



**AMD TRAINING LEAD
and BIOINFORMATICS
REGIONAL RESOURCE**

Demystifying Series: Microbial ID Using NGS Data

Joel R. Sevinsky, Ph.D.

WY PHL AMD Training Lead

MA DPH AMD Bioinformatics Resource

January 22, 2020

Outline for today's webinar

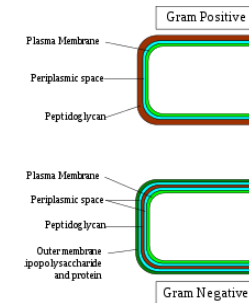
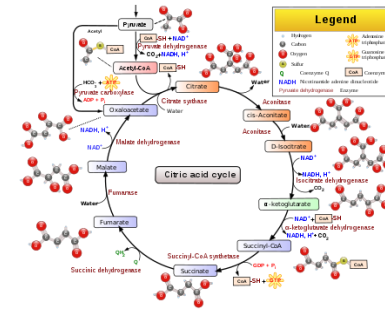
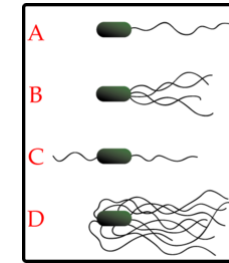
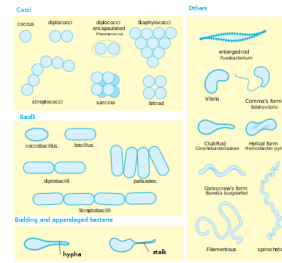
- Traditional microbial species identification
- Transition to genotyping methods
- Using NGS data to genotype pathogens
 - Average nucleotide identity (ANI)
 - MinHash dimensionality-reduction (Mash)

This webinar will be recorded and placed at <http://www.staphb.org>.



Phenotypic Characteristics of Taxonomic Value

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



Techniques of bacterial taxonomy
Yuhua Xin

China General Microbiological Culture Collection Center (CGMCC)
Institute of Microbiology
Chinese Academy of Sciences

