



# Taxonomic classifiers

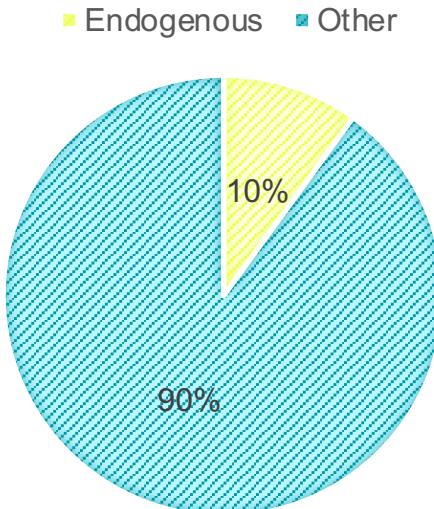
IMPRS Lecture – Maxime Borry





# Why do I care ?

## DNA

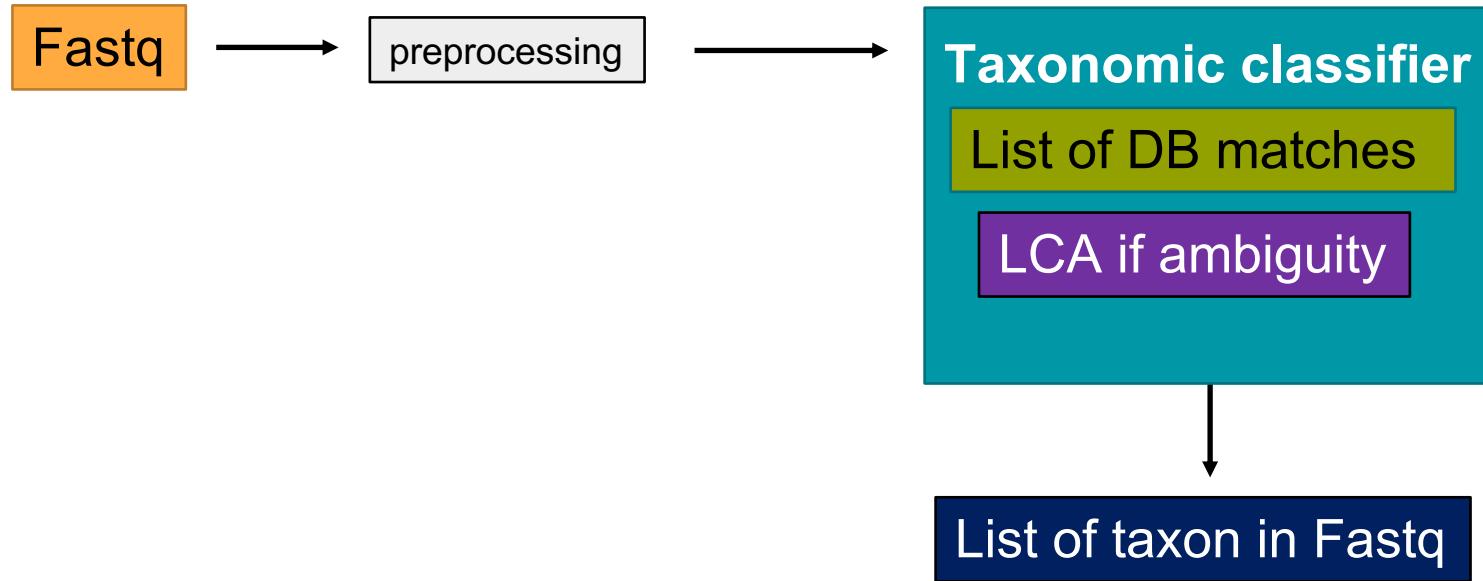


To extract information from the 90% "other", you need a taxonomic classifier to answer the question:

*Who is there ?*

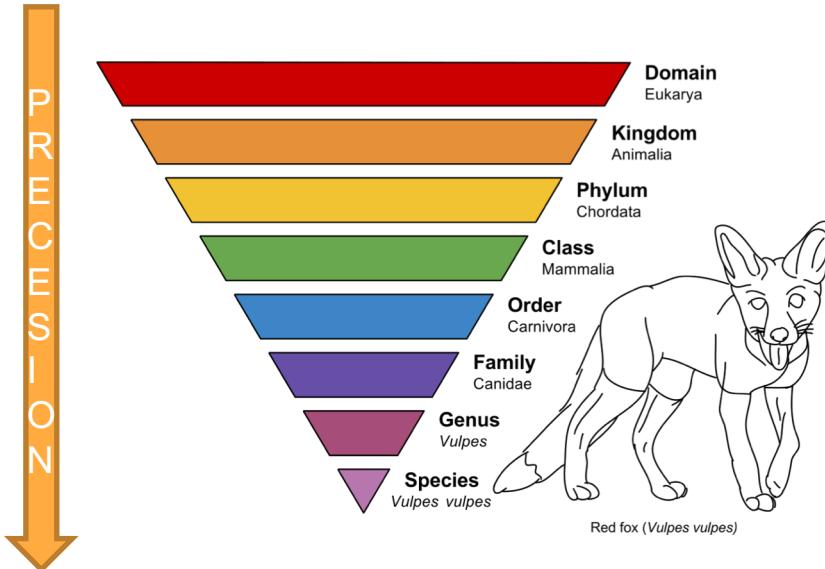


# What is a taxonomic classifier ?

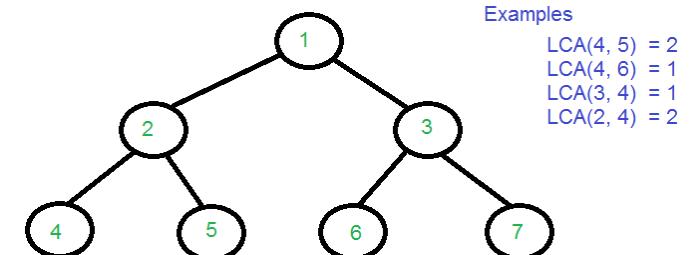




# Why not species classifier ?



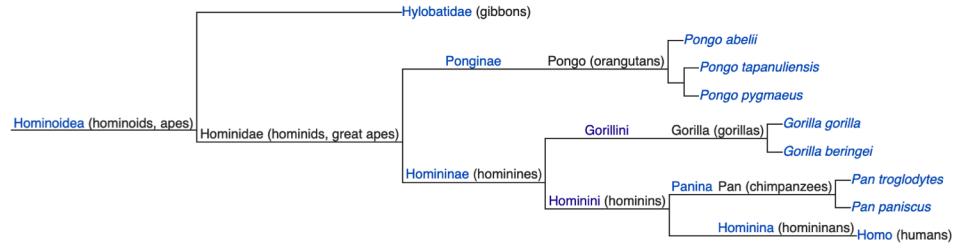
- Species level assignation is not always possible.
- Possibility of hits in different species
- Ambiguities solved by LCA (Lowest Common Ancestor) algorithm.





