

ライフサイエンスデータベース統合推進事業
統合化推進プログラム
平成24年度 進捗報告会

ゲノム・メタゲノム情報を基盤とした 微生物DBの統合

東京工業大学大学院生命理工学研究科
黒川 顕



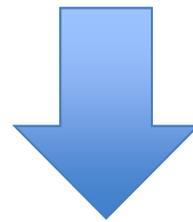
©2013黒川 顕(東京工業大学) licensed under CC表示2.1日本

微生物研究を取り巻く状況

- 微生物は環境と密接に関与し存在する
- 微生物研究はバイオ分野のみならず, 他の多くの分野と連携可能
- 微生物研究分野には多様なDBが多数存在する
- 環境との関連性を記述しているDBは未だ存在しない
- ほとんどのDBは専門知識を持っていないバイオ分野以外の人には利用困難

研究開発の目標・ねらい

ゲノム情報を核として様々な微生物学上の知識を統合し、幅広い分野での微生物学の発展に資することのできる「**微生物エンサイクロペディア**」の構築を目標とする。



**微生物学分野のオミックス研究の発展に寄与
データ駆動型研究による新しい仮説の提唱**

研究開発メンバー

東京工業大学

黒川 顕: 微生物DBにおける研究統括

森 宙史: ゲノム・メタゲノムDB、メタデータの構築、オントロジー構築

山田拓司: メタゲノムDBの構築

山本希: オントロジーの構築

吉野弘二, 竹原潤一: メタデータDBの構築、オントロジー構築

小西史一: スパコンにおける解析システムの開発および実装

国立遺伝学研究所

中村保一: 微生物アノテーションリファレンスの整備と共用化

藤澤貴智: TogoAnnotation構築、モデル微生物情報の高度化

菅原秀明: 微生物ゲノム基盤情報資源の共用化

神沼英里: TogoAnnotationの拡張

基礎生物学研究所

内山郁夫: 比較ゲノム解析に立脚した微生物ゲノム情報の統合化

千葉啓和: MBGDの統合化

統合データベースセンター(技術アドバイザー)

岡本忍, 片山俊明, 川島秀一, 川本祥子, 山本泰智: 技術協力

H24年度の当初計画

- 保存菌株情報(NBRC, JCM)のRDF化
- 培地情報オントロジーCMOの整備(w/ DBCLS)
- GTPSのRDF化および各オミックスデータの統合
- MBGDオーソログ遺伝子情報のRDF化
- 放線菌以外のモデル微生物ゲノムアノテーション高度化
- メタゲノムデータのRDF化およびGTPS、RefSeq等との統合
- 各種オントロジーの開発
- 各種アプリケーション、結果表示要素「Stanza」の標準化および開発(w/ DBCLS)

微生物における各種DBを統合化し、環境情報との連携を徹底的に記述した新たなDBの構築

H24年度の実施状況

- セマンティックウェブの技術を活用
- データ間をリンクするためのゲノム情報、オーソログ遺伝子情報、メタゲノム情報の整備
- 全データのRDF化、各データID間のリンク構築
- 各種オントロジーの開発、各データにマッピング
- アノテーション高度化システムの開発
- ユーザ認証システム
- データテンプレート「Stanza」の開発

RDF/OWLファイルを用いてFull RDFなDB「MicrobeDB.jp」の構築 ⁶

NBRC/JCMの菌株データのRDF化



NBRC No.	NBRC 102086
Scientific Name of this Strain	<i>Burkholderia multivorans</i> Vandamme et al. 1997
Synonymous Name	
Type Strain	
History	NITE <- Meijo Univ. (S. Ichihara, 25)
Other Culture Collection No.	
Other No.	NITE 02208=25
Rehydration Fluid	702
Medium	802
Cultivation Temp.	25 C
Source of Isolation	
Locality of Source	
Country of Origin	Japan
Biosafety Level	
Applications	Phenylacetic acid;degradation
Mating Type	
Genetic Marker	
Plant Quarantine No.	
Animal Quarantine No.	
Herbarium No.	
Restriction	
Comment	
References	
Sequences	16S rDNA

株数：約16,000株
単離元：1,627
培地情報：432種類



ゲノムデータおよび
メタゲノムデータ等と
統合を目標にRDF化

JCMも対象とする(約
14,000株)

NBRCおよびJCMデータ項目とその内容の概観

1. Strain_Number(株番号)
2. Other_Collection_Numbers(他機関での株番号)
3. Name(学名)
4. Organism_Type(生物分類)
5. History_of_Deposit(来歴)
6. Date_of_Isolation(分離日)
7. Isolated_from(分離源)
8. Geographic_Origin(地理)
9. Status(ステータス)
10. Optimum_Temperature_for_Growth(至適培養温度)
11. Maximum_Temperature_for_Growth(最高培養温度)
12. Minimum_Temperature_for_Growth(最低培養温度)
13. Medium(培地)
14. Application(利用法)
15. Literature(文献)

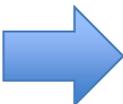
NBRC/JCMにおける総トリプル数

	NBRC株	JCM株
株番号	17,367	13,396
他機関での株番号	62,914	122,286
学名	48,791	51,769
生物分類	45,771	35,350
来歴	59,773	58,442
分離日	0	0
分離源	57,240	46,275
地理	15,338	0
ステータス	4,090	6,004
至適培養温度	16,611	13,396
最高培養温度	746	0
最低培養温度	746	0
培地	18,429	18,787
利用法	1,470	1,316
文献	21,600	15,646

菌株IDとNCBI Taxonomy IDの名寄せ

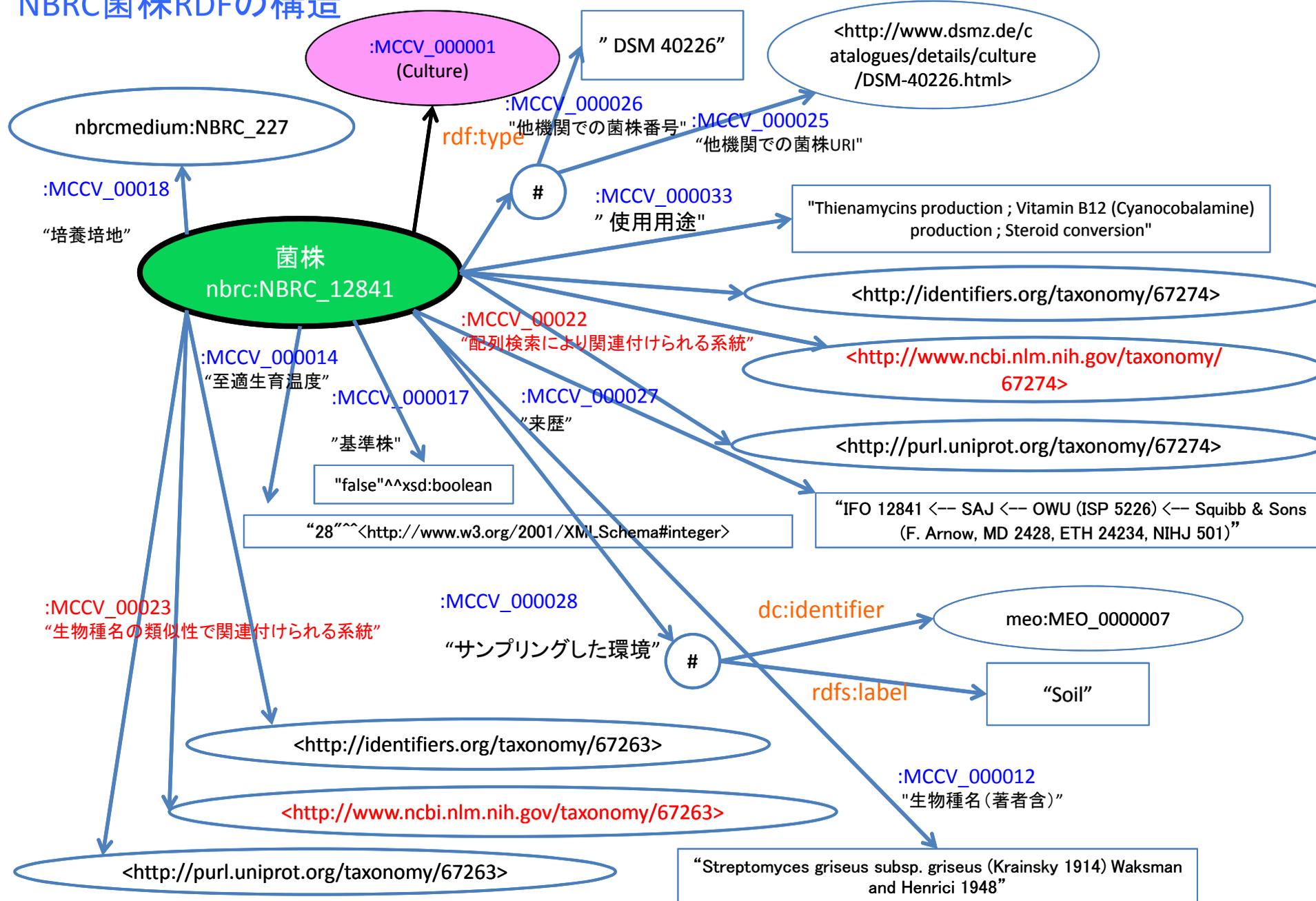
National Institute of Technology and Evaluation		独立行政法人 製品評価技術基盤機構
NBRC		Biotechnology Center
NITE Biological Resource Center		
NBRC No.	NBRC 12841	
Scientific Name of this Strain	<i>Streptomyces griseus</i> subsp. <i>griseus</i> (Krainsky 1914) Waksman and Henrici 1948	
Synonymous Name	Synonym: <i>Streptomyces griseus</i>	
Type Strain		
History	IFO 12841 <- SAJ <- OWU (ISP 5226) <- Squibb & Sons (F. Arnow, MD 2428, ETH 24234, NIHJ 501)	
Other Culture Collection No.	ATCC 11009=ATCC 23882=CBS 662.68=RIA 1168=ISP 5226=AS 4.1693=BCRC 11815=DSM 40226=JCM 4229=JCM 4623=KCTC 1742=LMG 5967=NCIMB 9625=NRRL B-1806=VKM Ac-747	
Other No.	IMET 43659	

現状のNBRC/JCMの菌株データは、配列データとは独立しており、両者の「糊しろ」となるNCBI Taxonomy IDも菌株データには付加されていない


 菌株データと配列データを統合するため、NBRC/JCM菌株IDとNCBI Taxonomy ID間の名寄せが必要
 名寄せ結果を含む菌株RDFで用いる様々な語彙のオントロジー、MCCV (Microbial Culture Collection Vocabulary)を設計し、これを用いて菌株RDFを作成

11

NBRC菌株RDFの構造



複数の異なる方法による菌株IDとNCBI Taxonomy ID間の名寄せ結果を、predicateで方法ごとに区別可能

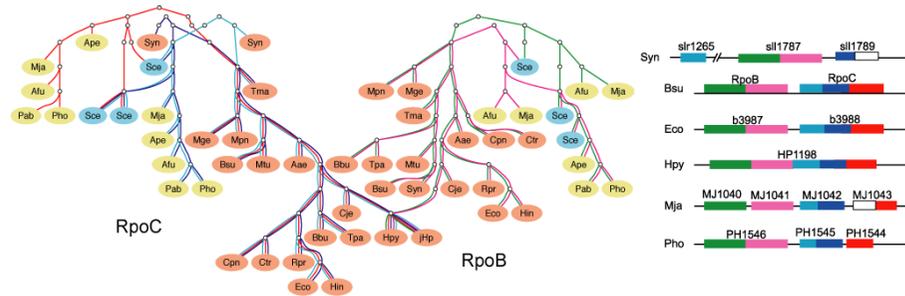


MBGDオーソログ遺伝子情報のRDF化

微生物比較ゲノムデータベースMBGD



オーソログ分類アルゴリズム DomClust



オーソロググループ

オーソログテーブル

MBGD Ortholog Cluster Table Overview

A Multiple sequence alignment (ClustalW MAP MAFFT)
M Multiple genome map comparison
H Find homologous clusters
P Similar phylogenetic pattern search (Correlation coefficient Hamming distance Mutual information)

Analyze the checked clusters
 Whole genome map comparison Merging multiple clusters

ClusterID	Name	#species	#genes	Description	Phylogenetic pattern (Set species color)
<input type="checkbox"/> O1149 A M H P	folE	328	377	GTP cyclohydrolase I	
<input type="checkbox"/> O1279 A M H P	folB	320	352	Dihydroneopterin aldolase	
<input type="checkbox"/> O1033 A M H P	ubiE	367	397	Ubiquinone/menaquinone biosynthesis methyltransferase	
<input type="checkbox"/> O990 A M H P	folK	371	406	2-amino-4-hydroxy-6-hydroxymethylidihydropteridine pyrophosphokinase	
<input type="checkbox"/> O952 A M H P	aroB	374	412	3-dehydroquinate synthase	
<input type="checkbox"/> O1077 A M H P	atpH	382	388	FOF1 ATP synthase subunit delta	
<input type="checkbox"/> O1399 A M H P	hemH	324	332	Ferrochelatase	
<input type="checkbox"/> O662 A M H P	sdhB	370	485	Succinate dehydrogenase / fumarate reductase protein, iron-sulfur subunit	
<input type="checkbox"/> O985 A M H P	yggS	386	405	Alanine racemase domain-containing protein	
<input type="checkbox"/> O315 A M H P	fabF	388	691	3-oxoacyl-(acyl carrier protein) synthase II	

Ortholog Cluster

Cluster	1149
Gene	folE
Title	GTP cyclohydrolase I
Size	328 species, 377 genes
Xref-COG	COG0302 GTP cyclohydrolase I [Equivalent, 34/35]
Xref-KEGG	K01495 GTP cyclohydrolase I [EC:3.5.4.16] [Equivalent, 249/258]
Xref-TIGR	TIGR00063 GTP cyclohydrolase I [Subgroup, 172/177]
Xref-GO	GO:0003934 GTP cyclohydrolase I activity [Supergroup, 305/762] GO:0016787 hydrolase activity [Supergroup, 303/173030] GO:0005737 cytoplasm [Supergroup, 287/174310] GO:0046654 tetrahydrofolate biosynthetic process [Supergroup, 287/604] GO:0006730 one-carbon compound metabolic process [Supergroup, 268/4398]

[Summary] [Gene List] [Clustering Tree]

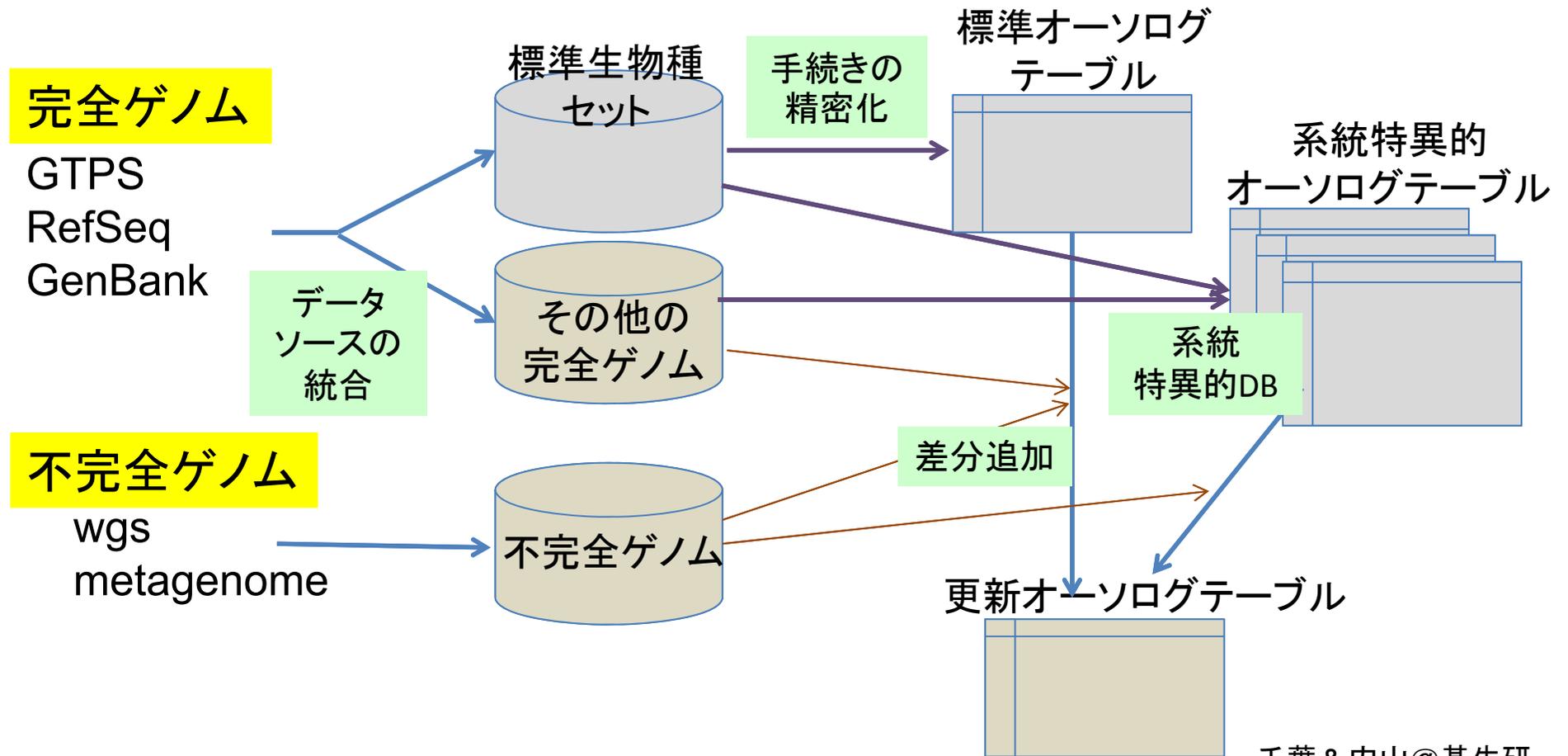
Redraw Display xref-ortholog Display xref-motif (cutoff 0.001)

