

# De novo metagenomics assembly

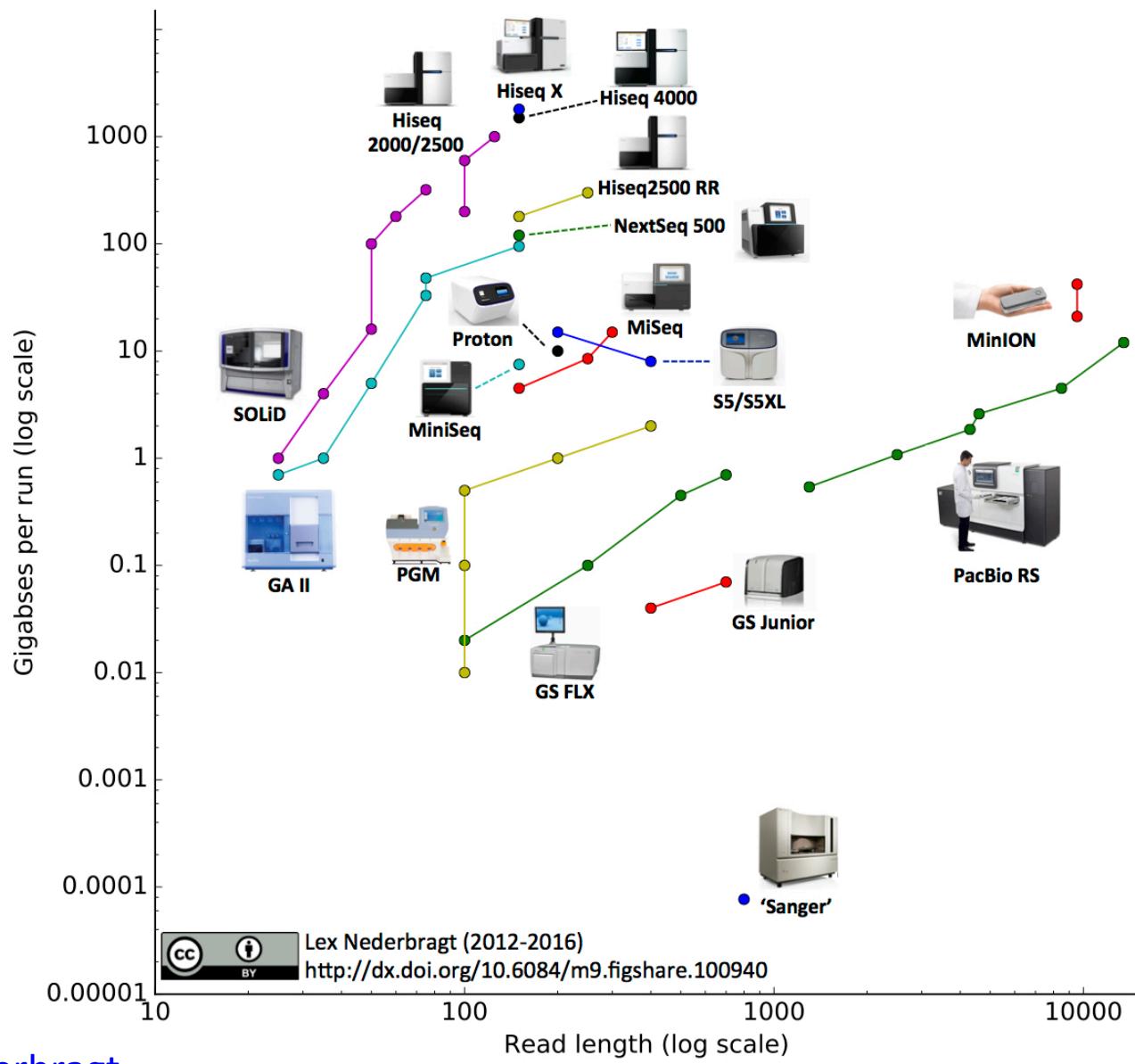
Christopher Quince  
Warwick Medical School

# Introduction

- What is de novo assembly?
- Why assemble?
- Reference based assembly simple but limited:
  - Map reads onto known genome
  - Requires closely related reference
  - Cannot find novel genes, plasmids, rearrangements

# Overview

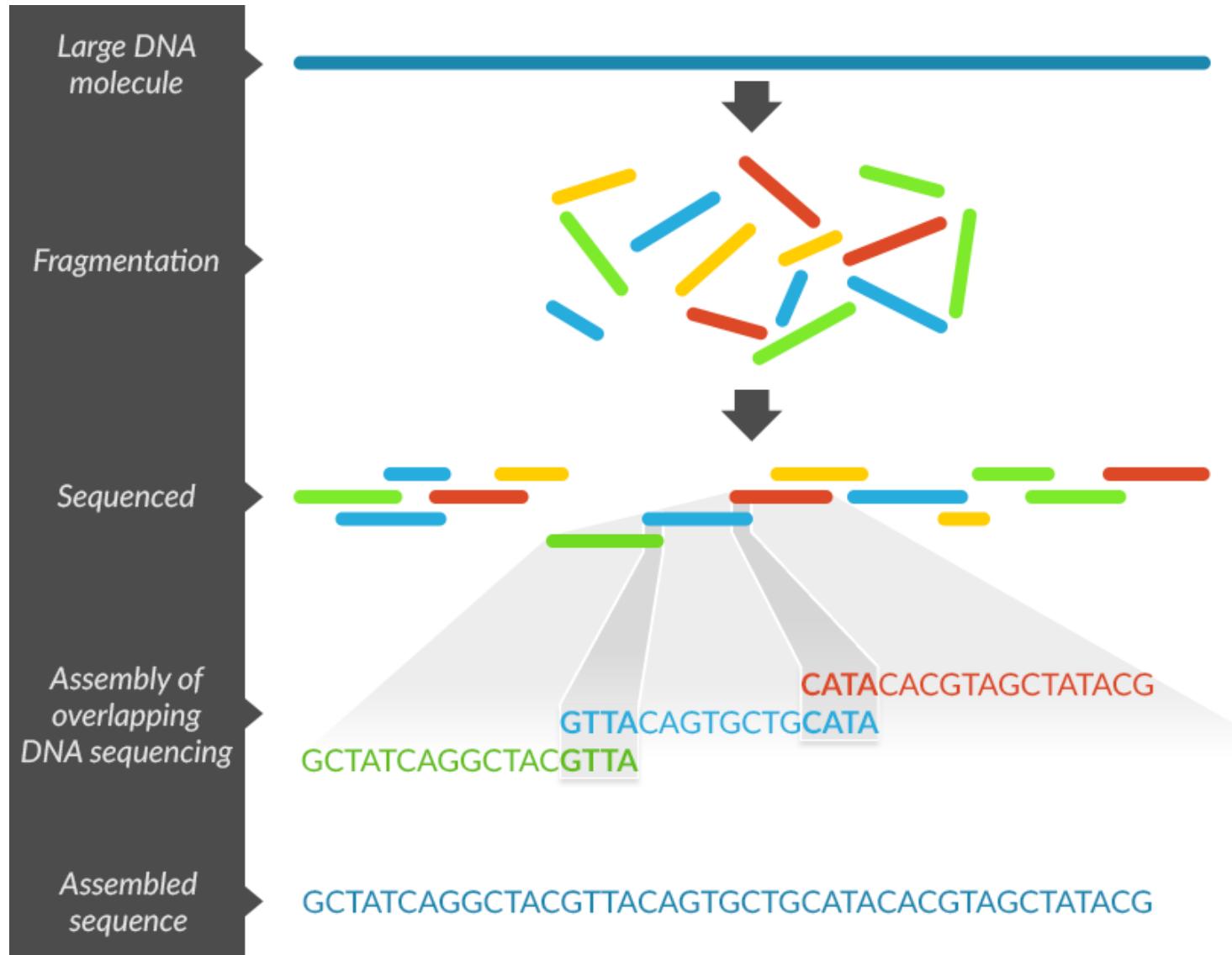
- Current state of sequencing technologies
- De novo assembly paradigms
- Coverage
- Repeats
- Metagenomics assembly:
  - Coassembly
  - Strains
- Assembly software



