

Decentralised Orchestration of Service-Oriented Workflows

Adam Barker and Rajkumar Buyya
Cloud Computing and Distributed Systems Laboratory
Department of Computer Science and Software Engineering
The University of Melbourne, Australia.

Abstract—Service-oriented workflows in the scientific domain are commonly composed as Directed Acyclic Graphs (DAGs), formed from a collection of vertices and directed edges. When orchestrating service-oriented DAGs, intermediate data are typically routed through a single centralised engine, which results in unnecessary data transfer, increasing the execution time of a workflow and causing the engine to become a performance bottleneck.

This paper introduces an architecture for deploying and executing service-oriented DAG-based workflows across a peer-to-peer proxy network. A workflow is divided into a set of vertices, disseminated to a group of proxies and executed without centralised control over a peer-to-peer proxy network. Through a Web services implementation, we demonstrate across PlanetLab that by reducing intermediate data transfer, end-to-end workflows are sped up. Furthermore, our architecture is non-intrusive: Web services owned and maintained by different institutions do not have to be altered prior to execution.

I. INTRODUCTION

Service-oriented architectures are an architectural paradigm for building software applications from a number of loosely coupled distributed services. Although the concept of service-oriented architectures is not novel, this paradigm has seen wide spread adoption through the Web services approach, which has a suite of simple standards (e.g., XML, WSDL and SOAP) to facilitate interoperability.

These core standards however do not provide the rich behavioural detail necessary to describe the role an individual service plays as part of a larger, more complex collaboration. Co-ordination of services is often achieved through the use of workflow technologies. As defined by the Workflow Management Coalition [13], a workflow is the automation of a business process, in whole or part, during which documents, information or tasks are passed from one participant (a resource, either human or machine) to another for action, according to a set of procedural rules.

Workflows in the scientific community are commonly modelled as Directed Acyclic Graphs (DAGs), formed from a collection of vertices (units of computation) and directed edges. The Genome Analysis and Database Update system (GADU) [18], the Southern California Earthquake Centre (SCEC) [11] CyberShake project, and the Laser Interferometer Gravitational-Wave Observatory (LIGO) [21] are examples of High Performance Computing applications composed using DAGs. DAGs present a *dataflow view*, here data is the primary concern, workflows are constructed from data processing

(vertices) and data transport (edges). In contrast, the Business Process Execution Language (BPEL) [22] is a *process-centric* language and is less well suited to modelling scientific scenarios [19].

Taverna [16] is an example of a popular graphical Web service composition tool used in the life sciences community in which workflows are represented as DAGs. Graph vertices can be one of a set of service types: WSDL Web services, BeanShell (lightweight scripting for Java) components, String constants etc. Services are given input and output ports which correspond to individual input and output variables. Edges are then formed by connecting services together by mapping output ports with input ports.

A. Motivating Scenario - Calculating Redshift

At this point, in order to put our motivation and problem statement into perspective, it is useful to consider an illustrative scenario. The Redshift scenario is taken from the AstroGrid [1] (UK e-Research project) science use-cases and involves retrieving and analysing data from multiple distributed resources. This scenario is representative of a class of large-scale scientific workflows, where data and services are made available through a Web service. It will be referenced throughout the remainder of this paper.

Photometric Redshifts use broad-band photometry to measure the Redshifts of galaxies. While photometric Redshifts have larger uncertainties than spectroscopic Redshifts, they are the only way of determining the properties of large samples of galaxies. This scenario describes the process of querying a group of distributed databases containing astronomical images in different bandwidths, extracting objects of interest and calculating the relative Redshift of each object.

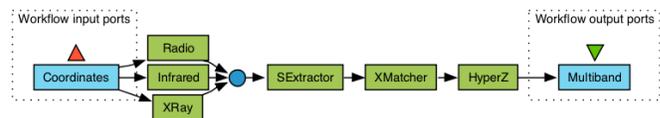


Fig. 1. AstroGrid scenario – Taverna representation. Workflow inputs are the RA and DEC coordinates, services are represented as rectangles, links correspond to the flow of data between services.

The scenario represents a workflow and begins with a scientist inputting the RA (right ascension) and DEC (declination)

coordinates into the system, which define an area of sky. These coordinates are used as input to three remote astronomical databases; no single database has a complete view of the data required by the scientist, as each database only stores images of a certain waveband. At each of the three databases the query is used to extract all images within the given coordinates which are returned to the scientist. The images are concatenated and sent to the SExtractor [5] tool for processing. SExtractor scans each image in turn and uses an algorithm to extract all objects of interest (positions of stars, galaxies etc.) and produces a table for each of the wavebands containing all the data. A cross matching tool is then used to scan all the images and produce one table containing data about all the objects of interest in the sky, in the five wavebands. This table is then used as input to the HyperZ¹ algorithm which computes the photometric Redshifts and appends it to each value of the table used as input. This final table consists of multi-band files containing the requested position as well as a table containing for each source all the output parameters from SExtractor and HyperZ, including positions, magnitudes, stellar classification and photometric Redshifts and confidence intervals; the final table is returned to the user. Figure 1 is a representation of the AstroGrid redshift scenario in Taverna.

