

Bacterial gene and genome annotation

13/06/2024

Dr. Luu Phuc Loi (luu.p.loi@googlemail.com)

Content

1. Bacterial genomes
2. Bacterial Genes
3. Genome Annotation
4. Prokka
5. Curating genomes
6. Bakta

Bacterial genomes: small genomes

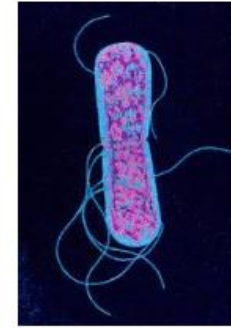


6,000,000,000
letters

30,000 genes



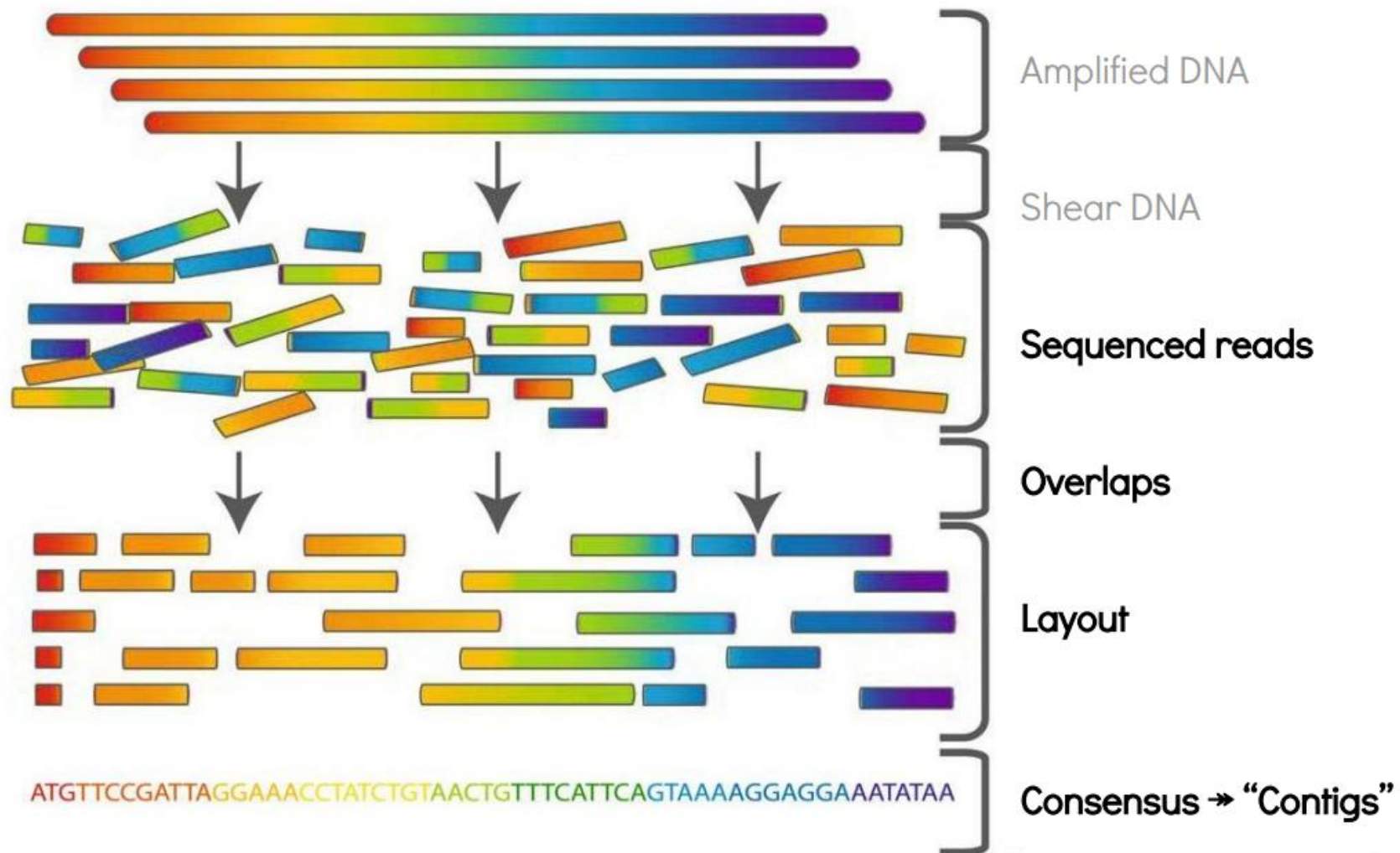
Genome
A T G C



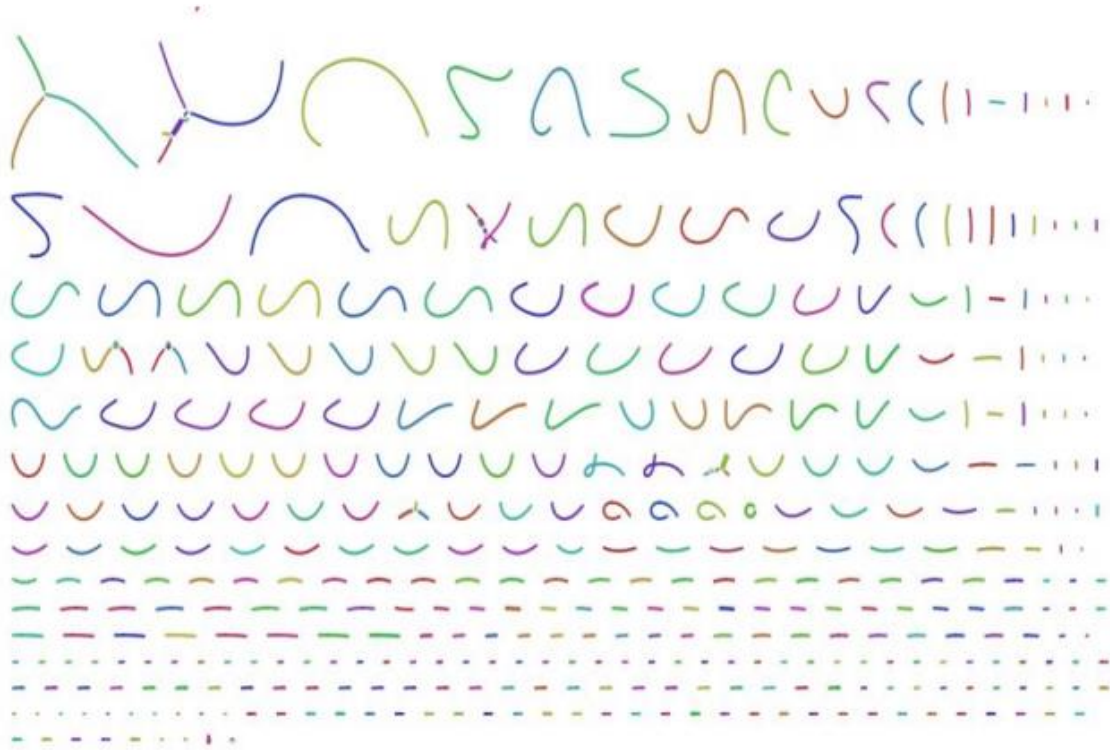
3,000,000
letters

3,000 genes

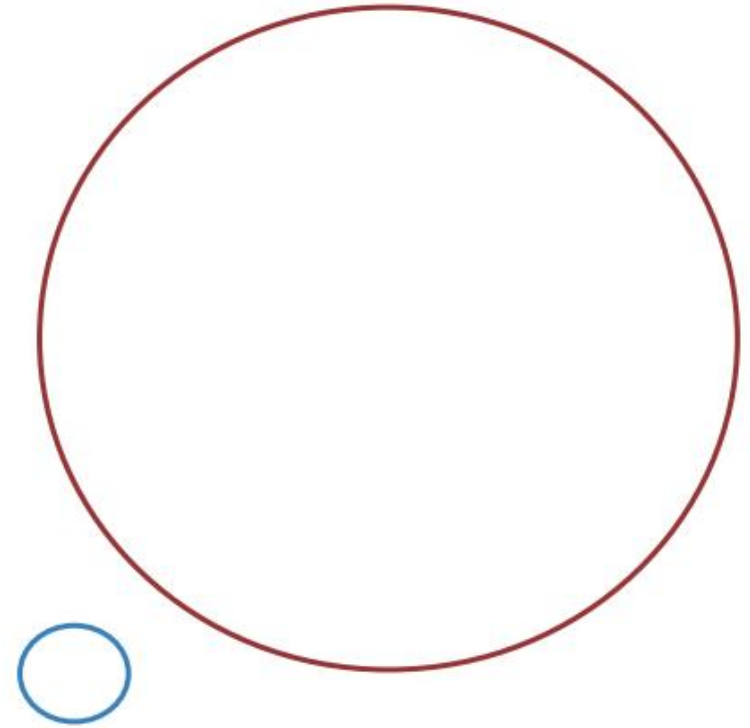
De novo assembly



Draft vs Finished genomes



Lots of contigs

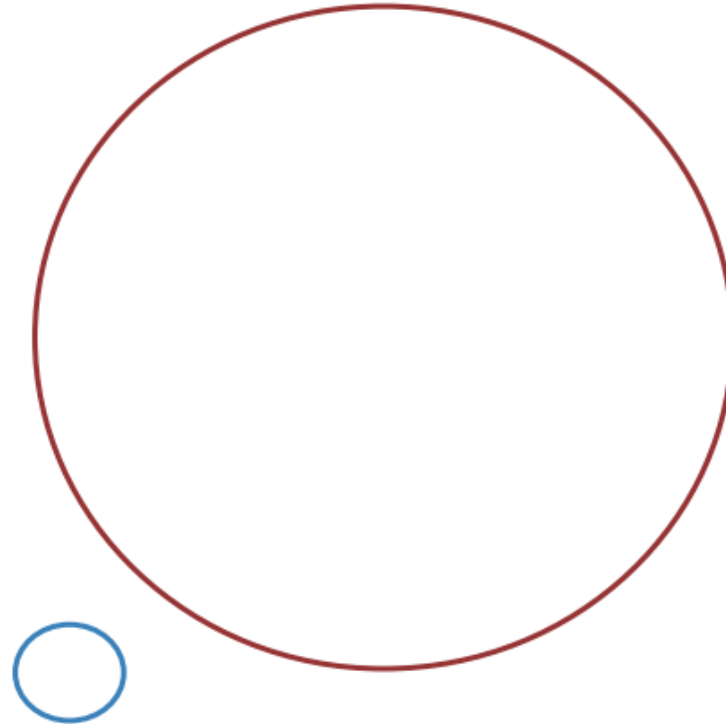


One contig per replicon

```

48541 agcccttcaa agaaatgttc tcagcaggca tggagcccag gacttgctcc ctttgggtag
48601 agagccgggt tgaaggtgac tgaagtgaat tgggacagta gaggcggggg ggggtggtgag
48661 ttcctggagg tggggggtgt gggaaacctgc tgtgtactga gatgcacccc tgccagttct
48721 gcctgaagat ttgaggcggg gggcaggggg gcggagtga gtcattttac tggtaaagtaa
48781 ttttaaacct tttaatatta aagcaaacgt ggatatgtaa tgaatgaaat tcattctgga
48841 atgaaaaatt cacgtgatgt tgaaaaataa cacggggcctt cagagaggac tttctggctg
48901 gcagcagact ccagattccc agggcccctg caccctcctc tgcccacagg gcaccttaat
48961 ggagaaggtg tgggaggaga gccaggccgg agtcagagca cactggtgac tccacatttg
49021 cagcgtgccc tgcctctctc ctgaggcttg gcaacgtgca atatgctaag caaactcccc
49081 ctgtccccgt ccagtttctg aggacaagag ccaccacctg tagcaaataa agaccagca
49141 accctttgac tcacttttgt gagtctcttg aatcagaggg tagccacatc gctgagagggt
49201 ggagtgaagc actcgggtga aaaggtacaa ggaagtcagg gacaggagtg tggggacatc
49261 acctagacaa tgacagagaa gaggggcaca gccgagttag gggagagggg ccggcagtc
49321 tacatcccct ggctgaagc acgctccagg gcagaaggaa aaacactgtc tttggggtcc
49381 aagagacctg agttcaaatt ctggctccac cactgaccac ctgtgtaacc ttgaactgct
49441 gctgcctgaa cctcagggtt cccttctaaa aatagaggag aaaaggatgc atttctcctt
49501 ccccctgtga gaacgaaatg gtgcaagcac caaggagcct cagcaaagggt cgggcctgcc
49561 cccgcctggc caaacctttc ctcttcagga ggccacggca accgtagttt gacagaagag
49621 cagcaccttg attaatgct tcccagcatg tgtccttgag caagtcacct aacctctctg
