

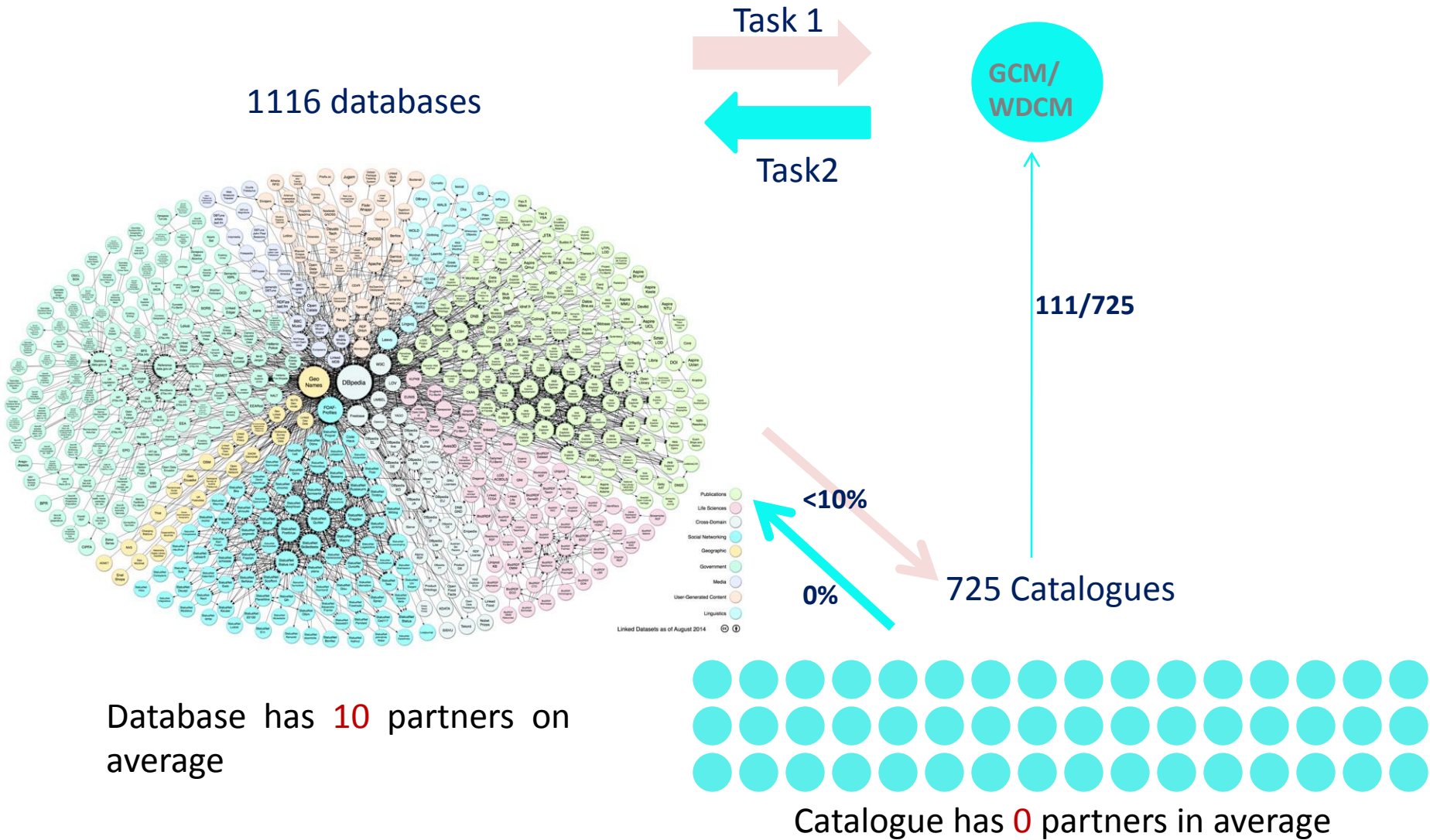
Data integration with Life Science databases: the technical aspects

Alexander Vasilenko, Svetlana Ozerskaya, Oleg Stupar, Lujdmila Evtushenko ,
Paolo Romano, Linhuan Wu, Juncai Ma

Potential help for the Life Science world from microbial culture collections (mBRCs)

1. Assurance of repeatability of experimental data
2. Resolving nomenclatural issues related to microorganisms
3. Strain-specific characters

