

Ocean Omics: Metagenome Data Processing

SIOB 278
Week 1: Introduction
April 04, 2017



DATE	TOPICS / ACTIVITIES
Apr 05	Introduction to metagenomics, research questions, and experimental design <i>AWS access; Command line operations; Local software needs; Selection of independent data sets</i>
Apr 12	Data types, expectations, and quality control <i>Sequencing technologies; File formats; Data quality control processing (FastQC, Trimmomatic)</i>
Apr 19	Assembly of metagenome sequence data <i>Assembly processing (IDBA-UD, MetaSPAdes); Assembly statistics & validation (Quast)</i>
Apr 26	Check point 1: Research questions, discussion, and papers
May 03	Gene prediction and functional annotation of assembled contigs/scaffolds <i>Annotation pipeline (Prokka); Web-based analysis portal (JGI-IMG); Annotation viewing (Artemis)</i>
May 10	Binning of assembled sequences I: DNA compositional metrics <i>Nucleotide frequency analysis (%GC); Coverage (mapping); Analysis & visualization (Anvi'o)</i>
May 17	Binning of assembled sequences II: taxonomy assignments <i>Homology searches (BLAST, DIAMOND, DarkHorse); Visualization of binned data (Anvi'o)</i>
May 24	Check point 2: Research questions, discussion, & papers
May 31	Ad-hoc analyses: population genome assembly, gene & pathway discovery
Jun 07	Presentation of independent research projects

Extent of Microbial Diversity

Environment	Average No. cells per volume	Total No. cells ($\times 10^{28}$)	Estimated No. taxa
Aquatic habitats			
Freshwater	700,000	2	200,000
Ocean:	100,000	6	2,000,000 (20,000 Archaea)
< 200 m	800,000	2	
> 200 m	100,000	4	
Marine subsurface			
0-10 cm	280,000,000	25	
>10 cm	14,000,000	330	
Terrestrial habitats			
Soil	800,000,000	26	4,000,000
Subsurface (> 10 m)	800,000	25-250	
Host-associated habitats			
Human gut	320,000,000,000	0.000039	1,000

**~ 5×10^{30} microbial cells globally
representing ~ 10^7 species**

*Data from Whitman WB *et al.* PNAS (1998), Karner MB *et al.* Nature (2001), Curtis TP *et al.* PNAS (2002), Curtis TP *et al.* Phil Trans B (2006)

ENVIRONMENTAL MICROBIAL BIOLOGY

outstanding questions relating microbial processes and activities in the environment

Physiology & metabolism

Genetic & metabolic potential of uncultivated organisms

Genotype : phenotype linkage

Activity in the environment

Genes of unknown function

Cultivating the uncultivated

Autecological phenomena

Environmental adaptation

Biogeographical patterning

Diversity & spatio-temporal dynamics

Environmental control of expression

Evolutionary processes

Modes and tempo of species evolution / selection

Species definition

Horizontal gene transfer dynamics

Evolution of organelles...

Discovery

Novel genes, metabolisms, physiologies

Novel diversity

Biotechnology / biomedicine products

Human health

Genomic approaches to characterize environmental microbial biology

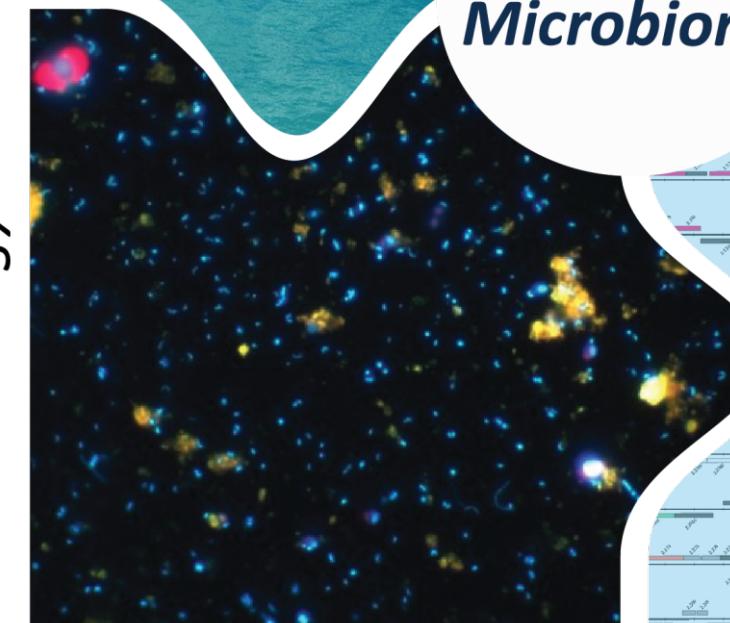
- 16S rRNA gene tag sequencing
- Metagenomics
- Metatranscriptomics
- Metaproteomics
- Metabolomics
- Isolate sequencing



