

StrainInfo.net: breaking down information barriers into holistic data integration scenarios using globally unique identifiers

Peter Dawyndt^{a,b,*}, Bernard De Baets^b, Xianhua Zhou^c, Juncai Ma^c, Jean Swings^{a,d}

^a Laboratory of Microbiology, Ghent University, Belgium
^b Department of Applied Mathematics, Biometrics and Process Control, Ghent University, Belgium
^c Information Network Center, Institute of Microbiology, Chinese Academy of Sciences, PR China
^d BCCMTM/LMG Bacteria Collection, Ghent University, Belgium
 *corresponding author: Peter.Dawyndt@UGent.be

中国科学院微生物研究所
 Institute of Microbiology, Chinese Academy of Sciences

Mission statement

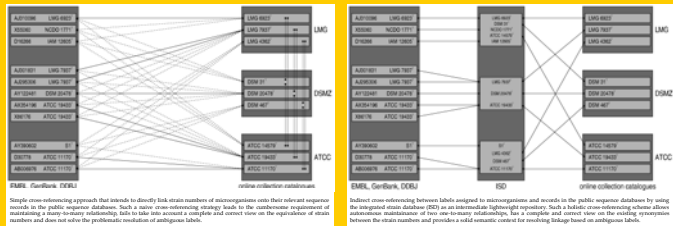
Once a microbial strain gets isolated from its natural habitat, a wealth of information about its genome, proteome, metabolism, clinical and ecological traits can be collected, on which basis it might eventually turn out to become important reusable material for scientific and industrial purposes. Perhaps a culture of the biological material gets deposited into a biological resource center (BRC) for long-term preservation and global dissemination among other BRCs or research institutions, its raw observational data are stored into private or public repositories to establish large-scale identification applications, it becomes commonly accepted as key reference material to support some artificial (human-conceived) taxonomic framework designed as a higher-level conception of biological diversity, it forms a cornerstone for the implementation of an industrial process that is protected by the patent law, or some conclusions drawn from the knowledge attained during scientific research activities wherein the microbial strain was involved are bundled into scientific publications.

Both the actual content of this downstream information on the microorganism and its location in private databases or on the World Wide Web are sensitive to modification over time. As science and technology are moving rapidly, thereby increasingly making use of the scientific merits of previous research results, instant and effortless visibility of this creative and scientific downstream information has become imperative for the realization of successful innovation chains that take full opportunity of the exploitation of biological resources.

The StrainInfo.net portal (www.straininfo.net) therefore envisages the establishment of a technology platform that can stimulate this movement towards using multi-perspective integrated information in a broadened biological and clinical context, as indeed we all would rather like to benefit from automated ICT technologies for keeping track of downstream information on biological resources than putting all our efforts into the tedious and error-prone compilation of relevant knowledge from the heterogeneous and autonomous data collections spread across the information jungle.

What genes are sequenced of a given microorganism ?

- strain-level searches in public sequence databases not accurate enough
- simple integration approach not aware of existing synonyms between labels
- simple integration approach hampered by homonymous strain numbers
- holistic integration scheme of StrainInfo.net performs accurate parallel searches with all known strain numbers in the Integrated Strain Database
- incorporating references to the biological resources from the public sequence databases establishes a true divide and conquer strategy



Advanced queries (workflows) crossing database barriers

More accurate statistics based on integrated information

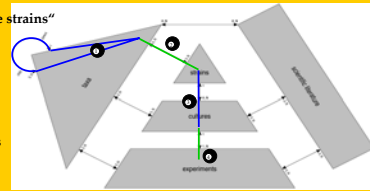
For a selection of biological resource centres (BRC), the table on the left gives a detailed overview of the number of isolates that each pair of BRCs has in common. Additionally the bottom of the table shows a reliable estimation of the number of strains that are unique within each collection^(*).

(*) acronyms of culture collections are taken from the World Data Centre of Microorganisms, taken into account the migration of culture collections and changes in culture collection acronyms.

Example query: "Find all 16S rRNA sequences of *Enterococcus* type strains"

4-step query resolution (right panel)

- 1) find all species of the genus *Enterococcus* (including synonyms)
- 2) find type strain for each species
- 3) find all synonym labels assigned to each type strain
- 4) find all 16S rRNA sequences linked to each of the synonym labels



Integrated Strain Database

Features & philosophy

- automatic assignment of globally unique identifiers to microbial strains
- maintenance of equivalence relation of strain numbers (machine learning)
- error detection/correction of inconsistencies within underlying sources
- provision of semantic context for disambiguation of label resolution
- dynamic integration platform for search results on microbial resources

Facts & Figures

- integrates information of 42 biological resource centres (culture collections)
 - organisms within scope: bacteria, archaea, filamentous fungi, yeasts
 - geographically distributed over all continents worldwide
 - ranging from niche specific to general purpose collections
- # organisms (isolates): 260,444
- # strain numbers (labels): 589,157
- # bacterial type strains affected by errors: 768/6137 (12.51%)

What collections do have a given microorganism in their holdings ?

