

<https://xkcd.com/1691/>

# Mash: fast genome and metagenome distance estimation

Don L. Armstrong

Institute for Genomic Biology, Computing Genomes for Reproductive Health, University of Illinois, Urbana-Champaign

June 23, 2016

Code and slides are here:



<http://dla2.us/p/mashminhash2016>

SOFTWARE

Open Access



# Mash: fast genome and metagenome distance estimation using MinHash

Brian D. Ondov<sup>1</sup>, Todd J. Treangen<sup>1</sup>, Páll Melsted<sup>2</sup>, Adam B. Mallonee<sup>1</sup>, Nicholas H. Bergman<sup>1</sup>, Sergey Koren<sup>3</sup> and Adam M. Phillippy<sup>3\*</sup>

## Abstract

Mash extends the MinHash dimensionality-reduction technique to include a pairwise mutation distance and  $P$  value significance test, enabling the efficient clustering and search of massive sequence collections. Mash reduces large sequences and sequence sets to small, representative sketches, from which global mutation distances can be rapidly estimated. We demonstrate several use cases, including the clustering of all 54,118 NCBI RefSeq genomes in 33 CPU h; real-time database search using assembled or unassembled Illumina, Pacific Biosciences, and Oxford Nanopore data; and the scalable clustering of hundreds of metagenomic samples by composition. Mash is freely released under a BSD license (<https://github.com/marbl/mash>).

**Keywords:** Comparative genomics, Genomic distance, Alignment, Sequencing, Nanopore, Metagenomics

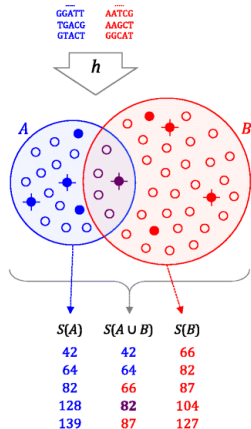
## Background

When BLAST was first published in 1990 [1], there were less than 50 million bases of nucleotide sequence in the

any problem where an approximate, global distance is acceptable, e.g. to triage and cluster sequence data, assign species labels, build large guide trees, identify

# The Problem

# MinHash Algorithm



- Decompose dataset into k-mers
- Hash k-mers (32/64bit)
- Estimate Jaccard Index  $J(A, B)$

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

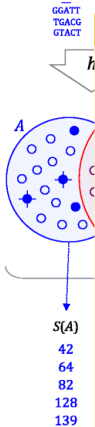
# MinHash Algorithm

## What about strandedness

- Take lexically lowest sequence
  - Given 7-mers 5'-ACTGCAC-3' and its reverse complement, 5'-GTGCAGT-3'
  - A G
  - Use ACTGCAC
- Because  $S(A \cup B)$  is a random sample of  $A \cup B$  the fraction of elements in  $S(A \cup B)$  which are shared by  $S(A)$  and  $S(B)$  is an unbiased estimate of  $J(A, B)$

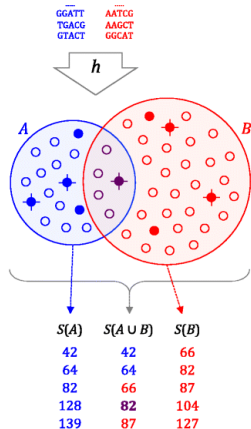
$$J(A, B) \approx \frac{|S(A \cup B) \cap S(A) \cap S(B)|}{|S(A \cup B)|}$$

into k-mers  
it)  
ex  $J(A, B)$



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

# MinHash Algorithm



- Decompose dataset into k-mers
- Hash k-mers (32/64bit)
- Estimate Jaccard Index  $J(A, B)$

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

# MinHash Algorithm

GGATT  
TGACG  
GTACT

## Estimating Jaccard Index

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

- Sample randomly from  $A$  and  $B$
- Because  $S(A \cup B)$  is a random sample of  $A \cup B$  the fraction of elements in  $S(A \cup B)$  which are shared by  $S(A)$  and  $S(B)$  is an unbiased estimate of  $J(A, B)$

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

into k-mers  
it)  
ex  $J(A, B)$   
y sample?  
!

