

Session 4.2 – Annotation

BU-ISCIII

Unidades Comunes Científico Técnicas – SGSAFI-ISCIII

11 al 15 Noviembre 2024

4ª Edición

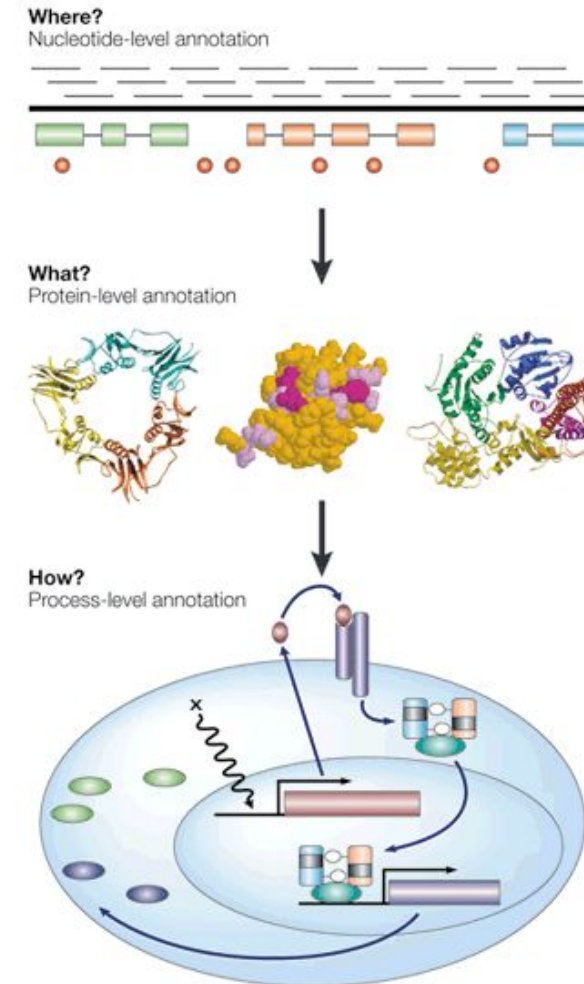
Programa Formación Continua, ISCIII

Annotation

Genome annotation is the process of **attaching biological (and positional) information to sequences**. It consists of three main steps:

- identifying portions of the genome that **do not code for proteins**
- Identifying coding elements on the genome, a process called **gene prediction**
- attaching **biological information** to these elements

<https://galaxyproject.github.io/training-material/topics/genome-annotation/tutorials/genome-annotation/tutorial.html>

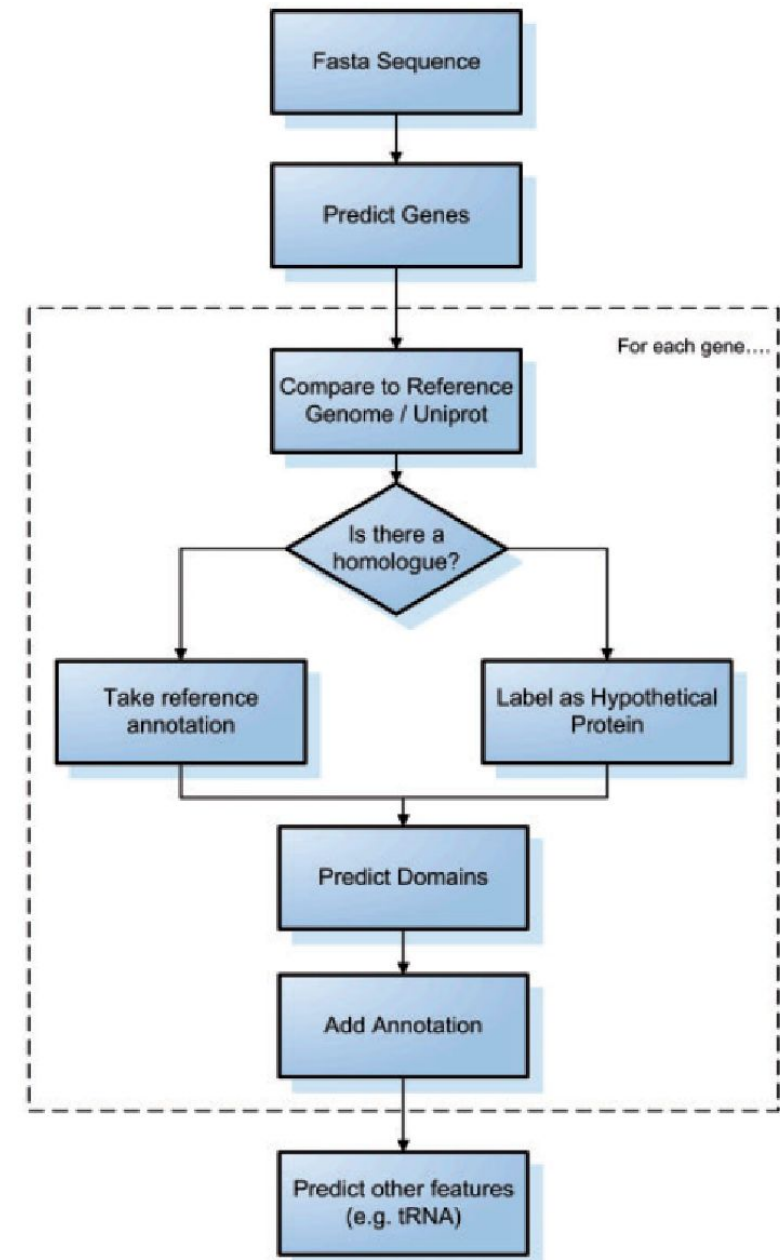


Main categories

- **Structural annotation** – Positions of genomic features along the genome. Finding genes and other biologically relevant sites with **specific locations but unknown function**
 - ORFs
 - Coding sequences(cds)
 - Promoters and regulatory regions
- **Functional annotation** – Assigning functions to features. Elements used in **database searches** to attach biologically relevant information to whole sequence and individual objects.

Automatic annotation

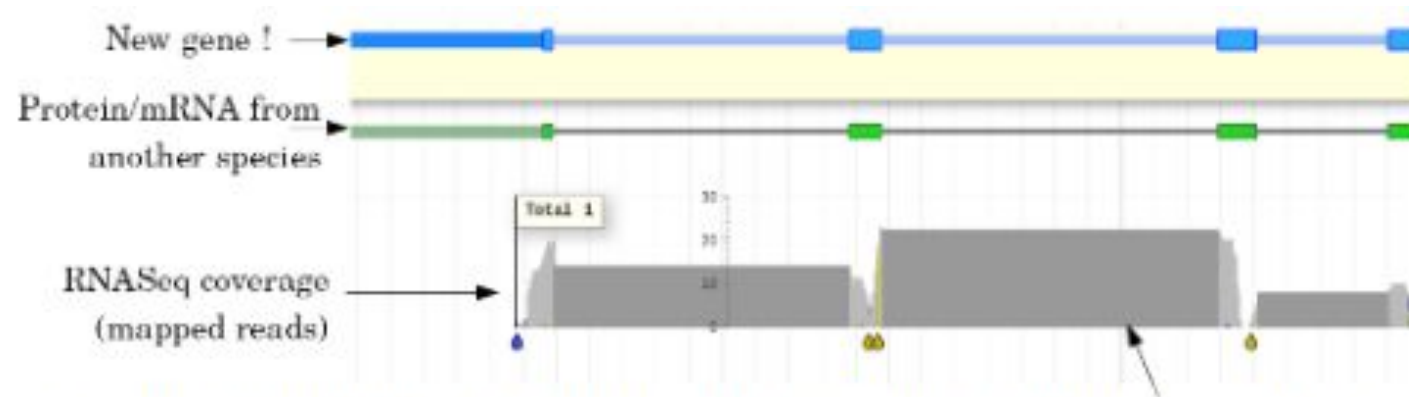
- Exponential submission of genomes
- Databases
 - Uniprot
 - RefSeq
 - Encyclopedia of DNA elements (ENCODE)
 - Entrez Gene
 - Ensembl
 - GENCODE
 - Gene Ontology Consortium (COGs)
 - GeneRIF
 - KEGG
 - Vertebrate and Genome Annotation Project (Vega)
 - Pfam
 - etc



