# Scientific and Grid Workflow Management

Cesare Pautasso
University of Lugano
http://www.pautasso.info

# Abstract

Grid workflow management systems coordinate multiple job submissions over heterogeneous Grid resources.

They feature visual programming environments to give scientist a high-level view over distributed computations composed of Grid services.

This brief introduction to the field of scientific and Grid workflows includes a survey of selected workflow management tools and outlines current research trends.

**Swiss Grid School SGS'09 @ GPC'09, Geneva, Switzerland**

Università della Svizzera italiana

Facoltà di scienze informatiche

# Cesare Pautasso

Ph.D. at ETH Zürich (2004)

Post-Doc at ETH Zürich in the Systems (IKS) Group

- Software: JOpera: Process Support for more than Web services

http://www.jopera.org/

Researcher at IBM Zurich Research Lab (2007)

Assistant Professor at the new Faculty of Informatics, University of Lugano (USI), Switzerland (since September 2007)

- USI Representative in the SwiNG Assembly (since 2007)
- Grid Workflow Working Group Lead (since 2007)

More Information: http://www.pautasso.info/

Follow me on: http://twitter.com/pautasso/

# Acknowledgements

Some material contained in this tutorial was adapted from slides originally published by:

Gustavo Alonso

Win Bausch

Ewa Deelman

Ian Foster

Yolanda Gil

Carole Goble

Roy Grønmo

Thomas Heinis

Rajesh Kalyanam

Francesco Lelli

Omer F. Rana

Heiko Schuldt

Frank Terpstra

# Outline

Why Workflow Management on the Grid?

Discussion: Scientific vs. Grid vs. Business Workflows

- Some Application Examples

Workflow Modeling Languages and Tools Overview

- Grid Workflow Language Patterns

Running Workflows on the Grid

- JOpera: Scientific Workflow for Eclipse
- Workflows and Provenance

# Why Workflow Management on the Grid?

# Kinds of Grid Computation

## One Job Submission

## Parameter Sweep

Scientific and Grid Workflow (Cesare Pautasso)

Scientific and Grid Workflow (Cesare Pautasso)

8

**Copy & Paste**

between different Websites

**Programming**

Java, C++, C#, Fortran...
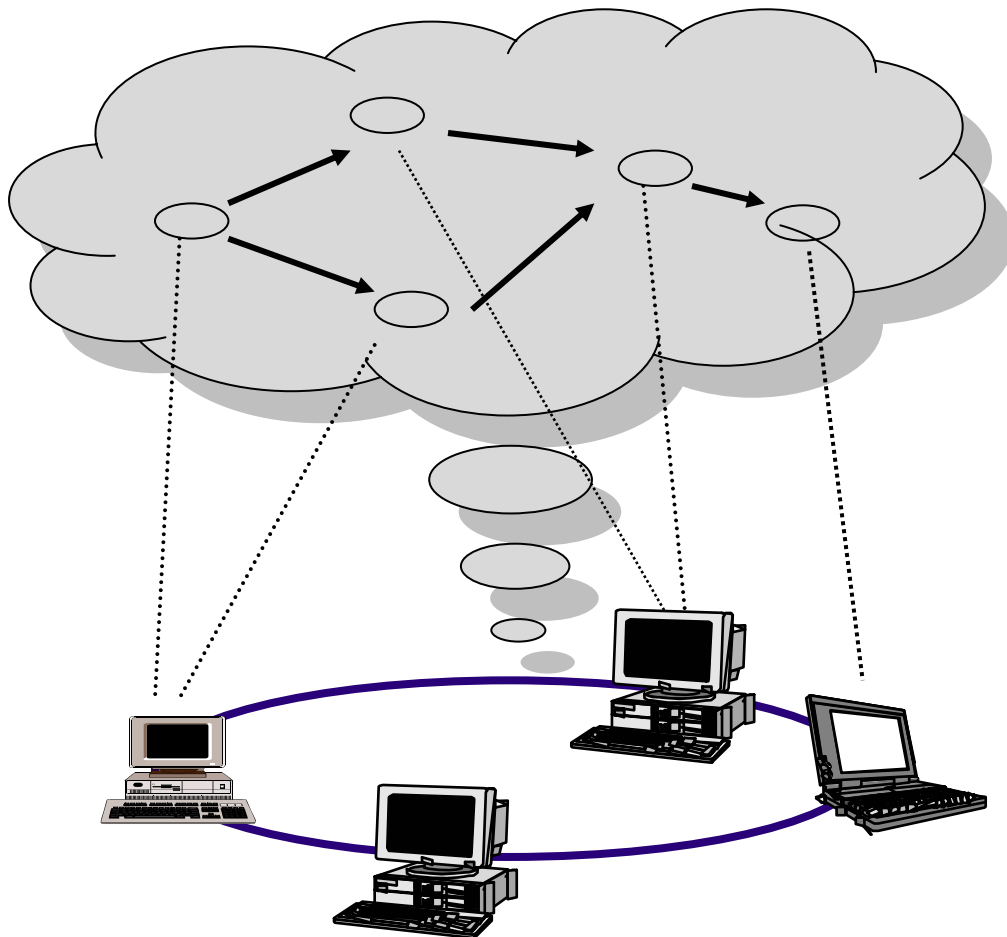
**(Shell) Scripts**

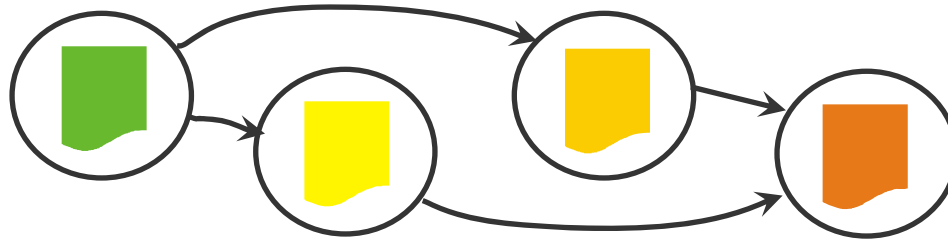Tcsh, Bash, Makefiles, Python, Perl...

**Workflows**

Graphical, Drag & Drop and Connect Environments
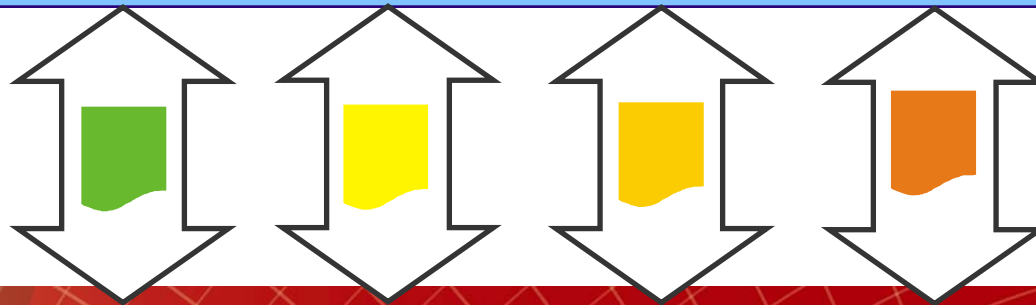
# Vision for Scientific and Grid Workflows

Make it easy to build Grid applications composed of **multiple jobs**



" Provide the scientist with a platform that takes care of all data handling and record keeping chores so that the user can concentrate on the *science* and not *computer science* "

# **Workflow**

# GRID

# Some (Scientific) Workflow Management Systems

Askalon

Bigbross Bossa

BioPipe

BPMN

Breeze

Carnot

Con:cern

DAGMan

DiscoveryNet

Dralasoft

GEL

GridAnt

Grid Job Handler

GWFE

GWES

ICENI

Inforsense

JIGSA

JOpera

Kepler

Karajan

Oakgrove's reactor

OSIRIS

OSWorkflow

OpenWFE

Pegasus

Pipeline Pilot

P-GRADE

PowerFolder

Ptolemy II

Savvion

Seebeyond

SCIRun

ScyFLOW

SDSC Matrix

SHOP2

Taverna

Teuta (UML)

Triana

Trident

Twister

Ultimus

Versata

Viztrails

wftk

XFlow

YAWL

Wildfire

WFEE

WS-BPEL

ZBuilder

# Scientific vs. Grid vs. Business Workflows

# The Origins: Business Process Management

# <u>who</u> has to do <u>what</u>, <u>when</u>

# The Origins: Business Process Management

- A business process describes key procedures within an organization. They involve:
  - multiple steps
  - numerous people
  - large amounts of resources

- In large business organizations there are many factors that increase the complexity of the business processes:
  - processes are not well documented
  - conformance to rules not guaranteed
  - people lack information about context
  - company lacks monitoring tools
  - steps, people and resources are not properly coordinated

- Workflow Management Systems try to address these problems by automating the coordination aspects of a business process: who has to do what, when, and with which software tools.

