

Recording the Context of Action for Process Documentation

Ian Wootten and Omer Rana

School of Computer Science, Cardiff University, UK

Abstract. In reviewing evidence about real world processes, being aware of the context in which activities within such processes are performed enables us to make more informed judgements. It is necessary to distinguish between the environment in which a process occurs, and the sequence of activities which form part of the description of that process. Each of these types of information is complementary to understanding the other and therefore making associations between them is also important. Our work has been exploring the use of *context* in documenting a process and working toward a solution which incorporates the two. We present an approach to automatically relating properties of workflow actors to the documentation of the process within which these actors are involved.

1 Introduction

Context plays a crucial role in support of evidence for a given argument. Statements which are taken ‘out of context’ could face criticism from those who note such omissions as a distortion of the original intended meaning. There are a number of definitions of context in distributed systems – Brown [1] defines context to be the elements of a user’s environment which the computer knows about. Dey and Abowd [2], refer to context as “any information that can be used to characterize the situation of entities (i.e., whether a person, place or object) that are considered relevant to the interaction between a user and an application, including the user and the application themselves”. Therefore the amount of information we have about context affects how we interpret the event we are given. Similarly, the way in which we interpret the recording of some data may be altered based on the context in which it is presented.

We use Groth et al.’s view on provenance [3] as the process which led to a piece of data, and that such processes may be described using evidence represented in the form of process documentation. In service oriented architectures we believe that the context of an action or sequence of actions may be provided by each actor involved in execution of a process. Records of context may exist outside of the notion of a process or the control of a client and this makes later interpretation based on review of both difficult. Our work has been exploring the use of context in the documentation and structuring of such evidence. Known relationships between a process and its context will be different depending on the application, and it may not always be possible to document both in open, loosely coupled architectures. In this paper we present a system which documents both the sequence of actions which describe a process and the situation in

which those actions took place, enabling both navigation of process documentation and prediction of future actor properties. The rest of this paper is organised as follows: in section 2 we review the relevant literature in the area of context as relating to provenance systems and describe our motivation. In section 3 we describe our model of the context of actions upon an actor, followed by a description of the architecture we have adopted in section 4 and a demonstration of its implementation in section 5. Finally in section 6 we conclude.

2 Background and Motivation

Provenance is important for scientists to be able to record information in order to, for example, ensure experiments are performed correctly and to be able to repeat processes which produced interesting results. As many disparate activities may be involved in such a process, without such recording it is difficult to determine precisely how results have been reached. Several solutions have been developed to capture provenance or enable applications to be *provenance aware*, e.g. in Bioinformatics[11] or Chemical Sciences[9]. The Oxford English dictionary describes context as: *the circumstances that form the setting for an event, statement or idea*. For actors in a service oriented system, documenting a process involves describing each of those steps which comprise it. Research to date has widely addressed documenting messages which are sent between actors[3] as these typically indicate invocation of some functionality. Other actions commonly include the actions of scientists who control such systems, which may be documented in a more ad-hoc manner. Recording single events does not however describe properties or conditions which hold true for the actor over a given period of time. We refer to such properties as a description of the *context* of the action, where the action is the function being performed by a particular entity. A universal agreement on the content of process documentation representing context submitted by actors has not yet been reached. Attempts to provide a generic schema for what is to be recorded as context have so far proved to be fruitless, with the Grid Provenance project¹ choosing not to adopt any formal structure for the state of actors during a process. This is due to the diversity in the types of use cases that are required to be satisfied for all those domains which have unique provenance requirements. Such diversity has led to scientists building a variety of tools able to capture specific contextual data. Recently, the Open Provenance Model has been developed to enable sharing of provenance data amongst the different systems that adopt it [8]. In this model, contextual pieces of data could be considered *artifacts* as they are immutable pieces of state. As yet, no systems are known to have implemented the model and no formal representation of the model has been specified. The model we present focuses less on what the content of this contextual data may be and more on how those elements are represented and recorded over time, to be of use during queries of process documentation.

