

SANS serif: alignment-free, whole-genome based phylogenetic reconstruction

Roland Wittler, Andreas Rempel, Marco Sohn

Bielefeld University

DSB 2021

Outline

SANS

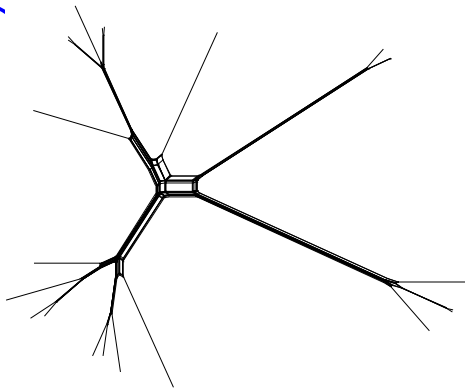
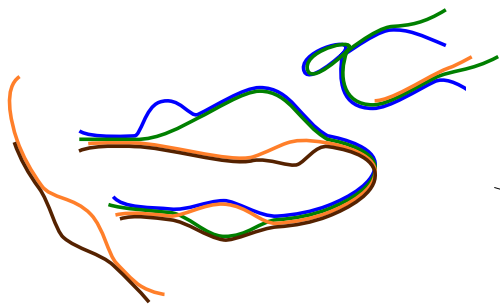
SANS serif

IUPAC Characters

Coding Sequences

Clustering MAGs

Outlook



Classical approaches

- ▶ mult. alignment of marker genes
- ▶ alignment to reference → SNPs
- ▶ pairwise distances or ML
- ▶ tree

Classical approaches

- ▶ mult. alignment of marker genes
- ▶ alignment to reference → SNPs
- ▶ pairwise distances or ML
- ▶ tree

Whole genome

- ▶ k -mers,
spaced k -mers, ...
- ▶ pairwise distances
- ▶ tree

Classical approaches

- ▶ mult. alignment of marker genes
- ▶ alignment to reference → SNPs
- ▶ pairwise distances or ML
- ▶ tree

Whole genome

- ▶ k -mers,
spaced k -mers, ...
- ▶ pairwise distances
- ▶ tree

SANS: Symmetric Alignment-free phylogeNomic Splits

- ▶ whole genome (no markers, no alignment, no reference)

Classical approaches

- ▶ mult. alignment of marker genes
- ▶ alignment to reference → SNPs
- ▶ pairwise distances or ML
- ▶ tree

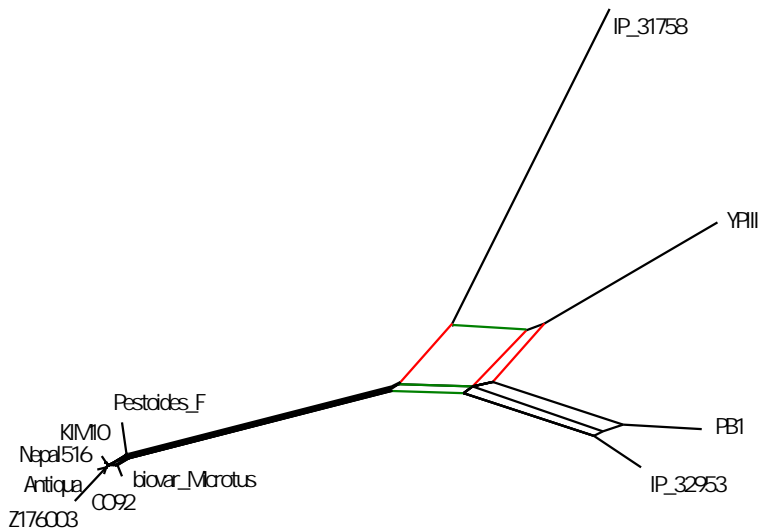
Whole genome

- ▶ k -mers,
spaced k -mers, ...
- ▶ pairwise distances
- ▶ tree

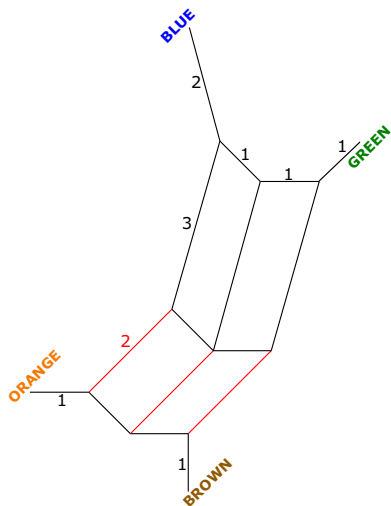
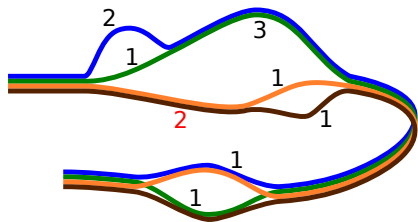
SANS: Symmetric Alignment-free phylogeNomic Splits

