# Workflows Accessibility in Bioinformatics

Konstantinos Karasavvas

Netherlands Bioinformatics Centre

Leiden University Medical Center

# Workflows / Pipelines

- Large data sets
  - e.g. NGS
- Data analysis
  - multiple steps
    - automation
    - flexibility
    - reusability
- User groups
  - biologists
  - bioinformaticians

# Simple Task

- List all enzymes catalyzing reactions involving a given compound

- Services on the Internet provide this functionality
  - e.g. KEGG services
  - http://www.genome.jp/kegg/

- Scripts
- Galaxy
- Taverna
- Taverna-Galaxy
- Taverna-Web

Workflows Accessibility in Bioinformatics

# Script/Program (1)

```ruby
#!/usr/bin/env ruby

require 'soap/wsdlDriver'

wsdl = "http://soap.genome.jp/KEGG.wsdl"
kegg_service = SOAP::WSDLDriverFactory.new(wsdl).create_rpc_driver

enzyme_classifications = []

begin
  reactions = kegg_service.get_reactions_by_compound(ARGV[0])
rescue => err
  puts err.message
end

reactions.each do |reaction|
  begin
    enzyme_classifications << kegg_service.get_enzymes_by_reaction(reaction)
  rescue => err
    puts err.message
  end
end

enzyme_classifications.uniq.each { |ec| puts ec }
```

LU MC

nbic  Netherlands Bioinformatics Centre

# Script/Program (2)

```
$ ./simpleGetEnzymesFromCompoundWkf.rb C15973

ec:2.3.1.12
ec:2.3.1.61
ec:2.3.1.61
ec:2.3.1.168
ec:2.3.1.168
ec:2.3.1.168
ec:1.8.1.4

$
```

Workflows Accessibility in Bioinformatics

# Script/Program (2)

```
$ ./simpleGetEnzymesFromCompoundWkf.rb C15973

ec:2.3.1.12
ec:2.3.1.61
ec:2.3.1.61
ec:2.3.1.168
ec:2.3.1.168
ec:2.3.1.168
ec:1.8.1.4

$
```

Who does what?

# Multiple Scripts

```ruby
#!/usr/bin/env ruby

require 'soap/wsdlDriver'

wsdl = "http://soap.genome.jp/KEGG.wsdl"
kegg_service = SOAP::WSDLDriverFactory.new(wsdl).create_rpc_driver

enzyme_classifications = []

begin
  reactions = kegg_service.get_reactions_by_compound(ARGV[0])
rescue => err
  puts err.message
end

reactions.each do |reaction|
  begin
    enzyme_classifications <<
kegg_service.get_enzymes_by_reaction(reaction)
  rescue => err
    puts err.message
  end
end

enzyme_classifications.uniq.each { |ec| puts ec }
```

superscript
is needed

```ruby
#!/usr/bin/env ruby

require 'soap/wsdlDriver'

wsdl = "http://soap.genome.jp/KEGG.wsdl"
kegg_service = SOAP::WSDLDriverFactory.new(wsdl).create_rpc_driver

genes = []

enzymes.each do |ec|
  begin
    enzyme_classifications << kegg_service.get_genes_by_enzymes(ec)
  rescue => err
    puts err.message
  end
end

genes.uniq.each { |gene| puts gene }
```

Workflows Accessibility in Bioinformatics

nbic  Netherlands Bioinformatics Centre

# Multiple Scripts

```ruby
#!/usr/bin/env ruby

require 'soap/wsdlDriver'

wsdl = "http://soap.genome.jp/KEGG.wsdl"
kegg_service = SOAP::WSDLDriverFactory.new(wsdl).create_rpc_driver

enzyme_classifications = []

begin
  reactions = kegg_service.get_reactions_by_compound(ARGV[0])
rescue => err
  puts err.message
end

reactions.each do |reaction|
  begin
    enzyme_classifications <<
kegg_service.get_enzymes_by_reaction(reaction)
  rescue => err
    puts err.message
  end
end

enzyme_classifications.uniq.each { |ec| puts ec }
```

superscript
is needed

```ruby
#!/usr/bin/env ruby

require 'soap/wsdlDriver'

wsdl = "http://soap.genome.jp/KEGG.wsdl"
kegg_service = SOAP::WSDLDriverFactory.new(wsdl).create_rpc_driver

genes = []

enzymes.each do |ec|
  begin
    enzyme_classifications << kegg_service.get_genes_by_enzymes(ec)
  rescue => err
    puts err.message
  end
end

genes.uniq.each { |gene| puts gene }
```

