sylph: metagenomic profiling and containment by statistical k-mer sketching

**Jim Shaw**[1] and Yun William Yu[2]

Introduction

k-mer sketching

Statistical adjustment by ZIP model

Profiling

ANI-based MWAS

Flexible database profiling

Conclusion

Supp. Figs

# sylph: metagenomic profiling and containment by statistical k-mer sketching

**Jim Shaw**[1] and Yun William Yu[2]

[1]University of Toronto [2]Carnegie Mellon University

Genome Informatics 2023

UNIVERSITY OF
**TORONTO**

# Shotgun sequencing of sample

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

**Metagenomics** - analyzing shotgun sequences of an environmental sample

# Focus of the talk

sylph: metagenomic profiling and containment by statistical k-mer sketching

**Jim Shaw**[1] and Yun William Yu[2]

Introduction

k-mer sketching

Statistical adjustment by ZIP model

Profiling

ANI-based MWAS

Flexible database profiling

Conclusion

Supp. Figs

- **Database approach**

# Focus of the talk

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

- **Database approach**
- What genomes are in my sample?

# Focus of the talk

- **Database approach**
- What genomes are in my sample?
- Two distinct but similar approaches: **profiling** and **containment**

# What is metagenomic profiling?

sylph: metagenomic profiling and containment by statistical k-mer sketching

**Jim Shaw**[1] and Yun William Yu[2]

Introduction

k-mer sketching

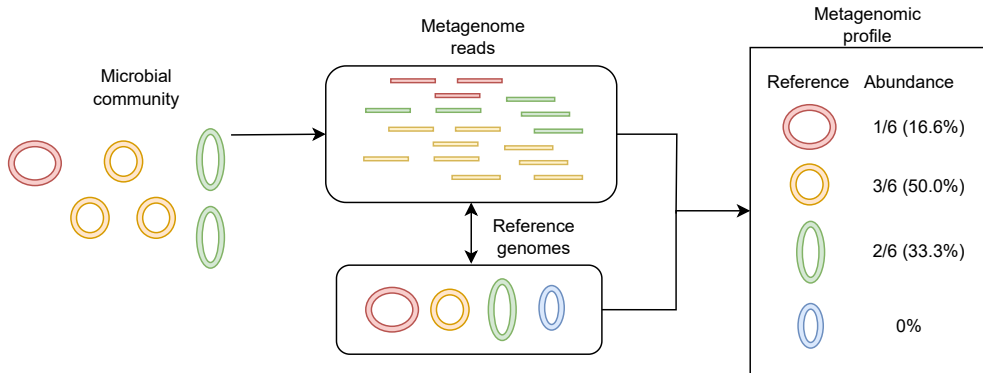Statistical adjustment by ZIP model

Profiling

ANI-based MWAS

Flexible database profiling

Conclusion

Supp. Figs

Profiling: **What taxa** are in the community and **how abundant** are they**?**

# What is containment estimation?

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

**Average nucleotide identity** (ANI): genome-to-genome similarity
**containment ANI**: genome-to-metagenome similarity (nearest neighbor)