ON OFF

Phylum	Species	Orf ID	Description
Actinobacteria <input type="checkbox"/> ON <input type="checkbox"/> OFF	<i>Acidimicrobium ferrooxidans</i> DSM 10331	<input type="checkbox"/> afo:AFER_1912	GTP cyclohydrolase I
	<i>Arcanobacterium haemolyticum</i> DSM 20595	<input type="checkbox"/> ahc:ARCH_0252	GTP cyclohydrolase I
	<i>Mobiluncus curtisii</i> ATCC 43063	<input type="checkbox"/> mcu:HMPREF0573_10367	GTP cyclohydrolase I
	<i>Catenulispora acidiphila</i> DSM 44928	<input type="checkbox"/> cai:CACL_2639	GTP cyclohydrolase I
		<input type="checkbox"/> cai:CACL_8425	GTP cyclohydrolase I
	<i>Corynebacterium glutamicum</i> ATCC 13032	<input type="checkbox"/> cgl:NCGL2602	GTP cyclohydrolase I
	<i>Gordonia bronchialis</i> DSM 43247	<input type="checkbox"/> gbr:GBRO_4016	GTP cyclohydrolase I
	<i>Mycobacterium tuberculosis</i> H37Rv	<input type="checkbox"/> mtu:RV3609C	GTP cyclohydrolase I
	<i>Nocardia farcinica</i> IFM 10152		

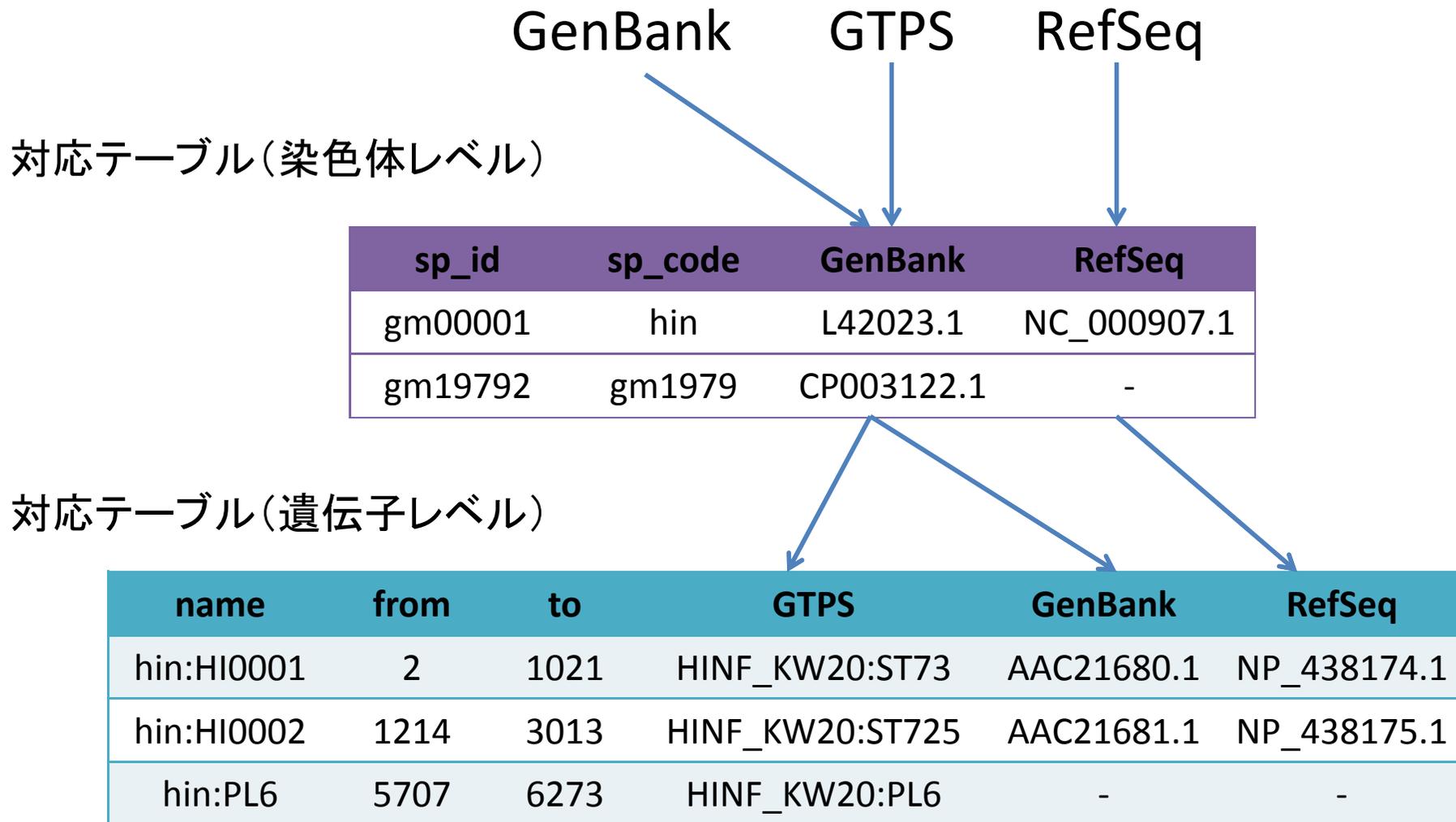


対象ゲノムデータの拡大と 効率的なオーソログ解析



千葉&内山@基生研

データソースの統合による拡張



千葉 & 内山 @ 基生研

系統特異的オーソログテーブルの作成

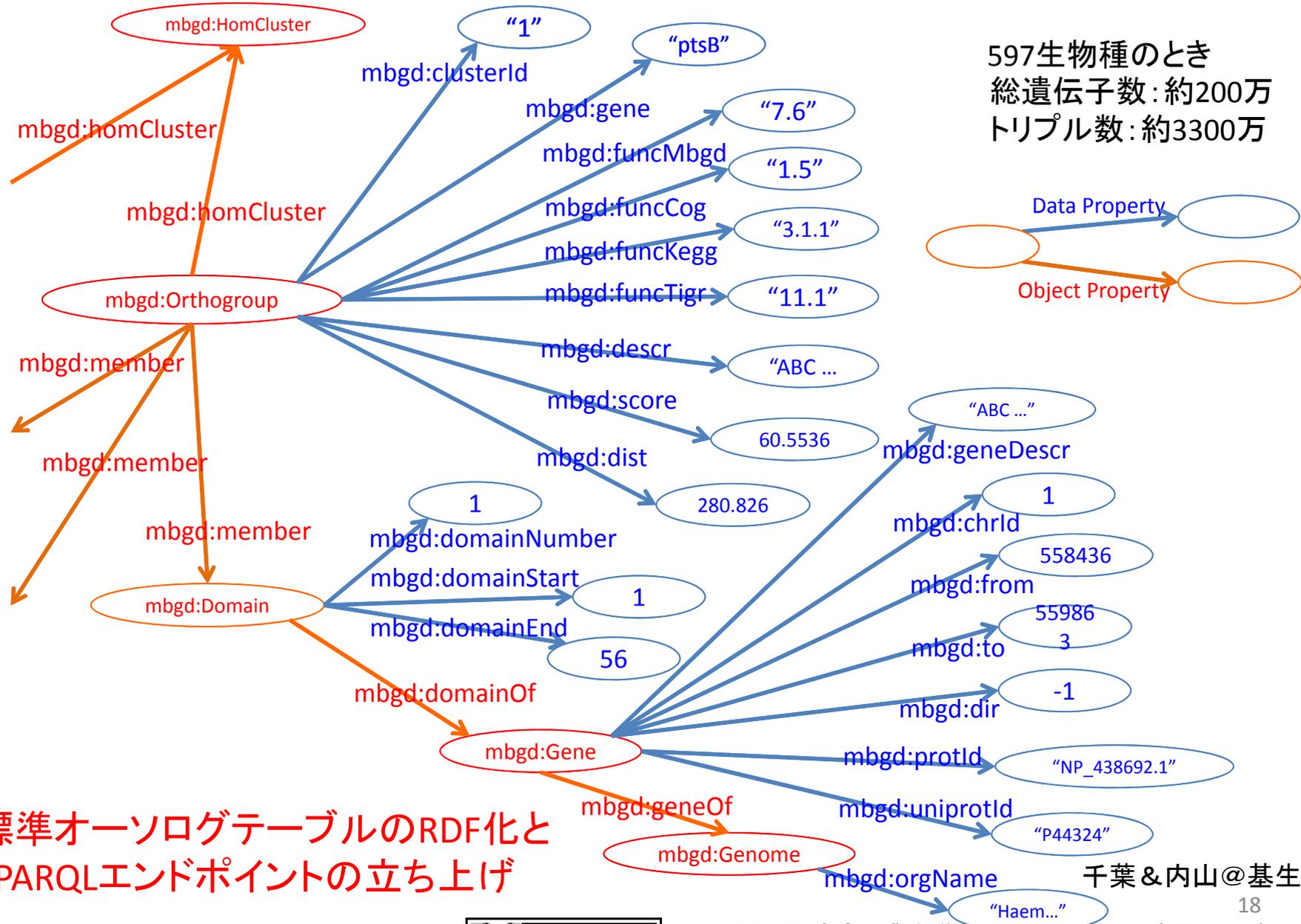
Target rank	Rank for representative selection	CoreAlign is available	Num. of taxa
all (default)	genus	No	1
all (extended)	genome	No	1
superkingdom	genus	No	2
phylum	genus	No	13
class	genus	No	18
order	species	No	50
family	species	Yes	58
genus	species	Yes	31
species	genome	Yes	21

代表生物種を6種以上含むタクサの数

千葉&内山@基生研

17

MBGDのRDF化



標準オーソログテーブルのRDF化と
 SPARQLエンドポイントの立ち上げ

千葉&内山@基生研



モデル微生物アノテーション高度化

文献情報に基づくモデル微生物 ゲノムデータベースの現状

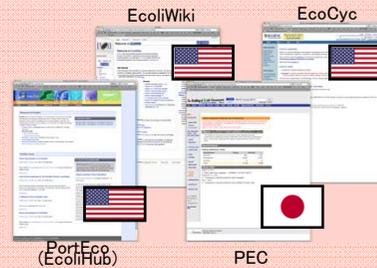
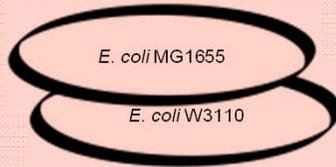
モデル微生物

リファレンス株

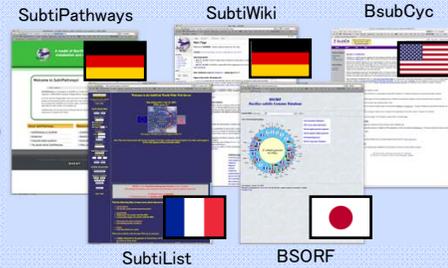
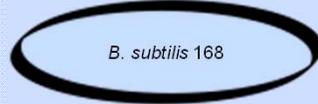
リファレンス株遺伝子の関連文献が
参照可能なデータベース

国内でゲノム解析された
病原性/産業有用株

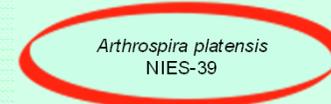
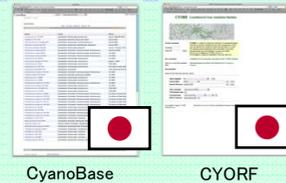
大腸菌



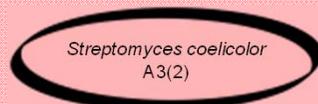
枯草菌



ラン藻



放線菌



藤澤&神沼&中村@遺伝研

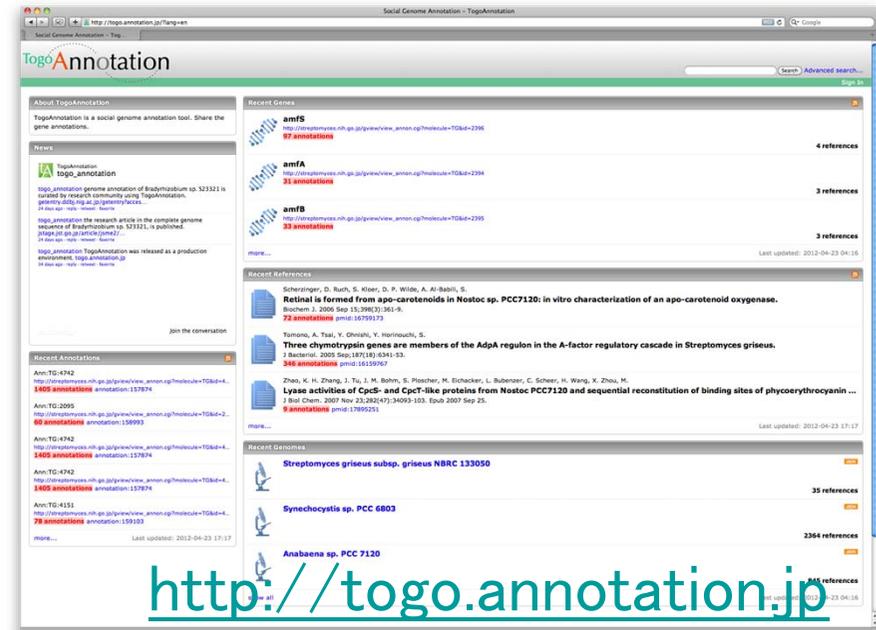


©2013黒川 顕(東京工業大学) licensed under CC表示2.1日本

TogoAnnotationの拡張整備

▪TogoAnnotation

ソーシャルブックマークシステムを基盤技術として構成様々なタイプのエンティティを容易にアノテーションすることが可能



☑ URL, サービス名称の変更 (H23)

KazusaAnnotationから改名 (H24.3)

☑ ゲノム、遺伝子、文献単位での再利用を目的としたAPI開発 (H23)

☑ モデル微生物情報(シアノバクテリア、放線菌)のマニュアルキュレーションによる文献リファレンス蓄積 (H23,H24)

☑ OpenID認証と連携したユーザ/グループ認証サービス開発・運用 (H24)

☑ モデル微生物情報(大腸菌、枯草菌)の文献リファレンス蓄積のための整備 (H24)

☑ TogoAnnotationのRDF化 (H24)

☑ モデル微生物情報のRDFを介した情報統合のユースケースプロトタイプ (H24)



