

A Practical Guide to Using Taverna

Biomed GRID Summer School 2007



Exercise 1 - Getting Started



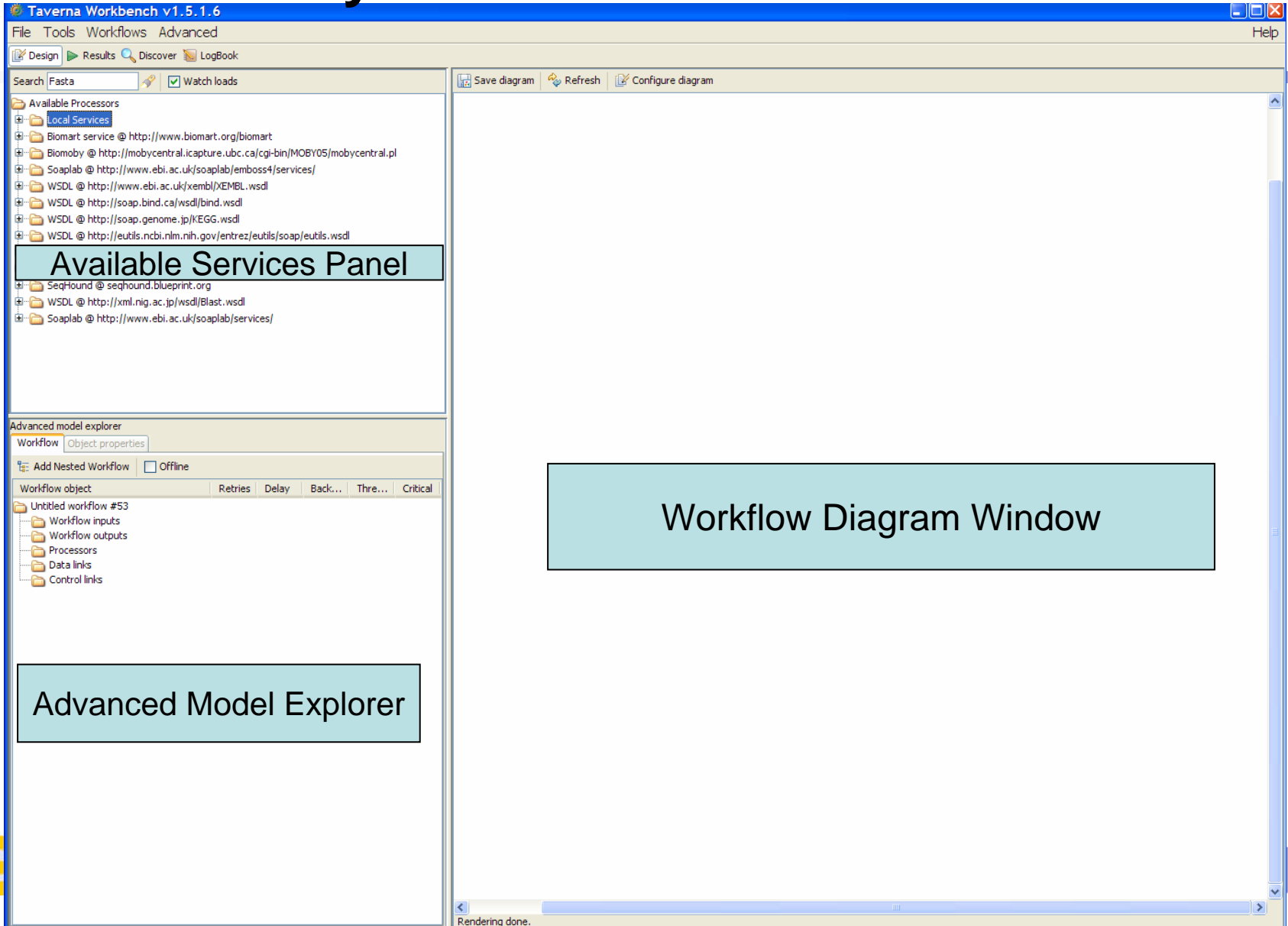
Exercise 1 - Running Taverna

- Taverna has already been installed for you
- To get Taverna started, click on the 'Taverna' icon, which is situated on the desktop