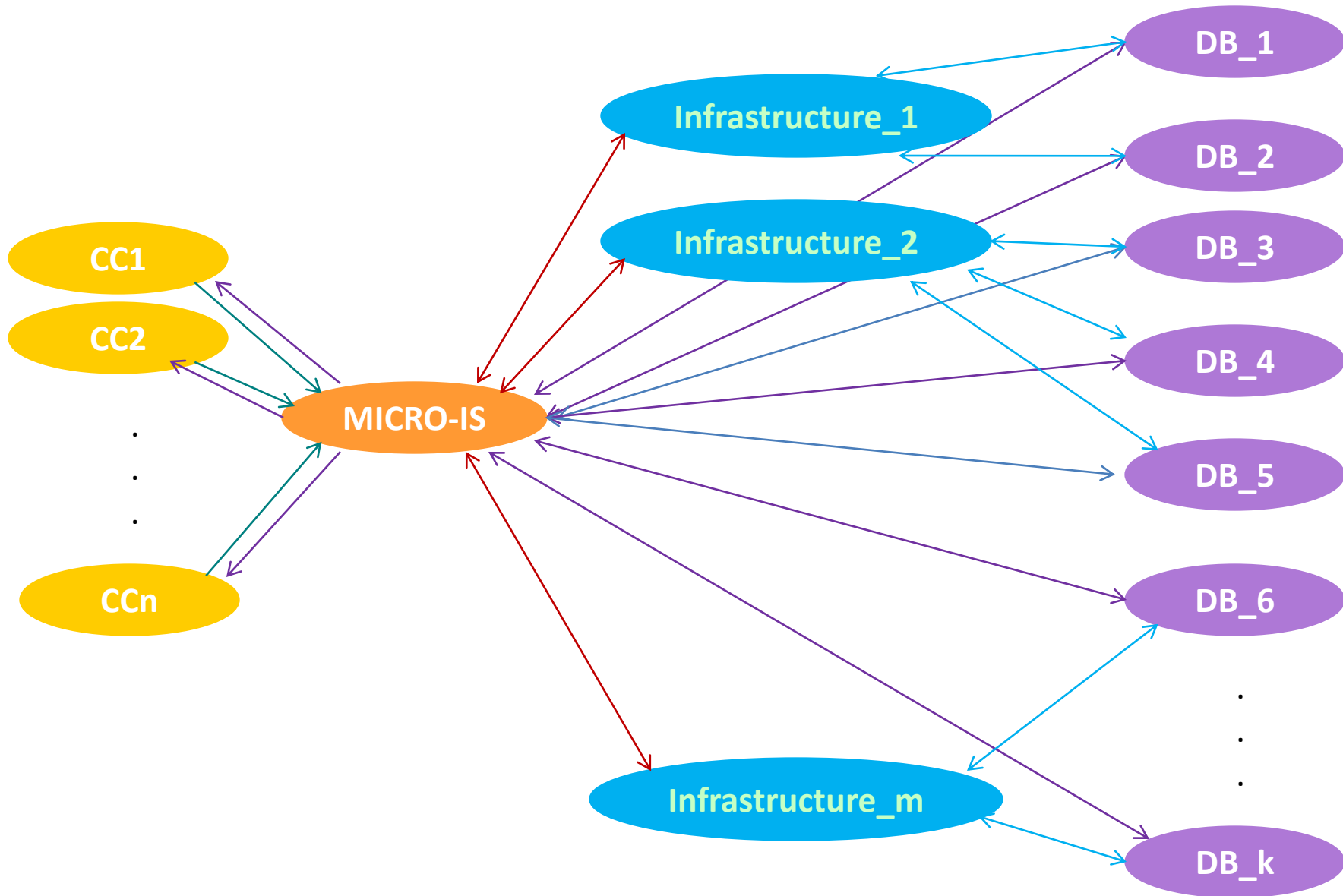
Microbial databases – CC interconnection



Potentially this data integration could mean the tasks:

1. To make CC data visible and accessible from partner Life Science databases,
2. To make partner database records visible and accessible from CC aggregated catalogue,
in the formats:
 - a. To give this data integration for human access,
 - b. To give this data integration for computer programs.

General integration schema



Life Science databases inspected

Total number of life science database names or references discovered in this study is more than 14 800. The total number of database references inspected manually is more than 5500 . The total number of life science databases collected visible online is 2660, the number of databases with microbial data collected is 1116 (plus 9318 bacterial databases in BioCyc).

The main sources:

- MB (1802 entries (<http://metadatabase.org>), 26.12.2015),
- Biosharing (724 databases, (<https://www.biosharing.org/>), 26.12.2015),
- BioMedBrigeds (814 Databases, 27.12.2015), (<http://wwwdev.ebi.ac.uk/fgpt/toolsui/>),
- Pathguide (363 database names, 2013), (<http://www.pathguide.org/>)
- ELIXIR list (579 entries, (<https://bio.tools/?q=database>), 28.1.2016)
- ExPASy (85 + 665 databases, (http://www.expasy.org/old_links), 12.2.2016)
- Bioinformatics Links Directory (621 databases)
- OBRC (<http://www.hsls.pitt.edu/obrc/>) 30.3.2017

Databases parameters collected (an example)

