Linking CABRI catalogues to the EMBL Data Library: lessons learned and future perspectives

Paolo Romano National Cancer Research Institute, Genova (paolo.romano@istge.it)

Outline of this presentation

- CABRI network services
- EBRCN project
- Methods for linking CABRI and EMBL
- Results
- Future perspectives

CABRI: the EU project

Common Access to Biological Resources and Information

Project funded by the EU from 1996 to 1999

Partners

- INSERM (coordinator L. Réchaussat)
- BCCM, CBS, DSMZ, ECACC, HGMP-RC, ICLC, NCCB (resources)
- HGMP-RC, IST, CERDIC (ICT)

<u>Resources</u>

- Microorganisms (bacteria, yeasts, fungi strains)
- Cells (animal and human cell lines, hybridomas, HLA typed B lines)
- Plasmids, phages, plant cells, plant cell viruses, DNA probes
- Overall, more than 90.000 items in catalogues

CABRI: the goals

- Setting Quality Management Guidelines
- Distributing biological resources of the highest quality
- Integrating searches and access to catalogues
 - A "one-stop-shop" for quality resources
 - An "ad hoc" search (CABRI Simple Search)
 - Shopping cart (pre-ordering facility)

CABRI: the SRS search engine

SRS was choosen because:

- it has a simple and effective interface
- it manages heterogeneous databases and it uses flat file format
- it offers both internal and external links and has a link operator
- it is easily expandible (new databases)
- it is widely used for many molecular biology database

CABRI: data structure

For each material, three data sets identified:

- Minimum Data Set (MDS): essential data, needed to identify individual resources
- Recommended Data Set (RDS): all data that are useful to describe individual resources
- Full Data Set (FDS): all data available on the resources
- For each information, data input and authentication guidelines, including:
- Detailed textual description of the information
- In-house reference lists of terms and controlled vocabularies
- Predefined syntaxes (e.g., Literature, scientific names)

CABRI: Data sets

Data set	Field label	Catalogues
MDS	Strain_number	All
MDS	Other_collection_numbers	All
MDS	Name	All
RDS	Race	All
MDS	Organism_type	All
MDS	Restrictons	All
MDS	Status	All
MDS	History	All
RDS	Misapplied_names	All
RDS	Substrate	All
RDS	Geographic_origin	All
RDS	Sexual_state	All
RDS	Mutant	All
FDS	Genotype	DSMZ

CABRI: Name field

Field	Name
Description	Full scientific and most recent name of the strain.
	It includes:
	 Genus name and species epithet
	•Subspecies
	Pathovar
	Authors of the name
	•Year of valid publication or validation
	 Approbation of the name
Input process	Enter full scientific name as given by depositor and confirmed (or changed) by collection. Names of authors of the name, year of valid publication or validation and approbation are included after a comma.
	Values for approbation:
	AL = approved list, c.f.r. IJSB 1980
	VL = validation list, in IJSB after 1980
	VP = validly published, paper in IJSB after 1980
	Reference list: DSMZ list of bacterial names
Required for	MDS

```
Strain number LMG 1(t1)
Other collection numbers CCUG 34964;NCIB 12128
Restrictions Biohazard group 1
Organism type Bacteria
Name Phyllobacterium rubiacearum, (ex Knösel 1962) Knösel 1984 VL
Infrasubspecific names -
Status Type strain
History <- 1973, D.Knösel
Conditions for growth Medium 1, 25C
Form of supply Dried
Isolated from Pavetta zimmermannia
Geographic origin Germany, Stuttgart-Hohenheim
Remarks Stable colony type isolated from LMG 1. See also Agrobacterium sp. LMG 1
   (t2)
Strain number LMG 1(t2)
Other collection numbers -
Restrictions Either Biohazard group 1 or Biohazard group 2
Organism type Bacteria
Name Agrobacterium sp.
Infrasubspecific names -
Status -
Other names Phyllobacterium rubiacearum, (ex Knösel 1962) Knösel 1984 VL
History <- D.Knösel (Phyllobacterium rubiacearum)</pre>
Conditions for growth Medium 16, 28C
Form of supply Dried
Isolated from Pavetta zimmermannia
Geographic origin Germany, Stuttgart-Hohenheim
Remarks One (t2) out of two stable colony types isolated from the original culture
   LMG 1.
```

