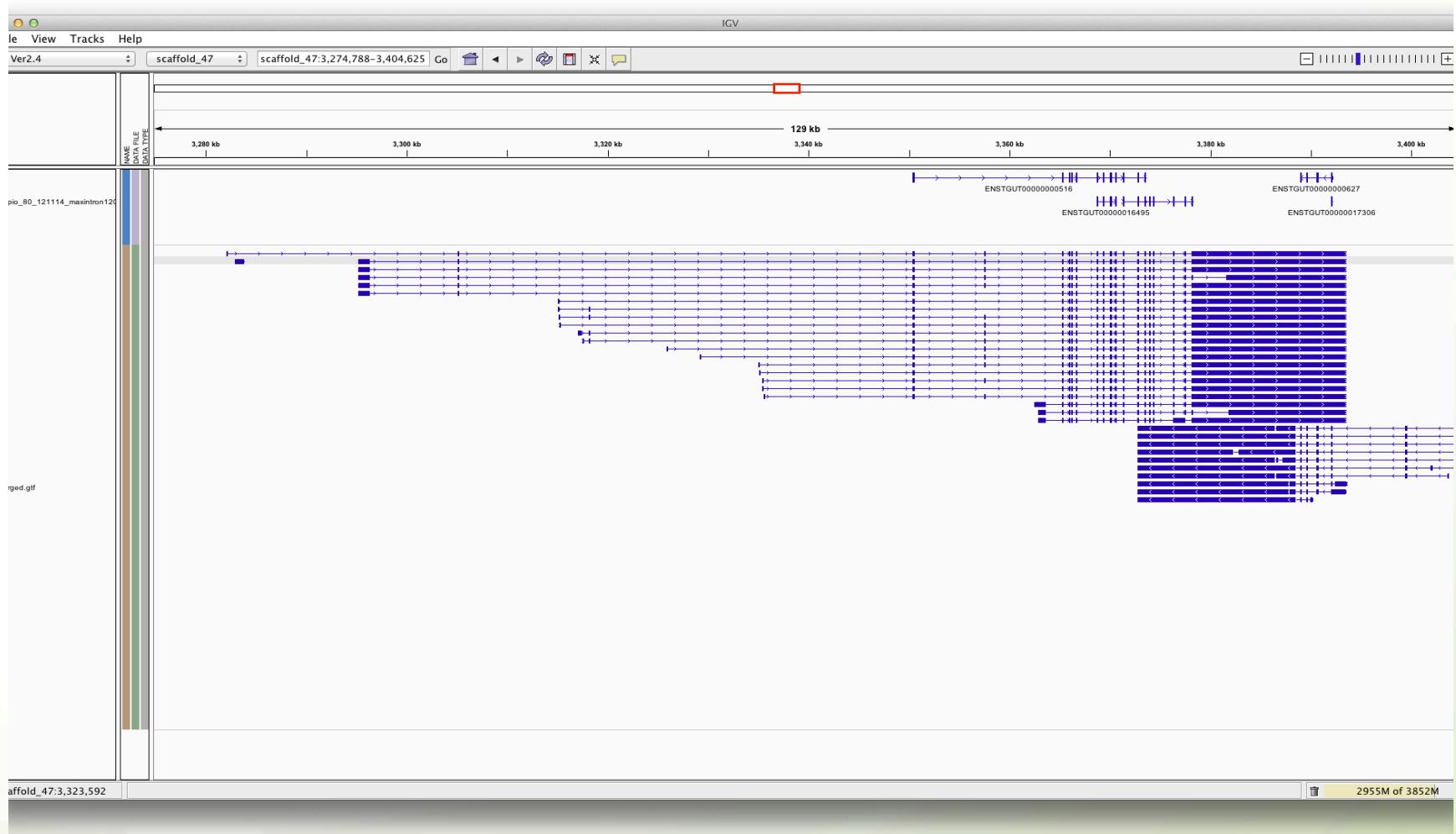
# Includes everything that is transcribed



