

Computational Methods in Metagenomics and Microbiome Research

0368-3116-01

Prof. Elhanan Borenstein
School of Computer Science
Semester B, 2019



Outline

A Few Updates about Seminar Logistics

“Metagenomics”

Genomes, Genes, and the Central Dogma

Real Metagenomics

Computational Challenges in Metagenomics

Outline

A Few Updates about Seminar Logistics

“Metagenomics”

Genomes, Genes, and the Central Dogma

Real Metagenomics

Computational Challenges in Metagenomics

Seminar Format

- ~15 students registered
- 1 paper per student, 1-2 students/papers per class
- Class structure:
 - 1:10 Start talk 1 (35 minutes!!)
 - 1:45-2:00 Feedback/Discussion
 - 2:00-2:10 Break
 - 2:10 Start talk 2 (35 minutes!!)
 - 2:45-3:00 Feedback/Discussion
- Paper selection:
 - List posted: Sunday, March 10, 9:00
 - Please select your paper by: Monday, 17:00
 - First talk: March 20



Presentations

- In your presentation:
 - Emphasize the main task of the paper
 - Cover required background
 - Focus more (but not only) on methods (may need to dig)
 - **Choose wisely what to cover and how deeply**
 - Summarize briefly
 - Add something original
- Remember the dos and don'ts of a good talk and of slide design
- Keep us interested

Other Tasks

- **As a speaker**

- Come early (or on a different day) to make sure your presentation works
- Print slides for everybody
- Send me slides to post online
- Prepare discussion points (2-3 points, NOT in paper)

- **As a listener**

- Read/skim the papers before class
- Listen, take notes, engage
- Participate in feedback/discussion
- Learn and have fun

Outline

A Few Updates about Seminar Logistics

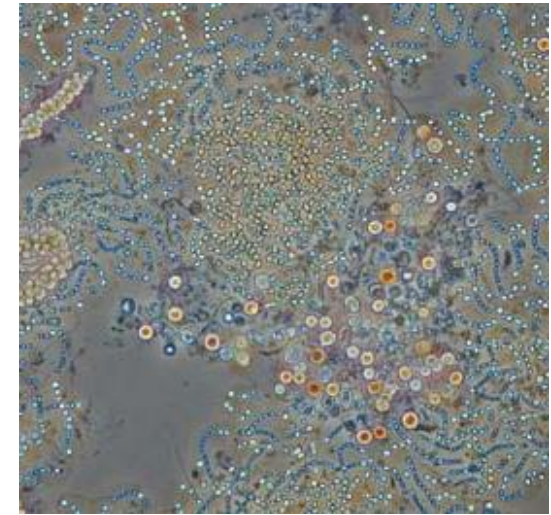
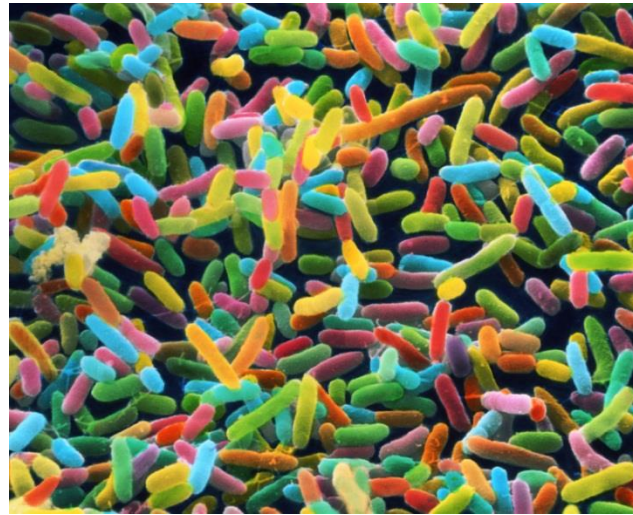
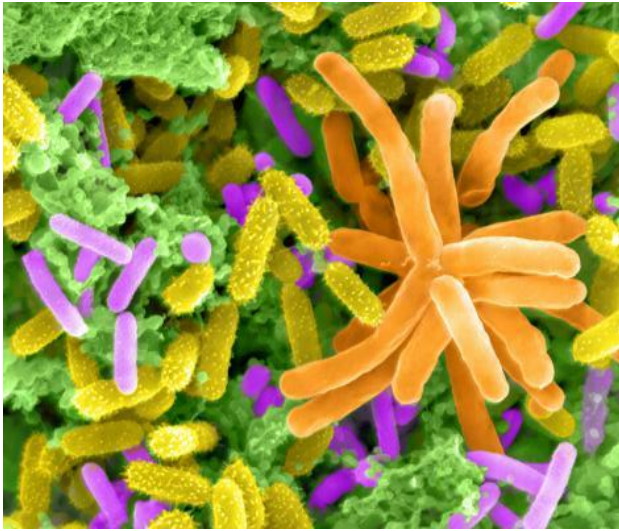
“Metagenomics”