# Business Workflows

" The **automation** of a business process where **documents**, **information** to be processed or **tasks** to be carried out are passed from one participant to another following a set of **procedural rules** "
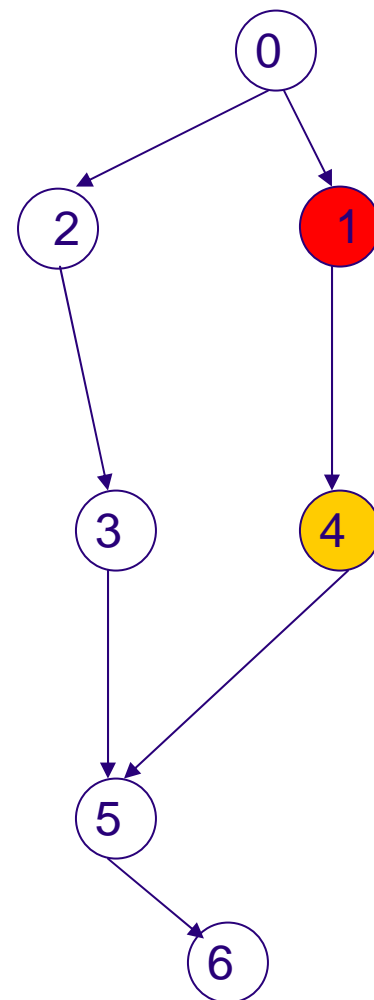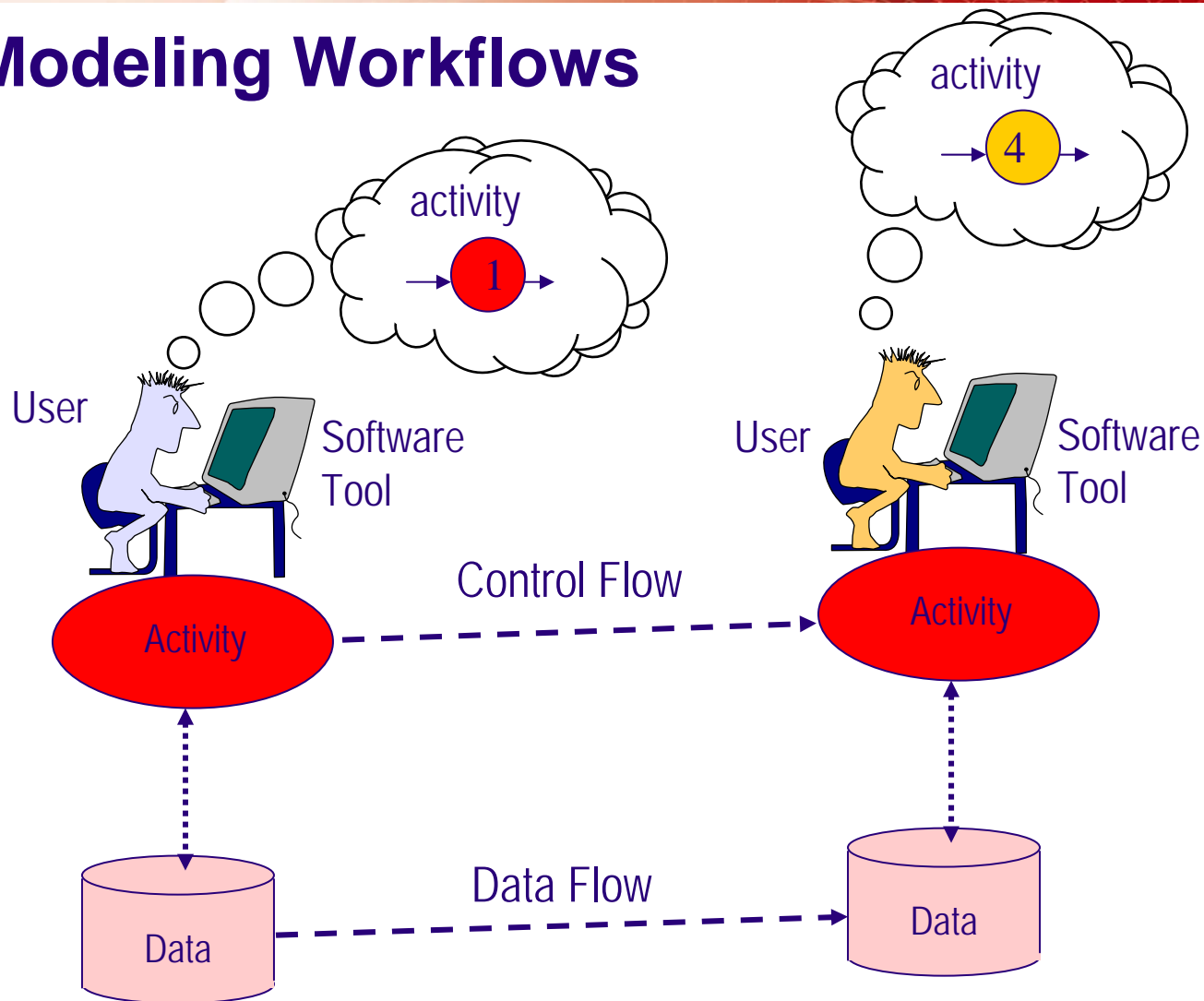
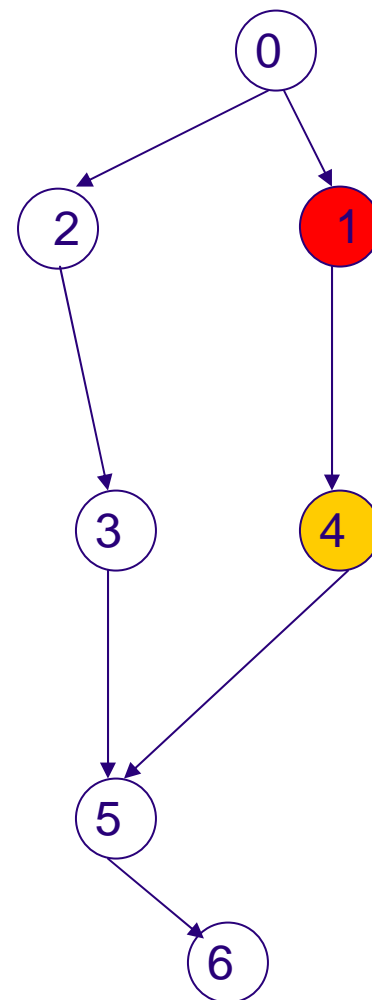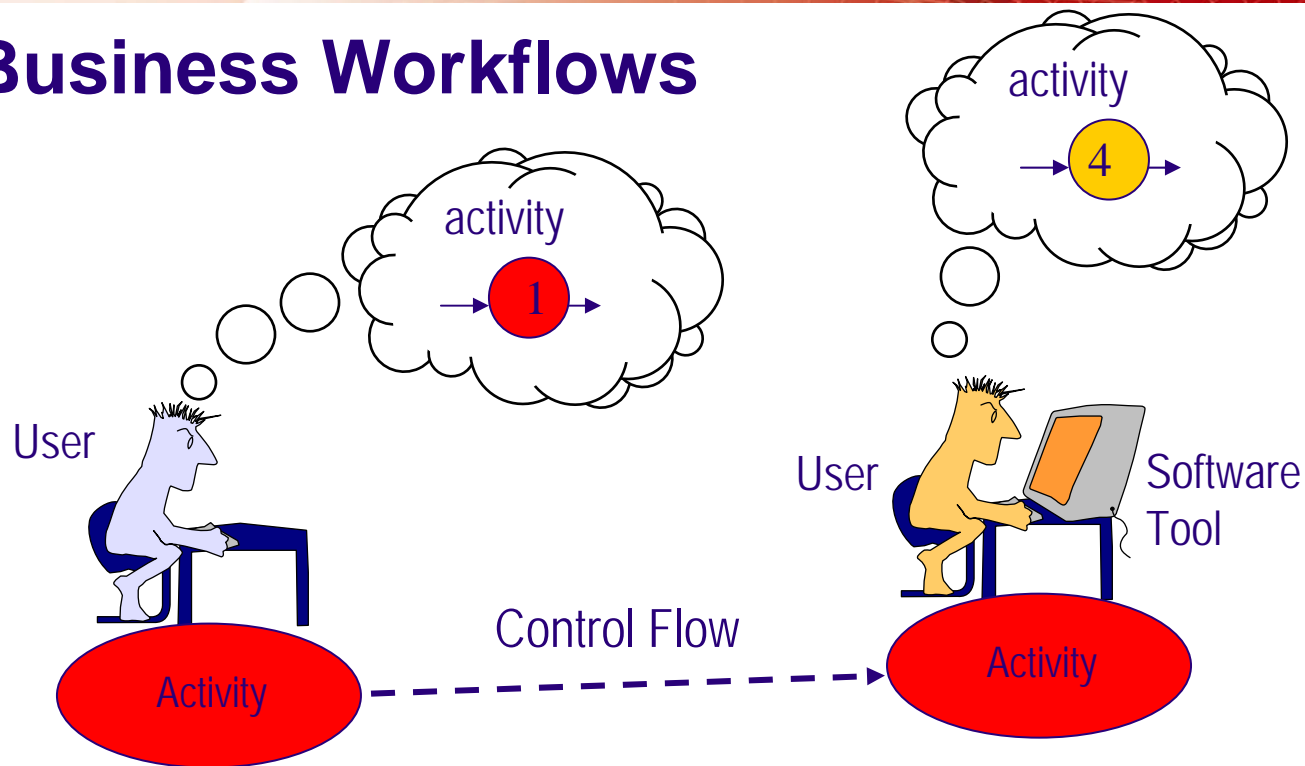Worfklow Management Coalition (WfMC, 1993)