49681 ggctgttcc tcattgggaa aatatggctg ccagtaaaac ctgccctgtc cacctcctgg
49741 ggcacttggc aaacagcaaa agagtccaaa tgtgcaggct gggccaggcg cagtggctca
49801 tgcctgtaat cccagcaatt taggaagcca aggtgggcgg atcacctgag gtcaggagtt
49861 tgagaccagc ctggccaaca tggtgaaacc ttgtctctac aaaaatacaa aaattagccg
49921 ggcatgatgg cgggtgcctg taatcccagt tactcgggag gctgaggcaa gagaatcgct
49981 tgaaccggga aggggaaggt tgcagtgagc caagattgtg ccactgcact ccagcctggg
50041 caacagagcg agactctgtc tcaaaaaaaa aaaaaaaaaa aaacaatgca gagctggctg
50101 tgtaaaaaac ctgttcact gcaggggcca gtgtccacca ggctgggggtg caggcctatg
50161 ggggtggggc ccagcatcag cctctcagca gccctgggag gcggggcgca tcccgtgcc
50221 ctctgtgtct ggtgtgttc tagcccaagt cctaggttac acctgccgtc gcctggcctc
50281 tcaggagagg cccagggtga ggaggagcat ggtaaagggt aagctgattg ggaagtcagc
50341 tgttgggaaa gcaactcctt gcacattgga ggaaccgaga aagactgacc ccgaggacag
50401 cagccagcat ggccttctt gggagcccat gttgggggat tcttgctgca gccaaggctc
50461 agcccttgtg gtcgcagggt ctgttccttg cctcttcccc tcccatgagc gagcacagga
50521 gagatggctt ctgaggacct gttgcagctg tggccctggg aatagatttg ccaggagct
50581 ttaaagcagc tgagtgtgtc atccagctaa gcctggggaa ggagcttggc tcaggctctg
50641 acagggtgta cagggatggg gactgggaag taagagatga aacctggct ggaggctgtg
50701 agcttcaca gccagcgtg gcagggagg gtccagatat acccactagt gccctacca

```



One contig per replicon

Genome Annotation: Adding biological info to sequences

```
48541 agcccttcaa agaaatgttc tcagcaggca tggagcccag gacttgctcc ctttgggttag
48601 agagccgggt tgaaggtgac tgaagtgaat tgggacagta gaggcggggg ggggtggttag
48661 ttccctggagg tggggggtgt gggaaacctgc tgtgtactga gatgcacccc tgccagttct
48721 gcctgaagat ttgaggcggg gggcaggggg gcgagtgaa gtcattttac ttgtaagtaa
48781 ttttaaacct ttttaattta aagcaaacgt ggaatgtaa tgaatgaaat tcattctgga
48841 atgaaaaaatt cacgtgatgt tgaataataa cagggggctt cagagaggac tttctggctg
48901 gcagcagact ccagattccc agggccctcg caccctctc tccccacagg gcaccttaat
48961 ggagaaggtg tgggaggaga gccaggccgg agtcagagca cactggtgac tccacatttg
49021 cagcgtgccc tgctctctc ctgaggcttg gcaacgtgca atatgctaag caaacctccc