**e.g. GAPSS pipeline**

LU MC

nbic Netherlands Bioinformatics Centre

# Multiple Scripts - Issues

- Difficult to manage
  - many scripts/programs
  - many arguments per script/program
  - many directories/files (scripts' I/Os)
- Difficult to reuse
  - due to low-level description
  - inconsistent program arguments

# In Galaxy Web Portal (1)



Workflows Accessibility in Bioinformatics

# In Galaxy Web Portal (1)



Workflows Accessibility in Bioinformatics

# In Galaxy Web Portal (1)



Workflows Accessibility in Bioinformatics

# In Galaxy Web Portal (2)



Workflows Accessibility in Bioinformatics

# In Galaxy Web Portal (3)

# In Galaxy Web Portal (3)



e.g. GAPSS pipeline

Workflows Accessibility in Bioinformatics

# In Taverna Workbench (1)

# In Taverna Workbench (2)

# In Taverna Workbench (2)



Who does what?

Workflows Accessibility in Bioinformatics

# In Taverna Workbench (3)

# In Taverna Workbench (4)
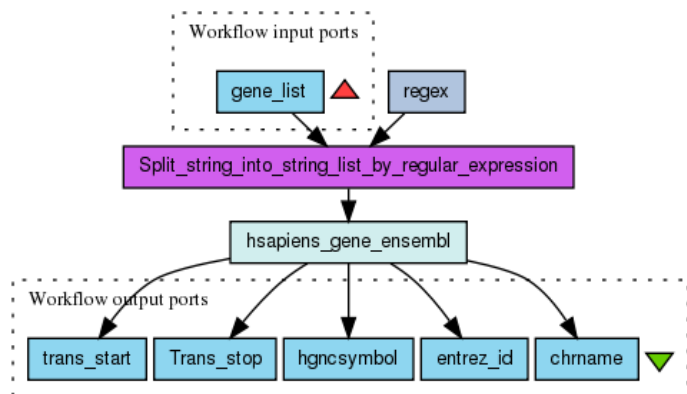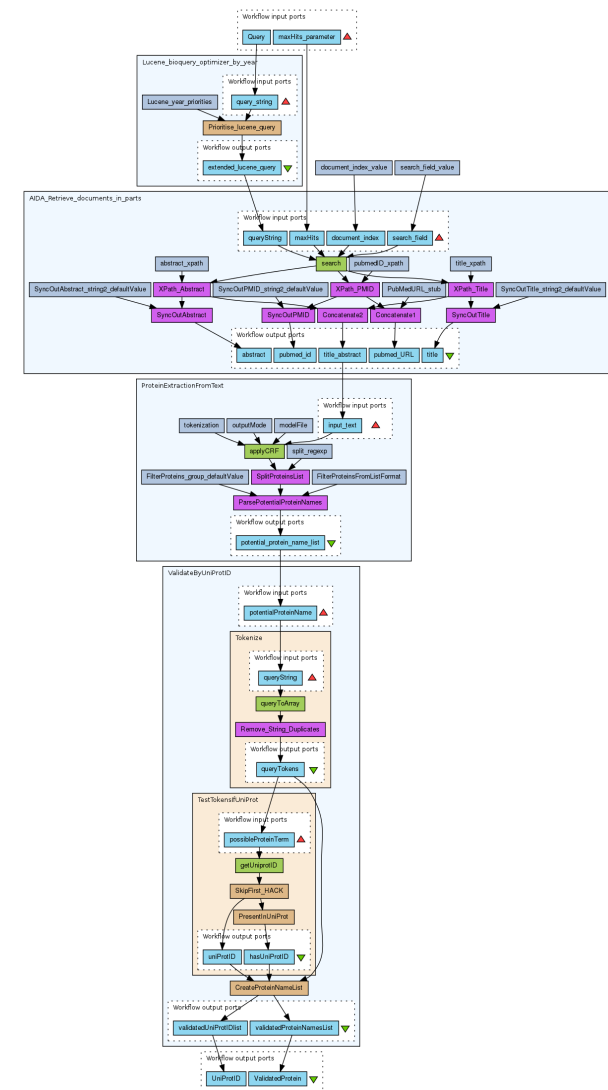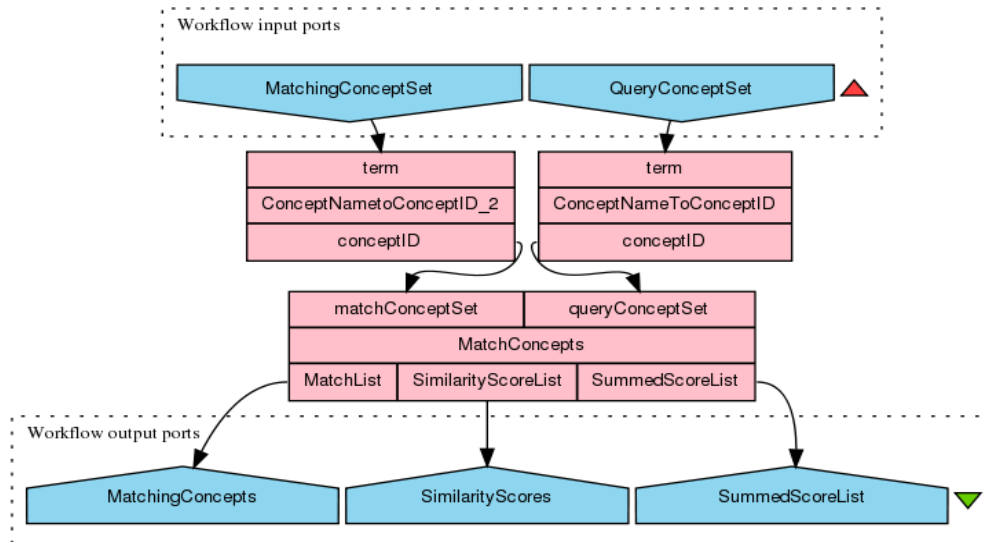
# Harish's Workflows (GWAS: Metabolic Syndrome)