AMD TRAINING LEAD *and* BIOINFORMATICS REGIONAL RESOURCE





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



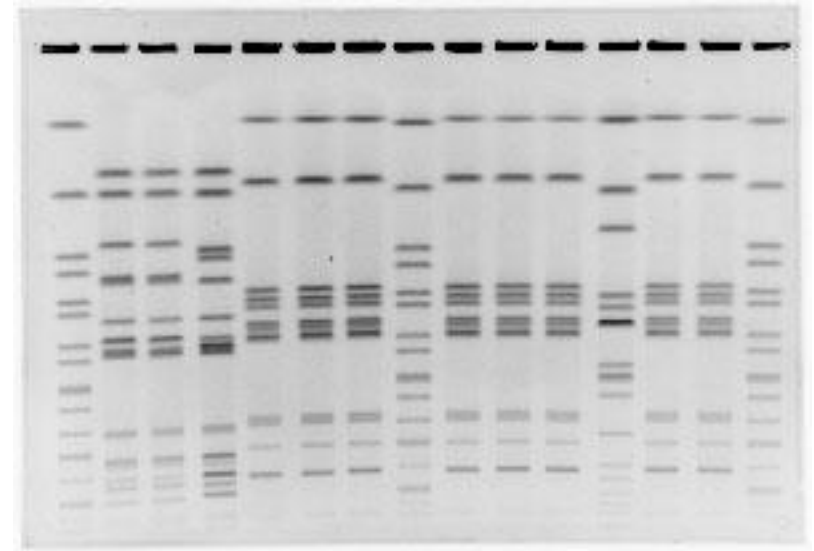
Genotypic Approach

Same genotypes, different phenotypes



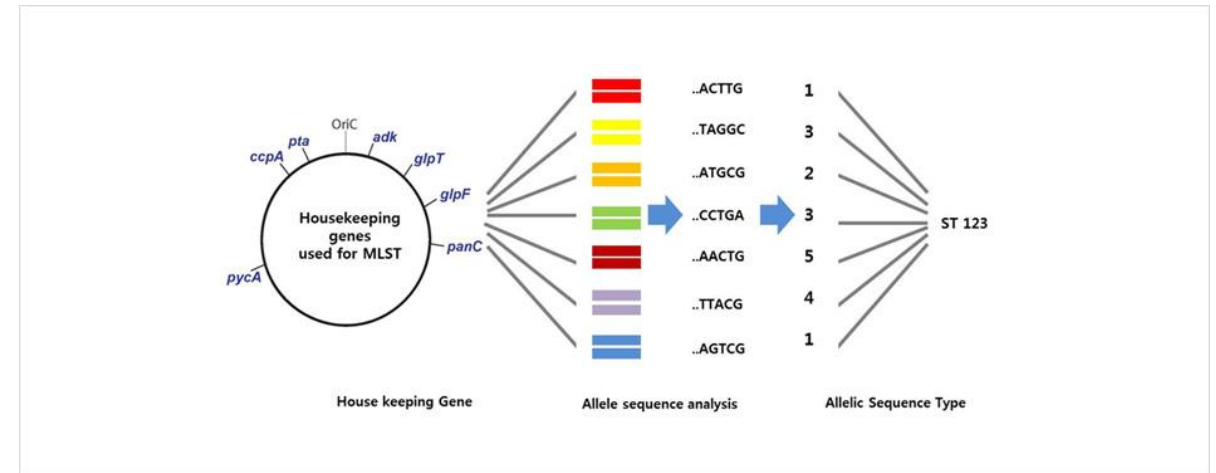
How to use genotype for microbial species ID

- PFGE – global examination of a subset of gDNA from a high level.
- 16S rDNA – specific examination of a small, highly discriminatory region of gDNA.
- MLST – specific examination of several small discriminatory regions of gDNA.
- NGS – global yet specific examination of almost all gDNA.



How to use genotype for microbial species ID

- PFGE – global examination of a subset of gDNA from a high level.
- 16S rDNA – specific examination of a small, highly discriminatory region of gDNA.
- **MLST – specific examination of several small discriminatory regions of gDNA.**
- NGS – global yet specific examination of almost all gDNA.

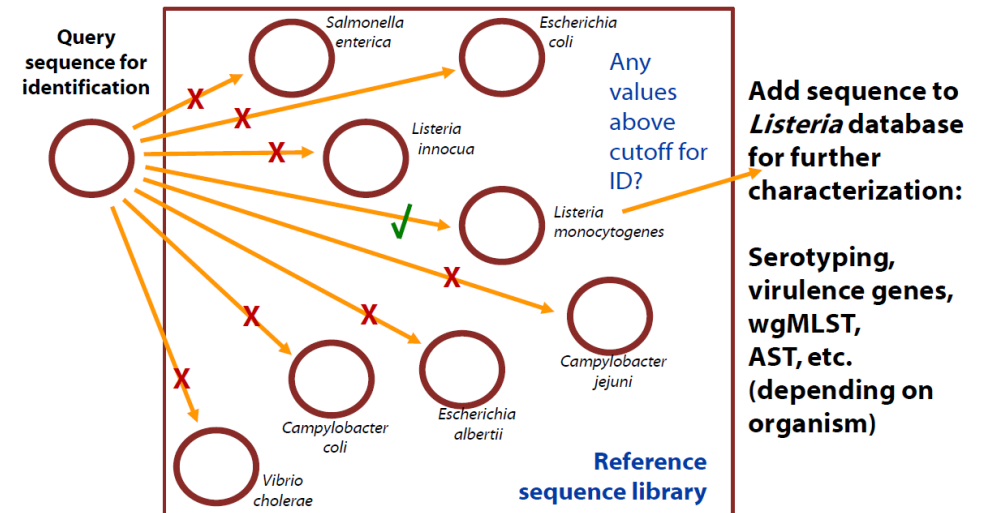


https://www.macrogen.com/en/business/ces_service4.php



How to use genotype for microbial species ID

- PFGE – global examination of a subset of gDNA from a high level.
- 16S rDNA – specific examination of a small, highly discriminatory region of gDNA.
- MLST – specific examination of several small discriminatory regions of gDNA.
- NGS – global yet specific examination of almost all gDNA.