B. Problem Statement

Although service-oriented workflows can be composed as DAGs using a dataflow model, in reality they are orchestrated from a single workflow engine, where intermediate data are typically routed through a centralised engine.

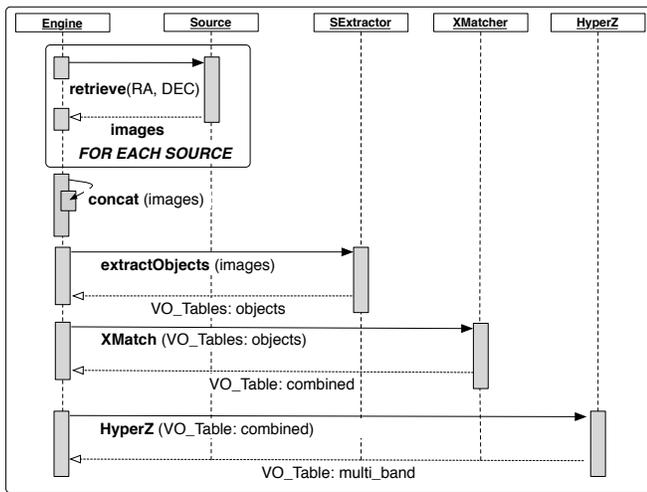


Fig. 2. UML Sequence diagram: AstroGrid scenario.

Figure 2 is a UML Sequence diagram displaying how the AstroGrid workflow is orchestrated. The initial RA and DEC coordinates are used as input to each of the three source databases: Radio, Infrared and XRay. Each source database then returns a set of images to the workflow engine. These

images are then combined and passed through the SExtractor, XMatcher and HyperZ services. Finally, the HyperZ service returns the Multiband table as output.

In the AstroGrid scenario, output from each of the source databases and processing services passes via the workflow engine, in order to be passed to the next service in the workflow chain. When one is orchestrating Web services from a tool such as Taverna, the workflow engine can become a bottleneck to the performance of a workflow. Standard WSDL Web services are designed with a request-response paradigm, therefore all intermediate data are routed via the workflow engine, which results in unnecessary data transfer, increasing the execution time of a workflow and causing the engine to become a bottleneck to the execution of an application.

Large-scale scientific applications (such as those described in Section I) can use third-party data transfers to move data between nodes in a workflow. However, nodes in these environments are usually closely coupled clusters, which have been specifically altered to support this kind of behaviour and do not tend to use Web services technology as a way of exposing application code. Choreography techniques describe service interactions from a global perspective, meaning that all participating services are treated equally, in a peer-to-peer fashion. Choreography techniques and infrastructure that supports third-party data transfers are discussed in relation to our research in Section IV.

C. Paper Contributions and Structure

This paper proposes a novel architecture for accelerating end-to-end workflows by deploying and executing a service-oriented workflow (composed as a DAG) across a peer-to-peer proxy network. Individual proxies are deployed at strategic network locations in order to exploit connectivity to remote Web services. By ‘strategic’ we could refer to one of many factors including network distance, security, policy etc. Proxies together form a proxy network which facilitate a number of functions such as service invocation, routing of intermediate data and data caching.

By breaking up a workflow and disseminating it across a peer-to-peer network, workflow bottlenecks associated with centralised orchestration are removed, intermediate data transfer is reduced and applications composed of interoperating Web services are sped up.

Importantly, our proposed architecture is a non-intrusive solution, Web services do not have to be redeployed or altered in any way prior to execution. Furthermore, a DAG-based workflow defined using a visual composition tool (e.g., Taverna) can be simply translated into our DAG-based workflow syntax and automatically deployed across a peer-to-peer proxy network.

A Java Web services implementation serves as the basis for performance analysis experiments conducted over Planet-Lab [8]. These experiments empirically demonstrate how our proposed architecture can speed up the execution time of a workflow when compared to standard orchestration. Although this paper focuses on Web services, the concept is generic

¹<http://webast.ast.obs-mip.fr/hyperz/> [26/02/2010]

and can be applied to other classes of application, i.e., High Performance Computing, Cloud Computing etc.

The remainder of this paper is structured as follows: Section II introduces the architecture of an individual proxy and multiple proxies which together form a proxy network. A syntax is introduced for service-oriented DAG-based workflows and the algorithms which divide a workflow, assign and enact a workflow across a proxy network are described. A Java Web services implementation, which serves as a platform for performance analysis is also discussed. Section III describes the performance analysis experiments. Related work is discussed in Section IV. Finally conclusions and future work are addressed in Section V.

II. PROXY ARCHITECTURE

A proxy is a middleware component that is deployed at a strategic location to a Web service or set of Web services. For the purposes of this paper, by strategic we mean in terms of network distance; as closely as possible to an enrolled service, i.e., on the same Web server or within the same domain, such that communication between a proxy and a Web service takes place over a local, not Wide Area Network. Proxies are considered to be volunteer nodes and can be arbitrarily sprinkled across a network, importantly not interfering with currently deployed infrastructure.

A proxy is generic and can concurrently execute any workflow definition. In order for this to be possible, the workflow definition is treated as an executable object which is passed between proxies in a proxy network. Proxies invoke Web services on behalf of the workflow engine, storing any intermediate data at the proxy. Proxies form peer-to-peer proxy networks and can route data directly to one another, avoiding the bottleneck problems associated with passing intermediate data through a single, centralised workflow engine.

