Computational Molecular Biology and Bioinformatics Metagenomics

Malay Bhattacharyya

Assistant Professor

Machine Intelligence Unit Indian Statistical Institute, Kolkata January, 2022

What is m	etagenomics?		
Introduction	Microbiome Analysis in Mass-transit Systems	Software and algorithms	Hands-on
●00000	0000000000	O	O

Metagenomics is defined as the direct genetic analysis of genomes (microbial communities) contained within a site (environmental, clinical, built-in, etc.).

Metagenomics allows for a detailed study of the diversity of communities, and therefore to clarify the mechanisms of their functioning, to determine the metabolic relationships.

Note: In Greek, meta means "beyond or above the range of normal or physical human experience".

What is mi	crobiome?		
Introduction	Microbiome Analysis in Mass-transit Systems	Software and algorithms	Hands-on
00000		O	O

• • = • • = •

What is mi	crobiome?		
Introduction	Microbiome Analysis in Mass-transit Systems	Software and algorithms	Hands-on
00000		O	O

• **Commensal:** An organism that uses food supplied in the internal or the external environment of the host, without establishing a close association with the host. E.g., Staphylococcus epidermidis found on human skin.

医子宫医子宫医

What is r	nicrohiome?		
00000	000000000	0	0
Introduction	Microbiome Analysis in Mass-transit Systems	Software and algorithms	Hands-on

- **Commensal:** An organism that uses food supplied in the internal or the external environment of the host, without establishing a close association with the host. E.g., Staphylococcus epidermidis found on human skin.
- **Symbiotic:** An organism that lives in beneficial association with the host. E.g., Bacteroides thetaiotaomicron found in human intestine.

伺 ト イ ヨ ト イ ヨ ト

What is	microhiome?		
00000	000000000	0	0
Introduction	Microbiome Analysis in Mass-transit Systems	Software and algorithms	Hands-on

- **Commensal:** An organism that uses food supplied in the internal or the external environment of the host, without establishing a close association with the host. E.g., Staphylococcus epidermidis found on human skin.
- **Symbiotic:** An organism that lives in beneficial association with the host. E.g., Bacteroides thetaiotaomicron found in human intestine.
- **Pathogenic:** An organism which is capable of causing diseases in a host. E.g., SARS-CoV-2 causing COVID-19 in human.

・吊り ・ ヨト ・ ヨト

 Introduction
 Microbiome Analysis in Mass-transit Systems
 Software and algorithms
 Hands-on

 000000
 0000000000
 0
 0
 0

What is microbiome?



イロト 不得 トイヨト イヨト 二日

Human	microhiome		
000000	000000000	0	0
Introduction	Microbiome Analysis in Mass-transit Systems	Software and algorithms	Hands-on

The human microbiome is the collection of microorganisms that reside on or within human tissues and biofluids corresponding to anatomical sites.

伺 ト イヨト イヨト

Microbiome Analysis in Mass-transit Systems

Software and algorithms

Hands-on 0

Human microbiome

The human microbiome is the collection of microorganisms that reside on or within human tissues and biofluids corresponding to anatomical sites.



Microbiome Analysis in Mass-transit Systems

Software and algorithms

Hands-on

The integrative Human Microbiome Project (iHMP)



Source: Nature, 569:641-648, 2019.

Introduction 00000● Microbiome Analysis in Mass-transit Systems

Software and algorithms 0 Hands-on O

Microbiome of the built environment

Microbiome of the built environment comprises the communities of microorganisms that reside in human constructed environments.

・ 同 ト ・ ヨ ト ・ ヨ ト

Microbiome Analysis in Mass-transit Systems

Software and algorithms 0 Hands-on 0

.⊒ →

Microbiome in mass-transit systems



Source: Cell, 184(13):3376-3393, 2021.

Microbiome Analysis in Mass-transit Systems

Software and algorithms 0 Hands-on 0

Availability of data

MetaSUB sequencing data is available at: https://pngb.io/metasub-2021

・ 同 ト ・ ヨ ト ・ ヨ ト

Microbiome Analysis in Mass-transit Systems

Software and algorithms 0 Hands-on O

Generating taxonomic profiles for samples

KrakenUniq (formerly KrakenHLL) is a novel metagenomics classifier that combines the fast k-mer-based classification of Kraken with an efficient algorithm for assessing the coverage of unique k-mers found in each species in a dataset.

On various test datasets, KrakenUniq gives better recall and precision than other methods and effectively classifies and distinguishes pathogens with low abundance from false positives in infectious disease samples.

<u>Note</u>: KrakenUniq requires a huge amount of primary memory (ideally 128-512 GB). For performing more memory efficient classification, consider using Centrifuge (ideally requiring 4-12 GB).