- ▶ whole genome (no markers, no alignment, no reference)
- ▶ k -mers → splits (no pairwise comparison)

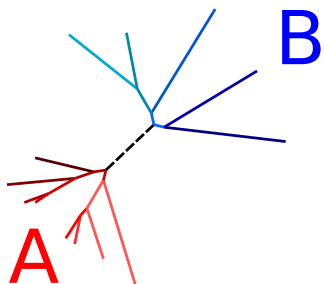
Introduction — SANS



Introduction — SANS



Introduction — SANS



A and B separated in phylogeny

⇒ mutations on that edge

⇒ sequences unique to A and
sequences unique to B

⇒ A and B k -mers

of about same total amount $w \approx w'$

⇒ $\sqrt{w \cdot w'}$ to downweight asymmetry

or $\sqrt{(w+1) \cdot (w'+1)}$ to not lose any

Outline

SANS

SANS serif

IUPAC Characters

Coding Sequences

Clustering MAGs

Outlook

First SANS version [DSB 2019, WABI 2019]:

1. build C-DBG (Bifrost, Melsted&Holley)
2. extract splits from unitigs and store in trie

Re-implementation:

- ▶ hash tables
- ▶ no dependencies
- ▶ space and time efficient
- ▶ includes filtering of splits (e.g. tree filter)

First SANS version [DSB 2019, WABI 2019]:

1. build C-DBG (Bifrost, Melsted&Holley)
2. extract splits from unitigs and store in trie

Space and time **Efficient Re-Implementation** incl. **Filters: SANS serif**

- ▶ hash tables
- ▶ no dependencies
- ▶ space and time efficient
- ▶ includes filtering of splits (e.g. tree filter)

```

>SAL_BA7171AA_AS_NODE_40
GGCAGTGGAGAATAAGTTGCACTGGCGGTT
GGATGTGGTGATGAATGAAGATGATTGCAGAATAAG
AAGAGGGAATGCTGCGGAATTGTTTTAGGGATTAG
>SAL_BA7171AA_AS_NODE_19
CTCGTTGATGGGGTAGTTATTGTGGAATGTCCACCG
GTGTGCCATCAAGAAAAATTTATCAGCATAGCGAG
TTGAAAAATTCATATTTATGAAGAACATAAGAAATT
TCTCCATCATTGCTCACATTGACCACGGTAAATCGA
CGCTGTCTGACCGTATTATCCAGATCTGCGGTGGCC
>SAL_BA7171AA_AS_NODE_13
GAAGATACAGGACTACATAAAGCACCAGCTTGAAGA
GGATAAAATGGGAGAGCAGTTATCGATCCCTTATCC
GGTAGCCCGTTTACGGGCCGTAAGTAGCGAAGTCT
GATGCAAATGTCAGATCGCGTGCCTGTTAGGGCG
CGGCTGGTAAGAGAGCCTTATAGCGCATCAGAAAA
ACCTCCGGCTATGCCGGAGGATATTTATTACATTCT
G GCG TA CA GGT T
  
```

<i>k</i> -mer	files
0b10010010100000	0b10000000000000
0b01110110111110	0b11000000000000
0b11000101111000	0b11000000000000
0b01001010000001	0b01011101100000
0b01000011000110	0b00111111111111
0b00010001101001	0b00111111111111
0b00011101010101	0b00111111111111
0b00011110101111	0b01011101100000
0b10100011101001	0b11000000000000
0b01100011111101	0b00001001111111
.....

SANS serif

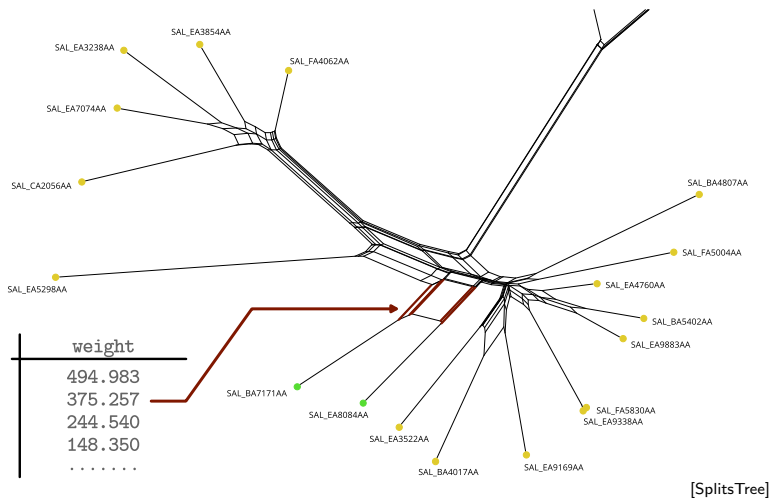
Andreas Rempel

<i>k</i> -mer	files		split	# <i>k</i> -mers
0b10010010100000	0b10000000000000		0b001011101000000	[499, 491]
0b01110110111110	0b11000000000000	}	0b11000000000000	[362, 389]
0b11000101111000	0b11000000000000		0b00011010110000	[260, 230]
0b01001010000001	0b01011101100000		0b00100111001000	[131, 168]
0b01000011000110	0b00111111111111	
0b00010001101001	0b00111111111111			
0b00011101010101	0b00111111111111			
0b00011111010111	0b01011101100000			
0b10100011101001	0b11000000000000			
0b01100011111101	0b00001001111111			
.....			

