

CAMI tutorial

assessing metagenomics software with the CAMI benchmarking toolkit

Jie Zhu

zhujie@genomics.cn

2021-03-15 Mon

BGI Research

Outline

Introduction

Challenges

Datasets

CAMI 1st

CAMI 2nd

Reference

Introduction

Challenges

Datasets

CAMI 1st

CAMI 2nd

Reference

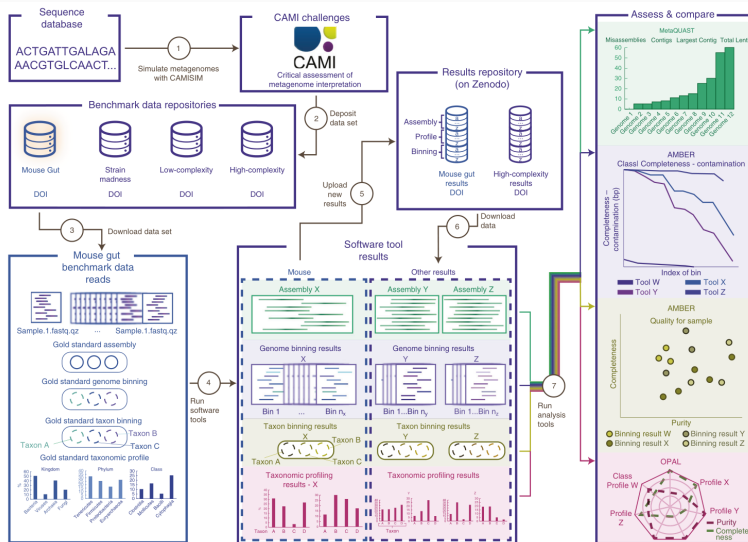
What is CAMI?

In just over a decade, metagenomics has developed into a powerful and productive method in microbiology and microbial ecology. The ability to retrieve and organize bits and pieces of genomic DNA from any natural context has opened a window into the vast universe of uncultivated microbes. Tremendous progress has been made in computational approaches to interpret this sequence data but none can completely recover the complex information encoded in metagenomes. A number of challenges stand in the way. Simplifying assumptions are needed and lead to strong limitations and



potential inaccuracies in practice. Critically, methodological improvements are difficult to gauge due to the lack of a general standard for comparison. Developers face a substantial burden to individually evaluate existing approaches, which consumes time and computational resources, and may introduce unintended biases. The Critical Assessment of Metagenome Interpretation (CAMI) is a new community-led initiative designed to help tackle these problems by aiming for an independent, comprehensive and bias-free evaluation of methods

Overview



Assembly challenge

What is an assembly method?

An assembly method returns longer nucleotide sequences derived by puzzling together individual sequencing reads. These sequences are assumed to represent contiguous stretches from one genome included in the microbiome sample that was sequenced.

What is a profiling method?

A profiling method returns an estimate for the frequencies of different taxa in a sequenced microbial community based on analysis of the sequence sample. The main output is a vector with relative abundances for the different sample taxa. The relative abundances of taxa from the same 'rank' of the taxonomy (e.g. superkingdom, including archaea, bacteria and eukaryotes) cannot sum up to more than 1.

What is a binning method?

A binning method assigns an identifier to every sequence in a sequence sample, where the total number of identifiers is ideally less than the total number of sequences. Thus the act of binning places the sequences into broader categories. A bin includes all the sequences with the same identifier. If these identifiers identify taxa from a taxonomy, the method is a taxonomic binning method.

Introduction

Challenges

Datasets

CAMI 1st

CAMI 2nd

Reference

1st CAMI Challenge Dataset 1: CAMI low diversity

summary

Platform	Hiseq
Number of samples	1
Total Size	15 Gbp
Read length	2x150 bp
Insert size mean	270 bp
Insert size stddev	27 bp

list

```
java -jar ~/.local/bin/camiClient.jar -l \  
https://openstack.cebitec.uni-bielefeld.de:8080/swift/v1/CAMI_I_L  
CAMI_low_RL_S001__insert_270_GoldStandardAssembly.fasta.gz  
RL_S001__insert_270.fq.gz  
gold_standard_low_single.fasta.gz  
goldstandard_low_1.filtered.profile  
goldstandard_low_1.profile  
gs_read_mapping.binning.gz  
gsa_mapping.binning  
novelty_complete.tsv  
source_genomes_low.tar.gz  
taxonomy.tar.gz  
unique_common.tsv
```

1st CAMI Challenge Dataset 2: CAMI medium diversity

summary

Platform	Hiseq
Number of samples	2
Total Size	40 Gbp
Read length	2x150 bp
Insert size mean	270 bp and 5kbp
Insert size stddev	10%

list

cami_list https://openstack.cebitec.uni-bielefeld.de:8080/swift/v1/CAMI_I_MEDIUM

1st CAMI Challenge Dataset 3: CAMI high diversity

summary

Platform	Hiseq
Number of samples	5
Total Size	75 Gbp
Read length	2x150 bp
Insert size mean	270 bp
Insert size stddev	10%
Diversity	High

list

`cami_list` https://openstack.cebitec.uni-bielefeld.de:8080/swift/v1/CAMI_I_HIGH

CAMI1 Toy Test Dataset: Low Complexity

summary

Platform	simulated from public genomes(30)
Number of samples	
Total Size	15 Gbp
Read length	2x100 bp
Insert size mean	180 bp
Insert size stddev	10%
Diversity	Low

list

`cami_list` https://openstack.cebitec.uni-bielefeld.de:8080/swift/v1/CAMI_I_TOY_LOW

