

November 29, 2019
統合データベース講習会 : AJACS世田谷

微生物統合データベース MicrobeDB.jpの活用法

森 宙史
Hiroshi Mori
国立遺伝学研究所
情報研究系

<http://microbedb.jp/>



Text Analysis Statistics

 Search

Environment: hot spring

Taxonomy: Enterococcus faecalis

Taxonomy: Streptomyces avermitilis

Gene: psbA

ID: 29

2011年から公開している
原核生物を主とした微生物の統合データベース(DB)

MicrobeDB.jp v.3 project members

National Institute of Genetics: (Genome, Metagenome, Ontology)

Ken Kurokawa, Yasukazu Nakamura, Hiroshi Mori,
Takatomo Fujisawa, Eli Kaminuma (TMDU), Koichi Higashi

National Institute of Basic Biology: (Ortholog)

Ikuo Uchiyama, Hirokazu Chiba (DBCLS), Hiroyo Nishide

Tokyo Institute of Technology: (Metagenome)

Takuji Yamada, Zenichi Nakagawa

Chiba University: (Fungal & Bacterial culture collection info.)

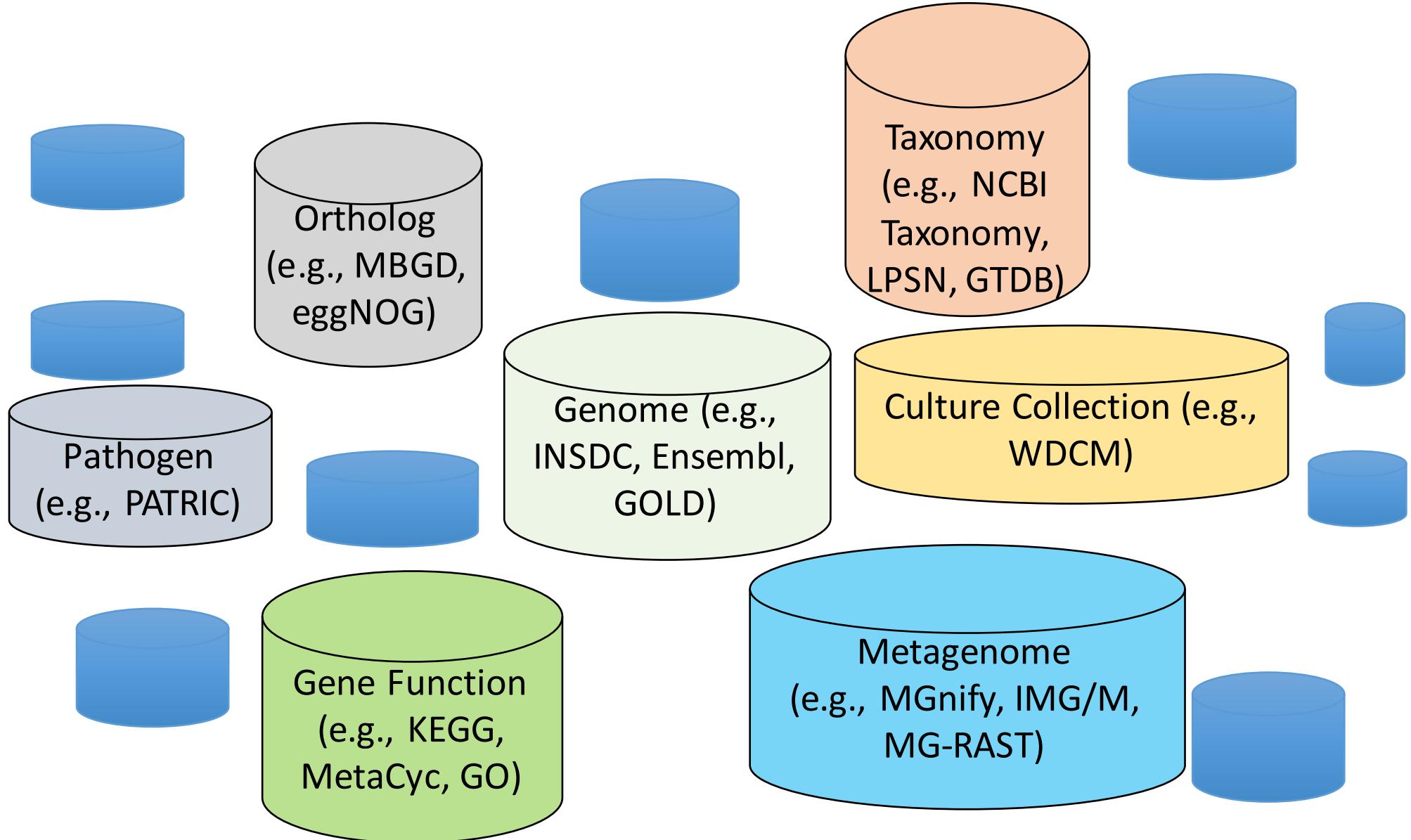
Hiroki Takahashi, Takashi Yaguchi

Technical adviser:

DBCLS (especially Shuichi Kawashima, Toshiaki Katayama)



微生物のDBは既に良い特化型のDBが多数存在



特化型のDBでは、異なる種類のデータ(ゲノムとメタゲノム等)を
関連付けて検索することは困難 → **統合DB**の出番

微生物統合DB MicrobeDB.jpの基本的な開発方針

- 遺伝子と系統のリンク
 - 系統が持つ遺伝子機能的な特徴を推定
- 遺伝子と表現型のリンク
 - 表現型を担う遺伝子を推定
- 遺伝子と環境のリンク
 - その環境で生息するために必要な遺伝子機能を推定
- 系統と環境のリンク
 - その環境に適応した系統を推定

既存のデータにこれらのリンクを付与し、リンクから微生物に関する新たな知見を得られるDBを目指す⁵



Figure 10.14 Microbiology: A Clinical Approach 2e © Garland Science 2016

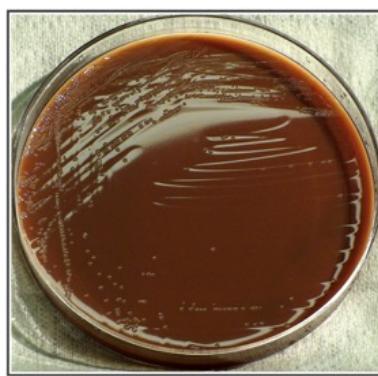
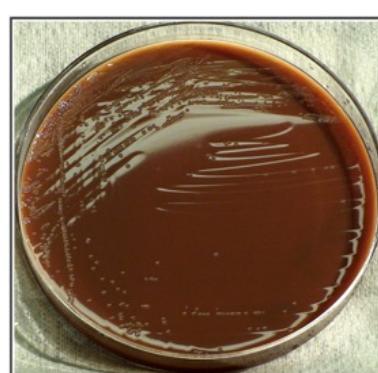


Figure 10.14 Microbiology: A Clinical Approach 2e © Garland Science 2016

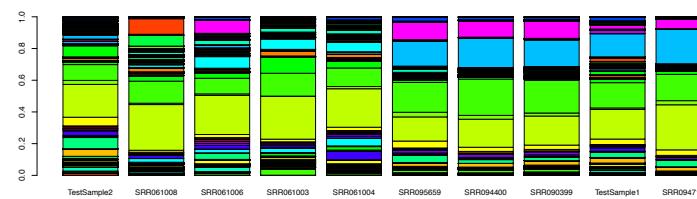
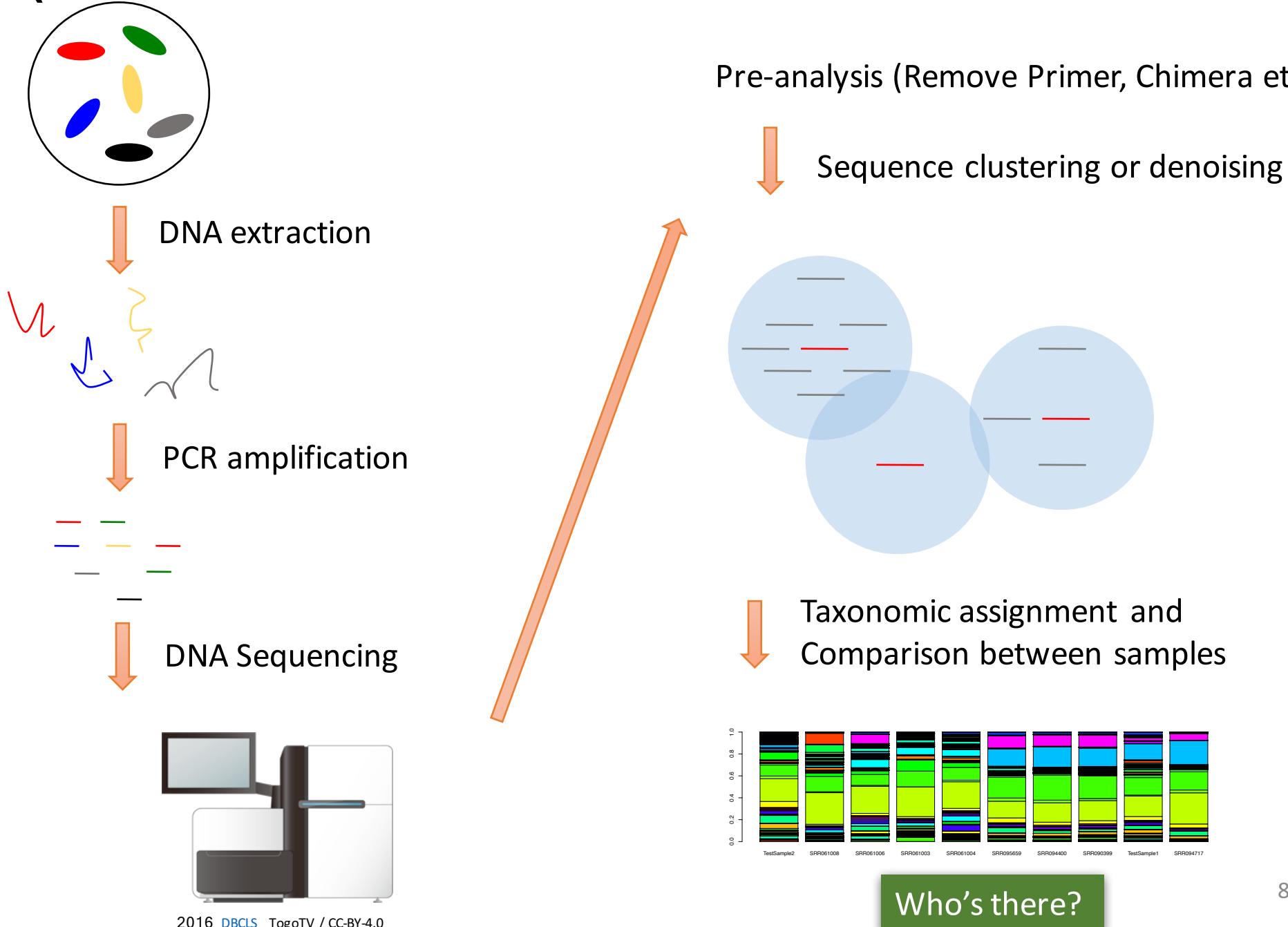


数%ぐらいの菌しか
培養できない

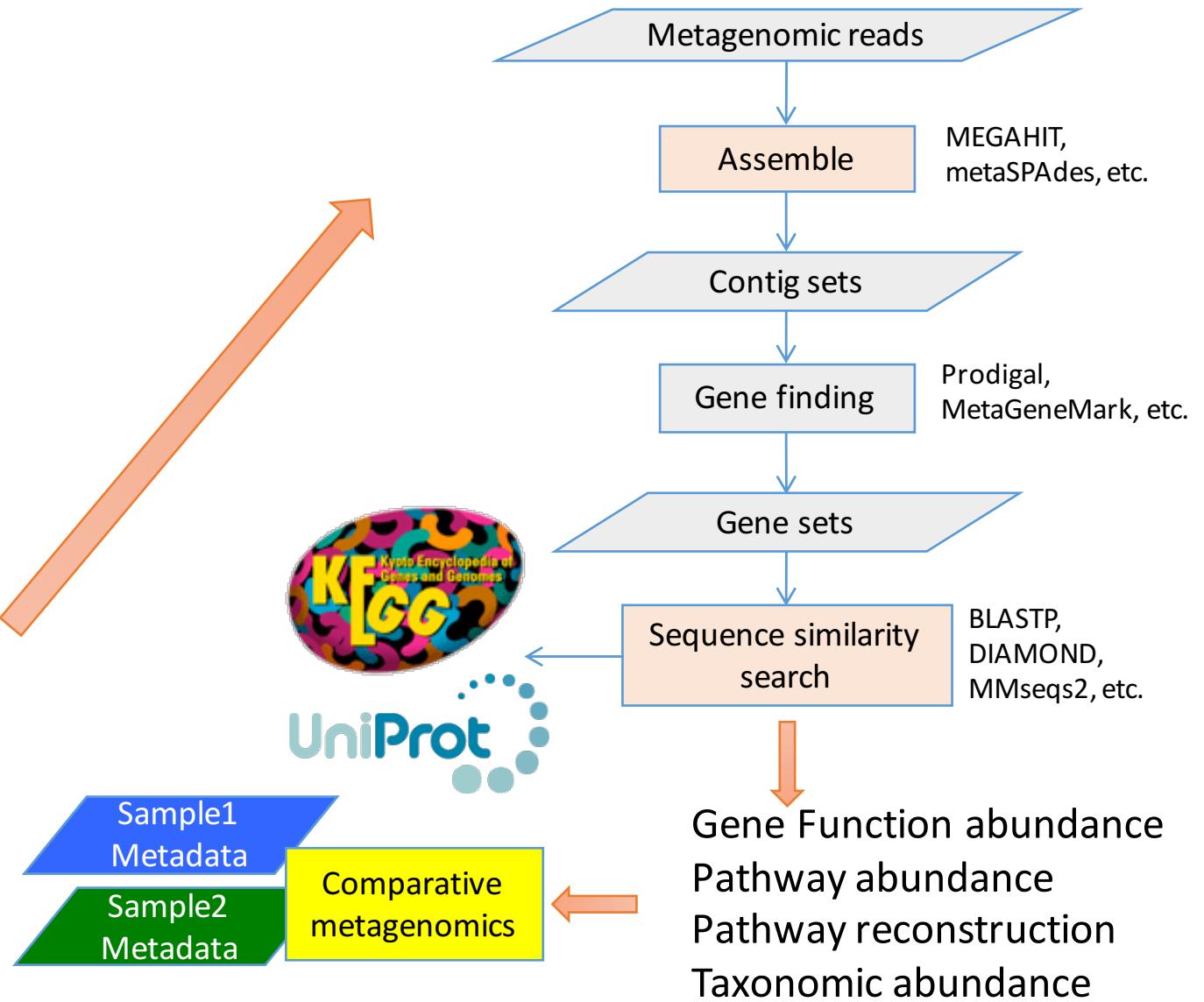
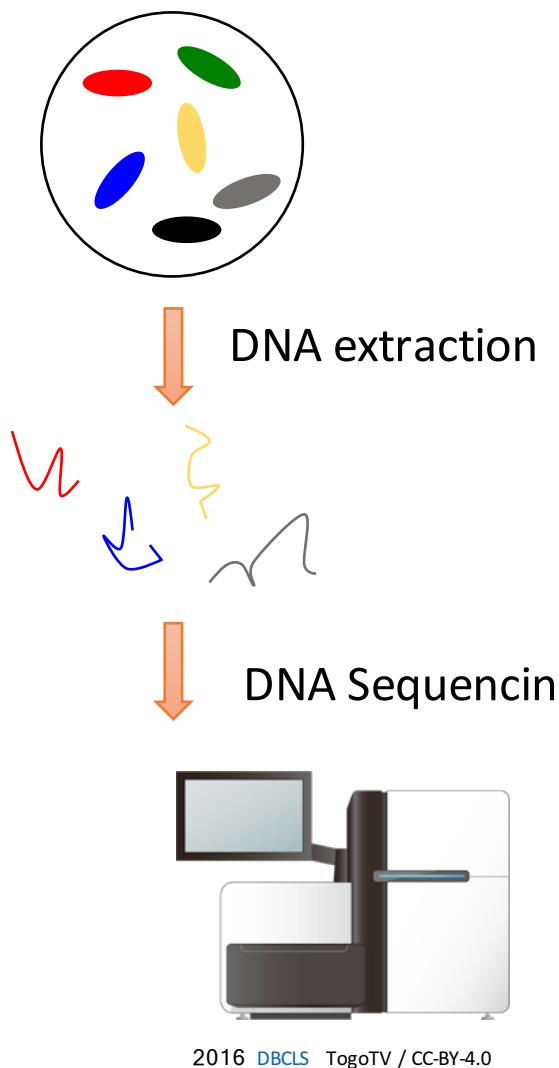
What's metagenomics?

