# Scientific workflow management a way to enable e-science on both Grids and Clouds
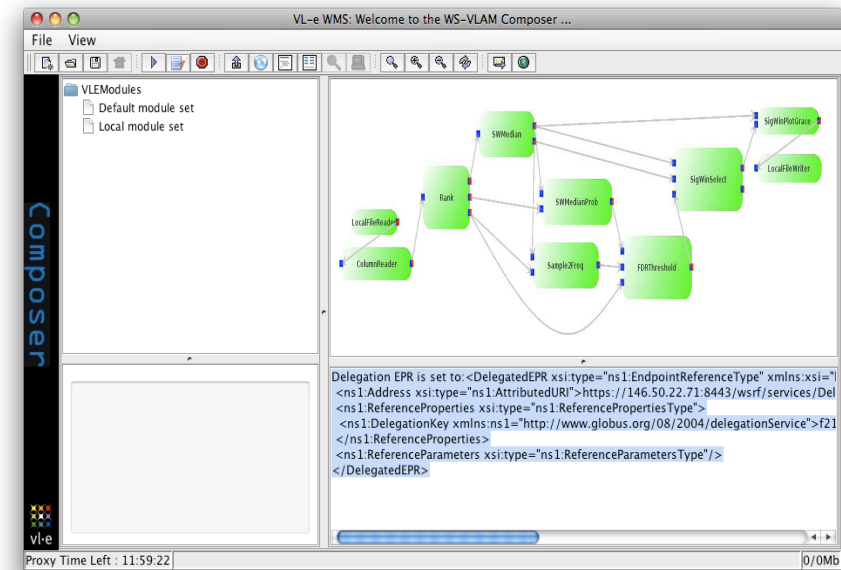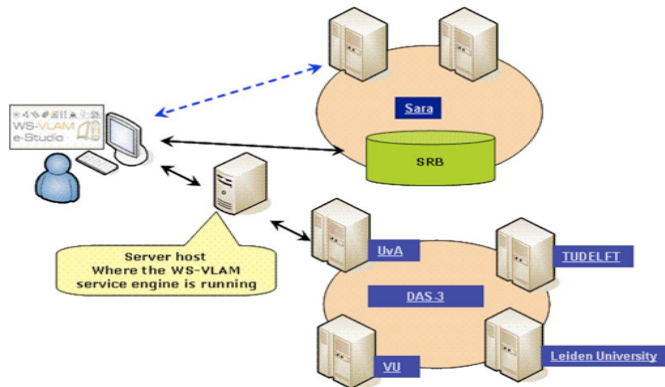
Adam Belloum

Institute of Informatics

University of Amsterdam

a.s.z.belloum@uva.nl

# Outline

- Introduction
- Lifecycle of an e-science workflow
- Workflow management Systems
- Scientific workflows Applications
- Provenance
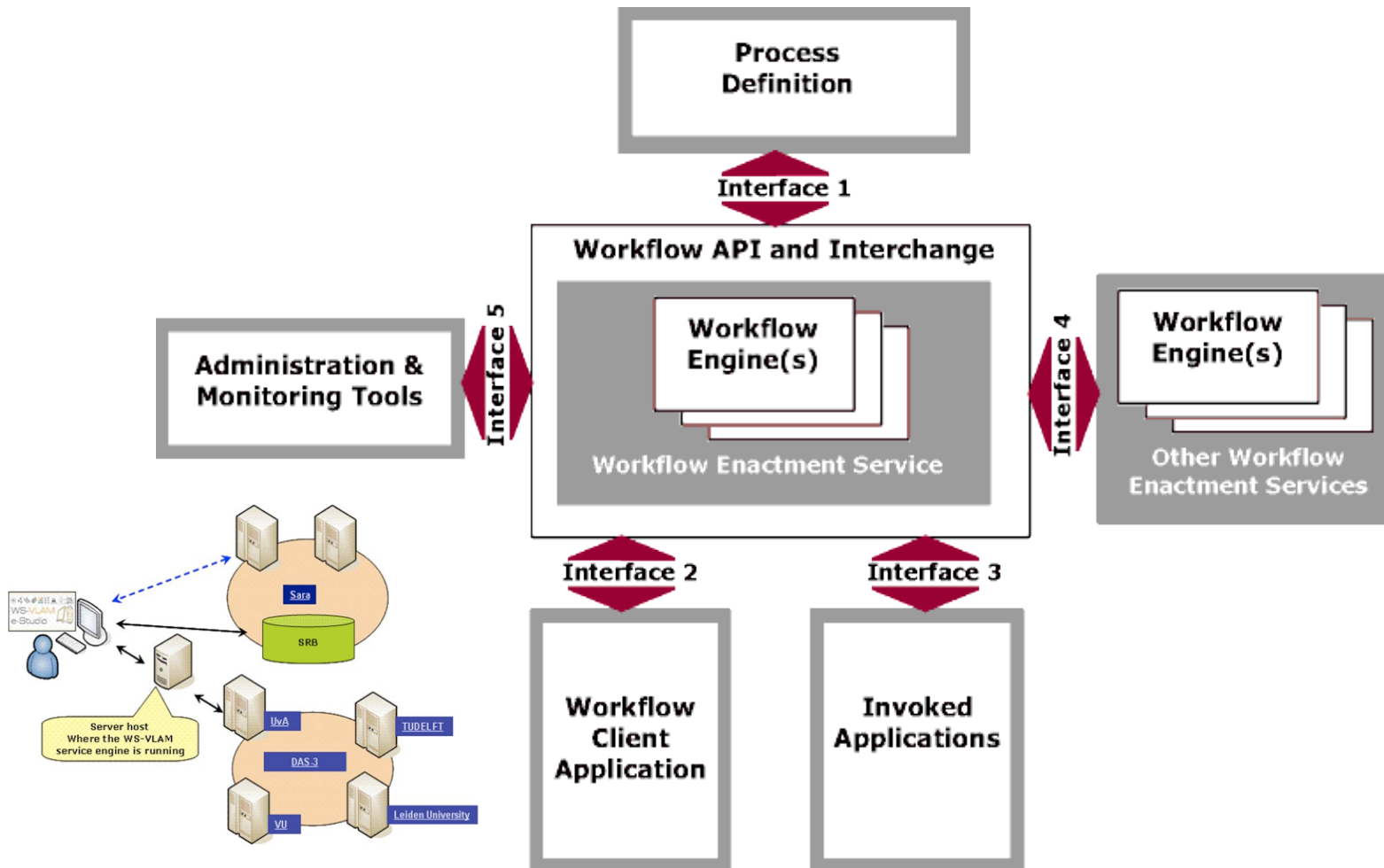- Examples of Scientific workflow managements

# Workflow management system

- **Workflow management system** is a computer program that manages the execution of a workflow on a set of computing resources.



*The user interface of the WS-VLAM* a workflow management system developed in the VL-e project to execute application workflow on geographically distributed computing resources

Deployed as service on Dutch super Computer (DAS3), and Dutch NGI (BigGrid) Clusters
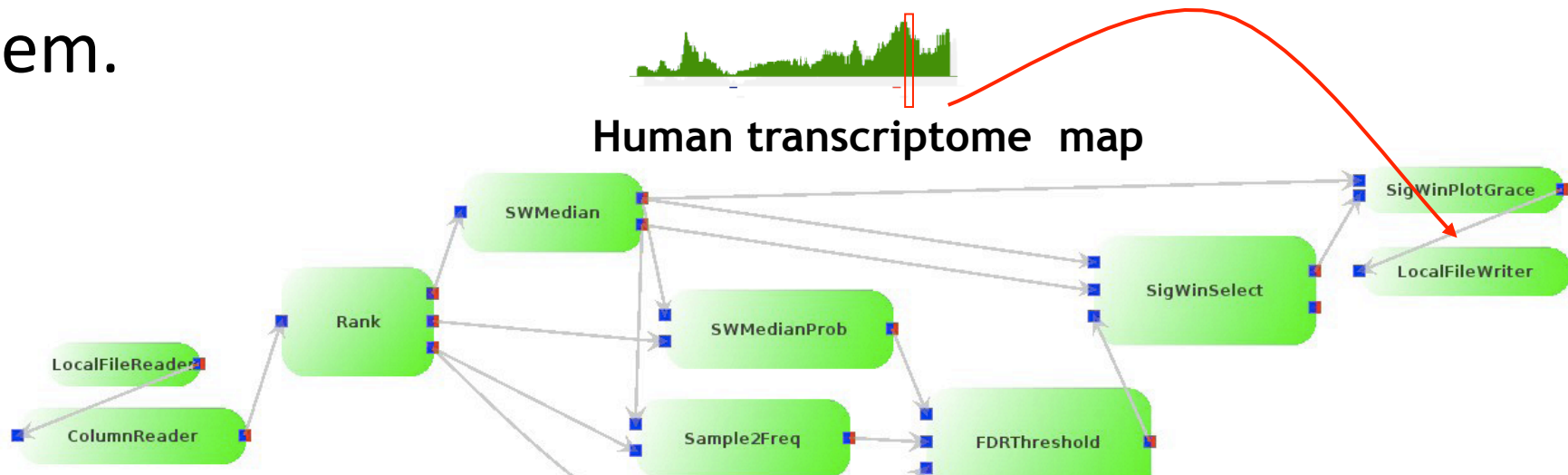
# Reference Model From WFMC



*The automation of a **business process**, in whole or parts, where documents, information or tasks are passed from one participant to another to be processed, according to a set of **procedural rules**. (WFMC definition of a Workflow)*

# Challenges of running workflows on e-infrastructure (grids and clouds)

- ***co-allocate*** *resources needed for workflow enactment across multiple domains*?

- *achieve **QoS** for data centric application workflows that have special requirements on network connections*?

- *achieve **Robustness** and fault tolerance for workflow running across distributed resources*?

- *increase **re-usability** of Workflow, workflow components, and refine workflow execution*?
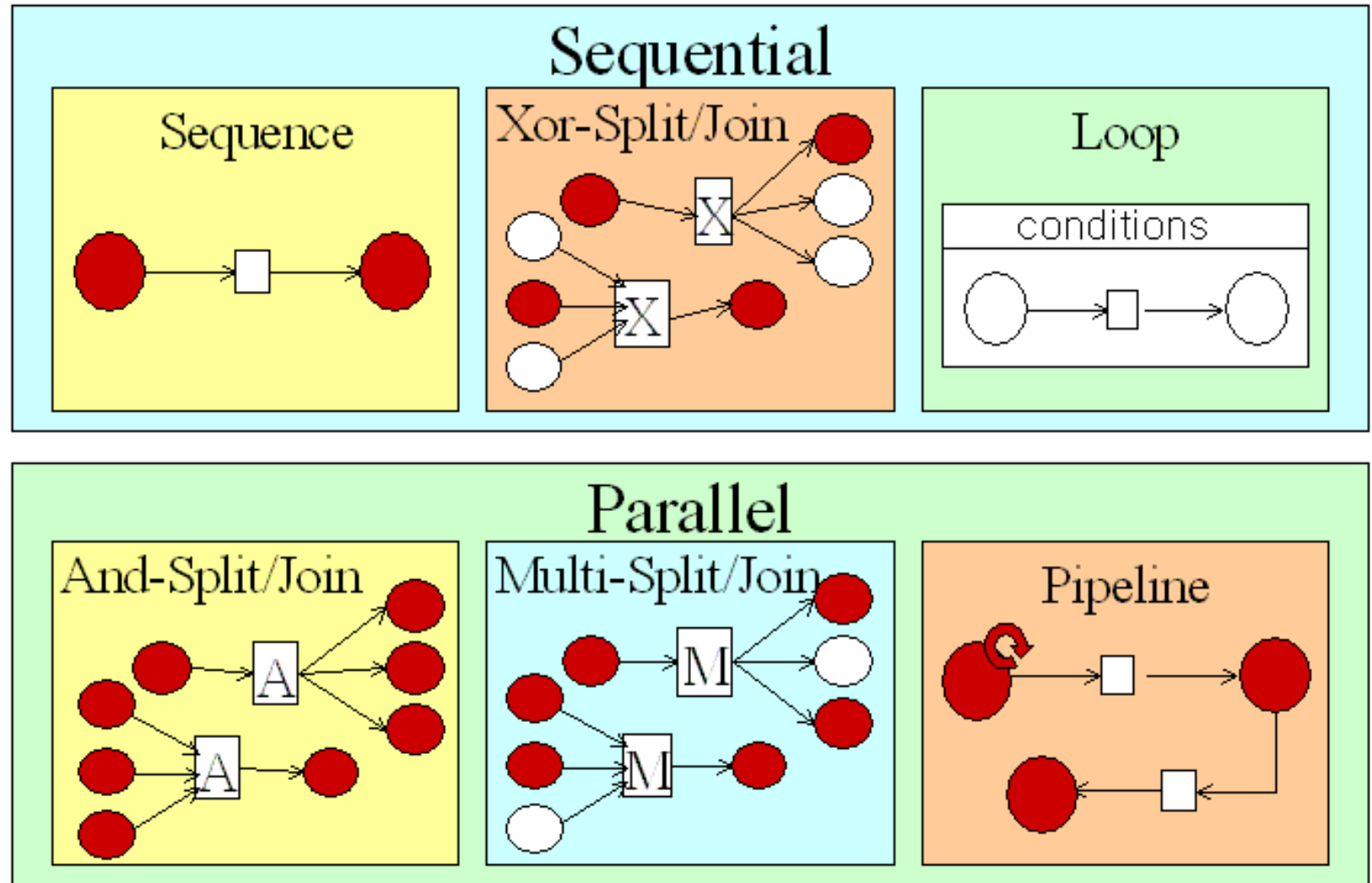
# Workflow

A workflow is a model to represent a **reliably repeatable sequence** of operations/tasks by showing explicitly the interdependencies among them.