ocean chemistry

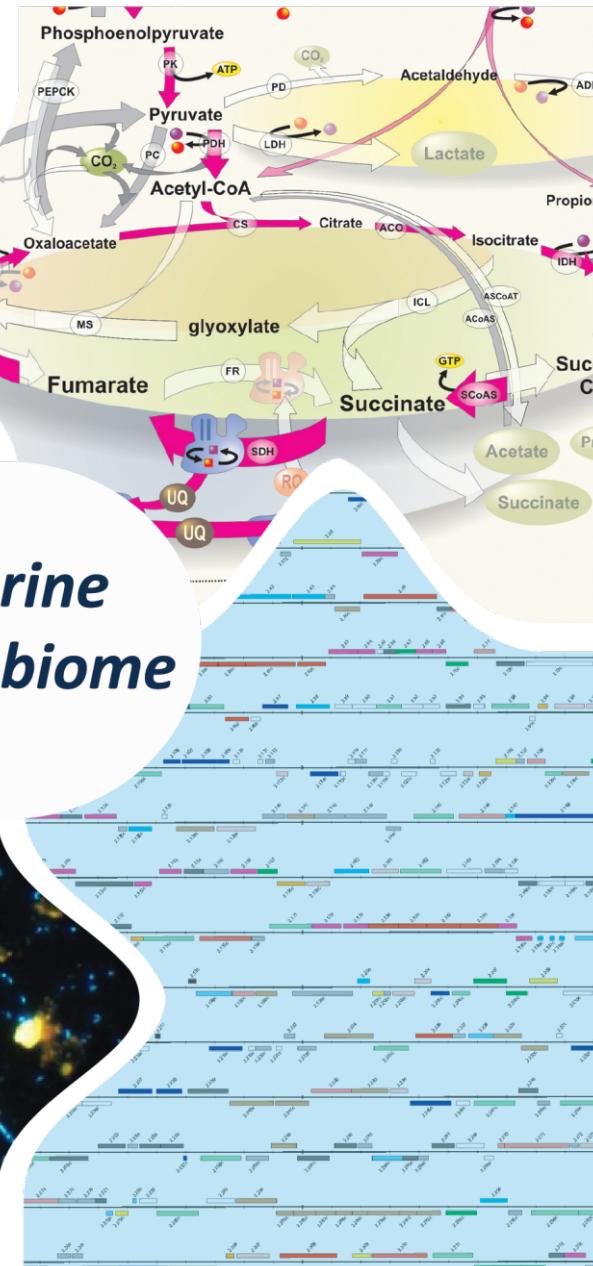


oceanography



microbial diversity

microbial impacts



genomics

metabolism

evolution

...the metagenomic approach...



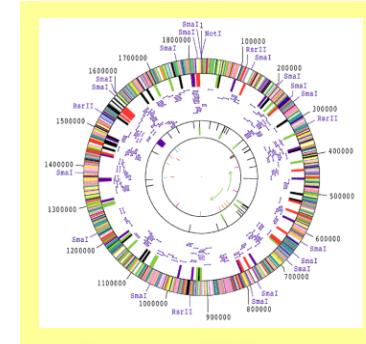
in situ

- Field collection
- Geochemistry
- Sample analysis
- Sample preservation



ex vivo

- DNA extraction
- PCR amplification
- Library construction
- RNA extraction
- cDNA synthesis
- DNA sequencing
- Protein extraction
- LC/MS-MS analysis



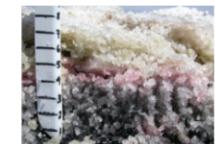
in silico

- Sequence processing
- Quality control filtering
- Sequence assembly
- Phylogenetic binning
- Functional annotation
- Comparative genomics
- Metabolism/physiology
- mRNA/protein expression

*in
situ*



Environmental Sample

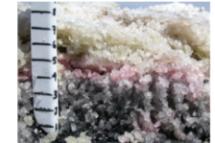
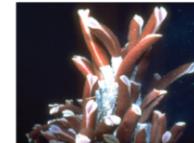


*ex
vivo*

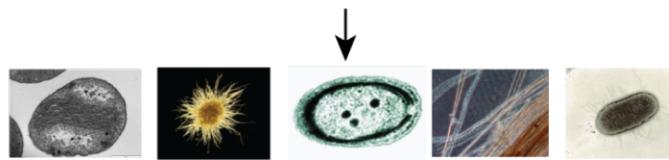
*in
silico*



Environmental Sample

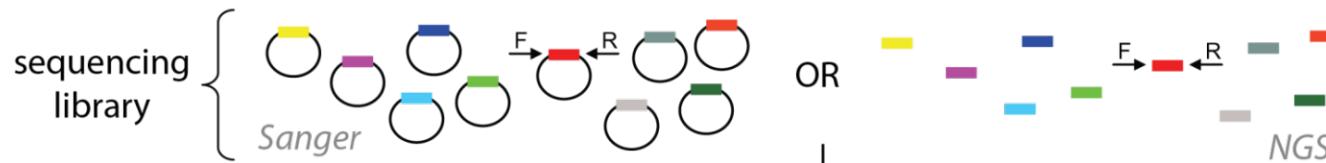


[Geochemistry / Hydrology / Mineralogy / Temperature / Pressure / Location / Time]



