

Module 1: Artemis

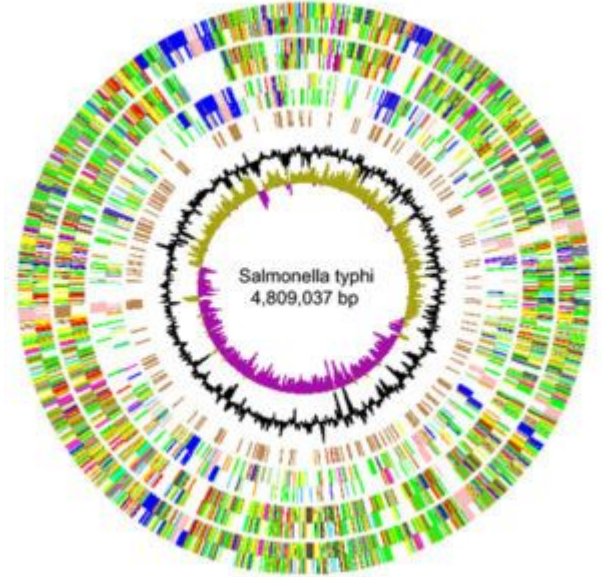
Lionel Uran Landaburu

Instituto de Investigaciones Biotecnológicas
lionel.u.l@iib.unsam.edu.ar

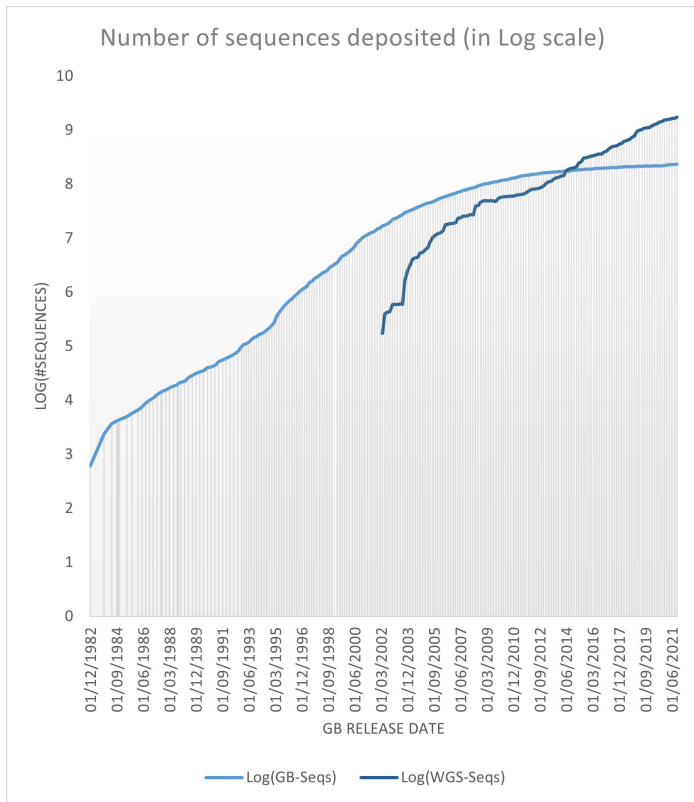
Rodrigo Quiroga

Instituto de Investigaciones Biotecnológicas
lionel.u.l@iib.unsam.edu.ar

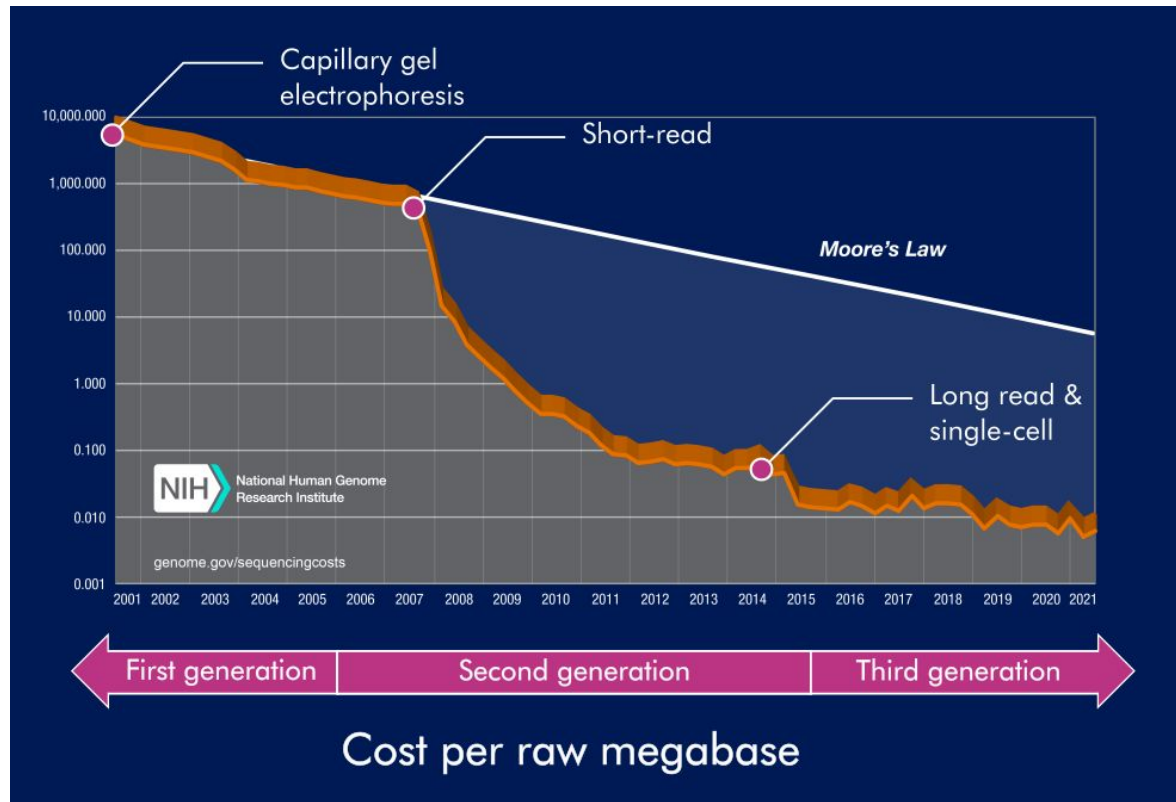
Working with Pathogen Genomes 7-11th, February 2022.



Next generation sequencing, a paradigm shift



Source: GenBank statistics
(<https://www.ncbi.nlm.nih.gov/genbank/statistics/>)



Adapted from genome.gov/sequencingcosts & Athina Gkazi (2021). An Overview of Next-Generation Sequencing. Technology Networks: Genomics Research.

We have sequences... now what?



See anything yet?

Entry: ☒ sequence.fasta

Nothing selected

I K G L Y L P R # Q T N Q L S I S C R S V L # T N F K I C V A V T R L H A + C T H A V # L I T N Y C R * Q D T S N S S I
L K V Y T F P G N K P T N F R S L V D L F S K R T L K S V W L S L G C M L S A L T Q Y N # # L I T V V D R T R V T R L S
. # R F I P S Q V T N Q P T F D L L + I C S L N E L # N L C G C H S A A C L V H S R S I I N N # L L S L T G H E # L V Y L