# Scientific Workflows

**"** are **networks** of analytical **steps** that may involve, e.g., database access and querying, data analysis and mining, and many other steps including computationally intensive **jobs** submitted to high performance clusters and **Grids** **"**

Bertram Ludäscher

# Modeling Workflows

# Business Workflows

# Scientific Workflows

# Similarities:
# Are scientists doing e-Business?

## Capturing knowledge/best practices

Capture business processes within a company

*Capture scientific experiments*

## Executable Models for Repeated Execution

Run a well defined procedure many times

*Ensure that an experiment can be reproduced*

## Incorporate human decision in the process

Can we always do straight-through processing?

*Hard to achieve full automation*

# Differences:
# Do scientists need business transactions?

## Rate of change

Changing business procedures requires management approval

*Exploratory scientific processes require high flexibility*

## Which kind of data?

Travel reservations, Loan applications

*Large protein sequence databases, Astronomy image catalogs*

## What is the ultimate goal?

Making profit

*Making science*

# Scientific vs. Grid Workflows

Scientific workflows emphasize the design of virtual experiments:

- Data flow models
- Reusable "scientific computing" component library
- Interactive debugging, monitoring and steering
- Data provenance and lineage tracking for reproducibility
- Model versioning for exploratory customization

Grid workflows focus on the large-scale execution of scientific workflows:

- Mapping and adaptation to a dynamic run-time environment
- Provide access to shared workflows as a Grid service
- Parameterized Execution
- Centralized vs. Distributed Execution Architectures
- Fault Tolerance
- Optimization

# Scientific Workflows on the Grid

- How can Scientific WF benefit from the Grid?
  - 1. Leverage underlying Grid middleware:
    - Resource Management
    - Job Scheduling
    - Large Data Transfers (GridFTP) between Activities
  - 2. Improved QoS based on the workflow model
    - Grid resource reservation
    - Data replication
    - Data placement
    - Fault Tolerance