- Unique identifier: **BIODBCORE-000438**
- Database acronym: **dbSNP**
- Database name: **The Database of Short Genetic Variation (single nucleotide polymorphism)**
- Database URL: **<http://www.ncbi.nlm.nih.gov/SNP/>**
- Access level: **Open**
- Practical domain: **health**
- Microbial level: **sp**
- Year of the last correction: **2017**
- Developer/Owner: **USA, NCBI; USA, National Library of Medicine, National Institutes of Health**
- Comment: **Escherichia coli**
- Orientation: **-**
- Properties: **DNA, gene, genome, proteomics, publications, RNA**
- Search by: **-**
- Ontologies list: **SO**
- Partner databases: **Assembly, BioProject, BioSample, ClinVar, dbGaP, dbMHC, dbSTS, dbVar, Ensembl, GenBank, Homologene, IGSR, MapViewer, NCBI Gene, Nucleotide, OMIM, PMC, Protein, PubChem Substance, PubMed, RDP, RefSeq, UniGene, UniProtKB**
- Program interface: **ELIXIR WEB UI, Entrez Programming Utilities (E-Utills)**

Property keywords in databases

826	Gene	226	Interactome
716	Proteomics	219	Taxonomy
625	Publications	195	Drugs
517	Image	195	Peptide
389	RNA	183	Molecules
395	DNA	166	Metabolite
355	Genome	151	Pathogen
361	Enzyme	167	Immunology
316	Cell	156	Toxicology
297	Chemistry	132	Lipid
270	Pathways	45	Microbiome
263	Disease		

Average number of keywords assigned to a database = 6,55

Databases with microbiome data

Biology Reference	HPMCD	PANGAEA
BioSamples	IMG	PLOS One
BioSystems	IMG/M	PMC
Bookshelf	IMG/VR	PSP
EMBL	MEDLINE	pubget
EMBL-EBI	MeSH	PubMed
ENA	Microbiome	PubMed Health
Espacenet	NARCIS	QIAGEN
Europe PMC	NCBI	RefSeq
ForestScience Current Database	NFSD	ScienceDirect
GO Database	NLM Catalog	SRA
GONUTS	Nowomics	TACONIC
GoPubMed	OMIM	UniProtKB
HGTree	OMIM (1)	VetMed Resource
HOMD	OReFiL	WikiGenes

Biggest database producer: BESC (BioEnergy Science Center)

BioCyc pathway/genome database: 9367 databases totally (<http://www.biocyc.org/biocyc-pgdb-list.shtml>)

Group 1 are 7 databases: EcoCyc, MetaCyc, HumanCyc, AraCyc, YeastCyc, LeishCyc, TrypanoCyc

Group 2 are 41 databases generated by program with curation done each is one strain:

Agrobacterium fabrum C58

Anopheles gambiae

Aurantimonas manganoxydans SI85-9A1

Bacillus anthracis Ames

Bacillus subtilis 168

Bacteroides thetaiotaomicron VPI-5482

Candidatus Cardinium hertigii

Candidatus Evansia muelleri

Candidatus Portiera aleyrodidarum BT-QVLC

Caulobacter crescentus CB15

Caulobacter crescentus NA1000

Chlamydomonas reinhardtii

Clostridium saccharoperbutylacetonicum ATCC 27021

Cryptosporidium hominis TU502

Cryptosporidium parvum Iowa

Drosophila melanogaster

Escherichia coli B str REL606

Escherichia coli CFT073

Escherichia coli K-12 substr W3110

Escherichia coli O157:H7 str EDL933

Eubacterium rectale ATCC 33656

Helicobacter pylori 26695

Listeria monocytogenes 10403S

Methylosinus trichosporium OB3b

Mus musculus

Mycobacterium tuberculosis CDC1551

Mycobacterium tuberculosis H37Rv

Penicillium chrysogenum Wisconsin 54-1255

Peptoclostridium difficile 630

Plasmodium berghei ANKA

Plasmodium chabaudi

Plasmodium falciparum 3D7

Plasmodium vivax Sal-1

Plasmodium yoelii 17XNL

Schistosoma mansoni

Shigella flexneri 2a str 2457T

Streptomyces coelicolor A3(2)

Synechococcus elongatus PCC 7942

Thalassiosira pseudonana CCMP1335

Toxoplasma gondii ME49

Vibrio cholerae O1 biovar El Tor str N16961

Group 3 are 9318 databases each database is one bacterial strain with no curation yet