Big Picture

- NGS genotyping methods attempt to match 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



AMD TRAINING LEAD *and* BIOINFORMATICS REGIONAL RESOURCE



What is Average Nucleotide Identity (ANI)

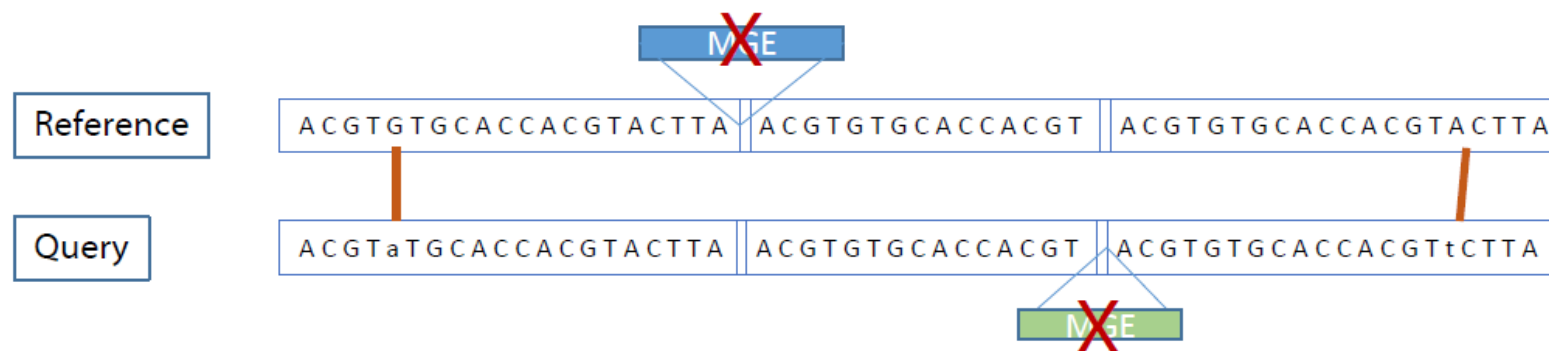
- A computation method to compare two genomes
 - Compares an unknown query sequence to a well-characterized reference genome.
 - Two calculations:
 - Compares the genetic similarity of shared sequences.
 - Determines the proportion of bases aligned.
- Closely mirrors comparisons by DNA-DNA hybridization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

Konstantinidis KT, Tiedje JM. Genomic insights that advance the species definition for prokaryotes. Proc Natl Acad Sci U S A. 2005 Feb 15;102(7):2567-72.

Richter M, Rosselló-Móra R. Shifting the genomic gold standard for the prokaryotic species definition. Proc Natl Acad Sci U S A. 2009 Nov 10;106(45):19126-31.



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

DNA-DNA hybridization values and their relationship to whole-genome sequence similarities

Johan Goris,^{1†} Konstantinos T. Konstantinidis,^{1‡} Joel A. Klappenbach,¹ Tom Coenye,² Peter Vandamme² and James M. Tiedje¹

¹Center for Microbial Ecology, Michigan State University, East Lansing, MI 48824, USA

²Laboratory for Microbiology, Ghent University, K. L. Ledeganckstraat 35, B-9000 Ghent, Belgium