JOpera
Powered by

http://www.jopera.org/

SwiNG
SWISS NATIONAL
GRID ASSOCIATION

# Example

PURDUE
UNIVERSITY
Information Technology at Purdue

ITaP

TeraGrid

# A Web Service-Enabled Workflow System for Climate Modeling Data Processing in TeraGrid

*Rajesh Kalyanam*

*Lan Zhao*
*Taezoon Park*
*Sebastien Goasguen*

# Architecture

# Portal

# Workflow Model

# Workflow Execution

# Workflow Results
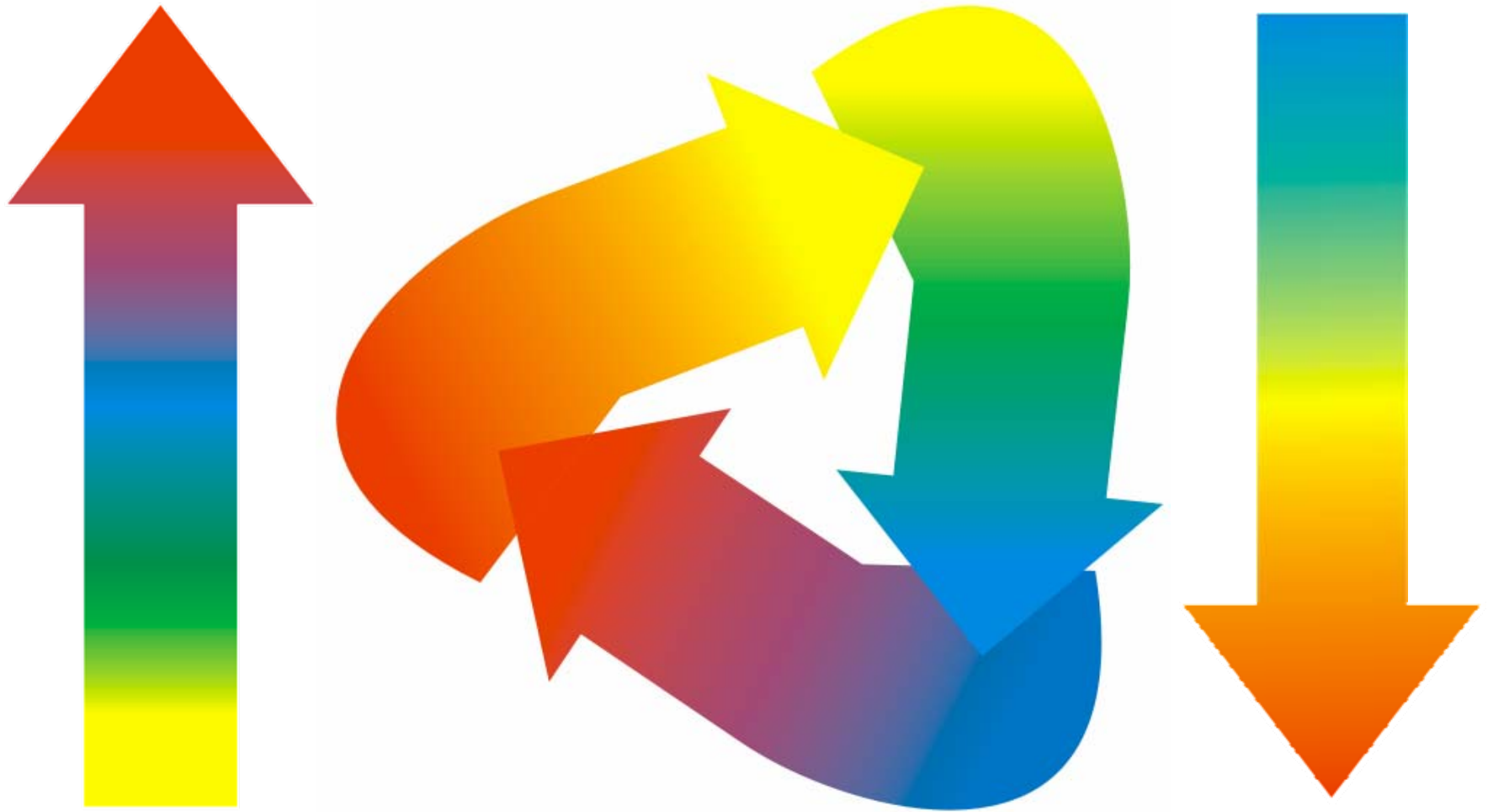


From Rajesh Kalyanam

# Workflow Lifecycle

# Workflow Lifecycle



From Gustavo Alonso

# Workflow Modeling Methodologies

# Bottom up Composition

4. Share and Publish it as Web Service

3. Run, Test, and Debug the execution **within the same modeling environment**

2. Build a workflow using a drag, drop and connect **modeling** environment

1. Select components from a **library**
    a. Lookup services in a public registry
    b. Import from external Web service (WSDL)
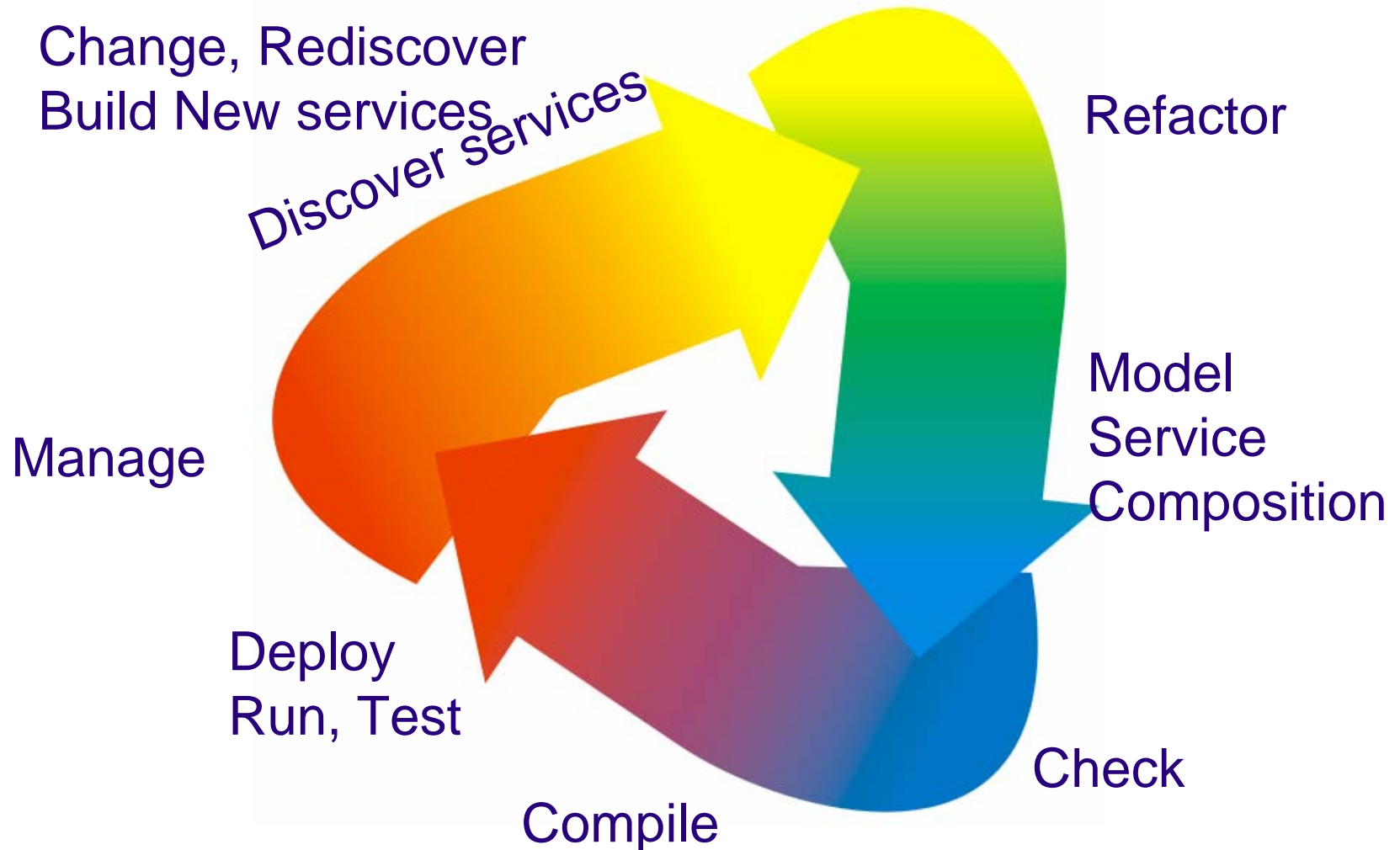    c. Search the standard library

# Top down Decomposition

1. Define a **goal** and Draw a *skeleton of the workflow* that satisfies it

2. Refine it and **Bind** services into it:
   - Search for existing matching services
   - Build missing services (if necessary)
   - Add required data transformations

3. Run, Test, and Debug the execution **within the same modeling environment**
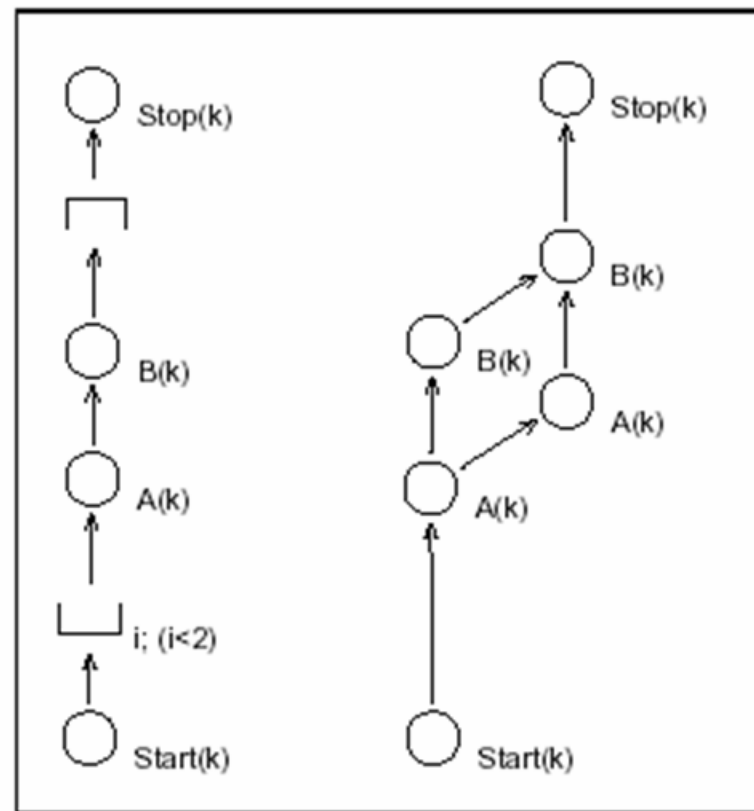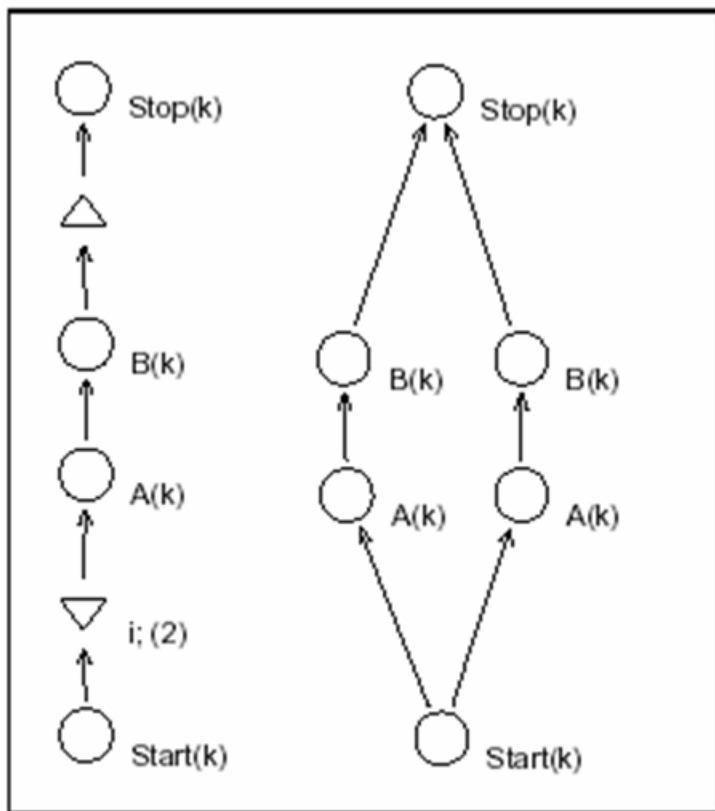
4. Share and Publish it as Web Service