Automatic annotation

Two strategies for identifying coding genes:

- Evidence: Sequence alignment to find known protein sequences in the contigs
 - transfer the annotation
 - will miss proteins not present in your database
 - may miss partial proteins



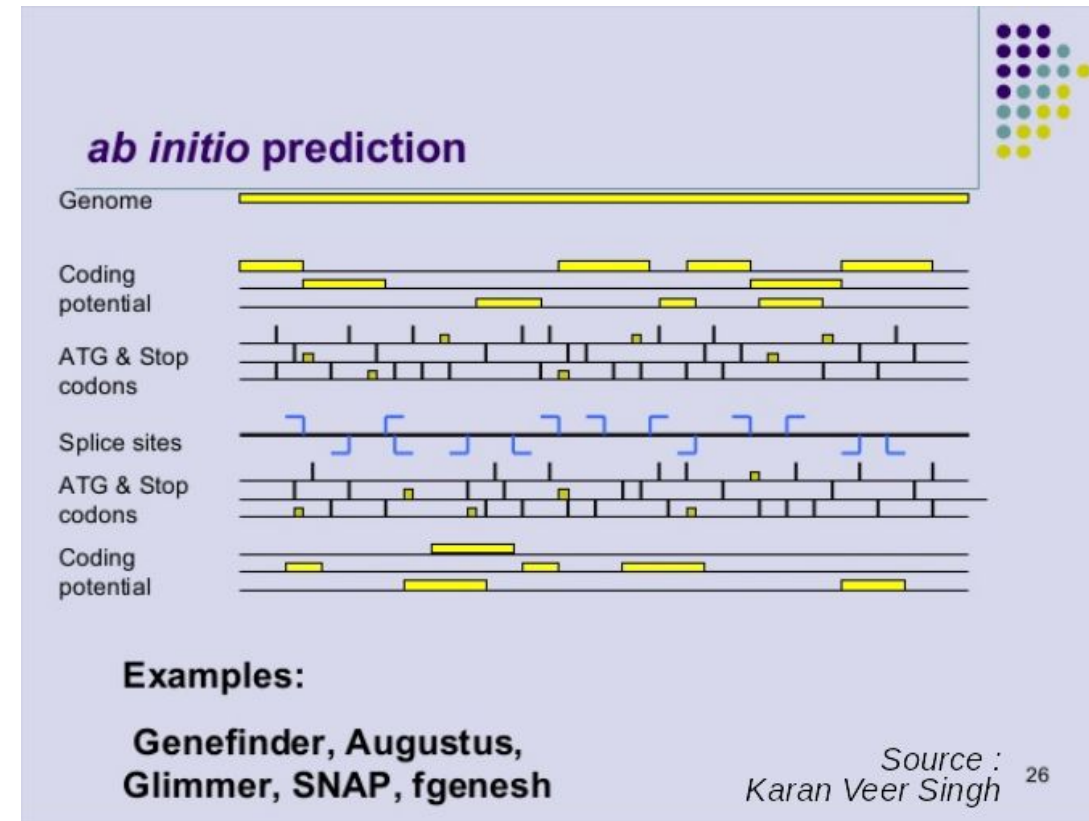
But data unavailable for novel or very distant genes, or unexpressed genes

[training.galaxyproject](https://training.galaxyproject.org/)

Automatic annotation

Two strategies for identifying coding genes:

- Ab initio gene finding o find candidate open reading frames:
 - Build model of ribosome binding sites
 - predict coding regions
 - may choose the incorrect start codon
 - may miss atypical genes, overpredict small genes



[training.galaxyproject](https://training.galaxyproject.org/)

Automatic annotation

- **tRNA:** easy to find and annotate: anti-codon
- **rRNA:** easy to find and annotate: 5s 16s 23s
- **CDS:** straightforward to find candidates
 - false positives are often small ORFs
 - wrong start codon o partial genes
 - Pseudogenes
 - assigning function is the bulk of the workload

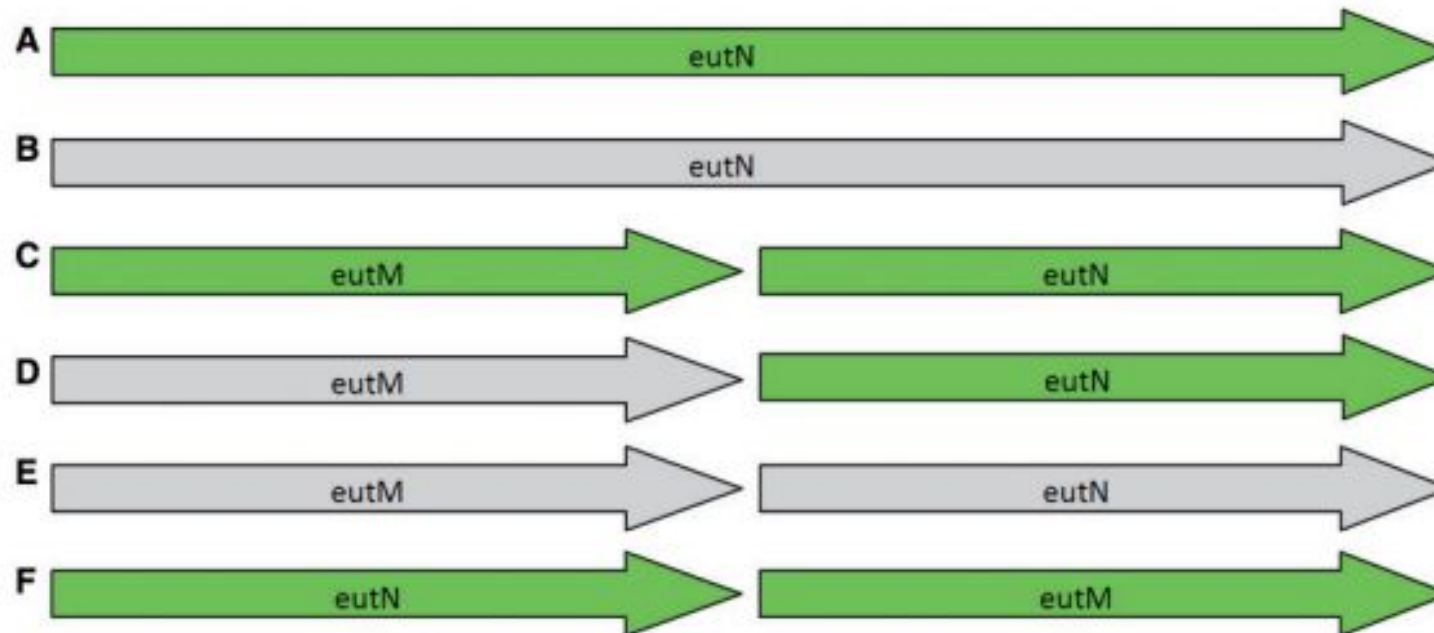
Automatic annotation: limitations

- If sequence homologous are found, may **not be functional** homologous
- If **no homology found**- limited information can be inferred
- Incorrect annotation can be **propagated** when similarity is over part on sequence not used in annotation
 - Multidomain proteins (HMM)
- Inconsistent annotation (**Different names, same protein**)
- Same **gene name, different product** name
- Spelling mistakes
- Looking for **new genes**, not present in DDBB
- Expression experiments / Manual annotation needed

Richardson and Watson. Briefings in Bioinformatics. 2012

Automatic annotation: limitations

Inconsistent annotation



Salmonella typhi CT18 (NC_003198) and *Salmonella typhi* Ty2 (NC_004631) there is a single ORF of 690 bp

Figure 2: The six different models present across 17 RefSeq entries for *Salmonella* species for the *eutM/eutN* locus. Green indicates normal gene/CDS features, lighter grey indicates gene features annotated as pseudogenes. (A) A single intact gene of 690 bp; (B) a single pseudogene of 690 bp; (C) two short intact genes ~300 bp in length; (D) one pseudogene and one intact gene, each ~300 bp in length; (E) two pseudogenes, each 300 bp in length; and (F) two intact genes with the order reversed.

