

Session 8: Intro to Sequence Similarity - BLAST, MASH, ANI, and Other Applications in Public Health Bioinformatics

Joel R. Sevinsky, Ph.D.
AMD Academy June 2020

Overview

- BLAST - Overview and brief description of algorithm
- BLAST exercise using NCBI interface
- Microbial identification - BLAST in ANI
- Microbial identification - beyond BLAST → MASH
- Other areas of sequence similarity searching

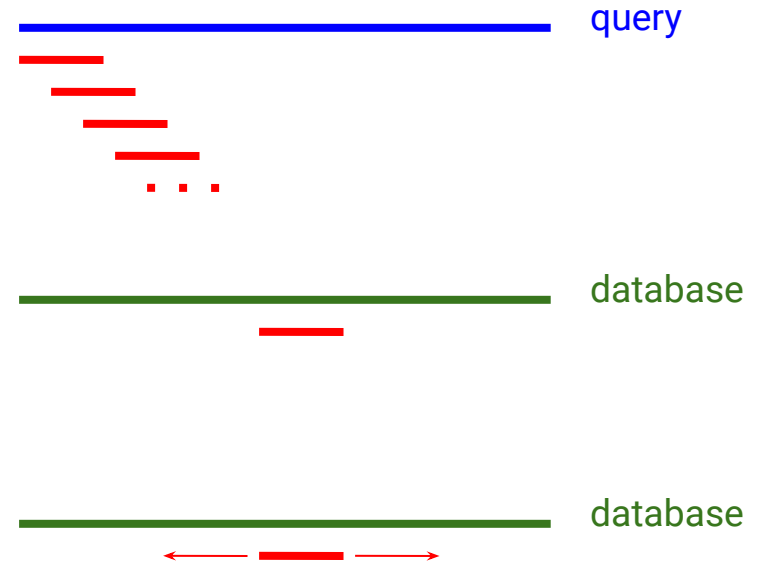
BLAST: Basic Local Alignment Search Tool

- Suite of tools for comparing a query sequence to a database of sequences
- You can compare DNA:DNA, DNA:protein, protein:DNA, and protein:protein
- “Seed and extend” heuristic (indexing based local alignment)
 - Seed
 - make starting words (k-mers) in assembly from your query sequence
 - compare query words to an indexed word list from your database
 - identify exact matches
 - Extend
 - once exact matches are found, try and extend the alignments with a scoring algorithm
 - different algorithms available for DNA vs protein
 - keep score as you extend beyond exact match, and maximize the score

BLAST: Basic Local Alignment Search Tool

Seed and extend process

- Create query dictionary
- Align/search database dictionary
- Extend query along target



BLAST: Basic Local Alignment Search Tool

```

g a g t a t t g a g g t c a g a t g g t c a c t g a - query (seed=10mer)
a t t c g c t g a g g t c a g a t g a c t t g a c t - database

```



```

+1+1+1+1+1+1+1+1+1+1 = 10
+1+1+1+1+1+1+1+1+1+1+1 = 12 --> Highest scoring
-2+1+1+1+1+1+1+1+1+1+1+1-2 = 8

```

- Initial, or seed match, is extended on each side until the aggregate alignment score falls below a predetermined threshold.

BLAST: Basic Local Alignment Search Tool

<https://blast.ncbi.nlm.nih.gov/Blast.cgi>

What Tredegar steps use sequence similarity?

```
(base) amd-academy@blast:~/data$ staphb-wf tredegar -o blast illumina/
/home/amd-academy/anaconda3/lib/python3.7/site-packages/staphb_toolkit/lib
3.7/site-packages/staphb_toolkit/workflows/tredegar//tredegar.nf -profile
blast/logs/20_06_14_17_40_41_Tredegar_trace.txt -with-report blast/logs/20
st/logs/work
Starting the Tredegar pipeline:
N E X T F L O W ~ version 20.01.0
Launching `/home/amd-academy/anaconda3/lib/python3.7/site-packages/staphb_
- revision: 679f403fdf
executor > local (8)
[ea/c591c8] process > preProcess      [100%] 6 of 6 ✓
[dc/decla6] process > mash_dist      [ 17%] 1 of 6
[-          ] process > mash_species -
[88/cacb86] process > trim           [  0%] 0 of 6
[-          ] process > cleanreads   -
[-          ] process > shovill      -
[-          ] process > quast        -
[-          ] process > cg_pipeline  -
[-          ] process > emmtype_finder -
[-          ] process > seqsero     -
[-          ] process > serotypefinder -
[-          ] process > results      -
```

1. mash
2. shovill
3. cg_pipeline
4. emmtype_finder
5. seqsero
6. serotypefinder

What Tredegar steps use sequence similarity?