# **Iterative Composition**

Change, Rediscover
Build New services

Discover services

Refactor

Model
Service
Composition

Manage

Deploy
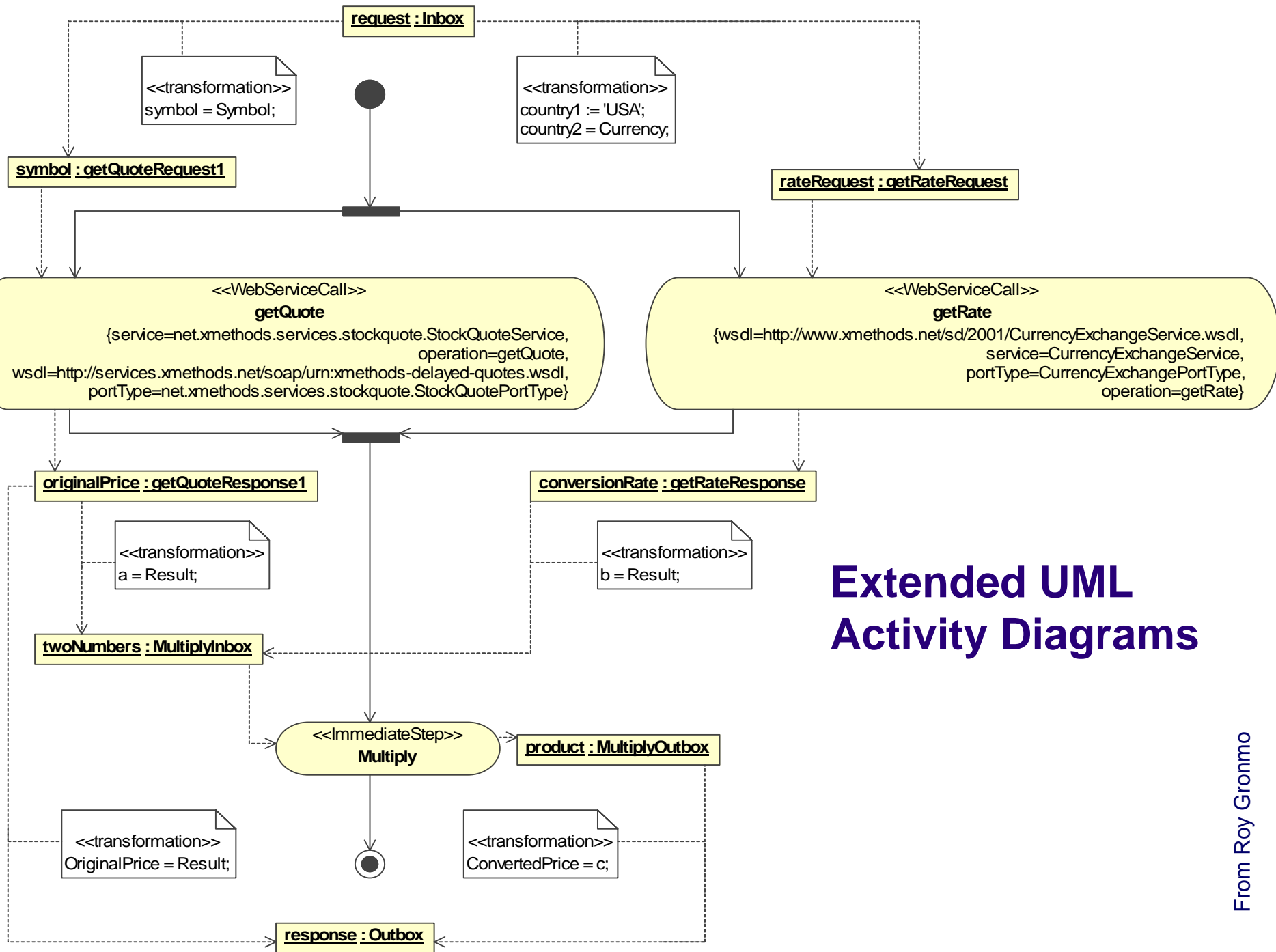Run, Test

Check

Compile

# Workflow Modeling Languages and Tools Overview

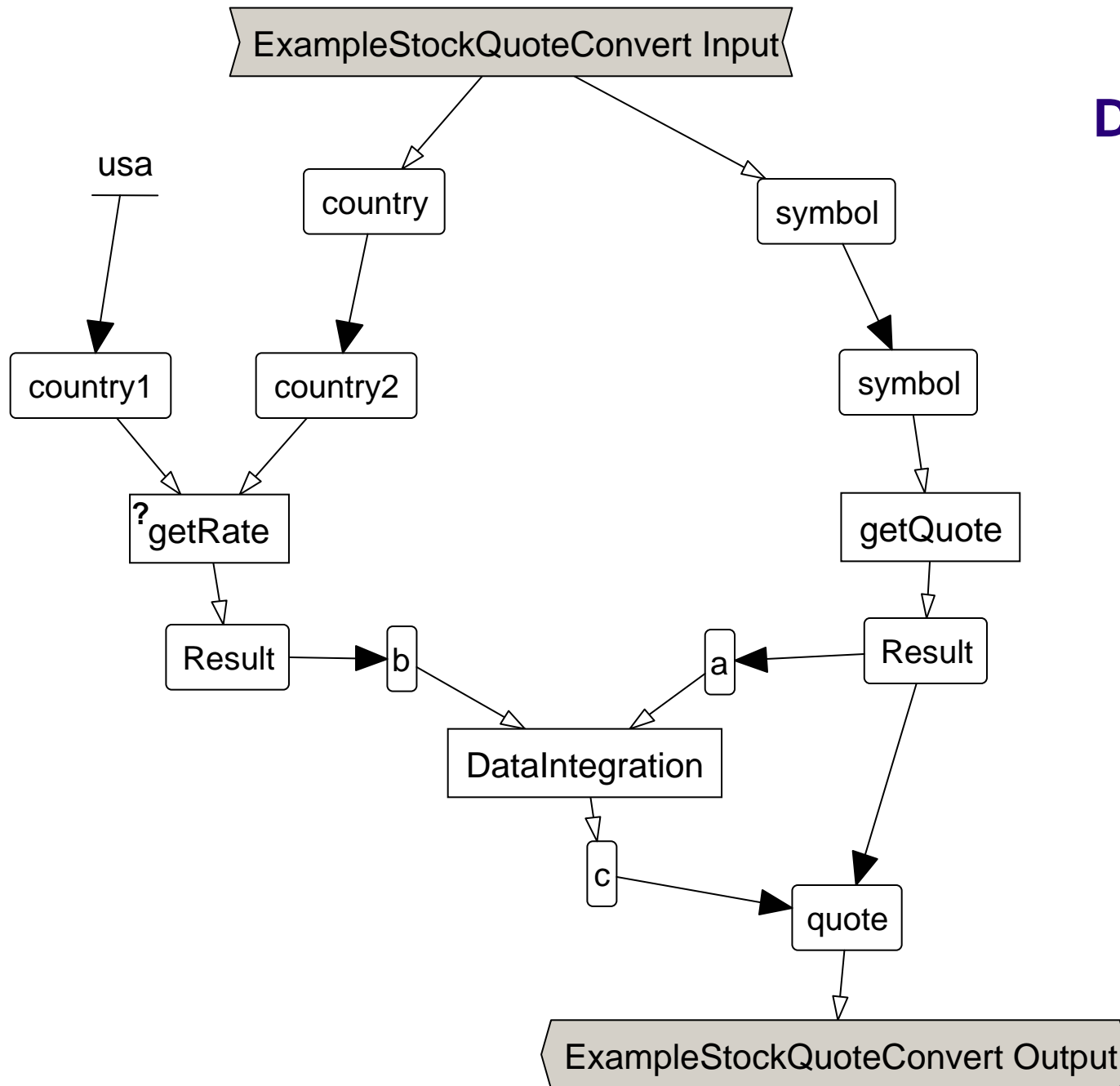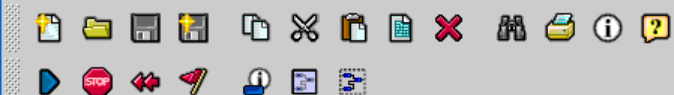# HeNCE - The Ancestor of Grid Workflows?



A. Beguelin, J. J. Dongarra, G. A. Geist, R. Manchek, V. S. Sunderam, Graphical Development Tools for Network-Based Concurrent Supercomputing, in: Proc. of the 1991 ACM/IEEE conference on Supercomputing, Albuquerque, New Mexico, 1991, pp. 435–444.

