# "Small Bugs, Big Data": Developing an integrated Database for Microbes with Semantic Web Technologies
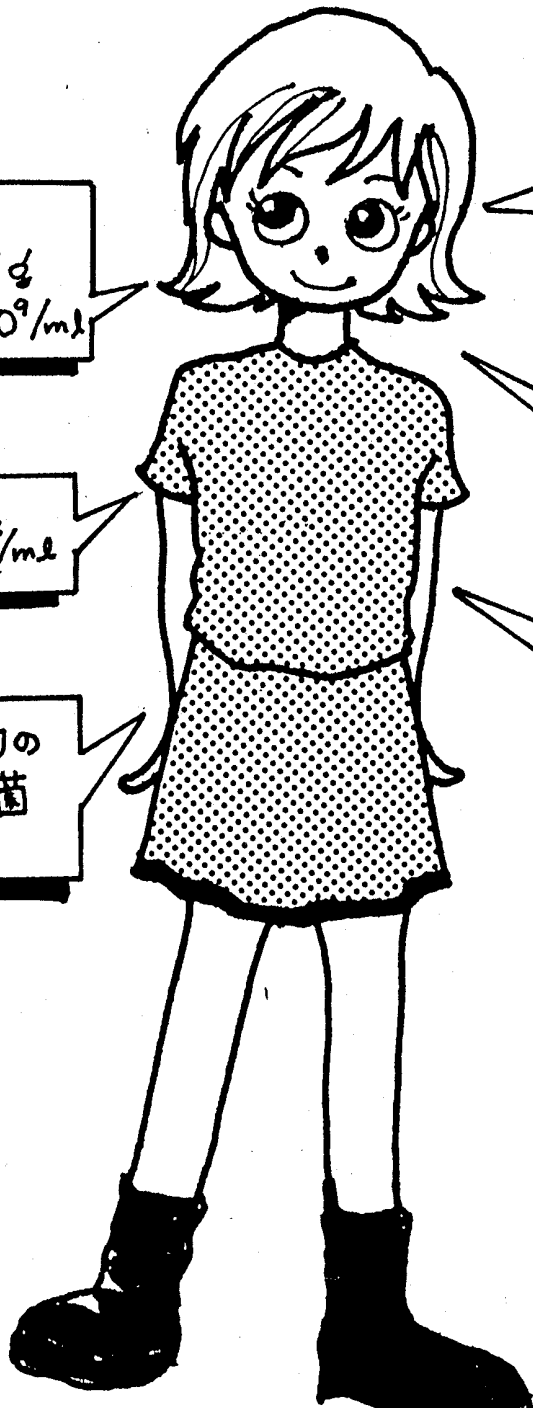
Ken Kurokawa (Earth-Life Science Institute, TITECH)