- Assembly
 - genome assembly, reference based mapping
- Identity
 - species identification
- Annotation
 - gene prediction, serotype, AMR
- Genetic Characterization
 - phylogenetic trees, SNPs

What Tredegar steps use sequence similarity?

- Assembly
 - genome assembly, reference based mapping
- Identity
 - species identification
- Annotation
 - gene prediction, serotype, AMR
- Genetic Characterization
 - phylogenetic trees, SNPs

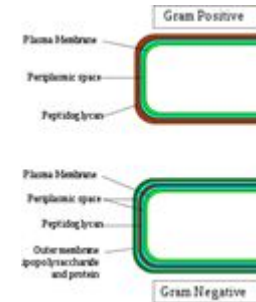
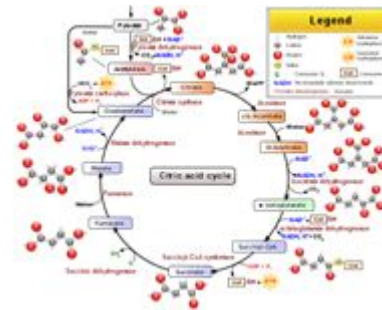
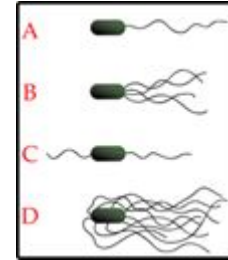
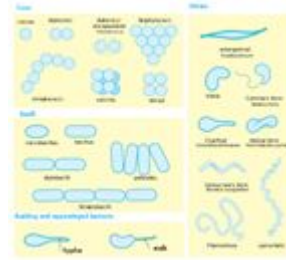
Microbial Identification Using Sequence Similarity

- Traditional microbial species identification
- Transition to genotyping methods
- Using NGS data to genotype pathogens
 - Average nucleotide identity (ANI)
 - MinHash dimensionality-reduction (Mash)



Phenotypic Characteristics of Taxonomic Value

- Morphology
- Motility
- Metabolism
- Physiology and Biochemical Data
- Cell Chemistry
- Others





Phenotypic Approach – Disadvantages

- Need experienced staff
 - Lots of validations, competencies, etc
- Can be a complicated process
 - Multiple tests and results necessary for interpretation
- Labor consuming
 - Hands on process for most tests
- Time consuming
 - Some testing needs to be sequential, often growth required

Genotypic Approach

Same genotypes, different phenotypes



Big Picture

- NGS genotyping methods attempt to match (sequence similarity) experimentally acquired DNA sequences with reference DNA sequences for microbial identification.
 - No reference sequence, no identification.
- Methods vary by:
 - How are the query and reference sequence represented?
 - Are the sequences converted to a new data structure or kept as a string?
 - How are the strings broken up?
 - Are the reads used natively or are they assembled first?
 - What is the algorithm for comparison?
 - Is BLAST used?
 - Is an algorithm optimized for the new data structure used?
 - Does the algorithm use any approximations?

Genotyping - Average Nucleotide Identity (ANI)

The next few slides borrow heavily from the presentations:

National Center for Emerging and Zoonotic Infectious Diseases



Whole Genome Sequence (WGS) of Enteric Bacteria using the BioNumerics RefID Database

Steven Stroika

PulseNet WGS Technical Lead

BioNumerics 7.6 Workshop for Analyzing WGS Data

May 2019



The Use of Average Nucleotide Identity (ANI) for Bacterial Identification

Patti Fields

for

Maryann Turnsek

Enteric Diseases Laboratory Branch (EDLB)