3 Modeling the Context of a Process

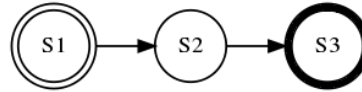
We consider a number of numerical variables which are all able to be measured over the same particular period of time upon an actor. Each of these is recorded at

¹ <http://www.gridprovenance.org>

some regular interval which differ between variables. We choose numeric variables as these are commonly able to be monitored in distributed systems, but our method is equally applicable to other primitive types, such as strings. We term the conceptual representation of these variables at any given time the *state* of the actor on which they were observed. It is possible for an actor to progress through a number of different states over a period of time. We assume that the primary cause for any transition between states is an instantaneous event in our model. Such events are considered to occur when any one of the measured variables changes value, which results in a transition in the state of that actor. These events may be internal and only visible to the actor in question, or external and able to be documented by others (such as receiving or responding to messages). We assume when state transitions are documented that an actor is only involved in one process at a time. In a set of actions invoked by dispatch of request messages, the request event is always assumed to be the cause of transition for the state changes in the interval which follows. An example of deriving a number of states from a given set of variables (v_1, v_2, v_3) is given in figure 1(a) with a finite state machine representing the states mined from the data shown in figure 1(b). Here we see that although some common values exist throughout measurement of each variable (such as the value of v_3) a change in any variable can lead to a state change.

State	Observation Time	v_1	v_2	v_3
s_1	1164277522	4.71	13084	2.56
s_2	1164282522	4.71	15698	2.56
s_3	1164287522	4.00	15698	2.56

(a) Deriving unique states from variables



(b) FSM of variable data

Fig. 1

In the scenario of a service based architecture, the most interesting states are those which occur within the interval when a request message was sent to a service and the associated response message was sent by that service. In this period, the actor's observed state may be documented as a part of any process documentation that is recorded for a process. Representing state using *time intervals* means it is possible to mine series which describe the same property as being true over a number of non-overlapping intervals. In figure 2(a) we demonstrate how thresholds are used to segment the measurements of a variable[4], to determine when it was within pre-defined ranges. By coinciding two series of this type, (as in figure 2(b)) we are able to determine the unique states of an actor over that period, when using those two variables as state components. Our resultant (coincidence) series therefore describes the periods over which each variable is described as high, medium or low.

The relationship type which exists between the states which occur during a process and any documented events is assumed by default to be causal. That is, states observed upon an actor are assumed to be an effect of the request (cause) event from which the observation interval begins. This is in order to support prediction of future actor properties - which would be impossible without causal

knowledge. Response events are assumed not to hold a direct relationship with a state, although these are also assumed to be an effect of the request event. Documenting both these relationships mean that it is possible to easily determine the request/response events which hold relevance for a sequence of states.

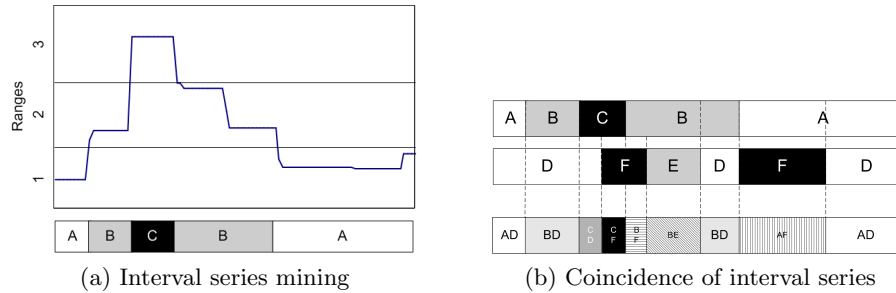


Fig. 2: Mining unique context series from numeric time series

To the best of our knowledge, modeling observations of context using an interval representation is a novel approach to gathering metadata about processes within the provenance community. Particular focus has previously been paid to documenting the set of events which comprise a process, without discussion of whether any one of these events may hold relationships with properties and conditions holding true for longer than a single point in time. By adopting a representation capable of representing intervals we hope to capture such knowledge.

