# Gene Sequence Analysis

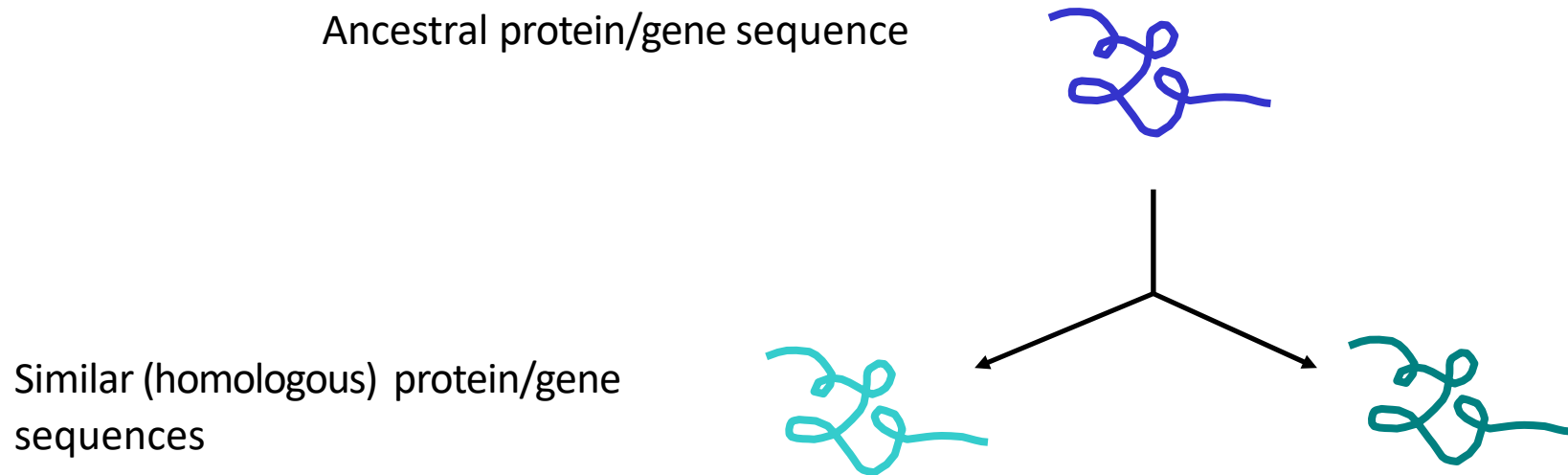# Lecture 2: **Pairwise Alignment**

30/05/2024

Phuc-Loi Luu, PhD
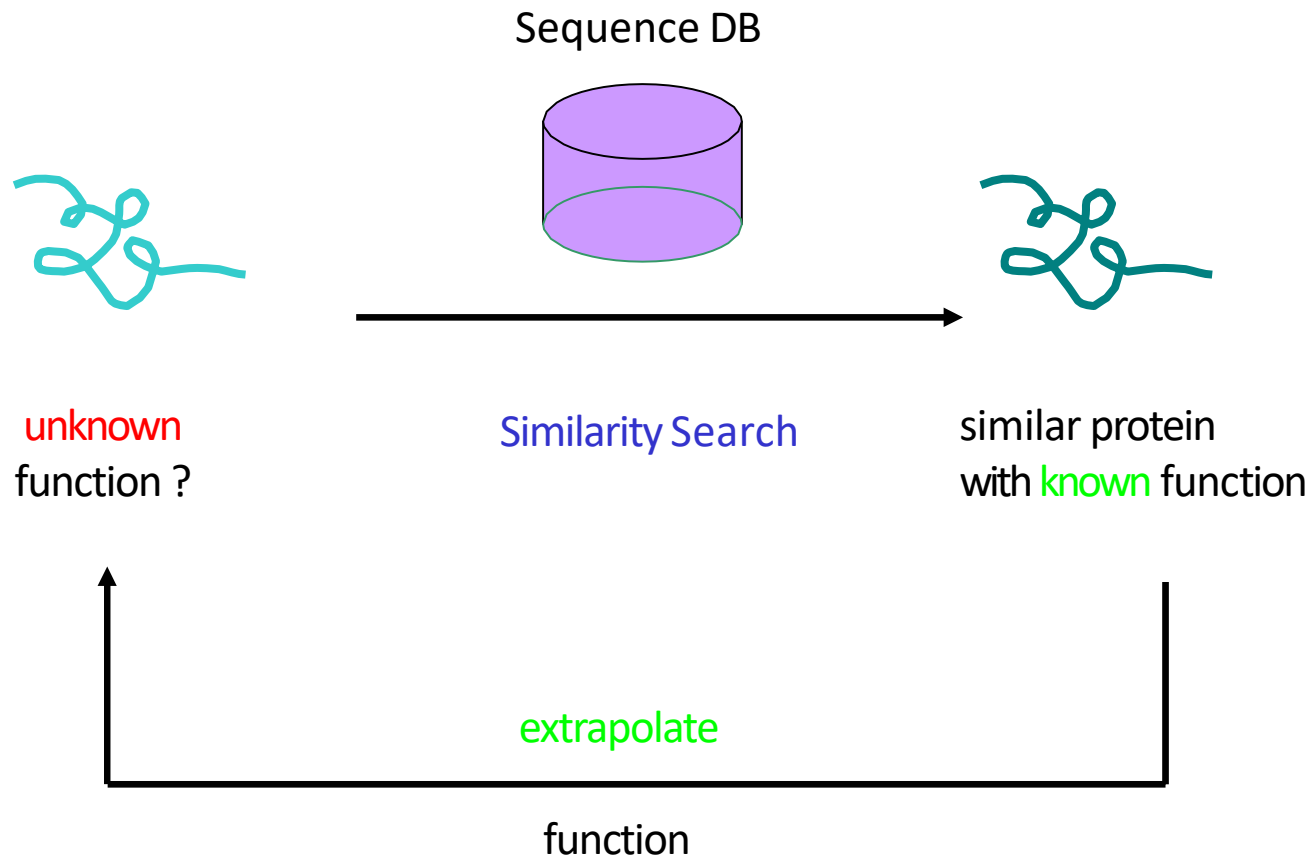
Adapted from Dr. Morgan Langille's Lecture

# Importance of Similarity

Ancestral protein/gene sequence

Similar (homologous) protein/gene
sequences

Similar sequences: probably have the same ancestor, share the same structure, and have a similar biological function

# Importance of Similarity

Sequence DB



unknown
function ?

Similarity Search

similar protein
with known function

extrapolate

function

# Importance of Similarity

Rule-of-thumb:
If your sequences are more than **100 amino acids** long (or 100 nucleotides long) you can considered them as homologues if **25%** of the **aa** are identical (**70%** of **nucleotide** for DNA). Below this value you enter the twilight zone.

Twilight zone = protein sequence similarity between ~0-20% identity: is not statistically significant, i.e. could have arisen by chance.

Beware:
- E-value (*Expectation value*)
- length of the segments similar between the two sequences
- The number of insertions/deletions

BLAST, FASTA, SSEARCH, and other commonly used similarity searching programs produce accurate statistical estimates that can be used to reliably infer homology. Searches with protein sequences (BLASTP, FASTP, SSEARCH,) or translated DNA sequences (BLASTX, FASTX) are preferred because they are 5- to 10-fold more sensitive than DNA:DNA sequence comparison. The 30% identity rule-of-thumb is too conservative; statistically significant [$E() < 10^{-6} - 10^{-3}$] protein homologs can share less than 20% identity. E()-values and bit scores (bits >50) are far more sensitive and reliable than percent identity for inferring homology.

