

Applying Software Engineering Technology to Support the Clear and Precise Specification of Scientific Processes

Leon J. Osterweil¹, Lori A. Clarke¹, Rodion Podorozhny², Alexander Wise¹, Emery Boose³, Aaron M. Ellison³, Julian Hadley³

ljo@cs.umass.edu, rp31@txstate.edu, aellison@fas.harvard.edu, clarke@cs.umass.edu, boose@fas.harvard.edu, wise@cs.umass.edu, jhadley@fas.harvard.edu

¹University of Massachusetts, Department of Computer Science, Computer Science Building, Amherst, Massachusetts 01003 USA

²Texas State University, San Marcos, TX USA

³Harvard University, Harvard Forest, 324 North Main Street, Petersham, Massachusetts 01366 USA

Abstract With the availability of powerful computational and communication capabilities, scientists now readily access large, complicated derived datasets and build on those results to produce, through further processing, yet other derived datasets of interest to themselves and others. These scientific processes must be clearly documented so that scientists can evaluate the soundness of the scientific processes used to create these datasets, reproduce these results, and build upon them in responsible and appropriate ways. Here, we present the concept of an *analytic web*, which defines the scientific processes employed and the details of the exact application of those processes in creating derived datasets. The work described here is similar to work often referred to as “scientific workflow”, but emphasizes the need for semantically richer, more rigorously defined process definition languages, such as those that were first developed to define software engineering processes. This paper illustrates the information that might comprise an analytic web for a scientific process for measuring and analyzing the flux of water through a forested watershed. This is a complex and demanding scientific process that illustrates the benefits of using a semantically rich executable language for defining the process, supporting automatic creation of process provenance data, assuring data reproducibility, and serving as the basis for analysis of the data’s scientific soundness.

1 Introduction

1.1 The Problem

Modern computation and communication systems have dramatically changed the way in which science is done. These systems enable scientists to work with datasets and to create models of their research subjects that are far larger and more detailed than were possible in the past. Faster computing speeds enable far more ambitious analyses of these models, leading to the production of far greater quantities of derived scientific datasets. Ever faster global networks make these datasets accessible to scientists around the world. While these new computational and communications capabilities have opened up the

possibility of exciting new research, they have also lead to new challenges and problems. When new datasets, often the products of extensive processing, are generated, they are most often promulgated without adequate documentation that describes their creation. If scientists are to make appropriate use of the datasets produced by others and avoid misuse by inappropriate application of subsequent processing, then it is imperative to know how such datasets were produced. Indeed, before scientific results can be accepted, they should be reproduced by other scientists, as such reproducibility is a fundamental requirement for results to be confirmed. Consequently, the very conduct of science requires that the way in which scientific results are produced be carefully documented for others.

An obvious way to address these problems is to associate as an *annotation* a precise description of a dataset along with each dataset itself . Such annotations, essentially data items that describe data, are called *metadata*. There have already been many calls for the use of metadata, which typically document such details as the date of generation of a dataset, the name of the investigator, and perhaps some specifications of the hardware and software systems used. We argue that it is necessary to go further. We suggest that a particular type of metadata annotation, *process provenance* metadata, be attached to all datasets, and when necessary to individual data items. The benefits of such process provenance metadata include facilitation of the reproduction of the data by others, expedited identification of data items and datasets of interest, and better understanding of which forms of subsequent processing should, and should not, be applied to data items and datasets.

It is our view that concepts drawn from the domain of software engineering can provide the basis for the generation and association of such process provenance metadata with the derived data items and datasets.

1.2 Analytic Webs

Scientific datasets can be viewed as products emerging from a distributed enterprise: input datasets may be stored and retrieved remotely, analytic services may be obtained from external sources, and the datasets that are the products of a scientist's work are increasingly likely to be made directly accessible. The totality of data and capabilities produced and consumed by a working scientific team in pursuit of a particular scientific objective can be thought of as a scientific (usually online) data processing enterprise, and we refer to it by the term *analytic web* [Boose, et.al. 2007; Ellison, et. al. 2006; Osterweil, et.al. 2005]. One view might be that the purpose of an analytic web is to expedite the participation of a scientific team in the marketplaces of scientific investigation. From this view there is then a need to provide strong support for two distinct activities that scientists engage in routinely, namely the *production* of datasets for such marketplaces and the *consumption* of datasets from such marketplaces, for use in their subsequent scientific investigation.

Production: As raw data are collected by individual investigators or by sensors, they may be pre-processed with data loggers, field computers, etc.; for example, many environmental measurements (*e.g.*, temperature, solar radiation, carbon

flux, water flow) are sampled at high frequencies but only hourly averages are stored. Individual researchers take these data and post-process them, entering them into data repositories such as spreadsheets, checking them for errors, and transforming them into datasets used for analysis by the investigators themselves or by others. These datasets are often stored on the investigator's personal computer, where they might be further analyzed (and in the process new datasets may be created) and condensed for publication. The producers of these datasets could use help in managing the execution of these increasingly intricate processes, in documenting exactly what processing was applied to which data items, and in reasoning about the soundness of the resulting data items and datasets. These data items and datasets, on which publications are often based, should be archived on institutional servers or within national repositories, where it is expected that they may be accessed and used by others: *i.e.*, dataset consumers.

Consumption: Datasets collected by other investigators increasingly are available via the Internet and may be reanalyzed to verify existing models and results or used to generate new models, hypotheses, and scientific insights. We refer to this use of previously existing datasets as dataset consumption. It is not atypical for a consumer who synthesizes datasets to subsequently become a producer of new (synthetic) datasets that are consumed by others (who then may become producers of further datasets).

The conceptualization of dataset production and consumption suggests specific ways in which an analytic web should support a scientific team's activities. Support for production should consist of facilities for generating and storing new data items and datasets. In recognition that others may consume these datasets, the datasets and when necessary the individual data items should be annotated with precise process provenance metadata, and the analytic web should provide facilities for accessing such metadata and evaluating its subsequent use. This may entail reproducing the dataset, evaluating its use in further scientific data processing, or using the process generation metadata as a guide to generation of other datasets.

To make this conceptual vision of an analytic web a useful reality, we propose that a specific analytic web be realized by a set of tools aimed at creating, analyzing, and managing two types of closely interrelated graph structures, namely *Dataset Derivation Graphs* (DDGs) and *Process Derivation Graphs* (PDGs). The purpose of a DDG is to organize datasets into a structure based upon the way in which the datasets are derived from each other. The purpose of a PDG is to define precisely the processes by which these derivations are performed. The PDG also serves as the vehicle for executing the process definition and generating subsequent data items and datasets. Moreover, execution of a PDG can result in the automatic creation of the DDG for each of these data items and datasets.

In this paper we provide a concrete example of the way in which an analytic web can support the activities of a specific scientific research team. This example illustrates that DDGs are likely to be large structures that need to be defined in detail and with great care. Thus, this example illustrates the importance of automatically executing the PDG and, during that execution, automatically generating the corresponding DDGs. A major objective of this example is to demonstrate a set of semantic features that seem essential

in a process definition language employed to define the PDGs used by actual working scientific teams. An existing process definition language, Little-JIL [Wise 2006; Wise, et.al. 2000; Cass, et.al. 2000], incorporates many of these semantic features. Key features of Little-JIL are highly effective in defining the PDGs needed to support our example scientific process, but some important semantic features are missing.

2 Motivating Example

Measuring and forecasting water flux and storage in the ecosystem (including ground water, soils, surface water, snow pack, vegetation, and atmospheric boundary layer) is of tremendous importance to society, of pressing interest to scientists, and a central focus of major scientific investigation efforts, such as the NEON project [NEON URL] and the Waters Initiative [WATERS URL] (formed from the synergies of CUASHI [CUASHI URL] and CLEANER [CLEANER URL]). Such forecasts require detailed hydrological measurements of natural and human-dominated ecosystems; these measurements increasingly are coming from vast networks of real-time sensor systems that may be subjected to elaborate real-time adjustments and considerable, perhaps iterative, post-processing over the ensuing months or years. Producers of such datasets could use help in systematically and correctly applying various processing and analysis tools. Consumers of these datasets will require support to help them understand the datasets, to reproduce them if desired, and to use them appropriately in further processing.

To address this scientific problem, a group of ecologists at the Harvard Forest Long-Term Ecological Research (LTER) [Harvard Forest URL] site is currently designing a real-time system for estimating water budgets for two small, forested watersheds at their site. Their system is being designed to calculate change in water storage using the water balance equation: $dS = P - ET - Q$, where the change in water storage (dS) is a linear function of precipitation (P), evapotranspiration (ET), and surface discharge or streamflow (Q). The complete system will include additional measurements of snow pack, soil moisture, ground water, *etc.* Equation inputs come initially from five streams of real-time data from three sources:

- **Precipitation (P)** – 15-minute precipitation totals ($P1$, $P2$) measured at two rain gauges. Two gauges are used in order to minimize gaps in the data.
- **Surface Discharge (Q)** – 15-minute average stream flow values measured at a stream gauge. Short gaps in Q caused by sensor failure, excessive ice build-up, etc. can be filled by modeling Q as a function of preceding measurements of P and Q .
- **Evapotranspiration (ET)** – 30-minute average ET values measured at an eddy-flux (micro-meteorological) tower.
- **Photosynthetically active radiation (PAR)** – 30-minute average PAR values measured at the eddy-flux tower. PAR is measured continuously because over short time spans (up to several weeks) PAR is the environmental variable most often highly correlated with ET . Thus, PAR can be used to estimate ET when it cannot be accurately measured due to instrument failure, moisture on the sonic anemometers

within the eddy-covariance system, very low air turbulence, or unfavorable interactions of wind with local topography.

- This system will incorporate three features that are typical of virtually any sensor network and raise challenging issues for dataset producers as well as dataset consumers:

1. Real-time quality control entails non-trivial processing, much of it determined on the fly, which may cause different data items in a dataset to have different process provenance. For example such systems may incorporate duplicate sensors (here *P1* and *P2*), real-time modeling (yielding estimates \hat{Q} and \hat{ET} of *Q* and *ET*, respectively), and rules for value selection. Thus, in this case the two precipitation measurements are to be compared and specific actions taken if the values differ by more than a specified amount. To check for instrument problems, *Q* and *ET* are to be compared to their modeled values (\hat{Q} and \hat{ET}) using respectively (i) verified runoff models and recent precipitation and flow values and (ii) regressions of recent *ET* and *PAR* values. Modeled values also may be substituted (imputed) for measured values when sensors fail or wind conditions do not support reliable eddy-flux measurements.

2. Regularly scheduled post-processing of data (e.g., after 30 days) is used so that individual imputed data values can be computed using models that take into account values coming from both preceding and subsequent measurements. This feature is necessary because real-time modeling is retrospective only (i.e., based on past data). Thus, post-processing is important especially for times when the ecosystem is undergoing rapid change (e.g., during spring leaf-out or when soils become saturated during heavy or prolonged precipitation).

3. Alternative past measurements also may become available and may then be included in additional datasets. This may be desirable, for example, to make use of recent measurements that did not arrive in time for real-time processing, corrections of earlier measurements as a result of later detection of sensor drift, or replacement or imputation of faulty or missing values with measurements from other sources. Note that the substitution of alternative measurements must be accompanied by post-processing of a sufficient number of preceding and subsequent values to ensure that any “ripple effects” – possible impacts on subsequent modeled values – are accounted for. This indicates another way in which it becomes possible for each data item of a dataset to be the product of a different process.

Figure 1 is a data flow graph (DFG) representation of the water budget process. In this data flow graph, rounded boxes represent types of tools or subprocesses and types of datasets are represented by boxes with a clipped corner. Edges connecting these boxes represent the flow of data or datasets into and out of the tools and subprocesses. In this figure, the **Real-time Selection Criteria** box represents the on-the-fly processing of the **Real-time Data** that is used to create new or updated **Models** and to update the model variables \hat{Q} and \hat{ET} . Similarly, the **Retrospective Selection Criteria** and the **Alternative Selection Criteria** boxes represent the reprocessing of the **Real-time** and **Retrospective Data**, respectively, and can also result in new **Models** and updated model

values. In subsequent sections of this paper we indicate shortcomings of this DFG representation and suggest another graph notation that seems to offer better facilities for representing complex processes such as this one.

3 Proposed Model of an Analytic Web

As noted above, we propose that the realization of an analytic web be centered on the production and management of two complementary graph structures, a *Dataset Derivation Graph* (DDG) and a *Process Derivation Graph* (PDG). Both types of graphs are illustrated for this example.