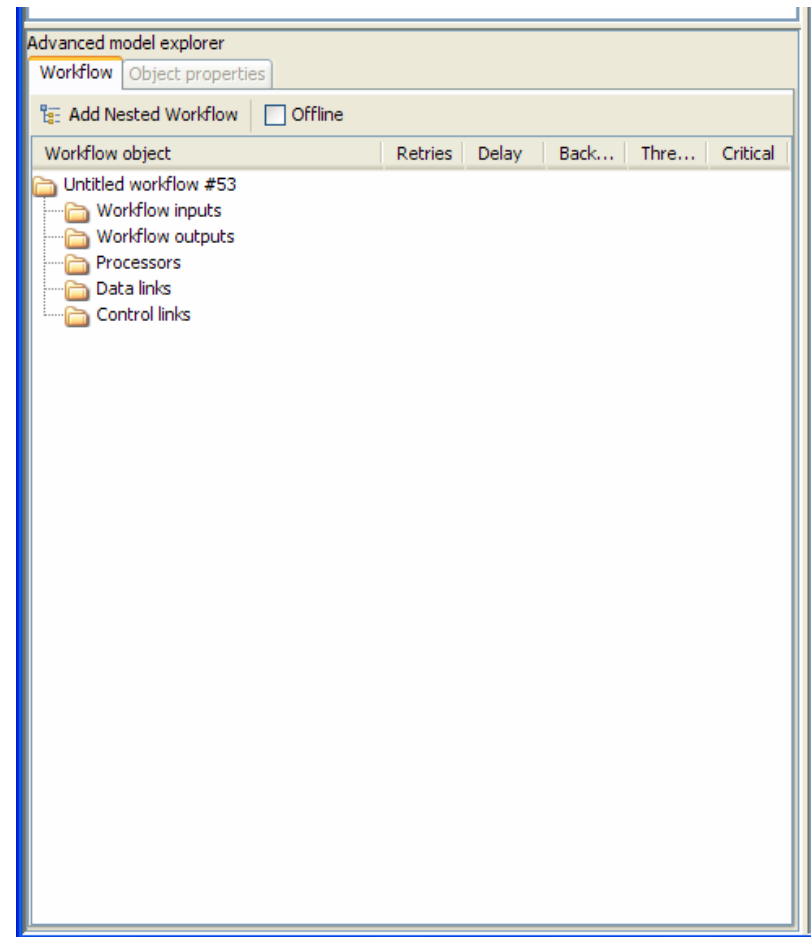
Taverna

Workbench Layout



AME – Advanced Model Explorer

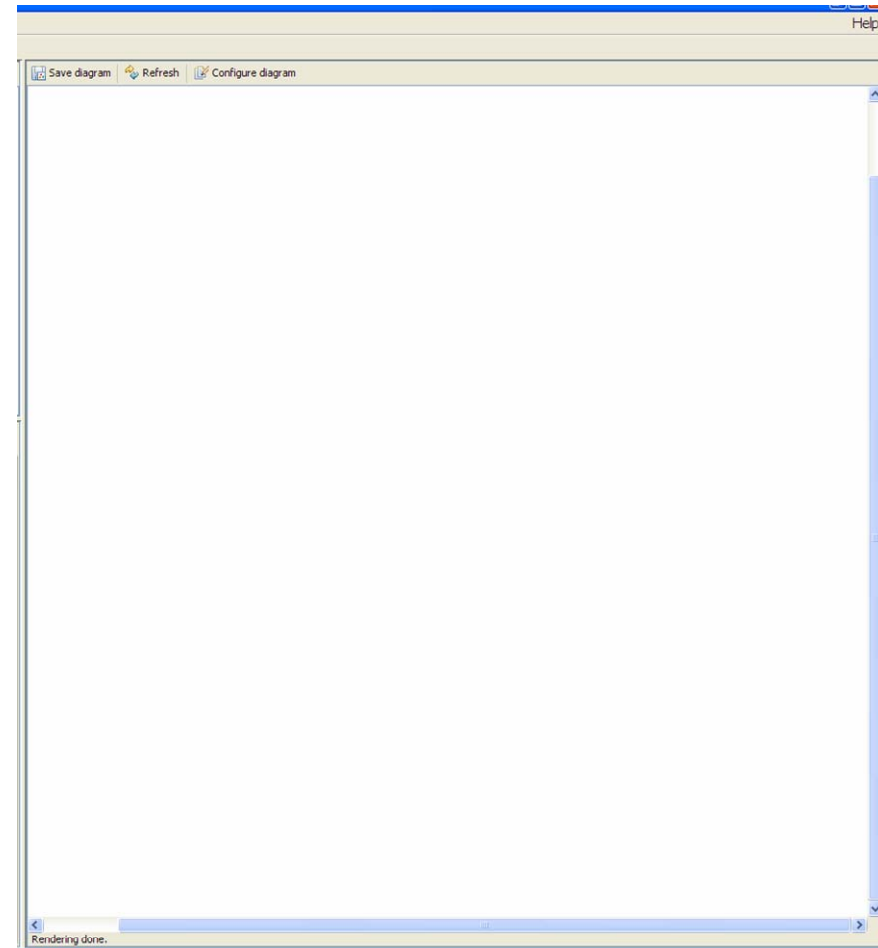
- The Advanced Model Explorer (AME) is the primary editing component within Taverna. Through it you can **build**, **load**, **save** and **edit** any property of a workflow.