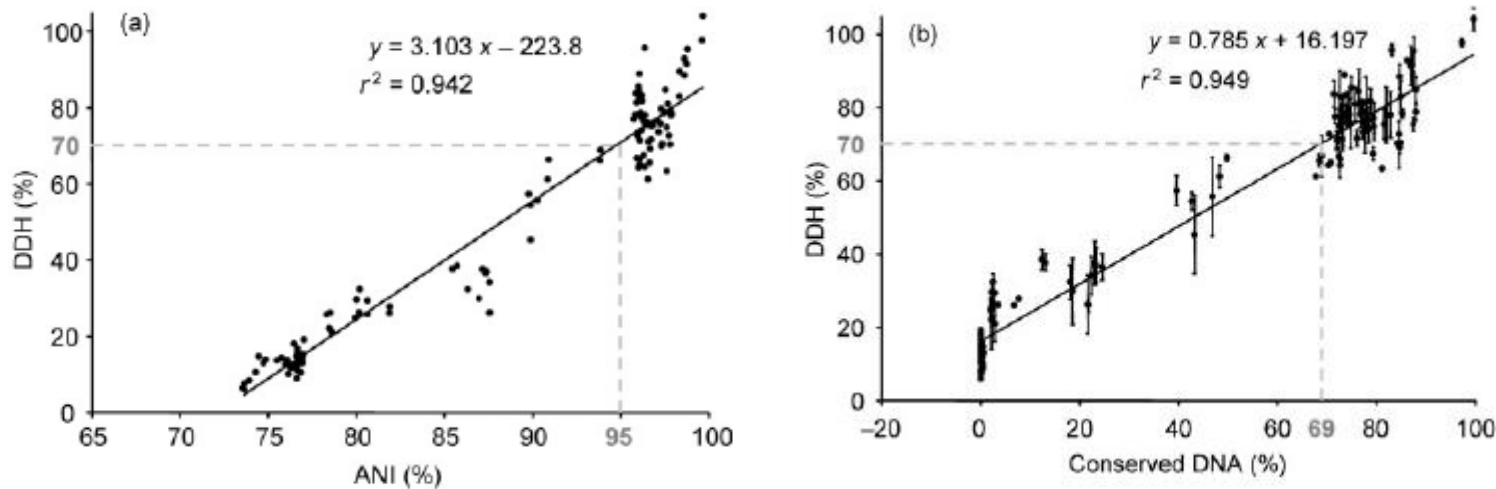


Fig. 1. Relationship between DDH values and genomic sequence identity and conservation. Each filled circle represents the value for DDH between two strains (y-axis), plotted against the ANI of the conserved genes between the strains (a) and the percentage of conserved DNA between the strains (b). The 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.

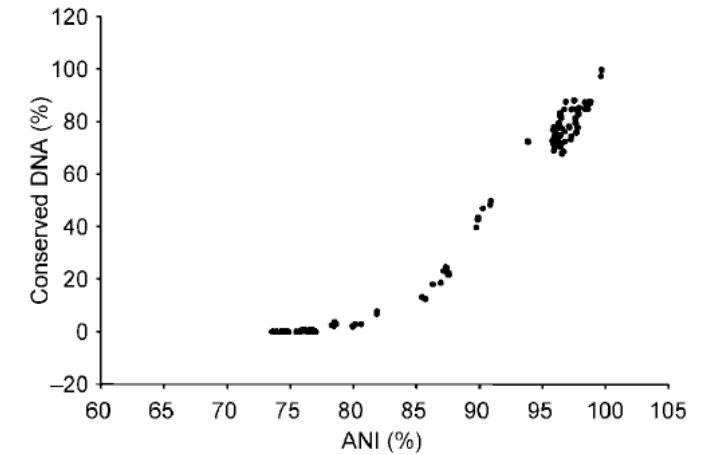


Fig. 2. Relationship between genomic sequence identity and conservation. Each filled circle represents the percentage of conserved DNA shared between two strains (determined at 90% nucleotide identity), plotted against the ANIs of their common genes.



ANI Algorithms

- ANI Blast used originally (ANiB)
- ANI MUMer used in BN (ANIm)