3.1 Dataset Derivation Graph

A DDG documents the specific data items or dataset instances created when a producer applies processes (e.g. perhaps defined using a PDG), using specific tools and

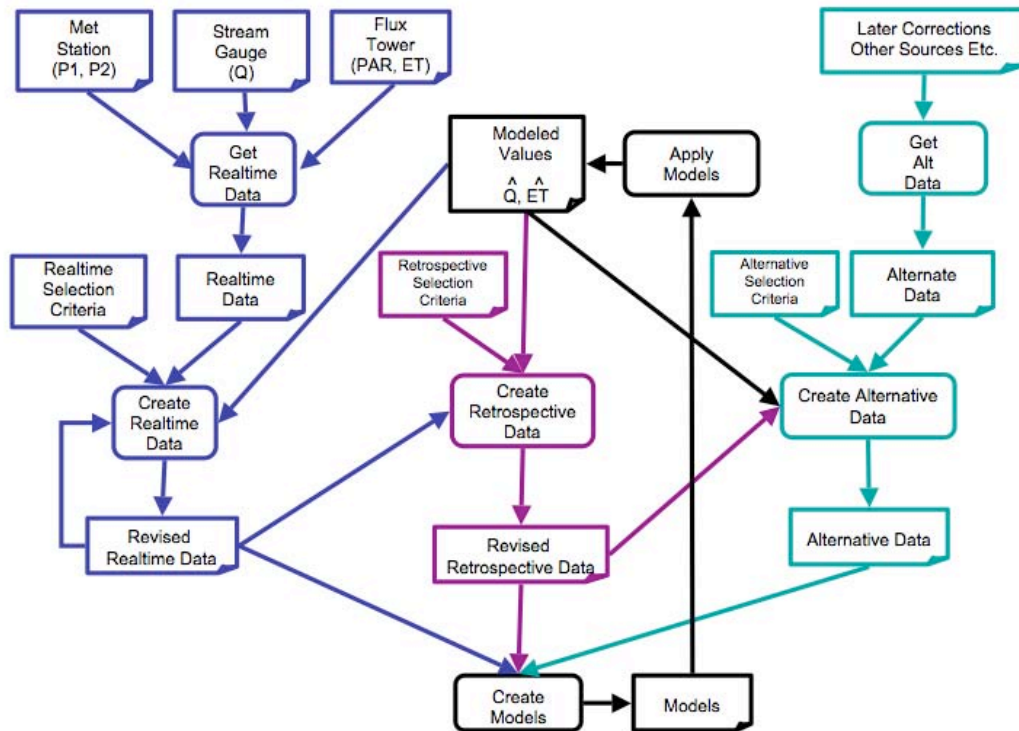


Figure 1: A data flow graph of the water budget process (To make the diagram easier to understand, each processing loop, real-time, retrospective, and alternative processing as well as model maintenance, is represented in a different color, with each edge colored by the same color as the node it emanates from.)

subprocesses on specific input data items or dataset instances. A DDG thus contains detailed metadata about the process sequences used to build all of its datasets. This metadata is thus an example of the *process provenance metadata* needed to inform consumers about how a dataset was generated and to support reproduction of the dataset by other scientists. A DDG, as depicted in Figure 2, uses rectangles to represent specific data items and dataset instances and ovals to represent specific tools or subprocess instances. There is an edge from each data item and dataset instance node to the process instance node from which it was derived (unless the data item or dataset instance represents raw data that was not previously processed). Each oval process instance node is connected by one or more edges to the data item(s) and dataset instance(s) that it used as input(s) to derive the indicated output data item or dataset instance. Each time a process is executed, a new set of data items and dataset instances is created, and these data items and dataset instances, as well as the process instances that created them, must be added as nodes of the evolving DDG. Each DDG node instance can be stored independently with a unique URL for identification.

Data items and dataset instances, such as those represented by DDG nodes, are the usual focus of scientific attention and thus are the objects intended to be documented with metadata such as specified by Ecological Metadata Language (EML) [EML URL]. For example, *Transpiration Data For 2/26/06* is a specific instance of type transpiration data and, thus is depicted in Figure 2 by a box. Figure 2 also shows that *Aaron's outlier rejection parameters* is another dataset instance. These two instances are taken as inputs by a specific processing tool, namely *Reject Outlying Data*, to produce the output dataset instance, *Cleaned 2/26/06 Data*. This dataset instance is, in turn, taken as input

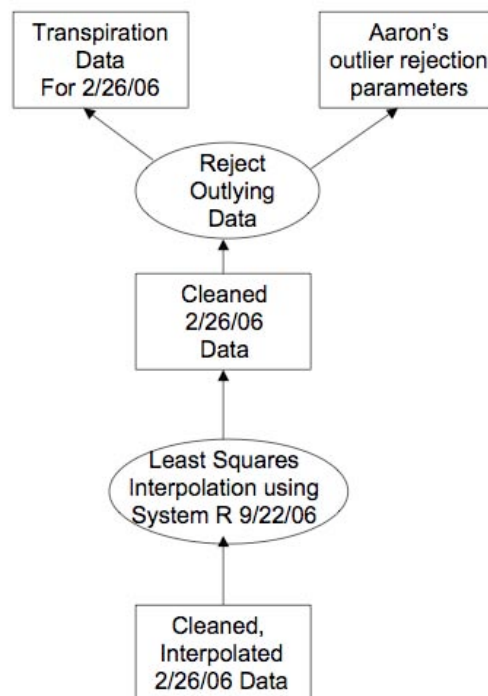


Figure 2: An example of a Dataset Derivation Graph

to the processing tool Least Squares Interpolation using System R 9/22/06, producing as output the dataset instance Cleaned, Interpolated 2/26/06 Data. Note that the specification of the specific interpolation tool, complete with version number and date of application, are important as they specify precisely which tool or system was applied to the dataset instance. Different versions of a tool may produce different results. Consumers who want to fully understand the provenance of a dataset, and producers who want to reproduce it exactly, require the documentation provided by a DDG: the specific datasets and tools that were *actually used*. The quantity and intricacy of this documentation is considerable, but it can be produced automatically with the aid of a suitable suite of tools, as described below.

To make this clearer, Figure 3 illustrates a DDG that could have been generated by executing a sequence of tools and subprocesses represented as a path through the boxes in the diagram shown in Figure 1. Note that this DDG shows the result of two iterations, represented by two traversals through a loop in the diagram in Figure 1, resulting in the creation of three instances of Revised Data, denoted here by Revised Data 1, Revised Data 2, and Revised Data 3. Revised Data 1 is simply the set of data items that resulted from filtering the initial real-time data stream by applying some specific filtering criterion. This generally results in a dataset where some data item values are missing, most often due to intentional deletion. Revised Data 2 results from applying initial \hat{Q} and \hat{ET} models to this dataset, thereby filling in missing data item values and replacing others. Revised Data 3 results from creating and applying a second pair of \hat{Q} and \hat{ET} models, (e.g., from regularly scheduled post-processing) and applying them to Revised Data 2. This DDG provides documentation of the exact processing steps that were taken to produce these datasets. It specifies what datasets must be generated to reproduce these processing steps, and thus the dataset that is their final result.

3.2 Process Derivation Graphs

A Process Derivation Graph (PDG) is a precise representation of the sequences of steps used to process the data items and datasets in performing a scientific process (and thereby can be used as the basis for generating DDGs). According to this definition, we note that a DFG, could in relatively simple cases be used as a PDG. But for complex scientific processes, DFGs usually do not accurately describe the allowed sequence of process steps since they often do not distinguish between those paths that are intended to be part of the scientific process and those that are not. For example, the DFG in Figure 1 does not preclude the retroactive processing loop from preceding the real-time processing loop for the same dataset. For this reason, such diagrams are typically unreliable for specifying which paths through the graph can be used safely to generate additional DDG nodes. The challenges indicated in Section 2 suggest that a PDG should instead be defined by means of a language that incorporates a diverse and powerful set of semantic features in order to suffice as an adequate representation of a scientific process. Among the desirable features are: hierarchical decomposition, abstraction mechanisms; concurrent processing; exception management and reaction to contingencies; and methods for dealing with streaming real-time data.

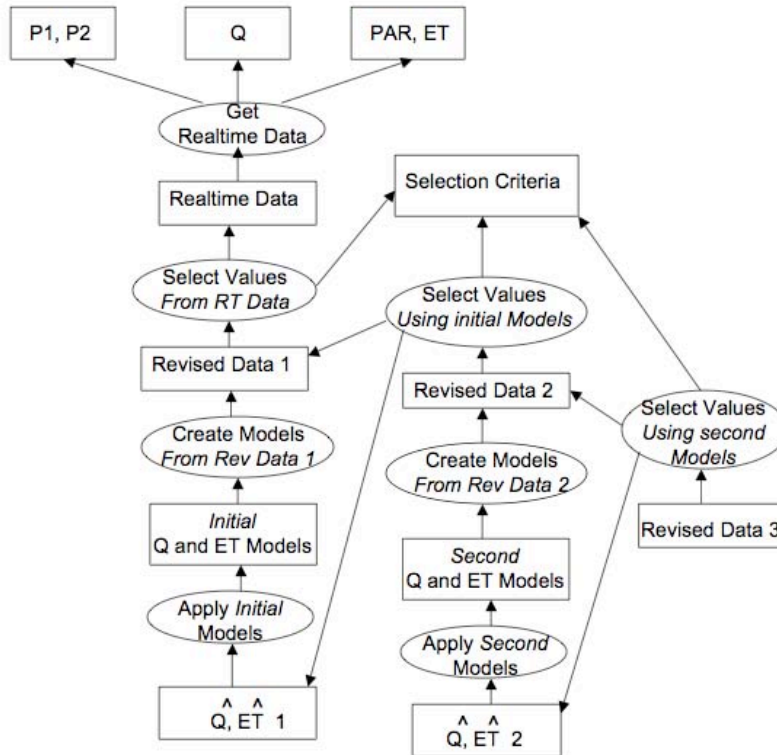


Figure 3: The DDG showing the complete provenance of Revised Data 3

Hierarchical Decomposition: Hierarchy is commonly used to simplify the representation of large and complex processes, which can become quite cumbersome, unless high-level abstract processes can first be introduced and their details subsequently and separately elaborated. Such a facility is typically effective for dealing with complexity.

Abstraction: It is not enough simply to illustrate processes hierarchically. In addition it seems necessary to take the modular processing abstractions that scientists think in terms of, and use actual process execution contexts to tailor and specialize the abstractions into specific processing steps that fit the specific execution context. Modern programming languages support the use of abstraction in this way as an effective approach for dealing with the substantial challenges of complexity and scaling. One key mechanism for using context information to tailor abstractions is by binding artifacts as inputs and outputs to individual subprocesses. This capability is tantamount to a capability for passing arguments to a procedure that can perform differently in different contexts. Further, different instantiations of a process may offer different exception handlers (see below), thereby establishing different execution contexts.

Concurrency: Many scientific processes and activities can occur simultaneously (in parallel). For example, data streams from various sensors may be processed in parallel. In addition, there may be a (much slower) parallel activity of generating new models and

evaluating both new and existing models. It is important to be able to represent accurately the concurrency of these various activities so that the right data items are handled by the right processing steps at the right times. There are, in addition, situations in which it is important to synchronize activities taking place in parallel. Both producers and consumers need a clear picture of such concurrency and synchronization to assure them that what has been defined is indeed what is desired.

Exception Management: In the real world, a broad range of contingencies can arise that require specialized reactions. In handling real-time data streams, for example, it is necessary to deal with sensor failures in ways that do not interfere with the continued flow of data. Different failure modes may dictate different responses. Suitably powerful process language features are needed to support the specification of such exception management.

The focus of the rest of this paper is on demonstrating how these semantic capabilities are needed by delving into the details of the Harvard Forest water budget estimation process. Our view is that these needs are projected most clearly by using a specific process definition language. We use Little-JIL, a visual process definition language originally developed for defining software engineering processes, as an example of a language that incorporates many of the features needed in order to define the PDG of an analytic web sufficiently completely and precisely. Little-JIL is a high-level language that can produce hierarchically structured definitions of scientific processes that can be invoked in different contexts. Little-JIL has precise, executable semantics (defined by finite-state automata). Among its key features are facilities for scoping to determine context, for specifying parallel processing, for procedural abstraction, for handling exceptional conditions, and for specifying and controlling iteration. Limitations of Little-JIL that arise in supporting the definition of the process in this example are also pointed out.

3.3 Semantics of Little-JIL

A process is defined in Little-JIL using hierarchically decomposed steps (Figure 4), where a step represents a task to be done by an assigned agent. Each step has a name and a set of badges to represent control flow among its sub-steps, its interface (a specification of its input and output data), the exceptions it handles, etc. A step with no substeps is called a *leaf step* and represents an activity to be performed by an agent, without any guidance from the process.

Resources and Agents: As part of its interface, each Little-JIL step contains a specification of the type of agent that is required to assume responsibility for the step's execution. The agent specification is a specification of a capability. It is assumed that this specification will be considered by a separate Resource Manager that maintains a repository of available resources, along with their capabilities. The Resource Manager is expected to identify a specific resource instance to be bound as the Agent in response to the step's need for a specified capability. Little-JIL agents may be either humans or automated devices. There are cases where either might be appropriate or

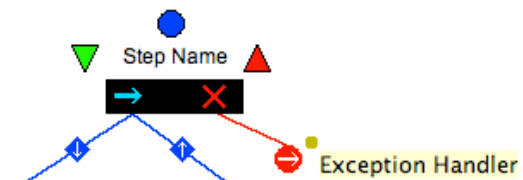


Figure 4: A Little-JIL step icon

where there are multiple capable agents and then the choice should be made by the Resource Manager, rather than being dictated by the process definition. A step may also specify the need for resources other than the agent.

Substep Decomposition: Little-JIL steps may be decomposed into substeps of two different kinds, ordinary substeps and exception handlers. The *ordinary substeps* define the details of just how the step is to be executed. The substeps are connected to the parent step by edges, which may be annotated by specifications of the artifacts that flow between parent and substep, and also by cardinality specifications. Cardinality specifications define the number of times the substep is to be instantiated and may be a fixed number, a Kleene * (for zero or more times), a Kleene + (for one or more times), or a Boolean expression (indicating whether the substep is to be instantiated or not). *Exception handlers* define the way in which exceptions thrown by the step's descendants are to be handled. The edge connecting an exception handler to its parent is annotated with the type of the exception being handled, as well as an indication of how execution is to continue after the exception has been handled.

Step sequencing: Every non-leaf step has a *sequencing badge* (an icon embedded in the left portion of the step bar; e.g., the right arrow in Figure 4), which defines the order in which its substeps execute. For example, a sequential step (right arrow) indicates that its substeps are to be executed sequentially from left to right and is only considered completed after all of its substeps have completed. A parallel step (equal sign) indicates that its substeps can be executed in any (possibly arbitrarily interleaved) order. It, too, is considered completed after all of its substeps have completed. A choice step (circle slashed with a horizontal line) indicates that the human executing the step is to make a choice between any of the step's substeps. A try step (right arrow with an X on its tail) mandates a sequence in which substeps are to be tried in any order until one completes successfully.

Artifacts and artifact flows: An *artifact* is an entity (e.g., a datum or dataset) that is used or produced by a step. Parameter declarations are specified in the interface to a step (circle atop the step bar) as lists of the artifacts used by the step (IN parameters) and the artifacts produced by the step (OUT parameters). Artifact flow through steps can be defined to take place in either of two different ways, 1) hierarchically, as the flow of artifacts between parent and child steps, and 2) by means of channels. The flow of artifacts between parent and child steps is (as noted above) indicated by attaching to the edges between parent and child identification of the artifacts as well as arrows indicating the direction of flow of each artifact.