# Containment is a continuous measure

sylph: metagenomic profiling and containment by statistical k-mer sketching

**Jim Shaw**[1] and Yun William Yu[2]

Introduction

k-mer sketching

Statistical adjustment by ZIP model

Profiling

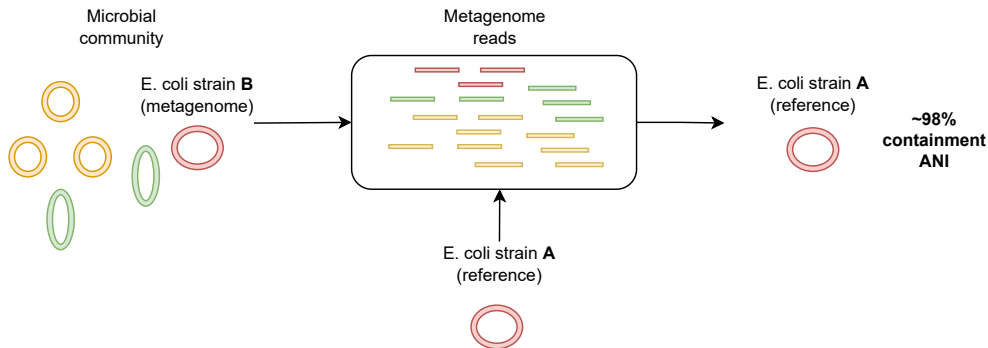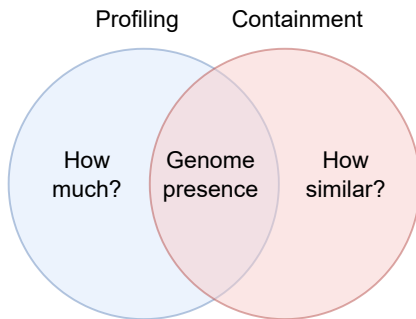ANI-based MWAS

Flexible database profiling

Conclusion

Supp. Figs

Containment: How **similar** is a genome to the genomes in the community?

Microbial community

E. coli strain **B** (metagenome)

Metagenome reads

E. coli strain **A** (reference)

~98% containment ANI

E. coli strain **A** (reference)

# Profiling vs containment

**Jim Shaw**[1] and Yun William Yu[2]

1. **Profiling** - relative abundances of genomes/taxa
2. **Containment** - nucleotide similarity of genome within metagenome

# sylph - containment and profiling

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS
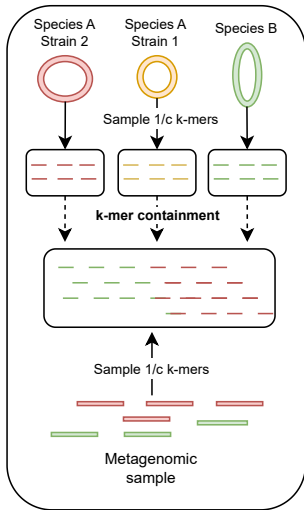
Flexible database
profiling

Conclusion

Supp. Figs

We present **sylph**, a new *profiler* with *containment* capabilities

# sylph part 1 - k-mer sketching

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Given a genome or a metagenome:

- Take all k-mers, sample only $1/c$ of them using FracMinHash (Irber et al., 2022). Default $c = 200$.

# sylph sketching

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

(1) Sketching k-mers and containment

Species A Strain 2 · Species A Strain 1 · Species B

Sample 1/c k-mers

**k-mer containment**

Sample 1/c k-mers

Metagenomic sample

# Estimating containment ANI

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
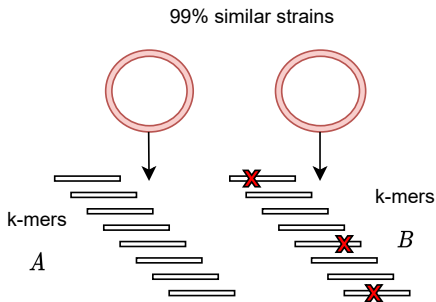adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

### k-mers differ between different genomes



99% similar strains

k-mers

$A$

k-mers

$B$

# Estimating containment ANI

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Fact:

$$ANI \approx \left( \frac{|A \cap B|}{|A|} \right)^{1/k}$$



99% similar strains

k-mers

$A$

k-mers

$B$

# Low abundance genomes

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

Problem: this model had no **read sampling**. k-mers are missing in low
abundance genomes due to read sampling

# Low abundance genomes

$|A \cap B|/|A|$ **underestimates** ANI when k-mers are missing!