Workflow Diagram Window

Visual representation of workflow

- Shows inputs / outputs, services and control flows
- Enables saving of workflow diagrams for publishing and sharing

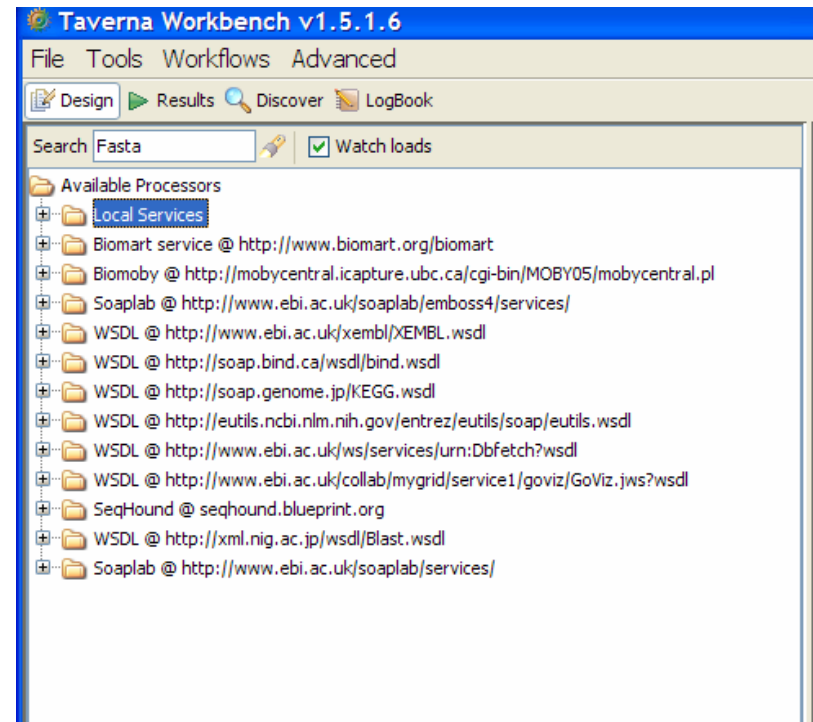


Available Services Panel

Lists services available by default in Taverna

- ~ 3000 services
 - Local java services
 - Simple web services
 - Soaplab services – legacy command-line application
 - Gowlab services
 - BioMart database services
 - BioMoby services

Allows the user to add new services or workflows from the web or from file systems



Plugins

- Additional plugins such as Feta and Logbook are available for Taverna. These can be added via the Tools menu.
- For this tutorial these have already been installed
- You should be able to see two options '*Discover*' and '*LogBook*' at the top of the screen.
- Feta is now available through the Discover tab
- To use the LogBook, you also need a ^{my}SQL database (we will not use this today)

Exercise 2 - Adding New Services

Exercise 2 - Adding a New Service

New services can be gathered from anywhere on the web

1. Go to <http://taverna.sourceforge.net/> and go to the documentation page.
2. Select '(partial) service list'

These services are not all included by default when Taverna opens.