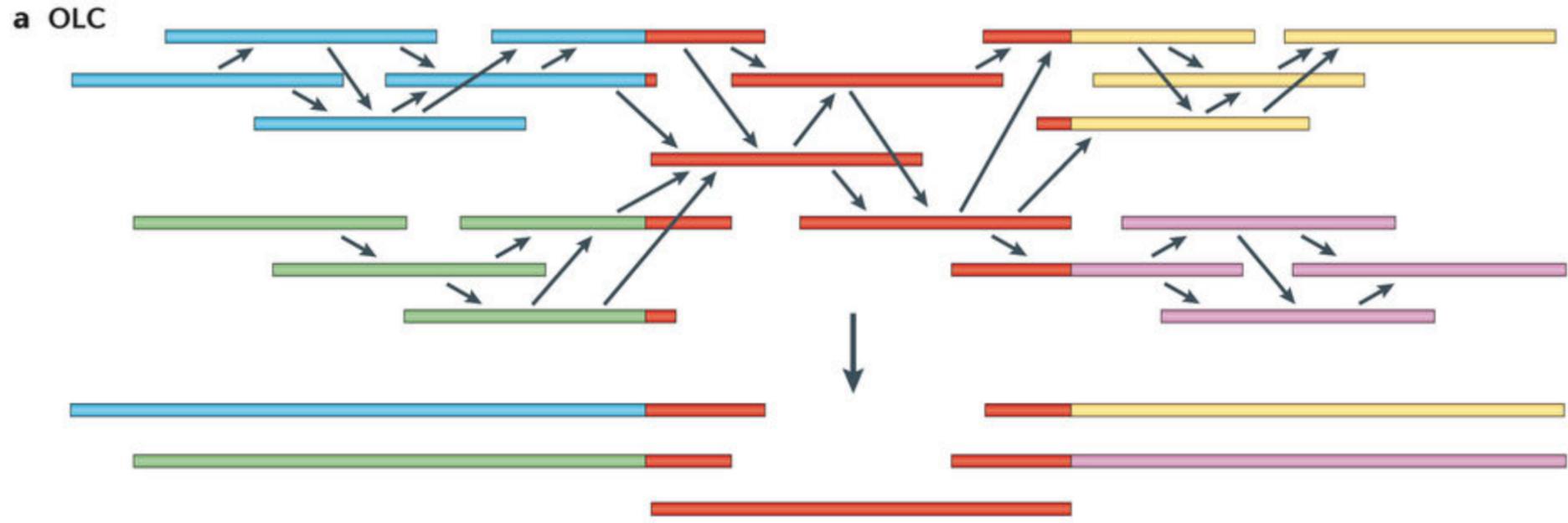
Lex Nederbragt

[https://figshare.com/articles/developments\\_in\\_NGS/100940](https://figshare.com/articles/developments_in_NGS/100940)

# What is de novo sequence assembly?



# Overlap layout consensus



- OLC slow because of pair-wise comparisons
- Renaissance with long read technologies

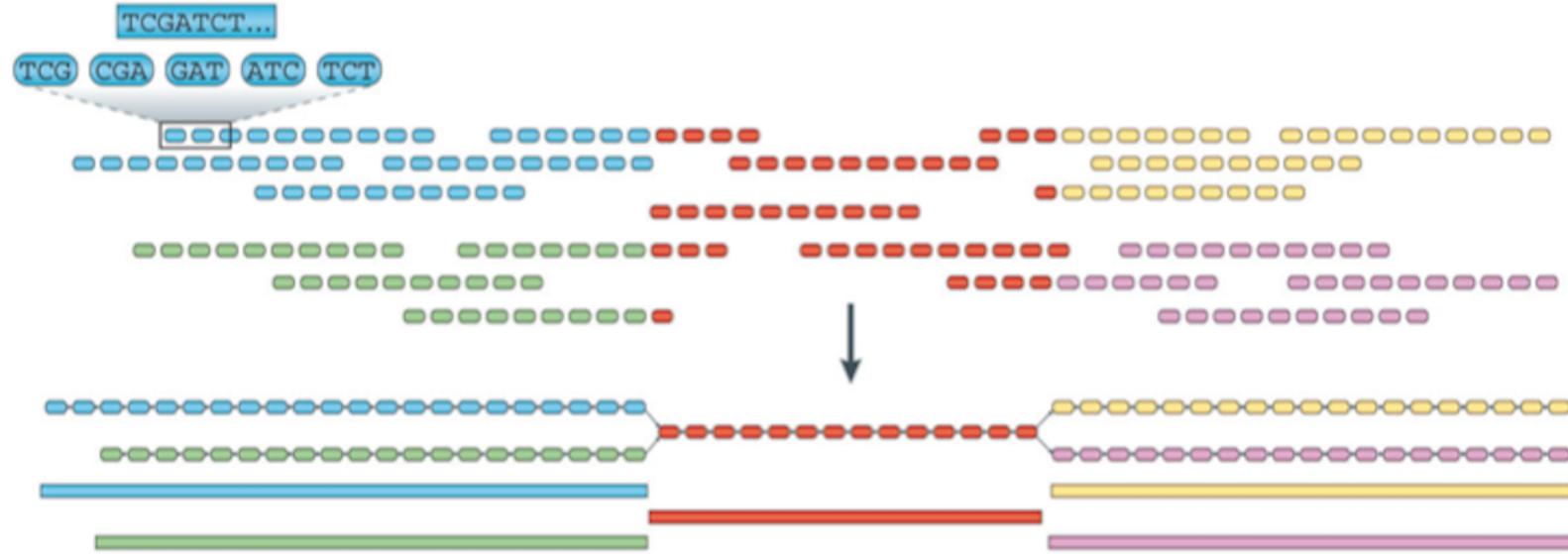
Genetic variation and the de novo assembly of human genomes

Mark J. P. Chaisson, Richard K. Wilson & Evan E. Eichler

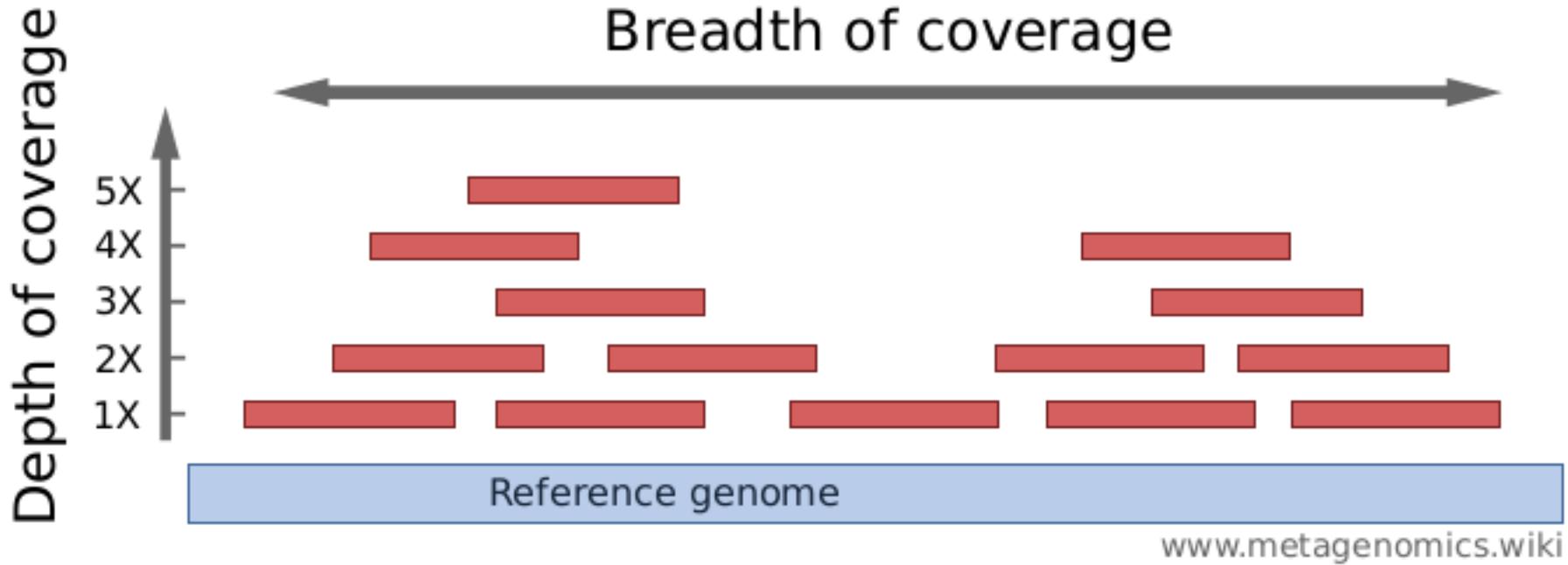
Nature Reviews Genetics 16, 627–640 (2015) doi:10.1038/nrg3933