49081 ctgtccctcg ccagtttctg aggacaagag ccaccacttg tagcaataaa agaccagca
49141 accctttgac tcattcttgt gactctctg aatcagaggg tagccacatc gctgagaggt
49201 ggagtgaagc actcgggtga aaaggtacaa ggaagtcagg gacagagtg tggggacatc
49261 acctagacaa tgacagagaa gaggggcaca gccgagtgag gggagagggg ccggcagttc
49321 tacatccctt ggcctgaagc acgtctcagg gcagaaggaa aaacactgtc tttgggttcc
49381 aagagacctg agttcaaat ctggctccac cactgaccac ctgtgtaacc ttgaactgct
49441 gctgcctgaa cctcaggttc ccttctctaa aatagaggag aaaaggtatg atttctctt
49501 gccctgtga gaacgaatg gtgcaagcac caaggagcct cagcaaaagt cgggcctgcc
49561 cccgcctggc caaaccttct ctcttcagga ggccacggca accgtagttt gacagaagag
49621 cagcaccctg atttaagtct tcccagcatg tgtccttgag caagtacact aacctctctg
49681 ggctgtcttc tcattgggaa aatatggctg ccagtaaaac ctgccctgtc caccctctgg
49741 ggcaacttggc aaacagcaaa agagtccaaa tgtgcaggct gggccaggcg cagtggctca
49801 tgcctgtaat cccagcaatt taggaagcca aggtgggctg atcacctgag gtcaggaggt
49861 tgagaccagc ctggccaaca tggtgaaacc ttgtctctac aaaaatacaa aaattagccg
49921 ggcatgatgg cgggtgcctg taatccagat tactcgggag gctgaggcaa gagaatcgct
49981 tgaaccggga aggggaaggt tgcagtgaac caagattgtg ccactgcact ccagcctggg
50041 caacagagcg agactctgtc tcaaaaaaaa aaaaaaaa aaacaatgca gagctggctg
50101 tgtaaaaaac ctgttcact gcaggggcca gtgtccacca ggctggggtg caggcctatg
50161 ggggtggggc ccagcatcag cctctcagca gccctgggag gcggggcgca tcccgtgccc
50221 ctgctggctc ggaatgttct tagcccaagt cctaggttac acctgcctc gctgcgctc
50281 tcaggagagg cccagggtga ggaggagcat ggtaaaggtg aagctgattg ggaagtgcg
50341 tgttgggaaa gcaactcctt gcacattgga ggaaccgaga aagactgacc ccgaggacag
50401 cagccagcat ggcccttctt gggagcccat gttgggggat tctgtctgca gccaaggctc
50461 agcccttgtg gtgcaggtg ctgtctctg cctcttccc tccatgcag gagcacagga
50521 gagatggctt ctgaggacct gttgcagctg tggccctggg aatagatttg ccagggagct
50581 ttaaagcagc ttagtgtgtc atccagctaa gcctggggaa ggagcttggc tcaggtcctg
50641 acaggtgtga cagggatggg gactgggaag taagagatga aacctggct ggaggtctg
50701 agcttcacca gccagcgtg gacaggagg gtcagatat acccactagt gccctacca
```