CDC

2017 APHL Annual Meeting

Providence, Rhode Island

June 14, 2017

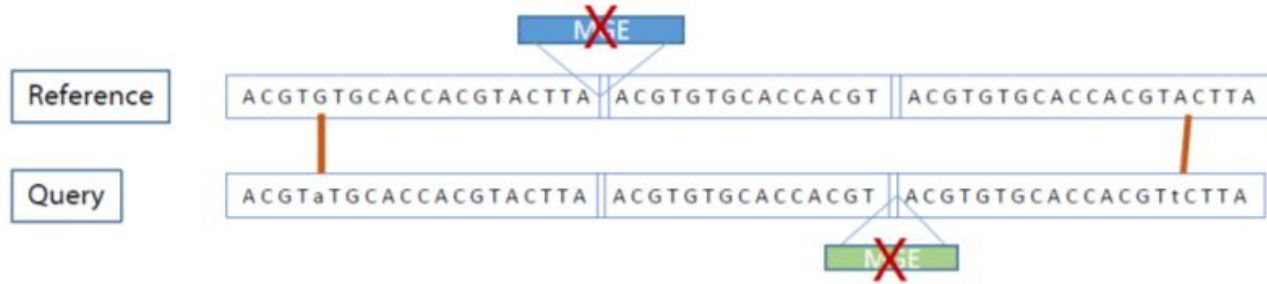
National Center for Emerging and Zoonotic Infectious Diseases
Division of Foodborne, Waterborne, and Environmental Diseases



What is Average Nucleotide Identity (ANI)

- A computation method to compare two genomes
 - Compares an unknown query sequence to a well-characterized reference genome.
 - Two calculations:
 - Compares the genetic similarity of shared sequences.
 - Determines the proportion of bases aligned.
- Closely mirrors comparisons by DNA-DNA hybridization
 - The traditional gold standard method for determining species boundaries.

How ANI Works



- Aligns shared sequences and calculates percent identical nucleotides
- Answers the question: Are these two genomes the same taxon? Yes or No
- In this example, 53/55 aligned bases = 96.4% identity
- The ANI “cutoff” value for % identity and % bases aligned is determined empirically for each taxon
 - Published values are on the order of 95% identity.

ANI vs DNA-DNA Hybridization

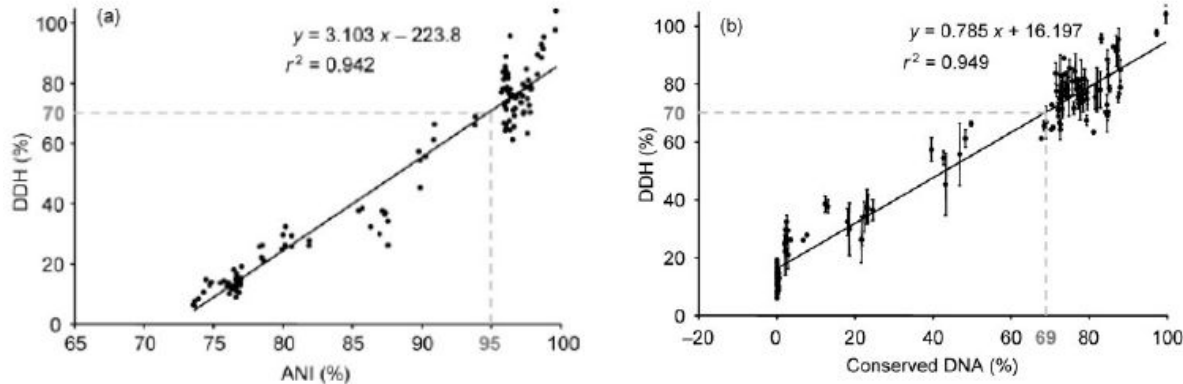


Fig. 1. Relationship between DDH values and genomic sequence identity and conservation. Each filled circle represents the value for DDH between two strains (y-axis), plotted against the ANI of the conserved genes between the strains (a) and the percentage of conserved DNA between the strains (b). The standard deviations for the DDH values, omitted from (a) for simplicity, are shown in (b). A linear trend line is shown, but other regression models were evaluated as well (see text). The horizontal broken lines denote the 70% DDH recommendation for species delineation, while the vertical broken lines denote the corresponding ANI (a) and percentage of conserved DNA (b) values for linear regression.