- **Microflora, microbiota, microbial community:** 微生物群集
Total collection of microorganisms within a community
- **Metagenome:** ある群集の遺伝情報の総体
Total genomic potential of a community
[Handelsman et al. 1998, Chem. & Biol.]
- **Microbiome:** マイクロバイオーム
Microbiota and metagenome in a microbial community

amplicon sequencing analysis (アンプリコン解析, 16S rRNA遺伝子のアンプリコン解析)

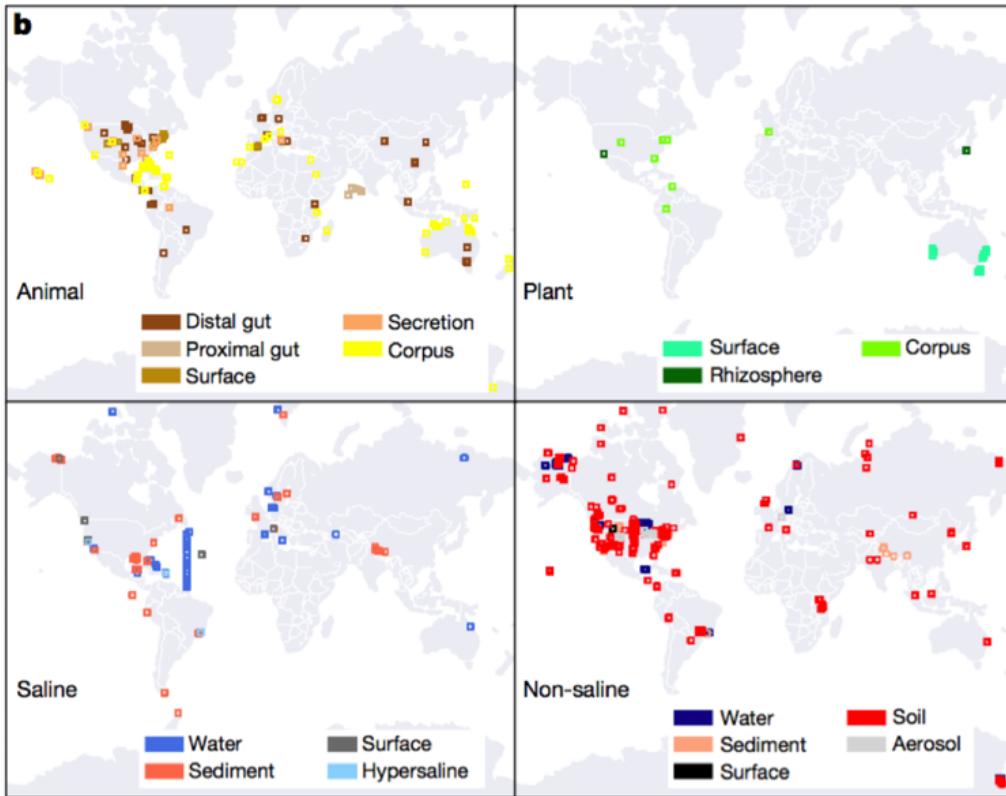


Metagenomic sequencing analysis (メタゲノム解析, ショットガンメタゲノム解析)



Who's there?

What are they doing?



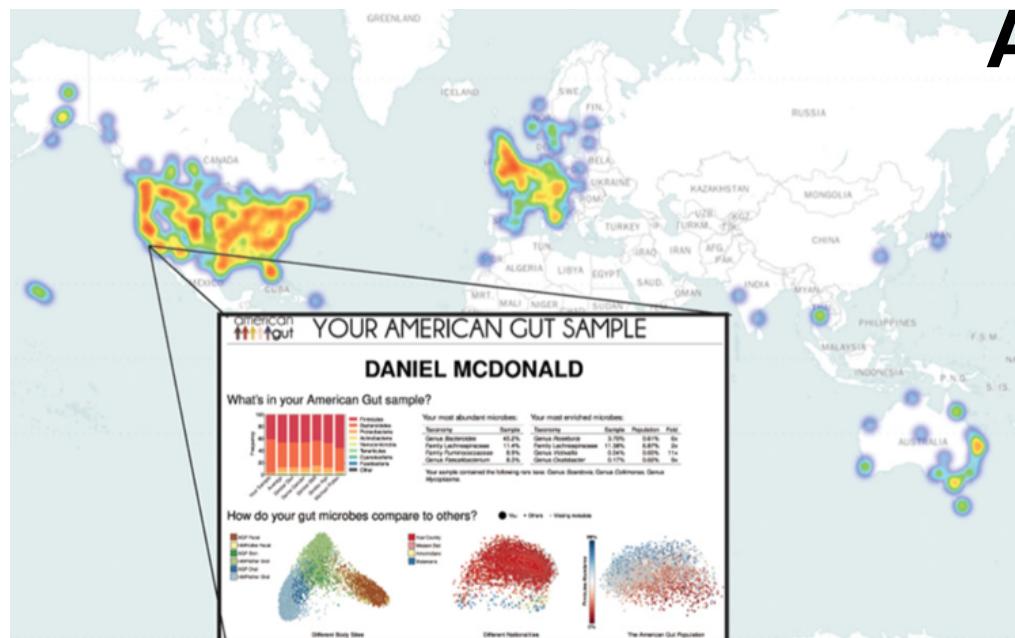
Earth Microbiome Project

Thompson et al. 2017, Nature

Soil, Sea, Pond, Animal, etc.

23,828 samples

Standardized experimental and bioinformatics procedure



American Gut Project

McDonald et al. 2018, mSystems

Human feces

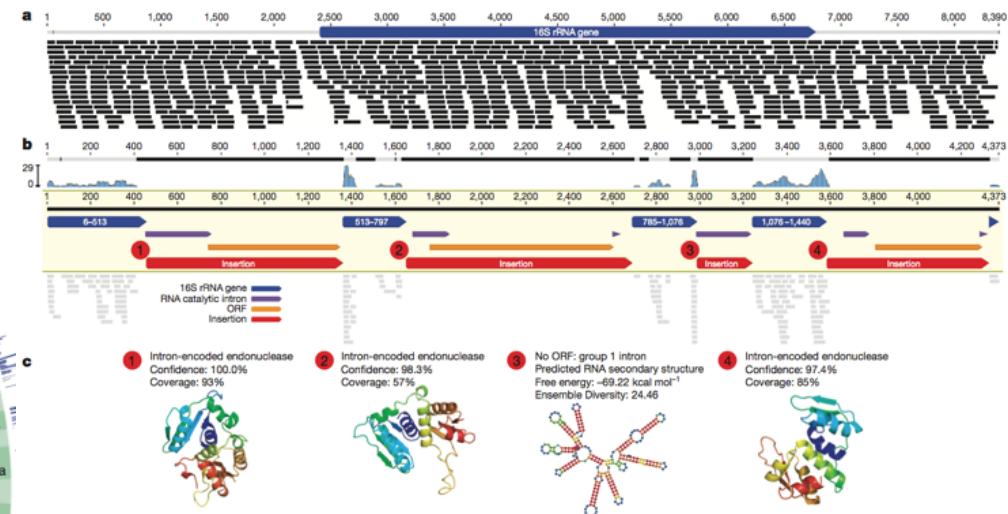
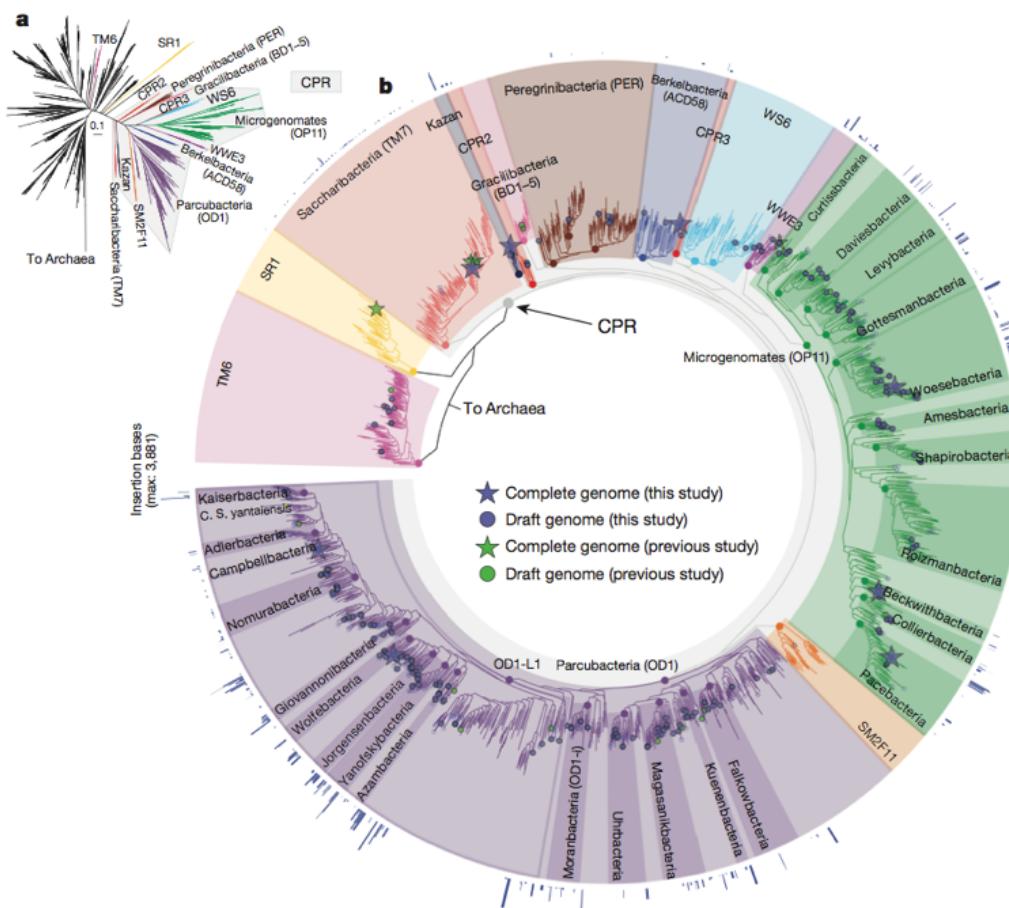
9,511 samples

Standardized experimental and bioinformatics procedure.

Cloud funding

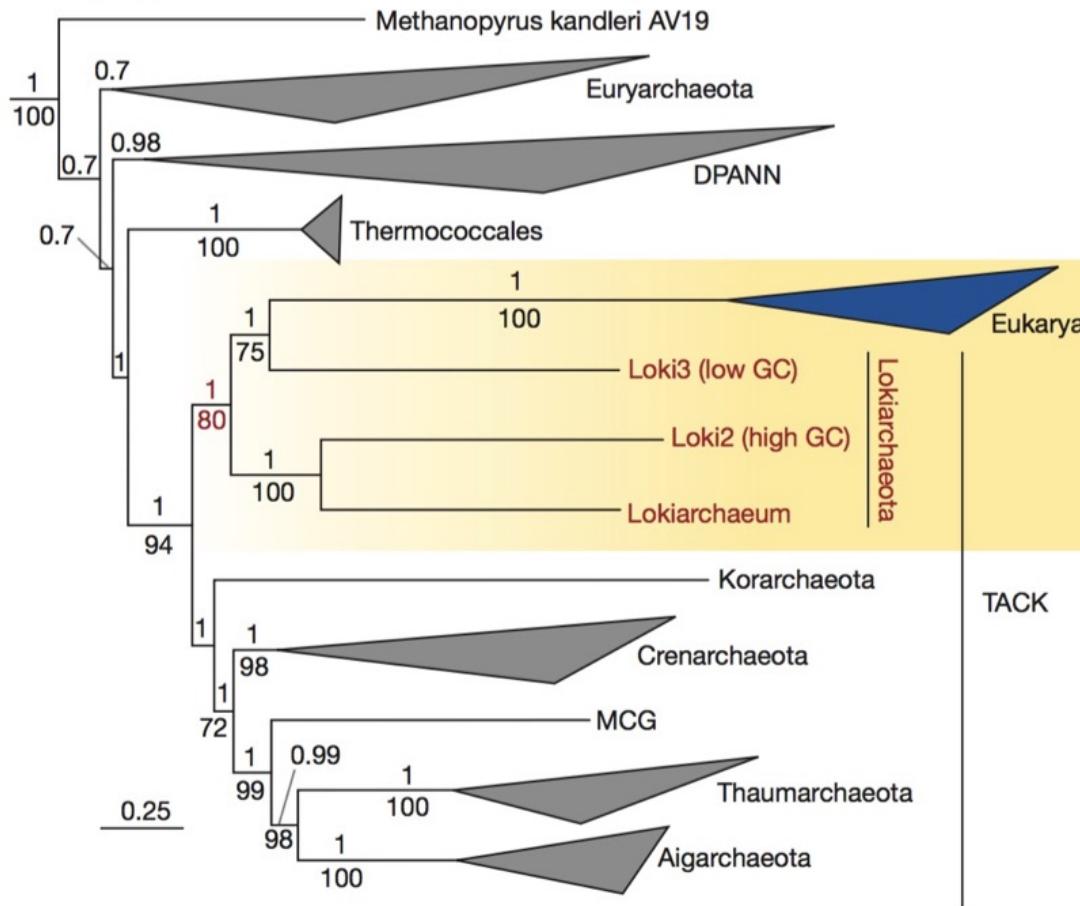
Unusual biology across a group comprising more than 15% of domain Bacteria

Christopher T. Brown¹, Laura A. Hug², Brian C. Thomas², Itai Sharon², Cindy J. Castelle², Andrea Singh², Michael J. Wilkins^{3,4}, Kelly C. Wrighton⁴, Kenneth H. Williams⁵ & Jillian F. Banfield^{2,5,6}