CABRI: Extra features

Simple Search:

- Search by ID(s), name(s), all other fields
- Search by name(s) with synonyms support

Shopping cart:

- Set of mixed javascripts and perl scripts
- Pre-order facility (email or fax)

Links to reference databases:

 SRS and HTML links to reference databases (media, synonyms, hazard, etc...) for all catalogues

CABRI: Synonyms' support

The synonyms' support allows to retrieve strains by using synonyms or alternate names

- Only allowed for micro-organisms
- Managed through a perl script
- Searched terms are first matched against synonyms' reference databases with getz
- Results are added to the initial search and a new search is carried out
- New results are then displayed
- Links to synonyms' databases are provided

CABRI: Contents updating

- New collections added after the end of the project:
 - CABI Biosciences (fungi, yeasts and bacteria)
 - CIP (bacteria)
 - NCIMB (bacteria and phages)
- No new catalogues since 2002
- Catalogues updated 1 to 2 times a year
- Now 28 catalogues for more than 110.000 items
- HyperCatalogue: hypertext for accessing records through a hyerarchical set of indexes

CABRI: HyperCatalogue



Romano, Linking CABRI and EMBL, Brussels 2005

CABRI: collections and catalogues

	DNA probes	Bacteria strains	Fungi & yeasts strains	Plasmids	Phages	Cell lines	Hybrid omas	HLA typed B lines	Plant cells	Plant cell viruses
BCCM		11,912	15,865	681						
CABI		243	11,677							
CBS		1,038	31,746							
CIP		7,038								
DSMZ		8,366	2,296	284	89	526			476	426
ECACC	261					956	397	238		
ICLC						237				
NCCB		4,609		396	69					
NCIMB		6,986			69					
Total	261	40,192	61,584	1,361	227	1,719	397	238	476	426

CABRI: some limitations

Closed system

Lacking of external links

Lacking of interoperability

Goals of the integration

• In biology, integration is needed in order to:

- Achieve a better and wider view of all available information
- Carry out analysis and/or searches involving more databases and softwares in one step only
- Carry out a real data mining
- For biological resources, integration is useful for:
 - Improving catalogues descriptions of the resources by retrieving new data from well curated databases
 - Retrieving detailed information on resources while searching molecular biology databases
 - Improving visibility of high quality collections of biological resources
 - Extending use of biological materials of certified quality

Integration methods

From syntactical to semantical links:

- Explicit (reciprocal) links (xrefs)
- Implicit links (e.g., names)
- Common contents (vocabularies)
- Shared data models and schemas
- Ontologies

EBRCN: A focal point for BRCs

European Biological Resource Centres Network

- Wp1 Co-ordinate European BRC policies, prepare a coordinated European response to international initiatives on biodiversity and become the European focal point for BRCs
- **Wp2** Develop new and maintain existing **quality standards** for European BRCs
- Wp3 Establish a framework to maximise **complementarity** and minimise **duplication** among European BRCs
- Wp4 Introduce new techniques in Information Technology to the EBRCN to add value to current catalogue information and enhance accessibility

Wp5 Collate and disseminate relevant information to the BRCs

EBRCN: Extending contents

Workpackage 4

"Introduce new techniques in information technology to the EBRCN to add value to current catalogue information and enhance accessibility"

Objective

Link catalogue data to external databases Improve accessibility of catalogues Create an interconnected biological resource database

EBRCN: links to external databases

For all catalogues:

- Links to Medline through Pubmed ID
- Links to representative EMBL records

For selected catalogues:

- Links to plasmids' maps (plasmids)
- Links to microscope images (microorganisms)
- Links to other dbs (nomenclature, acronyms, genes) under evaluation

EBRCN: links to Medline

```
Syntax change:
```

```
add [PMID: <number>] after bibliographic reference
```

```
Links in place (> 7000):

Plasmids: LMBP (375), NCCB (30)

Cell lines: ICLC (294), DSMZ (905)

Fungi: CBS (454)

Yeasts: CBS (1132)

Phages: NCCB (30)
```