# Stranded rna-seq

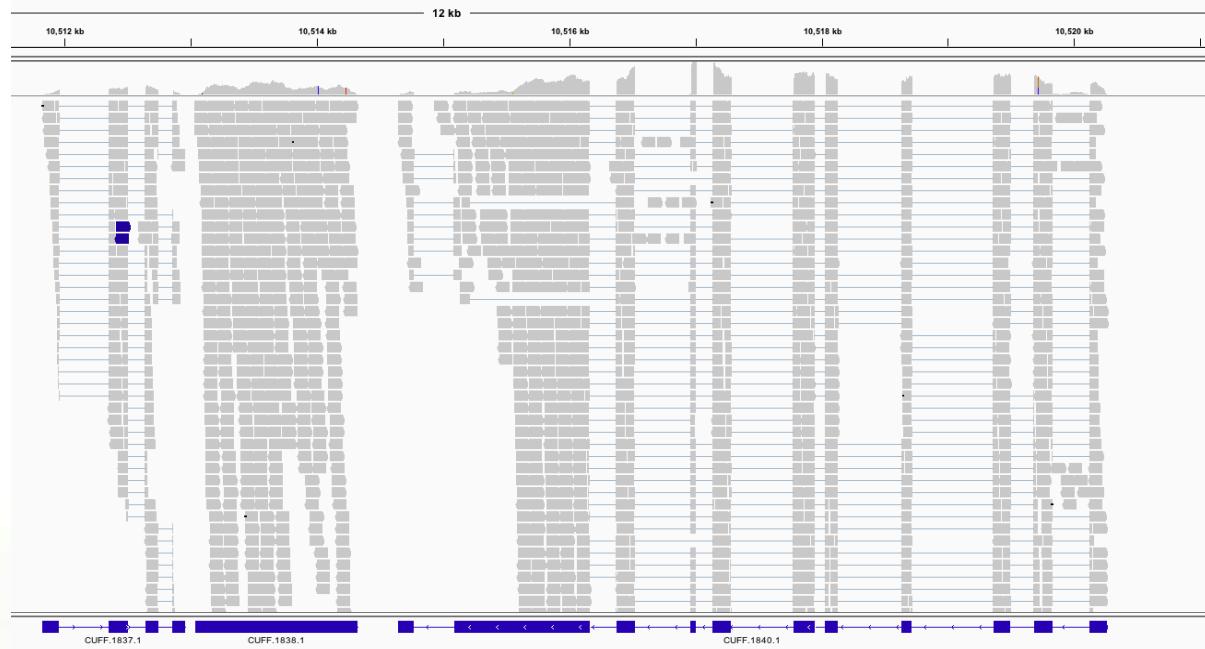


# Three-prime bias in polyA-selected rna-seq



# How to use RNA-seq

- Maker will align transcripts (ESTs), but these need to be assembled first.
- Cufflinks: mapped reads -> transcripts

