


Shotgun Metagenomic

SLU intro to genomics course 2020

 @metamoritz

 @moritzbuck

 @metamoritz

Before the computer

- **Get sample**
- **Extract DNA**
- **Prep library**
- **Send to Sequencing facility**
- **????**
- **PROFIT**

- **Load and load of reads**

3

Dream-quest of unknown metagenome

From an environment with the power of genomics we want to know:

- Who is in the environment? (taxonomy)
- How many of each? (abundance)
- What can they do they? (genetic potential)
- What are they doing? (expression analysis)
- What do they eat? (metabolism)
- Where do they come from? (evolution and ecology)

Two main roads

Assembly based :

- Assembling
- Binning
- Annotating
- Slow, tricky, more specific

Assembly-free :

- map/classify reads
- Direct annotation
- Heavily database reliant
- Faster, easier, more noisy



The tutorials

- **Let's assemble some metagenome!**
- **And get some genomes out of it!**
- **Check how good they are!**
- **And then let's compare some metagenomes directly assembly free**

Assembly based approaches

Turning the reads into genes and more! (aka The Contig):

- **Full length coding sequences**
- **Pieces of genomes**
- **Or anything else made of DNA**

Increases specificity of previous approaches, but might lose some sensitivity (e.g. not everything assembles)

Linking reads to contigs

The dark art of mapping.

Normally, quantitative information gets lost during assembly, we need to “map back”. Read based approaches “map” to databases

Mappers :

- **bowtie/bwa/bwa2**
- **Bbmap**
- **Squeeze...**

RNA/DNA euk or prok, all different!

=> pile-ups!

The gene atlas approach

If you only care about genes:

- **Predict genes and annotate genes:**

- Functionally
- Taxonomically

- **Gene clustering:**

- By similarity
- Into COGs
- Into families



- **Quantify and analyse diversity**

The king discipline

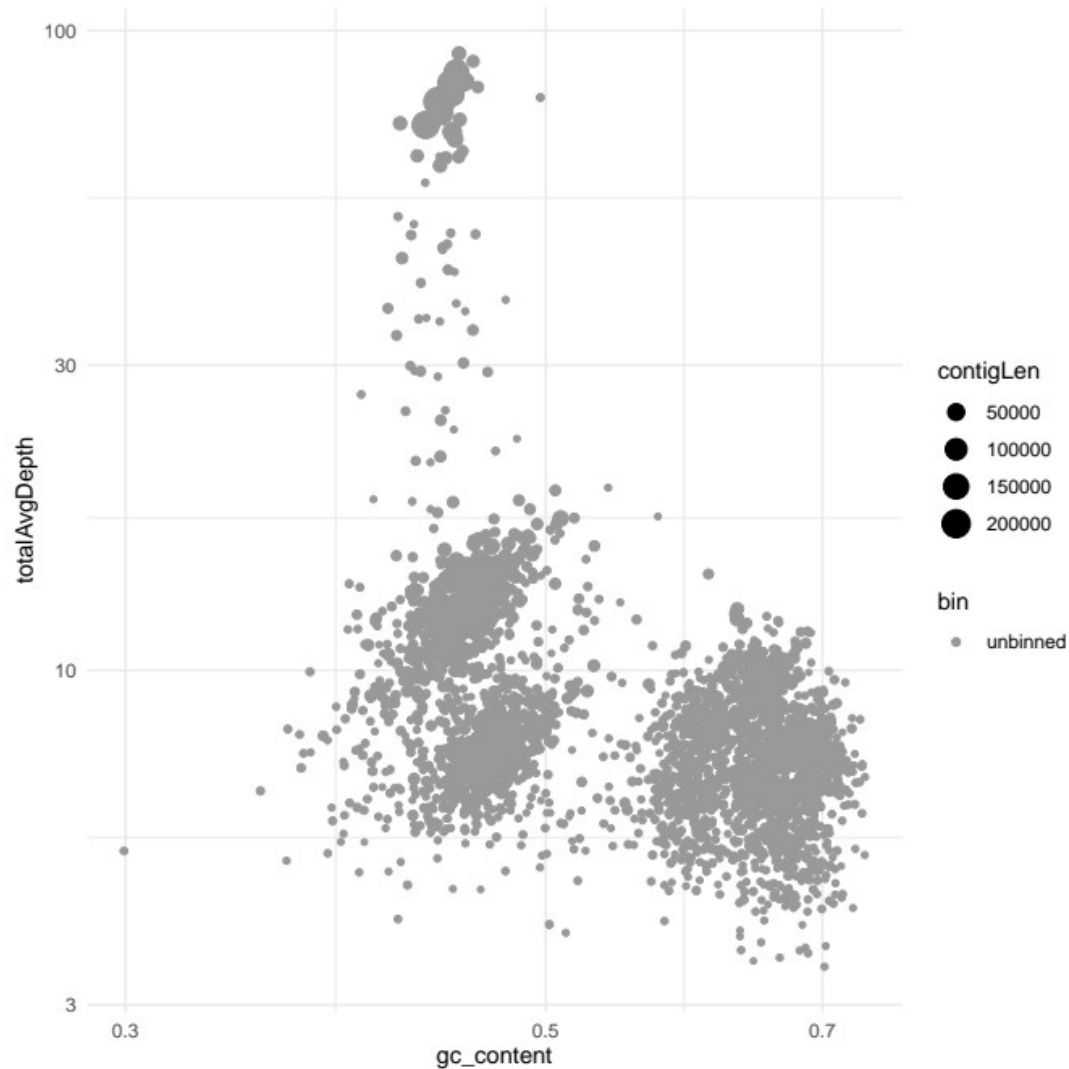
Genome-resolved metagenomics:

- **Cluster contigs into bins**
- **Identify possible genomes**
- **Quantify completeness/redundancy**
- **Assign taxonomy**
- **Have fun?**

So vamos!

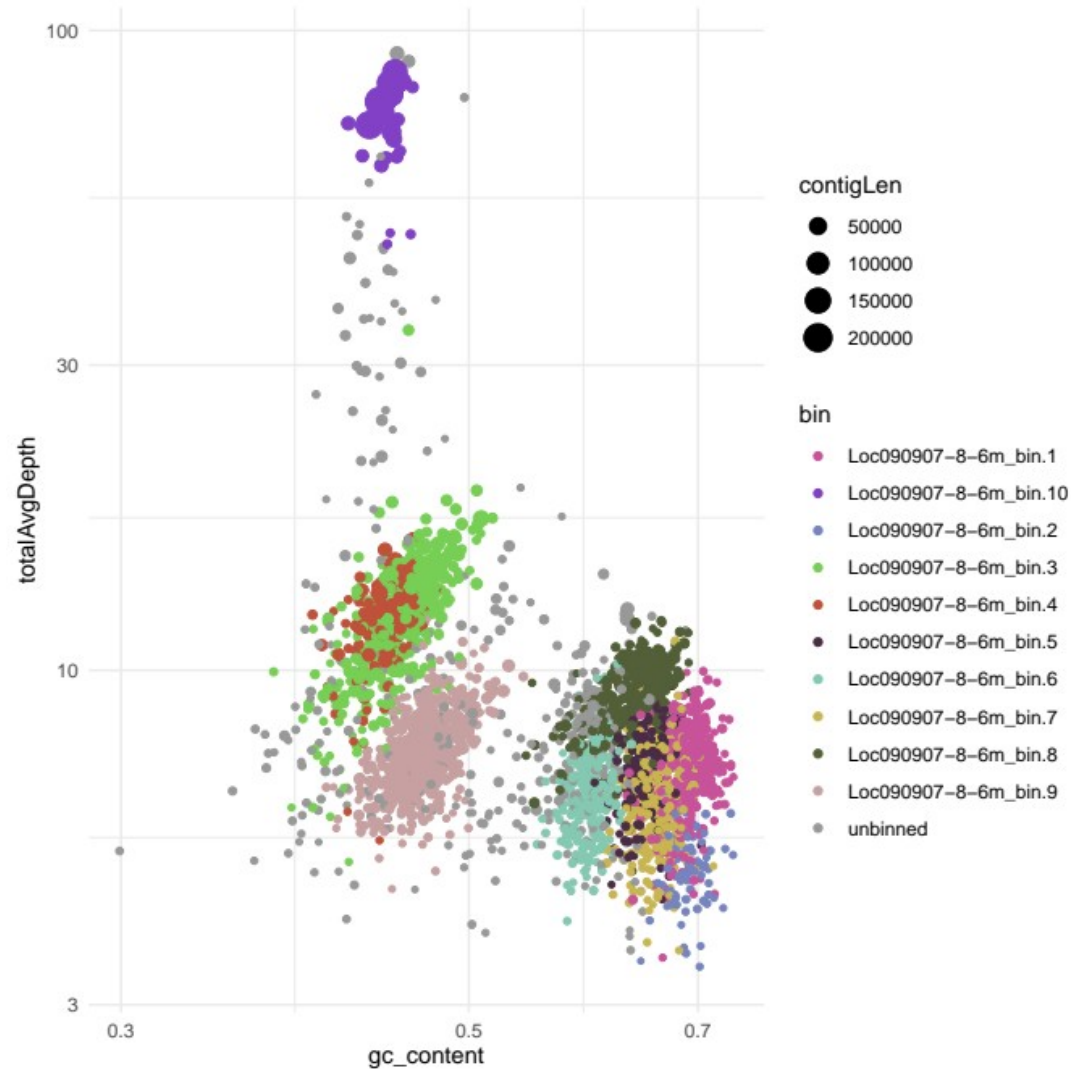
Clustering contigs

Binning:



Clustering contigs

Binning:



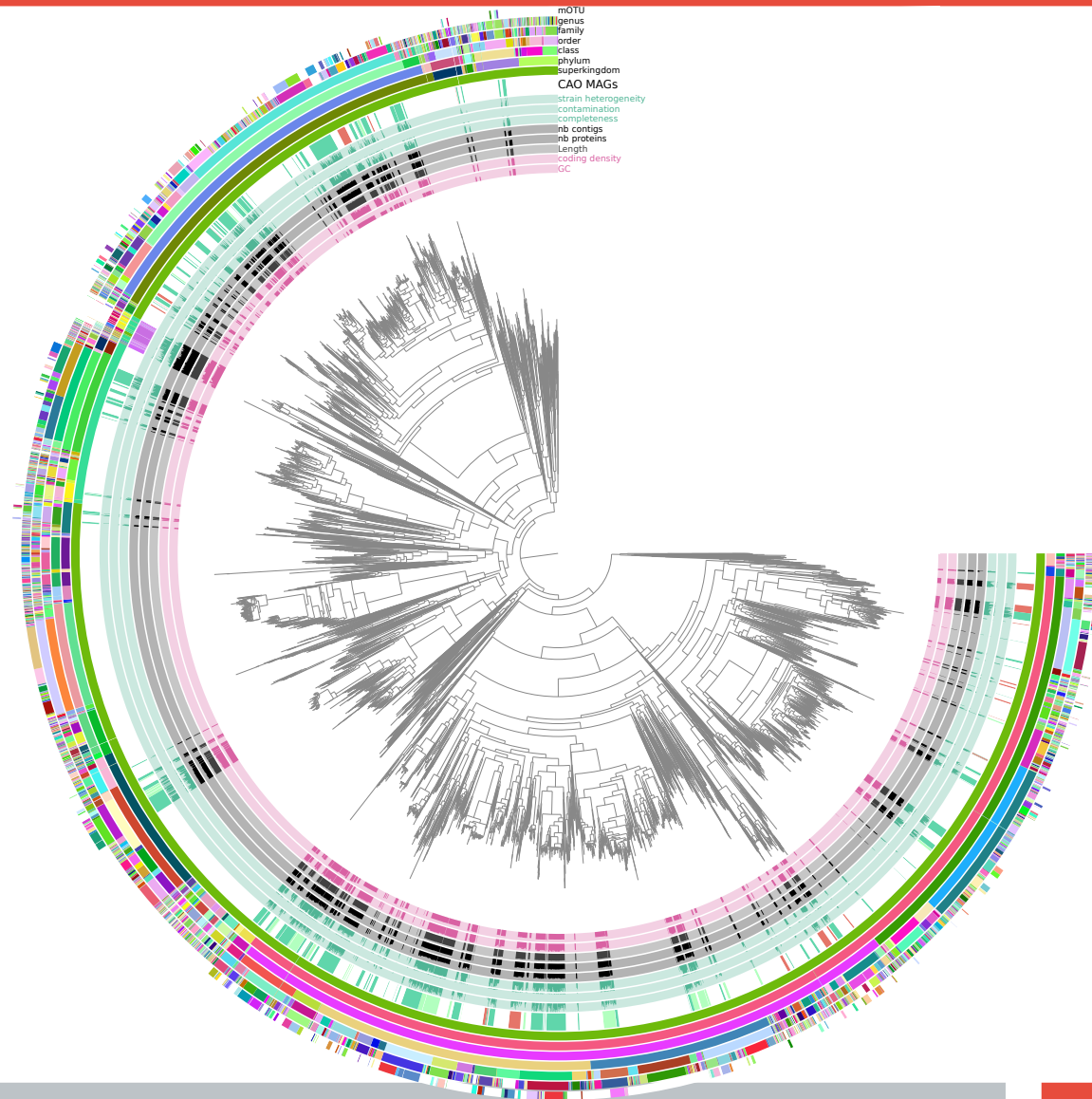
What is a genome?