**Human transcriptome map**

http://www.youtube.com/watch?v=R6bTFrzaR_w&feature=player_embedded

*SigWin-Detector workflow* has been developed in the VL-e project to detect ridges in for instance a Gene Expression sequence or Human transcriptome map, BMC Research Notes 2008, 1:63 doi: 10.1186/1756-0500-1-63.

# Workflow Pattern

# Business vs Scientific Workflows (Similarities)

- **Capturing knowledge/best practices**
  - Capture **business process** based on the company policy
  - Capture **best practices of scientist**, expert from a specific domain

- **Series of structured activities and computations**
  - Both involves **repeated execution** of certain procedures, and both describes tasks within this procedures.

- **Incorporate human decision in the process**
  - There are exceptional cases that can not be automated both in business and scientific workflow

http://www.csc.ncsu.edu/faculty/mpsingh/papers/databases/ workflows/sciworkflows.html
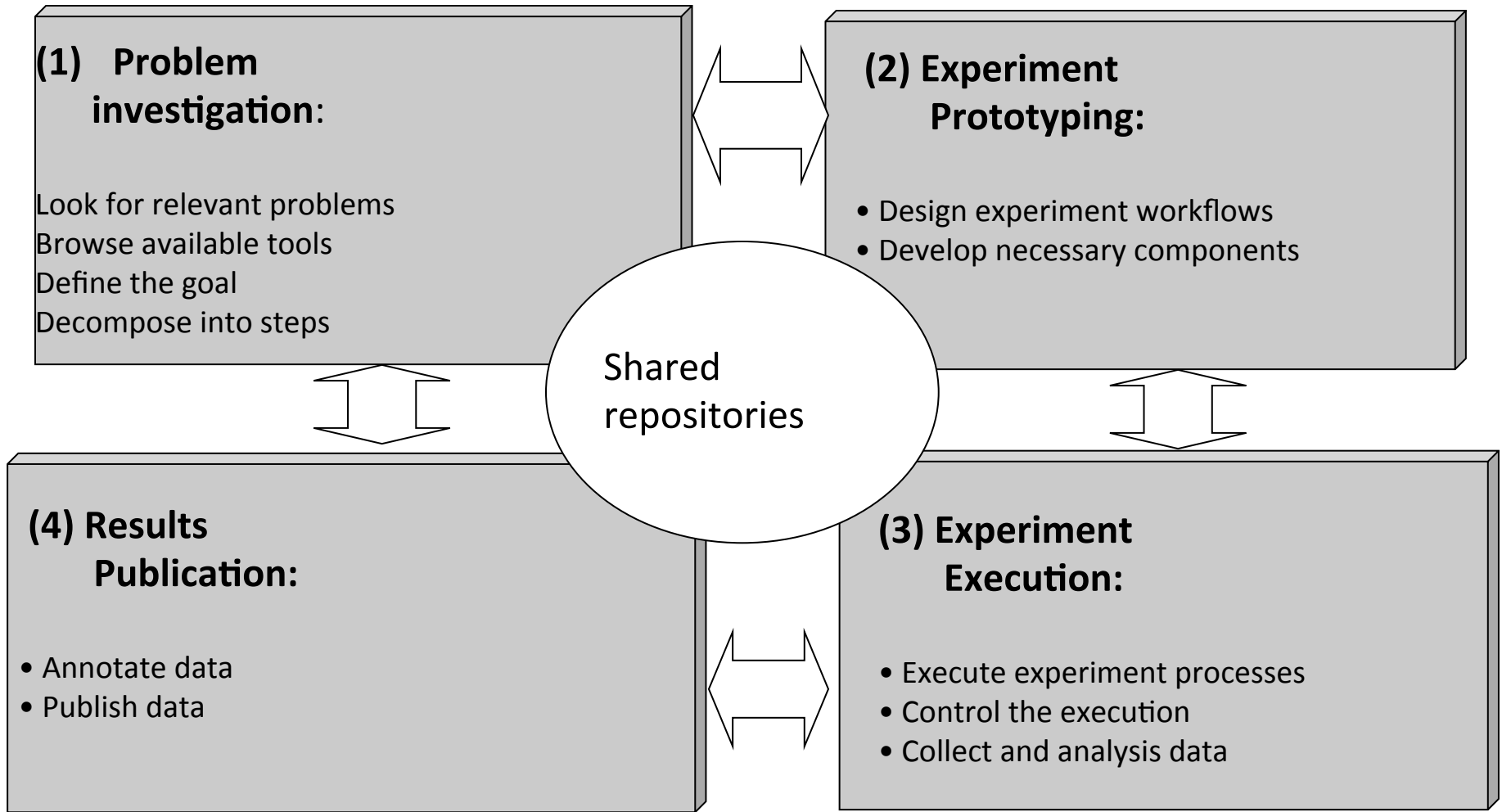
# Business vs Scientific Workflows (Differences)

- **Business Process**
  - **Information**, task, procedural rules of a certain company
  - Driven by business profit goals

- **Static Procedures**
  - Reflecting certain policy within a company
  - Rigid, any changes require approval from management

- **Closed Environment**
  - Managed own resources
  - Within company, actual organization

- **Documents, task descriptions**
  - Flight reservation, credit approval, supply chain, billing, resource planning

- **Scientific process**
  - **Data** analysis, experiment, data manipulation recipes
  - Driven by problem solving goal

- **Dynamic**
  - Exploratory and speculative
  - Flexible, scientist manage their own business (they are their own user/ manager).

- **Open Environment**
  - Non Centralized grid environment
  - Across boundary, Virtual Organizations

- **Large Data**:
  - High energy physics data, bioinformatics micro array/ genomic data etc.

# Scientific Workflow Specific Needs

- What makes scientific workflow different?
  - Need for **large** data flows support
  - Need to do **parameterized** execution of large number of jobs
  - Need to monitor and control workflow execution including **ad-hoc** changes
  - Need to execute in **dynamic** environment where resources are not known a priori and may need to adapt to changes
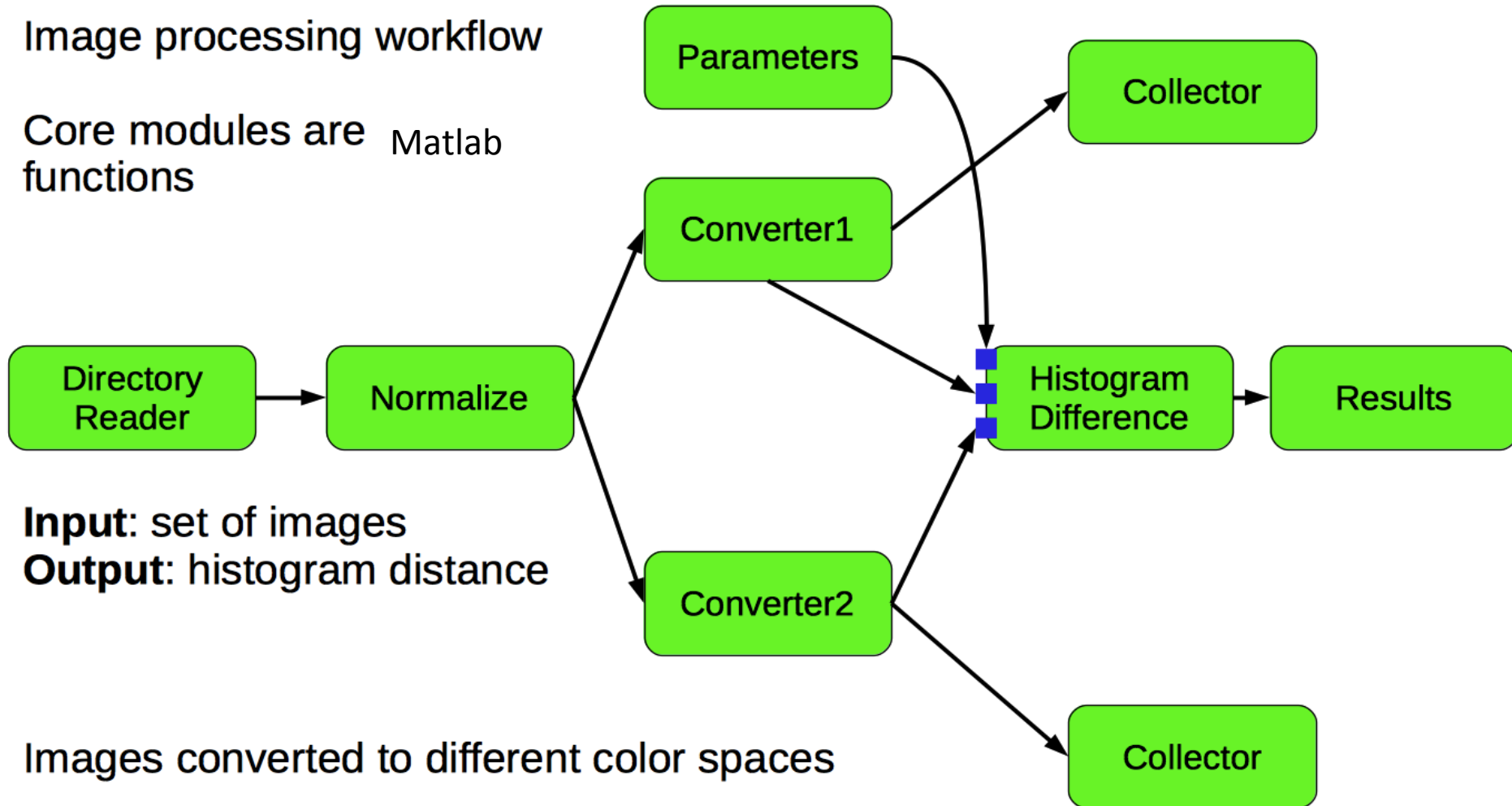  - Hierarchical execution with sub-workflows created and destroyed when necessary

# Complex Scientific experiments model



**(1) Problem investigation:**

Look for relevant problems
Browse available tools
Define the goal
Decompose into steps

**(2) Experiment Prototyping:**

- Design experiment workflows
- Develop necessary components

**Shared repositories**

**(4) Results Publication:**

- Annotate data
- Publish data

**(3) Experiment Execution:**

- Execute experiment processes
- Control the execution
- Collect and analysis data