$S(A)$

42  
64  
82  
128  
139

0/ 12/

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



# Hash functions and their properties

- A function which can map arbitrary sized data to output of fixed size
- They're used everywhere: caches, duplicate filtering, cryptography, finding similar substrings, finding similar records, etc.
- Good hash functions are deterministic, **uniform**, and rarely produce collisions.
- Cryptographic hash functions are non-invertible (SHA-512, etc.)
- Other hash functions may have continuity (IE, AAAB and AAAC will hash to nearby values)
- for example, the MD5 of "DON" is  
78ed35b3eb1f21c6179d824610d6d5c9 while "DO" is  
194199909f0994cd6244c4442642f3ff

# MurmurHash3 – properties

- Simple
- Good distribution
- Good collision resistance
- Good avalanche properties (single bit in the input changes the output significantly).

# MurmurHash3 – properties

- Simple
- Good distribution
- Good collision resistance
- Good avalanche properties (single bit in the input changes the output significantly).
  - 0 -> 25dd8f33676144d5e42427eff1c8546724d63d15
  - 1 -> 6c30934a0ea2c0473d37b6d8bb5b955b435a8bc1
  - 2 -> 315a5aa84aa6cfa4f3fb4b652a596770be0365e8

# MurmurHash3 – properties

- Simple
- Good distribution
- Good collision resistance
- Good avalanche properties (single bit in the input changes the output significantly).
- Fast
- Key questions:
  - 1 Does it map DNA and protein sequences uniformly into output?
  - 2 Does the non-random input distribution of DNA/protein sequences expose pathologic behavior of the hash function?
- Partially answered by the smhasher test suite which exposed issues with MurmurHash2 (which is ideally the DieHarder of hash testing.)

# Why is this quick

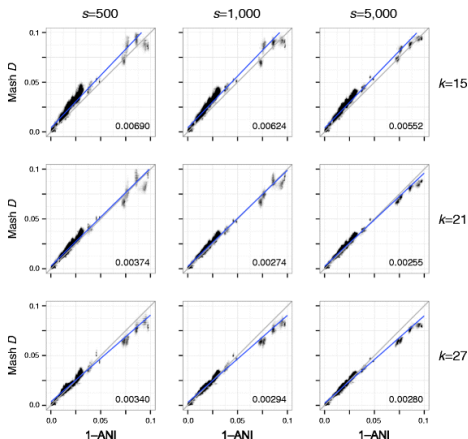
- Hashing is relatively fast
- Only read the data once
- Can be trivially parallized
- Only need to keep the  $k$  lowest hash values for each dataset
- Much faster than string alignment

**Table 2** Mash runtime and output size for all-pairs RefSeq computation using various sketch and k-mer sizes

Sketch size	k = 16				k = 21			
	Sketch (CPU h)	Dist (CPU h)	Size (Mb)	gzip (Mb)	Sketch (CPU h)	Dist (CPU h)	Size (Mb)	gzip (Mb)
500	26.4	8.4	120.1	89.7	31.3	9.0	229.8	201.8
1000	27.7	15.9	224.9	179.7	31.3	17.4	439.2	399.6
5000	26.4	74.5	1022.5	873.8	31.6	83.6	2034.5	1924.6
10,000	26.8	146.9	1961.8	1691.1	31.7	164.0	3913.0	3696.2

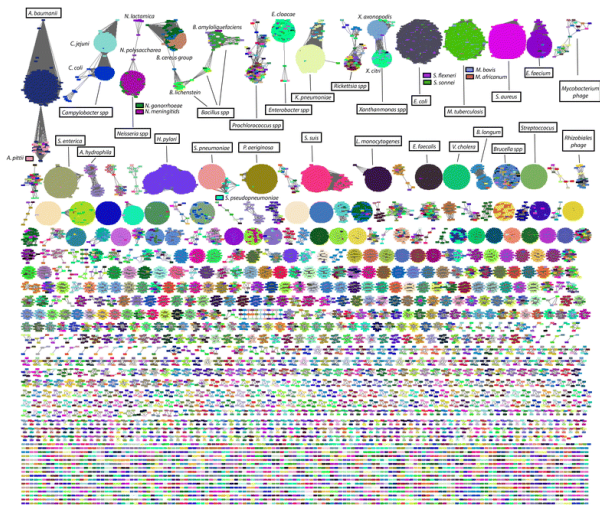
*Sketch*: CPU h required for the Mash *sketch* operation for all 54,118 RefSeq genomes. *Dist*: CPU h required for the Mash *dist* table operation for all pairs of sketches. *Size*: combined size of the resulting sketches in megabytes. *gzip*: combined size of the resulting sketches after gzip compression