東京工業大学
Tokyo Institute of Technology

wpi
World Premier International
Research Center Initiative

JST

ELSI
EARTH - LIFE SCIENCE INSTITUTE

口腔：
　歯垢 $10^{11}/g$
　唾液 $10^5\sim10^9/ml$

皮膚：
　$10^3\sim10^6/cm^2$

鼻腔・副鼻腔
咽喉：
　鼻汁 $10^4\sim10^7/ml$

胃：胃液 $0\sim10^3/ml$

十二指腸・空腸：
　ほぼ無菌

大腸：固型物の
　$\frac{1}{2}\sim\frac{1}{4}$は細菌
　$10^{12}/g$

# *The next generation…*
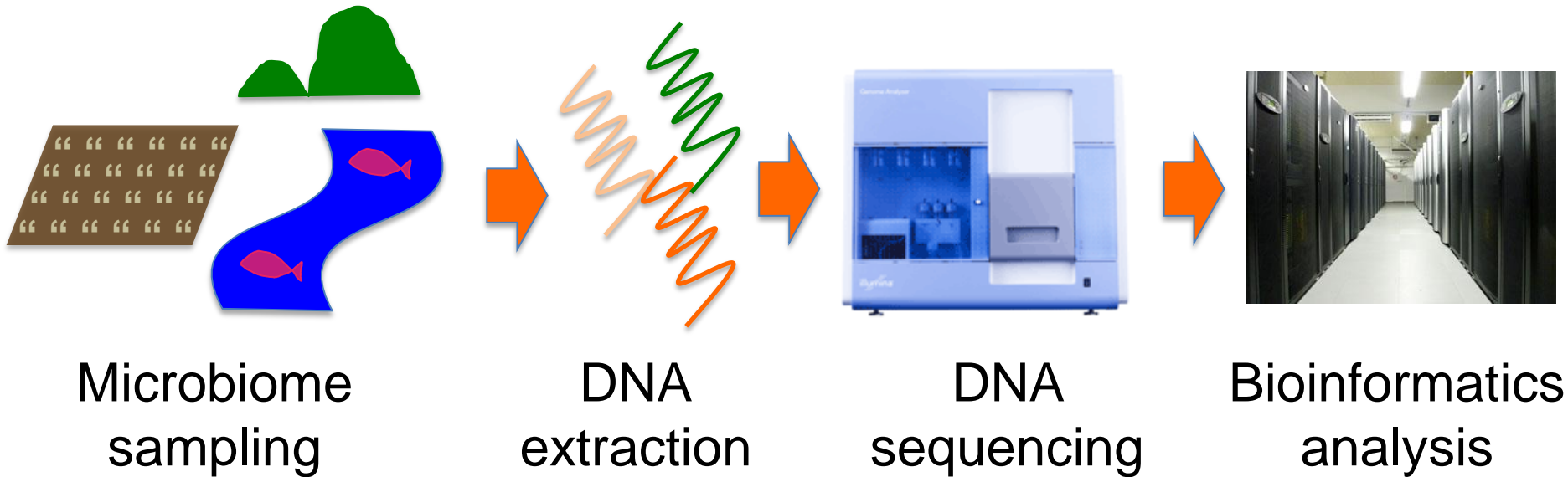
## Metagenomics

Bacterial community …

Natural environment (marine, river, soil..)
Human, animal (intestine, oral, skin..)

To elucidate the bacterial gene pools in environments, deeply sequencing the entire genomes extracted from bacterial community without cultivation

Genome analysis against "Microbiome"

# Metagenomic analysis



Microbiome
sampling

DNA
extraction

DNA
sequencing

Bioinformatics
analysis

# Both genomic and metagenomic data stored in public Databases

## Microbial genome data (NCBI RefSeq DB)

| Taxonomic division | Genomes |
|---|---|
| Archaea | 375 |
| Bacteria | 24,119 |

## Metagenomic data

| DB | Env. metagenomics | Human metagenomics |
|---|---|---|
| MG-RAST | 14,188 | 3,291 |
| JGI IMG/M | 1,694 | 840 |
| INSDC SRA | 23,214 | 18,108 |

Aprirl 2014

# The Importance of Metadata

Metadata are data about data (Gray et al. 2005). They are also about biology. Metadata are the descriptions of sampling sites and habitats that provide the context for sequence information. Metadata are of great importance for metagenomic sequence data for two reasons. First, only by fully describing the samples from which metagenomics sequences have been obtained can one have any possibility of replicating a study. Samples from environmental or biological sources can never be fully replicated, but it is important that samples be sufficiently well described for an independent researcher to have the possibility of resampling. Second, metadata are essential for the analysis of metagenomics sequence data. Metagenomic sequence data that lack an environmental context have no value.

Metadata are the description of sampling sites and habitats that provide the context for sequence information

Human microbiome:
   Age, Sex, BMI, Body habitat, Country, Diet, Disease stage, Family relationship …

Environmental microbiome:
   lat / long, pH, Depth, Dissolved oxygen, Wind speed, Total nitrogen, Temp. …

# Metagenomic analysis is performed in a variety of environments, and its massive data is stored in the public databases

**hot spring sediment**

| | |
|---|---|
| Identifiers | SRA: SRS152471 |
| Organism | human skin metagenome<br>unclassified sequences; metagenomes; organismal |