# De Bruijn graph assembly

b de Bruijn



- Fast but effectively fixed length exact overlaps
- Still default for short read next generation



Coverage depth of taxa n

$$C_n = \frac{\rho_n RL}{G_n}$$

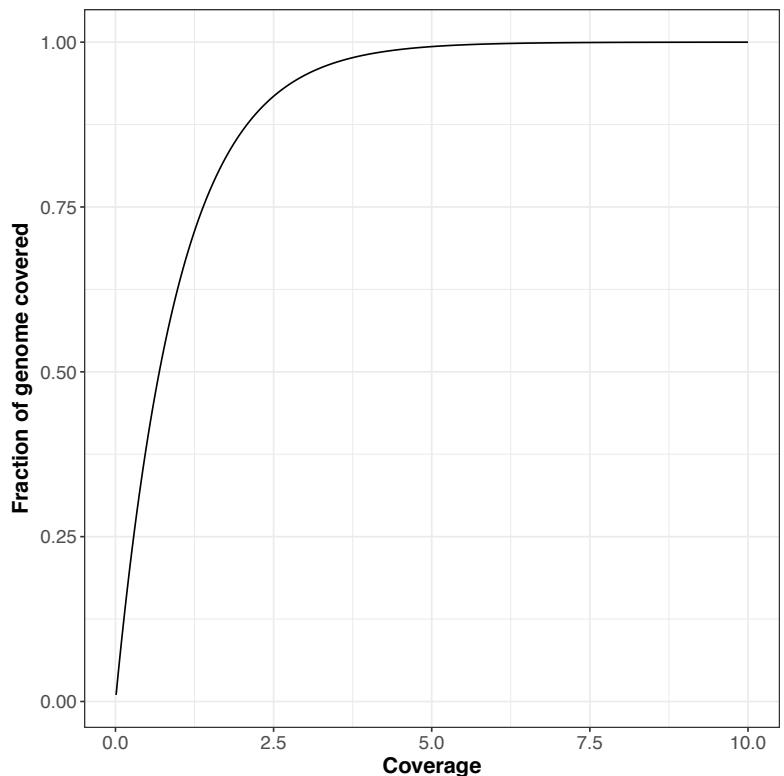
- Where  $\rho_n$  is genome relative frequency of taxa n
- R number of reads in library
- L read length
- $G_n$  is genome length of taxa n

# Lander-Waterman statistics

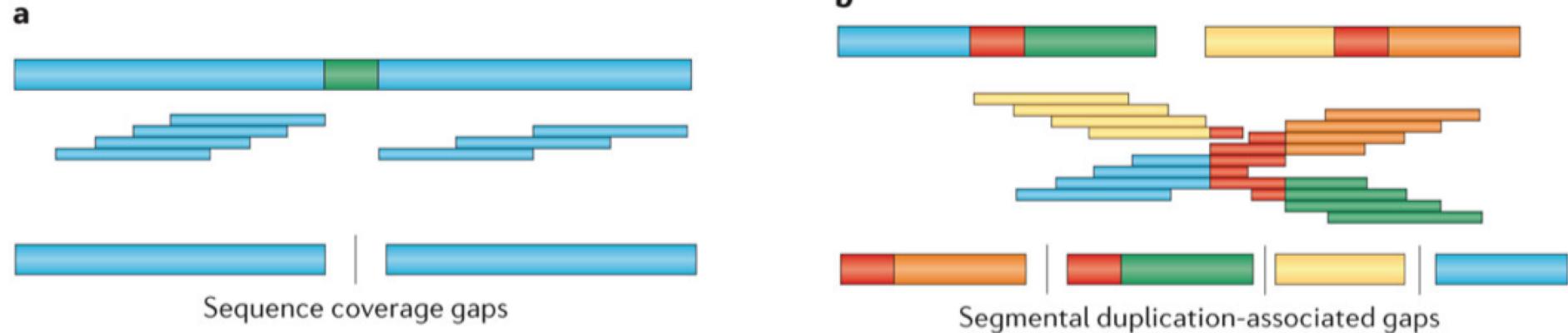
- Assume reads randomly distributed across genome
- Fraction of genome covered:
- Mean number of contigs:

$$Re^{-c_n}$$

$$1 - e^{-c_n}$$



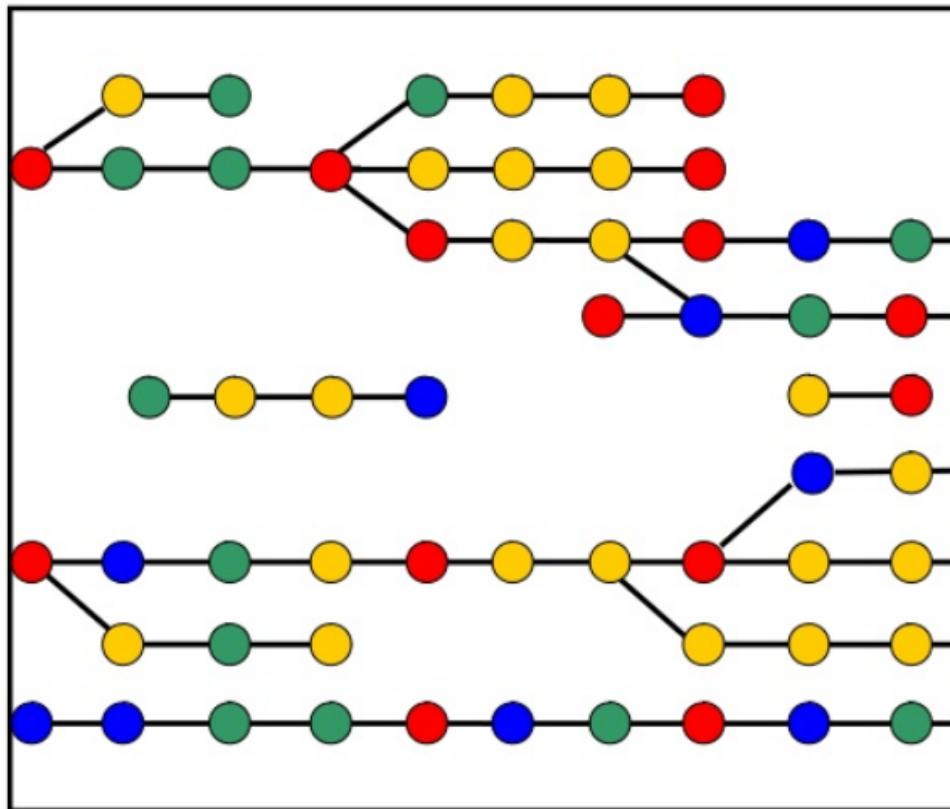
# Gaps in assembly



- A contig is just an unambiguously assembled genome fragment
- What is the optimal kmer length for de Bruijn graph assembly?

