

## **Semantic Metadata Generation for Large Scientific Workflows**

Jihie Kim<sup>1</sup>, Yolanda Gil<sup>1</sup>, and Varun Ratnakar<sup>1</sup>,

<sup>1</sup> Information Sciences Institute, University of Southern California  
4676 Admiralty Way, Marina del Rey CA 90292, United States  
{jihie, gil, varunr}@isi.edu

**Abstract.** In recent years, workflows have been increasingly used in scientific applications. This paper presents novel metadata reasoning capabilities that we have developed to support the creation of large workflows. They include 1) use of semantic web technologies in handling metadata constraints on file collections and nested file collections, 2) propagation and validation of metadata constraints from inputs to outputs in a workflow component, and through the links among components in a workflow, and 3) sub-workflows that generate metadata needed for workflow creation. We show how we used these capabilities to support the creation of large executable workflows in an earthquake science application with more than 7,000 jobs, generating metadata for more than 100,000 new files.

**Keywords:** metadata reasoning, workflow generation, grid workflows.

### **1 Introduction**

Scientists have growing needs to use workflows to manage large distributed computations [13, 5, 2, 24]. In recent years, uses of large workflows have been significantly increased. Often they adopt grid-based environments that enable efficient execution of workflows by making use of distributed shared resources [22]. In such cases, computations in scientific workflows are represented as grid jobs that describe components used, input files required, and output files that will be produced as well as file movements, and deposition to distributed repositories [4].

Metadata describe the data used and generated by workflow components. Semantic web techniques have been applied for metadata reasoning on workflows such as validation of input parameters based on provenance data using component semantics [23], representing and managing dependencies between data products [14], helping scientists relate and annotate data and services through ontology-based generation and management of provenance data [25], etc. However, most of the existing metadata reasoning approaches focus on analyses of provenance data that are created from execution [18] rather than generation of input and output file descriptions needed in the workflow before execution.

The metadata reasoning capabilities of existing systems focus on files and simple collections and cannot effectively handle constraints on nested collections. Existing

checks on files are limited to validation of inputs for individual components. However, often there are global constraints on inputs and outputs of multiple components, and the workflow should be validated against such constraints in order to prevent execution of invalid workflows and wasting of expensive computations. In addition, unnecessary execution of individual components or multiple components in the given workflow should be detected and avoided when datasets that are equivalent to the ones to be produced already exist.

The creation of large workflows in the domains we use required several novel metadata reasoning capabilities:

- Keeping track of constraints on datasets used (i.e. files and file collections), including global constraints among multiple components as well as local constraints within individual components.
- Describing datasets that are used or created by the workflow.
- Detecting equivalent datasets and prevent unnecessary execution of workflow parts when datasets already exist.
- Managing large datasets and their provenance.

This paper presents novel metadata reasoning capabilities that we have developed to support the creation of large workflows. They include 1) use of semantic web technologies in handling metadata constraints on file collections and nested file collections, 2) propagation and validation of metadata constraints from inputs to outputs in a workflow component, and through the links among components in a workflow, and 3) sub-workflows that generate metadata needed for workflow creation. We illustrate these novel capabilities to support the creation of large workflows in an earthquake science application.

## 2 Motivation

A computational workflow is a set of executable programs (called *components*) that are introduced and linked together to pass data products to each other. The purpose of a computational workflow is to produce a desired end result from the combined computation of the programs. We will call a computational workflow as a *workflow* in this paper for brevity. Whereas a workflow represents a flow of data products among executable components, a *workflow template* is an abstract specification of a workflow, with a set of *nodes* and *links* where each node is a placeholder for a *component* or *component collections* (for iterative execution of a program over a file collection), and each link represents how the input and output parameters are connected. For example, Figure 1-(a) shows a template that has been used by earthquake scientists in SCEC (Southern California Earthquake Center) in Fall 2005. The template has two nodes (seismogram generation and calculation of spectral accelerations), each one containing a component collection. The workflow created from the template is shown in Figure 1-(b). This workflow was used in estimating hazard level of a site with respect to spectral acceleration caused by ruptures and their variations over time.

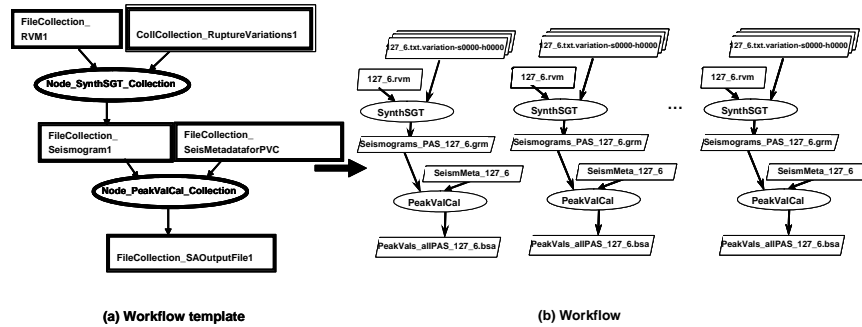


Figure 1: Workflow creation for seismic hazard analysis in Fall 2005.