**request : Inbox**

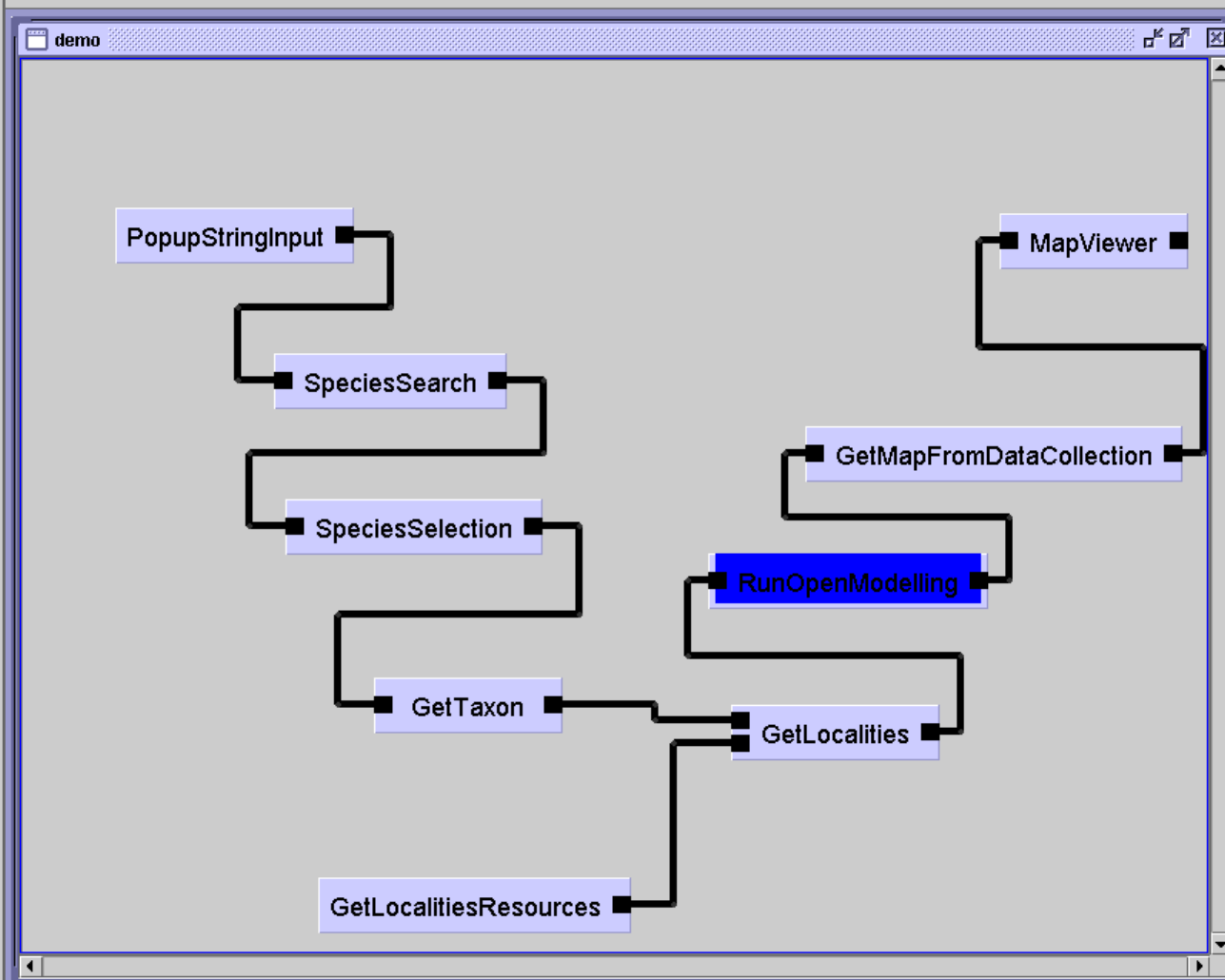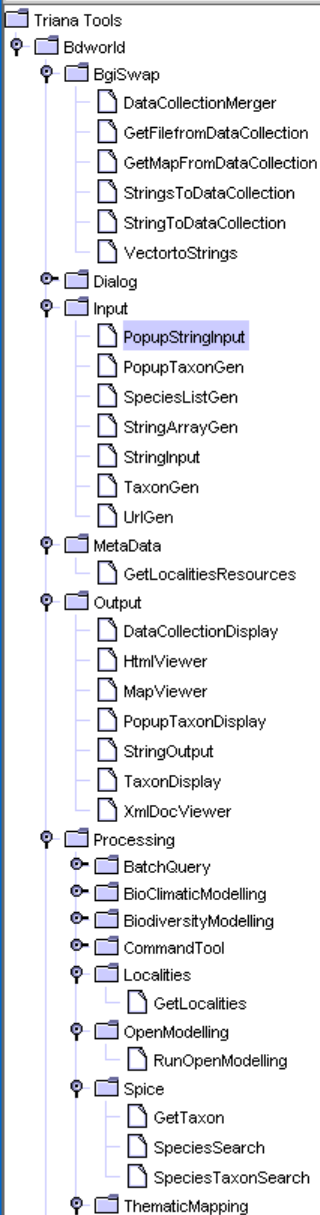<<transformation>>
symbol = Symbol;

<<transformation>>
country1 := 'USA';
country2 = Currency;

**symbol : getQuoteRequest1**

**rateRequest : getRateRequest**

<<WebServiceCall>>
**getQuote**
{service=net.xmethods.services.stockquote.StockQuoteService,
operation=getQuote,
wsdl=http://services.xmethods.net/soap/urn:xmethods-delayed-quotes.wsdl,
portType=net.xmethods.services.stockquote.StockQuotePortType}

<<WebServiceCall>>
**getRate**
{wsdl=http://www.xmethods.net/sd/2001/CurrencyExchangeService.wsdl,
service=CurrencyExchangeService,
portType=CurrencyExchangePortType,
operation=getRate}

**originalPrice : getQuoteResponse1**

**conversionRate : getRateResponse**

<<transformation>>
a = Result;

<<transformation>>
b = Result;

# Extended UML
# Activity Diagrams

**twoNumbers : MultiplyInbox**

<<ImmediateStep>>
**Multiply**

**product : MultiplyOutbox**

<<transformation>>
OriginalPrice = Result;

<<transformation>>
ConvertedPrice = c;

**response : Outbox**

From Roy Gronmo

**JOpera
Data Flow
Graph**

ExampleStockQuoteConvert Input

usa

country

symbol

country1

country2

symbol

**?**getRate

getQuote

Result

b

a

Result

DataIntegration

c

quote

ExampleStockQuoteConvert Output

**VizTrails**

Scientific and Grid Workflow (Cesare Pautasso)

# Grid Workflow Language Patterns

## Workflow Pattern          Variants

1. **Simple Parallelism**
   - Implicit
   - Explicit

2. **Data Parallelism**
   - Static
   - Dynamic

3. **Pipelining**
   - Best Effort
   - Blocking
   - Buffered
     - Hybrid
   - Superscalar
     - Synchronized
     - Out of Order
   - Streaming

# Modeling Simple Parallelism

Data Flow, Graph Based



SCIRun

Kepler

Triana

# Modeling Simple Parallelism

## Control Flow, Graph Based



JOpera GEL

UML

# Modeling Simple Parallelism

## Control Flow, Block Based



BPMN

WS-BPEL

# Modeling Data Parallelism

## Data Flow, Graph Rewriting



### Static or Dynamic

Triana

Taverna
JOpera

# Modeling Data Parallelism

Control Flow, Block Based, Dynamic



WS-BPEL

AGWL

Karajan

GEL

# Modeling Pipelined Execution

# Pipelining Semantics

# Best Effort Pipelined Execution



Drop data elements on pipeline collisions

Advantages:

- Simplified implementation
- Some applications may tolerate data loss

Problem:

- Downsampling is non deterministic

# Blocking Pipelined Execution



Tasks are blocked if successors are busy

Advantages:

- Avoid data loss in the pipeline

Problem:

- Pipeline speed limited by slowest task
- Data may be lost before it enters the pipeline

# Buffered Pipelined Execution



Tasks are decoupled by buffers

Advantages:

- Collisions are prevented
- Best applied to tasks having variable speed

Problem:

- Buffer capacity is limited
  (Blocking still needed – Hybrid semantics)

# Streaming Pipelined Execution



Tasks exchange data while running

Advantages:

- Suitable for a distributed (P2P) engine

Problems:

- Shifts complexity from the workflow engine to the tasks
- Tasks exchange data while running
- Workflow/Task interface more complex

# Running Workflows on the Grid

**Basic Architecture**

**Workflow Model**

Act 1 → Act 2
Act 1 → Act 3
Act 2 → Act 4
Act 2 → Act 5
Act 2 → Act 6
Act 3 → Act 5
Act 3 → Act 6
Act 4 → Act 7
Act 5 → Act 7

**Workflow Management System**

Adapters

Grid
Schedulers

Grid
Resources

# Standard APIs

From WFMC,
Workflow Reference Model, 1998

# Wrappers and Grid Applications

# Wrappers and Legacy Applications

- The workflow engine is also in charge of connecting the different scientific applications.

- These applications do not have to talk directly to each other, they do it through the workflow engine.

- Most engines target a service oriented applications for which they provide very good connectivity through standardized protocols. Otherwise, the interface adapters must be developed on a case by case basis (as a last resort manual integration may be required!)

- For legacy application, a wrapper must be built so that the workflow engine can communicate with the application. The wrapper can be a simple relay of commands and data, or a complete translation program implementing functionality not present in the legacy application.

- For most Grid applications, the interaction takes place through a Grid scheduler, which is responsible for managing the distributed execution of the applications.

# Run-time Abstraction Levels

# Run-time Abstraction Levels

- A design-time workflow model needs to be mapped across different abstraction levels in order to be executed at run time.

- User request the execution of a new workflow instance.