## Descriptions

**Sequences producing significant alignments:**

Select: All None  Selected:0

Alignments  Download ˅  GenBank  Graphics  Distance tree of results

| Description | Max score | Total score | Query cover | E value | Ident | Accession |
|---|---|---|---|---|---|---|
| Homo sapiens laminin subunit alpha 4 (LAMA4), transcript variant 2, mRNA | 11596 | 11596 | 100% | 0.0 | 99% | NM_002290.4 |
| Homo sapiens laminin, alpha 4, mRNA (cDNA clone MGC:74960 IMAGE:6164247), complete cds | 11533 | 11533 | 100% | 0.0 | 99% | BC066552.1 |
| Homo sapiens laminin subunit alpha 4 (LAMA4), transcript variant 3, mRNA | 11516 | 11516 | 100% | 0.0 | 99% | NM_001105207.2 |
| Homo sapiens laminin subunit alpha 4 (LAMA4), transcript variant 1, mRNA | 11498 | 11498 | 100% | 0.0 | 99% | NM_001105206.2 |
| laminin alpha 4 chain [human, fetal lung, mRNA, 6204 nt] | 11407 | 11407 | 98% | 0.0 | 99% | S78569.1 |
| PREDICTED: Gorilla gorilla gorilla laminin, alpha 4, transcript variant 1 (LAMA4), mRNA | 11383 | 11383 | 100% | 0.0 | 99% | XM_004044552.1 |
| PREDICTED: Homo sapiens laminin, alpha 4 (LAMA4), transcript variant X1, mRNA | 11367 | 11367 | 99% | 0.0 | 99% | XM_005266983.3 |
| PREDICTED: Gorilla gorilla gorilla laminin, alpha 4, transcript variant 2 (LAMA4), mRNA | 11311 | 11311 | 100% | 0.0 | 99% | XM_004044553.1 |
| PREDICTED: Gorilla gorilla gorilla laminin, alpha 4, transcript variant 3 (LAMA4), mRNA | 11287 | 11287 | 100% | 0.0 | 99% | XM_004044554.1 |
| PREDICTED: Homo sapiens laminin, alpha 4 (LAMA4), transcript variant X2, mRNA | 11236 | 11236 | 97% | 0.0 | 99% | XM_005266984.3 |
| H.sapiens mRNA for laminin alpha 4 protein | 11158 | 11158 | 95% | 0.0 | 100% | X91171.1 |
| Homo sapiens mRNA for LAMA4 variant protein, clone: hh01833 | 10966 | 10966 | 95% | 0.0 | 99% | AB210027.1 |
| PREDICTED: Pan paniscus laminin, alpha 4 (LAMA4), transcript variant X2, mRNA | 10924 | 10924 | 96% | 0.0 | 99% | XM_003822216.3 |
| PREDICTED: Pan troglodytes laminin, alpha 4 (LAMA4), transcript variant X2, mRNA | 10890 | 10890 | 96% | 0.0 | 99% | XM_518696.5 |
| PREDICTED: Pan paniscus laminin, alpha 4 (LAMA4), transcript variant X1, mRNA | 10852 | 10852 | 96% | 0.0 | 99% | XM_008975760.1 |
| PREDICTED: Pan paniscus laminin, alpha 4 (LAMA4), transcript variant X3, mRNA | 10831 | 10831 | 96% | 0.0 | 99% | XM_008975761.1 |
| PREDICTED: Pan troglodytes laminin, alpha 4 (LAMA4), transcript variant X1, mRNA | 10816 | 10816 | 96% | 0.0 | 99% | XM_009451861.1 |

---

Download ˅  GenBank  Graphics

**Homo sapiens laminin subunit alpha 4 (LAMA4), transcript variant 2, mRNA**
Sequence ID: ref|NM_002290.4|  Length: 7355  Number of Matches: 1

Range 1: 116 to 6412 GenBank  Graphics          ▼ Next Match  ▲ Previous Match

| Score | Expect | Identities | Gaps | Strand |
|---|---|---|---|---|
| 11596 bits(6279) | 0.0 | 6291/6297(99%) | 0/6297(0%) | Plus/Plus |

```
Query  1    AGAAGGTAAAAAGGGAGTGGTGAGAATGAATGTGAGAAGGAAGCCAGGACAGCGCAGTCC  60
            ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct  116  AGAAGGTAAAAAGGGAGTGGTGAGAATGAATGTGAGAAGGAAGCCAGGACAGCGCAGTCC  175

Query  61   CCAGTCCCGAACGGCCAGGGAGAGGAGGTGGCCTAGCGCTGGCGGGGCTCACCCCAATCC  120
            ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct  176  CCAGTCCCGAACGGCCAGGGAGAGGAGGTGGCCTAGCGCTGGCGGGGCTCACCCCAATCC  235

Query  121  GTCTGCCTTTTGATGCCGTACTCTGCTGGTTGCGCAGCCACCTCGGGATACTGCACACGG  180
            ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct  236  GTCTGCCTTTTGATGCCGTACTCTGCTGGTTGCGCAGCCACCTCGGGATACTGCACACGG  295

Query  181  AGAGGAGGGAAAATAAGCGAGGCACCGCCGCACCACGCGGGAGACCTACGGAGACCCACA  240
            ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct  296  AGAGGAGGGAAAATAAGCGAGGCACCGCCGCACCACGCGGGAGACCTACGGAGACCCACA  355

Query  241  GCGCCCGAGCCCTGGAAGAGCACTACTGGATGTCAGCGGAGAAATGGCTTTGAGCTCAGC  300
            ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct  356  GCGCCCGAGCCCTGGAAGAGCACTACTGGATGTCAGCGGAGAAATGGCTTTGAGCTCAGC  415
```
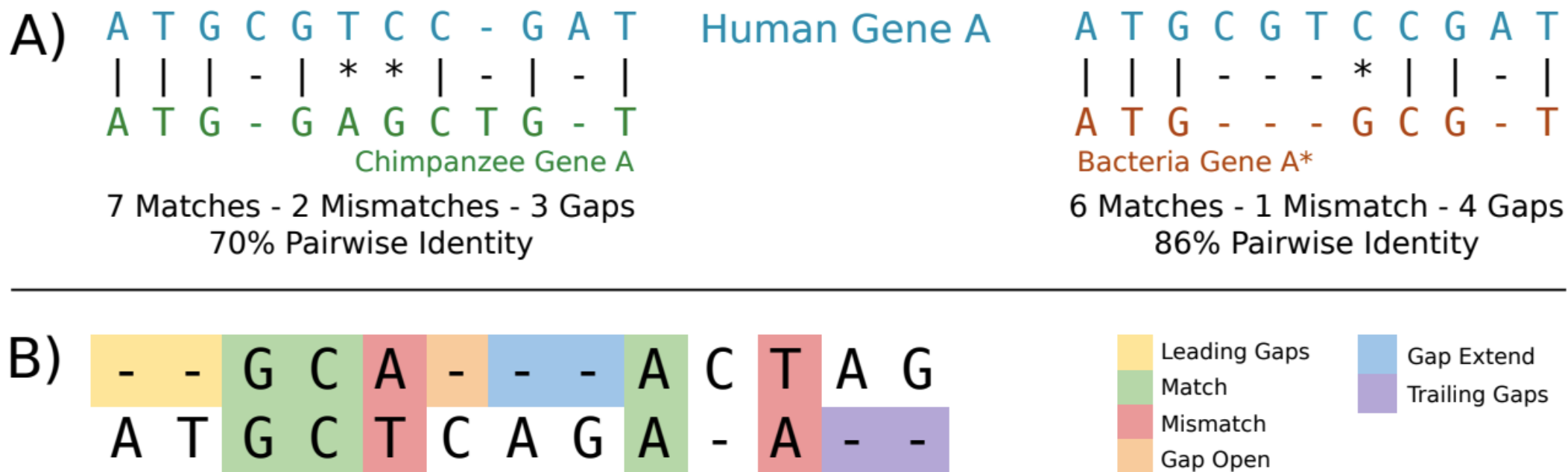
FIGURE 1.2: **Similarity versus homology and alignment scoring:** *A)* Two sequence alignments to Human Gene A highlighting how true homologs (Chimpanzee Gene A) can have similar similarity to non-homologous sequences (Bacteria Gene A*). *B)* A sequence alignment with different reward or penalties highlighted to illustrate how homology search algorithms score alignments based on sequence similarity.

Doolittle's twilight zone is between 15-25% and is a 'rule of thumb' for sequences of 100 amino acid length or greater (29). The 100 amino acids caveat is due to the greater probability of short sequences having the same amino acid sequence by chance and therefore the level of confidence in short alignments is not as easy to give general advice for (19, 29). No quantitative experiment was done to determine this

As a part of a study of structural ncRNAs that were used to benchmark multiple sequence alignment programs the twilight zone for nucleotide alignment accuracy was calculated (11). It was between 50-60% which is much higher then the 10-20% calculated by Thompson *et al.* 1999 for protein alignments (32). This zone represents

billion years ago (e.g. humans to bacteria). Moreover, DNA:DNA alignment statistics are less accurate than protein:protein statistics; while protein:protein alignments with expectation values $< 0.001$ can reliably be used to infer homology, DNA:DNA expecation values $< 10^{-6}$ often occur by chance, and $10^{-10}$ is a more widely accepted threshold for homology based on DNA:DNA searches. The most effective way to improve search sensitivity with DNA sequences is to use translated-DNA:protein alignments, such as those produced by BLASTX and FASTX, rather than DNA:DNA alignments.