Literature reference file: DSMZ (3818)

EBRCN: revised reference paper field

Field	Reference paper
Description	Original paper [if available]
Input process	New entries: JournalTitle Year; Volume(issue): beginning page#-ending page#
	The title is abbreviated following international standard rules (ISSN). Abbreviations are without dot. Authors and title of the article are not mentioned.
	The reference can be followed by the Pubmed ID enclosed within square brackets as follows:
Required for	

EBRCN: links to other databases

Plasmids' maps:

Syntax:

New FDS field: 'External_links map <name>'

Links in place:

Plasmids: LMBP (777)

Images of micro-organisms:

Syntax:

New FDS field: 'External_links image <name>'

Links in place:

None (waiting for next catalogues' update)

EBRCN: EMBL Data Library

- EMBL Data Library is the European database for DNA sequences
- It is updated daily and a coordination with NCBI and DDBJ ensures its completeness
- Static releases are issued every three months
- It is offered at EBI and other centers by means of SRS

EBRCN: EMBL Data Library size

• EMBL Data Library 74 (Mar 2003):

- Sequences: 23,234,788, Bases: 30,356,786,718
- EMBL Data Library 81 (Dec 2004):
 - Sequences: 40,696,839, Bases: 44,285,259,441
 - WGS sequences: 5,408,558, Bases: 34,986,041,399
- EMBL Data Library 82 (Mar 2005):
 - Sequences: 43,246,005, Bases: 46,927,070,905
 - WGS sequences: 6,228,397, Bases: 38,207,643,477

• Size: 7,3% more vs 81 (3 months), 112,9% vs 74 (24 months)

Test have been conducted to identify how to link "on-the-fly" to EMBL Data Library through SRS

Tests performed on: Bacteria and Archaea Animal and Human Cell Lines Fungi and Yeasts Plasmids Viruses

EBRCN: EMBL links variability

Links vary according to the materials

Links can be created on the basis of various EMBL fields:

- All-text (not very useful)
- Organism (for micro-organisms)
- Division (useful for viruses and plasmids)
- Feature Table data (it allows a correct definition of a source through Key, Qualifier, Description)

Key	Location	Qualifier	Value
source	1513	db_xref	taxon:56450
		mol_type	genomic DNA
		organism	Xanthomonas cassavae
		strain	LMG673
misc rna	1513	note	16S-23S ribosomal RNA intergenic spacer

EBRCN: EMBL links variability

Example:

set up a link from the CBS fungi catalogue to EMBL for CBS 100.20 strain

Fields and values:

Organism: fungi

Ft-Key: source

Ft-Qualifier: strain

Ft-Description: "cbs 100.20"

EBRCN: EMBL links variability

Annotation problems in EMBL:

CBS 100.20 can be annotated as: CBS 100.20 or CBS100.20 or CBS 10020 or CBS10020

CBS 112345 can be annotated as CBS12345

Indexing problems at EBI: CBS 100.20 is indexed as "CBS", "100" and "20" The dot is not included and is used as a separator CABRI unique index key is "CBS 100.20"

EBRCN: Linking to EMBL on-the-fly

Example of search "on-the-fly":

Searching for fil. fungi strain CBS 100.20 It involves: fungi & source & cbs 100.20

- ((([emblrelease-FtKey:source] &
 [emblrelease-FtQualifier:strain] &
 (([emblrelease-FtDescription:cbs] &
 [emblrelease-FtDescription:100]) |
 [emblrelease-FtDescription:cbs100]) &
 [emblrelease-FtDescription:20]))
 - < [emblrelease-Organism:fungi*])

Test for linking "on-the-fly" to EMBL Data Library through SRS, without explicit IDs, gave negative results

An alternative approach:

- Identify "a priori" all existing crossreferences between EMBL and CABRI catalogues, based on CABRI IDs
- Add properly formatted links to CABRI in EMBL and use these links when linking from CABRI

