

Portals and workflows: Taverna Workbench

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



- Information and data integration in biology
- Web Services and workflow management
- Taverna Workbench
- A workflow enactment portal

∥ℒℼ

Information in biology: well known facts

- Biomedical research produces an increasing quantity of new data and new data types
 - EMBL size: ver 82 7,3% more vs ver 81 (3 months), 112,9% vs ver 74 (24 months)
- Emerging domains, like mutation and variation analysy, polymorphisms, metabolism, as well as new technologies, e.g., microarrays, will contribute with even huger amounts of data
- Analysis softwares must interoperate with databases
 - Databases as input for softwares
 - Results as new data to record and analyze

∥ℒℿ

Heterogeneicity of databanks

- A few dbs are managed in an almost homogenous way (e.g., sequences at EBI, NCBI, DDBJ)
- Secondary databases are of high quality (good and extended annotation, quality control)
- Many databases are highly specialized, e.g. by gene, organism, disease, mutation, etc...
- Many databanks are created by small groups or by single researchers
- Databanks are distributed:
 - Different DBMS, data structures, query methods
 - Different information, semantics

Goals of the integration

In this context, data integration and work automation are needed to:

- Carry out analysis and/or searches involving more databases and softwares automatically
- Perform analysis involving large data sets effectively
- Achieve a better and wider view of all available information
- Carry out a real data mining

Data integration longevity

• Integration needs stability

- Standardization.....
- Good domain knowledge
- Well defined data
- Well defined goals

Integration fears:

- Heterogeneicity of data and systems
- Uncertain domain knowledge
- Fast evolution of data
- Highly specialized data
- Lacking of predefined, clear goals
- Originality, experimentalism ("let me see if this works")

11871

Integration of biological information

In biology:

- A pre-analysis and reorganization of the data is very difficult, because data and related knowledge change very quickly
- Complexity of information makes it difficult to design data models which can be valid for different domains and over time
- Goals and needs of researchers evolve very quickly according to new theories and discoveries

Integration must therefore be carried out by using flexible systems that are easy to adapt and to extend

11 22 11

Integration methods

From syntactical to semantical links:

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

Web Services in brief

- XML based network services
- Implement standard transport protocols (SOAP, HTTP)
- Standards also available for retrieval and identification (UDDI), description (WSDL) and composition (WSFL)
- Allow software applications to access data "intelligently": identification of contents, interpretation of semantics information
- Metadata needed
- Web Services already implemented by many Institutes and service centers

1187

Workflow management

"A computerized facilitation or automation of a business process, in whole or part". (Workflow Management Coalition)

Main goal is:

 the implementation of data analysis processes in standardized environments

Main advantages relate to:

- **effectiveness**: being an automatic procedure, it frees bioscientists from repetitive interactions with the web and it supports good practice,
- reproducibility: analysis can be replicated over time,
- reusability: intermediate results can be reused,
- traceability: the workflow is carried out in a transparent analysis environment where data provenance can be

checked and/or controlled.

∥ℒℼ

Workflow management software

Workflow management software for bioinformatics applications:

- Biopipe, an add-on to bioperl
- GPipe, an extension of the Pise interface
- Taverna (EBI), a component of the myGrid platform
- Wildfire (Bioinformatics Institute, Singapore)
- Pipeline Pilot (SciTegic)
- BioWBI, Bioinformatic Workflow Builder Interface, from IBM

They all require knowledge of the systems and skills and time for development of the workflows.

Workflow management: Taverna

Taverna Workbench

- builds complex analysis workflows
- is able to access both remote and local processors
- defines alternative processors
- runs workflows
- visualizes the results

It includes a an ontology for bioinformatics data

Requirements: java, Windows or Linux Open source: http://taverna.sourceforge.net/ Current version: 1.3.1 (stable, next version 2.0) 15876

Taverna: GUI

Graphical user interface (GUI) including:

- Advanced Model Explorer (AME)
- Workflow diagram
- Available services
- Run workflow
- Enactor invocation

Options: default services, workflow editor, debug Extra features: FETA search engine ∥ℒℼ

Taverna: available processors (services)

WSDL services

Web Service Description Language (WSDL) file: adds WSDL based service nodes

Soaplab servers

Soaplab server: adds a list of soaplab provided services

BioMOBY registries

Moby Central repository: determines hosts and their services

Workflows

• XScufl definition file: either adds the workflow as a sub-workflow or imports processors

Biomart databases

Biomart data warehouse: allows for searching all available data sets

Local processors

Simple list/string processors, r/w, extra special remote links, constant values, beanshell scripts

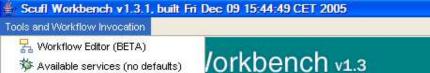
Other processors

Styx, Talisman, Seqhound, API Consumer

ℿℒℿ

	* Available services
erna Scufl Workbench v1.3	Search list
http://taverna.sf.net	Available Processors
djm,jf,mg,pl,ktg,mp,ms,tmo,mf,ek,pa, jb et al.	🔄 🦢 🗁 Local Services
actmen	
🔓 Advanced model explorer 💦 🕄 🔀	🖨 🗁 Local Java widgets
Workflow Object properties	ist ⊞ io
Load 😂 Load from web 🎱 Save 🔙 New subworkflow 🐲 🥅 Offline 🛛 Reset 💥	⊞ in metadata
	u missaus u mi
Workflow object Retries Delay Back Thre Critical	🖶 👘 💼 ui
Workflow model	🕀 🗁 💼 ncbi
Workflow outputs	🕀 💼 conditional
Processors	🕀 💼 net
	🛱 🗁 🦢 text
Control links	Byte[] to String
	String list union
Vorkflow diagram	Concatenate two strings
	String list difference
e as 📑 💽 🔬 🚍 📕 🕼 Configure diagram	Filter list of strings by regex
ə as 📑 😰 🚵 📮 📑 🕼 Configure diagram	
a as 📄 💽 🚵 📮 😫 Configure diagram	Split string into string list by regular expression
ə as 📄 😰 🕍 📑 📔 🕼 Configure diagram	
as 📄 💽 🛃 🚍 🕼 Configure diagram	Pad numeral with leading 0s Filter list of strings extracting match to a regex
e as 📄 🔊 🔊 📮 😰 Configure diagram	Split string into string list by regular expression Pad numeral with leading 0s Filter list of strings extracting match to a regex String list intersection
e as 📄 🔊 🔊 🗔 🥵 Configure diagram	Split string into string list by regular expression Pad numeral with leading 0s Filter list of strings extracting match to a regex String list intersection Diojava Diojava Diojava Diojava Diojava Diojava
e as 📄 💽 🛃 🚍 🕼 Configure diagram	Split string into string list by regular expression Pad numeral with leading 0s Filter list of strings extracting match to a regex String list intersection Diojava Diojava Diojava Diojava Diojava Diojava Diojava Diojava Diojava Diojava Diojava Diojava Diojava Diojava Diojava Diojava Diojava Diojava Diojava
e as 📄 🔊 🔊 🚍 🥵 Configure diagram	Split string into string list by regular expression Pad numeral with leading 0s Filter list of strings extracting match to a regex String list intersection Jobc Jobc Jobc Jobc Jobc Jobc Jobc Jobc Jobc Jobc Jobc
e as 📄 🔊 🔊 🗐 🥵 Configure diagram	Split string into string list by regular expression Pad numeral with leading 0s Filter list of strings extracting match to a regex String list intersection Group biojava Group base64 Group test Beanshell scripting host
e as 📄 🔊 🔐 🚍 🥵 Configure diagram	Split string into string list by regular expression Pad numeral with leading 0s Filter list of strings extracting match to a regex String list intersection Group biojava Group biojava Group base64 Group test Beanshell scripting host String Constant
e as 📄 🔊 🔐 📮 🖳 🕼 Configure diagram	Split string into string list by regular expression Pad numeral with leading 0s Filter list of strings extracting match to a regex String list intersection
e as 📄 🔊 🕼 📮 🥵 Configure diagram	Split string into string list by regular expression Pad numeral with leading 0s Filter list of strings extracting match to a regex String list intersection String list intersection String base64 String Constant String Constant String Constant String Constant String Constant String Mttp://www.ebi.ac.uk/~tmo/defaultMartRegist String WSDL @ http://www.ebi.ac.uk/collab/mygrid/service1/goviz/GoVi:
e as 📄 🔊 🕼 📮 🥵 Configure diagram	Split string into string list by regular expression Pad numeral with leading 0s Filter list of strings extracting match to a regex String list intersection Gravity biojava Gravity biojava Gravity Beanshell scripting host String Constant Biomart registry @ http://www.ebi.ac.uk/~tmo/defaultMartRegist Gravity Biological Scripting host String Constant String Constant String Constant String Constant String Constant String Constant String Constant String Constant Pad numeral with leading 0s String Constant String String Constant String String Constant String String St
e as 💽 🔊 🔊 🚰 🥵 Configure diagram	Split string into string list by regular expression Pad numeral with leading 0s Filter list of strings extracting match to a regex String list intersection String list intersection String list intersection String biojava String base64 String Constant String String Constant String String Stri
e as 💽 🔊 🔊 🖓 🥵 Configure diagram	Split string into string list by regular expression Pad numeral with leading 0s Filter list of strings extracting match to a regex String list intersection String list intersection Jobc Jobc Desce Beanshell scripting host String Constant String Constant WSDL @ http://www.ebi.ac.uk/~tmo/defaultMartRegist WSDL @ http://www.ebi.ac.uk/collab/mygrid/service1/goviz/GoViz Porttype: GoViz [RPC]

Rendering done.



lorer

Save 🔡

📝 Configure diagram

ies

http://taverna.sf.net

actmon

djm,jf,mg,pl,ktg,mp,ms,tmo,mf,ek,pa, jb et al.

New subworkflow 🏪

Delay

Retries

Offline

Thre...

Back...