Richardson and Watson. *Briefings in Bioinformatics*. 2012

Automatic annotation: limitations

These two regions are more than 97% identical at the nucleotide level; however, the annotation differs considerably.

While *E. coli* K12 MG1655 contains features with gene names *araA*, *araB* and *araC*, the equivalent features in *E. coli* O157:H7 Sakai do not have those gene names and have been assigned uninformative locus tags

Inconsistent annotation

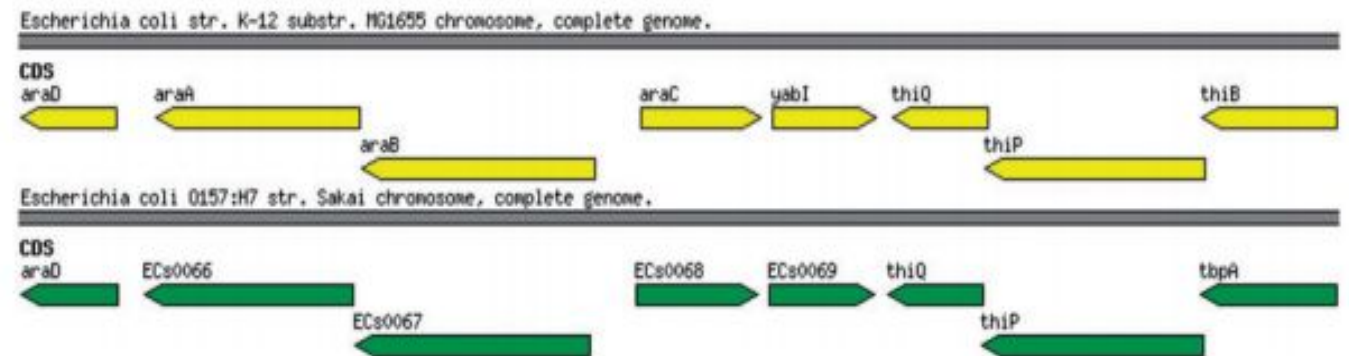


Figure 3: A syntenic block of genes showing inconsistent gene name annotations in *E. coli* K12 MG1655 and *E. coli* O157:H7 Sakai.

*Richardson and Watson. Briefings
in Bioinformatics. 2012*

Automatic annotation: limitations

- Spelling mistakes

- There are 128 proteins in UniProt that contain the word 'syntase', an incorrect spelling of the word 'synthase'
- If a user was to visit any of these databases and search for 'dihydrofolate synthase' the misspelled entries would be omitted from the search results

*Richardson and Watson. Briefings
in Bioinformatics. 2012*

Automatic annotation: limitations

- ‘Same gene name, different product name’
 - The NCBI validation software specifically highlights when this occurs intra-genomically with the description ‘Same gene name, different product name’

Table 1: Different product names assigned to features with the gene name 'int' across 17 different RefSeq entries for *Salmonella* species

Gene name	Product name	Accession
int	bacteriophage integrase	NC.003198, NC.004631, NC.015761
int	Gifsy-1 prophage Int	NC.006905
int	hypothetical protein	NC.006905
int	Integrase	NC.003198, NC.004631, NC.006511, NC.012125
int	integrase (fragment)	NC.003198
int	phage integrase family site specific recombinase	NC.006905
int	putative cytoplasmic protein	NC.006905
Int	Putative integrase	NC.003384
int	putative integrase protein	NC.006905
int	putative P4-type integrase	NC.006905
int	putative phage integrase protein	NC.006905
int	site-specific recombinase, phage integrase family	NC.012125

Richardson and Watson. Briefings
in Bioinformatics. 2012

Automatic annotation: limitations

Hypothetical proteins

- These may be real genes with no known function or they may be artifacts of the gene prediction process.
- Often there are features which are only orthologous to other hypothetical features and do not contain any domains. These could either be regions with no functionality, a relic of the feature prediction software or the domains present have not been discovered yet
- Whether or not to include them is often a decision made by the annotation team and varies between groups
- As experimental data becomes more ubiquitous evidence tags should play a larger role in annotation.

*Richardson and Watson. Briefings
in Bioinformatics. 2012*

Automatic annotation: limitations

Distinguishing orthologs from paralogs

orthologs tend to retain similar functions, whereas paralogs tend to diverge over time to perform different functions

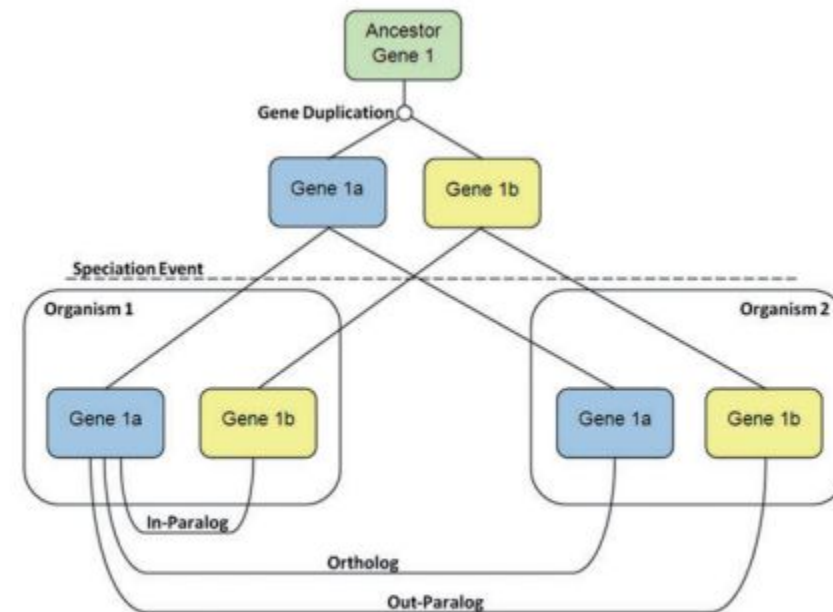


Figure 4: A diagram displaying the processes that can lead to, and define, orthologs and paralogs. Gene duplication and speciation events create complex evolutionary relationships between genes.

*Richardson and Watson. Briefings
in Bioinformatics. 2012*

Automatic annotation: limitations

- RefSeq is one attempt to standardize and improve the quality of genome annotation
 - WP_ prefix. All identical proteins regardless of species
 - Standard classification

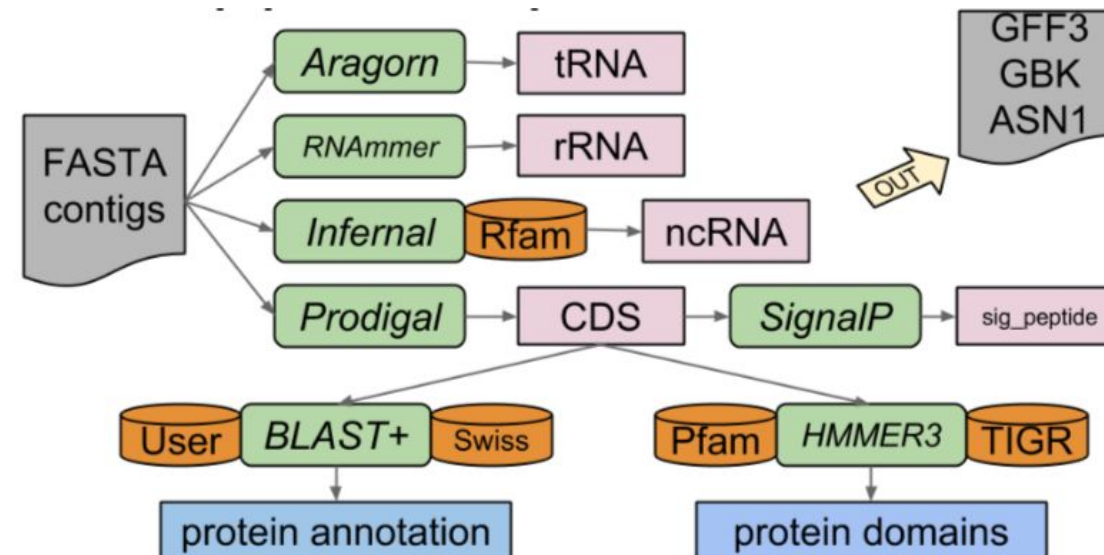
```

beta-lactamase (conceptual)
  class A beta-lactamase (HMM:NF033103)
  metallo-beta-lactamase (HMM:NF012229)
    subclass B1 metallo-beta-lactamase (HMM:NF033088)
      NDM family subclass B1 metallo-beta-lactamase (HMM:NF000259)
        subclass B1 metallo-beta-lactamase NDM-1 (allele)
        subclass B1 metallo-beta-lactamase NDM-2 (allele)
        subclass B1 metallo-beta-lactamase NDM-3 (allele)
      VIM family subclass B1 metallo-beta-lactamase (HMM:NF012100)
      SPM family subclass B1 metallo-beta-lactamase (HMM:NF012150)
    subclass B2 metallo-beta-lactamase (HMM:NF033087)
    subclass B3 metallo-beta-lactamase (HMM:NF033105)
  class C beta-lactamase (HMM:NF033085)
  class D beta-lactamase (conceptual)
    class D beta-lactamase (main branch) (HMM:NF012161)
    class D beta-lactamase (other branch) (HMM:NF000270)
  