4 Documenting Context in Service based Architectures

We use the PreServ software created at the University of Southampton to capture assertions of provenance to a repository known as a *provenance store* [3]. This software has been built in response to a large variety of requirements gathered from numerous domains such as bioinformatics, high energy physics and medicine [5]. PreServ breaks up process documentation into three sub-categories of assertion known as *p-assertions*. Interaction p-assertions document message exchange between services, relationship p-assertions document the causal dependencies between events or data items and actor state p-assertions document the state an actor is in at a given point in time during a process. Dividing documentation into these three types means that parts of it may be recorded by each of the actors which were involved in a process to a repository common to all, known as a *provenance store*. PreServ is suitable as a capture mechanism for assertions of state as it does not prescribe their contents, instead leaving it up to specific applications to define this. This leaves us free to specify our own XML representation of state and assert that to storage.

Our implementation of a system capable of automatically documenting state centres around use of a State Assertion Registry (StAR) co-located with a service [10]. StAR is implemented as a Java library and acts as a wrapper to the

service enabling it to dynamically record assertions of provenance according to a policy file with which the service is co-located. The state descriptions which are documented detail the intervals in which they occurred and a *pattern* unique to each state describing the range which each variable fell within. StAR represents a benefit to the scientist in capturing assertions automatically upon policy configuration, with the alternative being a manual instrumentation of each actor with the content for each assertion to be made.

Variable data is collected at discrete time instants and *segmented*[4] to one of a set of possible values using thresholds based on average values of the variables previously observed. The segmented value corresponds to an element within the pattern for a particular state. We use techniques from the Time Series Knowledge Representation (TSKR) [6] to determine the intervals in which the segmented series coincide with one another as shown in figure 2(b). Details of these series are then used as the content of an actor state p-assertion, along with a pattern description indicating which conditions hold over the series. Following recording of documentation, queries of the evidence submitted by each of the actors may be performed through a query to a provenance store.

A user can determine future states for a process based upon the states previously documented within a provenance store. For all states which are related to the same event, a transition table listing the probability of state transition between two past states (given that observed event) may be calculated. The most likely next state for an actor is that with the largest probability value corresponding to the current state in the table. In cases where the two states (predicted and actual) do not match we use a similarity measure to find how different those states are. It is a simple distance measure of each corresponding pattern value, shown in equation 1, where q and r are the two patterns being compared and p and t are the number of items in the patterns and the number of possible values for each of those items. The total number of possible states is p^t , though for any given process run not all states may be observed.

$$s = 1 - \frac{\sum_{n=0}^p |q_n - r_n|}{p \times (t - 1)} \quad (1)$$

The distance of two states therefore is the total measured error between each of the states pattern elements, divided by the maximum total distance possible for error on each of those elements. The total similarity of a process is calculated from the product of the similarities of each state against those in a comparison process. This similarity value gives us a single measure of how similar the conditions under which each of the actors involved were operating were for two processes.

5 Evaluation

We now demonstrate two uses of documenting actor state for a process; 1) Attempting to predict future actor properties for processes based on previously documented ones and 2) Assisting navigation of process documentation based on comparison of states over the intervals actors were involved in a process. The

workflow we use to demonstrate this was the subject of both the first and second provenance challenges [7]. It is used to create population-based brain atlases from high resolution anatomical data from the Functional Magnetic Resonance Imaging (fMRI) Data Center². We use StAR to automatically collect assertions of interaction and state to a provenance store for each of the services which comprise the workflow. We focus on the last two actions in the workflow, which convert an averaged brain image (determined from the average of intensities of MRI scans) gathered from a collection of high resolution anatomical data into graphics files showing slices of the brain. Actor state assertions identify the interval over which the actor is invoked (between request and response messages) based upon TSKR mined series from the segmented values of the three metric values for an actor: bytes in per second, one minute load average and the amount of buffered memory.