Reset 💢

Critical

🔚 Advanced model explorer

Dot text

🎋 Available services

🛃 Workflow diagram

💄 Workflow diagram

Save as

Run workflow

DEBUG - Workflow XML preview

DEBUG - Enactor event monitor

Workflow object

Workflow model

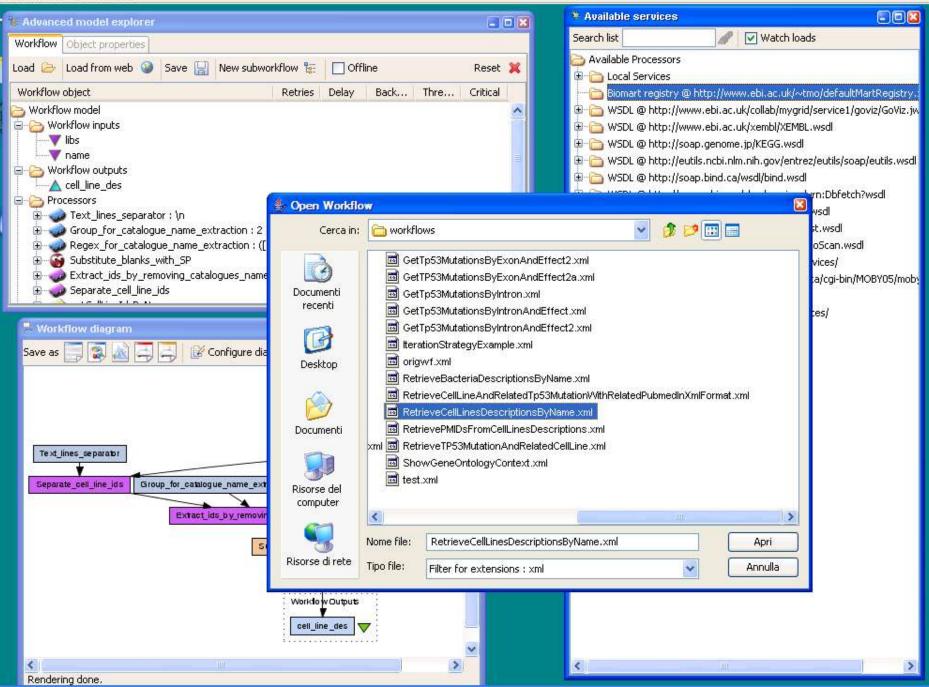
Processors

🛅 Data links 🛅 Control links

🛅 Workflow inputs

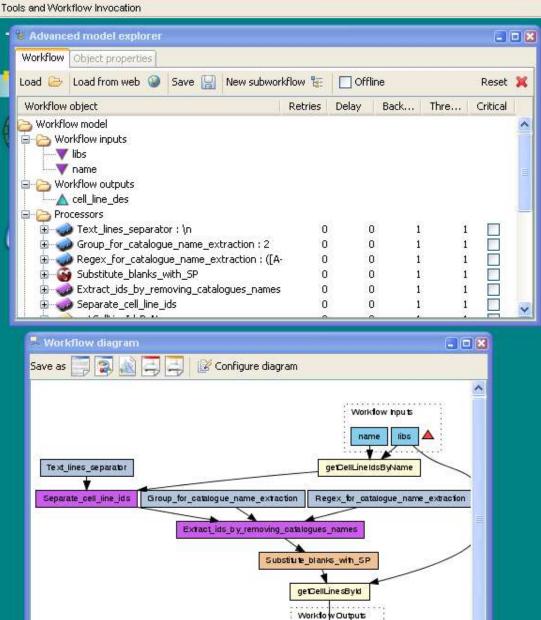
Workflow outputs





<

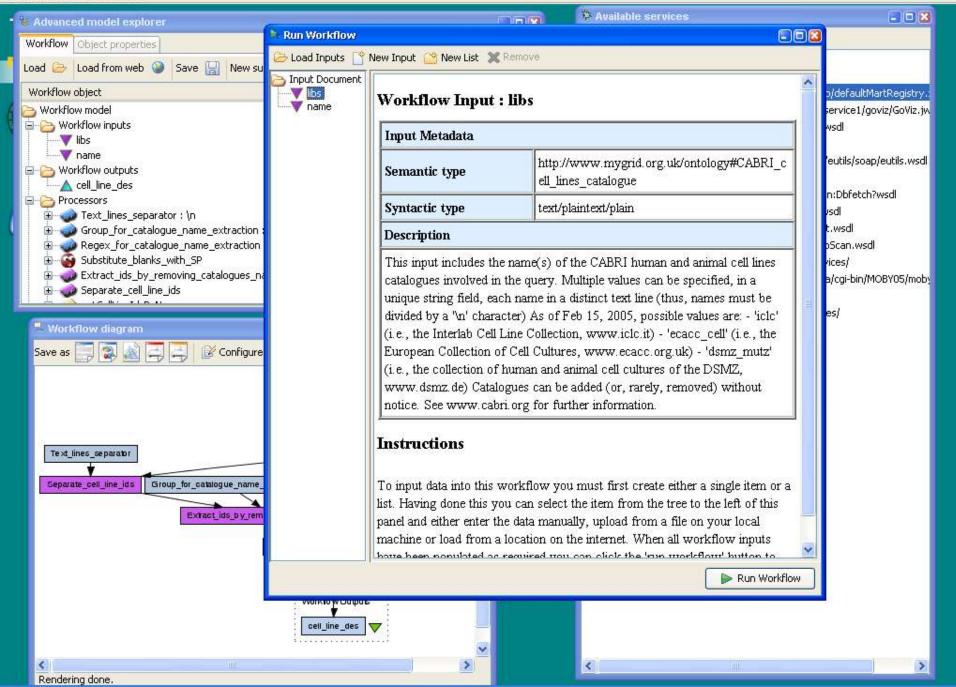
Rendering done.

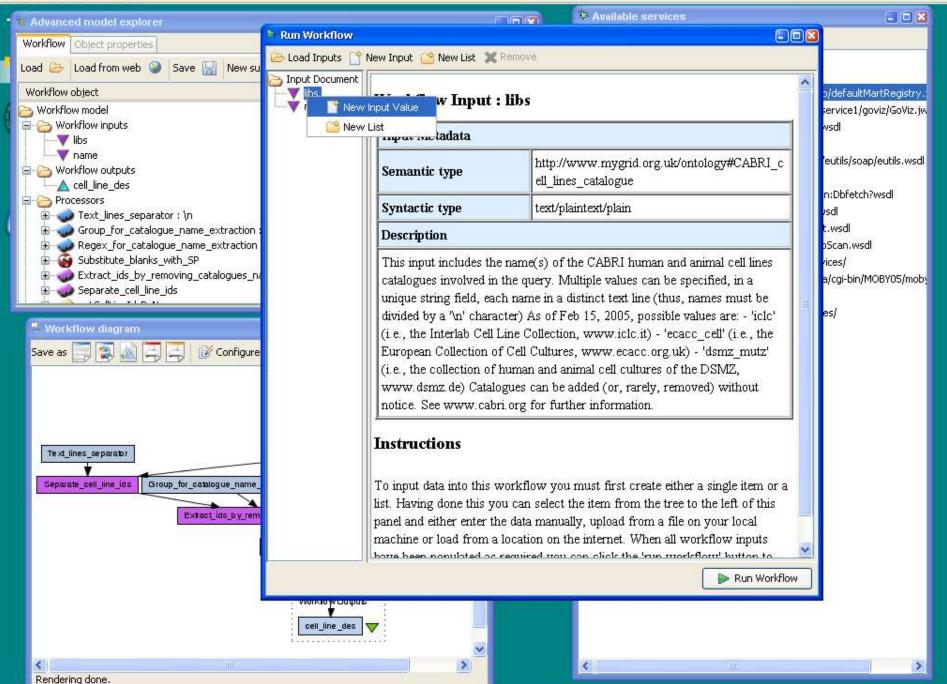


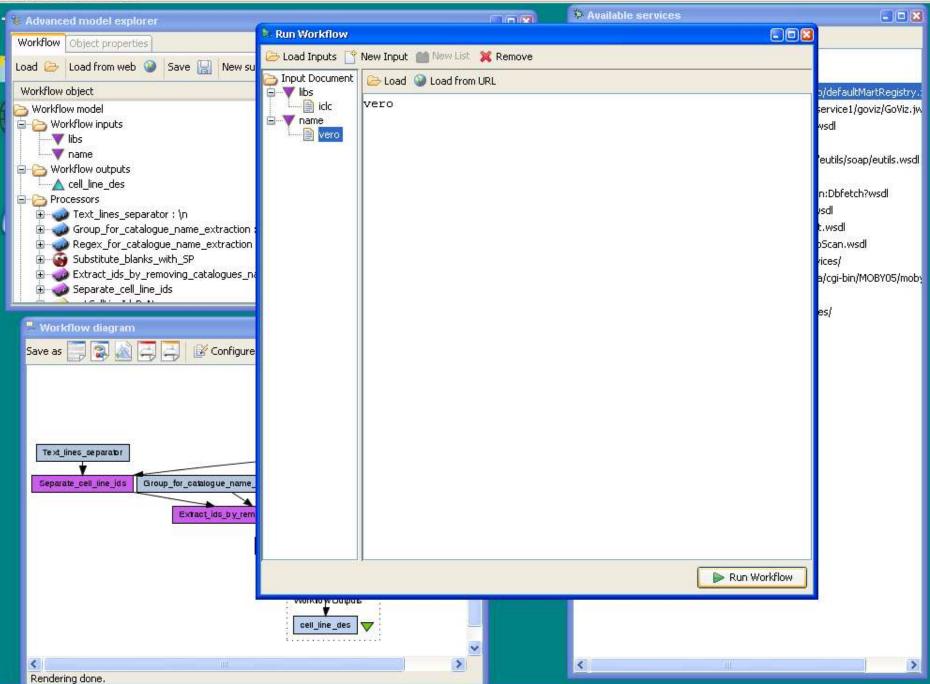
cell_line_des

>

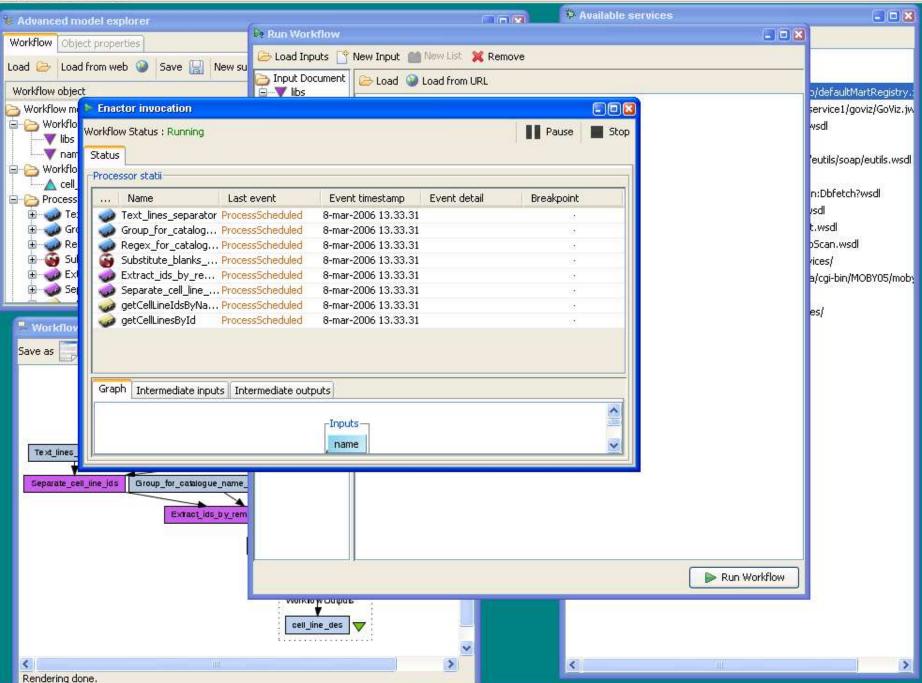
* Available services
Search list 🛛 🖉 🔽 Watch loads
🔁 Available Processors
🗄 🛅 Local Services
🖳 🛅 Biomart registry @ http://www.ebi.ac.uk/~tmo/defaultMartRegistry.
🐵 🛅 WSDL @ http://www.ebi.ac.uk/collab/mygrid/service1/goviz/GoViz.jw
😟 🛅 WSDL @ http://www.ebi.ac.uk/xembl/XEMBL.wsdl
🐵 🛅 WSDL @ http://soap.genome.jp/KEGG.wsdl
🐵 🛅 WSDL @ http://eutils.ncbi.nlm.nih.gov/entrez/eutils/soap/eutils.wsdl
🐵 🛅 WSDL @ http://soap.bind.ca/wsdl/bind.wsdl
🐵 🛅 WSDL @ http://www.ebi.ac.uk/ws/services/urn:Dbfetch?wsdl
🕀 🛅 WSDL @ http://www.ebi.ac.uk/ws/WSFasta.wsdl
🕀 🛅 WSDL @ http://www.ebi.ac.uk/ws/WSWUBlast.wsdl
🕸 🛅 WSDL @ http://www.ebi.ac.uk/ws/WSInterProScan.wsdl
🕀 🛅 Soaplab @ http://www.ebi.ac.uk/soaplab/services/
🕀 🛅 Biomoby @ http://mobycentral.icapture.ubc.ca/cgi-bin/MOBY05/mob
🕀 🛅 SeqHound @ seqhound.blueprint.org
🗄 📹 Soaplab @ http://www.o2i.it:8080/axis/services/



