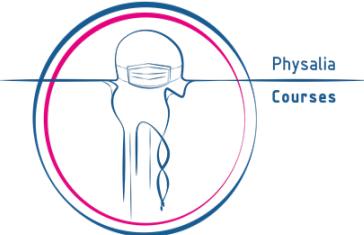
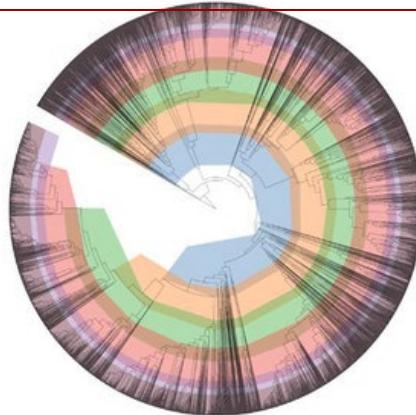


ENVIRONMENTAL METAGENOMICS

Physalia course, online, 11-15 November 2024

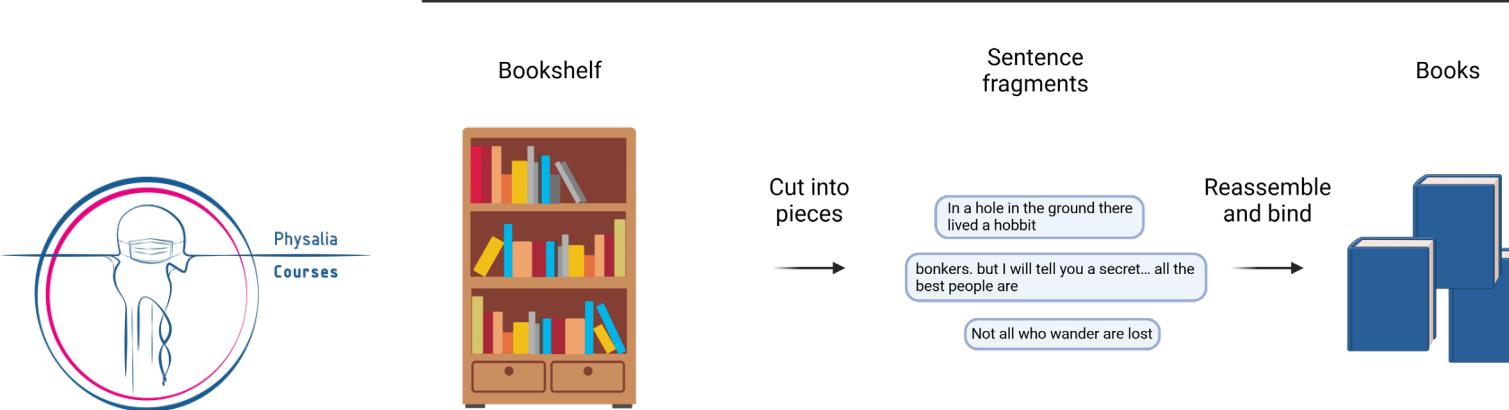
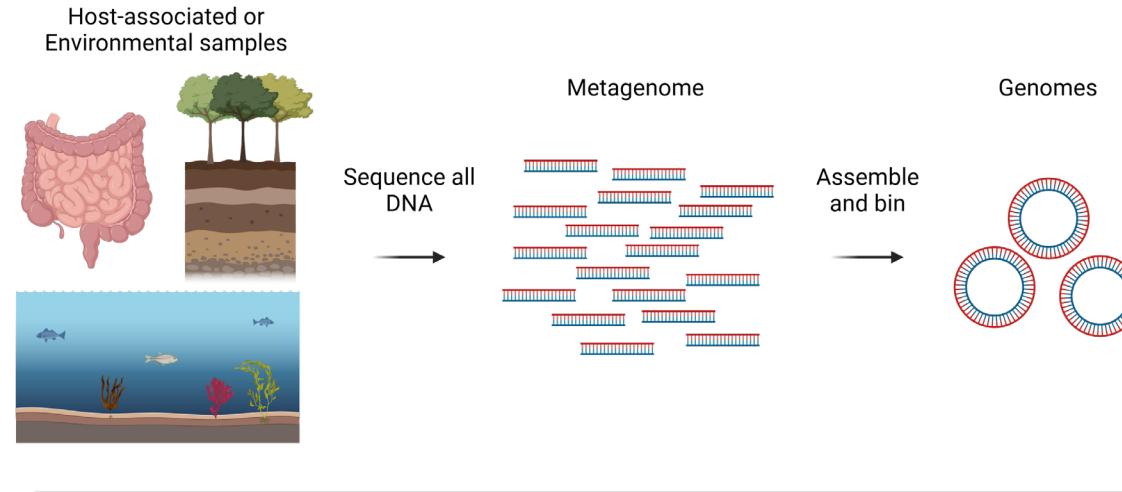
MAG binning