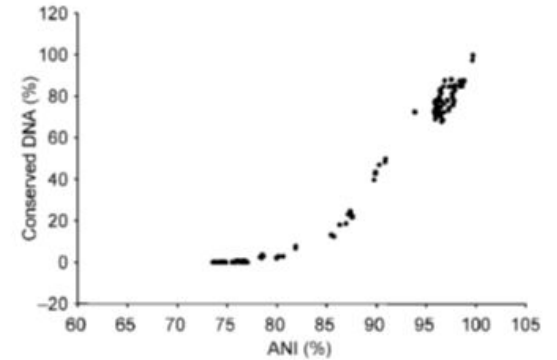


Fig. 2. Relationship between genomic sequence identity and conservation. Each filled circle represents the percentage of conserved DNA shared between two strains (determined at 90% nucleotide identity), plotted against the ANIs of their common genes.

ANI Algorithms

- ANI BLAST used originally (ANiB)
- ANI MUMer used in BN (ANiM)

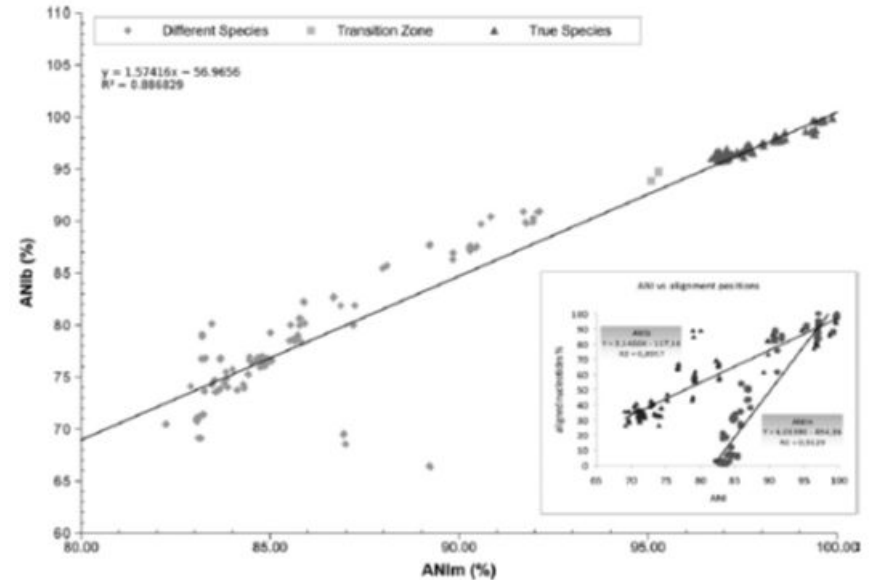


Fig. 1. Plotted results of ANiB versus ANiM. The triangles show those values that correspond to what taxonomists consider as “true” species according to the DDH values traditionally applied and that have previously been classified. *Inset* shows the regression lines of the pairwise comparisons of ANiB or ANiM values with their corresponding percentage of aligned stretches (percentage of nucleotides included in the study).