Each of the actions performed in the workflow is hosted separately on a IBM JS20 blade machine (2 x 2.1GHz, 1.5GB RAM) and the provenance store for each of them is a Sun x2100 (1 x 2.2GHz, 4GB of RAM). When the process executes, a single action is performed by each of the services used and a set of states are recorded for it. We perform the process 1000 times, delaying subsequent invocations to allow the systems to recover. For our state prediction evaluation, results are based on this experiment being performed twice.

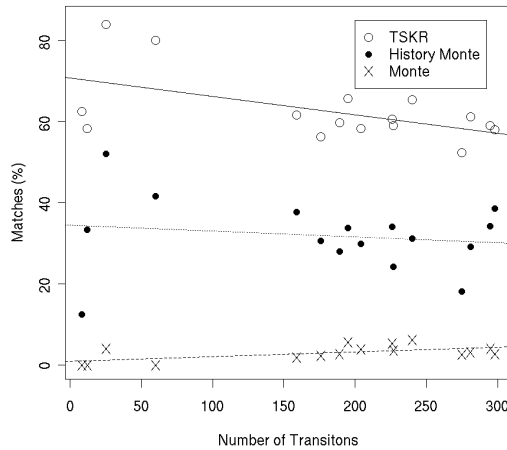


Fig. 3: Match rate of predicted states to those observed

Prediction of Future Actor Properties - A scientist is able to build a likely model for future properties using a transition table (as described in section 4) for each action. A state prediction is made for each actor after a single process is executed, based upon the last known state for that actor and the most common transition observed. Figure 3 shows the percentage of matches for

² <http://www.fmridc.org/>

predicted states to actual ones for our approach along with a history based and simple monte-carlo prediction. Each point represents the match rate for a single actor in the process. All machines consistently predicted states at a reasonably high success rate (50-85%), which was always above that of the monte-carlo predictions. The average trend indicates that as more transitions are observed, state becomes more difficult to predict in the future. This is due to the increased complexity of the model which is built when more transitions are found. It is likely that a more sophisticated analysis of the transition pattern leading to a state could further increase this success rate. Using this approach, the scientist executing the process is able to form a hypothesis detailing the most likely states to occur for each actor during future invocation of each action.

Assistance of Process Documentation Navigation - We demonstrate reduction of manual navigation by determining the similarity of a single state for each action within the process against that observed in a “comparison process”, with our results shown in figure 4(a). The total similarity of a process (as defined in section 4) is the product of multiplying each of these similarities for each action. The scientist may then use these distances as a filter to locate the most interesting processes from a large collection of process documentation. Our results in figure 4 show the distribution of process similarity values. For our scenario, we are able to see that the total documentation to be navigated is reduced dramatically when searching for either those processes with a high or low similarity (≥ 0.9 or ≤ 0.4). This corresponds to 12% and 8% of all of the documentation recorded. If we look at the lowest similarity processes (≤ 0.3), we can reduce this figure even further to 1% of all documentation. Figure 4(b) shows the same processes being compared but with an average similarity value for when multiple states are observed within each invocation. We reveal a further number of interesting processes within the 0.1-0.2 range and 0.7-0.8 ranges by doing this, including even smaller subsets of documentation. Without the documentation of actor states, it is perfectly feasible that the navigation of records of all 1000 processes (totalling 56MB’s worth of XML documentation to be queried in our own experiments) would have to be navigated manually.

6 Conclusion

In modeling actions performed by entities working as part of a process, strict event-based documentation may not be appropriate for documenting all process features. Instead, an interval based representation – such as the one presented in this paper, better represents observation of properties which hold true over a period of time. We have shown here that by documenting context along with process, it is possible to query provenance repositories to predict the future properties of actors or find other process traces which exhibit similarities to a model trace. Where vast collections of process documentation exist for the same process, being able to filter more interesting information for a scientist can present both time saving benefits and a reduction of queries of documentation. In our evaluation given, this was able to be reduced to as little as 1% of the overall documentation which was captured.