**Metadata**

| Attributes | | |
|---|---|---|
| | calcium | 3.3 |
| | chloride | 13.9 |
| | magnesium | 3.33 |
| | nitrate | 0.16 |
| | ph | 4.98 |
| | potassium | 5.85 |
| | sulfate | 1856 |
| | temp | 73.5 |

| Extra attributes | | |
|---|---|---|
| | biological_specimen | hot spring sediment |
| | env_biome | hot spring |
| | env_feature | spring |
| | env_matter | sediment |
| | latitude | 44.76 |
| | longitude | 110.43 |
| | sample_name | WB09-2 |

| | |
|---|---|
| Description | black sediment |
| Submission | Colorado School of Mines, Chuck Pepe-Ranney; 2011-01-12 |
| ID: 190846 | |

Temperature
pH
Anion & Cation
Long/Lat
:
:

## Use of both metagenome sequence data and its metadata will enable the large-scale comparative metagenomic analysis

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

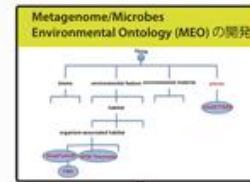Microbe DB .JP
http://microbedb.jp/

Microbe DB.jp
MicrobeDB.jp プロジェクトでは様々な微生物学上の知識を、ゲノム情報を核として遺伝子、系統、環境の３つの軸に沿ってセマンティックウェブの技術を駆使して整理統合し、幅広い分野での微生物学の発展に資することの出来るデータベースの構築を目標としています。

Metagenome/Microbes
Environmental Ontology (MEO) の開発

**Ontology**
オントロジー: 検索タームの柔軟化＆明確化

Gene ⟷ Taxon ⟷ Environment

Ortholog: MBGD
オーソログデータ

Genome: GTPS/RefSeq
オミックスデータ

Annotation:
TogoAnnotation
モデル微生物の高品質アノテーションデータ

Taxonomy:
NCBI Taxonomy
系統分類データ

Culture Collection:
NBRC/JCM
菌株データ

Metadata:
INSDC SRA
環境のメタデータ

Metagenome
INSDC SRA
メタゲノムデータ

# Integration of microbe's data centers around genome information

**Phylogenetic information**

**Environmental information**

Strain data
taxonomy
optimal temp.
medium…

Genome data
taxonomy
gene number
Isolation source…

Metagenome data

Meta data
pH environment
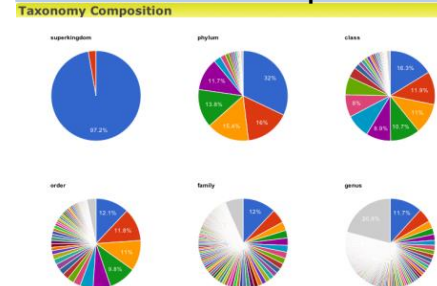temperature…

**Genome information**

Gene A

Taxonomic composition

Ortholog data
organism
gene name
gene function

Gene annotation
gene name
gene function…

Functional compostion

Gene A'

GTPS/RefSeq

Togo Annotation

Accurate annotation
from Model-organisms

**Genetic information**

# RDF is a standard data model of Semantic Web technology

RDF (Resource Description Framework)
 Data model which uses Triples
 (**S**ubject – **P**redicate – **O**bject)



| S | P | O |
|---|---|---|
| <URI> | <URI> | <URI>/Literal |
| gtps:Gene1 | rdfs:label | "16S rRNA gene" |

URI node can be linked to other nodes



**RDF**



Gene1 — has Function — KO:03043 / GO:0003700

Genome1 — organism — *Escherichia coli*

Organism1 — has Genome — Genome1

Organism1 — inhabit — Lake

**Ontology**

**Triple store**

**SPARQL**

Search

To prepare data in RDF,
the database management system automatically recognize same resources.

# An example of RDF relationships for *E. coli* K-12 genome data



Various types of data are integrated within a genome sequence, and we can retrieve all information about *E. coli* K-12 by following these graphs.

# How to integrate the data from two different DBs?

**DB 1**

**DB 2**



1. When two DBs use same URI, already two DB's data are integrated.
2. If not, you can integrate two DB's data by adding one Triple (db1:A owl:sameAs db2:B)

   You don't need to place all of these data in one DB management system.

   How can we discriminate whether two DB's resources are same or not?

# You should describe your resource by using some Ontologies

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

For example, to answer: What is soil? What is a relationship between soil and sand?