99% similar strains

k-mers in ref.
genome

$A$

k-mers from read
sampling

$B$

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

# Zero-inflated Poisson model

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching
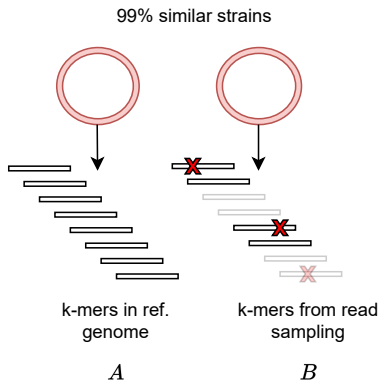
Statistical
adjustment by ZIP
model

Profiling

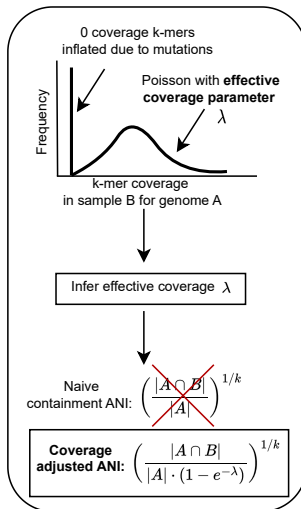ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

(2) k-mer ANI with coverage adjustment

99% similar strains

k-mers in ref.
genome

$A$

k-mers from read
sampling

$B$

0 coverage k-mers
inflated due to mutations

Poisson with **effective
coverage parameter**
$\lambda$

Frequency

k-mer coverage
in sample B for genome A

Infer effective coverage $\lambda$

Naive
containment ANI: $\left(\dfrac{|A \cap B|}{|A|}\right)^{1/k}$

**Coverage
adjusted ANI:** $\left(\dfrac{|A \cap B|}{|A| \cdot (1 - e^{-\lambda})}\right)^{1/k}$

# Coverage adjustment

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

1. Estimate the true coverage parameter $\lambda$

# Coverage adjustment

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

1. Estimate the true coverage parameter $\lambda$
2. $\hat{\lambda} = \frac{\text{\# k-mers with multiplicity } a+1}{\text{\# k-mers with multiplicity } a} \cdot (a+1)$ (similar to Skmer, Sarmashghi et al., 2019)

# Coverage adjustment

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

1. Estimate the true coverage parameter $\lambda$
2. $\hat{\lambda} = \frac{\# \text{ k-mers with multiplicity } a+1}{\# \text{ k-mers with multiplicity } a} \cdot (a+1)$ (similar to Skmer, Sarmashghi et al., 2019)
3. Coverage adjusted ANI:

$$\left( \frac{|A \cap B|}{|A| \cdot (1 - e^{-\hat{\lambda}})} \right)^{1/k}$$

# Coverage adjustment

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

1. Estimate the true coverage parameter $\lambda$
2. $\hat{\lambda} = \frac{\text{\# k-mers with multiplicity a+1}}{\text{\# k-mers with multiplicity a}} \cdot (a+1)$ (similar to Skmer, Sarmashghi et al., 2019)
3. Coverage adjusted ANI:

$$\left( \frac{|A \cap B|}{|A| \cdot (1 - e^{-\hat{\lambda}})} \right)^{1/k}$$

4. Intuition: small coverage $\implies$ denominator is small, pushes ANI upwards

# Coverage adjusted ANI - synthetic experiment

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
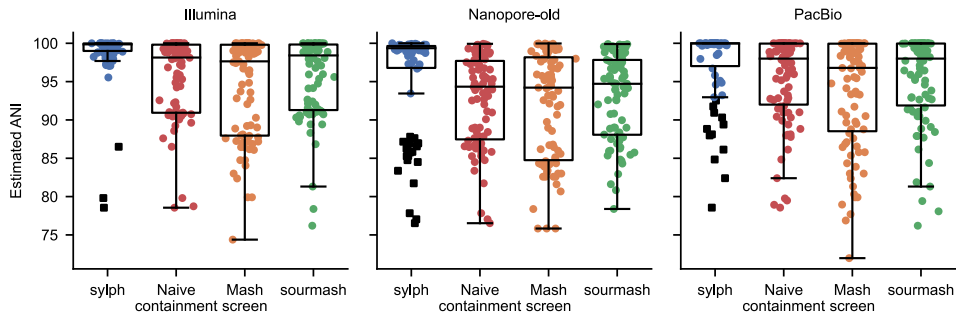profiling

Conclusion

Supp. Figs

Figure: Containment ANI against synthetic reads from a *Klebsiella pneumoniae* genome.