Channels: *Channels* are named entities that directly (without the need for hierarchical flow) deliver artifacts produced by specifically identified source step(s) as arguments to specific destination step(s). A Little-JIL channel is defined at present as a FIFO queue. Steps that use the channel either *write* to the end of the channel, *take* data from the front of the channel, or *read*, without removing, data from the front of the channel. The channel construct can be used as a vehicle to coordinate and synchronize steps executing in parallel. On the other hand, as the example below shows, it is restrictive to use only a

FIFO queue to model the channel's handling of data. A more complete understanding of what is needed for the clear and precise definition of models of scientific data processes is emerging from examples such as this.

Requisites: *Requisites* are optional and enable the checking of a specified condition either as a precondition for step execution or as a postcondition check to assure that the step execution has been completed acceptably. A downward arrowhead to the left of the step bar represents a prerequisite, and an upward arrowhead to the right of the step bar represents a post-requisite. If a requisite fails, an exception is triggered.

Exception Handling: A step in Little-JIL can define *exceptional conditions* when some aspects of the step's execution fail (*e.g.*, one of the step's requisites is violated). This violation triggers the execution of a matching exception handler associated with the parent of the step that throws the exception. An exception handler is represented as a step attached by an edge to an X on the right of the step bar (as shown in Figure 4). Of particular importance is the Little-JIL facility for specifying how execution proceeds after completion of the exception handler. Little-JIL currently supports four different ways that execution can proceed after the successful execution of an exception handler.

Scoping and Recursion: The parent step and all of its descendants represent a *scope*, specifying what data is considered local to that scope. Little-JIL also supports *recursive execution* of steps, which specifies the iterative application of a process step to specified inputs.

4 Using Little-JIL to define and execute a PDG and to create a DDG

The purpose of the Water Budget process is to provide estimates of the rate of change of water storage dS over various time intervals on the basis of time-ordered sequences of readings from sensors that measure various parameters. While this process may at first glance appear to be relatively straightforward, there are a number of aspects of the fully elaborated process that are challenging to define precisely. One aspect of this complexity is that the process is long-running because the time intervals can be as long as months or years. In addition, a number of contingencies can arise during execution. For example, malfunctioning sensors may produce large numbers of incorrect readings or may not produce readings at all over extended periods of time. Researchers continually develop new models to use to interpolate when sensor data is unavailable or considered to be unacceptable, and they continually attempt to validate these new models on previously collected datasets. This leads to a proliferation of derived datasets, each incorporating interpolated values that come from the applications of different combinations of models. Clearly all such datasets must be distinguished from each other, and the differences in the ways in which the various data items in each dataset were derived also must be carefully documented. One way to do this is to distinguish the different data items, and the datasets containing them, based upon differences in their derivation histories. This problem bears very strong similarities to problems in software configuration management, where, for example, executables can be created by loading different combinations of modules and

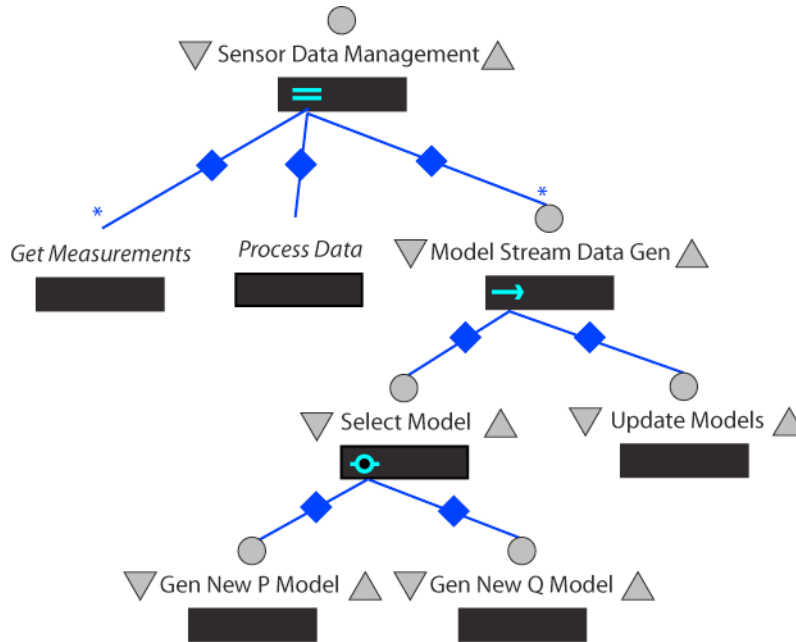
libraries, each of which might have been created by applying different compiler versions to different source code versions.

Below we describe selected aspects of a Little-JIL definition of the Water Budget process that illustrate some of the difficulties that arise in defining scientific processes and some of the useful semantic features that are desirable in a process definition language used to define the PDG in an analytic web. Although simplified for this presentation, the process defined in this section is representative of the process actually used at Harvard Forest. The process step definitions are illustrated here by step diagrams that have been created using the Little-JIL editor. In the interest of reducing clutter, the editor does not depict all the details and annotations of all step definitions unless specifically requested by the user. The user may obtain complete detailed information about any aspect of a process definition by moving the cursor over the appropriate icon. We also append to each figure a legend that contains additional explanatory information. Information about a step begins with the name of the step, followed by specifications of the input and output parameter types for that step, and when appropriate the channels that are used by the step and the types of data items that are carried by these channels. For some steps there is also a specification of exceptions that may be thrown by the step and/or exceptions that may be handled by the step. For edges, the information begins with a specification of the pair of steps that define the edge, followed by information about parameter flow between parent and child steps. Note that there may also be comments that informally describe the general purpose or approach of a step or edge. Much of the information being provided by these legends could be specified through extension of existing annotation schemes, such as EML. Thus, a process language such as Little-JIL can provide a useful and natural complement to existing approaches to the specification of scientific processes and the provenance of their datasets.

4.1 The Sensor Data Management step: use of decomposition and concurrency

We begin by describing the step **Sensor Data Management** shown in Figure 5. **Sensor Data Management**, the root of the Water Budget process, is a parallel step consisting of the execution of three substeps; one, **Get Measurements**, collects data from the sensors, another, **Model Stream Data Gen**, creates new predictive models for this data, and the third, **Process Data**, processes the data for publication. The collected data and created models are communicated to **Process Data** via channels that are declared in **Sensor Data Management** and are accessible to all of its substeps. The `sensorStream` channel is declared to be a FIFO channel, so that data is removed in the order in which it is put into the channel. The `modelStream` channel is a singleton channel, meaning that at most one item can be in the channel at a time.

The first two substeps are described in Sections 4.2 and 4.5 below. The third substep, **Model Stream Data Gen**, represents the generation of new models. The annotation indicates that this step is done by human experts, each of whom executes a different instantiation (as indicated by the Kleene * on the edge leading into this step) of this step, thereby allowing for the parallel and asynchronous generation of new models. The use of a channel allows for the possibility that new models can be dynamically placed into the



Step: Sensor data management

FIFO channel `SensorStream sensorStream = new SensorStream();`
 Singleton channel `ModelStream modelStream = new ModelStream();`

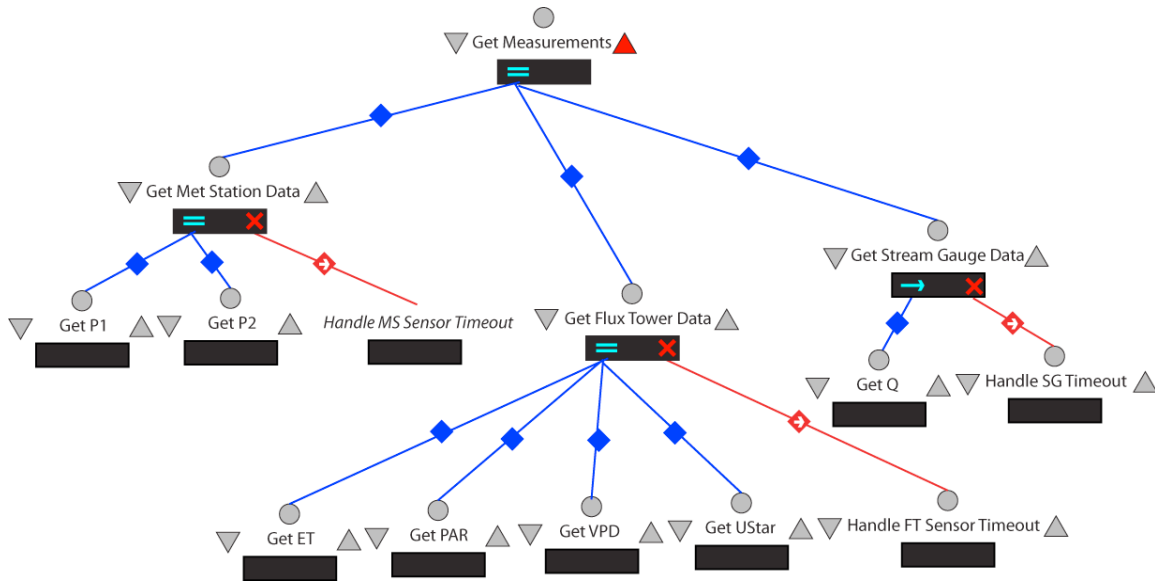
Step: Model Stream Data Gen

Agent: human expert

Figure 5: The Sensor Data Management step

channel at any time during the execution of this process. Doing so makes the model available for application to subsequent data items. Note also that there is a different step for the generation of each type of model. `Gen New P Model` is the step that is responsible for the generation of new P models, while `Gen New Q Model` is responsible for the generation of new Q models. The example shows only two such substeps for simplicity, but the complete process definition would require a substep for each model type.

Because each of these steps can be executed in parallel, each can proceed at its own pace. Sensor data streams in as it is generated and new models are derived as needed by human scientists, presumably at a far slower pace. The processing of the data streams (including application of the models as needed) is driven by interval timers and human curiosity. Little-JIL clearly represents this parallelism and indicates which models are applied to which data items as part of which datasets.



Step: Get Measurements
 Outputs: ETreading et; Qreading q; PrecReading p1, p2; PARreading par;
 VPDreading vpd; UStarReading uStar; // written to the channel sensorStream
 Output assertions: all output values are marked with flag set to measured

Step: Get Met Station Data
 Outputs: PrecReading p1, p2;
 Exception caught: "Sensor Down" with "Handle MS Sensor Timeout"

Step: Get P1
 Outputs: PrecReading p1;
 Exception raised: "Sensor Down"

Step: Get Flux Tower Data
 Outputs: ETreading et; PARreading par; VPDreading vpd; UStarReading uStar;
 Exception caught: "Sensor Down" with "Handle FT Sensor Timeout"

Step: Get Stream Gauge Data
 Outputs: Qreading q;
 Exception caught: "Sensor Down" with "Handle SG Timeout"

Figure 6: The Get Measurements step

4.2 The Get Measurements step: support for multiple data streams

The step Get Measurements, elaborated in Figure 6, reads and processes the values from the sensors and sends the results to the real time stream. This is done by having a different subprocess, Get Met Station Data, Get Flux Tower Data, and Get Stream Gauge Data, take responsibility for examining the three types of data sources. Note, that each of these subprocesses can execute independently and in parallel with the other two, and each may throw a different type of exception if difficulties arise with their sensors.

The **Get Met Station Data** step has two substeps, **Get P1** and **Get P2**, each of which is responsible for dealing with measurements coming from one of the two precipitation measurement gauges. Similarly, the **Get Flux Tower Data** step has four substeps, each of which is responsible for handling data items coming from each of the various sensors on the flux tower and **Get Stream Gauge Data** has one substep for handling surface discharge data.

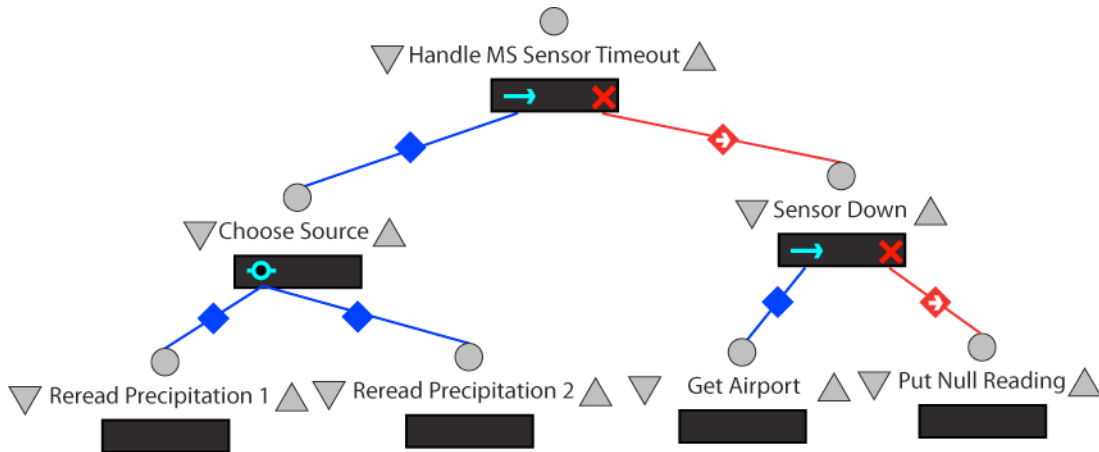
The substeps that access the individual data items also are responsible for annotating each data item with some provenance information. At present this consists of attaching very rudimentary provenance information, specifically a date/time stamp and a quality flag. Shortly we describe the quality flag that is currently used and indicate how a DDG improves significantly upon the weaknesses of the current quality flag.