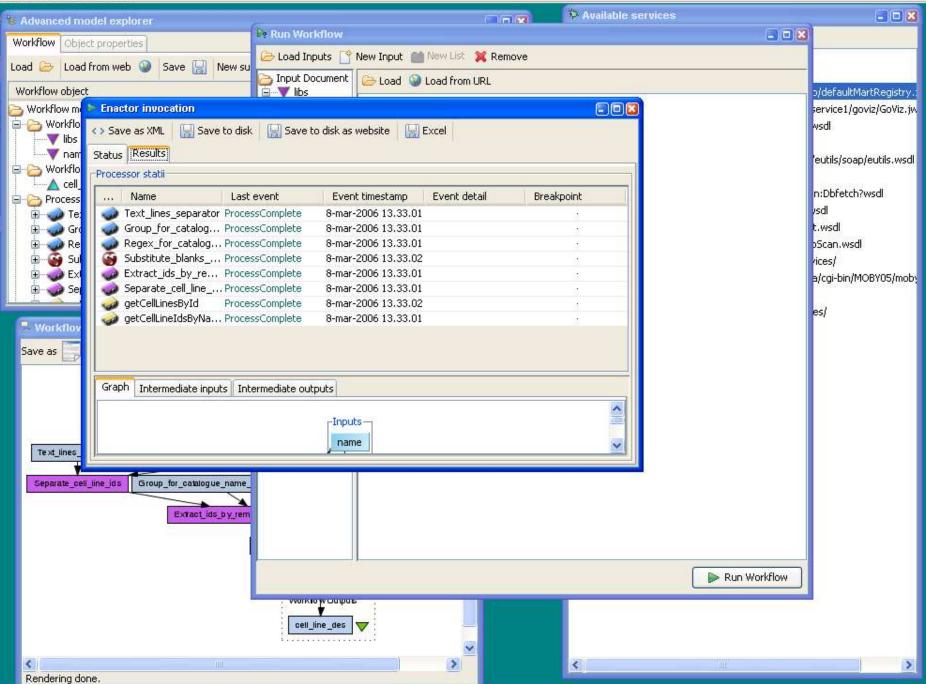




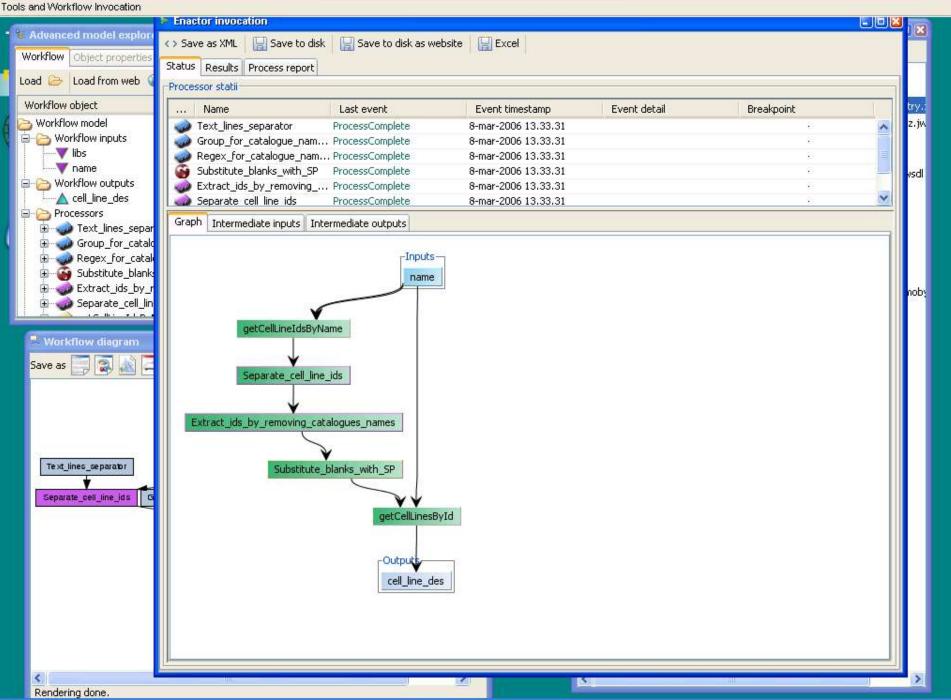
Tools and Workflow Invocation

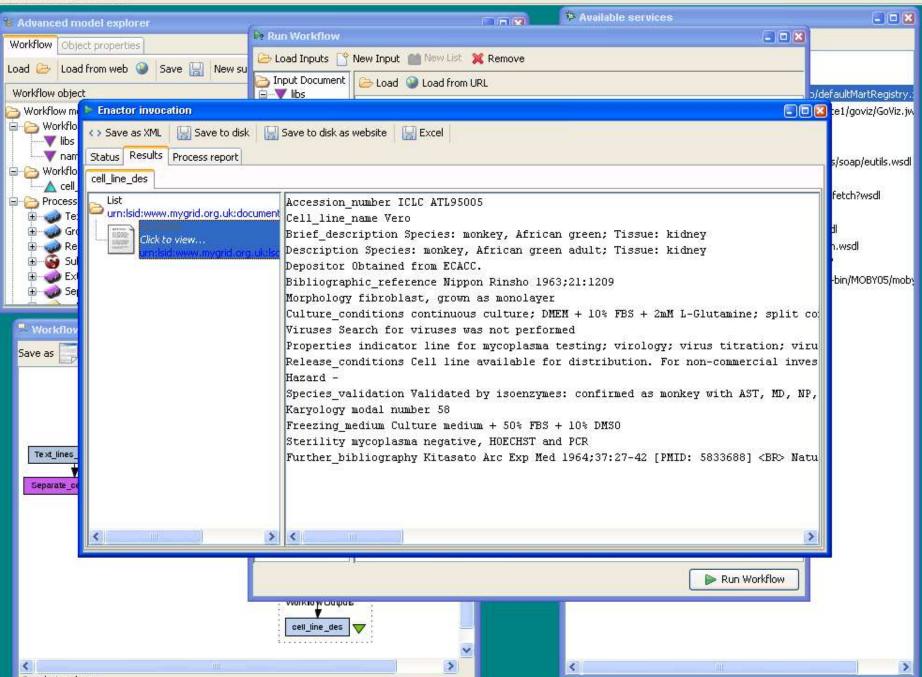


- 72











and Workflow Invocation						
	Enactor invocation					
Advanced model explore	Save as XML Save to dis	sk 🛛 🔛 Save to disk as websit	e 🔚 Excel			
/orkflow Object properties	Status Results Process report					
oad 📴 🛛 Load from web 🧃	Processor statii					
Uerlifleur abiaet	7530		N NS 5360 50	N 22 2742-025320		
Workflow object	Name	Last event	Event timestamp	Event detail	Breakpoint	<u> </u>
Workflow model	Regex_tor_catalogue_nam		8-mar-2006 13.33.31		1. State 1.	~
Workflow inputs	Substitute_blanks_with_SP Extract_ids_by_removing		8-mar-2006 13.33.31 8-mar-2006 13.33.31		K.	
name	Separate_cell_line_ids	ProcessComplete	8-mar-2006 13.33.31		100 - 100 -	
- Workflow outputs	getCellLineIdsByName	ProcessComplete	8-mar-2006 13.33.31		2. *:	
▲ cell_line_des	getCellLinesById	ProcessComplete	8-mar-2006 13.33.32		· ·	~
Processors					<i>K</i> .	
🚡 🥥 Text_lines_separ	r Graph Intermediate inputs In	termediate outputs				
🗄 🥥 Group_for_catalo	result					
🗄 🧼 Regex_for_catal						
🗄 🎯 Substitute_blank			ICLC:ICLC ATL95005			
🗄 🧼 Extract_ids_by_r	urn:lsid:www.myarid.org.	uk:Isdocument:7Y4DDO713113				
🗄 🧼 Separate_cell_lin	1					
📥 🗻 ie lie et e e	4					
The second second second second						
Save as 🧾 🕵 🚵 🚄	4					
Text_lines_separator						
Separate_cell_line_ids G						
Separate_sen_ine_ius	4					
	J					
<						

Watch loads

jp/KEGG.wsdl

wsdl/bind.wsdl

k/ws/WSFasta.wsdl

k/ws/WSWUBlast.wsdl

.uk/soaplab/services/

3080/axis/services/

int.org

>

k/ws/WSInterProScan.wsdl

.nih.gov/entrez/eutils/soap/eutils.wsdl

l.icapture.ubc.ca/cgi-bin/MOBY05/moby

k/ws/services/urn:Dbfetch?wsdl

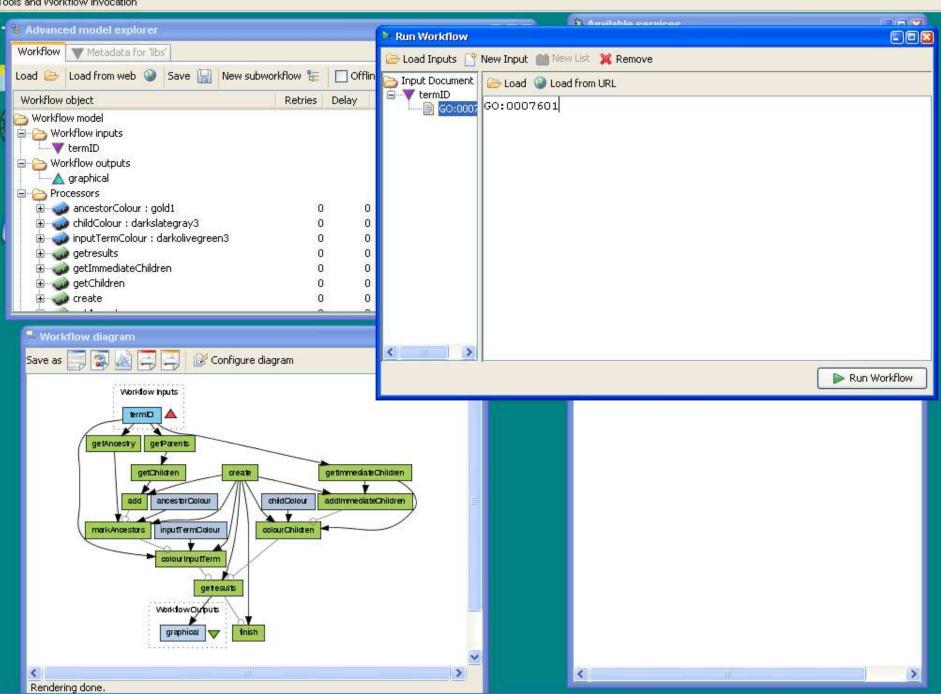
Scull Workbench v1.3.1, built Fri Dec 09 15:44:49 CET 2005 **Tools and Workflow Invocation** 🍄 Available services Advanced model explorer Search list Workflow Wetadata for 'libs' Available Processors Load 🗁 🛛 Load from web 🕥 🛛 Save 🔛 New subworkflow 🏣 Offline Reset 💢 🗄 🛅 Local Services Critical Workflow object Delay Back ... Thre... Biomart registry @ http://www.ebi.ac.uk/~tmo/defaultMartRegistry.. Retries WSDL @ http://www.ebi.ac.uk/collab/mygrid/service1/goviz/GoViz.jw Workflow model ÷---Contraction inputs WEDL @ http://www.obi.ac.uk/xembl/XEMBL.wsdl 🛅 Workflow outputs **Open Workflow** 흨 Processors ¥ 🤌 📂 📖 🚍 workflows Cerca in: 🛅 Data links 🛅 Control links GetTp53MutationsByExonAndEffect2.xml 0 GetTP53MutationsByExonAndEffect2a.xml GetTp53MutationsByIntron.xml Documenti recenti GetTp53MutationsByIntronAndEffect.xml GetTp53MutationsByIntronAndEffect2.xml B IterationStrategyExample.xml 🔟 origwf.xml Desktop RetrieveBacteriaDescriptionsByName.xml 💄 Workflow diagram RetrieveCellLineAndRelatedTp53MutationWithRelatedPubmedInXmlFormat.xml RetrieveCellLinesDescriptionsByName.xml Save as 🧾 🕵 an 4 RetrievePMIDsFromCellLinesDescriptions.xml Documenti xml 🖾 RetrieveTP53MutationAndRelatedCellLine.xml ShowGeneOntologyContext.xml Risorse del computer < Nome file: ShowGeneOntologyContext.xml Apri Risorse di rete Tipo file: Annulla Filter for extensions : xml

>

<

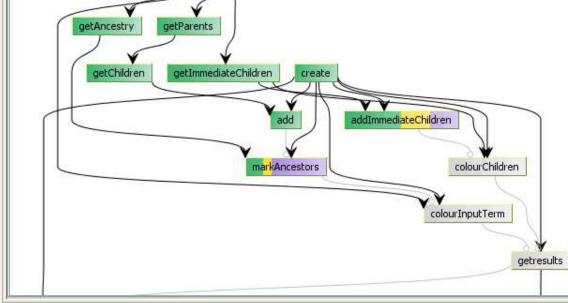
Rendering done.

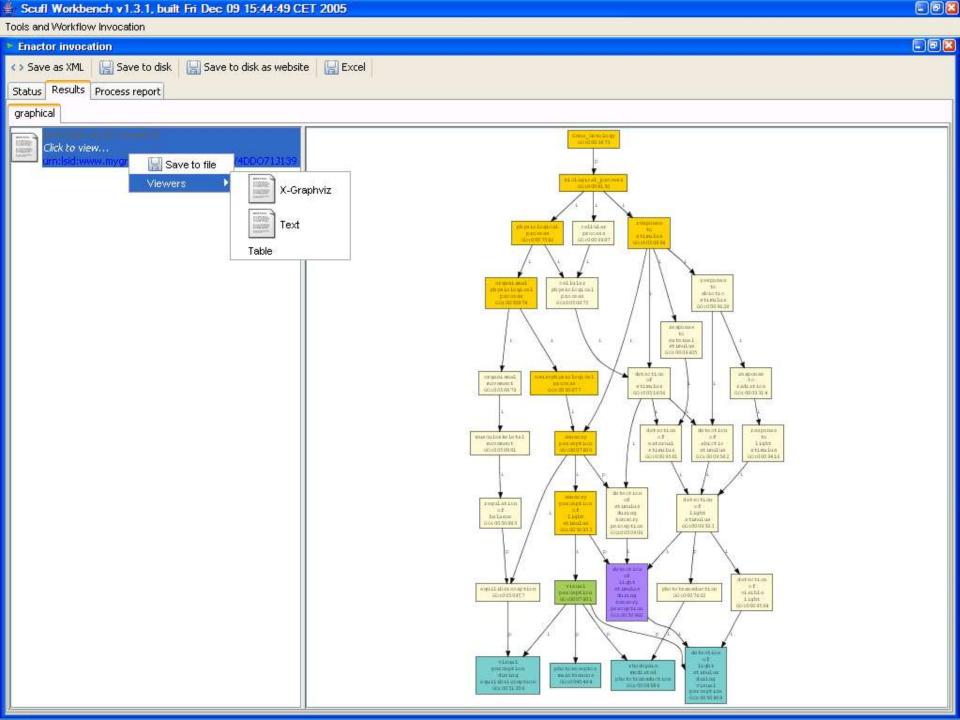
Tools and Workflow Invocation



nactor invocation					. 0
kflow Status : Running				Pause	Sto
atus					
ocessor statii					
Name	Last event	Event timestamp	Event detail	Breakpoint	-
ancestorColour	ProcessComplete	8-mar-2006 14.20.21		1 0	1
🍺 childColour	ProcessComplete	8-mar-2006 14.20.21		N	
inputTermColour	ProcessComplete	8-mar-2006 14.20.21		N:	
🧼 getresults	ProcessScheduled	8-mar-2006 14.20.21		N	
🍺 getImmediateChildren	ProcessComplete	8-mar-2006 14.20.24		N	
🍺 getChildren	ProcessComplete	8-mar-2006 14.20.37		N)	
🌶 create	ProcessComplete	8-mar-2006 14.20.24		193 193	
🤌 getAncestry	ProcessComplete	8-mar-2006 14.20.24		1 9	_
🌶 finish	ProcessScheduled	8-mar-2006 14.20.21		N	
markAncestors	InvokingWithIteration	8-mar-2006 14.20.38	IterationNumber='3' IterationTotal='9' ActiveWorkers		
🍺 add	ProcessComplete	8-mar-2006 14.20.38		**	3