# LCA example

Hit	Identity
<i>Pan paniscus</i>	97%
<i>Pan troglodytes</i>	96%
<i>Homo sapiens</i>	92%
<i>Gorilla gorilla</i>	87%



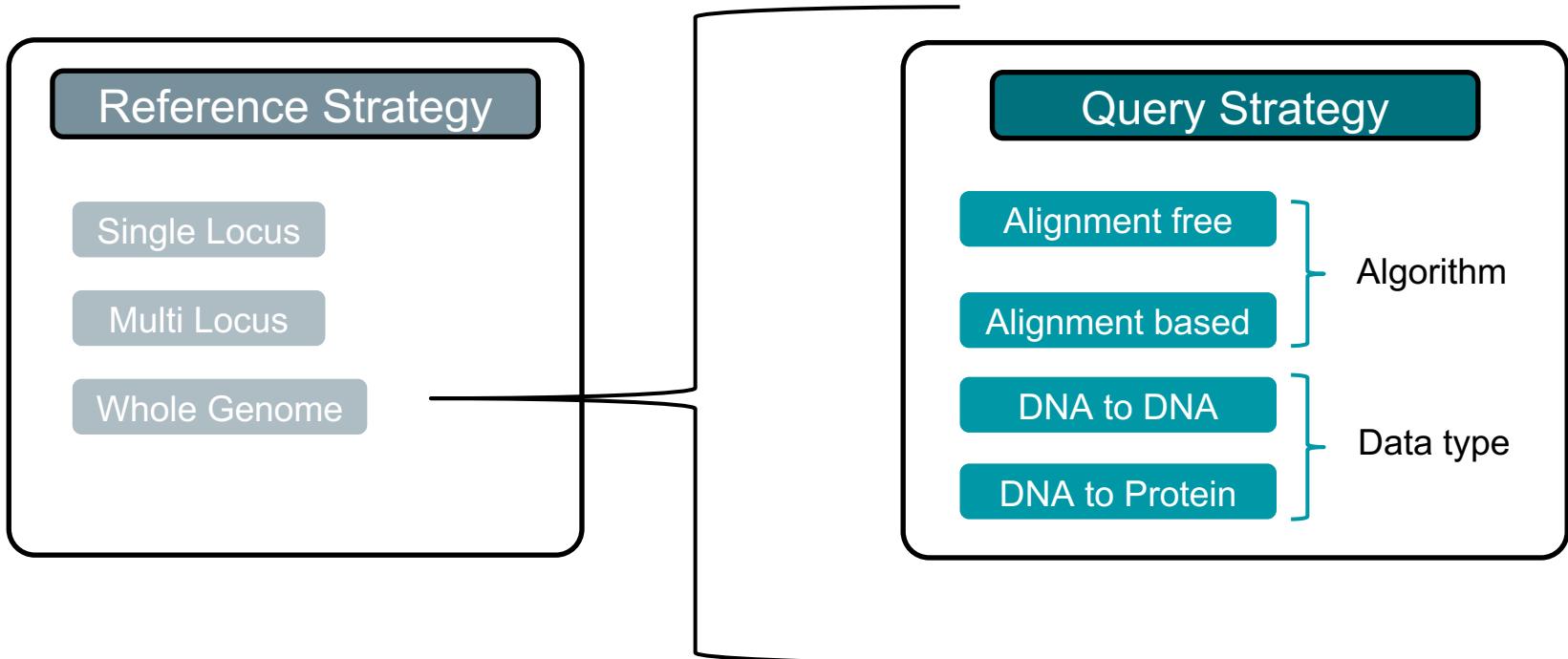
LCA 95%      Pan

LCA 90%      Hominini

LCA 85%      Homininae



# Taxonomic classifiers overview



# Reference strategy: single locus

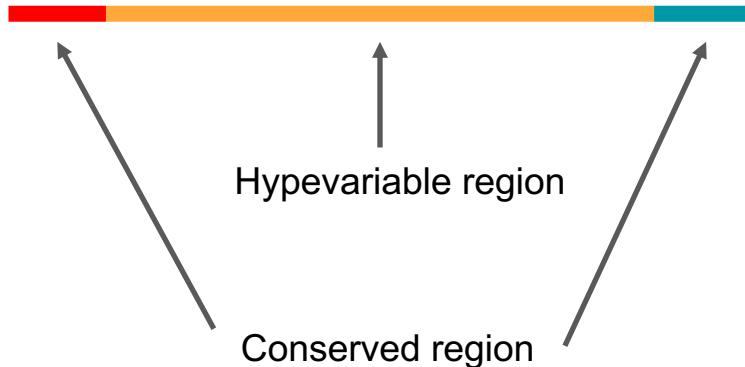
PCR

WGS



(Also known as Amplicon metataxonomics, Phylotyping, Metabarcoding)

Targeted amplification and deep sequencing of clade-universal gene



- Amplification of locus by PCR with primers targeting conserved regions or directly from WGS
- (Deep) Sequencing of amplicons
- Comparison to reference marker database



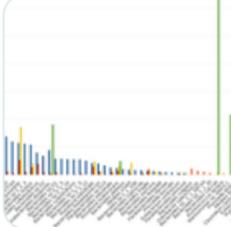
# Side note: vocabulary matters ! (to some)

## Metataxonomics vs Metagenomics

 **Jonathan Eisen**   
@phylogenomics 

Well, this software might be useful for **#metagenomics** but drives me crazy when people refer to 16S PCR as **#metagenomics**  
[plosone.org/article/info:d...](http://plosone.org/article/info:d...)

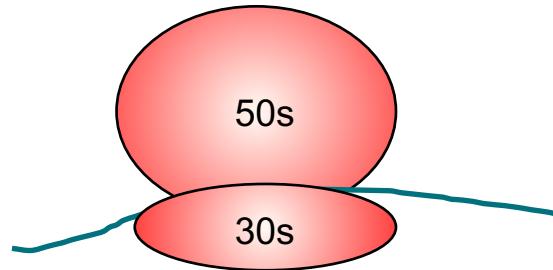
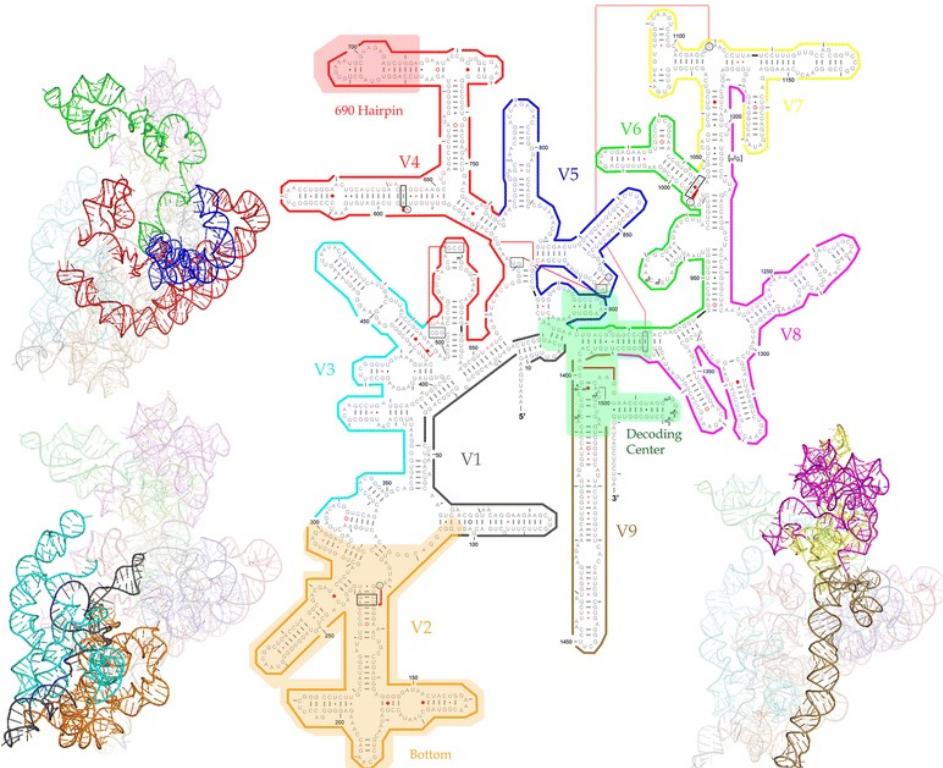
 20 8:22 AM - Aug 22, 2012 



**Genometa - A Fast and Accurate Classifier f...**  
Summary Metagenomic studies use high-throughput sequence data to investigate microbial communities *in situ*. However, considerable  
[journals.plos.org](http://journals.plos.org)

 See Jonathan Eisen's other Tweets >

# Single locus marker genes for bacteria: 16s



## 16s rRNA

- Part of the 30s prokaryotic ribosome subunit
- Stems are more stable -> conserved
- Loops are mutating faster -> variable

Yang, B., Wang, Y., & Qian, P. Y. (2016). Sensitivity and correlation of hypervariable regions in 16S rRNA genes in phylogenetic analysis. *BMC bioinformatics*, 17(1), 135.