# Outline

- PSSMs/PSI-BLAST
- HMMs/HMMer RNA Alignments
- Genome Alignments
- Assemblers  Mappers

# Different tools for homology searching

- Searching for protein families

- Aligning genomes

- Looking for RNA genes

- Combining overlapping sequences (assemblers)

- Finding the position of a sequence in a genome

# One tool does not do it all

- Blast may give you an answer

    - BUT you could find the answer much quicker or with more precision by using the right tool!

# Typical BLAST output



**Histone H1** (residues 120-180)

```
HUMAN  KKASKPKKAASKAPTKKPKATPVKKAKKKLAATPKKAKKPKTVKAKPVKASKPKKAKPVK
CHIMP  KKASKPKKAASKAPTKKPKATPVKKAKKKLAATPKKAKKPKTVKAKPVKASKPKKAKPVK
MOUSE  KKAAKPKKAASKAPSKKPKATPVKKAKKKPAATPKKAKKPKVVKVKPVKASKPKKAKTVK
RAT    KKAAKPKKAASKAPSKKPKATPVKKAKKKPAATPKKAKKPKIVKVKPVKASKPKKAKPVK
COW    KKAAKPKKAASKAPSKKPKATPVKKAKKKPAATPKKTKKPKTVKAKPVKASKPKKTKPVK
       ***.*********.************** :.**** **.*********** * **
```

NON-CONSERVED
AMINO ACIDS

Conservative    Conservative    Non-conservative    Conservative    Non-conservative    Semi-conservative    Conservative    Non-conservative

https://en.wikipedia.org/wiki/Sequence_alignment
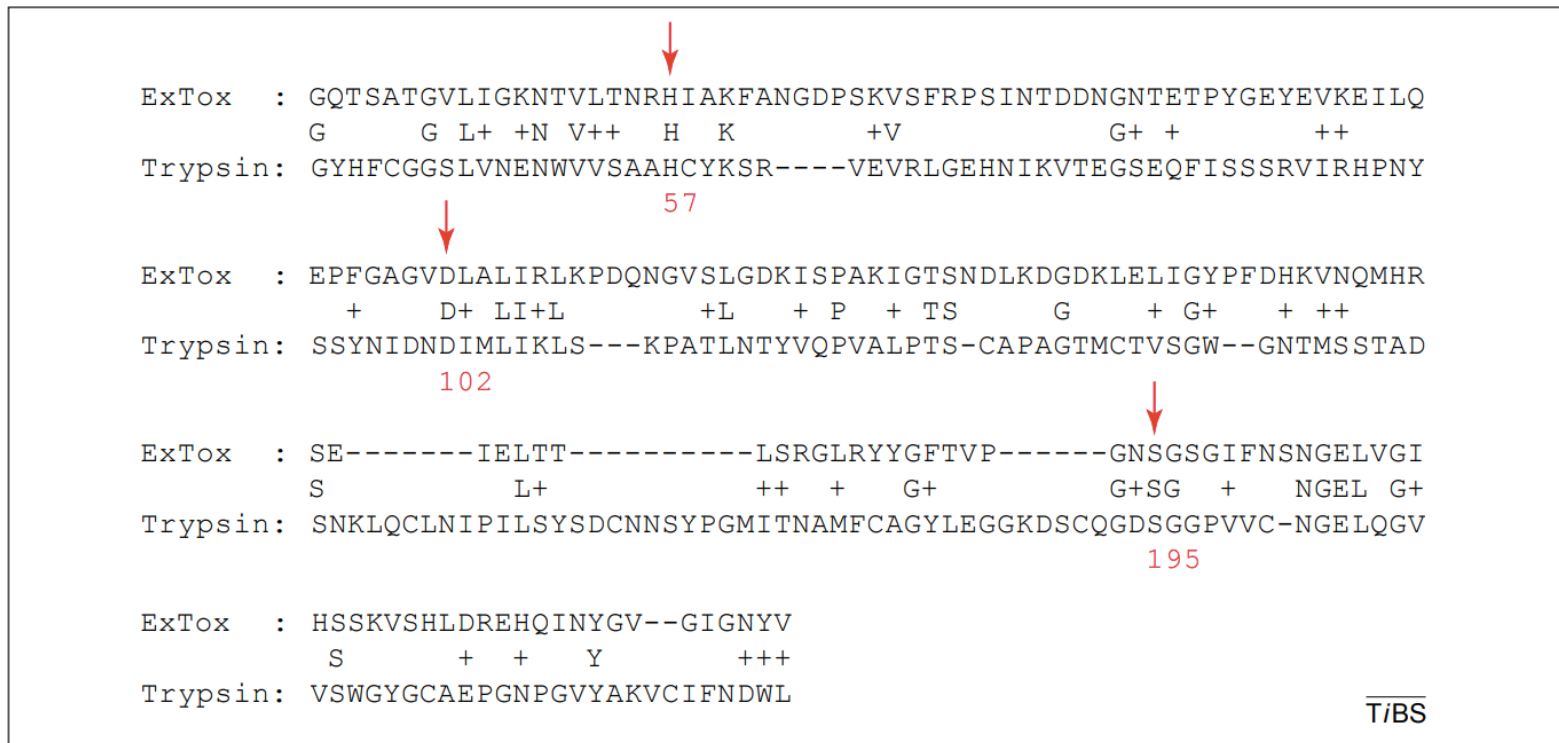
# Distance related proteins

# PSI-BLAST



**Fig. 1.** Third-iteration PSI–BLAST result from querying exfoliative toxin (ExTox) A from *Staphylococcus aureus* (BAA97652.1) against the non-redundant database. His, Asp and Ser residues are indicated with arrows and numbered as for trypsin [anionic; complexed with the inhibitor benzamidine (1bit)]. The alignment has 15% identity (32/206) and the $E$-value $= 6 \times 10^{-21}$. The threshold $E$-value for inclusion in the profile was 0.005 and the effective search space was 22 926 875 677.

10.1016/S0968-0004(01)02039-4

# PSI-BLAST



```
gi|2622094 (AE000872) conserved protein [Methanobacterium thermoautotrophicum]
             Length = 143

Score = 84.7 bits (206), Expect = 4e-16
Identities = 56/156 (35%), Positives = 81/156 (51%), Gaps = 16/156 (10%)

                  ──────────────────────────────►
Query: 4    MYKKILYPTDFSETAEIALKHVKAFKTLKAEEVILLHVIDEREIKKRDIFSLLLGVAGLN 63
            MY KIL PTD S+ A   A +H              E+I L V++        S L+G+
Sbjct: 1    MYSKILLPTDGSKQANKAAEHAIWIARESGAEIIALTVMET---------SSLVGLPA-- 49

Query: 64   KSVEEFENELKNKLTEEAKNKMENIKKELEDVGFKVKDIIVV--GIPHEEIVKIAEDEGV 121
               ++      L+  L EEA    +E +KK +E+ G   +K    G P E I++  E EGV
Sbjct: 50   ---DDLIIRLREMLEEEASRSLEAVKKLVEESGADIKLTVRTDEGSPAEAILRTVEKEGV 106

Query: 122  DIIIMGSHGKTNLKEILLGSVTENVIKKSNKPVLVV 157
            D+++MG+ GK L    LLGSV E V++ +    PVLVV
Sbjct: 107  DLVVMGTSGKHGLDRFLLGSVAEKVVRSAGCPVLVV 142
```
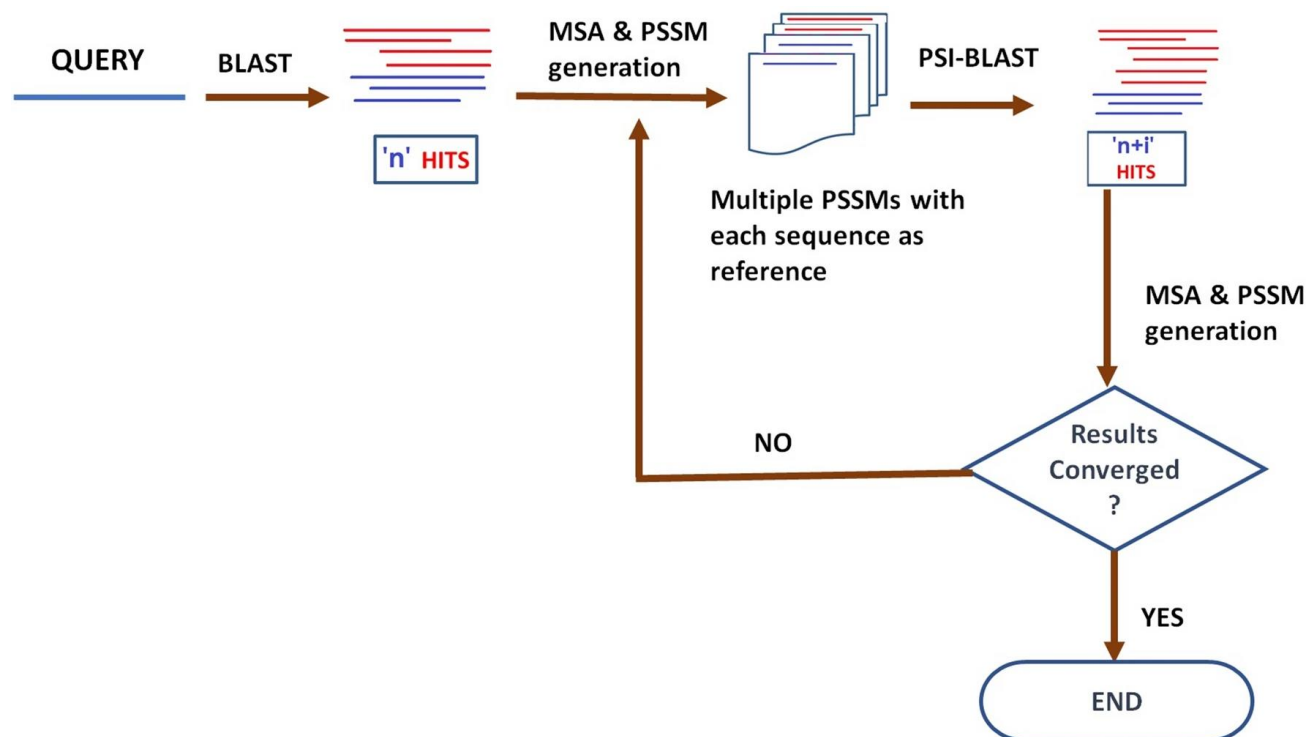
An example of high-scoring segment pairs (HSP) found by PSI-BLAST. The first peptide pairs as marked by the box are similar, and we assign the secondary structure element of each amino acid in MYKKILY to its counterpart in MYSKILL.

