

A Scientific Workflow Model Designer based on Scientific Information Control Nets

Minjae Park*, Hyukjun Na*, Hyun Ahn**, Kwanghoon Pio Kim**

*BISTel, Inc., Seoul, Korea

**Dept. of Computer Science, KYONGGI UNIVERSITY, Suwon Kyonggido, Korea

{mjpark, hjna}@bistel-inc.com, {hahn, kwang}@kgu.ac.kr

Abstract—In this paper, we design and implement a scientific workflow process designer with a conceptual building block depicting its architectural structure. The designer is theoretically designed from the scientific information control net[1], and it is graphically implemented by expanding the standardized BPMN(business process modeling notations)[2]. In particular, the designer is able to automatically transform a BPMN-based graphical form into an XML-based textual form of the scientific workflow process model. Finally, we illustrate two captured-screens of the designer, as an operational example, which are corresponding to a simple scientific workflow process model and its XML-based representation, respectively.

Keywords-scientific workflow; information control net; data intensive workflow; formal workflow description; scientific workflow model designer

I. INTRODUCTION

Traditionally, the goal of a workflow management system was to describe, control, and monitor the enactments of business procedures in a workflow-supported organization. According to flourish on the almost all industries through the various types of process-driven business models, such as process portals, POD (Process-on-Demand), process choreography and orchestration, process collaboration, and inter-organizational workflows, the applicability of the workflow management system has been changing, broadening, and scaling up, too, along with dramatic progressions of workflow and business process technologies. In recent, what we see happening now in the workflow literature is that the applicability reaches up to the scientific knowledge discovery arena, which is so-called “Scientific Workflow[3][4],” that is able to not only support the whole stage of scientific experiments and simulations but also automate exploratory processes for discovering scientific knowledge involving cycles of observation, hypothesis formation, experiment design and execution. At this point, we would faithfully follow the conceptual definition of the scientific workflow stated in [3],[4], and [5], the like of which the scientific workflow implies the large scale and data intensive workflow, in terms of describing and deploying the overall idea and the basic concepts throughout the paper.

In this paper, we particularly focus on the quality of service issues in scientific workflow technologies, in terms of efficient load-balancing strategies in designing a large scale scientific workflow enactment architecture, and effective exception handling and recovery strategies in implementing a data inten-

sive scientific workflow system. In order to cope with these quality of service issues, we need a well-defined scientific workflow meta-model that is surely applicable to the scientific workflow model as well as the architectural functionality of the underlying scientific workflow system. As we know, there exist tens of proposals[6][7][8][9] and rationales[5][10] about defining the concept of scientific workflow models and systems, so far. However, we, in this paper, try to newly define a scientific workflow meta-model rather than choosing one of those existing scientific workflow models. We also extend the conventional information control net methodology[11] to graphically and formally represent the scientific workflow model spawned from the defined meta-model, which is dubbed “Scientific Information Control Net” that is abbreviated to **sciCN**, from now on.

In organizing the paper, we start from defining the basic concept of the scientific workflow meta-model with a series of functional definitions of the scientific activity types. The next section formalizes the scientific information control net and applies its graphical and formal representation to a pseudo scientific workflow model with describing the implications of the proposed mathematical formalism for large scale and data intensive scientific workflows. And next section describes the structure of the proposed scientific workflow designer. With describing an implementation of nodes which is support a scientific workflow model that is the information control nets, we finalize the paper.

II. SCIENTIFIC WORKFLOW MODEL

In this section, we define a scientific workflow meta-model for modeling scientific experimental processes characterized by the large scale and data intensive properties. Based upon the meta-model, we are able to instantiate sciCN-based scientific workflow models that are conceptually extended from the information control net methodology[11][12].