ribosome
binding site

delta toxin
PubMed: 15353161

```
ACCGGCCGAGACAGCGAGCATATGCAGGAAGCGGCAGGAATAAGGA
AAAGCAGCCTCCTGACTTTCCTCGCTTGGTGGTTTGAGTGGACCTC
CCAGGCCAGTGGCGGGCCCCCTCATAGGAGAGGAAGCTCGGGAGGTG
GCCAGGCGGCAGGAAGGCGCACCCCCCAGCAATCCGCGCGCCGGG
ACAGAATGCCCTGCAGGAACCTTCTTCTAGAAGACCTTCTCCTCCTG
CAAATAAAACCTCACCCATGAATGCTCACGCAAGTTTAATTACAGA
CCTGAAACAAGATGCCATTGTCCCCCGGCCTCCTGCTGCTGCTGCT
CTCCGTCCGTCCGTGGGCCACGGCCACCGCTTTTTTTTTTTGCC
```

transfer RNA
Leu-(UUR)

tandem repeat
CCGT x 3

homopolymer
10 x T

What's in an annotation?

- Location

- which sequence? *chromosome 2*
- where on the sequence? *100..659*
- what strand? *-ve*

- Feature type

- what is it? *protein coding gene*

- Attributes

- protein product? *alcohol dehydrogenase*
- enzyme code? *EC:1.1.1.1*
- subcellular location? *cytoplasm*
- note? *beer processing*

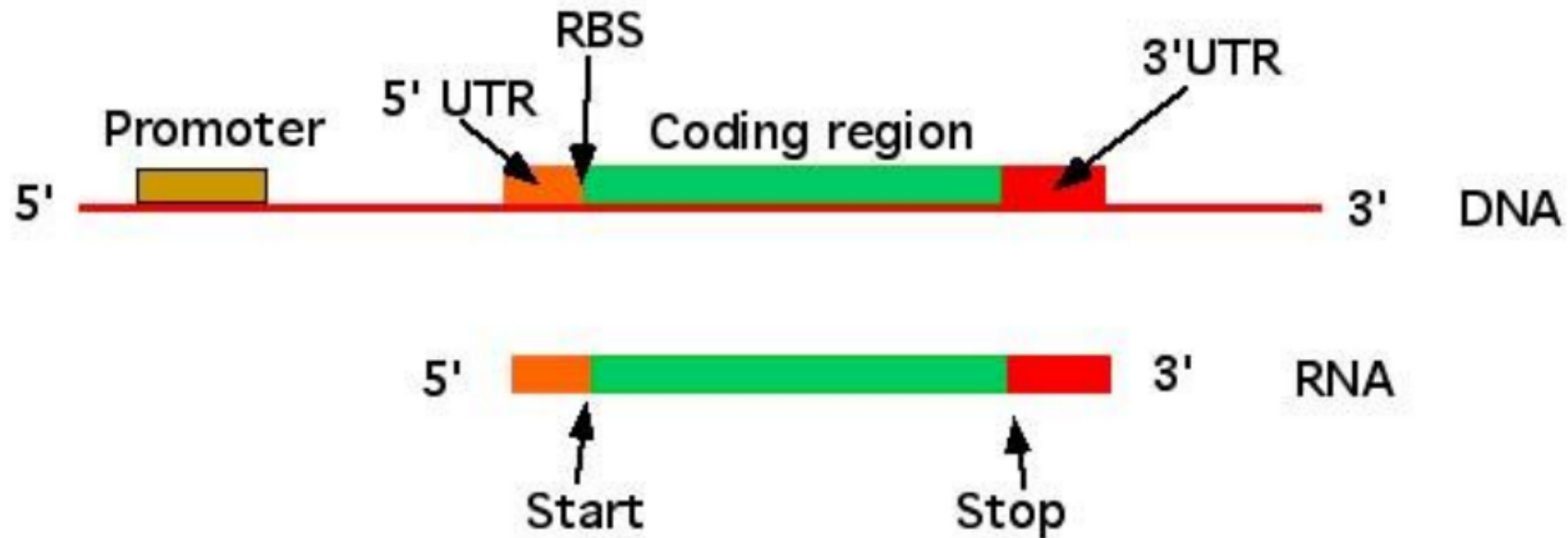
Bacterial feature types

- protein coding genes
 - promoter (-10, -35)
 - ribosome binding site (RBS)
 - coding sequence (CDS)
 - signal peptide, protein domains, structure
 - terminator
- non coding genes
 - transfer RNA (tRNA)
 - ribosomal RNA (rRNA)
 - non-coding RNA (ncRNA)
- other
 - repeat patterns, operons, origin of replication, ...