# How does it stack up against ANI?



- Sketch size  $s$
- $k$ -mer size  $k$
- Mash D is the mash distance
- ANI is the average nucleotide identity between genomes
- Increasing  $s$  improves match;  $k$  seems to max near 27.
- ANI only considers core genome, so at some point ANI and Mash will diverge

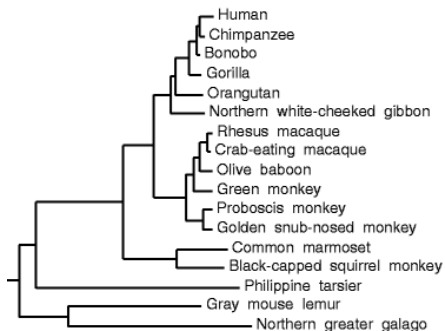




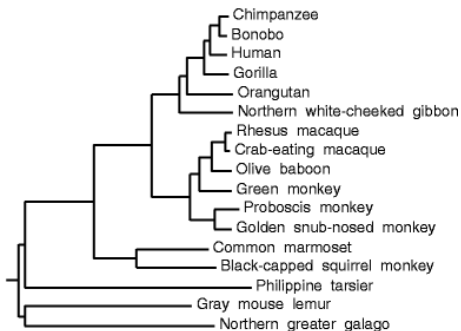
- *De novo* clustering of all RefSeq genomes
- Each node is a genome
- Nodes are connected if Mash distance  $\leq 0.05$  and  $p \leq 10^{-10}$



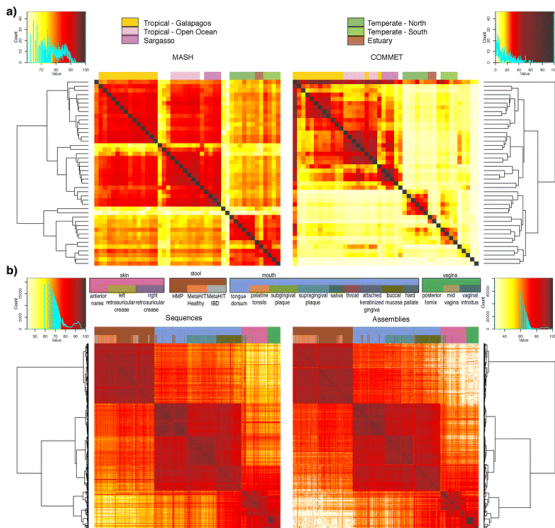
a) UCSC genome browser



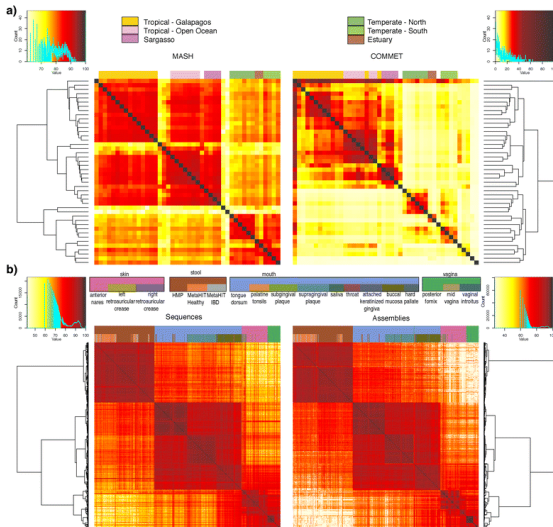
b) Mash



Primate species trees calculated by mash (right) in comparison to UCSC genome browser trees (left)



- Comparison between mash and COMMET using Global ocean survey clustering
- Mostly reproduces the clusters seen in the original study



- Comparison of sequences and assemblies from 888 sequencing runs and 879 assemblies from MHP and MetaHit projects
- Mostly reproduces the clusters seen in the original study
- COMMET would have required 140,000 hours to run this analysis; mash did it in 25 hours or so.

# Conclusions

- Compare genomes on a massive scale which would be impossible for typical techniques
- Good for rapid triaging of sequencing data to identify outliers
- Metagenomics and identification of samples
- Caveats
  - Based on k-mer sketch; possibility of batch effects
  - Phylogeny reconstructions are approximate, and do not track mutation rates, etc.
  - What about deconvolution of mixed samples? Probably need different techniques.

# References