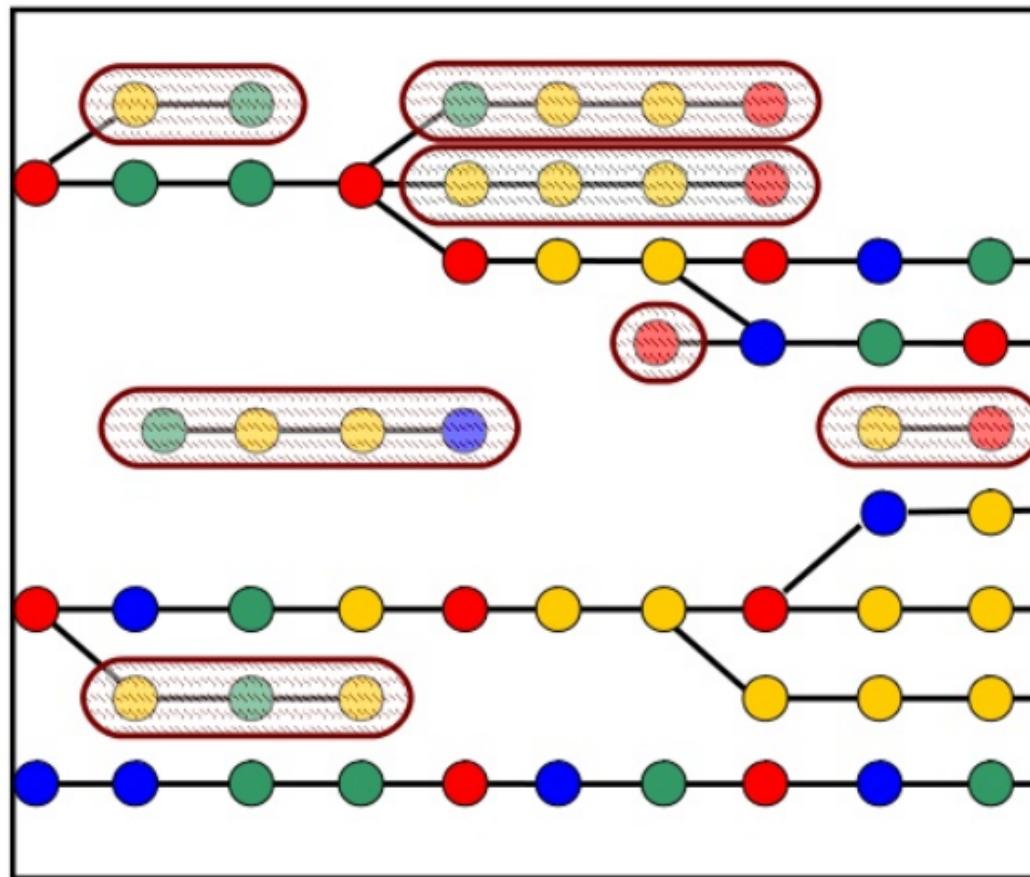
# Genome assembly in practice

- Read errors generate tips that are pruned:



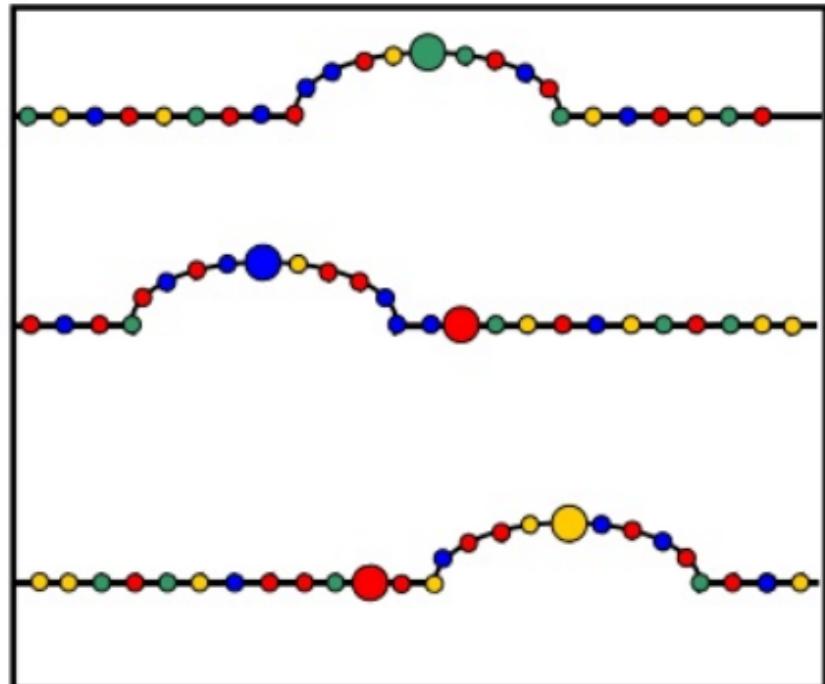
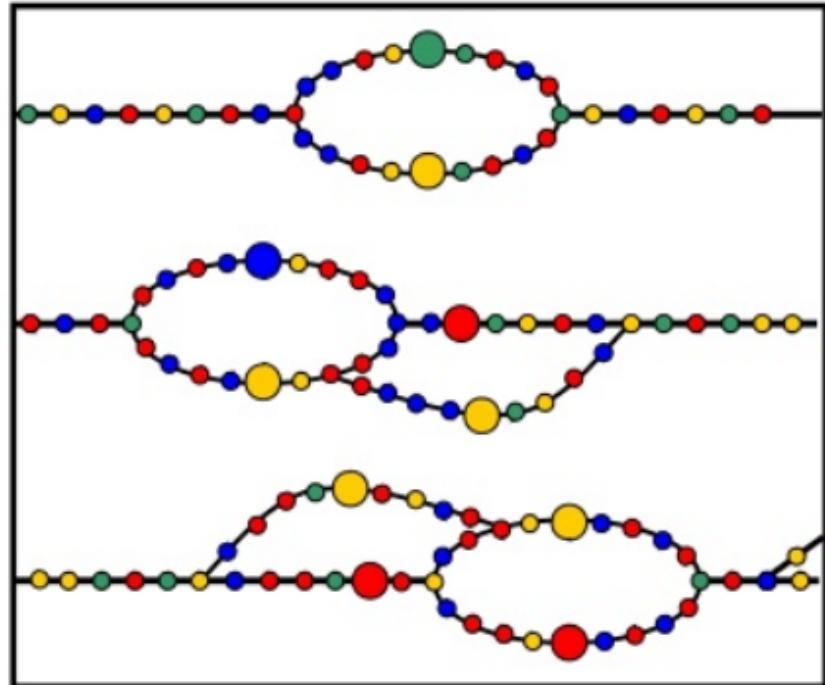
# Genome assembly in practice

- Read errors generate tips that are pruned:



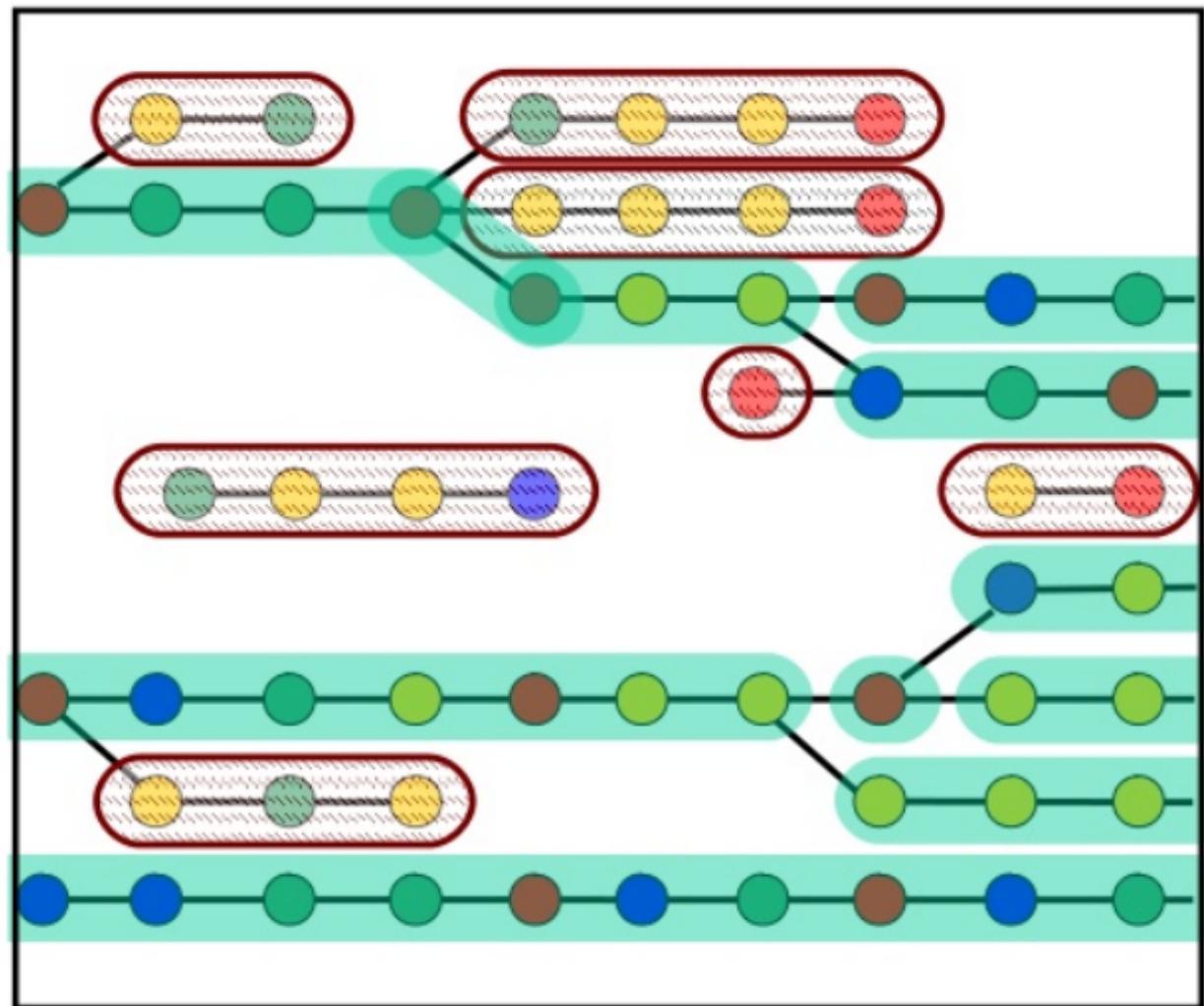
# Bubbles

- Errors can also generate bubbles but these can also be real variants
- These are popped
- **Removing** real variants from assembly



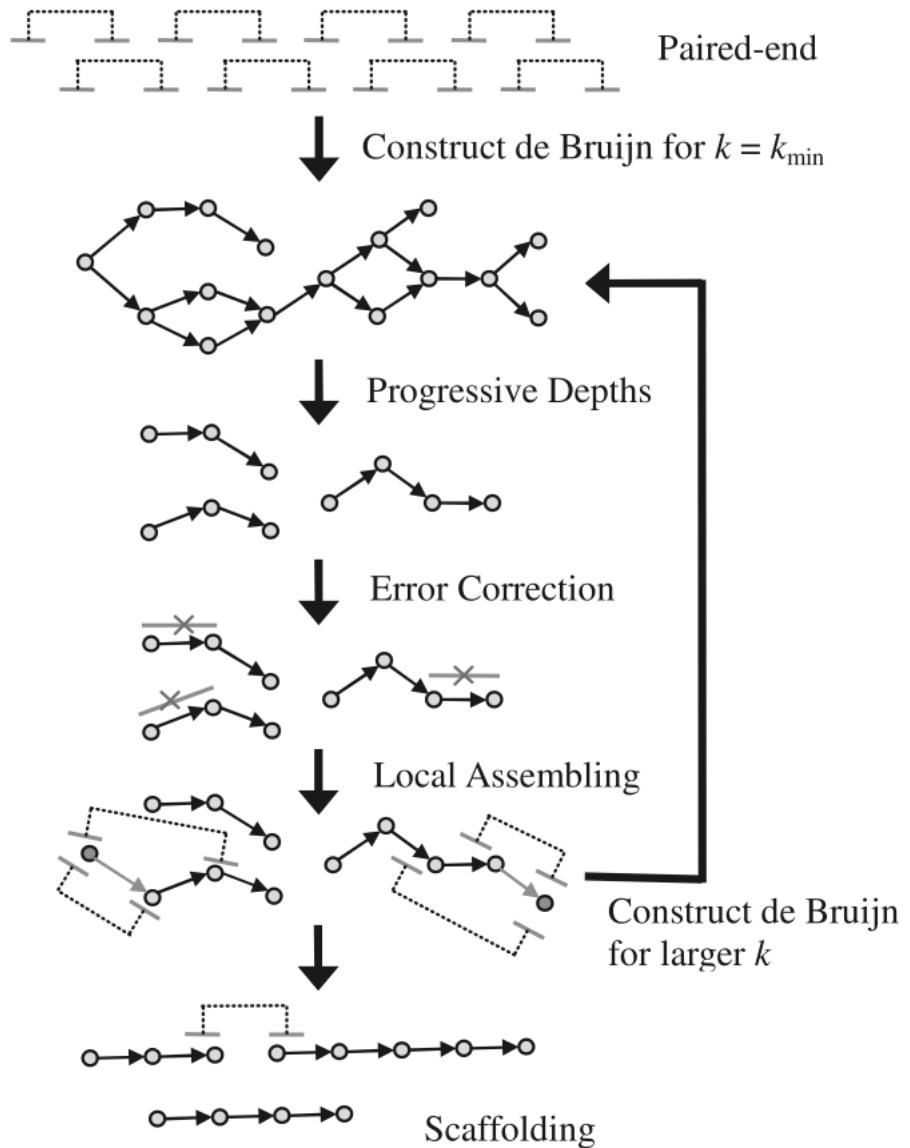
# Contigs

- Remove ambiguous edges  
output linear portions of graph:



# Variable k-mer length dBG assemblers

- Small  $k$  more branches
- Larger  $k$  more gaps
- Approach pioneered in IDBA is to iterate through  $k$ mer lengths to get best of both
- Used in most dBG assemblers now: idba, idba\_ud, megahit, minia