CABRI IDs included in EMBL data library and distributed with it

EBRCN: procedure

Work carried out with EMBL 80 and CABRI catalogues version 2004.1

Common SRS site for CABRI and EMBL established Modified indexing -> common keys format SRS internal links established for micro-organisms Automatically identified links downloaded and analysed List of dubious references sent to collections Some comments returned Data submitted to EBI and included in EMBL

EBCRN: Common site established

iet)	Quick Search
Search Options Select the databanks you want to search Enter your search terms in the Quick Search box, or choose a query form from below Standard Query Form Extended Query Form U can browse through the entries in any stabanks. "st, select the databanks u want to browse, then ck: Browse Entries Tips Dookmark this link to curn to your project	Available Databanks • Expand all • Collapse all • Sequence databanks - complete • Sequence databanks - subsections all • EMBL (Release) □ EMBL (Updates) • Animal and Human Cell Lines • HLA typed B Cell Lines • Hybridomas • Bacteria all • BCCM/LMG Bacteria Catalogue □ CABI BACT □ CIP BACT □ CBS BACT □ DSMZ BACC □ NCIMB BACT □ NCCB BACT □ BCCM IHEM □ BCCM MUCL □ CBS FIL □ CBS YEAST □ CABI FIL □ CABI YEAST □ DSMZ FUNGI • Plasmids • Plant Cells • Plant Cells • DNA Probes

EBCRN: Common format index keys created

CABRI indexing: by whole ID

CBS 100.20 -> 'CBS 100.20'

EMBL indexing: by single words

CBS 100.20 -> 'CBS' + '100' + '20'

CBS100.20 -> 'CBS100' + '20'

Common indexing: name (only letters), possibly followed by space, followed by string (including letters, numbers, dot, dash), punctuation removed

CBS 100.20 -> 'CBS10020' CBS100.20 -> 'CBS10020'

Special case (not currently managed): NCCB LMD and Phabagen bacteria catalogues

EMBL (before)

EMBL (after)

```
| /[^= \n]+/ {$Uniq:name}('=' (/[a-zA-Z0-9]+/ {$Uniq:word} | /./)*)?) ~
```

CABRI (before)

CABRI (after)

EBCRN: SRS links EMBL - CABRI

\$Link: [from:\$EMBLRELEASE_DB
 to:\$BCCM_LMG_DB
 fromField:\$DF_FtDescription
 toField:\$DF CABRI Strain number]

\$Link: [from:\$EMBLRELEASE_DB
 to:\$CBS_BACT_DB
 fromField:\$DF_FtDescription
 toField:\$DF_CABRI_Strain_number]

EBRCN: Automatic identification of links

5	From Databank	Entries Linked	To Databank	Entries Linked	Links Total	Indexing Date
	EMBLRELEASE	0	EMBLRELEASE features	0	0	
	EMBLRELEASE	0	EMBLRELEASE reference	0	0	
	EMBLRELEASE	0	EMBLRELEASE counter	0	0	
	EMBLNEW	249034	EMBLRELEASE	<u>249034</u>	249034	06-Sep-2004
	EMBLRELEASE	858	BCCM LMG	<u>659</u>	864	06-Sep-2004
	EMBLRELEASE	0	CABI BACT	0	0	06-Sep-2004
	EMBLRELEASE	0	CBS BACT	0	0	
	EMBLRELEASE	<u>3910</u>	DSMZ BACT	1713	3910	06-Sep-2004
	EMBLRELEASE	335	NCIMB BACT	204	335	06-Sep-2004
	EMBLRELEASE	1	NCCB BACT	1	1	06-Sep-2004
	EMBLRELEASE	717	CIP BACT	388	717	06-Sep-2004
	EMBLRELEASE	26	BCCM IHEM	22	26	06-Sep-2004
	EMBLRELEASE	162	BCCM MUCL	133	162	06-Sep-2004
	EMBLRELEASE	3148	CBS FIL	2043	3155	06-Sep-2004
	EMBLRELEASE	<u>55781</u>	CBS YEAST	953	55789	06-Sep-2004
	EMBLRELEASE	233	CABI FIL	<u>187</u>	233	06-Sep-2004
	EMBLRELEASE	0	CABI YEAST	0	0	06-Sep-2004
	EMBI RELEASE	47	DSMZ FUNGI	34	47	06-Sep-2004