# Eleni's Workflows (Huntington's Disease)

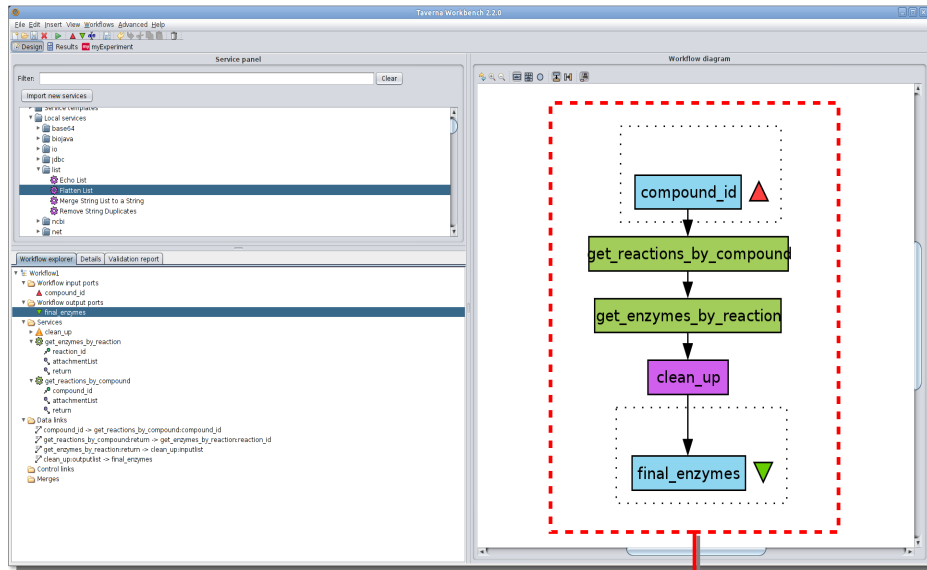# Marco's workflows (Semantic Text Mining)

# Taverna and Galaxy workflows

- The systems have a different focus
  - Some overlapping functionality but different strengths
  - Different fan clubs!