AGCGCGCGCATACGAGCGC AGCGATAGCAGACTCGATTAGCT
GAGTGCGATA TAGC AGCGATAGCAGACTCGATTAGCT GAGTGCGATA TAGC
AGCGATAGCAGAGTCGATTAGCT AGCGCGCGCATACGAGCGC GAGTGCGATA TAGC
AGCGCGCGCATACGAGCGC AGCGATAGCAGAGTCGATTAGCT AGCGCGCGCATACGAGCGC
AGCGATAGCAGAGTCGATTAGCT AGCGATAGCAGAGTCGATTAGCT

whole genome shotgun sequencing



end-sequencing



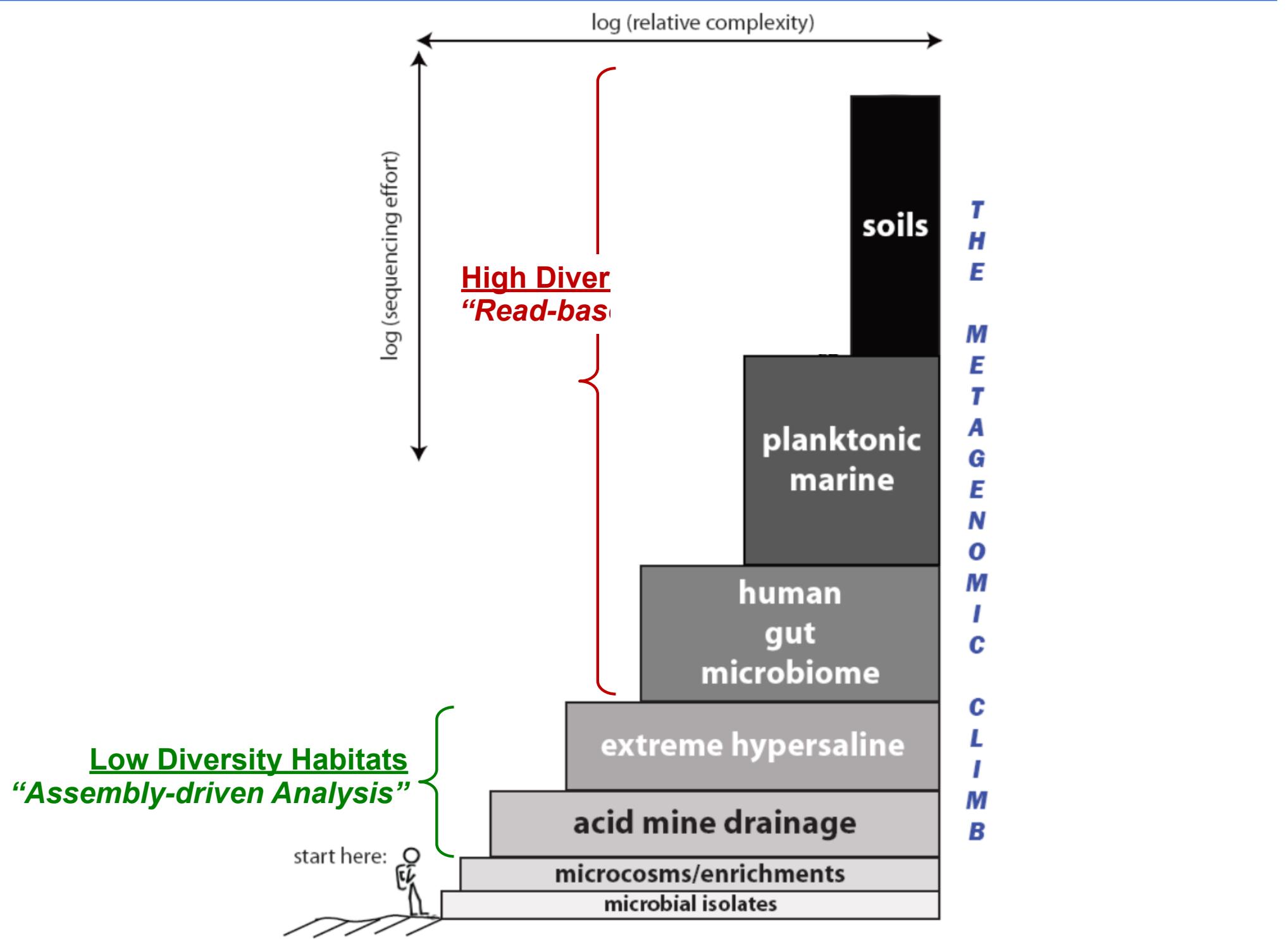
”Read-based” analysis

assembly

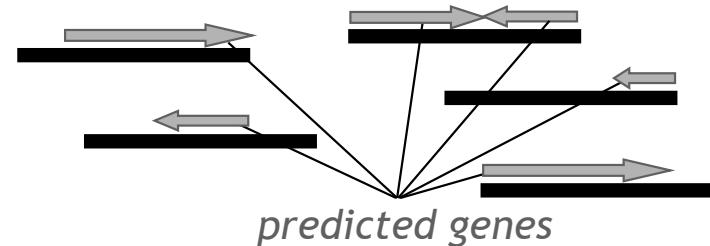
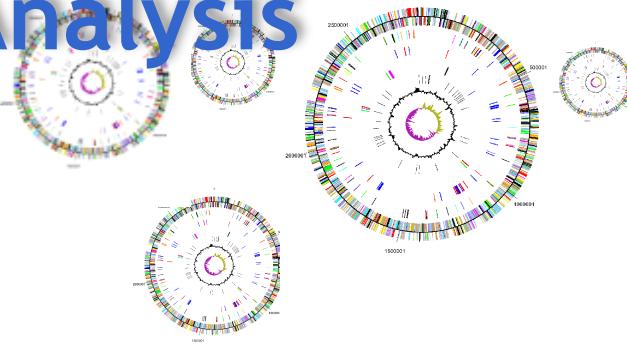