A T T A A A G G T T T A T A C C T T C C C A G G T A A C A A A C C A C C A C T T T C G A T C T C T T G T A G A T C T G T T C T C T A A A C G A A C T T T A A A A T C T G T G G C T G T C A C T C G G C T G C A T G C T T A G T G C A C T C A C G C A G T A T A A T T A A T A A C T A A T T A C T G T C G T T G A C A G G A C A C G A G T A A C T C G T C T A T C T A A T T T C C A A A T A T G G A A G G G T C C A T T G T T T G G T T G G T G A A A G C T A G A G A A C A T C T A G A C A A G A G A T T T G C T T G A A A T T T A G A C A C A C C G A C A G T G A G C C G A C G T A C G A A T C A C G T G A G T G C G T C A T A T T A A T T A T T G A T T A A T G A C A G C A A C T G T C C T G T G C T C A T T G A G C A G A T A G A
N F T # V K G P L L G V L K R D R T S R N E L R V K F D T H S D S P Q M S L A S V C Y L # Y S I V T T S L V R T V R R D .
L N I G E W T V F W G V K S R K Y I Q E R F S S # F R H P Q * E A A H K T C E R L I I L L + N S D N V P C S Y S T + R
L P K Y R G L Y C V L W S E I E Q L D T R + V F K L I Q T A T V R S C A # H V * A T Y N I V L # Q R Q C S V L L E D I H

How about now? Does it ring a bell?