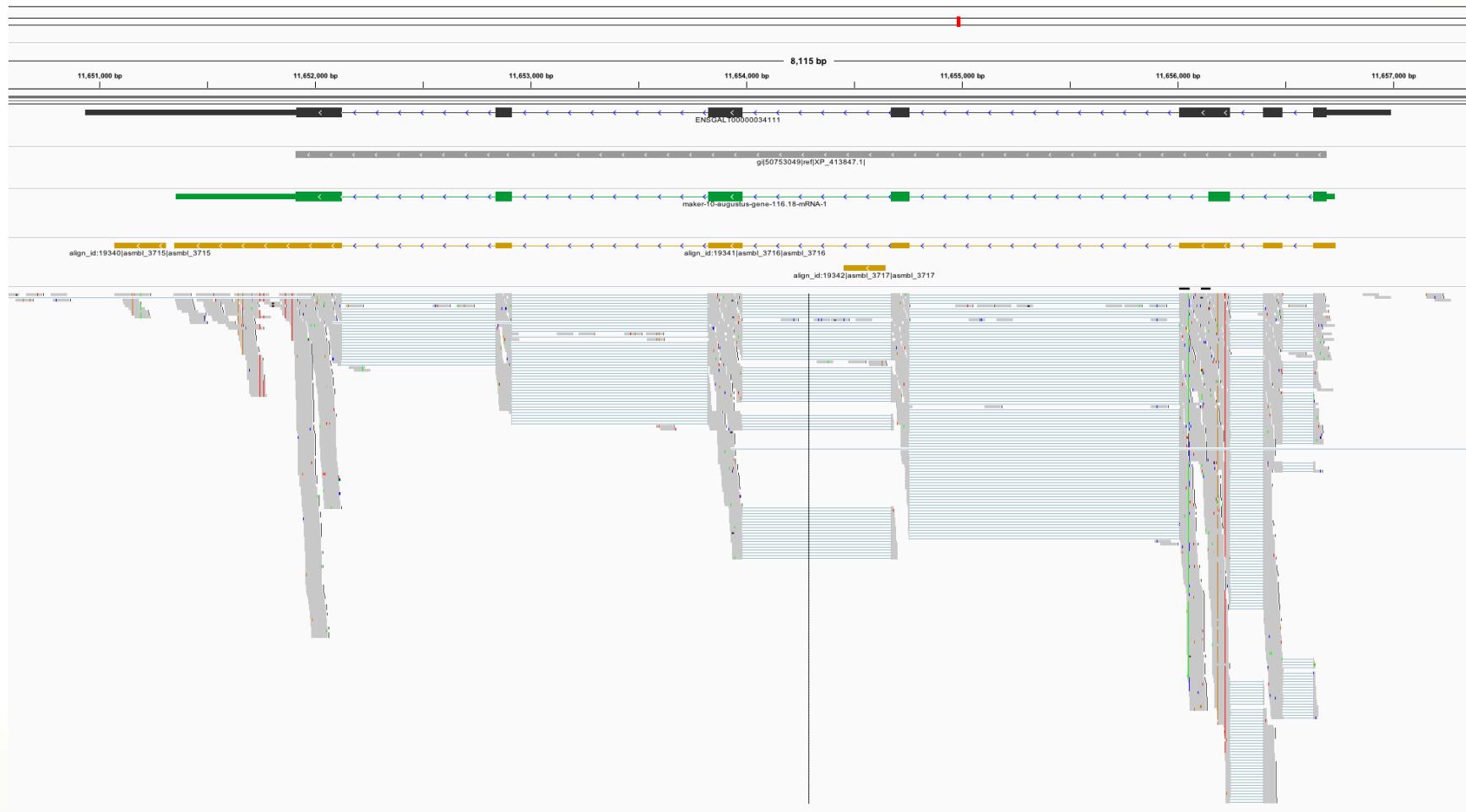


# How to use RNA-seq

- Maker will align transcripts (ESTs), but these need to be assembled first.
- Cufflinks: mapped reads -> transcripts
- Trinity: assembles transcripts without a genome

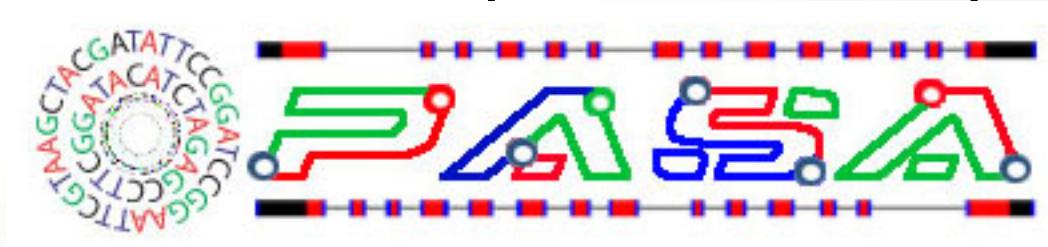


# Mapped Trinity-assembled transcripts



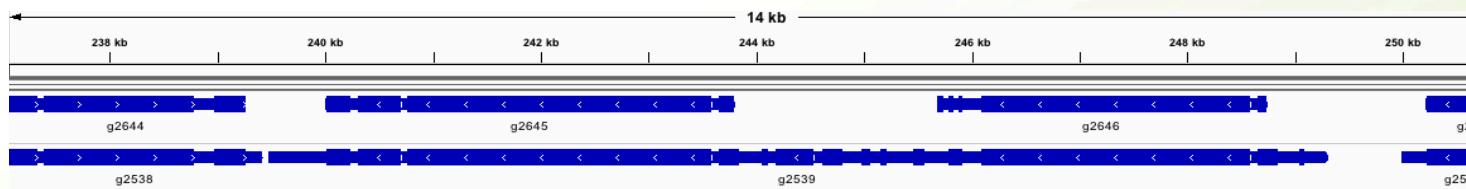
# How to use RNA-seq

- Maker will align transcripts (ESTs), but these need to be assembled first.
- Cufflinks: mapped reads -> transcripts
- Trinity: assembles transcripts without a genome
- PASA can be used to improve transcript quality