# Coverage adjusted ANI - real experiment

Figure: **Real** reads for mock metagenome from Meslier et al. (2022) with known references; black squares have uncorrected ANI.

# sylph can do containment... but profiling?

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

- containment ANI - doesn't say **how abundant** a microbe is

# sylph can do containment... but profiling?

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
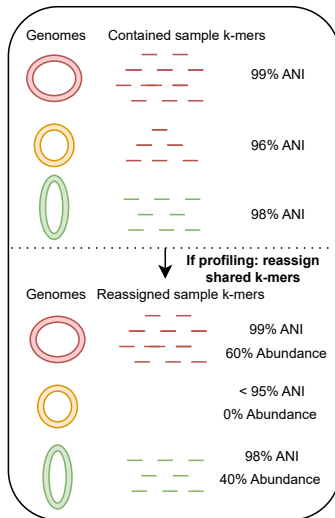model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

- containment ANI - doesn't say **how abundant** a microbe is
- **Problem**: k-mers are shared between genomes... which genome does a k-mer belong to?

# Reassigning k-mers for profiling

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
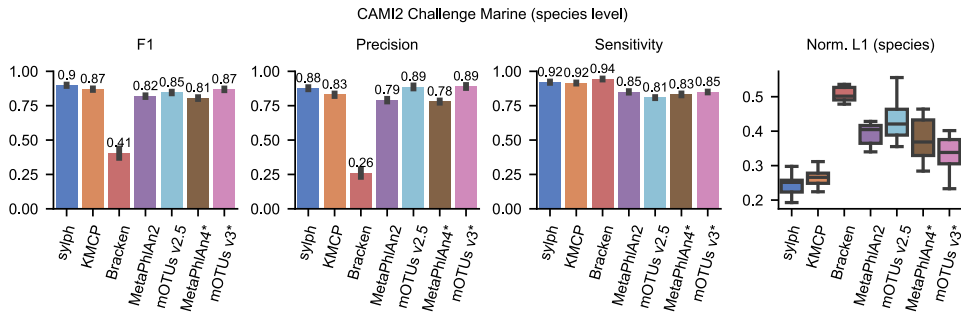profiling

Conclusion

Supp. Figs

(3) Top - ANI querying (sylph **query**)
Bot - Taxonomic profiling (sylph **profile**)

# Synthetic metagenome - CAMI2 Marine

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Figure: CAMI2 marine metagenome profiling challenge.

# Synthetic metagenome - varying ANI

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS
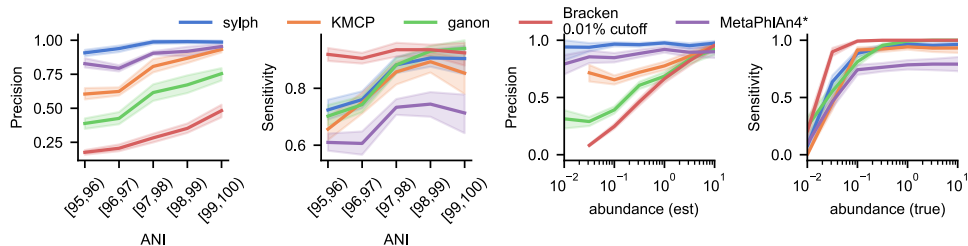
Flexible database
profiling

Conclusion

Supp. Figs

Figure: sylph retains high precision for lower ANI and abundance microbes (species level classification).