Key bacterial features

- 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
 - partial genes, remnants
 - pseudogenes
 - assigning function is the bulk of the workload

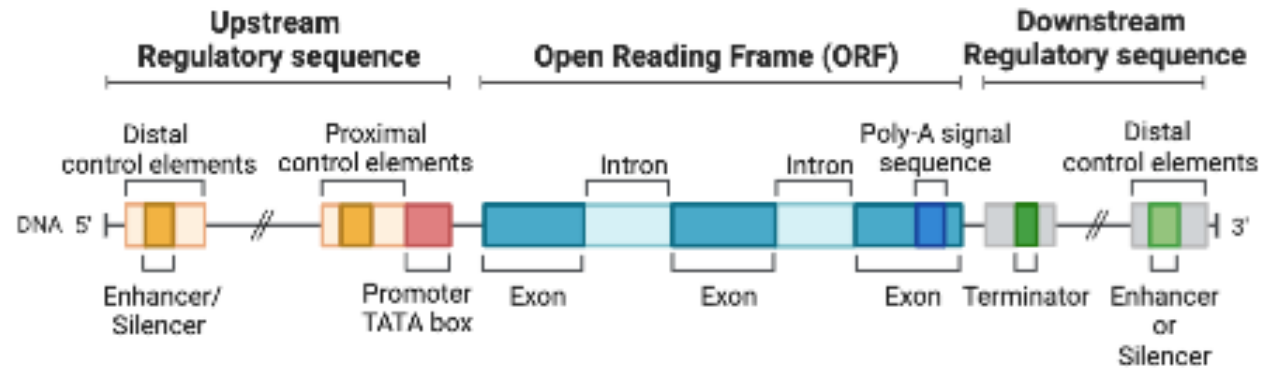
Bacterial genes: no introns!



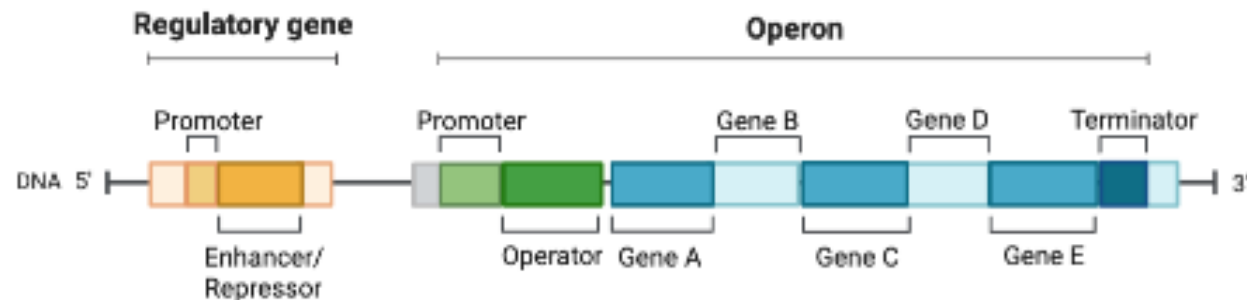
- have ≥ 3 potential start codons (species dependent)
- haploid, but lots of horizontal gene transfer
- methylation used as primitive immune system
 - restriction modification system against phage

Prokaryotic vs eukaryotic Genes

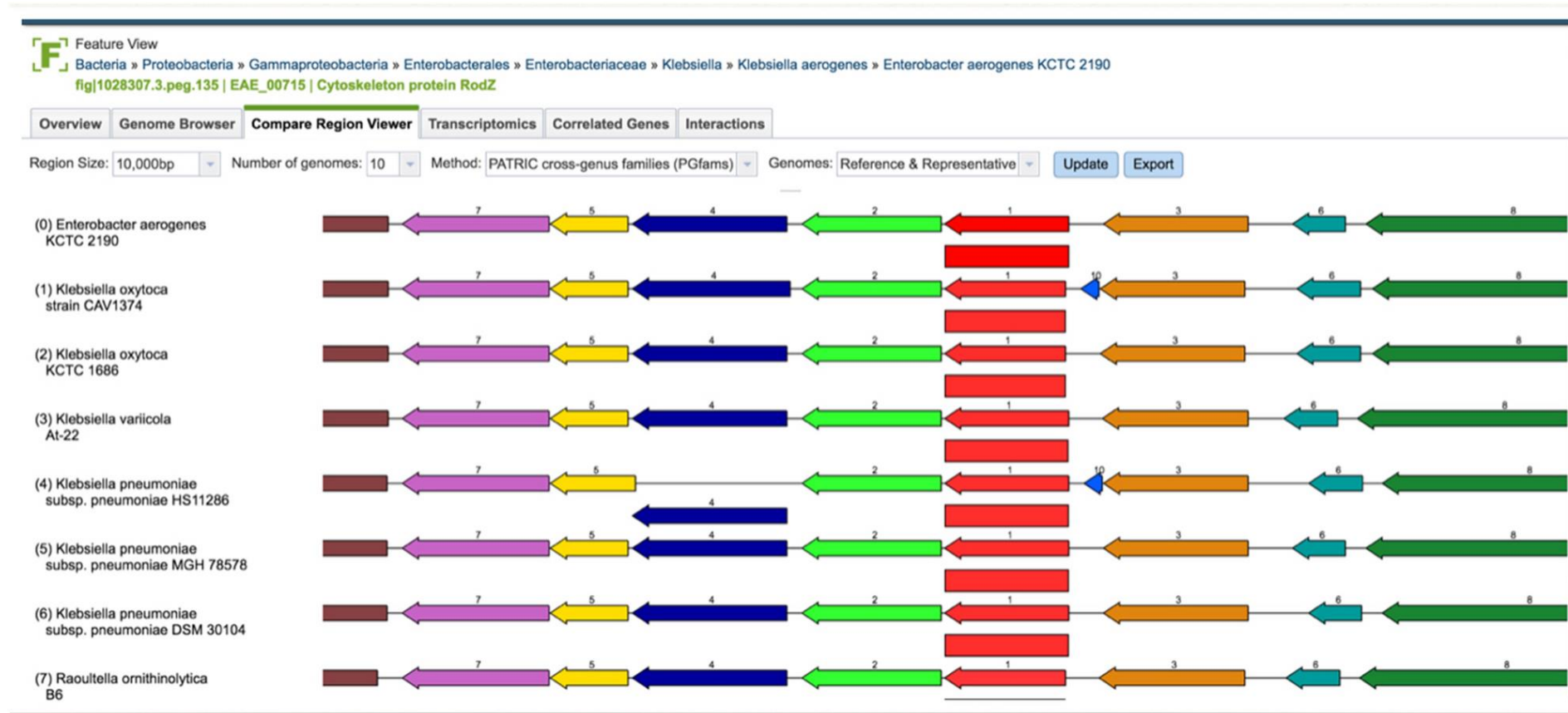
Eukaryotic Gene Structure



Prokaryotic Gene Structure



Detecting operons in bacterial genomes via visual representation learning



Snapshot of the Compare Region Viewer service provided by PATRIC (<https://www.patricbrc.org>). The image shows a genomic region of the query genome (first row) aligned against a set of other genomes, anchored at the focus gene (represented as a red arrow). The service starts with finding other genes that are of the same family as the focus gene, and then aligns their flanking regions accordingly.