Nikolay Oskolkov, Lund University, NBIS SciLifeLab
Samuel Aroney, Queensland University of Technology



NB: original course material courtesy:
Dr. Antti Karkman, University of Helsinki
Dr. Igor Pessi, Finnish Environment Institute (SYKE)
As. Prof. Luis Pedro Coelho

Binning: collating contigs into collections



Standard binning strategies

- Nucleotide composition
 - Microbes tend to have similar %GC and codon choice across genome

A

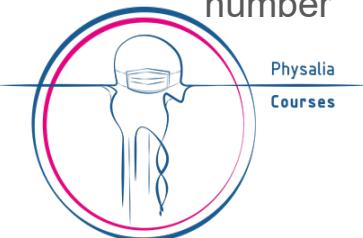
Nucleotide composition

Contig 1: AGCTGTCGCTCGGAGTCCAGG...	AGCT	AAAA - 5557	Contig 1	0.0321	0.0144	0.0231	0.0011	0.0031	...
	GCTG	AAAC - 2345	Contig 2	0.0362	0.0152	0.0211	0.0012	0.0029	...
	CTGT	AAAG - 1034	Contig 3	0.0033	0.0083	0.0344	0.0651	0.0005	...
	TGTC	AAAT - 6345	Contig 4	0.0021	0.0079	0.0339	0.0649	0.0040	...
	GTCG	AACA - 2493	Contig 5	0.0031	0.0141	0.0220	0.0011	0.0033	...
	Contig 6	0.0009	0.0022	0.0016	0.0014	0.0001	...
			Contig 7	0.0012	0.0110	0.0141	0.0284	0.0261	...
			...						
			Contig n	0.0229	0.0091	0.0045	0.0001	0.0042	...
			Contig n	0.0229	0.0091	0.0045	0.0001	0.0042	...

B

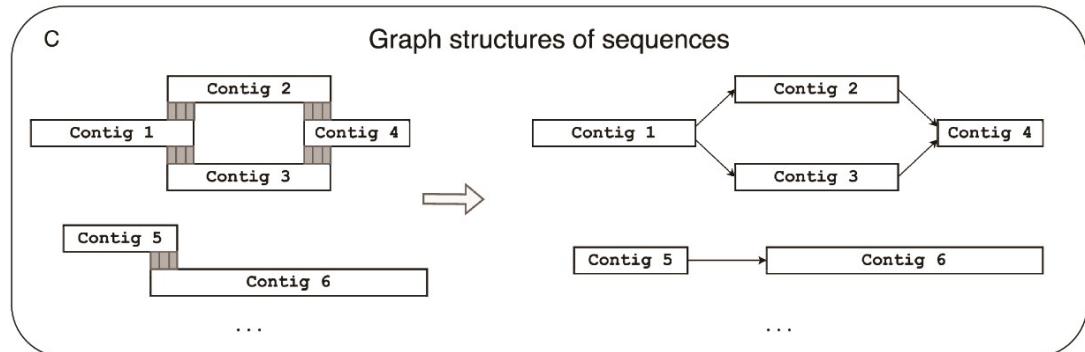
Abundance

Contig 1: AGCTGTCGCTCGGAGTCCAGG...	AGCTGTCGCT	GGAGTCCAGG	Sample 1	10	Sample 2	...	Sample m	22	
	GCTGTCGCTC	GAGTCCAGG	Contig 2	11	49	...	19	...	
	CAGCTGTCGCT	GGAGTCCAGG	Contig 3	47	33	...	65	...	
	CTGTCGCTCG	AGTCCAGGA	Contig 4	0	17	...	63	...	
	AGCTGTCGCG	CGGAGTCCAG	Contig 5	83	42	...	17	...	
	...		Contig 6	0	18	...	58	...	
			Contig 7	52	44	...	0	...	
			...						
			Sample m	GCTGTCGCTC	GAGTCCAGG	Contig n	22	71	...
				CAGCTGTCGCT	GGAGTCCAG				30
				...					