EBRCN: EMBL records citing CABRI strains

Reset	Query "(EMBLNEW < BCCM_LMG)" found 892 entries	next
Apply Options to:	EMBL (Updates):LM191	
	EMBL (Updates):LM192	
O selected results only	EMBL (Updates):LM195	
	EMBL (Updates):LM366	
C unselected results only	EMBL (Updates):LM367	
	EMBL (Updates):LM368	
Result Options	EMBL (Updates):LM369	
·	EMBL (Updates):LM370	
Link to related information:	EMBL (Updates):LM371	
Link	EMBL (Updates):LM372	
Covo requita:	EMBL (Updates):LM373	
Save results. Save	EMBL (Updates):LM374	
	EMBL (Updates):LM375	
Display Options	EMBL (Updates):CEY51B9A	
	EMBL (Updates):LMFL2719	
View results using:	EMBL (Updates):LMFL4768	
SeqSimpleView	EMBL (Updates):A37476	
	EMBL (Updates):AX007138	
Sort results by:	EMBL (Updates):AX008203	
unsorted 💌	EMBL (Updates): AB032353	
ascending	EMBL (Updates): AB045098	
	EMBL (Updates): AB071140	
o descending	EMBL (Updates):AB071141	1

EBRCN: CABRI records cited in EMBL



EBRCN: EMBL custom views

Reset Query "(emb	olrelea	ise < cip_bact)" found 717 entr	ies <u>next</u>	2			
Apply Options to:		<u>EMBL (Release)</u>	<u>Organism</u>		CIP BACT	Other collection numbers	<u>Name</u>
C selected results only		EMBL (Release): CMDNAADE1	Candida maltosa	<u>CIP</u>	BACT:CIPA10		Pseudomonas aeruginosa
unselected results only		EMBL (Release): CMADNAPSS	Candida maltosa	<u>CIP</u>	BACT:CIPA101	-	Corynebacteriur diphtheriae
Result Options	Γ	EMBL (Release): AB008684	Acinetobacter baumannii	<u>CIP</u>	BACT:CIP7034	ATCC 19606; CCUG 19096; DSM 30007; IAM 12088; LMG 1041; NCTC 12156; NCIMB 12457	Acinetobacter baumannii
Link to related information:		EMBL (Release): AB008686	Acinetobacter sp.	<u>CIP</u>	BACT:CIP7029	ATCC 19004	Acinetobacter genomospecies 3
Save results: Save	Π	EMBL (Release): AB008687	Acinetobacter sp.	<u>CIP</u>	BACT:CIPA165	ATCC 17979	Acinetobacter genomospecies 6
Display Options		EMBL (Release): AB008688	Acinetobacter sp.	<u>CIP</u>	BACT:CIP7031	ATCC 9957	Acinetobacter genomospecies 9
shortEmblCip	Г	EMBL (Release): AB008689	Acinetobacter sp.	CIP	BACT:CIP7012	ATCC 17924; LMG 1003; NCIMB 9019	Acinetobacter genomospecies 10
Sort results by:	Γ	EMBL (Release): AB008690	Acinetobacter sp.	CIP	BACT:CIP6346	ATCC 11171; DSM 590; NCIMB 8250	Acinetobacter genomospecies 11
• ascending • descending	Г	EMBL (Release): AB008691	Acinetobacter haemolyticus	CIP	BACT:CIP643	ATCC 17906; CCM 2358; CCUG 888; LMG 996; NCTC 10305; NCIMB 12458	Acinetobacter haemolyticus
Show 30 results	Г	EMBL (Delease): AB008692	Acinetobacter	CIP	BACT:CIP646	ATCC 17909; CCUG 19095;	Acinetobacter