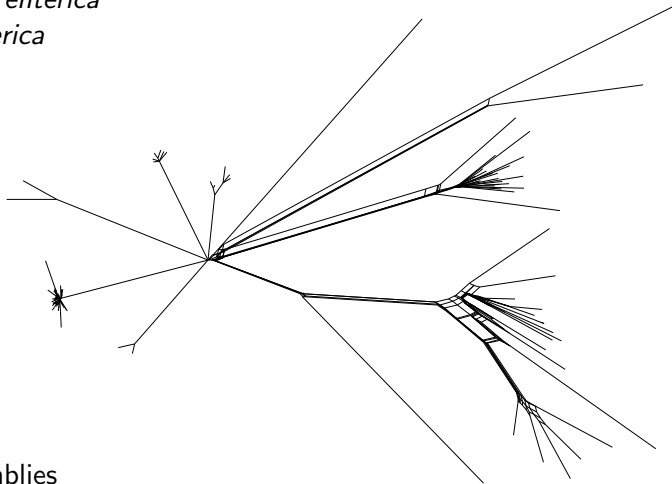
split	# <i>k</i> -mers		weight
0b001011101000000	[499, 491]		494.983
0b110000000000000	[362, 389]	→	375.257
0b00011010110000	[260, 230]		244.540
0b00100111001000	[131, 168]		148.350
.....

SANS serif

Andreas Rempel



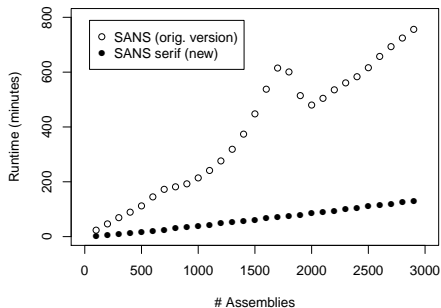
Salmonella enterica
subsp. *enterica*



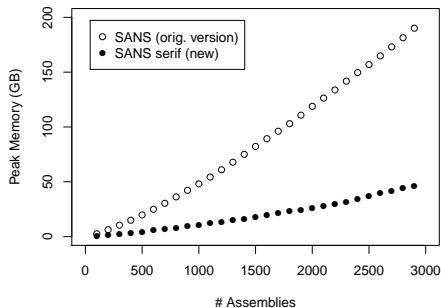
2964 assemblies
[Zhou et al., 2017]

SANS serif

Andreas Rempel



250 assemblies: [andi](#): 110 min, [Co-phylog](#): 9 h, [FSWM](#): 50 h, ...



Outline

SANS

SANS serif

IUPAC Characters

Coding Sequences

Clustering MAGs

Outlook

IUPAC Characters

Andreas Rempel

Consider all variants:

R → A or G

Y → C or T

S → G or C

W → A or T

K → G or T

M → A or C

B → C or G or T

D → A or G or T

H → A or C or T

V → A or C or G

N → any base

AARCGYA ⇒ $\left\{ \begin{array}{l} \text{AAACGCA} \\ \text{AAACGTA} \\ \text{AAGCGCA} \\ \text{AAGCGTA} \end{array} \right.$

IUPAC Characters

Andreas Rempel

R → A or G

Y → C or T

S → G or C

W → A or T

K → G or T

M → A or C

B → C or G or T

D → A or G or T

H → A or C or T

V → A or C or G

N → any base

Consider all variants:

AARCGYA ⇒ $\left\{ \begin{array}{l} \text{AAACGCA} \\ \text{AAACGTA} \\ \text{AAGCGCA} \\ \text{AAGCGTA} \end{array} \right.$

True variant:

supports / is supported by further k -mers

False variant:

cancelled out due to missing inverse split

IUPAC Characters

Andreas Rempel

Simulated phylogeny, 100 leaf genomes, length ~ 96 kb, 5 PAM to the root simulated with ALF [Dalquen et al., 2012]

	original	with 0.1 % N's	
		skipped	replaced
<i>unweighted</i>			
precision	0.90	0.63	0.78
recall	0.70	0.38	0.58
<i>weighted</i>			
precision	0.98	0.95	0.98
recall	0.88	0.68	0.88

Outline

SANS

SANS serif

IUPAC Characters

Coding Sequences

Clustering MAGs

Outlook

Coding Sequences

Marco Sohn

whole genome \Rightarrow coding sequences
DNA \Rightarrow amino acids

-a, --amino Consider amino acids:

--input provides amino acid sequences

Implies --norev and a default k of 10

-c, --code Translate DNA: --input provides coding sequences

Implies --norev and a default k of 10

optional: ID of the genetic code to be used

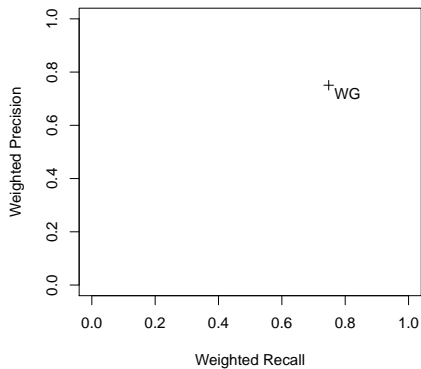
Default: 1 (The Standard Code)

Use 11 for Bacterial, Archaeal, and Plant Plastid Code

Coding Sequences

Marco Sohn

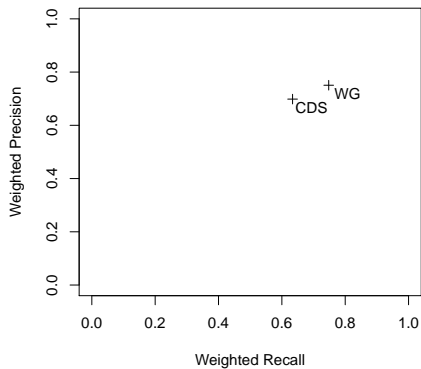
90 Pseudomonas genomes



Coding Sequences

Marco Sohn

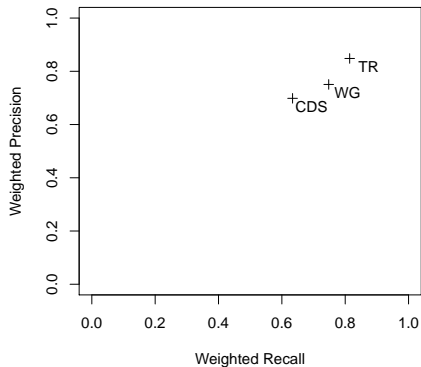
90 Pseudomonas genomes



Coding Sequences

Marco Sohn

90 Pseudomonas genomes



Outline

SANS

SANS serif

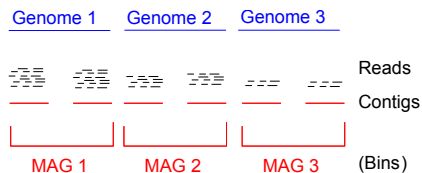
IUPAC Characters

Coding Sequences

Clustering MAGs

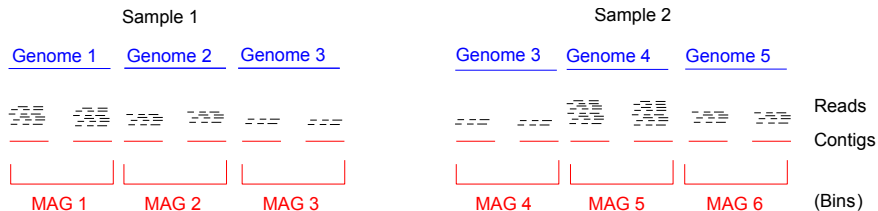
Outlook

Clustering MAGs

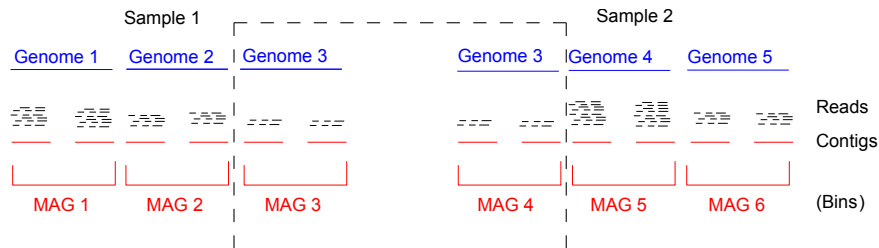


Metagenome Assembled Genomes

Clustering MAGs



Clustering MAGs



Given: Many MAGs from many samples.

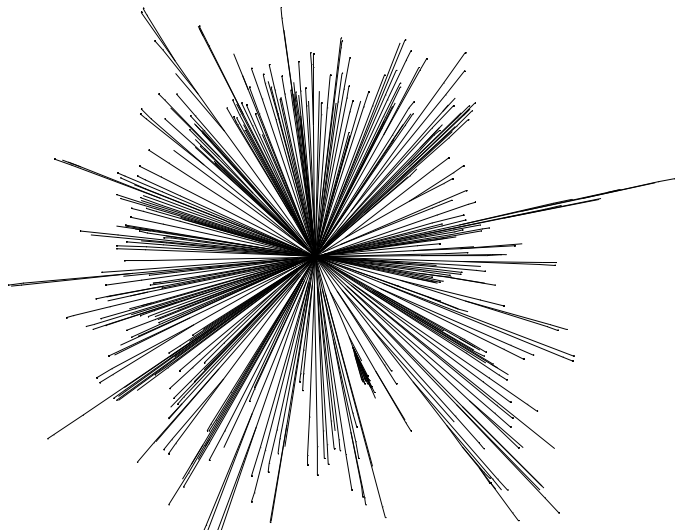
Problem: Cluster MAGs such that: cluster = genome