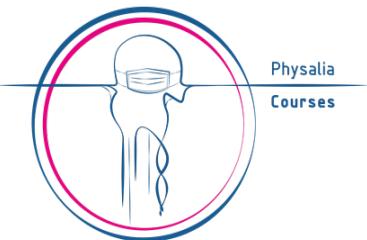
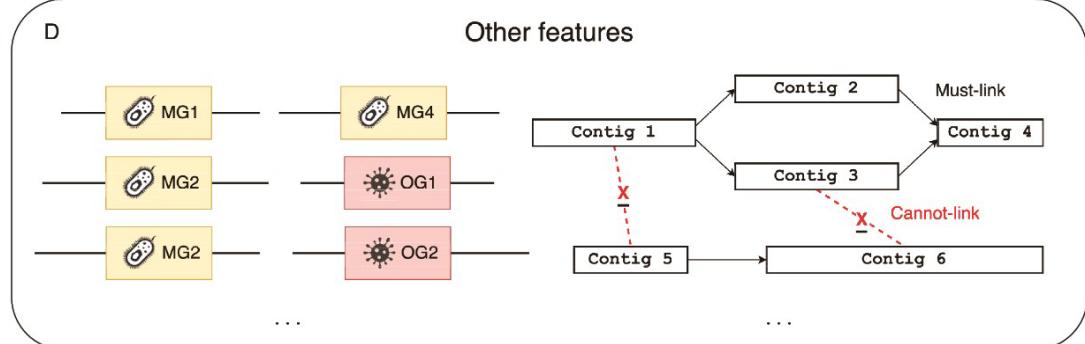


More binning strategies

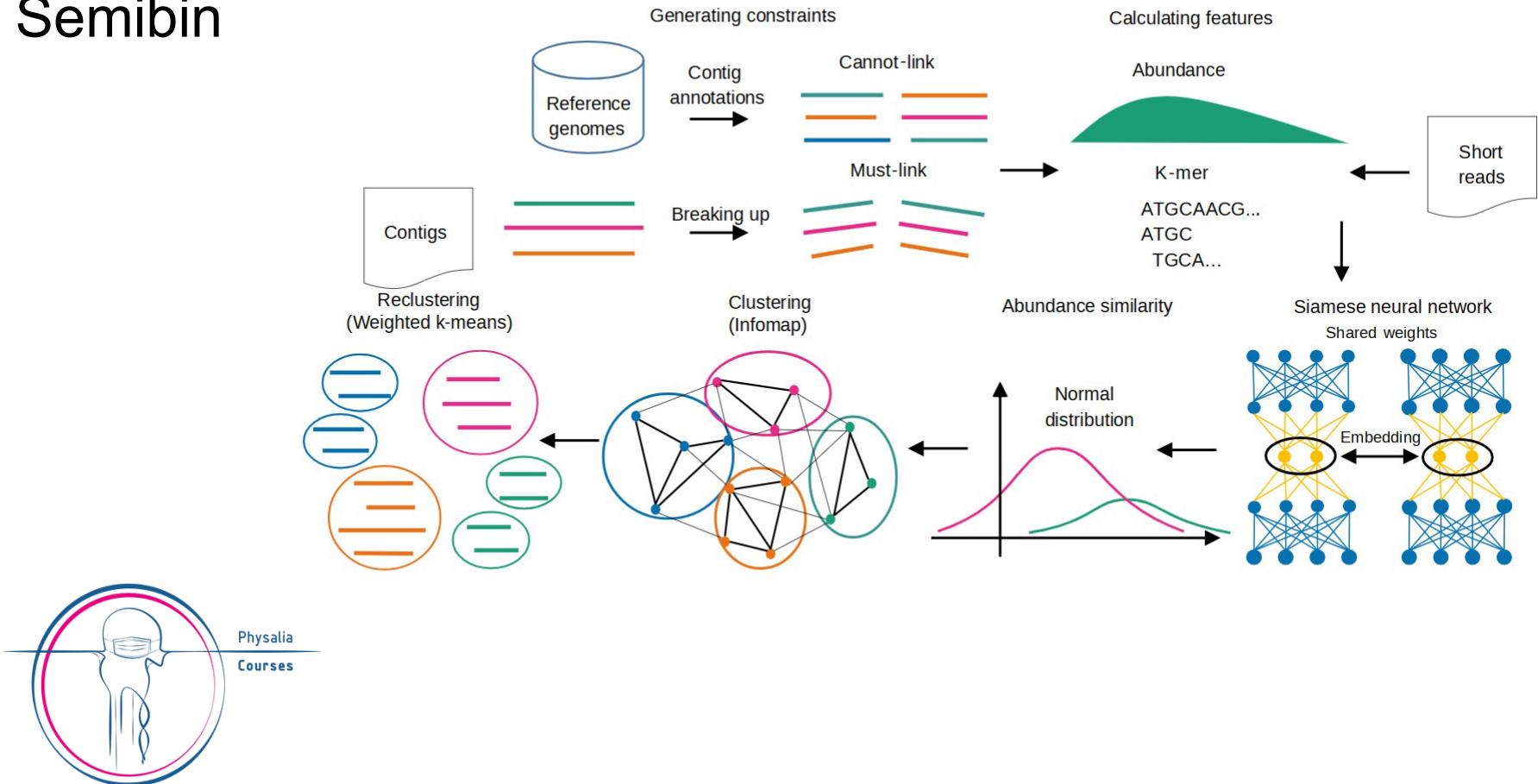
- Graph structures
 - Assembly graph structures can be reused