CAMI1 Toy Test Dataset: Medium Complexity

summary

Platform	simmlater from public genomes(225)
Number of samples	2
Total Size	
Read length	
Insert size mean	180 bp and 5kb
Insert size stddev	
Diversity	Medium

list

cami_list https://openstack.cebitec.uni-bielefeld.de:8080/swift/v1/CAMI_I_TOY_MEDIUM

CAMI1 Toy Test Dataset: High Complexity

summary

Platform	simulated from public genomes(450)
Number of samples	5
Total Size	75Gbp
Read length	2x100 bp
Insert size mean	180bp
Insert size stddev	18bp
Diversity	High

list

`cami_list` https://openstack.cebitec.uni-bielefeld.de:8080/swift/v1/CAMI_I_TOY_HIGH

CAMI2: Rhizosphere challenge

short reads

Platform	simulated Illumina HiSeq metagenome data
Number of samples	21
Total size	105 Gb
Read length	2x150 bp
Insert size mean	270 bp
Insert size s.d.	20 bp

long reads

- Pacific Bioscience

Platform	simulated Pacific Bioscience metagenome data
Number of samples	21
Total size	105 Gb
Average read length	3,000 bp
Read length s.d.	1,000 bp

- Oxford Nanopore

Platform	simulated Oxford Nanopore metagenome data
Number of samples	21
Total size	105 Gb
Average read length	1,610 bp
Read length s.d.	~3,000 bp

taxonomy database

[https://openstack.cebitec.uni-bielefeld.de:](https://openstack.cebitec.uni-bielefeld.de:8080/swift/v1/CAMI_DATABASES/taxdump_cami2_toy.tar.gz)

[8080/swift/v1/CAMI_DATABASES/taxdump_cami2_toy.tar.gz](https://openstack.cebitec.uni-bielefeld.de:8080/swift/v1/CAMI_DATABASES/taxdump_cami2_toy.tar.gz)

CAMI2: Clinical pathogen detection challenge

CAMI2: Toy Human Microbiome Project Dataset

Simulated endogenous data from five different body sites of the human host, namely gastrointestinal tract, oral cavity, armpits, skin and urogenital tract.

Platform	simulated Illumina HiSeq metagenome data
Number of samples	68 (10 ES tract, 32 oral cavity, 10 uterine, 10 skin, 9 congenital tract)
Total size	248 Gbp
Read length	2x100 bp
Insert size mean	270 bp
Insert size stdev	20 bp

Platform	simulated Pacific Bionessence monitoring data
Number of samples	40 (10 CI tract, 10 real cavity, 10 airways, 10 skin, 9 ungeriatric tract)
Total size	288 Gbp
Average read length	3000 bp
Read length	2000 bp

[illegible][illegible]

In every sample folder there are three subfolders, bam, contigs and reads

folded is a harm file for every genome for which at least one read was produced, which is uniquely indicated by a combination of OTU and a running ID number for the number of genomes included in that OTU in the sample. OTU_{id}.

The contigs folder contains the gold standard assembly for that particular sample. It contains two files, the gold standard in fasta format: `anonymous22,0222,20.fa`. And the mapping for each contigs to its reference genome. In this case: `anonymous22,0222,20.map`.

The reads folder contains the sorted reads for this sample. It contains two files, one with the `read` identifiers, containing both ends for paired-end sequencing and with anonymous names: `anonymous_{taskid}.id.gz`. And the second one is a mapping of every single read to the genome it originated from and the original read ID (you anonymized): `reads_{taskid}.map.gz`.

every data set contains one abundance file per sample mapping OTUs to genomic abundance[†] but

every data set contains the pooled gold standard assembly over all samples in the whole AnonymizedGenomicData set

```

reverting (the state set) config.ini
reverting (the state set) config.ini

```

(`bin/genome/ucsc`) genomes. This folder contains all the fasta files of the downloaded genomes.
`genome_human`

Since the contigs are anonymized, a file mapping each contig to its genome/baculovirus id and position in the reference sequence is provided separately to the user.

For each input GTU (genetic transferable unit), two columns are suggested: One of the first on which the GTU was mapped to the NCBI and one to the specifically downloaded genome, contains a nearby *usage* column in case new genomes are provided, otherwise this column is "unknown" and can be ignored (see table 1a).

In the top folder "Hybrid" there are assembly and mapping gold standards created from both the short and long read data sets. For every sample as well as all samples pooled, the bam files of short and long read simulators (as described for the "bam" subfolder above) are merged and the gold standards calculated the same way as for the individual short or long read samples.

CAMI2: Toy Mouse Gut Dataset

Introduction

Challenges

Datasets

CAMI 1st

CAMI 2nd

Reference