How to recognize a genome:

- Size and coding-density
- Similarity to known genome
- Universal single copy genes

```
>k141_3
GCCATCCGAGGTGGGCCGGCCGCTCCGCGGGCGGCGGATTGGAGAGAGAGGCGTGGCCGG
CGAAGCCGCGCGCGGAGCGCACGGCCACCCGGGGACACCGGACCCTGCACGAACTGAG
ATGCGCCACCCTCACGCAGCGGGCGACCCGTGGCGAAATCAGCGATCGAACGGCGGGCG
TTTTTACGCGCCCATCTCAGGTGCGGCGTTGGGAAACCGCCGGTCTCTCCGCCGCT
CACAGCGGCCGTTGGAGAGGACGACGTCTCCACGAACAGCGTCTGCGCCTGATTGTGG
CCGTTGTTGGCGGTGAGTCCGACCTGCATCCGGTCGTAGACCGTGCGGGAGCGCGGCAGC
GTTTTGCCGGCTTCGTGAGCACCTTGACGCCGTCTGCCAGACTTGATCCGGCCGCTCC
GCTCCTTCGAGAGGCGCAGGTGGATTTTACCTCGACCCACTTGTCTTGGGGAACCTC
GGCGCACCGACGACCGCGCGGAATTTTTCGCCGTCCACCCTTGCCAGATCGGAGGCG
ACCATTTGCGCGTTCTGCAGAAAGAGGCGGCGCGGGCGTGTGCGCAGCGACGTGAC
TCCAGATCCCAGAGAAAAATCGAGGACGATTCCGTGCCGCCTTCCAGCCAGAACCGGCG
CTGAACCACACATGGTCGCCTTTCACGAACCGCAATCCCTCGCGCTCGATATCTGCCTCG
GAGGCCCCGCGGCCGTGCAACGGGGCGCGACGCATTTAGAGCGTTGGTTCCCGAGTGA
ACCGGTTGCGTGGTCAGCGTAACCGAATTGGACGCGAGACCCGTCGTGGACTCGATTTCC
AGGCCCTGCCACCGGCTCGAATCCGCGGGGAACAAATCTTCGGTCTTCGCCGAACCTTCG
AACCATCGCGAAAGACGAATGGCGCTGCCGGGGCGCGCAGCGGCCGCGCCACCCAAAAC
CACCACGCGCGCCCCGCACTCACGCCCAAGGCGGAAATGCTCCCCGGCGACCGACGCAA
CTACCCAGCGGGCGTTGCTCGCGGACTCGGGGAGGGCGGAGTGGTCGGCGTCATCAGG
GTATGGGGTCTTAGGAACAAAAGCGCTGGGCGGAAATCGCAACAAGTGTCTGCGGGC
CATCGCAGCTCTTTCTGTTTTCCCGGCTTGCGCCGCGCGCGGCTCTTTGCGGCCATGTCT
CACGCATGAGAATACCGGACTCGCCCTCGTTCGCGCGCGACCGCGCTGACCTCCCGC
AGGTCACTCCCGAGCTCCTCGCGTCCGTTCTCGCCCGTTACTCCCGCAGCAATCAGGGCA
TCGCCGCCATCCTGTCCAAGGTCGATCTCGCCAACCCCGACGCTTCGATCGACCGCATCC
TCAGTTTCGTGACTACGGCCATGCGTCCATCGGCGGCCTCACCGGCGGTCTCGCGGTGG
CCCTCGACGATGTGTGATGTGGCTGGCCTACAAAATTTTTGAAATCGCCACCATGGCCG
ACGGCCAGGAGTCGAGCACCCGCTACATCAGATGGACGCGGCGAACCTCCCGACCGCCG
CCGAGTTGGGGATCCCCGACGATCTGGCGCGCGCTGGTCCGCTGTGATGGCCAAGTCTCT
TCGCCGCCTACCACGCGGAGTATACCCGGCTCGATGCCCTCGCGACCGCCACCCCGGCC
TCGTCCGCCTCCCGCCGACGCCAAACCCGCGCTCGTCACGCGCTCCGCAAGAACTACG
CGCTCGACCGCGCCCGCTACTTCATCCCGCTCGCCACGCGCACCAATGTGCGGCTCGTGC
AGTCGTGCGCATGTGGGCGATGACGGTGAAGCACCTCGATTGCTGCCGACCCCGAGG
CCCGGGCGGCGCGCGGTGATCCGCGACGAACTGCTGAAGCTTTCCCGCGGCTGATGC
GCCACAGTTCGGCGGAACAGTCTTACACGGCCAGGCGGCCAAGAGCTGGCGACCTCGT
```