## MEO (Microbes Environmental Ontology)

- Thing
  - atmosphere
    - 'aerobic environment'
    - aerosol
    - 'anaerobic environment'
  - geosphere
    - 'geographic feature and biome'
    - 'rock, sand and soil'
    - sediment
  - 'human activity association'
    - 'artificial natural environment'
    - food
    - fuel
    - 'large-scale artifact'
    - 'liquid artifact'
    - 'small artifact'
    - 'waste treatment'
  - hydrosphere
    - 'aquatic feature and biome'
    - ice
    - water
  - 'organism association'
    - 'animal associated'
    - biofilm
    - 'excrement and secretion'
    - 'fungi associated'
    - 'organic feature and biome'
    - 'plant associated'
    - 'symbiotic microbe'

## PDO (Pathogenic Disease Ontology)

- 'Disease involving body sites'
  - 'Breast disease'
  - 'Cardiovascular disease'
  - 'Digestive system disease'
  - 'Immune system disease'
  - 'Musculoskeletal system disease'
  - 'Nervous system disease'
  - 'Reproductive system disease'
  - 'Respiratory system disease'
  - 'Skin disease'
  - 'Systemic disease'
  - 'Urinary system disease'
- 'Disease involving unidentified body site'

## MCCV (Microbial Culture Collection Vocabulary)

## MSV (Metagenome Sample Vocabulary)

## MPO (Microbial Phenotype Ontology)

## MBGD Ortholog Ontology

Most of them can be obtained from

BioPortal

# More than 1 billion Triples!

| グラフ名 | 説明 | 作成元 | トリプル数 |
|---|---|---|---|
| refseq | RefSeq Prokaryoteゲノムデータ | DBCLS | 550,273,744 |
| mbgd | MBGD Orthologデータ | 基生研 | 291,714,037 |
| gtps | GTPSゲノムデータ | 遺伝研 | 197,069,932 |
| taxonomy | SPARQLthonで作成したNCBI Taxonomyオントロジー改良版版 | DBCLS,遺伝研,東工大 | 10,183,714 |
| meta16S | 各SRSメタ16Sの系統組成データ | 東工大 | 9,831,600 |
| gazetteer | 地理オントロジー | 外部機関 | 7,062,536 |
| srs_metadata | SRSメタ16S・メタゲノムの様々なメタデータ | 東工大 | 4,982,739 |
| srs_ortholog | 各SRSメタゲノムのMBGD Ortholog組成 | 東工大,基生研 | 2,026,746 |
| go | Geneオントロジー | 外部機関 | 1,211,571 |
| brc | JCM/NBRC菌株データ with NCBI Taxonomy ID | 遺伝研,東工大,DBCLS | 903,319 |
| gold | GOLDの個別ゲノムのMEO等へのオントロジーマッピングデータ | 東工大,DBCLS | 150,899 |
| srs | SRSメタ16S・メタゲノムのMEO等へのオントロジーマッピングピングデータ | 東工大 | 53,691 |
| so | Sequenceオントロジー | 外部機関 | 43,060 |
| pdo | 感染症オントロジー + 症状オントロジー + ゲノムへのオントオントロジーマッピングデータ | 東工大 | 8,809 |
| meo | 微生物の生息環境オントロジー | 東工大 | 4,975 |
| msv | SRSメタ16S・メタゲノムのメタデータオントロジー | 東工大 | 1,601 |
| mpo | 微生物フェノタイプオントロジー | DBCLS | 734 |
| mccv | 菌株オントロジー | 東工大,DBCLS | 293 |
| その他中間データ | いくつかのデータ集計系のSPARQLクエリは遅いため、MSSMSSが集計結果のデータを作成 | | 440,773 |
| 合計 | | | 1,075,964,773 |

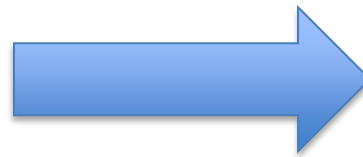# http://microbedb.jp/

# Stanza Development

To obtain biological knowledge from low data (sequence and metadata), we developed a variety of "Stanza", which is a compact, modular, and reusable application for data analysis.