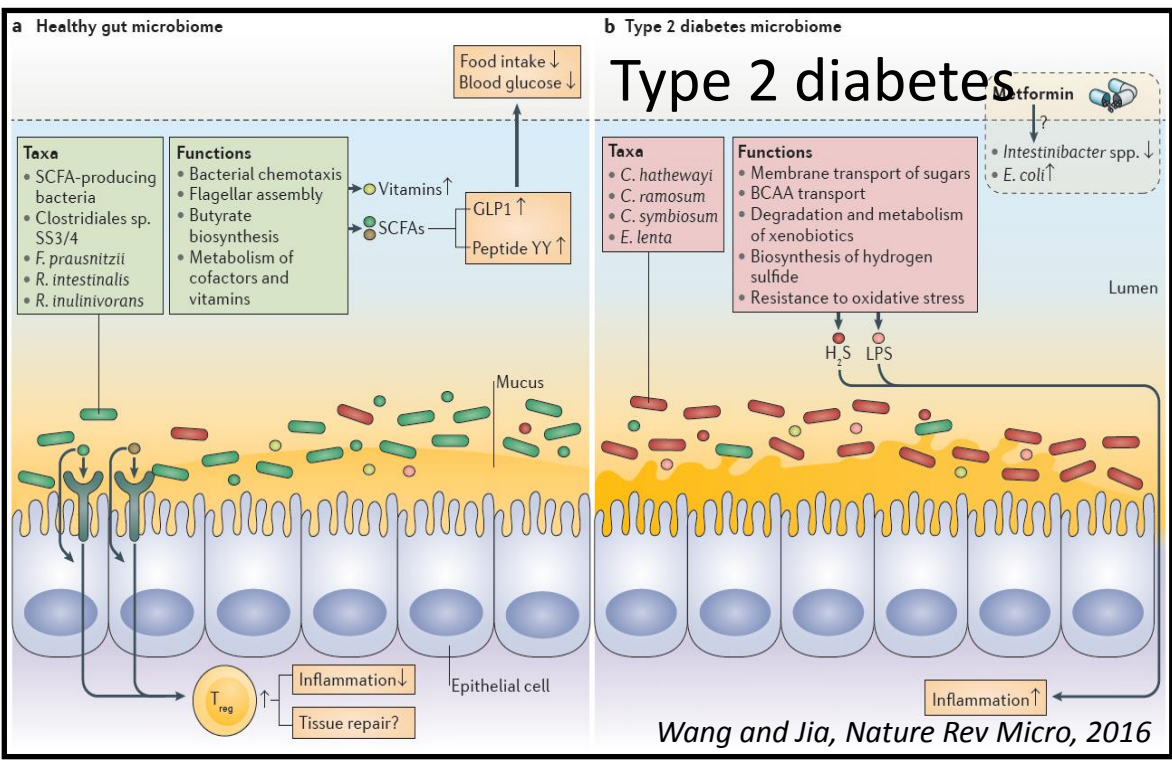
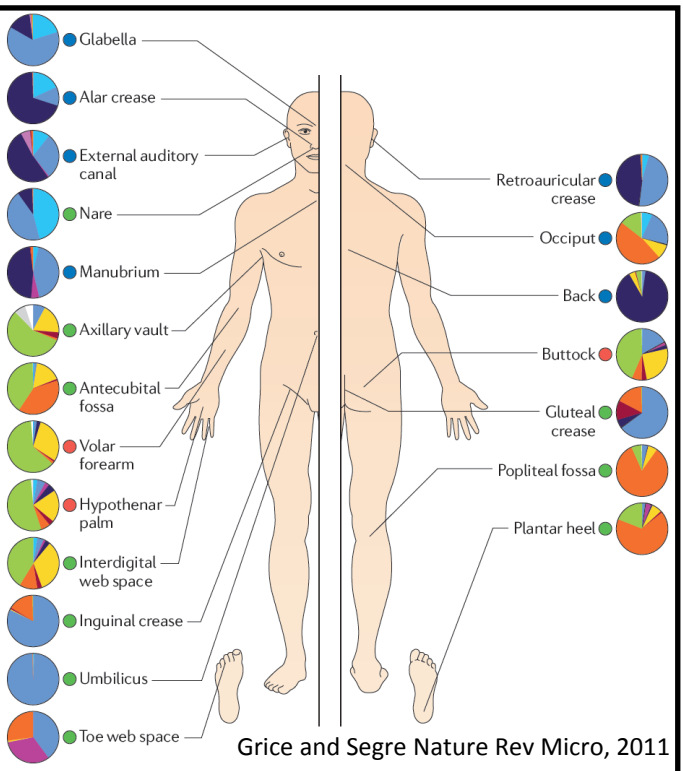
Genomes, Genes, and the Central Dogma

Real Metagenomics

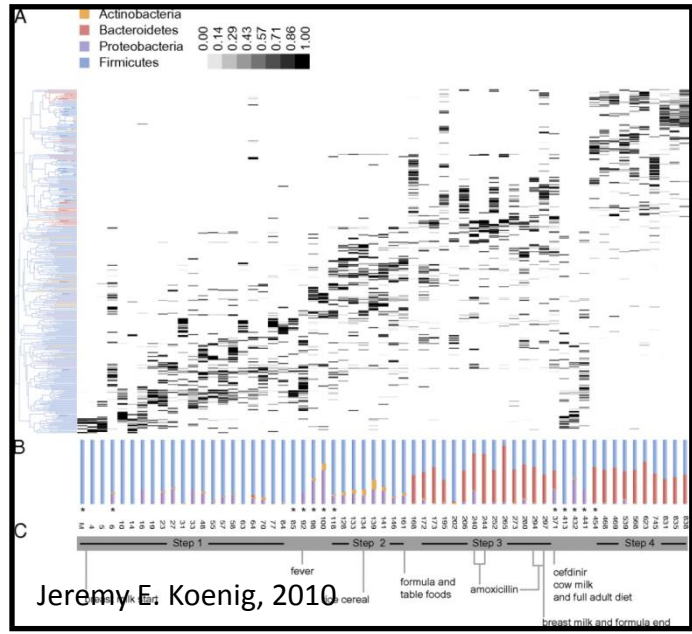
Computational Challenges in Metagenomics

Microbial communities





- **Hundreds of species!**
- **100 trillion microbes!**
(weighing ~3-4lb)
- **150x more genes**
(~3,300,000)

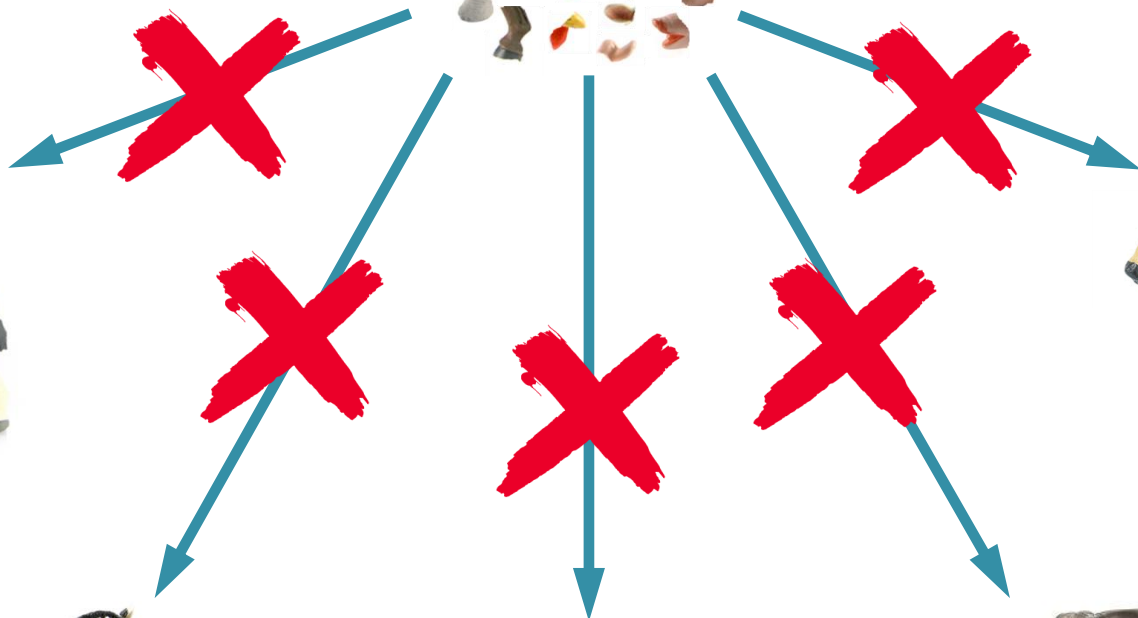


metagenomics

“The study of genetic material recovered directly from environmental samples”

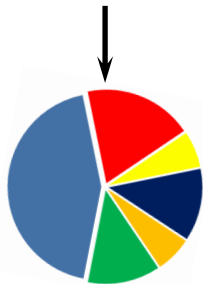
(the term was first used by Jo Handelsman in 1998)







Who's there?