3. Scroll down the page to DDBJ services. You will see a list of available DDBJ services.
4. Click on the DDBJ blast service (<http://xml.nig.ac.jp/wsdl/Blast.wsdl>) and copy the web page address

Exercise 2 - Adding a New Service

5. Go to the '*Available services*' panel and right-click on '*Available Processors*'.

For each type of service, you are given the option to add a new service, or set of services.

6. Select '*Add new wsdl scavenger*'.

A window will pop-up asking for a web address

7. Enter the Blast Web service address

8. Scroll down to the bottom of the '*Available Services*' panel and look at the new DDBJ service that is now included.

Exercise 3 - Finding and Invoking a Single Service

Exercise 3.1 - Finding a Service

1. Go to the 'Available Services' Panel
2. Search for *Fasta* in the 'search list' box at the top of the panel (we will start with simple sequence retrieval)
3. You will see several services highlighted in red
4. Scroll down to '*Get Protein FASTA*'

This service returns a Fasta sequence from a database if you supply it with a sequence identifier

Exercise 3.2 - Invoking a Single Service

5. Right click on the '*Get Protein FASTA*' service and select '*Invoke service*'
6. In the pop-up '*Run workflow*' window add a protein sequence GI by selecting ID and right-clicking. Select '*new input value*' and enter a value (e.g., 1220173) in the box on the right

Note: GI is a genbank gene identifier (you don't need the gi: just the number, for example, the MAP kinase phosphatase sequence 'GI:1220173' would be entered as '1220173')

7. Click '*Run workflow*' and the service is invoked

Exercise 3.3 - Viewing Results

8. Click on '*Results*'

The fasta sequence is displayed on right when you select *click to view*

9. Click on '*Process Report*'

Look at processes.

This shows the experiment provenance – *where* and *when* processes were run

10. Click on '*Status*'

Look at options

As workflows run, you can monitor their progress here.



Exercise 3 - Conclusion

The processes for running and invoking a single service are the basics for any workflow and the tracking of processes and generation of results are the same however complicated a workflow becomes

In the next few exercises, we will look at some example workflows and build some of our own from scratch

Exercise 4 - Finding and Using Workflows



Exercise 4.1 - Finding and Using Workflows

1. Select 'Load Workflow' from the File menu at the top of the workbench.

You will see a selection of .xml files in an examples directory. These are workflow definition files

2. Select '*ConvertedEMBOSSTutorial.xml*' and a pre-defined workflow will be loaded

3. View the workflow diagram - you will see services of in different colours

Exercise 4.2 - Workflow Documentation

4. Find out what the workflow does by reading the workflow metadata
5. In the AME – click on the name of the workflow – in this case '*A workflow version of the EMBOSS tutorial*' and then select the '*workflow metadata*' tab at the top of the AME.

You will see a text description of the workflow, its author and its unique LSID. When publishing workflows for others, this annotation is useful information and allows the acknowledgement of intellectual property



Exercise 4.3 - Workflow Features

6. Run the workflow by selecting 'run workflow' from the file menu
7. Watch the progress of the workflow in the '*enactor invocation*' window. As services complete, the enactor reports the events. If a service fails, the enactor reports this also

Loading workflows from the Web

- Go to the webpage <http://www.cs.man.ac.uk/~katy/taverna>
- Select '*conditional_control.xml*' and copy the web address
- Go back to the taverna workbench and select 'Open from web'
- Run the workflow (using the same GI as input '1220173')
- You will see at least one of the services fail. These are conditionals – fail of false or fail if true. These are useful operators for controlling progress of your workflow based on intermediate results
- You will see black arrows and white circles – black arrows show the flow of the data and white circles are *control links*.

Exercise 5 - Building a Simple Workflow

Exercise 5.1 - Adding a Service

1. Import the '*Get Protein FASTA*' service into a new workflow model

First, you will need to close the current workflow from the file menu, then find the '*Get Protein Fasta*' service again in the '*Available services*' panel.