```


Automatic annotation: Prokka (Rapid prokaryotic genome annotation)

Seeman, Bioinformatics 2014

Tool (reference)	Features predicted
Prodigal (Hyatt 2010)	Coding sequence (CDS)
RNAmmmer (Lagesen et al. , 2007)	Ribosomal RNA genes (rRNA)
Aragorn (Laslett and Canback, 2004)	Transfer RNA genes
SignalP (Petersen et al. , 2011)	Signal leader peptides
Infernal (Kolbe and Eddy, 2011)	Non-coding RNA
BLAST+ (Camacho <i>et al.</i> , 2009)	Specific function or name Personal database



<https://galaxyproject.github.io/training-material/topics/genome-annotation/tutorials/annotation-with-prokka/slides.html#8>

Automatic annotation: Prokka

- Optional **user-provided** set of annotated proteins
- All bacterial proteins in **UniProt**
- All proteins from finished bacterial genomes in **RefSeq**
- Hidden Markov model profile databases, **Pfam** and **TIGRFAMs**
- Hypothetical protein

Prokka uses this method, but in a hierarchical manner, starting with a **smaller trustworthy database**, moving to medium sized but **domain-specific databases**, and finally to **curated models of protein families**

Automatic annotation: Prokka

- Facts
 - searching against smaller databases is faster
 - searching against similar sequences is faster
- Idea
 - start with small set of close proteins
 - advance to larger sets of more distant proteins
- Prokka
 - your own custom "trusted" set (optional)
 - core bacterial proteome (default)
 - genus specific proteome (optional)
 - whole protein HMMs: PRK clusters, TIGRfams
 - protein domain HMMs: Pfam

Prokka uses this method, but in a hierarchical manner, starting with a **smaller trustworthy database**, moving to medium sized but **domain-specific databases**, and finally to **curated models of protein families**

Automatic annotation: Prokka output

Suffix	Description of file contents
.fna	FASTA file of original input contigs (nucleotide)
.faa	FASTA file of translated coding genes (protein)
.ffn	FASTA file of all genomic features (nucleotide)
.fsa	Contig sequences for submission (nucleotide)
.tbl	Feature table for submission
.sqn	Sequin editable file for submission
.gbk	Genbank file containing sequences and annotations
.gff	GFF v3 file containing sequences and annotations
.log	Log file of Prokka processing output
.txt	Annotation summary statistics

Viral genome annotation

PROPERTIES

- DNA, ssDNA, dsDNA, RNA, ssRNA, fragmented RNA
- Non-coding ORF
- Coding ORF
- Overlapping reading frames
- Non-standard nomenclature for viral gene products
- RNA editing (the RNA polymerase co-transcriptionally adds one or two nucleotides that are not on the template, including multiple proteins in a single gene. Annotated protein sequence does not match the expected translated nucleotide sequence)
- Ribosome slippage (Allow viruses to produce two proteins from a single mRNA transcript by having the ribosome 'slip' one or two nucleotides along the mRNA transcript, thus changing the reading frame.)
- Viral sequence variability

Viral genome annotation

APPROACHES

- Identification hallmark genes conserved within known virus families
- Detection of short nucleotide sequences believed to be enriched in viruses (DeepVirFinder: reference-free and alignment-free machine learning method, for identifying viral sequences in metagenomic data using deep learning. Ren et al., Quan Biol 2020)
- Tools for specific virus (i.e. Influenza)

Viral genome annotation

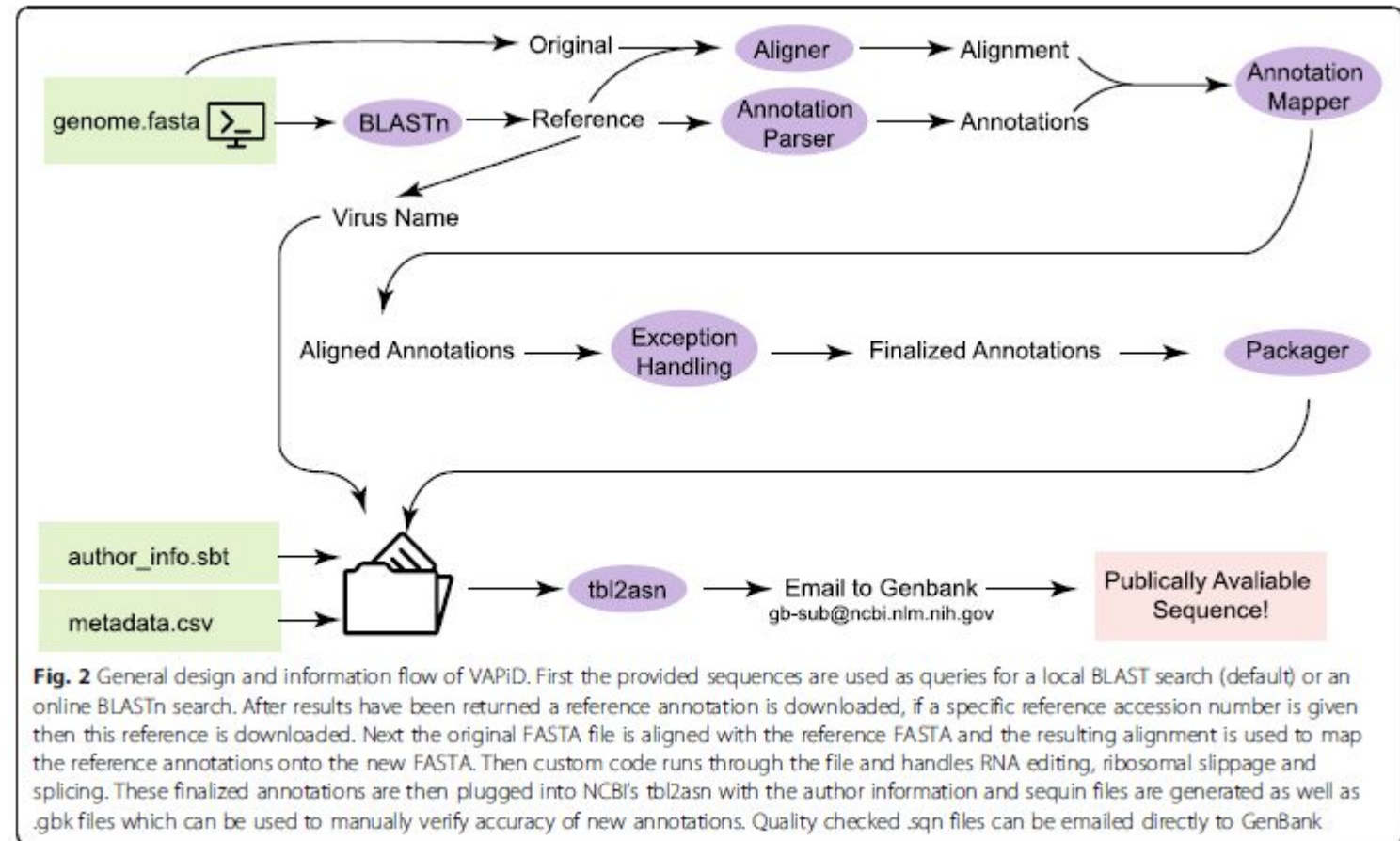
LIMITATIONS

- Pitfalls that can lead to false-positives or false negatives
- Some tools are limited by minimum sequence length
- Detection of a limited range of virus families.
- High diversity of DNA and RNA viruses presents a challenge for development of a universal annotator