Clustering MAGs

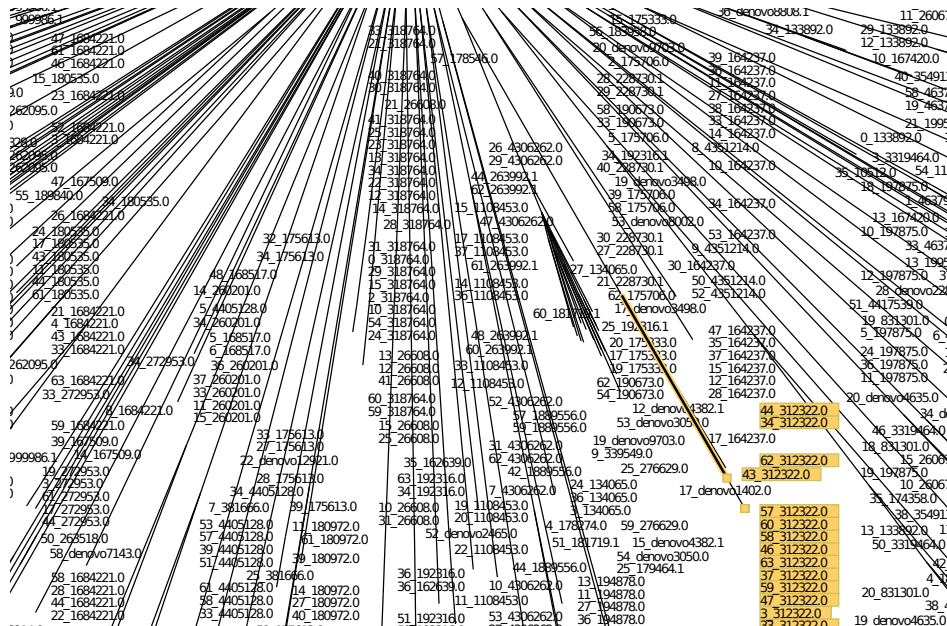
1. infer phylogeny using SANS
2. chop tree into clusters

Clustering MAGs

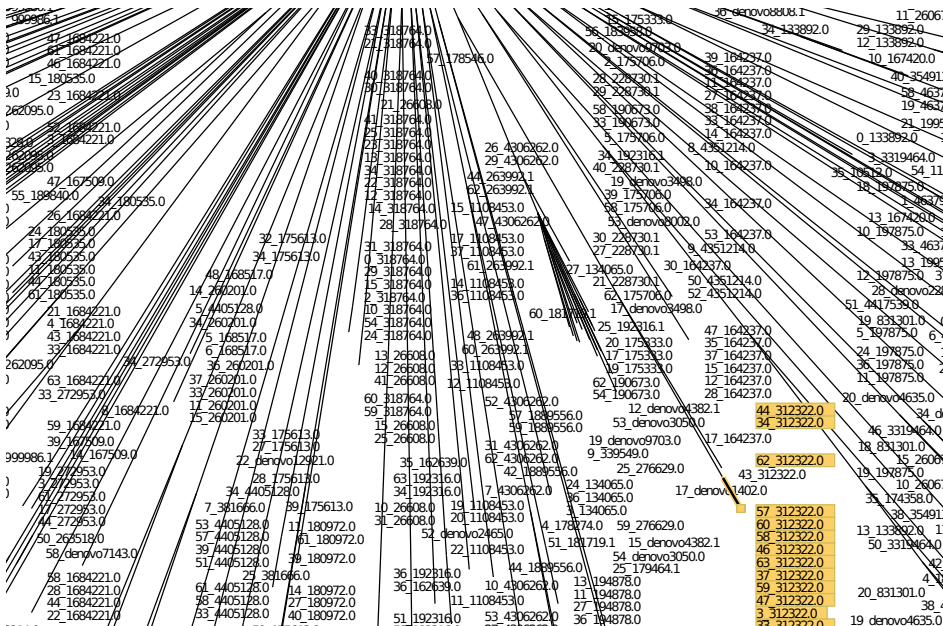
1. infer phylogeny using SANS
2. chop tree into clusters



