# Bacterial (Microbial) Genomics

Introduction to Computational & Quantitative Biology

Anne-Catrin Uhlemann, MD, PhD
Department of Medicine / Division of Infectious Diseases

November 15th 2022

Metazoa

Archaea

Bacteria

Excavata
Amoebozoa
SAR

Archaeplastida
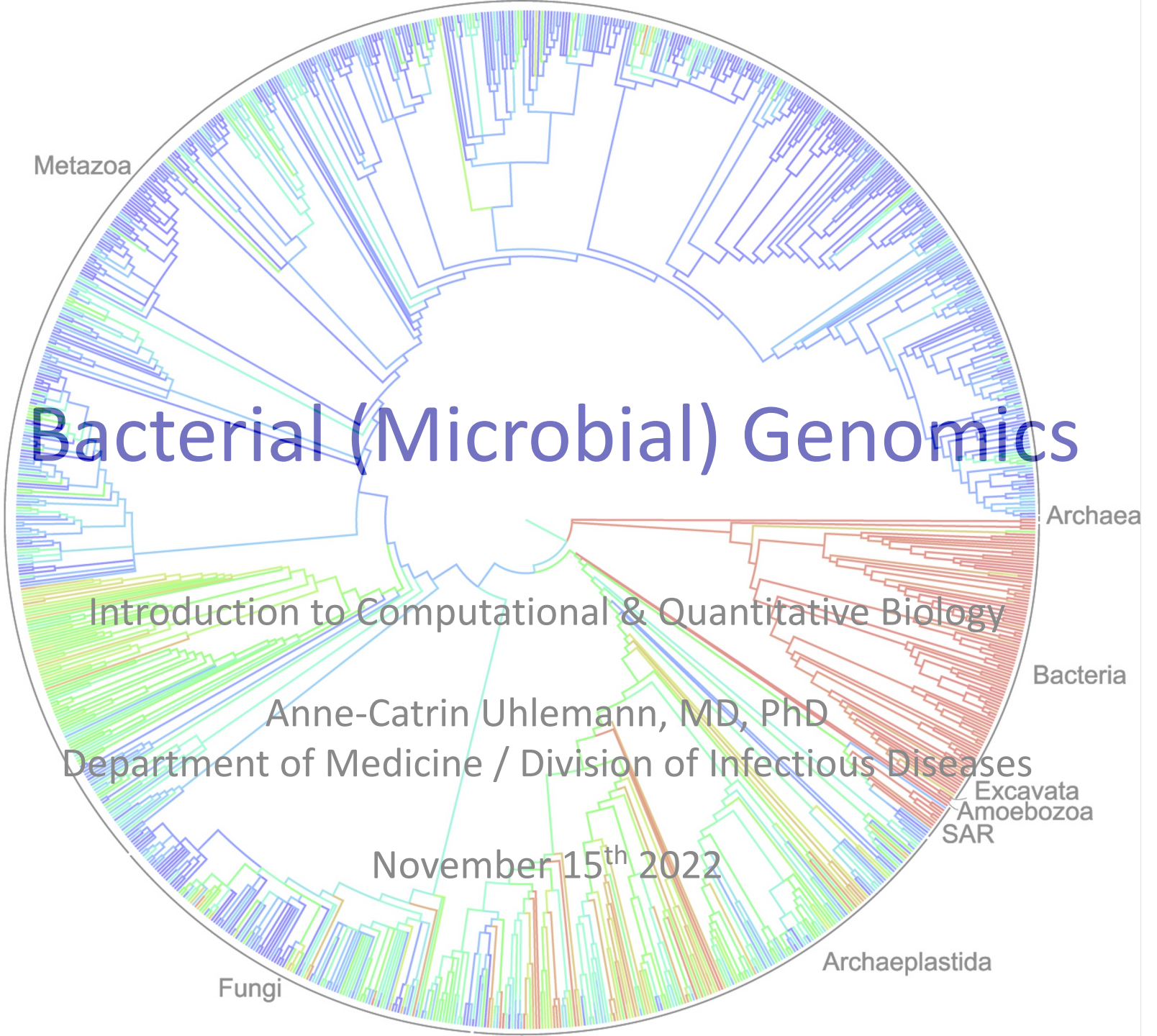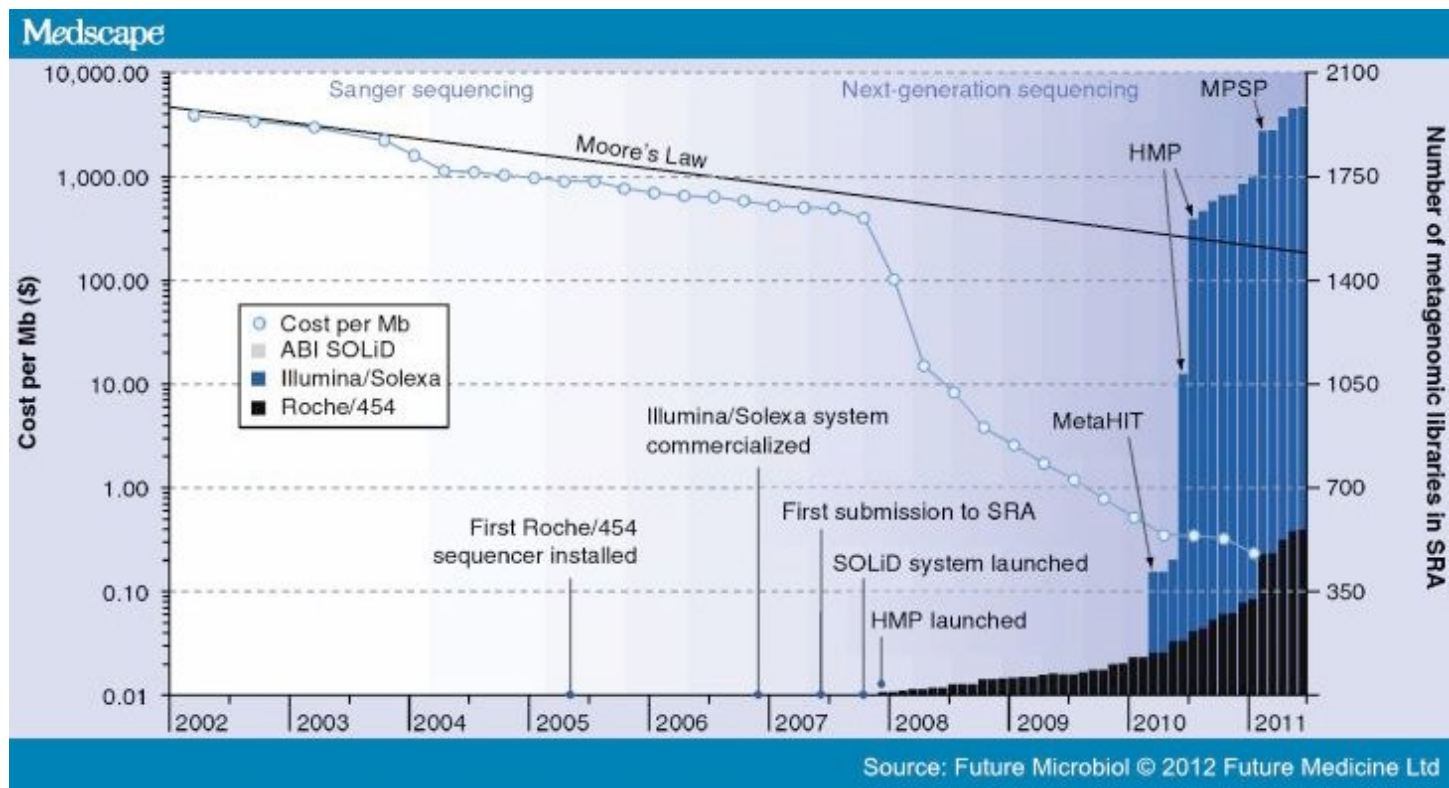
Fungi

# Overview

- Sequencing technologies
- Analysis approaches

  - Bacterial whole genome sequencing

  - Microbiome analyses

- Some practical examples
- Maybe a few words on COVID-19 at the end

# Recap: Sanger sequencing

- Dideoxynucleotide sequencing
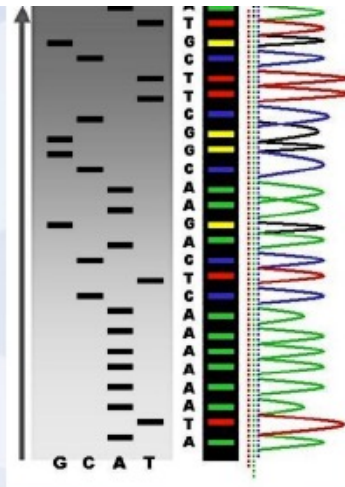- Chain-termination method