A. Meta-Model

In order to clearly define a scientific workflow model, it is needed to design a meta-model that is used to be signified by a set of entity types and their relationships. Furthermore, we have to take account of the conceptual definition of scientific workflows to conceive a reasonable and applicable meta-model. In general, the scientific workflow model is

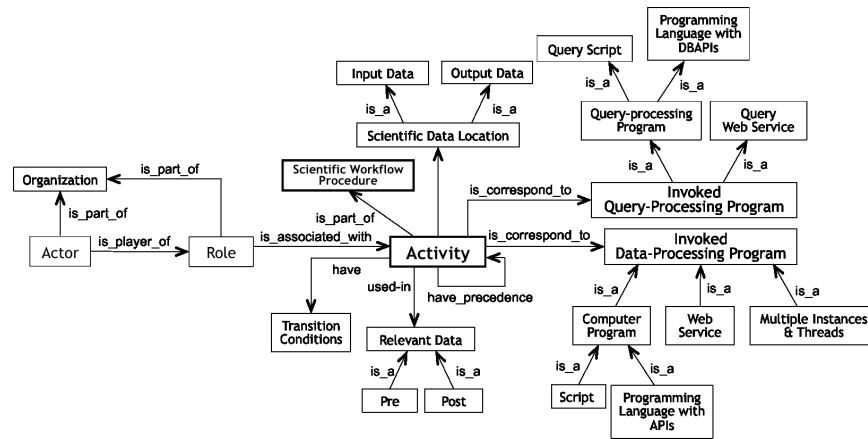


Fig. 1. Meta-model for Modeling Scientific Workflows

a mathematical methodology to specify various phases of the large and complex science process involving cycles of modeling and automation of computational experiments, data analysis, and data management. The well-described conceptual definition of the scientific workflows was introduced in the project DAKS[5] as followings:

“A scientific workflow is the description of a process for accomplishing a scientific objective, usually expressed in terms of tasks and their dependencies. Typically, scientific workflow tasks are computational steps for scientific simulations or data analysis steps. Common elements or stages in scientific workflows are acquisition, integration, reduction, visualization, and publication (e.g., in a shared database) of scientific data. The tasks of a scientific workflow are organized (at design time) and orchestrated (at runtime) according to data-flow and possibly other dependencies as specified by the workflow designer. Workflows can be designed visually, e.g., using block diagrams, or textually using a domain-specific language.”

We design a scientific workflow meta-model by faithfully reflecting the above definition as much as possible. Figure 1 depicts the scientific workflow meta-model drawn from in-depth discussions in the authors’ collaborative research groups. The primitive entity types constituting the scientific workflow model are Activity and Tasks (Invoked Query-processing program and Invoked Data-processing program) entity types; the primitive entity types has also certain associations with their own supplementary entity types, like Role, Data Location, and Relevant Data. The followings are the basic definitions of the primitive entity types:

- An **Activity** is a conceptual entity of the basic unit of work (computational task or step), and the activities in a scientific workflow model have precedence relationships, each other, according to their execution sequences based upon data-flow dependencies and/or a temporal ordering of computational steps. Also, the activity entity can be precisely specified by one of the elementary entity types, G-type, D-type, Q-type, and DQ-type, as depicted in Figure 2.
 - The elementary activity implies a computational step that can be realized by either data-processing

(D-type), query-processing (Q-type), or both (DQ-type) applications implemented through computer programs, threads, transactions, query-scripts, or web services. Especially, we suppose that the D-type and Q-type activities can be realized by applying the concepts of multiple instance patterns[5], like Multiple Instances with a Priori Design Time Knowledge and Multiple Instances with a Priori Run Time Knowledge implemented in Kepler[7][13]. Also, the D-type data processing activity implies a kind of computational steps for scientific simulations or data analysis tasks performing acquisition, integration, reduction, visualization, or publication (e.g., in a shared database) of scientific data.

- Particularly, the gateway (G-type) activity means a pair of split-join pseudo step that is used to controlling execution sequences of elementary/compound computational steps. Also, we assume that each of the gateway steps forms in the structured properties, *matched pair* and *proper nesting*. There are three sorts of G-type activities, such as conjunctive gateway (paralleling), disjunctive gateway (alternating), and iterative (looping) gateway. Particularly, the disjunctive gateway needs to be set some specific **Transition Conditions** in order to select one of the possible transition paths during the execution time. The transition condition itself can be defined by using certain input/output relevant data registered on the **Relevant Data** repository.
- Additionally, there may be a special computational step, the compound activity, that signifies a computational step containing another scientific workflow model, which is so-called **Subprocess**.
- A **Scientific Workflow Procedure** is defined by a pre-defined or intended set of tasks or computational steps, called Activities, which reflects the activities’ temporal ordering of executions and/or data-flow dependencies. A system, so-called scientific workflow management system, helps to organize, control, execute, and monitor such defined scientific workflow models. Conclusively, a