ANI in Public Health Bioinformatics



WGS data



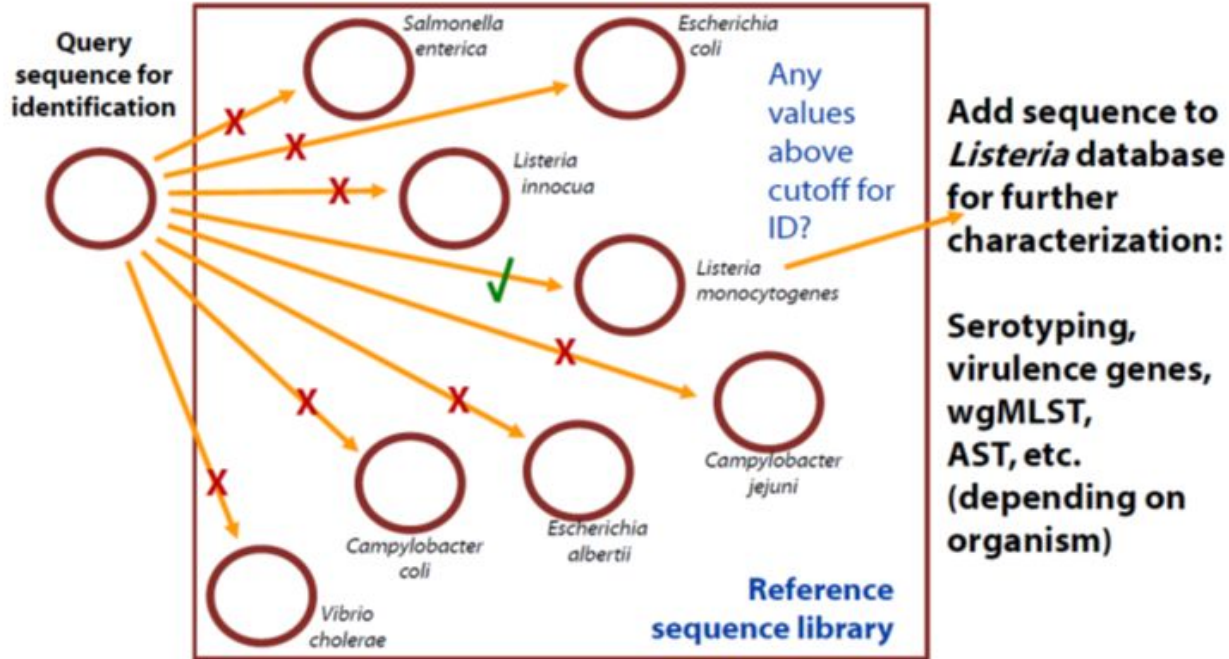
**Reference
Identification Database**
(ANI: Genus, Species)



**Organism-
specific Databases**

Further characterization:
MLST, serotype, lineage, AST,
virulence, plasmids, etc...

How searching works



“...determined empirically for each taxon.”

Tableau visualization of ANI values between and within species.



Empirical Values Used in BN

Genera	Species	ANI value (%)	Genome size (MB)
<i>Campylobacter</i>	<i>coli</i>	≥92	1.4-2.2
	<i>fetus</i>		
	<i>jejuni</i>		
	<i>lari</i>		
	<i>upsaliensis</i>		
	<i>hyointestinalis*</i>		
<i>Escherichia</i>	<i>albertii*</i>	≥95	4.5-5.5
	<i>coli</i> and <i>Shigella</i>		
	<i>fergusonii*</i>		
<i>Listeria</i>	<i>innocua*</i>	≥92	2.7-3.2
	<i>ivanovii*</i>		
	<i>marthii*</i>		
	<i>monocytogenes</i>		
	<i>seeligeri*</i>		
	<i>welshimeri*</i>		

Genera	Species	ANI value (%)	Genome size (MB)
<i>Salmonella</i>	<i>bongori</i>	≥93	4.5-5.0
	<i>enterica</i>		
<i>Vibrio</i>	<i>cholerae</i>	≥95	4.0-5.0
	<i>Parahaemolyticus</i>		
	<i>vulnificus</i>		
	<i>alginolyticus*</i>		
	<i>cidicii*</i>		
	<i>cincinnatiensis*</i>		
	<i>fluvialis*</i>		
	<i>furnissii*</i>		
	<i>garveyi*</i>		
	<i>metoecus*</i>		
	<i>metschnikovii*</i>		
<i>mimicus*</i>			
<i>navarrensii*</i>			

Pros and Cons for ANI

- Pros

- Replicates species determinations by DNA-DNA hybridization
- Very rapid: Compare two genomes in seconds
- Very robust: Reliable answer with 5X sequence coverage (based on down-sampling experiment)
- Relatively easy to interpret with clear cut off values

- Cons

- Definitive identification requires representative genome is in the Reference Sequence Library
 - New or unrepresented species cannot be identified
- Useful for comparing closely related bacteria only
 - Distantly related => No Match
- **As reference library gets bigger, computation time gets longer**

Genotyping – Mash (MinHash)



- Current ANI database for Bionumerics contains ~40 reference genomes.
- Current NCBI Pathogen Detection Browser (January 2020) contains ~500,000 isolates.

If you want to start dramatically increasing your reference database size, use reads rather than assemblies, and keep your search times short, you will need to reduce the dimensionality of your data and make some assumptions.