# Single locus taxon assignation

**Spoiler:** Reads are clustered by sequence identity

The example of **QIIME**: OTU picking

**OTU (Operational Taxonomic Unit):** Cluster of Organism grouped by sequence identity (usually 97%)

## **De Novo OTU picking:**

Reads are clustered against one another without any external reference sequence collection

## **Closed Reference OTU picking**

Reads are clustered against a reference sequence collection and any reads which do not hit a sequence in the reference sequence collection are excluded

## **Open Reference OTU picking:**

Reads are clustered against a reference sequence collection and any reads which do not hit the reference sequence collection are subsequently clustered de novo.



# Single locus metaxomics tools and databases

## Tools



Qiime 2



Mothur



Dada 2

## Databases



**GREEN**GENES



Maxime Borry

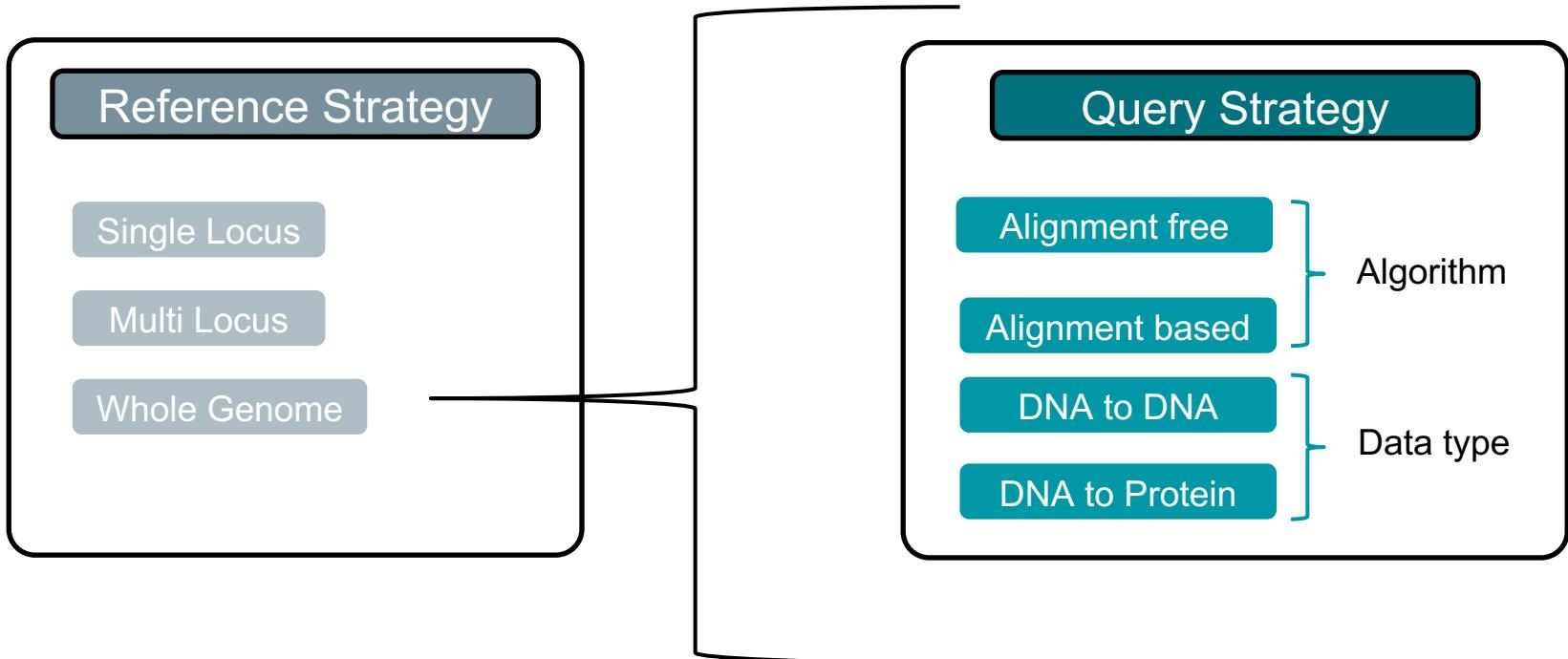


# 16s single locus limitations

- Not specific enough for sub-species classification
- Amplification bias: primer binding and GC-content
- If recovered from WGS data, very low amount (~ 0.2 - 0.6 %)

## Limitations specific to ADNA

- Hypervariable regions length greater than typical aDNA fragment length



# Reference strategy: multi-locus

WGS



## Use a set of single-copy housekeeping genes

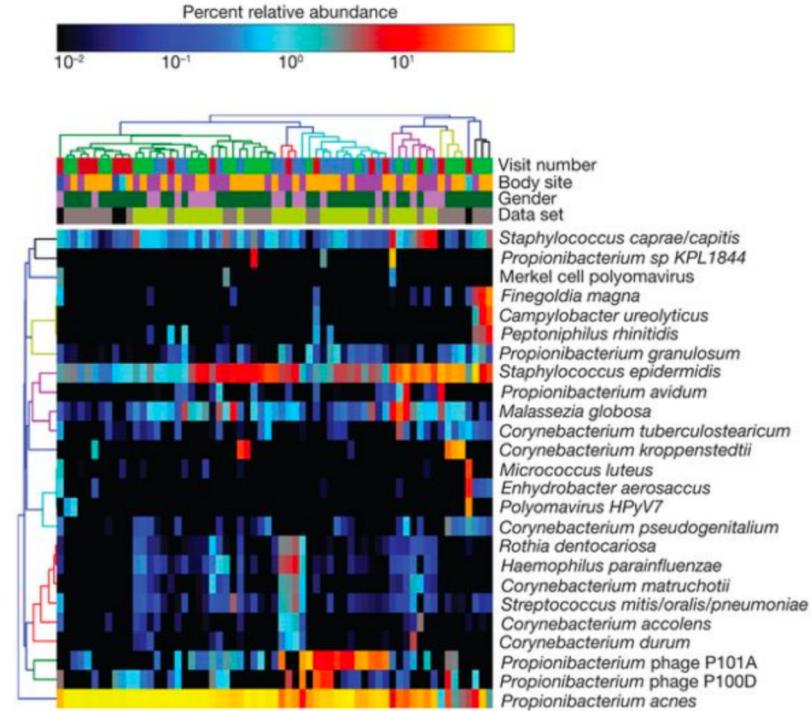
- Sequence divergence greater than 16s
- Single-copy: access to quantification