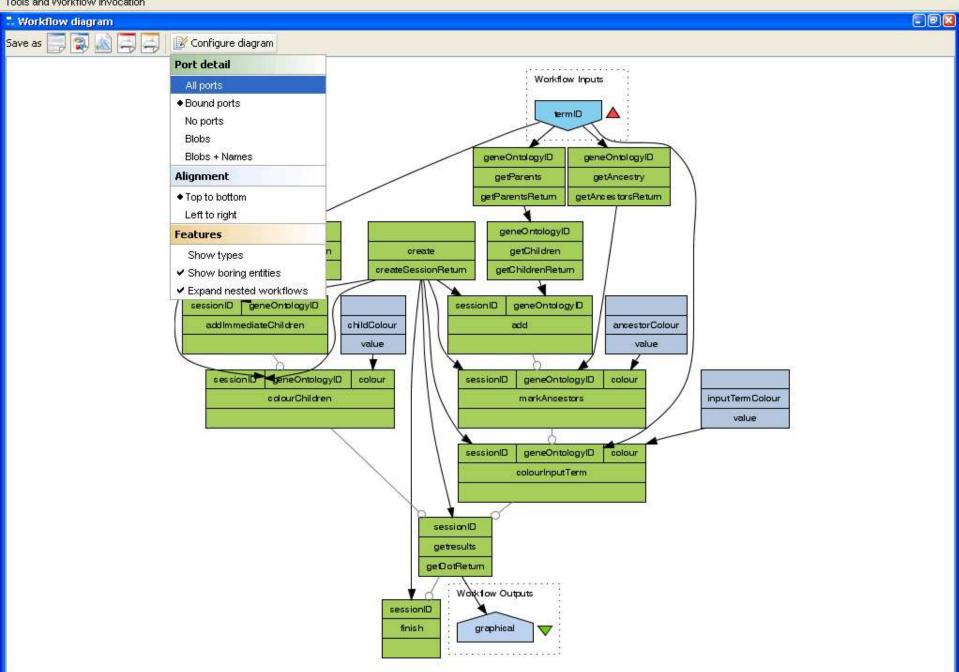
¥

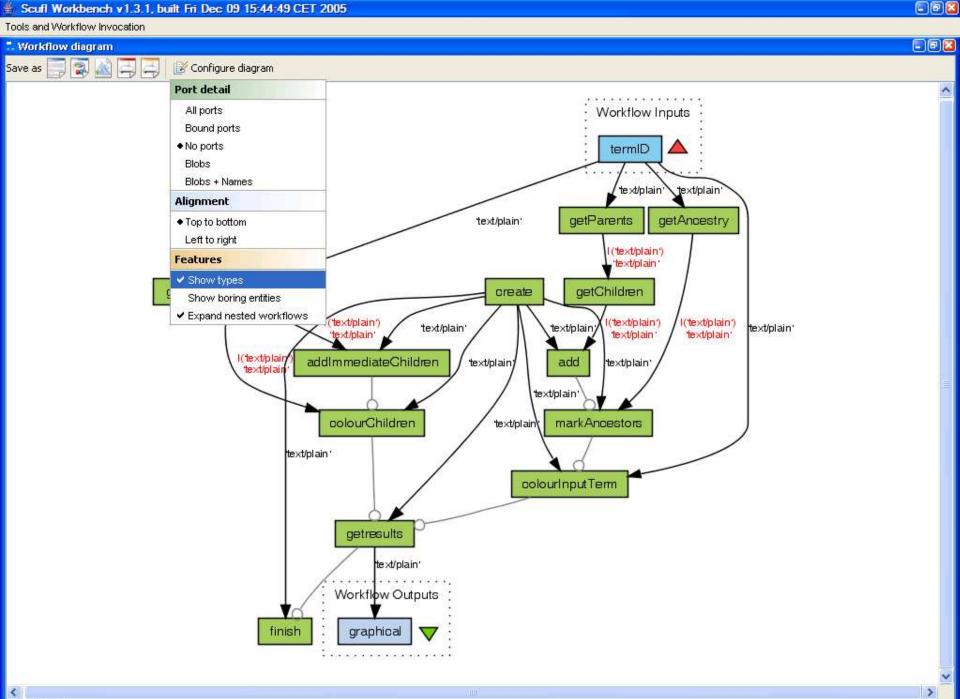


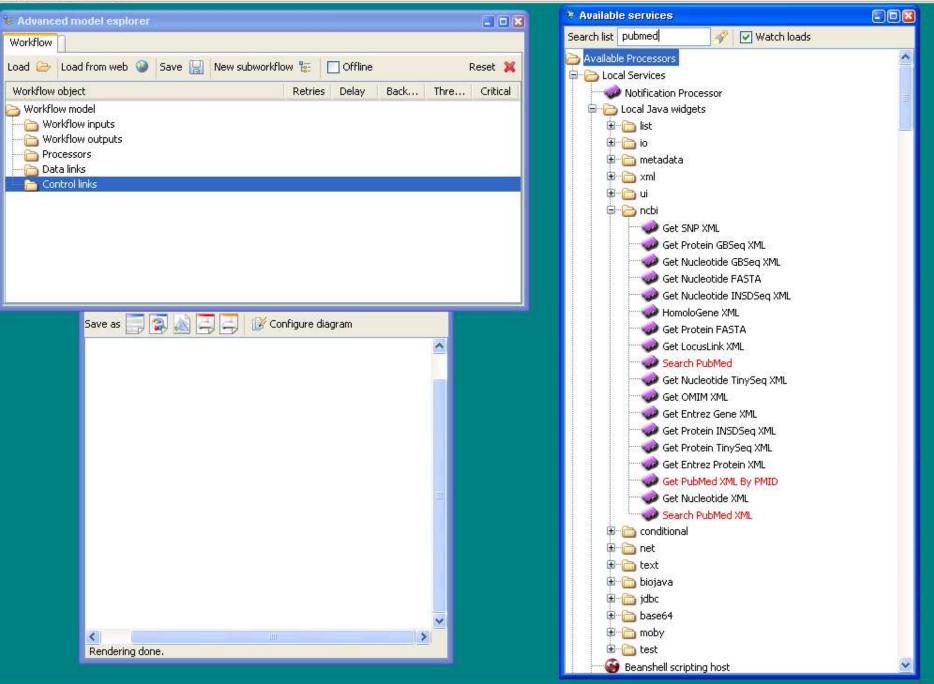


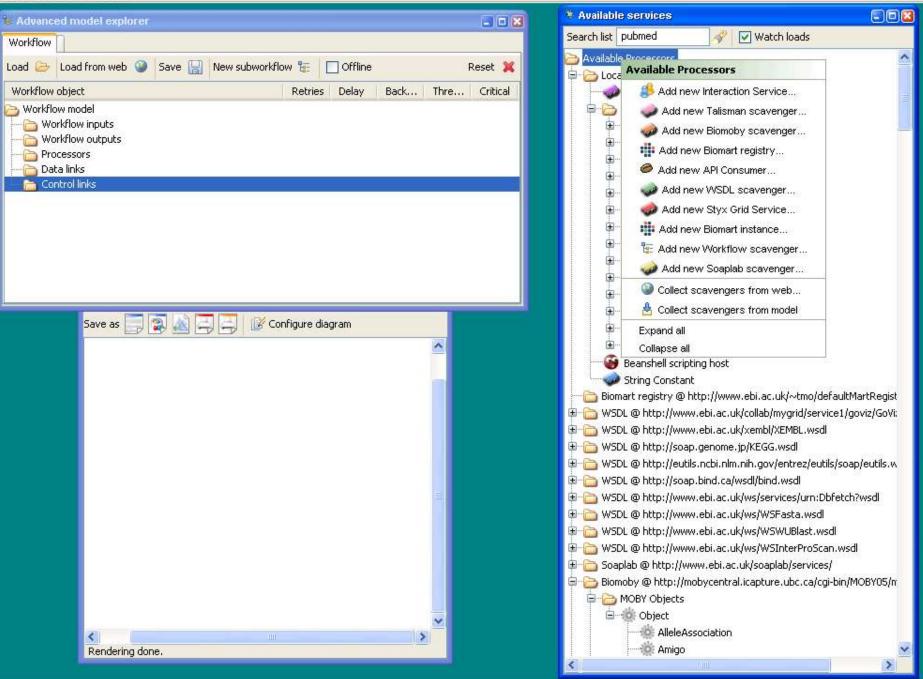
Tools and worknow invocation				
🕺 Workflow diagram				
Save as 📑 😰 🚵 🚍	📝 Configure diagram	_0		
Worklow hputs	Port detail			
getAncesty geParent getChildren add ances b	Blobs Blobs + Names Alignment Top to bottom	mediateChildren		
mark Ancestors Tinputt	Left to right			
colour	Features			
	Show types Show boring entities			
Workdo	w Kerner			
	hical. 🗸 İnish			

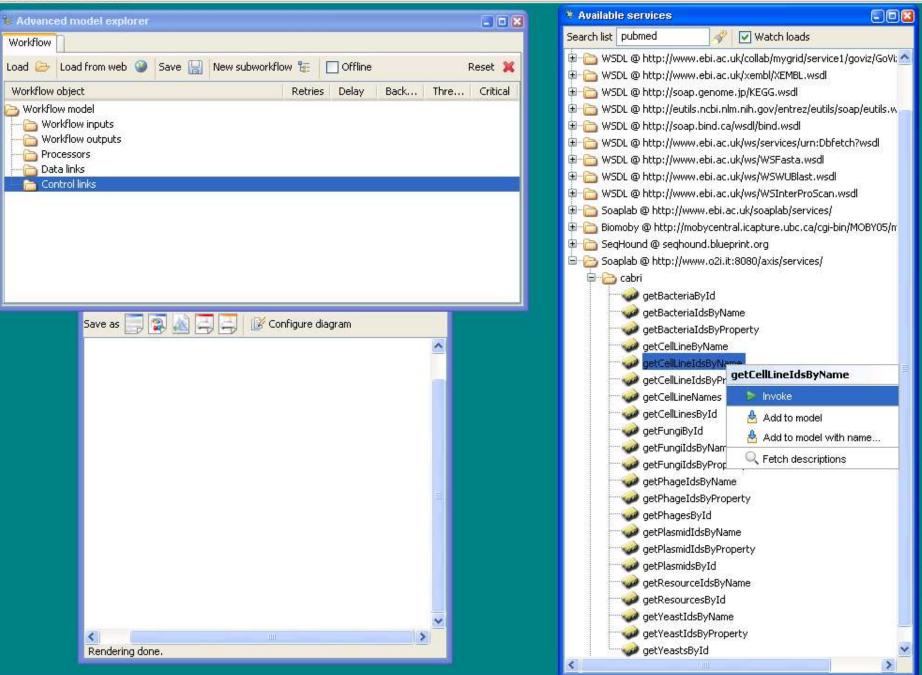






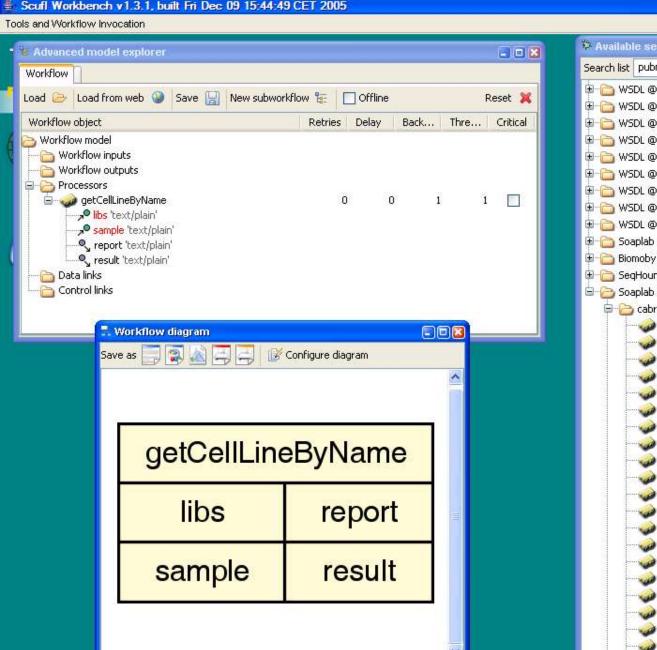




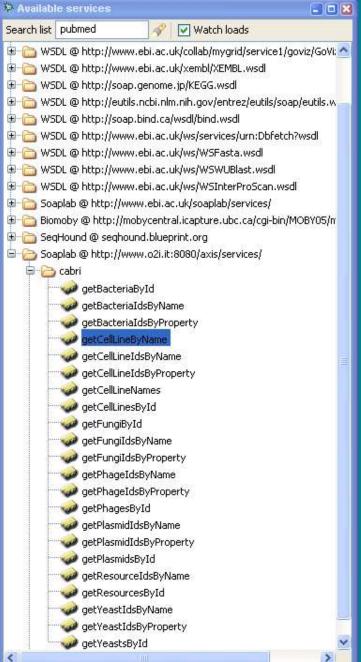


<

Rendering done.