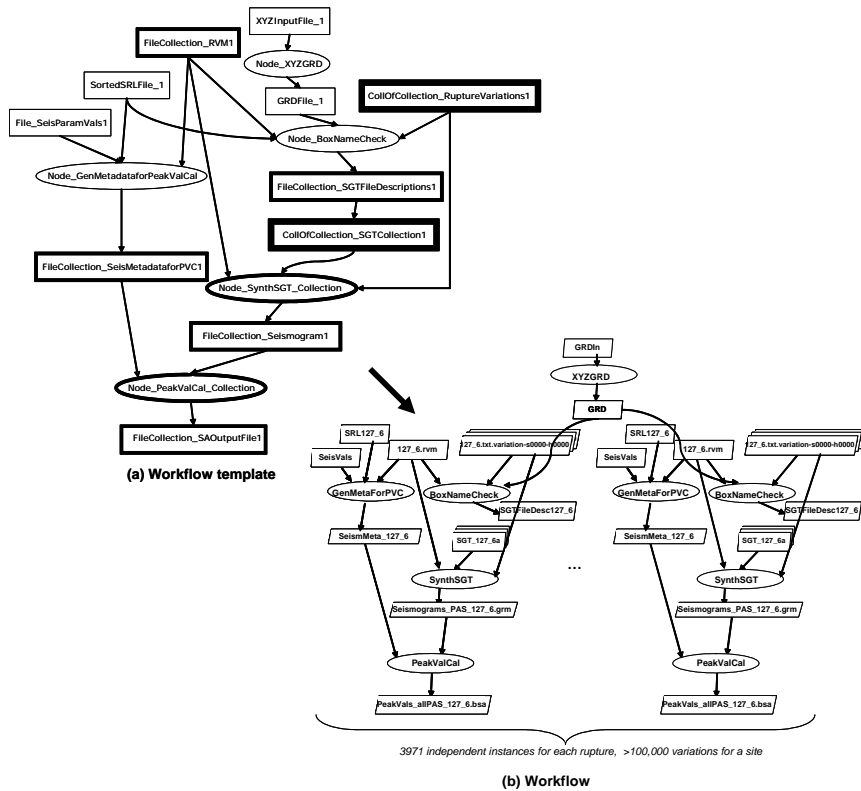


Figure 2: Workflow creation for seismic hazard analysis in Spring 2006.

The workflow was generated from manually created scripts that specify how to bind files to input parameters of the components and what are the expected output file names. An important feature of the workflow is that their data products are stored in files, often organized in directory structures that reflect the structure of the workflows. The names of the files and the directories follow conventions to encode

metadata information in the names such as the creation date or the relative area covered by the analysis. Therefore, the scripts that generate the workflow must orchestrate the creation of very particular data identifiers, namely file names that comply with those conventions and are instantiated to the appropriate constants. For example, a file containing the points for a hazard curve would be named using the rupture id and the fault id that were used in the simulation of the wave, as well as the lat-long of the location for the curve. The script included calls to functions or other scripts that generate information needed by the workflow (e.g. seismic parameter values). These manual ‘seam’ steps were not a part of the workflow. Most of the validation checks on the files and the collections were done by hand.

Figure 2-(a) shows an extension in the template in Spring 2006. This extension was needed to include strain green tensors (SGTs) as additional data input for seismogram generation. As the workflow template and descriptions of components become more complex, the script based approach becomes infeasible. First of all, there are more manual seam steps to handle. For example, since the SGT files that should be used in the workflow are unknown, the function that generates appropriate SGT file names should be executed beforehand. Validation of the workflow requires more checks. For example, now we need to check whether the SGTs use in generating seismogram are consistent with the rupture variations used for calculating peak values. If the seismogram generation step uses ruptures for Pasadena and their corresponding SGTs but the peak value calculation uses a rupture variation map for LA, the execution of the workflow will fail. When there exists a dataset that is equivalent to the expected output from executions of some components (e.g. SGT name datasets for Pasadena already exist), scientists had to identify them by hand.

In summary, generation of large workflows for this type of applications requires flexibility in adding or changing components to the template, systematic identification of files that are needed and generated by the workflow, incorporation of manual ‘seam’ steps into the workflow (making them a part of the workflow), and automatic validation of files and collections that are input to the workflow.

### 3 Approach

In developing new metadata reasoning capabilities for workflow creation, we use a workflow creation framework called Wings [6]. Wings takes a workflow template and initial input file descriptions, and creates an abstract workflow called DAX (DAG XML description). A DAX is transformed into an executable concrete workflow through a mapping that assigns available grid resources for execution by Pegasus [4]. Wings uses OWL-DL for representing files and collections, components, workflow templates, and workflows [6]. Currently Jena supports the reasoning.

In this work, Wings was extended to support metadata reasoning and generation. We developed an approach for representing metadata constraints on files and collections, and supporting metadata reasoning capabilities. Figure 3 shows an overview of the relevant components in the system, described in the following subsections. Although the descriptions rely on earthquake science examples, the same approach is used for other applications [6].