EBRCN: EMBL custom views (ii)

y:		/pseudo	
g	EMBL (Release): AB008684	/db_xref="taxon:470" /mol_type="genomic DNA" /organism="Acinetobacter baumannii" /strain="CIP 70.34; ATCC 19606"	CIP BACT:CIP7034
results y view		/db_xref="GOA: <u>Q43907</u> " /db_xref="UniProt/Swiss-Prot:Q43907" /transl_table=11 /EC_number="5.99.1.3" /gene="gyr8" /product="DNA gyrase subunit 8" /product="DNA gyrase subunit 8" /protein_id=" <u>BAA75401</u> .1" /translation="DNSYKVSGGLHGVGVSVVNALSSKLHLTIYRAGQIHEQEYHHGDP QYPLRVIGETDNTGTTVRFWPSAETFSQTIFNVEILARRLRELSFLNAGVRIVLRDERI NLEHVYDYEGGLSEFVKYINEGKNHLNEIFHFTADADNGIAVEVALQWNDSYQENVRCF TNNIPQKDGGTHLAGFRAALTRGLNQYLENENILKKEKVNVTGDDAREGLTAIISVKVP DPKFSSQTKEKLVSSEVKPAVEQAMNKEFSAYLLENPQAAKSIAGKIIDAARARDAARK AREMTRRKSALDIAGLPGKLADCQEKDPALSELYLVEGDSAGGSAKQGRNRKMQAILPL KGKILNVERARFDKMISSQEVGTLITALGCGIGREEYNPDKLRYHKII"	
	EMBL (Release): AB008686	/db_xref="taxon:472" /mol_type="genomic DNA" /organism="Acinetobacter sp." /strain="CIP 70.29; ATCC 19004"	CIP BACT:CIP7029
		/db_xref="GOA: <u>Q44270</u> " /db_xref="UniProt/Swiss-Prot:Q44270" /transl_table=11 /EC_number="5.99.1.3" /gene="gyrB" /product="DNA gyrase subunit B" /protein_id=" <u>BAA75403</u> .1" /translation="DNSYKVSGGLHGVGVSVVNALSSKLHLIINRAGQVHEQEYHHGDP QYPLRVIGETDSSGTTVRFWPSELTFSQTIFNVEILARRLRELSFLNAGVRIVLRDERI NLEHVYDYEGGLSEFVKYINEGKTHLNEIFHFTADTENGIGVEVALQWNESYQENVRCF TNNIPQKDGGTHLAGFRAALTRGLNQYLENENILKKEKVNVTGDDAREGLTAIISVKVP DPKFSSQTKEKLVSSEVKPAVEQAMNKEFSAYLLENPQAAKSIAGKIIDAARARDAARK AREMTRRKSALDIAGLPGKLADCQEKDPALSELYLVEGDSAGGSAKQGRNRKMQAILPL KGKILNVERARFDKMISSQEVGTLITALGCGIGREEYNPDKLRYHKII"	
	EMBL (Release): AB008687	/db_xref="taxon:472" /mol_type="genomic DNA" /organism="Acinetobacter sp." /strain="CIP A165; ATCC 17979"	CIP BACT:CIPA165
-		/db_xref="GOA: <u>Q44271</u> " /db_xref="UniProt/Swiss-Prot:Q44271"	

EBRCN: CABRI custom views

BCCM LMG	Strain number	Name	EMBL (Release)	AccNumber	Organism		
BCCM LMG:LMG383	LMG 383	Rhizobium radiobacter, (Beijerinck and van Delden 1902) Young, Kuykendall, Martínez-Romero, Kerr and Sawada 2001 VP	<u>EMBL</u> (Release):ARA130719	AJ130719	Agrobacterium tumefaciens		
BCCM LMG:LMG471	LMG 471	Xanthomonas sacchari, Vauterin,	EMBL (Release):AF209762	AF209762	Xanthomonas sacchari		
		and Swings 1995	EMBL (Release):XSY10766	Y10766	Xanthomonas sacchari		
BCCM LMG:LMG568	LMG 568	Xanthomonas campestris pv.	EMBL (Release):AF209755	AF209755	Xanthomonas campestris		
		(Pammel 1895) Dowson 1939	EMBL (Release):AY227424	AY227424	Xanthomonas campestris pv. campestris		
BCCM LMG:LMG673	LMG 673	Xanthomonas cassavae, (ex Wiebe and	EMBL (Release):AF209756	AF209756	Xanthomonas cassavae		
			Wiehe and Dowson 1953) <u>E</u> Vauterin, Hoste, (1	Dowson 1953) Vauterin, Hoste,	EMBL (Release):AY227425	AY227425	Xanthomonas cassavae
		Kersters and Swings 1995 VP	EMBL (Release):XCY10762	Y10762	Xanthomonas cassavae		
BCCM LMG:LMG689	LMG 689	Xanthomonas arboricola pv. corylina, (Miller, Bollen, Simmons, Groes and Barros	EMBL (Release):AY227394	AY227394	Xanthomonas arboricola		