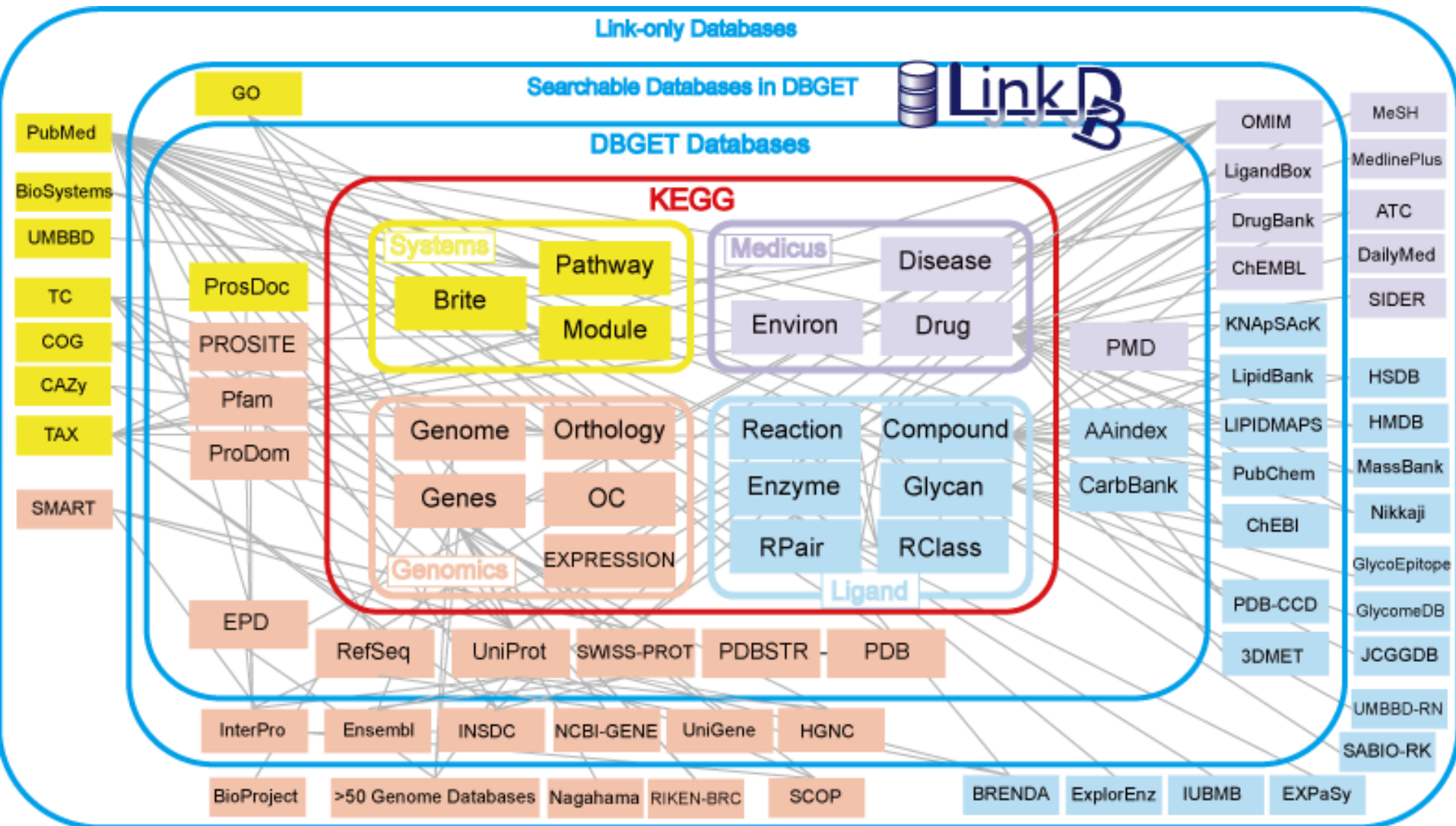
EMBL-EBI (97_{db})

ArrayExpress	e!EnsemblChimpanz	Ensembl	logRECOORD	PICR
ASD	ee	Enzyme Portal	MACiE	PomBase
ASTD	e!EnsemblCow	Enzyme Structures	MEROPS	PRIDE
ATD	e!EnsemblDog	EVA	MetaboLights	PROCOGNATE
BioModels	e!EnsemblFugu	Expression Atlas	Metal MACiE	Reactome
BioSamples	e!EnsemblFungi	FunTree	MicroCosm	RECOORD
Cellular Phenotype Db	e!EnsemblGenomes	GeneDB	MIRIAM collection	Rfam
ChEBI	e!EnsemblGorilla	GWAS Catalog	MTBLS	RNAcentral
ChEMBL	e!EnsemblHorse	HGNC	NRNL1	SAS
CluSTR	e!EnsemblMetazoa	HipSci	NRNL2	SRS@EMBL-EBI
CSA	e!EnsemblMouse	IGSR	NRPL1	SureChEMBL
DGVa	e!EnsemblPig	IMEx	NRPL2	TreeFam
DNAtraffic	e!EnsemblPlants	IMGT/HLA	OLDERADO	UniChem
DrugPort	e!EnsemblProtists	IntAct	PANDIT	UniProt-GOA
e!Ensembl	e!EnsemblRabbit	IntEnz	PDBe	UniSave
e!Ensembl S. cerevisiae	e!EnsemblZebrafish	InterPro	PDBe EM Resources	VASCO
e!EnsemblBacteria	EGA	IPD	PDBeChem	VectorBase
e!EnsemblCat	EMBL	IPD-ESTDAB	PDBsum	
e!EnsemblChicken	EMBL-EBI	IPD-HPA	Pfam	
	EMDB	IPD-KIR	Pfam	
	ENA	IPD-MHC	PhenoDigm	