Complex archaea that bridge the gap between prokaryotes and eukaryotes

Anja Spang^{1*}, Jimmy H. Saw^{1*}, Steffen L. Jørgensen^{2*}, Katarzyna Zaremba-Niedzwiedzka^{1*}, Joran Martijn¹, Anders E. Lind¹, Roel van Eijk^{1†}, Christa Schleper^{2,3}, Lionel Guy^{1,4} & Thijs J. G. Ettema¹



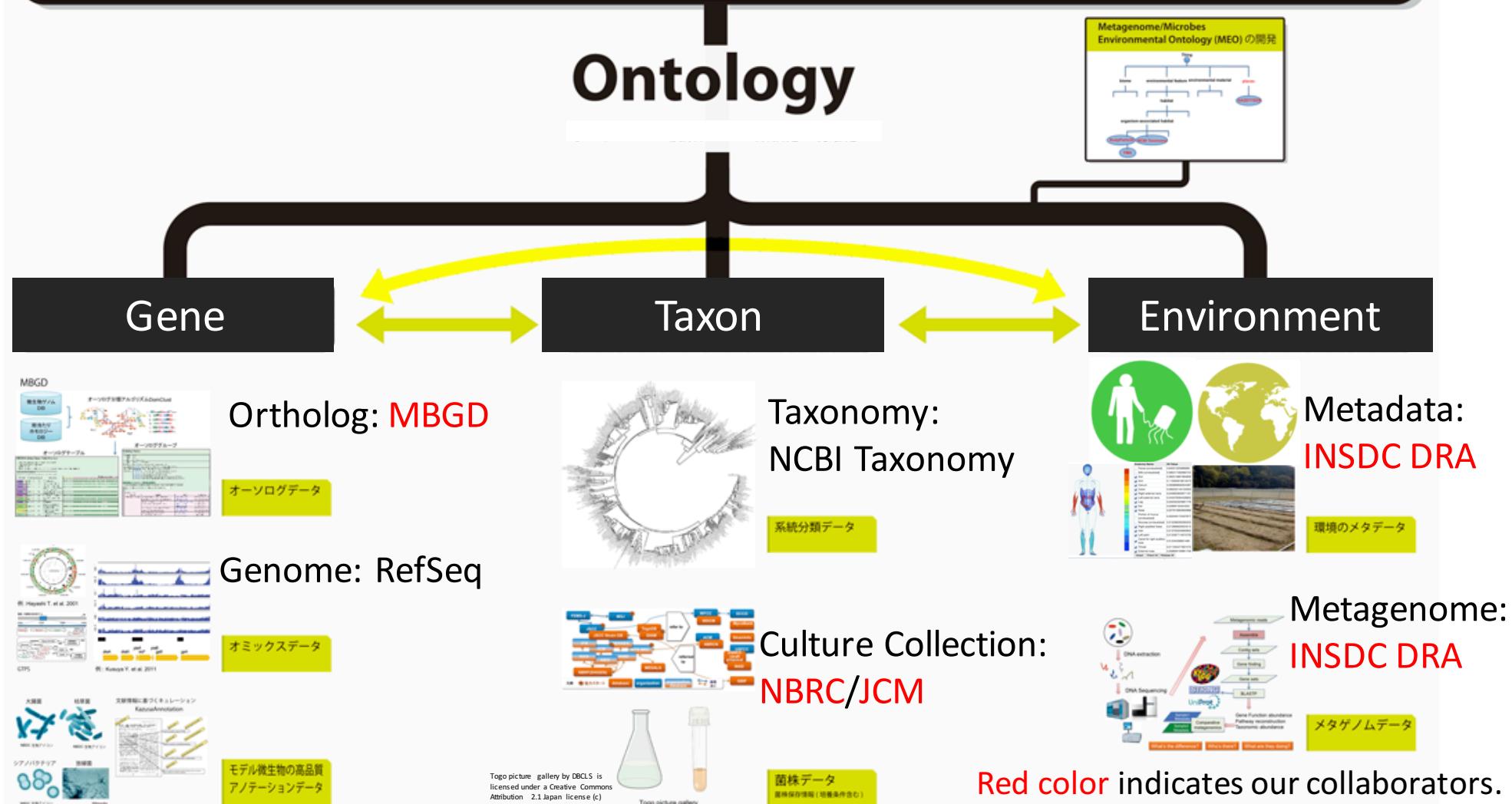
メタゲノム解析の情報解析手法については、
<https://github.com/AJACS-training/AJACS79>

Microbiome data in public DB

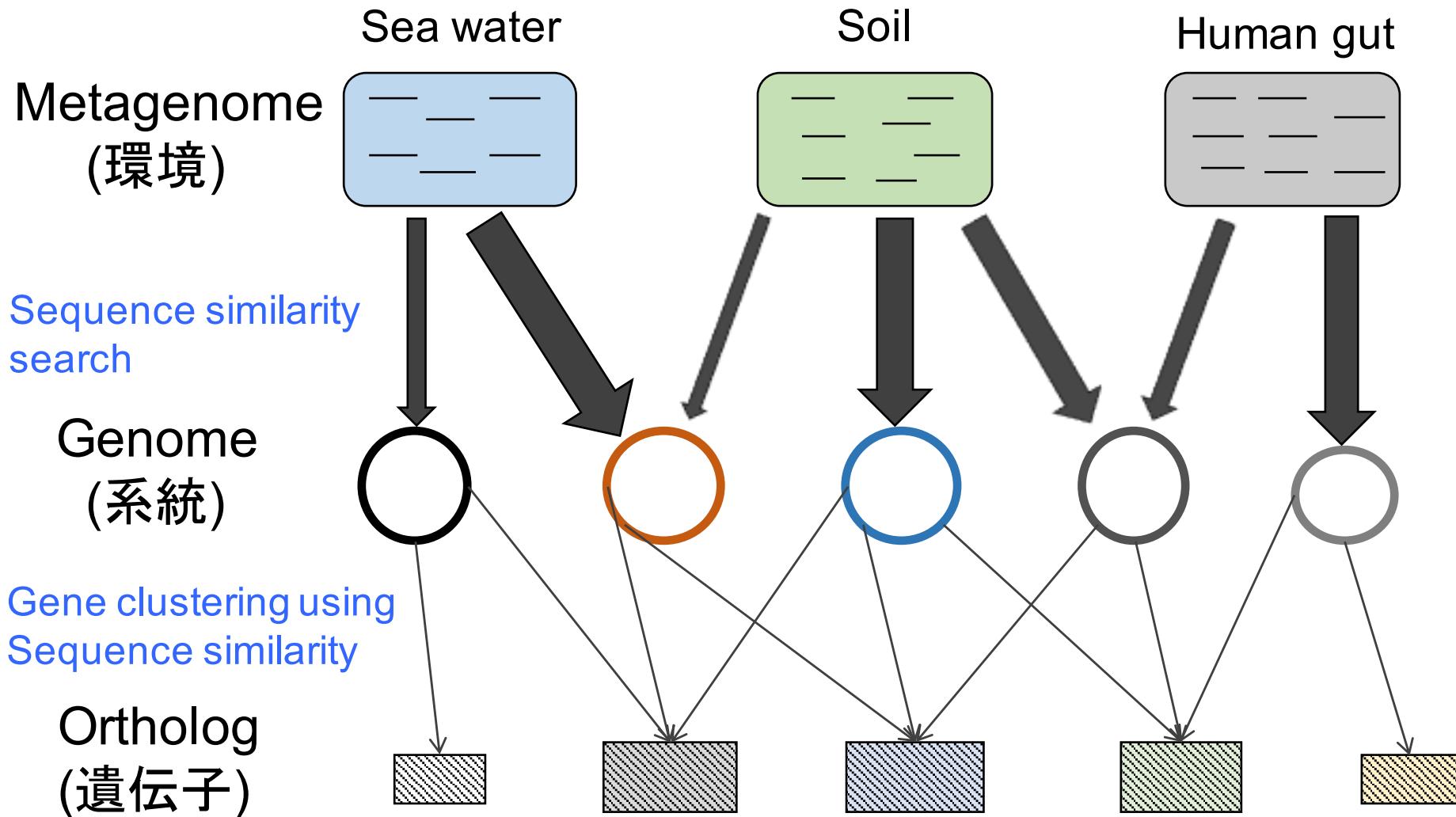
	Admin	Database URL	Sequence data	Separate amplicon and shotgun?	Taxa	Function	Number of samples in Sep. 2019
INSDC DRA/ERA/SRA	Japan, EU, USA		○	×	△	×	>1,600,000
GOLD	JGI, USA	https://gold.jgi.doe.gov/	×	×	×	×	50,821
IMG/M	JGI, USA	https://img.jgi.doe.gov/cgi-bin/m/main.cgi	○	×	○	○	16,170
MG-RAST	Chicago U. USA	http://www.mg-rast.org/	○	○	○	○	392,514
MGnify	EBI, EU	https://www.ebi.ac.uk/metagenomics/	×	○	○	○	176,360
MicrobeDB.jp v.2	NIG, Japan	http://microbedb.jp/MDB/	×	○	○	○	60,551

1次データベースであるINSDC DRA/ERA/SRAには
系統組成・遺伝子機能組成のデータが無い。
メタデータもまともに構造化されておらず整理されていない

 **MicrobeDB**.JP integrates lots of data related to microbes.
Especially, we integrates the microbial data that can be linked to **genomes**. since 2011



配列類似性 (進化的類縁性) を基にしたデータの統合化



概念 (Ontology) を基にしたデータの統合化

Ontology is a structured controlled vocabulary to describe properties and types of resources.

例: 森とは何か? 林と何が違うのか?

データの由来が異なれば、同じ語彙でも意味が異なる場合がある
それをOntologyを用いて統一する

MEO (Microbes/Metagenomes Environmental Ontology)

MSV (Metagenome Sample Vocabulary)

MCCV (Microbial Culture Collection Vocabulary)

MPO (Microbial Phenotype Ontology)

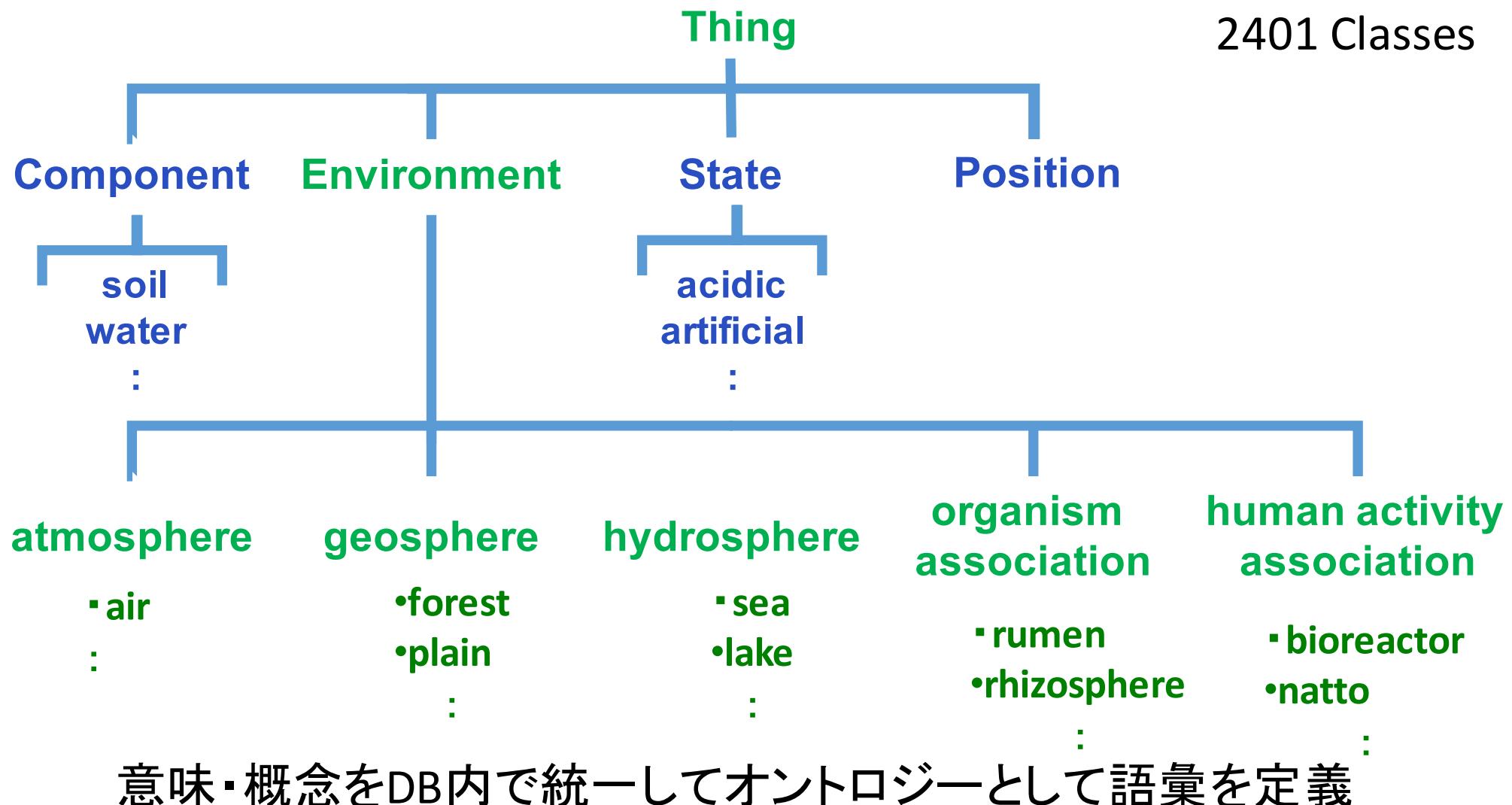
MBGD Ontology

PDO (Pathogenic Disease Ontology)