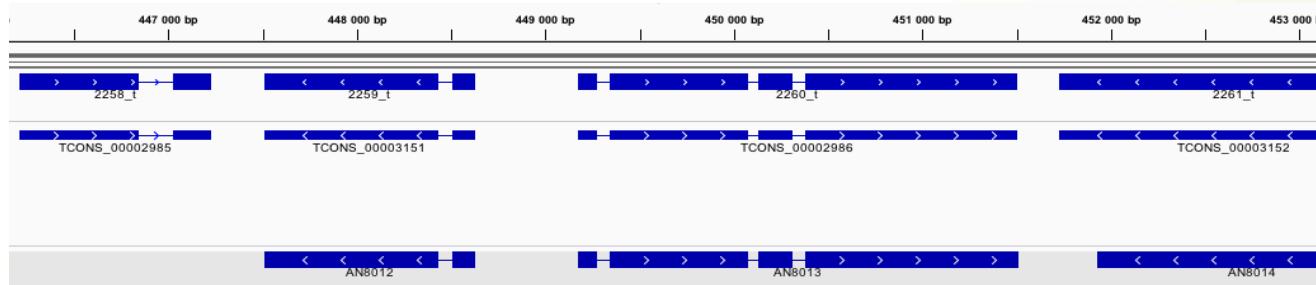
# Ab initio gene finders are used in Maker

- Commonly used programs: Augustus, Snap, Genemark-ES, FGENESH, Genscan, Glimmer-HMM,...
- Uses HMM-models to figure out how introns, exons, UTRs etc. are structured
- These HMM-models need to be trained!



# Liftovers are very useful for orthology determination

- Kraken
- Align the two genomes (Satsuma) and then transfer annotations between aligned regions



# General recommendations

- Always combine different types of evidence!
- One single method is not enough!
- Use Maker!