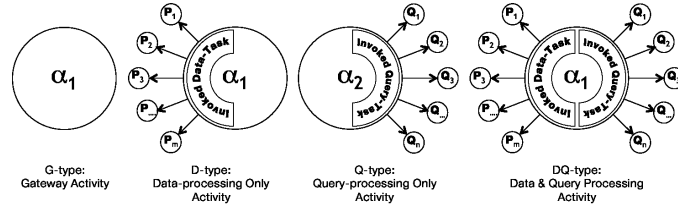


Fig. 2. Activity Types of Scientific Workflows

scientific workflow model can be described by a temporal order of the associated computational steps through the combinations of sequential logics, conjunctive logics (after activity A, do activities B and C, simultaneously), disjunctive logics (after activity A, do either activity B or C, alternatively), and iterative logics, as illustrated in Figure 3.

B. Scientific Information Control Net

In this paper, we would extensively adopt the original information control net (ICN) methodology[11] as the scientific workflow modeling methodology, which is dubbed to scliCN-based scientific workflow model upon the meta-model described in the previous subsection. Based upon the scientific information control net, we are going to dig up a possible way of analyzing activity dominancies on a scientific workflow model in the next section. The scientific workflow model needs to be described by a formal representation so that it is able to provide a means to eventually specify the model either in textual language or in database format, or in both. The following definition is the formal representation of the scientific information control net:

Definition 1: scliCN: Scientific Information Control Net of the scientific workflow model. A basic scliCN is 10-tuple $\Gamma = (\delta, \gamma, \lambda, \theta, \varepsilon, \pi, \vartheta, \kappa, \mathbf{I}, \mathbf{O})$ over a set \mathbf{A} of activities (including a set of compound/elementary/gateway computational steps), a set \mathbf{T} of transition conditions, a set \mathbf{R} of relevant data, a set \mathbf{D} of invoked data-task applications, a set \mathbf{Q} of invoked query-task applications, a set \mathbf{P} of scientific roles, a set \mathbf{L} of scientific data locations, and a set \mathbf{C} of participants, where $\varphi(\cdot)$ is the Power-set function.

- \mathbf{I} is a finite set of initial input data locations, assumed to be loaded with information by some external scliCNs before execution of the corresponding scliCN;
- \mathbf{O} is a finite set of final output data locations, perhaps containing information used by some external scliCNs after execution of the corresponding scliCN;
- $\delta = \delta_i \cup \delta_o$
where, $\delta_o : \mathbf{A} \rightarrow \wp(\mathbf{A})$ is a multi-valued mapping function of a computational step to its set of (immediate) successors,
and $\delta_i : \mathbf{A} \rightarrow \wp(\mathbf{A})$ is a multi-valued mapping function of a computational step to its set of (immediate) predecessors;
- $\gamma = \gamma_i \cup \gamma_o$
where $\gamma_o : \mathbf{L} \rightarrow \wp(\mathbf{A})$ is a multi-valued mapping function of a computational step to its set of output data

locations,

and $\gamma_i : \mathbf{L} \rightarrow \wp(\mathbf{A})$ is a multi-valued mapping function of a computational step to its set of input data locations;