Entry: ☒ sequence.gb

Selected feature: bases 3822 amino acids 1273 S (/gene="S" /locus tag="GU280 gp02" /gene synonym="spike glycoprotein" /codon start=1 /product="sur")

The genomic map displays the S gene with various Open Reading Frames (ORFs) and features. The main ORF is labeled 'S' and is highlighted in red. Other ORFs include ORF1ab, ORF1, ORF2, ORF3, ORF4, ORF5, ORF6, ORF7a, ORF7b, ORF8, ORF9, ORF10, and ORF10 loop. The map also shows the 5' UTR, 3' UTR, and various other features like ORF1ab, ORF1, ORF2, ORF3, ORF4, ORF5, ORF6, ORF7a, ORF7b, ORF8, ORF9, ORF10, and ORF10 loop.

5'UTR: ab ab ORF1ab ORF1ab ORF1 ORF2 ORF3 ORF4 ORF5 ORF6 ORF7a ORF7b ORF8 ORF9 ORF10 ORF10 loop

2200 4400 6600 8800 11000 13200 15400 17600 19800 22000 24200 26400 28600

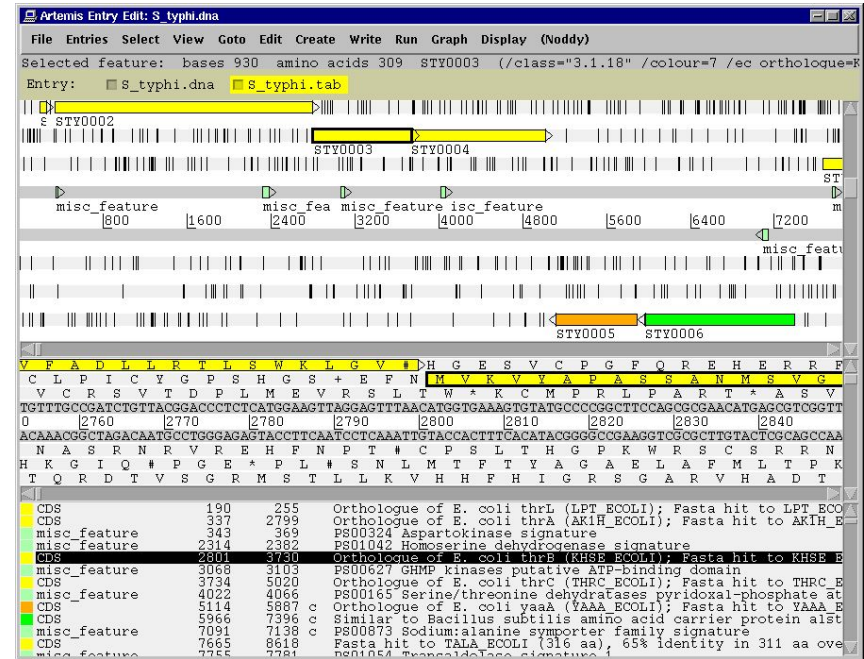
I K G L Y L P R # Q T N Q L S I S C R S V L # T N F K I C V A V T R L H A + C T H A V # L I T N Y C R * Q D T S N S S I F C R L L
L K V Y T F P G N K P T N F R S L V D L F S K R T L K S V W L S L G C M L S A L T Q Y N # # L I T V V D R T R V T R L S S A G C
R F I P S Q V T N Q P T F D L L + I C S L N E L # N L C G C H S A A C L V H S R S I I N N # L L S L T G H E # L V Y L L Q A A
ATTAAAGGTTTATACCTTCCCAGGTAAACAAACCAACCAACTTTCGATCTCTTGAGATCTGTTCTCTAAACGAACCTTTAAAATCTGTGTGGCTGTCACTCGGCTGCATGCTTAGTGCACTCACGCAGTATAATTAATAACTAATTACTGTCGTTGACAGGACACGAGTAACCTCGTCTATCTTCTGCAGGCTGCT
TAATTTCCAAATATGGAAGGGTCCATTGTTGGTTGGTTGAAAGCTAGAGAACATCTAGACAAGAGATTTGCTTGAAATTTAGACACACCGACAGTGAGCCGACGTACGAATCACGTGAGTGCGTCATATTAAATTATTGATTAAATGACAGCAACTGTCTGTGCTCATTGAGCAGATAGAAAGACGTCGACGA
NFT # V K G P L L G V L K R D R T S R N E L R V K F D T H S D S P Q M S L A S V C Y L # Y S I V T T S L V R T V R R D E A P Q K
L N I G E W T V F W G V K S R K Y I Q E R F S S # F R H P Q * E A A H K T C E R L I I L L + N S D N V P C S Y S T + R R C A A
L P K Y R G L Y C V L W S E I E Q L D T R + V F K L I Q T A T V R S C A # H V * A T Y N I V L # Q R Q C S V L L E D I K Q L S S