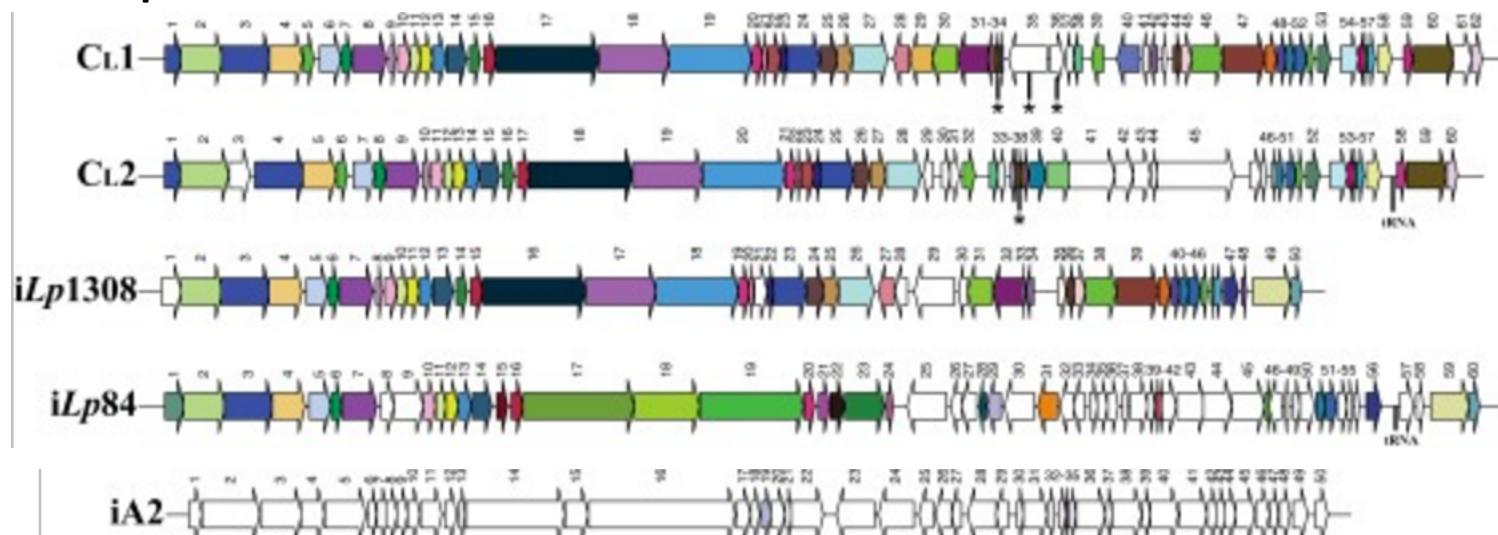
**Fig. 1.** Flowchart of IDBA-UD

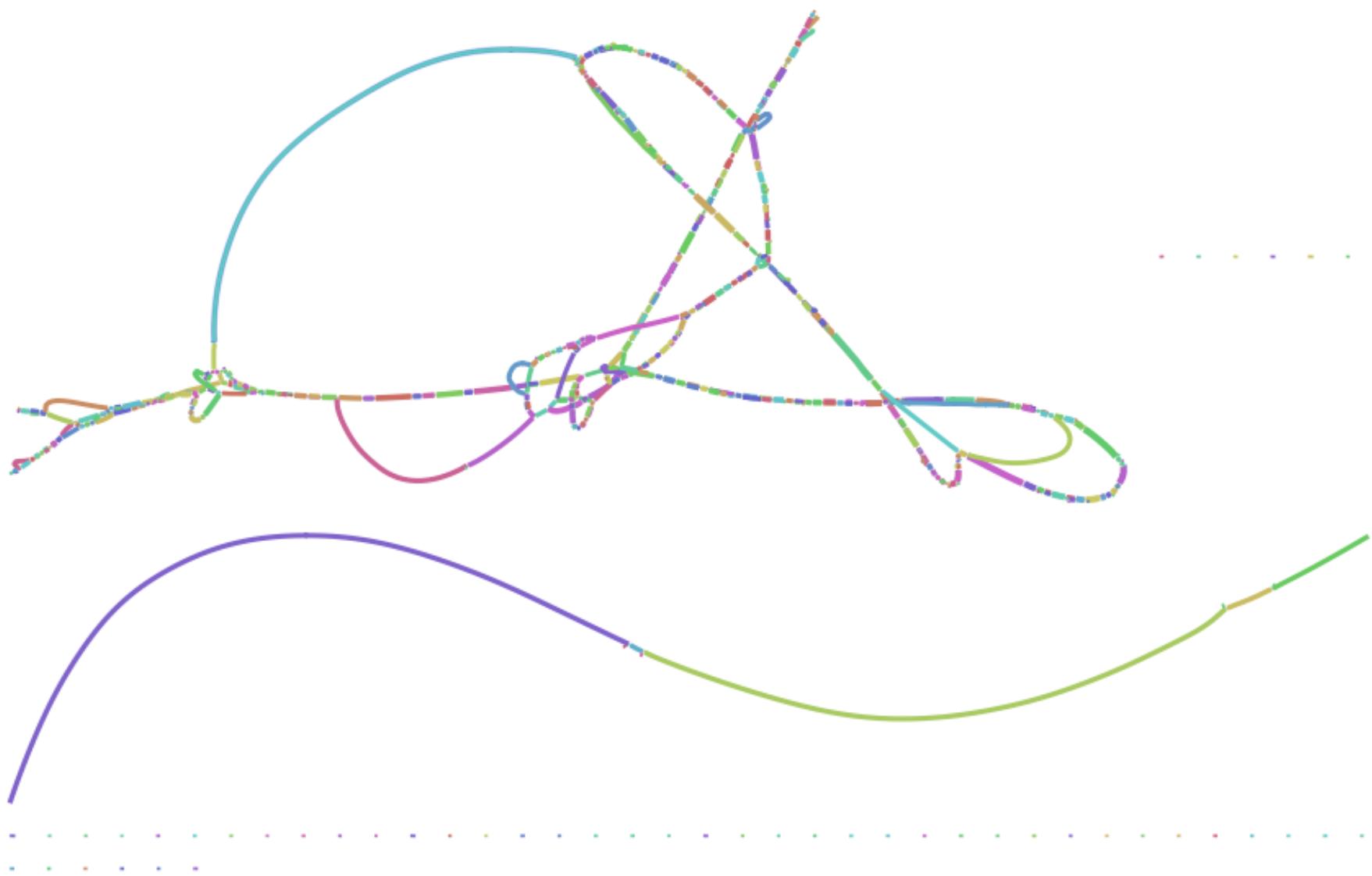
# Evaluating assemblies

- Total size of assembly
- Number of contigs
- N50 – minimum contig length in which at least 50% of bases are contained
- But contigs can be chimeric (Quast)

# Why is metagenomics assembly hard?

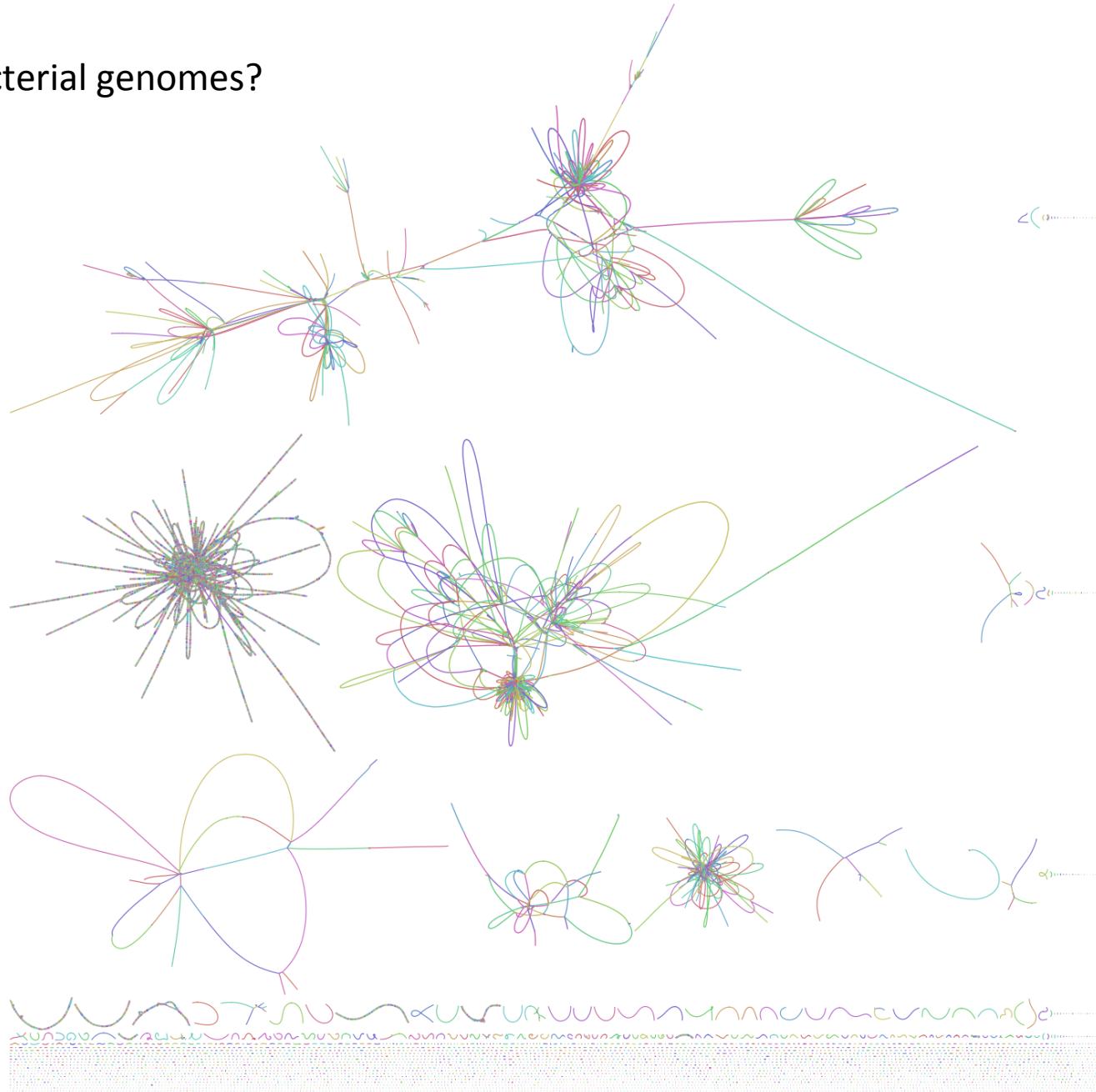
- Took five *Lactococcus paracasei* temperate phage genomes: CL1, CL2, iLp84, iLp1308, and iA2
- The genome lengths ranged from 34,155 bp (iA2) to 39,474 bp (CL1)
- Generated 10,000 synthetic 2X150bp reads for 16 samples





De Bruijn graph Bcalm assembly with kmer length 71

20 bacterial genomes?