NCBI (70_{db})

Assembly	Dengue virus database	MedGen	PubChem Compound
BioProject	ECRbase	MEDLINE	PubChem Substance
BioSample	Genbank	MeSH	PubMed
BioSystems	Gene	MMDB	PubMed Health
Bookshelf	Genetic Codes	NCBI	RefSeq
CCDS	Genome	NCBI taxonomy	RefSeqGene
CDD	GEO	NCBI Trace Archives	Retroviruses
ClinGen	GEO DataSets	NLM Catalog	SKY/M-FISH and CGH
ClinVar	GEO Profiles	Nucleotide	SPARCLE
Clone DB	GSS	OMIM	SpliceInfo
COGs	GTR	Organelle genomes	SRA
dbEST	Histone	PMC	Structure
dbGaP	HIV-1	PopSet	TPA
dbMHC	Homologene	Probe	UniGene
dbProbe	IBIS	Protein	UniVec
dbSNP	Influenza Virus Resource	Protein Clusters	Viral genomes
dbSTS	MapViewer	PubChem	Virus Variation
dbVar		PubChem BioAssay	

Databases interconnection example



Interconnection matrix

	3D Lectin	3D RIBOSC	5S RNA Da	A.pernix	UniPROBE	UniProt-Gl	UniProtKB	UniRef	Y MPL	YPM	Zif-BASE		AN
2P2Idb	0	0	0	0	0	0	0	0	0	0	0		1
3D Genome Browser	0	0	0	0	0	0	0	0	0	0	0		1
3D Lectin	0	0	0	0	0	0	0	0	0	0	0		0
3D RIBOSOMAL MODIFICATIO	0	0	0	0	0	0	0	0	0	0	0		3
3DBIONOTES	0	0	0	0	0	0	0	0	0	0	0		0
3Dee	0	0	0	0	0	0	0	0	0	0	0		0
3DEM Benchmark	0	0	0	0	0	0	0	0	0	0	0		0
3D-footprint	0	0	0	0	0	0	0	0	0	0	0		1
3DID	0	0	0	0	0	0	0	0	0	0	0		6

PubChem	0	0	0	0	0	0	0	0	0	0	0		57
PubChem BioAssay	0	0	0	0	0	0	0	0	0	0	0		16
PubChem Compound	0	0	0	0	0	0	0	0	0	0	0		31
PubChem Substance	0	0	0	0	0	0	0	0	0	0	0		25
pubget	0	0	0	0	0	0	0	0	0	0	0		2
Public Data Portal	0	0	0	0	0	0	0	0	0	0	0		1
PubMed	0	0	0	1	0	0	1	0	0	0	1		565
PubMed Health	0	0	0	0	0	0	0	0	0	0	0		6

YPM	0	0	0	0	0	0	0	0	0	0	0		1
ZFIN	0	0	0	0	0	1	1	0	0	0	0		20
Zif-BASE	0	0	0	0	0	0	0	0	0	0	0		0
ZOBODAT Vespoidea	0	0	0	0	0	0	0	0	0	0	0		1
ZorapteraSF	0	0	0	0	0	0	0	0	0	0	0		1
ИИПС	0	0	0	0	0	0	0	0	0	0	0		0
CN	5	8	4	10	6	148	161	15	2	3	5		

UniProtKB database has integration parameters CN=161, AN=501.

In CC catalogues both values are mostly 0.

Databases with more than 26 partners (CN)

364	Pathguide	55	dbProbe	38	NCBI taxonomy	31	PANDORA
161	UniProtKB	55	NCBI	38	OpenHelix	30	FungiDB
159	iProClass	53	PiroplasmsDB	38	ViralZone	30	JPGV
153	COL	52	e!EnsemblGenomes	37	GPCRs	30	Reactome
148	UniProt-GOA	49	EMBL	37	T3DB	30	STRING
143	OReFiL	49	ENA	36	EcoCyc	29	BioModels
129	EcoliWiki	49	SBKB	36	Nucleotide	29	CAZy
128	GeneCards	46	EcoGene	35	BacMet	29	ConsensusPathDB
91	PIR	46	SGD	35	ChEBI	29	EuPathDB
73	UCD 2D-PAGE	45	Gene	35	PRODORIC	29	MACiE
63	Hits	44	InterMitoBase	34	ESTHER	29	RefSeq
61	SWISS-2DPAGE	43	EMBL-EBI	33	Genome	28	EcoProDB
57	PubChem	43	OMIM	33	gpmDB	28	RNAcentral
57	PubChem BioAssay	43	OMIM (1)	32	Ebolavirus	28	ThaleMine
57	PubChem Compound	42	MetaCyc	32	HOGENOM	28	YeastMine
57	PubChem Substance	40	Guide to Pharmacology	32	STITCH	27	BioSystems
		40	MalaCards	31	eNet	27	OrthoDB
				31	HPIDB		