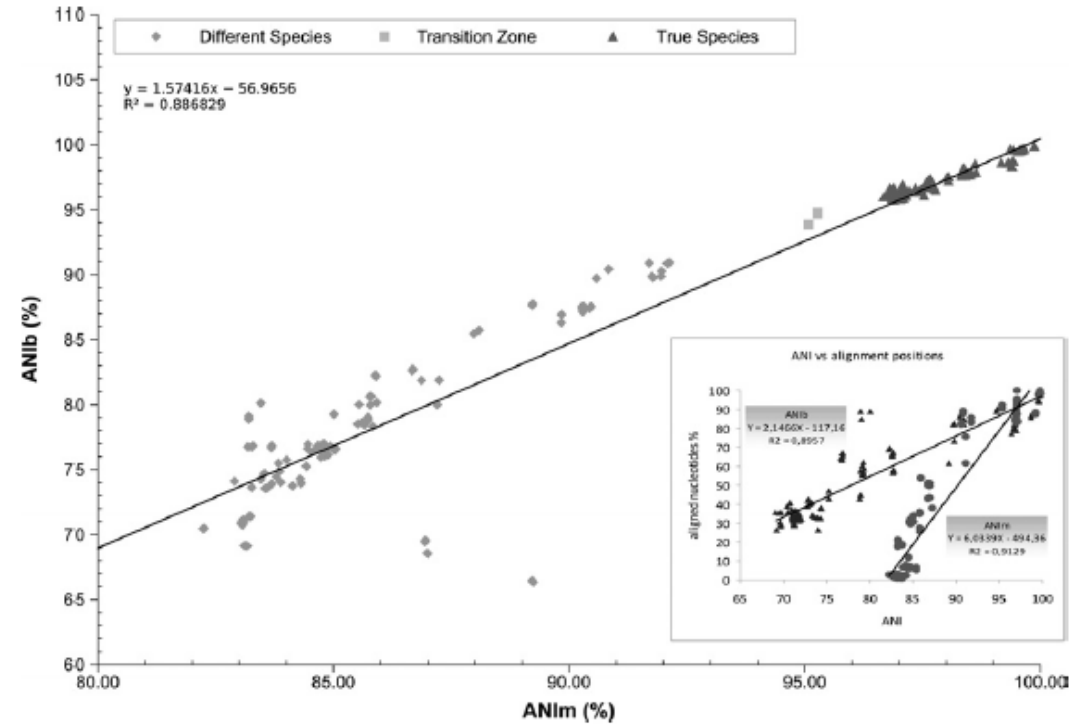
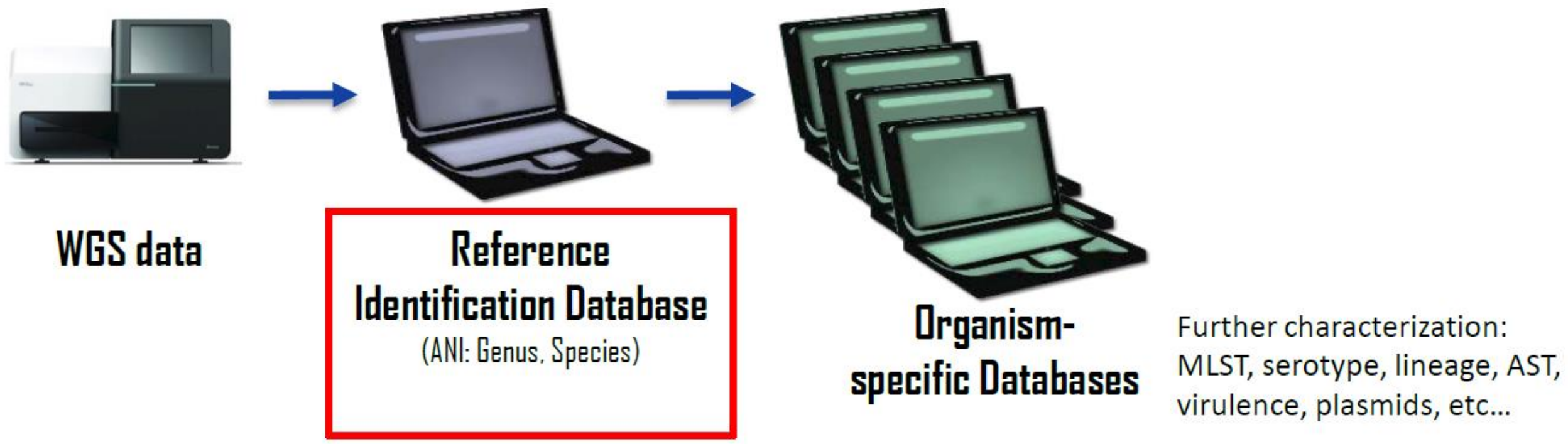


Fig. 1. Plotted results of ANiB versus ANIm. The triangles show those values that correspond to what taxonomists consider as “true” species according to the DDH values traditionally applied and that have previously been classified. *Inset* shows the regression lines of the pairwise comparisons of ANiB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