2. Right-click on '*Get Protein Fasta*' and import it into the workbench by selecting '*Add to Model*'
3. Go to the AME and expand the [+] next to the newly imported '*Get Protein Fasta*' service. You will see:
 - 1 input (green arrow pointing up)
 - 1 output (purple arrow pointing down)

Exercise 5.2 - Adding Input

1. Define a new workflow input by right-clicking on '*Workflow Input*' and selecting 'create new Input'
2. Supply a suitable name e.g., '*genelIdentifier*'
3. Connect this new input to the '*Get Protein Fasta*' service by right-clicking on '*genelIdentifier*' and selecting '*getFasta ->id*'

You always build workflows with the flow of data

Exercise 5.3 - Adding Output

4. Define a new workflow output by right-clicking on '*workflow output*' and selecting 'create new output'
5. Supply a suitable name e.g. '*fastaSequence*'
6. Connect this new output to the '*Get Protein Fasta*' service. remembering to build with the flow of data
You have now built a simple workflow from scratch!
7. Run the workflow by selecting 'run workflow' from the '*Tools and Workflow Invocation*' menu at the very top of the workbench

You will again need to supply a GI – for later exercises, please use a protein GI – e.g. 1220173

Exercise 6 – Stringing Services Together

Exercise 6 - Stringing Services Together

We have used '*Get Protein Fasta*' to retrieve a sequence from the genbank database. What can we do with a sequence?

- Blast it?
- Find features and annotate it?
- Find GO annotations?

Blast it?

The first thing you need to do is find a service which performs a blast. For this, we are going to use the ***Feta Semantic Discovery*** Tool

Feta is a tool to semantically describe services. Instead of the user needing to know exactly what a service provider has called their services, the user can search by the biological tasks that are performed by the services, or by properties of the service, for example, the types of inputs it requires/outputs it produces

Exercise 6.1 - Finding Blast using Feta

1. Select the *'Discover'* tab and select 'uses method from the first drop down menu
2. When you select it, 'bioinformatics algorithm' will appear in the adjoining box. Scroll down this list to find 'Similarity search algorithm', and then the subclass of this, 'BLAST'
3. Select 'BLAST' and click *'Find Service'*

The results are all the annotated services that perform blast analyses (there may be more un-annotated ones!)

Exercise 6.1 - Finding Blast

4. Select '*searchSimple*' from the list and look at the details
5. Look at the service description
This tells you what the service does and what each input/output is expecting/produces. It also tells you where the service comes from. For this example, we are using BLAST from the DNA Databank in Japan
6. Right-click on '*searchSimple*' in the Feta results list and select 'add to model'
This adds the service to your current workflow in the 'Design Window'
7. Before you go back to the Design window, go back to search services and experiment with other ways of finding services – e.g., by task, input/output, resource etc

Exercise 6.2 - Adding Blast

1. Go back to the Design window. '*Search simple*' will have been imported into your model
2. In the AME expand the [+] for the '*search simple*' service and view the input/output parameters

This time, you will see three inputs and two outputs.

For the workflow to run, each input must be defined. If there are multiple outputs, a workflow will usually run if at least one output is defined.



Exercise 6.2 - Adding Blast

3. Create an **output** called '*blast_report*' (in the same way we did before)
4. The sequence input for the Blast service will be the output from the '*Get Protein Fasta*' service. Connect the two together, from '*fastaSequence*' (output from *GetProteinSequence* service) to '*search simple query*'
5. Create two more inputs called '*database*' and '*program*' and connect them to the '*database*' and '*program*' inputs on '*search simple*' service

Exercise 6.3 - Running the 'Blast it' Workflow

1. Once more select '*run workflow*' from the '*Tools and Workflow Invocation*' menu. You will see a run workflow window asking for 3 input values
2. Insert a GI (e.g., 1220173), a program (blastp for protein-protein blast), and a database (SWISS for swissprot)
3. Click 'run workflow'. This time you will see a blast report and a fasta sequence as a result

Exercise 6.4 - Changing Parameters

For parameters that do not change often, you will not wish to always type them in as input. In this example, the database and blast program may only change occasionally, so there is an alternative way of defining them

1. Go back to the AME and remove the '*database*' and '*program*' inputs by right-clicking and selecting '*remove from model*'

Exercise 6.5 - String Constants

1. Select '*string constant*' from '*Available Services*'
2. Right-click and select 'add to model with name...'
3. Insert '*program*' in the pop-up window
4. Select '*string constant*' for a second time and repeat for a string constant named '*database*'
5. In the AME, right-click on '*program*' and select '*edit me*'
6. Edit the text to 'blastp'. Repeat for 'database' and enter 'SWISS' for the swissprot database
7. Run the workflow – it runs in the same way
8. Save the workflow by selecting the 'save' icon at the top of the AME.