Clustering MAGs



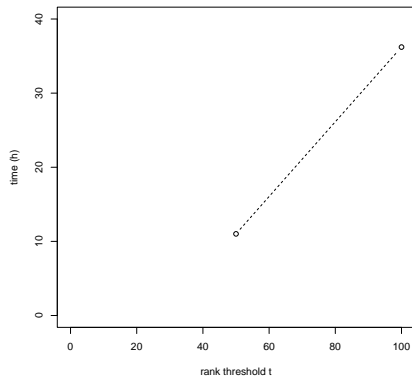
Clustering MAGs



Clustering MAGs

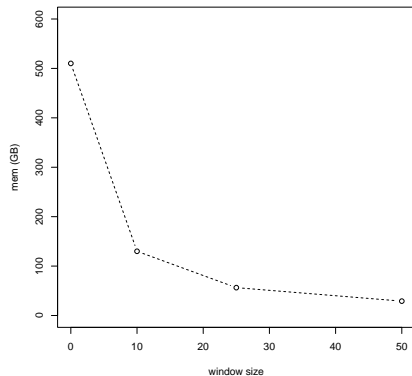
Rank threshold

MIMAG medium, w = 10 / 25



Minimizers *[Andreas Rempel]*

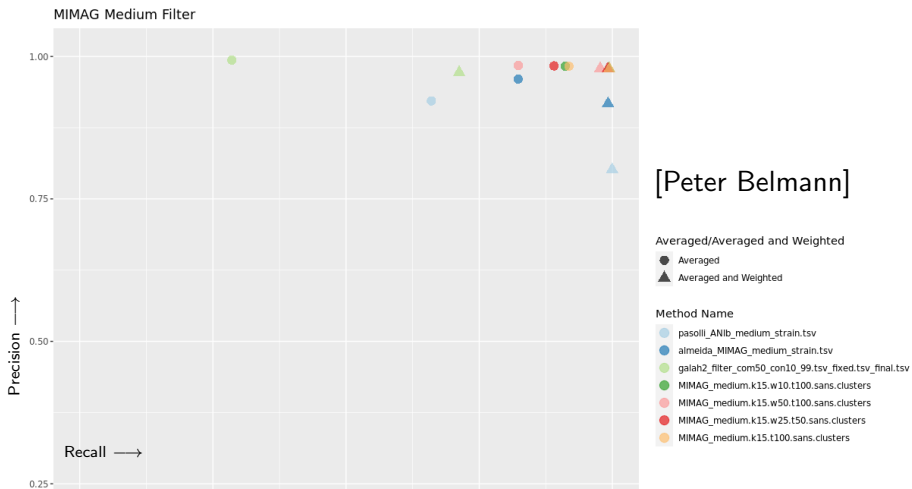
MIMAG medium, top 100



Clustering MAGs

CAMI mouse gut metagenome [Sczyrba et al.]

5 786 MAGs, 686 genomes (simulated)



[Peter Belmann]

Outline

SANS

SANS serif

IUPAC Characters

Coding Sequences

Clustering MAGs

Outlook

- ▶ parallelization [Fabian Kolesch]
in particular reading the input files

- ▶ parallelization [Fabian Kolesch]
in particular reading the input files
- ▶ counting k -mers [Ann-Cathrin Groba]
1 1 0 0 vs. 9 9 0 0 vs. 9 9 1 0

- ▶ parallelization [Fabian Kolesch]
in particular reading the input files
- ▶ counting k -mers [Ann-Cathrin Groba]
1 1 0 0 vs. 9 9 0 0 vs. 9 9 1 0
- ▶ estimating k [Rebecca Pfeil]
Shannon entropy(?), sampling

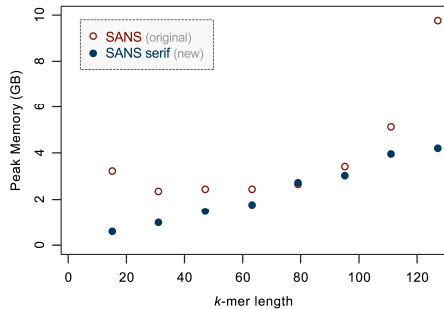
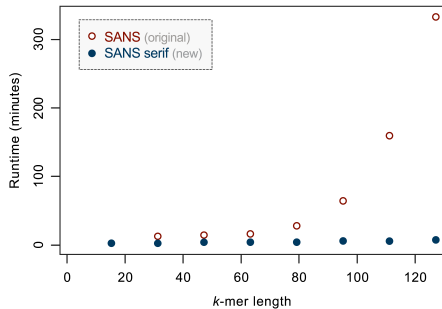
- ▶ parallelization [Fabian Kolesch]
in particular reading the input files
- ▶ counting k -mers [Ann-Cathrin Groba]
1 1 0 0 vs. 9 9 0 0 vs. 9 9 1 0
- ▶ estimating k [Rebecca Pfeil]
Shannon entropy(?), sampling
- ▶ reference free quality measure [Marco Sohn, Rebecca Pfeil]
tree-likeness

- ▶ parallelization [Fabian Kolesch]
in particular reading the input files
- ▶ counting k -mers [Ann-Cathrin Groba]
1 1 0 0 vs. 9 9 0 0 vs. 9 9 1 0
- ▶ estimating k [Rebecca Pfeil]
Shannon entropy(?), sampling
- ▶ reference free quality measure [Marco Sohn, Rebecca Pfeil]
tree-likeness

Thank you!