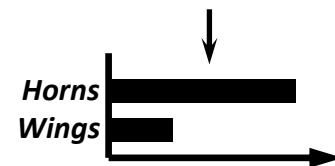


What are they doing?

Horns



Wings



Outline

A Few Updates about Seminar Logistics

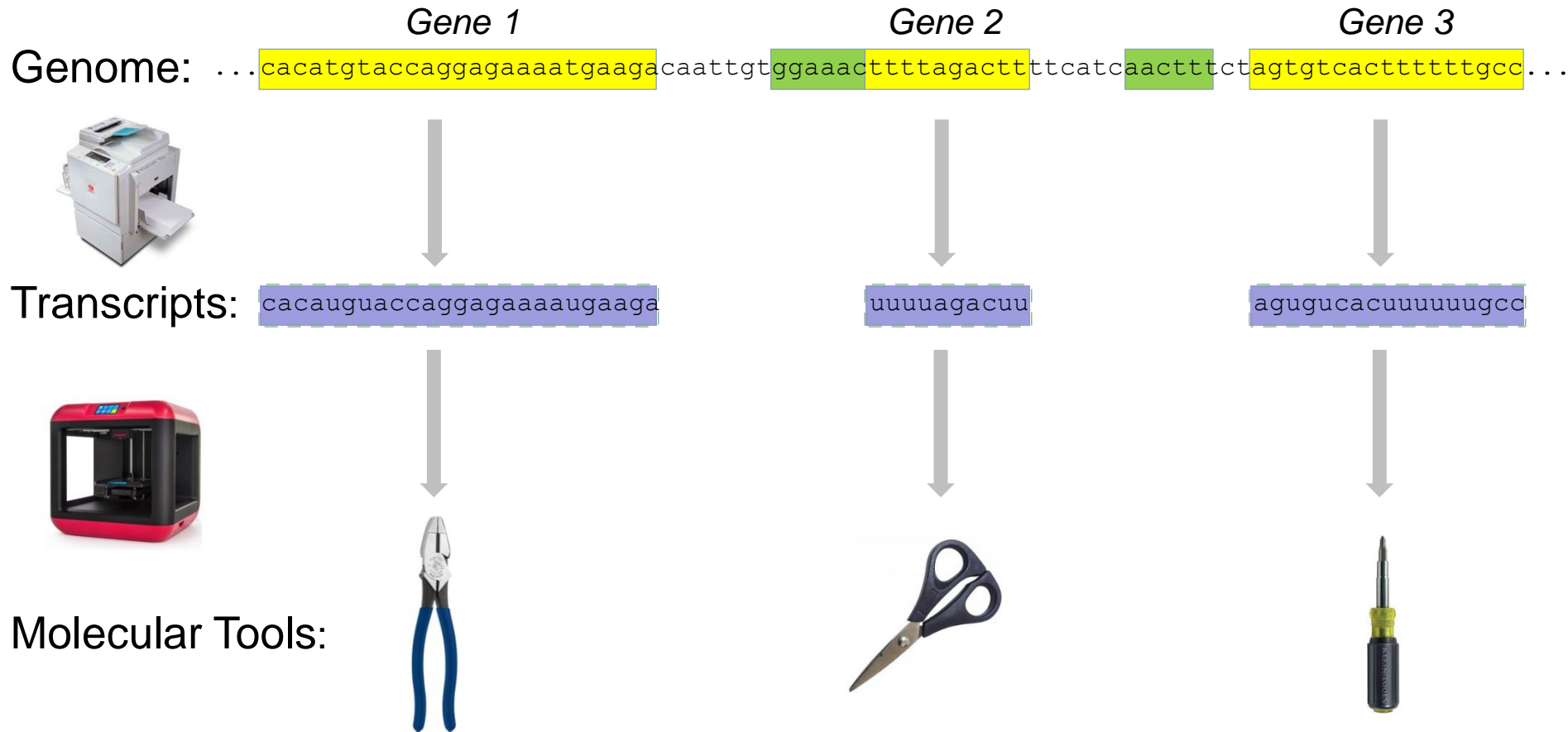
“Metagenomics”