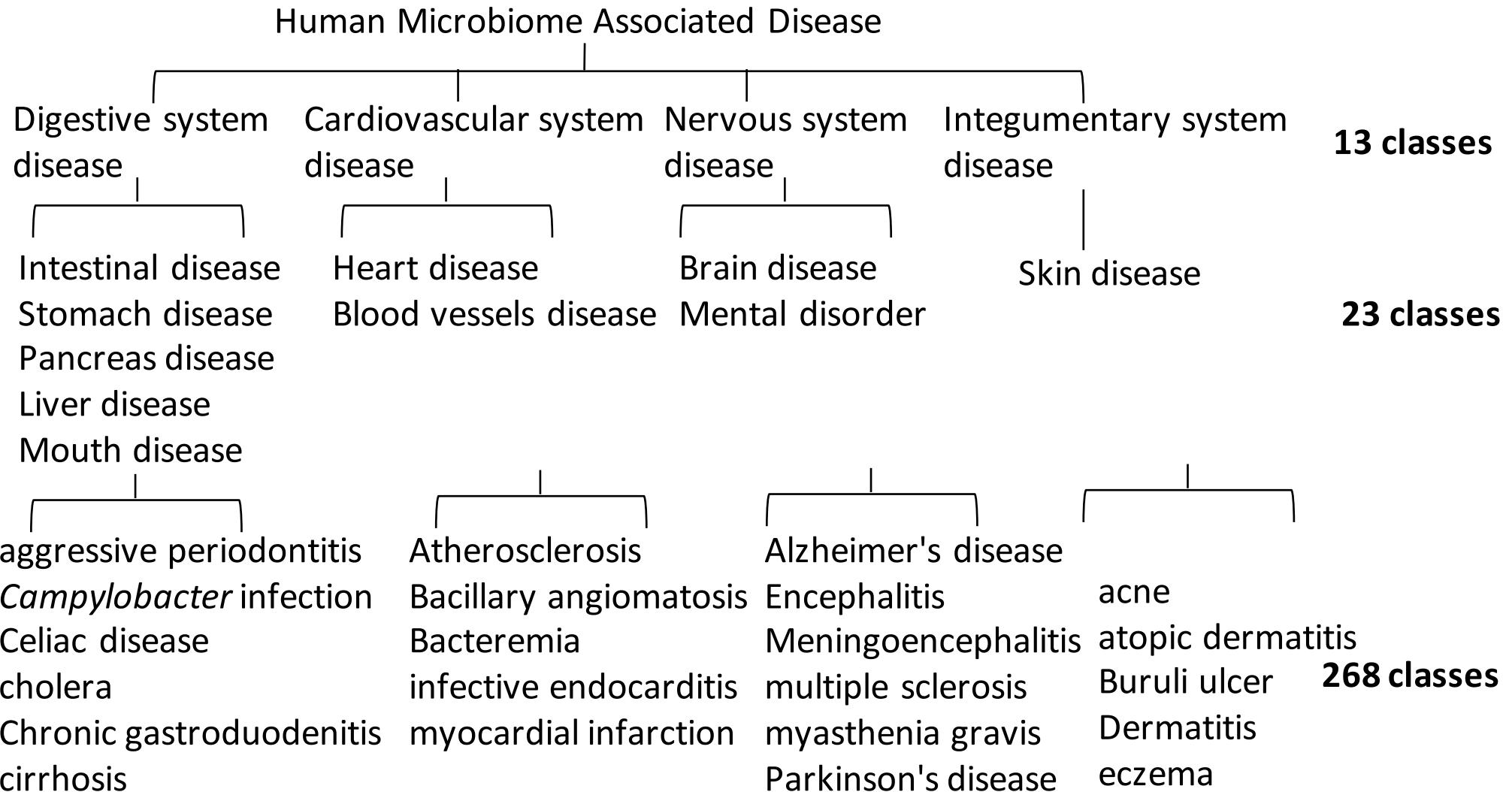
Most of them can be obtained from



Metagenome/Microbes Environmental Ontology (MEO) Ver. 0.9.2



Human Microbiome Associated Disease Ontology (HMADO)



感染症なのか否かなど、微生物群集が関係するヒトの病気の分類

MicrobeDB.jp version 3 data

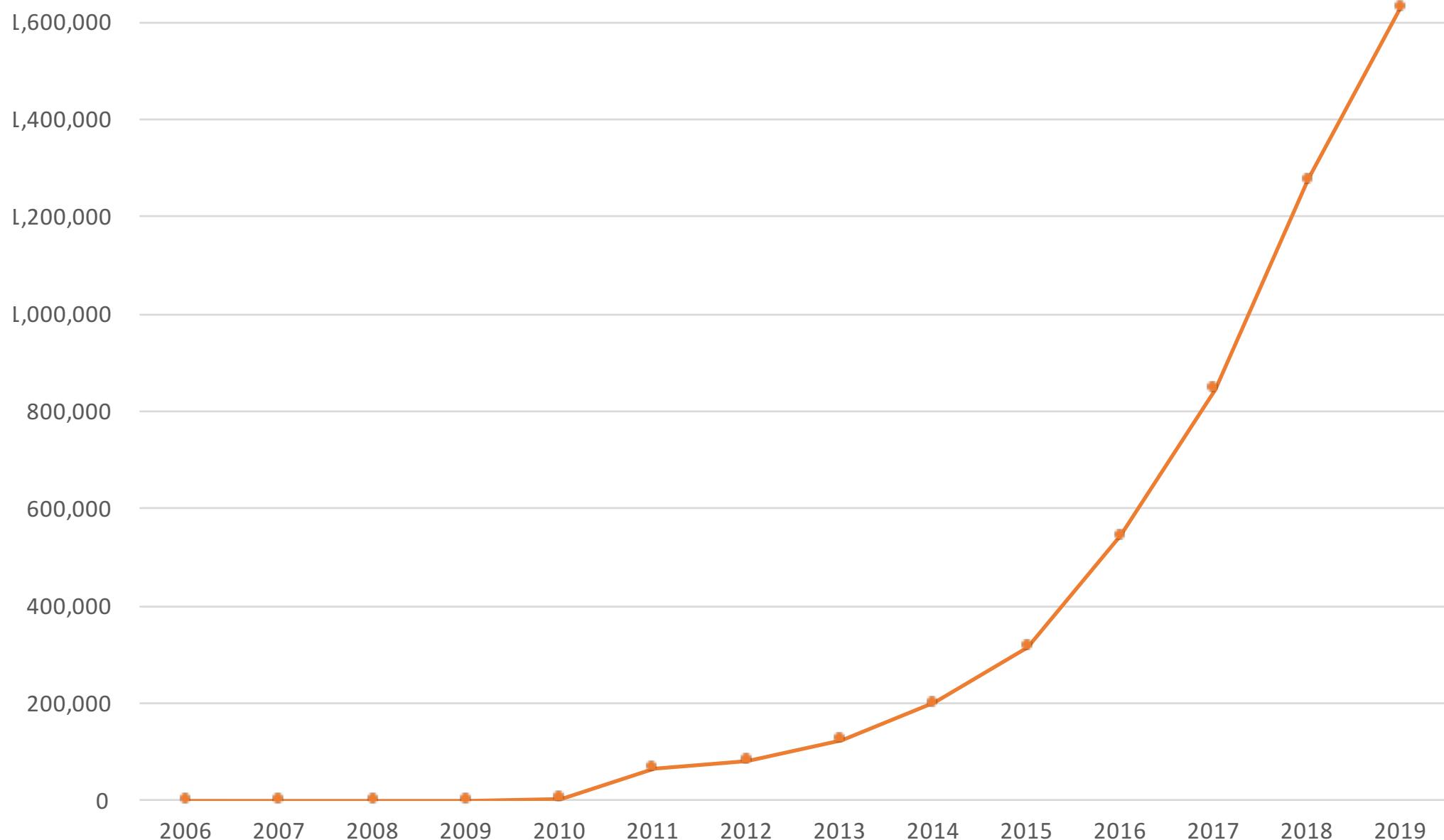
Data category	Data sources	Ontologies
Genome Metadata	INSDC BioSample, NCBI Assembly report	MEO, HMADO, NCBITAX
Ortholog	MBGD	ORTH
Culture collection	JCM (RDF-Portal)	MCCV, MPO, MEO, PDO, CSSO
Culture collection	NBRC (RDF-Portal)	MCCV, MPO
Metagenome	INSDC DRA	NCBITAX, KEGG Orthology
Metagenome metadata	INSDC BioSample	MEO, HMADO, SIO

MicrobeDB.jp version 3 data

Data category	Number of entry
Genome metadata (from RefSeq)	290,208 genomes
Ortholog cluster data (from MBGD)	375,228 clusters
Culture collection strain data from JCM/NBRC (RDF-Portal)	38,414 strains
Microbiome metadata (from INSDC DRA)	1,631,611 samples
Microbiome taxonomic composition data	96,766 samples
Microbiome functional composition data	4,784 samples

Ver. 2(2014)と比べてゲノムは約20倍、メタゲノムは約10倍の数に

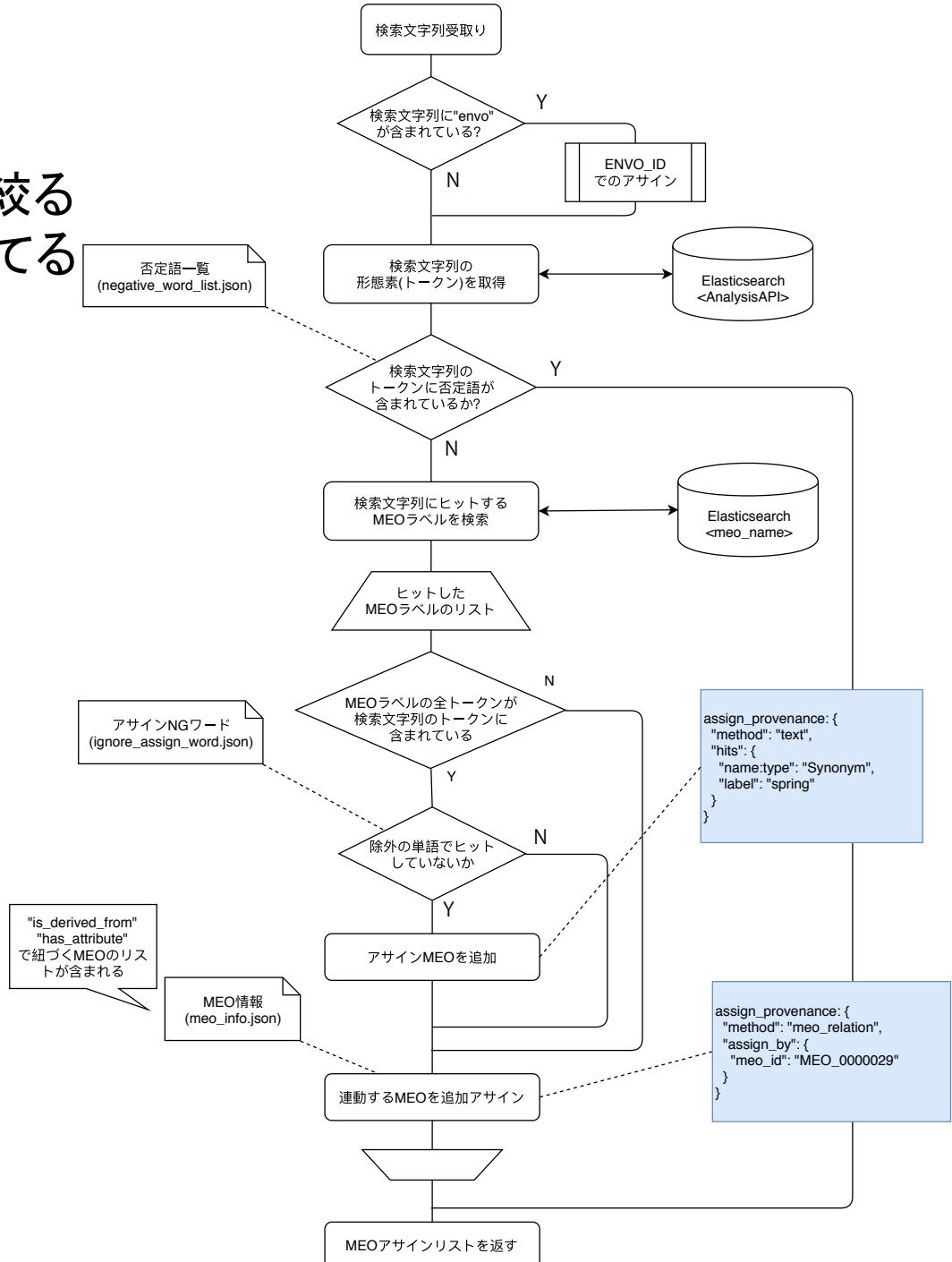
INSDC DRA/ERA/SRAで公開されたマイクロバイオームサンプル数 (積算・2019年9月時点)



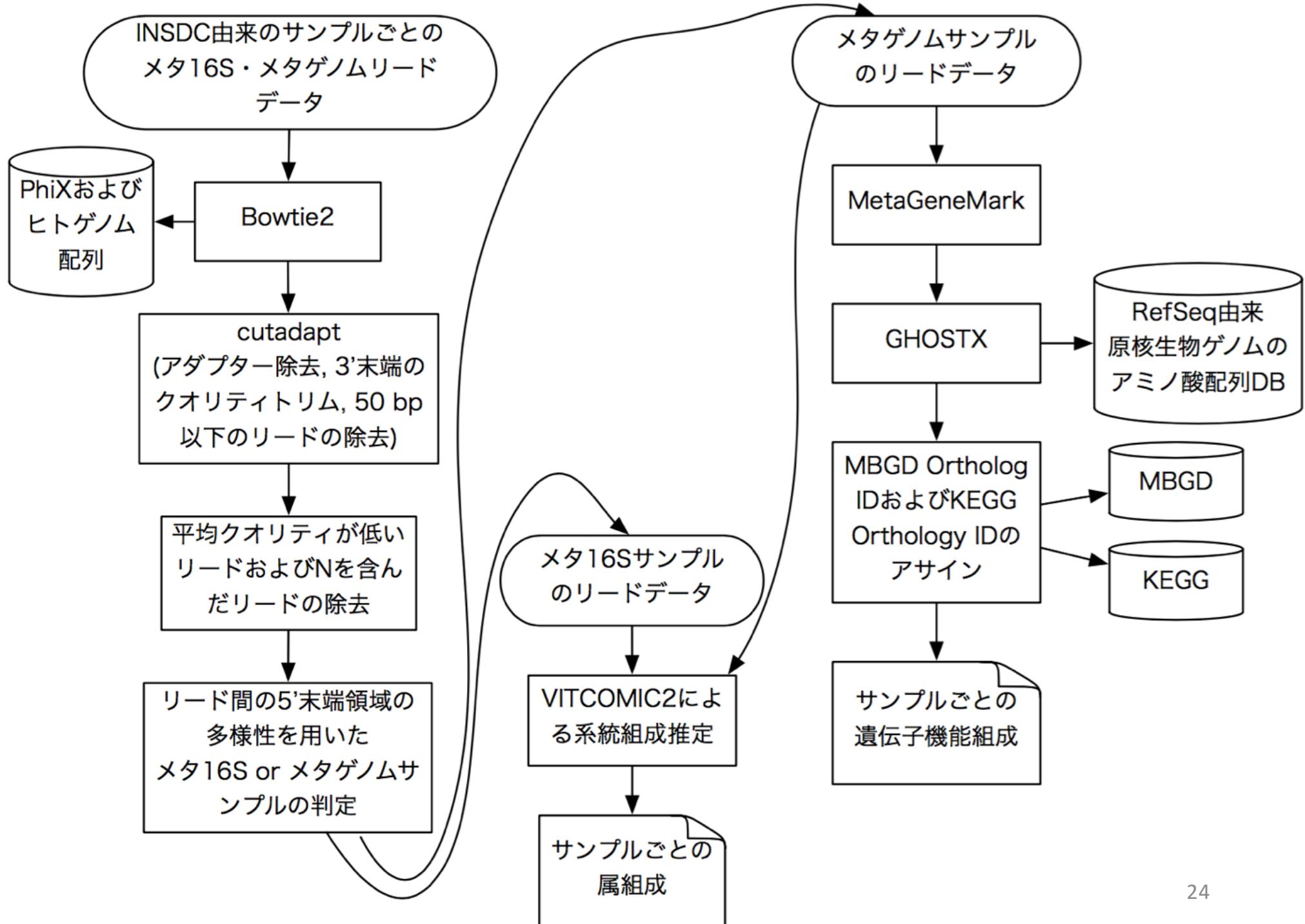
非構造TextデータからのMEOの自動アサイン

- 数的にマイナーなものは諦める
- attribute name(harmonized name)を絞る
- オントロジーのラベル、synonymをあてる
- 否定語対策
- stemmingの有無
- 例外ワードリストの作成
- より下位概念のアサイン