Use contig binning to **cluster** contigs back into strain/species genomes

# Critical assessment of metagenome interpretation

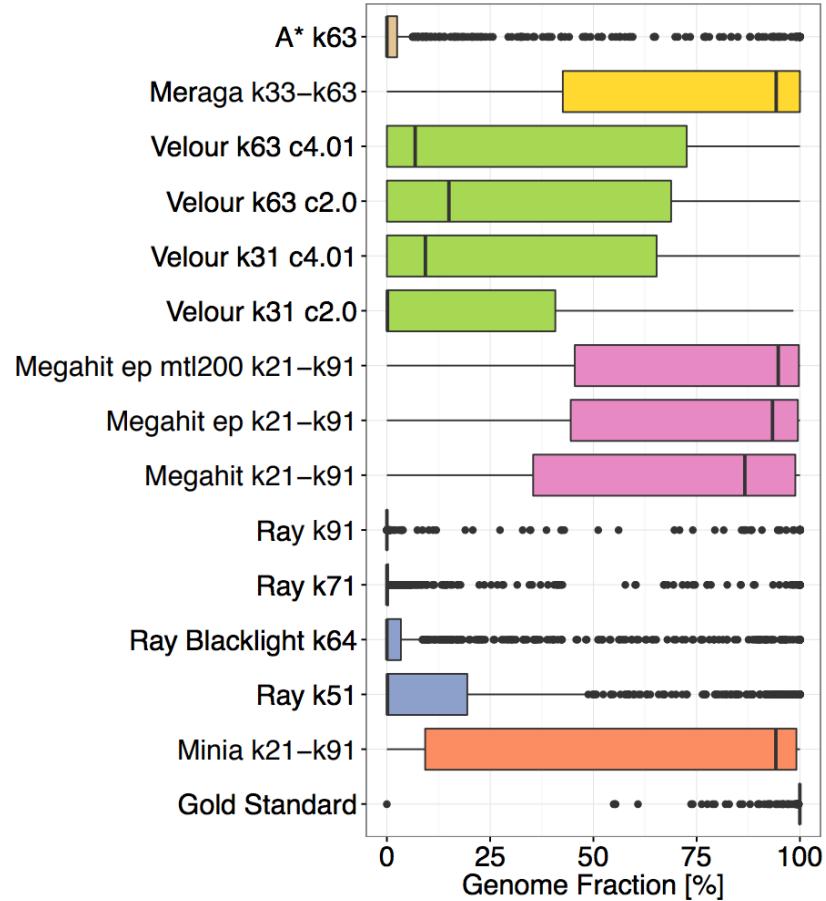
(<https://www.biorxiv.org/content/biorxiv/early/2017/01/09/099127.full.pdf>)

- Three datasets:
  - 15 Gbp – single sample low complexity (40 genomes and 20 circular)
  - 40 Gbp – differential abundance dataset with two samples of a medium complexity community (132 genomes and 100 circular elements)
  - 75 Gbp time series dataset with five samples from a high complexity community (596 genomes and 478 circular elements).

Software	No Spades	Description
<b>Assemblers</b>		
Megahit v.0.2.2		Metagenome assembler using multiple k-mer sizes and succinct de Bruijn graphs
Ray Meta v2.3.2		Distributed de Bruijn graph metagenome assembler
Meraga v2.0.4		Meraculous + Megahit
Minia 2 and Minia 3		De Bruijn graph assembler based on a Bloom filter
A*		OperaMS Scaffolder using SOAPde novo2 on medium complexity and Ray assemblies on low and high complexity data sets
Velour		De Bruijn graph genome assembler