Databases with big attraction number (AN)

565 PubMed	76 PROSITE	54 COGs	38 CATH	16
501 UniProtKB	71 CAS	53 Genome	38 MetaCyc	
275 NCBI taxonomy	71 IntAct	51 DIP	36 CDD	
255 RCSB PDB	70 Reactome	51 STRING	36 SUPERFAMILY	
239 Genbank	69 ChEBI	50 BioCyc	35 dbSNP	
229 Gene	63 FlyBase	50 ENZYME	35 EcoCyc	
199 KEGG	61 MEDLINE	49 KEGG PATHWAY	34 ChEMBL	
193 RefSeq	61 UniGene	48 HPRD	34 PDBsum	
187 EC	59 BioGRID	48 WormBase	33 PLOS One	
182 Pfam	57 PubChem	47 BioProject	33 PRINTS	
160 InterPro	56 GEO	47 BRENDA	31 Homologene	
157 Protein	56 NCBI	47 DDBJ	31 ProDom	
152 Ensembl	56 SMART	47 TAIR	31 PubChem Compound	
125 Nucleotide	55 DrugBank	47 TIGRFAMS	31 wwPDB	
109 OMIM	55 MGI	46 MINT	30 PDBe	
99 SGD	55 PIR	45 MeSH	30 UCSC Archaeal Genome Browser	
80 ENA	55 PMC	43 PANTHER		
80 HGNC	55 SCOP	40 GeneCards		

Integration candidates

	Producer	Databases	Total attraction	% of maximal
1	NCBI	70	2909	33
2	EMBL-EBI	97	1209	14
3	SIB	37	762	9
4	Kioto University	19	348	4
5	Instute Paster	18	148	2
6	BioCyc	9378	133	2
7	InterMine	16	20	0
	1+2+3+4	133	5228	59
	1+3+4	92	4019	45
	2+3+4	78	2319	26
	Total	1115	8870	100

MICRO-IS database partners in efficient solution

ArrayExpress, Assembly, BioModels, BioProject, BioSample, BioSamples, BioSystems, Bookshelf, CDD, Cellosaurus, ChEBI, ChEMBL, COGs, CSA, dbEST, dbProbe, dbSNP, Dengue virus database, DNATraffic, DrugPort, e!Ensembl, e!Ensembl Saccharomyces cerevisiae, e!EnsemblBacteria, e!EnsemblFungi, e!EnsemblGenomes, e!EnsemblProtists, EMBL, EMBL-EBI, EMDB, ENA, Ensembl, ENZYME, Enzyme Structures, EPD, EVA, Expression Atlas, Genbank, Gene, GeneDB, Genetic Codes, Genome, GEO, GEO DataSets, GEO Profiles, GSS, HAMAP, Hits, HIV-1, Homologene, IMEx, Influenza Virus Resource, IntAct, InterPro, KEGG, KEGG BRITE, KEGG DISEASE, KEGG GENES, KEGG GENOME, KEGG GLYCAN, KEGG LIGAND, KEGG MEDICUS, KEGG MODULE, KEGG Organisms, KEGG ORTHOLOGY, KEGG PATHWAY, MACiE, MapViewer, MedGen, MEDLINE, MEROPS, MeSH, MetaboLights, MIAPEGelDB, MMDB, MTBLS, NCBI, NCBI taxonomy, NCBI Trace Archives, neXtProt, NLM Catalog, Nucleotide, OMA, OMIM, OpenFlu, Organelle genomes, PathComp, PathPred, PathSearch, PaxDB, PDBe, PDBe EM Resources, PDBsum, Pfam, PICR, PMC, PMP, PomBase, PopSet, PRIDE, Probe, PROSITE, Protein, Protein Clusters, Protein Spotlight, Proteomes, PubChem, PubChem BioAssay, PubChem Compound, PubChem Substance, PubMed, PubMed Health, Reactome, RefSeq, RefSeqGene, Retroviruses, Rfam, Rhea, RNACentral, SPARCLE, SpliceInfo, SRA, Structure, SugarBind, SWISS-2DPAGE, SWISS-MODEL, SwissVar, UniGene, UniProt-GOA, UniProtKB, UniRef, Viral genomes, ViralZone, Virus Variation

Task 1a: Strains algorithm

NCBI Resources How To

Probe (Penicillium) AND "Penicillium bialowiezense"