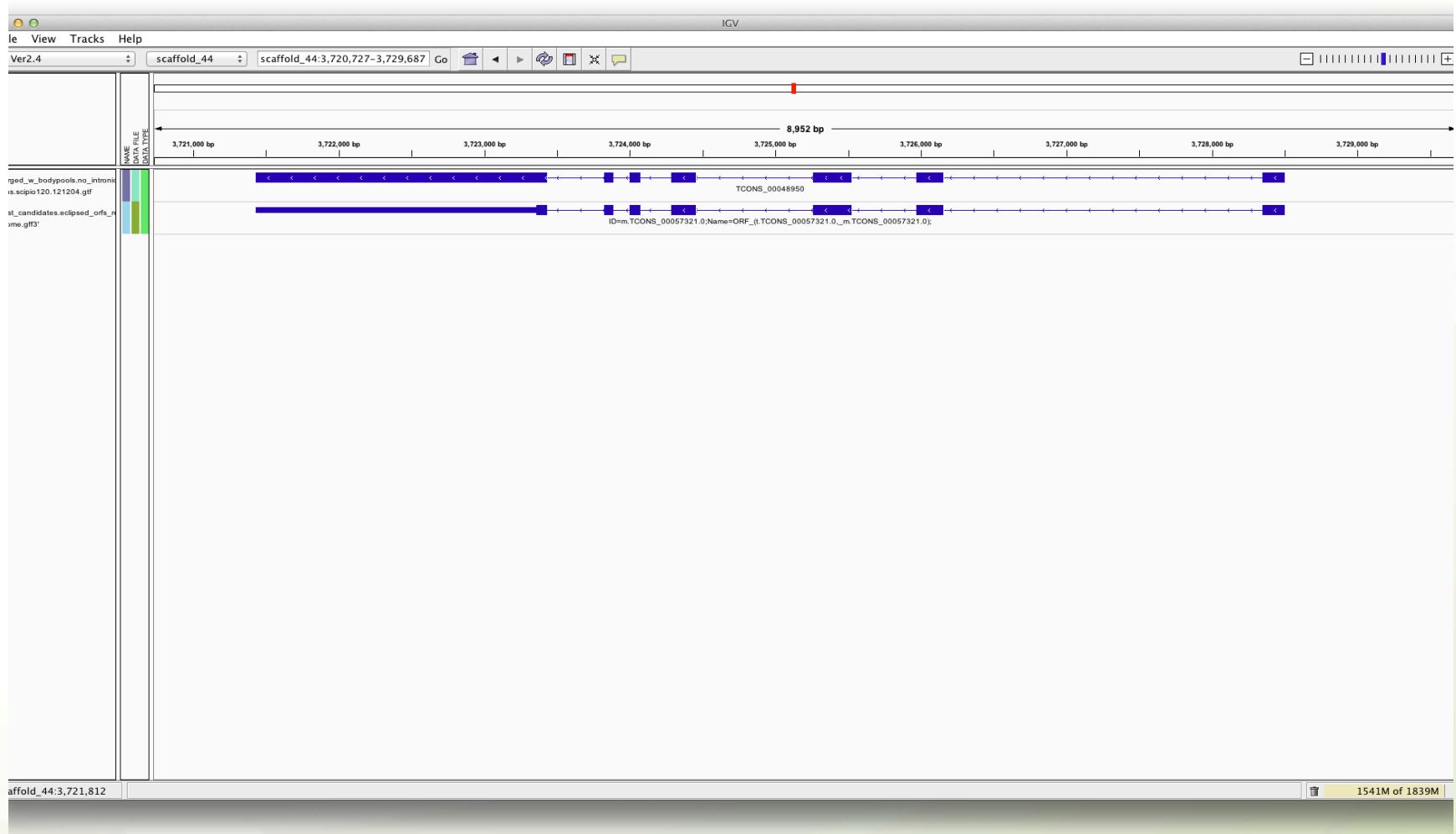


MAKER  
Annotate this!

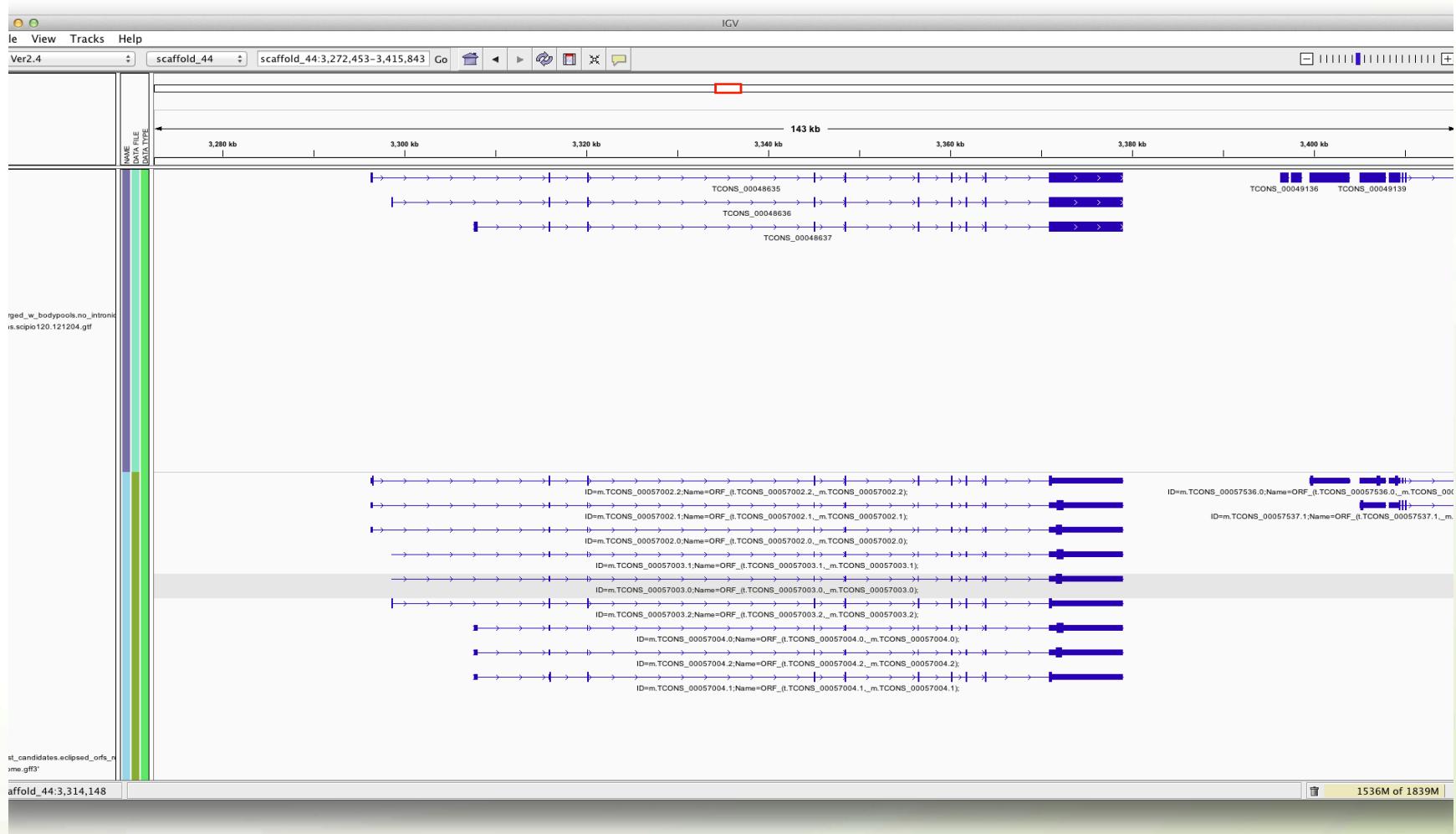
# Transcript annotation

- Here the transcript is already defined. The challenge is to find where the coding regions starts and stops
- Transdecoder