4.3 The Handle MS Sensor Timeout step: use of exception management

Each substep of the **Get Measurements** step is responsible for acquiring data from its sensors. And, each is also responsible for defining the subprocess to be employed if exceptional conditions arise. Here we describe how **Get Met Station Data** does this exception handling. The other two substeps of the **Get Measurements** step, **Get Flux Tower Data** and **Get Stream Gauge Data**, carry out their responsibilities similarly to **Get Met Station Data**. The final output of the **Get Measurements** step is exactly one datum of each of the types **P1**, **P2**, **ET**, **Q**, **PAR**, **VPD**, and **UStar**, each accompanied by provenance annotations.

The **Get Met Station Data** step deals with situations where some, or all, of the expected data items do not arrive by throwing an exception. For example, the substep **Get P1** attempts to obtain a precipitation reading from meteorological sensor 1. If this access succeeds, then the value is passed as **P1**, annotated with the observation date and time and with the *measured* quality attribute. The **Get P1** step is also responsible, however, for determining when new **P1** data is not available and then, subsequently, for throwing an exception that is to be handled by **Handle MS Sensor Timeout** step, passing the identifier of this sensor (namely sensor #1) as a parameter.

In this context, the **Handle MS Sensor Timeout** step, shown in Figure 7, is specified as the handler for the **Sensor Timeout** exception that can be thrown in the **Get Met Station Data** step. The **Handle MS Sensor Timeout** begins its processing by choosing, based on the value of the input parameter, **Reread Precipitation 1** or **Reread Precipitation 2**. We do not show the details of either **Reread Precipitation** step here, but include these steps only to illustrate that it is possible for the developers of this process to decide that unusual (and presumably expensive) measures might be taken under these circumstances to attempt to extract the desired data directly from the sensor, rather than from the data stream. If this direct request to the sensor succeeds then output of this step is annotated with the date and time and quality flag. If however the selected **Reread Precipitation (1 or 2)** step fails (e.g. because the sensor is inoperative), a **Sensor Down** exception is thrown (inside of the **Handle MS Sensor Timeout** exception handler itself). The **Sensor Down** handler executes the **Get Airport** step, which attempts to obtain the desired reading by getting it from a local airport. Here again,



Step: Handle MS Sensor Timeout

Inputs: PrecipitationSourceID precSourceID;

Outputs: PrecipitationReading p;

Exceptions caught: “Sensor Down” with “Sensor Down Handler”

Step: Choose Source

Outputs: PrecipitationReadingFlagged measuredP;

Exceptions raised: “Sensor Down” (if selected alternative does not succeed)

Step: Sensor Down

Outputs: PrecipitationReadingFlagged measuredOrMissingP;

Exceptions caught: “Airport Data Read Failure” with “Put Null Reading”

Step: Get Airport

Outputs: PrecipitationReadingFlagged measuredOrMissingP;

Exceptions raised: “Airport Data Read Failure”

Step: Put Null Reading

Outputs: NullReadingFlagged nullReading;

Figure 7: The Handle MS Sensor Timeout exception handler step

nested exception handling capabilities are used to define the handling of still further failure. Thus, if the **Get Airport** step succeeds, then the resulting data item is annotated with the date and time, source, and the *measured* quality attribute. If the **Get Airport** step also fails (*e.g.* because the airport is also unable to provide the precipitation data), then it must also throw an exception, which would then be caught by the **Put Null Reading** step of the **Get Airport** exception handler, to produce as output a null value for P1 (or P2) and a quality attribute “missing”. A “missing” value will not take place unless all three exception-handling alternatives have been explored in sequence. Little-JIL’s exception-handling capabilities make the way in which these alternatives are explored relatively easy to follow. Certainly a variety of other sequences of exception handling could have been defined equally well.

4.4 Using the Little-JIL process definition to create the DDG

The above process results in the creation of two data items, P1 and P2, each of which might have been arrived at in a number of different ways. Specifically, each of the two measurements might be arrived at by any of the following sequences of process steps:

1. Pi arrives in a timely fashion and is recorded.
2. Pi does not arrive in time, a timeout exception is thrown, Reread Precipitation is executed, and Pi is obtained.
3. Pi does not arrive in time, a timeout exception is thrown, Reread Precipitation is executed, Pi still does not arrive, a Sensor Down exception is thrown, the Airport Read step is executed, and Pi arrives.
4. Pi does not arrive in time, a timeout exception is thrown, Reread Precipitation is executed, Pi still does not arrive, a Sensor Down exception is thrown, the Airport Read step is executed, Pi does not arrive, an Airport Data Read Failure exception is thrown, the Put Null Value step is executed, providing a null value for Pi.

The differences among these four possibilities are important to subsequent process steps. Accordingly the current process specifies that a quality flag, mentioned previously, be attached to Pi. This quality flag is currently a simple annotation, at this point having the following values:

missing - no measured value is available (e.g. because the sensor is down);

measured – the measured value is available, but it is unknown whether or not it is within prespecified bounds;

Unfortunately, the current process specifies that Pi is annotated as *measured* in all of the first three of the cases we have just enumerated, making it impossible to distinguish among them. In the fourth case, the data item is annotated as “missing”; again leaving out all of the details describing what alternatives had been tried. Clearly this simple annotation flag is inadequate. Instead, we note that each of the four different sequences of process steps can be thought of as a different trace through the process, illustrating the importance of annotating each value with process provenance information, as provided in the DDG.

Examples of the DDGs that represent cases #3 and #4 are provided in Figure 8. Note that the boxes in this figure represent actual data instances, namely the actual data values that are bound at execution time to the type specifications in the process definition. Thus, for example, one of the boxes at the top of Figure 8 is annotated by `sensor1 null @"try1 time"`, indicating that this box represents the actual (null) value that was delivered at the specific time, “try 1 time”. The adjacent box likewise represents the specific (null) value that was delivered at the specific time “try 2 time”.

The ovals in Figure 8 represent the process step instances that were the actual producers and consumers of the actual data items. Thus, there is an oval labeled `Get Airport` that indicates that an instance of the `Get Airport` step was used to generate the data item in the box shown below this oval. This step instance represents the instance of `Get Airport` that was invoked as the process’s response to the two null readings. Two arrows from this

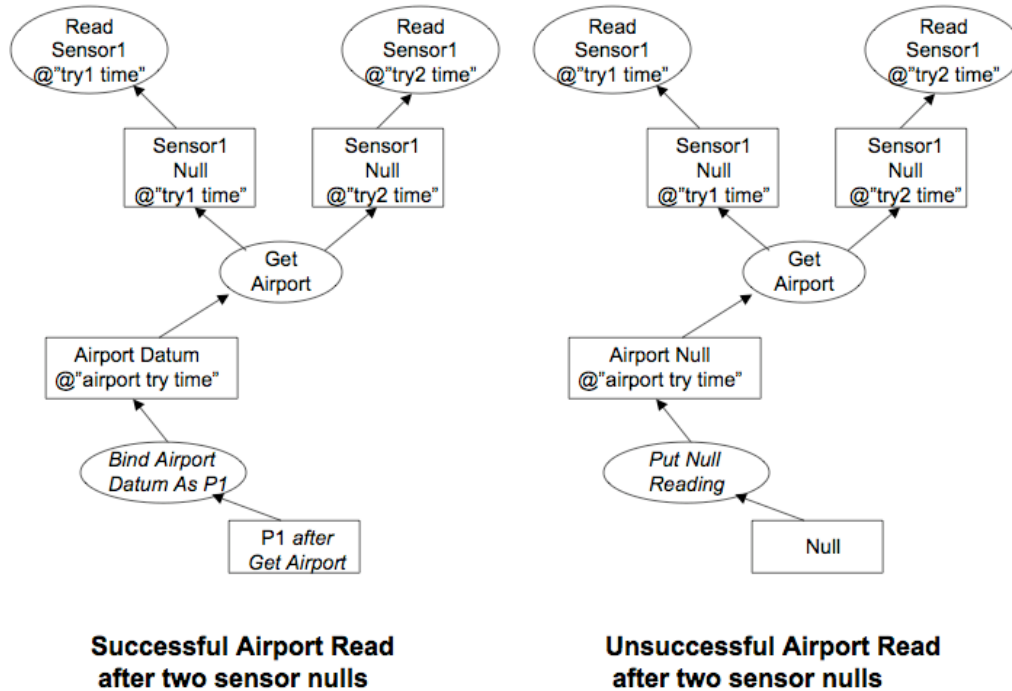
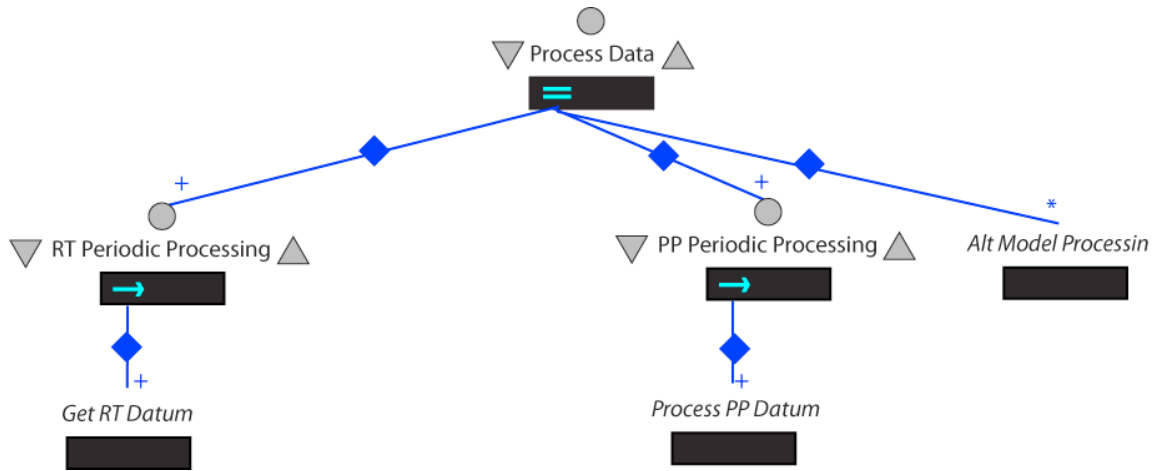


Figure 8: The DDG of the original acquisition of P data

oval connect it to two boxes, representing the fact that the values represented by these boxes were used as inputs to the step represented by that oval. In this case, the use that is made of these data items is simply to note that they are both Null, causing the *Get Airport* step to be executed to produce this output. Other ovals may make more substantive use of their inputs in generating their outputs. Thus, for example, in the left-hand DDG of Figure 8, the result of the execution of the instance of the *Get Airport* step is an actual value, annotated with date and time information, which is taken as the final value of P_i . In this case, no actual step is used to generate that value, and instead the DDG indicates that the value is produced as a consequence of the parameter binding operation that occurs as an integral part of the execution of every step. The fact that this oval does not represent an actual step is indicated by the use of italics in its annotation. In the other case, a null reading is obtained, and a null value is then the final value of P_1 . Clearly such DDGs provide much more useful information about the provenance of the resulting value of P_1 than a mere annotation, *measured* or *missing*.

Although the structures shown in Figure 8 seem to be large annotations, they can be represented efficiently as sequences of pointers to process definition steps, parameter values, and the actual data item instances. Further, these pointers can be generated quite easily as side effects of the actual execution of the process that is represented by the PDG shown in the preceding figures. In many cases the DDGs can be generated at a modest incremental cost of the execution of the process and, with the appropriate encodings, can be represented relatively efficiently in those cases where the intermediate datasets are already being stored.



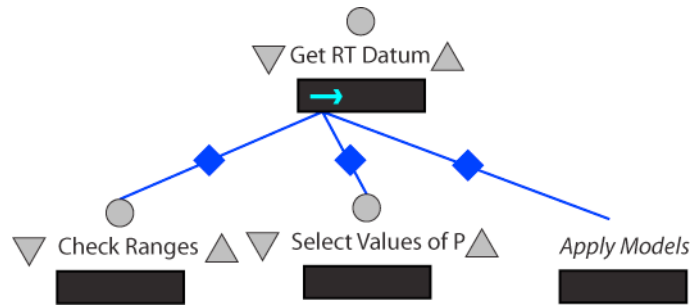
Step: RT Periodic Processing
 Comments: the channels sensorStream and modelStream are accessed from an enclosing scope
 Outputs: seqofdS <dS>;

Figure 9: The Process Data step

4.5 The Process Data step: alternative approaches to processing the data

The second subprocess, **Process Data**, in the root **Sensor Data Management** process, shown in Figure 5, is described and elaborated in Figure 9. This step takes input from sensors and investigators, as described above, and produces a time-ordered sequence denoted <dS>, where each dS estimates change in water storage over the previous time interval. Note that <dS> (also dS) is a type and that <dS> instances are generated by each of three different subprocesses of **Process Data**, namely **RT Periodic Processing**, **PP Periodic Processing**, and **Alt Model Processing**. **Process Data** indicates that these three subprocesses can be executed in parallel, but whether the actual executions of these steps overlap in time will depend on decisions of the agents bound to perform these substeps and, most importantly, upon the availability of the input data.

At a high-level, **RT Periodic Processing** produces a <dS> dataset every 24 hours, based on the data that was collected over the proceeding 24 hour period. **PP Periodic Processing** also produces a <dS> dataset every 24 hours, but it is for a 24 hour period 30 days before. This subprocess has the advantage of doing interpolation that uses data obtained both prior and subsequently to the data in need of interpolation. **Alt Model Processing** allows investigators to experiment with alternative models and time periods. Here we describe each of these subprocesses in turn. Note that in the DDG, it is important to know which subprocess was responsible for generating each of the different <dS> instances and their different individual data items. As illustrated above, the process definitions offer a natural basis for generating DDGs that document these differences.