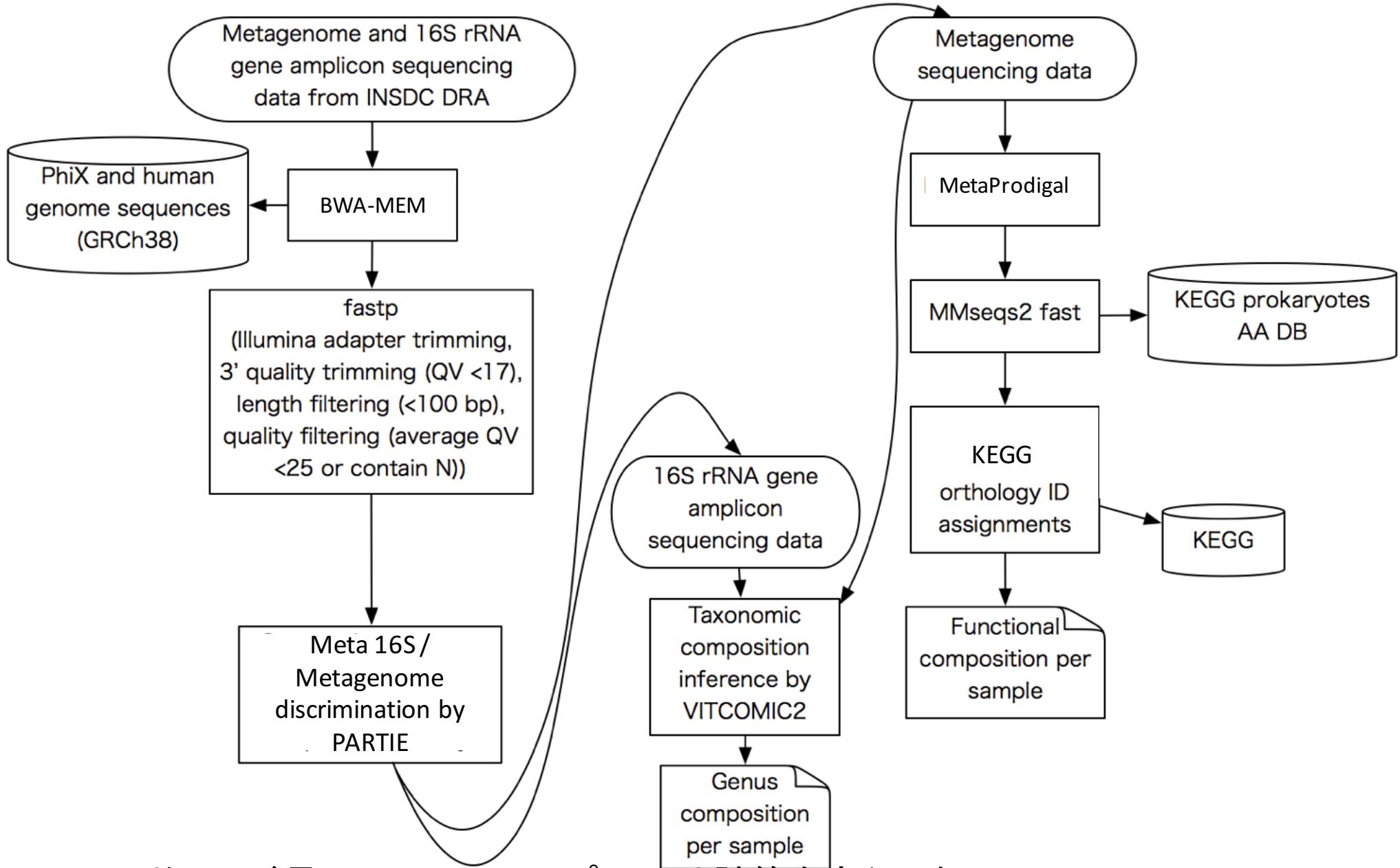
6割のゲノム、
9割のメタゲノムのBioSampleに
MEOがアサインされた



MeGAP2 (MicrobeDB.jp ver. 2のパイプライン)



MeGAP3 (MicrobeDB.jp ver. 3のパイプライン)



平均リード長 <100 bp のサンプルでは計算を実行しない。

16S rRNAアンプリコンについては、結果のファイルサイズでフィルタリング、
メタゲノムについては、アサインされたKOの種類数 <100 のサンプルを除外

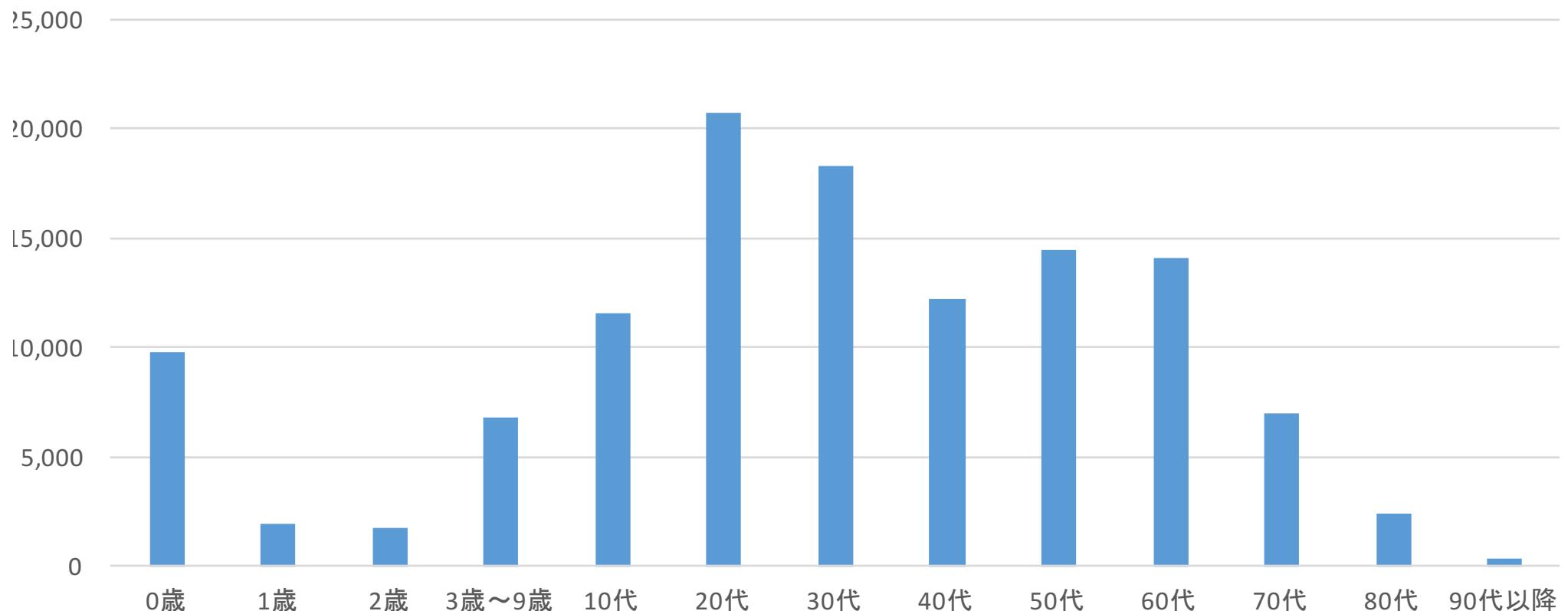
INSDC DRA/ERA/SRAで公開されたマイクロバイオームサンプル数 (2019年9月時点)

	サンプル数
マイクロバイオーム	約 1,600,000
ヒト	約 420,000
マウス	約 95,000
土壤	約 240,000
水環境	約 200,000
人工環境	約 140,000
植物共生	約 120,000

160万サンプルはヒトや土壤等マイクロバイオームサンプルの総数

ヒトマイクロバイオームについて、特に詳細にメタデータをキュレーション

	サンプル数
ヒト	約 420,000
ヒト (年齢情報あり)	約 120,000



ヒトマイクロバイオームについて、特に詳細にメタデータをキュレーション

ヒトの体の部位	サンプル数
gut (feces)	243,958
oral cavity	55,990
respiratory system	45,655
skin	41,381
vagina	21,761

- 抗生物質投与(約2-3万サンプル)
- Probiotics投与(約5千サンプル)
- Ethnicity(約5千サンプルを、Ethnicity ontologyとGAZETTEERでアノテーション)

INSDC DRA/ERA/SRAで公開された 疾患患者のマイクロバイオーム サンプル数(2019年9月時点)

サンプル数が100以上ある疾患は
約50疾患

サンプル数が1000以上ある疾患が
右表の18疾患
(肥満と未熟児を除く)

大多数が16S rRNA遺伝子の
アンプリコン解析

疾患名 (HMADO)	サンプルサイズ
IBD	37,094
皮膚炎	13,219
にきび	7,488
喘息	6,178
乾癬	5,151
大腸炎	4,453
囊胞性線維症	4,040
アレルギー	3,902
下痢	2,684
大腸がん	2,085
細菌性膿炎	2,030
アテローム性動脈硬化	1,968
赤痢	1,918
パーキンソン病	1,760
腺がん	1,742
統合失調症	1,712
歯周炎	1,555
虫歯	1,370

MicrobeDB.jpの使い方

詳しくは、

<https://togotv.dbcls.jp/20190902.html>



https://beta.microbedb.jp



Home

Document

Analysis ▾

e.g. hot spring, Enterococcus faecalis, psbA

Search



MicrobeDB.jp

Integrating and representing genome, metagenome, taxonomy resources and the analysis datasets with Semantic Web Technologies.

Database statistics

Total number of Metagenomic samples (SRA/SRS):	173,359 samples
- with taxonomic analysis results:	60,551 samples
- with functional analysis results:	4,048 samples
Total number of Assembled Genomes (RefSeq/Genbank):	16,983 taxa
Total number of Strains (JCM/NBRC):	16,671 strains
Total number of Environmental terms in ontology (MEO):	2,381 terms



Search id or term...

Public/Private

 metagenome_public 3729

hasMetagenomeAnalysis

 taxonomy 3729 function 609

hasMEO (Text)

 gut

hasMEO: Component

Component for environment 389

hasMEO: Env

Environment for microbes 3729

hasMEO: Position

Position toward environment 5

hasMEO: State

Metagenomic samples 3729 results found in 112mshasMetagenomeAnalysis: taxonomy hasMEO (Text): gut

Clear all filters

Previous

1

2

3

4

...

Next

10

Select All

Deselect All

Select	MDB SampleID	msv:sampleTitle	msv:scientificName	msv:hasTaxonID	msv:hasBioProjectID	msv:hasBioSar
<input type="button" value="Remove"/>	SRS551059	Content of the intestinum from animals fed one meal of heparinized sheep blood	gut metagenome	749906	PRJNA237098	SAMN026145
<input type="button" value="Remove"/>	SRS551061	Content of the intestinum of animals fed two meals of heparinized sheep blood four weeks apart	gut metagenome	749906	PRJNA237098	SAMN026145
<input type="button" value="Remove"/>	SRS452599	Environmental/Metagenome sample for mouse gut metagenome	mouse gut metagenome	410661	PRJNA209582	SAMN022138
<input type="button" value="Add"/>	SRS367344	Pooled 16S rRNA gene sequences	Bacteria	2	PRJNA209582	SAMN017587
<input type="button" value="Add"/>	SRS369251	Pooled 16S rRNA gene sequences		2	PRJNA209582	SAMN017656



Taxonomic composition (bar)

Taxonomic composition (heatmap)

Diversity index

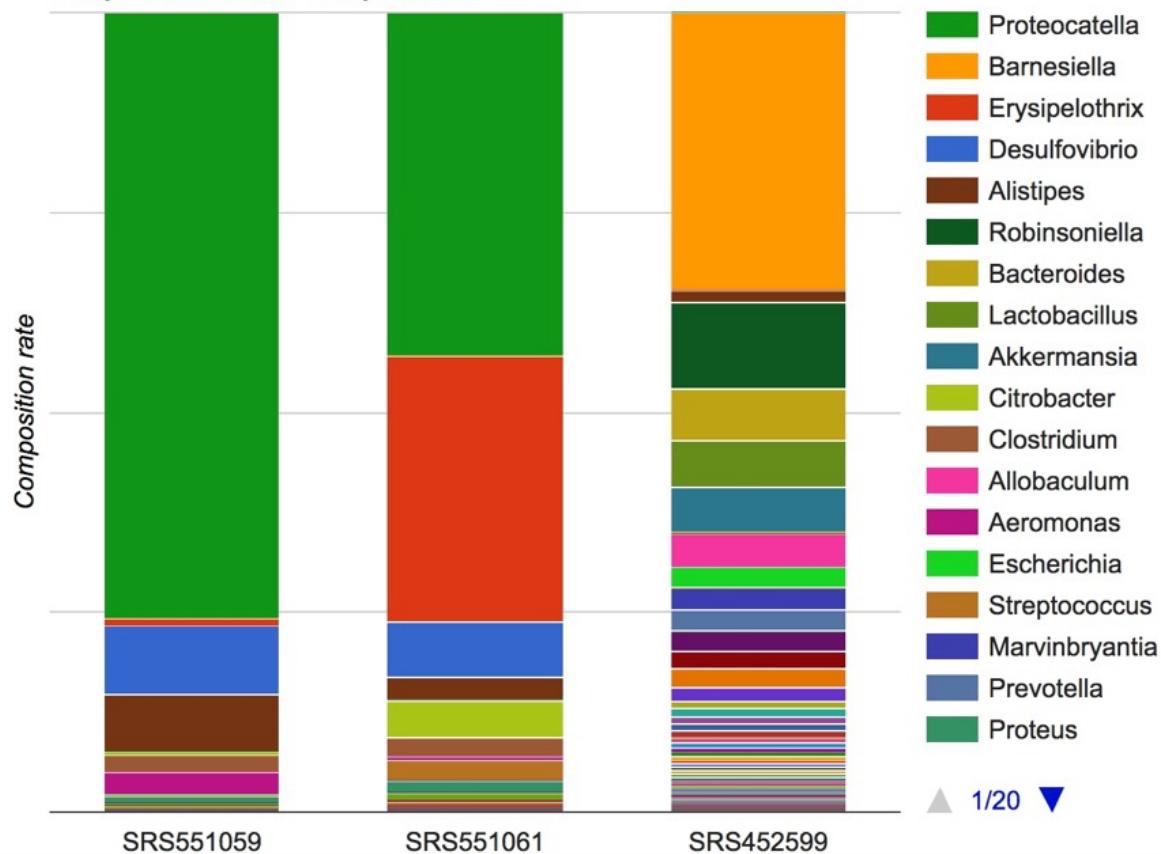
Hierarchical clustering

PCoA

Functional composition (bar)

Functional composition (heatmap)

Samples Taxonomic Composition



ここからversion 3の新機能
(12月中に切り替わります)

MEO(soil)由来マイクロバイオームサンプルの絞り込み検索(これはver. 2でもできました)

Home Document Analysis e.g. hot spring, Enterococcus faecalis, psb. Search

Index

facet_sample 247853

hasMetagenomeAnalysis

taxonomy 8500

function 246

Search id ...

attribute name

Search attribute name ...