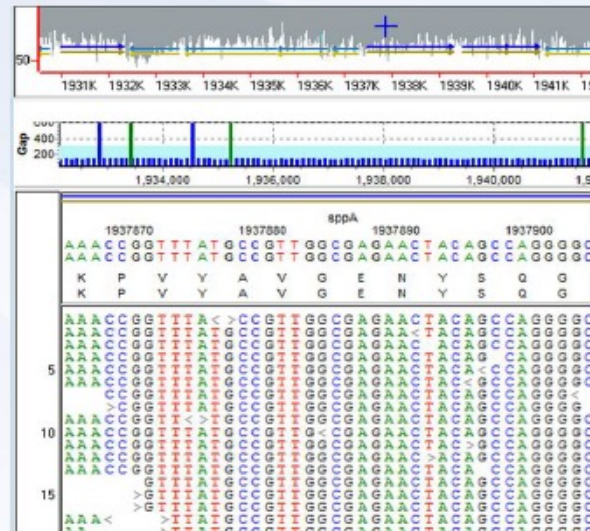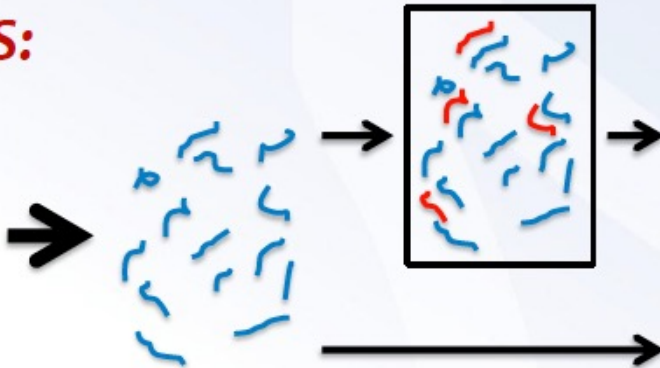# Microbial genomics – byproducts of the race for the human genome

# Traditional versus
# Next-Generation Sequencing

# High-throughput Next generation sequencing by synthesis



Genomic DNA → Cut DNA → Add Linkers

Input library → Flow cell → In Situ PCR → Sequencing → An image of hundreds of extended molecules

blood

# Read length



Ex. Illumina's sequencing by synthesis

# of sequencing cycles = Read length

Reads

Chemistries limit read length, are constantly being improved
- short         < 50 consecutive bases
- mid-length   51 - <400
- long          > 400 (< 1000)

# Depth of coverage



Numbers can be misleading!

# Some terminology

- FASTQ: text-based format storing sequence data and quality scores

- FASTA file: sequence in text format

- SAM file: tab-delimited text file that contains sequence alignment data

- BAM file: binary version of a SAM file

# Bacterial genomics overview

# Bacterial sequencing applications

16S rRNA sequencing        Metagenomics        Single cell sequencing

# Single isolate bacterial sequencing

- Comparative sequencing
  - SNPs / indels that determine virulence
  - evolution
  - outbreak investigations
  - compare presence / absence mobile elements
- *De novo* assembly
  - only way to determine new gene content
  - not always optimal for variant calling

**Figure 1. An example workflow for high-throughput whole genome sequencing in bacteria.**

# Typing schemes of bacterial DNA

Bacterial DNA

Plasmids



Multi-locus sequence typing (housekeeping genes)



..GCTTG..   2
..TAGGC..   3
..ATGCG..   1
..CGCTG..   1
..TGATC..   4
..TAAGG..   4
..ACTGA..   3

ST239

Plasmid PCR typing
- *rep* gene
- *bla*$_{KPC}$ gene mutations

MLST reflects clonal type quite well

# Workflows for comparative analyses

# K-mers

- Sequence of K base calls (DNA that is k long)
    - ATGC = 4-mer, ATGCTG = 6-mer
    - all of a sequence's subsequences of length
- Only consecutive bases are used
- Reads with high sequence similarity must share K-mers in overlapping regions
- Shared K-mers are easier to find then overlaps
- Fast detection of shared K-mer content reduces computational cost / time
- Disadvantage: lower sensitivity in overlap regions

# De Bruijn graph assembly

- AACCGGTTA

- GGTTATAC

AACC

ACCG

CCGG

CGGT

GGTT

GTTA

GGTT

GTTA

TTAT

TATA

ATAC

AACCGGTTATAC

Spades: uses k=31 ->127

**a**

**b**

CGTGCAA

TGCAATG

ATGGCGT

GGCGTGC

CAATGGC

ATGGCGT
GGCGTGC
CGTGCAA
TGCAATG
CAATGGC
ATGGCGT

Genome: ATGGCGTGCAATGGCGT

Short-read sequencing

Vertices are *k*-mers
Edges are pairwise alignments

Vertices are (*k*–1)-mers
Edges are *k*-mers

**c**

ATG
TGG
GGC
GCG
CGT
GTG
TGC
GCA
CAA
AAT

**Hamiltonian cycle**
Visit each vertex once
(harder to solve)

*k*-mers from vertices

Genome: ATGGCGTGCAATG

*k*-mers from edges

ATG
TGG
GGC
GCG
CGT
GTG
TGC
GCA
CAA
AAT
ATG

**d**

AA
AT
TG
GG
GC
CG
GT
CA

AAT
ATG
TGG
GGC
GCG
CGT
GCA
GTG
TGC
CAA

**Eulerian cycle**
Visit each edge once
(easier to solve)

Compeau et al. Nat Biotechnol 2011

# Bacterial sequencing application: Evolution of *S. aureus* USA300

# California: Initial documentation of CA-MRSA epidemic / USA300



CA-MRSA:
- No hospitalization past 6 months
- Not nursing home
- Culture + < 48hrs
- Not dialysis
- Not homeless
- Skin & Soft tissue infections
- Invasive ~5% cases

# USA300 genome composition



**Genomic & biological features:**

Core genome
- Increased expression of core virulence genes PSM, $\alpha$-toxin

Mobile genetic elements
- Small SCC*mec* IV
- φ2 / PVL toxin
- SaPi5
- ACME I

    detoxifies host anti-microbials
- derived from USA500?

Diep et al. The Lancet 2005

# Whole genome sequencing
# 387 isolates



**Sequencing:**
- Large dataset of infectious and household isolates
- Mate-paired libraries 100 bp paired-end
- Illumina Hi-Seq
- Coverage 100 to 170 fold

**Mapping**:
- Reference genome FPR3757
- Exclusion unmapped reads, MGEs
- Repeat Scout
- SNP calling

**Phylogenetic tree**
- Core genome
- Concatenated SNPs
- RAxML

# Phylogenetic tree

## All life is related by common ancestry.



Phylogeny:  pattern of historical relationships

Tree: mathematical structure used to model the evolutionary history  of a group of organisms

# Tree Notation

# Genome composition of *S. aureus*
# Mobile elements do not follow tree like evolution

Chromosomal Genome:
1. Stable core
   - MLST

2. Variable core
- Surface proteins
    spa-types
- Some virulence factors

3. Mobile genetic elements
- Integrated Pro-phages
- Pathogenicity islands
- Transposons and Insertion sequences

Plasmids:
- Resistance genes

# SNP calling (dataset n = 375)

- Mapped to reference

- Mapped genome ~90%

- MGEs excluded (need to analyze separately)