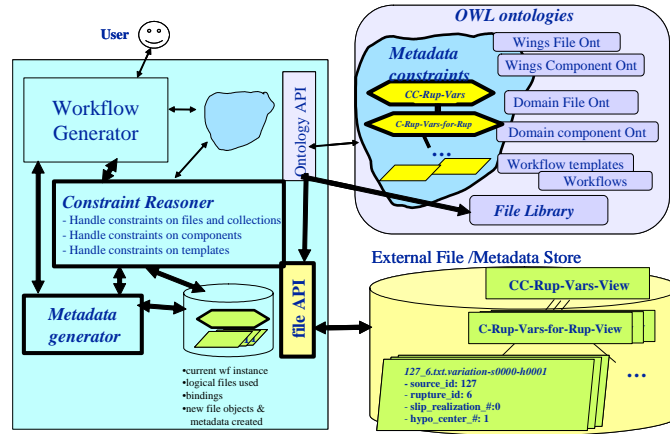


Figure 3: Metadata reasoning for workflow creation

### 3.1 Representing metadata constraints

One of the novel capabilities addresses the issue of keeping track of constraints on individual files, constraints on collections and their elements, constraints on inputs and outputs of each component, and global constraints among multiple components.

#### 3.1.1 Metadata constraints on individual files

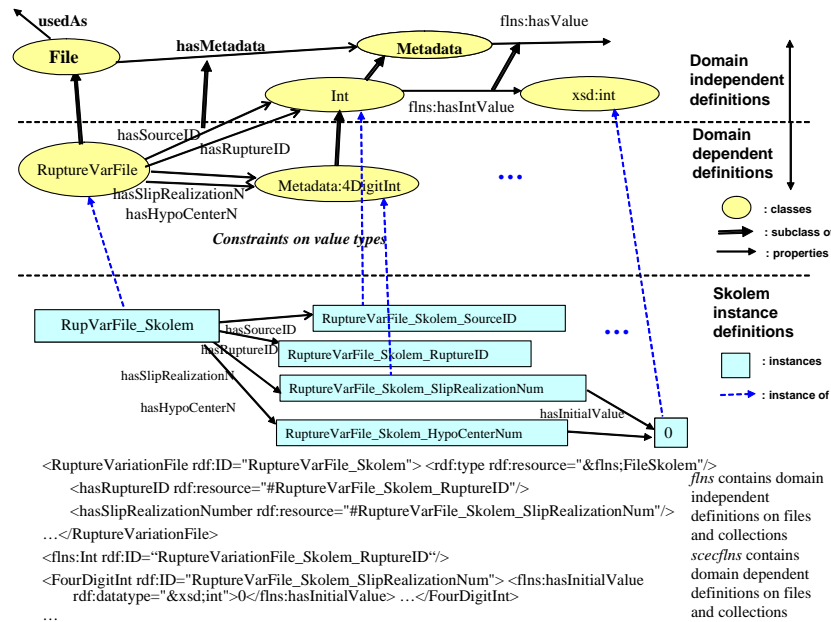


Figure 4: Metadata constraints on individual files

Each file class can have one or more metadata properties associated with it. In representing metadata constraints of a file class, we use a *skolem* instance (e.g., *RupVarFile\_Skolem*) that represents prototypical instances of the class. The metadata can describe what the file contains, how it was generated, etc. For example, a rupture variation file can have Rupture ID, SourceID, SlipRealizationN, and HypoCenterN that represent what it contains. Each metadata property has value ranges and can have some initial values. Other workflow generation functions such as how to derive filenames from metadata can be represented using the skolem instance. The actual metadata property values of file instances can be used in checking constraints on input and output files/collections used in the workflow, as described below.