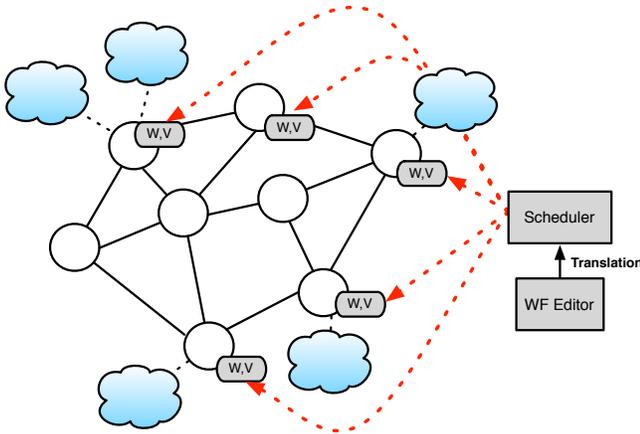


Fig. 3. Proxy architecture: Web services represented by clouds, proxies by circles, the workflow definition and vertex (W,V) by a rounded rectangle.

Figure 3 shows a high level architectural diagram of a proxy network. A user designs a service-oriented DAG-based workflow using a visual workflow editor such as Taverna.

A scheduling service assigns workflow vertices to proxies, the unique identifier of each proxy is then spliced into the workflow definition. Each proxy is passed an entire copy of the workflow definition. Once all proxies have received the workflow definition, each proxy executes its assigned set of vertices and passes any intermediate data directly to one another according to the directed edge definition. Once deployed a DAG-based workflow is executed without any centralised control over a peer-to-peer proxy network.

The following subsections describe in detail how a user designs a workflow, how the workflow definition is divided and assigned to a set of proxies, deployed across a proxy network and enacted.

A. Workflow Definition

A workflow is specified as a DAG according to the syntax displayed in Figure 4. We have taken inspiration from the Taverna SCUFL language [16], our syntax is a simplified version which does not support the additional processor types such as BeanShell etc.

Workflow	::=	$ID_w, \{Vertex\}, \{Edge\}, ID_s$
Vertex	::=	$vertex(ID_v, Processor)$
Processor	::=	$WS \mid Input \mid Output$
WS	::=	$s(Config^{(k)}, \{Inport\}, \{Outport\}, ID_p)$
Input	::=	$input(value, Outport)$
Output	::=	$output(value, Inport)$
Inport	::=	$in(ID_{in}, Type, [digit])$
Outport	::=	$out(ID_{out}, Type)$
ID	::=	$ID_w \mid ID_s \mid ID_v \mid ID_p \mid ID_{in} \mid ID_{out}$
Type	::=	XML RPC Types
Config	::=	$\langle name, value \rangle$
Edge	::=	$ID_v:ID_{out} \rightarrow \{ID_v:ID_{in}\}$

Fig. 4. Workflow definition syntax.

A workflow is labelled with a unique identifier ID_w and consists of a set of vertices $\{Vertex\}$ and a set of edges $\{Edge\}$. A vertex is given a globally unique identifier ID_v and can consist of one of a set of processor types. As this paper focuses on service-oriented workflows, processors are Web service definitions WS or input and output variables.

A Web service is defined firstly as a list of configuration pairs $Config^{(k)}$ which are simply $\langle name, value \rangle$ pairs and define the information necessary to invoke an external Web service: WSDL location, operation name etc. Secondly, a set of input ports $\{Inport\}$; each inport is given a unique identifier within a processor ID_{in} and a Type definition, types map to the standard set of XML RPC Types². The final parameter (optional) defines how many inputs are expected at a given inport. A set of output ports $\{Outport\}$; each outport is given a unique identifier within a processor ID_{out} and a Type definition. The final parameter of a processor is ID_p ,

²<http://ws.apache.org/xmlrpc/types.html> [26/02/2010]

which represents the globally unique identifier (mapping to an individual IP address) of a proxy which is executing a given vertex; these are initially null and spliced in before the workflow definition is disseminated to a set of proxies. Proxy identifiers are included in the workflow definition so that individual proxies can communicate with one another when executing a workflow. ID_s , the final parameter of a Workflow definition is the location (IP address) of the scheduling service, this is also spliced in before the workflow is disseminated so that the final output from a workflow can be passed back to the user.

The final processor types supported are Input (used as input to a Web service) and Output variables (used as output from a Web service). An Input variable is defined as a value, which is the actual value assigned to the variable and an output definition. An Output variable is defined as a value, which is initially a wildcard and an inport definition. The Types supported are the same set of XML RPC types.

In order to complete the workflow, a set of directed edges are formed which constitute a dataflow mapping from processor inports to processor outports. This is specified by providing the following mapping: $ID_v:ID_{out} \rightarrow \{ID_v:ID_{in}\}$. The types of the output to inport mappings must match and are enforced by the workflow editor.