VAPiD: a lightweight cross-platform viral annotation pipeline and identification tool to facilitate virus genome submissions to NCBI GenBank

- Users can provide a specified reference from which to annotate all viruses
- Provide their own BLASTn database
- Force VAPiD to search NCBI's NT database

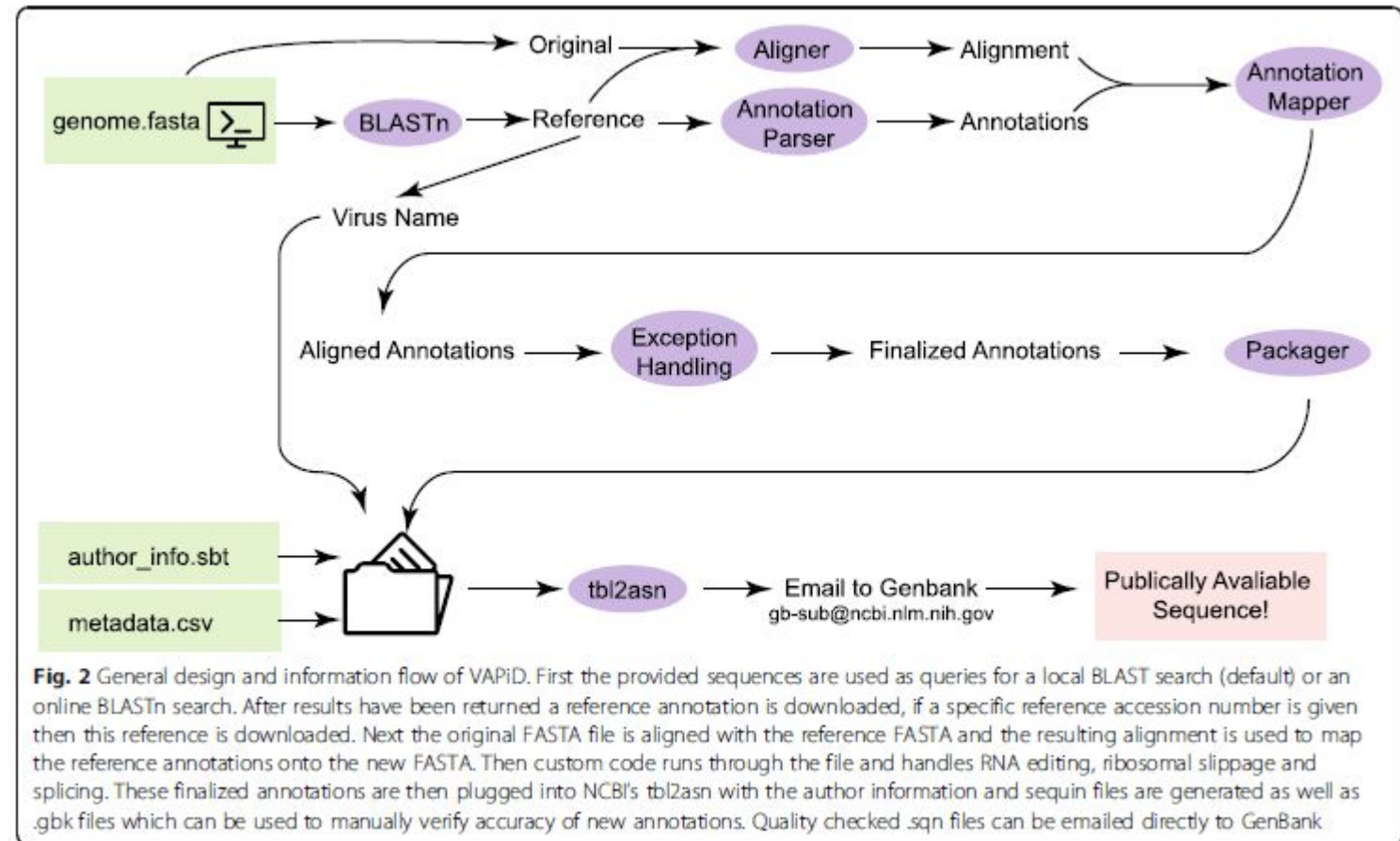
<https://github.com/rcs333/VAPiD>



VAPiD: a lightweight cross-platform viral annotation pipeline and identification tool to facilitate virus genome submissions to NCBI GenBank

ALGORITHM STEPS:

1. Find the correct reference sequence.
2. Gene locations are stripped from the reference
3. Pairwise nucleotide alignment between the reference and the submitted sequence is generated using MAFFT
4. The relative locations of the genes on the reference sequence are then mapped onto the new sequence
5. Gene names are taken from the annotated reference sequence
6. Spellchecking
7. RNA editing
8. Ribosome slippage
9. Genbank file generation



VAPiD: a lightweight cross-platform viral annotation pipeline and identification tool to facilitate virus genome submissions to NCBI GenBank

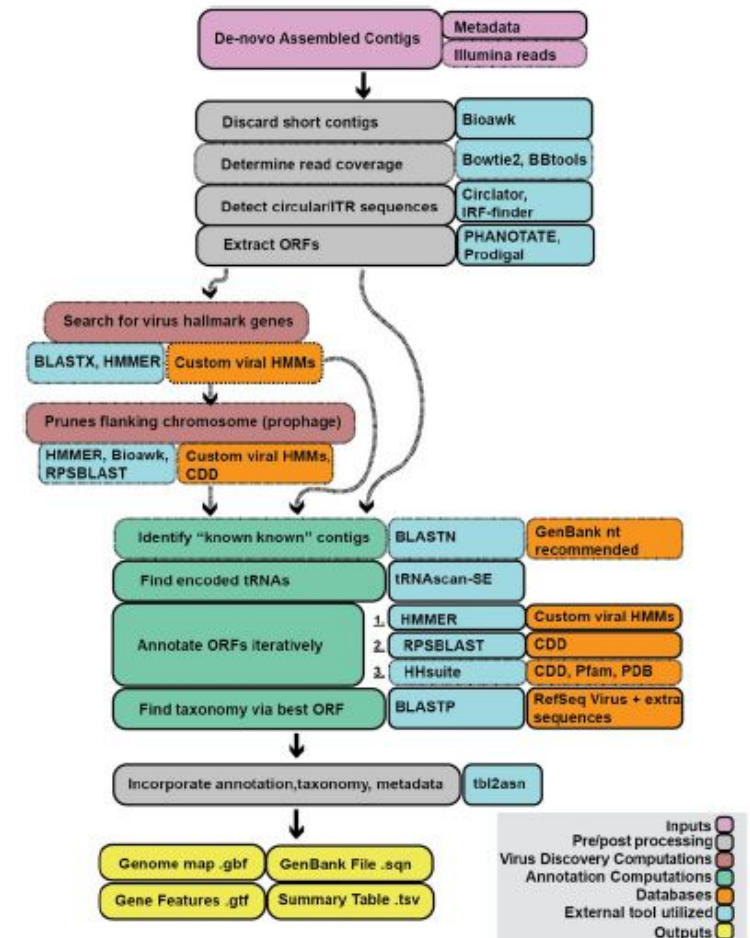
LIMITATIONS

- VAPiD is not the preferred annotation tool for novel or extremely divergent viral species
- Not perform ab initio gene annotation
- Any errors that are in the downloaded reference will be transferred to the new genome (i.e. misspelling)
- VAPiD performs best on high-quality and accurate reference sequences

Cenote-Taker 2

Discovery and annotation of viruses in DNA and RNA genomes of multiple data types (genomic, metagenomic, transcriptomic, etc.)