attribute value

Search attribute value ...

hasMEO (Text)

soil

hasMEO: Component

Component for environment 237292

hasMEO: Env

Environment for microbes 171749

taxonomy (Text)

Search taxonomy terms ...

taxonomy

root 247853

hasHostTaxonomy (Text)

Search HostTaxonomy...

hasHostTaxonomy

root 30348

pH

Add SAMN02898402 MIMS Environmental/Metagenome sample from soil metagenome soil metagenome 410658 PRJNA431674 SRS666697 SRR6516229 SAMN02898402 2014-07-08T00:00:00Z

Add SAMN08408074 Metagenome or environmental sample from soil metagenome soil metagenome 410658 PRJNA431674 SRS2872044 SRR6516229 SAMN08408074 2018-01-27T00:00:00Z

Add SAMN08408073 Metagenome or environmental sample from soil metagenome soil metagenome 410658 PRJNA431673 SRS2872043 SRR6516228 SAMN08408073 2018-01-27T00:00:00Z

Add SAMN07281069 crop soil soil metagenome 410658 PRJNA392052 SRS2317230 SRR5760457 SAMN07281069 2017-06-27T00:00:00Z

Add SAMN07281068 crop soil soil metagenome 410658 PRJNA392052 SRS2317229 SRR5760458 SAMN07281068 2017-06-27T00:00:00Z

Add SAMN07281067 crop soil soil metagenome 410658 PRJNA392052 SRS2317207 SRR5760480 SAMN07281067 2017-06-27T00:00:00Z

Add SAMN07281066 crop soil soil metagenome 410658 PRJNA392052 SRS2317228 SRR5760459 SAMN07281066 2017-06-27T00:00:00Z

Add SAMN07281065 crop soil soil metagenome 410658 PRJNA392052 SRS2317227 SRR5760460 SAMN07281065 2017-06-27T00:00:00Z

Add SAMN07281064 crop soil soil metagenome 410658 PRJNA392052 SRS2317226 SRR5760461 SAMN07281064 2017-06-27T00:00:00Z

Add SAMN07281063 crop soil soil metagenome 410658 PRJNA392052 SRS2317205 SRR5760482 SAMN07281063 2017-06-27T00:00:00Z

Add SAMN07281062 crop soil soil metagenome 410658 PRJNA392052 SRS2317206 SRR5760481 SAMN07281062 2017-06-27T00:00:00Z

Add SAMN07281061 crop soil soil metagenome 410658 PRJNA392052 SRS2317214 SRR5760473 SAMN07281061 2017-06-27T00:00:00Z

Add SAMN07281060 crop soil soil metagenome 410658 PRJNA392052 SRS2317215 SRR5760472 SAMN07281060 2017-06-27T00:00:00Z

Add SAMN07281059 crop soil soil metagenome 410658 PRJNA392052 SRS2317212 SRR5760475 SAMN07281059 2017-06-27T00:00:00Z

Add SAMN07281058 crop soil soil metagenome 410658 PRJNA392052 SRS2317213 SRR5760474 SAMN07281058 2017-06-27T00:00:00Z

Add SAMN07281057 crop soil soil metagenome 410658 PRJNA392052 SRS2317210 SRR5760477 SAMN07281057 2017-06-27T00:00:00Z

Add SAMN07281056 crop soil soil metagenome 410658 PRJNA392052 SRS2317211 SRR5760476 SAMN07281056 2017-06-27T00:00:00Z

Add SAMN07281055 crop soil soil metagenome 410658 PRJNA392052 SRS2317208 SRR5760479 SAMN07281055 2017-06-27T00:00:00Z

HMADO(cancer)由来マイクロバイオームサンプルの絞り込み検索

MicrobeDB.jp

Home Document Analysis e.g. hot spring, Enterococcus faecalis, psb. Search Sign Up Sign In

Index
facet_sample 3587

hasMetagenomeAnalysis
taxonomy 41

Search id ...

attribute name
Search attribute name ...

attribute value
Search attribute value ...

hasMEO (Text)
Search MEO terms ...

hasMEO: Env
Environment for microbes 3587

taxonomy (Text)
Search taxonomy terms ...

taxonomy
root 3587

hasHostTaxonomy (Text)
Search HostTaxonomy...

HMADO (Text)
cancer

HMADO (Text): cancer x Clear all filters

Metagenomic samples 3587 results found in 43ms

Previous 1 2 3 4 ... Next

10 Select All Deselect All

Select	MDB SampleID	title	organism.name	organism.identifier	BioProjectID	SRAID	SRRID	BioSampleID	publishedDate
Add	SAMN07452485	MIMS Environmental/Metagenome sample from human gut metagenome	human gut metagenome	408170	PRJNA397219	SRS2409846	SRR5903357 SRR5903369 SRR5903386 SRR5903748	SAMN07452485	2017-08-06T00:00:00.00
Add	SAMN07452484	MIMS Environmental/Metagenome sample from human gut metagenome	human gut metagenome	408170	PRJNA397219	SRS2409841	SRR5903342 SRR5903368 SRR5903385 SRR5903757	SAMN07452484	2017-08-06T00:00:00.00
Add	SAMN07452483	MIMS Environmental/Metagenome sample from human gut metagenome	human gut metagenome	408170	PRJNA397219	SRS2409840	SRR5903336 SRR5903367 SRR5903384 SRR5903756	SAMN07452483	2017-08-06T00:00:00.00
Add	SAMN07452482	MIMS Environmental/Metagenome sample from human gut metagenome	human gut metagenome	408170	PRJNA397219	SRS2409834	SRR5903329 SRR5903366 SRR5903383 SRR5903761	SAMN07452482	2017-08-06T00:00:00.00
Add	SAMN07452481	MIMS Environmental/Metagenome sample from human gut metagenome	human gut metagenome	408170	PRJNA397219	SRS2409835	SRR5903330 SRR5903365 SRR5903382 SRR5903760	SAMN07452481	2017-08-06T00:00:00.00
Add	SAMN07452480	MIMS Environmental/Metagenome sample from human gut metagenome	human gut metagenome	408170	PRJNA397219	SRS2409832	SRR5903327 SRR5903345 SRR5903381 SRR5903759	SAMN07452480	2017-08-06T00:00:00.00
Add	SAMN07452479	MIMS Environmental/Metagenome sample from human gut metagenome	human gut metagenome	408170	PRJNA397219	SRS2409833	SRR5903328 SRR5903346 SRR5903380 SRR5903758	SAMN07452479	2017-08-06T00:00:00.00
Add	SAMN07452478	MIMS Environmental/Metagenome sample from human gut metagenome	human gut metagenome	408170	PRJNA397219	SRS2409830	SRR5903325 SRR5903375 SRR5903379 SRR5903765	SAMN07452478	2017-08-06T00:00:00.00
Add	SAMN07452477	MIMS Environmental/Metagenome sample from human gut metagenome	human gut metagenome	408170	PRJNA397219	SRS2409831	SRR5903326 SRR5903337 SRR5903378 SRR5903764	SAMN07452477	2017-08-06T00:00:00.00
Add	SAMN07452476	MIMS Environmental/Metagenome sample from human gut metagenome	human gut metagenome	408170	PRJNA397219	SRS2409828	SRR5903323 SRR5903338 SRR5903377 SRR5903763	SAMN07452476	2017-08-06T00:00:00.00

Metagenome sample comparison analysis x Compare 0 Please select items within range 2 - 100

Japanese由来マイクロバイオームサンプルの絞り込み検索

Home Document Analysis e.g. hot spring, Enterococcus faecalis, psb. Search Sig

Index
facet_sample 1539

hasMetagenomeAnalysis
taxonomy 1385

Search id ...

attribute name
Search attribute name ...

attribute value
Search attribute value ...

hasMEO (Text)
Search MEO terms ...

hasMEO: Component
Component for environment 4

hasMEO: Env
Environment for microbes 1539

taxonomy (Text)
Search taxonomy terms ...

taxonomy
root 1539

hasHostTaxonomy (Text)
Search HostTaxonomy...

hasHostTaxonomy
root 1539

HMADO (Text)

Metagenomic samples

1539 results found in 147ms

HostEthnicity: japan | [x](#) | [Clear all filters](#)

Previous **1** 2 3 4 ... Next

10 | [Select All](#) | [Deselect All](#)

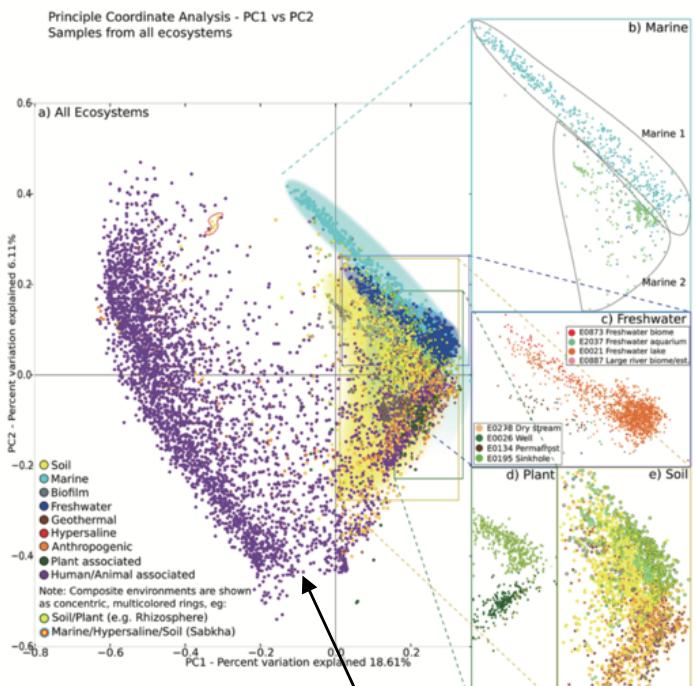
Select	MDB SampleID	title	organism.name	organism.identifier	BioProjectID	SRAID	SRRID	BioSampleID	publishedDate
Add	SAMD00079966	Japanese469	human gut metagenome	408170	PRJDB4360	DRS057897	DRR092102	SAMD00079966	2018-01-19T00
Add	SAMD00024579	Yms34	human gut metagenome	408170	PRJDB3417	DRS020590	DRR028771	SAMD00024579	2015-08-24T00
Add	SAMD00058608	Microbiota of the fecal sample from subject 516	human gut metagenome	408170	PRJDB4998		DRR068403	SAMD00058608	2016-11-17T22
Add	SAMD00035679	NA40	human gut metagenome	408170	PRJDB4064		DRR041845	SAMD00035679	2016-06-30T22
Add	SAMD00043213	Japanese371	human gut metagenome	408170	PRJDB4360		DRR049363	SAMD00043213	2016-05-30T13
Add	SAMD00036348	TS-41	human gut metagenome	408170	PRJDB3601		DRR042663	SAMD00036348	2016-03-08T11
Add	SAMD00079965	Japanese468	human gut metagenome	408170	PRJDB4360	DRS057896	DRR092101	SAMD00079965	2018-01-19T00
Add	SAMD00079964	Japanese467	human gut metagenome	408170	PRJDB4360	DRS057895	DRR092100	SAMD00079964	2018-01-19T00
Add	SAMD00079963	Japanese466	human gut metagenome	408170	PRJDB4360	DRS057894	DRR092099	SAMD00079963	2018-01-19T00
Add	SAMD00079962	Japanese465	human gut metagenome	408170	PRJDB4360	DRS057893	DRR092098	SAMD00079962	2018-01-19T00

Metagenome sample comparison analysis | [Compare 0](#) | Please select items within range 2 - 100

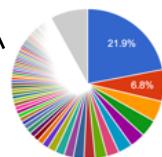
もう一つの検索アプローチ

配列データで判断する

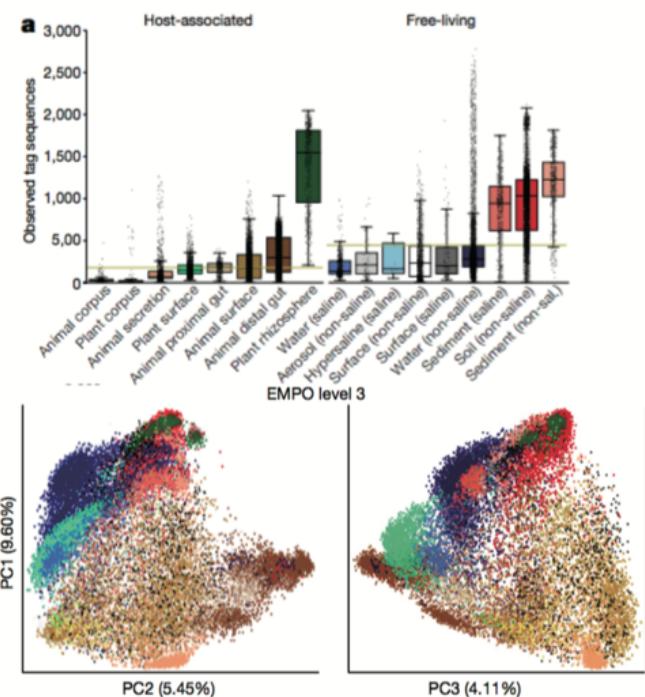
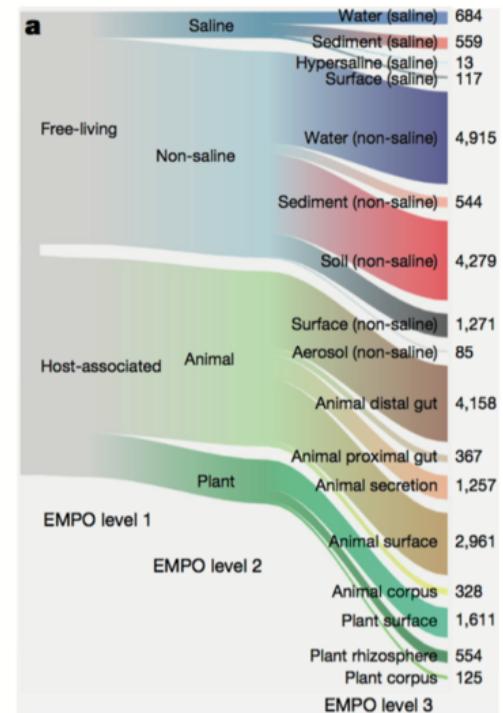
Global diversity of microbial community structures



Henschel, A., Anwar, M. Z., & Manohar, V. (2015).
PLoS computational biology, 11(10), e1004468.



Community structure
(Bacterial & Archaeal taxonomic abundance)



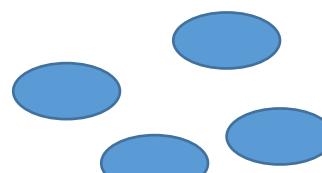
Thompson, Luke R., et al. Nature 551.7681 (2017).

Design “environmental labels” to be compared
 => Assign environmental labels to samples
 => Evaluate similarities or differences between labeled groups
 But, can “environment” be discretely categorized?

All samples in DB has metadata (documents described by researchers).
Using these documents and community structure data analyzed by MicrobeDB.jp,
Correspondence between microbes and natural languages can be learned by a
hierarchical Bayesian model.

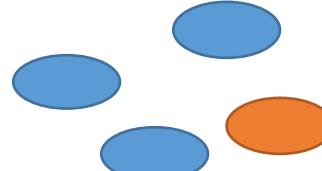
Sample #1

river microbiome.
thames river.
water samples collected from ...