Exercise 7 - Completing the Protein Annotation Workflow



Exercise 7.1 - Adding More Services

How can we use Taverna to annotate our protein with function descriptions?

1. In the 'available services' panel, find the emboss soaplab services and find the '*protein_motifs*' section
2. *Hint: use the simple text search at the top of the panel*
3. Find out which of these services enable searching of the Prosite and Prints databases by fetching the service descriptions. To do this *right-click* on '*protein_motifs*' and select '*fetch descriptions*'
4. Import both services into the workflow model.

Exercise 7.1 - Adding More Services

5. Connect these services up to the workflow so that you can find prints and prosite matches in the query sequence returned from '*Get Protein Fasta*' – you will see that soaplab services have many input values

Soaplab services have many input parameters, but many have default values so may not always need to be altered. In this case, you can run the services by simply adding the query sequence. Go to the EMBOSS home page to find out which input(s) relate to the query sequence.

This extra searching is impractical – the Feta Semantic Discovery tool is designed to combat this problem

Exercise 7.2 - Running the Protein Annotation Workflow

1. Run the workflow – now you have blast results and protein domain/motif matches

How else can you annotate your protein? As an advanced exercise, you might want to search for other ways of characterising your sequence e.g. structural elements, GO annotation?

Exercise 8 - Saving Results

Taverna provides several options for saving data.

- a) Individual data items can be saved by right-clicking on them
- b) All data can be saved to disk
- c) Textual/tabular data can be saved to excel

Save all the data from your workflow

Advanced Exercises

The previous exercises have covered the basics of myGrid workflows. The following demos and exercises cover more advanced features, such as rendering output, configuring BioMart services, dealing with service failure and iterating over datasets. You may not reach the end of these exercises, but they will provide a some examples to take home

Exercise 9 - Defining Output Formats

So far, most of the outputs we have seen have been text, but in bioinformatics, we often want to view a graph, a 3D structure, an alignment etc. Taverna is able to display results using a specific type of renderer if the workflow output is configured correctly.

1. Reset the workbench and load *'convertedEMBOSSTutorial'* from the *'examples'* directory
2. Look at the workflow diagram and read the workflow metadata to find out what the workflow does
3. Run the workflow

Exercise 9 - Defining Output Format

4. Look at the results.

For *'tmapPlot'* and *'outputPlot'*, you will see the results are displayed graphically. This is achieved by specifying a particular mime type in the output.

5. Go back to the AME and look at the metadata for *'tmapPlot'* and *'outputPlot'*.

6. Select MIME Types. As you can see, each has the image/png mime type associated with it. If you wish to render results in anything other than plain text, you **MUST** specify the mime-type in the workflow output

Taverna MIME-Types

The following mime-types are currently used by Taverna

text/plain=Plain Text

text/xml=XML Text

text/html=HTML Text

text/rtf=Rich Text Format

text/x-graphviz=Graphviz Dot File

image/png=PNG Image

image/jpeg=JPEG Image

image/gif=GIF Image

application/zip=Zip File

chemical/x-swissprot=SWISSPROT Flat File

chemical/x-embl-dl-nucleotide=EMBL Flat File

chemical/x-ppd=PPD File

chemical/seq-aa-genpept=Genpept Protein

chemical/seq-na-genbank=Genbank Nucleotide

chemical/x-pdb=Protein Data Bank Flat File

chemical/x-mdl-molfile

Exercise 9 - Taverna MIME types

The ‘chemical/’ mime-types are rendered using SeqVista to view formatted sequence data

7. Reset the workbench and load ‘*seqVistaRendering*’ from the ‘*examples*’ directory for a demo
8. Run the workflow and have a look at results – you should see that seqVista has enabled the results to be viewed graphically as the mime type was set to ‘*chemical/x-embl-dl-nucleotide*’. To see how the results can be viewed in different ways, right click on the ‘click to view’ icon in the results pane. You should see a list of the results can be rendered.