モデル微生物のマニュアルキュレーションによる 文献リファレンス蓄積: 遺伝子情報

ver. 2013.01.15

モデル微生物	生物種	対象文献数	アノテーション 文献数	アノテーション 遺伝子数	アノテーション ブックマーク数	アノテーション ブックマーク数 (H23-)
大腸菌	<i>Escherichia coli</i> K-12					
枯草菌	<i>Bacillus subtilis</i> 168					
ラン藻	<i>Synechocystis</i> sp. PCC 6803	2524	2244	3003	74363	7543
	<i>Anabaena</i> sp. PCC 7120	974	929	2685	28031	8196
	<i>Synechococcus elongatus</i> PCC 7942	900	787	763	16356	2105
	<i>Synechococcus</i> sp. PCC 7002	299	251	251	3943	236
	<i>Thermosynechococcus elongatus</i> BP-1	329	252	2534	6659	304
	<i>Chlorobium tepidum</i> TLS	180	143	751	5532	11
	<i>Nostoc punctiforme</i> ATCC 29133	138	143	751	3242	870
	<i>Anabaena variabilis</i> ATCC 29413	130	113	246	1712	121
	<i>Prochlorococcus marinus</i> MED4	66	59	383	2119	335
	<i>Gloeobacter violaceus</i> PCC 7421	49	46	4486	5588	4754
	<i>Prochlorococcus marinus</i> MIT9313	41	40	235	885	166
	<i>Prochlorococcus marinus</i> SS120	40	34	127	527	70
	<i>Arthrospira platensis</i> NIES-39	8	7	258	785	391
	<i>Synechococcus</i> sp. WH8102	1	2	3	7	1
	<i>Synechococcus elongatus</i> PCC 6301	1	1	1	4	4
	<i>Trichodesmium erythraeum</i> IMS101	1	1	1	2	2
放線菌	<i>Streptomyces coelicolor</i> A3(2)	2143	21	266	2958	2958
	<i>Streptomyces griseus</i> IFO 13350	169	112	299	12814	12814

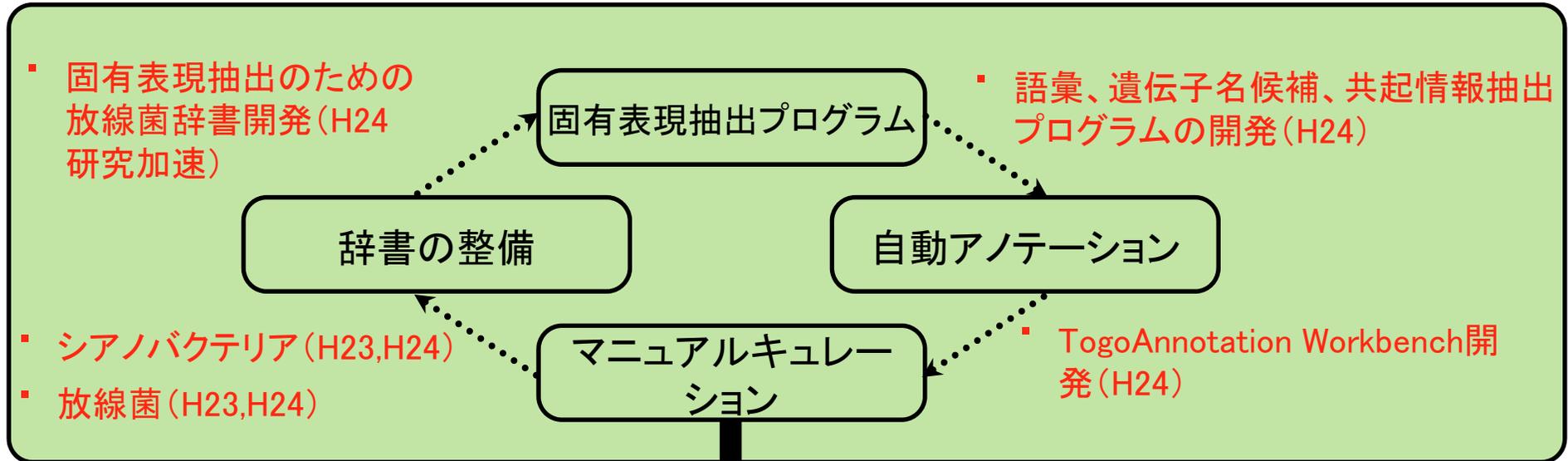
TogoAnnotationへモデル微生物17043遺伝子に対して合計165,527アノテーション(H23-24: 8,869遺伝子40,881アノテーション)のブックマークを集積した。H24において、新規に放線菌リファレンスゲノムである *Streptomyces coelicolor* A3(2) のマニュアルキュレーションを開始した。

藤澤&神沼&中村@遺伝研



©2013黒川 顕(東京工業大学) licensed under CC表示2.1日本

モデル微生物アノテーションリファレンス集積加速のための整備

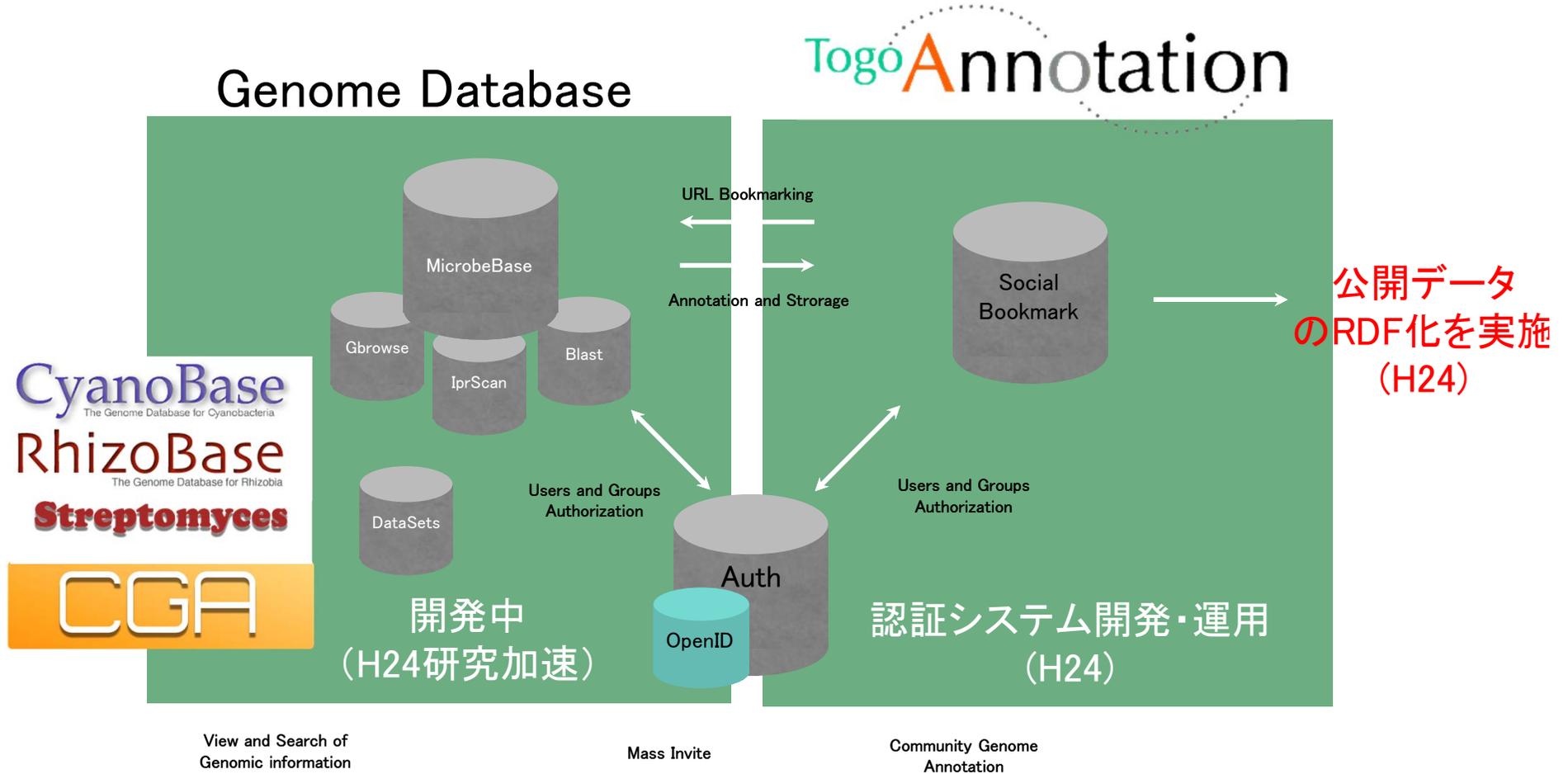


TogoAnnotationへのモデル微生物
アノテーションリファレンス集積

ゲノムデータベースへのマッピング

H24までに各処理系の開発を実施した。H25に実データでの運用を開始する。

ゲノムアノテーションプラットフォームの整備



Community

Bradyrhizobium sp. S23321 (H23論文公開)
シアノバクテリア (H24~)
放線菌

藤澤&神沼&中村@遺伝研



TogoAnnotationキュレーション情報とNGS多型解析情報の統合

放線菌ゲノム情報統合の事例

streptomyces/A3-2 - A3(2):SCO2958 編集

- reference genome
- gene annotation
- pumbed links

Sequence Read Archive + DDBJ pipeline

Example: SRA002840

S.coelicolor A3(2)[REF]	All genes	7824
S.lividans TK24	Unmapped genes	599
	Mapped genes with variants(SNP/INDEL)	5795

S.lividans TK24	SNP	27630
	INS	402
	DEL	377

- target strain
- new variants

ゲノム情報

キュレーション情報

多型解析情報

gene_id	gene_product	gene_symbol	annotation_id	curated_gene_symbol	reference_count	annotation_count	variation_id	variation_type	seqid	location	allele
SCO0613	arginine deiminase		http://togo.annotation.jp/annotations/193815	arcA	1	1	SITK24_000002495	SNP	chr	652849	C
SCO0674	endo-1,4-beta-xylanase		http://togo.annotation.jp/annotations/197580	xysA	1	4	SITK24_000002703	SNP	chr	713507	C
SCO0674	endo-1,4-beta-xylanase		http://togo.annotation.jp/annotations/197580	xysA	1	4	SITK24_000002702	SNP	chr	712924	T
SCO0674	endo-1,4-beta-xylanase		http://togo.annotation.jp/annotations/197580	xysA	1	4	SITK24_000002704	SNP	chr	713822	G
SCO0713	lipase		http://togo.annotation.jp/annotations/191409	lipA	1	1	SITK24_000002874	SNP	chr	756500	G
SCO0713	lipase		http://togo.annotation.jp/annotations/191409	lipA	1	1	SITK24_000002875	SNP	chr	756716	A
SCO0713	lipase		http://togo.annotation.jp/annotations/191409	lipA	1	1	SITK24_000002876	SNP	chr	756797	C
SCO0713	lipase		http://togo.annotation.jp/annotations/191409	lipA	1	1	SITK24_000002877	SNP	chr	757109	C
SCO1483	carbamoyl phosphate synthase large subunit	carB	http://togo.annotation.jp/annotations/193838	pyrA	1	5	SITK24_000005845	SNP	chr	1584905	G
SCO1483	carbamoyl phosphate synthase large subunit	carB	http://togo.annotation.jp/annotations/193838	pyrA	1	5	SITK24_000005847	SNP	chr	1586465	C
SCO1483	carbamoyl phosphate synthase large subunit	carB	http://togo.annotation.jp/annotations/193838	pyrA	1	5	SITK24_000005841	SNP	chr	1584152	A
SCO1483	carbamoyl phosphate synthase large subunit	carB	http://togo.annotation.jp/annotations/193838	pyrA	1	5	SITK24_000005842	SNP	chr	1584191	G
SCO1483	carbamoyl phosphate synthase large subunit	carB	http://togo.annotation.jp/annotations/193838	pyrA	1	5	SITK24_000005843	SNP	chr	1584359	A
SCO1483	carbamoyl phosphate synthase large subunit	carB	http://togo.annotation.jp/annotations/193838	pyrA	1	5	SITK24_000005844	SNP	chr	1584890	A
SCO1483	carbamoyl phosphate synthase large subunit	carB	http://togo.annotation.jp/annotations/193838	pyrA	1	5	SITK24_000005846	SNP	chr	1586213	G
SCO1484	carbamoyl phosphate synthase small subunit		http://togo.annotation.jp/annotations/193839	pyrAa	1	2	SITK24_000005848	SNP	chr	1588303	T
SCO1486	dihydroorotase	pyrC	http://togo.annotation.jp/annotations/193846	pyrC	1	1	SITK24_000005849	SNP	chr	1589698	G
SCO1486	dihydroorotase	pyrC	http://togo.annotation.jp/annotations/193846	pyrC	1	1	SITK24_000005850	SNP	chr	1589974	G
SCO1486	dihydroorotase	pyrC	http://togo.annotation.jp/annotations/193846	pyrC	1	1	SITK24_000005851	SNP	chr	1590196	C
SCO1487	aspartate carbamoyltransferase catalytic subunit	pyrB	http://togo.annotation.jp/annotations/193841	pyrB	1	8	SITK24_000005852	SNP	chr	1590420	G
SCO1487	aspartate carbamoyltransferase catalytic subunit	pyrB	http://togo.annotation.jp/annotations/193841	pyrB	1	8	SITK24_000005853	SNP	chr	1590854	A
SCO1513	GTP pyrophosphokinase		http://togo.annotation.jp/annotations/194993	relA	1	3	SITK24_000005907	SNP	chr	1618230	G
SCO1513	GTP pyrophosphokinase		http://togo.annotation.jp/annotations/194993	relA	1	3	SITK24_000005908	SNP	chr	1618512	C
SCO1513	GTP pyrophosphokinase		http://togo.annotation.jp/annotations/194993	relA	1	3	SITK24_000005909	SNP	chr	1618563	G



オミックスデータとゲノム情報との統合イメージ

研究コミュニティ提供の *S. griseus* オミックスデータ

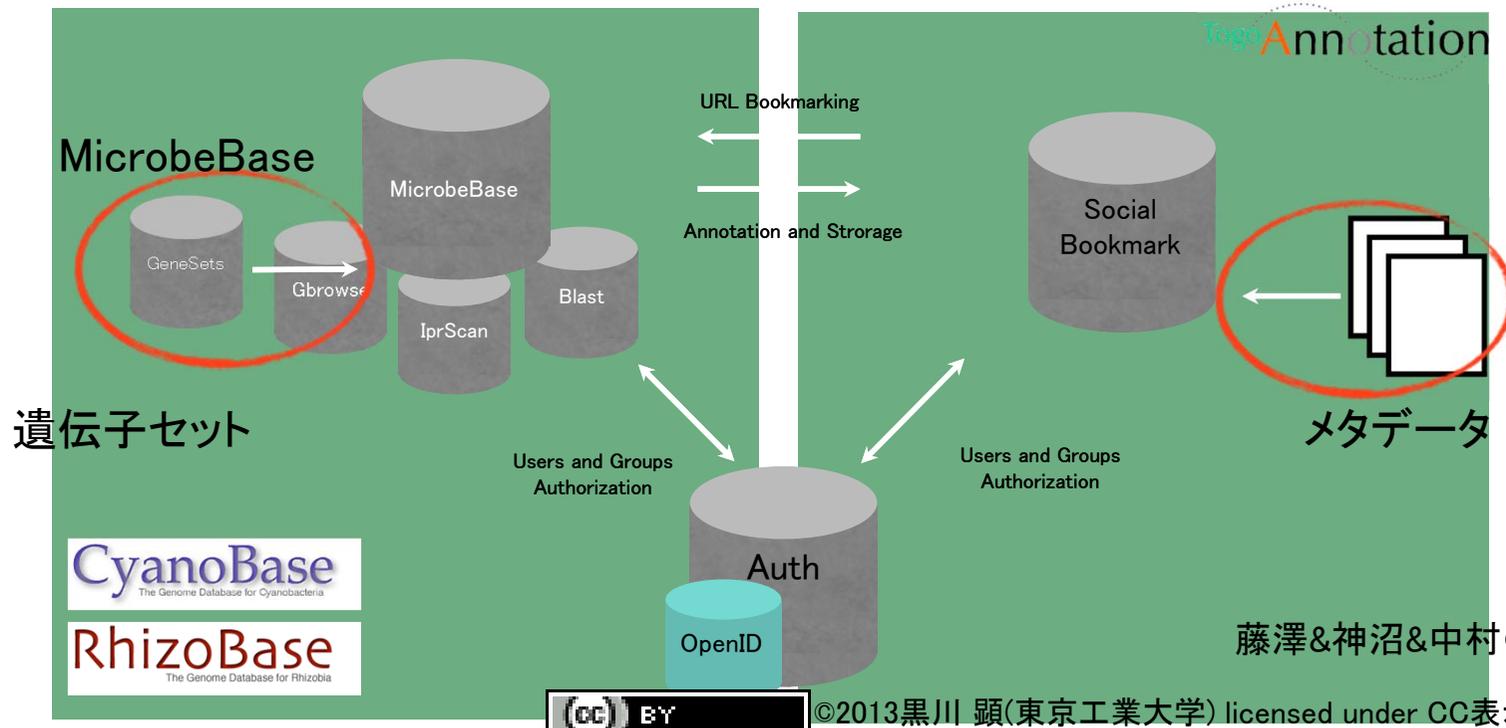
H24研究加速により実施

マイクロアレイ	Gbrowse
ChIP-seq解析	新規
プロテオーム解析 (MS)	新規
プロテオーム解析 (二次元電気泳動)	KazusaWiki
RNA-seq	新規
転写解析	新規
遺伝子破壊株	TogoAnnotation
small RNA	Gbrowse
ゲルシフトアッセイとフットプリント	新規
文献情報(遺伝子グループ)	新規