- 12,451 SNPs

- Coverage 3-fold per SNP base needed

# Validation of SNP calls: Comparison of multiple isolates per person

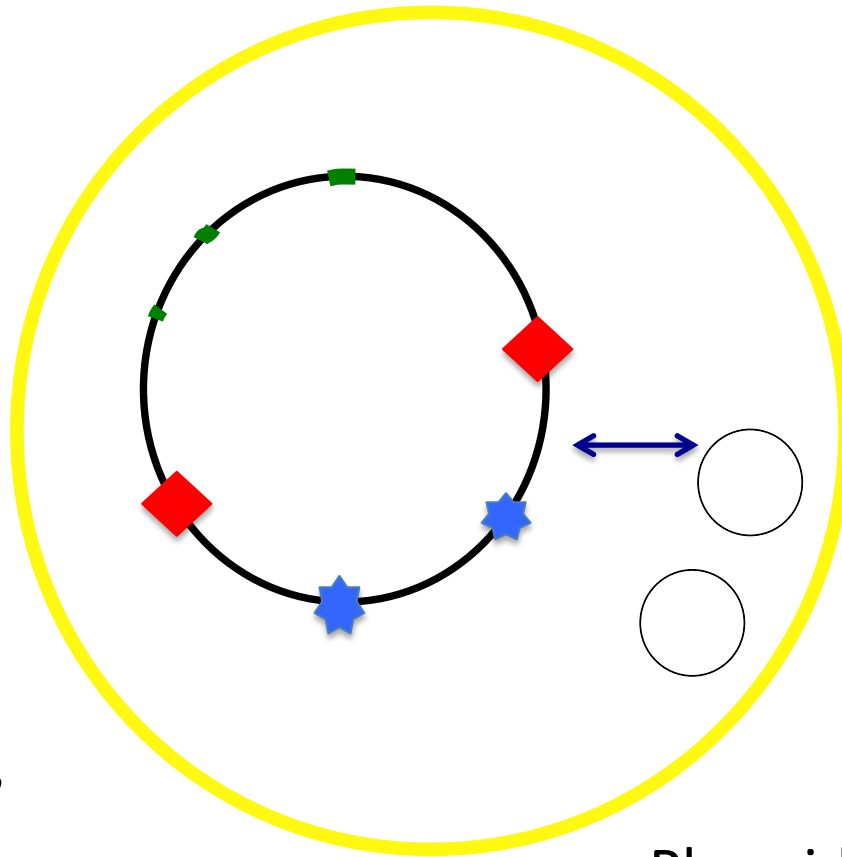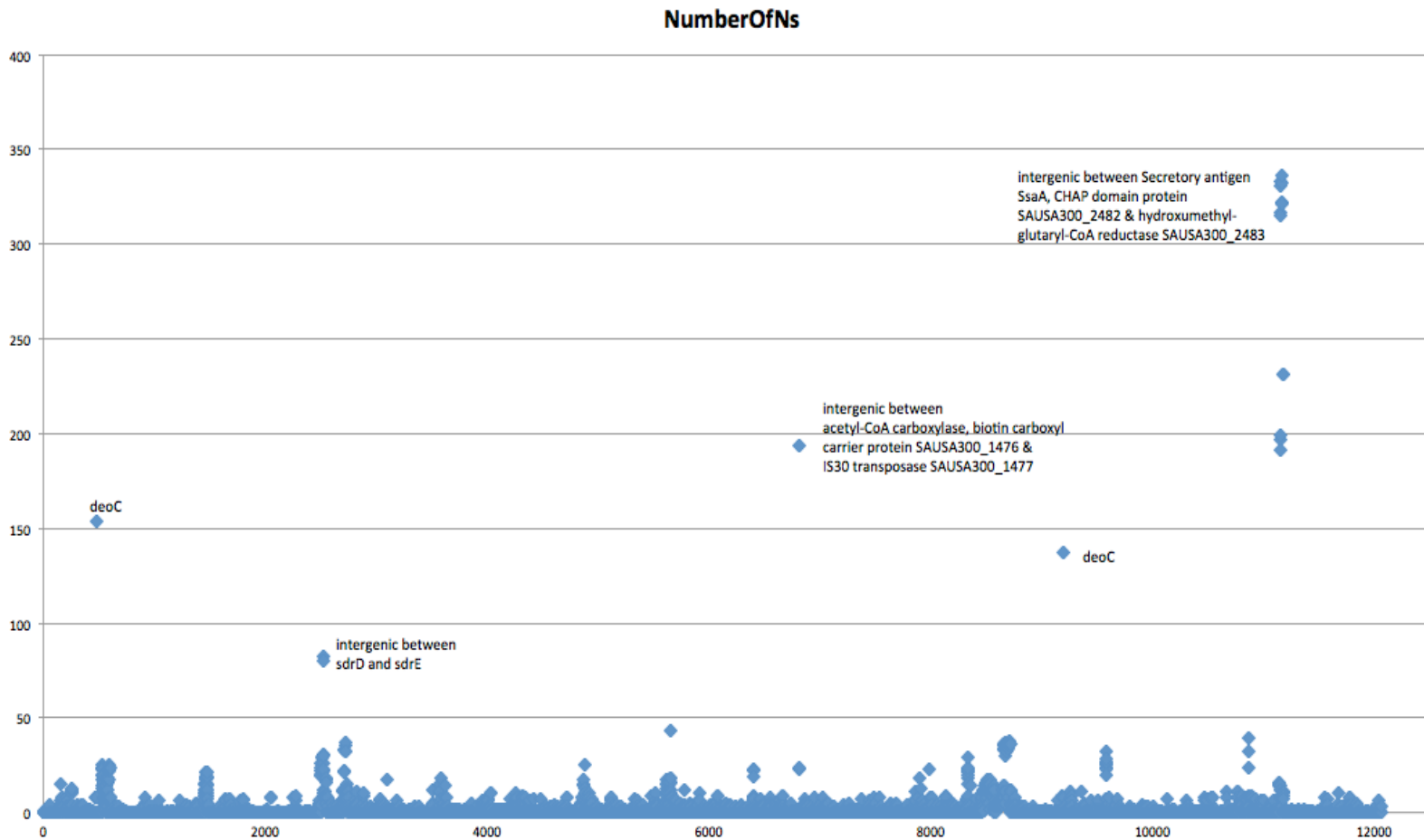| Position_in_ | CDS/rRNA/tR | CDS_name | product | Synonymous | Ref_base | SNP_base | Total | 7748_4#94 | 9266_1#14 | 9266_1#15 | 7748_4#95 | 9266_1#16 | 7790_7#31 | 7748_4#96 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | | | | | 2079_1 | 2079_1_ob1 | 2079_1_ob2 | 2079_2 | 2079_2_ob | 2079_3 | 2079e |
| 108951 | Intergenic | - | - | - | G | T | 6 | T | T | T | T | N | T | T |
| 349157 | CDS | SAUSA300_0297 | putative lipoprotein | S | C | T | 6 | T | T | T | T | T | N | T |
| 349159 | CDS | SAUSA300_0297 | putative lipoprotein | N | C | T | 6 | T | T | T | T | T | N | T |
| 1074221 | Intergenic | - | - | - | G | T | 21 | T | N | T | T | T | T | T |
| 2217199 | Intergenic | - | - | - | C | T | 1 | T | . | . | . | . | . | . |
| 110608 | CDS | SAUSA300_0100 | staphylococcal tandem lipopr | N | C | T | 5 | N | T | T | T | T | N | T |
| 110610 | CDS | SAUSA300_0100 | staphylococcal tandem lipopr | S | A | T | 5 | N | T | T | T | T | N | T |
| 110613 | CDS | SAUSA300_0100 | staphylococcal tandem lipopr | S | A | T | 5 | N | T | T | T | T | N | T |
| 110709 | CDS | SAUSA300_0100 | staphylococcal tandem lipopr | S | C | T | 3 | N | T | N | N | T | N | N |
| 2031077 | Intergenic | - | - | - | T | G | 72 | N | N | N | G | G | G | G |
| 349259 | CDS | SAUSA300_0297 | putative lipoprotein | S | T | C | 2 | N | N | C | N | C | N | N |
| 110614 | CDS | SAUSA300_0100 | staphylococcal tandem lipopr | N | A | G | 5 | N | G | G | G | G | N | G |
| 110616 | CDS | SAUSA300_0100 | staphylococcal tandem lipopr | S | A | G | 5 | N | G | G | G | G | N | G |
| 110620 | CDS | SAUSA300_0100 | staphylococcal tandem lipopr | N | A | G | 5 | N | G | G | G | G | N | G |
| 110634 | CDS | SAUSA300_0100 | staphylococcal tandem lipopr | N | A | C | 3 | N | C | C | N | C | N | N |
| 110718 | CDS | SAUSA300_0100 | staphylococcal tandem lipopr | S | T | A | 3 | N | A | N | N | A | N | N |
| 110628 | CDS | SAUSA300_0100 | staphylococcal tandem lipopr | N | T | A | 3 | N | A | A | N | A | N | N |
| 349265 | CDS | SAUSA300_0297 | putative lipoprotein | S | G | A | 3 | N | A | A | N | A | N | N |
| 2179424 | Intergenic | - | - | - | G | A | 40 | N | A | A | N | A | N | N |
| 1630757 | Intergenic | - | - | - | T | A | 215 | A | N | N | N | N | N | A |
| 1074250 | Intergenic | - | - | - | G | A | 21 | A | A | A | A | A | N | A |
| 960514 | Intergenic | - | - | - | C | A | 13 | A | A | A | A | N | A | A |
| 1433531 | CDS | SAUSA300_1302 | ATPase family protein | 2 | G | A | 1 | . | . | A | . | . | . | . |