Advanced Features

- Spotlight on BioMart
- BioMoby Services
- Iteration
- Control Flow
- Substituting Services and fault tolerance

Exercise 10 - Spotlight on Biomart

Biomart enables the retrieval of large amounts of genomic data e.g., from Ensembl and Sanger, as well as Uniprot and MSD datasets

1. After saving any workflows you want to keep, reset the workbench in the AME
2. Load the workflow 'BiomartAndEMBOSSAnalysis.xml' from the '*examples*' directory
3. Run the Workflow

What does the workflow do?

This Workflow Starts by fetching all gene IDs from Ensembl corresponding to human genes on chromosome 22 implicated in known diseases and with homologous genes in rat and mouse.

For each of these gene IDs it fetches the 200bp after the five prime end of the genomic sequence in each organism and performs a multiple alignment of the sequences using the EMBOSS tool 'emma' (a wrapper around ClustalW).

It then returns PNG images of the multiple alignment along with three columns containing the human, rat and mouse gene IDs used in each case.

Exercise 10.1 - Configuring Biomart

1. Right-click on the service and select '*configure BioMart query*'
2. By selecting '*filters*' – change the chromosome from 22 to 21 – now the workflow will retrieve all disease genes from chromosome 21 with rat and mouse homologues
3. Run the workflow and look at the results
4. See how the '*disease gene*' filter was configured and the '*sequence exports*' were configured on the other Biomart queries for mouse and rat

Exercise 10.2 - Adding Extra Information

Find out which diseases the known diseases are on your chosen chromosome by adding a new Biomart query process

1. Select '*hsapiens_gene_ensembl*' from the available services panel and select 'invoke with name....' (as there is already a service with that name!)
2. Call the service '*hsapiens_disease*'
3. Configure '*hsapiens_disease*' by selecting an '*ensembl gene IDs*' filter under the '*gene*' tab
4. Configure the output attribute '*disease description*' under the '*gene*' tab in the attributes section

Exercise 10.2 - Adding Extra Information

5. Connect the input to the '*hsapiens_gene_ensembl*' service via the '*gene_stable_id*'
6. Create a new workflow output for the '*disease_description*' output
7. Re-run the workflow and view which diseases are associated with your chromosome

Exercise 11 - Spotlight on BioMoby

The process of adding a BioMoby service is different from other services. BioMoby services need to be defined using terms from the Moby Object ontology

1. Reset the workflow and load the '*blast-biomoby.xml*' workflow from

<http://www.cs.man.ac.uk/~katy/taverna/>

Exercise 11 - Spotlight on BioMoby

2. Run the workflow and look at the results
As the workflow name suggests, a blast search is performed on a sequence
3. Look at the workflow diagram
Instead of simply giving the blast service a fasta sequence, there is a '*Fasta*' sequence object defined
4. Look at the inputs for '*Fasta*'
5. Read the metadata for the '*Fasta*' object in the AME window

Exercise 11 - Spotlight on BioMoby

The Fasta object is defined by

- a) The sequence (as a plain string)
- b) The namespace (i.e. the database the sequence came from)
- c) A unique identifier for the sequence
- d) A name

These extra definitions take time for the user to define, but they have other advantages

Exercise 11 - Spotlight on BioMoby

7. Right-click on the '*Fasta*' object in the AME and select '*Moby Object Details*'
8. A pop-up window will show you what BioMoby services a '*Fasta*' sequence is produced by and what services it can feed into
9. Right-click on the '*getDragonBlastText*' service and select '*Moby Object Details*'. This tells you what the service requires as inputs and what it produces as output

Spotlight on BioMoby

- The BioMoby services are annotated using terms from the Moby ontology to enable semantic searching for services.
- BioMoby services are specialist kinds of service from a closed community. The object model, ontology and annotations have been agreed by the BioMoby service providers.
- Semantic discovery queries over other ^{my}Grid services are also possible using the ^{my}Grid ontology and the Feta Semantic discovery component.
- The ^{my}Grid ontology and the BioMoby ontology both share the same service ontology, so feta can search both types of service

Exercise 11 - Iteration

Taverna has an implicit iteration framework. If you connect a set of data objects (for example, a set of fasta sequences) to a process that expects a single data item at a time, the process will iterate over each sequence

1. Reload the BiomartandEMBOSSAnalysis.xml workflow from the examples directory
2. Watch the progress report. You will see several services with 'Invoking with Iteration'

Exercise 11 - Iteration

The user can also specify more complex iteration strategies using the service metadata tag

1. Reset the workflow and load the '*IterationStrategyExample.xml*'
2. Read the workflow metadata to find out what the workflow does
3. Select the '*ColourAnimals*' service and read the metadata for that service. Under the description is the iteration strategy
4. Click on '*dot product*'. This allows you to switch to cross product

Exercise 11 - Iteration

5. Run the workflow twice – once with '*dot product*' and once with '*cross product*'.
6. Save the first results so you can compare them – what is the difference? What does it mean to specify dot or cross product?