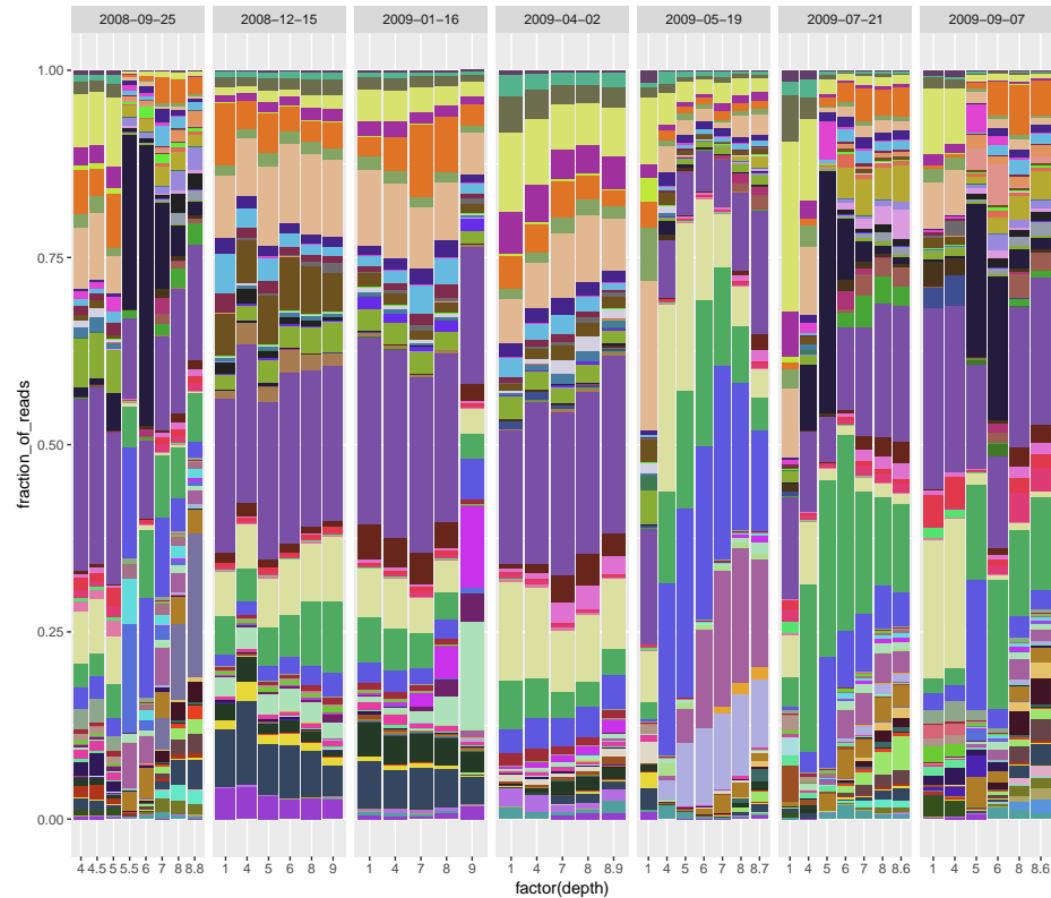
Putting it into phylogenomics context!



Quantifying it

Mapping read to the bins:

- Statistics
- Diversity metrics
- Network analysis



What are those reads?

The inventory approach:

- **Annotating the reads directly**
- **Heavily relying on databases**
- **Good for well known microbiomes (western European poop, agricultural soils)**
- **Dangerous for others as methods tend to overclassify**


What are those reads?

- **Taxonomically :**

- Kraken
- Kaiju

- **Functionally :**

- Blast ...
- Megan
- Humann

Domain	Eukarya	
Kingdom		
Phylum		
Class		
Order		
Family		
Genus		
Species		