- Viral genomes annotation:
 - Viral genomes with divergent ORFs,
 - Diamond, Hmmer, BLAST
- Viral discovery:
 - de novo assembly (Megahit, SOAPdenovo2.)
 - Identification of viral structural and replication genes using reference gene databases.
- Prophage Pruning Module
 - Processes bacterial genomes to identify and prune induced prophages.
 - Reads aligned to reference genomes using Bowtie2.
 - Visualization of prophage coverage with the Integrative Genomics Viewer (IGV).



Tisza et al., Virus Evolution 2021

SnpEff

SnpEff (SNP effect) multi-platform open source variant effect predictor program. It annotates variants and predicts the coding effects of genetic variations (SNPs, INDELs, MNPs).

Properties:

- Speed
- Flexibility
- Integration with Galaxy
- Multi-species
- Compatible with GATK
- Non-coding annotation

https://pcingola.github.io/SnpEff/adds/SnpEff_paper.pdf

SnpEff

Steps:

1. Database build: reference genome .fasta + annotation .gtf/.gff
2. Effect calculation: Interval forest algorithm
 - Hash of interval trees indexed by chromosome. Each node has five elements
 - Querying an interval tree
 - Effect prediction
3. Output: annotated vcf

Column	Notes
Chromosome	Chromosome name (usually without any leading 'chr' string)
Position	One based position
Reference	Reference
Change	Sequence change
Change type	Type of change (SNP, MNP, INS, DEL)
Homozygous	Is this homozygous or heterozygous (Hom, Het)
Quality	Quality score (from input file)
Coverage	Coverage (from input file)
Warnings	Any warnings or errors.
Gene_ID	Gene ID (usually ENSEMBL)
Gene_name	Gene name
Bio_type	BioType, as reported by ENSEMBL
Transcript_ID	Transcript ID (usually ENSEMBL)
Exon_ID	Exon ID (usually ENSEMBL)
Exon_Rank	Exon number on a transcript
Effect	Effect of this variant. See details below
old_AA/new_AA	Amino acid change
old_codon/new_codon	Codon change
Codon_Num(CDS)	Codon number in CDS
Codon_degeneracy	Codon degeneracy
CDS_size	CDS size in bases
Custom_interval_ID	If any custom interval was used, add the IDs here (may be more than one)

SnpEff

Table 2. Detailed effect list from SnpEff

Effect	Note
INTERGENIC	The variant is in an intergenic region
UPSTREAM	Upstream of a gene (default length: 5K bases)
UTR_5_PRIME	Variant hits 5'UTR region
UTR_5_DELETED	The variant deletes and exon which is in the 5'UTR of the transcript
START_GAINED	A variant in 5'UTR region produces a three base sequence that can be a START codon
SPLICE_SITE_ACCEPTOR	The variant hits a splice acceptor site (defined as two bases before exon start, except for the first exon)
SPLICE_SITE_DONOR	The variant hits a Splice donor site (defined as two bases after coding exon end, except for the last exon)
START_LOST	Variant causes start codon to be mutated into a non-start codon
SYNONYMOUS_START	Variant causes start codon to be mutated into another start codon
CDS	The variant hits a CDS
GENE	The variant hits a gene
TRANSCRIPT	The variant hits a transcript
EXON	The variant hits an exon
EXON_DELETED	A deletion removes the whole exon
NON_SYNONYMOUS_CODING	Variant causes a codon that produces a different amino acid
SYNONYMOUS_CODING	Variant causes a codon that produces the same amino acid
FRAME_SHIFT	Insertion or deletion causes a frame shift
CODON_CHANGE	One or many codons are changed
CODON_INSERTION	One or many codons are inserted
CODON_CHANGE_PLUS_CODON_INSERTION	One codon is changed and one or many codons are inserted
CODON_DELETION	One or many codons are deleted
CODON_CHANGE_PLUS_CODON_DELETION	One codon is changed and one or more codons are deleted
STOP_GAINED	Variant causes a STOP codon
SYNONYMOUS_STOP	Variant causes stop codon to be mutated into another stop codon
STOP_LOST	Variant causes stop codon to be mutated into a non-stop codon
INTRON	Variant hits an intron. Technically, hits no exon in the transcript
UTR_3_PRIME	Variant hits 3'UTR region
UTR_3_DELETED	The variant deletes and exon which is in the 3'UTR of the transcript
DOWNSTREAM	Downstream of a gene (default length: 5K bases)
INTRON_CONSERVED	The variant is in a highly conserved intronic region
INTERGENIC_CONSERVED	The variant is in a highly conserved intergenic region

Sub-field	Notes
Effect	Effect of this variant. See details below
Codon_Change	Codon change: old_codon/new_codon
Amino_Acid_change	Amino acid change: old_AA/new_AA
Warnings	Any warnings or errors
Gene_name	Gene name
Gene_BioType	BioType, as reported by ENSEMBL
Coding	[CODING NON_CODING]. If information reported by ENSEMBL (e.g., has 'protein_id' information in GTF file)
Transcript	Transcript ID (usually ENSEMBL)
Exon	Exon ID (usually ENSEMBL)
Warnings	Any warnings or errors (not shown if empty)

The information is added to the INFO fields using an tag 'EFF'. The format for each effect is "Effect (Effect_Impact | Codon_Change | Amino_Acid_change | Gene_Name | Gene_BioType | Coding | Transcript | Exon [| ERRORS | WARNINGS])"

Annotation format: gff3

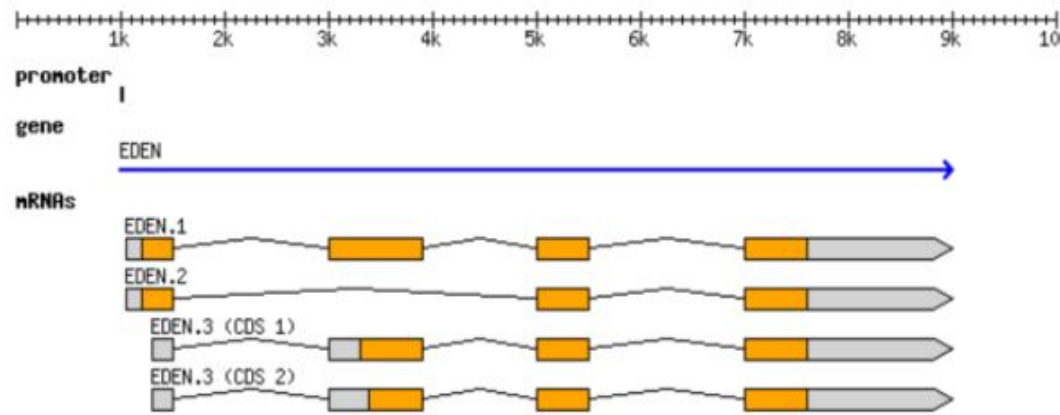
1. Seqid - name
2. Source - program
3. Type - term or SOFA sequence ontology
4. Start
5. End
6. Score
7. Strand - (+/-)
8. Phase - (0/1/2)
9. Attributes
 - Name
 - Alias
 - Parent
 - Target
 - Gap
 - Derives_from
 - Note
 - Dbxref
 - Ontology_term

```
##gff-version 3.2.1
##sequence-region ctg123 1 1497228
ctg123 . gene 1000 9000 . + . ID=gene00001;Name=EDEN
ctg123 . TF_binding_site 1000 1012 . + . ID=tfbs00001;Parent=gene00001
ctg123 . mRNA 1050 9000 . + . ID=mRNA00001;Parent=gene00001;Name=EDEN.1
ctg123 . mRNA 1050 9000 . + . ID=mRNA00002;Parent=gene00001;Name=EDEN.2
ctg123 . mRNA 1300 9000 . + . ID=mRNA00003;Parent=gene00001;Name=EDEN.3
ctg123 . exon 1300 1500 . + . ID=exon00001;Parent=mRNA00003
ctg123 . exon 1050 1500 . + . ID=exon00002;Parent=mRNA00001,mRNA00002
ctg123 . exon 3000 3902 . + . ID=exon00003;Parent=mRNA00001,mRNA00003
ctg123 . exon 5000 5500 . + . ID=exon00004;Parent=mRNA00001,mRNA00002,mRNA00003
ctg123 . exon 7000 9000 . + . ID=exon00005;Parent=mRNA00001,mRNA00002,mRNA00003
ctg123 . CDS 1201 1500 . + 0 ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
ctg123 . CDS 3000 3902 . + 0 ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
ctg123 . CDS 5000 5500 . + 0 ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
ctg123 . CDS 7000 7600 . + 0 ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
ctg123 . CDS 1201 1500 . + 0 ID=cds00002;Parent=mRNA00002;Name=edenprotein.2
ctg123 . CDS 5000 5500 . + 0 ID=cds00002;Parent=mRNA00002;Name=edenprotein.2
ctg123 . CDS 7000 7600 . + 0 ID=cds00002;Parent=mRNA00002;Name=edenprotein.2
ctg123 . CDS 3301 3902 . + 0 ID=cds00003;Parent=mRNA00003;Name=edenprotein.3
ctg123 . CDS 5000 5500 . + 1 ID=cds00003;Parent=mRNA00003;Name=edenprotein.3
ctg123 . CDS 7000 7600 . + 1 ID=cds00003;Parent=mRNA00003;Name=edenprotein.3
ctg123 . CDS 3391 3902 . + 0 ID=cds00004;Parent=mRNA00003;Name=edenprotein.4
ctg123 . CDS 5000 5500 . + 1 ID=cds00004;Parent=mRNA00003;Name=edenprotein.4
ctg123 . CDS 7000 7600 . + 1 ID=cds00004;Parent=mRNA00003;Name=edenprotein.4
```