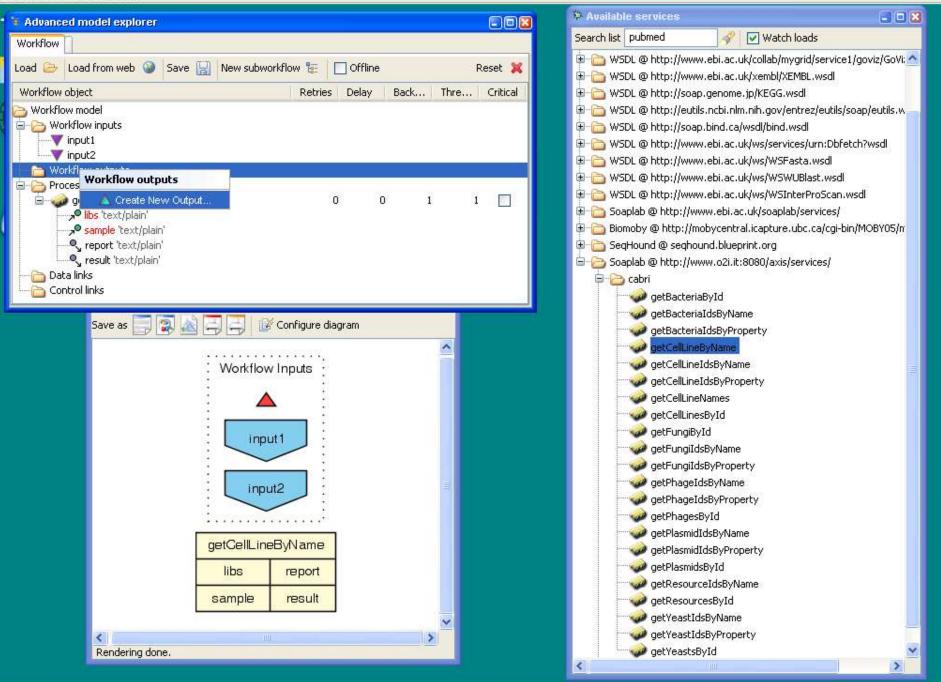


>



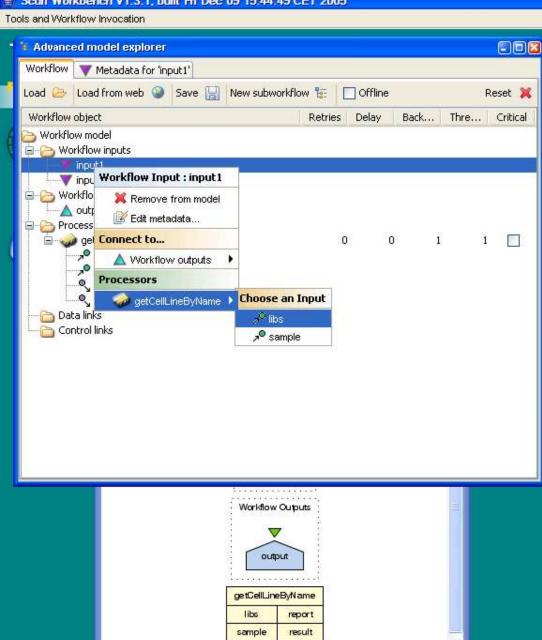
- 7

Tools and Workflow Invocation

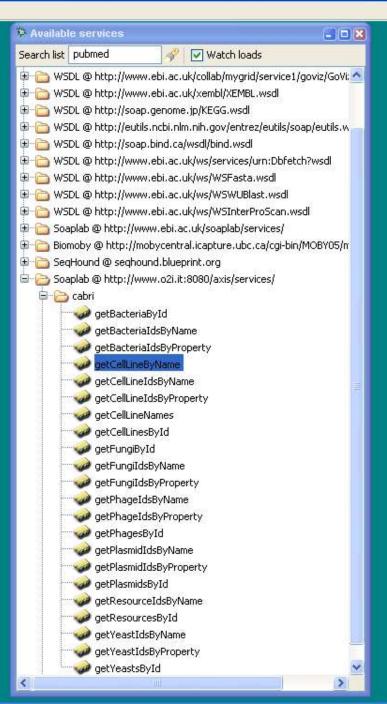


<

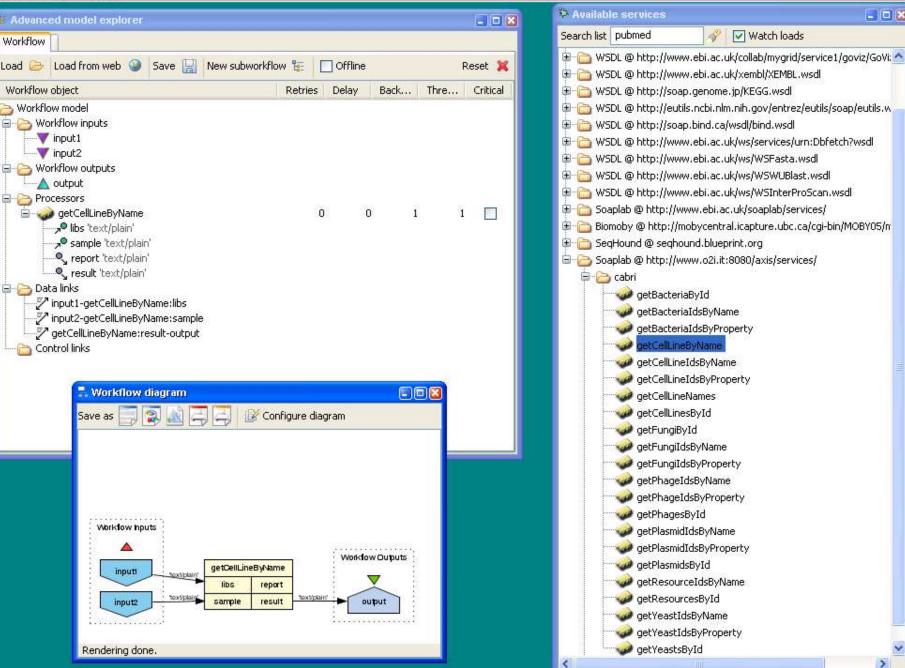
Rendering done.



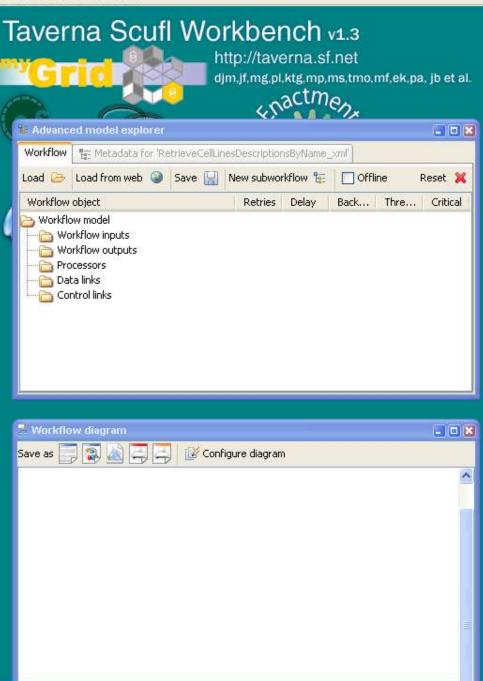
>

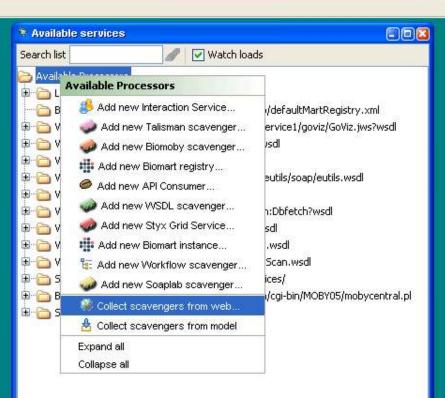


Tools and Workflow Invocation 🐮 Advanced model explorer Workflow Workflow object



Tools and Workflow Invocation



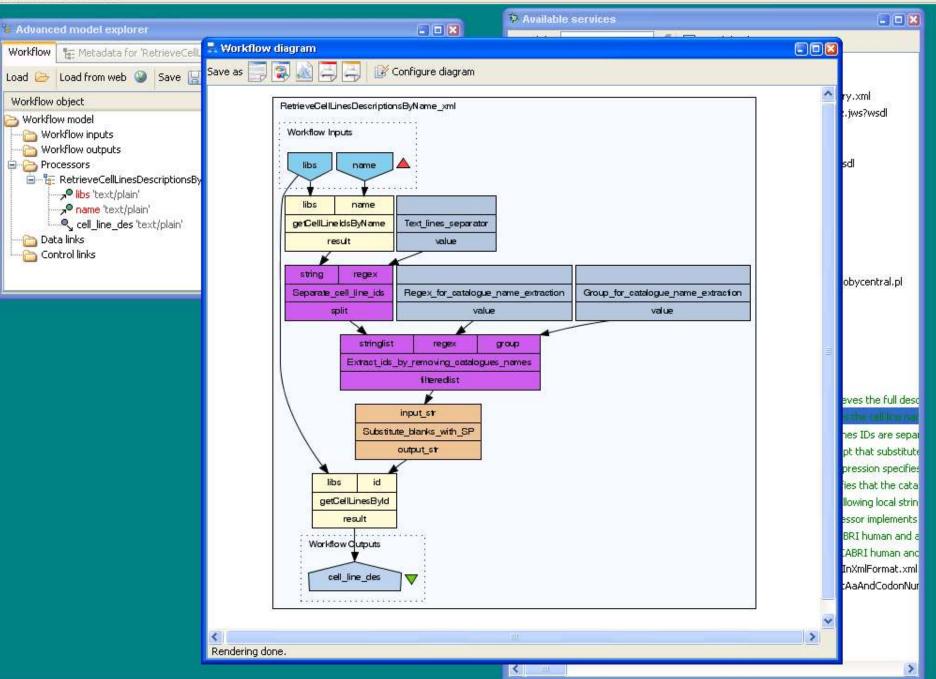


Tools and Workflow Invocation



Available services		
earch list	🥒 🗹 Watch loads	
Available Processors Local Services Biomart registry @ http WSDL @ http://www.el WSDL @ http://www.el WSDL @ http://soap.ge WSDL @ http://soap.ge WSDL @ http://soap.bi	://www.ebi.ac.uk/~tmo/defaultM bi.ac.uk/collab/mygrid/service1/go bi.ac.uk/xembl/XEMBL.wsdl enome.jp/KEGG.wsdl cbi.nlm.nih.gov/entrez/eutils/soap nd.ca/wsdl/bind.wsdl bi.ac.uk/ws/services/urn:Dbfetch?	viz/GoViz.jws?wsdl /eutils.wsdl
Soaplab @ http://www Soaplab @ http://www Biomoby @ http://moby SeqHound @ seqhound Web crawl @ http://ww The GetMutationsByE	bi.ac.uk/ws/WSInterProScan.wsdl .ebi.ac.uk/soaplab/services/ vcentral.icapture.ubc.ca/cgi-bin/Mi .blueprint.org w.o2i.it/workflows/index2.html	OBY05/mobycentral.pl
CetMutationsByI CetMutationsByI CetMutationsByI CetTP53Mutation CetTP53Mutation CetTeveBacterial CetTeveCellLines CetTeveCellLines CetTeveCellLines	Example 2 Examp	w retrieves the full desc w takes the cell line nar t cell lines IDs are separ
Group_for_cal Separate_cell Extract_ids_b getCellLinesBy getCellLineIds to the RetrieveCellLineAnd to the RetrieveCellLineAnd	talogue_name_extraction - This re talogue_name_extraction - This de line_ids - This processor implemen y_removing_catalogues_names - 1 rId - Get cell line descriptions by Id ByName - Get cell lines' IDs by nan dRelatedTp53MutationWithRelated unctionEntriesAndTP53CellLineEntr unctionEntriesByExon.xml -	ata specifies that the cata its the following local strin This processor implements from CABRI human and a ne from CABRI human and dPubmedInXmlFormat.xml