# Example of Scientific workflow

Image processing workflow

Core modules are Matlab functions

Directory Reader → Normalize → Converter1 → Parameters → Collector

**Input**: set of images
**Output**: histogram distance

Converter2 → Histogram Difference → Results

Collector

Images converted to different color spaces

Histogram difference is calculated between color spaces

# Example of Scientific workflow
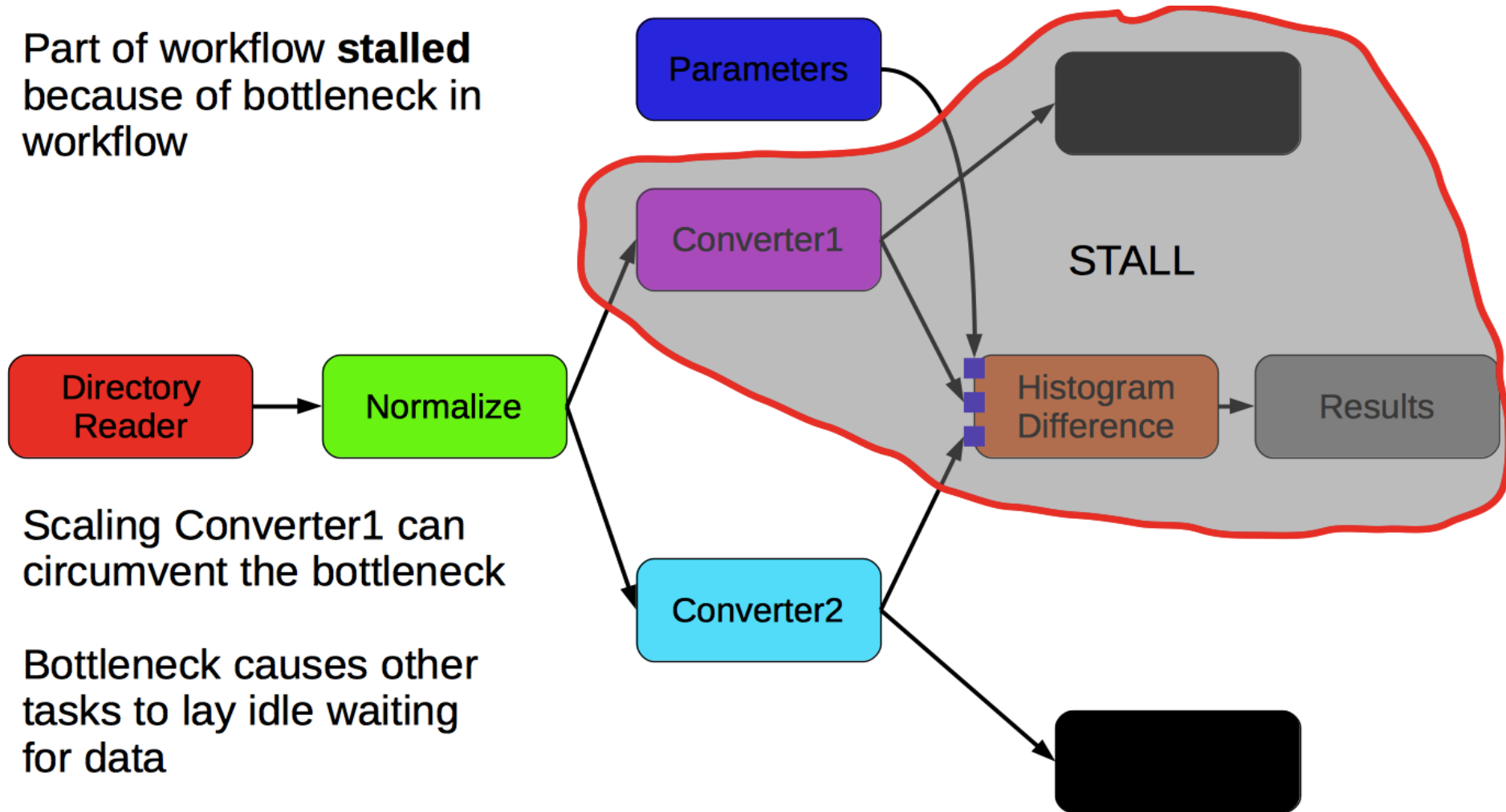


Color coded workflow to better understand the result graphs

Parameters

Converter1

Directory Reader

Normalize

Histogram Difference

Results

Converter2

# Workflow Without Scaling



Slow task causing a **bottleneck** in the workflow

# Example of Scientific workflow (1)



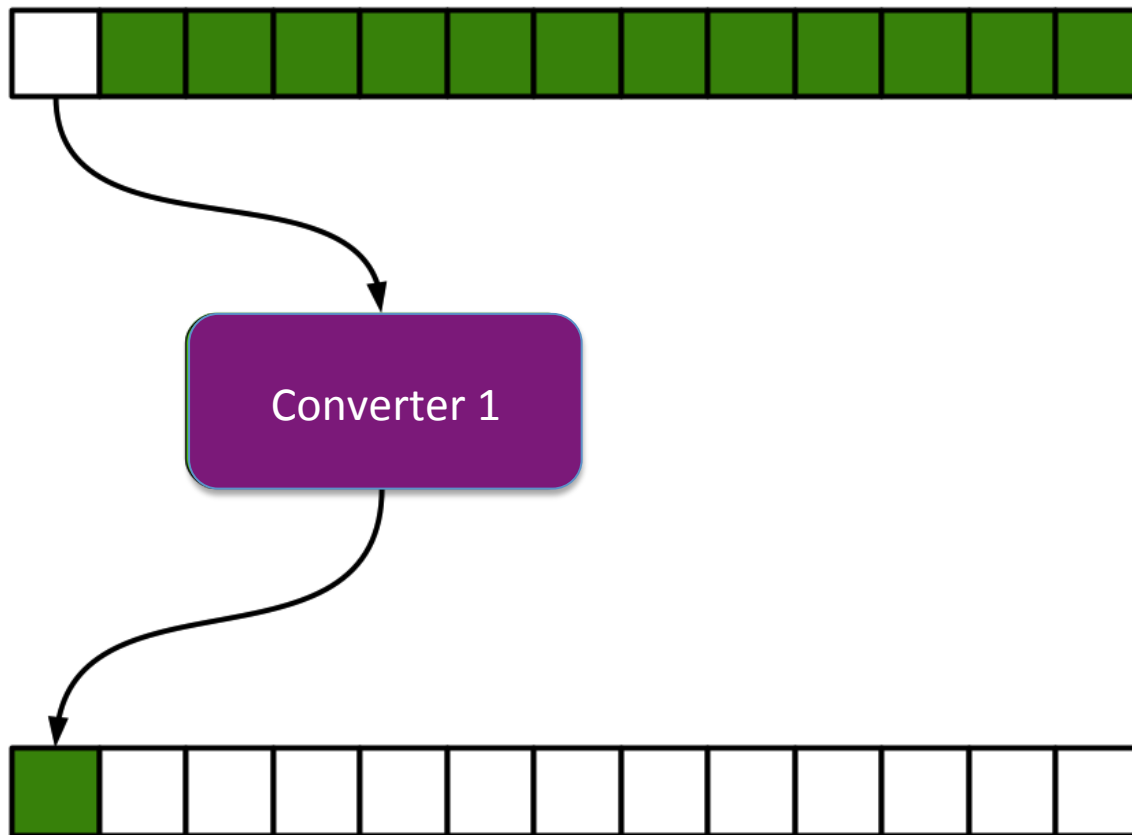Part of workflow **stalled** because of bottleneck in workflow

Scaling Converter1 can circumvent the bottleneck

Bottleneck causes other tasks to lay idle waiting for data

Parameters

Converter1

STALL

Directory Reader

Normalize

Histogram Difference

Results

Converter2

# Zooming into the Task Converter 1



Data organized in atomic parcels(messages)

Converter 1

Part of workflow **stalled** because of bottleneck in workflow

Scaling Converter1 can circumvent the bottleneck

Bottleneck causes other tasks to lay idle waiting for data

Parameters

Converter1

STALL

Directory Reader

Normalize

Histogram Difference

Results

Converter2

# Zooming into the Task Converter 1



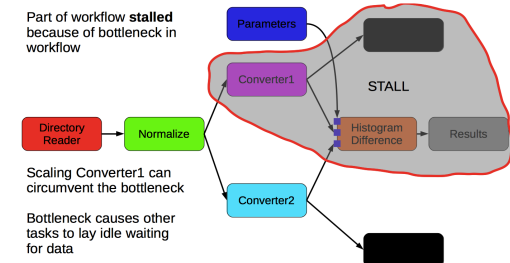Data organized in atomic parcels(messages)

Task processes data sequentially

Part of workflow **stalled** because of bottleneck in workflow
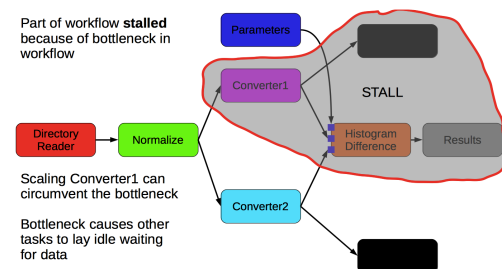
Scaling Converter1 can circumvent the bottleneck

Bottleneck causes other tasks to lay idle waiting for data

Parameters

Converter1

STALL

Directory Reader

Normalize

Histogram Difference

Results

Converter2

# Scaling Concepts

Data organized in atomic parcels(messages)

Task processes data sequentially

**Converter 1**

Part of workflow **stalled** because of bottleneck in workflow

Scaling Converter1 can circumvent the bottleneck

Bottleneck causes other tasks to lay idle waiting for data

Parameters

Converter1

STALL

Directory Reader → Normalize

Histogram Difference → Results

Converter2

# Zooming into the Task Converter 1

Data organized in atomic parcels(messages)

Task processes data sequentially

Converter 1

Part of workflow **stalled** because of bottleneck in workflow

Scaling Converter1 can circumvent the bottleneck

Bottleneck causes other tasks to lay idle waiting for data

Parameters

Converter1

STALL

Directory Reader

Normalize

Histogram Difference

Results

Converter2

# Zooming into the Task Converter 1



Data organized in atomic parcels(messages)

Tasks processes data **concurrently**

Adding more tasks increases **message consumption** rate

Converter 1

Part of workflow **stalled** because of bottleneck in workflow

Scaling Converter1 can circumvent the bottleneck

Bottleneck causes other tasks to lay idle waiting for data

Parameters

Converter1

STALL

Directory Reader

Normalize

Histogram Difference

Results

Converter2

# Zooming into the Task Converter 1



Data organized in atomic parcels(messages)

Task processes data sequentially

Adding more tasks increases **message consumption** rate

**Challenge:** How many tasks to create?

Too **many** and tasks get stuck on queues. Too **few** and optimal performance not achieved
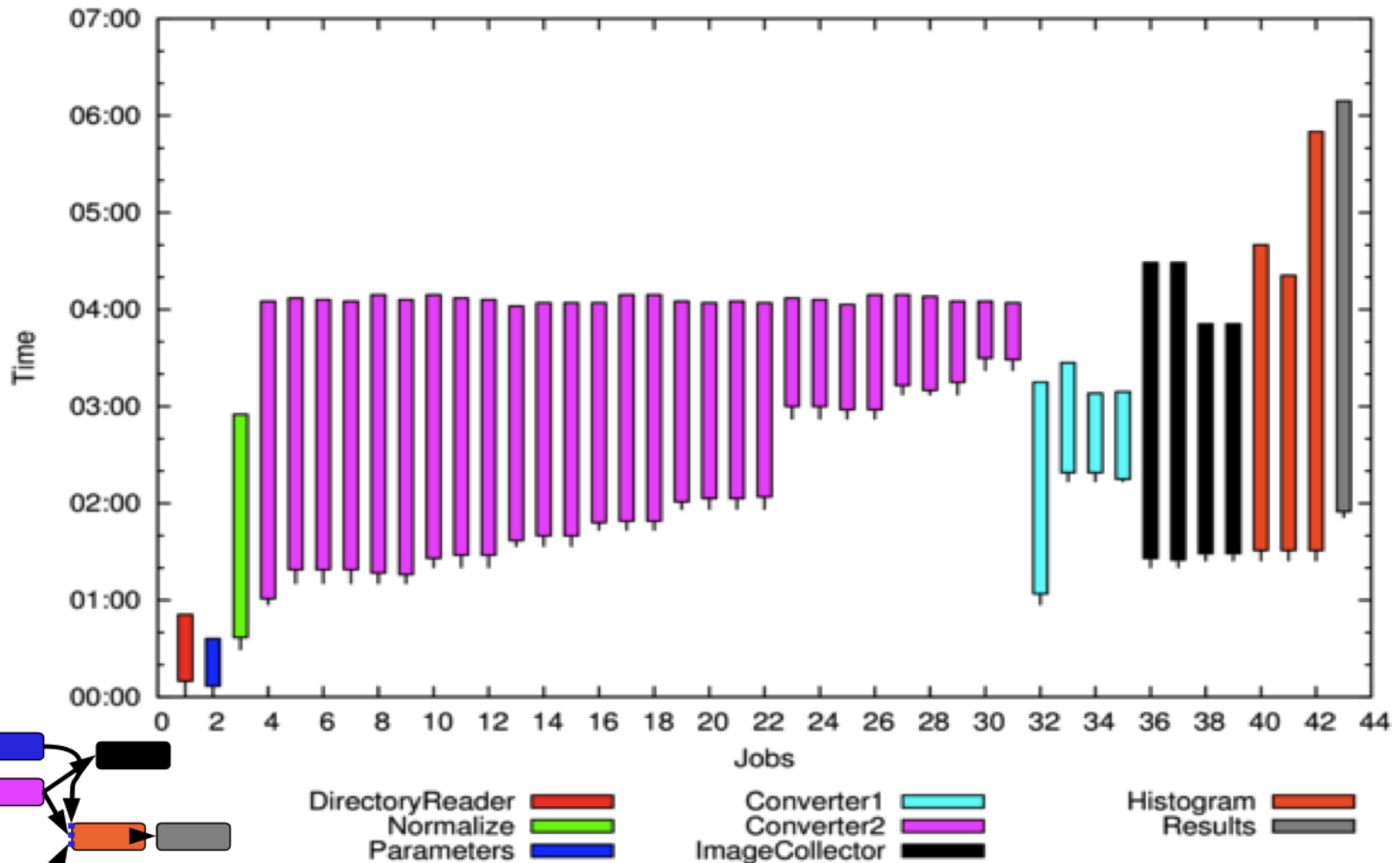
# Load Prediction