Formato del fichero GFF (General feature format)

- Formato standard para describir genes o transcritos.... 9 col, tab-delimited, plain text files

The Canonical Gene



seqid	source	type	start	end	score	strand	phase	atributes
0	##gff-version	3.1.26						
1	##sequence-region	ctg123	1	1497228				
2	ctg123	gene	1000	9000	.	+	.	ID=gene00001;Name=EDEN
3	ctg123	TF_binding_site	1000	1012	.	+	.	ID=tfbs00001;Parent=gene00001
4	ctg123	mRNA	1050	9000	.	+	.	ID=mRNA00001;Parent=gene00001;Name=EDEN.1
5	ctg123	mRNA	1050	9000	.	+	.	ID=mRNA00002;Parent=gene00001;Name=EDEN.2
6	ctg123	mRNA	1300	9000	.	+	.	ID=mRNA00003;Parent=gene00001;Name=EDEN.3
7	ctg123	exon	1300	1500	.	+	.	ID=exon00001;Parent=mRNA00003
8	ctg123	exon	1050	1500	.	+	.	ID=exon00002;Parent=mRNA00001,mRNA00002
9	ctg123	exon	3000	3902	.	+	.	ID=exon00003;Parent=mRNA00001,mRNA00003
10	ctg123	exon	5000	5500	.	+	.	ID=exon00004;Parent=mRNA00001,mRNA00002,mRNA00003
11	ctg123	exon	7000	9000	.	+	.	ID=exon00005;Parent=mRNA00001,mRNA00002,mRNA00003
12	ctg123	CDS	1201	1500	.	+	0	ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
13	ctg123	CDS	3000	3902	.	+	0	ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
14	ctg123	CDS	5000	5500	.	+	0	ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
15	ctg123	CDS	7000	7600	.	+	0	ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
16	ctg123	CDS	1201	1500	.	+	0	ID=cds00002;Parent=mRNA00002;Name=edenprotein.2
17	ctg123	CDS	5000	5500	.	+	0	ID=cds00002;Parent=mRNA00002;Name=edenprotein.2
18	ctg123	CDS	7000	7600	.	+	0	ID=cds00002;Parent=mRNA00002;Name=edenprotein.2
19	ctg123	CDS	3301	3902	.	+	0	ID=cds00003;Parent=mRNA00003;Name=edenprotein.3
20	ctg123	CDS	5000	5500	.	+	1	ID=cds00003;Parent=mRNA00003;Name=edenprotein.3
21	ctg123	CDS	7000	7600	.	+	1	ID=cds00003;Parent=mRNA00003;Name=edenprotein.3
22	ctg123	CDS	3391	3902	.	+	0	ID=cds00004;Parent=mRNA00003;Name=edenprotein.4
23	ctg123	CDS	5000	5500	.	+	1	ID=cds00004;Parent=mRNA00003;Name=edenprotein.4
24	ctg123	CDS	7000	7600	.	+	1	ID=cds00004;Parent=mRNA00003;Name=edenprotein.4

<https://github.com/The-Sequence-Ontology/Specifications/blob/master/gff3.md>

Annotation format: gbk

- LOCUS – Annotated sequence
- DEFINITION
- ACCESSION
- FEATURES
 - source
 - gene
 - CDS
 - Locus tag
 - function
 - Product
 - protein_id
 - Translation (sequence)

```

LOCUS       AF068625                200 bp    mRNA    linear    ROD 06-DEC-1999
DEFINITION  Mus musculus DNA cytosine-5 methyltransferase 3A (Dnmt3a) mRNA,
            complete cds.
ACCESSION   AF068625 REGION: 1..200
VERSION     AF068625.2  GI:6449467
KEYWORDS    .
SOURCE      Mus musculus (house mouse)
  ORGANISM  Mus musculus
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Mammalia; Eutheria; Euarchontoglires; Glires; Rodentia;
            Sciurognathi; Muroidea; Muridae; Murinae; Mus.
REFERENCE   1 (bases 1 to 200)
  AUTHORS   Okano,M., Xie,S. and Li,E.
  TITLE     Cloning and characterization of a family of novel mammalian DNA
            (cytosine-5) methyltransferases
  JOURNAL   Nat. Genet. 19 (3), 219-220 (1998)
  PUBMED    9662389
REFERENCE   2 (bases 1 to 200)
  AUTHORS   Xie,S., Okano,M. and Li,E.
  TITLE     Direct Submission
  JOURNAL   Submitted (28-MAY-1998) CVRC, Mass. Gen. Hospital, 149 13th Street,
            Charlestown, MA 02129, USA
REFERENCE   3 (bases 1 to 200)
  AUTHORS   Okano,M., Chijiwa,T., Sasaki,H. and Li,E.
  TITLE     Direct Submission
  JOURNAL   Submitted (04-NOV-1999) CVRC, Mass. Gen. Hospital, 149 13th Street,
            Charlestown, MA 02129, USA
REMARK      Sequence update by submitter
COMMENT     On Nov 18, 1999 this sequence version replaced gi:3327977.
FEATURES             Location/Qualifiers
     source            1..200
                       /organism="Mus musculus"
                       /mol_type="mRNA"
                       /db_xref="taxon:10090"
                       /chromosome="12"
                       /map="4.0 cM"
     gene              1..>200
                       /gene="Dnmt3a"
ORIGIN
1 gaattccggc ctgctgccgg gccgcccgc ccgccgggcc acacggcaga gccgcctgaa
61 gccacgcgt gaggctgcac tttccgagg gcttgacatc agggctcatg ttttaagtctt
121 agctcttgct tacaagacc acggcaattc cttctctgaa gccctcgag cccacagcgc
181 ccctcgagc cccagcctgc
//