### 3.1.2 Handling constraints on nested collections

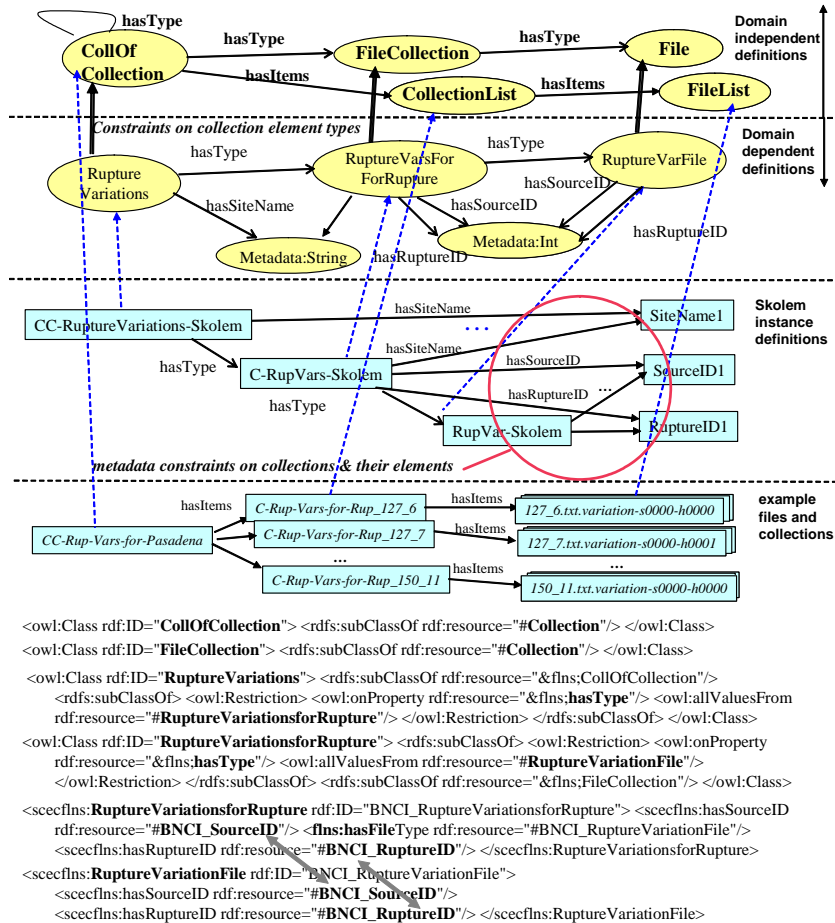
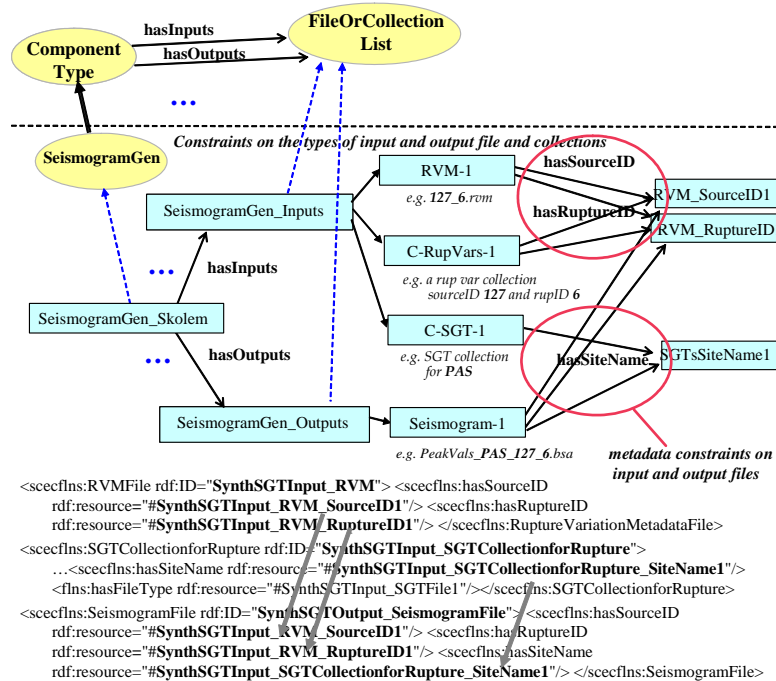


Figure 5: Nested file collections and their metadata constraints.

In general, for a given site (e.g. Pasadena), several ruptures are used in performing the hazard analysis. According to rupture dynamics of earthquakes that depend on hypocenter and slip values, each temporal variation of the stress is described in a rupture variation file. That is, rupture variations for a site are naturally structured as a collection of file collections. In our ontology, the concept *collection* represents both simple file collections and nested collections. Each collection should specify the type for the collection element using the ‘hasType’ property. There can be constraints between a collection and its elements. For example, for a rupture variation collection for a rupture, the SourceID and the RuptureID of individual rupture variation file should be the same as the rupture’s SourceID and RuptureID. That is, if the rupture variation collection for a rupture has SourceID 127 and RuptureID 6, each element (a rupture variation file) should have SourceID 127 and RuptureID 6. Figure 5 shows how these constraints on collections and nested collections are represented with skolem instances.

### 3.1.3 Constraints on components: constraints on input and output files and collections



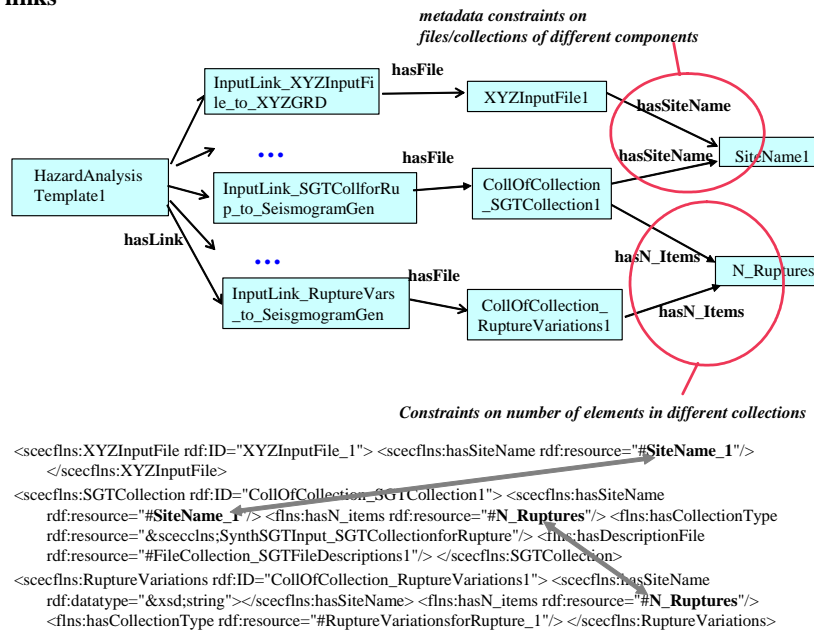
**Figure 6:** Constraints on metadata properties of input/output files or collections.