**Simplified Load** = $6x/k$ time slots

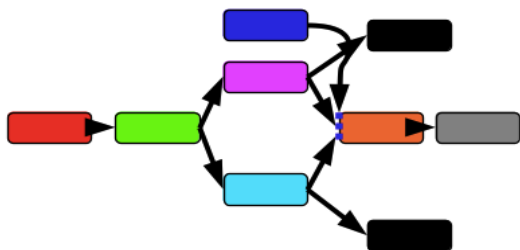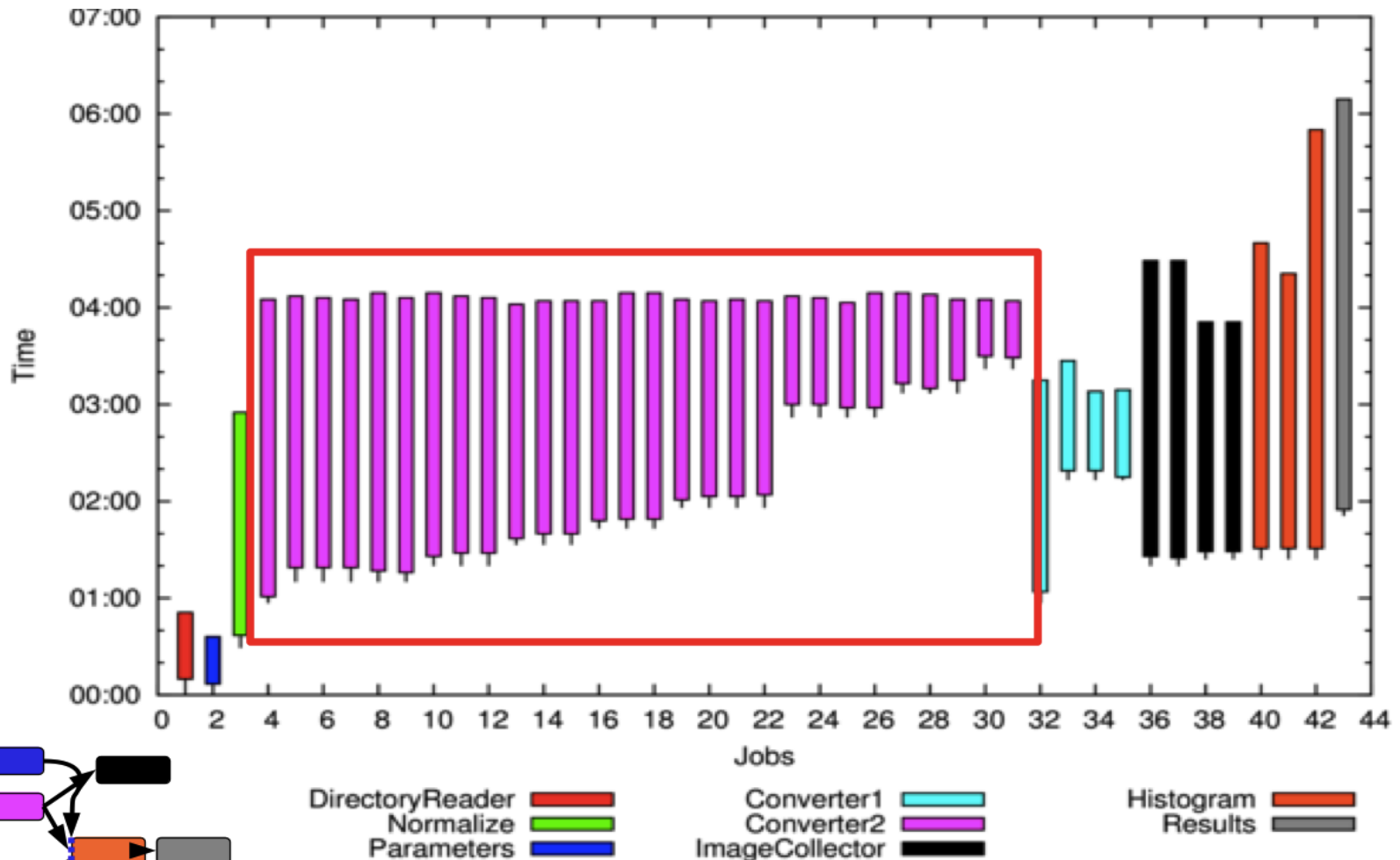**Assumption**: Size of data directly proportional to computation time. May not always be the case
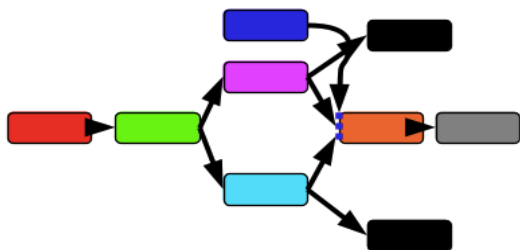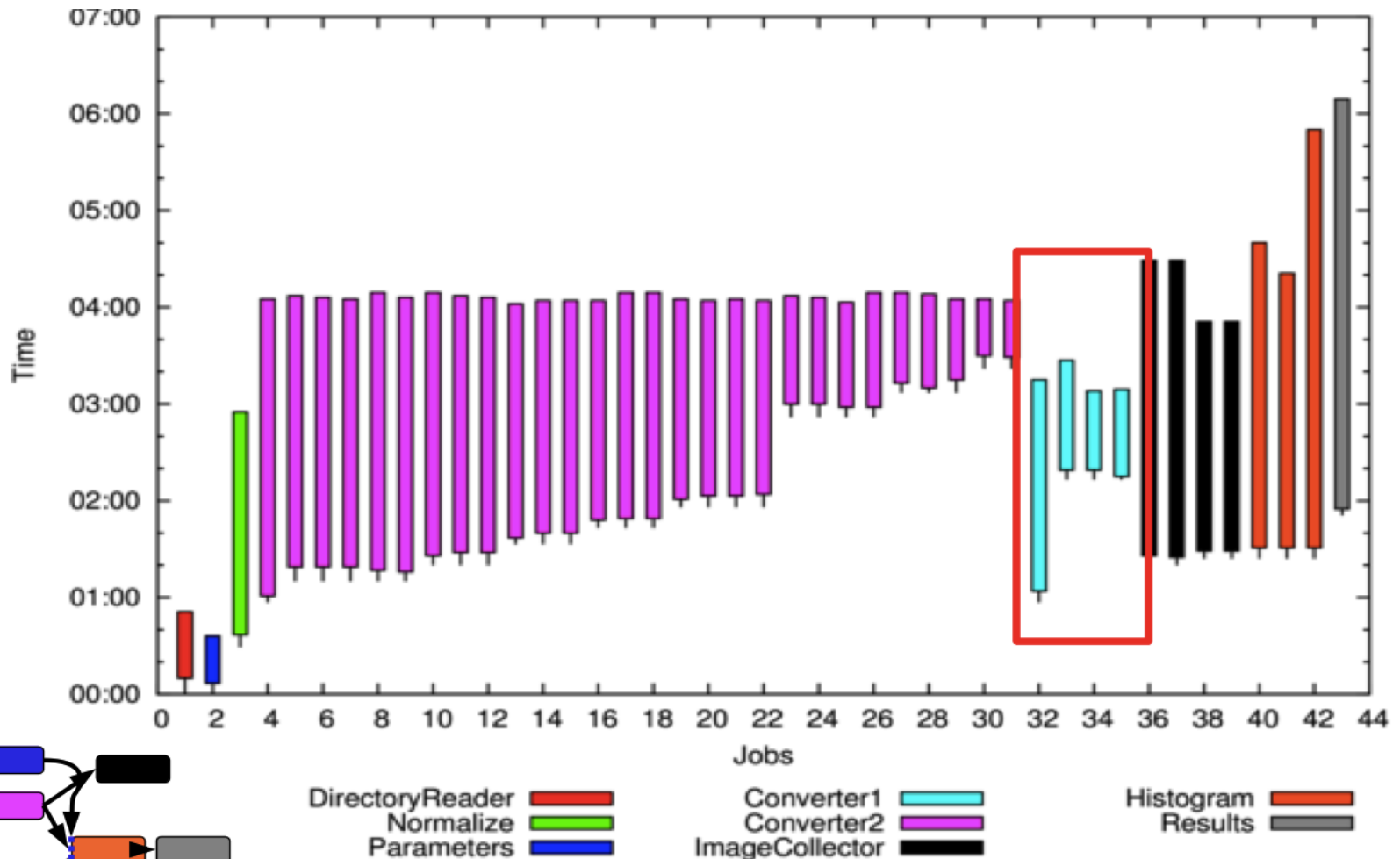
# Workflow execution with Scaling

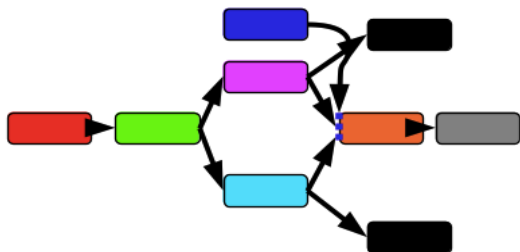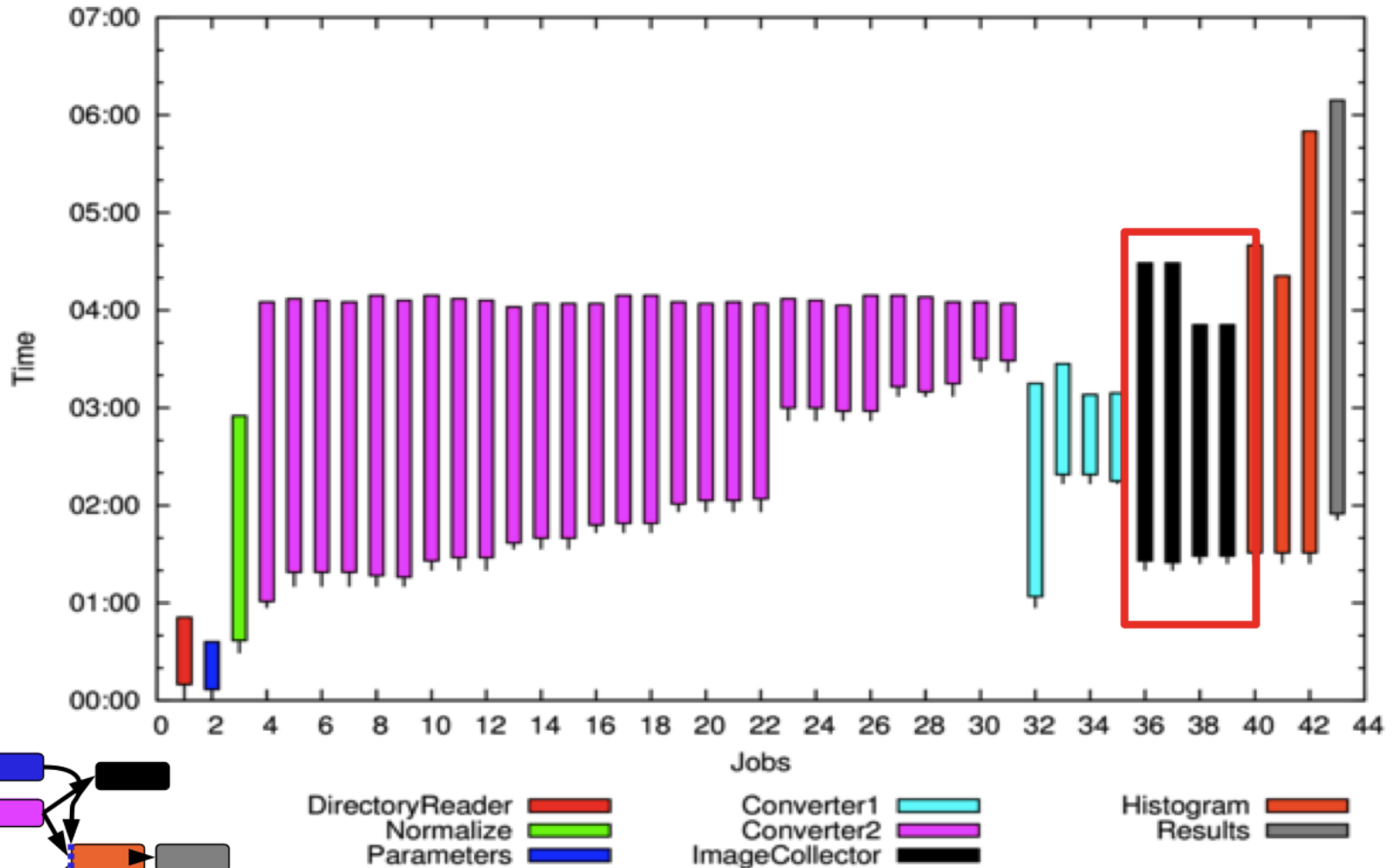# Workflow execution with Scaling

# Workflow execution with Scaling Task -1
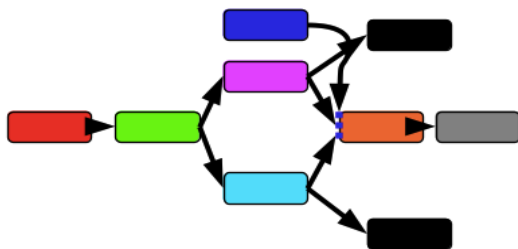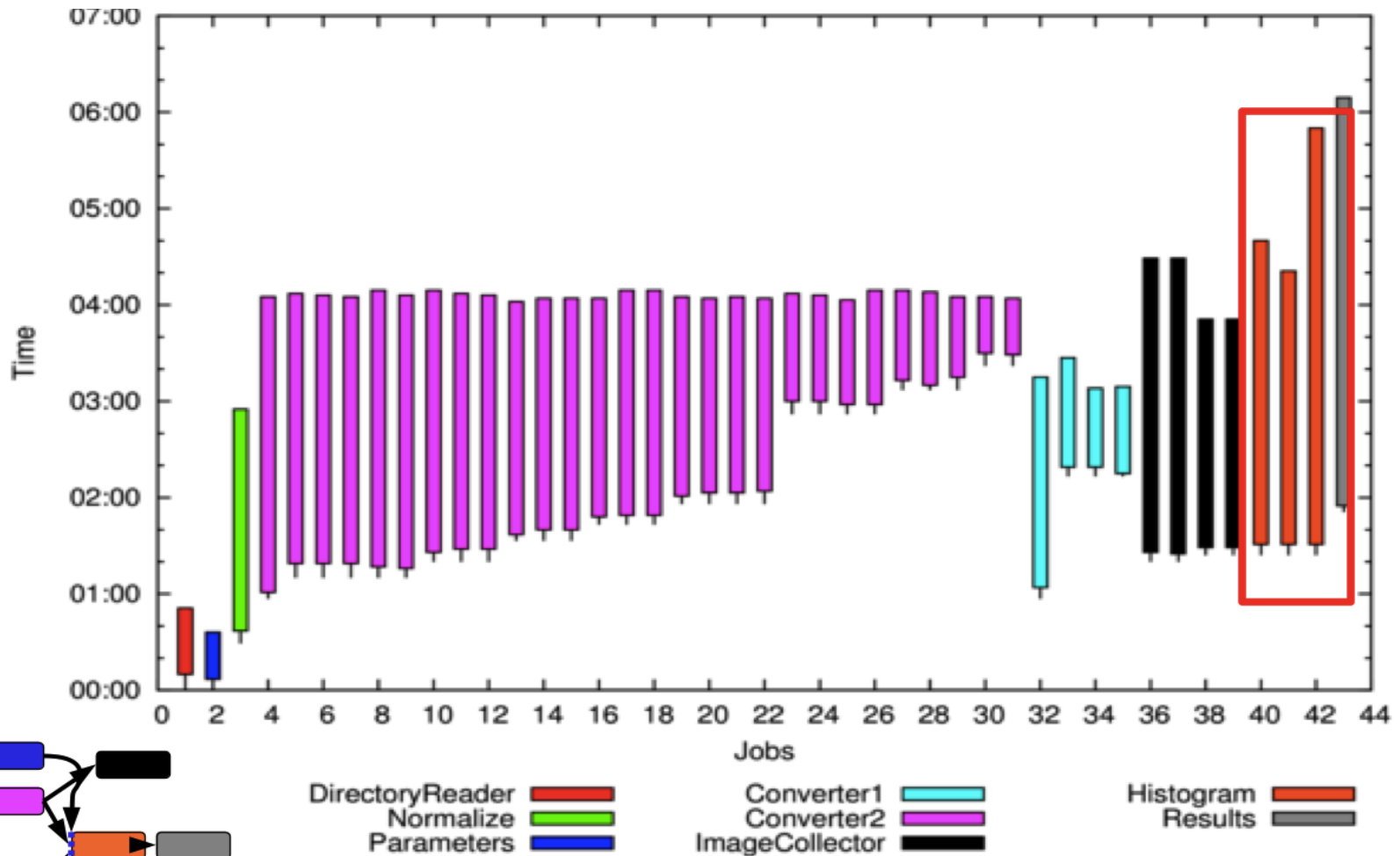
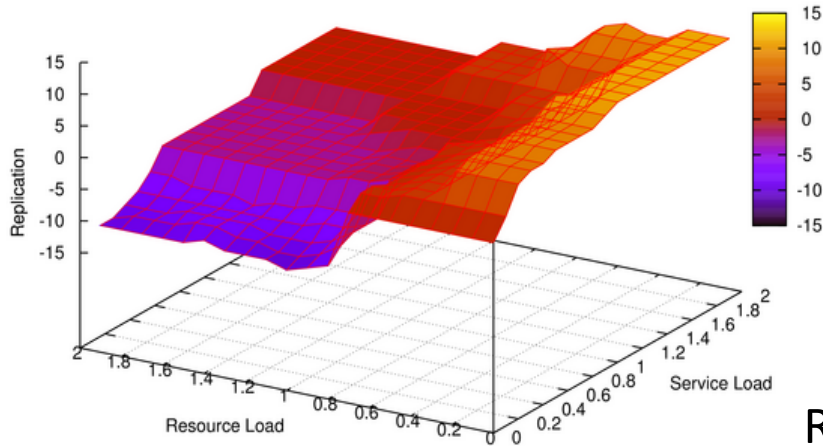# Workflow execution with Scaling Task -2