Feature	Start	End	Description
CDS	21563	25384	structural protein; spike protein
gene	25393	26220	
CDS	25393	26220	
gene	26245	26472	
CDS	26245	26472	ORF4; structural protein; E protein

Genome browser and annotation tool



- Visualization of sequence
 - DNA
 - six frame translation
 - Panoramic and sequence view
- Annotation
 - Features
 - Mapped and listed
 - Editable
 - In layers (entry)
- Perform and view analysis
 - basic analysis
 - Basic stats & index can be plotted
 - import and view the results of other searches/analysis
 - Different lines of evidence can be seen together



Can be used simply as a sequence viewer allowing the visualization of sequence and annotation taken directly from EMBL or GeneBank.

Files in Artemis



Sequence	Annotation
.fasta	.tab
.seq	
.dna	
	.embl

It can read several file formats (FASTA, EMBL, GENE BANK, GFF).

EMBL

Two-character line code indicates the type of information contained in the line

```
ID      ECRSMA      standard; DNA; PRO; 500 BP.
XX
AC      L40173;
XX
SV      L40173.1
XX
DT      10-AUG-1995 (Rel. 44, Created)
DT      04-MAR-2000 (Rel. 63, Last updated, Version 4)
XX
DE      Erwinia carotovora repressor (rsmA) gene, complete cds.
XX
KW      repressor; rsmA gene.
XX
OS      Pectobacterium carotovorum
OC      Bacteria; Proteobacteria; Gammaproteobacteria; Enterobacteriaceae;
OC      Pectobacterium.
XX
RN      [1]
RP      1-500
RA      Cui Y., Chatterjee A., Liu Y., Dumenyo C.K., Chatterjee A.K.;
RT      "Identification of a global repressor gene, rsmA, of Erwinia carotovora
RT      subsp. carotovora that controls extracellular enzymes,
RT      N-(3-oxohexanoyl)-L-homoserine lactone, and pathogenicity in soft-rotting
RT      Erwinia spp";
RL      J. Bacteriol. 177(17):0-0(1995).
XX
DR      GOA; Q47620; Q47620.
DR      SWISS-PROT; Q47620; CSRA_ERWCA.
```

Feature Key

Qualifier

Key	Location/Qualifiers
source	1..500 /db_xref="taxon:554" /organism="Pectobacterium carotovorum" /strain="71" /sub_species="carotovora"
-10_signal	107..112 /gene="rsmA"
RBS	235..239 /gene="rsmA"
CDS	246..431 /codon_start=1 EIQYRIQAEKSQPTSY"

```
XX
SQ      Sequence 500 BP; 140 A; 101 C; 120 G; 139 T; 0 other;
      ggatccggcga agcaggatag aaagtgtgtt accttcagat attctgaagc ttatcatgct      60
      cagttcttgtt gttgtgataa caaaagcaca agctactgat atcgactaaa ctaacaagta      120
      gtgacaaacc ggagttgtgat ggtgtgtgta taccatcgctc taggtttacg ttttcacagc      180
      acatgatgga taatggcggg gagacagaga gaccggactc tttataatct ttcaagggagc      240
      aaagaatgct tattttgact cgtcagagtg gcgaaacccct catcatcgcg gatgaggttaa      300
      cggttaccgt attaggagtg aaaggcaacc aggtgcgctat tgggtgtaat gcacctaaag      360
      aggtttctgt ccaccgtgaa gagatctatc agcgtattca ggccgaaaaa tctcaaccaa      420
      cgtcatattg attgacaatg cgtctcgtgt tcgggggacg caattgttat ttccggtttt      480
      tccccacac atttctcgat      500
```

Genbank

Header

```
LOCUS      ERWRSMA                      500 bp    DNA       linear   BCT 19-AUG-1995
DEFINITION Erwinia carotovora repressor (rsmA) gene, complete cds.
ACCESSION  L40173
VERSION    L40173.1   GI:927031
KEYWORDS   repressor; rsmA gene.
SOURCE     Pectobacterium carotovorum
            ORGANISM   Pectobacterium carotovorum
            Bacteria;  Proteobacteria; Gammaproteobacteria; Enterobacteriaceae;
            Pectobacterium.
REFERENCE  1 (bases 1 to 500)
AUTHORS    Cui,Y., Chatterjee,A., Liu,Y., Dumenyo,C.K. and Chatterjee,A.K.
TITLE      Identification of a global repressor gene, rsmA, of Erwinia
            carotovora subsp. carotovora that controls extracellular enzymes,
            N-(3-oxohexanoyl)-L-homoserine lactone, and pathogenicity in
            soft-rotting Erwinia spp
JOURNAL    J. Bacteriol. 177(17) (1995) In press
COMMENT     Original source text: Erwinia carotovora (strain 71, sub_species
            carotovora) DNA.
```