Automatic Genome Annotation

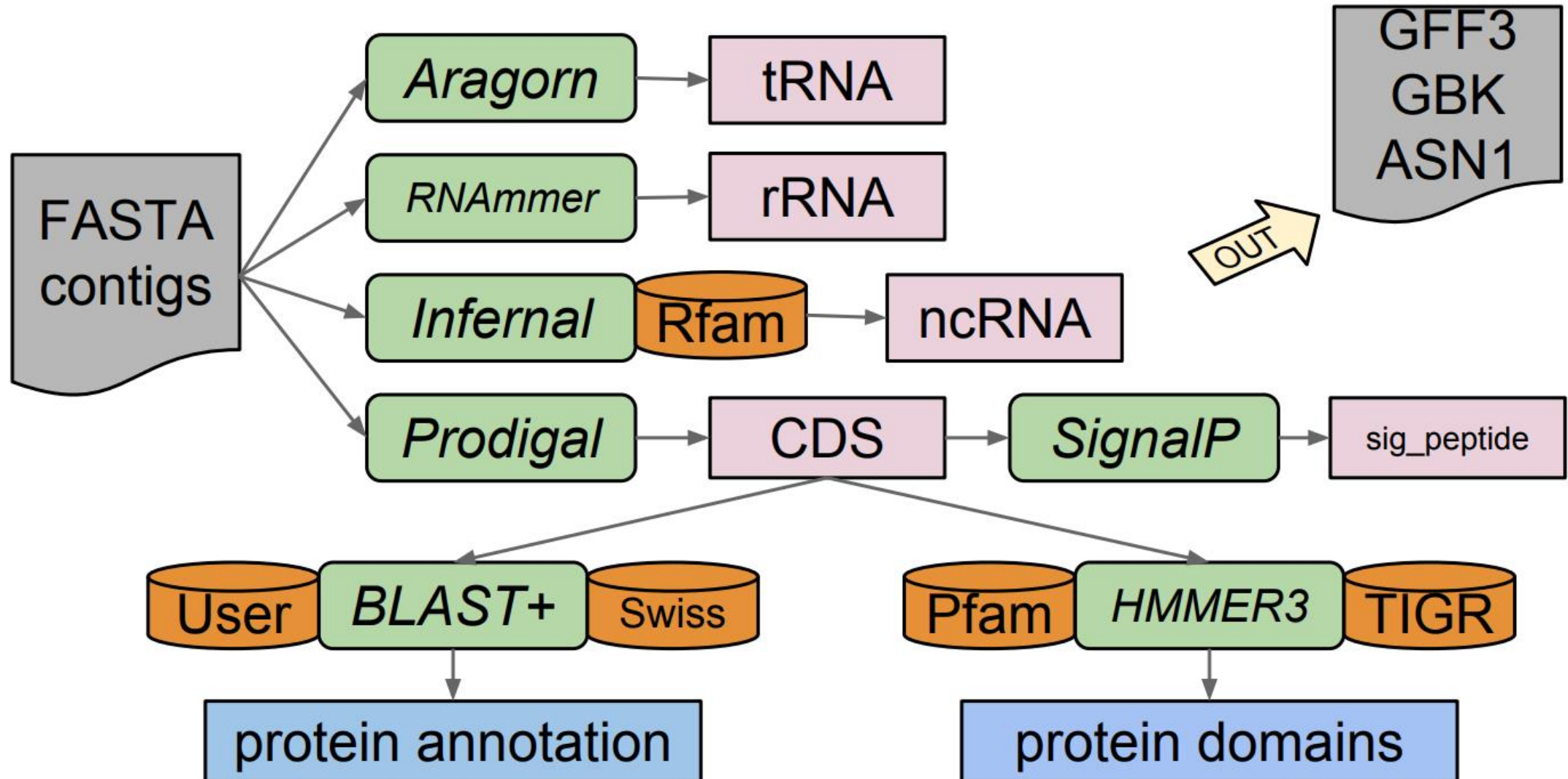
Two strategies for identifying coding genes:

- **sequence alignment**
 - find known protein sequences in the contigs
 - transfer the annotation across
 - will miss proteins not in your database
 - may miss partial proteins
- ***ab initio* gene finding**
 - 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

Some good existing tools

Software	<i>ab initio</i>	align- ment	Availability	Speed
RAST	yes	yes	web only	12-24 hours
xBASE	yes	no	web only	>4 hours
BG7	no	yes	standalone	>10 hours
PGAAP (NCBI)	yes	yes	email / we	>1 month

Prokka pipeline (simplified)



Prokka pipeline (simplified)

Bioinformatics Advance Access published March 18, 2014

Genome Analysis

Prokka: rapid prokaryotic genome annotation

Torsten Seemann^{1,2,*}

¹ Victorian Bioinformatics Consortium, Monash University, Melbourne, Australia.

² Life Sciences Computation Centre, Victorian Life Sciences Computation Initiative, Melbourne, Australia.

Associate Editor: Prof. Alfonso Valencia

ABSTRACT

Summary: The multiplex capability and high yield of current day DNA sequencing instruments has made bacterial whole genome sequencing a routine affair. The subsequent *de novo* assembly of reads into contigs has been well addressed. The final step of annotating all relevant genomic features on those contig can be achieved slowly using existing web and email-based systems, but these are not applicable for sensitive data or integrating into computational pipelines. Here we introduce Prokka, a command line software tool to fully annotate a draft bacterial genome in about ten minutes on a typical desktop computer. It produces standards-compliant output files for further analysis or viewing in genome browsers.

Availability and Implementation: Prokka is implemented in Perl and is freely available under an open source GPLv2 license from <http://vicbioinformatics.com/>.