# Other Scaled Task -1



ImageCollector was set to a **fixed** amount (4)

# Auto Scaling Task -2
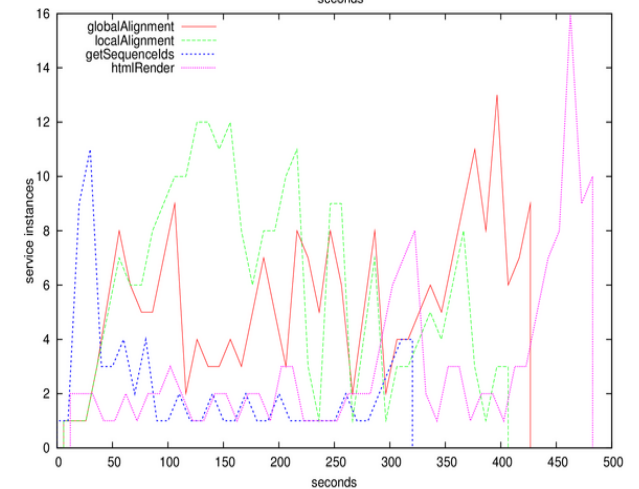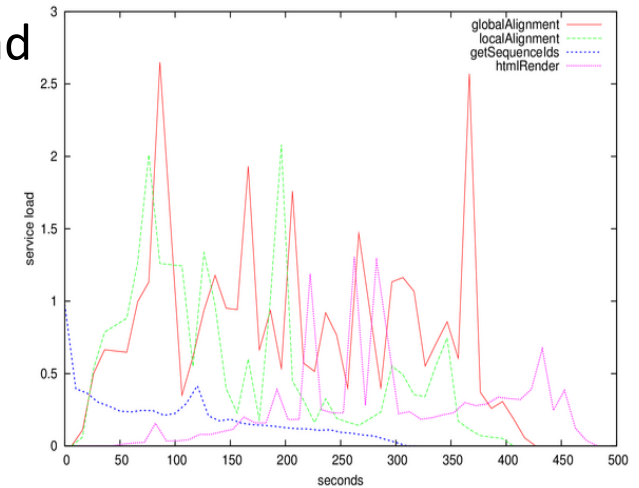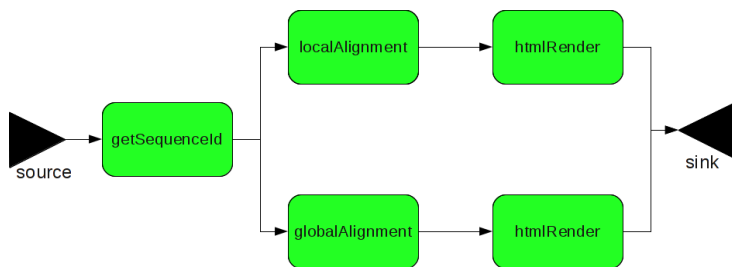


HistogramDifference was set to **one-to-one** scaling.
Each parameter generates a new task

# Example of Scientific workflow (2)

Service Load



Running Service instances



Reginald Cushing, Spiros Koulouzis,  Adam S. Z. Belloum, Marian Bubak, **Dynamic Handling for Cooperating Scientific Web Services**, 7th IEEE International Conference on e-Science, December 2011, Stockholm, Sweden

# Workflow as a Service (WFaaS)

- Once a workflow is initiated on the resources it stays alive and process data/jobs continuously

- Reduce the scheduling overhead



Reginald Cushing, Adam S. Z. Belloum, V. Korkhov, D. Vasyunin, M.T. Bubak, C. Leguy
ECMLS'12, June 18, 2012, Delft, *Workflow as a Service: An Approach to Workflow Farming*, The Netherlands

# Workflow Issues

- **Workflow description**

  How to capture **knowledge** of expert while still **hiding** complexity of underlying system.

  - **Workflow Models:** allow to **model** the tasks and **dependencies** between them (DAG, Petri Net)

  - **Workflow languages:** provide the required support to express the workflow model.

- **Workflow Enactment:** The functions provided by enactment are scheduling, fault management and data movement.
  - In the context of Grid environment workflow enactment service can be built on the top of low level Grid middleware

# Workflow Enactment

- Workflow Refinement
  - Modification from the workflow description
  - Reduction of workflow if some data already exist
  - Additional data movement preparation if needed

- Mapping to actual resource
  - Resource discovery, allocation and management
  - Bind to real computing resource

- Workflow Fault Tolerance & Monitoring of Execution
  - Two level failure recovery techniques
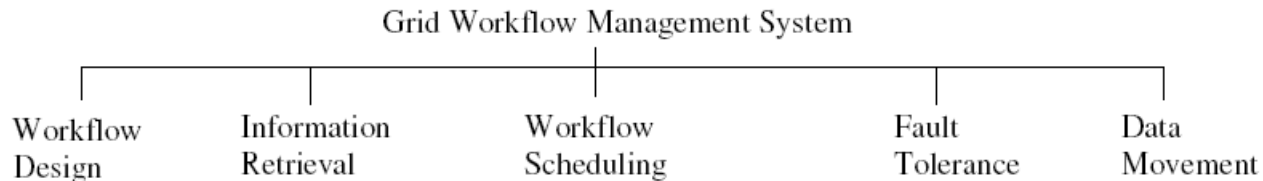    - Task Level
    - Workflow Level

# Model of computation

- Model of computation: stream-based process network.
  - Engine **co-allocates** all workflows.
  - Components waste time idling.
  - Co-allocation difficult.

- Communication: time **coupled**
  - Assumes components are running
  - Simultaneously
  - Synchronized p2p
  - Fixed TCP/IP

V. Korkhov et al. VLAM-G: Interactive data driven workflow engine for Grid-enabled resources, Scientific Programming 15 (2007) 173–188 173 IOS Press
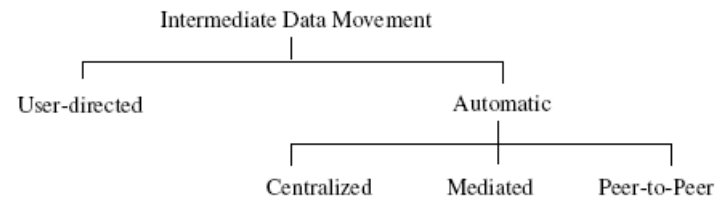
# Model of computation

- Model of computation: dataflow network
  - components **scheduled** depending on data
  - components **only activated** when data is available
  - **no need** for **co-allocation**

- Communication: time decouples
  - messaging communication system.
  - components not synchronized
  - communication not strictly TCP/IP

# Workflow Taxonomy

Grid Workflow Management System

Workflow Design | Information Retrieval | Workflow Scheduling | Fault Tolerance | Data Movement

Intermediate Data Movement

User-directed | Automatic

Centralized | Mediated | Peer-to-Peer

- For Grid workflow applications,
  - the input files of tasks need to be staged to a remote site before processing the task.

  - Similarly, output files may be required by their children tasks which are processed on other resources.

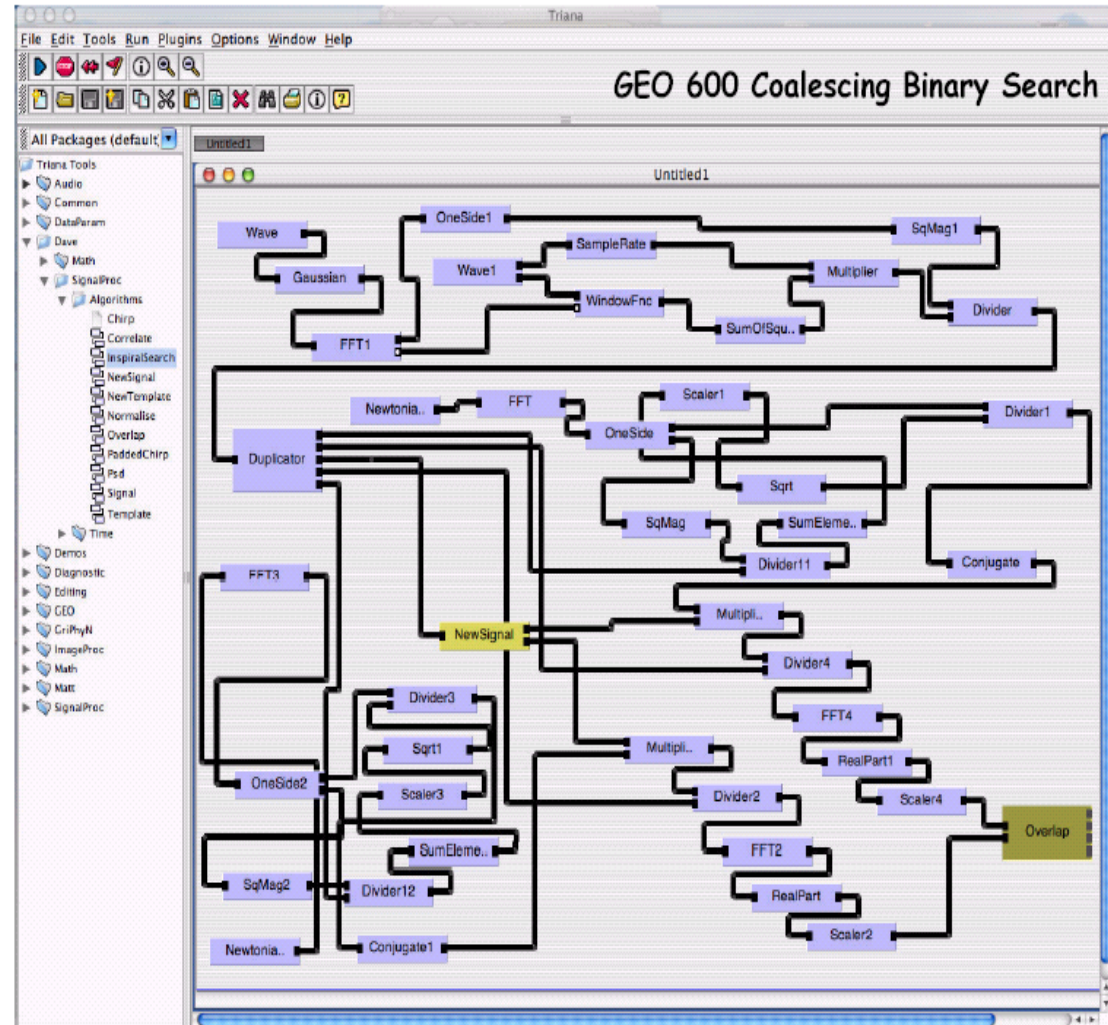- The intermediate data has to be staged out to the corresponding Grid sites.