Annotation

FEATURES	Location/Qualifiers
source	1..500 /organism="Pectobacterium carotovorum" /strain="71" /sub_species="carotovora" /db_xref="taxon:554"
gene	107..431 /gene="rsmA"
-10_signal	107..112 /gene="rsmA"
RBS	235..239 /gene="rsmA"
CDS	246..431 /gene="rsmA" /function="global repressor" /note="putative" /codon_start=1 /transl_table=11 /protein_id="AAA74502.1" /db_xref="GI:927032" /translation="MLILTRRVGETLIIGDEVTVLVGVKNQVRIGVNAPKEVSVHR EETIYRIQAEKSQPTSY"

Sequence

BASE COUNT	140 a	101 c	120 g	139 t
ORIGIN	1	61	121	181
	ggatccggcga agcaggatag aaagtgtgtt accttcagat attctgaagc ttatcatgct	gttggtgataa caaaagcaca agctactgat atcgactaaa ctaacaagta	gtgacaaacc ggagttgtgat ggtgtgtgta taccatcgctc taggtttacg ttttcacagc	aatgatgga taatggcggg gagacagaga gaccggactc tttataatct ttcaagggagc
	aaagaatgct tattttgact cgtcagagtg gcgaaacccct catcatcgcg gatgaggttaa	cggttaccgt attaggagtg aaaggcaacc aggtgcgctat tgggtgtaat gcacctaaag	aggtttctgt ccaccgtgaa gagatctatc agcgtattca ggccgaaaaa tctcaaccaa	cgtcatattg attgacaatg cgtctcgtgt tcgggggacg caattgttat ttccggtttt
	tccccacac atttctcgat			

Artemis panels & navigation

Drop Down Menus

Entry Button Line

Main Sequence View Panel

Magnified Sequence View Panel

Feature Menu

CDS	190	255	Orthologue of E. coli thrL (LPT_ECOLI); Fasta hit to LPT_ECOLI
CDS	337	2799	Orthologue of E. coli thrA (AK1H_ECOLI); Fasta hit to AK1H_ECOLI
misc_feature	343	369	P800324 Aspartokinase signature
misc_feature	2314	2382	P801042 Homoserine dehydrogenase signature
CDS	2801	3730	Orthologue of E. coli thrB (KHSE_ECOLI); Fasta hit to KHSE_ECOLI
misc_feature	3068	3103	P800627 GHMP kinases putative ATP-binding domain
CDS	3734	5020	Orthologue of E. coli thrC (THRC_ECOLI); Fasta hit to THRC_ECOLI
misc_feature	4022	4066	P800165 Serine/threonine dehydratases pyridoxal-phosphate at
CDS	5114	5887	Orthologue of E. coli yaaA (YAAA_ECOLI); Fasta hit to YAAA_ECOLI
CDS	5966	7396	c Similar to Bacillus subtilis amino acid carrier protein alst
misc_feature	7091	7138	P800873 Sodium:alanine symporter family signature
CDS	7665	8618	Fasta hit to TALA_ECOLI (316 aa), 65% identity in 311 aa over
misc_feature	7755	7781	P801054 Transaldolase signature

Sliders

A word on genome browsers

	Artemis	GIVE	IGB	IGV	Jbrowse	Tablet	UCSC
LOCAL INSTALLATION							
Native app	●		●	●	●	●	
Web app		●●		●●	●●		●
Creation date	1999	2017	2001	2008	2009	2007	2014
Development status (2019)	stalled	early	stalled	mature/ early for web	active	mature	active
Software License	GPL3	Apache 2.0	Common Public License 1.0	MIT License	GNU LGPL v2.1	BSD-2 Clause	Copyright © 2001 UC Regents
PUBLIC WEB INSTANCE							
Creation date		2017		2018			2000
Development status (2019)		early		early			active

Features of the Artemis module

Exercise 1: Familiarize with Artemis

- Loading sequences and annotation files
- Changing the view
- Searching and getting around

Exercise 2: Find insight within a genome

- Graphs and plots

Exercise 3: Analyze a genome

- Basic analysis
- Generating features

Exercise 4: Feature editing

- Adding and Modifying annotations
- Finding evidence : Database searches



**HANDS-ON
LEARNING**