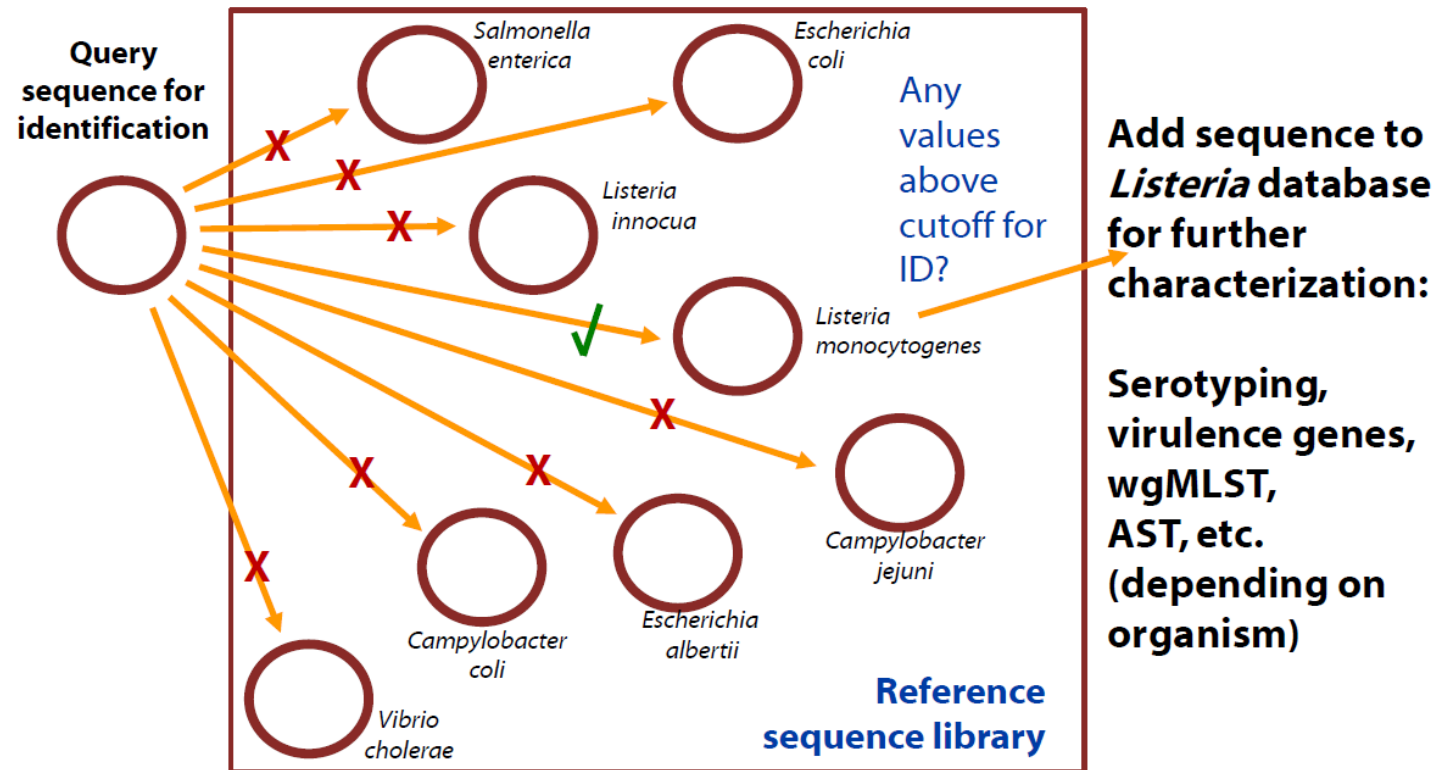




ANI in Public Health Bioinformatics



How searching works



“...determined empirically for each taxon.”

1. 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>navarrensis*</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)

Ondov et al. *Genome Biology* (2016) 17:132
DOI 10.1186/s13059-016-0997-x

Genome Biology

SOFTWARE

Open Access

Mash: fast genome and metagenome
distance estimation using MinHash



Brian D. Ondov¹, Todd J. Treangen¹, Páll Melsted², Adam B. Mallonee¹, Nicholas H. Bergman¹, Sergey Koren³
and Adam M. Phillippy^{3*}

- 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

5' -AGGGCGGTTTAATAATCTACGGCTTATTGTTGAACGA-3'