<p>Step: Get RT Datum Inputs: SetOfET {et}; SetOfQ {q}; SetOfP {p}; SetOfPAR {par}; SetOfVPD {vpd}; SetOfUStar {uStar} ET et; Q q; P p; PAR par; VPD vpd; UStar uStar; // taken from the channel sensorStream Outputs: SetOfETout {et}; SetOfQout {q}; SetOfPout {p}; SetOfPARout {par}; SetOfVPDout {vpd}; SetOfUStarOut {uStar};</p> <p>Step: Select Values of P Inputs: PrecReading p1, p2; Outputs: PrecReading p;</p>
--

Figure 10: The Get RT Datum step

4.6 The RT Periodic Processing step: handling of real-time streaming data

The RT Periodic Processing step is connected to its parent by an edge annotated with a Kleene +, prescribing that the Little-JIL interpreter will keep instantiating a new instance of this child step indefinitely (the decision to terminate iteration is made by the agent bound to the iterated child step). The agent for this step is an interval timer that initiates step execution every 24 hours. The **Get RT Datum** step specifies the actual processing of the real-time data and the rate of execution is determined by the availability of data from the sensors and the resources required for the processing of that data.

Observe that good software engineering concepts have contributed to the clarity of this process definition: the Kleene + annotation provides a clear specification of the desired iteration, an agent specification defines periodicity, and hierarchical decomposition is used to hide details of data gathering until there is interest in seeing them.

4.7 The Get RT Datum step: dealing with individual data items

The **Get RT Datum** step, shown in Figure 10, defines the heart of the Water Budget process, specifying the way in which the real-time data items are subjected to cleaning. This step retrieves the readings collected from the sensors as they become available, filters the readings to ensure each reading is within predefined bounds, tags each reading to indicate how it was obtained, and appends readings to corresponding time-ordered sequences. This step is defined as the sequential execution of its substeps and results in

an instance of dS, which has data items of type ET, Q, P, PAR, VPD and UStar as its components. The set dS is then appended to the time ordered sequence <dS> as each new set of data items is processed.

The first and second substeps of **Get RT Datum** are relatively straightforward, but each adds further useful provenance information to data items. The **Check Ranges** step takes as input a measurement artifact that was output by a sensor and determines whether it is beyond preset bounds. While the details of the range checking are not particularly difficult or interesting as process features, the specification of the details of the check are important provenance information that should be associated with the data item.

Figure 11 shows the DDG that results from the application of such a filter. The box labeled **Filter Range Data** represents the specific values used in applying the filtering to the data item, and the oval labeled **Check Ranges** represents the application of that step to the value represented by the box labeled **P1**, using the values specified by **Filter Range Data**. This results in the creation of a new instance that is bound to the variable **P1**. Note that this DDG documents the exact quality criterion that was applied in checking the plausibility of the resulting data item. If an investigator wanted to experiment with a more stringent, or more relaxed, filtering criterion, it could be reapplied to the previous data items available from this provenance documentation, thus potentially allowing for the application and evaluation of alternative scientific models

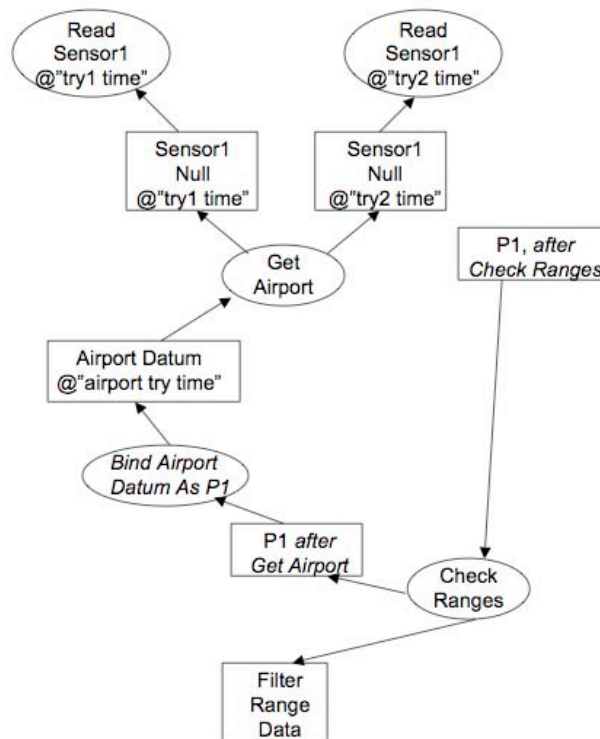


Figure 11: The DDG representing the application of range filtering

using previously gathered (and alternatively processed) models.

Similarly, the **Select Values of P** step selects one of the two values P1 and P2 and assigns it as its output artifact, namely P. This example suggests the importance of using a PDG language that supports the specification of the agent. In this case, this selection is done by specifying “Expert P Selector” as the type of the agent that is to be assigned to this step. Note that this specification does not preclude the possibility that the agent might be a human or an automated agent having expertise in selecting the preferred value. The choice of which agent to bind to this step is made by the scientist or it could be relegated to an automated resource management system, which would presumably use other execution state data as the basis for this decision. Little-JIL would document this choice in the DDG, as illustrated in Figure 12. Note that this figure shows the case where Sensor 2 had succeeded in delivering a datum on the first request and where the human agent decided to accept this value as the final value of P. The figure contains an italicized annotation indicating that the chosen value was the one provided by P2, but the actual DDG need only contain the value. The fact that this value had been provided by P2 would be inferable by inspecting the two input values to the **Select Values of P** step, both of which would be available as nodes in the DDG.

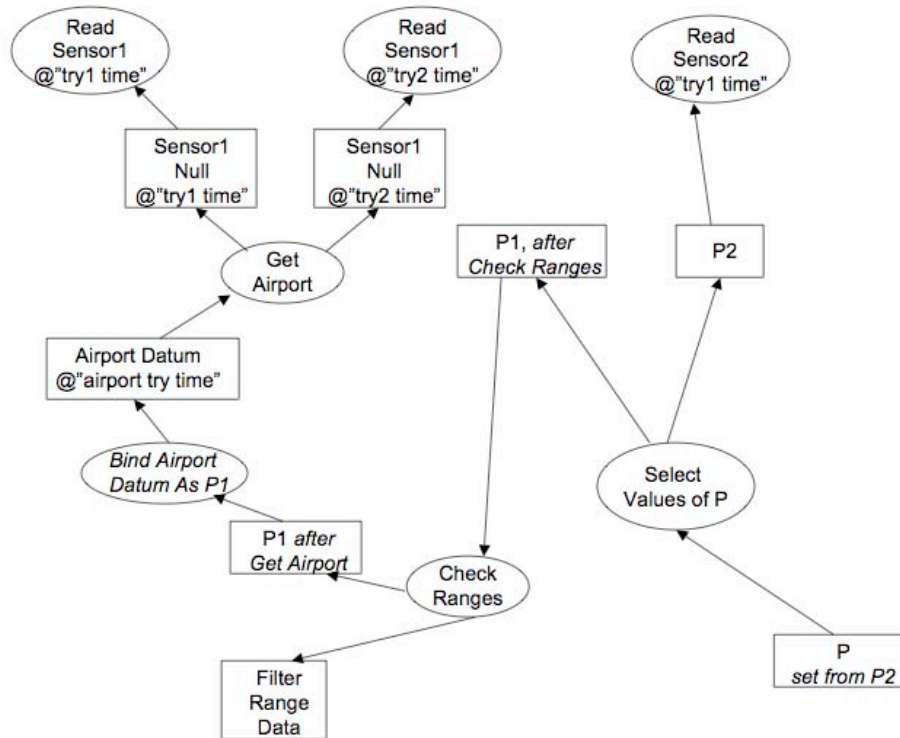


Figure 12: The DDG representing the selection of P as a choice between P1 and P2

The comparison of measured and modeled values (like the comparison of redundant sensors) also provides a form of real-time quality control and may provide an early warning of a sensor problem. The essence of the **Apply Models** step, in the **Get RT Datum** step of Figure 10, is to effect the use of an empirical formula, called a *model*, to fill in data readings in such cases. As noted above, these models have been created by scientists for the purpose of interpolating values in place of those that are missing or suspect. The models are designed to provide replacements for such values, by using other data items and formulas that are believed to accurately capture the relations of these data items to the missing or suspect data. There are generally a number of alternative models that can be applied in these circumstances, as these models are the subjects of considerable research. The model selected and applied to a data item is important provenance information that should be attached to that item. We indicate the way in which this additional provenance information is generated as part of the execution of the **Apply Models** step of the Water Budget example.

The models themselves are often structured into sets of models of different types, for example, different types of regression models (*e.g.* linear or second-order) or different probabilistic models. The suitability of a model is then evaluated by a human, perhaps supported by various statistical tools. Eventually a model is chosen and then used to

create a data value, which is substituted for the original data value. As this sort of synergy between humans and tools in the evaluation and application of models seems to be at the heart of many scientific activities, it is important to detail it here and to demonstrate the challenges of representing such a process and documenting the human choices made.

4.8 The Apply Models step: use of abstraction

The purpose of the **Apply Models** step is to replace the data readings obtained from the sensors when the quality attribute of the reading is either missing or determined to be out of range (by the **Check Ranges** step). Although models can be applied to any number of data items, Figure 13 shows the application of models to only two types of data, **P** and **Q**. The definition of this step illustrates the value of abstraction in defining processes, since here the **Apply A Model** step is simply instantiated once for each type of data requiring the application of a model. Each **Apply A Model** substep has a prerequisite (not shown here) whose purpose is to examine the incoming data item and determine whether it has a quality attribute of “missing” or “out of range”. Either attribute causes the prerequisite to be satisfied and the main body of the step to be executed.

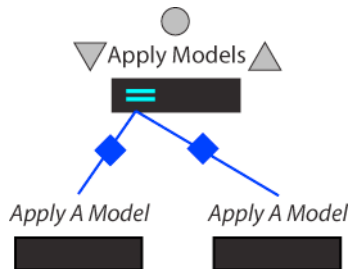
Also note that the output of this step is a triple of items, namely the original **P** value (**P.original**), a replacement **P** value that has been produced by the model that has been selected (**P.modeled**), and an identifier (**PModelID**) indicating which model was used to generate the replacement value. In fact, the first and third elements of this triple are redundant with information contained in the DDG, which provides more precise detail about the way in which **P** was derived. The triple specified here is included as an indication that such annotations might be derived and included to help users by providing such information as documentation. Thus, this step provides the basis for distinguishing among different data items and different data streams that have been produced by different applications of different models by different scientists at different times. Note, in addition, that all three of the steps are carried out by agents who are required to be of type “HumanExpert”. Each of the two substeps is defined to have as “local data” a set of readings that can be used by a selected model in computing the model-generated output of this step. Thus it is important that the language used to define the PDG for this process allows for the possibility of such “local data”. Further the example illustrates the value of incorporating scoping semantics into the language in order to support specification of how such “local data” can be collected and held in the defined local scope, presumably in any way that the model and its agent decides.

4.9 The Apply A Model step: synergizing the efforts of humans and tools, and representing a history of human decision-making

The purpose of this step is to consider iteratively each of a set of models selected from the pool of models available, and then to select and apply the model that appears to be most effective in replacing the datum that has been identified for replacement. The step is comprised of two substeps, executed in sequence. The step also incorporates an exception handler to deal with the case in which no model is selected and applied.

The first substep of Apply A Model, Eval Models is shown in Figure 14 and is the iterative consideration of the set of candidate models. Note that the agent for this step is identified as being of type “Model selector”. The Resource Repository used to support the execution of this process is responsible for identifying the agent instance to be bound to this step. It is conceivable that this agent might be a human or might be an automated system that is capable of performing preliminary evaluations of models. This process does not mandate either, but allows the runtime agent binding capability to make that choice from among the alternatives that the process designer has established.

Eval Models iteratively evaluates each model in the collection of currently available models by first checking the model’s applicability (represented by the pre-requisite) and then consulting an expert to evaluate the model and the results that it gives. Note that Eval Selected Model is to be executed by an agent of type “human expert”, indicating that human expertise is required. We do not show the details of this step, but it might entail the use of statistical tools to study various aspects of how well the data suggested by the model seems to fit with the other data being collected. In this case, although a



Step: Apply Models

Outputs: ModeledPtriple (P.original, P.modeled, PModelID);
 ModeledQtriple (Q.original, Q.modeled, QModelID);

Step: Apply A Model

Inputs: PrecReading p;
 Outputs: ModeledPtriple (P.original, P.modeled, PModelID);
 Comments: the modelStream from Sensor Data Management is used in this step

Step: Apply A Model

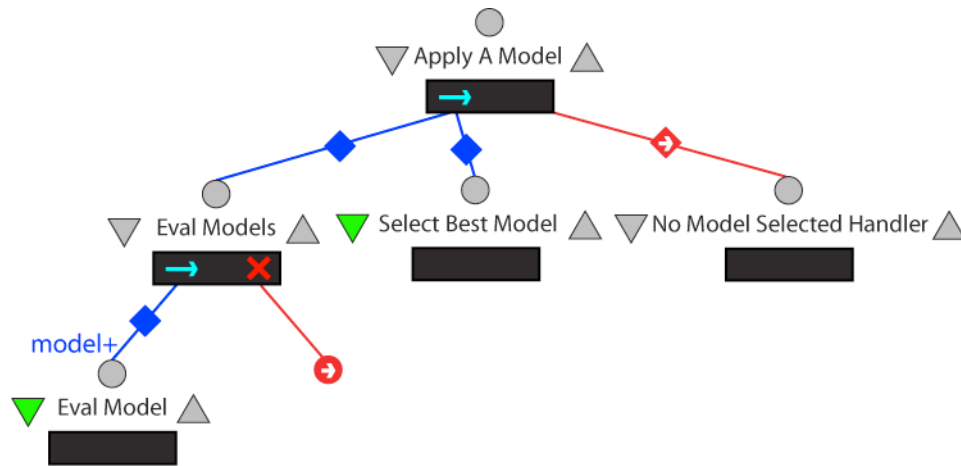
Inputs: Qreading q;
 Outputs: ModeledQtriple (Q.original, Q.modeled, QModelID);
 Comments: the channel modelStream from Sensor Data Management is used in this step

Figure 13: The Apply Models step

human expert executes **Eval Model**, some of its substeps might well be executed by automated agents (*e.g.* statistical tools).