- Other features
 - Methylation patterns
 - Single-copy marker genes
 - Machine learning inputs



Semibin



Learning a new model is resource intensive

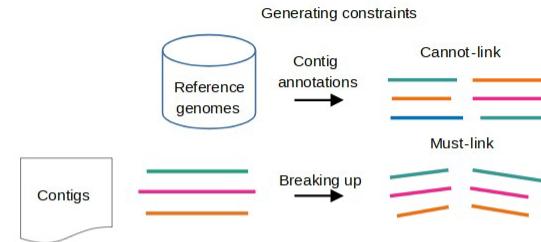
- SemiBin uses a learned model
- Generating the data for training is expensive (mmseqs2 is slow)
- Learning is slow (a GPU helps, but few people have a GPU cluster)

		Single-sample		
		Human gut	Dog gut	Ocean
Computing features	Time(min)	5	9	12
	Memory(MB)	923	736	1,675
Generating cannot-link	Time(min)	88	81	114
	Memory(MB)	39,070	37,904	46,091
Training(CPU)	Time(min)	181	209	222
	Memory(MB)	2,497	2,373	3,211
Training(GPU)	Time(min)	34	36	45
	Memory(MB)	4,487	4,355	5,222
Binning	Time(min)	2	2	3
	Memory(MB)	4,501	3,641	7,622



From Semi- to self-supervised learning

- In SemiBin1, must-links were self-supervised
- In SemiBin1, cannot-links were from a reference taxonomy



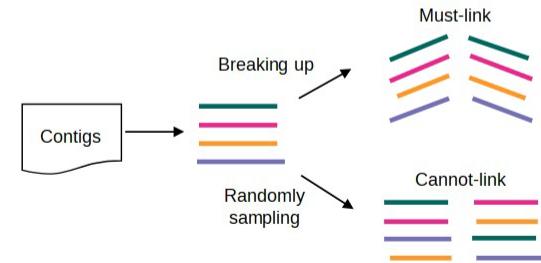
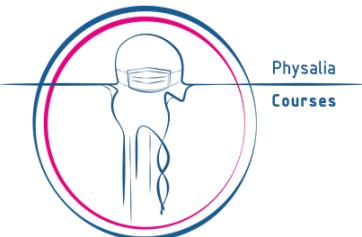
Use random sampling to generate cannot-links