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

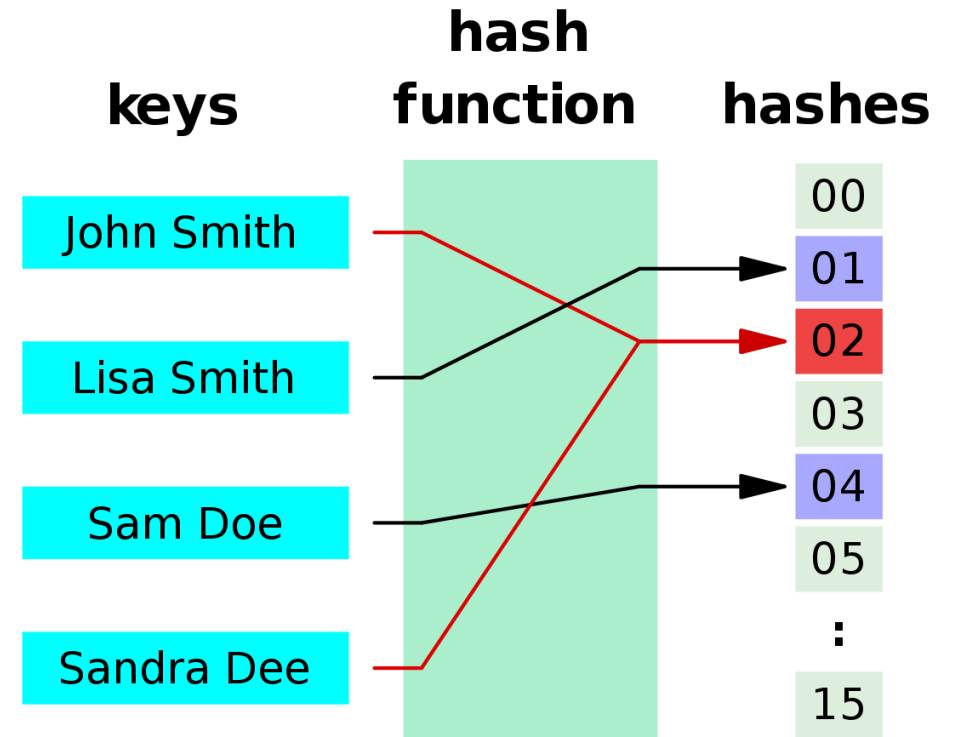
- DNA sequence $L=37$
- k-mer size $k=21$
- # of k-mers = $L - k + 1 = 17$ k-mers
- n^k possible k-mers
- ~17+ billion unique k-mers



Hash Tables

Convert a string into a number in a reproducible way.

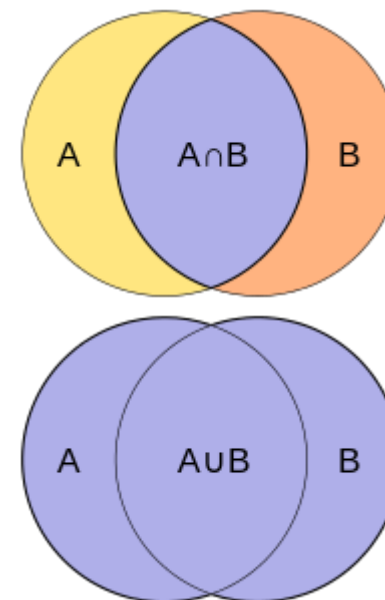
Numbers are faster to compare than text.



Jaccard Index (Similarity Coefficient)

Jaccard Index

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



Mash

- Uses subsampling

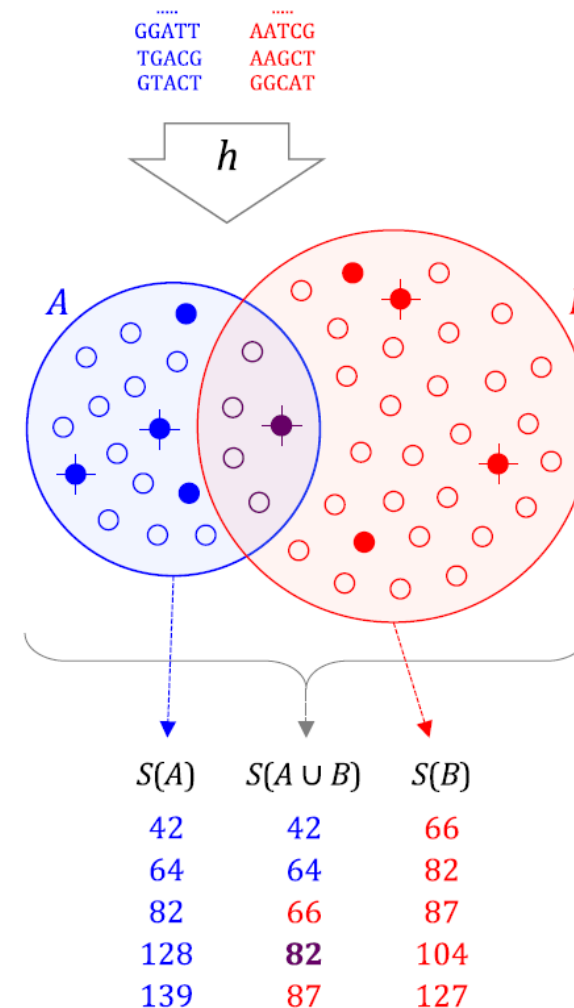
$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} \approx \frac{|S(A \cup B) \cap S(A) \cap S(B)|}{|S(A \cup B)|}$$