Concern: high number of SNPs in isolates samples at the same time from same person based on "N" (i.e. inability to call sequence) in "non-mutant"

Suspicious: clustering in one gene/region

# Distribution of N's across SNPs



**NumberOfNs**

# Possible explanations for N's

- Repetitive sequences
- Does not overtly match deletions/insertions
- Difficulty of mapping repeats (read length only 100bp – disadvantage of Illumina)
- Duplications or recombination?
- Identified new ORF in "intergenic region" transposase type (published in other genomes)

# How are N's addressed in the literature?

- Usually no mention!
- Between the lines:

"Unmapped reads and sequences that were not present in all genomes were not considered as part of the core genome, and therefore SNPs from these regions were not included in the analysis… as were SNPs falling in high-density SNP regions, which could have arisen by recombination. The core genome was curated manually to ensure a high-quality data set…"

# SNP matrices, distances and trees

Multiple alignment

1  AGGCCAAGCCATAGCTGTCC
2  AGGCAAAGACATACCTGACC
3  AGGCCAAGACATAGCTGTCC
4  AGGCAAAGACATACCTGTCC

Distance matrix

|   | 1 | 2 | 3 | 4 |
|---|---|------|------|------|
| 1 | – | 0.20 | 0.05 | 0.15 |
| 2 |   | –    | 0.15 | 0.05 |
| 3 |   |      | –    | 0.10 |
| 4 |   |      |      | –    |

- Once we compute the distances, how do we find a good tree?

- There are several methods.

# Trees are like mobiles

- The same tree can be represented in different ways, by permuting the branches.

# Different trees

- Alignment of homologous sequences

  -> concatenated SNPs

- Topology (no lengths):
  - Cladogram: relative common ancestry without specifying lengths.

- Topology + lengths:
  - Additive trees: incorporate the length of the branch representing the amount of evolutionary change.

NJ 2271 sites J-C 100 repl.

A/shorebird/Delaware/246/2006H1N1
A/NewYork/4290/2009H1N1
A/swine/Ohio/23/1935H1N1
A/swine/Tennessee/7/1976H1N1
A/NewJersey/8/1976H1N1
A/Brevig
A/Melbourne/35H1N1
A/HongKong/117/77H1N1
A/SouthAustralia/45/2000H1N1
A/Denmark/16/04H1N1
A/Malaysia/54H1N1

# Methods of constructing trees

1. Distance methods
   - Minimal Evolution
   - Least Squares
   - UPGMA
   - Neighbor-Joining
2. Parsimony
3. Likelihood. PHYLIP (Felsenstein)
4. Bayesian methods

# Maximum likelihood estimation

- Principle: Choose the tree which makes the data most probable
- Each position evolves independently
- Accommodates time structure of temporally-spaced sequences
- Tips have isolation date; internal nodes are unknown -> arbitrary starting times (order on tree)
- Substitution rate used to scale times into units of expected number of substitutions per site
- Likelihood of the model; standard multi-dimensional optimization ->maximum likelihood

# Maximum likelihood (II)

- Allows hypothesis testing and model comparison via likelihood ratio test

- Test if one hypothesis provides better fit (nested hypotheses)

- Problem: can be time / computationally intensive

# Bootstrapping

- How much do we trust a tree that we have constructed?

- A simple method for parsimony, distance or ML is bootstrapping.
  - Select some random positions, with repetition.
  - Construct another tree with the bootstrapped data.
  - Repeat many times.
  - Check the consistency of the results.

- In Bayesian methods, one can estimate the confidence by looking at posterior probabilities.

# Phylogeny of ST8 and the emergence of USA300



**Phylogenetic tree:**
- 433 isolates
- Additional isolates
  2005 / 2006 study
  California SSTIs
- 12,212 SNPs in core genome
- Maximum likelihood tree with 1000 bootstraps
- Homoplasy index 0.007

Cases
Case contact
Control
California isolates

Rooted to midpoint

# Mobile genome analysis - PVL in USA300 core lineage differs from other ST8
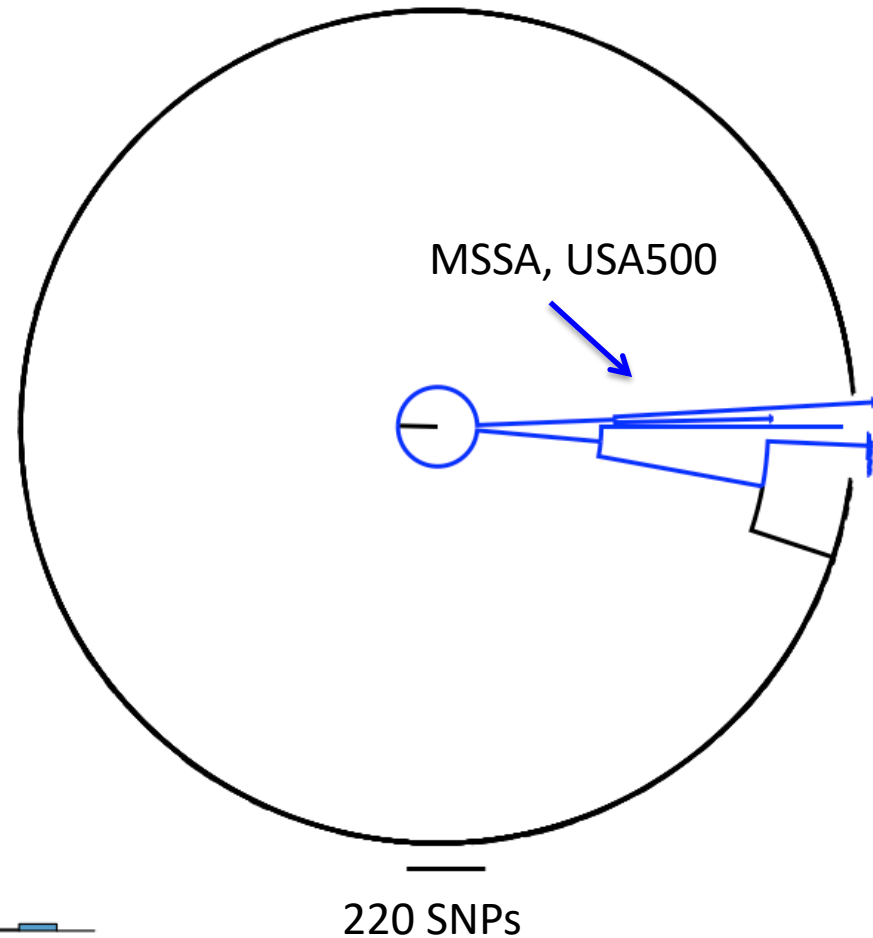


Distribution of pairwise distances

MSSA, USA500

220 SNPs

# Mapping of mobile genetic elements matches core phylogeny



A

ACME & SaPI5
- red: present
- black: absent

SCCmec
- blue: IVa
- lime: IVb
- green: IVc
- orange: IVg
- black: MSSA

PVL-prophage (Sa2int)
- orange: subtype 1 (USA500)
- green: subtype 2
- magenta: subtype 3
- light blue: subtype 4
- red: subtype 5
- black: absent

~ 53 SNPs

B  USA300 FPR3757* (SCCmecIVa SaPI5+ ACME+)

SCCmec    ACME

SCCmec    ACME remnant (70 bp)

7748_4#59* (SCCmec IVa SaPI5+ ACME-)

# What determines "strain similarity"

- Substitution rate

    - Root-to-tip analysis

    - Bayesian reconstruction (subset of isolates)

- Pairwise SNP distance

# Root-to-tip linear regression

- First estimates rooted phylogeny

  - matrix pairwise genetic distance using empiric model of substitution

  - matrix used for neighbor-joining tree

- Second linear regression between time of sampling of each tip and genetic distance (sum of reconstructed branch length)

  $$E[d_{root,i}] = m(t_i - t_{root}) = mt_i - mt_{root}$$

- Root of tree picked to maximize $R^2$ value of regression

- Advantage: fast visualization

- Not the final model!

