

Evaluating Data Integrity in the Cloud using the UPPAAL

Sachi Ishida¹ and Yoshiyuki Shinawa²

¹Fujitsu FSAS Inc., 13-2 Nakamaruko, Nakahara-ku, Kawasaki, Kanagawa, Japan

²Graduate School of Science and Technology, Ryukoku University, 1-5 Seta Oe-cho Yokotani, Otsu, Shiga, Japan

Keywords: Cloud Computing UPPAAL Transaction Processing Data Integrity

Abstract: There are several considerations when implementing a transaction processing system in cloud environments like Google App Engine (AE). One of the most critical ones is the data integrity since the cloud provides us with limited capability for it. Therefore we need to evaluate the applications and the cloud platform carefully from the data integrity viewpoint. This paper presents a model-based data integrity evaluation method using the UPPAAL model checker. In order to make the model reusable we built it as a set of application-independent functional modules. On the other hand, the application unique functionalities are totally included in the model as UPPAAL functions written by the C-like UPPAAL language. The data integrity evaluation is performed in two different ways. One is a simulation-based method in which the model is executed by the UPPAAL simulator to obtain the resultant variable values. The other is a verification-based method in which the given integrity constraints are examined by the UPPAAL verifier using full state space search of the model.

1 INTRODUCTION

Data integrity is one of the most critical concerns for distributed and concurrent systems especially for those in cloud environments e.g. Google App Engine (AE) Sanderson 2009 Amazon Web Services AWS van Lier and Paganelli 2011 or IBM Bluemix IBM 2015. One of the typical systems is database transaction processing and the data integrity becomes a crucial issue to make such systems robust. Ishida and Shinawa 2014.

Therefore the evaluation of the data integrity from both application and platform viewpoints for transaction processing in the cloud seems important to the spread of cloud computing. However, there are several difficulties in evaluating and validating this data integrity for transaction processing. The above difficulties are mainly caused by the different principles of the data integrity from traditional transaction processing which is adopted by the cloud.

This new and different principle is referred to as ASE standing for **B**asically **A**vailable **S**oft **s**tate and **E**ventually **c**onsistent Pritchett 2008. The basic differences between the ASE and the traditional principle ACID¹ Gray and Reuter 1993 both of which are a set of properties to be satisfied in order to

¹ ACID is standing for **A**tomicity **C**onsistency **I**solation and **D**urability.

guarantee the data integrity in transaction processing are

- 1 While the ACID restricts the concurrent data access accesses within a critical section, the ASE allows arbitrary concurrent data access accesses from any transaction by **Basically Available** property.
- 2 While the ACID postulates the transparent replication of the data bases, the ASE tolerates the non-transparent replication by **Soft state** property.
- 3 While the ACID aims at the data integrity at every instant, the ASE tries for achieving the data integrity within some duration by **Eventually consistent**.

According to the above differences between these two principles namely ASE and ACID we need a different approach to evaluating the data integrity in the cloud.

Since this evaluation must be performed before system implementation we need a precise model that reflects the cloud platform mechanism implementing the ASE principle along with the detailed application logic that determines the data values. The reason why is that the data integrity of transaction processing is affected by both of the

however, most modeling tools are specialized to a specific aspect of a system e.g. software specification languages like Spivey 2008. D

Fitgerald et al. 2004 and so on which are specialised to the functional aspect system modeling tools like finite state machines Petri nets² Reisig 1985 SDL Thiel 2001 and so on which are specialised to the behavioral aspect and architecture oriented modeling tools like UML class diagrams local diagrams and so on which are specialised to the structural aspect

On the other hand it is desirable to express multiple aspects of a system simultaneously in a single model for accurate evaluation of the data integrity. For this purpose we use the UPPAAL model checker David et al. 2015 as a modeling and evaluation tool since it can express the behavior of a system as a set of timed automata connected through communication channels along with the functional and data structure specifications using a C-like language provided by the tool.

The rest of the paper is organized as follows. In section 2 we introduce the basic concepts of the data integrity in the cloud along with the transaction behavior following the ASE principle. Section 3 shows how transaction processing in the cloud is modeled using the UPPAAL. Section 4 discusses an evaluation and validation method for the data integrity using the UPPAAL.

2 TRANSACTION PROCESSING IN THE CLOUD

Data integrity in transaction processing has been hitherto relying on the ACID principle that is guaranteed by a transaction processing monitor (TPM) under which they are running. One of the advantages of the ACID is that the serialised execution of transactions always maintains the data integrity. Therefore at the implementation level the TPM isolates and serialises the critical sections of each transaction by *locking* mechanism. In addition the ACID implicitly presumes the transparent replication or synchronous replication of data bases to realise *Consistency* property of it.

This approach could cause the reduction of data base availability along with the performance degradation of transaction processing. In cloud computing the ACID principle becomes a burden too much to guarantee the high availability scalability and stable performance of a system. Therefore more light-weight mechanisms to maintain the data integrity is desired in the cloud.

²Except for higher order Petri nets like Coloured Petri nets Jensen and Kristensen 2009

The ASE principle is a newly introduced principle to compromise the conflicting requirements that is availability and integrity in the cloud. In order to improve the availability the ASE principle does not serialise the critical sections of each transaction and allows the non-transparent or asynchronous replication. For maintaining the data integrity in such an environment a TPM following the ASE principle provides us with version information instead of a locking mechanism in order to determine whether the referred data are valid. If some of the referred data are invalid the relevant transaction aborts the database updates. This mechanism is known as *optimistic locking*.