# PSI-BLAST

- Position Specific Iterated – BLAST A
- cycling/iterative method
  - Gives increased sensitivity for detecting **distantly related proteins**
  - Can give insight into functional relationships
  - Very refined statistical methods
- Fast – still based on BLAST methods
- Simple to use

# PSI-BLAST

- Essentially we are using intermediate sequences to infer similarity between two sequences that are too dissimilar to link directly.

# Profiles & PSSMs Need Multi-sequence Alignment

```
          1                                                    50
P43871-1  .......... ...IKKLDSN SIHAIISDIP YGIDYDDWDI LHSNTNSALG
S18997-1  LMSKIYQMDA VDWLKTLENC SVDLFITDPP YESL.EKYRQ IGTTTRLKES
P23192-1  EINKIHQMNC FDFLDQVENK SVQLAVIDPP YNL....... ..........
P29538-1  MDQRLICSNA IKALKNLEEN SIDLIITDPP YNLG.KDY.. ..........
P14751-1  TRHVYDVCDC LDTLAKLPDD SVQLIICDPP YNI....... ..........
P34721-1  KNFNIYQGNC IDFMSHFQDN SIDMIFADPP YFLS.NDG.L TFKNSIIQ..
P50178-1  ENAILVHADS FKLLEKIKPE SMDMIFADPP YFLS.NGG.M SNSGGQIV..
P20590-1  FLNTILKGDC IEKLKTIPNE SIDLIFADPP YFMQ.TEGKL LRTNGDEF..
S43876-1  GPETIIHGDC IEQMNALPEK SVDLIFADPP YNLQ.LGGDL LRPDNSKV..
P28638-1  EAKTIIHGDA LAELKKIPAE SVDLIFADPP YNIG.KNF.. ..........
P23941-1  DLGKLYNGDC LELFKQVPDE NVDTIFADPP FNLD.KEY.. ..........
P14230-1  RSCKIIVGDA REAVQGLDSE IFDCVVTSPP YWGL.RDY.. ..........
P14243-1  NGATLFEGDA LSVLRRLPSG SVRCIVTSPP YWGL.RDY.. ..........
Q04845-1  LNNMLLQGNC AETLKKLPDE SVNLVFTSPP YY........ ..........
S53866-1  WVNDIHEGDA EEVLAELPES SVHMVMTSPP YFGL.RDY.. ..........
P29568-1  .......... ...MNELKDK SINLVVTSPP YPMV.EIWDR LFSELNPKIE

Signature Sequence:                     DPP Y
```
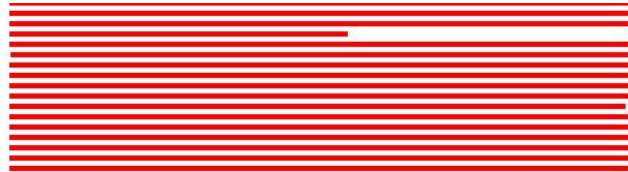
# How does PSI-BLAST work?

1) First, a standard blastp is performed

2) The highest scoring hits are used to generate a multiple alignment

3) A Position Specific Scoring Matrix (PSSM) is generated from the multiple alignment.
   - Highly conserved residues get high scores
   - Less conserved residues get lower scores
   - The PSSM describes the sequence similarity between your query and all significant blastp hits

4) Another similarity search is performed, this time using the new PSSM as the query sequence.
   - This PSSM (scoring matrix) is now customized to find sequences that are related to your original query
   - Steps 2-4 can be repeated until convergence
   - Convergence occurs when no new sequences appear after iteration

# PSI-BLAST Example



## Descriptions

Legend for links to other resources: U UniGene  E GEO  G Gene  S Structure  M Map Viewer  PubChem BioAssay

**NEW** - alignment score below the threshold on the previous iteration

● - alignment was checked on the previous iteration

Run PSI-Blast iteration 2 with max [500]  [Go]

### Sequences producing significant alignments with E-value BETTER than threshold

| | Accession | Description | Max score | Total score | Query coverage | △ E value | Max ident | Links |
|---|---|---|---|---|---|---|---|---|
| NEW | ☑NP_625085.1 | tetracycline resistance protein [Streptomyces coelicolor A3(2)] | 1172 | 1172 | 100% | 0.0 | 98% | G |
| NEW | ☑ZP_07295505.1 | tetracycline resistance protein TetP [Streptomyces hygroscopicus ATCC | 650 | 650 | 97% | 0.0 | 61% | |
| NEW | ☑YP_003342575.1 | unnamed protein product [Streptosporangium roseum DSM 43021] | 637 | 637 | 97% | 0.0 | 58% | G |
| NEW | ☑YP_004811081.1 | small GTP-binding protein [Streptomyces violaceusniger Tu 4113] | 634 | 634 | 97% | 0.0 | 58% | G |
| NEW | ☑ZP_06532924.1 | tetracycline resistance protein [Streptomyces lividans TK24] | 587 | 587 | 48% | 0.0 | 99% | |
| NEW | ☑YP_004920166.1 | unnamed protein product [Streptomyces cattleya NRRL 8057] | 586 | 586 | 97% | 0.0 | 58% | G |
| NEW | ☑YP_003118884.1 | small GTP-binding protein [Catenulispora acidiphila DSM 44928] | 560 | 560 | 97% | 0.0 | 54% | G |
| NEW | ☑YP_003486610.1 | unnamed protein product [Streptomyces scabiei 87.22] | 472 | 472 | 95% | 2e-156 | 50% | G |
| NEW | ☑ZP_06921723.1 | translation elongation factor G [Streptomyces sviceus ATCC 29083] | 468 | 468 | 95% | 3e-155 | 50% | |
| NEW | ☑YP_003383525.1 | small GTP-binding protein [Kribbella flavida DSM 17836] | 467 | 467 | 97% | 4e-155 | 48% | G |
| NEW | ☑ZP_04163174.1 | GTP-binding elongation factor protein, TetM/TetO [Bacillus mycoides Roc | 462 | 462 | 96% | 4e-153 | 38% | |
| NEW | ☑ZP_07276632.1 | translation elongation factor G [Streptomyces sp. AA4] | 462 | 462 | 97% | 5e-153 | 47% | |
| NEW | ☑ZP_04151746.1 | GTP-binding elongation factor protein, TetM/TetO [Bacillus pseudomycoi | 459 | 459 | 96% | 9e-152 | 38% | |
| NEW | ☑ZP_04157524.1 | GTP-binding elongation factor protein, TetM/TetO [Bacillus mycoides Roc | 457 | 457 | 96% | 5e-151 | 38% | |
| NEW | ☑ZP_04198058.1 | GTP-binding elongation factor protein, TetM/TetO [Bacillus cereus AH60: | 457 | 457 | 96% | 6e-151 | 38% | |
| NEW | ☑ZP_09405061.1 | putative tetracycline resistance protein [Streptomyces sp. W007] | 454 | 454 | 97% | 1e-149 | 50% | |
| NEW | ☑ZP_04097136.1 | GTP-binding elongation factor protein, TetM/TetO [Bacillus thuringiensis | 453 | 453 | 97% | 1e-149 | 38% | |
| NEW | ☑ZP_04228518.1 | GTP-binding elongation factor protein, TetM/TetO [Bacillus cereus Rock3 | 451 | 451 | 96% | 1e-148 | 38% | |
| NEW | ☑ZP_07308418.1 | translation elongation factor G [Streptomyces viridochromogenes DSM 4 | 449 | 449 | 97% | 6e-148 | 47% | |