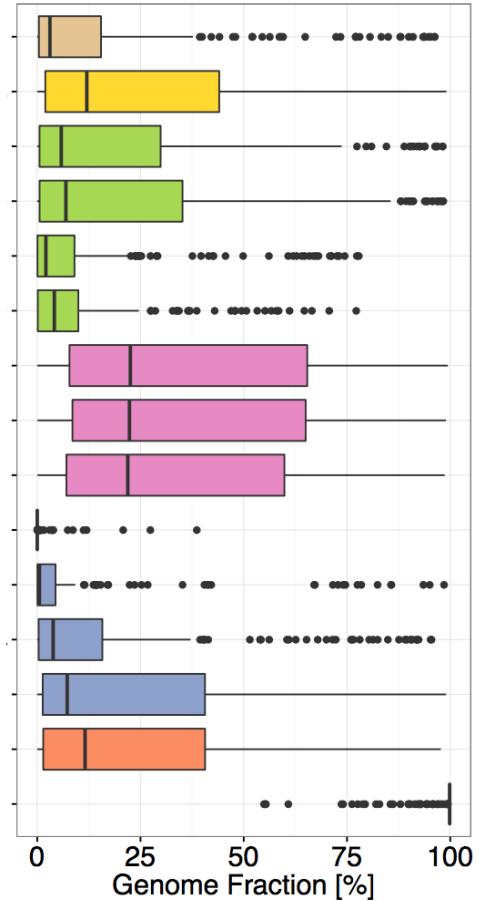
## All Genomes

a



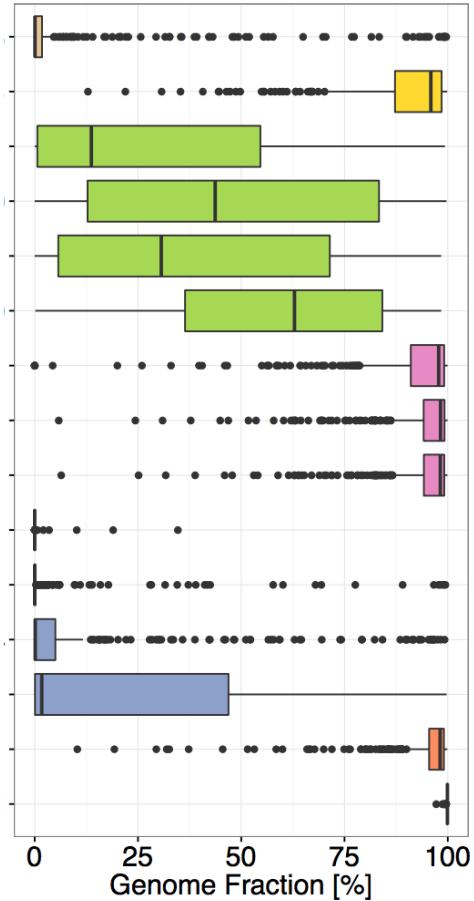
## ANI > 95%

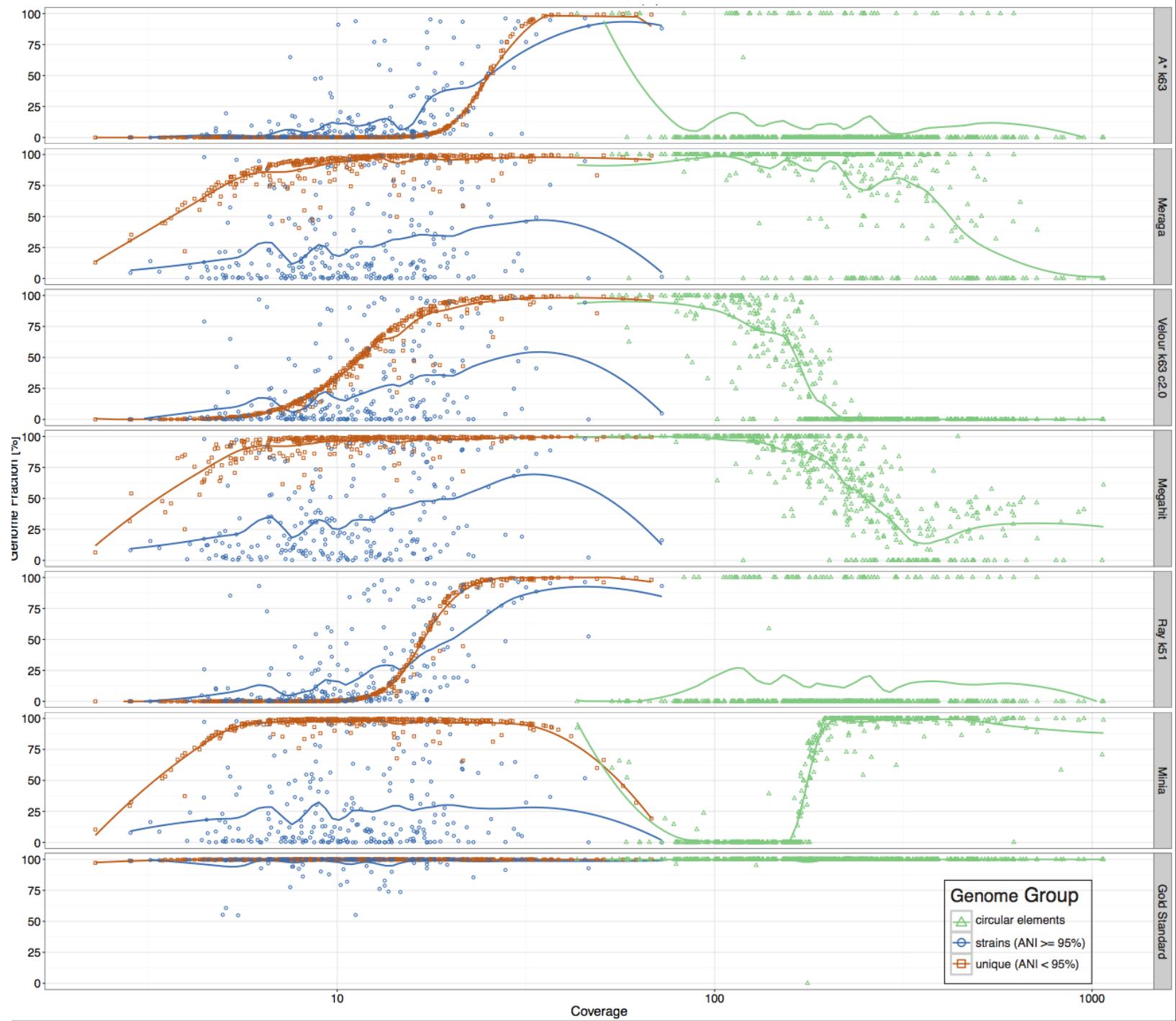
b



## All Genomes ANI < 95%

c

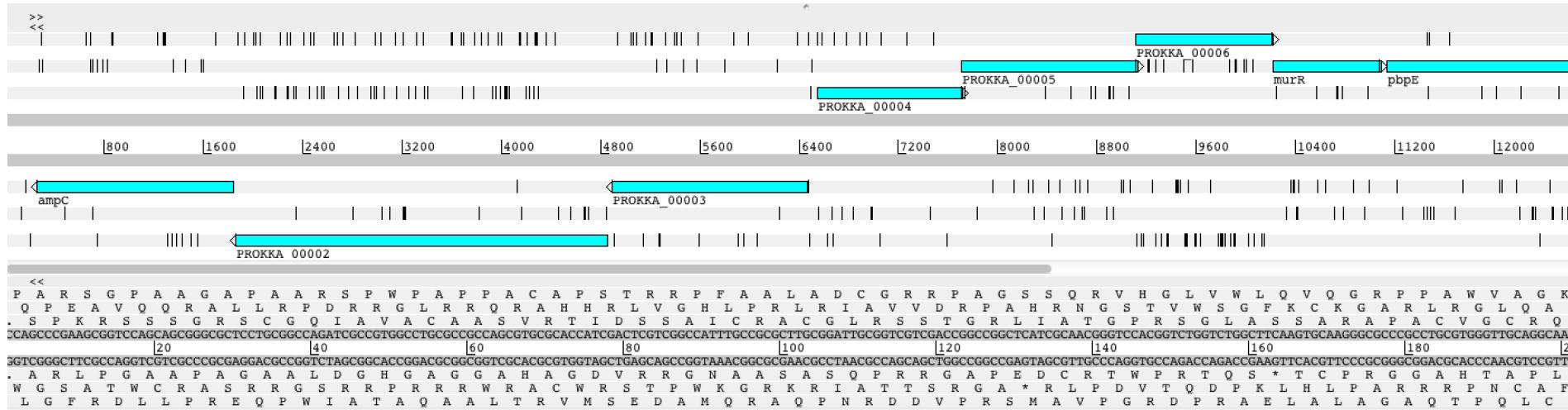




# Coassembly

- In a metagenome study we may have multiple studies from similar environments with the same organisms:
  - Time series e.g. reactor
  - Horizontal sampling
- Combining samples increases coverage of rare organisms and enables binning
- Cost is increased strain confusion
- Are contigs useful?

# Can look at context in contig e.g. AMR



<>  
P A R S G P A A G A P A A R S P W P A P P A C A P S T R R P F A A L A D C G R R P A G S S Q R V H G L V W L Q V Q G R P P A W V A G K  
Q P E A V Q R A L L R P R R G L R R Q R A H R N G S T V W S G F K C K G A R R L R G L Q A  
. S P K R S S S G R S C G Q I A V A C A A S V R T I D S S A I C R A C G L R S S T G R L I A T G P R S G L A S S A R A P A C V G C R Q  
..CAGGCCGAAGGGTCCAGCGCCGCTCTCGCCGAGATCGCCGTGGCCTCGCCGCCAGCGTGCACCATCGTCGCCATTGCGCCCTTGCGATTGCGTCACGGTCTGGCTCAAGTCGAAGGGCGCCGCTGGCTGGTTGCAAGGCA  
[20] ..GTCCGGGCTTCCCCAGGTCTCGCCCCCGAGGACGGCGGCTAGCCGGCACCGGCGGCTCCACGGCTGGTAGCTGAGCAGCGGTAAACGCCCGAACGCCAACGCCCTGGCCGGCGAGTAGCGTIGGCCAGGTGCCAGCCAGACCGGA  
[40] ..GAGTAAACGCCAACGCCCTGGCCGGCGAGTAGCGTIGGCCAGGTGCCAGCCAGACCGGAAGTTCACGTTCCCGCGGGCGACGCCAACCAAACGTCGGT  
[60] ..A R L P G A A P A G A A L D G H G A G G A H A G D V R R G N A A A S A S Q P R R G A P E D C R T W P R T Q S \* T C P R G G A H T A P L  
[80] ..W G S S A T W C R A S R R G S R P R R R W R A C W R S T P W K G R K R I A T T S R G A \* R L P D V T Q D P K L H L P A R R R P N C A F  
[100] ..L G F R D L L P R E Q P W I A T A Q A A L T R V M S E D A M Q R A Q P N R D D V P R S M A V P G R D P R A E L A L A G A Q T P Q L C  
[120] ..  
[140] ..  
[160] ..  
[180] ..

<<  
CDS 255 1835 c  
CDS 1859 4846 c  
CDS 4890 6455 c  
CDS 6549 7706  
CDS 7712 9118  
CDS 9115 10212  
CDS 10217 11080  
CDS 11138 12868  
CDS 12921 13277 c

>PROKKA\_00001 Beta-lactamase  
>PROKKA\_00002 hypothetical protein  
>PROKKA\_00003 hypothetical protein  
>PROKKA\_00004 hypothetical protein  
>PROKKA\_00005 hypothetical protein  
>PROKKA\_00006 L-Ala-D/L-Glu epimerase  
>PROKKA\_00007 HTH-type transcriptional regulator MurR  
>PROKKA\_00008 Penicillin-binding protein 4\*  
>PROKKA\_00009 hypothetical protein

Top blast hit Lysobacter antibioticus strain ATCC 29479 genome 78% 75%

# Summary

- De novo assembly of genomes is a powerful way to extract biological information from data set
- We will end up with very fragmented assemblies
- These are still useful and can be made even more useful with binning