# Root-to-tip analysis to estimate date of ancestry



Correlation coefficient 0.4853
$R^2$: 0.2355

Substitution rate/site/year: $1.56 \times 10^{-6}$
Time most common recent ancestor: ~1995

Uhlemann et al. PNAS 2014

# USA300 substitution rate comparable to other MRSA clones



Corresponds to ~ 4 SNPs per year

# SNPs for ruling in / out transmission

# Applying SNPs to estimate transmission between community households



**Transmission within households**

- ▲ Yes
- ▼ No
- ● No – singleton

**Linkage between households**

- → Sequence only
- → Sequence and epidemiology
- → Epidemiology only

# Bayesian interference of evolutionary rate

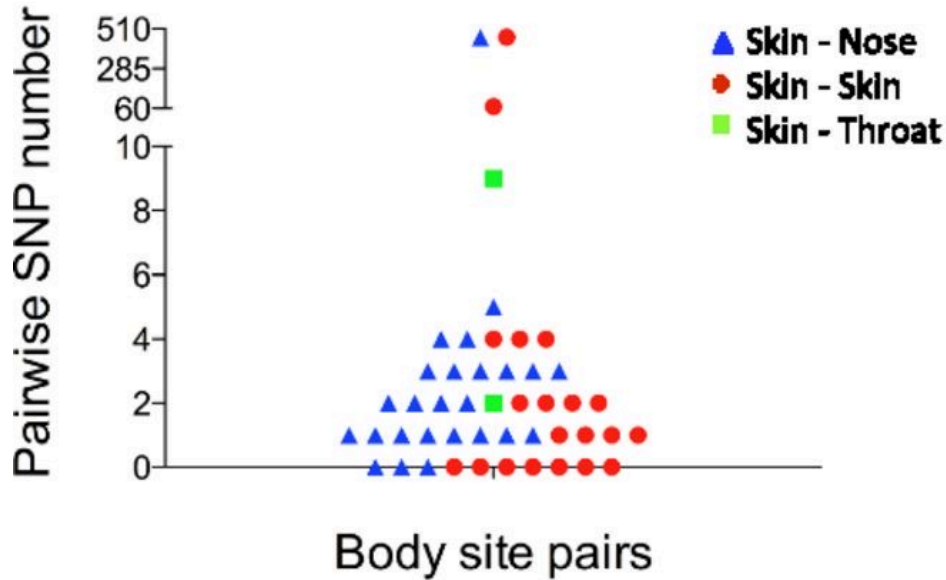- Phylogeny as how to assign probability to different trees given that we observed some sequences.
  - We can think that we do not know the right history but a few histories can be compatible.
  - P(T|D): probability of a tree given the data.
    - that is the inverse of likelihood: P(D|T).
- Uses Markov chain Monte Carlo
- To estimate substitution rates includes:

   - tree topology

   - times of ancestral nodes

   - substitution are (?)

   - substitution parameters (transition/transversion)

# Phylogeographic reconstruction



Support for root in Fort Washington neighborhood (site of CUMC)

Uhlemann et al. PNAS 2014

# Have unique genomic USA300 subpopulations emerged?

Stop codons:



**ebh**
- 10 isolates, 6 households
- ECM binding

**wbrA**
- 29 isolates, 13 households, 20 months apart
- Tryptophan-repressor binding protein, NADP(H)-quinone-oxidoreductase, oxidative stress response?

0.0030

# Expansion of Fluoroquinolone-resistant clone (*gyrA / grlA* SNPs)



CDC survey

    - Decrease in FQ-susceptibility: 63% to 45% from 2004 – 2008

National prescription data overlap with FQ-R prevalence

Additional 15 non-synonymous SNPs associated with gyrA/grlA

# Time scaled evolution of USA300



1927
PP = 1

SaPI5

ACME,
Sa2int-5

1993
PP = 1

1995
PP = 0.92

*

USA300 subclade

927      1941      1954      1968      1982      1995      2009

ST8              SaPi5              ACME & PVL
                                    53 non-synon SNPs

# Summary – part 1

- Comparative bacterial short read sequencing informs

  - outbreak info

  - geographic spread / number of introductions

  - evolutionary history

    - acquisition of MGEs

    - acquisition of drug resistance

- Enabled by different analytical approaches from same dataset

- Beware of limitations in sequence included in analyses (~90% of genome for most bacterial species)

# Microbiome analyses

# ~2,000,000 bacterial genes



MICROBIAL CELLS
~100 TRILLION
(~70-90%)

HUMAN CELLS
~30 TRILLION

MICROBIAL GENES
~2,000,000
(~99%)

HUMAN GENES
~23,000

Many bacterial species previously not recognized because unculturable with current methods.

# Microbiota? Or Microbiome?



Microbiota
16S rRNA
Taxonomic identification

Metagenome
Genes and genomes
of microbiota

Microbiome
Genes, genomes,
products, host
proteins

# 16s rRNA sequencing

- 16S rRNA gene present in all bacterial species
- Highly conserved and variable sequences
- Variable = "molecular fingerprint"
- Amplification with degenerate primers targeting conserved regions
- Large public database for comparisons

# Taxonomy assignment

- Challenges:

  multiple matches

  no match (new OUT)

- Some species may share

  >97% similarity, no resolution
  at species level

# Output taxa distribution

- Bar chart

- Heat map

# Alpha diversity

- Diversity within a sample

  - taxon based

  - phylogeny based

- Richness – number of species present

  - Chao-index

- Evenness – abundance of different species

  - Shannon index

# Beta diversity

- Comparisons of samples to each other
- How different are types present?
- Measure of distance / dissimilarity between sample pair
- UniFrac (weighted, unweighted)

# UniFraq example



Red vs Yellow

Red vs Blue
shared=Grey

Yellow vs Blue

PCoA

PC 2 (25%)

PC 1 (75%)

|   | R | Y | B |
|---|---|---|---|
| R | 0 | .5 | .7 |
| Y | .5 | 0 | 1 |
| B | .7 | 1 | 0 |

Distance Matrix

Hierarchical
Cluster

R
Y
B

Lozupone & Knight 2005 AEM 71:8228

# Principal Coordinate Analysis

- Visualization of beta diversity matrix
- Transform distance matrix into new set of orthogonal axes
- 2D or 3D

# QIIME / QIITA



- Open-source bioinformatics platform
  - data analysis from raw reads to figures
- Qiita: online data repository / data analysis platform

Sequencing output
(454, Illumina, sanger, ...)

fasta/qual/sff files; fastq files;
trace files

Sample metadata

Pre-processing
e.g. remove primer(s), demultiplex,
quality filter

OTU table
(i.e. per sample OTU
counts)