The edge leading into **Eval Model** has a Kleene + annotation, indicating that this step is to be instantiated as many times as there are models. Combined with the pre-requisite, the net effect of this is to instantiate a step for the consideration of only those models that the step’s agent deems worthwhile. It is important to note that **Eval Models** also has access to a cache of data values collected and stored locally to this step that can be used to help in deciding the suitability of each candidate model.

The output of the **Eval Models** step is a set of pairs, where each pair consists of a model and the output that it produces. Once this stream of sets of ordered pairs has been created



Step: Apply a Model

Outputs: ModeledDatumTriple modeledDatum =
 new ModeledDatumTriple (sensorDatum.original, sensorDatum.modeled, modelID);
 Exceptions caught: “No Model Selected” with “No Model Selected Handler”

Step: Eval Models

Inputs: SensorReading sensorDatum; Cache readingsCache;
 Outputs: SetOfModelIDandDatumPairs {datumPairs}; //
 where datumPairs=new DatumPairTuple(ModelID, SensorDatum.modeled);

Step: Select Best Model

Agent: human expert
 Inputs: Cache readingsCache; SensorReading sensorDatum.original;
 SetOfModelIDandDatumPairs datumPairs //
 where datumPairs=new DatumPairTuple(ModelID, SensorDatum.modeled);
 Outputs: ModeledDatumTriple modeledDatumTriple =
 new ModeledDatumTriple(sensorDatum.original, sensorDatum.modeled, modelID);
 Exceptions raised: No Model Selected

Step: Eval Model

Inputs: SensorReading sensorDatum; Cache readingsCache; ModelType model;
 Outputs: DatumPairTuple = new DatumPairTuple(ModelID, SensorDatum.modeled);

Figure 14: The Apply a Model step

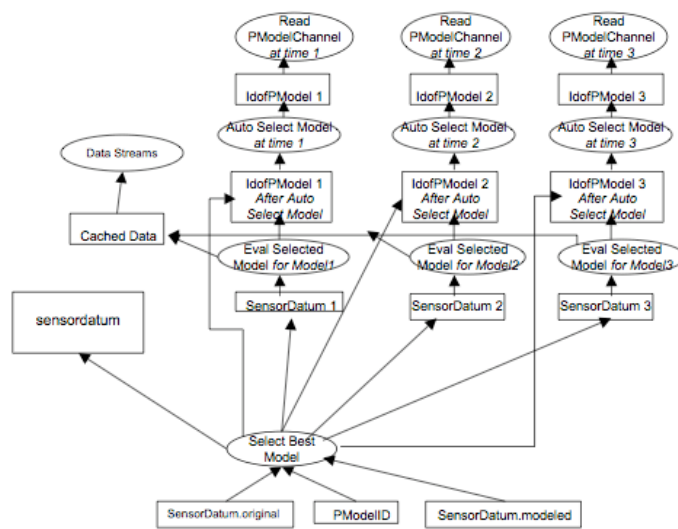


Figure 15: The DDG representing the selection and application of a model to generate a modeled value for P

the next substep of Apply A Model, **Select Best Model**, considers all of these candidate models and selects the one deemed best. Note that a “human expert” is required to perform the **Select Best Model** step, requiring a scientist to use expert judgment to select the model that seems to produce the best value. In this example, this value then becomes a third component of the ordered triple, namely (**sensordatum.original**, **sensordatum.modeled**, Model ID), that is the final output of the Apply A Model step. Again, the first and third elements of this triple are redundant with the more precise derivation information contained in the DDG and should be regarded as documentation intended to be of value to the user.

Figure 15 provides an example of a DDG that might be generated from one possible execution of this step. In this case, it is assumed that three different models (shown as PModel1, PModel2, and PModel3) are selected from the channel containing the accesses to all possible P Models. Each of the three is subsequently evaluated by the **Eval Selected Model** step, using **Cached Data**, the set of data instances that has been cached by the process, thereby producing three instances of **Sensor Datum** (namely **Sensor Datum 1**, **Sensor Datum 2**, and **Sensor Datum 3**). These three instances of **Sensor Datum**, as well as the three models that produced them are then shown to be inputs to the **Select Best Model** step, which then produces as its output the ordered triple (**sensordatum.original**, **sensordatum.modeled**, Model ID), where the original sensor datum, which was an input parameter to this step, is now referred to as **sensordatum.original**.

Note that the two DDGs described in Figure 12 and Figure 15 are both necessary to provide the complete provenance of **SensorDatum.modeled**. The DDG in Figure 15 represents the evolution of **SensorDatum.modeled** from the datum that was selected as

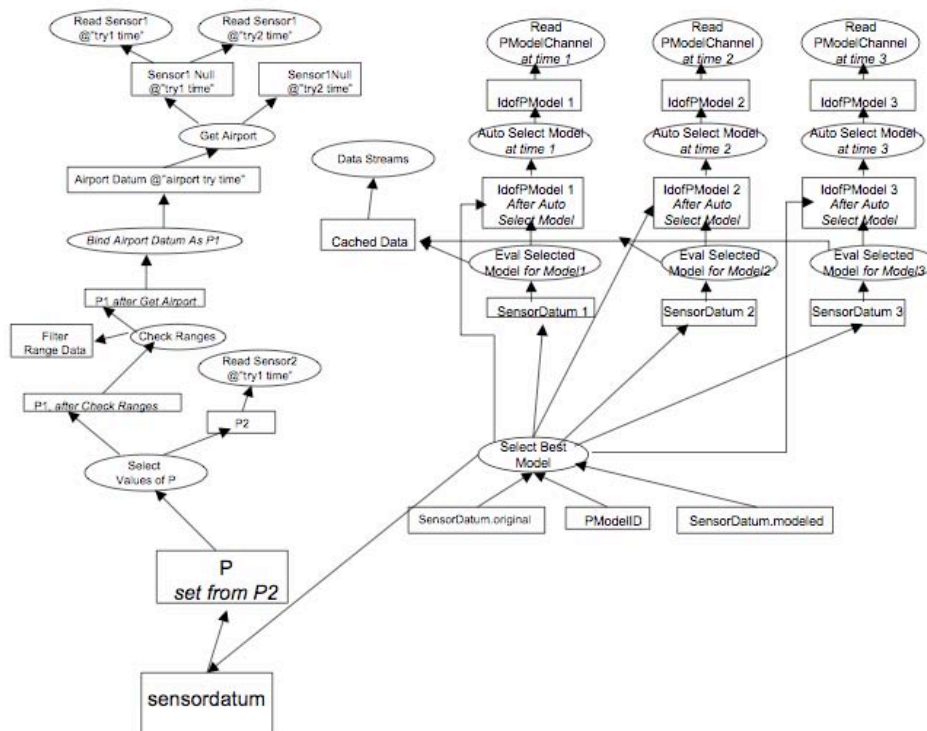


Figure 16: The DDG formed by combining the DDGs of Figure 12 and Figure 15

the result of the process depicted in Figure 14. This combined DDG is depicted in Figure 16. This DDG accurately documents that it is this **P** value that is bound as the actual datum taken as the subject for consideration for replacement by the alternative models.

After **Apply Models** has executed, the new values that have been obtained from all of the sensors are appended to the end of the data streams that are being accumulated by the **RT Periodic Processing** step (Figure 9). As described above, some of these newly appended data items are modeled, rather than measured. Here again each datum in the data stream is shown to be a packet that contains information that is redundant with the DDG. This redundant information should be thought of as documentation intended to be helpful to the user. In this example, the packet consists of **SensorDatum.original** along with appropriate information about its provenance, such as the date and time of the original measurement, the **sensordatum.modeled**, and the model used to generate **sensordatum.modeled** along with the date and time that model was applied.

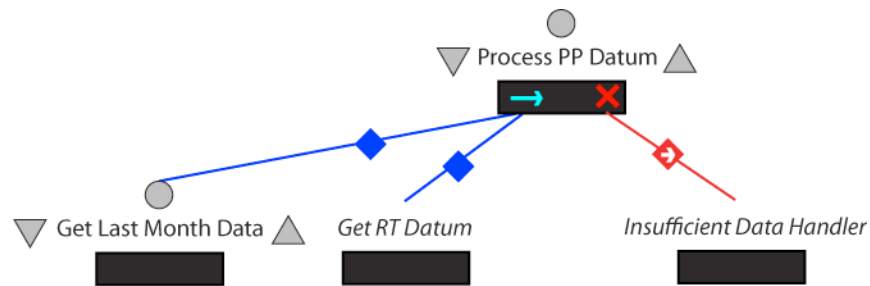
This stream of ordered triples is the real-time output of this process and is presumably made available via the Internet in real time in addition to being collected up into the 24-hour datasets. We assume that such data streams will also be labeled by a unique process identifier and the DDG representing the precise trace through the process, as well as with further annotations, indicating the place, dates, and times at which the data streams were collected. As this information consists largely of pointers into the DDG, most of the cost

of carrying this provenance information lies in the cost of the DDG itself, which should be relatively efficiently represented as a collection of pointers.

4.10 The PP Periodic Processing and Alt Model Processing steps: abstraction supporting process reuse

Recall from the Process Data step shown in Figure 9 that the real-time data stream is not the only data stream that is produced as the product of this process. The Water Budget process calls for automatic post-processing of newly collected data streams exactly 30 days later. As noted above, the correction of missing or out of bounds data during real-time processing must rely only upon retrospective data, and we have shown that this can be done by allowing for the caching of such data locally in the Apply A Model step. But post-processing after 30 days allows data gathered both before and after the selected data item to be taken into consideration. This post processing is done in the PP Periodic Processing substep of Process Data. In addition, consideration of further models is made perhaps every few weeks or months, resulting in the creation of new models and their application to various historical data streams. The way this is done is defined in the Alt Model Processing substep of Process Data. Note that the Process Data definition indicates that PP Periodic Processing is invoked every 24 hours, using an interval timer mechanism that works analogously to the way in which a timer is used to trigger processing in the RT Periodic Processing step.

The Process PP Datum step is defined in Figure 17 and consists of the reinstatement of the step structure that supported the Process RT Datum step, described previously. In particular, Process PP Datum consists of a first substep, Get Last Month Data, followed sequentially by Get RT Datum. The purpose of Get Last Month Data is to identify the data stream produced exactly 30 days earlier and to output the information



<p>Step: Process PP Datum Outputs: SetOfET {et}; SetOfQ {q}; SetOfP {p}; SetOfPAR {par}; SetOfVPD {vpd}; SetOfUStar {uStar}</p> <p>Step: Get Last Month Data FIFO Channel SensorStream sensorStream = new SensorStream(); // to be populated with "last month" data items</p> <p>Step: Get RT Datum Comments: "Augmented model set" is a set of models that also includes models that can benefit from prospective data, now available, as well as retrospective data</p>
--

Figure 17: The Process PP Datum step

italics), which is now made available for possible future consideration. Note that subsequent comparisons of `SensorDatum.modeled` with the results of other models would be represented as successive elaborations of the DDG shown in Figure 18.

To be completely fair, the reused subprocess steps of `Get RT Datum` must be designed to be applicable in both the `RT Periodic Processing` and `PP Periodic Processing` contexts. Thus, for example, the data measurements that are cached for use in the `Apply A Model` steps of `Apply Models` must include prospective data (acquired by some sort of lookahead) in the `PP Periodic Processing` context, while they will have only retrospective data in the `RT Periodic Processing` context.

The reuse of modular capabilities is an inherent abstraction capability in programming languages that usually requires careful design choices. Thus, it should not be a surprise that the reuse of process definitions in a process language also requires some care.

Finally, note that the `Alt Model Processing` step is defined in a manner that is analogous to the definition of the `PP Periodic Processing` step. This step defines the way in which scientists can recall historical datasets and apply new models to them, thus offering new ways to replace missing or out of bounds data. The step is to be executed at any time that a scientist wishes to reanalyze a particular dataset. In all other ways, however, this process is virtually identical to the `PP Periodic Processing` step. Thus, in particular, the first substep of `Alt Model Processing` is a step whose only role is to identify historical datasets to be reanalyzed and to pass the information needed to access them as handles to channels. The sequentially executed next step is again `Get RT Datum`.

5 Evaluation

The Water Budget example has provided a vehicle for demonstrating how an analytic web provides provenance information needed by both producers and consumers of scientific datasets. More specifically this example has demonstrated that a broad range of semantic capabilities, including hierarchical decomposition, abstraction, concurrency, exception management, and complex data handling facilities, are needed to define modern scientific processes. We now present a more detailed analysis of what the need for these capabilities tells us about desiderata for the semantic features of languages that are to be used as the basis for defining an analytic web's PDG.

5.1 Strengths and weaknesses of Little-JIL

Little-JIL is not simply a vehicle for supporting hierarchical decomposition (as is the case in many other process languages), but is better thought of as a vehicle for implementing abstraction. The difference is that a Little-JIL step is accurately thought of as the definition of an abstract concept, capable of being made a concrete specification by its bindings to concrete artifacts and placement in a specific execution context. A step defines a scope, and thus establishes a context. One key mechanism for context definition is the binding of artifacts as a step's inputs and outputs. This capability is tantamount to a capability for passing arguments to a procedure. By varying the argument stream to and

from a step, the step is made to perform somewhat differently in different contexts. Steps also provide different contexts by providing different exception handling capabilities. Every step may define a set of handlers for the various types of exceptions that may be raised in its scope. Different instantiations of a step may offer different exception handlers, thereby establishing different execution contexts.

The Water Budget example made interesting use of Little-JIL's facilities for abstraction, for example, by its reuse of the **Get RT Datum** step. Reuse of this step emphasized the strong similarities in the ways in which **RT Periodic Processing**, **PP Periodic Processing** and **Alt Model Processing** perform their work. This reuse shortens the process definitions and clarifies the reader's (e.g. the consumer of a analytic web's PDG) understanding of these processes. Another application of this concept was seen in the **Apply Models** step, which consisted of two different instances of the **Apply A Model** step. This step definition emphasized the iterative nature of the step, yet left little doubt about the differences between the two invocations of its substeps, namely the differences in their arguments.