- The abstract workflow is mapped to an executable instance by:
  - Finding suitable service implementations and binding them to the tasks
  - Rewriting the workflow graph based on a set of refinement rules
  - Planning required data staging, registration, placement, replication and transfer operations

- Each task of the resulting executable workflow is then submitted to a Grid resource manager so that it can be scheduled on suitable resources

- The mapping can be done:
  - when the workflow is started at instantiation time (statically)
  - incrementally as the workflow runs (adaptive execution with dynamic late binding)

# Example: Binding with WS-BPEL

# Workflow Binding Lifecycle

- Library Registration time *(classification)*
- Modeling time (***static early binding***)
- Compilation time *(blacklisting)*
- Deployment time *(customization)*
- Startup time *(testing)*
- Task Execution time (***dynamic late binding***)
- Failed invocation time *(rebind on retry)*

http://www.jopera.org/

# JOpera
# Scientific Workflow for Eclipse

- **High Level Workflow Language**

  - Data and Control Aspects (Visual Representation)
  - Recursion, Iteration, Parallelism and Pipelining

- **Open and Extensible Component Model**

  - Run existing code without changes
  - Synchronous, Asynchronous, and Streaming interaction
  - Web services support (Axis, WSIF)
  - Secure access to remote file systems and hosts (SSH)
  - Easy to integrate with existing schedulers (e.g. Condor)

- **High Level Workflow Language**
  - Data and Control Aspects (Visual Representation)
  - Recursion, Iteration, Parallelism and Pipelining

- **Open and Extensible Component Model**
  - Run existing code without changes
  - Synchronous, Asynchronous, and Streaming interaction
  - Web services support (Axis, WSIF)
  - Secure access to remote file systems and hosts (SSH)
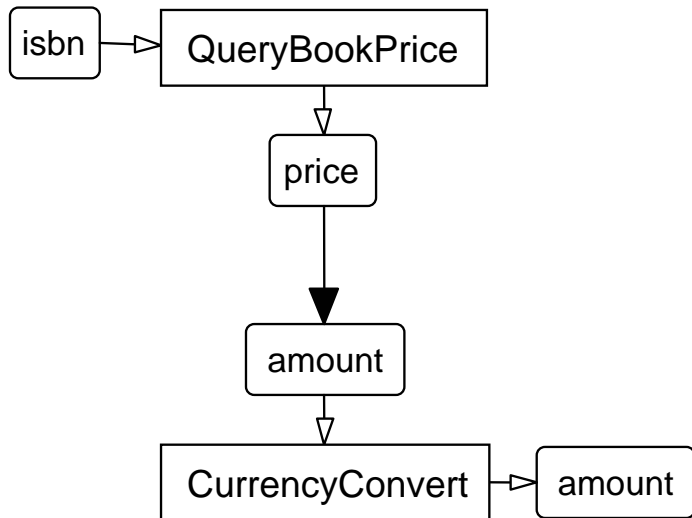  - Easy to integrate with existing schedulers (e.g. Condor)

- **Strong Eclipse Foundation**
  - Platform Independent (Eclipse/Java)
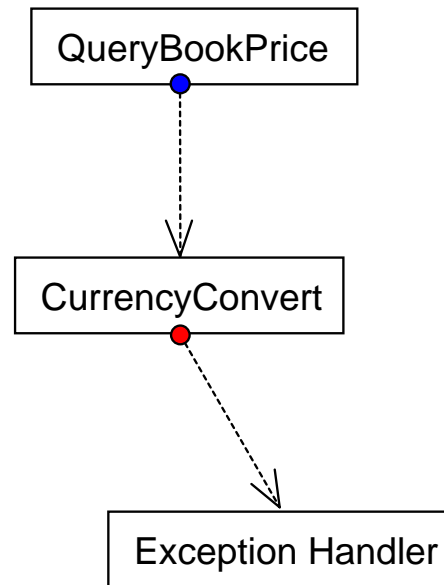  - Flexible, Extensible, Modular and Embeddable

# JOpera Visual Composition Language

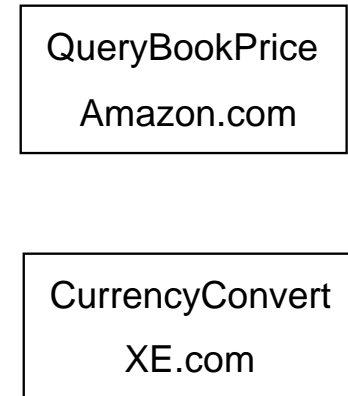## Workflows are modeled using multiple viewpoints:

### 1. Data Flow Graph
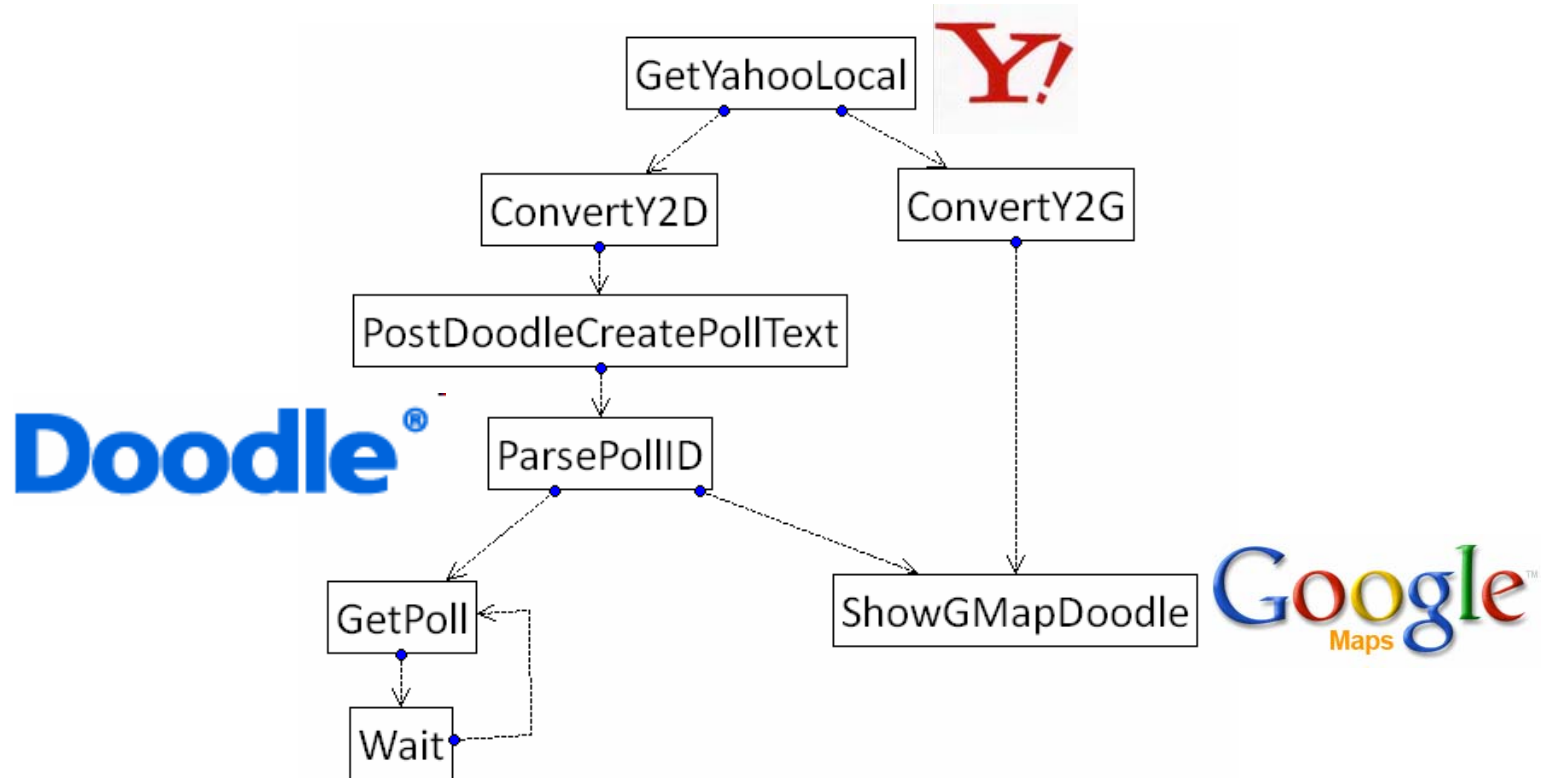


### 2. Control Flow Graph
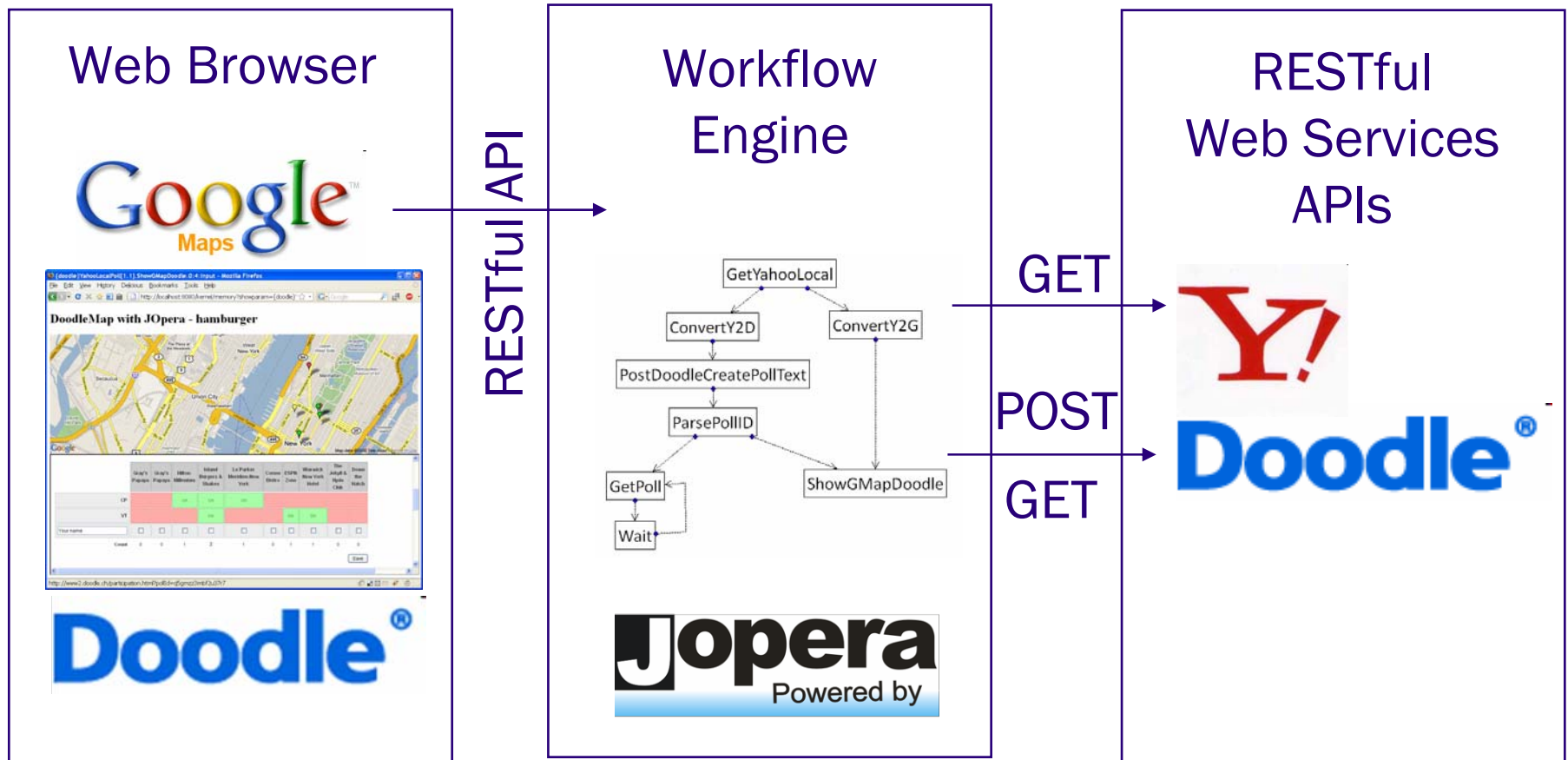


### 3. Service Bindings

# JOpera Example: Doodle Map Mashup

Setup a Doodle with Yahoo! Local search and visualize

the results of the poll on Google Maps

# Doodle Map Mashup Architecture

# Extensible JOpera Component Model

Combine in the same workflow jobs implemented using an open and extensible set of technologies

**JOpera Workflow**

| WSDL | Java | Human | XML | SQL | SSH | Condor |

- Snippets
- Methods

- XSLT
- XPath

# Sharing Workflows as a Service

JOpera processes are automatically published to clients using a variety of access protocols

| Web Clients | WS Clients | Eclipse RCP Clients |
|---|---|---|

REST          WSDL          Java

## JOpera Workflow

| WSDL | Java | Human | XML | SQL | SSH | Condor |
|---|---|---|---|---|---|---|

SWISS NATIONAL GRID ASSOCIATION

**JOpera ARC Integration Demo**

# Workflows and Provenance

# Lineage in Scientific Workflows

Scientists consider the "capture and generation of provenance information as a critical part of the workflow-generated data"

"Sharing workflows is an essential element of education, and acceleration of knowledge dissemination."

Ewa Deelman *et al.*

# Where does this picture come from?



**METADATA**

This photo was taken July 21, 1981, when the Voyager 2 spacecraft was 33.9 million km from the Saturn planet

**METADATA**

**Date:** 16.4.2005
**Dimension:** 640x480
**Colors:** 32bits
**Size:** 1.2MB
**Format:** JPEG

FOR SALE

**Title:** White Arabian Horse

# Would you buy a horse without this?

# Lineage in Spreadsheets

# Lineage in Spreadsheets

# Lineage in Spreadsheets

# Lineage in Spreadsheets

# Lineage in Databases

What is the relationship between these tuples?



SQL

Problem:
**Query Inversion**

# Lineage in Software Development

What's in a Makefile?

```
CC = gcc
CFLAGS = -Wall -g

program: main.o input.o output.o logic.o
        $(CC) $(CFLAGS) main.o input.o output.o logic.o -o program

main.o: main.c input.h output.h logic.h
        $(CC) $(CFLAGS) -c main.c
input.o: input.c input.h
        $(CC) $(CFLAGS) -c input.c
output.o: output.c output.h
        $(CC) $(CFLAGS) -c output.c
logic.o: logic.c logic.h
        $(CC) $(CFLAGS) -c logic.c
```

# Lineage in Software Development

Where does my program come from?

```
CC = gcc
CFLAGS = -Wall -g

program: main.o input.o output.o logic.o
        $(CC) $(CFLAGS) main.o input.o output.o logic.o -o program

main.o: main.c input.h output.h logic.h
        $(CC) $(CFLAGS) -c main.c
input.o: input.c input.h
        $(CC) $(CFLAGS) -c input.c
output.o: output.c output.h
        $(CC) $(CFLAGS) -c output.c
logic.o: logic.c logic.h
        $(CC) $(CFLAGS) -c logic.c
```

# Lineage in Scientific Workflows

**Input Data** Theory

**Output Data** Published Paper

Scientific Workflow

**Input Data** Observation

*An ideal scientific workflow should document all of the steps linking the original observations with the final published results so that the process can be reproduced*

# Data Provenance



Where does this output document come from?

# Change Propagation

What to recompute if this input changes?

# Conclusion

**Reuse**

**Modeling**

**Data Products**

**Workflow and Component Libraries**

**Data, Metadata Catalogs**

**Adapt, Modify**

**Workflow Template**

**Populate with data**

**Data, Metadata, Provenance Information**

**Workflow Instance**

**Execution**

**Execute**

**Executable Workflow**

**Map to available resources**

**Resource, Application Component Descriptions**

**Compute, Storage and Network Resources**

**Distributed**

**Mapping**

From Ewa Deelman

# e-Science as Workflow?



Executed
Executing
Executable
Not yet executable

Provenance Query

What I Did

What I Am Doing

Model

What I Want to Do

...

Execution environment

Schedule

From Ian Foster

# Some References

Workflows for
e-Science

Ian J. Taylor
Ewa Deelman
Dennis B. Gannon
Matthew Shields (Eds)

Scientific Workflows for Grids

Springer

Gil, Y. *et al.*; Examining the Challenges of Scientific Workflows. IEEE Computer, Dec 2007

Taylor, I.J.; Deelman, E.; Gannon, D.B.; Shields, M. (Eds.) Workflows for e-Science: Scientific Workflows for Grids, Springer 2007

Yu, J.; Buyya, R.: A taxonomy of workflow management systems for grid computing, Journal of Grid Computing, 3(3–4):171–200 (2005)

Pautasso, C.; Alonso, G.: Parallel Computing Patterns for Grid Workflows, Proc. Of WORKS@HPDC06, Paris, France, 2006

OGF Workflow Research Group
http://www.isi.edu/~deelman/wf–rg/

Download This Tutorial Material

http://www.pautasso.info/lectures/sgs09workflow.pdf

# Free Download

**Jopera** Powered by

# http://www.jopera.org/