```

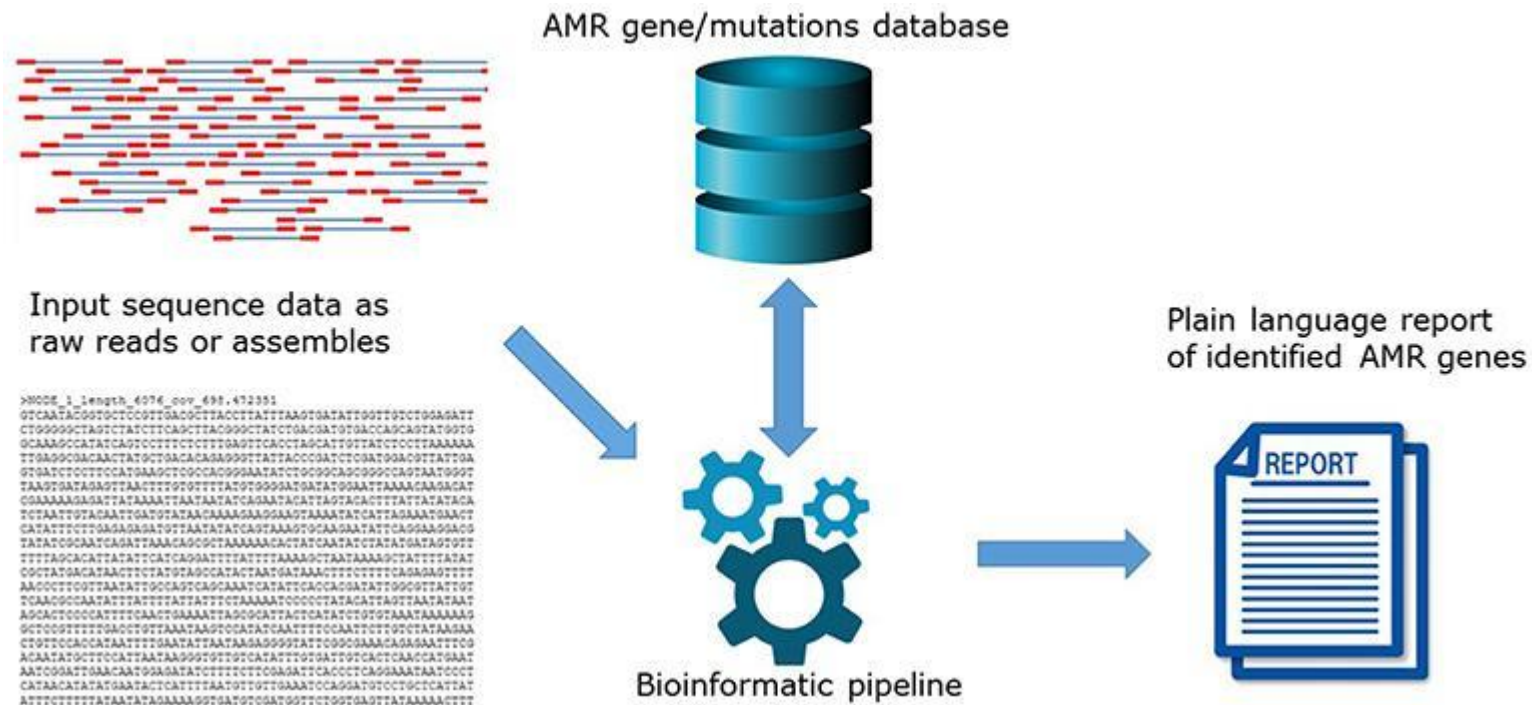

Annotation format: gbk

- LOCUS – Annotated sequence
- DEFINITION
- ACCESSION
- FEATURES
 - source
 - gene
 - CDS
 - Locus tag
 - function
 - Product
 - protein_id
 - Translation (sequence)

FEATURES	Location/Qualifiers
source	1..381113 /organism="Klebsiella pneumoniae subsp. pneumoniae SA1" /mol_type="genomic DNA" /strain="SA1" /sub_species="pneumoniae" /db_xref="taxon:1379688" /note="contig LPSB1_2557_Contig_49"
gene	415..1536 /locus_tag="KPST86_490001"
CDS	415..1536 /locus_tag="KPST86_490001" /inference="ab initio prediction:AMIGene:2.0" /note="Evidence 4:Homologs of previously reported genes of unknown function" /codon_start=1 /transl_table=11 /product="conserved hypothetical protein" /protein_id="C0I25656.1" /translation="MAYQLININWPEFLEKYWQKQPVVVKNAFPDFVDPITPDELGLA MEPEVDSRLVSLKNGKNQASNGPFEHFDGLGETGWSLLAQAVNHWMPAAELVRPFRV LPDWRLLDLMISFSVPGGGVGHIDQYDFIQQWIGSRRWRVGDKLPHRQFCPPALL HVDPPFPIIDEDLQPGDILYIPPGFPHDGIHETALNYSVGFGRGNRDLISSFADYV LENDLGDHYSDPDLTCREHPGRVEEYELRLRTHMIDMIRQPEDFKQWFGSFVTTTPR HELDIAPAEPPYEEVEVLDALLGGEKLSRLSGLRVLHIGDSFFVHSEQLDITDAEALD ALCRYTSLGQEELGSLQNPFAVSELRLINQGYNYFEE"
gene	complement(1584..2117) /locus_tag="KPST86_490002"
CDS	complement(1584..2117) /locus_tag="KPST86_490002" /inference="ab initio prediction:AMIGene:2.0" /note="Evidence 4:Homologs of previously reported genes of unknown function" /codon_start=1 /transl_table=11 /product="conserved hypothetical protein" /protein_id="C0I25658.1" /translation="MEQQLTIEMIADAFSYDITGFDGCEALNTFLKEHLKRQHDGQI LRGVALVSGDTPRLLGYITLGSFCFERGMLPSKTQQKKIPYQNPVTLGRLAIDKS VQQQGWGEMLVAAHMRVWVGASKAVGIYGLFVEALNEKAKAFYLRGLFIQLVDENSNL LFYPTKSIEQLFTDDES"
gene	complement(2128..2394) /locus_tag="KPST86_490003"
CDS	complement(2128..2394) /locus_tag="KPST86_490003" /inference="ab initio prediction:AMIGene:2.0" /note="Evidence 4:Homologs of previously reported genes of unknown function"

Resistance prediction using WGS

Hendrisken et al. *Frontiers in Microbiology*. 2019.



Resistance prediction using WGS

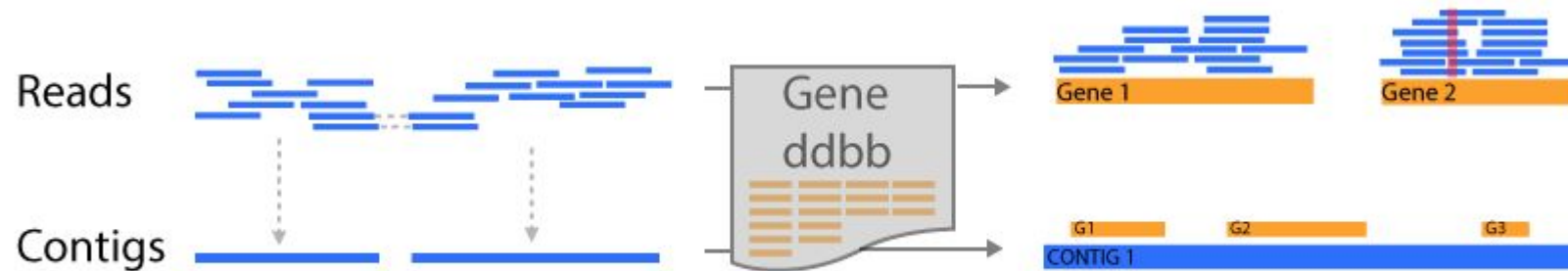
Hendrisken et al. *Frontiers in Microbiology*. 2019.

- Huge list here:
https://www.frontiersin.org/files/Articles/478239/fpubh-07-00242-HTML/image_m/fpubh-07-00242-t002.jpg

Software	Type
SRST2	Mapping
Ariba	Mapping + assembly
ABRICATE	Assembly
ResFinder	Assembly

Mapping vs Assembly

- **Functional annotation based on mapping (srst2)**
 - Pro: more resolute / high quality ddbb
 - Con: Unable to locate genes / no ab initio annotation
- **Functional annotation based on assembly (Resfinder)**
 - Pro: genes are located / related
 - Depend on assembly (close to repetitive regions)



Manual annotation: Artemis

Artemis is a DNA sequence viewer and annotation tool that allows visualisation of sequence features and the results of analyses within the context of the sequence, and its six-frame translation.



Thanks for your attention!