A Taxonomy of Workflow Management Systems for Grid Computing
Jia Yu and Rajkumar Buyya, http://www.cloudbus.org/reports/GridWorkflowTaxonomy.pdf

# Component Based Workflow Description: Triana
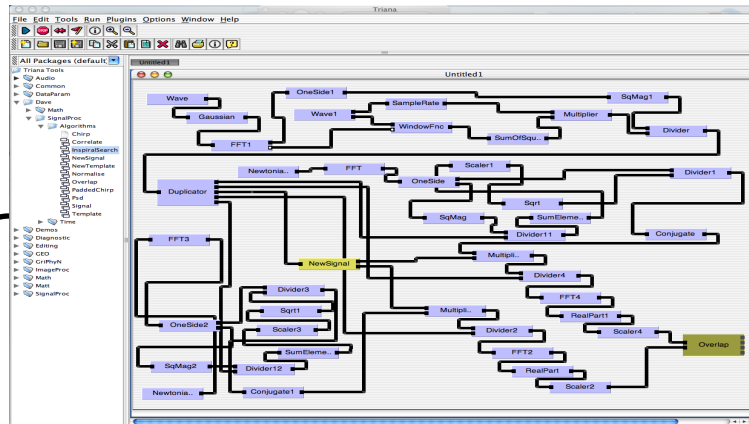
- Workflow design
  - workflow structure: Non-DAG
  - workflow model: concrete
  - workflow composition: user-directed: Graph-based Modeling: User-defined component
- Information retrieval
  - GAT
- Scheduling
  - Architecture: decentralized
  - Decision Making: local
  - Planning scheme: just-in-time
  - Scheduling strategy: GAT
- Performance estimation: N/A
- Fault tolerance :GAT
- Data movement: P2P

# Triana, the GAT and the GAP
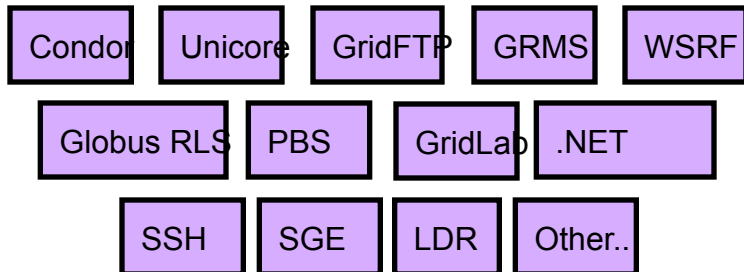


**Grid Computing:**

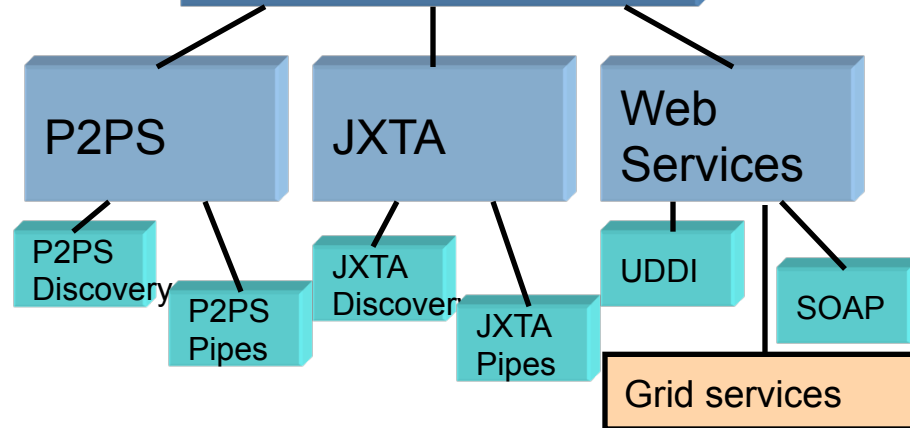Job Submission, data services On top of a number of Grid Middleware

**Service Based Computing:**

Deployment, discovery and communication with distributed services e.g. P2P and (GSI) Web services

## GAT Interface
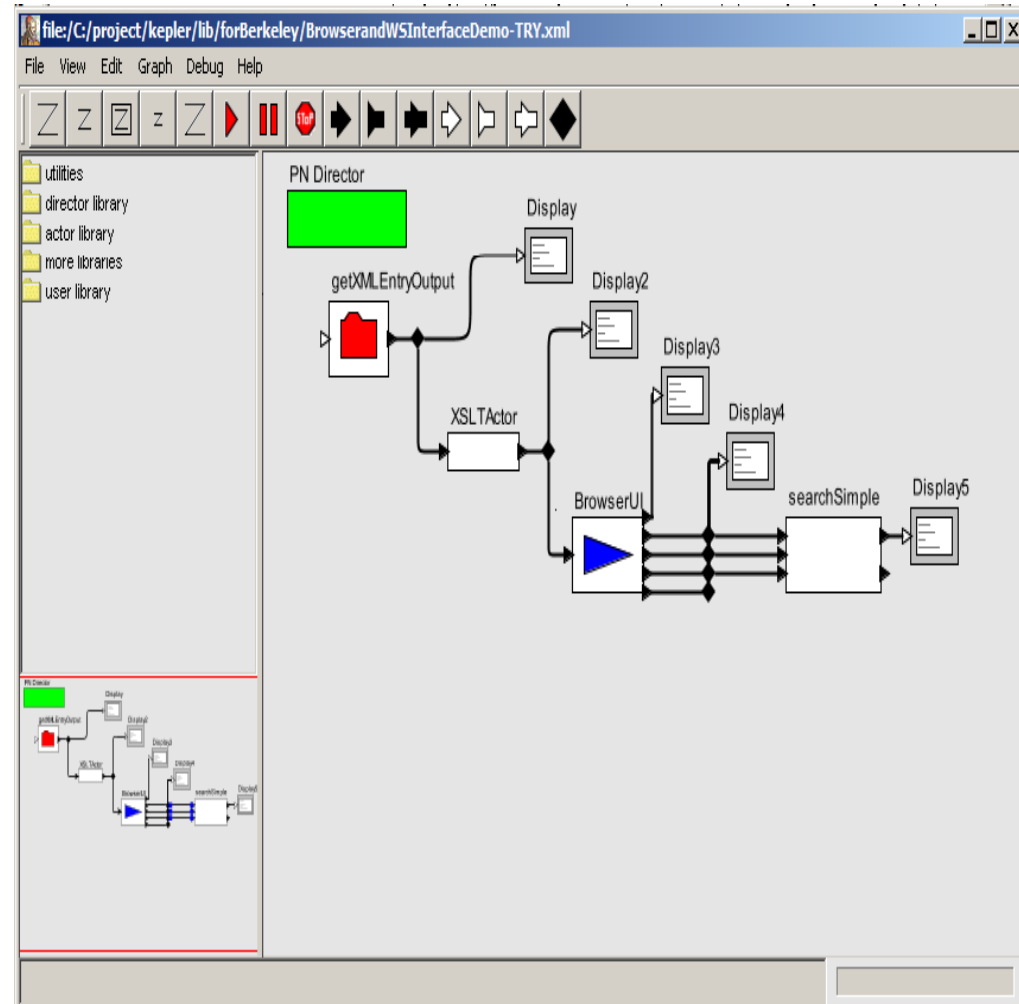
| Condor | Unicore | GridFTP | GRMS | WSRF |

| Globus RLS | PBS | GridLab | .NET |

| SSH | SGE | LDR | Other.. |

## GAP Interface

| P2PS | JXTA | Web Services |

P2PS Discovery
P2PS Pipes
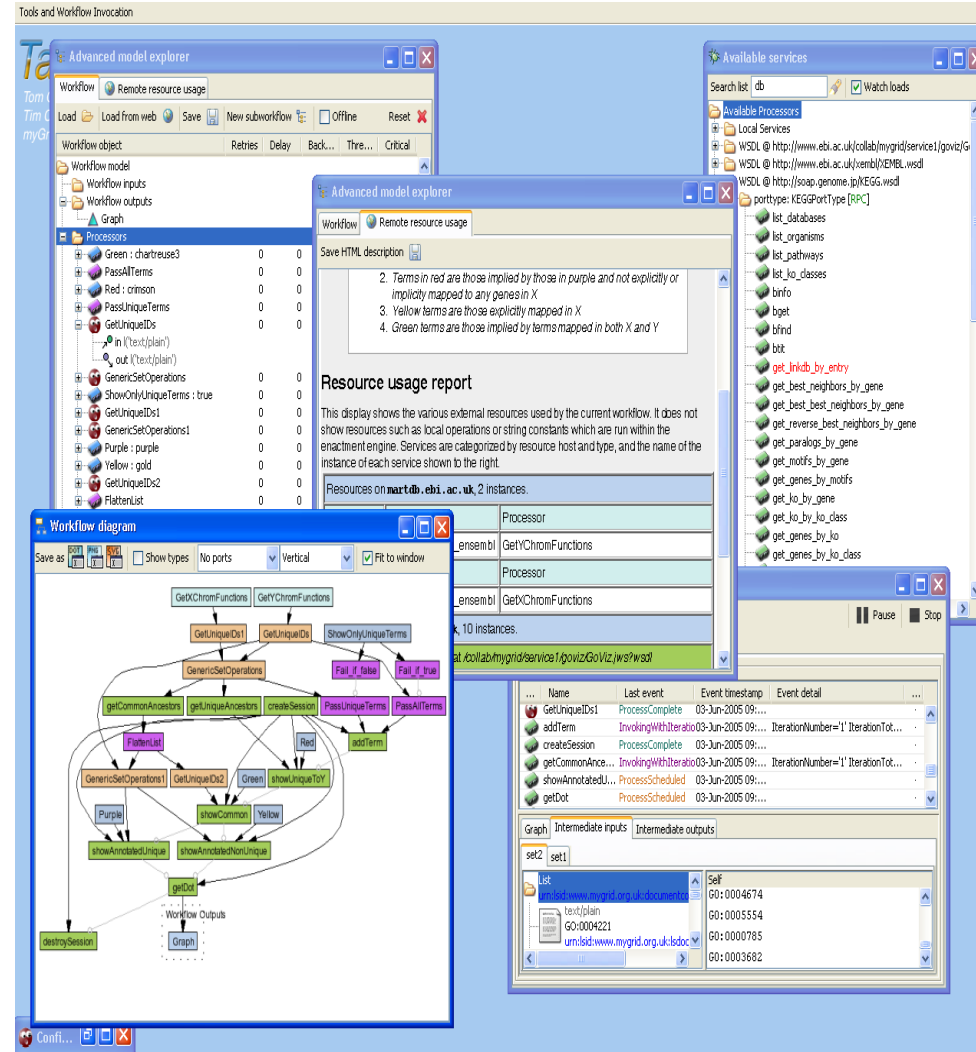
JXTA Discover
JXTA Pipes

UDDI
SOAP

Grid services

# Component Based Workflow Description: Kepler

- Workflow design
  - workflow structure: Non-DAG
  - workflow model: concrete
  - workflow composition: user-directed: Graph-based Modeling: User-defined component

- Scheduling
  - architecture: centralized
  - Decision Making: local
  - Planning scheme: static
  - Scheduling strategy: Performance-driven

- Performance estimation: N/A

- Fault tolerance: N/A

- Data movement: P2P

# Program/Application: workflow Based: Taverna
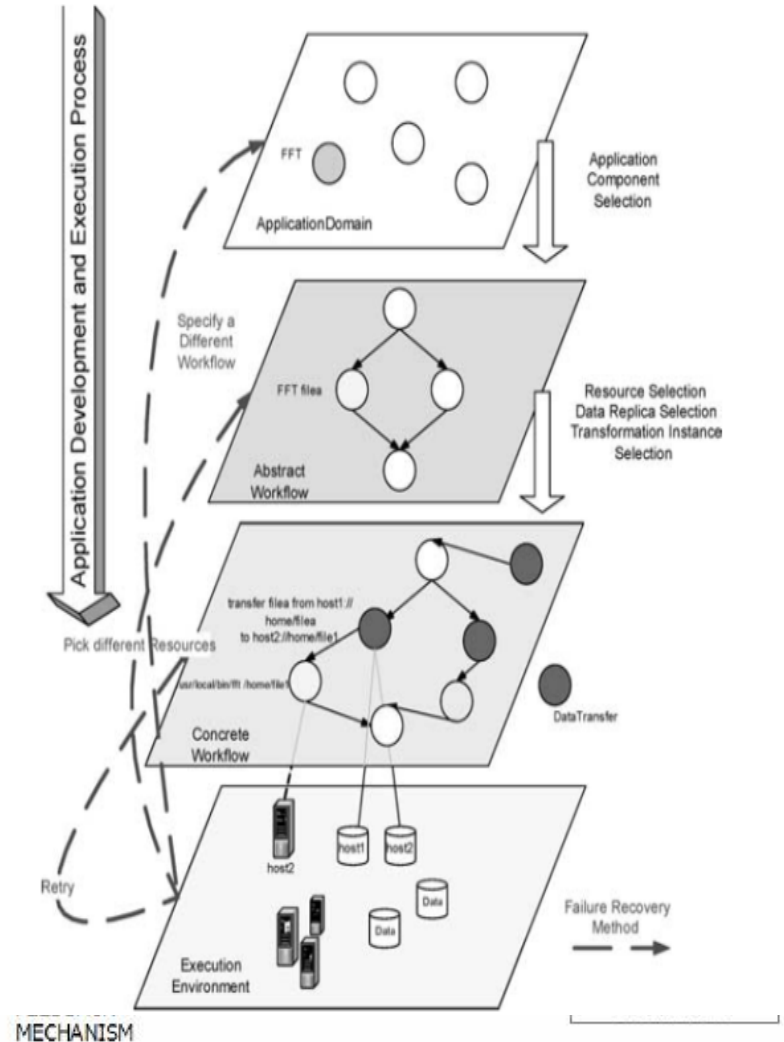
- Workflow design
  - workflow structure: DAG
  - workflow model: concrete
  - workflow composition: user-directed: Graph-based Modeling: User-defined component

- Information retrieval
  - Dynamic: execution related

- Scheduling
  - Architecture: centralized
  - Decision Making: local
  - Planning scheme: just-in-time
  - Scheduling strategy: N/A
  - Performance estimation: N/A

- Fault tolerance: Task-level (retry, Alternate)

- Data movement: Centralized

# Pegasus (GriPhyN)

- Pegasus Planner and Condor G Manager

- Pegasus convert abstract workflow into concrete Workflow
  - Prepare for data movements
  - Reduce workflow if data exist
  - Assign Resources to Processes

- Condor G Manager



The process of developing data intensive applications for Grid environments.