# Transdecoder



# Transdecoder



## Or get help - BILS assembly and annotation team

- Five people working with assembly and annotation
- Deliver high quality annotations
- Enable visualization and manual curation through a web interface
- Also available for consultation
- [support@bils.se](mailto:support@bils.se)

# Biosupport.se

LATEST 144 OPEN 1 RNA-SEQ 14 CHIP-SEQ SNP 9 ASSEMBLY 2 TUTORIALS TOOLS JOBS FORUM PLANET ALL »

Welcome to Biosupport.se! about • faq • rss

**SBS** Swedish Bioinformatics Support Community User Login New Post

Live search: start typing... or

Limit to: all time ▾ <prev • 144 results • page 1 of 6 • next > Sort by: update ▾

0 votes	2 answers	355 views	Translation of metagenomic DNA to protein sequence	metagenomics	written 12 weeks ago by Jutta • 110
2 votes	2 answers	2.6k views	mapping metagenomic reads to human reference genome using bowtie2	bowtie2 metagenomics referencemapping	written 2.4 years ago by Stef- • 210
0 votes	1 answer	152 views	analysing tri-allelic loci in a gwas	analysis plink association gwas	written 4 weeks ago by niclas • 110
1 vote	2 answers	1.4k views	Step by step instructions to carry out Principal Component Analysis of a Molecular Dynamics trajectory of a protein	pca amber	written 4 weeks ago by SuchetanaG • 110
1 vote	2 answers	126 views	Cuffmerge: merged.gtf correct?	cuffmerge cufflinks	written 4 weeks ago by Christina • 110
0 votes	1 answer	174 views	Making forest plot using R	stata r plot python	written 6 weeks ago by Stenemo- • 110
2 votes	2 answers	193 views	DNA methylation	methylation dna course	written 7 weeks ago by BigBen • 150
0 votes	1 answer	168 views	How to plot single peak for ChIP-seq data?	coverage chipseq plot	written 7 weeks ago by walker • 140
3 votes	3 answers	1.6k views	detect small RNA from normal RNA-seq data?	genome rna-seq	written 24 months ago by walker • 140
0 votes	1 answer	357 views	Tool for miRNA enrichment analysis	mirna rnaseq	written 4 months ago by Mary • 170
4 votes	3 answers	2.0k views	Cloud services for bioinformatics?	service cloud-computing	written 23 months ago by AndersW- • 210
1 vote	3 answers	2.2k views	SLURM output for automated processing	slurm automation	written 23 months ago by AndersW- • 210

**Recent Votes**

**Recent Locations • All »**

- Stockholm, Sweden, 53 minutes ago
- European Union, 2 hours ago
- Sweden, 15 hours ago
- Sweden, 15 hours ago

**Recent Awards • All »**

- Teacher ☺ to daho ++ 3.4k
- Teacher ☺ to daho ++ 3.4k
- Teacher ☺ to daho ++ 3.4k
- Rising Star ★ to daho ++ 3.4k
- Popular Question ☺ to daho ++ 3.4k
- Popular Question ☺ to daho ++ 3.4k

**Recent Replies**

- A: analysing tri-allelic loci in a gwas by DagAhrén • 2.5k  
Hi Niclas! I am no expert on GWAS, but I have done some searching and reading to try to come up w...
- A: Step by step instructions to carry out Principal Component Analysis of a Molecular Dynamics trajectory of a protein by WoA- • 110  
Check this tutorial from Bio-3D package in R:  
<http://thegrantlab.org/bio3d/tutorials/trajectory-a...>
- A: Step by step instructions to carry out Principal Component Analysis of a Molecular Dynamics trajectory of a protein by wes • 210  
Can't provide step-by-step instructions but maybe a few information resources will help get you s...
- C: Cuffmerge: merged.gtf correct? by Christina • 110  
Just an update: the analysis worked fine following your suggestions. Thanks again.