Tools and Workflow Invocation



Tools and Workflow Invocation

Т

		🎋 Available services	
Advanced model explorer		Search list	/ Watch loads
Workflow TetrieveCellLinesDescriptionsByName_xml		Available Processors	
Load 🗁 Load from web 🥹 Save 🔚 New subworkflow 🏣 🔲 Offline	Reset 💢	🗊 🛅 Local Services	
Workflow object Retries Delay Back Thre	. Critical		://www.ebi.ac.uk/~tmo/defaultMartRegistry.xml
 Workflow model Workflow inputs Workflow outputs Processors RetrieveCellLinesDescriptionsByNam 0 1 <li< th=""><th>1</th><th> Image: Solution of the second state of the second st</th><th>enome.jp/KEGG.wsdl hcbi.nlm.nih.gov/entrez/eutils/soap/eutils.wsdl hd.ca/wsdl/bind.wsdl hi.ac.uk/ws/services/urn:Dbfetch?wsdl hi.ac.uk/ws/WSFasta.wsdl hi.ac.uk/ws/WSFuBlast.wsdl hi.ac.uk/ws/WSInterProScan.wsdl</th></li<>	1	 Image: Solution of the second state of the second st	enome.jp/KEGG.wsdl hcbi.nlm.nih.gov/entrez/eutils/soap/eutils.wsdl hd.ca/wsdl/bind.wsdl hi.ac.uk/ws/services/urn:Dbfetch?wsdl hi.ac.uk/ws/WSFasta.wsdl hi.ac.uk/ws/WSFuBlast.wsdl hi.ac.uk/ws/WSInterProScan.wsdl
Save as Ref		SeqHound @ seqhound. SeqHound @ seqhound. Web crawl @ http://www GetMutationsByExor GetMutationsByIntro	ww.o2i.it/workflows/index2.html on.xml - rronAndEffect.xml - rron.xml - ByIntronAndEffect2.xml - escriptionsByName.xml - This workflow retrieves the full de
RetrieveCellLinesDescriptionsByName_xml libs cell_line_des name		Substitute_blar Regex_for_cat Group_for_cat Separate_cell_ Extract_ids_by getCellLinesBy1 getCellLineIdsB	anks_with_SP - This is a trivial beanshell script that substitu atalogue_name_extraction - This regular expression specifi italogue_name_extraction - This data specifies that the cal _line_ids - This processor implements the following local str y_removing_catalogues_names - This processor implement yId - Get cell line descriptions by Id from CABRI human and ByName - Get cell lines' IDs by name from CABRI human and dRelatedTp53MutationWithRelatedPubmedInXmlFormat.xr unctionEntriesAndTP53CellLineEntriesByMutAaAndCodonN
Rendering done.			

<

🐇 Scufl V	/orkbench v1.3.1, built Fri Dec 09 15:44:49 CET 2005	
Tools and V	/orkflow Invocation	
Advance	d model explorer	
Workflow	V Metadata for 'libs'	
-		
Ontology	Description MIME Types	
-Pick from	ontology	
	Find from regex :	
🗅 Availa	ble ontologies :	
	oot:Types	
	task	
	() filtering	
	etrieving	
	grouping	
	() manipulating	
	Calculating	
	erging	
	o parsing	
	ioining	
	e summarising	
	© removing	
	e splitting	
	O inserting O translating	
	e cransiaulig e searching	
	isplaying	
	e distinguishing	
	E · () aligning	
	bioinformatics_concept	
	O bioinformatics_application	
	🗄 🙆 bioinformatics_data	
	🗄 🙆 bioinformatics_metadata	
	🗄 🙆 bioinformatics_database	
	B obioinformatics_algorithm	
	B ioinformatics_diagram	
	🕒 3D_plot_of_Gene_Ontology_lattice	
	🕒 nucleotide_sequence_feature_diagram	
	ABI_graph_plot	
	e sequence_alignment_dot_plot	
	2D_alignment_quality_graph_plot	
Select fro	n ontology or manually edit term below	
	ww.mygrid.org.uk/ontology#CABRI_cell_lines_catalogue	
Tuch III	uuuu ausua ausua ausua ausua _con _nuos_ca canadao	

- 7

Tools and Workflow Invocation

Advance	d model explorer	
Workflow	Remote resource usage	
Save HTML	description 🔛	

Workflow information

This report applies to the workflow titled Search CABRI human and animal cell lines catalogues by cell line name and retrieve full cell line descriptions' authored by Paolo Romano, IST, Genova, Italy (paolo.romano@istge.it)' and with LSID 'urn:lsid:www.mygrid.org.uk:operation:US38XCRXJA0'. The textual description, if any is shown below :

This workflow takes the cell line name and the catalogue(s) name(s) as input and retrieve the full cell line description(s) by first retrieving the cell lines' unique IDs associated with the input (done via a call to the getCellLineIdsByName web service) and then using IDs for retrieving the full cell lines descriptions (done via a call to the getCellLineIdsByName web service) and then using IDs for retrieving the full cell lines descriptions (done via a call to the getCellLineIdsByName web service) and then using IDs for retrieving the full cell lines descriptions (done via a call to the getCellLineIdsByName web service).

Resource usage report

This display shows the various external resources used by the current workflow. It does not show resources such as local operations or string constants which are run within the enactment engine. Services are categorized by resource host and type, and the name of the instance of each service shown to the right.

Resources on www.o2i.it, 2 instances.					
Soaplab	Service rooted at /axis/services				
	App category and name Processors				
	/axis/services/cabri::getCellLineIdsByName getCellLineIdsByNa				
	/axis/services/cabri::getCellLinesById getCellLinesById				

Oncology over Internet (O₂I)

We designed a web system that:

- allows for the carrying out of a set of predefined workflows (of oncology interest)
- supports workflows annotation by using a simple ontology for bioinformatics processors (domain, task, i/o)
- implements search of workflows on the basis of their annotation
- supports retrieval of workflows based on users' registration and profiling
- allows storing and retrieval of workflows' executions and related results

11871

Oncology over Internet (O₂I)

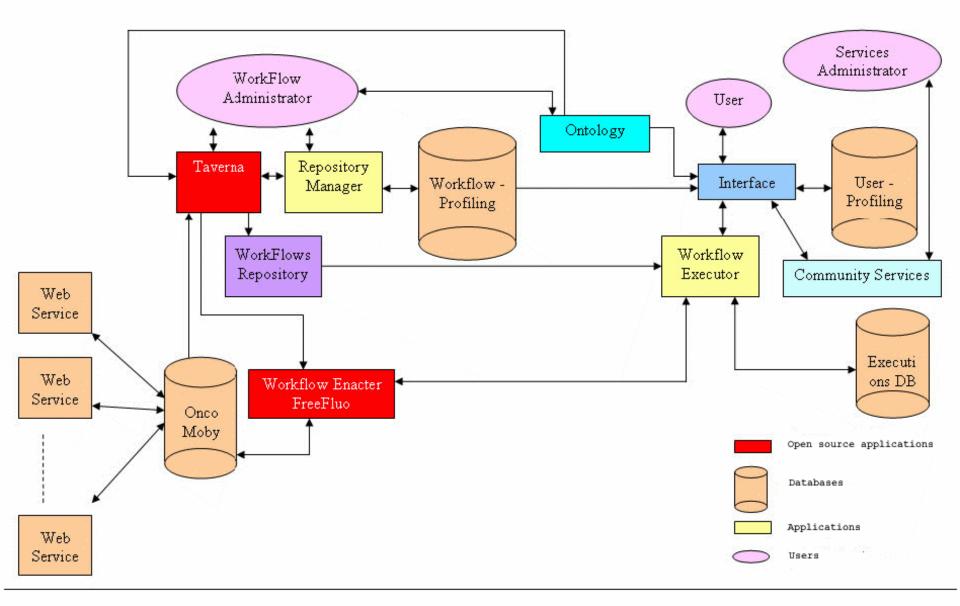
We designed a web system that:

- makes access to and retrieves data from Web Services and registries of Web Services
- stores workflows using the Simple conceptual unified flow language (Scufl) format
- is partially based on open source tools (Taverna WB, FreeFluo and mySQL)

Prototype available on-line: http://www.o2i.it:8080/portal/

11 22 11

O₂I architecture



ĪST

Predefined workflows

Workflows are:

- created by internal staff using Taverna
- stored in Scufl format
- maintained (workflow vs version)
- submitted by:
 - users
 - service providers

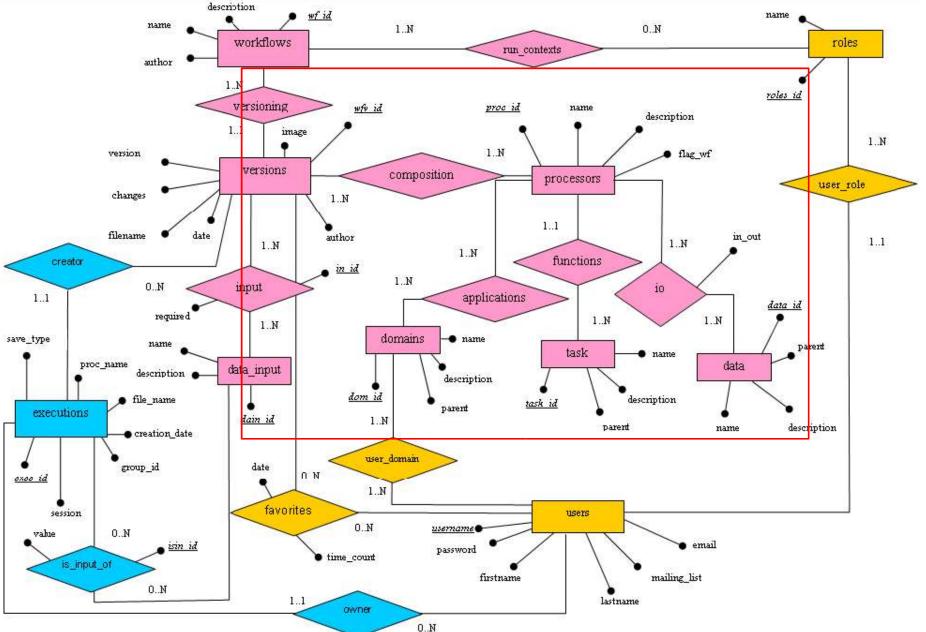
∥ℒℿ

Annotation of workflows

Workflows are annotated on the basis of:

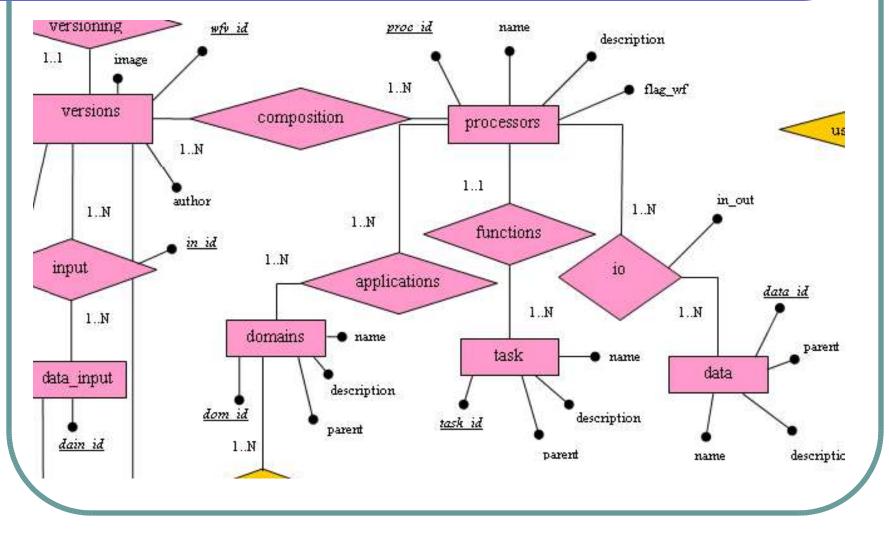
- a simple ontology for bioinformatics processors:
 - application domains
 - task
 - inputs/outputs
- ontology derived from Taverna:
 - new structure
 - some additions (biological resources, images, ...)
 - under further development

O₂I ERA schema



ĪST

O₂I workflows annotation



ĪST

Users' registration and profiling

Users are profiled on the basis of:

- role in their organization
 - computer scientist / physician / researcher / patient / journalist / ...
- domains of interest
- past workflows' executions

∥ℒℼ

O2I Project - Microsoft Internet Explorer		
File Modifica Visualizza Preferiti Strumenti ?		🦧 👖
Indirizzo inttp://www.o2i.it:8080/o2i/		Vai
00	O ₂ I (Oncology over Internet) Project Your personalized project research web site.	
Please login		
username: password: login		
New user? Please <u>register</u> !		
Please install the new <u>library</u> . <u>(Instructions)</u> Applets digital <u>certificate</u> .		
	O ₂ I project	
9/03/2006	P. Romano, BIOINFOGRID Initial Course	55

🖄 O2I Project - Microsoft Internet Exp	lorer				×
File Modifica Visualizza Preferiti Stru	menti ?		Merce .		1
Indirizzo 🙆 http://www.o2i.it:8080/o2i/chang	e_profile.jsp			💌 >	Vai
00			ver Internet) Project oject research web site.		^
USER: PaoloR				Clone Window logout	I
All workflows list My last executed My domains workflows	Username: Password:	PaoloR	(Must be six characters or more)	Choosing your password You will use this information to access O ₂ I each time. Capitalization matters for your password!	
My role most popular My role last executed Search by ontology	Re-type Password:	•••••			
All available results Unsaved results	First Name: Last Name:	Paolo Romano		Info about you Please give us some information about you.	
Temporary saved results Persistently saved results Edit your profile	Role:	computer scientist - bio computer scientist - un journalist patient physician (clinical onc	specified		
	E-Mail Address:	paolo.romano@istge.i	t		
	Domain:	Regulatory affairs and Research Scientific divulgation Translational research Virology			
	Do you want	to receive news through Yes O No	n your e-mail address?		