Genomes, Genes, and the Central Dogma

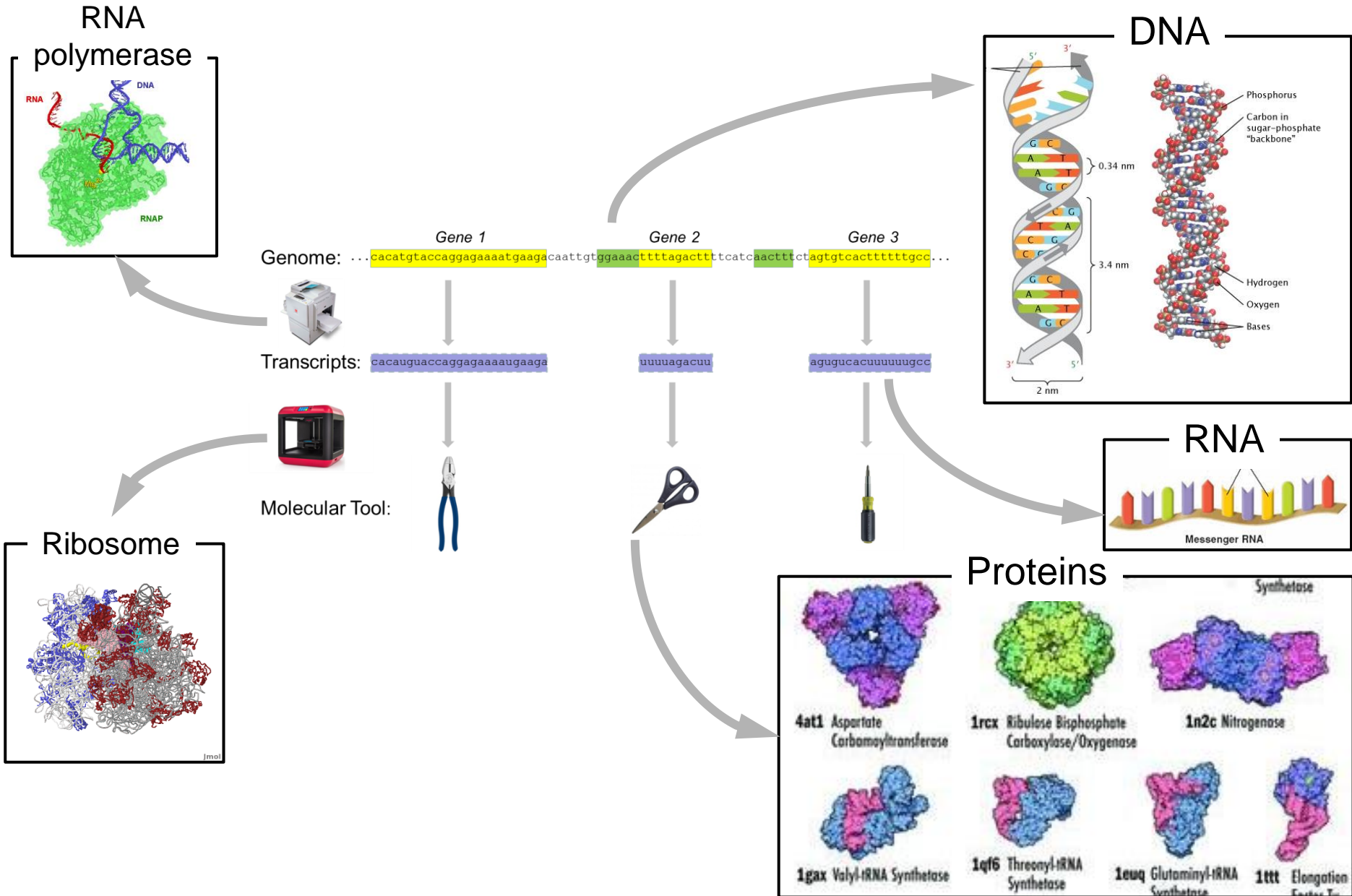
Real Metagenomics

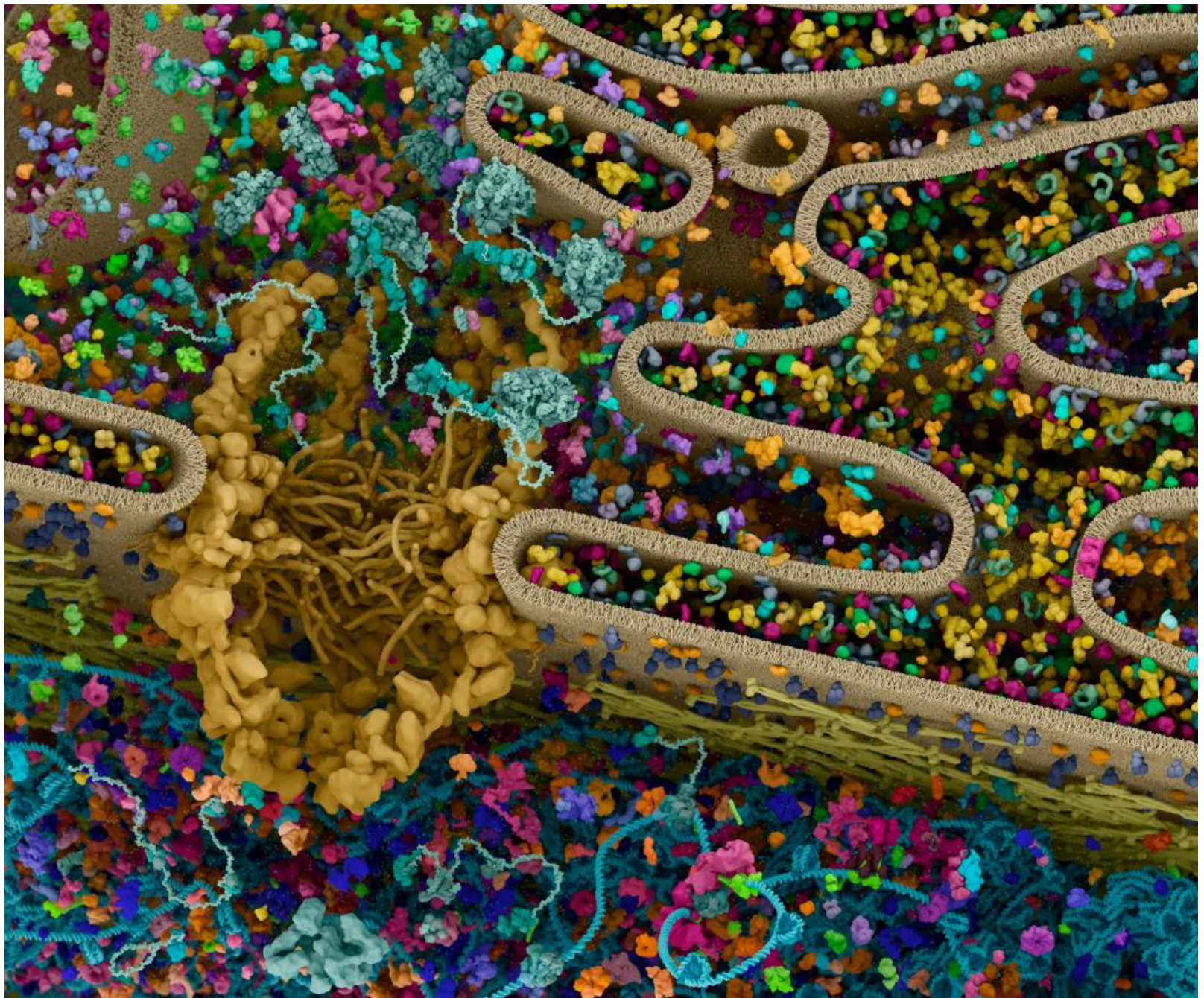