Each workflow component is described in terms of its input and output data types. In Figure 6, the SeismogramGen component has three inputs: an RVM (rupture variation map) file, a rupture variation collection, and a SGT file collection. Each RVM file has a SourceID and a RuptureID of the rupture that it represents. In order to create valid

results, their values should be the same as the RuptureID and SourceID of the input rupture variation collection. The input SGT collection should have a site name associated with it. Given these inputs, the SeismogramGen component produces a seismogram file.

The metadata for the generated seismogram file depends on the metadata of the inputs. In the above example, the site name of the SGT collection (PAS), and the SourceID and RuptureID of the RVM file (127 and 6) are propagated to corresponding metadata properties of the output seismogram file. The procedure for metadata validation and propagation during workflow creation is described in Section 3.2.

### 3.1.4 Global constraints on templates: constraints among different nodes and links



**Figure 7:** Global constraints on metadata properties among files and collections used by different components in a template

There are additional validation checks that should be made in order to create a valid workflow. First of all, the components should use seismic data for the same site (e.g. PAS) in performing hazard analysis. In Figure 7, the site name of the XYZinput file used in generating a mesh for simulation should be the same as the site name of the SGT collection of collections. (We also use a isSameAs property in representing equalities of metadata.) In addition, the components should use the same number of ruptures throughout the workflow. For example, the number of elements in a collection of collection rupture variations indicates the number of ruptures used in modeling the site. This number (i.e. the number of ruptures) should be the same as the number of elements (SGT collections) in the collection of collection SGTs that are



used. If the specific number of ruptures is known, the value can be given for the `N_Ruptures` using the `flns:hasValue` property. Figure 7 shows the current representations. In representing these global constraints, we make use of *link skolems*. Each link skolem is a placeholder for a file or collection that is bound to the input and output parameters of the components associated with the link during workflow creation. If more than one link skolems in a template share the same metadata objects, when the bindings for the links are created their corresponding metadata values should be the same. These constraints are used by metadata reasoner in creating consistent and correct workflows. The details of metadata based validation are described below.

### 3.2 Metadata propagation and validation

**Table 1:** Steps for propagating metadata and checking constraints during workflow creation.

<p><b>Bind&amp;ValidateWorkflow</b> (WorkflowTemplate wt, InputLinks ILinks)</p> <ol style="list-style-type: none"> <li>1. Assign ILinks to LinksToProcess.</li> <li>2. While LinksToProcess is not empty             <ol style="list-style-type: none"> <li>2.1. Remove one from LinksToProcess and assign it to L1.</li> <li>2.2. Let F1 be the link skolem for binding files or collections to L1.</li> <li>2.3. <i>If metadata for F1 should be generated from an execution of a component</i> <ol style="list-style-type: none"> <li>2.3.1. <i>if the execution results are not available, continue.</i> <ul style="list-style-type: none"> <li>;; i.e. exclude this link in the sub-workflow</li> </ul> </li> </ol> </li> <li>2.4. If any metadata of F1 depends on a link L2 that is not bound yet,               <ol style="list-style-type: none"> <li>2.4.1. Mark L1 as a dependent of L2 and continue.</li> </ol> </li> <li>2.5. If L1 is an input link,               <ol style="list-style-type: none"> <li>2.5.1. Get metadata of the file from the user or a file server</li> <li>2.5.2. <b>Check consistencies</b> with links that L1 depends on</li> <li>2.5.3. <b>Check consistencies</b> with existing bindings based on template-level constraints</li> <li>2.5.4. If any metadata are inconsistent, report inconsistency and return.</li> <li>2.5.5. <b>Bind</b> file/collection name and metadata to F1.</li> <li>2.5.6. If the file type for F1 is a collection, recursively get the metadata of its elements</li> </ol> </li> <li>2.6. Else (i.e. L1 is InOutLink or OuputLink)               <ol style="list-style-type: none"> <li>2.6.1. Generate file names and metadata base on the definition of the depending links.                   <ul style="list-style-type: none"> <li>;; <b>metadata propagation</b></li> </ul> </li> </ol> </li> <li>2.7. For each link L2 that is dependent on l1,               <ol style="list-style-type: none"> <li>2.7.1. if all the links that L2 is depending on are bound, put L2 in LinksToProcess.</li> </ol> </li> <li>2.8. If L1 is an output link, continue.</li> <li>2.9. Else (L1 is InputLink or InOutLink)               <ol style="list-style-type: none"> <li>2.9.1. If all the inputs to the destination node (i.e. the component that L1 provides an input to) have been bound,                   <ol style="list-style-type: none"> <li>2.9.1.1. Add all the OutputLinks and InOutLinks from the destination node to the LinksToProcess.</li> </ol> </li> </ol> </li> </ol> </li> </ol>
--