# Fast and efficient

sylph: metagenomic profiling and containment by statistical k-mer sketching

**Jim Shaw**[1] and Yun William Yu[2]

Introduction

k-mer sketching

Statistical adjustment by ZIP model

**Profiling**

ANI-based MWAS

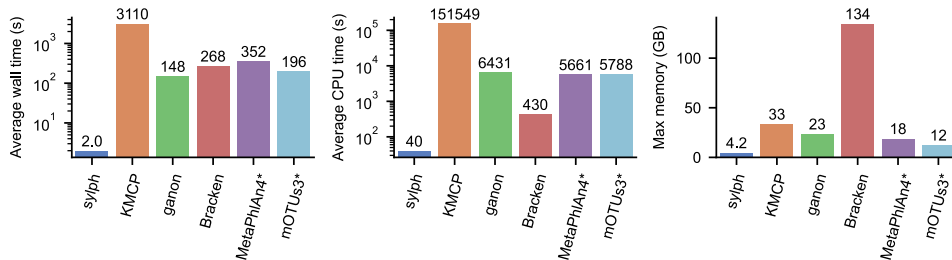Flexible database profiling

Conclusion

Supp. Figs

Figure: Runtime/memory on 200 genome synthetic community (3 Gbp, 2x150bp). 50 threads.

# Why is it so fast?

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

**Profiling**

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

- sylph does not classify reads (unlike Kraken et al.)
- sylph does not align reads (unlike MetaPhlAn, mOTUs)
- sylph shares one database for multi-sample profiling (unlike Kraken)
- Engineering (uses AVX2 instructions for k-mer sketching, etc)

# Applications

sylph: metagenomic profiling and containment by statistical k-mer sketching

**Jim Shaw**[1] and Yun William Yu[2]

Introduction

k-mer sketching

Statistical adjustment by ZIP model

Profiling

ANI-based MWAS

Flexible database profiling

Conclusion

Supp. Figs

**Applications**

# Wallen et al. Parkinson's Disease Metagenomics

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

- Wallen et al. (2022) performed gut metagenome wide association study (MWAS) for 490 Parkinson's Disease (PD) and 234 controls

# Wallen et al. Parkinson's Disease Metagenomics

**Jim Shaw**[1] and Yun William Yu[2]

- Wallen et al. (2022) performed gut metagenome wide association study (MWAS) for 490 Parkinson's Disease (PD) and 234 controls
- They used differential abundance testing → p-value

# Wallen et al. Parkinson's Disease Metagenomics

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

- Wallen et al. (2022) performed gut metagenome wide association study (MWAS) for 490 Parkinson's Disease (PD) and 234 controls
- They used differential abundance testing → p-value
- **What we did**: differential containment ANI → p-value

# Wallen et al. Parkinson's Disease Metagenomics

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

Jim Shaw[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

- Wallen et al. (2022) performed gut metagenome wide association study (MWAS) for 490 Parkinson's Disease (PD) and 234 controls
- They used differential abundance testing $\rightarrow$ p-value
- **What we did**: differential containment ANI $\rightarrow$ p-value
- Queried *289,232* genomes (UHGG) against 5.5 tb of reads; took a $\approx$ 3 hours with 40 threads

# Containment ANI MWAS

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

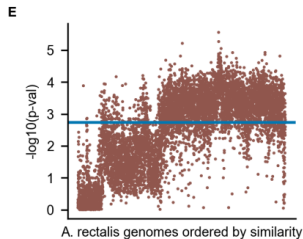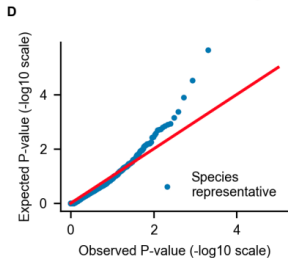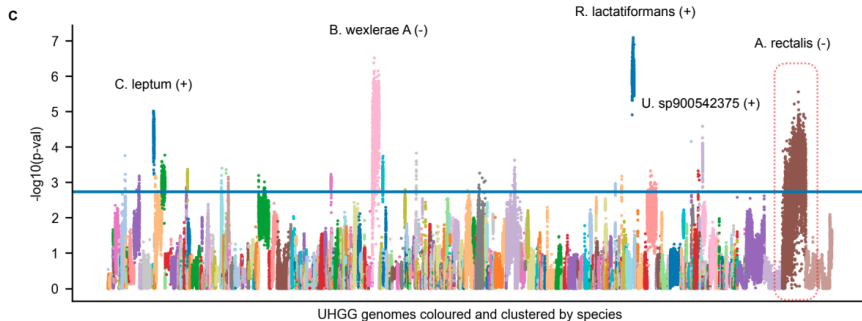k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

# Results

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

- Results are concordant with Wallen et al.