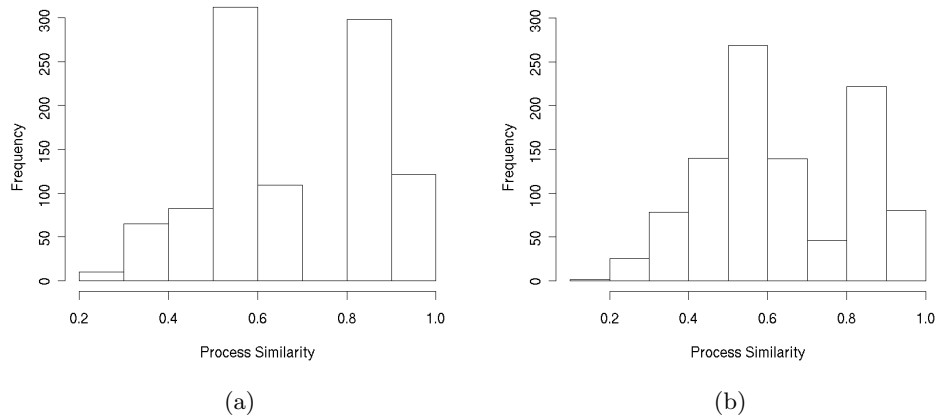


Fig. 4: Distribution of process similarity values when compared to model process

References

1. P. J. Brown. The Stick-e Document: A Framework for Creating Context-aware Applications. *Electronic Publishing - Origination, Dissemination, and Design*, 8(2/3):259–272, June/September 1995.
2. Anind K. Dey and Gregory D. Abowd. Towards a Better Understanding of Context and Context-Awareness. In *Workshop on The What, Who, Where, When, and How of Context-Awareness*, ACM Press, New York, April 2000.
3. Paul Groth, Simon Miles, and Luc Moreau. PReServ: Provenance Recording for Services. In *Proceedings of the UK OST e-Science second All Hands Meeting 2005 (AHM'05)*, 2005.
4. E. Keogh, S. Chu, D. Hart, and M. Pazzani. *Segmenting Time Series: A Survey and Novel Approach*, 1993.
5. Simon Miles, Paul Groth, Miguel Branco, and Luc Moreau. The requirements of using provenance in e-Science experiments. January 01 2006.
6. Fabian Moerchen. Algorithms for time series knowledge mining. In *KDD '06: Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 668–673, New York, NY, USA, 2006. ACM.
7. Luc Moreau and Bertram Ludäscher Editors. The First Provenance Challenge. *Concurrency and Computation: Practice and Experience*, Volume 20, Issue 5:409 – 418, 2007.
8. Luc Moreau, Juliana Freire, Joe Futrelle, Robert McGrath, Jim Myers, and Patrick Paulson. The Open Provenance Model, December 2007. [Online] <http://eprints.ecs.soton.ac.uk/14979/>.
9. James D. Myers, Carmen M. Pancarella, Carina S. Lansing, Karen L. Schuchardt, Brett T. Didier, and C N. Ashish, Goble. Multi-scale Science: Supporting Emerging Practice with Semantically Derived Provenance, March 06 2006.
10. Ian Wootten, Shrija Rajbhandari, and Omer Rana. Automatic Assertion of Actor State in Service Oriented Architectures. *ICWS07*, pages 655–662, 2007.
11. Jun Zhao, Carole A. Goble, and Robert Stevens. An Identity Crisis in the Life Sciences. In Luc Moreau and Ian T. Foster, editors, *IPAW*, volume 4145 of *Lecture Notes in Computer Science*, pages 254–269. Springer, 2006.