”Assembly-driven” analysis

Metagenomic Adventures



“Assembly-driven” vs “Read-based” Analysis



- Improved phylogenetic binning of DNA fragments
- Insight into operon structure & genomic islands
- Improved metabolic reconstruction for organisms
- Provides reference sequences for transcriptomic, proteomic, & read-based studies
- Rapid prediction of genes & functional annotation
- Provides a census of taxa & metabolic processes
- Allows comparative analysis within & between samples
- Broad view of community structure & function

*The level of environmental sequencing must be commensurate with the expected diversity of the sample
undersampling complicates interpretation due to incomplete data sets*

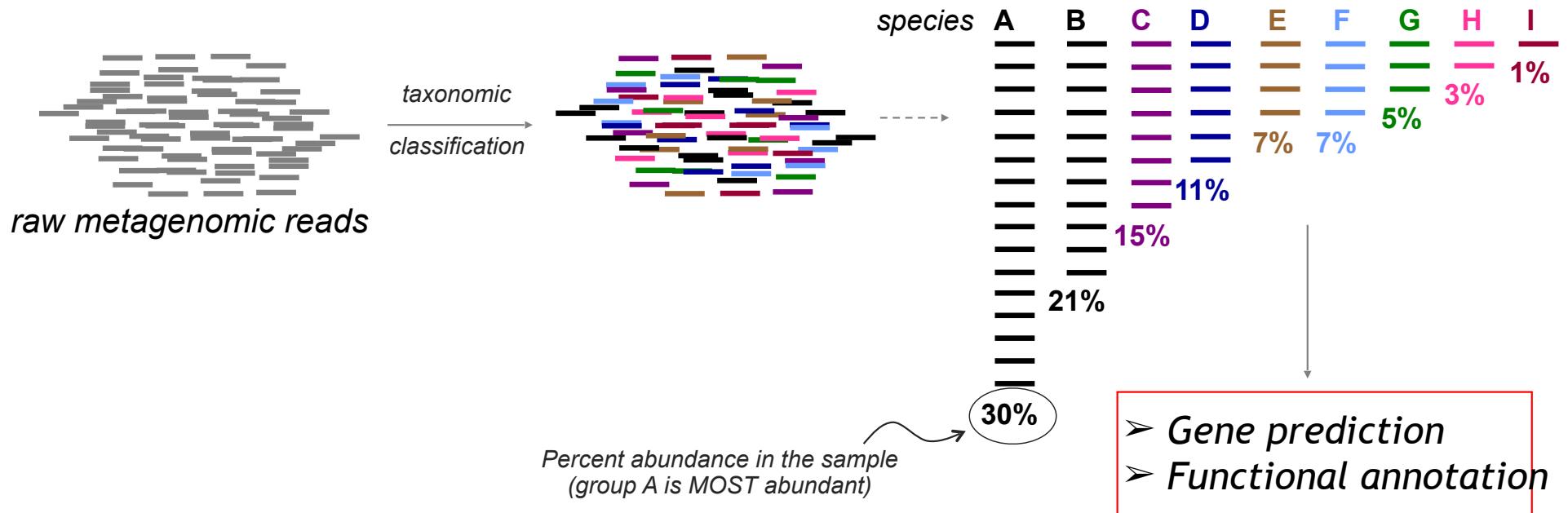
“Assembly-driven” vs “Read-based” Analysis

- Read-based analyses look at the genes encoded on individual metagenomic sequencing reads
- Assembly-driven analyses seek to assemble complete genomes from metagenomic reads

Having only a few pieces of the puzzle complicates interpretation but there are still bits of information in these small pieces (metagenomic reads)