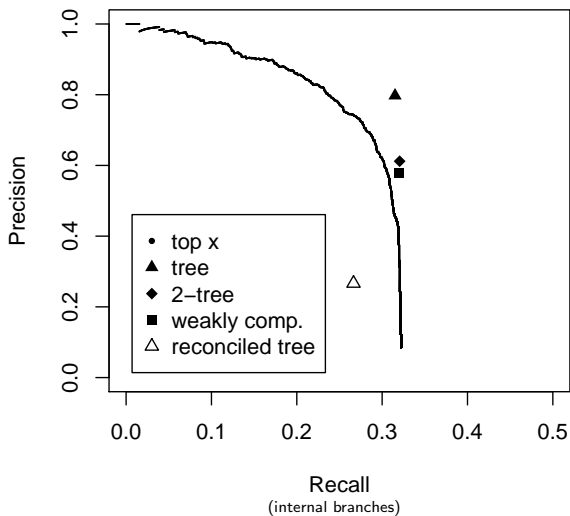
Appendix

Appendix



(100 assemblies)

Appendix



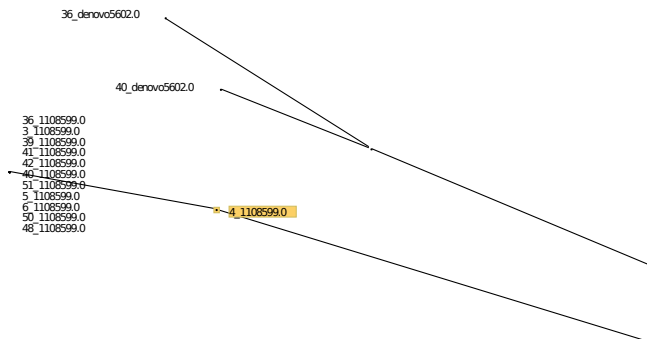
Appendix

90 *Pseudomonas* genomes

	Whole Genome	Coding Sequences	
	DNA	DNA	AA
k	25	25	8
time	17min	6.9s	1.5s
memory	12G	140M	17M

Appendix

1. SANS, filter for tree, re-root
2. post-order traversal:
 - 2.1 ignore non-branching node
 - 2.2 get clusters from sub-trees (recursively)
 - item if edge \geq parent edge:
 - remove found clusters from current leaf set
 - remaining leaf set =: new cluster



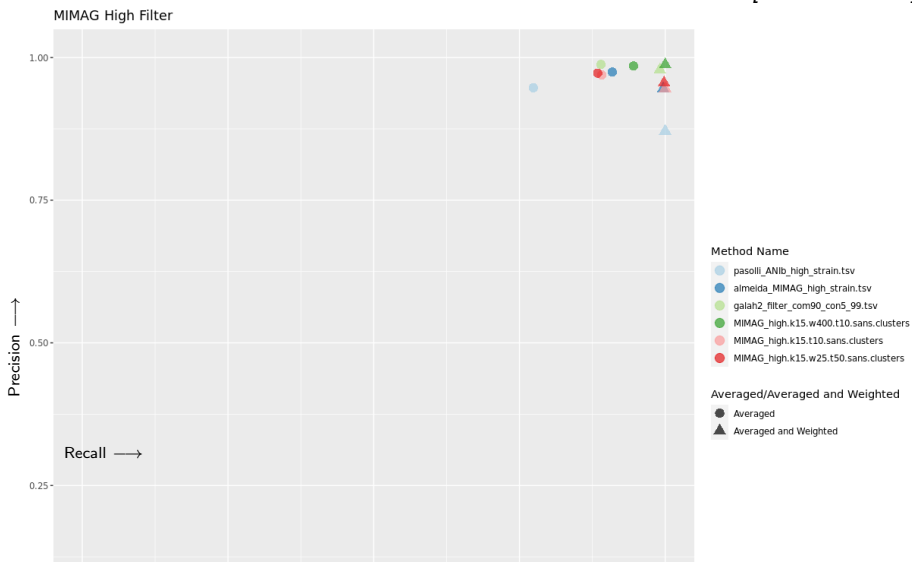
CAMI mouse gut metagenome [Sczyrba et al.]

Simulated dataset representing 64 metagenome samples.

	All	MIMAG medium	MIMAG high
contamination		< 10 %	< 5 %
completeness		≥ 50 %	> 90 %
# MAGs	11 602 (23 GB)	5 786	2 510
# genomes	791	686	349
# species	509	448	271