B. Example Definition

Figure 5 is a representation of the AstroGrid scenario in the workflow definition syntax, Figure 6 is a corresponding diagrammatic representation. With reference to Figure 5, within the scope of the workflow identifier `calculate_redshift`, eight vertices are defined: `ra`, `dec` represent the workflow input parameters, the variables are defined by the physical values which are transferred via the outputs `ra_output` and `dec_output`. As the workflow output needs to be written back to a user's desktop, the vertex `wf_o` represents the final workflow output which will eventually be written to the inport `multi_band`; this output will then be passed back to the scheduling service which initiated the workflow.

The remaining vertices are WS definitions, `radio`, `infra` and `xray` are the distributed data sources, containing data from each of the required wave lengths; `tools` represents the co-located services `SExtractor` and `XMatcher`; finally `z` is the `HyperZ` processing service. Each of the service definitions contains typed inport and output definitions; note the 3 in the `tools` inport `images` states that 3 inputs (which will be merged) are required, for simplicity `config` represents the concrete details of individual Web services. The wildcard at the end of each service definition is the unique identifier of the proxy ID_p which is spliced in before the workflow definition is disseminated to a set of proxies. A set of Edge definitions connect vertex outports with vertex inports according to the flow of data in the AstroGrid scenario.

```

calculate_redshift,

//RA, DEC and Output Vertices
{vertex(ra, input(100, {out(ra_output, String)})),
vertex(dec, input(50, {out(dec_output, String)})),
vertex(wf_o, output(_, {in(multi_band, Object[])})),

//Source Vertices
vertex(radio, s(config, {in(ra_input, String),
in(dec_input, String)}, {out(image_set, byte[]), _}),
vertex(infra, s(config, {in(ra_input, String),
in(dec_input, String)}, {out(image_set, byte[]), _}),
vertex(xray, s(config, {in(ra_input, String),
in(dec_input, String)}, {out(image_set, byte[]), _}),

//Processing Vertices
vertex(tools, s(config, {in(images, byte[], 3)},
{out(combined, Object[]), _}),
vertex(z, s(config, {in(combined, Object[]),
{out(multi_band, Object[])})), _),

//Edge definitions
{ra:ra_output -> radio:ra_input, infra:ra_input,
xray:ra_input,
dec:dec_output -> radio:dec_input, infra:dec_input,
xray:dec_input,
radio:image_set -> tools:images,
infra:image_set -> tools:images,
xray:image_set -> tools:images,
tools:combined -> z:combined,
z:multi_band -> wf_o:multi_band}, _

```

Fig. 5. AstroGrid scenario workflow definition.

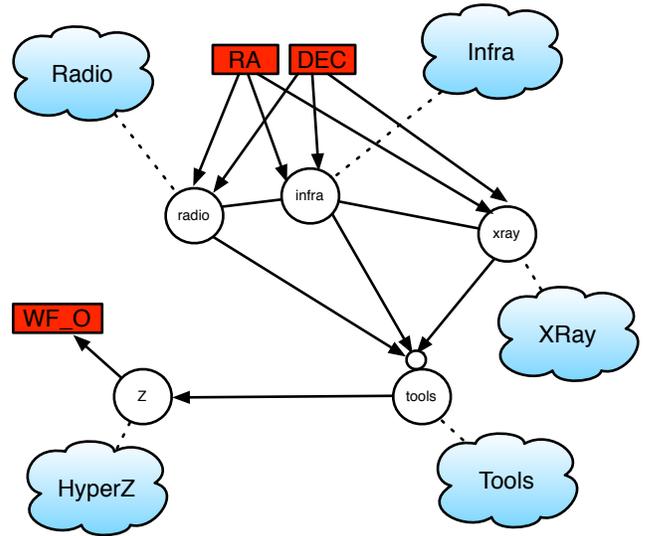


Fig. 6. Possible proxy configuration for the AstroGrid scenario: Edges are directed and show dataflow between proxies. Workflow inputs (RA, DEC) and outputs (WF_O) are also labelled. For simplicity, tools represents the co-located services `SExtractor` and `XMatcher`. Workflow engine and scheduling service not shown.

C. Web Services Implementation

The proxy architecture is available as an open-source toolkit implemented using a combination of Java (version 1.6) and Apache Axis (version 2) Web services [2], it consists of the following core components:

- **Registry service.** When a proxy is deployed it is automatically enrolled with the registry service, which contains global knowledge of distributed proxies. The registry service logs data of previous successful interactions and proxies are polled to ensure availability.

- **XML Syntax.** The workflow syntax displayed in Figure 4 is encoded in XML, allowing the registry service to splice in the proxy identifiers and proxies. Type checking between outport and inport definitions is enforced at the syntax level.

- **Translation.** A translation component automatically converts workflows defined in the Taverna SCUFL dataflow language into our workflow specification syntax. Translations from other languages are possible, we have chosen Taverna SCUFL as it is a widely accepted platform, particularly in the life sciences community.

- **Scheduling service.** Once a user has designed a workflow and it has been translated, the scheduling service (a local component) takes as input the workflow definition and consults the registry service, splicing in a unique proxy identifier for every vertex in a workflow definition. The scheduling service's IP address is spliced into the workflow definition, so that final output can be sent back to the user.