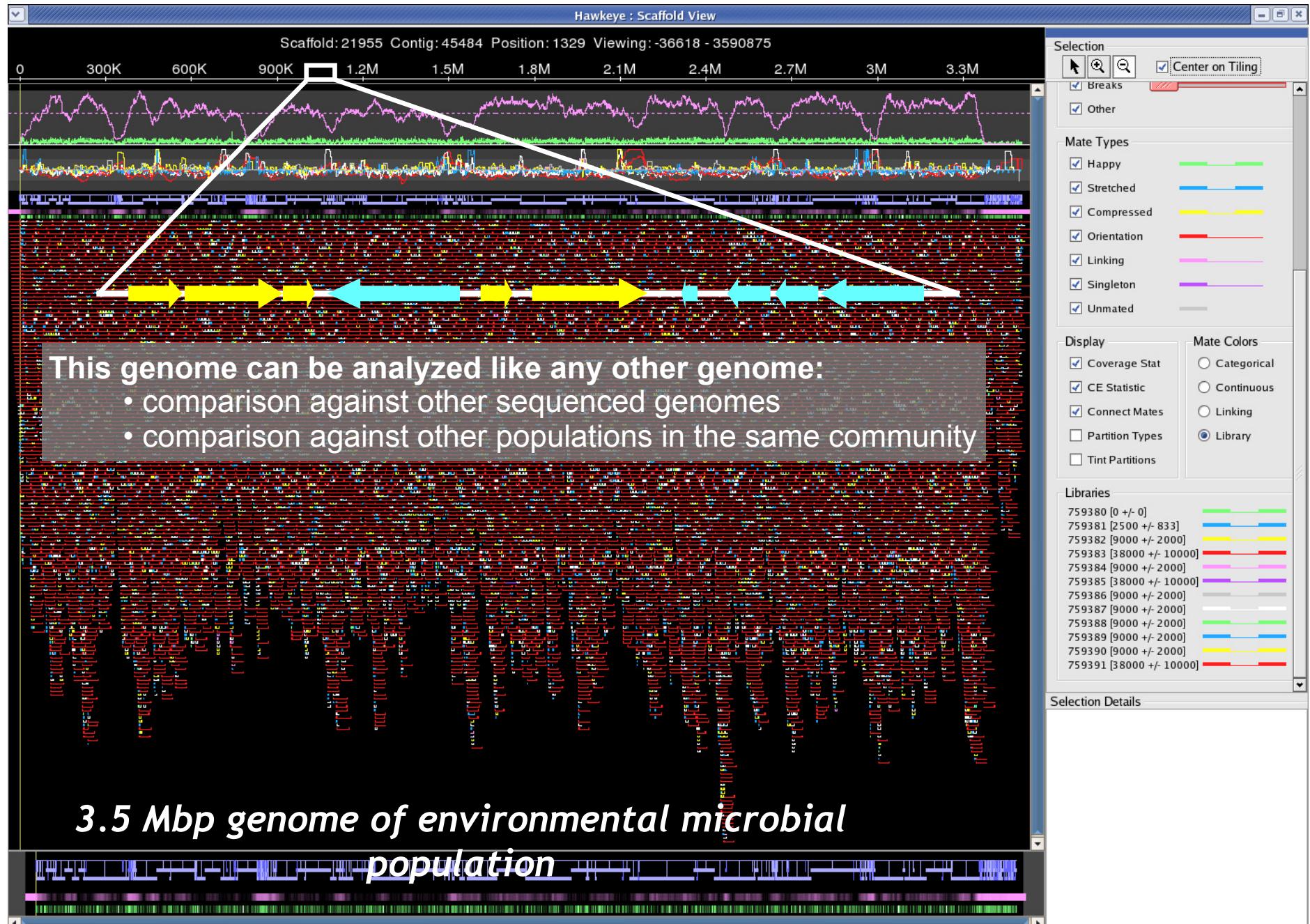


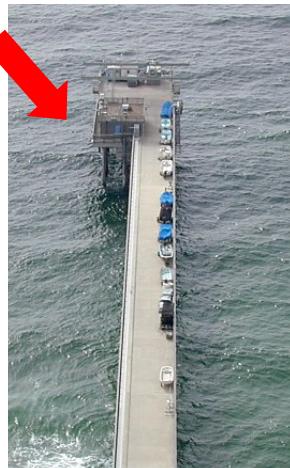
“Read-based” Sequence Analysis



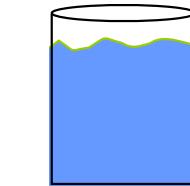
No assembly required to provide an inventory of the taxonomic diversity and functional capabilities within the community sampled