# Results

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

- Results are concordant with Wallen et al.
- Butyrate-producing bacteria (*F. prausnitzii*, *A. rectalis*) depleted in PD

# Results

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

- Results are concordant with Wallen et al.
- Butyrate-producing bacteria (*F. prausnitzii*, *A. rectalis*) depleted in PD
- Previous study (Becker et al. 2022) showed:

$$\frac{\text{B. wexlerae abundance}}{\text{R. lactatiformans abundance}} \propto \text{fecal butyrate concentration}$$

# Flexible database profiling

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

- sylph allows for flexible database choice

# Flexible database profiling

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

- sylph allows for flexible database choice
- Viruses and eukaryotes can be profiled too

# Virome profiling and customized databases vs RefSeq

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

Comprehensive databases:

- ▸ GTDB-R214: 85,205 prokaryotic species genomes
- ▸ IMG/VR4: 2,917,521 species genomes

RefSeq:

- ▸ RefSeq representative prokaryotic: 18,325 genomes
- ▸ RefSeq viral: 14,993 genomes
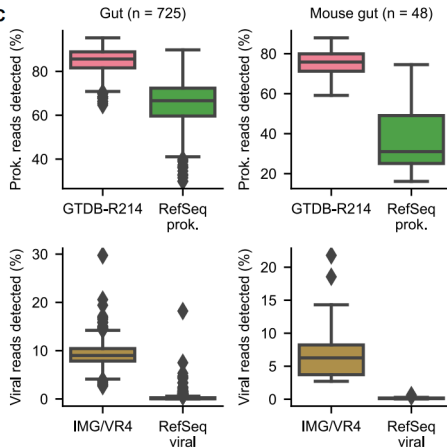
# Virome profiling and customized databases vs RefSeq

sylph: metagenomic profiling and containment by statistical k-mer sketching

**Jim Shaw**[1] and Yun William Yu[2]

Introduction

k-mer sketching

Statistical adjustment by ZIP model

Profiling

ANI-based MWAS
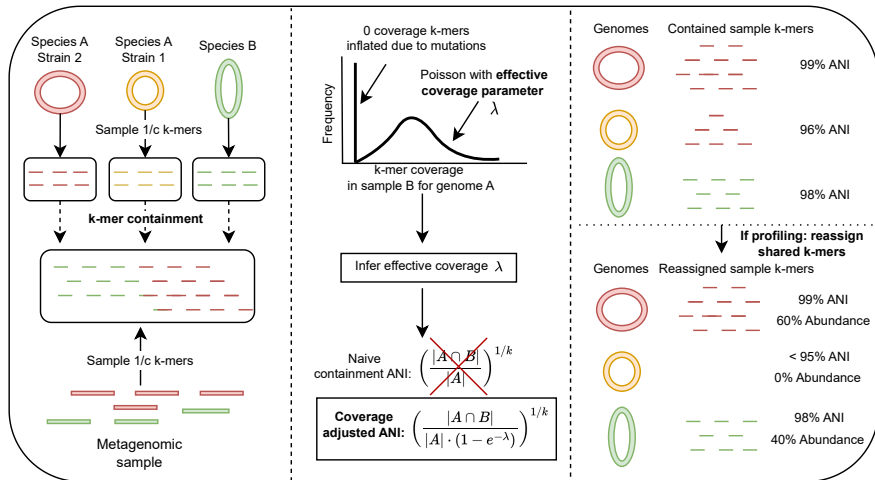
**Flexible database profiling**

Conclusion

Supp. Figs

# Conclusion

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Presented **sylph**

- Statistical k-mer sketching approach for **containment** and **profiling**



Profiling     Containment

How much? | Genome presence | How similar?

# Algorithm recap

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

# Funding and acknowledgements

**Metagenome profiling and containment estimation through abundance-corrected k-mer sketching with sylph** now available on bioRxiv.