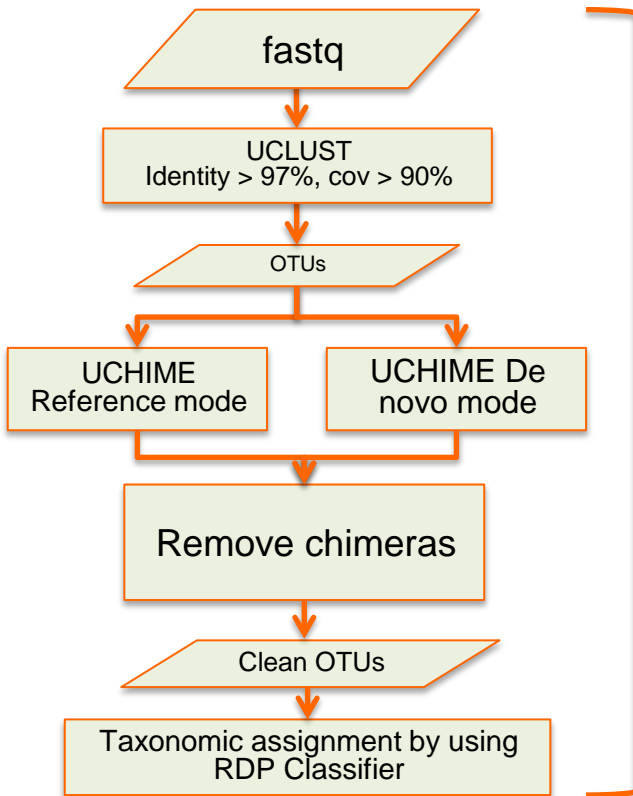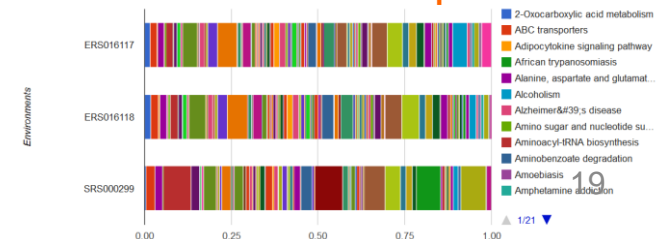Correlation analysis between gene abundance and metadata

Analyze data by using the Stanza

Comparison of taxonomic composition

# Stanza categories in MicrobeDB.jp



**Genes**

Gene Definition
Gene Publication
Ortholog Definition

Sample Function
Mapping to Environment (Chromosome)
Mapping to Environment (Plasmid)

Gene Annotation
Ortholog Group Members
Ortholog Cluster
Genome Information
GTPS Gene/Genome Feature
RefSeq Gene/Genome Feature

Ortholog Abundance among Environments
Ortholog Abundance in Environment

**Taxon**

**Environment**

GTPS Genome
GTPS Genome Definition
Other Collection Numbers
Pathogen Information
Phenotype Information
RefSeq Genome
RefSeq Genome Definition
Strain Definition
Strain Genome
Strain Reference
Taxon Definition
Taxon Hierarchy

Disease Definition
Environment Definition
MEO Hierarchy
MEO Ontology View
Meta16S Sample List
Metagenome Sample List
Numeric Metadata Histogram
Sample Definition
Sample Metadata
SRS Cross Reference
Symptom List

Genome-Sequenced Strains
Sequenced Genome List
Strain List
Taxonomic Composition of Genomes
Taxonomic Composition of Meta 16S
Human Meta Body Mapping
Strain Metadata

Cost per Raw Megabase of DNA Sequences

Data from the NHGRI Genome Sequencing Program (GSP)
http://www.genome.gov/sequencingcosts/

# MicrobeDB.jp Project Team

・ <u>Ken Kurokawa (Tokyo Institute of Technology)</u>
   Hiroshi Mori, Junichi Takehara, Koji Yoshino, Shinya, Suzuki, Nozomi Yamamoto, Takuji Yakada

・ <u>Yasukazu Nakamura (National Institute of Genetics, DDBJ)</u>
   Takatomo Fujisawa, Eri Kaminuma, Hideaki Sugawara

・ <u>Ikuo Uchiyama (National Institute for Basic Biology)</u>
   Hirokazu Chiba, Hiroyo Nishide

Advisor (DataBase Center for Life Science)
   Shinobu Okamoto, Shuichi Kawashima, Toshiaki Katayama, Yasunori Yamamoto, Shoko Kawamoto

# Funding