Display Settings: Summary, 20 per page

Results: 15

- Penicillium bialowiezense** microarray element probe Pen_COX1_11g
1. Accession: Pr010275422 ID: 10275422
Name: Pen_COX1_11g
Type: microarray element
Application: genotyping
Target organism: [Penicillium bialowiezense](#)
- Penicillium bialowiezense** microarray element probe Pen_COX1_11f
2. Accession: Pr010275421 ID: 10275421
Name: Pen_COX1_11f
Type: microarray element
Application: genotyping
Target organism: [Penicillium bialowiezense](#)
- Penicillium bialowiezense** microarray element probe Pen_COX1_11e
3. Accession: Pr010275420 ID: 10275420

StrainInfo BETA

4 search results for query 'taxon = 'Penicillium bialowiezense''

Species Name	Strain Numbers
<i>Penicillium bialowiezense</i>	NRRL 865 , Thom 5010.5
<i>Penicillium bialowiezense</i>	MUCL 46394
<i>Penicillium bialowiezense</i>	MUCL 46540
<i>Penicillium brevicompactum</i>	CBS 227.28 T , FRR 3574 T , IBT 23044 T , IMI 092237 T ,

WDCM 133 Centraalbureau voor Schimmelcultures Filamentous fungi and Yeast Collection, Netherlands

WDCM 18, Food, Science, Australia, Ryde

WDCM 758, IBT, Culture Collection of Fungi, Denmark

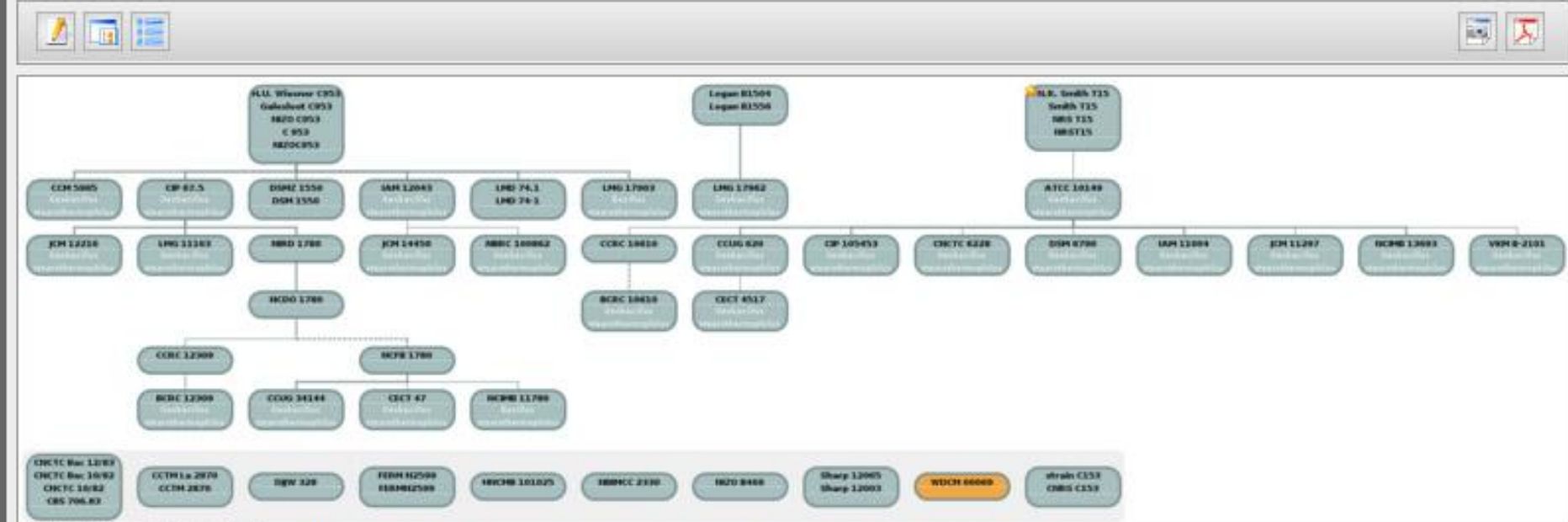
WDCM 214, CABI, Genetic Resource Collection, UK

StrainInfo strain exchange tree

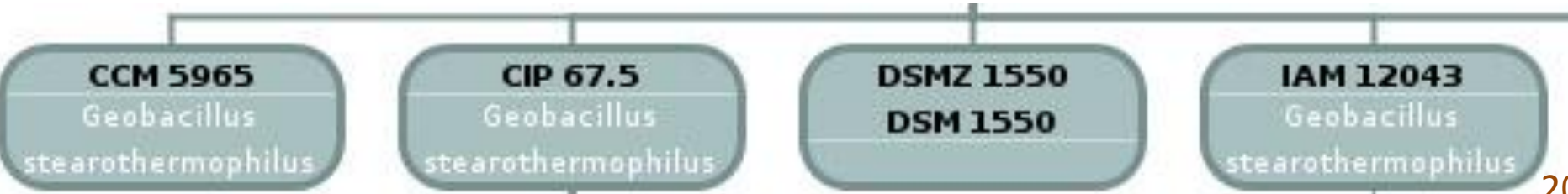
Strain Passport

WDCM 00069 *Geobacillus stearothermophilus*

history



This Histri was manually curated (last saved on 2011-04-12).