🚰 021 Project - Microsoft Internet Explorer		
File Modifica Visualizza Preferiti Strumenti ?		TT, 🌾
Indirizzo i http://www.o2i.it:8080/o2i/main.jsp?modo=0		Vai
00	O ₂ I (Oncology over Internet) Project Your personalized project research web site.	
USER: PaoloR		Clone Window logout

All workflows list My last executed

My domains workflows

My role most popular My role last executed

Thy fold labe excedeed

Search by ontology

All available results Unsaved results

Temporary saved results

Persistently saved results

Edit your profile

All workflows:

Workflow			Description	Version	
Conditional Branch Choice		run	This is a demo workflow distributed with Taverna Workbench (see <u>taverna site</u>). If the input is true then the string 'foo' is emited, if false then 'bar'. Just a simple example to show how the conditional branch processor works.	1.0	
Retrieve Cell Lines Descriptions By Name	details	run	This workflow takes the cell line name and the catalogue(s) name(s) as input and retrieve the full cell line description(s) by first retrieving the cell lines' unique IDs associated with the input (done via a call to the getCellLineIdsByName web service) and then using IDs for retrieving the full cell lines descriptions (done via a call to the getCellLinesByIds web service). Both these web services are available at the soaplab system at http://www.o2i.it:8080/axis/services A number of string or string list local elaborations are required: - returned IDs are in a string and this must be transformed in a list (done by the 'Separate_cell_line_ids' processor, that is implemented by using a Split_string_into_string_list_by_regular_expression local processor) - returned IDs include catalogues' names and this must be removed before their utilization for further processing (done by the 'Extract_ids_by_removing_catalogues_names' processor, that is implemented by using a Filter_list_of_strings_extracting_match_to_a_regex local processor) - returned IDs include a blank character and this must be substituteb by a '_SP_' characters string before submissing the data to the 'getCellLinesByIds' web service (done by a trivial beanshell script). Special requirements on input data are: - one or more of the following catalogues names can be specified: 'iclc', 'ecacc_cell', 'dsmz_mutz'. Other names may lead to errors, - when specifying more than one catalogue names, they must be in a unique input string but on distinct text lines, - cell lines names can only be made by a single word, excluding special characters as '/',-' and '*', - cell lines names are case insensitive. Example of valid cell lines names are: - vero - hela - a172 - calu6	1.0	
Retrieve decriptions of	details	run	Retrieve full descriptions of bacteria strains from CABRI catalogues (see <u>CABRI site</u>) by their scientific name (genus and species only). Inputs of the workflow arethe name of the involved CABRI catalogues (text/plain string with one catalogue name per line) and the scientific name of the desired bacteria strain (a text/plain string including genus and species separated by a blank space).	1.0	

P. Romano, BIOINFOGRID Initial Course

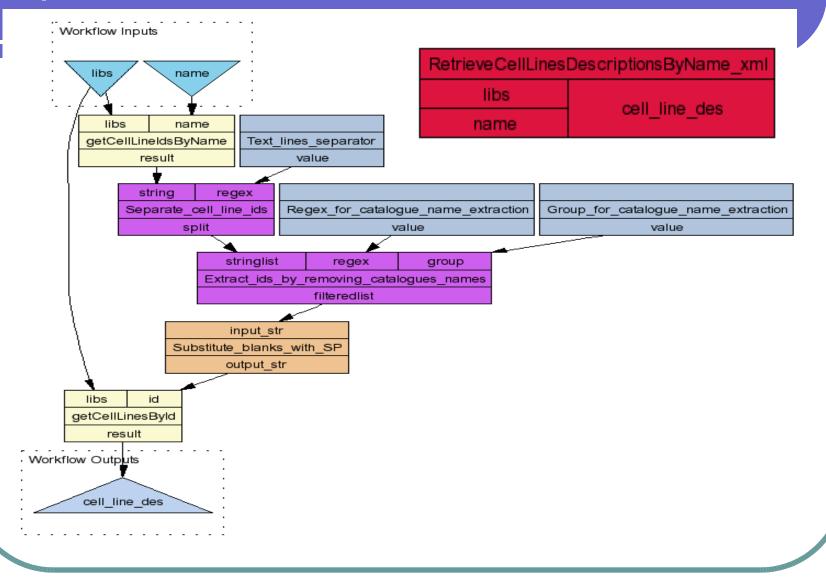
🗿 021 Project - Microsoft Internet	Explorer				
File Modifica Visualizza Preferiti	Strumenti ?	Salfana -			🥂 T
Indirizzo i http://www.o2i.it:8080/o2i/in	ndex.jsp			✓ →	Vai
00		(Oncology over Interne rsonalized project resea			
USER: PaoloR				Clone Window logou	t
All workflows list My last executed	My role most popular workf	lows:	- pi		
My domains workflows	Workflow	Description	Version	Executed	
My role most popular	No workflows corresponds to your	selected criteria.			
My role last executed					
Search by ontology					
All available results					
Unsaved results					
Temporary saved results					
Persistently saved results					
Edit your profile					
		O ₂ I project			2
9/03/2006	P. Rom	nano, BIOINFOGRID Initial	Course	58	

Ø 021 Project - Microsoft Internet Explorer		
File Modifica Visualizza Preferiti Strumenti ?		🥂 'T
Indirizzo i http://www.o2i.it:8080/o2i/main.jsp?modo=4		Vai 🔁
00°	O ₂ I (Oncology over Internet) Project Your personalized project research web site.	
USER: PaoloR		Clone Window logout

All workflows list My application domains workflows: My last executed Workflow Version Description My domains workflows This workflow takes the cell line name and the catalogue(s) name(s) as input and My role most popular retrieve the full cell line description(s) by first retrieving the cell lines' unique IDs associated with the input (done via a call to the getCellLineIdsByName web service) My role last executed and then using IDs for retrieving the full cell lines descriptions (done via a call to the getCellLinesByIds web service). Both these web services are available at the soaplab Search by ontology system at http://www.o2i.it:8080/axis/services A number of string or string list local elaborations are required: - returned IDs are in a All available results string and this must be transformed in a list (done by the 'Separate cell line ids' processor, that is implemented by using a Unsaved results Split string into string list by regular expression local processor) - returned IDs Retrieve Cell include catalogues' names and this must be removed before their utilization for further Lines Temporary saved results details run processing (done by the 'Extract ids by removing catalogues names' processor, that 1.0 Descriptions is implemented by using a Filter list of strings extracting match to a regex local Persistently saved results By Name processor) - returned IDs include a blank character and this must be substituteb by a '_SP_' characters string before submissing the data to the 'getCellLinesByIds' web service (done by a trivial beanshell script). Edit your profile Special requirements on input data are: - one or more of the following catalogues names can be specified: 'iclc', 'ecacc cell', 'dsmz mutz'. Other names may lead to errors, - when specifying more than one catalogue names, they must be in a unique input string but on distinct text lines, - cell lines names can only be made by a single word, excluding special characters as '/','-' and '*', - cell lines names are case insensitive. Example of valid cell lines names are: - vero - hela - a172 - calu6 Retrieve full descriptions of bacteria strains from CABRI catalogues (see CABRI site) by their scientific name (genus and species only). Inputs of the workflow arethe name of the involved CABRI catalogues (text/plain string with one catalogue name per line) and the scientific name of the desired bacteria strain Retrieve decriptions of (a text/plain string including genus and species separated by a blank space). details run 1.0 bacteria Data is retrieved from CABRI Web Services in two steps. First, all bacteria strains IDs are retrieved by using the getBacteriaIdsByName method, and after descriptions are strains retrieved by using the getBacteriaById method. Some list/text elaboration is required to remove catalogue names from returned IDs

P. Romano, BIOINFOGRID Initial Course

Simple demo workflow



IST

File Modifica Visualizza Preferiti Strumenti ?

Indirizzo i http://www.o2i.it:8080/o2i/details.jsp?wfv_id=2



O₂I (Oncology over Internet) Project Your personalized project research web site.

USER: PaoloR

1 P

🗸 🛃 Vai

All workflows	list
---------------	------

My last executed

My domains workflows

My role most popular

My role last executed

Search by ontology

All available results

Unsaved results

Temporary saved results

Persistently saved results

Edit your profile

Workflows details

Name: Retrieve Cell Lines Descriptions By Name

Description: This workflow takes the cell line name and the catalogue(s) name(s) as input and retrieve the full cell line description(s) by first retrieving the cell lines' unique IDs associated with the input (done via a call to the getCellLineIdsByName web service) and then using IDs for retrieving the full cell lines descriptions (done via a call to the getCellLinesByIds web service). Both these web services are available at the soaplab system at http://www.o2i.it:8080/axis/services

A number of string or string list local elaborations are required: - returned IDs are in a string and this must be transformed in a list (done by the 'Separate_cell_line_ids' processor, that is implemented by using a

Split_string_into_string_list_by_regular_expression local processor) - returned IDs include catalogues' names and this must be removed before their utilization for further processing (done by the 'Extract_ids_by_removing_catalogues_names' processor, that is implemented by using a Filter_list_of_strings_extracting_match_to_a_regex local processor) - returned IDs include a blank character and this must be substituteb by a '_SP_' characters string before submissing the data to the 'getCellLinesByIds' web service (done by a trivial beanshell script).

Special requirements on input data are: - one or more of the following catalogues names can be specified: 'iclc', 'ecacc_cell', 'dsmz_mutz'. Other names may lead to errors, - when specifying more than one catalogue names, they must be in a unique input string but on distinct text lines, - cell lines names can only be made by a single word, excluding special characters as '/','-' and '*', - cell lines names are case insensitive.

Example of valid cell lines names are: - vero - hela - a172 - calu6

Author: Paolo Romano, IST, Genova, Italy (paolo.romano@istge.it)

Roles list:

- researcher molecular biologist
- researcher cellular biologist
- researcher structural biologist
- researcher microbiologist
- researcher immunologist
- researcher virologist

Version details

Version: 1.0

Date: 18:05 - 27/09/2005

P. Romano, BIOINFOGRID Initial Course

嶜 021 Project - Microsoft Internet	Explorer								
File Modifica Visualizza Preferiti	Strumenti ?							1	J
Indirizzo i http://www.o2i.it:8080/o2i/d							- Ð	Vai	
	Filename: RetrieveC Changes: Original ve Image: show in a new Input list: • biological resou • cell line name i Output list: • CABRI human Domains list: • Microbiology • Cellular Biologi	2005 no, IST, Genova, Italy (paolo.roman cellLinesDescriptionsByName.xml ersion w window urce database (required) (required) and animal cell lines record	9274 201004	They are not orde	ered.			~	
	Name	Description	Task	Domains	Inputs	Outputs			
	Get cell lines id by name	Retrieve CABRI cell lines' IDs after a search in CABRI Web Services by cell lines' name	biological resource retrieval	Microbiology Cellular Biology	biological resource database cell line name	biological resource identifier			