# Workflow Refinement

- Example of a simple abstract workflow in which
  - the logical component *Extract* is applied to an input file with a logical filename F.a.
  - The resulting files F.b1 and F.b2, are used as inputs to the components identified by logical filenames *Resample* and *Decimate*.
  - Finally, the results are *Concatenated*

- If we assume that F.c2 is already available
  1. Reduces the workflow to 3 components, namely *Extract*, *Resample*, and *Concat*.
  2. Adds the transfer nodes for transferring F.c2 and F.a from their current locations.
  3. Adds transfer nodes between jobs that will run on different locations.
  4. Adds output transfer nodes to stage data out and registration nodes if the user requested that the resulting data be published and made available at a particular location.
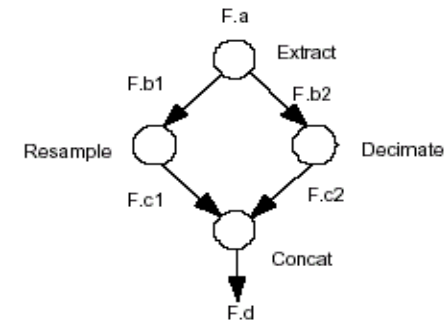


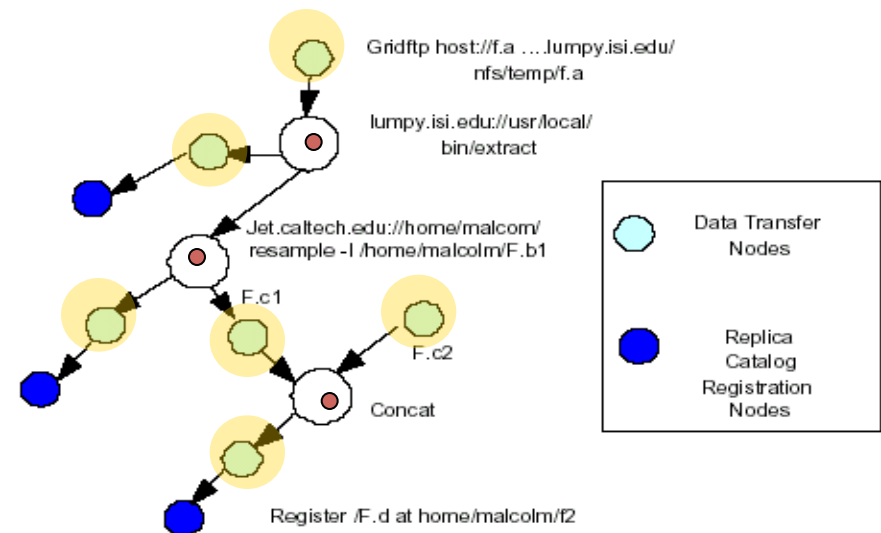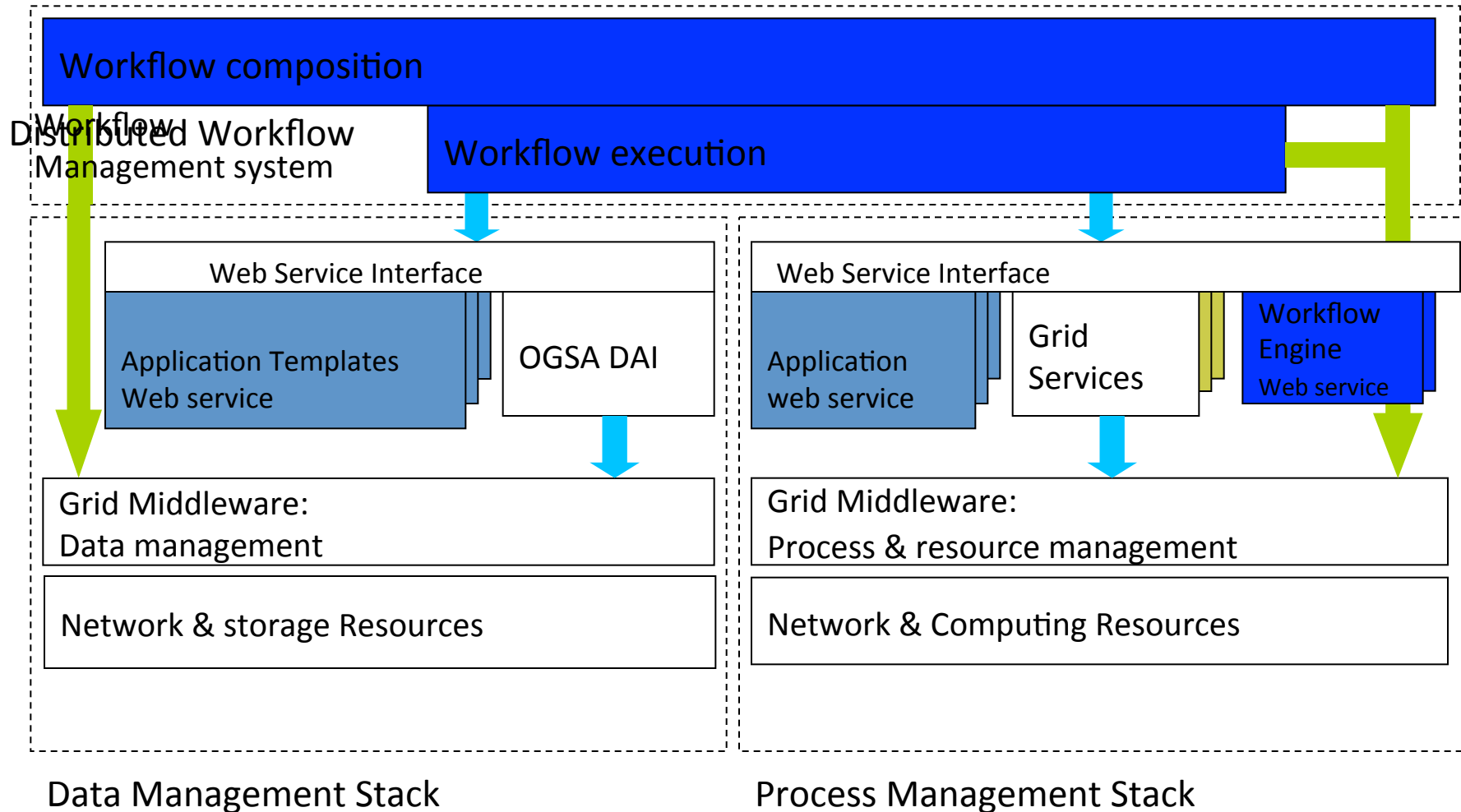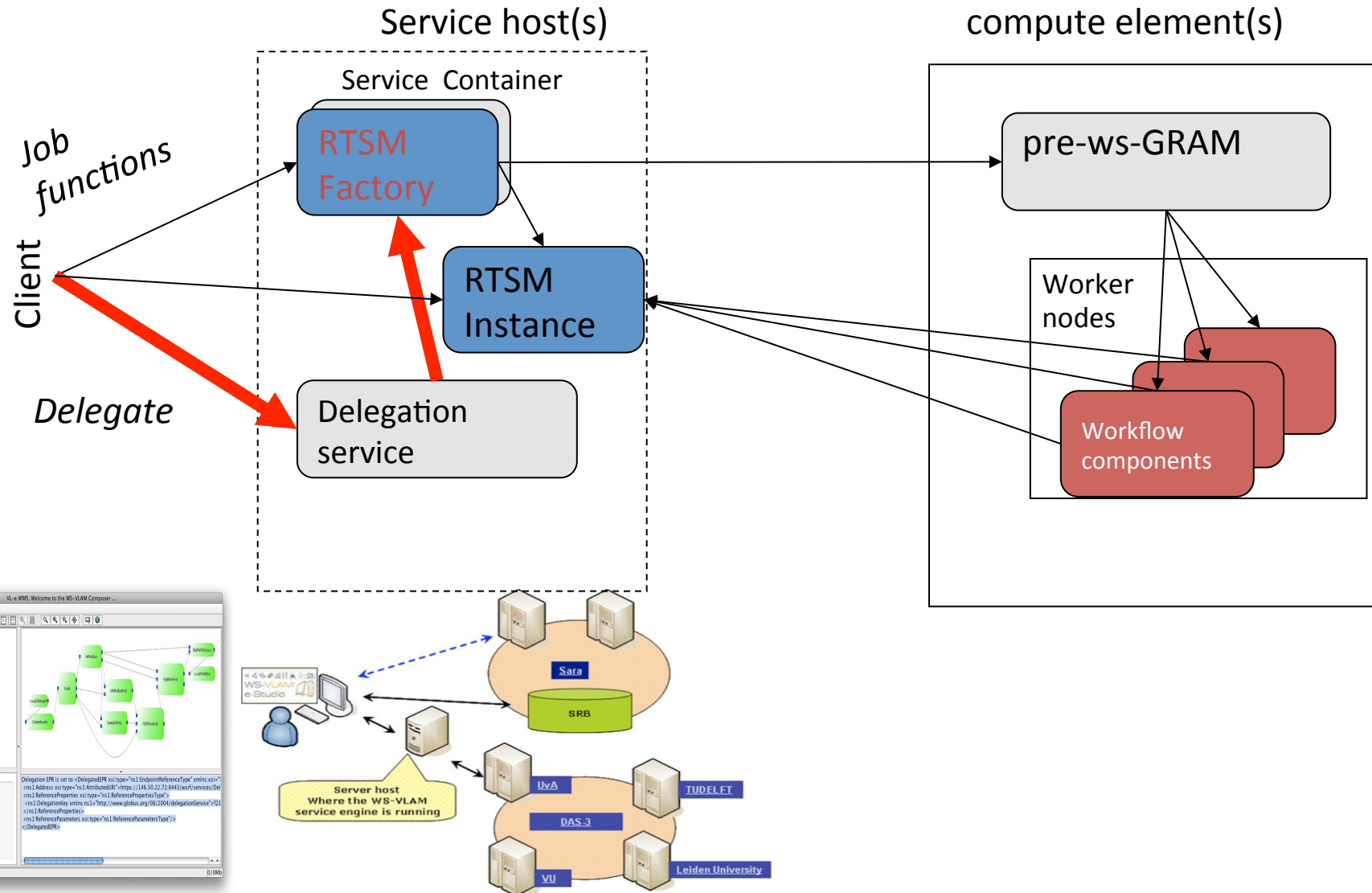*Figure 1.2.* An example abstract workflow.
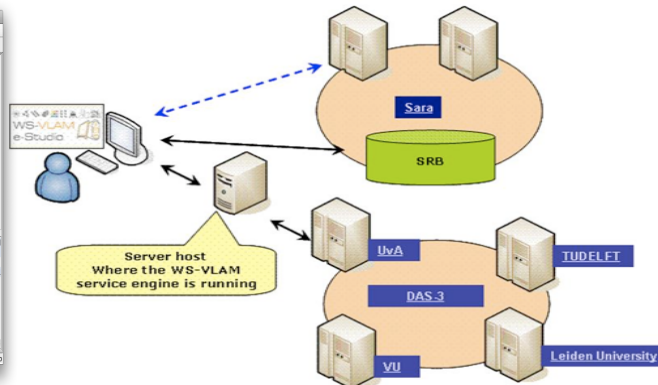


*Figure 1.3.* An example reduced, concrete workflow.