Mash (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. Ondov¹, Todd J. Treangen¹, Páll Melsted², Adam B. Mallonee³, Nicholas H. Bergman¹, Sergey Koren³ and Adam M. Phillippy^{2*}



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.

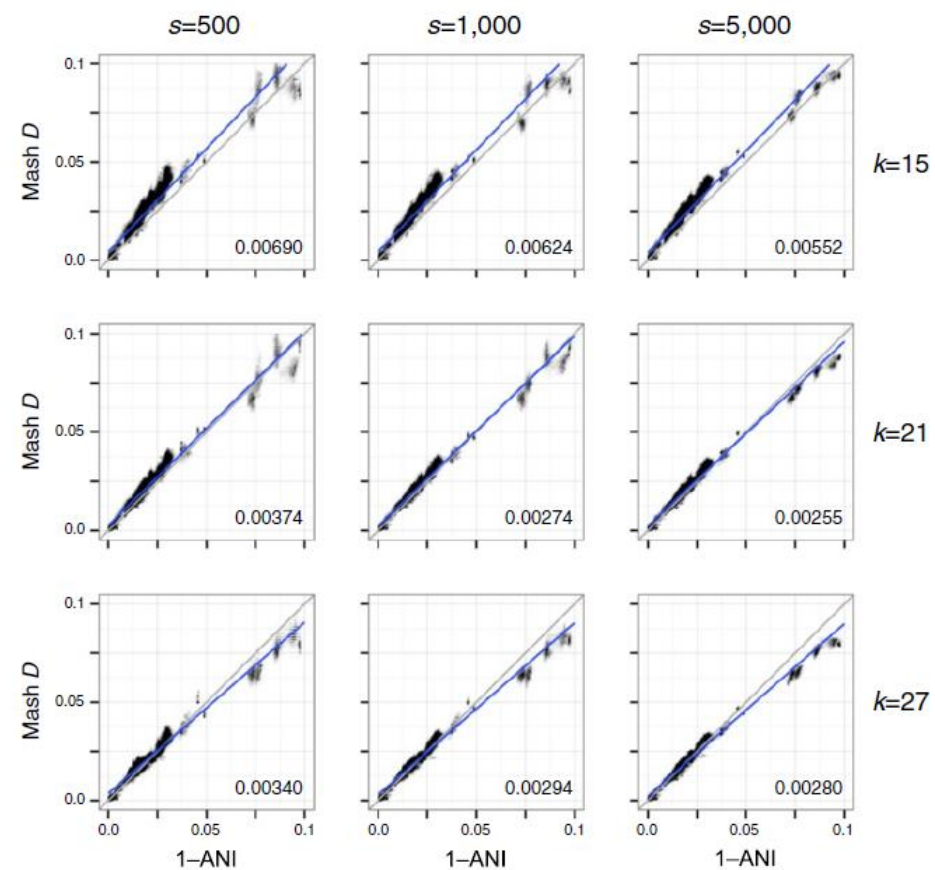


Tableau Analysis of Mash

1. Mash output – command line
2. Mash output - Tableau
3. Escherichia spp. analysis



Why is all this important?

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



Questions?

AMD TRAINING LEAD *and* BIOINFORMATICS REGIONAL RESOURCE