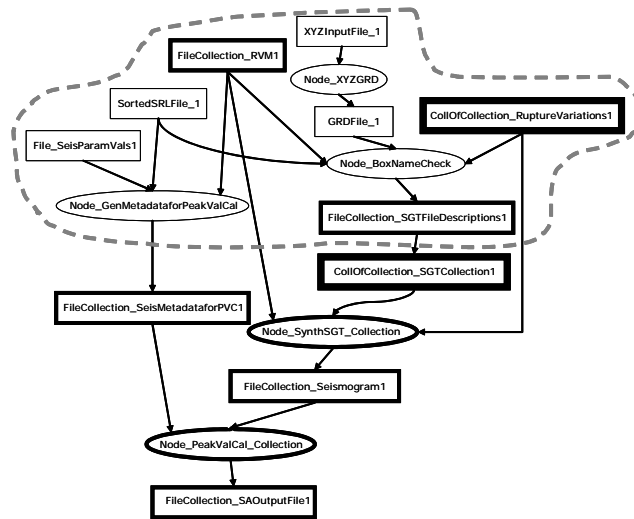
Table 1 shows the procedure for propagating metadata constraints and validating workflows created using metadata constraints. The procedure significantly extends the existing Wings algorithm by including steps for metadata propagation and validation checks. It traverses links in the workflow template and generates consistent bindings for link skolems. There are three classes of links: InputLink, InOutLink, and OutputLink. An InputLink is a link from an initial input file or collection to a node. Each InOutLink represents a connection from an output parameter of a node to an input parameter of another node. An OutputLink represents an end result from a node. The procedure specifies how the system starts with the input links of a template, identifies dependencies among the links based on definitions of metadata constraints, binds link skolems to files or collections, propagates and checks constraints of the

bindings based on metadata constraints, and traverses the next unbound links based on the dependencies.

A link I1 is dependent on I2 if some of the metadata of I1 needs to be filled in based on some metadata of I2. For example, in Figure 6, the metadata of the output of SeismogramGen step depends on the metadata of the RVM file and the SGT collection. The input link for a rupture variation collection depends on the input link for an RVM, if the SourceID and the RuptureID of the rupture variations are derived from the values in the RVM file. We assume that there are no cyclic dependencies in the definition of metadata constraints.

The file names and the metadata for initial input files or collections can be given from the user or existing file library (in OWL) through a file API. The metadata of the initial inputs can also be retrieved from other external file stores using the same API. Currently we use a web repository, but we are exploring uses of grid catalogs such as MCS (Metadata Catalog Service) [19]. The italicized steps handle sub-workflows, which are explained in the next section.

### 3.3 Sub-workflows for generating metadata needed for workflow creation



**Figure 8:** Identifying and executing sub-workflows for full workflow creation.

As described in Section 2, creation of workflows needed manual ‘seam’ steps that call functions that generate information needed by the workflow, such as file names and parameter values. In order to minimize such manual steps, we have created new workflow components that model such steps. For example, in Figure 8 (an enlargement of Figure 2-(a)), individual files in CollOfCollection\_SGTCollection1 are unknown initially and the file names should be generated by executing the BoxNameCheck component. Previously, the execution of BoxNameCheck was done manually. We represent such components as workflow components, and link them to the depending component inputs or outputs (e.g. SGT files needed by SynthSGT) in the template.

In generating grid workflows, for each execution of a component, the names of the inputs and output files for the component should be specified beforehand. That is, what data are created, and what data are staged in and out of the computation should be known before execution. Names (or descriptions) of some of the files in the workflow are not given initially, and their names should be automatically generated from metadata of other files.

As shown in Table 1, our Bind&ValidateWorkflow procedure checks these dependencies, and generates a ‘sub’-workflow that includes only the parts that can be instantiated with the currently available data. For the template in Figure 8, a sub-workflow with bindings for input and output links of the three components (XYZGRD, GenMetaForPeakValCal and BoxNameCheck), highlighted with dotted lines, is generated. The resulting sub-workflow is mapped to grid resources through Pegasus [5] and executed in a grid environment. The execution of a sub-workflow provides results for dependent input/output links, such as file names needed for component inputs or outputs. The metadata for these new file names are generated and added to our file repository by the metadata generator (shown in Figure 3) so that they can be used in creating an expanded workflow. The creation and execution of sub-workflows can be interleaved until the complete workflow is generated. The above workflow template needs only one iteration of sub-workflow creation and execution.