Phylogenetic tree
Evolutionary relationship
between OTUs

Denoise 454 data
AmpliconNoise,
denoiser

Database submission
Submit sequences and
metadata to database

α-diversity and rarefaction
e.g. Phylogenetic
Diversity, Chao1,
observed species

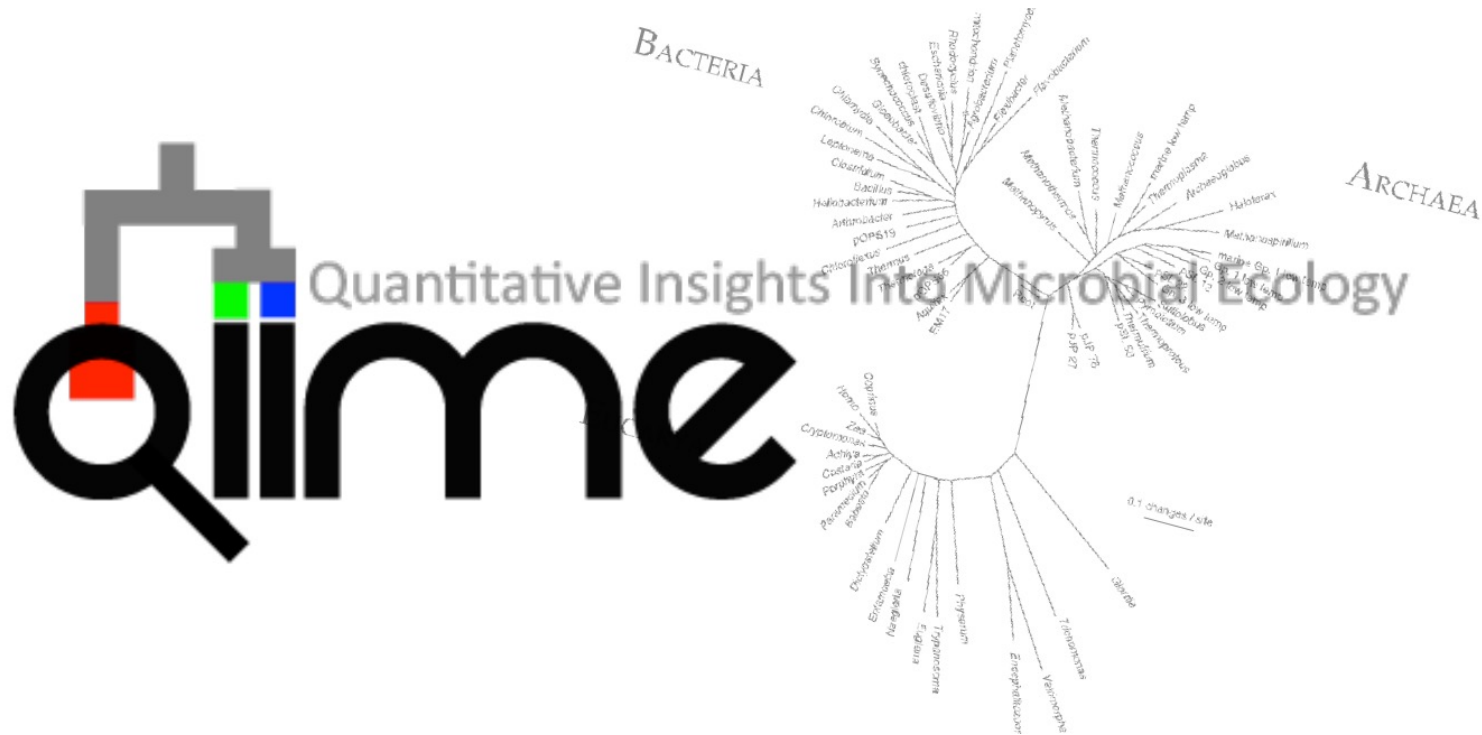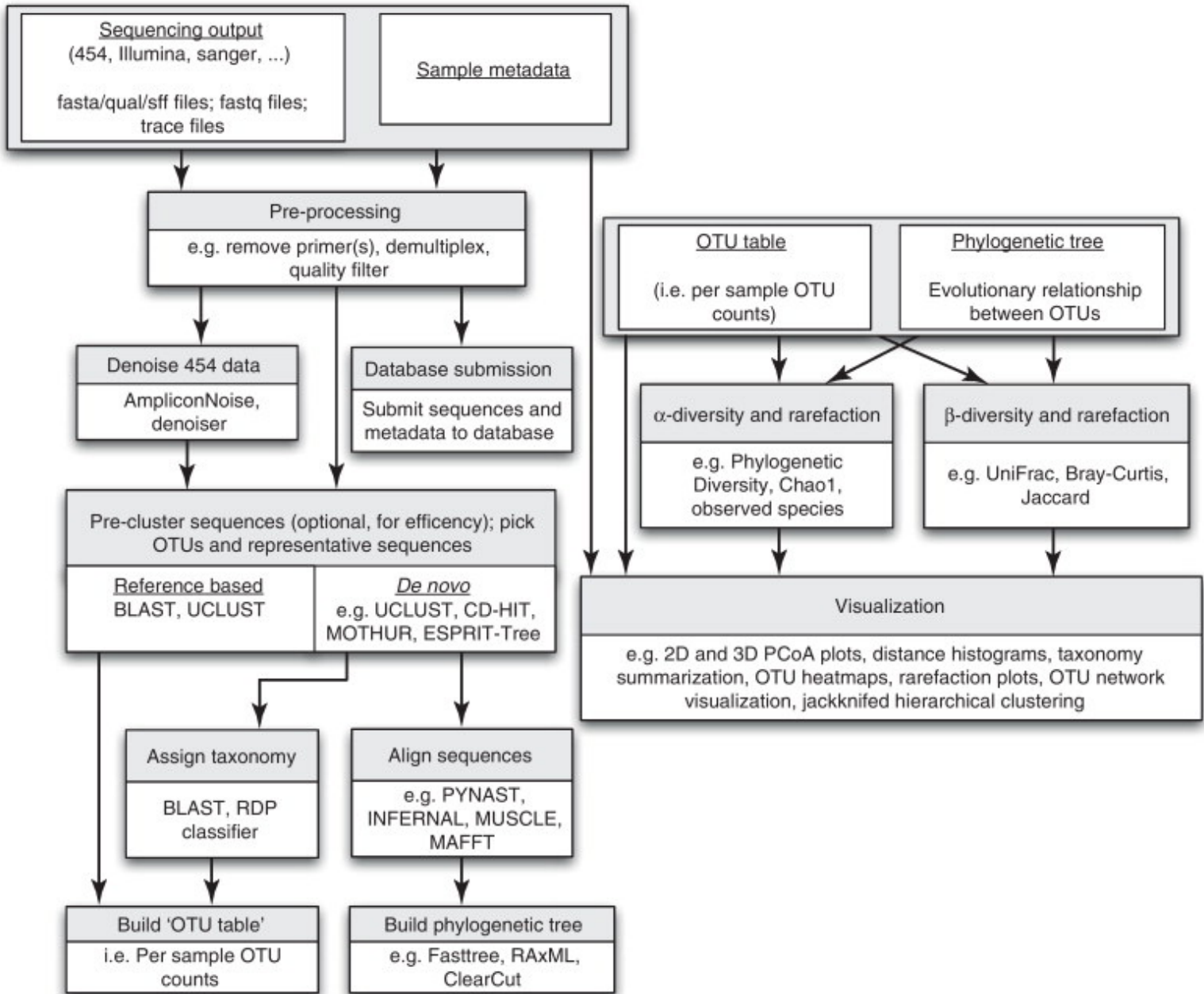β-diversity and rarefaction
e.g. UniFrac, Bray-Curtis,
Jaccard

Pre-cluster sequences (optional, for effiency); pick
OTUs and representative sequences

Reference based
BLAST, UCLUST

De novo
e.g. UCLUST, CD-HIT,
MOTHUR, ESPRIT-Tree

Visualization
e.g. 2D and 3D PCoA plots, distance histograms, taxonomy
summarization, OTU heatmaps, rarefaction plots, OTU network
visualization, jackknifed hierarchical clustering