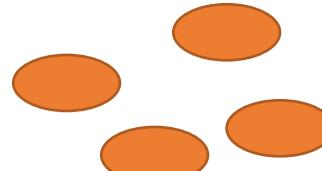
Sample #2

river microbiome.
wastewater contaminated river.



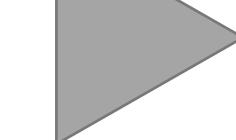
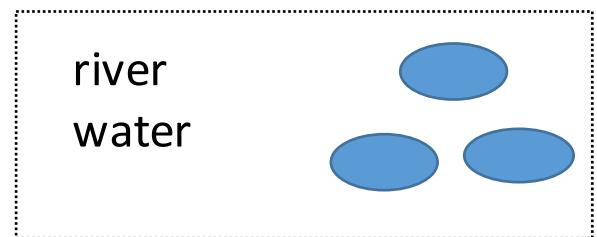
Sample #3

Activated sludge.
wastewater.
collected from the sewage
treatment facility

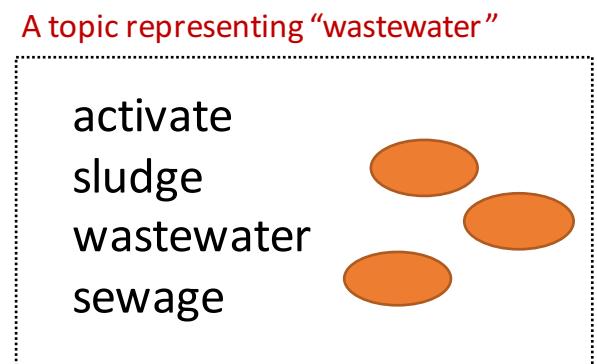


A combination of words and microbes
can be extracted as a base variable of
whole data set

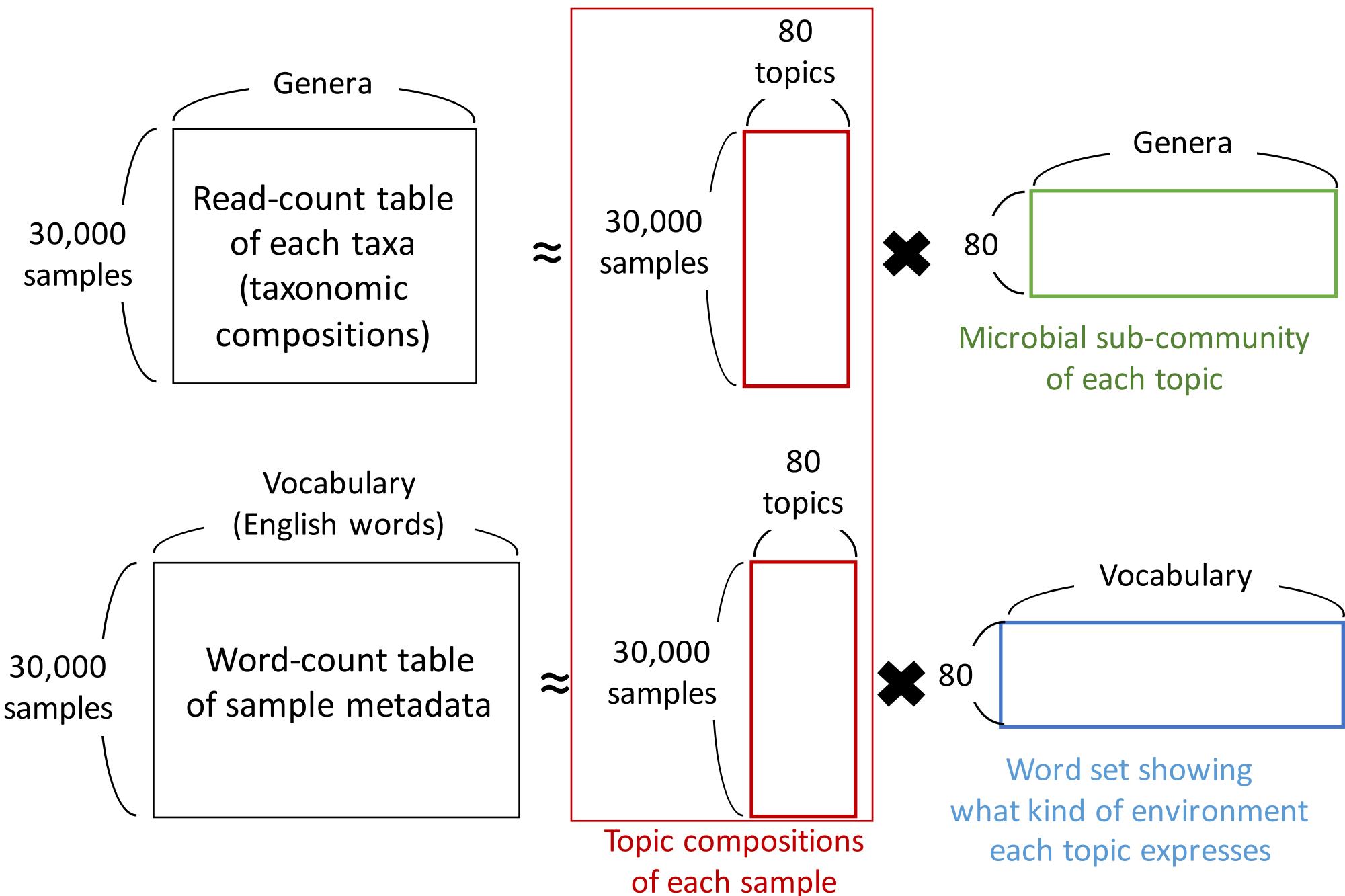
A topic representing “river”



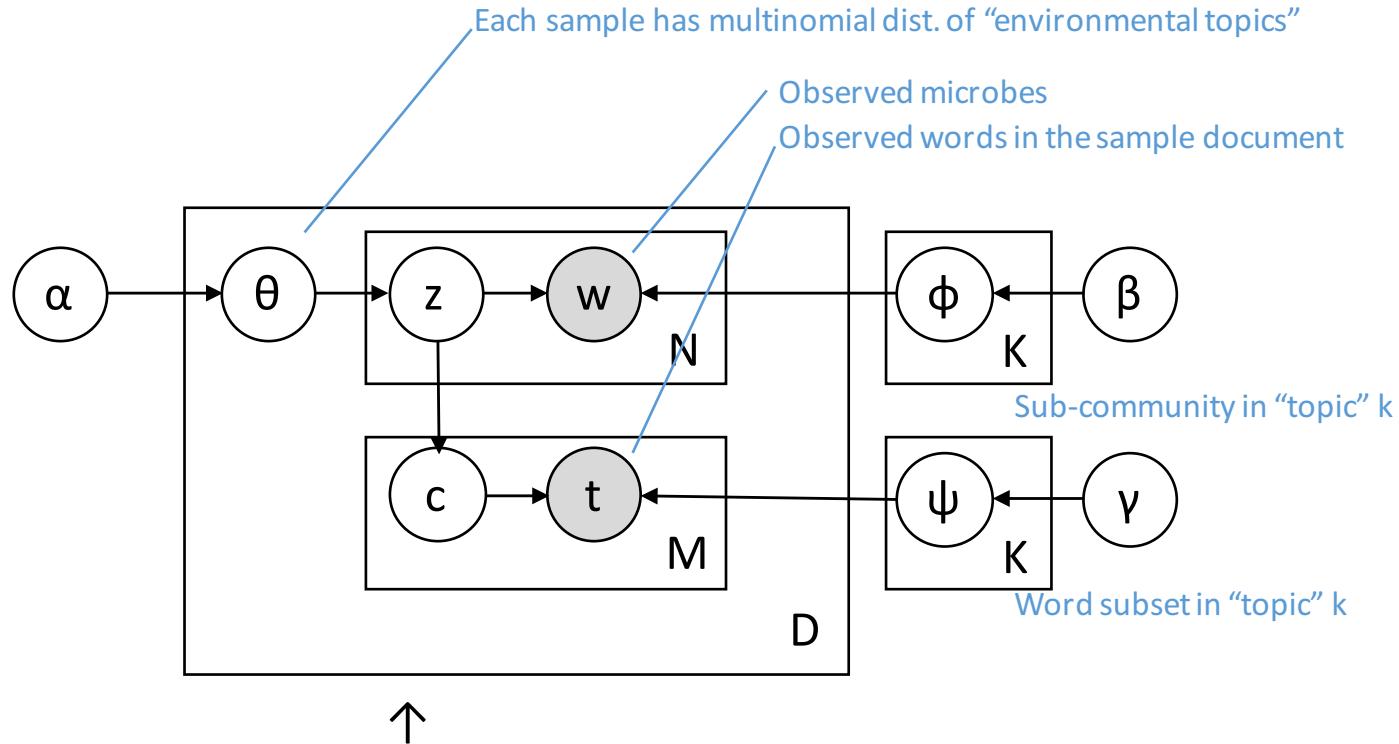
Machine
learning
(Topic model)



Matrix factorization for {read | word}-count tables (MicrobeDB.jp ver. 2 data)



Generative model of Latent environment allocation (LEA)
 based on the **topic model** approach (Latent Dirichlet Allocation-model)
 (cf. Admixture model of population genetics)



$$P(W, T, Z, C | \alpha, \beta, \gamma) = P(Z|\alpha)P(W|Z, \beta)P(C|Z)P(T|C, \gamma)$$

Parameter optimization by Gibbs sampling

$$P(z_{dn} = k | W, T, Z_{\setminus dn}, C) \propto (N_{kd \setminus dn} + \alpha_k) \frac{N_{kw_{dn} \setminus dn} + \beta}{N_{k \setminus dn} + \beta W} \left(\frac{N_{kd \setminus dn} + 1}{N_{kd \setminus dn}} \right)^{M_{kd}}$$

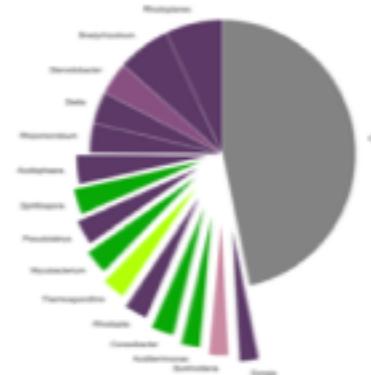
$$P(c_{dm} = k | W, T, C_{\setminus dm}, Z) \propto N_{kd} \frac{M_{kt_{dm} \setminus dm} + \gamma}{M_{k \setminus dm} + \gamma T}$$

Examples of estimated core “environmental topics”

Topic #17



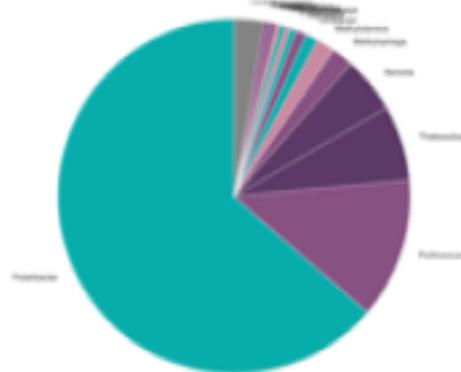
temperate
peat
fungus
ant
spruce
SOIL
forest
garden
mineral
tropical
finland
sweden
mountain
coniferous
wood
nest
active
natural



Topic #40



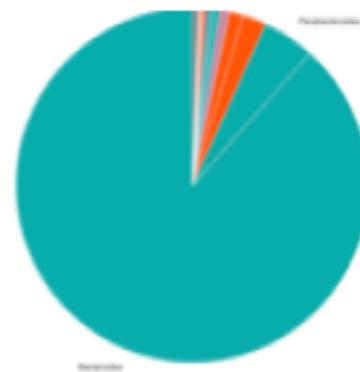
sea saline
southern ocean
water
glass marine
dark station
station light



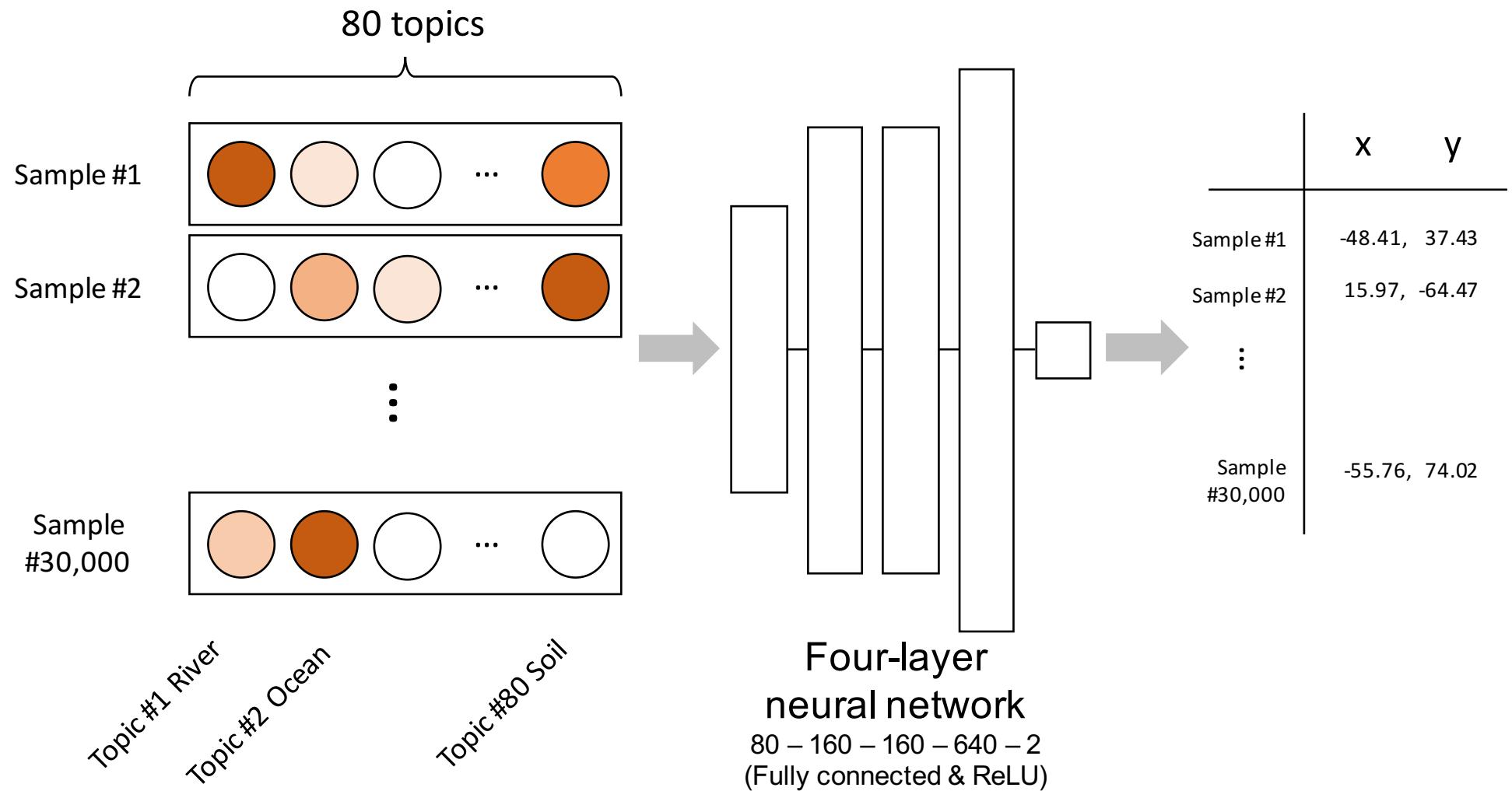
Topic #79



female patient adult
color new
stool child disease
feces male
york america intestinal