## 4 Results

**Table 2:** Number of files and OWL instances created during workflow generation

	Workflow creation time	Number of file instances created for the workflow	Number of OWL individuals created
A sub workflow for hazard analysis	7 minutes, 59 seconds	15,888	322,473
A full workflow for hazard analysis	22 minutes, 52 seconds	117,379	2,001,972

The above metadata reasoning capabilities are used in creating workflows for seismic hazard analysis. In creating a workflow for an LA site with the template in Figure 8, there were about 3,971 ruptures and 97,228 variations of ruptures to take into account. As the number of files and file collections become large, many OWL objects that represent file and collections and their metadata should be created and queried. The number of files in the workflow we have represented was 117,379, as shown in Table 2. The number of OWL individuals created was over two million. (We excluded the anonymous individuals that are created as a by-product of rdf:list in the count, so the actual number is larger.) For the full workflow, the DAX included 7,945 jobs. Large workflows pose challenges on computational resources (CPUs and memory) used during workflow creation. Currently it takes about 8 minutes to create a sub workflow and about 23 minutes to generate the full workflow on a Pentium 4 3.0GHz with 1GB of RAM. The current system is being used for other applications including statistical natural language processing tasks where parallel processing of a large corpus is needed.

In order to efficiently perform the required metadata reasoning with many objects, we split a workflow into multiple independent workflow parts and create each separately. In splitting, we make use of metadata properties that can divide collections into independent sub-collections. For example, separate sub-trees in Figure 2-(b) can be independently generated. We currently use the SourceID to split rupture file collections into sub-collections. Other collections such as rupture variation collections are divided using the same set of metadata properties. Currently we select such metadata properties by hand, but we are investigating an automatic approach that takes into account sizes of file collections. The independent workflow parts are accumulated in the workflow generator and are automatically merged in the end, creating a complete workflow.

Using the same collection splitting approach described above, we can store the resulting files and collections into separate file library entities. The objects can be selectively loaded and used in creation of new workflows. Equivalent files or collections can be identified using metadata, which enables detection of unnecessary execution of components or workflow parts that will produce equivalent datasets.

## **5 Related Work**

Semantic web techniques have been used in supporting many e-science workflow systems [10, 7]. Applications include semantic description of web services, resource discovery, data management, composition of workflow templates [17, 20, 1], etc. Our work complements existing work by supporting creation of large workflows needed for data and/or compute intensive scientific applications.

Recently various data management and provenance techniques have been developed for e-Science applications [18,8]. Most of the existing work focuses on pedigree or lineage metadata that describes the data resources and the processes used in generating data products. These provenance metadata are often used in qualifying data products and supporting data management and reuse. Our current work focuses on metadata reasoning that support workflow creation and validation. The metadata that are generated during workflow creation can be used in combination with other provenance metadata for supporting file reuse. Our work extends existing approaches for validating workflows in that we take into account constraints on nested collections and global constraints among multiple components as well as constraints on inputs within individual components [23,14]. Another difference is that we make use of metadata in generating valid workflows before execution instead of validating already executed workflows with provenance data, also enabling detection of unnecessary jobs before execution.

## **5 Conclusion and Future Work**

We presented a semantic metadata generation and reasoning approach that supports creation of large workflows. Given the metadata of initial input files, the system propagates metadata constraints from the inputs to the outputs, and through the links

among the components during workflow creation. Both global constraints among multiple components and local constraints are used for workflow validation. The files that will be produced from workflow execution as well as the input files are identified during the metadata propagation and validation process. Some of the metadata are generated through creation and execution of sub-workflows when the metadata need to be computationally generated. Because we are able to identify data collections and their properties before the workflow is executed, we can detect whether the data has been generated before by querying an existing data repository. This is important for optimizing execution performance: If some intermediate data product already exists then there is no need to re-execute the portion of the workflow that produces it. We also use the metadata in managing large collections and their provenance.

We are currently working on extensions of the workflow template shown in Figure 2-(a) and they will use more datasets for seismic analysis of different sites in Southern California. In order to further improve the efficiency of the workflow creation and metadata reasoning, we are considering several extensions to our system. One area of improvement is creating a scalable metadata repository. Currently we can store metadata in multiple OWL file libraries, but we are planning to explore its integration with MCS that can store metadata of data products (such as files) published on the Grid [19]. With this approach, when there are new files and metadata added to MCS by a different workflow or system, we will be able to use them in creating new workflows. In order to perform iterative sub-workflow generation and execution more efficiently, we are investigating a client-server style approach where our system can call a workflow execution server with a newly generated sub-workflow, and the execution results can be notified to our system (a client). The newly generated metadata during workflow creation can be used in combination with other metadata for data provenance applications. For example, the metadata can tell whether the two files (or collections) contain the same kind of information, even when they are generated from different workflows. We are exploring various uses of metadata in relating datasets used in scientific workflows.

**Acknowledgments.** We thank David Okaya, Philip Maechling, Scott Callaghan, Hunter Francoeur, and Li Zhao in the Southern California Earthquake Center (SCEC) for valuable discussions on seismic hazard analysis workflows. We would also like to thank Gaurang Mehta and Ewa Deelman for their help in executing workflows with Pegasus. This research was funded by the National Science Foundation (NSF) with award number EAR-0122464. The SCEC contribution number for this paper is 1016.