Figure: Github QR code, Jim Shaw (PhD student), Yun William Yu (Advisor)

Natural Sciences and Engineering Research Council of Canada

Conseil de recherches en sciences naturelles et en génie du Canada

Canadá

UNIVERSITY OF TORONTO

Carnegie Mellon University

# Synthetic metagenome - CAMI2 Marine

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

Figure: CAMI2 marine metagenome challenge profiling.
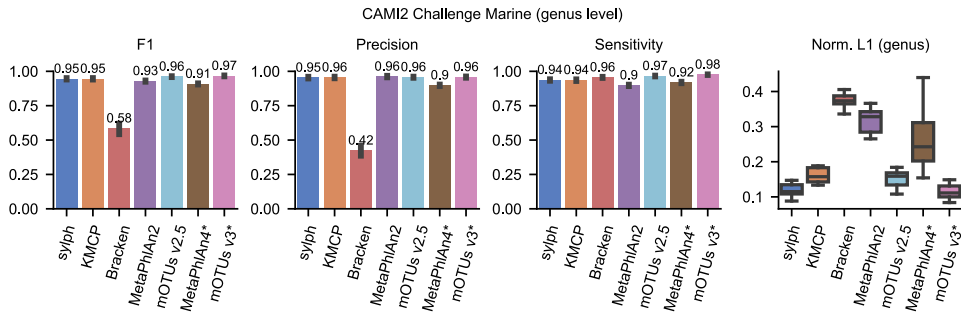
# Synthetic metagenome - CAMI2 Strain Madness

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

Figure: CAMI2 strain madness challenge.

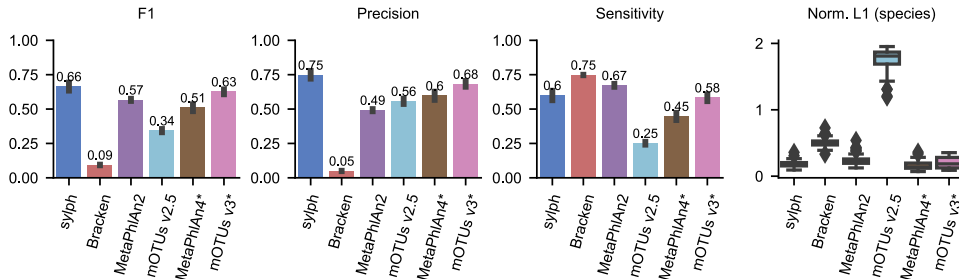# Synthetic metagenome - CAMI2 Strain Madness

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

Figure: CAMI2 strain madness metagenome challenge profiling.

# Profiling real reads

sylph:
metagenomic
profiling and
containment by
statistical k-mer
sketching

**Jim Shaw**[1] and
Yun William Yu[2]

Introduction

k-mer sketching

Statistical
adjustment by ZIP
model

Profiling

ANI-based MWAS

Flexible database
profiling

Conclusion

Supp. Figs

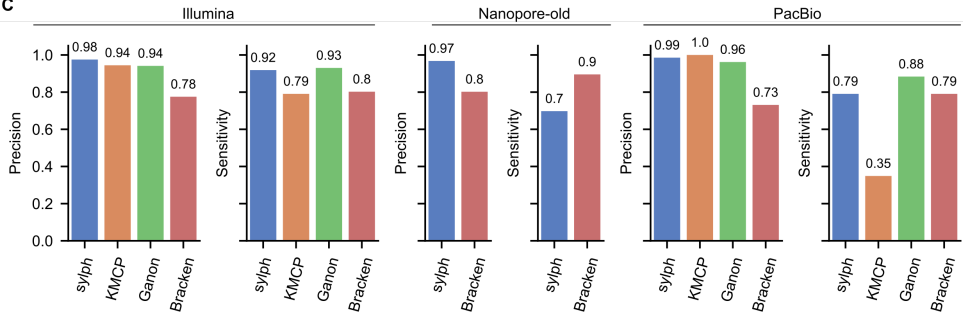Figure: Mock community profiling from Meslier et al. (2022)