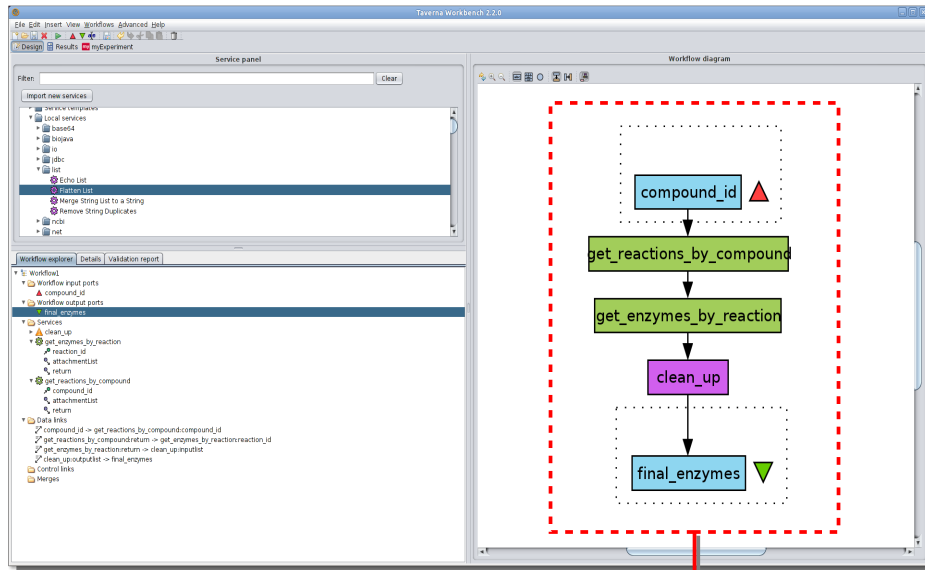| Galaxy | Taverna |
| --- | --- |
| straightforward workflows | very expressive workflows |
| exposing existing scripts | exposing existing web services |
| typically local tools | typically remote tools |

- Do we have to make a choice?
- Does that limit our potential users?
- How to make them more interoperable?

Workflows Accessibility in Bioinformatics

Netherlands Bioinformatics Centre

# Taverna Workflows in Galaxy (1)



Workflows Accessibility in Bioinformatics

# Taverna Workflows in Galaxy (1)



Who does what?

Workflows Accessibility in Bioinformatics

# Taverna Workflows in Galaxy (2)



Workflows Accessibility in Bioinformatics

# Taverna Workflows on a Web Browser (1)



**Workflows Accessibility in Bioinformatics**

# Taverna Workflows on a Web Browser (1)



Who does what?

Workflows Accessibility in Bioinformatics

# Taverna Workflows on a Web Browser (2)



Workflows Accessibility in Bioinformatics

# Taverna Workflows on a Web Browser (3)



Workflows Accessibility in Bioinformatics

# Summary

- Workflow example using several approaches

- Taverna workfkows can be accessed in Galaxy
  - Bioinformatician *creates/finds* appropriate workflow
  - He uses Taverna-Galaxy to create new tool and installs it
  - … biologist will see the new tool in the Galaxy server
    - The taverna workflow can now take part in a Galaxy workflow

- Taverna workflows can be accessed via the web
  - Bioinformatician *creates/finds* appropriate workflow
  - … sends the URL to biologist

Workflows Accessibility in Bioinformatics

Netherlands
Bioinformatics
Centre

- More information
  - http://galaxy.psu.edu/
  - http://www.taverna.org.uk/
  - https://trac.nbic.nl/elabfactory/wiki/eGalaxy

- Questions?
  - kostas.karasavvas@nbic.nl