## Metaphlan 2

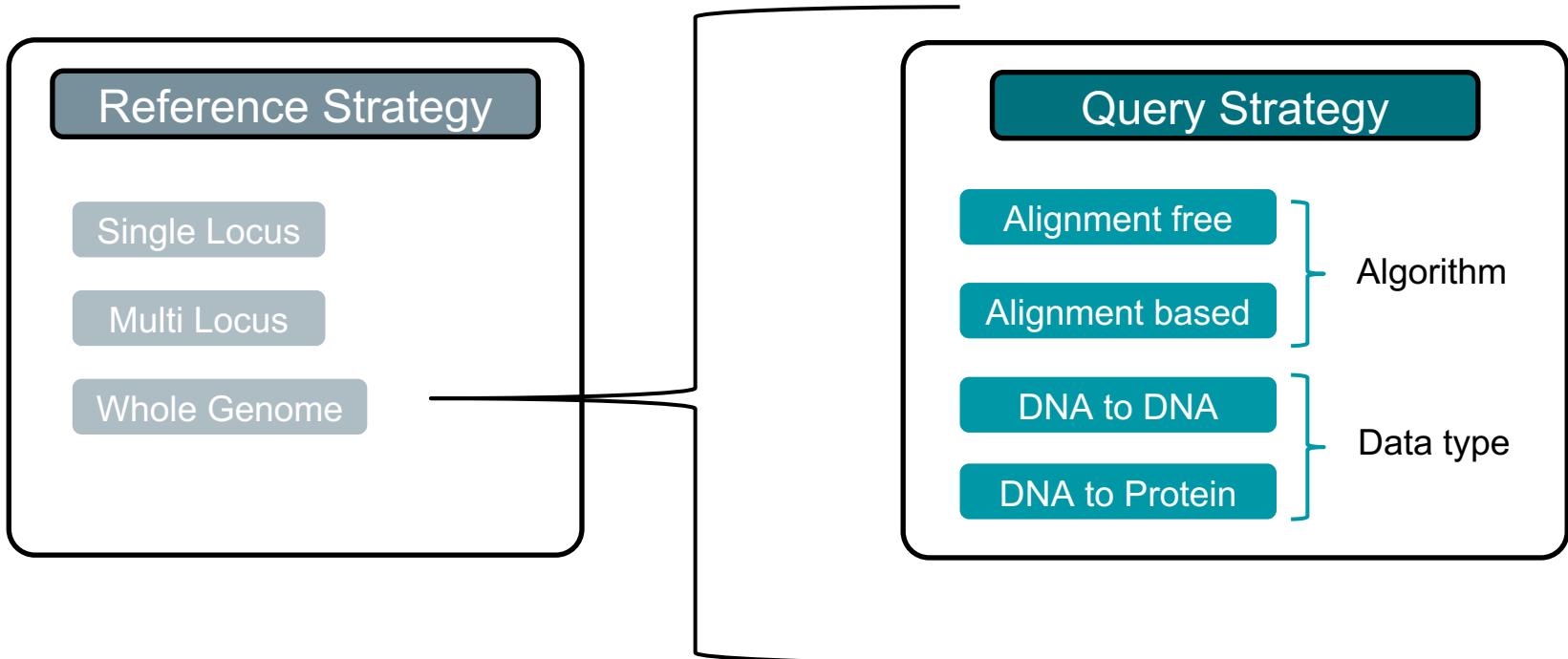
- 1 million marker genes
- ~17,000 reference genomes
- ~13,500 bacterial and archaeal
- ~3,500 viral
- ~110 eukaryotic
- Uses Bowtie2 for mapping

## Limitations

- Database not updated: good specificity but poor sensibility



Truong, Duy Tin, et al. "MetaPhiAn2 for enhanced metagenomic taxonomic profiling." *Nature methods* 12.10 (2015): 902.



# Reference strategy: whole genome

WGS

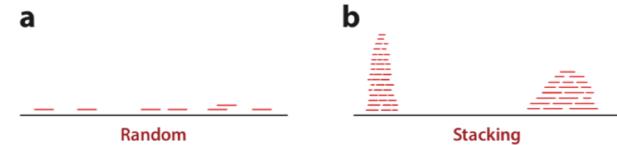


## Use entire genomes as reference database:

- Greatest sequence diversity
- Beneficial for aDNA when only traces of ancient organism
- Need heuristic to efficiently search the reference genome database

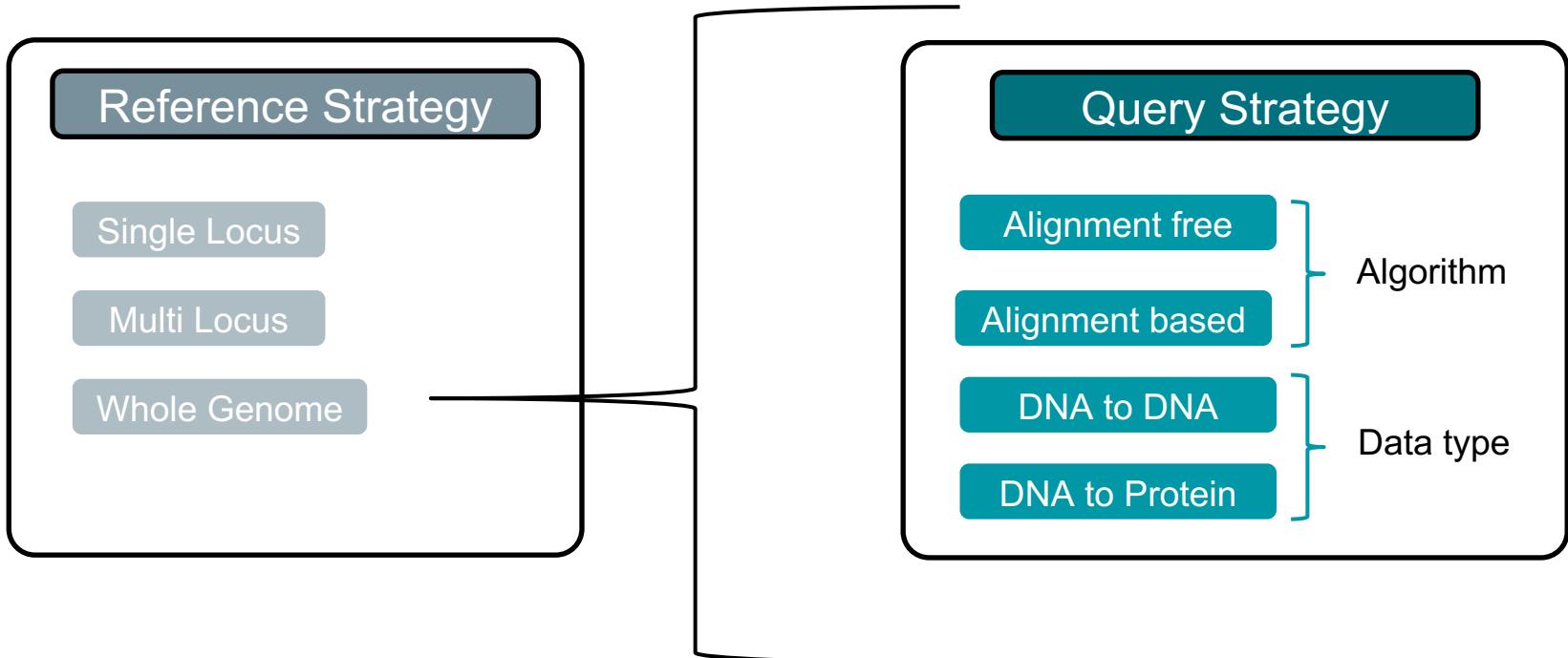
## Limitations:

- Horizontal gene transfer, mobile elements, recombination decrease precision (b)
- Variable bacterial genome sizes skew DNA proportion estimation
- Sparse database compared to 16s lead to more false assignments



More false negative (not in DB)