- $\lambda = \lambda_a \cup \lambda_p$
where $\lambda_p : \mathbf{D} \rightarrow \wp(\mathbf{A})$ is a single-valued mapping function of a computational step to its invoked data-task application with multiple threads,
and $\lambda_a : \mathbf{A} \rightarrow \wp(\mathbf{D})$ is a multi-valued mapping function of an invoked data-task application with multiple threads to its set of associated computational steps;
- $\theta = \theta_a \cup \theta_p$
where $\theta_p : \mathbf{Q} \rightarrow \wp(\mathbf{A})$ is a single-valued mapping function of a computational step to its invoked query-task application with multiple threads,
and $\theta_a : \mathbf{A} \rightarrow \wp(\mathbf{Q})$ is a multi-valued mapping function of an invoked query-task application with multiple threads to its set of associated computational steps;
- $\varepsilon = \varepsilon_a \cup \varepsilon_p$
where $\varepsilon_p : \mathbf{P} \rightarrow \wp(\mathbf{A})$ is a single-valued mapping function of a computational step to one of the scientific roles,
and $\varepsilon_a : \mathbf{A} \rightarrow \wp(\mathbf{P})$ is a multi-valued mapping function of a scientific role to its sets of associated computational steps;
- $\pi = \pi_p \cup \pi_c$
where, $\pi_c : \mathbf{C} \rightarrow \wp(\mathbf{P})$ is a multi-valued mapping function of a scientific role to its set of associated scientists,
and $\pi_p : \mathbf{P} \rightarrow \wp(\mathbf{C})$ is a multi-valued mapping function of a scientist to its sets of associated scientific roles;
- $\vartheta = \vartheta_i \cup \vartheta_o$
where $\vartheta_o : \mathbf{R} \rightarrow \wp(\mathbf{A})$ is a multi-valued mapping function of a computational step to its set of output relevant data,
and $\vartheta_i : \mathbf{R} \rightarrow \wp(\mathbf{A})$ is a multi-valued mapping function of a computational step to its set of input relevant data;
- $\kappa = \kappa_i \cup \kappa_o$
where $\kappa_i(\alpha) : \mathbf{T}$: sets of control-transition conditions, \mathbf{T} , on each arc, $(\delta_i(\alpha), \alpha), \alpha \in \mathbf{A}$;
and $\kappa_o(\alpha) : \mathbf{T}$: sets of control-transition conditions, \mathbf{T} , on each arc, $(\alpha, \delta_o(\alpha)), \alpha \in \mathbf{A}$;
where the set $\mathbf{T} = \{default, or(conditions), and(conditions)\}$.

In terms of the graphical representation, we principally adopt the original notations of information control net[11][12]. As depicted in Figure 3, there are possible primitive computation-flows and their related gateway G-type activities. That is, the conjunctive (or parallel) computation-flow type with a pair of conjunctive gateway G-type activities, split-AND and join-AND, is represented by solid dots(\bullet), mean-

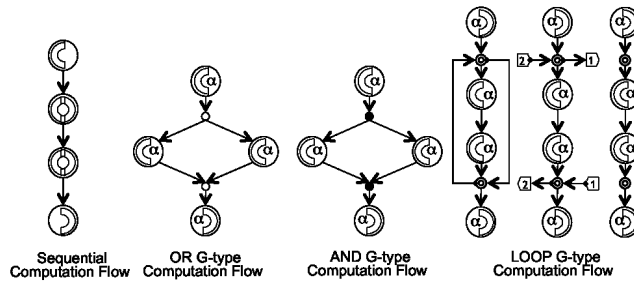


Fig. 3. Logics of Temporal Orders with Gateway Activities in Scientific Workflows

while the disjunctive (or decision) computation-flow type with a pair of disjunctive gateway G-type activities, split-OR and join-OR, is represented by hollow dots(\circ). Also, the iterative (loop) computation-flow type with a pair of loop gateway G-type activities, split-LOOP and join-LOOP, is represented by double hollow dots.

Besides, in order to be syntactically safe, it is very important for these gateway G-type activities to keep the structured properties—proper nesting and matched pair properties. Therefore, in specifying a scientific information control net, not only each of the gateway G-type activities always keeps matched pair with split and join types, but also multiple sets of the gateway G-type activities keep in a properly nested pattern.

III. SCIENTIFIC WORKFLOW MODEL DESIGNER

In order to define a scientific workflow model, it is needed to implement a designer that is a scientific workflow model designer supporting scICN. So, we have implemented a scientific workflow model designer.

The proposed scientific workflow designer is based on standard workflow model designer which support plug-in extension. As follow workflow/BPM standard that is BPMN, it is easy integrate and execute a workflow/BPM engine like Activiti[14] which is an open source software.

In scICN designer, nodes are able to define three type activities, Q-Type, D-Type and DQ-Type, use an extensible function for workflow standard modeling notation that is the BPMN[2]. But, it could not define 'G-Type' activity, instead as it is 'Gateway' nodes of BPMN. That gateway nodes could describe the conjunctive or parallel logic of scientific workflow model.