-
- ・新規データの格納方法
 - ・サンプル情報などのメタデータ
 - ・遺伝子セットデータ (.gmt形式)



ゲノムデータベースへのマッピングを可能に



藤澤&神沼&中村@遺伝研



©2013黒川 顕(東京工業大学) licensed under CC表示2.1日本

メタゲノムデータのRDF化および GTPS、RefSeq等との統合

メタゲノムメタデータの集計結果

	サンプル数	メタデータの カテゴリ数	メタデータカテゴリーの例
ヒト共生細菌群集	72,236 (18,224)	85	Age , Sex ,Disease stage , Country , Body Habitat , Diet 等
環境細菌群集	6,356 (5,827)	627	pH , Temperature , Wind Speed , Dissolved Oxygen 等

()内は実際に配列が存在したエントリ (2012年6月5日時点)

	Age	Body Habitat	Body Site	Collection Date	Country	Disease Stage
Sample 1	22	Feces		2008		Obese
Sample 2					Japan	
Sample 3			Scalp			
Sample 4			Skin		USA	
Sample 5	1years		Gut	2011/8/8		Healthy

登録されているカテゴリーや値の語彙は登録者によってバラバラ



微生物の生息環境についてのオントロジー-MEOを利用して、メタデータの記述を行った

28



©2013黒川 顕(東京工業大学) licensed under CC表示2.1日本

配列取得可能なSRS IDの内訳

ヒト共生細菌群集 (計 18,224)

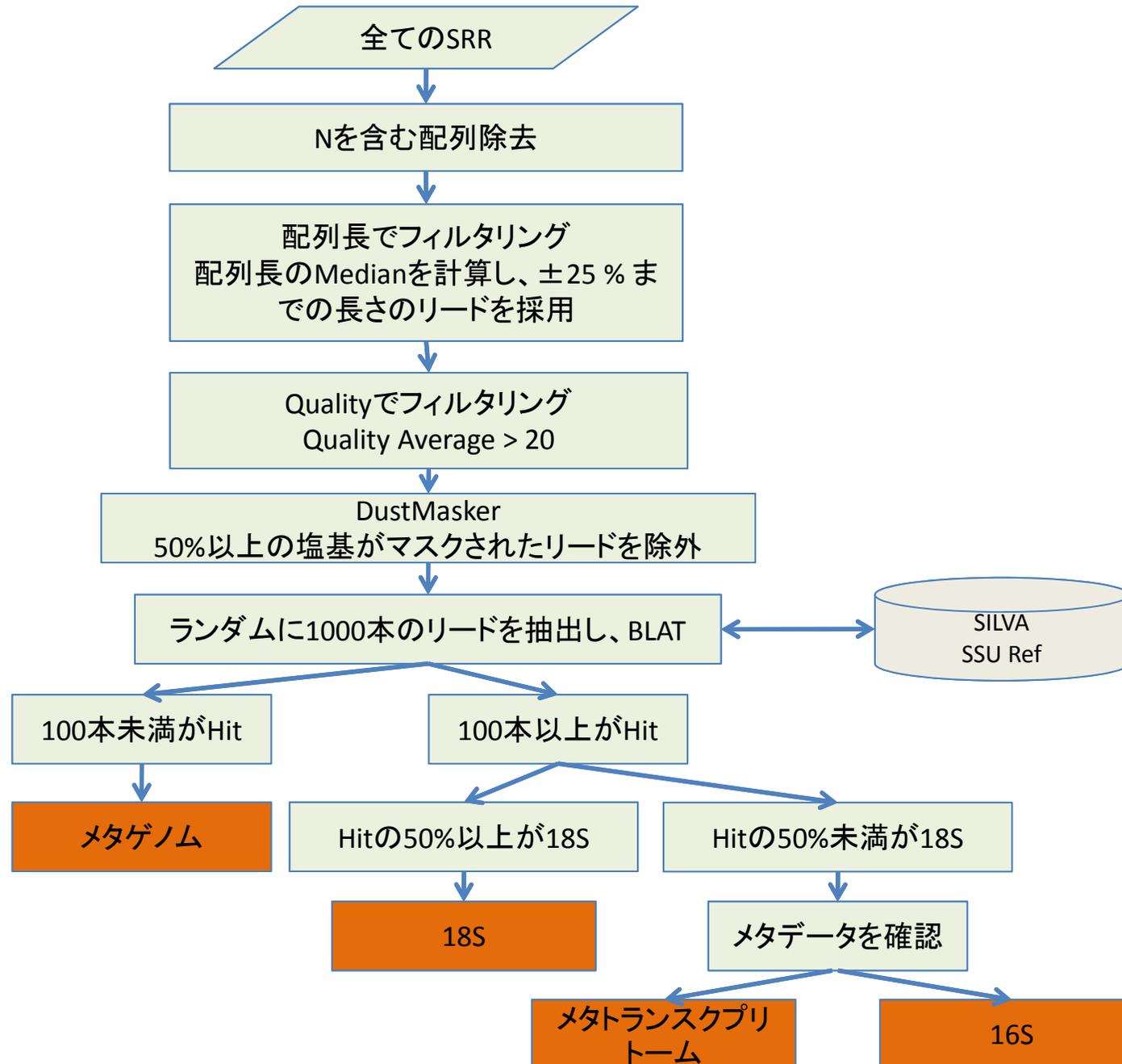
サンプルの種類	SRS IDの数
human metagenome	16,185
human gut metagenome	766
human oral metagenome	504
Homo sapiens	405
human skin metagenome	363
human lung metagenome	1
合計	18,224

配列取得可能なSRS IDの内訳

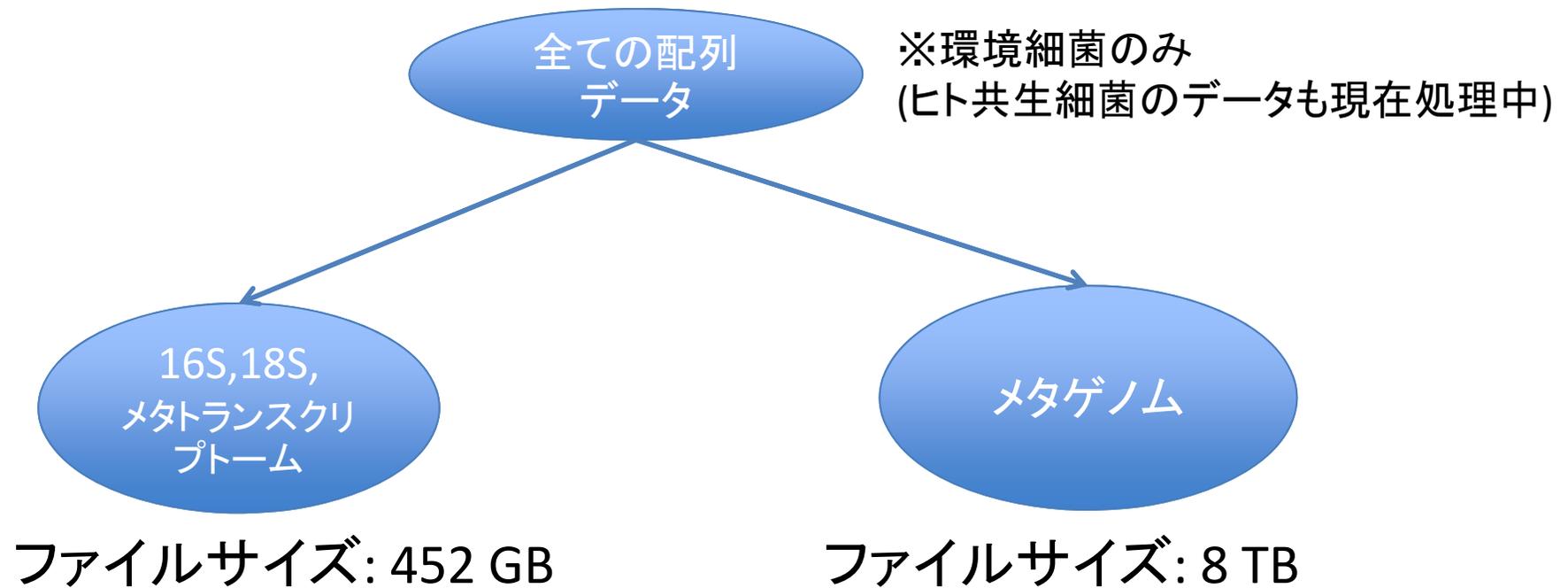
環境細菌群集 (計 5,827)

サンプルの種類	SRS IDの数
marine metagenome	1,212
mouse gut metagenome	985
hydrocarbon metagenome	480
soil metagenome	469
freshwater metagenome	419
gut metagenome	395
organismal metagenomes	210
metagenome sequence	198
bioreactor sludge metagenome	132
marine sediment metagenome	98
sediment metagenome	79
:	:
合計	5,827

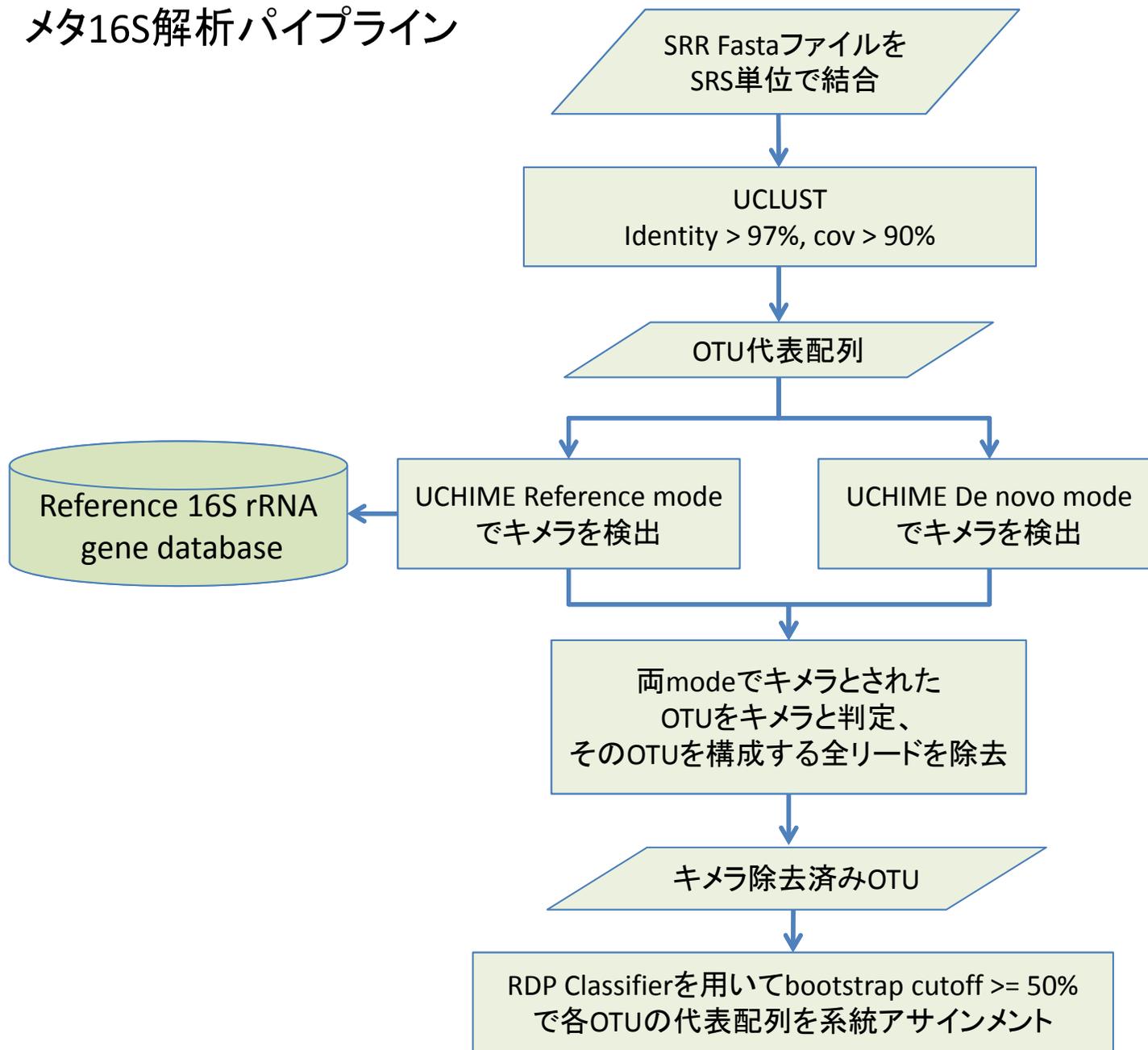
配列の前処理 (目的: メタゲノムとメタ16Sデータの区別 & 高精度配列データの抽出)



前処理によって区別された配列データの内訳



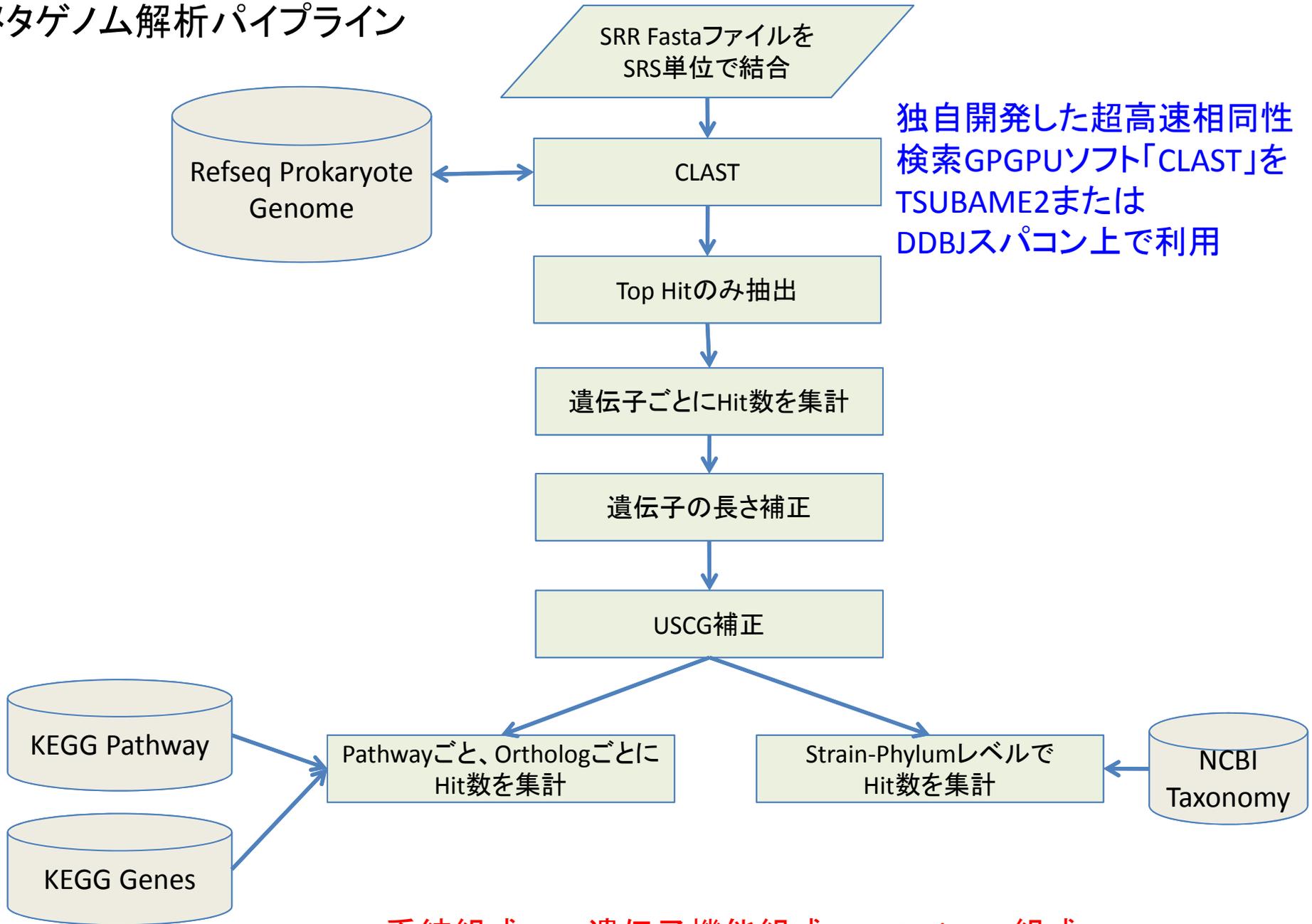
メタ16S解析パイプライン



Genus-Phylumレベルの系統組成



メタゲノム解析パイプライン



系統組成

遺伝子機能組成

Pathway組成

文献からの環境メタデータの抽出

SRAへの配列登録時、大部分のメタデータの記述は任意であるため、論文には記述されているが、SRAのXMLには記述されていないメタデータが大量に存在



生息環境の情報と微生物の遺伝子および系統組成の情報を統合して解析し、より多くの知見を得るためには、論文に記述されているメタデータの追加が必要



まず、論文のpubmed IDとSRS IDの対応付けを行う必要がある

2つの方法でpubmed IDとSRS IDの対応付けを行った

- ①SRAs (<http://sra.dbcls.jp/cgi-bin/publication.cgi>)の利用
- ②NCBIのftpサイトからダウンロードしたSRAのXMLファイルの利用