Random sampling advantages

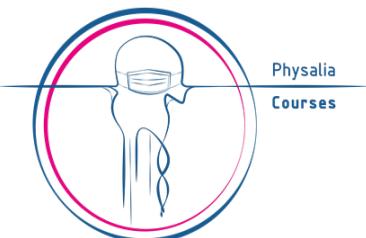
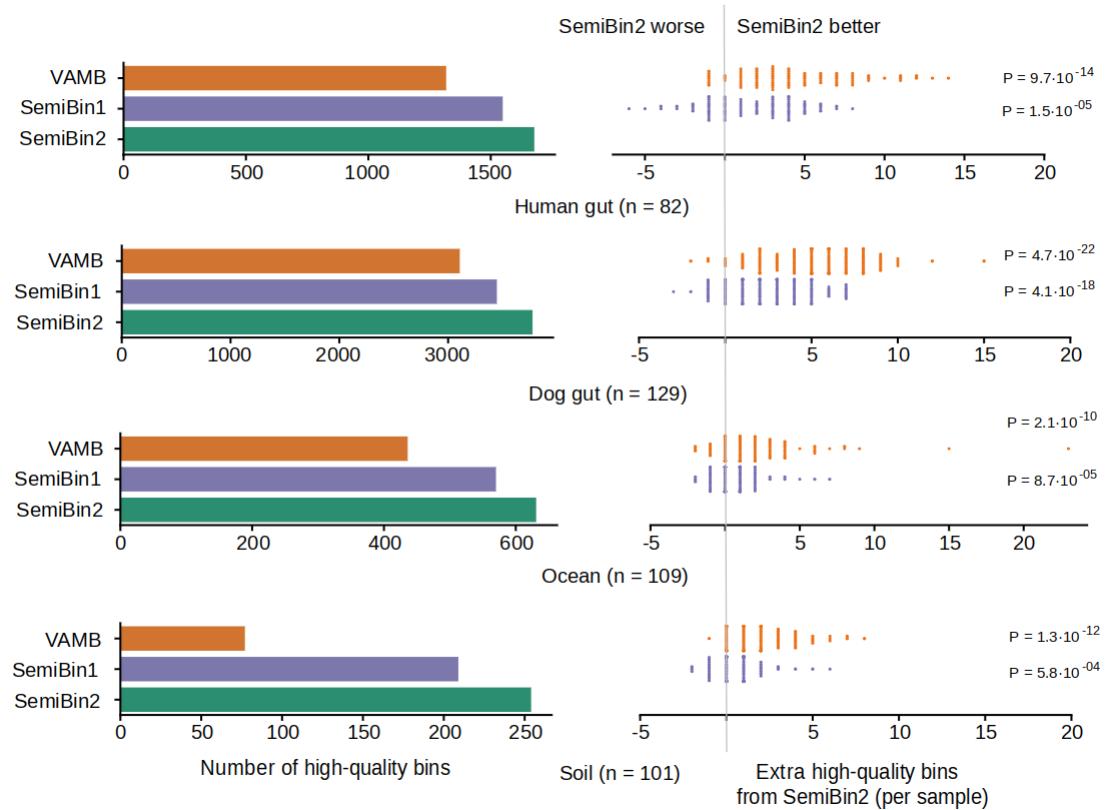
1. Random sampling is faster than MMSeqs2
2. It has no database biases

Random sampling disadvantages

1. Random sampling introduces errors (naturally)

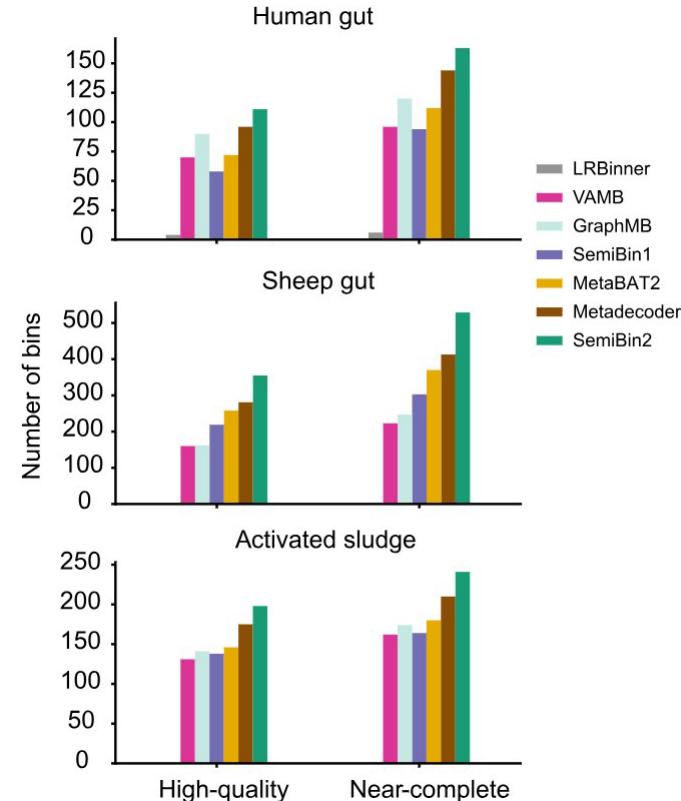


Semibin vs Semibin2



SemiBin2 surpasses alternatives for long-read binning

- Three datasets: human gut, sheep gut, and activated sludge
- Some of the alternatives (like SemiBin1) are not designed for long-reads
- We used [CheckM2](#) & [GUNC](#) to evaluated the results
- SemiBin2 is also best using [checkM1](#), but the two tools share some of the same marker genes (so evaluation is partly circular)