41

# Distributed enabled workflow engines



Workflow composition

Distributed Workflow
Management system

Workflow execution

**Data Management Stack**

Web Service Interface

Application Templates
Web service

OGSA DAI

Grid Middleware:
Data management

Network & storage Resources

**Process Management Stack**

Web Service Interface

Application
web service

Grid
Services

Workflow
Engine
Web service

Grid Middleware:
Process & resource management

Network & Computing Resources
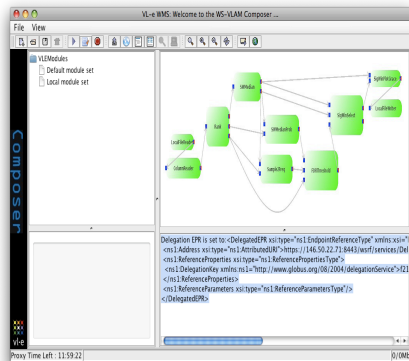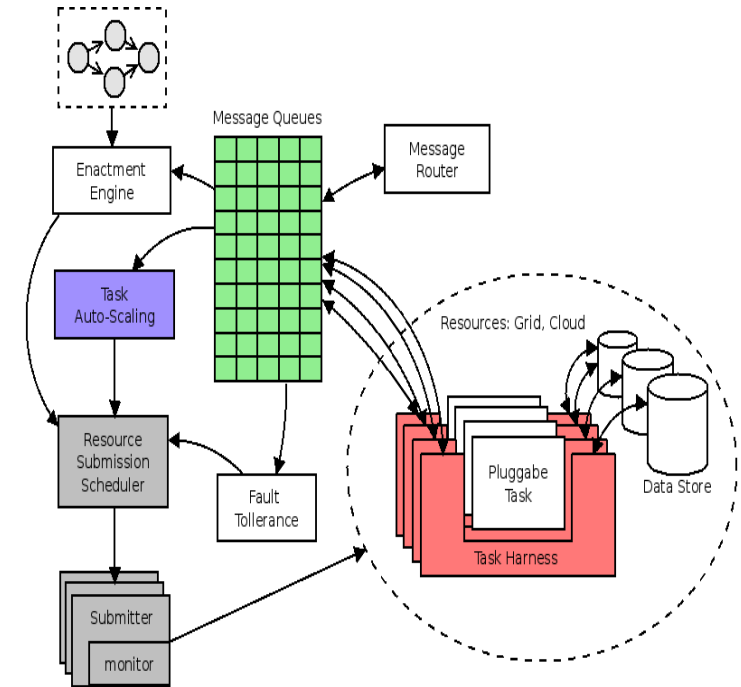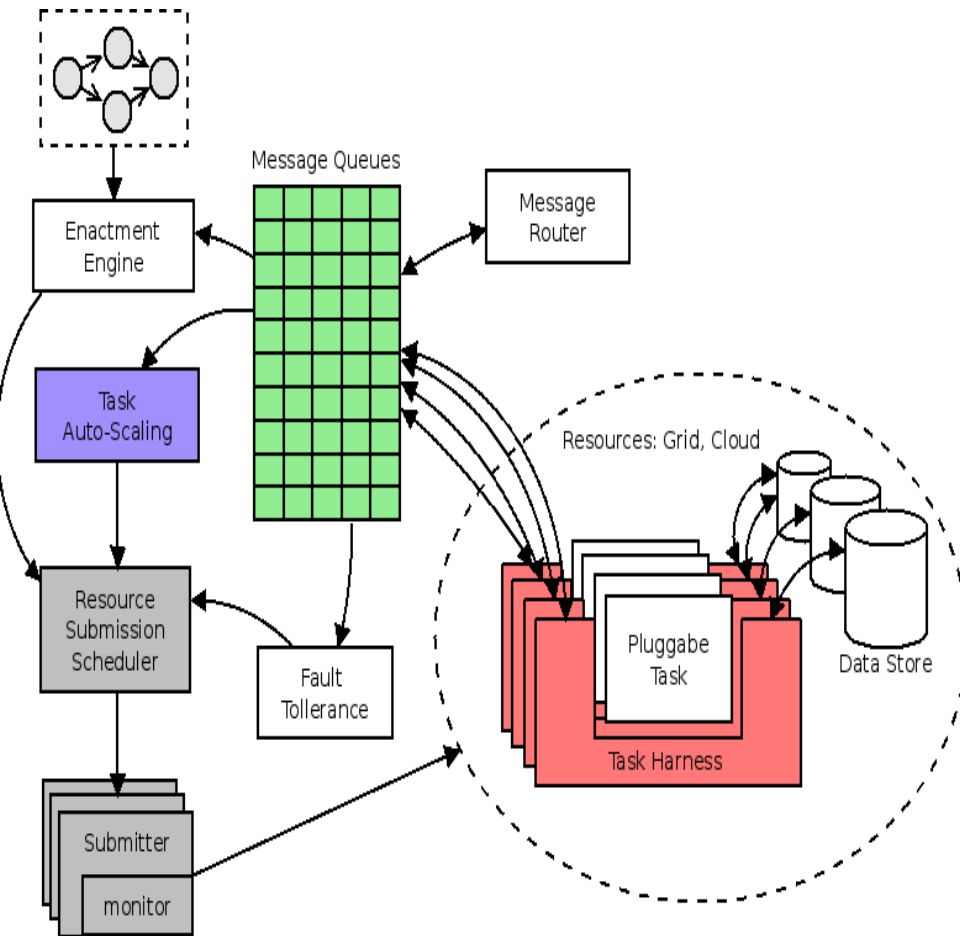
# WS-VLAM

# WS-VLAM

# DataFluo engine



- Automatic scaling of workflow components based
  - Resource load
  - Application load
  - provenance data

- Scaling across various infrastructures
  - desktop
  - Grids
  - Clouds

Reginald Cushing, Spiros Koulouzis,  Adam S. Z. Belloum, Marian Bubak, **Prediction-based Auto-scaling of Scientific Workflows**, Proceedings of the 9th International Workshop on Middleware for Grids, Clouds and e-Science, ACM/IFIP/USENIX December 12th, 2011, Lisbon, Portugal

# Usage of Web Services in e-science

- WS offer interoperability and flexibility in a large scale distributed environment.

- WS can be **combined** in a **workflow** so that more complex operations may be achieved,

- but any workflow implementation is potentially **faced** with a *data transport problem*
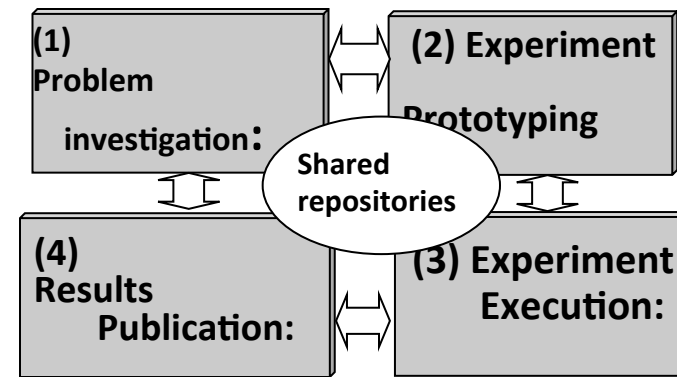
# Resource management

- Within a single workflow services are **competing** for resources.

- Scaling one service without any regard to the whole workflow may starve parts of the workflow and hamper progress

- It would be ideal to have a mechanism to **greedily** consume resources if **no one** is using them but **donate back** resources once they are requested.

Fuzzy controller tries to do just that.

# Outline

- Introduction
- Lifecycle of an e-science workflow
- Different approach to workflow scheduling
  - Workflow Process Modeling & Management In Grid/Cloud
  - Workflow and Web services (intrusive/non-intrusive)
- provenance

# Provenance/ reproducibility

- "A complete provenance record for a data object allows the possibility to reproduce the result and reproducibility is a critical component of the scientific method"

- Provenance: The recording of metadata and provenance information during the various stages of the workflow lifecycle
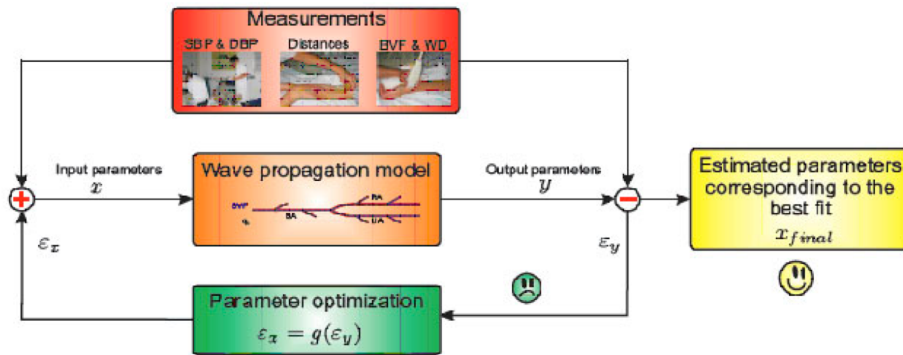
Workflows and e-Science: An overview of workflow system features and capabilities Ewa Deelmana, Dennis Gannonb, Matthew Shields c, Ian Taylor, Future Generation Computer Systems 25 (2009)

# History-tracing XML (FH Aachen)

- provides data/process provenance following an approach that
  - maps the workflow graph to a layered structure of an XML document.
  - This allows an intuitive and easy processable representation of the workflow execution path,
  - which can be, eventually, electronically signed.

```
<patternMatch>
  <events>
    <PortResolved> provenance
data</PortResolved>
    <ConDone>provenance data
            </ConDone>
    ...
  </events>
  <fileReader2>
    <events> ... </events>
    <sign-fileReader2> ...
        </signfileReader2>
  </fileReader2>
  <sffToFasta>
    Reference
  </sffToFasta>
  <sign-patternMatch> ...
        </sign-patternMatch>
</patternMatch>
```
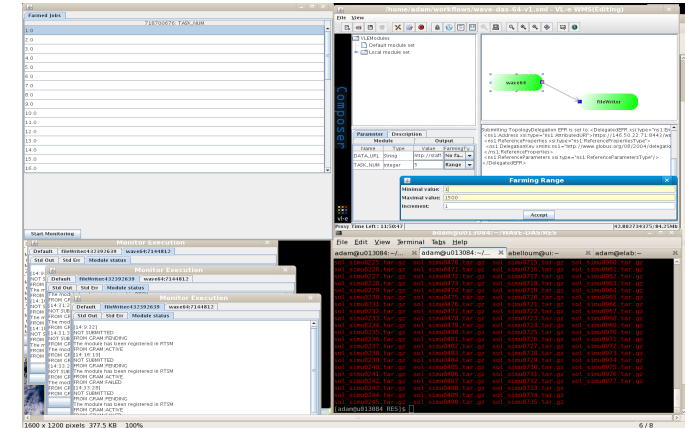
M. Gerards, Adam S. Z. Belloum, F. Berritz, V. Snder, S. Skorupa, **A History-tracing XML-base Proveannce Framework for workflows**, WORKS 2010, New Orleans, USA, November 2010

# wave propagation model applications



*[Biomedical engineering Cardiovascular biomechanics group TUE])*

*wave propagation model of blood flow in large vessels using an approximate velocity profile function:*

a biomedical study for which **3000 runs** were required to perform a global sensitivity analysis of a blood pressure wave propagation in arteries
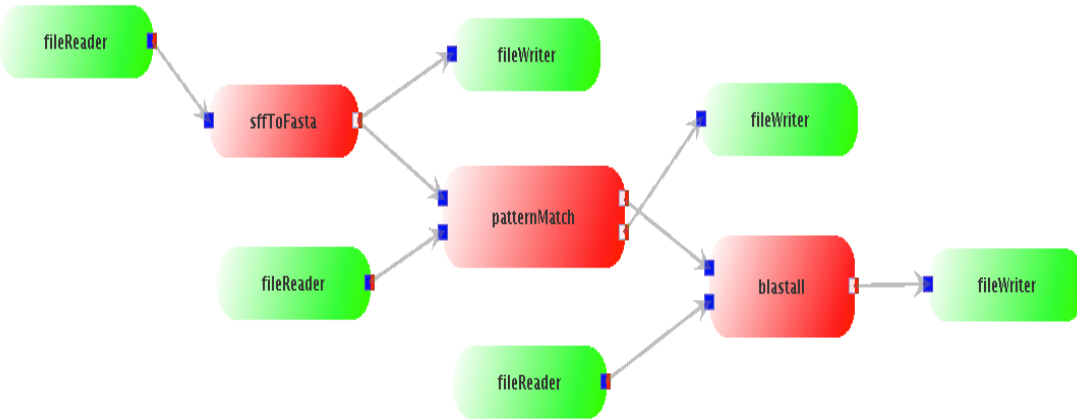


User Interface to compose workflow (top right), monitor the execution of the farmed workflows (top left), and monitor each run separately (bottom left) data
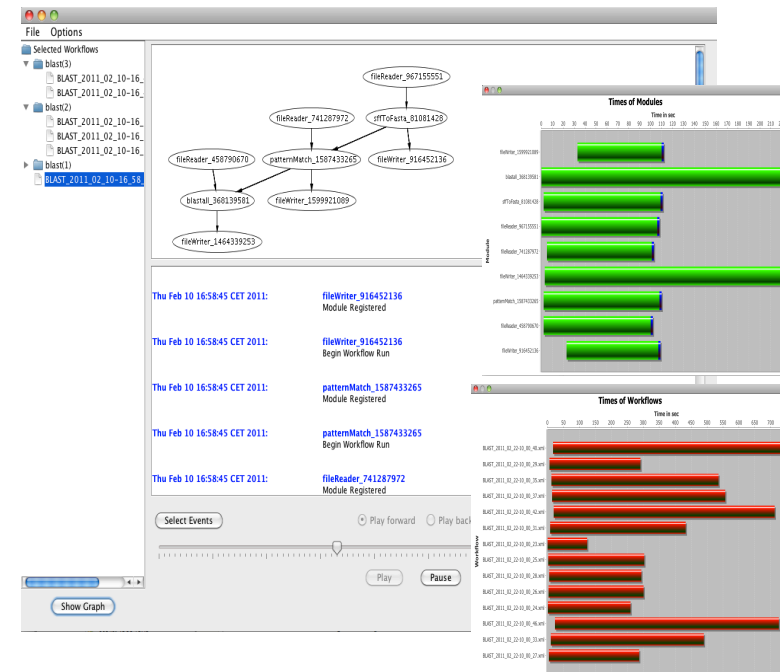


Query interface for the provenance data collected from 3000 simulations of the "*wave propagation model of blood flow in large vessels using an approximate velocity profile function"*

# Blast Application



*[Department of Clinical Epidemiology, Biostatistics and Bioinformatics (KEBB), AMC ]*



The aim of the application is the **alignment of DNA sequence** data with a given reference database. A workflow approach is currently followed to run this application on distributed computing resources.

For Each workflow run

The provenance data is collected an stored following the XML-tracing system

User interface allows to reproduce events that occurred at runtime (replay mode)

User Interface can be  customized (User can select the events to track)

User Interface show resource usage

on-going work UvA-AMC-fh-aachen

# More References

1. A.S.Z. Belloum, V. Korkhov, S Koulouzis, M. A Inda, and M. Bubak *Collaborative e-Science experiments: from scientific workflow to knowledge sharing* JULY/AUGUST, IEEE Internet Computing, 2011

2. Ilkay Altintas, Manish Kumar Anand, Daniel Crawl, Shawn Bowers, Adam Belloum, Paolo Missier, Bertram Ludascher, Carole A. Goble, Peter M.A. Sloot, Understanding Collaborative Studies Through Interoperable Workflow Provenance, IPAW2010, Troy, NY, USA

3. A. Belloum, Z. Zhao, and M. Bubak Workflow systems and applications , Future Generation Comp. Syst. 25 (5): 525-527 (2009)

4. Z. Zhao, A.S.Z. Belloum, et al., Distributed execution of aggregated multi domain workflows using an agent framework The 1st IEEE International Workshop on Scientific Workflows, Salt Lake City, U.SA, 2007

5. Zhiming Zhao, Adam Belloum, Cees De Laat, Pieter Adriaans, Bob Hertzberger Using Jade agent framework to prototype an e-Science workflow bus Authors Cluster Computing and the Grid, 2007. CCGRID 2007

# Summary

- Workflow research especially in the grid environments are rapidly growing research subject

- VOs in Grid can benefits from the experience of workflows in the business community

- Scientific Workflow in Grid Environment have their own characteristics that need to be dealt with new approach

- Scientific Workflow research is highly related with various other research topics: resource management, fault tolerance, application performance, ontology.