Contact: torsten.seemann@monash.edu

2 DESCRIPTION

2.1 Input

Prokka expects pre-assembled genomic DNA sequences in FASTA format. Finished sequences without gaps are the ideal input, but it is expected that the typical input will be a set of scaffold sequences produced by *de novo* assembly software. This sequence file is the only mandatory parameter to the software.

2.2 Annotation

Prokka relies on external feature prediction tools to identify the coordinates of genomic features within contigs. These tools are listed in Table 1, and all of them, except for Prodigal, provide co-

Example from Prokka

Feature Type:

tRNA

Location:

contig000341 @ 655..730 +

Attributes:

/gene="tRNA-Leu (UUR) "

/anticodon=(pos:678..680,aa:Leu)

/product="transfer RNA-Leu (UUR) "

/inference="profile:Aragorn:1.2"

Provenance

Recording where an annotation came from

Prokka uses Genbank "evidence qualifier" tags:

Wet lab

```
/experiment="EXISTENCE:Northern blot"
```

Dry lab

```
/inference="similar to DNA sequence:INSD:AACN010222672.1"
```

```
/inference="profile:tRNAscan:2.1"
```

```
/inference="protein motif:InterPro:IPR001900"
```

```
/inference="ab initio prediction:Glimmer:3.0"
```

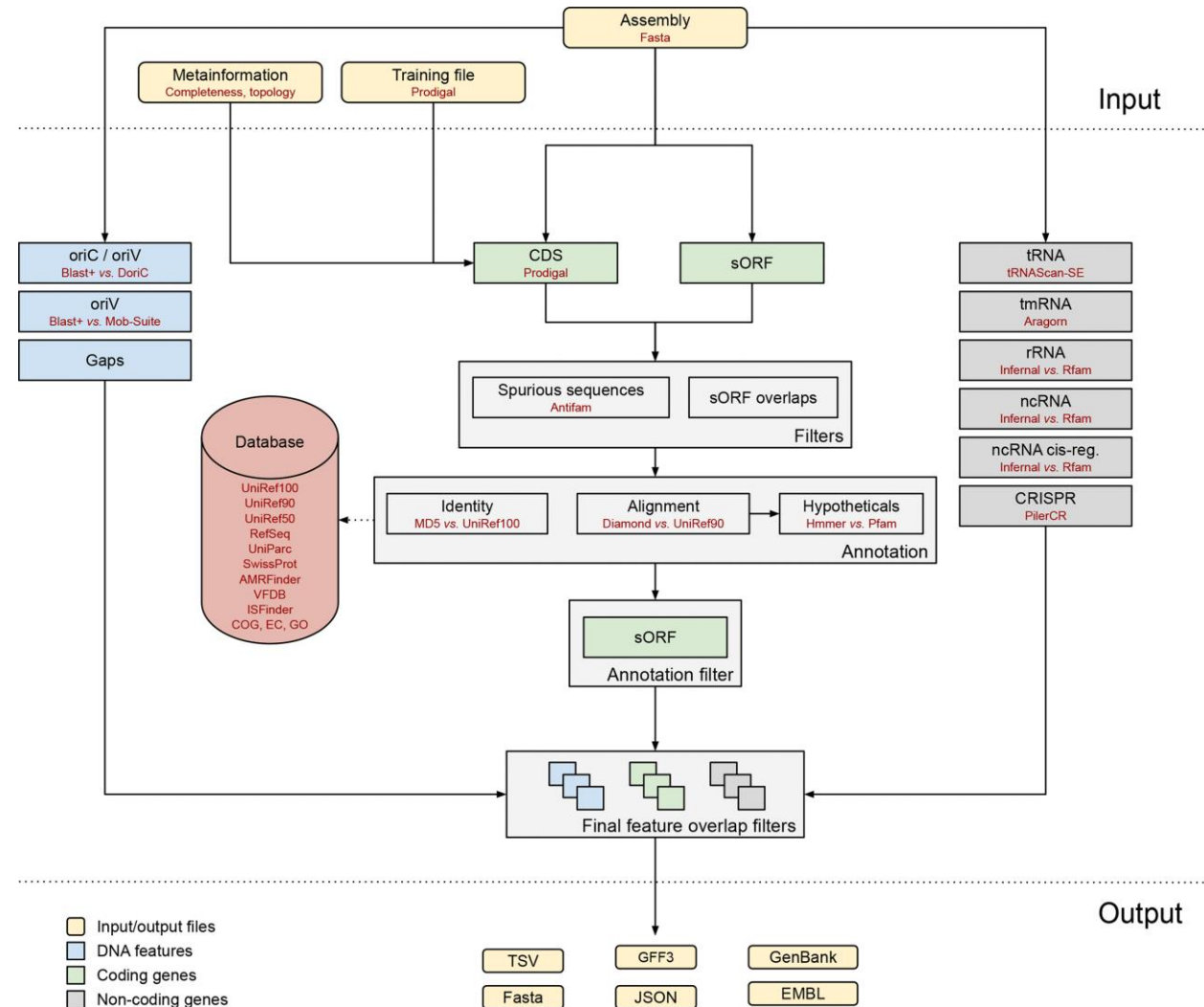
Curating genomes: Improving annotations

- Some annotations are wrong
 - False annotation
 - Missing annotation
 - Partially wrong annotation
- Curation
 - Manual effort to improve annotations
 - Community curation

Web Apollo: a web-based genomic annotation editing platform



Bakta is a tool for the rapid & standardized annotation of bacterial genomes and plasmids from both isolates and metagenome-assembled genomes (MAGs)




Overview of the Bakta annotation workflow ([Schwengers et al. 2021](#))

Further structural annotation

- **PlasmidFinder** ([Carattoli and Hasman 2020](#)), a tool for the identification and typing of plasmid sequences in Whole-Genome Sequencing.
- **IntegronFinder** ([Néron et al. 2022](#)), a tool for detecting integrons
- An **integron** is minimally composed of:
 - ❑ a gene encoding for a site-specific recombinase (intI)
 - ❑ a proximal recombination site (attI), which is recognized by the integrase and at which gene cassettes may be inserted
 - ❑ a promoter (P_c) which directs transcription of cassette-encoded genes
- **IS (Insertion Sequence)** elements: a short DNA sequence that acts as a simple transposable element
- **ISEScan** ([Xie and Tang 2017](#)) for detection of IS

Exercise 1: Using RAST to annotate bacterial genome sequence as following

<https://www.youtube.com/watch?v=plCJ4Gmqq1U>




RAST

Rapid Annotation using
Subsystem Technology

version 2.0

The NMPDR, SEED-based, prokaryotic genome annotation service.
For more information about The SEED please visit theSEED.org.