Computational Challenges in Metagenomics

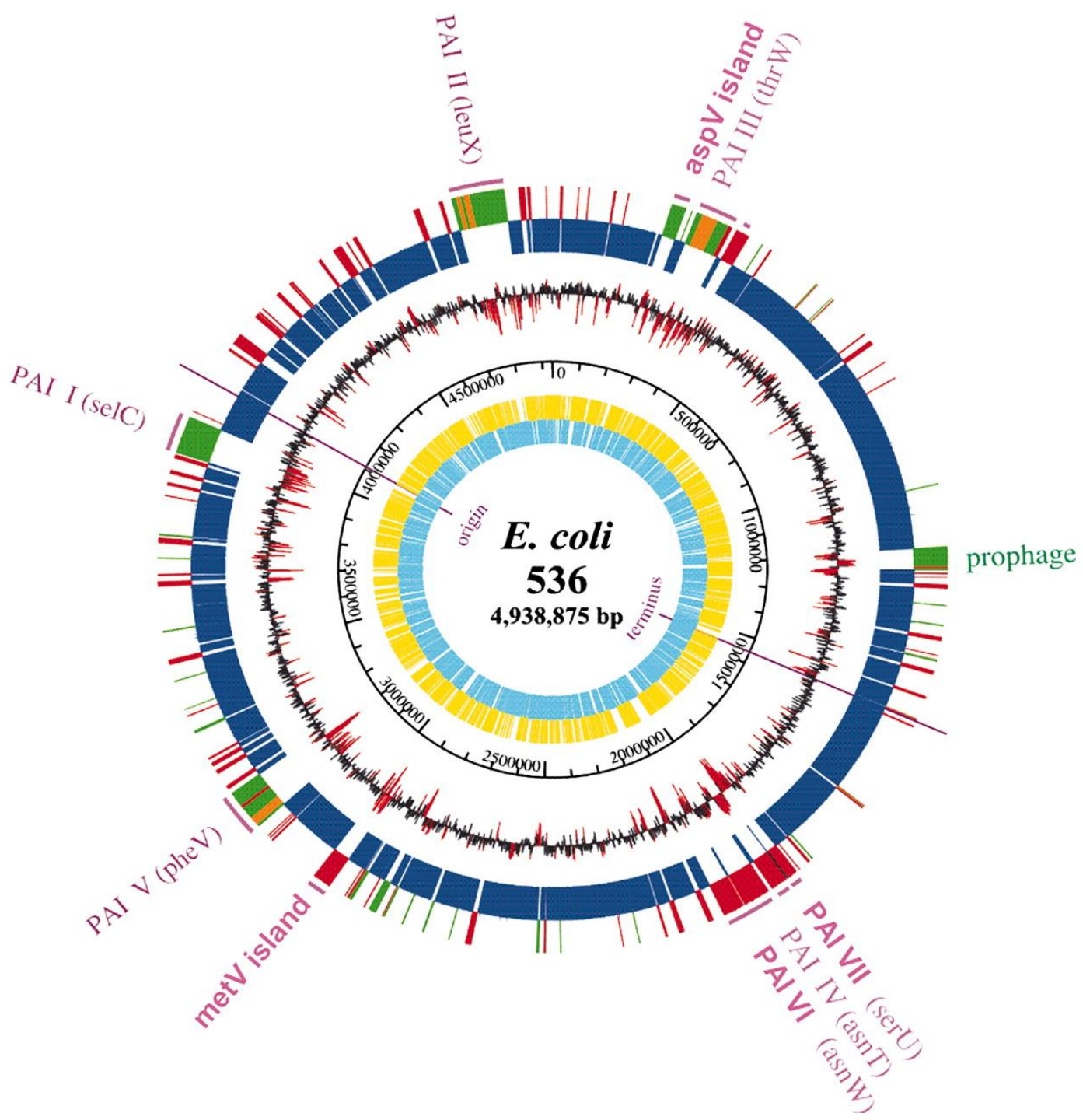
Genomes and Genes: A two minute conceptual view