EBRCN: Download of data

	delicatum	/mol_type="genomic DNA" /organism="Aquaspirillum delicatum" /strain="LMG 4328"		delicatum, (Leitson 1962) Hylemon, Wells, Krieg and
		/product="16S ribosomal RNA"		Jannasch 1973 AL
<u>ites):AF085226</u>	Wautersia paucula	/db_xref="taxon:82633" /mol_type="genomic DNA" /organism="Wautersia paucula" /strain="LMG 3413"	BCCM LMG:LMG3413	Ralstonia paucula, Vandamme, Goris, Coenye, Hoste, Janssens, Kersters
		/product="16S ribosomal RNA"		1999 VP
ites):AF090993	Lactobacillus fermentum	/db_xref="taxon:1613" /mol_type="genomic DNA" /note="conserved RAPD fragment UF" /organism="Lactobacillus fermentum" /strain="LMG 6902"	BCCM LMG:LMG6902	Lactobacillus fermentum, Beijerinck 1901 AL
		/codon_start=2 /db_xref="GOA: <u>Q9XDB9</u> " /db_xref="UniProt/TrEMBL:Q9XDB9" /transl_table=11 /gene="hda" /product="D-2-hydroxyacid dehydrogenase" /protein_id=" <u>AAD41656</u> .1" /translation="RALNARVLGYDLKPREEMEGIVEYVSKEELLRQSDVVSLHVDLNP TSEGLLTKAEFDQMKPGASLVNASRGPVVNTADLINALKDGPLAFAALDTVEGEGPIFN ADHRQDGLKGEPLVEELHAMDNVILTPHIAFFTNIAVQNMVDISLDDVITILNGQPSEH EVAN"		
<u>ites):AF090994</u>	Lactobacillus fermentum	/db_xref="taxon:1613" /mol_type="genomic DNA" /note="conserved RAPD fragment LF" /organism="Lactobacillus fermentum" /strain="LMG 6902"	BCCM LMG:LMG6902	Lactobacillus fermentum, Beijerinck 1901 AL
		/codon_start=1 /db_xref="UniProt/TrEMBL:Q9XDB8" /note="putative" /transl_table=11 /gene="ap1" /product="unknown" /protein_id=" <u>AAD41657</u> .1" /translation="TRPYIGSSSTNSVLNLAFGYNGTQRLFGQTTGTGGTFSGMGTKAS GKTKTGAPTGGTKQGGMKKPAGGAPTGQAKPTSGMKKGTTNGGKPSGKSLSSMPKTTKG RPTGMKGGAGQGGNGIFNVGTAGVTRLFQTALRRQISWFLPLALFGMVAAYLAAYDRKK KWWQTNRKQQDVLYWAGWLIPVGAFFSVANFFHPYYTIMLAPPLAVLAALAITSGLKRP		
<				

Validation of crossreferences

- A custom view including organism name, as specified by EMBL, and scientific name, as specified by CABRI, was created and data downloaded
- For each catalogue, three lists were created including:
 - Valid links: records having identical strain numbers and the same organisms/scientific name
 - **Dubious links**: records where organisms and name were different, but could be checked for synonyms of names, alternate names and previous names
 - **Invalid links**: records where names were different and clearly not coherent (e.g., bacteria versus human sequences)

For bacteria strains validity was further validated with data provided by P. Dawyndt (Univ. Ghent) Romano, Linking CABRI and EMBL,

Submission of crossreferences

- Validated links have been submitted to EBI and included in EMBL version 81
- Collections are invited to return validated data
- Analysis carried out every two EMBL versions