# HMMs & HMMer

- The more powerful way to search for protein families than PSSMs



Hidden Markov Models

# Hidden Markov Models in Bioinformatics

- **Used extensively in gene prediction**

- **Used to create Sequence Profiles and to classify sequences into families**

- **Used in Multiple Sequence Alignment**

# HMMER 3

- Suite of sequence analysis programs based on HMMs

- Used to build the Pfam database

- Available for free download at
  - http://hmmer.org/

# HMMER 3

- HMMER 2 was used for many years
  - Biggest draw back was always speed

- HMMER 3 released in 2011
  - Very fast with comparable speeds to BLAST
  - 100X faster than v2

# HMMER Programs

- **hmmbuild** – build a HMM from multiple sequence alignment

- **hmmscan** – searches a query sequence(s)  against a database of HMMs (used by PFAM)

- **hmmsearch**– searches a query HMM against a database of sequences (e.g. like psi-blast)

- **phmmer** – search a protein sequence vs a sequence database (e.g. like blastp)

# HMMER Search & Software

- http://hmmer.janelia.org/search

- PFAM

  - http://pfam.sanger.ac.uk/

# RNA Alignments

- RNA alignments are "special"

- RNA genes often have secondary structures that allow improved searching

- Improved searching is needed since

  - Must search in DNA space (less complex then protein sequences)

  - Often shorter length then proteins

# Infernal (RNA Search)

- Infernal is like HMMER
  - Includes use of secondary structure information
  - Uses profile "stochastic context-free grammar"
    - SCFGs vs HMMs
  - "consensus RNA secondary structure profiles"
- Infernal is slow!
- Infernal can be used to search RFAM

# RNAseq alignment with Magic-BLAST

## Magic-BLAST, an accurate RNA-seq aligner for long and short reads

Grzegorz M. Boratyn, Jean Thierry-Mieg, Danielle Thierry-Mieg, Ben Busby & Thomas L. Madden ✉

## Abstract

### Background

Next-generation sequencing technologies can produce tens of millions of reads, often paired-end, from transcripts or genomes. But few programs can align RNA on the genome and accurately discover introns, especially with long reads. We introduce Magic-BLAST, a new aligner based on ideas from the Magic pipeline.

### Results

Magic-BLAST uses innovative techniques that include the optimization of a spliced alignment score and selective masking during seed selection. We evaluate the performance of Magic-BLAST to accurately map short or long sequences and its ability to discover introns on real RNA-seq data sets from PacBio, Roche and Illumina runs, and on six benchmarks, and compare it to other popular aligners. Additionally, we look at alignments of human idealized RefSeq mRNA sequences perfectly matching the genome.