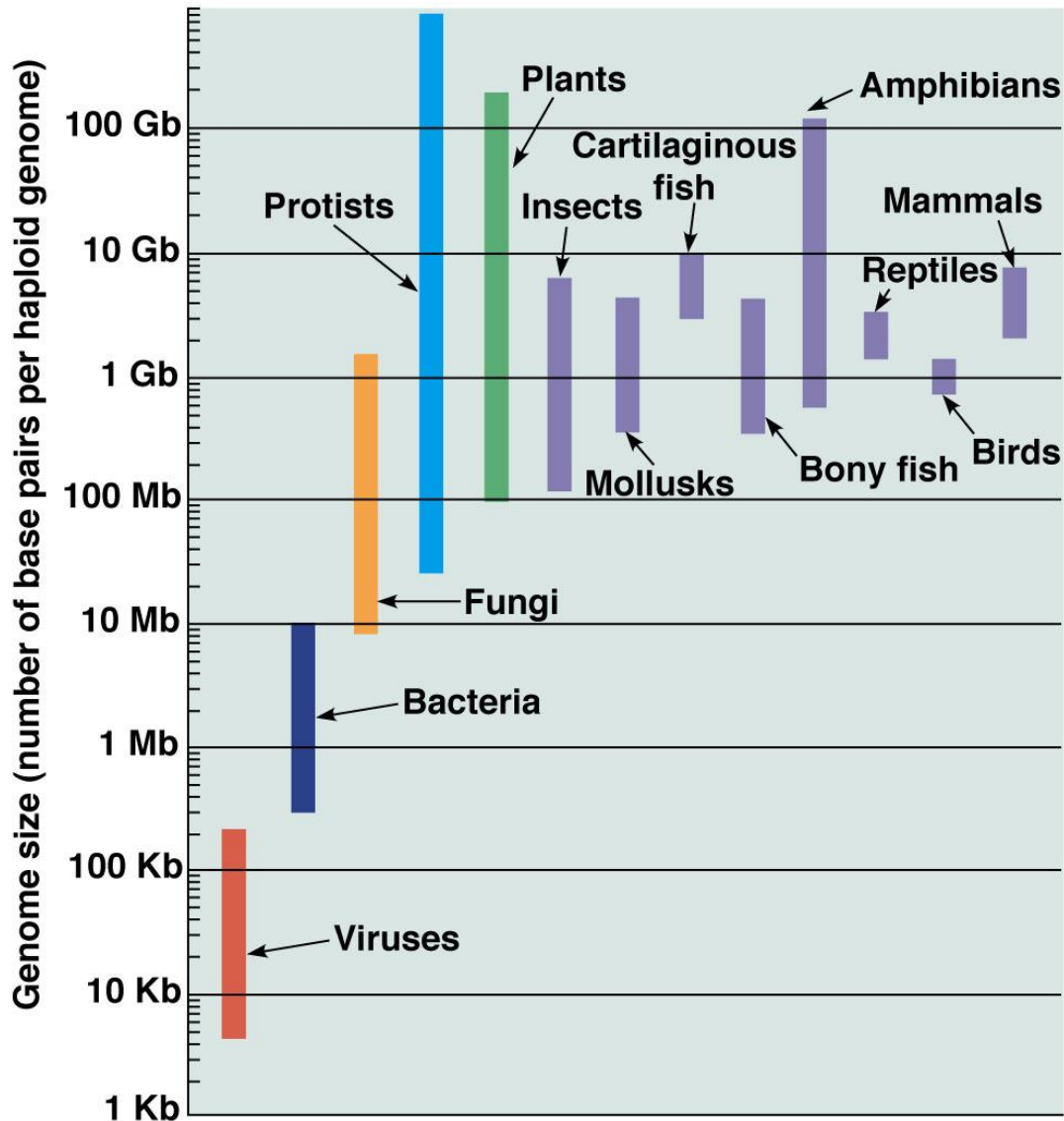
Genomes and Genes: In Real Life







Variation in Genome Length and # Genes



Variation in Genome Length and # Genes

organism	genome size (base pairs)	protein coding genes	number of chromosomes
model organisms			
model bacteria <i>E. coli</i>	4.6 Mbp	4,300	1
budding yeast <i>S. cerevisiae</i>	12 Mbp	6,600	16
fission yeast <i>S. pombe</i>	13 Mbp	4,800	3
amoeba <i>D. discoideum</i>	34 Mbp	13,000	6
nematode <i>C. elegans</i>	100 Mbp	20,000	12 (2n)
fruit fly <i>D. melanogaster</i>	140 Mbp	14,000	8 (2n)
model plant <i>A. thaliana</i>	140 Mbp	27,000	10 (2n)
moss <i>P. patens</i>	510 Mbp	28,000	27
mouse <i>M. musculus</i>	2.8 Gbp	20,000	40 (2n)
human <i>H. sapiens</i>	3.2 Gbp	21,000	46 (2n)
bacteria			
<i>C. ruddii</i> (smallest genome of an endosymbiont bacteria)	160 kbp	182	1
<i>M. genitalium</i> (smallest genome of a free living bacteria)	580 kbp	470	1
<i>H. pylori</i>	1.7 Mbp	1,600	1
Cyanobacteria <i>S. elongatus</i>	2.7 Mbp	3,000	1
methicillin-resistant <i>S. aureus</i> (MRSA)	2.9 Mbp	2,700	1
<i>B. subtilis</i>	4.3 Mbp	4,100	1
<i>S. cellulosum</i> (largest known bacterial genome)	13 Mbp	9,400	1
eukaryotes - multicellular			
pufferfish <i>Fugu rubripes</i> (smallest known vertebrate genome)	400 Mbp	19,000	22
poplar <i>P. trichocarpa</i> (first tree genome sequenced)	500 Mbp	46,000	19
corn <i>Z. mays</i>	2.3 Gbp	33,000	20 (2n)
dog <i>C. familiaris</i>	2.4 Gbp	19,000	40
chimpanzee <i>P. troglodytes</i>	3.3 Gbp	19,000	48 (2n)
wheat <i>T. aestivum</i> (hexaploid)	16.8 Gbp	95,000	42 (2n=6x)
marbled lungfish <i>P. aethiopicus</i> (largest known animal genome)	130 Gbp	unknown	34 (2n)
herb plant <i>Paris japonica</i> (largest known genome)	150 Gbp	unknown	40 (2n)

A Bit More About Genes (≈ 1 course)

- Different species differ in the set of genes they encode and in the exact sequence of each gene
- Homology:
 - Homologues genes are genes that derive from a common ancestor
 - Orthologues genes are homologous in different species (arise via speciation)
 - Paralogues are homologous genes in the same species (arise via gene duplication)
- As species evolve, their genomes (and genes) gradually diverge in sequence
- The closer the sequences of two genes are, the more likely it is the proteins they encode function similarly

Sequencing, Alignment, Assembly (≈ 1 course)

...cacgcttgcagetaccaggagaaaatgaacttttcatcaacttttctagtgtcacttttttgcc...

Replicate

Shred

Sequence

tttcatcaacttt tcaactgggtg tcactttacggg tcaaacccttttg
acttttcatc tgtaccaggagaaa tttcatcaacttt acttttcatc
caggagaaaat tcaaacccttttg tcactttacggg caggagaaaat tcactata

Align

Assemble

aacttttttg
gagaaaactt
aggagaaaac

cgcttgcag
tgcageta

aacttttttg
gagaaaactt
aggagaaaac

cgcttgcag
tgcageta
...cacgcttgcagetaccaggagaaaacttttttgcc...

cgcttgcageta aggagaaaacttttttg

Outline

A Few Updates about Seminar Logistics

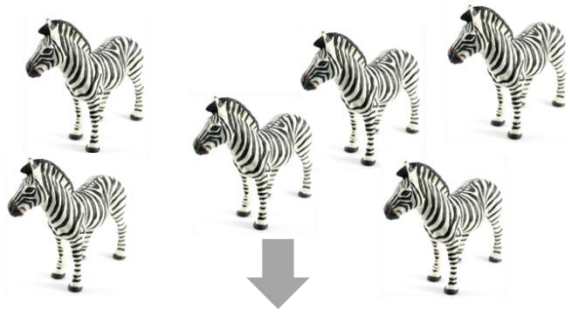
“Metagenomics”

Genomes, Genes, and the Central Dogma

Real Metagenomics

Computational Challenges in Metagenomics

“Genomics”



Real Genomics

...cacgcttgagetaccaggagaaaatgaacttttcatcaacttttctagtgtcacttttttgcc...

Replicate

Shred

Sequence

tttcatcaacttt tcaactgggtg tcaactttacggg tcaaacccttttg
acttttcatc tgtaccaggagaaa tttcatcaacttt acttttcatc
caggagaaaat tcaaacccttttg tcactttacggg caggagaaaat tcactata

Align

Assemble

aacttttttg
gagaaaactt
aggagaaaac

aacttttttg
gagaaaactt
aggagaaaac

cgcttgag
tgcageta

cgcttgag
tgcageta

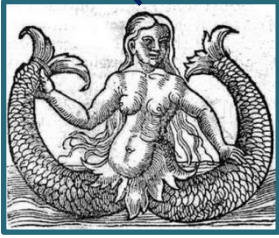
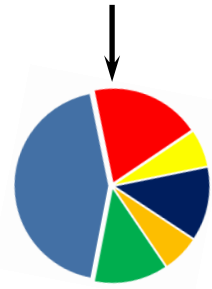
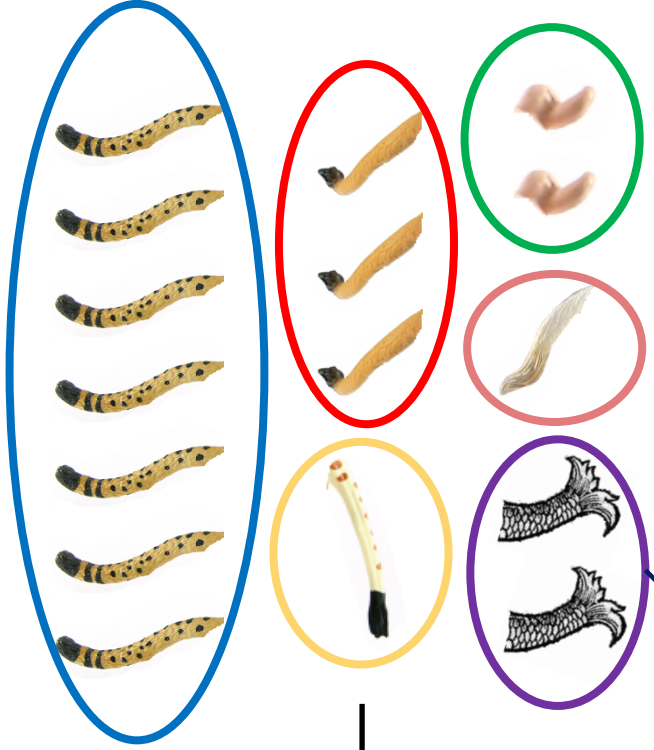
...cacgcttgagetaccaggagaaaacttttttgcc...

cgcttgageta aggagaaaacttttttg

“Metagenomics”



Who's there?

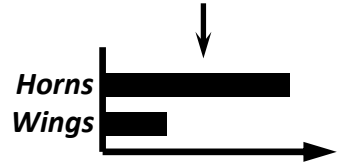


What are they doing?