- **Proxy.** Proxies are made available through a standard WSDL interface, allowing them to be simply integrated into any standard workflow tool. As discussed further in Section II-E, a proxy has two remote methods: `initiate` to instantiate a proxy with a workflow and a vertex, and `data`, allowing proxies to pass data to one another, which in our implementation is via SOAP messages over HTTP. Proxies are simple to install and configure and can be dropped into an AXIS container running on an application server, no specialised programming needs to take place in order to exploit the functionality.

D. Workflow Deployment and Vertex Assignment

For simplicity in the algorithm definition, we assume reliable message passing and service invocation, however, fault tolerance has been built into the corresponding Web services implementation. A proxy is generic and can concurrently execute any vertex from any workflow definition. In order for this to be possible, the workflow definition is treated as an executable object, which is passed between proxies in a proxy network. The workflow definition is passed to the scheduling service which needs to assign proxies to vertices. This process is formally defined by Algorithm 1.

All proxies are enrolled with the registry service, which is a global directory of each proxy within the system. For each ID_v in $\{\text{Vertex}\}$ a suitable proxy must be located, if the processor type is a service definition, i.e., not an input or output variable. In our existing implementation, the registry service selects 'suitable' proxies which are deployed with the

Algorithm 1 Vertex assignment

```
1: for each Vertex  $ID_v$  where  $ID_v \in \{\text{Vertex}\}$  do
2:   if (Processor = WS) then
3:      $ID_p \leftarrow \text{locate}(ID_v, \text{WS})$ 
4:      $\text{WS.ID}_p \leftarrow ID_p$ 
5:      $\{\langle ID_v, ID_p \rangle\} \leftarrow \{\langle ID_v, ID_p \rangle\} + ID_v, ID_p$ 
6:   end if
7: end for
8: for each Vertex  $ID_v$  where  $ID_v \in \{\langle \text{proxy}, ID_v \rangle\}$  do
9:    $\text{initiate}(\text{Workflow}, ID_v)$ 
10: end for
```

same network domain as the WS it will eventually invoke. However, we are investigating optimisation techniques which will be addressed by further work, discussed in more detail in Section V.

This suitability matching is performed by the scheduling service which in turn consults with the registry service. The scheduling service invokes the `locate` method on the registry service, which takes as input a vertex ID_v and a WS definition and returns the unique identifier of a proxy which will enact a given vertex ID_v . This identifier is then spliced into the processor definition; before the assignment process begins all ID_p definitions are wildcards, each vertex (multiple vertices can be assigned to the same proxy) is then assigned before the workflow is disseminated. The proxy identifier is added so that proxies can communicate throughout the system globally.

The proxy identifier along with the vertex identifier are added to a set. Once the proxy assignment process is complete, the workflow definition and vertex a proxy is to assume ID_v is sent to each proxy in the set $\{\langle ID_v, \text{proxy} \rangle\}$. The remote method `initiate` is invoked on each proxy.

E. Workflow Execution

A proxy can concurrently execute any vertex from any workflow. With reference to Algorithm 2, in order to initiate a workflow, the remote method `initiate` is invoked on each proxy, given a workflow definition and ID_v . The vertex definition ID_v is extracted from the workflow. If the vertex relies on data from other vertices, it must wait until all inports have been resolved. Therefore, each inport ID_{in} in $\{\text{Inport}\}$ must be resolved before execution of ID_v can begin. This is achieved through the internal method `resolve_in` which checks if data for a given inport has been received; if the inport vertex definition is simply an input variable then the corresponding value is retrieved.

Once all inport dependencies have been resolved, given the unique workflow identifier ID_w and ID_v , the input data set $\{\text{input}\}$ is retrieved through the internal method `get_input`. The proxy takes the WSDL location, operation name and other parameters defined in $\{\text{config}\}$ and dynamically invokes the service using $\{\text{input}\}$ as input to the Web service. Results are then stored locally at the proxy.

In order to determine where (i.e., which proxy) to send the output of a given service invocation, the $\{\text{Edge}\}$ set is

Algorithm 2 Vertex enactment

```
1: initiate(Workflow, IDv)
2: for each Inport IDin where IDin ∈ {Inport} do
3:   resolve_in(IDin)
4: end for
5: {input} ← get_input(IDw, IDv)
6: results ← invoke({config}, {input})
7: for each Outport IDout where IDout ∈ {Outport} do
8:   {IDv:IDin} ← resolve_out(IDv:IDout, {Edge})
9:   for each Vertex IDv where IDv ∈ {IDv:IDin} do
10:    if (Processor = WS) then
11:      IDp ← WS.IDp
12:      data(IDw, IDin, results)
13:      delete(results)
14:    else if (Processor = Output) then
15:      value ← results
16:      data(IDw, IDs, results)
17:    end if
18:  end for
19: end for
```

inspected which contains mappings from a vertex output to a set of vertex inports. The set of inports which map to a corresponding output is returned by the internal `resolve_out` method. In order to determine which proxy to send these data to, each vertex in this set is traversed and the location of the proxy, ID_p is retrieved from the workflow definition.

The remote method `data` is invoked on the proxy ID_p , using the workflow identifier ID_w , the inport identifier ID_{in} and the result data as input. Once received (confirmed by an acknowledgement) by the recipient proxy, these data are stored according to ID_w and ID_{in} and deleted from the sender proxy. If the output corresponds to a `output`, this variable is written back to the scheduling service ID_s , which is running on a user’s desktop. This process is repeated for each output.