More false positives (database bias) because assigned to closest relative



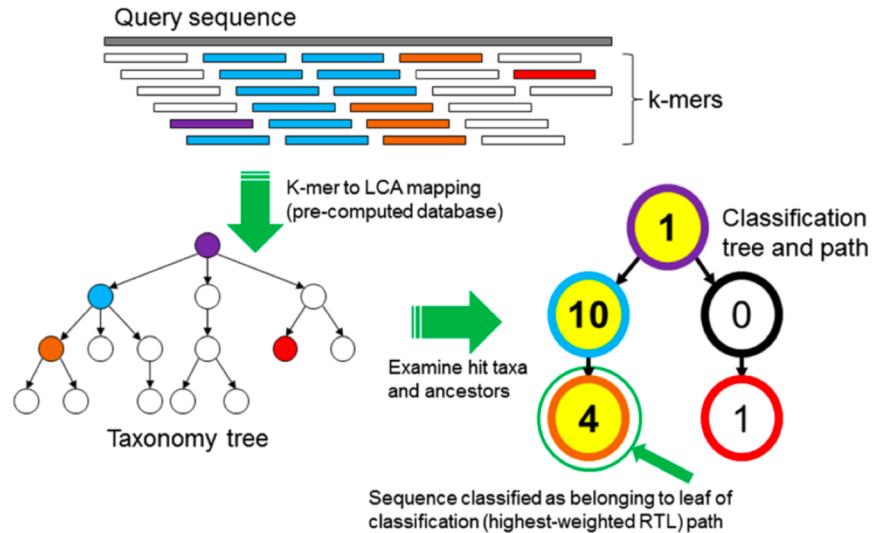


# Query strategy: Alignment free

Only look for exact matches between query and database

- Slice the reference genomes and the query into kmers of fixed size and save them into hash table
- Match kmers in query and reference
- Map the kmers to the taxonomy

kmer: substring of size k





# Alignment free methods

**Kraken**

Taxonomic Sequence Classification System

**Kraken 2**

Taxonomic Sequence Classification System

**Centrifuge**

Classifier for metagenomic sequences

**KrakenUniq: confident and fast metagenomics classification using unique k-mer counts**

**CLARK**

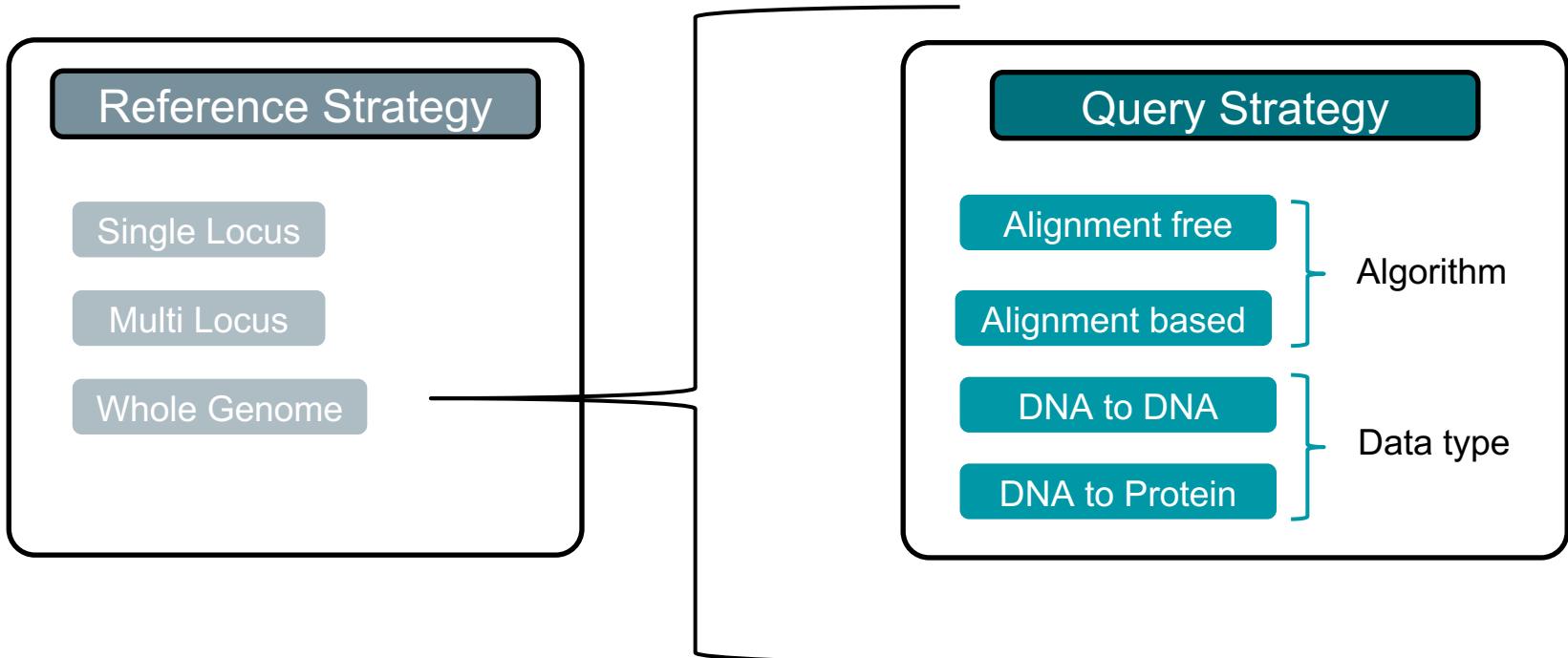
Fast, accurate and versatile sequence classification system



# Alignment based advantages and limitations

## Very quick

- Sensitivity and specificity depend on k:
  - Low k: more sensitive, less specific
  - High k: less sensitive, more specific
- No alignment: can't check for DNA damages or functions





# Query strategy: alignment based

## Local vs Global Alignments

**Local:** Start and end alignment at any location in sequence. Algorithm: *Smith-Waterman*

**Global:** Align every nucleotide in every sequence. Algorithm: *Needleman-Wunsch*

## Seed and Extend

1. Start with exact (or near exact) kmer matching – **Very fast**
2. Extend match on both sides if not too many mismatches then trigger alignment – **Slower**



# The Smith-Waterman algorithm



Temple Ferris Smith

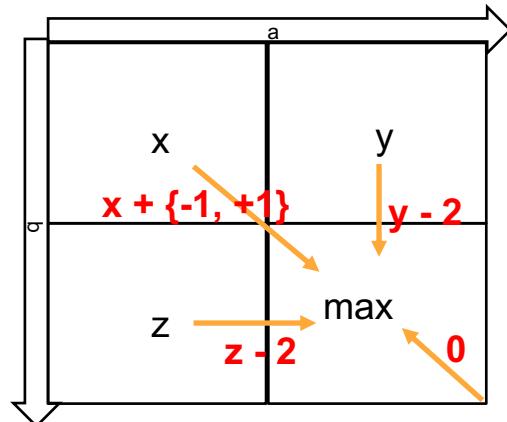


Michael Waterman

Smith, T. F., & Waterman, M. S. (1981). Identification of common molecular subsequences. *Journal of molecular biology*, 147(1), 195-197.

How to align these two sequences ? **ATGC**      **ATCC**

$$\begin{aligned}s_{\text{match}} &= +1 \\ s_{\text{mismatch}} &= -1 \\ s_{\text{gap}} &= -2\end{aligned}$$



Maxime Berry

	$A_{i=1}$	$T_{i=2}$	$G_{i=3}$	$C_{i=4}$
0				
$A_{j=1}$				
$T_{j=2}$				
$C_{j=3}$				
$C_{j=4}$				

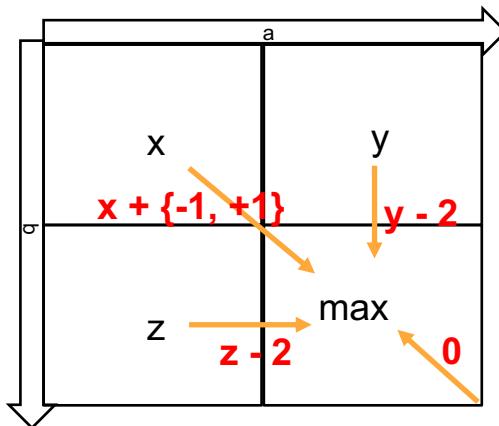


# The Smith-Waterman algorithm

## Scoring

$$\begin{aligned}s_{\text{match}} &= +1 \\ s_{\text{mismatch}} &= -1 \\ s_{\text{gap}} &= -2\end{aligned}$$

## Rule



## Matrix

		$A_{i=1}$	$T_{i=2}$	$G_{i=3}$	$C_{i=4}$
	0	0	0	0	0
$A_{j=1}$	0	1	0	0	0
$T_{j=2}$	0	0	2	0	0
$C_{j=3}$	0	0	0	1	0
$C_{j=4}$	0	0	0	0	2



# The Smith-Waterman algorithm

Traceback from (all) max scores, back to cell that gave maximum until reaching a 0