And, this designer should be generate executable language for execution of the defined scientific workflow model, executable language would be an extension of BPMN and BPEL.

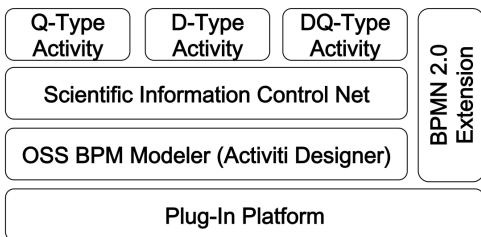


Fig. 6. The Conceptual Building Block of the Scientific Workflow Model Designer

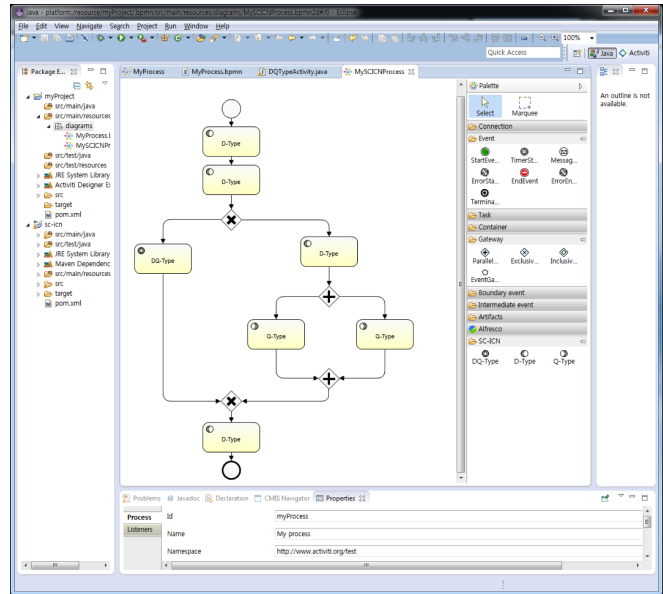


Fig. 4. The scICN Scientific Workflow Model Designer

```
<?xml-stylesheet href="http://www.omg.org/spec/BPMN/20100524/MODEL" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
xmlns:activiti="http://activiti.org/bpmn" xmlns:bpmndi="http://www.omg.org/spec/BPMN/20100524/DI"
xmlns:omgdc="http://www.omg.org/spec/DC/20100524/DC" xmlns:omgd="http://www.omg.org/spec/DO/20100524/DI"
type="text/xml" xmlns:xsd="http://www.w3.org/2001/XMLSchema" expressionLanguage="http://www.w3.org/1999/XMLSchema"
targetNamespace="http://www.activiti.org/test">
<process id="MyProcess" name="My process" isExecutable="true">
<startEvent id="startevent1" name="Start"/>
<endEvent id="endevent1" name="End"/>
<serviceTask id="servicetask1" name="D-Type" activiti:class="ctrl.sc.icn.runtime.DTypeJavaDelegat ion"
activiti:extensionId="ctrl.sc.icn.servicetasks.DTypeActivity"/>
<serviceTask id="servicetask2" name="D-Type" activiti:class="ctrl.sc.icn.runtime.DTypeJavaDelegat ion"
activiti:extensionId="ctrl.sc.icn.servicetasks.DTypeActivity"/>
<exclusiveGateway id="exclusivegateway1" name="Exclusive Gateway"/>
<serviceTask id="servicetask3" name="DQ-Type" activiti:class="ctrl.sc.icn.runtime.DQTypeJavaDelegat ion"
activiti:extensionId="ctrl.sc.icn.servicetasks.DQTypeActivity"/>
<serviceTask id="servicetask4" name="D-Type" activiti:class="ctrl.sc.icn.runtime.DTypeJavaDelegat ion"
activiti:extensionId="ctrl.sc.icn.servicetasks.DTypeActivity"/>
<parallelGateway id="parallelgateway1" name="Parallel Gateway"/>