文献からの環境メタデータの抽出

環境細菌群集のメタゲノム・メタ16S解析についての**113**報の論文を取得
SRAの**2,286**/5,827サンプル (39.2%)と紐付けられた。

取得した論文の内訳

	論文数
marine metagenome	29
soil metagenome	11
freshwater metagenome	6
mouse gut metagenome	5
sediment metagenome	4
marine sediment metagenome	4
:	:

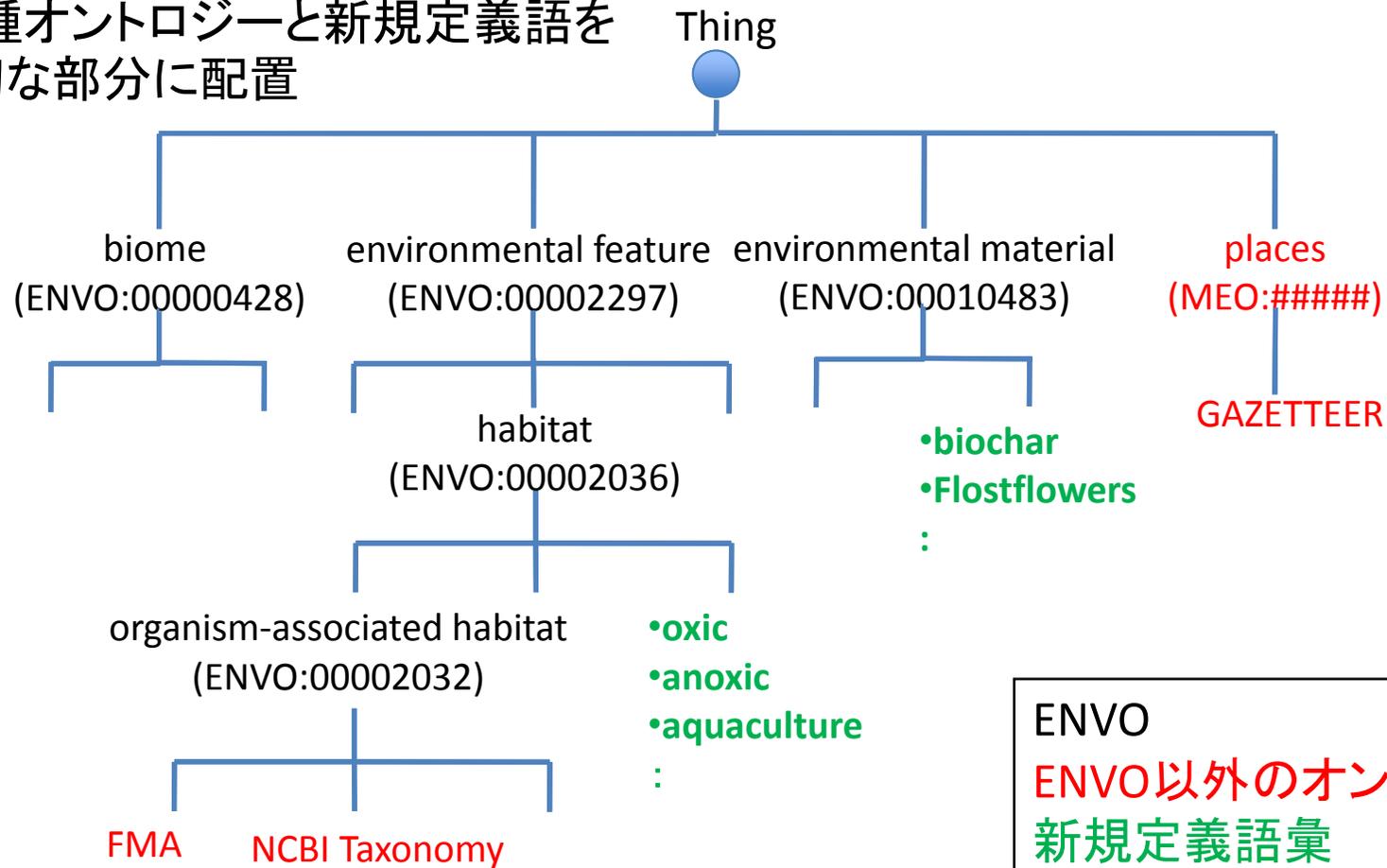
ヒト共生細菌群集のメタゲノム・メタ16S解析についての**30**報の論文を取得
SRAの**12,365**/18,224サンプル (67.9%)と紐付けられた。

各種オントロジーの開発

Metagenome/Microbes Environmental Ontology (MEO)

微生物の生息環境メタデータのオントロジー

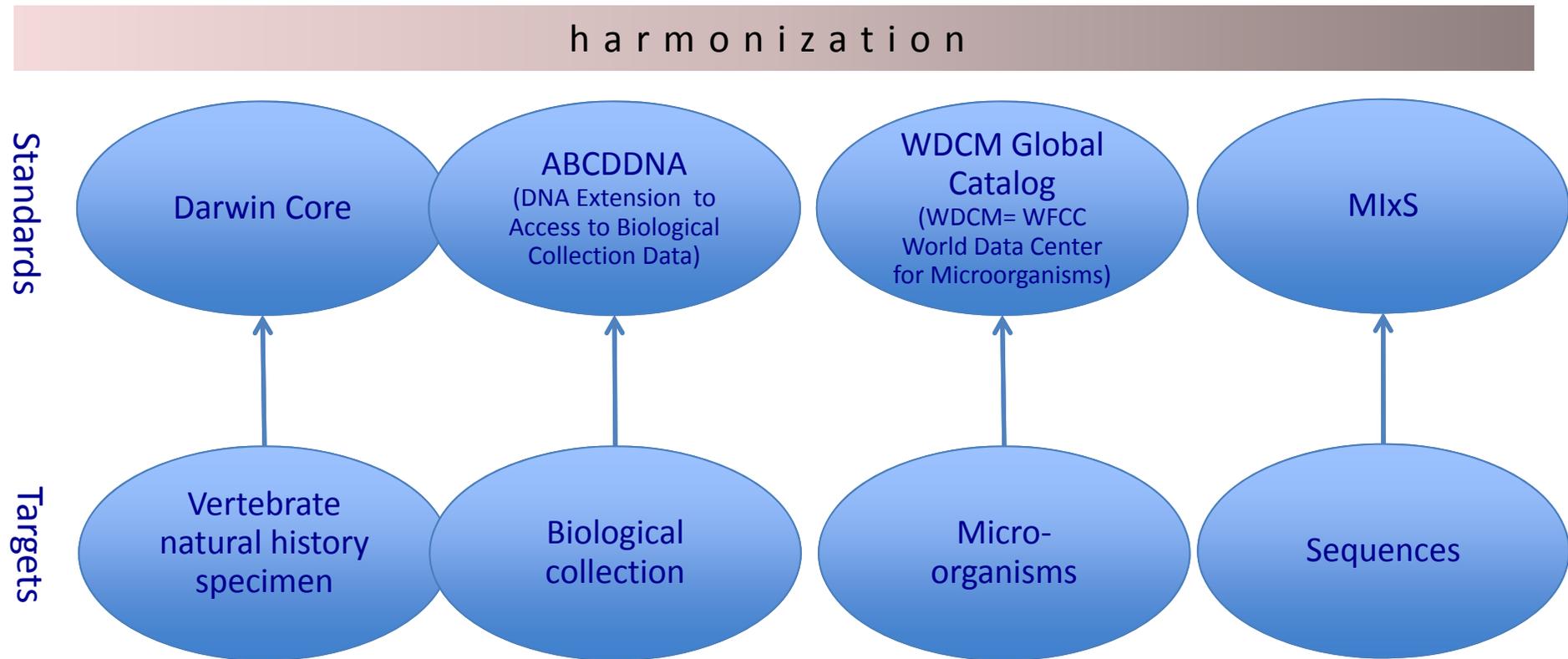
- ENVOを基礎とした構造
- 各種オントロジーと新規定義語を適切な部分に配置



BioPortal <http://bioportal.bioontology.org/> 及び
プロジェクトページ <http://mdb.bio.titech.ac.jp/meo/> でMEOを公開



Further Extension of GBIF standards



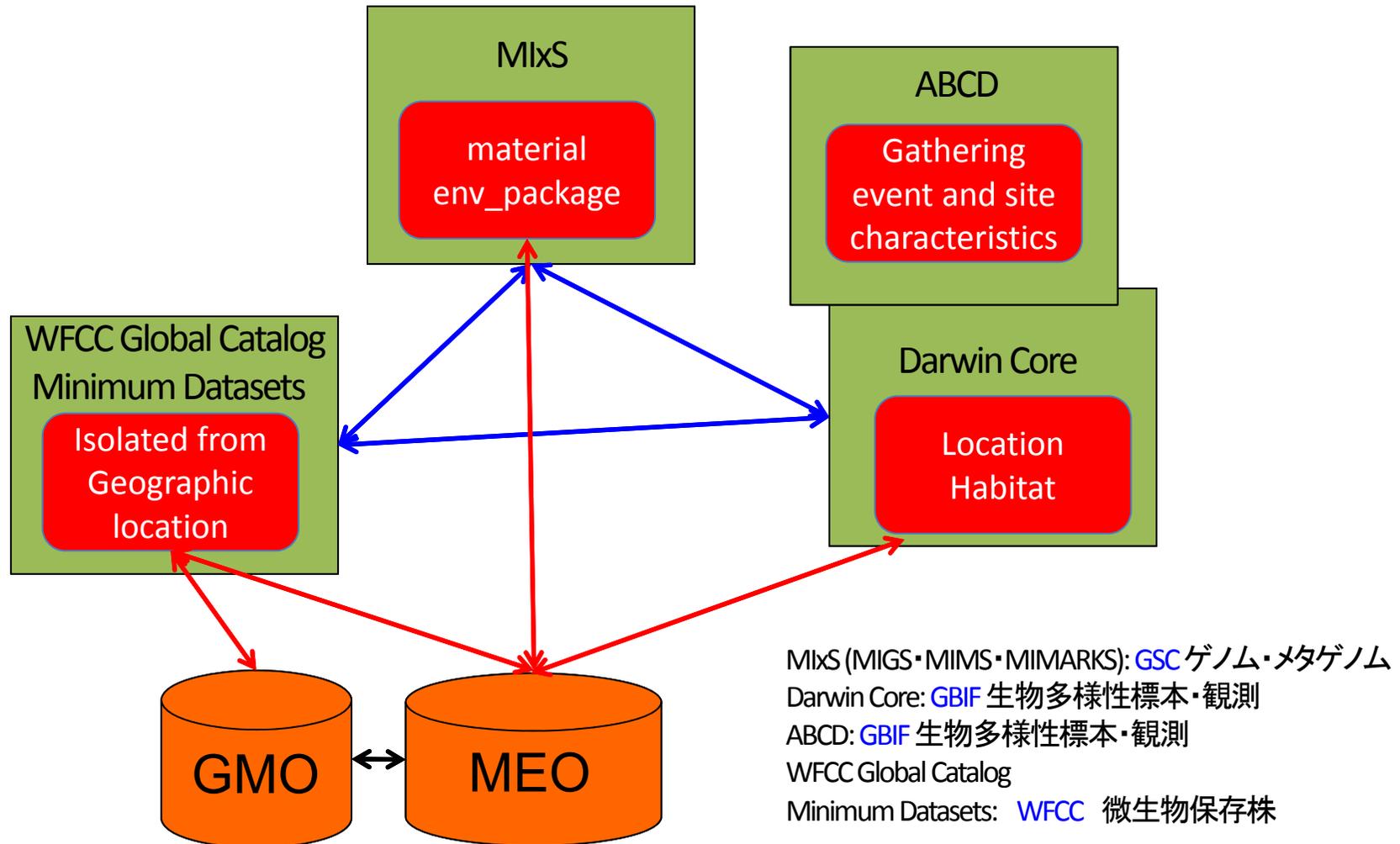
Sequences, sequences and sequences, e.g.,

--- Each GridION node and cartridge is initially designed to deliver tens of Gb of sequence data per 24 hour period, - snip -- Nanopore technology and the GridION system are well suited to miniaturisation of the sensing device. The MinION is a small., portable, single-molecule analysis device that is compatible with Oxford Nanopore's DNA sequencing and protein analysis techniques.--
 --- Oxford Nanopore intends to commercialise GridION and MinION directly to customers within 2012 ---



生物資料のゲノム情報 = その生物自体のゲノム + 共生細菌のメタゲノム

国際データ標準化グループとの連携



DwC* - MlXs** alignment workshop

27 – 29 February 2012

Oxford e-Research Centre, 7 Keble Road, Oxford, OX1 3QG

*DwC = Darwin Core

**MlXs = Minimum Information of {Genome | Metagenome | Marker gene} Sequence





Standards in Genomic Sciences

An Open Access Journal of the Genomic Standards Consortium

Standards in Genomic Sciences (2012) 7:166-170

DOI:10.4056/sigs.3166513

Meeting Report: Hackathon-Workshop on Darwin Core and MIxS Standards Alignment (February 2012)