		$A_{i=1}$	$T_{i=2}$	$G_{i=3}$	$C_{i=4}$
	0	0	0	0	0
$A_{j=1}$	0	1	0	0	0
$T_{j=2}$	0	0	2	0	0
$C_{j=3}$	0	0	0	1	0
$C_{j=4}$	0	0	0	0	2

ATGC  
| | \* |  
ATCC



# Your turn !

		G	T	T	G	A	C
	0						
G							
T							
T							
A							
C							



# Solution

		G	T	T	G	A	C
	0	0	0	0	0	0	0
G	0	1	0	0	1	0	0
T	0	0	2	1	0	0	0
T	0	0	1	3	1	0	0
A	0	0	0	1	2	2	0
C	0	0	0	0	0	0	3

## Solution 1

GTT  
| | |  
GTT

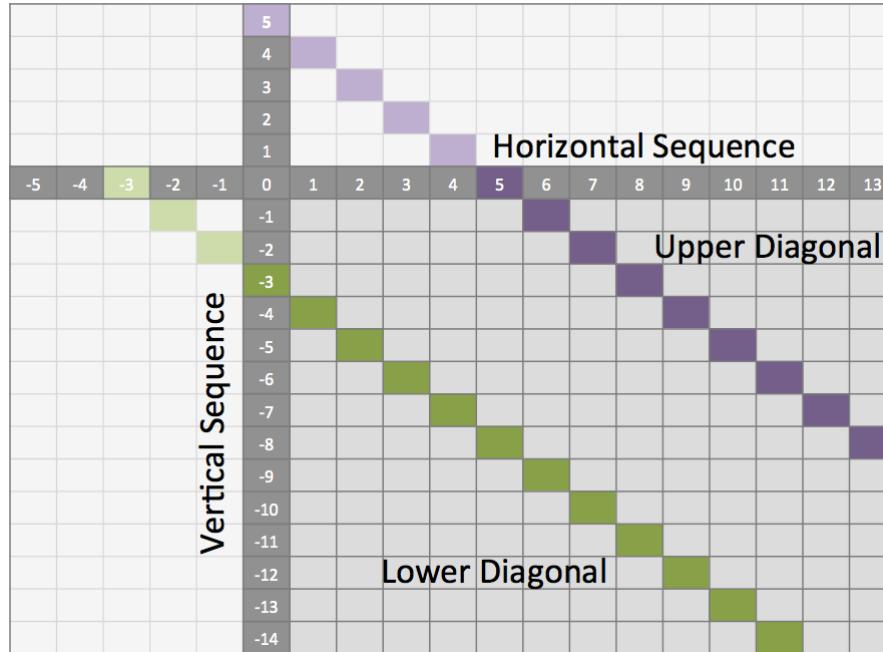
## Solution 2

GTTGAC  
| | | | |  
GTT-AC

Interactive Smith-Waterman:  
[rna.informatik.uni-freiburg.de/Teaching/index.jsp?toolName=Smith-Waterman](http://rna.informatik.uni-freiburg.de/Teaching/index.jsp?toolName=Smith-Waterman)