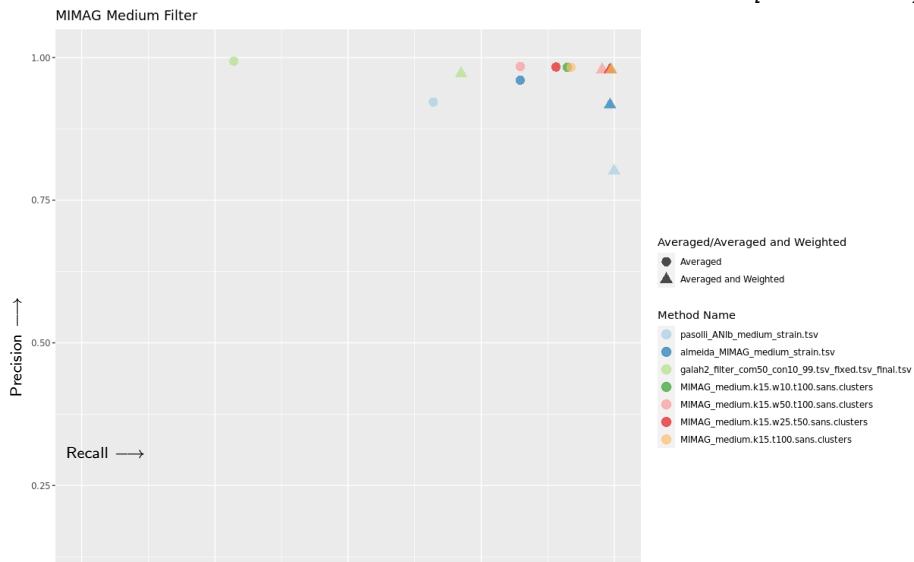
Clustering MAGS – Strain level

[Peter Belmann]



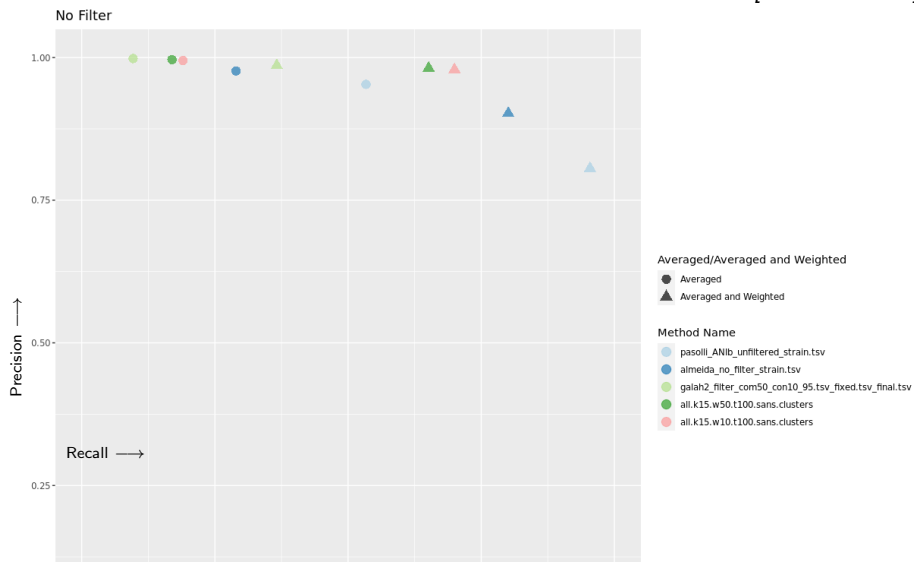
Clustering MAGS – Strain level

[Peter Belmann]



Clustering MAGS – Strain level

[Peter Belmann]



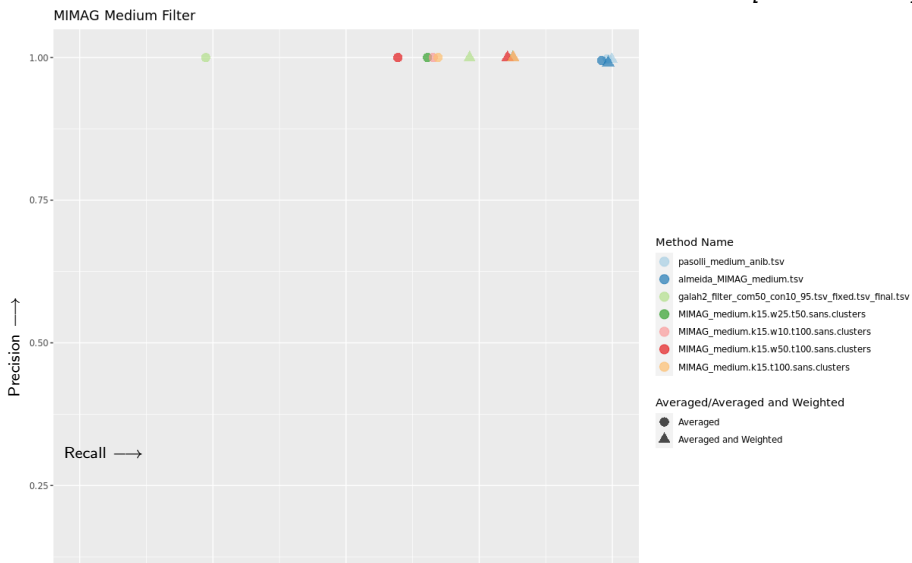
Clustering MAGS – Species level

[Peter Belmann]



Clustering MAGS – Species level

[Peter Belmann]



Clustering MAGS – Species level

[Peter Belmann]