Some of the complexity in the Water Budget process is attributable to the way in which different activities occur in parallel. Data streams from various sensors are processed in parallel, and the data must be processed in real-time as data are gathered and transmitted concurrently. Simultaneously, the much slower activity of generating new models and evaluating both new and existing models occurs. Little-JIL's parallel step is effective in defining what activities are executed in parallel with each other. Channels in Little-JIL are effective at defining data streams between steps that were distant from each other in the architecture of the Water Budget process. Channels also supported synchronization. For example, the parallelism defined by the top step in **Sensor Data Management** clearly depicted the way in which sensor data was generated and processed.

Little-JIL at present has a relatively simple communication model. Channels are limited to first-in/first out queues and parameters are passed by copy. The Water Budget process however requires more than this. For example, the **Eval Model** step accesses a more or less static collection of models but is required to retrieve a new copy at each access in case any of the existing models have been modified or new models have been added. Transaction-like semantics might better support the implementation of a step of this kind. Transactions would also support the manipulation of data at different levels of granularity. For example **Get RT Datum** produces a single result, but is used by **PP Periodic Processing** to produce a dataset containing all of the results from a 24-hour period. Transaction-like semantics could permit **Get RT Datum** to release individual results, but allow **PP Periodic Processing** to control the visibility until a complete dataset has been constructed.

The exception management facilities in Little-JIL enabled us to define features of the Water Budget process that contributed to its reliability and robustness. Thus these features seem important to include in a language used to define PDGs. For example, the **Get Met Station Data** step indicated both how to identify and how to respond to the lack of needed sensor readings. Requisites are particularly clear devices for showing where missing data can be detected, and exception handlers (such as **Handle MS Sensor Timeout**) were placed to clearly indicate where responsibility for responding to such

contingencies was located. This example also showed the importance of dealing with exceptions that occur during the handling of exceptions themselves. The nested exception handling in the *Get Met Station Data* step provided an example of this. Such situations also emphasize the importance of providing facilities for specifying how to continue execution upon completion of exception handling.

The example also showed the importance of supporting the late-binding of resources to steps to permit flexible reactions to contingencies (*e.g.* by the run-time selection of agents to retry an execution of a step). For example, in the elaboration of the *Handle MS Sensor Timeout* step, the process defined the need to execute a step to obtain a data value. The required agent was specified only as an agent having the capability to provide a precipitation reading. The choice of the agent was left to a Resource Manager having a repository of agents, some of which offered this capability. The indicated facility for specifying a needed capability, rather than a specific agent, enables the late-binding of any of a number of possible agents to this step. It enables a separate Resource Manager to keep track of the agents that are able to deal with a request at any given time, and thus supports the real-time selection of one that is actually able to satisfy the request.

In summary, we were able to use the Little-JIL language to support the specification of the PDG for a large and complex scientific process. Such processes are needed to handle and process contemporary scientific data. While Little-JIL offered many specification language features that seemed quite useful and effective, some important deficiencies were also noted. All of this has created a very useful picture of what semantic features seem to be needed in a language that can support the definition of the PDG in an analytic web. There is a striking similarity between these needs of the scientific community and what is generally provided by modern programming languages. This supports our intuition that scientific processes bear some strong similarities to computer software, and thus the challenges of defining them have strong parallels with the challenges of programming complex software systems. Thus, it is not surprising to find that a process language needs to incorporate the salient control features of modern programming languages. It is, moreover, not surprising to also find that modern capabilities for dealing with data, such typing mechanisms, also seem important to the precise specification of processes. Indeed, another weakness of Little-JIL as a language for specifying processes seems to be its relatively weak support for defining data objects,

5.2 DDG Evaluation

The example of the Water Budget process also showed how DDGs can be built incrementally as the execution of a PDG proceeds and can be defined as traces through the PDG. Dataset Derivation Graphs grow as DAGs, increasing in depth as PDG execution proceeds; iteration of processing steps is manifest as additional levels in the DDG DAG (*e.g.*, Figure 18 is an elaboration of the leaves of Figure 16, which in turn elaborates DDGs shown in earlier figures). Each iteration of the steps of a PDG creates a new scope, and such scopes are root nodes in successive DDGs. This emphasizes the role of these steps in establishing scopes, and the DDG clearly illustrates this role.

DDGs that are derived from the executions of lengthy processes seem large and cumbersome. But it is the pictorial depiction of an entire DDG that is large. Their internal representations are typical tree-like structures that are amenable to terse internal representation. In addition, the depiction of the entire DDG is not likely to be of interest in most cases; elided versions would probably suffice in many cases. Tools for allowing viewers of DDGs to tailor their views through devices such as elision seem necessary. A number of web-based viewing tools for hierarchical viewing of this kind already exist.

Clear depiction of the features of the DDG of greatest interest to dataset producers and consumers may prove to be a challenge. Long and complex process executions will yield DDGs whose depictions are indeed large and potentially confusing. Figure 18 may seem daunting, certainly at first. But DDGs will be produced automatically from executing PDGs and have the virtues of being precise and accurate. We do worry, however, about their clarity, and thus DDG presentation is an important area that requires future research.

6 Related Work

There are numerous other scientific workflow projects, many of which have been presented at meetings such as [Cooper, et.al. 2006; SAG 2004; SSDBM 2004]. Most of these projects (e.g. Kepler [Ludäscher, et.al. 2006; Altintas, et.al. 2005; Altintas, et. al 2004], Taverna [Wolstencroft, et. al. 2005; Oinn, et.al. 2004], and JOpera [Pautasso, et.al. 2005; Heinis, et.al 2006]) base their specification of process flow upon the use of various kinds of data flow graphs (DFGs; e.g., Fig. 1). Indeed, Kepler, perhaps the most noteworthy of these projects, is based upon Ptolemy II [Edwards, et.al. 2003; Baldwin, et.al. 2004; Ptolemy URL] which uses a powerful and flexible DFG structure to specify how datasets can be moved between processing capabilities. Kepler integrates a broad range of support tools that help with such key activities as specification, execution, and visualization of scientific data processes. It seems particularly effective in supporting the processing of streaming data, such as data produced by sensors and intended for real-time processing. Chimera [Foster, et.al. 2003; Foster, et.al. 2002; Deelman, et. al. 2003] was one of the earliest scientific workflow systems. It emphasized the use of pictorial visualizations to represent scientific processes. Chimera pictorial representations depicted a form of a DFG. Taverna [Wolstencroft, et. al. 2005; Oinn, et.al. 2004] is a more recent system that seems to focus on supporting the integration of web services, particularly for the creation of bioinformatics applications. Taverna's integration mechanism is a workflow notation that is also based upon a DFG formalism. More recently JOpera [Pautasso, et.al. 2005; Heinis, et.al 2006] has suggested the use of XML to specify scientific workflows as plugins that could be integrated using Eclipse. JOpera workflows also are based upon the use of a DFG formalism to represent scientific processes. Teuta [Fahringer, et.al. 2005a; Fahringer, et. al. 2005b] represents scientific processes through UML diagrams that offer some features, such as limited forms of concurrency, that go beyond the semantic features of a basic DFG.

We believe that the reliance of all of these other scientific workflow systems upon the DFG as their basis for the specification of processes is a major drawback. The DFG-based systems described above, for example, make it very difficult to support

specification of any but the most straightforward kinds of looping. The example presented in this paper indicates some ways in which complex iterative control is essential to scientific inquiry. Moreover DFG-based process definitions complicate the clear depiction of such semantic features as exceptions and abstraction, whose importance was also demonstrated by the example in this paper. Indeed, as noted above, Figure 1 is a data flow diagram of the Water Budget process, but it is not a satisfactory PDG, as it lacks the ability to specify the details of all the different cases that can arise during execution of a scientific process. On the other hand, a DFG can often provide an intuitive depiction of how information flows through a system. Thus, a DFG might be useful as a complementary representation to the PDG for defining an initial, high-level view of the process. But difficulties arise in trying to maintain consistency between the PDG and DFG, especially at lower levels of hierarchical elaboration, where typical data or control flow oriented graphs have been found to be clumsy for representing more detailed flow. Indeed, we have developed a tool that translates Little-JIL-based PDGs into internal representations that are essentially equivalent to DFGs. Our experimentation with this tool has shown that even modest amounts of use of concurrency and exception handling can cause the number of nodes in the generated DFG to be hundreds or thousands of times as large as the number of steps in the original PDG.

There is also a substantial amount of work aimed at supporting the documentation of the provenance of scientific datasets. Many of the approaches to provenance documentation are summarized in [Simmhan, et.al 2005a; Simmhan, et.al. 2005b]. Indeed, these approaches have been compared to each other more formally in [Moreau, et.al 2007]. These approaches seem to fall generally into two categories. In one approach (e.g. [Buneman, et.al 2001; Lanter 1991; Aiken, et.al. 1996] each data artifact generated by execution of a scientific process is annotated with detailed information about the tool or system used to create the artifact, along with precise specification of the input artifacts used, and the output artifacts created. Each such annotation is then stored in a database. The complete documentation of the provenance of an artifact can then be obtained by recursively querying the database for the annotations that describe the activities that produced as outputs the artifacts used as inputs to the query. The second approach entails building a derivation graph on the fly as execution of the scientific process proceeds (e.g., [Foster, et.al 2002] and Kepler [Altintas, et.al. 2006]). We note that these two approaches are essentially equivalent to each other. Both collect provenance information by documenting the execution trace that has led to the creation of the data artifact being documented. In the former, the provenance structure is stored implicitly and is created upon demand by database queries. In the latter, the derivation structure is built incrementally during execution.

Our own approach falls into the latter category, entailing the on-the-fly construction of a derivation structure, namely the DDG. What distinguishes our work from the prior efforts is that our DDG depicts the progress of execution through our PDG, a process definition structure that can define and depict more complex semantic structures such as concurrency, exception handling, non-trivial iteration, and abstraction. Our DDGs offer depictions of how these semantic features are used to contribute to the development of data artifacts. Thus, for example, artifacts produced on different iterations through a given activity are shown as the roots of distinctly different subgraphs of a DDG, where

the context of each activity execution is provided by the DDG. Our work is strongly reminiscent of earlier work on the Odin project [Clemm, et.al. 1990] that documented the ways in which collections of software tools are applied to produce software products. As in the case of the work described in this paper, Odin maintained two coordinated structures, a type structure, showing which types of software objects can be generated through the applications of which software tools, and an instance structure that recorded the specific software artifacts generated by a specific sequence of applications of tools.

The Odin Project was aimed at supporting the clear and precise documentation of how various software artifacts resulted from somewhat different applications of somewhat different versions of various tools. It could use that documentation to make smart decisions about what data to store and what data to rederive, as well as to determine when to automatically do the derivation based on desired outputs. It thus extended earlier work on software configuration management (SCM), such as Make [Feldman 1979] and SCCS [Rochkind 1975]. It seems significant, therefore, to note that fundamental problems in documenting scientific data artifact provenance bear a striking resemblance to fundamental problems in SCM [Estublier, et.al. 2005]. In both cases there is a need to document and communicate a clear and precise understanding of how artifacts of interest have come into existence. In both cases, the derivation history may be quite complex, and must be maintained despite such complications as changes in versions of tools, reworking of artifacts, concurrent activities by diverse agents, etc. Thus, it is not surprising that the solution approach suggested in this paper is strongly reminiscent of approaches taken in early work in SCM. Indeed, we note that other recent work in scientific data provenance has also started to recognize the problem of documenting provenance in situations where the scientific process is evolving [Altintas, et.al. 2006; Callahan et.al. 2006] and we suspect that future work in the area of scientific data provenance documentation is likely to follow closely the progress of the SCM field.

7 Future Work

This paper describes the concept of an analytic web as the integration of two complementary graph types. Preliminary evaluation of this concept suggests that it is promising, but also suggests the need for considerable amounts of subsequent development and evaluation.

While we believe that the semantic features in Little-JIL present a useful starting point for considering the features that should be incorporated into languages that are used as the basis for defining the PDG, we have also noted a number of shortcomings. Further investigation of the essential requirements for the semantics of an analytic web PDG is needed. Similarly we believe that the DDG described in this paper offers clear value as the basis for the process metadata needed by scientists. But specific details concerning the DDG require further evaluation. For example, we need to evaluate various internal representations of the DDG to determine how to store the DDG efficiently while still supporting efficient creation of needed visual representations. Moreover, we have suggested that the datasets represented by the nodes of the DDG might be either regenerated from scratch, or might be cached to expedite generation of subsequent

datasets. Specific strategies for determining when and what to cache should be the subject of future research.

In addition, we have suggested that the DDG be used as the basis for the creation and attachment of process metadata to datasets, but further research is needed to determine how this is done best. We have noted that metadata standards such as EML are starting to appear, and have suggested that process metadata might be most usefully seen as an augmentation of standard annotations of this sort.

The value of an analytic web will be greatly enhanced by the availability of a tool set that supports such capabilities as the creation of the PDG, the execution of the PDG, the automatic creation of the DDG, viewing and querying these graph structures, and reasoning about the soundness of the scientific processes defined. We have indeed begun the creation of such a prototype toolset, called *SciWalker*. Our existing prototype provides weak and preliminary support for dataset producers, and virtually no support for dataset consumers. Future work will entail integration of our existing Little-JIL language support into *SciWalker* and capture of the dataset products of Little-JIL execution to support DDG creation and management. Extensive work on the development and evaluation of user interfaces to these tools, particularly emphasizing the sorts of visual depictions of the PDG and DDG, is clearly required as well.

Of particular interest is the possibility of using PDGs as vehicles for integrating other systems that support scientific workflow. We have noted above that there are a number of other projects that have developed capabilities for the specification and execution of scientific processes. Our view is that the basis of most of these systems upon a DFG model of such processes is severely limiting, and that the effective representation of such processes requires more expressive semantic features, such as those supported by the PDG presented here. But we also note that our concept of a PDG supports the idea that PDG steps can be performed by different agents, either human or automated. We suggest that existing scientific workflow systems, such as Kepler, might be used to define lower level scientific processes (e.g., those not entailing complex iteration or exception management), and that those process fragments might then be considered to be the agents responsible for performing PDG steps. Such an approach could make good use of the better developed performance features of established systems such as Kepler and the needed semantic features offered by our PDGs.