All the binners

Binning metagenomic contigs by coverage and composition

Johannes Alneberg^{1,8}, Brynjar Smári Bjarnason^{1,8}, Ino de Brujin^{1,2}, Melanie Schirmer³, Joshua Quick^{4,5}, Umer Z Ijaz³, Leo Lahti^{6,7}, Nicholas J Loman⁴, Anders F Andersson^{1,9} & Christopher Quince^{3,9}

MetaBAT, an efficient tool for accurately reconstructing single genomes from

MetaBAT 2: an adaptive binning algorithm for robust and efficient genome reconstruction from metagenome assemblies

Dongwan D. Kang¹, Feng Li², Edward Kirton¹, Ashleigh Thomas¹, Rob Egan¹, Hong An² and Zhong Wang^{1,3,4}

Effective binning of metagenomic contigs using contrastive multi-view representation learning

Received: 28 June 2023

Ziye Wang^①, Ronghui You¹, Haitao Han^①, Wei Liu¹, Fengzhu Sun^② & Shafeng Zhu^{③,4,5,6}

Accepted: 7 December 2023

A de novo binning method for metagenomic datasets across different environments

Shaojun Pan^{①,2}, Chengkai Zhu^{1,2,3}, Xing-Ming Zhao^{①,2,4,5} & Luis Pedro Coelho^{①,2}

MaxBin 2.0: an automated binning algorithm to recover genomes from multiple metagenomic datasets

Yu-Wei Wu^{1,2,*}, Blake A.

Improved metagenome binning and assembly using deep variational autoencoders

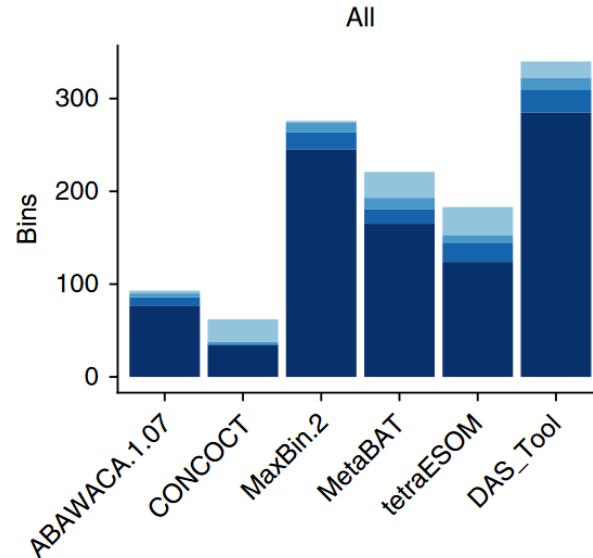
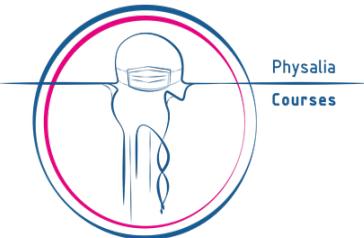
Jakob Nybo Nissen^{1,2}, Joachim Johansen^②, Rosa Lundbye Allesøe², Casper Kaae Sønderby¹, Jose Juan Almagro Armenteros^①, Christopher Heje Grønbæk^{3,4}, Lars Juul Jensen^②, Henrik Bjørn Nielsen^⑤, Thomas Nordahl Petersen⁶, Ole Winther^{3,4,7} and Simon Rasmussen^{②,8}

SemiBin2: self-supervised contrastive learning leads to better MAGs for short- and long-read sequencing

Shaojun Pan^{①,2}, Xing-Ming Zhao^{①,2,3,4,*}, Luis Pedro Coelho^{①,2,*}

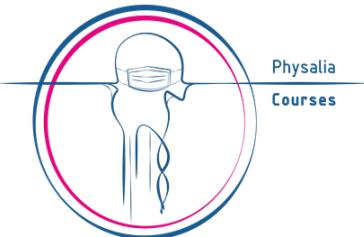
Ensemble binning

- Combining bins from multiple binners
- DAS Tool
- MetaWRAP
- MetaBinner
- BASALT



Aviary: Our Assembly + Binning pipeline

- Snakemake pipeline (automatically submit jobs to a cluster)
- Assembly: metaSPAdes or MegaHit
 - Particularly good at hybrid short+long read
- Binning
 - MetaBAT1/2, VAMB, SemiBin1/2, CONCOCT, MaxBin2, Rosella
 - DAS Tool



All the binners - comparison

semibin = Semibin2