https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-019-2996-x

# RNAseq alignment with Magic-BLAST

# Genome Alignment
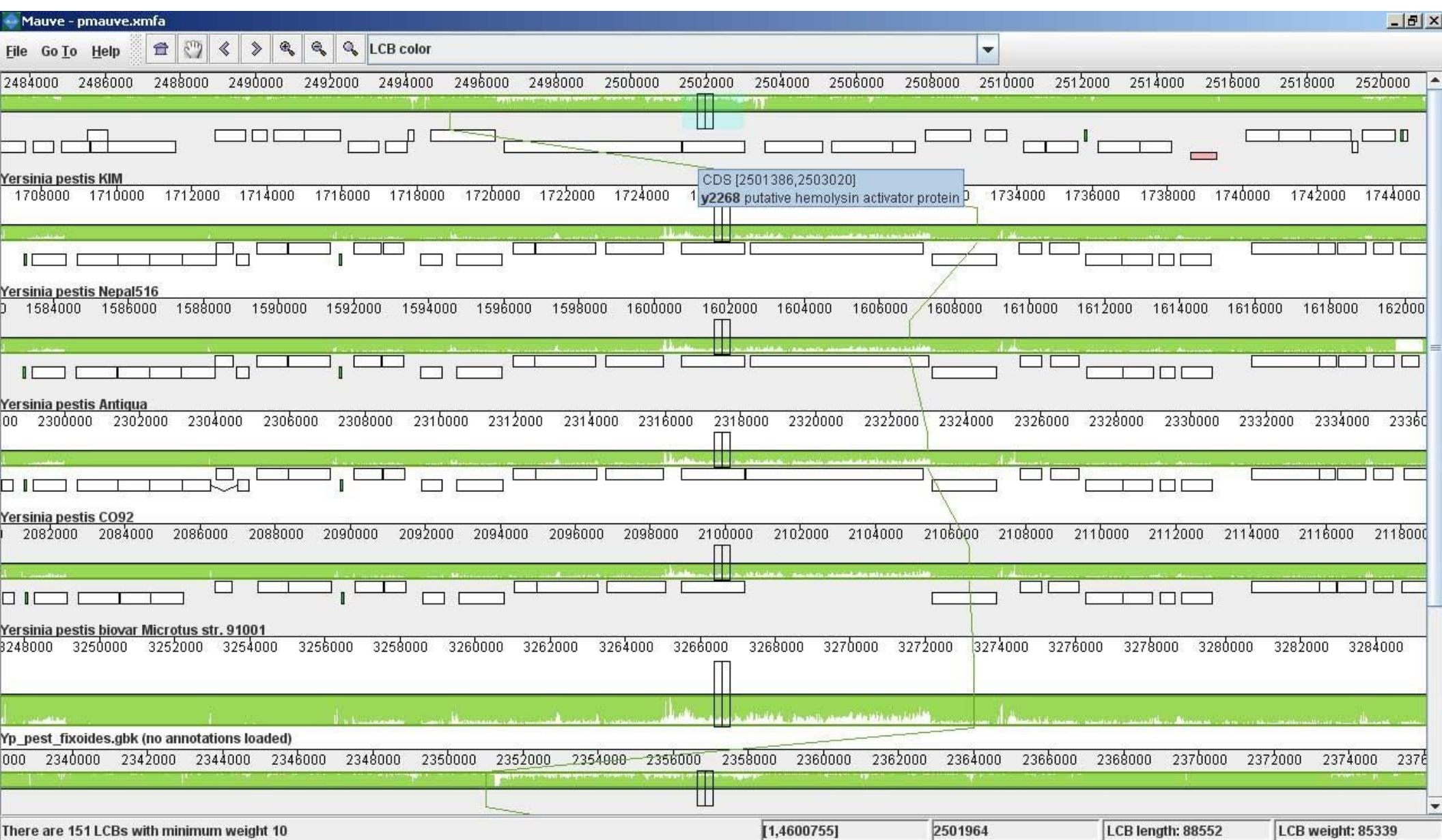
# Genome Alignment

- Genome alignment useful for

  - Visualizing genome

    - Rearrangements

    - Insertions/deletions

    - Inversions

  - Annotating genomes

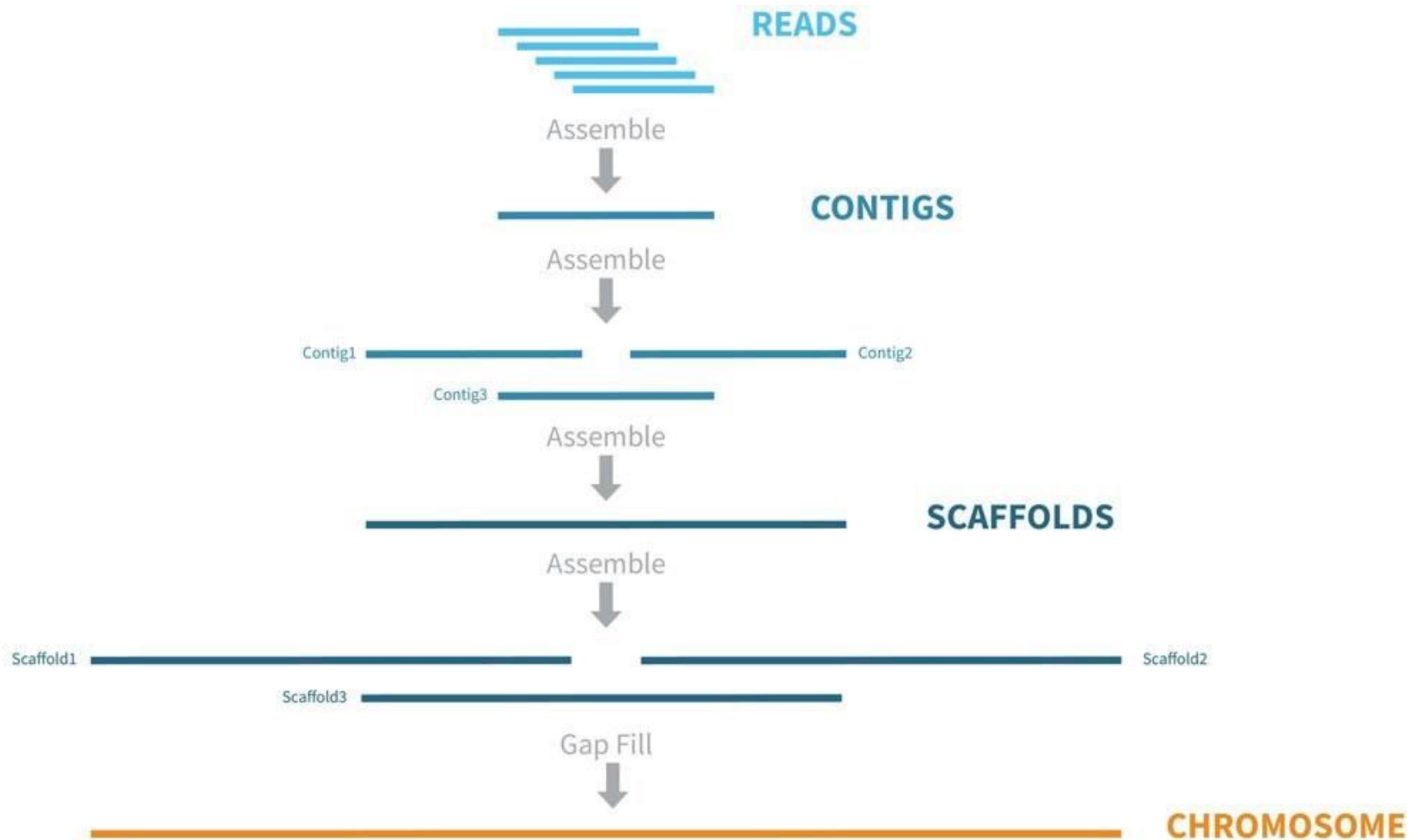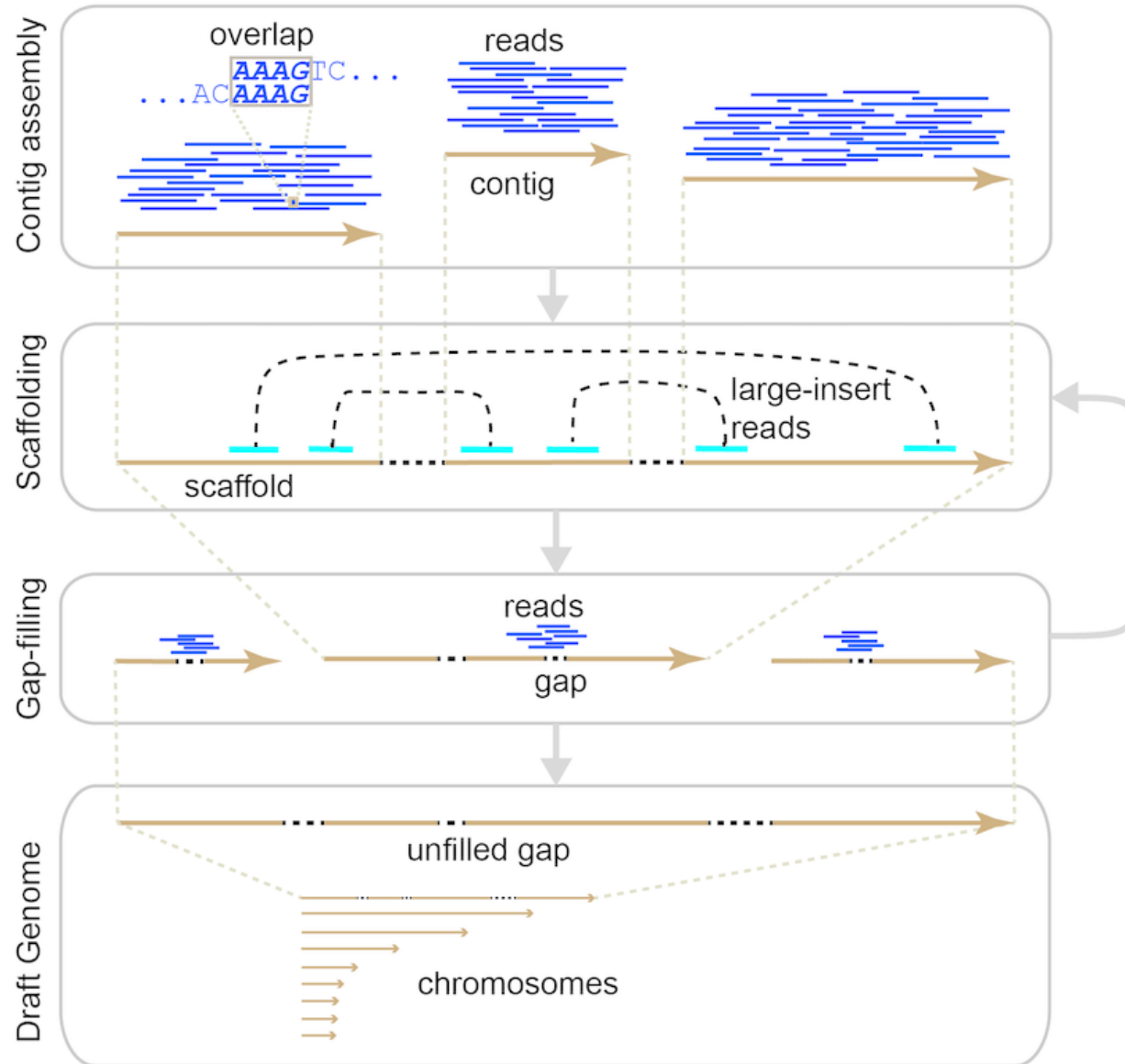    - Comparing gene annotations across species