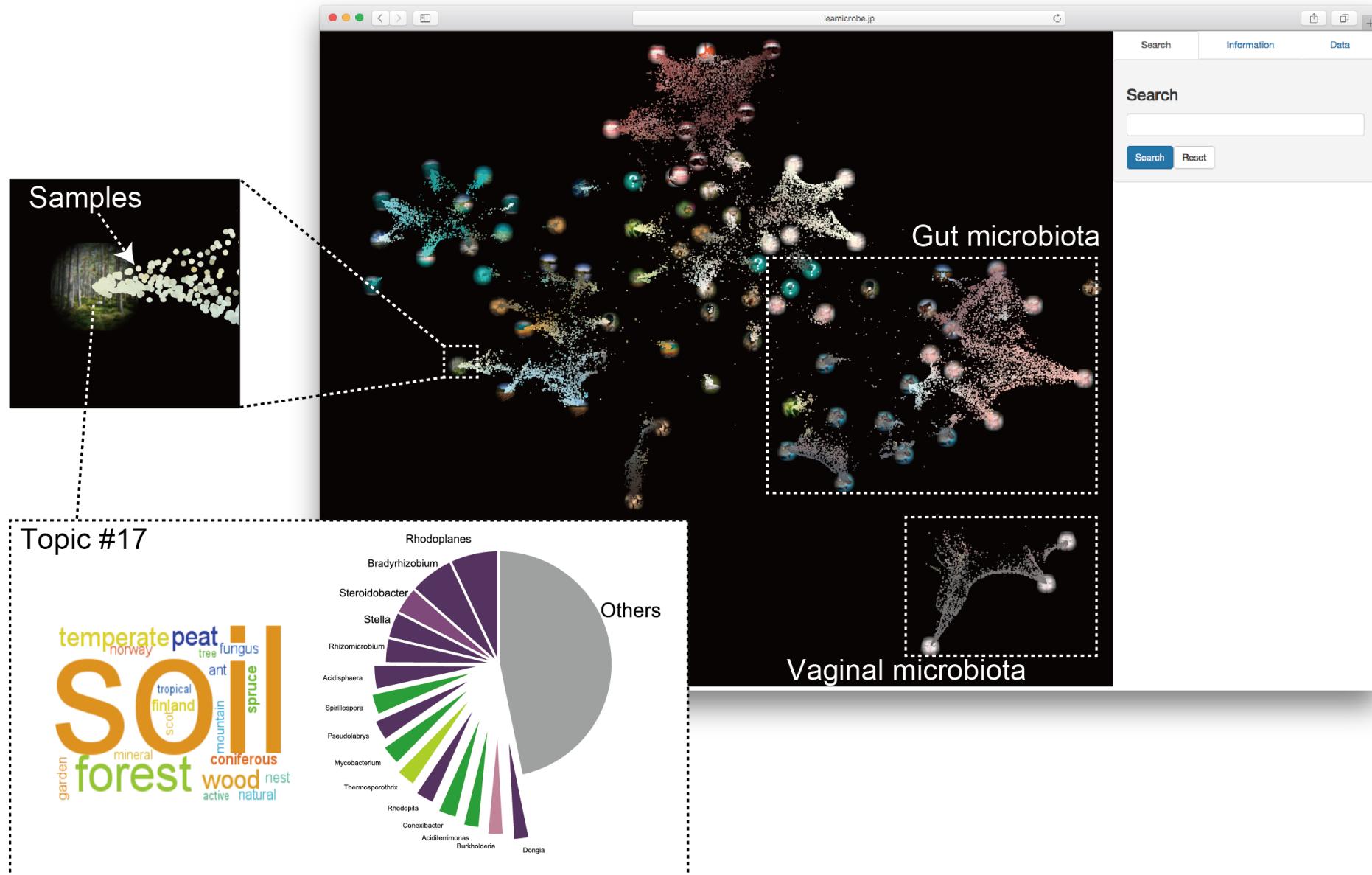


We extracted **80 topics** using >30,000 microbiome samples
 To visualize topic-based samples,
 we applied **parametric t-SNE** method as dimensionality reduction



LEA (Latent Environment Allocation)

<http://leamicrobe.jp/>



6 continents and some isolated islands

“Ocean” area



“Oral cavity” area



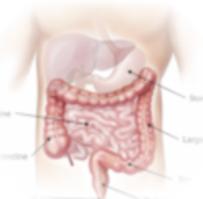
“Skin” area



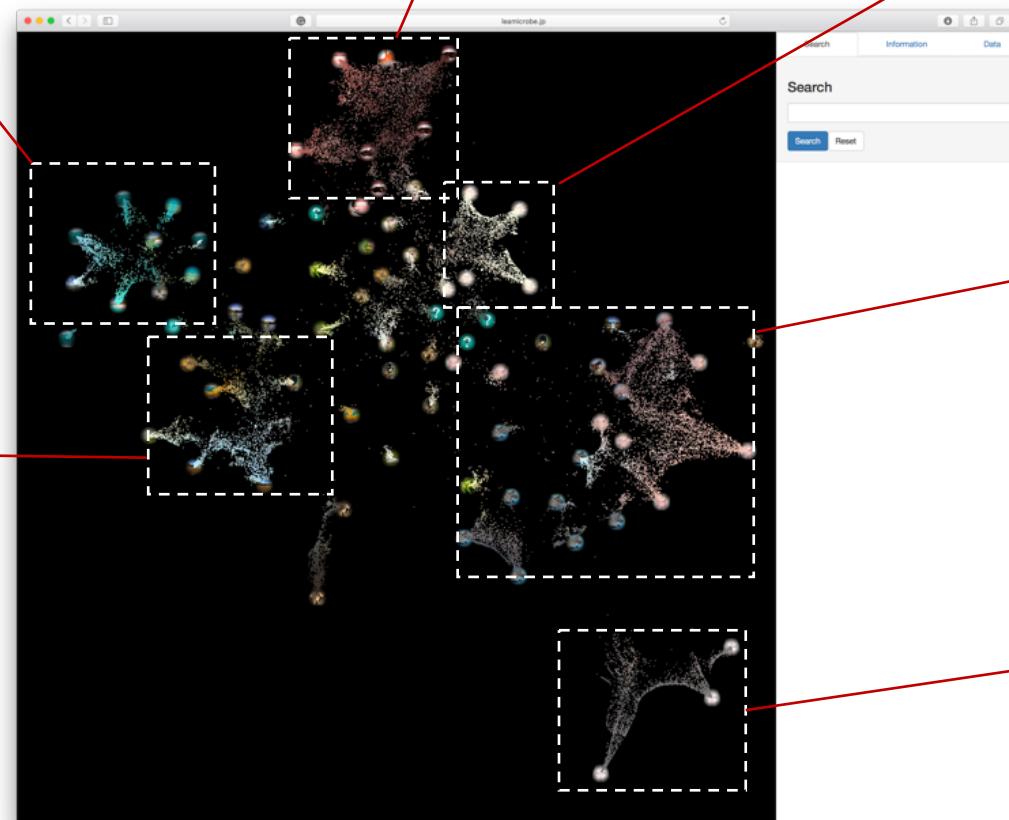
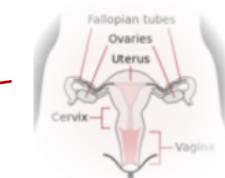
“Soil” area



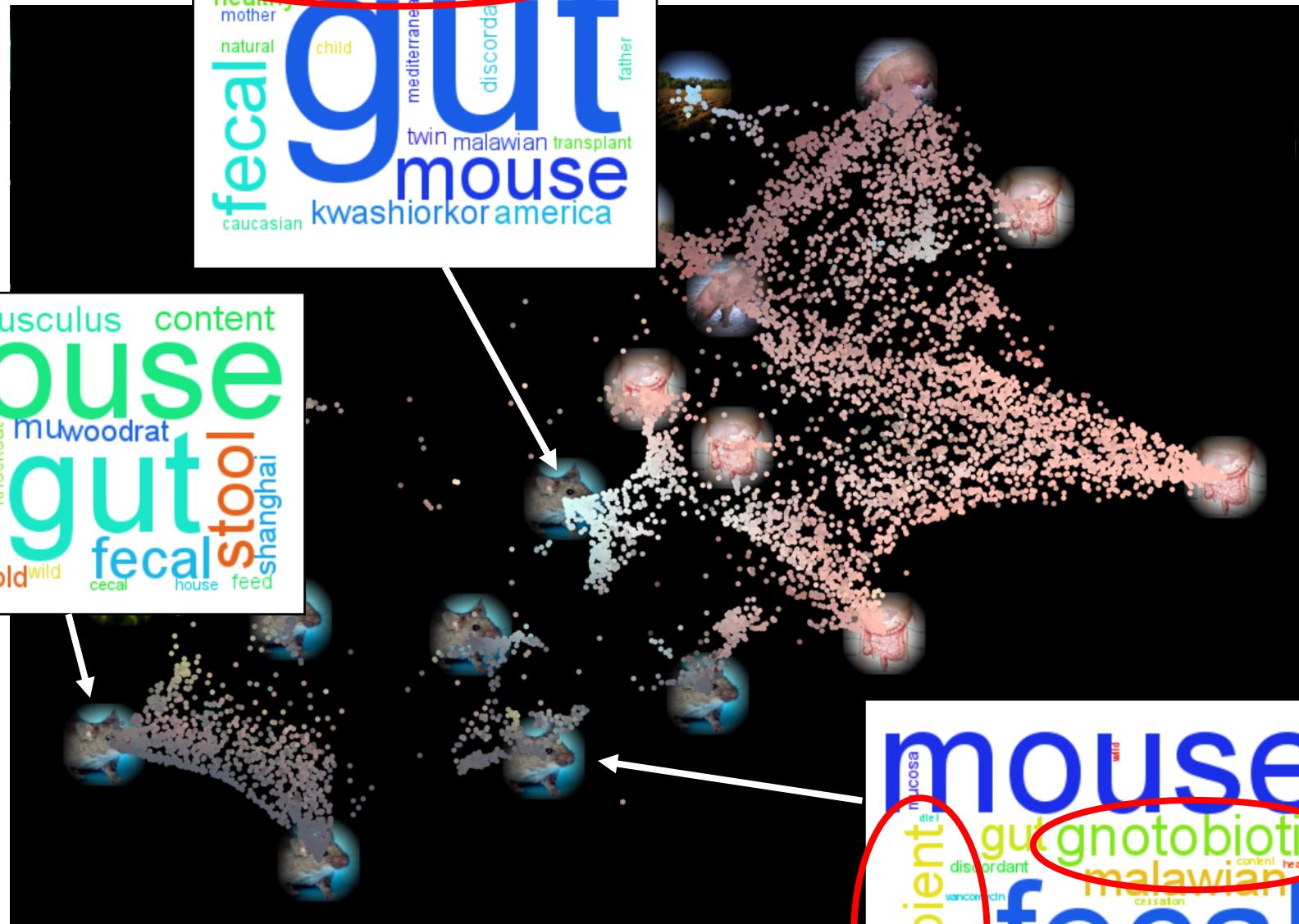
“Gut” area



“Vagina” area



mouse
gut
feces
diet
old
wild
field
knockout
cecum
inventory
musculus
content
muwoodrat
woodrat
shanghai
house
feed
cecral
twin

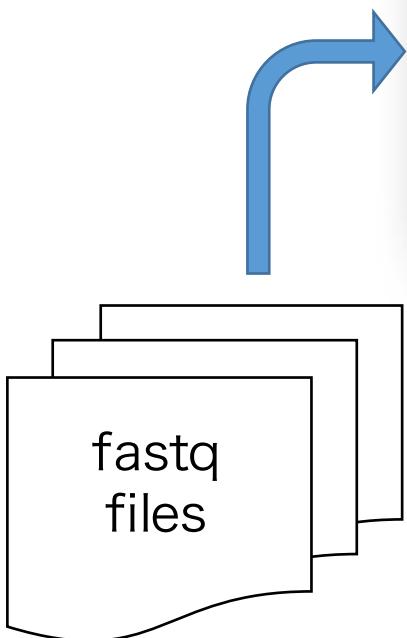


recipient
gnotobiotic
healthy
mother
natural
child
mediterranean
discordant
twin
malawian
transplant
caucasian
kwashiorkor
america

mouse
gut
gnotobiotic
recipient
twin
malawian
fecal
kwashiorkor
american
knockout
vancomycin
cessation
cessation
content
healthy
mucosa
ileal
disorder
twin
american
kwashiorkor
stomach
knockout
transplant
recipient

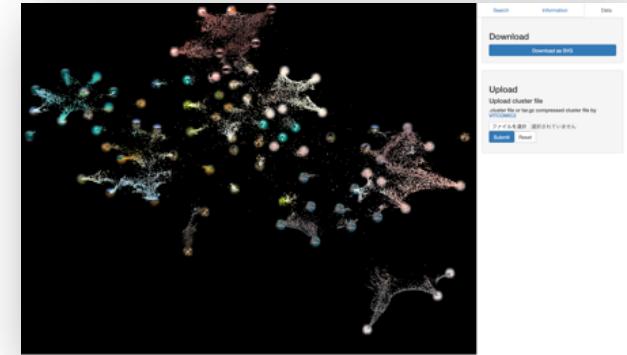
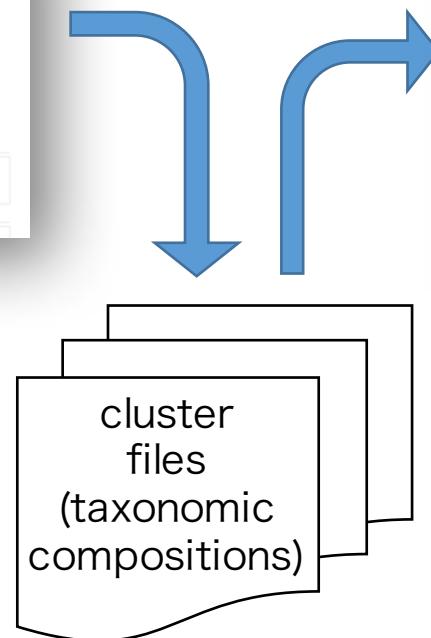
How to map your data on LEA

VITCOMIC2
(<http://vitcomic.org/>)



~ 6 min. /
a sample
(100,000
sequences)

LEA
(<http://leamicrobe.jp/>)



~ 0.5 sec. /
a sample

MicrobeDB.jp ver. 3の特徴

- かなりマイクロバイオームに特化した統合DB
- 環境や病気のオントロジーはメタデータから自動的にアノテーション
- 160万サンプル以上のマイクロバイオームデータをメタデータ等で検索可能
- MeGAP3による計算が終わったサンプルから、数ヶ月に一度、MicrobeDB.jpで検索できるように
- 現在主にセキュリティ面等をテスト中であり、12月中には<http://microbedb.jp/>から公開