K-mers and Hash Tables

5' -AGGGCGGTTTAATAATCTACGGCTTATTGTTGAACGA-3'

```
AGGGCGGTTTAATAATCTACG
GGGCGGTTTAATAATCTACGG
GGCGGTTTAATAATCTACGGC
GCGGTTTAATAATCTACGGCT
CGGTTTAATAATCTACGGCTT
GGTTTAATAATCTACGGCTTA
GTTTAATAATCTACGGCTTAT
TTTAATAATCTACGGCTTATT
TTAATAATCTACGGCTTATTG
TAATAATCTACGGCTTATTGT
AATAATCTACGGCTTATTGTT
ATAATCTACGGCTTATTGTTG
TAAATCTACGGCTTATTGTTGA
AATCTACGGCTTATTGTTGAA
ATCTACGGCTTATTGTTGAAC
TCTACGGCTTATTGTTGAACG
CTACGGCTTATTGTTGAACGA
```

DNA sequence $L=37$

k-mer size $k=21$

#k-mers = $L - k + 1$

17 k-mers

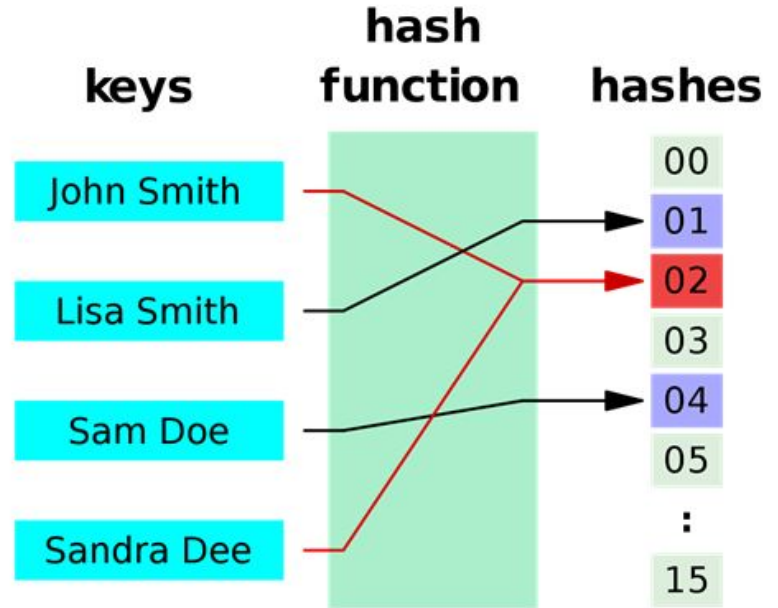
n^k possible k-mers

~17+ billion unique k-mers

K-mers and Hash Tables

Convert a string into a number in a reproducible way.

Numbers are faster to compare than text.



Mash: fast genome and metagenome distance estimation using MinHash

Brian G. Chikhi¹, Todd J. Treangen¹, Pat Medved², Adam S. Mallick³, Nicholas H. Bergman⁴, Sergey Korem⁵ and Adam M. Phillippy^{1*}

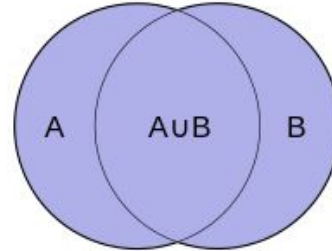
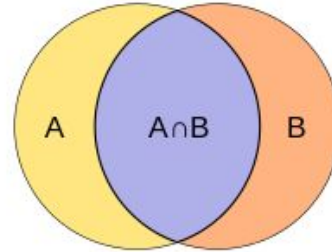
Mash

Jaccard Index

- Compute the ratio of the shared elements over all elements.