CBS Yeasts catalogue (huge records)	51103
CBS Filamentous fungi catalogue (huge records)	7116
DSMZ Bacteria catalogue	2618
CBS Filamentous fungi catalogue	2551
CBS Yeasts catalogue	1880
BCCM/LMG Bacteria catalogne	799
CIP Bacteria catalogne	546
NCIMB Bacteria catalogne	336
BCCM/MUCL Fungi and yeasts catalogue	119
DSMZ Fungi catalogne	31
BCCM/IHEM Fungi and yeasts catalogue	18
Total	67528

Links added to EMBL and set up at EBI

<u>File M</u> odifica <u>V</u> isualizza <u>P</u> referiti <u>S</u> trun	nenti <u>?</u>			//
	Pubmea	11931143		^
	Position	1-513		
2. 0	Goncalves,E.R.;	Rosato,Y.B.; 16S-23S	rDNA spacer regionUnpublished.	
	Position	1-513		
.د	Goncalves,E.R.; databases. CBM Campinas, SP 1	Kosato,Y.B.; Submitted IEG, UNICAMP, Cidade 3083-970, Brazil	Universitaria, Zeferino Vaz, P.O. Box 6109,	
1	Position	1-513		
Database Cross-	references			
CABRI	LMG 673.			
Features				
Key	Location	Qualifier	Value	
source	1513	db_xref	taxon:56450	
		mol_type	genomic DNA	
		organism	Xanthomonas cassavae	
		strain	LMG673	
misc rna	1513	note	16S-23S ribosomal RNA intergenic spacer	
Sequence				
Characteristics I	Length: 513 BP	, A Count:115, C Coun	t:121, G Count:144, T Count:133, Others Co	unt:0
Sequence	<pre>ggctggatca cctccttttg agcatgacgt cattgttcct gtcgggcgtc ctcacaaatt acctgcattc agagattcat accggcacag gtcggtatgc gaagtccctt catggggcct tagctcagct gggagagcac ctgctttgca agcagggggt cgtcggttcg atcccgacag gctccgccat agtgagtgaa aagacttcgg gtctgtagct cagtgggtta gagcgcgcc ctgataaggg tgaggtcggt agttcgagtc tacccagacc caccactctg aatgtagtgc acacttaaga atttatatgg ctcagccgtt gaggctgaga catgtcttt ataacttgg acgtagcgag cgtttgagat atctatctaa acgtgtcgtt gaggctaagg cgggacttca agtccctaaa taattgagtc gtatgtccgt tggtggcatt ggtaccccac acaacacgga ctgatatgga cctgaggcaa cttgggggtt ata //</pre>			60 120 180 240 300 360 420 480 513
Go to: <u>General</u> De	escription Refe	rences <u>Cross-referer</u> phts Reserved. <u>Terms of I</u>	i <mark>ces Features Sequence</mark> Use Feedback	

Some lessons learned

- Although SRS is a powerfull indexing and search engine, the EMBL size and complexity actually limits the possibility of performing complex searches on it
- Linking to EMBL is best done by using direct explicit links
- Since annotation of strain numbers is not coherent, an a priori identification of crossreferences is needed and it can be effectively carried out by using SRS
- A complete validation can only be achieved by manually analysing dubious data
- Comparison with crossreferences identified by using different methods is helpful for reducing manual intervention

Future perspectives

- The procedure should be further refined and made as automatic as possible
- Any effort should be done for requesting the insertion of strain numbers in a proper format and in a separate field
- The implementation of links from CABRI to EMBL should be done by adding a purpose information in those CABRI records that are actually cited in the EMBL Data Library
- This link can simply be implemented by searching the external references field of EMBL with the CABRI strain number through SRS

I wish to thank for their contributions...

National Cancer Research Institute (Genoa, IT) Federico Malusa, Domenico Marra, Francesca Piersigilli

LMG, Gent, BE Peter Dawyndt

BCCM, Brussels, BE François Guissart **DSMZ, Braunschweig, DE** Manfred Kracht

CBS, Utrecht, NL Gerrit Stegehuis

This work was partially supported by the EU within the projects Common Access to Biological Resources and Information (CABRI) and European Biological Resource Centers Network (EBCRN)