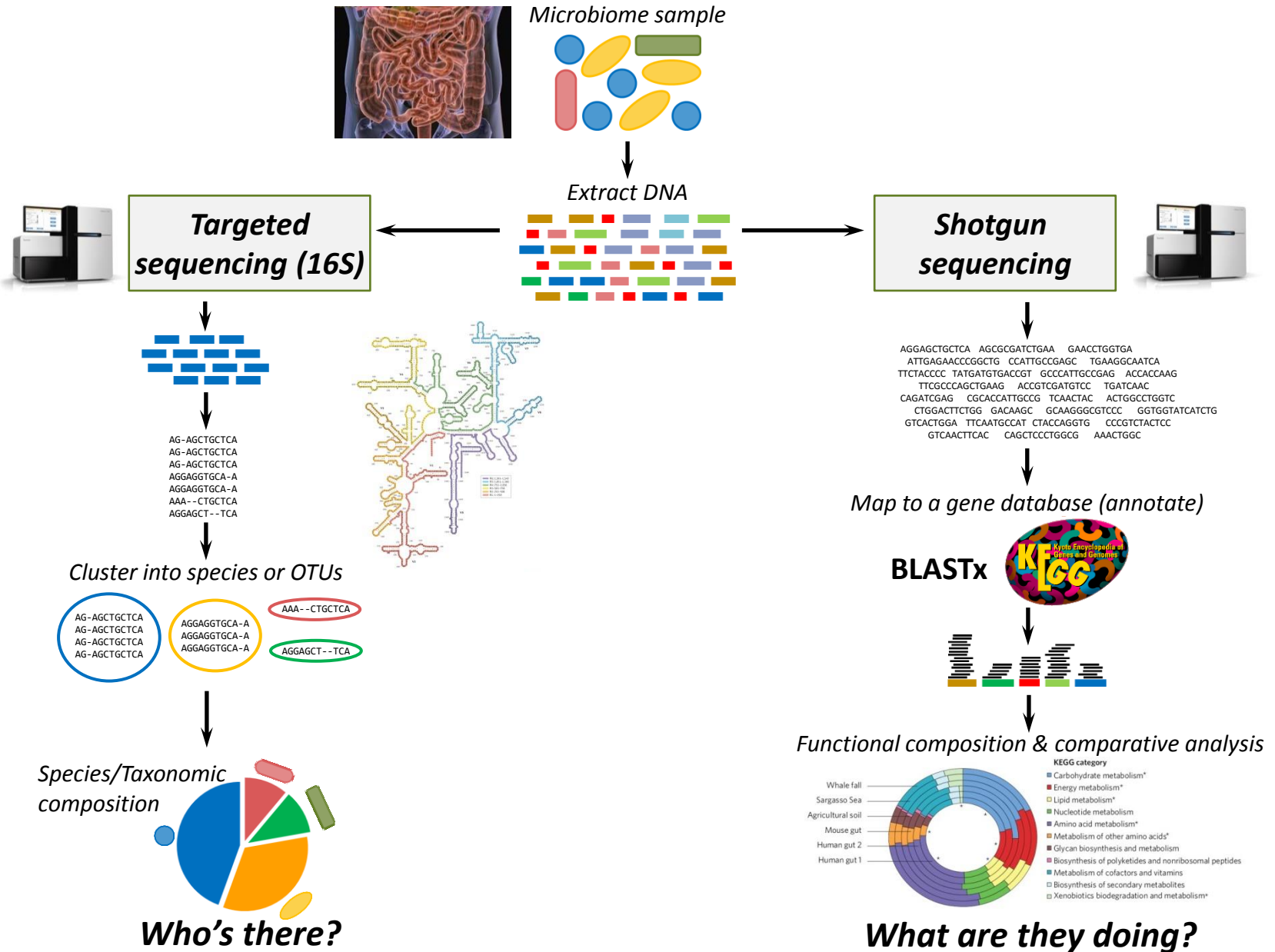
Horns



Wings



'Real' Metagenomics



Outline

A Few Updates about Seminar Logistics

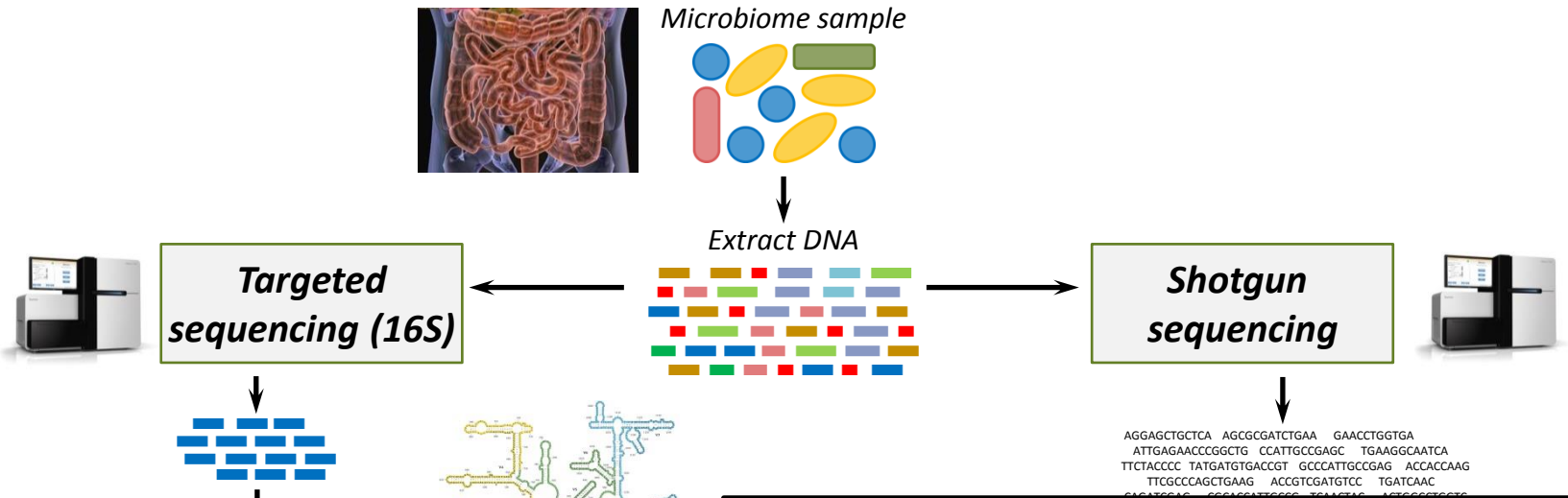
“Metagenomics”

Genomes, Genes, and the Central Dogma

Real Metagenomics

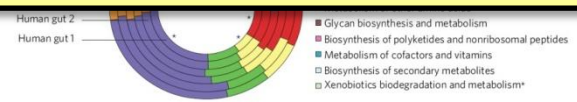
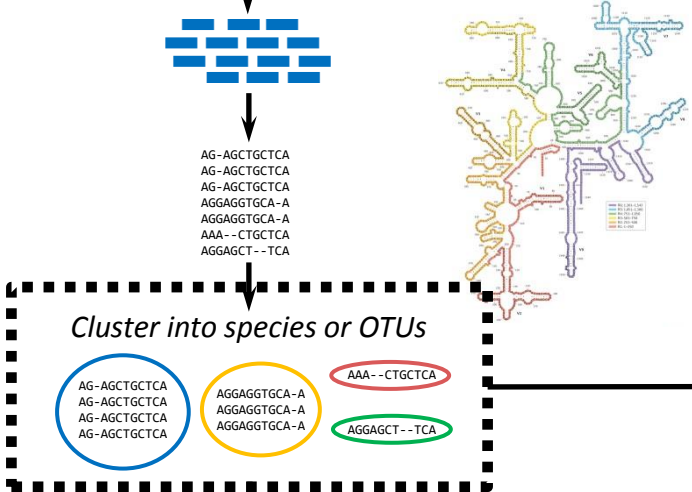
Computational Challenges in Metagenomics