Mash

- Uses subsampling



$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} \approx \frac{|S(A \cup B) \cap S(A) \cap S(B)|}{|S(A \cup B)|}$$

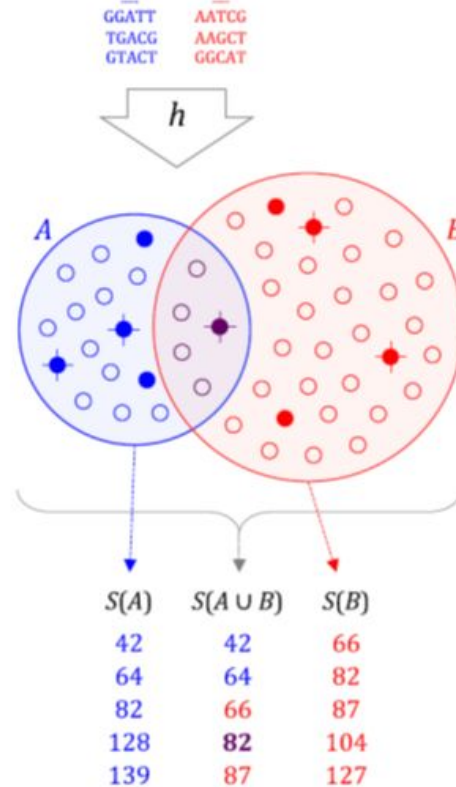
Mash: fast genome and metagenome distance estimation using MinHash

Brian D. Ondrejka¹, Todd J. Treangen¹, Pat Medved², Adam S. Malmer², Nicholas H. Bergman², Sergey Koren² and Adam M. Phillippy^{1*}

Mash (MinHash)

Mash will:

1. Create a hash sketch from k-mers of user defined size (15, 17, 19, 21, 23, ...)
2. Grab the X smallest hash values, where X is user defined (usually 500-1,000)
3. Compare these subsets as an estimate of similarity/dissimilarity and produce a Mash distance.



Mash: fast genome and metagenome distance estimation using MinHash

Brian D. Ondrejka¹, Todd J. Treangen¹, Pat Medved¹, Adam S. Mallick¹, Nicholas H. Bergman¹, Sergey Korem¹ and Adam M. Phillippy^{1*}

Mash

The Mash distance correlates well with ANI (or correctly 1-ANI), especially at high levels of similarity.

Not so good for distantly related species.

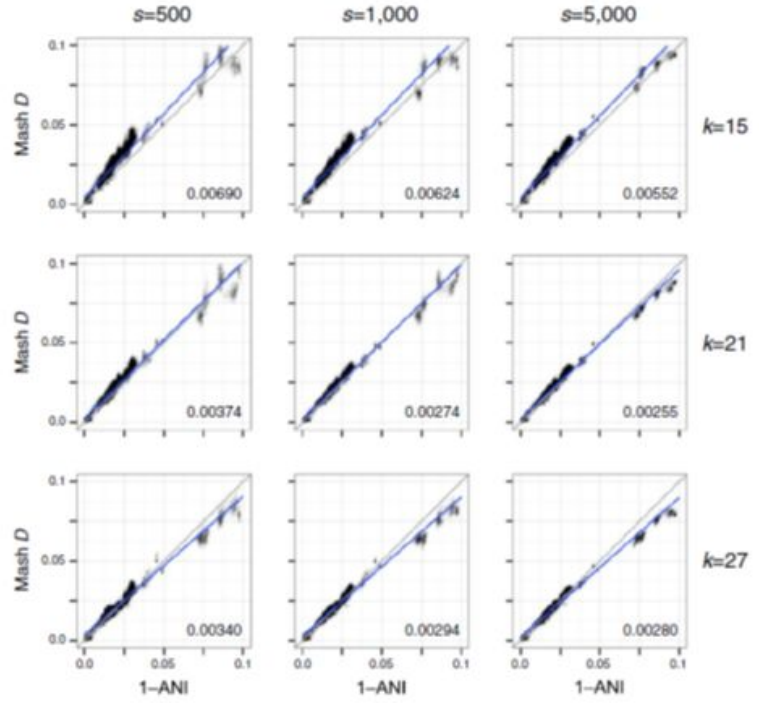


Tableau Analysis of Mash

1. Mash output – command line

Why is all this important?

1. These sequence similarity tools are becoming embedded in our workflows as NGS adoption continues.
2. These tools will need to be validated, and a deeper understanding of how they work, along with parameter optimization, is needed.
3. As NGS data increases, algorithms that use data reduction, subsampling, approximation, etc, to accelerate sequence similarity searching will become more and more necessary in order to take advantage of the wealth of data available.

Additional Resources

<http://www.staphb.org/training/pairwise/>