## References

1. Altintas, I., Berkley, C., Jaeger, E., Jones, M., Ludäscher, B., Mock, S.: Kepler: Towards a Grid-Enabled System for Scientific Workflows. The Workflow in Grid Systems Workshop in GGF10 - The Tenth Global Grid Forum, Berlin, Germany (2004)
2. Campobasso, M., Giles, M.: Stabilization of a Linear Flow Solver for Turbomachinery Aeroelasticity Using Recursive Projection Method. AIAA Journal, 42(9) (2004)

3. Churches, D., Gombas, G., Harrison, A., Maassen, J., Robinson, C., Shields, M., Taylor, I., Wang, I.: Programming Scientific and Distributed Workflow with Triana Services. Grid Workflow Special Issue of Concurrency and Computation: Practice and Experience (2004)
4. Deelman, E., Blythe, J., Gil, Y., Kesselman, C., Mehta, G., Patil, S., Su, M., Vahi, K., Livny, M.: Pegasus: Mapping Scientific Workflows onto the Grid. Across Grids Conference (2004)
5. Deelman, E., Blythe, J., Gil, Y., Kesselman, C.: Workflow Management in GriPhyN. The Grid Resource Management, Kluwer (2003)
6. Gil, Y., Ratnakar, V., Deelman, E., Spraragen, M., Kim, J.: Wings for Pegasus: A Semantic Approach to Creating Very Large Scientific Workflows. Internal project report (2006)
7. Goble, C.: Using the Semantic Web for e-Science: Inspiration, Incubation, Irritation. Lecture Notes in Computer Science 3729:1-3, (2005)
8. Goble, C.: Position Statement: Musings on Provenance, Workflow and (Semantic Web) Annotations for Bioinformatics. Workshop on Data Derivation and Provenance (2002)
9. Guo, Y., Pan'Z., Heflin, J.: An Evaluation of Knowledge Base Systems for Large OWL Datasets. Proc. of the Third International Semantic Web Conference (2004)
10. Hendler, J.: Science and the Semantic Web. Science 299 (2003) 520-521
11. Hustadt, U., Motik, B., Sattler, U.: Data Complexity of Reasoning in Very Expressive Description Logics. Proc. of the 19th International Joint Conference on AI (2005)
12. Kim, J., Spraragen, M., Gil, Y.: An Intelligent Assistant for Interactive Workflow Composition. Proceedings of the Intl. Conference on Intelligent User Interfaces (2004)
13. Maechling, P., Chalupsky, H., Dougherty, M., Deelman, E., Gil, Y., Gullapalli, S., Gupta, V., Kesselman, C., Kim, J., Mehta, G., Mendenhall, B., Russ, T., Singh, G., Spraragen, M., Staples, G., Vahi, K.: Simplifying Construction of Complex Workflows for Non-Expert Users of the Southern California Earthquake Center Community Modeling Environment. ACM SIGMOD Record, special issue on Scientific Workflows, 34 (3) (2005)
14. Myers, J., Pancerella, C., Lansing, C., Schuchardt, K., Didier, B.: Multi-scale Science: Supporting Emerging Practice with Semantically-Derived Provenance. Semantic Web Technologies for Searching and Retrieving Scientific Data Workshop (2003)
15. openRDF, 2006: <http://www.openrdf.org/> (2006)
16. OWL Web Ontology Language, 2006: <http://www.w3.org/TR/owl-features/> (2006)
17. Sabou, M., Wroe, C., Goble, C., Mishne, G.: Learning Domain Ontologies for Web Service Descriptions: an Experiment in Bioinformatics. Intl. Conf. on World Wide Web. (2005)
18. Simmhan Y., Plale B., Gannon, D.: A Survey of Data Provenance in e-Science. SIGMOD Record, vol. 34, 2005, pp. 31-36 (2005)
19. Singh, G., Bharathi, S., Chervenak, A., Deelman, E., Kesselman, C., Manohar, M., Patil, S., Pearlman, L.: A Metadata Catalog Service for Data Intensive Applications. SC (2003)
20. Sirin, E., Parsia, B., Hendler, J.: Filtering and selecting semantic web services with interactive composition techniques. IEEE Intelligent Systems, 19(4) (2004)
21. Sycara, K., Paolucci, M., Ankolekar, A., Srinivasan, N.: Automated Discovery, Interaction and Composition of Semantic Web services. Journal of Web Semantics 1(1) (2003)
22. TeraGrid 2006. NSF Teragrid Project, <http://www.teragrid.org/> (2003)
23. Wong, S., Miles, S., Fang, W., Groth, P., Moreau, L.: Validation of E-Science Experiments using a Provenance-based Approach. Proc. of 4th Intl. Semantic Web Conference (2005)
24. Wroe, C., Goble, C., Greenwood, M., Lord, P., Miles, S., Papay, J., Payne, T., Moreau, L.: Automating Experiments Using Semantic Data on a Bioinformatics Grid. IEEE Intelligent Systems special issue on e-Science (2004)
25. Zhao, J., Goble, C., Stevens R., Bechhofer, S: Semantics of a Networked World: Semantics for Grid Databases. Proc. of the First International IFIP Conference, ICSNW (2004)