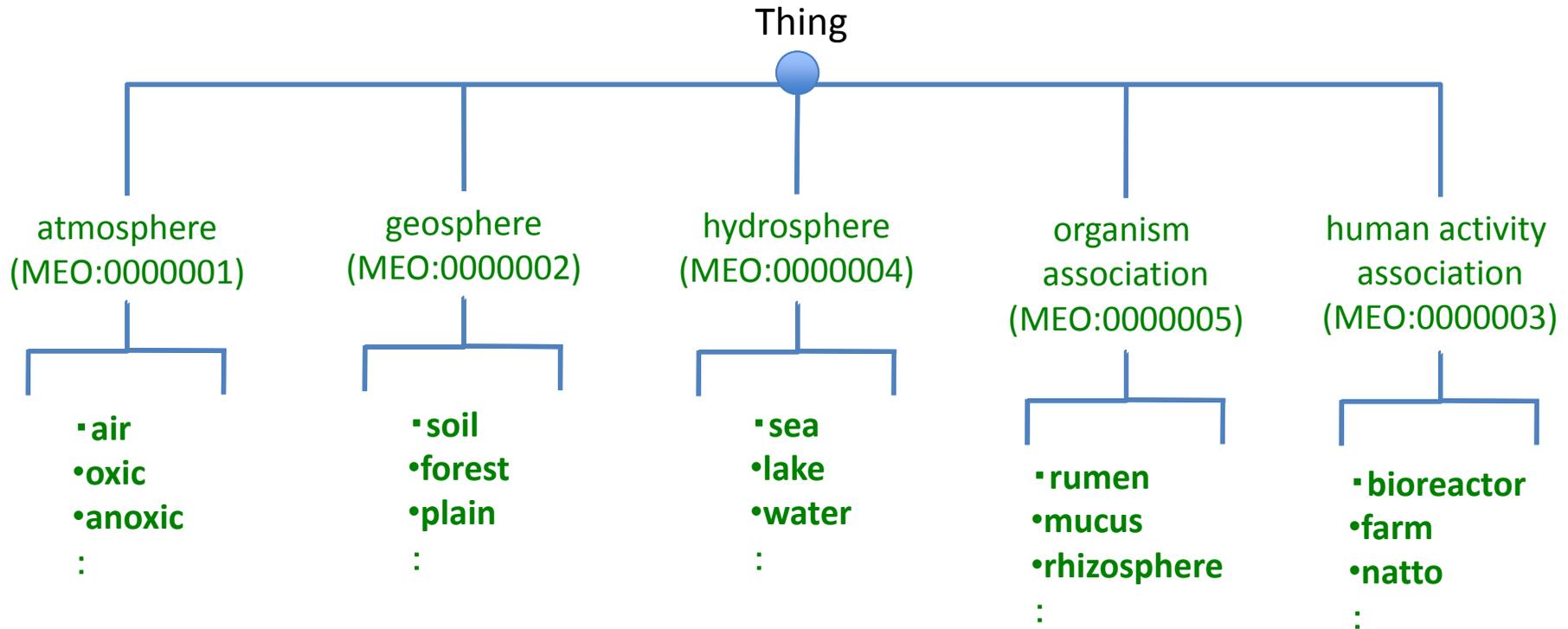
Éamonn Ó Tuama¹, John Deck², Gabriel Dröge³, Markus Döring⁴, Dawn Field⁵, Renzo Kottmann⁶, Juncai Ma⁷, Hiroshi Mori⁸, Norman Morrison^{9,10}, Peter Sterk¹¹, Hideaki Sugawara¹², John Wiczorek¹³, Linhuan Wu¹⁴, Pelin Yilmaz¹⁵

5. The MicrobeDB.jp project including MEO (http://mdb.bio.titech.ac.jp/meo/about_meo) coordinates microbial genomic and metagenomic data, and is supported by NBDC (National Bioscience Database Center, <http://biosciencedbc.jp/nbdc.cgi?lng=en>), Japan.

<http://standardsingenomics.org/index.php/sigen/article/view/sigs.3166513/779>

Metagenome/Microbes Environmental Ontology (MEO) Ver. 0.3

MEO is an application ontology for MicrobeDB.jp, not a domain ontology (BioHackathon 2012)



177 terms

MEO version 0.3 contains terms of GOLD complete genome metadata and SRA metagenome metadata by manually picking up terms related to microbial environment

微生物が引き起こす感染症・症状についての オントロジー PDO-SYMP の構築

感染症：病原体（主に微生物）の侵入により引き起こされる病気

既知のゲノム・メタゲノムデータには感染症関連のデータが多数存在



- 原因不明疾患の原因菌推定、微生物・ヒトor動植物間相互作用の解明等の目的にMicrobeDB.jpを利用できるようにしたい
- メタゲノムデータにも対応させたい

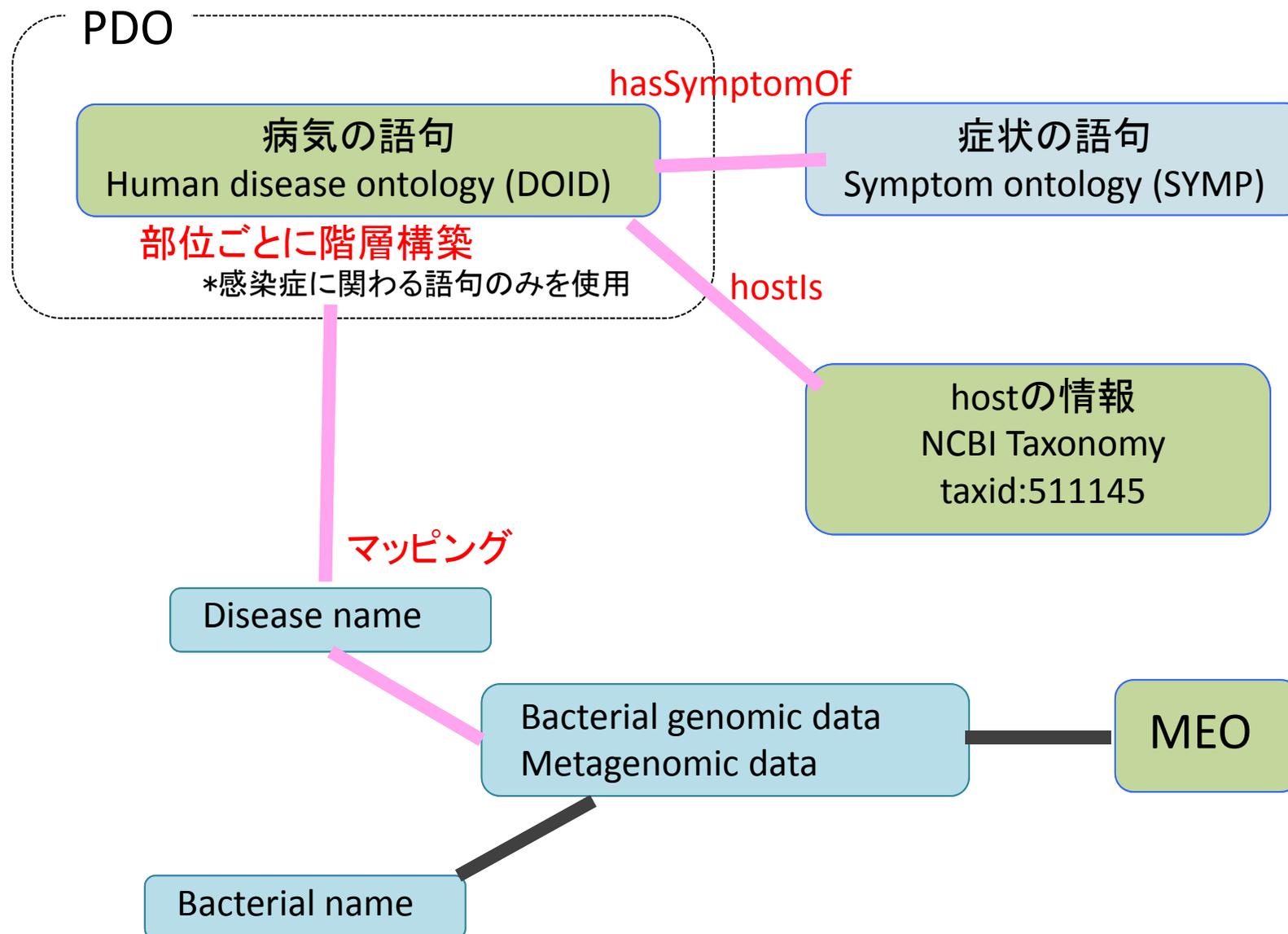


MicrobeDB.jpを感染症研究に活用可能にするためには...

感染症に関連する語句（原因菌、病名、症状）
の意味的関連性を明確にする→**オントロジーの構築**

オントロジー構造(全体)

名称: Pathogenic Disease Ontology (PDO)



ゲノムアノテーション標準のためのINSDCオントロジー開発

- ・ INSDC配列エントリー (GenBank, ENA, DDBJ 形式) のRDF化を目的としたオントロジー開発
- ・ 国内版BH12.12においてDBCLSと連携してINSDC Feature Table Definition (version 10.2)をOWLで定義
- ・ Genomic Standards Consortium (GSC15; April 22–24, 2013. NIH, USA) にて発表予定

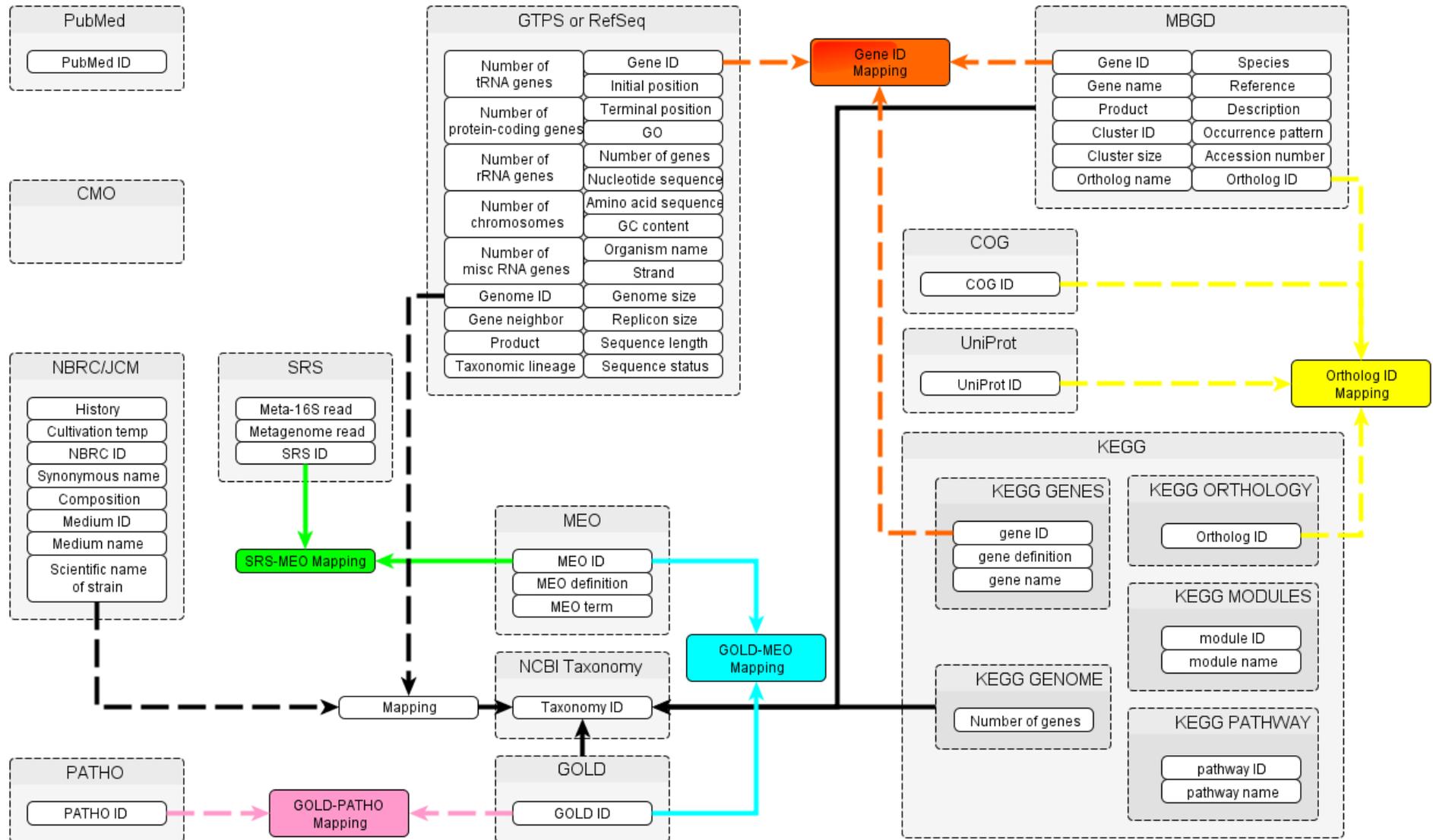
藤澤&神沼&中村@遺伝研

構築したオントロジー

- FALDO (Feature Annotation Location Description Ontology)
 - ゲノム中の各featureの位置情報を記述するためのオントロジー(DBCLS)
- INSDC Ontology
 - INSDCエントリのfeatureとqualifierのターム記述のためのオントロジー
- MCCV (Microbial Culture Collection Vocabulary)
 - 菌株データを記述するためのオントロジー
- MEO (Metagenome/Microbe Environmental Ontology)
 - 細菌の生息環境を記述するためのオントロジー
- PDO-SYMP (Pathogenic Disease Ontology with Symptom)
 - 細菌が引き起こす感染症の情報および感染症の症状を連結したオントロジー
- GMO (Growth Media Ontology)
 - 細菌の培地情報を記述するためのオントロジー(DBCLS)

MicrobeDB.jpにおけるDB間のID統合 およびStanzaの開発

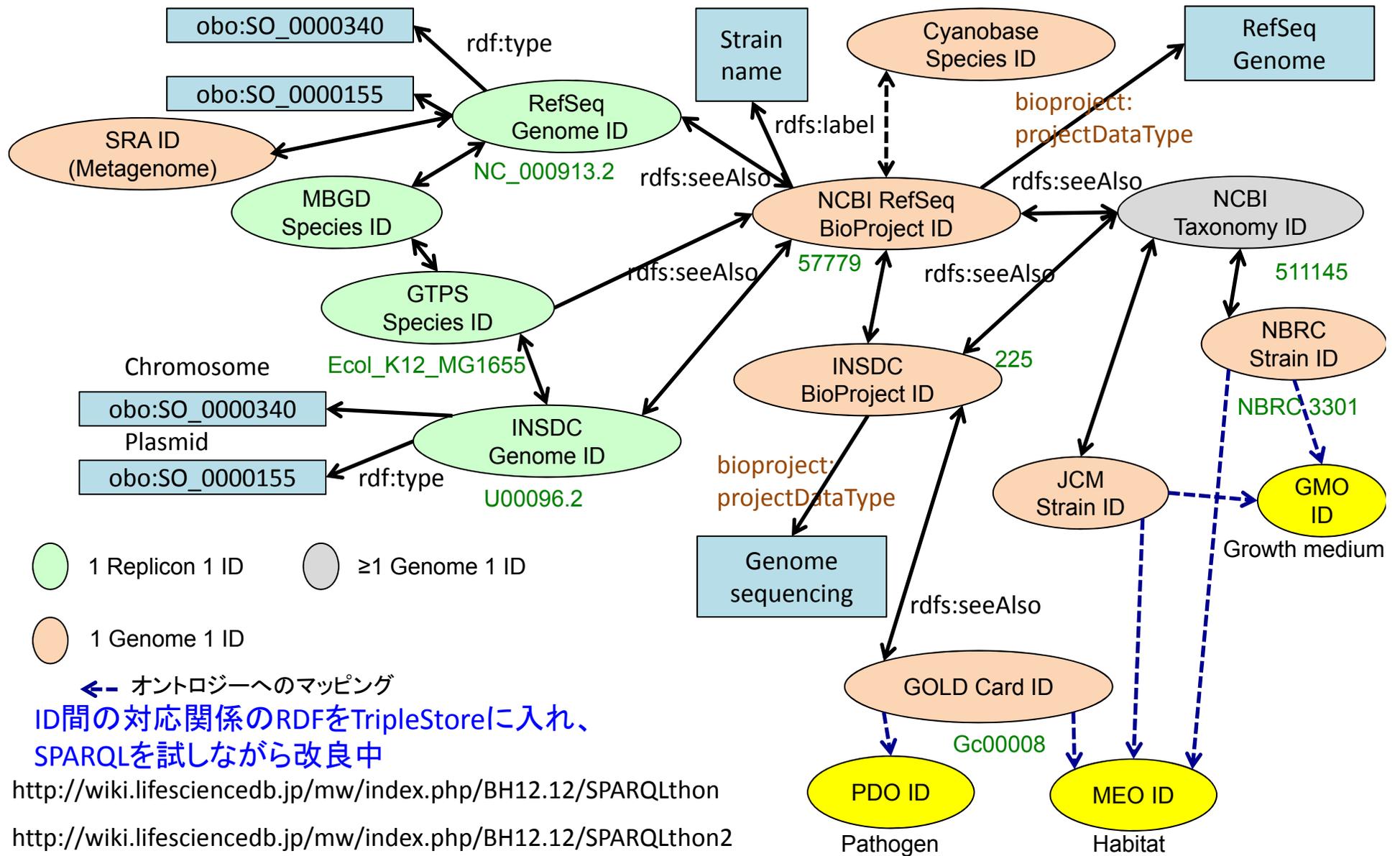
単一DB or 複数DB由来のStanza



複数のDB由来のデータからStanzaを作るには、DB間のID対応関係を整理する必要がある

Integration of DB resources in MicrobeDB.jp

Genomes (example: *Escherichia coli* str. K-12 substr. MG1655)



○ 1 Replicon 1 ID ○ ≥1 Genome 1 ID

○ 1 Genome 1 ID

←- オントロジーへのマッピング

ID間の対応関係のRDFをTripleStoreに入れ、
SPARQLを試しながら改良中

<http://wiki.lifesciencedb.jp/mw/index.php/BH12.12/SPARQLthon>

<http://wiki.lifesciencedb.jp/mw/index.php/BH12.12/SPARQLthon2>

<http://wiki.lifesciencedb.jp/mw/index.php/BH12.12/SPARQLthon3>



Stanza: data templates for a genome DB

- reusable
- high modularity
- easy to integrate

- TogoStanazaの仕様をDBCLSと共に策定
- MicrobeDBとCyanobaseをユースケースとして開発
- 重複開発を避ける必要あり
- H24年度までに約15種類のStanzaを開発(ゲノム基盤部分はDBCLSが開発)

The image displays a collection of data visualization templates for a genome database, organized into four main categories: **All**, **Gene**, **Taxon**, and **Environment**.