<serviceTask id="servicetask5" name="Q-Type" activiti:class="ctrl.sc.icn.runtime.QTypeJavaDelegat ion"
activiti:extensionId="ctrl.sc.icn.servicetasks.QTypeActivity"/>
<serviceTask id="servicetask6" name="Q-Type" activiti:class="ctrl.sc.icn.runtime.QTypeJavaDelegat ion"
activiti:extensionId="ctrl.sc.icn.servicetasks.QTypeActivity"/>
<parallelGateway id="parallelgateway2" name="Parallel Gateway"/>
<serviceTask id="servicetask7" name="D-Type" activiti:class="ctrl.sc.icn.runtime.DTypeJavaDelegat ion"
activiti:extensionId="ctrl.sc.icn.servicetasks.DTypeActivity"/>
<exclusiveGateway id="exclusivegateway2" name="Exclusive Gateway"/>
<sequenceFlow id="flow1" sourceRef="startevent1" targetRef="servicetask1"/>
<sequenceFlow id="flow2" sourceRef="servicetask1" targetRef="servicetask2"/>
<sequenceFlow id="flow3" sourceRef="servicetask2" targetRef="exclusivegateway1"/>
<sequenceFlow id="flow4" sourceRef="exclusivegateway1" targetRef="servicetask3"/>
<sequenceFlow id="flow5" sourceRef="servicetask3" targetRef="servicetask4"/>
<sequenceFlow id="flow6" sourceRef="servicetask4" targetRef="parallelgateway1"/>
<sequenceFlow id="flow7" sourceRef="parallelgateway1" targetRef="servicetask5"/>
<sequenceFlow id="flow8" sourceRef="servicetask5" targetRef="parallelgateway2"/>
<sequenceFlow id="flow9" sourceRef="servicetask5" targetRef="parallelgateway2"/>
<sequenceFlow id="flow10" sourceRef="servicetask6" targetRef="parallelgateway2"/>
<sequenceFlow id="flow11" sourceRef="parallelgateway2" targetRef="exclusivegateway2"/>
<sequenceFlow id="flow12" sourceRef="servicetask3" targetRef="exclusivegateway2"/>
<sequenceFlow id="flow13" sourceRef="exclusivegateway2" targetRef="servicetask7"/>
<sequenceFlow id="flow14" sourceRef="servicetask7" targetRef="endevent1"/>
</process>
```