[»Home](#) [»Your Jobs](#) [»Tutorials](#) [»Help](#)

 Ala

Info: **RAST Access Problems**

A number of users have recently reported problems accessing their data. There have been two primary causes.

First, they have been going to <http://rast.theseed.org/>; please note this is not RAST's canonical address, it is an alias. RAST's canonical address is <https://rast.nmpdr.org/>, and it is this canonical address that will provide you with the most reliable access to RAST.

Second, we have recently converted RAST from using unencrypted HTTP to encrypted HTTPS, and many users appear to have not yet updated their bookmarks. So it is possible that you are having problems due to using a stale URL or bookmark.

So as a first step, we recommend trying the following:

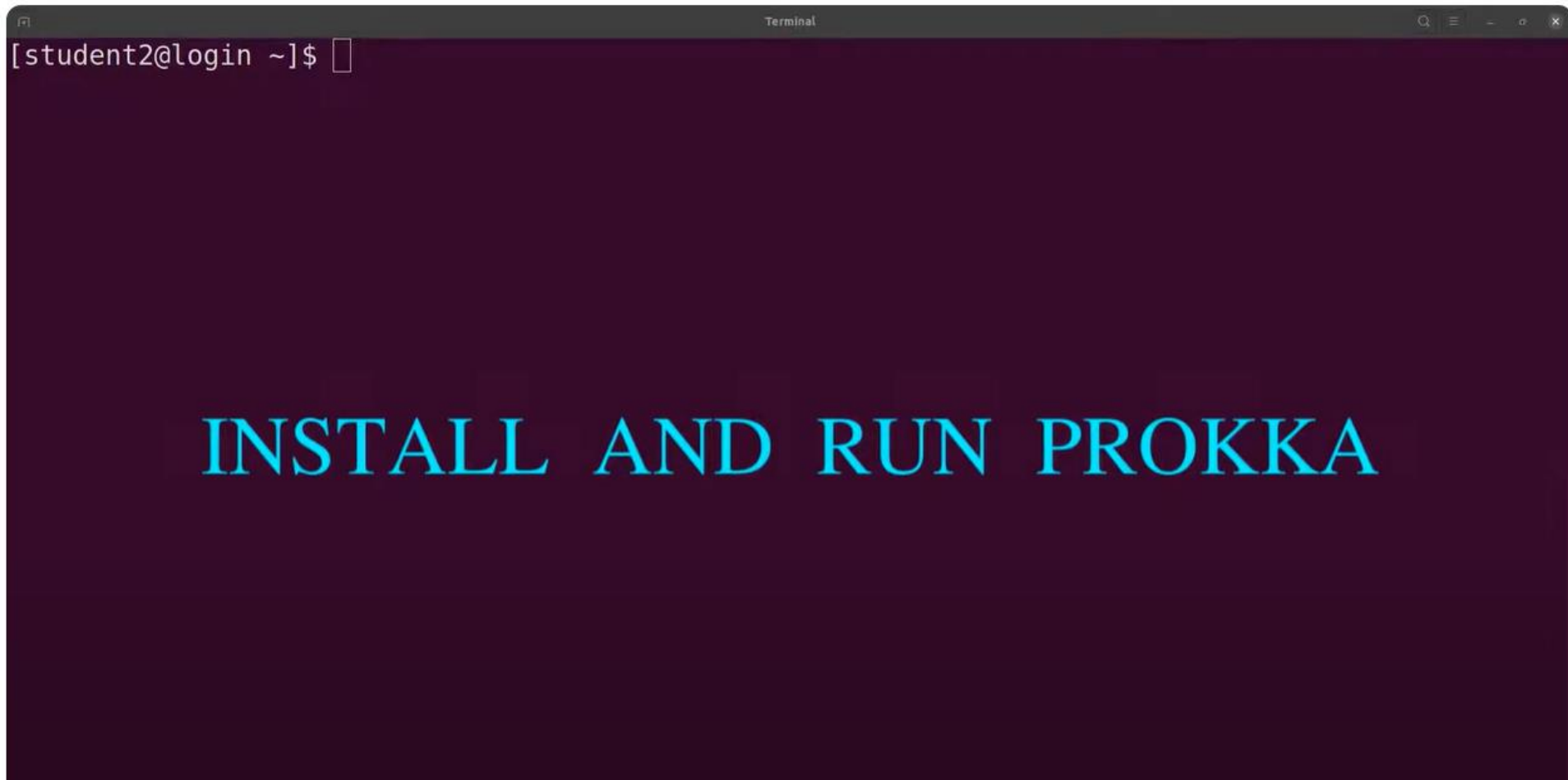
1. Make sure that you are logged out of RAST and that all RAST webpages are closed.
2. Update your RAST bookmark if you have one to the new secure encrypted URL: <https://rast.nmpdr.org/>. Again, please note that the new secure RAST URL now begins with 'https', not 'http'!
3. Clear your web browser's web cache and cookies, as they may contain stale copies of the RAST webpages
4. Repeat your login attempt.

If the above procedure does not resolve your problem, please contact us via the RAST Administrator address, rast@mcs.anl.gov, and we will investigate further.

To monitor RAST's load and view other news and statistics for RAST and the SEED, please visit "[The Daily SEED](#)."

As of Fri Jun 19 05:00:18 2020, there are 14 jobs in the RAST queue

Exercise 2: Using PROKKA to annotate bacterial genome sequence as following
<https://www.youtube.com/watch?v=NYbYxGf-8NM&list=PLe1-kjuYBZ04YmORKfXBCOqGDDYJtudVA&index=8>



Thank you for your attention!