III. PERFORMANCE EVALUATION

In order to validate the hypothesis that our architecture can reduce intermediate data transfer and speed up the execution time of a workflow, a set of performance analysis experiments have been conducted. Our architecture has been evaluated across Internet-scale networks on the PlanetLab framework. PlanetLab is a highly configurable global research network of distributed servers that supports the development of new network services.

The AstroGrid scenario described throughout this paper serves as the basis for our performance analysis. This scenario is representative of a class of large-scale scientific workflows and has been configured as follows:

- **PlanetLab Deployment.** Data sources are a Web service which take as input an integer representing how much data the service is to return; the service then returns a Byte array of the size indicated in the input parameter. Analysis services are simulated via a sleep and return a set of data representative

of the given input size. These data sources and analysis services were deployed over the PlanetLab framework. The exact domains are displayed in Table I, bold domains indicate that multiple servers were selected from the same domain.

- **Workflow engine.** In order to benchmark our architecture two configurations of the AstroGrid scenario were set up: the first was executed on the completely centralised Taverna workflow (version 1.7.2) environment, the second was the same representation executed across a peer-to-peer proxy network according to the implementation described in Section II-C.

- **Speedup.** The mean speedup is calculated by dividing the mean time taken to execute the workflow using standard orchestration (i.e., non-proxy, fully centralised) and dividing it by the mean time taken to execute the workflow using our proxy architecture, e.g., a result of 1.5 means that the proxy architecture executes 50% faster than standard orchestration.

- **Proxy configurations.** Three different proxy configurations are shown: “same machine”, here a proxy is installed on the same physical machine as the Web service it is invoking, “same domain”, the proxy is installed on a different machine within the same network domain, and “distributed” where a proxy is installed on a node within the same country as the Web service it is invoking. In each configuration one proxy is responsible for one service.

- **Graphs.** Each configuration was executed 50 times across the PlanetLab framework. In each graph, the x -axis displays the mean speedup ratio (along with 95% confidence intervals from each of the 50 runs per configuration) and the y -axis displays the total volume of data flowing through the workflow. The number of services involved is independent of the mean speedup ratio as we have taken the mean ratio across a set of scaling experiments: we have scaled the initial data sources from 2 to 20 and repeated this while executing the AstroGrid DAG in reverse order. To prevent the data processing from influencing our evaluation, it has not been accounted for in the performance analysis experiments.

TABLE I
PLANETLAB NODE SELECTION.

France	Germany	USA
inria.fr	uni-goettingen.de	mit.edu
inisa.fr	uni-konstanz.de	brown.edu
utt.fr	uni-paderborn.de	poly.edu
	fraunhofer.de	umd.edu
	tu-darmstadt.de	byu.edu
		postel.org
		iit-tech.net

A. Experiment 1

In this experiment, each of the data sources, analysis services and the workflow engine were installed on separate PlanetLab nodes in the USA. As one can see from Figure 7(a) in all configurations our architecture outperforms the centralised workflow configuration. If one calculates the average across all data points for each of the experiments, the “same machine” configuration observes a speedup of 75%, the “same domain” configuration 49% and the “distributed” configuration 30%.

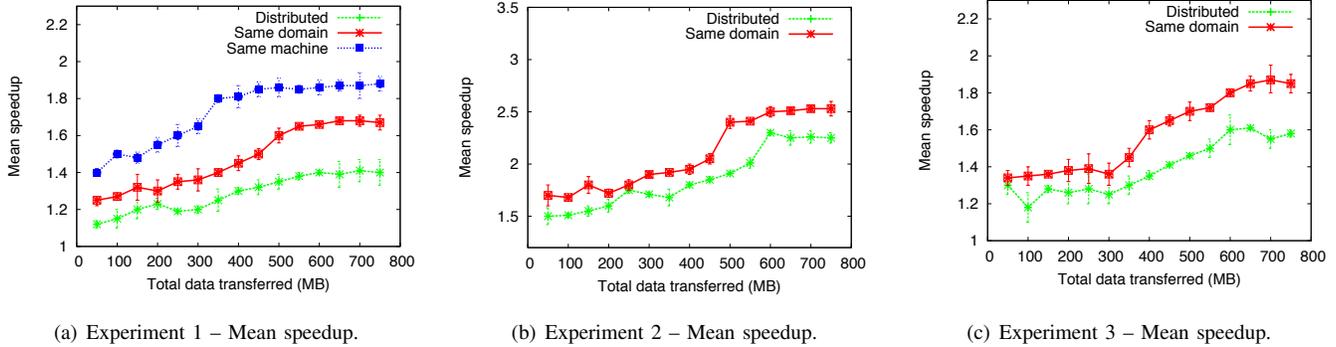


Fig. 7. PlanetLab performance evaluation.

As the results demonstrate, the speedup is greatest when a proxy is deployed as closely as possible to the back-end Web service, i.e., on the same machine. The cost of the proxy to service data movement increases as the proxy moves further away from the service it is invoking, in the “same machine” configuration, the input and output of a service invocation is flowing over a LAN. However, in the most extreme case, the “distributed” configuration an average speedup of 30% is observed over all runs.