Solutions: task 1a: Name processing *

Penicillium cyaneofulvum

Summary:	<i>Penicillium cyaneofulvum</i> Biourge, La Cellule 33: 174 (1923)
Synonymy:	<ul style="list-style-type: none">=<i>Penicillium brunneorubrum</i> Dierckx, Annales de la Société Scientifique de Bruxelles 25 (1): 88 (1901)=<i>Penicillium griseoroseum</i> Dierckx, Annales de la Société Scientifique de Bruxelles 25 (1): 86 (1901)=<i>Penicillium chrysogenum</i> Thom, U.S.D.A. Bureau of Animal Industry Bulletin 118: 58 (1910)=<i>Penicillium baculatum</i> Westling, Svensk Botanisk Tidskrift 4: 139 (1910)=<i>Penicillium notatum</i> Westling, Arkiv för Botanik 11 (1): 95 (1911)• • •=<i>Penicillium fluorescens</i> Laxa, Zentralblatt für Bakteriologie und Parasitenkunde Abteilung 2 86 (5-7):160-165 (1932)=<i>Penicillium camerunense</i> R. Heim, Bull. Acad. R. Belg. Cl. Sci.: 42 (1949)=<i>Penicillium aromaticum</i> f. <i>microsporum</i> Romankova, Uchenn. Zap. Leningr. Univ. Zhadanov: 102 (1955)=<i>Penicillium harmonense</i> Baghd., Novosti Sistematiki Nizshikh Rastenii 5: 102 (1968)
Current name:	<i>Penicillium chrysogenum</i> Thom, U.S.D.A. Bureau of Animal Industry Bulletin 118: 58 (1910)
Classification:	Fungi, Ascomycota, Pezizomycotina, Eurotiomycetes, Eurotiomycetidae, Eurotiales, Trichocomaceae, <i>Penicillium</i>
Facultative or heterotypic synonyms:	<ol style="list-style-type: none">1. <i>Penicillium aromaticum</i> f. <i>microsporum</i> Romankova, Uchenn. Zap. Leningr. Univ. Zhadanov: 102 (1955)2. <i>Penicillium baculatum</i> Westling, Svensk Botanisk Tidskrift 4: 139 (1910)3. <i>Penicillium brunneorubrum</i> Dierckx, Annales de la Société Scientifique de Bruxelles 25 (1): 88 (1901)4. <i>Penicillium camerunense</i> R. Heim, Bull. Acad. R. Belg. Cl. Sci.: 42 (1949)5. <i>Penicillium chlorophaeum</i> Biourge, La Cellule 33: 271 (1923)6. <i>Penicillium chrysogenum</i> Thom, U.S.D.A. Bureau of Animal Industry Bulletin 118: 58 (1910)• • •

TAXONOMY IN MICROBIAL DATABASES

Taxonomy	*	**
NCBI	115	267
GBIF	16	16
IF	7	9
COL	6	6
LPSN	2	2
MycoBank	2	5

* References in databases with taxonomical data

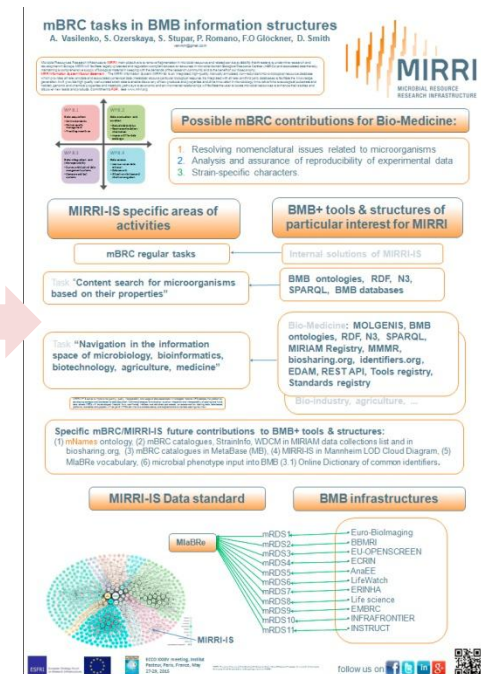
** In all the databases with microbial data

229 - Total number of databases with taxonomical data

94 - No reference to specific taxonomy database

Task 2: solution tools

Producer	Interface
NCBI	Entrez Programming Utilities (E-Utills)
EMBL-EBI	RESTful Web Services interface, Semantic WEB, RDF, SPARQL endpoint
SIB	RESTful Web Services interface, Semantic WEB
Kioto University	KEGG API, LinkDB



Thank you