ヘロト ヘ河ト ヘヨト ヘヨト

Microbiome Analysis in Mass-transit Systems

Software and algorithms 0 Hands-on

Overview of the KrakenUniq algorithm

A Read k-mers are looked-up in the database and assigned to taxa:



B For each taxon a data sketch records its k-mers for cardinality estimation



K-mer count and coverage in taxonomic report show evidence behind classifications:

reads	kmers	dup	cov	taxID	rank	name		Bad classification
122	112	144	0.0004	11855	species	Clostridioides difficile	-	with few k-mers
9650	7129	74.5	0.192	10632	species	Human polyomavirus 2		- Cood classification
15	1570	1	0.0002	7643	species group	Mycobacterium tb complex		reads cover genome
	T		+					reads cover genome

Number of distinct k-mers for taxon, and coverage of the taxon's k-mers

Source: Genome Biology, 19:198, 2018.

イロト イポト イラト イラト

Microbiome Analysis in Mass-transit Systems

Software and algorithms

Taxonomic tree for metagenome-assembled genomes (MAGs)



Microbiome Analysis in Mass-transit Systems

Software and algorithms

Hands-on

Dimensionality Reduction (UMAP vs t-SNE vs PCA)

UMAP	t-SNE	PCA	
1. Uniform Manifold	1. t-distributed	1. Principal Compo-	
Approximation and	Stochastic Neighbour-	nent Analysis	
Projection	hood Embedding		
2. Unsupervised	2. Unsupervised	2. Unsupervised	
3. Non-linear	3. Non-linear	3. Linear	
4. Not deterministic	4. Not deterministic	4. Deterministic	
5. Captures the global	5. Captures the local	5. Captures the global	
structure of data	structure of data	structure of data	
6. Relatively faster	6. Relatively slower	6. Relatively faster	

Analysis with UMAP

UMAP is a dimension reduction technique that can be used for visualisation of data in a lower dimension. The algorithm is founded on three assumptions about the data:

- The data is uniformly distributed on Riemannian manifold
- The Riemannian metric is locally constant (or can be approximated as such)
- The manifold is locally connected

Analysis with UMAP

UMAP is a dimension reduction technique that can be used for visualisation of data in a lower dimension. The algorithm is founded on three assumptions about the data:

- The data is uniformly distributed on Riemannian manifold
- The Riemannian metric is locally constant (or can be approximated as such)
- O The manifold is locally connected

From these assumptions, it is possible to model the manifold with a fuzzy topological structure. The embedding is found by searching for a low dimensional projection of the data that has the closest possible equivalent fuzzy topological structure.

・ロト ・同ト ・ヨト ・ヨト -

Microbiome Analysis in Mass-transit Systems

Software and algorithms 0 Hands-on 0

Analysis with UMAP



Introduction Microbiome Analysis in Mass-transit Systems Software and algorithms Hands-on o

Microbial signatures



Schematic of GeoDNA representation using indexing on graphs

Malay Bhattacharyya Computational Molecular Biology and Bioinformatics

Microbiome Analysis in Mass-transit Systems

Software and algorithms

Hands-on

Prediction of features



Prediction accuracy of a random forest model for a given feature

Microbiome Analysis in Mass-transit Systems

Software and algorithms 0 Hands-on 0

Antimicrobial resistance

Using the MegaRES ontology and alignment software, one can map reads to known antibiotic resistance genes.



Co-occurrence of AMR genes (left), AMR genes by city (right)

Malay Bhattacharyya

Computational Molecular Biology and Bioinformatics

Software and algorithms

Tool	Availability
AdapterRemoval v2.17	https://github.com/mikkelschubert/adapterremoval
Bowtie2 v2.3.0	https://sourceforge.net/projects/bowtie-bio/files/
	bowtie2/2.3.0
BLASTn	https://ftp.ncbi.nlm.nih.gov/blast/executables/
	blast+/LATEST
KrakenUniq v0.3.2	https://github.com/fbreitwieser/krakenuniq
Centifuge	https://github.com/infphilo/centrifuge
MASH v2.1.1	https://github.com/marbl/Mash
HUMAnN2	https://pypi.org/project/humann2
DIAMOND v0.8.36	https://github.com/bbuchfink/diamond
metaSPAdes v3.8.1	https://github.com/ablab/spades/releases/tag/v3.8.1
MegaRes v1.0.1	https://megares.meglab.org/download/index.php
MetaBAT2 v2.12.1	https://anaconda.org/ursky/metabat2
CheckM v1.0.13	https://github.com/Ecogenomics/CheckM
dnadiff v1.3	https://github.com/mummer4/mummer
GTDB-Tk v1.0.2	https://github.com/jianshu93/GTDB_Tk
FastTree v2.1.10	https://anaconda.org/bioconda/fasttree
iTOL v5.5	https://itol.embl.de

Malay Bhattacharyya Computational Molecular Biology and Bioinformatics

・ロト ・ 御 ト ・ 臣 ト ・ 臣 ト …

æ

Hands-on

- Read the following paper and access its data.
 - 0
 - D. Danko et al., A global metagenomic map of urban microbiomes and antimicrobial resistance. Cell, 184(13), pp.3376-3393, 2021.

Link: https://www.sciencedirect.com/science/ article/pii/S0092867421005857.

・ 同 ト ・ ヨ ト ・ ヨ ト