- All:** Includes a comparative heatmap, an ontology hierarchical viewer, a numerical data graph, and a references table.
- Gene:** Features a gene general overview, module general overview, nucleotide attributes, ortholog general overview, pathway general overview, protein attributes, genome browser, genome general overview, and replicon general overview.
- Taxon:** Includes taxon general overview, pathogen visualizer, taxonomic composition (with pie charts), and growth medium overview.
- Environment:** Features environmental distribution (with pie charts and maps), environment general overview, and a sample list.

Each template typically contains a table of data and a corresponding visualization (e.g., heatmap, pie chart, map, or hierarchical tree).



Stanzaの例 (遺伝子・ゲノム)

- Gene Annotation Stanza

Feature

[http://purl.obolibrary.org/obo/SO_0000704]

dbxref	http://www.ncbi.nlm.nih.gov/gene/897644
feature_gene	polC
feature_locus_tag	TM0576
location	605923..610026
isPartOf	http://genome.db/uuid/b4d48cd7-00ef-4e03-9adb-fda7de39e078
type	http://purl.obolibrary.org/obo/SO_0000704
label	TM0576

[http://purl.obolibrary.org/obo/SO_0000316]

dbxref	http://www.ncbi.nlm.nih.gov/gene/897644
dbxref	http://www.ncbi.nlm.nih.gov/nuccore/15643342
exons	nodeID://b71582

- Ortholog list Stanza

Ortholog

ID	Genome	Description	Protein	UniProt	GTFS	RefSeq
aac:AACT_1427	aac	DNA polymerase III subunit alpha	YP_003184842.1	C8WW12	AACT_ACIDOCALDARIUSDSM446:ST2344	NC_013205.1
aar:ACEAR_1599	aar	DNA polymerase III catalytic subunit, PolC type	YP_003828170.1		AARA_DSM5501:ST105	
acl:ACL_0247	acl	DNA polymerase III subunit alpha	YP_001620249.1	A9NEU3	ALAI_PG8A:ST588	NC_010163.1
af:AFLY_1700	afI	DNA polymerase III PolC	YP_002316046.1	B7GG80	AFLA_WK1:ST2505	NC_011567.1
afn:ACFER_1370	afn	DNA polymerase III subunit alpha	YP_003399045.1		AFER_DSM20731:ST1519	NC_013740.1
amt:AMET_2678	amt	DNA polymerase III subunit alpha	YP_001320489.1	A6TRL2	AMET_QYMF:ST2214	NC_009633.1

- Genome Information Stanza

Genome

length	1860725
location	1..1860725
molecularType	genomic DNA
organism	Thermotoga maritima MSB8
sequence	http://genome.db/uuid/b4d48cd7-00ef-4e03-9adb-fda7de39e078.fasta
start	1
stop	1860725
strain	MSB8
version	NC_000853.1
modified	2012-02-13
type	http://purl.obolibrary.org/obo/SO_0000340
type	http://purl.obolibrary.org/obo/SO_0000988
comment	Thermotoga maritima MSB8 chromosome, complete genome.

Stanzaの例 (系統・菌株)

- Taxonomic Hierarchy Stanza

Taxonomy Tree

superkingdom	Bacteria
phylum	Firmicutes
class	Bacilli
order	Lactobacillales
family	Enterococcaceae
genus	Enterococcus
species	Enterococcus faecalis

- Strain Metadata Stanza

Metadata

Medium	http://www.nbrc.nite.go.jp/NBRC2/NBRCMediumDetailServlet?NO=227
Strain number	NBRC 12841
Application	Thienamycins production ; Vitamin B12 (Cyanocobalamine) production ; Steroid conversion
Isolated from	Soil
Strain name	Streptomyces griseus subsp. griseus (Krainsky 1914) Waksman and Henrici 1948
History of deposit	IFO 12841 <-- SAJ <-- OWU (ISP 5226) <-- Squibb & Sons (F. Arnow, MD 2428, ETH 24234, NIHJ 501)
Taxonomy	http://purl.uniprot.org/taxonomy/67263
Temperature for growth	28

- Related Strain in Other Culture Collection Stanza

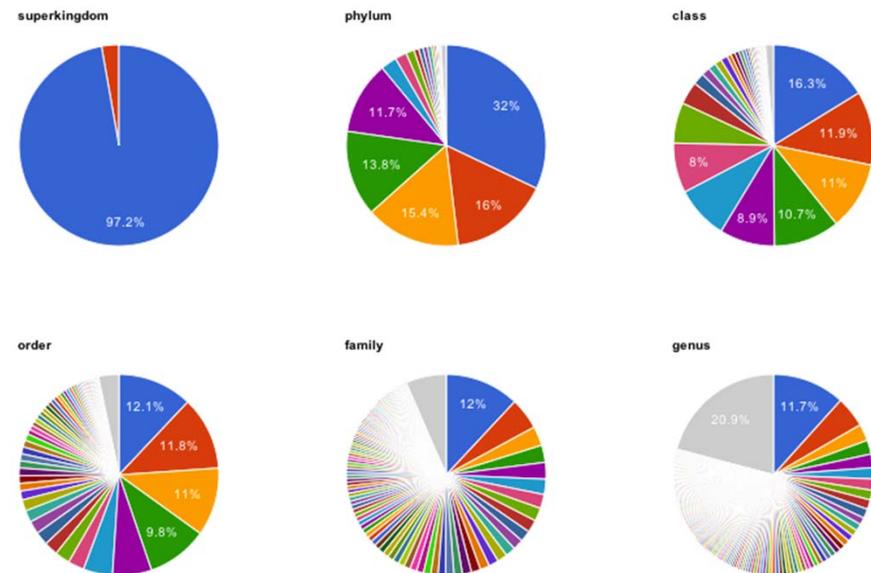
Other Collection Numbers

AS 4.1693
ATCC 11009
ATCC 23882
BCRC 11815
CBS 662.68
DSM 40226
ISP 5226
JCM 4229
JCM 4623
KCTC 1742
LMG 5967
NCIMB 9625
NRRL B-1806

Stanzaの例 (環境)

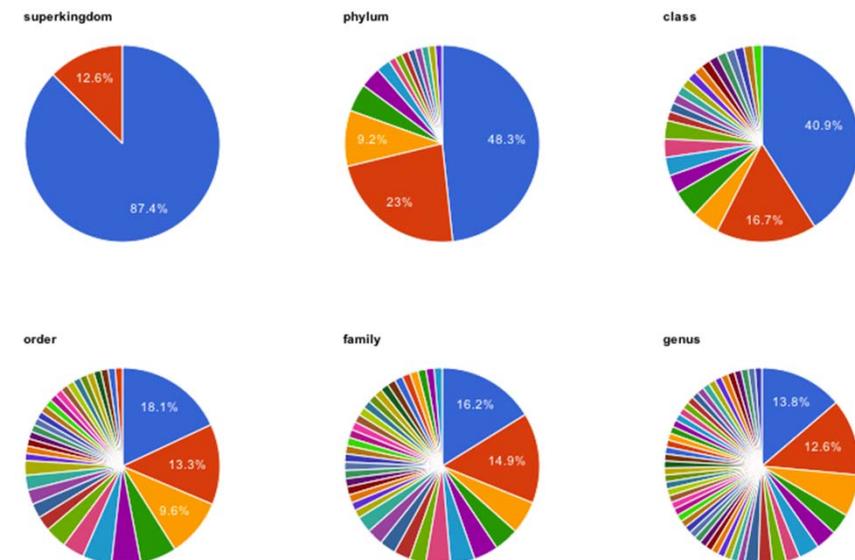
- Taxonomic Composition of the Environment from the Metagenome Stanza

Taxonomy Composition



- Taxonomic Composition of the Environment from the Genome-sequenced Strains Stanza

Taxonomy Composition via GOLD



DBCLSが開発中のTogoStanzaと仕様をあわせて特に微生物特有のStanzaを中心に開発中



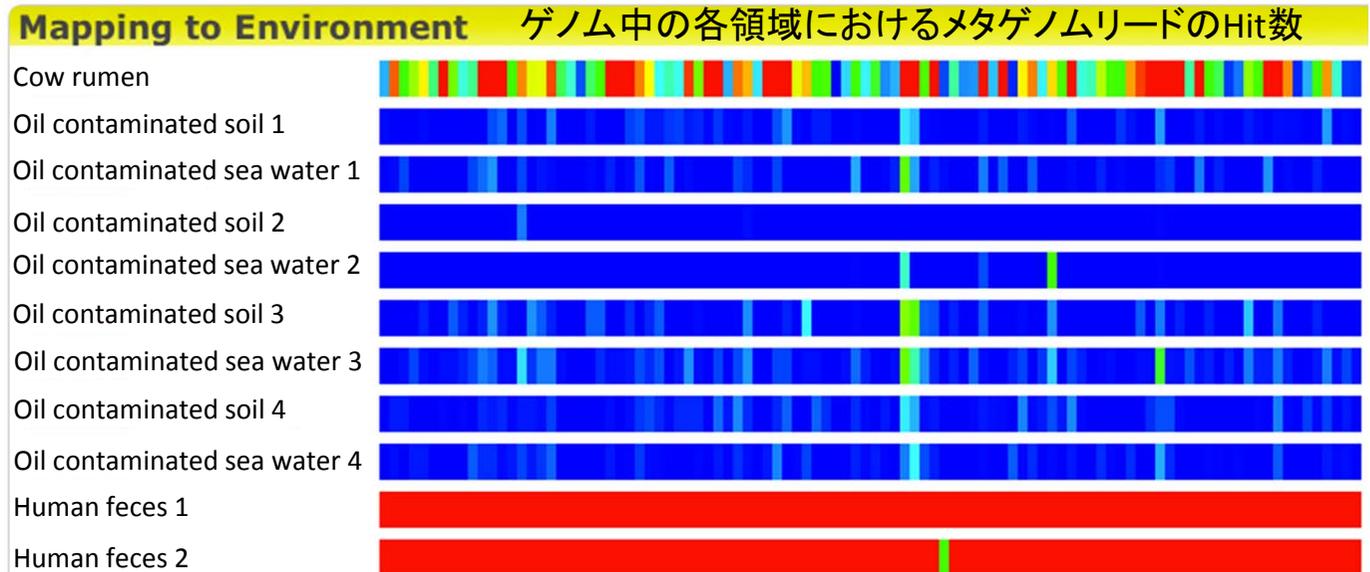
©2013黒川 顕(東京工業大学) licensed under CC表示2.1日本

Stanzaの例 (環境)

- Gene Abundance of the Genome from Several Metagenomes Stanza

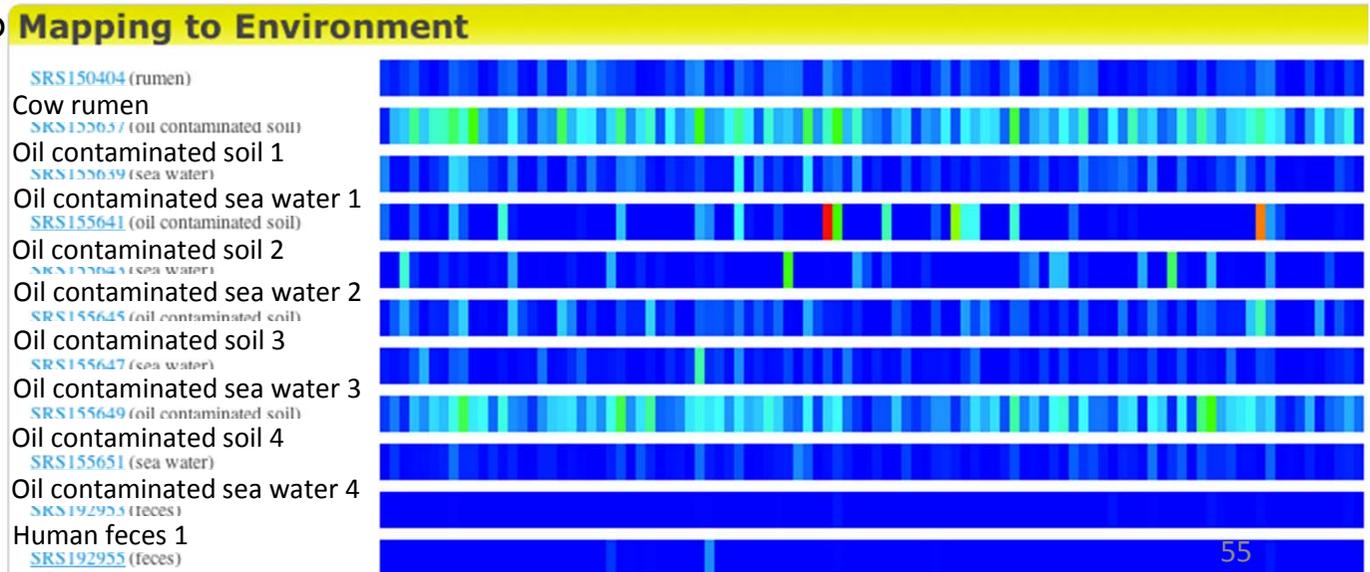
例1: *Bacteroides thetaiotamicron* VPI-5482