Retrieves cell lines' descriptions by CABRI id biological

resource

retrieval

Get cell lines

descriptions by

id

CABRI human and animal cell

lines record

biological resource database

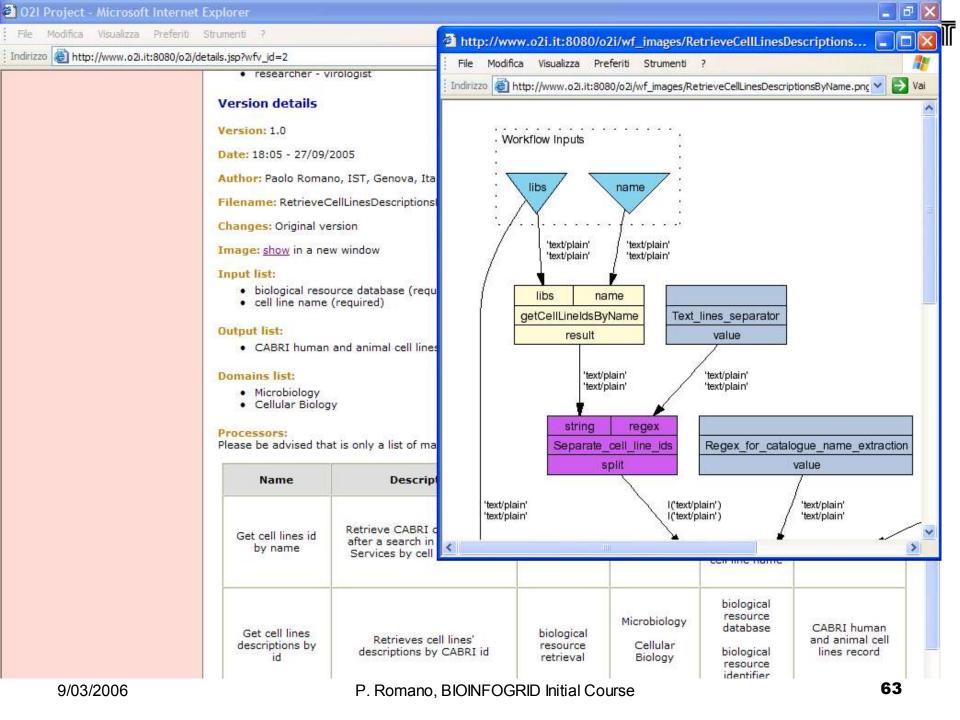
biological

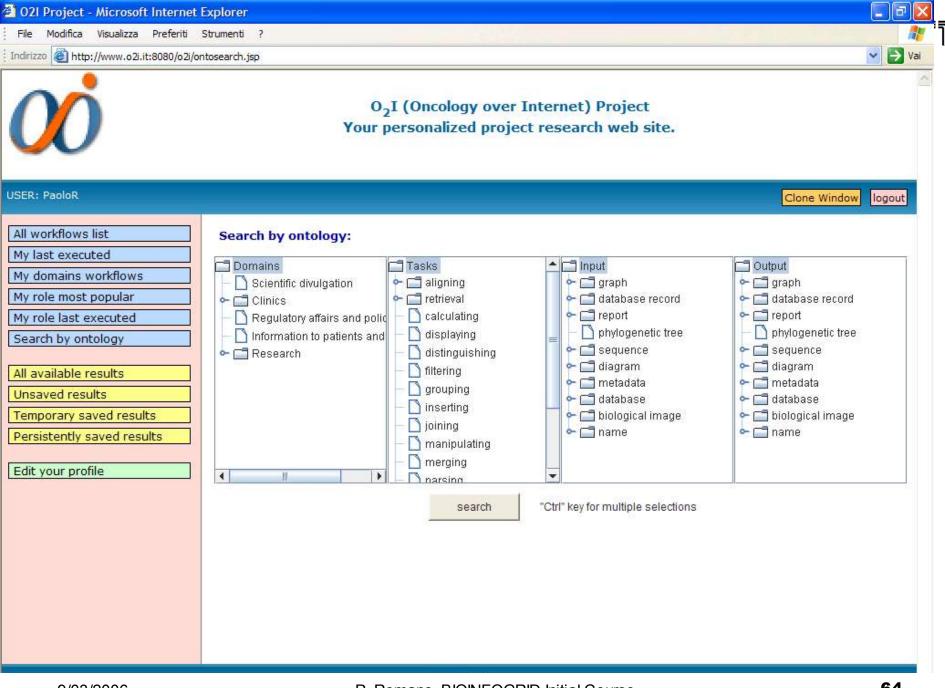
resource identifier

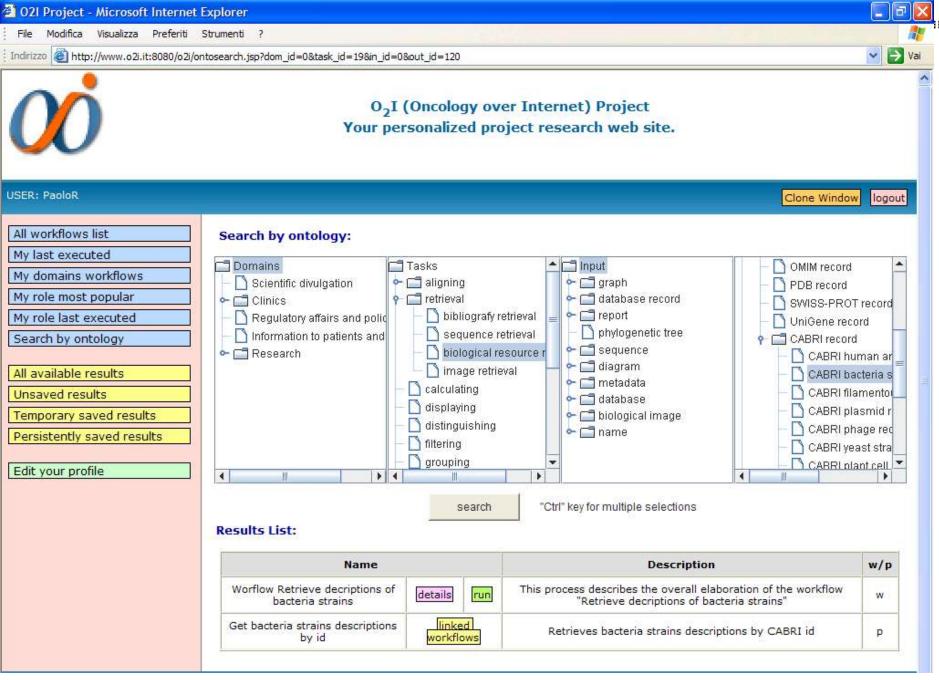
Microbiology

Cellular

Biology





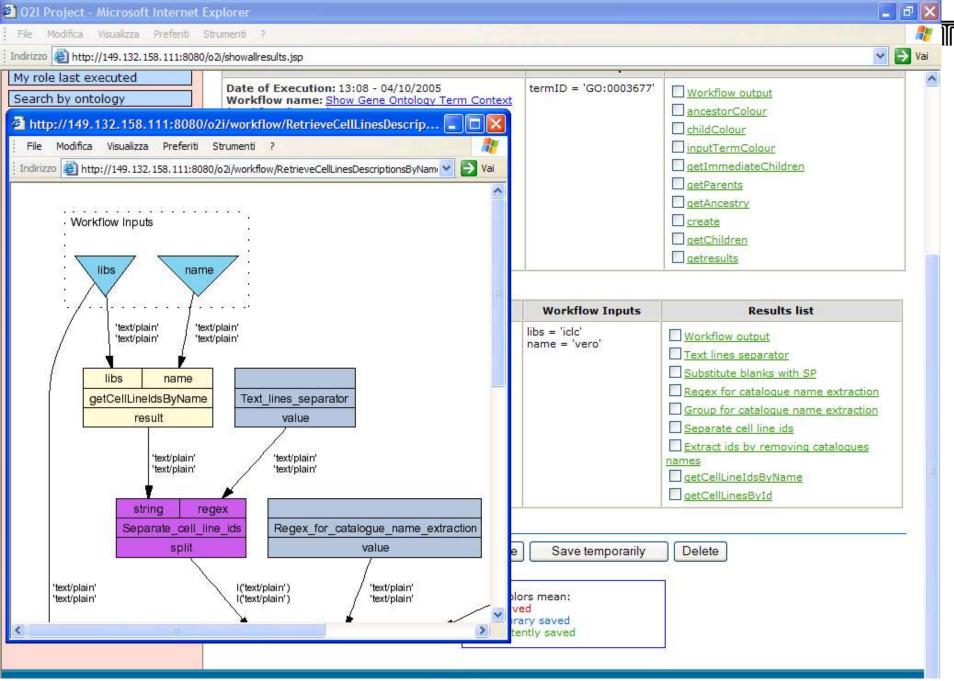


🚰 021 - Microsoft Internet Explorer			
<u>File M</u> odifica <u>V</u> isualizza <u>P</u> referiti (<u>S</u> trumenti <u>?</u>		🦧 ']
Indirizzo 🙆 http://www.o2i.it:8080/o2i/inp	put_insert.jsp?wfv <mark>_i</mark> d=2		💌 🄁 Vai
00	Υοι	O ₂ I (Oncology over Internet) Project r personalized project research web site.	
USER: PaoloR			Clone Window logout
All workflows listMy last executedMy domains workflowsMy role most popularMy role last executedSearch by ontologyAll available resultsUnsaved resultsTemporary saved resultsPersistently saved results	Multiple values can be sp '\n' character). As of Sep 15, 2005, poss - 'iclc' (i.e., the Interlab C - 'ecacc_cell' (i.e., the Eu - 'dsmz_mutz' (i.e., the c	ncludes the name(s) of the CABRI human and animal cell lines catalog cified, in a unique string field, each name in a distinct text line (thus,	names must be divided by a
Edit your profile	service (see the <u>http://ww</u> name. Moreover:		

O2I Project - Microsoft Internet Explorer _ 0 File Modifica Visualizza Preferiti Strumenti ? Indirizzo http://www.o2i.it:8080/o2i/applet_page.jsp?file_name=fin20050523123547686668.xml O₂I (Oncology over Internet) Project Your personalized research web site. USER: PaoloRomano Clone Window Llogout Result Viewer All workflows list **Results visualization** <>> Save to disc My last executed cell line des My domains workflows Execution state: Load com > List Accession_number ICLC ATL95005 My role most popular text/plain Cell_line_name Vero Click to view. My role last executed Brief description Species: monkey, African green; Tissue: kidney Description Species: monkey, African green adult; Tissue: kidney Search by ontology Depositor Obtained from ECACC. Bibliographic_reference Nippon Rinsho 1963;21:1209 All available results Morphology fibroblast, grown as monolayer Culture_conditions continuous culture; DMEM + 10% FBS + 2mM L-Glutamine; s Unsaved results Viruses Search for viruses was not performed Temporary saved results Properties indicator line for mycoplasma testing; virology; virus titration; virus rep Release_conditions Cell line available for distribution. For non-commercial inve Persistently saved results Hazard -Species validation Validated by isoenzymes: confirmed as monkey with AST, MD Edit your profile Karyology modal number 58 Freezing_medium Culture medium + 50% FBS + 10% DMSO Sterility mycoplasma negative, HOECHST and PCR Further_bibliography Kitasato Arc Exp Med 1964;37:27-42 [PMID: 5833688] <BR 4

021 Project - Microsoft Internet	Explorer		
File Modifica Visualizza Preferiti	Strumenti ?		A1
irizzo http://www.o2i.it:8080/o2i/sh	nowsaved.jsp?type=p		Vai
			1
R: PaoloRomano			Clone Window logout
workflows list / last executed / domains workflows	Your persistently saved results: All the results in the same table have been produced by	the same workflow executi	ion.
y role most popular	Execution Details	Workflow Inputs	Results list
v role last executed earch by ontology available results meaved results emporary saved results ersistently saved results dit your profile	Date of Execution: 13:08 - 04/10/2005 Workflow name: <u>Show Gene Ontology Term Context</u> (Workflow diagram)	termID = 'GO:0003677'	Workflow output ancestorColour childColour inputTermColour getImmediateChildren getParents getAncestry create getChildren getresults
	Execution Details	Workflow Inputs	Results list
	Date of Execution: 14:44 - 23/05/2005 Workflow name: <u>Retrieve Cell Lines Descriptions By</u> <u>Name</u> (Workflow diagram)	libs = 'iclc' name = 'vero'	Workflow output Text lines separator Substitute blanks with SP

Execution Details	Workflow Inputs	Results list
ate of Execution: 14:44 - 23/05/2005 orkflow name: <u>Retrieve Cell Lines Descriptions By</u> ame /orkflow diagram)	libs = 'iclc' name = 'vero'	Workflow output Text lines separator Substitute blanks with SP Regex for catalogue name extraction Group for catalogue name extraction Separate cell line ids Extract ids by removing catalogues names getCellLineIdsByName getCellLineSById



Some acknowledgements...

IST, Genoa Paolo Romano, Ulrich Pfeffer, Domenico Marra, Valentina Mirisola, M. Assunta Manniello

ISMAC, CNR, Genoa Patrizio Arrigo, Matteo Fattore

ITB, CNR, Milan Luciano Milanesi **DISCo, University of Milan Bicocca** Guglielmo Bertolini, Flavio De Paoli, Giancarlo Mauri

DIST, University of Genoa Ivan Porro, Silvia Scaglione

DMI, University of Camerino (MC) Emanuela Merelli, Ezio Bartocci

This work has partially been supported by the Italian Ministry for Education, University and Research (MIUR), projects "Oncology over Internet– O_2I " and "Laboratory for Interdisciplinary Technologies in Bioinformatics - LITBIO"

∥ℒℿ

... and an announcement!



Workshop NETTAB 2006 http://www.nettab.org/2006/

Distributed Applications, Web Services, Tools and GRID Infrastructures for Bioinformatics

July 10 - 13, 2005, Sardinia, Italy



Call for Papers deadline is next May 5 Take a brochure!

∥ℒℼ