Assign taxonomy
BLAST, RDP
classifier

Align sequences
e.g. PYNAST,
INFERNAL, MUSCLE,
MAFFT

Build 'OTU table'
i.e. Per sample OTU
counts

Build phylogenetic tree
e.g. Fasttree, RAxML,
ClearCut

Cell
PRESS

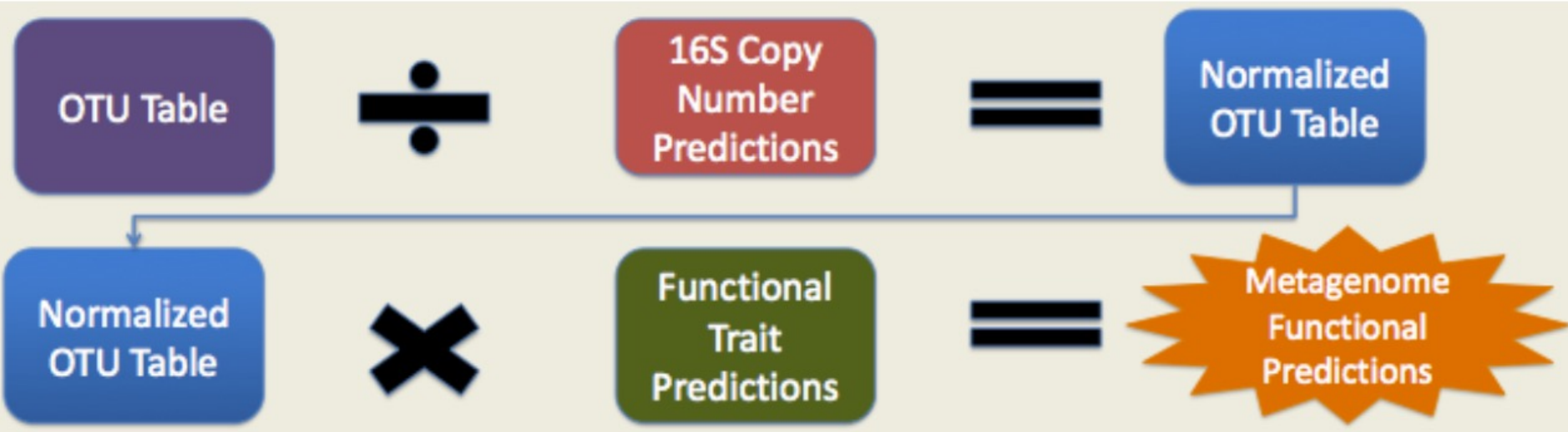TRENDS in Ecology & Evolution

Terms and Conditions
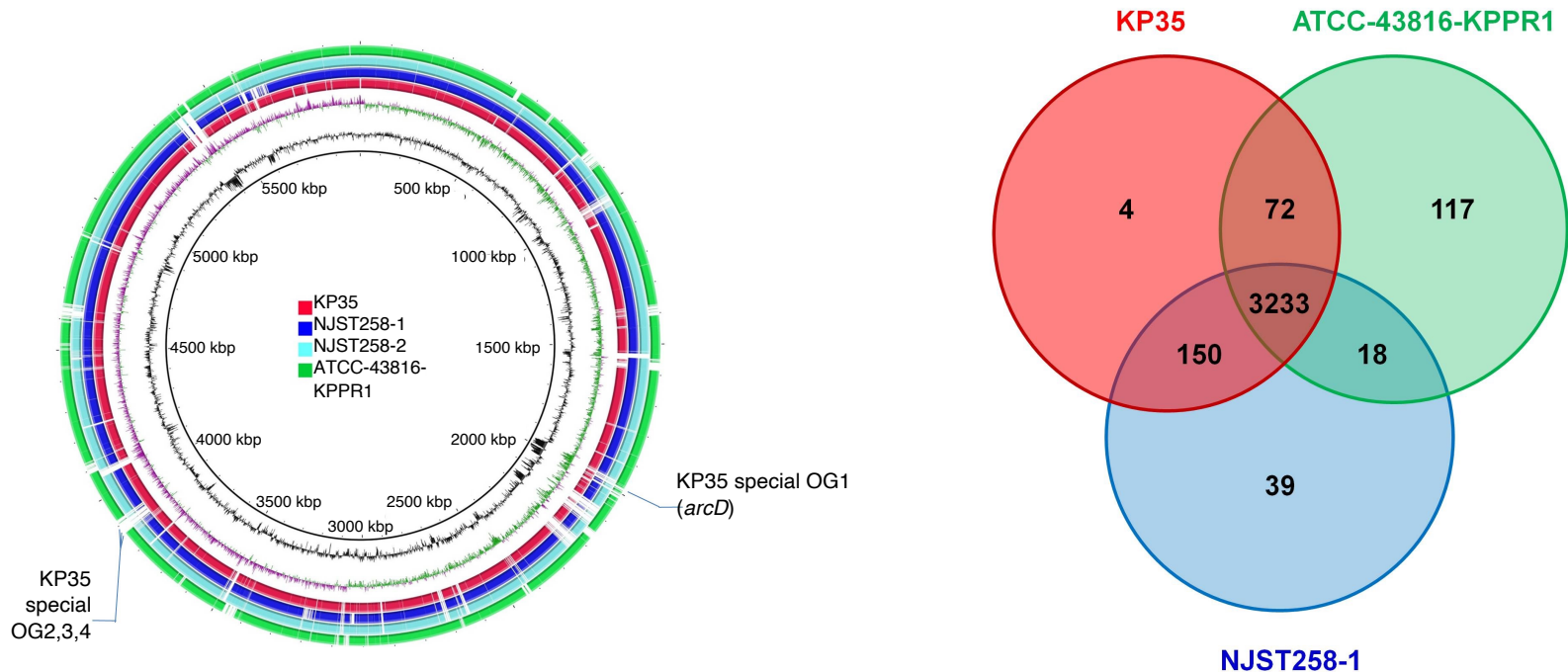
# Additional thoughts on 16S rRNA

- Uniqueness of variable region determines taxonomic resolution

    - V3 or V3-V4 or V1-V2

    - length of variable region

- Optimal resolution depends on sample composition

- Extraction methods (Gram-positive versus Gram-negative!) may play an important role in full recovery of species

# Predictive functional profiling of microbial communities by 16S rRNA genes

- PICRUSt software
- Validated using HMP data

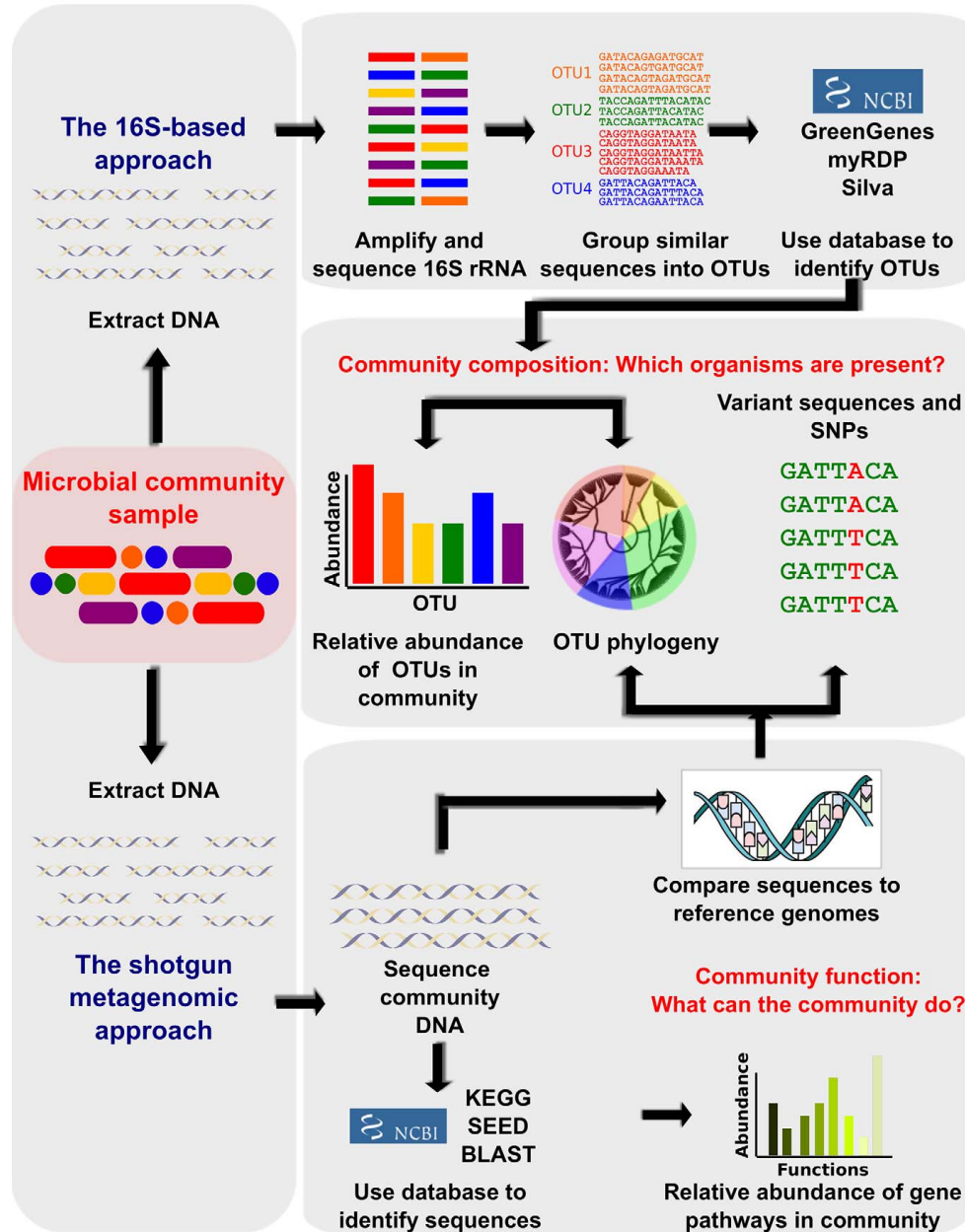# Simulation cannot predict clonal variants within species and genetic content in mobile genetic elements