B. Experiment 2

In this configuration, each of the data sources and analysis services were deployed on separate PlanetLab nodes across the USA. However, the workflow engine was now even further distributed from the services, running from a desktop machine in Melbourne. As one can see from Figure 7(b), as the cost (network distance) increases from the workflow engine to the workflow services, the hop any intermediate data has to make increases in cost. As the cost of intermediate data transport increases, the benefit of using our architecture grows as intermediate data transport is optimised. To quantify, this speedup ranged from 68% to 153% speedup for the “same domain” configuration and 51% to 125% for the “distributed” configuration.

C. Experiment 3

In order to distribute the services further, the data sources were deployed on PlanetLab nodes in the USA, the analysis services deployed on nodes in Europe (France and Germany) and the workflow engine was running from a desktop machine in Melbourne. With reference to Figure 7(c), speedup ranged from 34% to 85% for the “same domain” configuration and 30% to 58% for the “distributed” configuration. In this experiment one can observe an increased cost in distributing the workflow definition to each of the proxies prior to enactment, demonstrated by the lack of increase in mean speedup at lower data sizes.

Other experiments (not included in this paper due to space limitations) demonstrate similar trends.

IV. RELATED WORK

A. Choreography Languages

The research proposed by this paper is focused on decentralised orchestration rather than pure choreography: taking a workflow specification that is typically enacted by a centralised workflow engine, breaking it up and executing it through a set of decentralised peers.

There are relatively few languages targeted specifically at service choreography, the most prevalent being WS-CDL [4], BPEL4Chor [10] and Let’s Dance [23]. There are even fewer complete implementations of choreography languages, this means that choreography techniques are rarely deployed in practice. For example, there are only two documented prototype implementations of the WS-CDL specification. WS-CDL+, an extended specification [14] has been implemented in prototype form, although only one version, version 0.1 has been released. A further partial implementation [12] of the WS-CDL specification is currently in the prototype phase. The other widely known implementation is pi4soa³, an Eclipse plugin that provides a graphical editor to compose WS-CDL choreographies and generate from them compliant BPEL.

B. Techniques in Data Transfer Optimisation

Service Invocation Triggers [6] is an architecture for decentralised execution. Using the Triggers architecture, before execution can begin the input workflow must be deconstructed into sequential fragments, these fragments cannot contain loops and must be installed at a trigger; this is a rigid and limiting solution and is a barrier to entry for the use of proxy technology. In contrast with our proxy approach nothing in the workflow has to be altered prior to enactment.

The *Flow-based Infrastructure for Composing Autonomous Services* or FICAS [15] is a distributed data-flow architecture for composing software services. FICAS is intrusive to the application code as each application that is to be deployed needs to be wrapped with a FICAS interface.

In [7], an architecture for *decentralised orchestration* of composite Web services defined in BPEL is proposed. Our research deals with a set of challenges not addressed by

³<http://sourceforge.net/projects/pi4soa> [26/02/2010]

this architecture: how to optimise service-oriented DAG-based workflows, how to automatically deploy a workflow across a set of volunteer proxy nodes given a workflow topology, where to place proxies in relation to Web services, how these concepts operate across Internet-scale networks.

In our previous work [3] we proposed Circulate, a proxy-based architecture based on a centralised control flow, distributed data flow model. This paper has focused on executing DAG-based workflows without centralised control and explored a richer set of proxy functionality.

C. Third-party Data Transfers

This paper focuses primarily on optimising service-oriented workflows, where services are: not equipped to handle third-party transfers, owned and maintained by different organisations, and cannot be altered in anyway prior to enactment. For completeness it is important to discuss engines that support third-party transfers between nodes in task-based workflows.

Directed Acyclic Graph Manager (DAGMan) [9] submits jobs represented as a DAG to a Condor pool of resources. Intermediate data are not transferred via a workflow engine, instead they are passed directly from vertex to vertex. DAGMan removes the workflow bottleneck as data are transferred directly between vertices in a the DAG. *Triana* [20] is an open-source problem solving environment. It is designed to define, process, analyse, manage, execute and monitor workflows. Triana can distribute sections of a workflow to remote machines through a connected peer-to-peer network.

V. CONCLUSIONS AND FURTHER WORK

Through a motivating scenario, this paper has introduced an architecture for deploying and executing a service-oriented workflow (composed as a DAG) across a peer-to-peer proxy network. This architecture avoids workflow bottlenecks associated with centralised orchestration and reduces intermediate data transfer between interoperating services in a workflow. Importantly our proposed architecture is non-intrusive, Web services do not have to be altered in anyway prior to execution. Furthermore, users can continue to work with popular service-oriented DAG-based composition tools; our architecture translates DAG-based workflows (we have used Taverna SCUFL) into our workflow specification syntax, vertices are assigned, disseminated and enacted by an appropriate set of proxies.

A Web services implementation was introduced which formed the basis of our performance analysis experiments conducted over the PlanetLab framework. Performance analysis demonstrated across various network configurations that by reducing intermediate data transfer end-to-end workflows are sped up, in the best case from 68% to 192%.

Further work includes the following research challenges:

- **Expression of workflows.** This paper has focused on DAG-based workflows. Further work will address aligning the architecture with business process notations.