Fig. 5. The Executable Language for scICN Scientific Workflow

IV. CONCLUSIONS

So far, this paper has implemented a scientific workflow model designer with expatiation with the formal scientific information control net and its related definitions by extending the original information control net methodology. Based upon these formal and graphical definitions, we are able to construct a scientific workflow model. In recent, the literature needs various, advanced, and specialized scientific workflow modeling methodologies that are well-suitable to the domain of data intensive and very large scale scientific computing environments. We so strongly believe that this work might be one of those impeccable attempts and pioneering contributions for improving and advancing the large scale and data intensive and scientific workflow management technology.

ACKNOWLEDGEMENTS

This research was supported by the Technology Innovation Program (10045913, Development of Big Data based Analysis and Control Platform for Semiconductor Manufacturing Plants) funded By the Ministry of Trade, industry & Energy(MOTIE, Korea)

REFERENCES

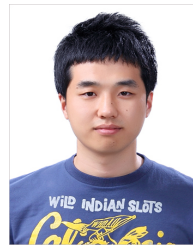
- [1] H. Ahn, M. Park, and K. P. Kim, "scicn: Scientific information control nets," in *Proceedings of the International Conference on Advanced Communications Technology*. IEEE, February 2014, pp. 1163–1166.
- [2] T. Allweyer, *BPMN 2.0: introduction to the standard for business process modeling*. BoD–Books on Demand, 2010.
- [3] A. Barker and J. van Hemert, "Scientific workflow: A survey and research directions," *Lecture Note in Computer Science*, vol. 4967, pp. 746–753, 2008.
- [4] J. Qin and T. Fahringer, *SCIENTIFIC WORKFLOWS: Programming, Optimization, and Synthesis with ASKALON and AWDL*. Springer, 2012.
- [5] U. Yildiz, A. Guabtini, and A. H. Ngu, "Business versus scientific workflow: A comparative study," *University of California Davis, Department of Computer Science, Research Report, Project DAKS*, vol. n 2009-3, pp. 1–18, 2009.
- [6] N. H. T. M. Anne H.H. Ngu, Shawn Bowers and T. Critchlow, "Flexible scientific workflow modeling using frames, templates, and dynamic embedding," *Lecture Note in Computer Science*, vol. 5069, pp. 566–572, 2008.
- [7] B. Ludascher and *et al.*, "Scientific workflow management and the kepler system," *Journal of Concurrency and Computation: Practice & Experience*, vol. 18, no. 10, pp. 1039–1065, August 2006.
- [8] P. Yang, Z. Yang, and S. Lu, "Formal modeling and analysis of scientific workflows using hierarchical state machines," in *Proceedings of IEEE International Conference on e-Science and Grid Computing*. IEEE, December 2007, pp. 619–626.
- [9] E. Ogasawara, D. de Oliverira, and P. Valduriez, "An algebraic approach for data-centric scientific workflows," in *Proceedings of the 37th International Conference on Very Large Data Bases*. VLDB Endowment, September 2011, pp. 619–626.
- [10] Y. Gil and *et al.*, "Examining the challenges of scientific workflows," *IEEE Computer*, vol. 40, no. 12, pp. 24–32, December 2007.
- [11] C. A. Ellis, "Information control nets: A mathematical model of information flow," in *Proceedings of Conference on Simulation, Modeling, and Measurement of Computer Systems*. ACM, 1979, pp. 225–240.
- [12] K. P. Kim and C. A. Ellis, *Section II / Chapter VII. An ICN-based Workflow Model and Its Advances, Handbook of Research on BP Modeling*. IGI Global, ISR, 2009.
- [13] V. Curcin and M. Ghanem, "Scientific workflow systems - can one size fit all?" in *Proceedings of the Cairo International Biomedical Engineering Conference*. IEEE, December 2008, pp. 1–9.
- [14] T. Rademakers, *Activiti in Action: Executable business processes in BPMN 2.0*. Manning Publications Co., 2012.



Minjae Park Minjae Park is a senior member of research staff at the solution R&D research center of BISTel, Inc., South Korea. He received B.S., M.S., and Ph.D. degrees in computer science from Kyonggi University in 2004, 2006, and 2009, respectively. His research interests include groupware, workflow systems, BPM, CSCW, collaboration theory, process warehousing and mining, workflow-supported social networks discovery and analysis, and process-aware factory automation systems.



Hyukjun Na Hyukjun Na is a senior member of R&D center of BISTel, Inc., South Korea. He received the B.S. degree in industrial engineering from Hansung University, Seoul, South Korea and the M.S. degree in industrial engineering from Korea University, Seoul, South Korea, in 1998 and 2000, respectively. His current research interests are data mining, clustering, fault detection, statistical process control, run to run control, workflow system, and BPM.



Hyun Ahn Hyun Ahn is a full-time Ph.D. student of computer science department and a graduate member of the collaboration technology research laboratory at Kyonggi University, South Korea. He received B.S. and M.S. degrees in computer science from Kyonggi University in 2011 and 2013, respectively. His research interests include workflow systems, BPM, scientific workflow systems, workflow-supported social and affiliation networks discovery, analysis, and visualization.



Kwanghoon Pio Kim Kwanghoon Pio Kim is a full professor of computer science department and the founder and supervisor of the collaboration technology research laboratory at Kyonggi University, South Korea. He received B.S. degree in computer science from Kyonggi University in 1984. And he received M.S. degree in computer science from Chungang University in 1986. He also received his M.S. and Ph.D. degrees from the computer science department at University of Colorado Boulder, in 1994 and 1998, respectively. He had worked as researcher and developer at Aztek Engineering, American Educational Products Inc., and IBM in USA, as well as at Electronics and Telecommunications Research Institute (ETRI) in South Korea. In present, he is a vice-chair of the BPM Korea Forum. He has been in charge of a country-chair (Korea) and ERC vice-chair of the Workflow Management Coalition. He has also been on the editorial board of the journal of KSII, and the committee member of the several conferences and workshops. His research interests include groupware, workflow systems, BPM, CSCW, collaboration theory, Grid/P2P distributed systems, process warehousing and mining, workflow-supported social networks discovery and analysis, process-aware information systems, data intensive workflows, and process-driven Internet of Things.