In future work, we also propose to add an important new dimension of support for dataset consumers by integrating powerful analyzers into *SciWalker*. One form of analysis that seems particularly important to dataset consumers is finite-state verification, which is capable of examining a PDG in order to determine whether or not it is possible to execute worrisome sequences of functional capabilities [Avrunin, et.al. 2000; Dwyer, et.al. 2004; Cobleigh, et.al. 2000]. We have noted that faulty scientific inferences can result for example when a dataset consumer applies certain types of interpolation models to datasets that been smoothed in certain ways by the dataset producer. This worrisome sequence of events may occur only for certain combinations of executions of the producer's process with the consumer's process. Such potential combinations can be detected by finite-state analysis of PDGs representing both processes. We suggest that it

is important to investigate how best to integrate such analysis capabilities into a toolkit such as SciWalker

Finally, we believe that the best way to make the progress needed in developing the ideas just outlined is to continue to create analytic webs to represent scientific processes of various kinds. Our work with ecological processes is encouraging, yet preliminary. We hope that there will be much more work of this sort, not just in ecology, but also in the representation of processes in a wide range of other sciences. This work should shed important light on the nature of languages needed to represent such processes, and tools needed to make them accessible to working scientists.

8 Acknowledgements

This material is based upon work supported by the National Science Foundation under Award No. CCR-0205575. Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

We are also grateful to many colleagues who supported this work and contributed key ideas that have led to our analytic web concept and our SciWalker prototype tool. In particular, we wish to thank Ed Riseman, Al Hanson, David Jensen, Paul Kuzeja, Howard Schultz, Bert Rawert, George Avrunin, and Mohammed Raunak for their advice, support, encouragement, and many stimulating conversations.

References

[Aiken, et.al. 1996]

A. Aiken, J. Chen, M. Stonebraker, and A. Woodruff, "Tioga-2: A Direct Manipulation Database Visualization Environment". Intl. Conf. on Data Engineering, 1996.

[Altintas, et. al 2004]

I. Altintas, C. Berkley, E. Jaeger, M. Jones, B. Ludäscher, S. Mock, system demonstration, "Kepler: An Extensible System for Design and Execution of Scientific Workflows", in 16th Intl. Conf. on Scientific and Statistical Database Management (SSDBM'04), 21-23 June 2004, Santorini Island, Greece.

[Altintas, et.al. 2005]

I. Altintas, A. Birnbaum, K. Baldrige, W. Sudholt, M. Miller, C. Amoreira, Y. Potier, and B. Ludäscher. "A Framework for the Design and Reuse of Grid Workflows". in Intl. Workshop on Scientific Applications on Grid Computing (SAG'04), LNCS 3458, Springer, 2005.

[Altintas, et.al. 2006]

I. Altintas, O. Barney, and E. Jaeger-Frank. "Provenance collection support in the Kepler Scientific Workflow System". in International Provenance and Annotation Workshop (IPAW), LNCS, Provenance and Annotation of Data, 4145: 118-132, 2006

[Avrunin, et.al. 2000]

G.S. Avrunin, J.C. Corbett, and M.B. Dwyer. "Benchmarking Finite-State Verifiers", **International Journal on Software Tools for Technology Transfer**, 2 (4), 2000, 317-320.

[Baldwin, et.al. 2004]

P. Baldwin, S. Kohli, E.A. Lee, X. Liu, and Y. Zhao. "Modeling of Sensor Nets in Ptolemy II. in Proc. of Information Processing in Sensor Networks", (IPSN), April 26-27, 2004, pp.359-368.

[Boose, et.al. 2007]

E. R. Boose, A. M. Ellison, L. J. Osterweil, L.A. Clarke, R. Podorozhny, J. L. Hadley, A. Wise, D. R. Foster. "Ensuring Reliable Datasets for Environmental Models and Forecasts", **Ecological Informatics**, ECOINF84, Vol. 2, No. 3, October 2007, Elsevier, pp. 237-247.

[Buneman, et.al 2001]

P. Buneman, S. Khanna, and W.C. Tan. "Why and Where: A Characterization of Data Provenance". in J. Van den Bussche and V. Vianu, editors, International Conference on Database Theory, pp. 316-330. Springer, LNCS 1973, 2001.

[Callahan, et.al. 2006]

S. Callahan, J. Freire, E. Santos, C. Scheidegger, C. Silva, and H. Vo. "Managing the Evolution of Dataflows with VisTrails". in Proceedings of the IEEE Workshop on Workflow and Data Flow for Scientific Applications (SciFlow 2006).

[Cass, et.al. 2000]

A. G. Cass, B.S. Lerner, E. K. McCall, L. J. Osterweil, S. M. Sutton, Jr., Alexander Wise. "Little-JIL/Juliette: A Process Definition Language and Interpreter". 22nd International Conference on Software Engineering (ICSE 2000), Limerick, Ireland, pp. 754-757, June 2000. (UM-CS-2000-066

[CLEANER URL]

<http://cleaner.ncsa.uiuc.edu/home/>

[Clemm, et.al. 1990]

G.M. Clemm and L.J. Osterweil. "A Mechanism for Environment Integration." **ACM Transactions on Programming Languages and Systems** (1990) 12(1): 1-25.

[Cobleigh, et.al. 2000]

J.M. Cobleigh, L.A. Clarke, and L.J. Osterweil. "Verifying Properties of Process Definitions". ACM SIGSOFT Intl. Symp. on SW Testing & Analysis. Portland, OR: ACM Press, 2000:96-101.

[Cooper, et.al. 2006]

B.F Cooper, R.S Barga. "Report on SciFlow 2006: the IEEE international workshop on workflow and data flow for scientific applications". SIGMOD Record, Volume 35, Number 3, p.54-56 (2006)

[CUASHI URL]

<http://www.cuahsi.org/>

[Deelman, et. al. 2003]

E. Deelman, J. Blythe, Y. Gil, and C. Kesselman. "Workflow Management in GriPhyN". in 14th International Conference on Scientific and Statistical Database Management (SSDBM'02), J. Nabrzyski, J. Schopf, and J. Weglarz editors, Kluwer, 2003.

[Dwyer, et.al. 2004]

M.B. Dwyer, L.A. Clarke, J.M. Cobleigh, and G. Naumovich. "Flow Analysis for Verifying Properties of Concurrent Software Systems". **ACM Trans. on Software Engineering and Methodology** 2004;13(4):359-430.

[Edwards, et.al. 2003]

S.A. Edwards and E.A. Lee. "The Semantics and Execution of a Synchronous Block-Diagram Language". **Science of Computer Programming**, Vol. 48, no. 1, July 2003, pp. 21-42.

[Ellison, et. al. 2006]

A. M. Ellison, L. J. Osterweil, J. L. Hadley, A. Wise, E. Boose, L. A. Clarke, D. R. Foster, A. Hanson, D. Jensen, P. Kuzeja, E. Riseman, H. Schultz. "Analytic Webs Support the Synthesis of Ecological Data Sets". **Ecology** V. 87, No. 6, pp. 1345-1358, June 2006.

[EML URL]

<http://knb.ecoinformatics.org/software/eml/>

[Estublier, et.al. 2005]

J. Estublier, D. Leblang, A. Van der Hoek, R. Conradi, G. Clemm, W. Tichy, and D. Wiborg-Weber. "Impact of Software Engineering Research on the Practice of Software Configuration Management", *ACM Trans. on Software Engineering and Methodology*, v. 14, #4 (2005) pp. 383-430.

[Fahringer, et.al. 2005a]

T. Fahringer, A. Jugravu, S. Pillana, R. Prodan, C. Seragiotto, and H. Truong. "ASKALON: A Tool Set for Cluster and Grid Computing". **Concurrency and Computation: Practice and Experience**, 17(2-4):143-169, 2005.

[Fahringer, et. al. 2005b]

T. Fahringer, R. Prodan, R. Duan, F. Nerieri, S. Podlipnig, J. Qin, M. Siddiqui, H.-L. Truong, A. Villazon, and M. Wiczorek. "ASKALON: A Grid Application Development and Computing Environment". In 6th IEEE/ACM International Workshop on Grid Computing, Seattle, USA, November 2005. IEEE Computer Society Press.

[Feldman 1979]

S.I. Feldman. "Make- a Program for Maintaining Computer Programs." **Software - Practice and Experience**. (1979) 9(4): 255-265.

[Foster, et.al. 2002]

I.T. Foster, J.-S. Voeckler, M. Wilde, Y. Zhao. "Chimera: A Virtual Data System for Representing, Querying, and Automating Data Derivation". in 14th International Conference on Scientific and Statistical Database Management (SSDBM'02) 2002.

[Foster, et.al. 2003]

I. T. Foster, J.-S. Vöckler, M. Wilde, and Y. Zhao, "The Virtual Data Grid: A New Model and Architecture for Data-Intensive Collaboration". in Conference on Innovative Data Systems Research, 2003.

[Harvard Forest URL]

<http://harvardforest.fas.harvard.edu/neon/neon.html>

[Heinis, et.al 2006]

T. Heinis, C. Pautasso, G. Alonso. "Mirroring Resources or Mapping Requests: Implementing WS-RF for Grid Workflows". *CCGRID 2006*, pp. 497-504.

[Lanter 1991]

Lanter, D.P., "Design of a lineage-based meta-data base for GIS", in **Cartography and Geographic Information Systems**, 18(4) pp. 255-261, 1991

[Ludäscher, et.al. 2006]

B. Ludäscher, I. Altintas, C. Berkley, D. Higgins, E. Jaeger-Frank, M. Jones, E. Lee, J. Tao, Y. Zhao, "Scientific Workflow Management and the Kepler System", in **Concurrency and Computation: Practice & Experience**, 18(10), pp. 1039-1065, 2006.

[Moreau, et.al 2007]

L. Moreau, B. Ludäscher, I. Altintas, R.S. Barga, S. Bowers, S. Callahan, B. Clifford, S. Cohen, S. Cohen-Boulakia, S. Davidson, E. Deelman, L. Digiampietri, I. Foster, J. Freire, J. Frew, J. Futrelle, T. Gibson, Y. Gil, C. Goble, J. Golbeck, P. Groth, D.A. Holland, S. Jiang, J. Kim, D. Koop, A. Krenek, T. McPhillips, G. Mehta, S. Miles, D. Metzger, S. Munroe, J. Myers, B. Plale, N. Podhorszki, V. Ratnakar, E. Santos, C. Scheidegger, K. Schuchardt, M. Seltzer, Y.L. Simmhan, C. Silva, P. Slaughter, E. Stephan, R. Stevens, D. Turi, H. Vo, M. Wilde, J. Zhao, and Y. Zhao. "The First Provenance Challenge", in **Concurrency and Computation: Practice and Experience**, Wiley InterScience, 2007

[NEON URL]

<http://www.neoninc.org/>

[Oinn, et.al. 2004]

T. Oinn, M. Addis, J. Ferris, D. Marvin, M. Senger, M. Greenwood, T. Carver, K. Glover, M. R. Pocock, A. Wipat and P. Li. "Taverna: A tool for the composition and enactment of bioinformatics workflows". **Bioinformatics Journal** 20(17) pp 3045-3054, 2004

[Osterweil, et.al. 2005]

L. J. Osterweil, A. Wise, L.A. Clarke, A. M. Ellison, J. L. Hadley, E. Boose, D. R. Foster. "Process Technology To Facilitate the Conduct of Science". Software Process Workshop (SPW2005), Beijing, China, Springer-Verlag Lecture Notes in Computer Science, Vol. 3840, pp. 403-415, May 25-27, 2005

[Pautasso, et.al. 2005]

C. Pautasso, G. Alonso, "The JOpera visual composition language". **Journal of Visual Language Computation** 16, pp. 119-152 (2005).

[Ptolemy URL]

<http://ptolemy.eecs.berkeley.edu/ptolemyII/>

[Rochkind, 1975]

M.J. Rochkind. "The Source Code Control System." **IEEE Transactions on Software Engineering** (1975) SE-1: 364-370.

[SAG 2004]

Scientific Applications of Grid Computing: First International Workshop (SAG 2004), Beijing, China, September 20-24, 2004, Lecture Notes in Computer Science, # 3458 (3), pp 119-132, ISBN3-540-25810-8.

[Simmhan, et.al 2005a]

Y.L. Simmhan, B. Plale, and D. Gannon. "A Survey of Data Provenance Techniques". Technical report 612, Comp. Sci. Dept., Indiana University, 2005

[Simmhan, et.al. 2005b]

Y.L. Simmhan, B. Plale, D. Gannon, "A survey of data provenance in e-science". in SIGMOD Rec. 34(3) pp. 31-36, 2005.

[SSDBM 2004]

16th Intl. Conf. on Scientific and Statistical Database Management (SSDBM'04), 21-23 June 2004, Santorini Island, Greece.

[WATERS URL]

<http://www4.nas.edu/webcr.nsf/MeetingDisplay1/WSTB-U-05-0A?OpenDocument&ExpandSection=1>

[Wise, et.al. 2000]

A. Wise, A. G. Cass, B. S. Lerner, E. K. McCall, L. J. Osterweil, S. M. Sutton, Jr. "Using Little-JIL to Coordinate Agents in Software Engineering". Automated Software Engineering Conference (ASE 2000), Grenoble, France, pp. 155-163, September 2000. (UM-CS-2000-045)

[Wise 2006]

A. Wise. "Little-JIL 1.5 Language Report". Department of Computer Science, University of Massachusetts, Amherst, MA 01003, October 2006. (UM-CS-2006-51)

[Wolstencroft, et. al. 2005]

K. Wolstencroft, T. Oinn, C. Goble, J. Ferris, C. Wroe, P. Lord, K. Glover, R. Stevens. "Panoply of Utilities in Taverna", in Proc E-Science 2005, 1st IEEE Intl Conf on e-Science and Grid Technologies, Melbourne, Australia, 5-8 December 2005