Information resource	country	WDCM	strain acronyms	records	homepage
American Type Culture Collection	U.S.A.	ATCC/ATCC BAA/ATCC MYA	44,975	www.atcc.org	
BIVTEC Culture Collection	Thailand	WDC	1,100	www.bivtec.org.th	
Bioresource Collection and Research Center	Taiwan	BCRC	10,035	www.brc.gov.tw	
Centre for Research and Biotechnology	Italy	CCM4	174	www.cnr.it/istituto/istituto.htm	
Canadian Collection of Fungal Cultures	Canada	CCM3	10,062	www.cbc.calgary.ca	
Centre for Microorganisms and Cell Culture	Czech Republic	CCM3/CCM F	3,477	www.mikrobiologie.cz	
Microorganism Collection of the University of Microorganisms	Spain	CCM5	763	www.ccm5.com	
Centre for Culture Type	Brazil	CCCT	2,189	www.cct.org.br	
Colección Nacional de Cultivos Microbianos	Spain	CCM2	10,035	www.ccm2.com	
Collection Nationale de Cultures Microbiennes et Cultures Cellulaires	Mexico	CCM8	10,035	www.ccm8.com	
Centre for Microbiology and Cell Culture	France	CCM1	4,254	www.ccm1.com	
Collection de Cultures de Microorganismes et de Cellules	France	CCM7/CCP A	7,812	www.ccm7.com	
Deutscher Sammlungszentrum für Mikroorganismen und Zellkulturen	Germany	DSMZ	14,028	www.dsmz.de	
Food Science Australia, Yeasts	Australia	FAO	2,949	www.foodscience.gov.au	
TEAMBI Culture Collection	Finland	TEAMBI	2,089	www.teambi.com	
International Collection of Micro-Organism Cultures and Cell Cultures	Japan	ICM	714	www.icm.ac.jp	
IAM Culture Collection	Japan	IAM	3,120	www.iam.ac.jp	
International Collection of Microorganisms and Plant	Spain	ICM3	14,035	www.icm3.com	
BCCM TM /BIM - Biomedical Culture and Yeasts Collection	Belgium	402/BIM	4,842	www.bim.ac.be	
Japan Collection of Microorganisms	Japan	CMC	7,914	www.jcm.ac.jp	
Korean Agricultural Culture Collections	Korea	KACC	12,039	www.kacc.or.kr	
Korean Collection for Type Cultures	Republic of Korea	KCTC	7,018	www.kctc.or.kr	
BCCM TM /LMG Bacteria Collection	Belgium	LMG	13,882	www.lmg.ac.be	
BCCM TM /MCC - Microbiology of Université catholique de Louvain	Belgium	308/MCC	11,970	www.mcc.ucl.ac.be	
Microorganism Collection of Mexico	Mexico	CCM6	113	www.ccm6.com	
National Bank for Industrial Microorganisms and Cell Cultures	Belgium	138/NIMCC	1,466	www.nimcc.com	
NITE Biological Resource Center	Japan	185/NITE	12,039	www.nite.or.jp	
National Collection of Agricultural and Industrial Microorganisms	Finland	38/NAMIB/NACAM/NCAM Y	2,033	www.namib.com	
National Collection of Industrial Microorganisms	India	NCIM	2,002	www.ncim.org.in	
National Collection of Industrial Food and Marine Bacteria	U.K.	NCIMB	7,904	www.ncimb.ac.uk	
National Collection of Yeast Cultures	U.K.	NCYC	2,706	www.ncyc.ac.uk	
Former Culture Collection of Cyprobiotechnia	France	481/CC	452	www.cccp.fr	
Iranian Type Culture Collection	Iran	124/PTCC	322	www.ptcc.ac.ir	
University of Alberta Microbiology Collection and Herbarium	Canada	78/UMS	452	www.umicrobiology.com	
All-Russian Collection of Microorganisms	Russian Federation	542/VKM B/VKM A/VKM F/VKM Y	6,120	www.vkm.ru	
Common Access to Biological Resources and Information	Taiwan	CABI	30,020	www.cabi.org	