“Assembly-driven” Analysis





Environment



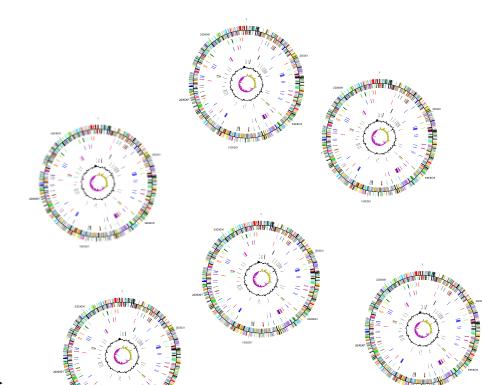
Seawater

Make genomic
sequencing library

Harvest cells
&
Extract DNA

Sequence
Sequence
Sequence

Assemble
&
Analyze

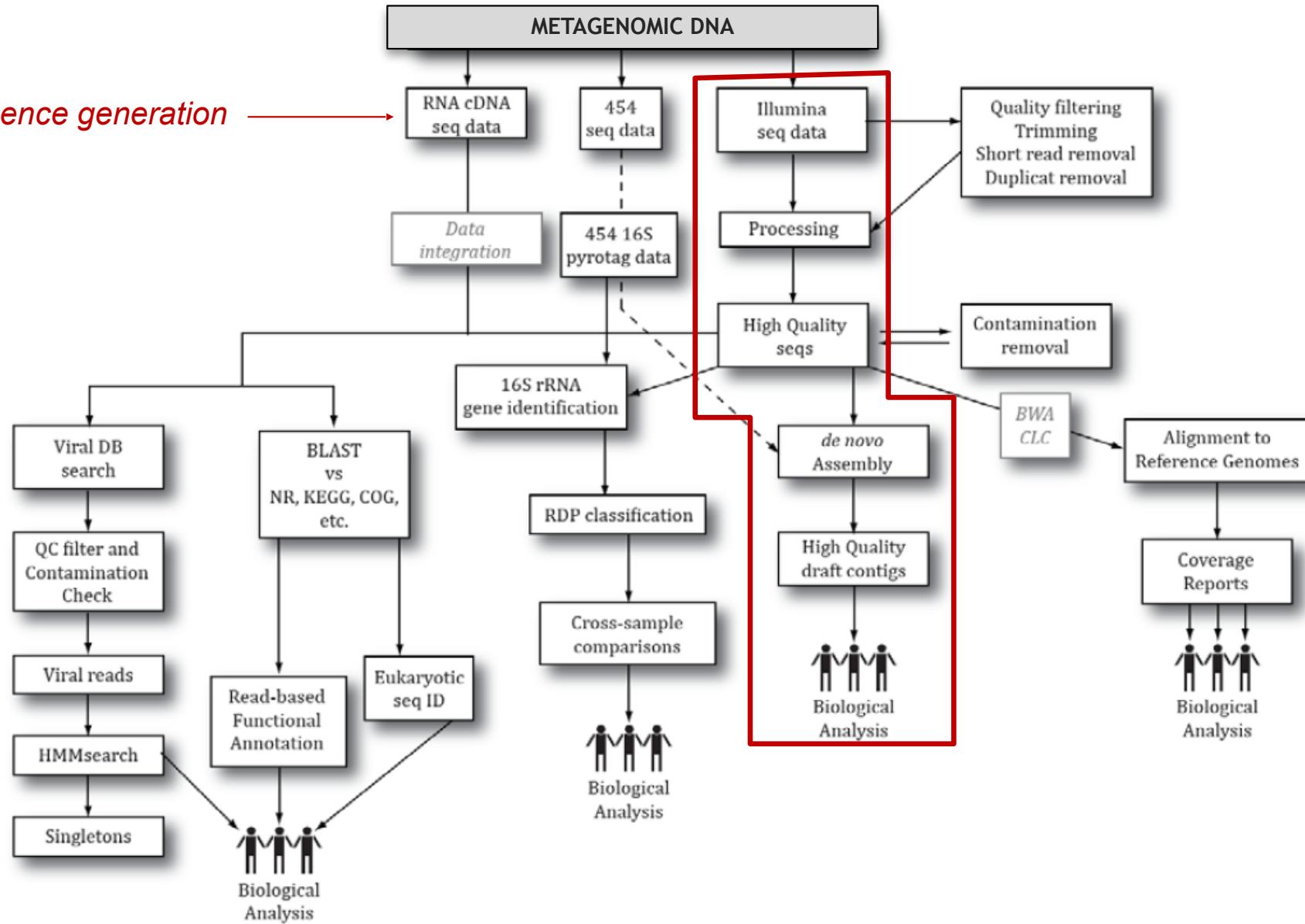


Genomes from
my environment

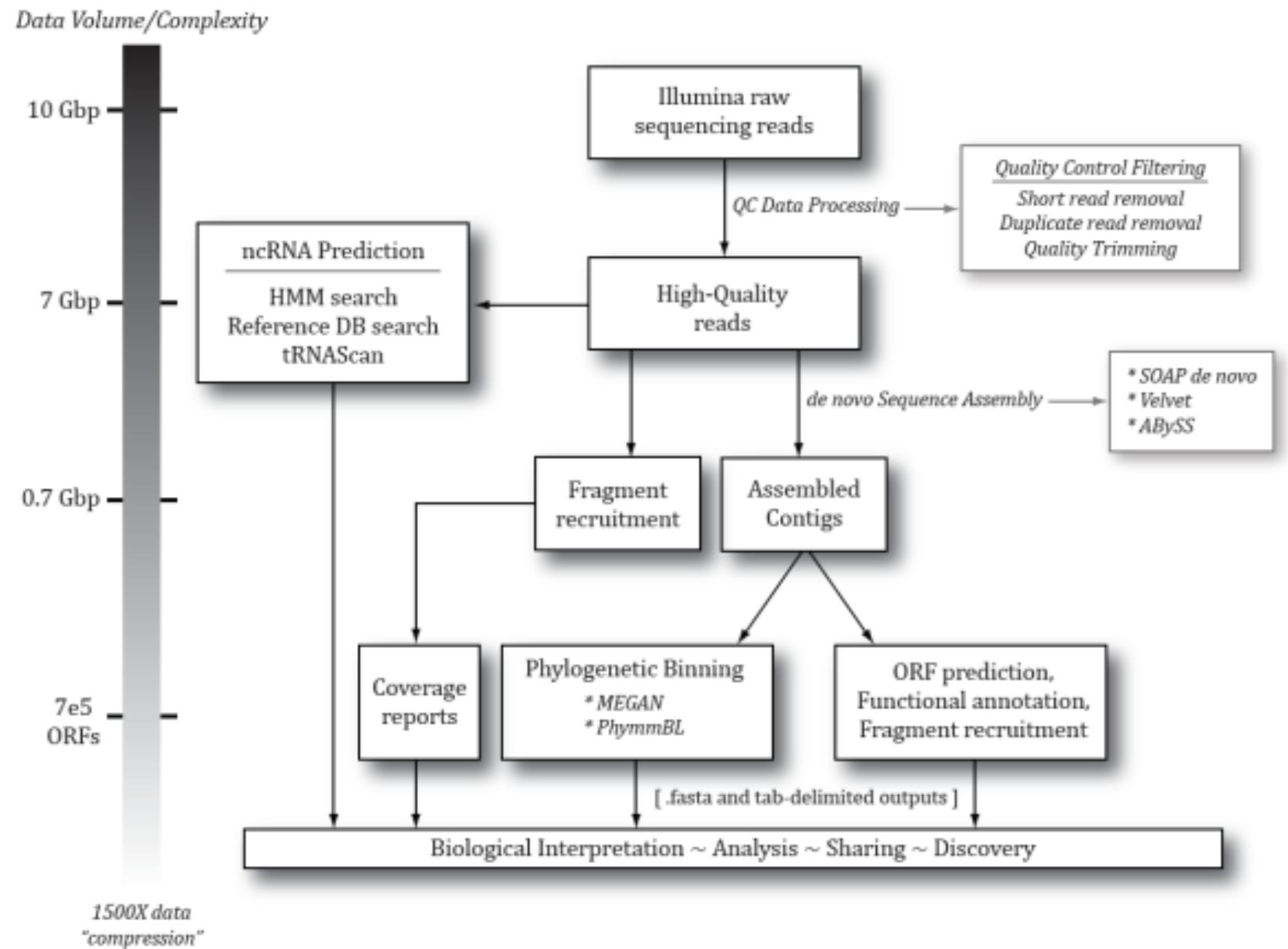
Schematic of simple metagenomic workflow