ヒト腸内における優占種



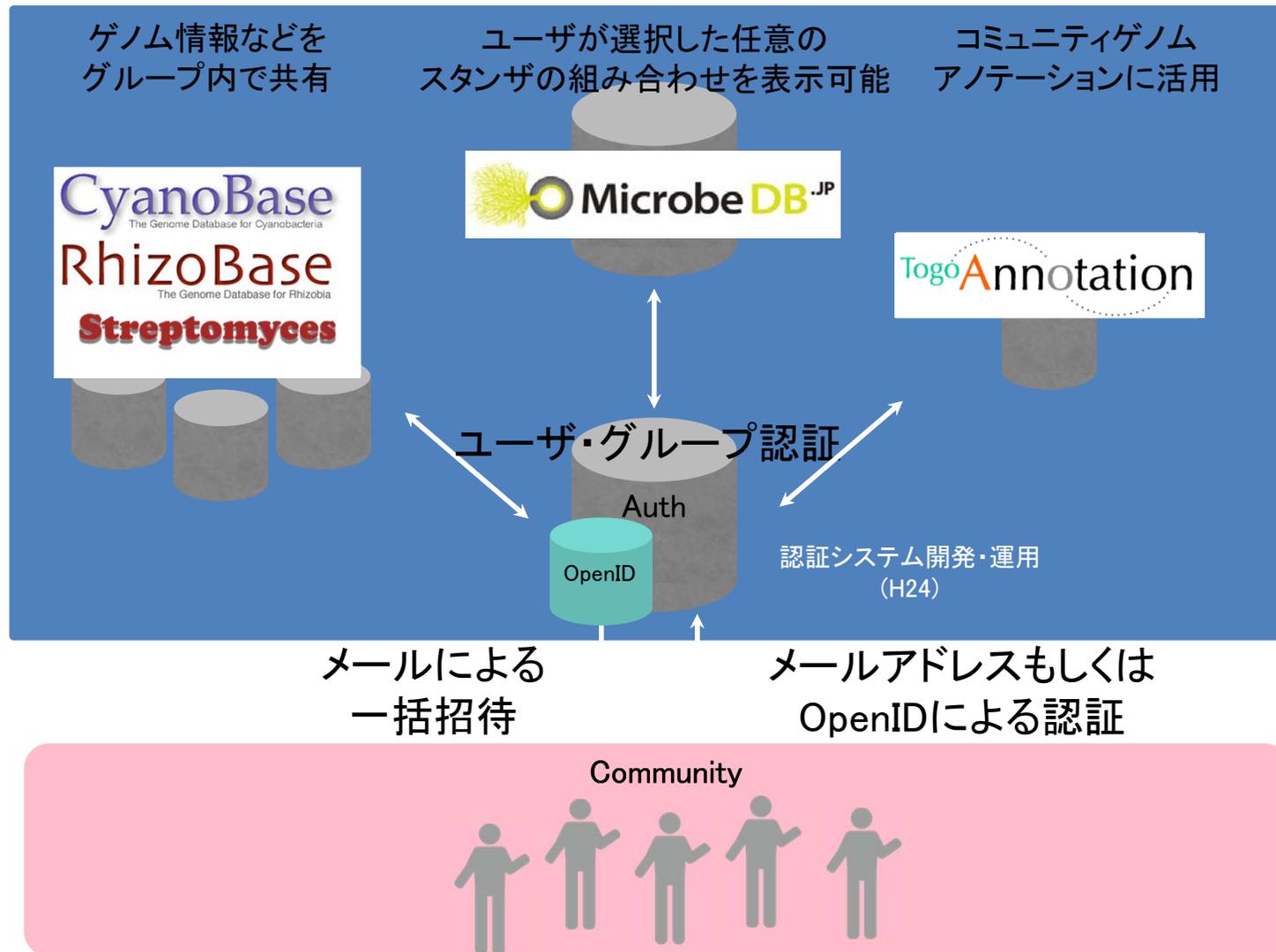
例2: *Burkholderia multivorans* ATCC 17616

芳香族分解関連遺伝子を多数持つ



OpenID-TAuthによる MicrobeDB.jpのユーザ認証

データベース連携に向けたユーザ認証サービス開発・運用



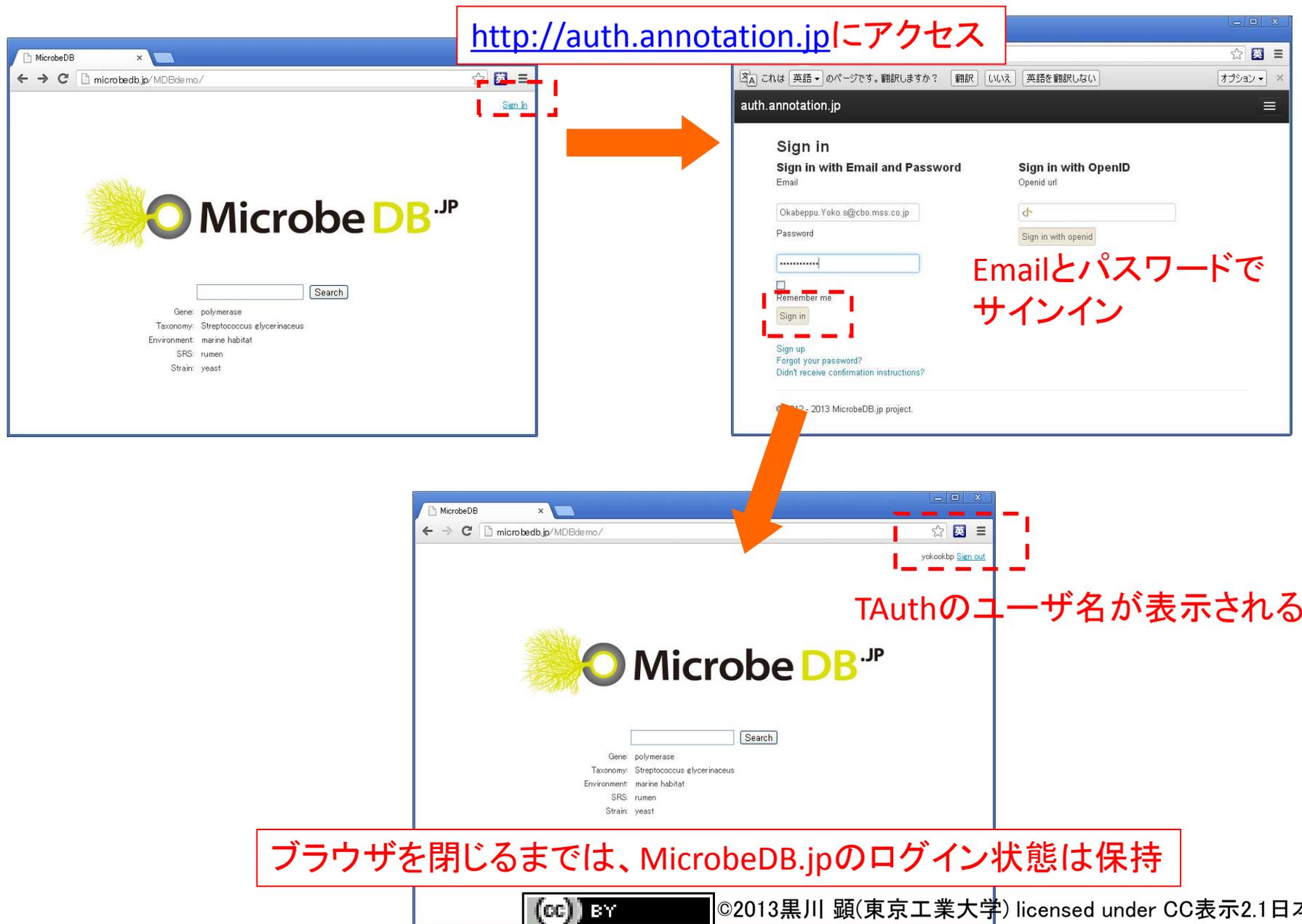
H24年度にOpenID認証と連携したグループ情報を含むユーザ認証サービスを開発し運用を開始した。この機能により、同一認証サービスを利用したMicrobeDB.jpプロジェクト運用データベース間でデータ共有等の連携が可能となった。

藤澤&神沼&中村@遺伝研

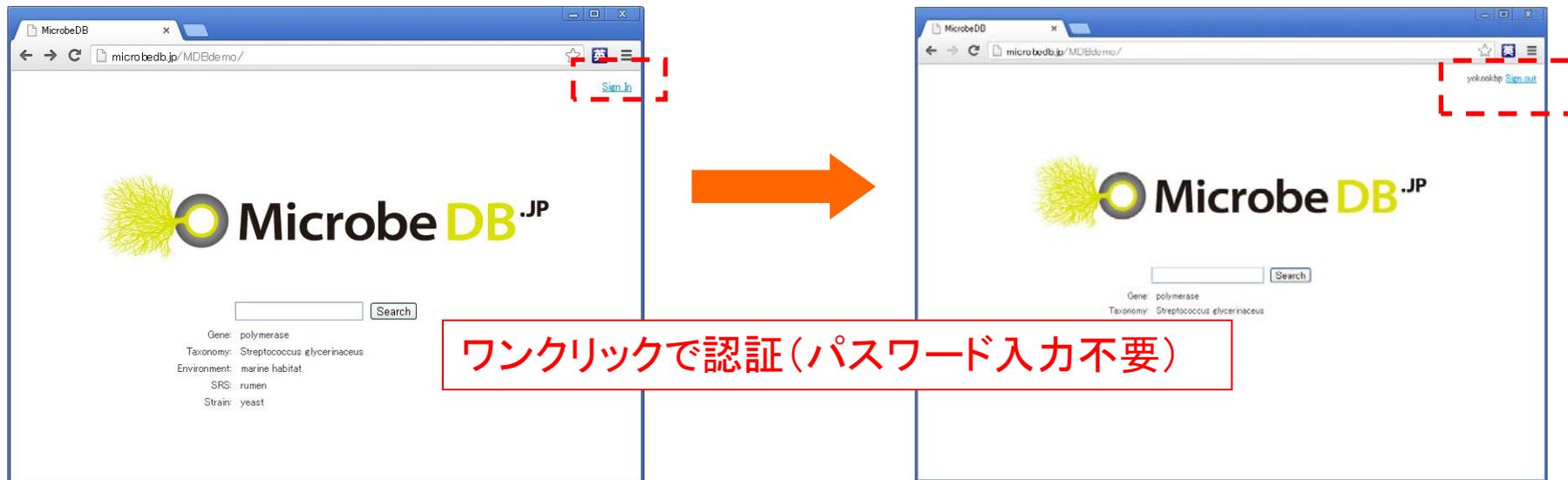


©2013黒川 顕(東京工業大学) licensed under CC表示2.1日本

① TAuthにログインしていない状態で、ログイン処理を行う場合



②TAuthにログインしている状態で、ログイン処理を行う場合



ワンクリックで認証(パスワード入力不要)

ユーザごとのStanzaの表示位置のカスタマイズ

生息環境についてのDefaultページ例

Microbe DB^{JP}

Definition
marine habitat

Sample List

ID	Description	Biome	EnvFeature	Taxonomy	EnvMaterial	Place
SRS009536	Bacterial diversity study from the Western English Channel	marine biome	marine habitat		sea water	English Channel
SRS009536	Bacterial diversity study from the Western English Channel	marine biome	marine channel		sea water	English Channel

GOLD List

ID	Description
G05313	saline water , marine habitat
G01989	saline water , marine habitat
G00426	saline water , marine habitat
G03565	saline water , marine habitat
G09627	saline water , marine habitat , sea water
G00960	saline water , marine habitat
G09622	saline water , marine habitat
G01400	saline water , marine habitat
G05262	saline water , marine habitat
G00841	marine habitat
G03569	saline water , marine habitat
G00513	saline water , marine habitat
G00384	saline water , marine habitat
G01406	saline water , marine habitat
G01408	saline water , marine habitat
G00864	saline water , marine habitat
G00834	saline water , marine habitat
G00912	saline water , marine habitat
G01661	saline water , marine habitat , sea water
G00523	saline water , marine habitat
G02959	saline water , marine habitat
G02146	saline water , marine habitat
G01097	saline water , marine habitat
G00421	saline water , marine habitat
G00166	saline water , marine habitat
G01444	saline water , marine habitat
G01696	saline water , marine habitat
G01918	saline water , marine habitat
G01338	saline water , marine habitat
G00891	saline water , marine habitat
G010761	saline water , ocean , marine habitat
G04086	saline water , marine habitat
G00853	saline water , marine habitat
G02962	saline water , marine habitat

Taxonomy Compisition via GOLD

生息環境について関連するゲノム・メタゲノムサンプルの情報のリスト表示に特化するようにカスタマイズしたページ

Microbe DB^{JP}

GOLD List

ID	Description
G05313	saline water , marine habitat
G01989	saline water , marine habitat
G00426	saline water , marine habitat
G03565	saline water , marine habitat
G09627	saline water , marine habitat , sea water
G00960	saline water , marine habitat
G09622	saline water , marine habitat
G01400	saline water , marine habitat
G05262	saline water , marine habitat
G00841	marine habitat
G03569	saline water , marine habitat
G00513	saline water , marine habitat
G00384	saline water , marine habitat
G01406	saline water , marine habitat
G01408	saline water , marine habitat
G00864	saline water , marine habitat
G00834	saline water , marine habitat
G00912	saline water , marine habitat
G01661	saline water , marine habitat , sea water
G00523	saline water , marine habitat
G02959	saline water , marine habitat
G02146	saline water , marine habitat
G01097	saline water , marine habitat
G00421	saline water , marine habitat
G00166	saline water , marine habitat
G01444	saline water , marine habitat
G01696	saline water , marine habitat
G01918	saline water , marine habitat
G01338	saline water , marine habitat
G00891	saline water , marine habitat
G010761	saline water , ocean , marine habitat
G04086	saline water , marine habitat
G00853	saline water , marine habitat
G02962	saline water , marine habitat

Sample List

ID	Description	Biome	EnvFeature	Taxonomy	EnvMaterial	Place
SRS009536	Bacterial diversity study from the Western English Channel	marine biome	marine habitat		sea water	English Channel
SRS009536	Bacterial diversity study from the Western English Channel	marine biome	marine channel		sea water	English Channel
SRS009536	Bacterial diversity study from the Western English Channel	marine biome	Aerobic		sea water	English Channel
SRS025487	no title	marine biome	marine habitat		sea water	
SRS005798	no title	aquatic biome	marine habitat		sea water	
ERS013839	no title	marine biome	saline water habitat		sea water	
SRS005855	no title	marine biome	marine habitat		sediment	
SRS025470	no title	marine biome	aquatic habitat		sea water	
ERS013826	no title	marine biome	saline water habitat		coastal water	
SRS172669	marine sediment metagenome Methanogenic sediments metagenome Genomic DNA sample DV-GN	ocean biome	marine habitat		marine sediment	
SRS005786	no title	aquatic biome	marine habitat		water	
SRS005910	no title	marine biome	marine habitat		sea water	Black Sea
SRS006225	no title	aquatic biome	marine habitat		sea water	
SRS025491	no title	aquatic biome	aquatic habitat		water	
SRS005881	no title	aquatic biome	marine habitat		water	
SRS005814	no title	marine biome	marine habitat		sea water	
SRS006286	no title	marine biome	marine habitat		marine sediment	
SRS005922	no title	aquatic biome	aquatic habitat		sea water	
SRS025469	no title	aquatic biome	marine habitat		sea water	
SRS006119	no title	marine biome	marine habitat		water	
SRS005786	no title	marine biome	aquatic habitat		sea water	
SRS006204	no title	marine biome	marine habitat		sea water	
SRS005664	no title	marine biome	marine habitat		water	
SRS006084	no title	marine biome	marine habitat		sea water	
SRS006082	no title	aquatic biome	marine habitat		water	
SRS025475	no title	aquatic biome	marine habitat		hydrothermal fluid	
SRS025458	no title	aquatic biome	marine		water	

ユーザがStanzaの表示位置をカスタマイズした場合、ログアウト後もそのユーザのStanza表示位置は記憶され、次のログイン後に再設定は必要無い



<http://microbedb.jp/MDBdemo/>

MicrobeDB

microbedb.jp/MDBdemo/

[Sign In](#)

Search

Gene: polymerase
Taxonomy: Streptococcus glycerinaceus
Mapping: Streptosporangium roseum
Environment: marine habitat
SRS: rumen
Strain: yeast

51

©2013黒川 顕(東京工業大学) licensed under CC表示2.1日本

H25年度の実施計画

- H23、H24年度に引続き、各DBの高度化および「MicrobeDB.jp」の開発
- GTPSのTogoGenomeを用いたRDF化
- TogoStanzaの仕様確定とMicrobeDB.jpで用いる各種Stanzaの開発 (w/ DBCLS)
- TogoAnnotationの各種オミックスデータへの対応
- MBGDの高度化およびバージョンアップ問題への対応
- メタゲノム解析パイプラインをTSUBAME&DDBJスパコンへの対応
- 生息環境以外の様々なメタデータを記述するオントロジーの構築
- オントロジーマッピング、文献情報抽出の自動化
- MicrobeDB.jpにおけるユーザ管理システムの構築
- 培地オントロジー「GMO」の開発 (w/ DBCLS)
- 各種オミックスデータの統合化
- 微生物自動アノテーションプログラム「MiGAP」との連携

H24年度 主な活動状況

- 微生物統合DB全体会議(4回開催)
- RDF会議(4回開催)
- オントロジー会議(6回開催)
- 菌株会議(1回開催)

- 微生物データベースの将来に関するフォーラム (5/28@東京)
- BioHackathon 2012 (9/2-7@富山)
- The 14th workshop of the GSC (9/17-21@Oxford)
- トーゴーの日シンポジウム 2012 (10/5@東京)
- BioJapan 2012 (10/10-12@横浜)
- SPARQLthon (10月,11月,12月,1月,計4回@DBCLS)
- 情報・システム研究機構シンポジウム 2012 (11/9@東京)
- 分子生物学会 (12/11-14@福岡)
- BioHackathon 12.12 (12/19-21@東京)
- ゲノム微生物学会 (3/8-10@長浜) 予定
- 農芸化学会 (3/24-28@仙台) 予定