# Optimization: banded alignments



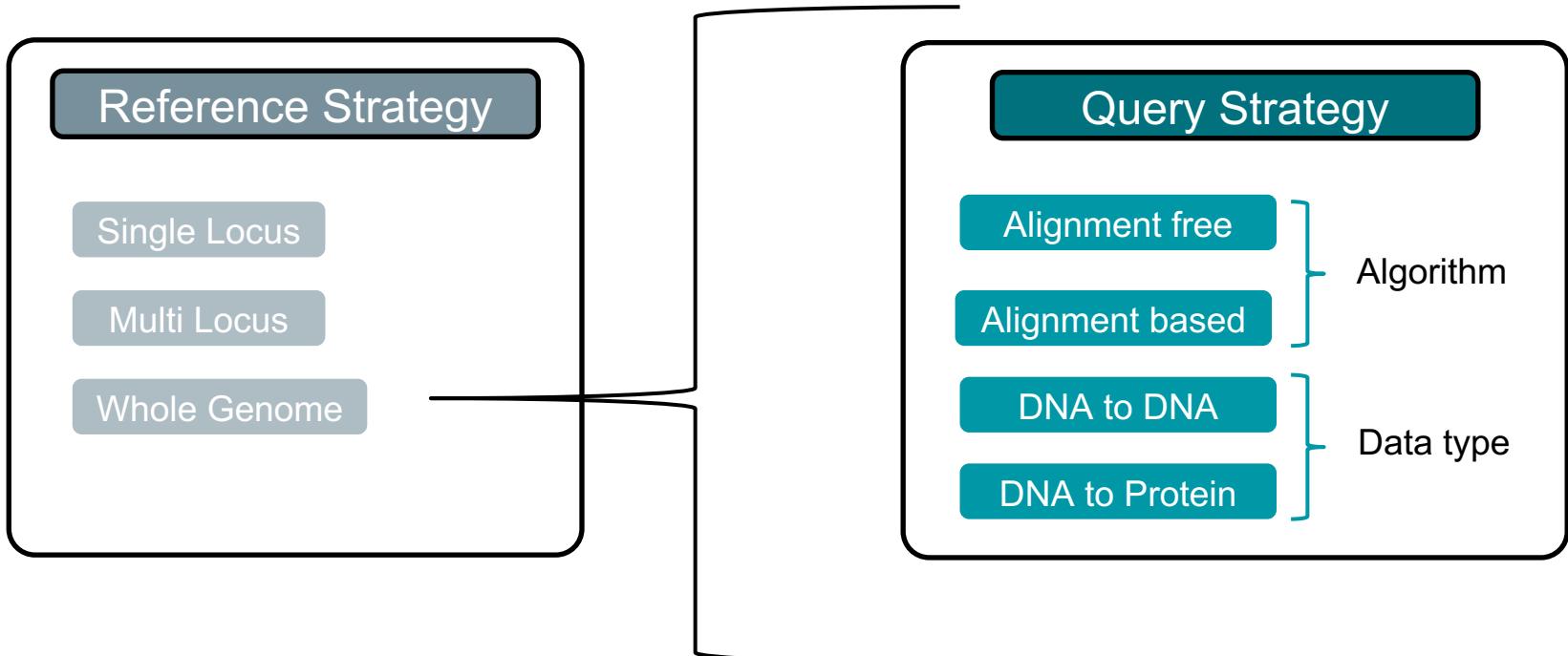


# Alignment based: advantages and limitations

## Alignment:

- Check for DNA damages
- Check for functions

Much (much) slower



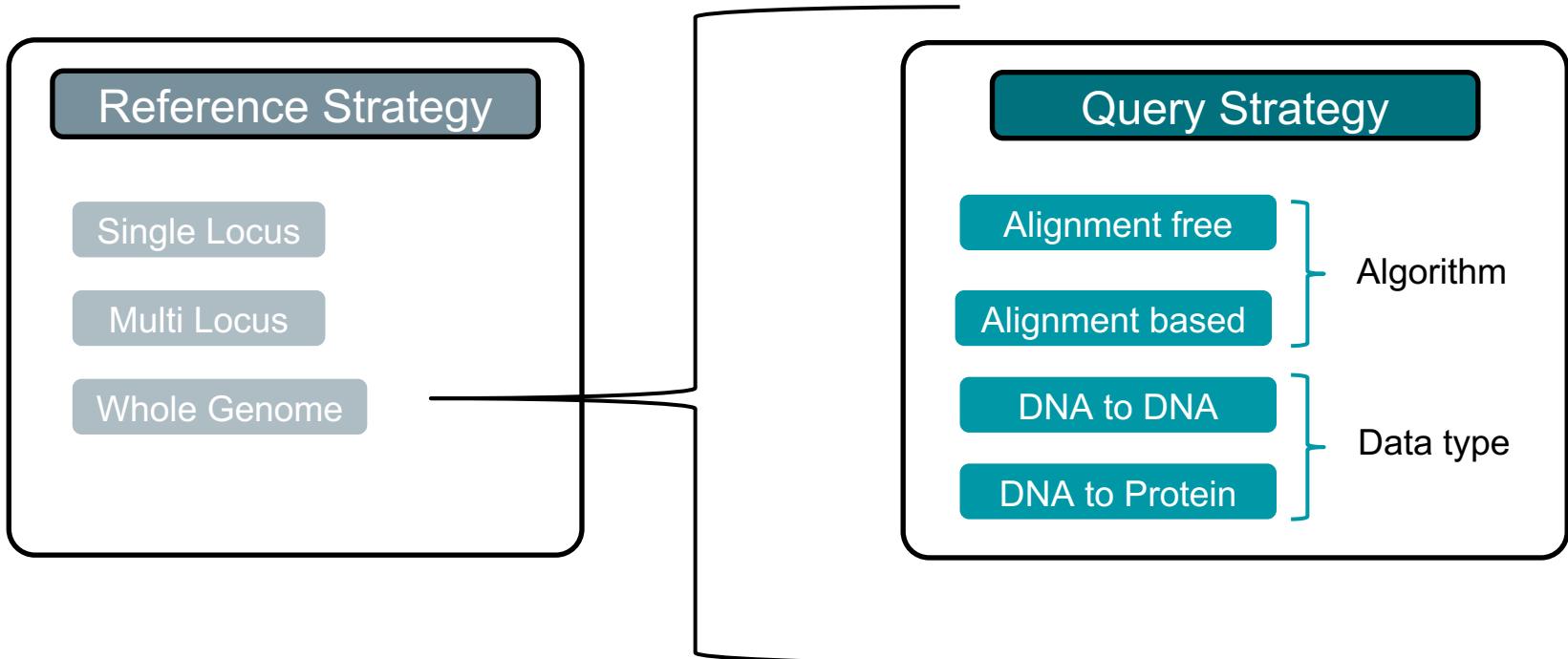


# DNA based alignment

## MALT: MEGAN ALignment Tool:

- Near exact seeds (spaced seeds)
- Mismatch tolerant extension (x-drop)
- Smith-Waterman (semi-global: read is aligned end-to-end)
- LCA





# Protein based alignment



<b>A</b>	+1			
<b>T</b>	-1	+1		
<b>G</b>	-1	-1	+1	
<b>C</b>	-1	-1	-1	+1
	<b>A</b>	<b>T</b>	<b>G</b>	<b>C</b>

**→**

Ala	4
Arg	-1 5
Asn	-2 0 6
Asp	-2 -2 1 6
Cys	0 -3 -3 -3 9
Gln	-1 1 0 0 -3 5
Glu	-1 0 0 2 -4 2 5
Gly	0 -2 0 -1 -3 -2 -2 6
His	-2 0 1 -1 -3 0 0 -2 8
Ile	-1 -3 -3 -3 -1 -3 -3 -4 -3 4
Leu	-1 -2 -3 -4 -1 -2 -3 -4 -3 2 4
Lys	-1 2 0 -1 -3 1 1 -2 -1 -3 -2 5
Met	-1 -1 -2 -3 -1 0 -2 -3 -2 1 2 -1 5
Phe	-2 -3 -3 -3 -2 -3 -3 -1 0 0 -3 0 6
Pro	-1 -2 -2 -1 -3 -1 -1 -2 -2 -3 -3 -1 -2 -4 7
Ser	1 -1 1 0 -1 0 0 0 -1 -2 -2 0 -1 -2 -1 4
Thr	0 -1 0 -1 -1 -1 -2 -2 -1 -1 -1 -1 -2 -1 1 5
Trp	-3 -3 -4 -4 -2 -2 -3 -2 -2 -3 -2 -3 -1 1 -4 -3 -2 11
Tyr	-2 -2 -2 -3 -2 -1 -2 -3 2 -1 -1 -2 -1 3 -3 -2 -2 2 7
Val	0 -3 -3 -3 -1 -2 -2 -3 -3 3 1 -2 1 -1 -2 -2 0 -3 -1 4
Ala Arg Asn Asp Cys Gln Glu Gly His Ile Leu Lys Met Phe Pro Ser Thr Trp Tyr Val	

Same strategy as DNA based alignment, but using Amino acid substitution matrix.

## BLOSUM: BLOcks SUbstitution Matrix

**See also:** Optimization with reduced alphabet (grouping amino acids by chemical properties)



# Protein based alignment tools

## **MALT:**

Protein to Protein Database  
DNA to Protein Database

## **Diamond:**

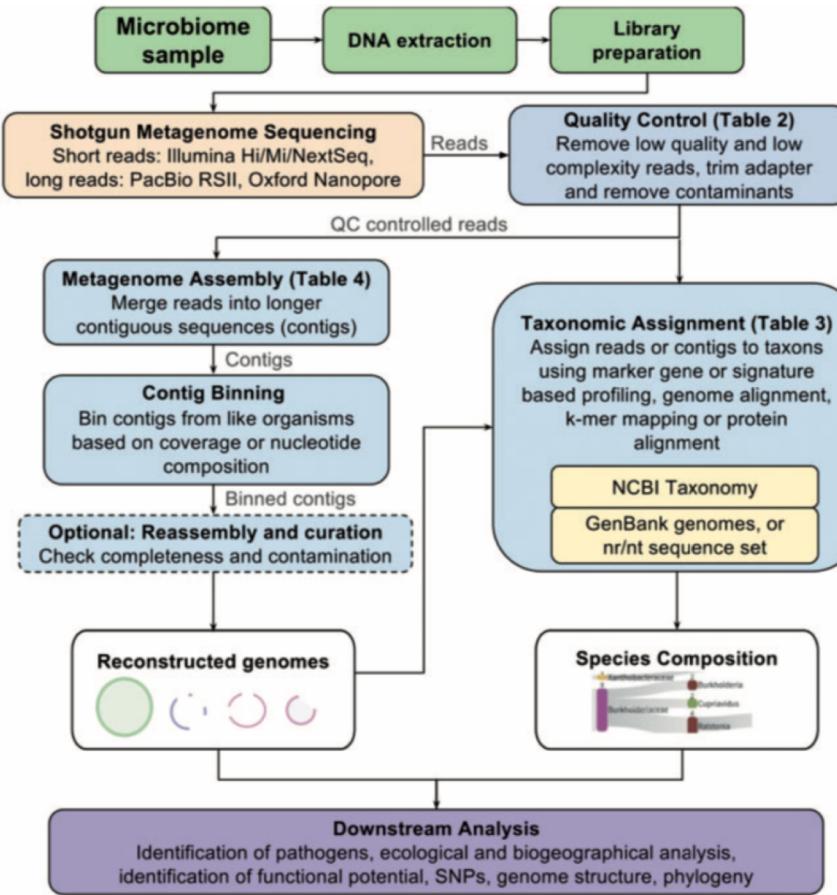
Protein to Protein Database  
DNA to Protein Database  
No LCA (Megan)