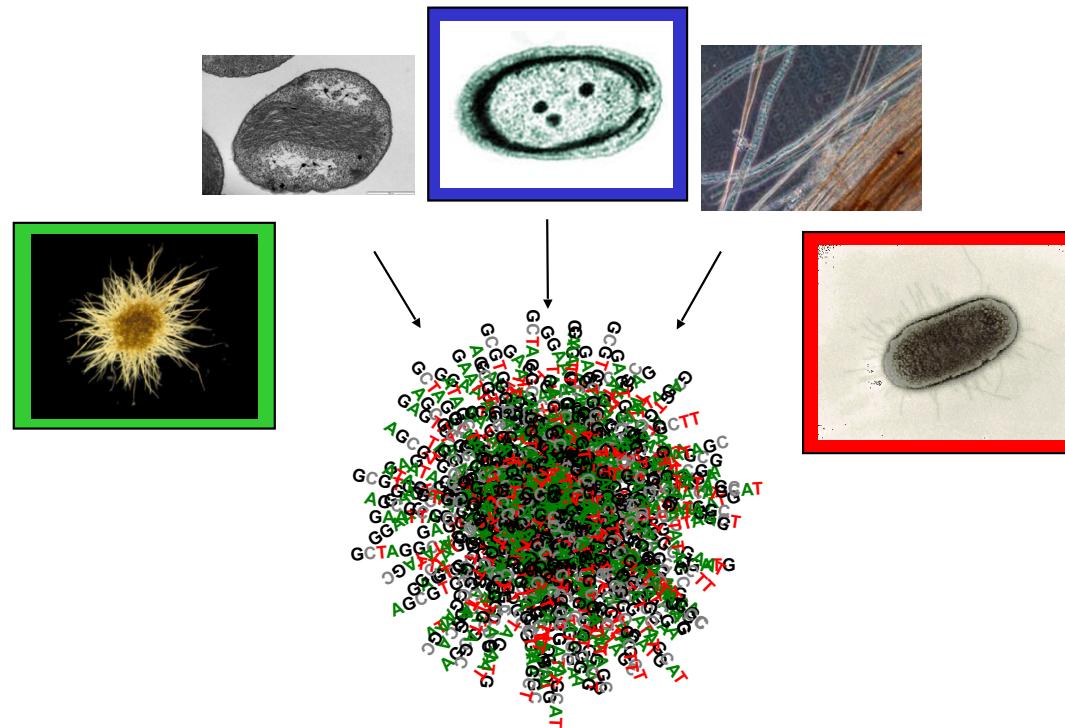
Analysis scheme for processing metagenomic data sets

Sequence generation

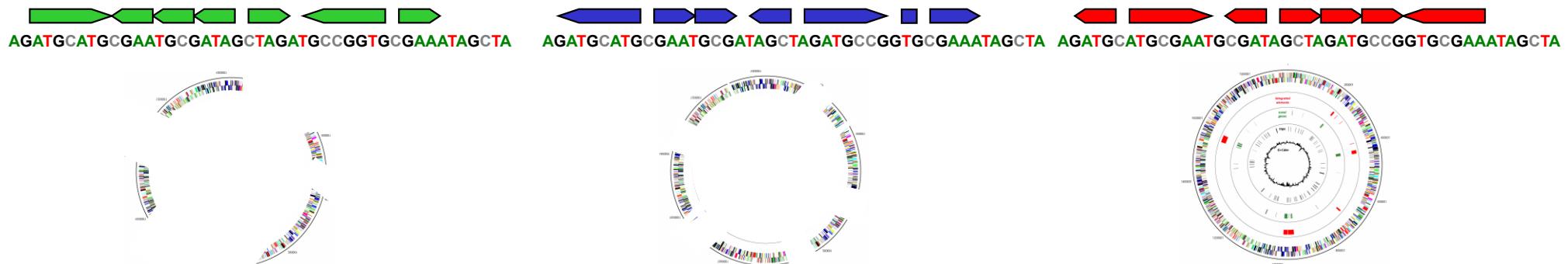


Metagenome *de novo* assembly workflow

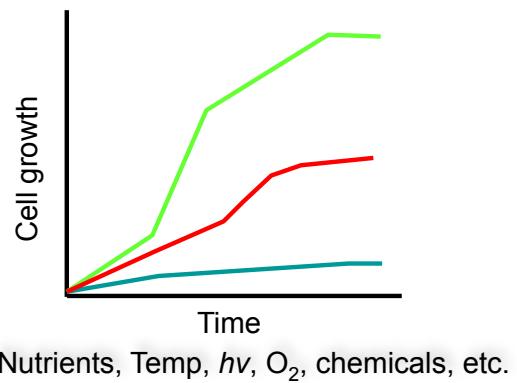
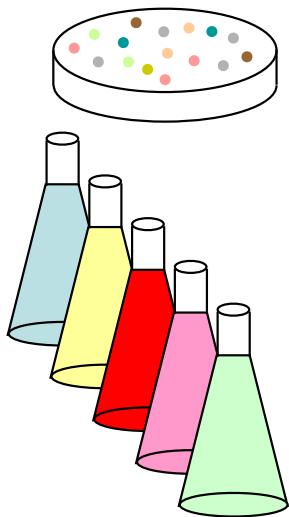