Exercise 12 - Substituting Services and Fault Tolerance



Exercise 12.1 – Substituting Services

Taverna does not own many of the bioinformatics services it provides. This means that it cannot control their reliability. Instead, Taverna provides strategies for dealing with services being unavailable

1. Reload the '*convertedEMBOSSTutorial.xml*' from the '*examples*' directory.
2. Look at the metadata for the '*emma*' service. It is an implementation of clustalw
3. Find the DDBJ clustalw service and add it to your workflow

HINT: use the Feta discovery tool

Exercise 12.1 - Substituting Services

4. When you have added this service to your workflow, right-click on it and select '*add as alternate*'
5. In the resulting menu select '*emma*'
The DDBJ version of the clustalw service is now added as an alternative to emma in the AME. It will be called '*alternate1*'
6. Select '*alternate1*' and look at the inputs and outputs. These need to be mapped to the correct inputs and outputs in emma

Exercise 12.1 - Substituting Services

7. Right-click on the '*query*' input in alternate1 and map it to '*sequence_direct_data*'. In both services, these inputs expect a set of fasta sequences.
8. Right-click on the '*result*' output and map it to '*outseq*' in emma in the same way.

Now you have a workflow which will run using emma when it is available – but will substitute it for DDBJ clustalw if emma fails!

Exercise 12.2 - Fault Tolerance

Taverna also allows the user to specify the number of times a service is retried before it is considered to have failed. Sometimes network traffic is heavy, so a working service needs to be retried

1. Select 'tmap' from the same workflow. To the right of the service name are a series of 0s and 1s. By simply typing the numbers, the user can specify the number of retries and the time between the retries
2. Change it to 3 retries for 'tmap' and set the status to 'critical' using the final tickbox. Now it is critical, it means the whole workflow will be aborted if '*tmap*' fails after 3 retries. Failures in non-critical services will not abort the workflow run.

Exercise 13 - Shim Services

This exercise highlights the services that do not perform biological functions, but are vital for running life science workflows



Exercise 13.1 - Finding Genes

1. Load the workflow entitled `genscan_shim_example.xml` from the page <http://www.cs.man.ac.uk/~katy/taverna>
2. Look at the workflow metadata – what does the workflow do?
3. Run the workflow.
4. For an input file, load `example_input.txt` from the same web page

What happens?

Did all the services return results?

Why did some fail?

Exercise 13.1 - Finding Genes

5. Load the workflow entitled `genscan_shim_example2.xml` from the page <http://www.cs.man.ac.uk/~katy/taverna>
6. Look at the workflow metadata – what does the workflow do? How is it different from the previous one?
7. Run the workflow (using the same input) – what happens this time?

Genscansplitter is a shim service – it performs no biological function, it simply parses a results file.

Exercise 13.2 - Other Shims

- There are many myGrid shim services. These are currently being described in a shim library, but for now, a small collection are documented here
<http://www.cs.man.ac.uk/~hulld/shims.html>

From the list,

1. Find a shim that will return a genbank DNA file from an id. Load the example workflow and run it in Taverna
2. Find a shim that will translate DNA

HINT: these services might be in the feta registry

Exercise 13.3 - Other Shims

1. Load the CompareXandYFunctions.xml workflow from the examples directory
This workflow contains several shims. Some are beanshell scripts
2. Select the 'GetUniqueIDs' service in the AME and right-click
3. Look at the script and see if you can work out what it is doing

Beanshell scripts allow users to write small, bespoke java scripts to allow incompatible service to work together

Exercise 13.4 - Other Shims

The emboss suite of programs have a subdivision – edit

1. All the edit services are shims
2. Experiment with the edit services
3. Find a service that will remove gaps from sequences