Key Challenge 1



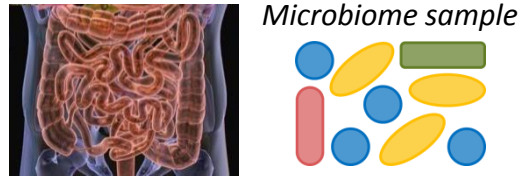
From 16S Sequences to Taxa Composition

- Clear clusters are not always feasible
- True variation vs. sequencing errors?
- Low resolution
- A fixed threshold doesn't always reflect the same phylogenetic closeness
 - Different species clustered together
 - Strains clustered separately



What are they doing?

Key Challenge 2

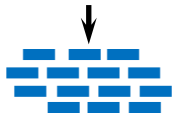


Extract DNA



Targeted sequencing (16S)

Shotgun sequencing



AG-AGCTGCTCA
AG-AGCTGCTCA
AG-AGCTGCTCA
AGGAGGTGCA-A
AGGAGGTGCA-A
AAA--CTGCTCA
AGGAGCT--TCA

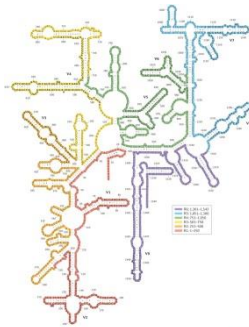
Cluster into species or OTUs



Species/Taxonomic composition



Who's there?

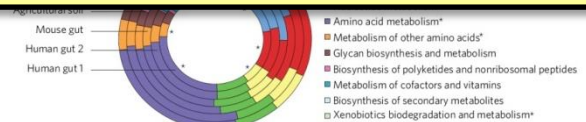


```

AGGAGCTGCTCA AGCGCGATCTGAA GAACCTGGTGA
ATTGAGAACCCGGCTG CCATTGCCGAGC TGAAGGCAATCA
TTCTACCCC TATGATGTGACCGT GCCCATTGCCGAG ACCACCAAG
TTGCCCCAGCTGAAG ACCGTCGATGCC TGATCAAC
CAGATCGAG CGCACCATTGCCG TCAACTAC ACTGGCCTGGTC
CTGGACTCTGG GACAAGC GCAAGGGCTGCC GGTGGTATCATCTG
GTCACCTGGA TTCAATGCCAT CTACCAGGTG CCGCTCTACTCC
GTCAACTTCAC CAGCTCCCTGGCG AAATCGGC
    
```

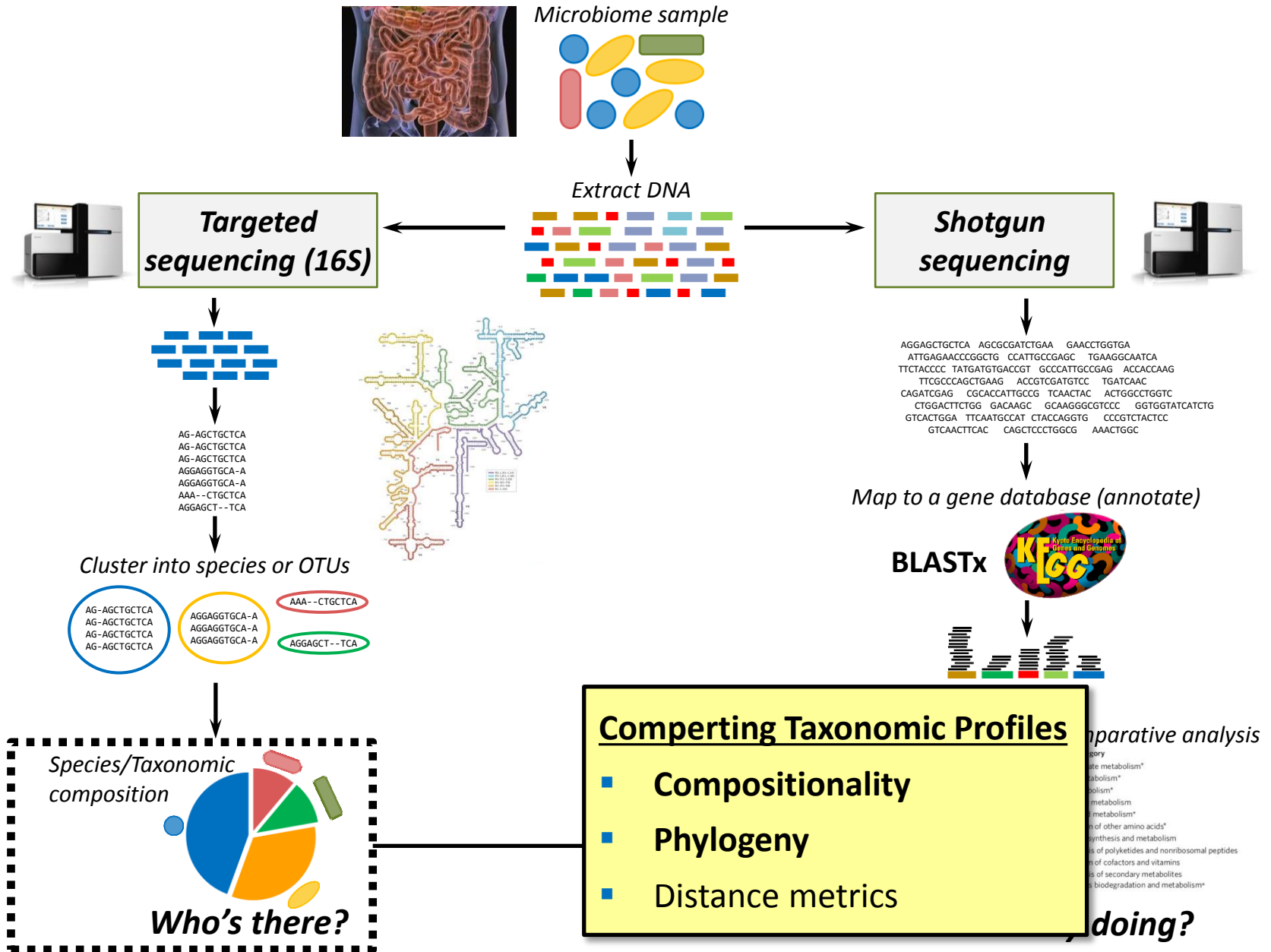
From Shotgun Sequencing to Species- and Strain-Level Profile

- Different species encode different genes
- Low abundance strains
- Shared variation



What are they doing?

Key Challenge 3



Key Challenge 4