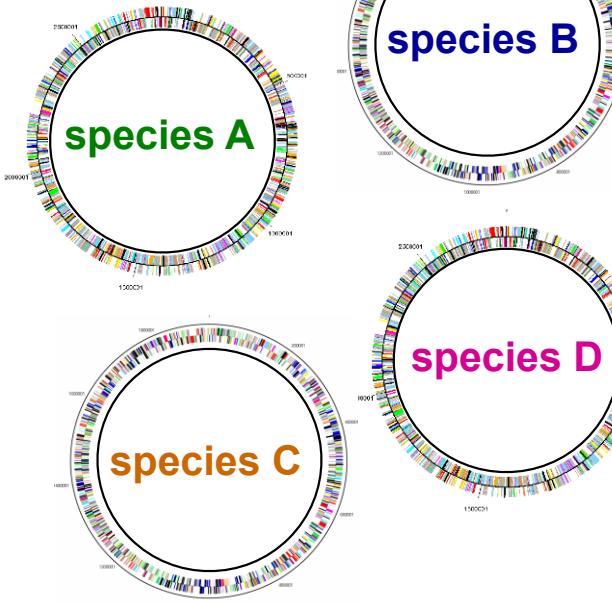
Metagenomic Assembly, Phylogenetic Binning, Functional Annotation



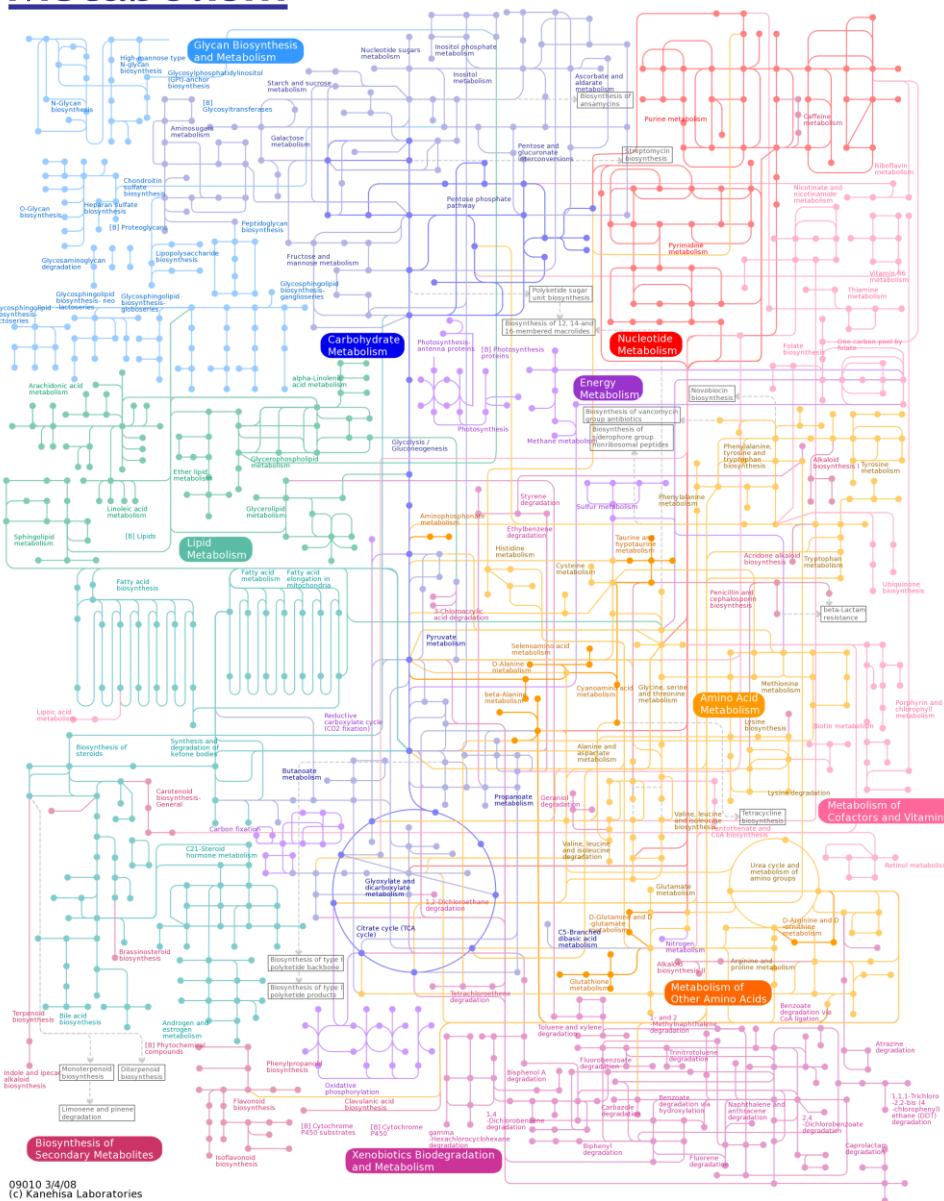
Experimentation



Genomes



Metabolism



...how an organism makes a living...