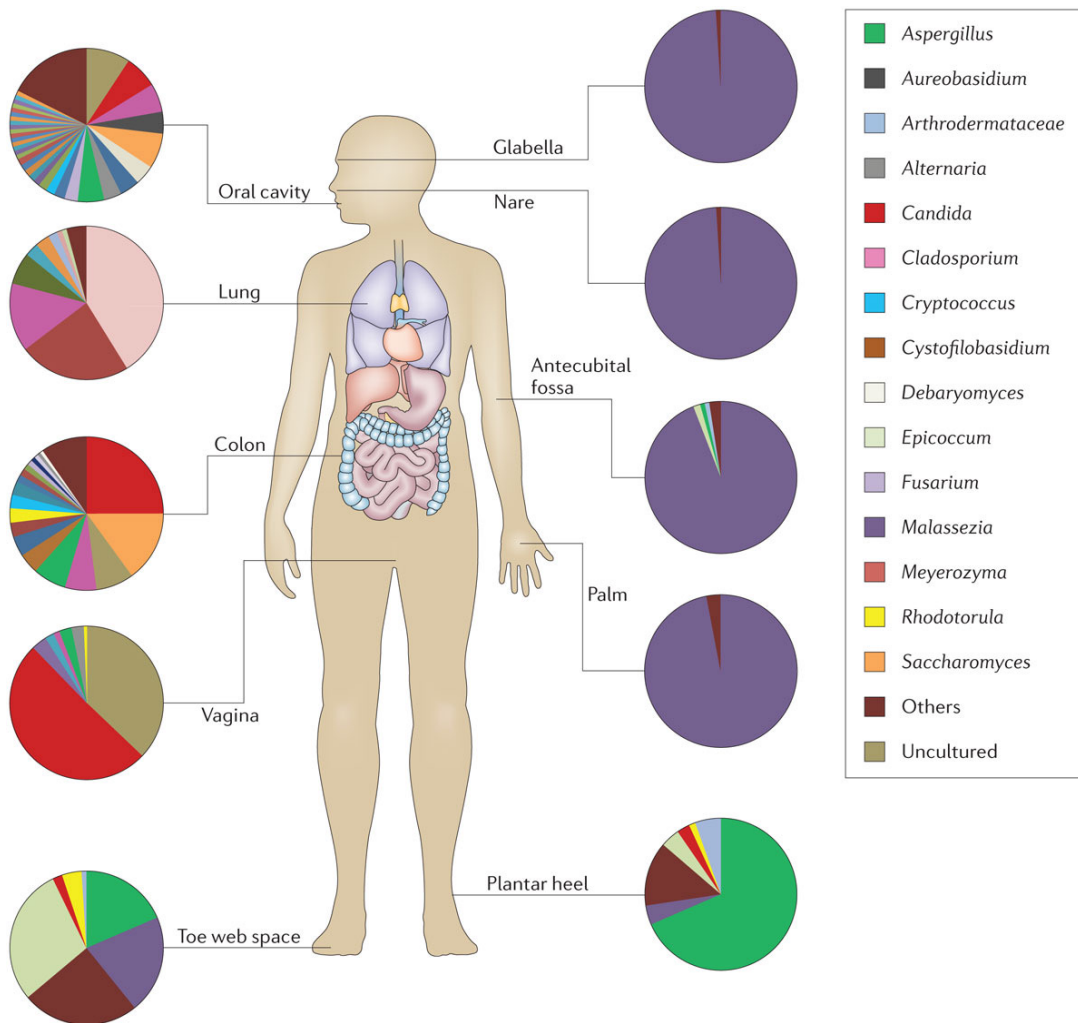Presence of arcD (arginine metabolism) results in decreased virulence, immune evasion

Ahn et al. JCI Insight.

# Bioinformatic methods for functional metagenomics



Morgan & Huttenhower, PLOS Comp. Biol., 2012

# Metagenomics: challenges for high-throughput remain

- Much more starting material required
- Higher sequencing cost (~3-4:1)
    - depends on depth of coverage
- Large amount of data generated, results in high demand on computing infrastructure for data processing and storage
- May still not allow assignment of mobile genetic elements and reliable identification of SNPs
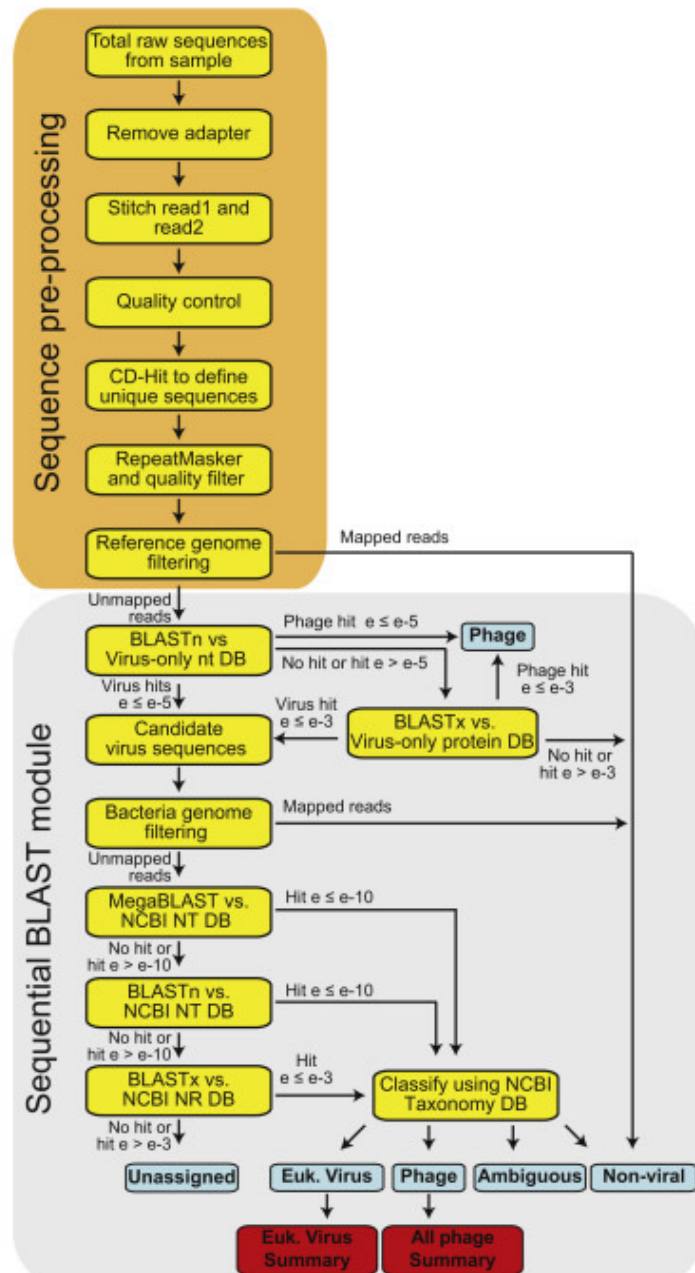
# "Other" microbes: mycobiota



- Sequencing of
  18S ss rDNA
  ITS region
- Longer reads needed
- Growing databases
  UNITE
- Interaction with innate and adaptive immune system

Legend:
- Aspergillus
- Aureobasidium
- Arthrodermataceae
- Alternaria
- Candida
- Cladosporium
- Cryptococcus
- Cystofilobasidium
- Debaryomyces
- Epicoccum
- Fusarium
- Malassezia
- Meyerozyma
- Rhodotorula
- Saccharomyces
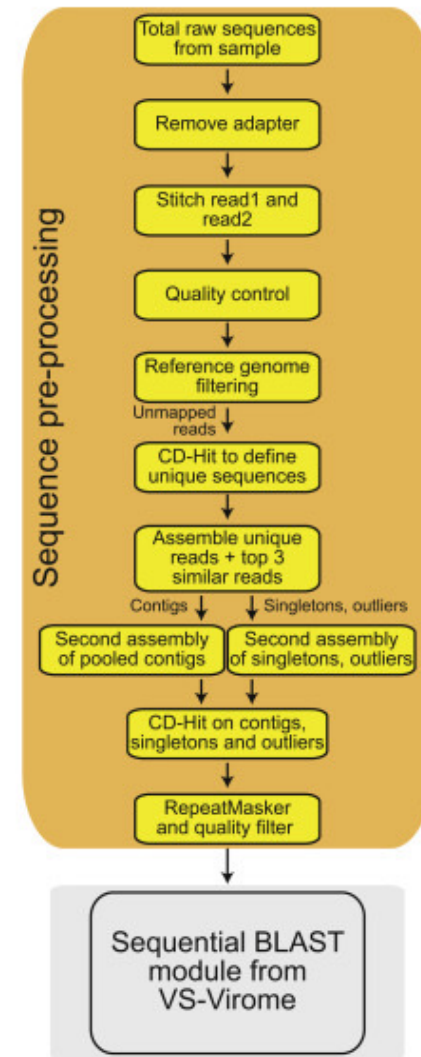- Others
- Uncultured

Underhill, Iliev Nat Rev Immun 2014

# "Other" microbes - Virome

- Most abundant and fastest mutating genetic elements
- Previously difficult to sequence / analyze given high diversity
- Difficult to extract (enrichment from filtrates, lysis of bacterial and human cells)
- Different types!
  - Eukaryotic, Bacterial, Archaeaic viruses
  - Integrated elements in human host DNA
- Trans-kingdom interaction
- Direct interaction with host / immune signaling
- Phages regulate bacterial content

H. Virgin Cell 2015

**VirusSeeker, a computational pipeline for virus discovery and virome composition analysis**
Virology, Volume 503, 2017, 21–30

# Microbiome summary

- Fingerprint of bacterial communities
- Relatively affordable and fast
- Does not provide information on unique functional features (MGEs…)
- Metagenomics will be more comprehensive but currently still expensive / data intensive, limiting widespread use. Difficulties assigning plasmids to organisms