# Mauve

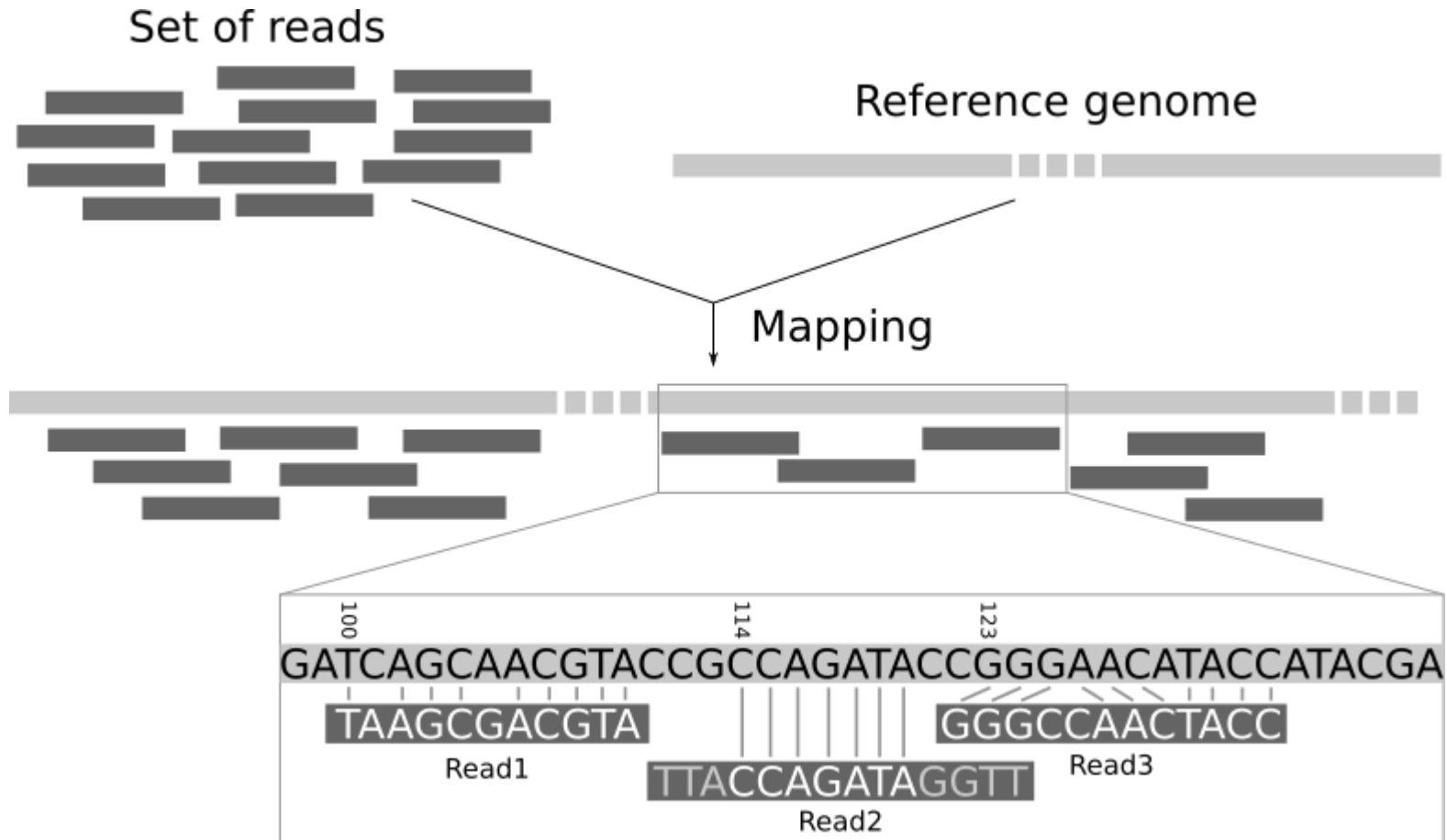# Mauve (zoomed in)

# Assemblers
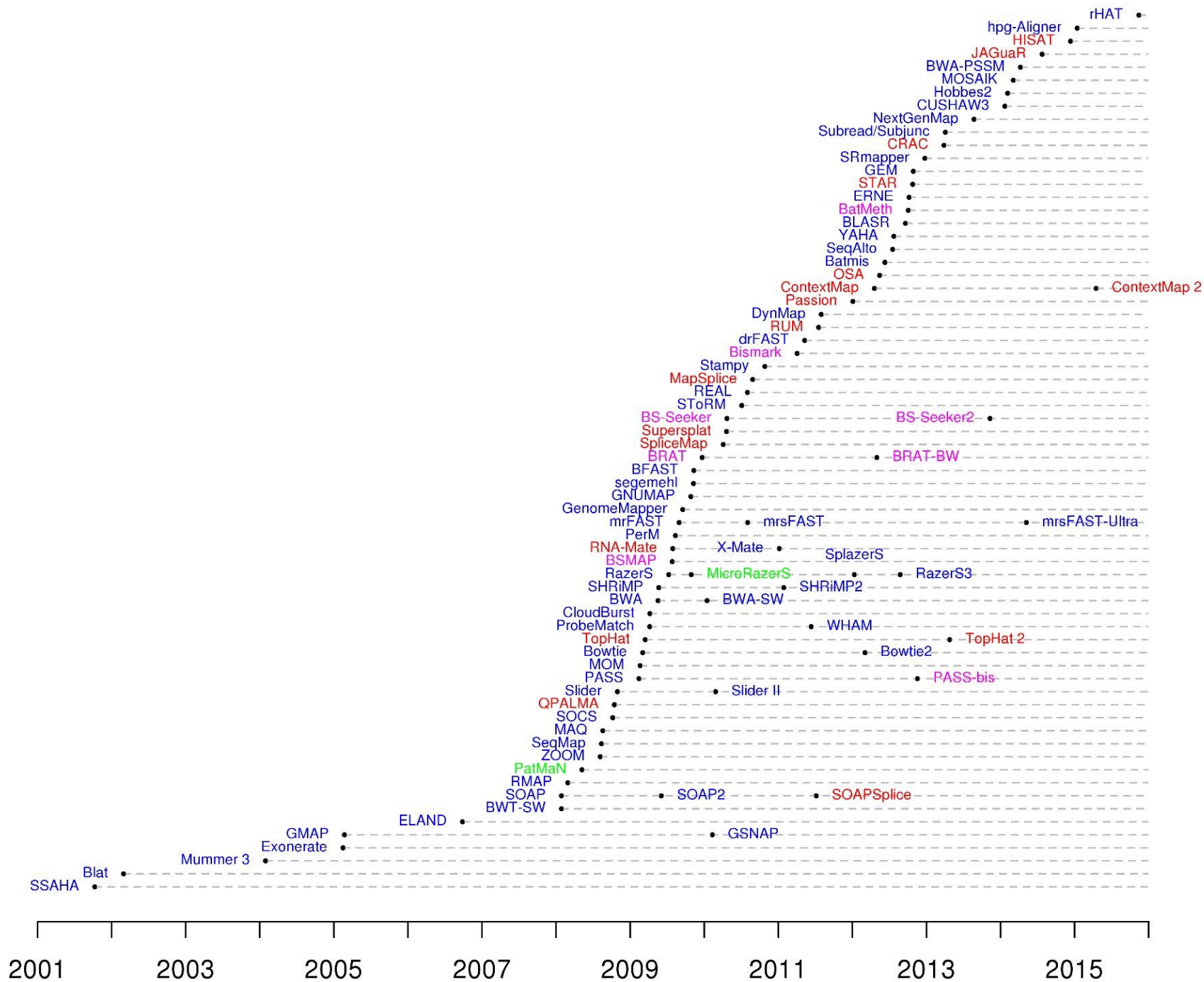
# Assemblers

# Assemblers

- Assemblers job is to make longer sequences from shorter ones.