No DNA damage estimation with protein alignment



# Other strategies

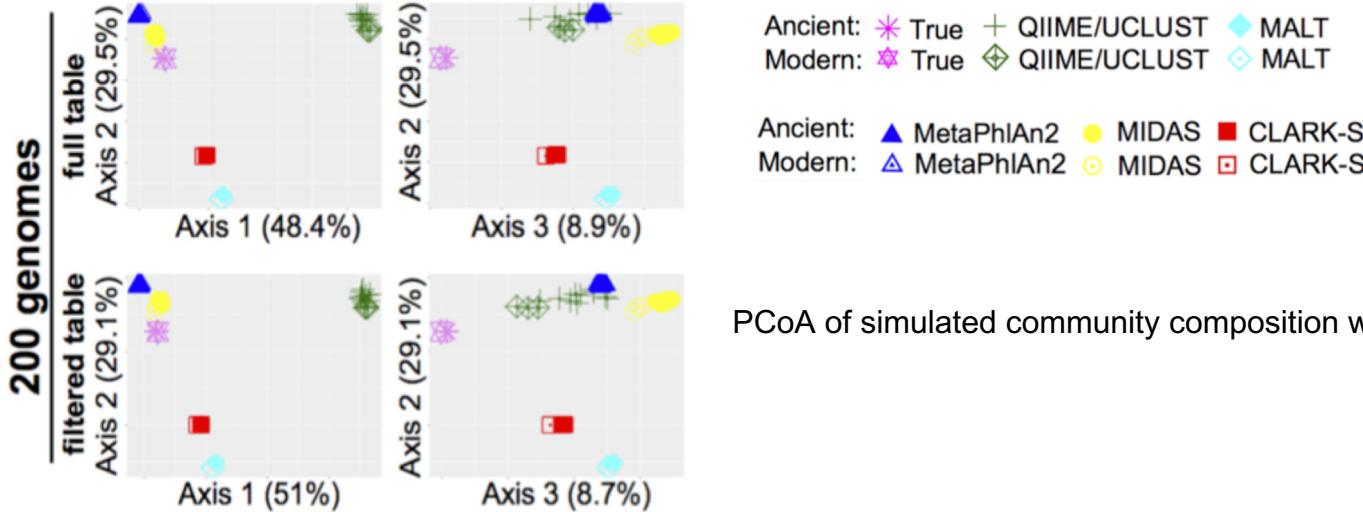




# Taxonomic classifiers for aDNA



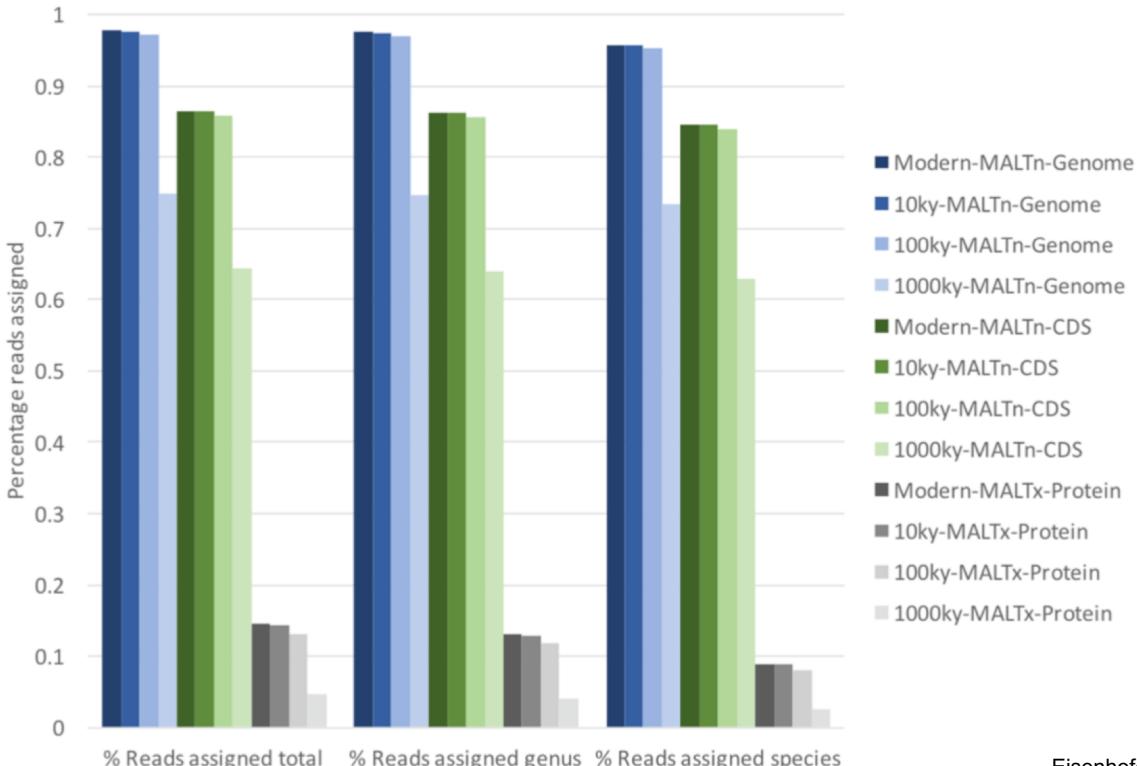
# Damage isn't really an issue



aDNA damage doesn't really affect inferred community composition



# The choice of database

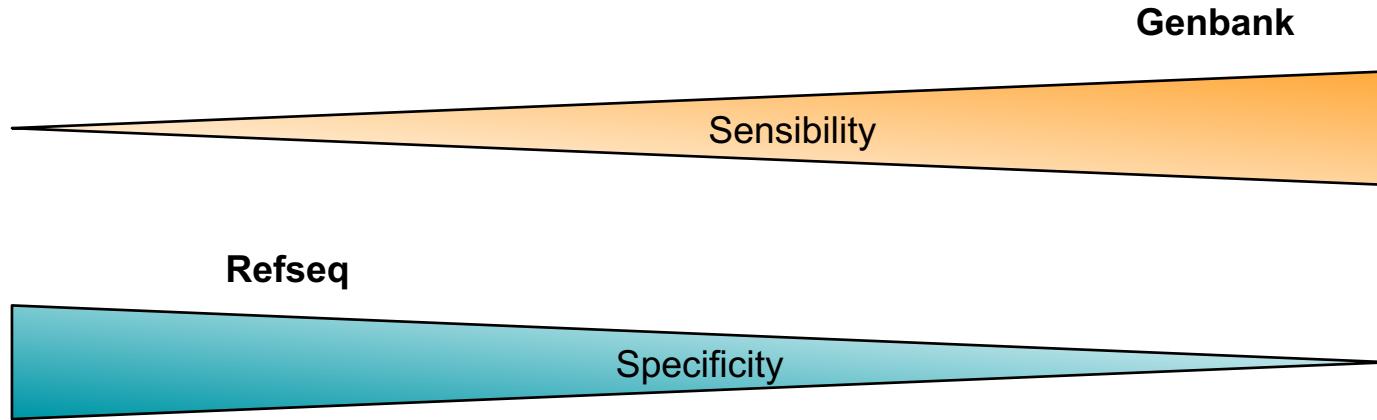


- DNA database is better for short aDNA (<60 bp) than protein databases
- Translation of shorter sequences gives too small peptides to align
- **Database composition matters**

Eisenhofer, R., & Weyrich, L. S. (2019). Assessing alignment-based taxonomic classification of ancient microbial DNA. *PeerJ*, 7, e6594.



# Genbank vs Refseq





# The definition of species

- **Taxonomic classifiers work with Taxonomy**
- Wanna be a species ? Better be culturable !
- Not culturable but otherwise well characterized : included as *Candidatus Genus species*
- Otherwise: **unclassified bacteria**

Lineage: Cellular Organisms; Bacteria; unclassified bacteria

-> Problem with LCA

# Species vs ANI



15%

- ANI: Average Nucleotide Identity
- **Gold standard:** 95% ANI for species
- 15% of Genomes in Genbank from the same species have ANI < 93%



# Choosing the right tool

First define the question:

Am I only asking **who is there** ?

Yes: Locus based or Alignment free method

No: Whole genome DNA alignment based method

Do I have a lot of computing time to spare (and/or not so many samples at once) ?:

Yes: Whole genome DNA alignment based method

No: Locus based or Alignment free method

Do I have WGS sequencing data ?

Yes: Prefer avoiding single locus (open debate)

No: Single locus it is

Do I have long fragments and feel adventurous ?

Yes: Try whole genome protein alignment based method, or even Assembly

No: Stick to the rest



Worst protein alignment ever (score = 2)