- **Peer-to-peer registry.** The registry service is currently centralised. Peer-to-peer techniques utilising Chord [17] are being investigated, with the view to improving scalability.

- **Security.** Although our system is only a prototype, to allow live deployment, security issues will have to be addressed.

REFERENCES

- [1] P. Allan, B. Bentley, and et al. AstroGrid. Technical report, Available at: www.astrogrid.org, 2001 [26/02/2010].
- [2] Apache Axis. <http://ws.apache.org/axis> [26/02/2010].
- [3] A. Barker, J. B. Weissman, and J. van Hemert. Orchestrating Data-Centric Workflows. In *The 8th IEEE International Symposium on Cluster Computing and the Grid (CCGrid)*, pages 210–217. IEEE Computer Society, 2008.
- [4] A. Barros, M. Dumas, and P. Oaks. A Critical Overview of the Web Services Choreography Description Language (WS-CDL). BPTrends Newsletter 3, 2005.
- [5] E. Bertin and S. Arnouts. Sextractor: Software for source extraction, Astronomy and Astrophysics, Suppl. Ser., 117:393–404, 1996.
- [6] W. Binder, I. Constantinescu, and B. Faltings. Decentralized Orchestration of Composite Web Services. In *Proceedings of ICWS'06*, pages 869–876. IEEE Computer Society, 2006.
- [7] G. B. Chaffe, S. Chandra, V. Mann, and M. G. Nanda. Decentralized orchestration of composite web services. In *Proceedings of the 13th international World Wide Web conference on Alternate track papers & posters*, pages 134–143. ACM, 2004.
- [8] B. Chun, D. Culler, T. Roscoe, A. Bavier, L. Peterson, M. Wawrzoniak, and M. Bowman. Planetlab: an overlay testbed for broad-coverage services. *SIGCOMM Comput. Commun. Rev.*, 33(3):3–12, 2003.
- [9] Condor Team. www.cs.wisc.edu/condor/dagman [26/02/2010].
- [10] G. Decker, O. Kopp, F. Leymann, and M. Weske. BPEL4Chor: Extending BPEL for Modeling Choreographies. In *Proceedings of ICWS'07*, pages 296–303. IEEE Computer Society, 2007.
- [11] E. Deelman and et al. Managing Large-Scale Workflow Execution from Resource Provisioning to Provenance tracking: The CyberShake Example. In *Proceedings of the Second IEEE International Conference on e-Science and Grid Computing*, 2006.
- [12] L. Fredlund. Implementing WS-CDL. In *Proceedings of the second Spanish workshop on Web Technologies (JSWEB 2006)*, 2006.
- [13] D. Hollingsworth. *The Workflow Reference Model*. Workflow Management Coalition, 1995.
- [14] Z. Kang, H. Wang, and P. C. Hung. WS-CDL+: An Extended WS-CDL Execution Engine for Web Service Collaboration. In *Proceedings of ICWS'07*, pages 928–935. IEEE Computer Society, 2007.
- [15] D. Liu, K. H. Law, and G. Wiederhold. Data-flow Distribution in FICAS Service Composition Infrastructure. In *Proceedings of the 15th International Conference on Parallel and Distributed Computing Systems*, 2002.
- [16] T. Oinn and et al. Taverna: a tool for the composition and enactment of bioinformatics workflows. *Bioinformatics*, 20(17):3045–3054, November 2004.
- [17] I. Stoica, R. Morris, D. Karger, M. F. Kaashoek, and H. Balakrishnan. Chord: A scalable peer-to-peer lookup service for internet applications. In *SIGCOMM '01: Proceedings of the 2001 conference on Applications, technologies, architectures, and protocols for computer communications*, pages 149–160, New York, NY, USA, 2001. ACM.
- [18] D. Sulakhe, A. Rodriguez, M. Wilde, I. T. Foster, and N. Maltsev. Interoperability of GADU in Using Heterogeneous Grid Resources for Bioinformatics Applications. *IEEE Transactions on Information Technology in Biomedicine*, 12(2):241–246, 2008.
- [19] W. Tan, P. Missier, R. Madduri, and I. Foster. Building Scientific Workflow with Taverna and BPEL: A Comparative Study in caGrid. In *ICSOC 2008 Workshops*, pages 118–129, 2008.
- [20] I. Taylor, M. Shields, I. Wang, and R. Philp. Distributed P2P Computing within Triana: A Galaxy Visualization Test Case. In *17th International Parallel and Distributed Processing Symposium (IPDPS 2003)*, pages 16–27. IEEE Computer Society, 2003.
- [21] I. J. Taylor, E. Deelman, D. B. Gannon, and M. Shields, editors. *Workflows for e-Science: Scientific Workflows for Grids*. Springer-Verlag, 2006.
- [22] The OASIS Committee. Web Services Business Process Execution Language (WS-BPEL) Version 2.0, 2007.
- [23] J. M. Zaha, A. Barros, M. Dumas, and A. ter Hofstede. Let's Dance: A Language for Service Behavior Modelling. In R. Meersman and T. Z. editors, *OTM Conferences (1)*, volume 4275 of LNCS, pages 145–162. Springer, 2006.