- Nothing like homology searching

- Must efficiently compare and join billions of sequences

- Soap-Denovo: http://soap.genomics.org.cn/soapdenovo.html

- Amos: http://sourceforge.net/apps/mediawiki/amos/index.php?title=AMOS

- Many, many, more

# Mappers (Aligners)
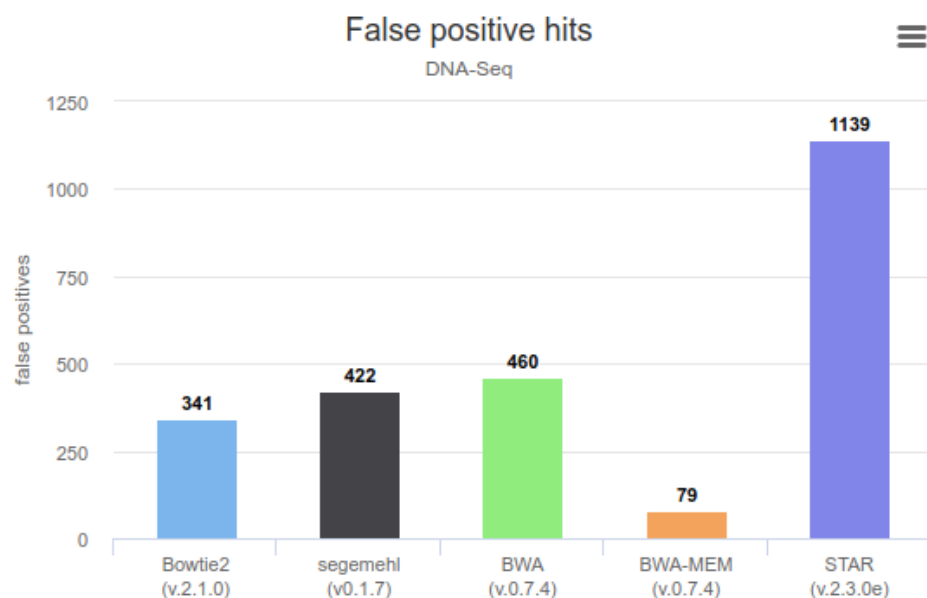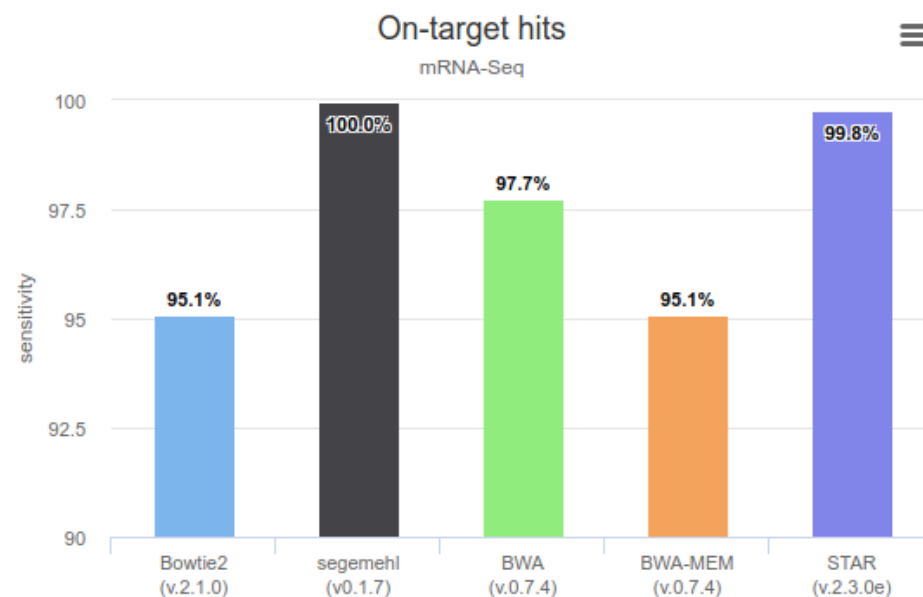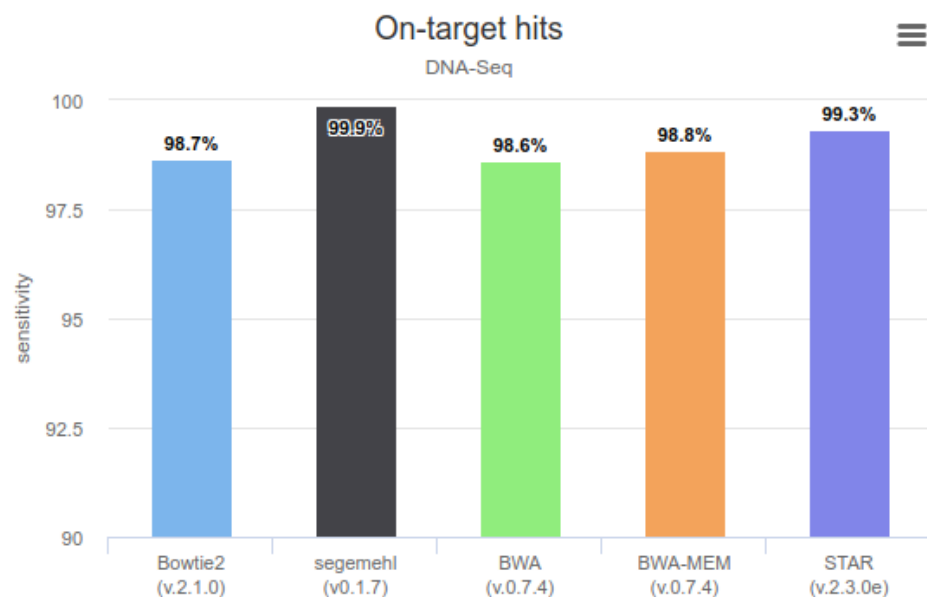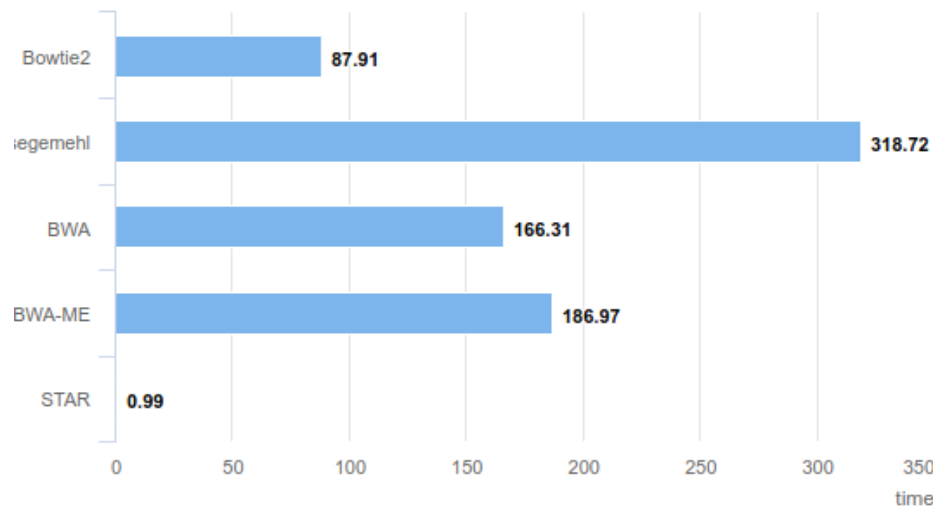
# Mappers

- These map a read to a reference genome

- Useful for assembly when a reference genome is already known

    - (think assembly of personal human genomes)

- Identifying SNPs within the same species Very Fast!

- BWA: https://github.com/lh3/bwa

- Bowtie: http://bowtie-bio.sourceforge.net/index.shtml

- Stampy: http://www.well.ox.ac.uk/project-stampy

- Many Others

-

Years

## On-target hits
### DNA-Seq

Bowtie2 (v.2.1.0): 98.7%
segemehl (v0.1.7): 99.9%
BWA (v.0.7.4): 98.6%
BWA-MEM (v.0.7.4): 98.8%
STAR (v.2.3.0e): 99.3%

## On-target hits
### mRNA-Seq

Bowtie2 (v.2.1.0): 95.1%
segemehl (v0.1.7): 100.0%
BWA (v.0.7.4): 97.7%
BWA-MEM (v.0.7.4): 95.1%
STAR (v.2.3.0e): 99.8%

## False positive hits
### DNA-Seq

Bowtie2 (v.2.1.0): 341
segemehl (v0.1.7): 422
BWA (v.0.7.4): 460
BWA-MEM (v.0.7.4): 79
STAR (v.2.3.0e): 1139

## False positive hits
### mRNA-Seq

Bowtie2 (v.2.1.0): 122
segemehl (v0.1.7): 7
BWA (v.0.7.4): 46
BWA-MEM (v.0.7.4): 266
STAR (v.2.3.0e): 167

Otto C, Stadler PF, Hoffmann S: 'Lacking alignments? The next generation sequencing mapper segemehl revisited', Bioinformatics. 2014 Jul 1;30(13):1837-43 (2014)

## User time [s]
### DNA-Seq

| Mapper | Time |
|--------|------|
| Bowtie2 | 87.91 |
| segemehl | 318.72 |
| BWA | 166.31 |
| BWA-ME | 186.97 |
| STAR | 0.99 |

## User time [s] (mRNA-Seq)
### mRNA-Seq

| Mapper | Time |
|--------|------|
| Bowtie2 | 32.63 |
| segemehl | 150.8 |
| BWA | 32.35 |
| BWA-ME | 12.58 |
| STAR | 3.31 |

## Memory consumption [GB]
### DNA-Seq

| Mapper | Memory |
|--------|--------|
| Bowtie2 | 3.77 |
| segemehl | 70.05 |
| BWA | 3.85 |
| BWA-ME | 5.75 |
| STAR | 28.12 |

## Memory consumption [GB]
### mRNA-Seq

| Mapper | Memory |
|--------|--------|
| Bowtie2 | 3.76 |
| segemehl | 70.05 |
| BWA | 3.73 |
| BWA-ME | 5.66 |
| STAR | 28.12 |

Otto C, Stadler PF, Hoffmann S: 'Lacking alignments? The next generation sequencing mapper segemehl revisited', Bioinformatics. 2014 Jul 1;30(13):1837-43 (2014)

# Mappers



IGV snapshot with reads in bam file