Before discussing the data integrity of ASE transactions³ we need to define the concept of data integrity rigorously in order to evaluate it effectively. The term integrity or data integrity is used differently in various contexts. For example it focuses on the relationships between directories and file allocation information e.g. i-node in the case of UNIX at the operating system level while it means the referential integrity that requires the existence of specific key values at the *Database Management System* level.

On the other hand at the application level there are no commonly recognised definitions since it depends on the semantics of the data rather than their structure. Therefore it seems more difficult to express the data integrity at this level than the former two levels. In order to determine whether an application can be performed in the cloud in the form of a transaction we have to evaluate the data integrity at the application level in this circumstance.

Consequently we first need to define rigorously the concept of data integrity at the application level using a unified notation. The data integrity at the application level can be defined as a set of constraints or rules on database occurrences. One of the ways to express these constraints is to use predicate logic formulae Shinjima 2012. In order to compose these logic formulae we first have to define the language \mathcal{L} and the structure \mathcal{S} to provide the syntax and semantics of the formulae.

The language \mathcal{L} stipulates the usage of symbols regarding constants variables functions predicates and logical operators. In the data integrity evaluation the \mathcal{L} deals with database related matters. Therefore each symbol for a variable or constant represents an entity or its value in the databases. As for functions and predicates there are two kinds of them that is database oriented and application oriented. The for-

³Transactions to be run under the control of a TPM implementing the ASE principle

er ones are the functions or predicates defined in a data base manipulation language like SQL. On the other hand, the latter ones are those used in a specific application domain e.g. production control, product management or customer management applications.

Therefore, we are to prepare the \mathcal{L} as composed of two parts: namely, the application independent part and application dependent part. While the former part can be reused among the different application domains, the latter need to be rebuilt every time a new application is dealt with. On the other hand, the structure \mathcal{S} consists of the domain of discourse \mathcal{D} and the interpretation I . All the objects that are referred to from the functions and predicates or assigned to variables and constants must be the elements of the above \mathcal{D} . In our case, this \mathcal{D} includes

- 1 all the data base instances D_i
- 2 all the data base records $r_j^{(i)}$ in each D_i and
- 3 all the attribute values $a_k^{(ij)}$ in each $r_j^{(i)}$

The interpretation I maps each symbol in the \mathcal{L} to an actual entity defined over the \mathcal{D} . Some of functions and predicates are predefined in a data base manipulation language e.g. SQL. Other symbols in the \mathcal{L} are defined during the modeling process discussed in the succeeding sections.

Using the above language \mathcal{L} and the structure \mathcal{S} , each constraint to be expressed as an integrity rule is represented by a standardized logic formula PCFP (Pre-negation Conjunctive Normal Form).

$$Q_1 \cdots Q_n \left(\bigvee_j \bigwedge_i P_{ij}(t_1^{(ij)} \cdots t_{m_{ij}}^{(ij)}) \right)$$

where Q_i is a variable with the quantifier \forall e.g. $\forall x_i$; P_{ij} is a predicate and $t_k^{(ij)}$ is a term composed of variables, constants and functions (Schoening 2008).

There are several kinds of constraints regarding data integrity e.g. restrictions on data values, existence of a record with some specificity or constraints on the values derived from a set of records.

However, any kind of those constraints can be expressed by the above predicate logic formulae in the form of PCFP.

3 MODELING THE TRANSACTION PROCESSING WITH THE BASE PRINCIPLE

Once the rules or constraints for data integrity are expressed in the form of predicate logic formulae, the next step is to model the transaction processing with

the ASE principle which updates the databases in the cloud. For this modeling, we use the UPPAAL model checker or the UPPAAL in short. The UPPAAL represents a system as a set of finite timed automata with variables along with the functions that manipulate the

Each timed automaton consists of states (*locations* in terms of the UPPAAL) and arcs (*edges* in terms of the UPPAAL) that represent the state transitions. Boolean expressions with clock type variables can be used as time constraints which are associated with any above stated *location* or *edge*. These timed automata are defined as parameterizable templates and must be instantiated by the system definition.

In order to reuse the models, these templates should be appropriately modularized. In our approach, the behavior of the transaction processing in the cloud is categorized into five types: namely, Initialization, Scheduling, Thread, Database and Replication. Each module works as follows:

- 1 The Initialization module sets up the databases to be used during the simulation. The databases are expressed as three-dimensional integer arrays. The first dimension represents the replication number, the second represents the record or row number, and the third represents the attributes in the database schema.
- 2 The Scheduling module sends a transaction to one of the Thread instances to process it. A transaction is expressed in the form of integer array, each element of which represents an argument or parameter to the transaction. These integer arrays compose a two-dimensional transaction list⁴.
- 3 The Thread module performs the functionality of each transaction. The functionality is determined by the transaction type and the specified arguments in the transaction list. The database update requests from a transaction are routed to the Database module through a UPPAAL channel.
- 4 The Database module is to be instantiated as many as database replications. Each instance reads and updates a specific replication of a database expressed in the form of an integer array.
- 5 The Replication module tries to keep the replicated databases identical in an asynchronous way.

⁴Since the UPPAAL allows only fixed size for arrays and each transaction type could require the different number of arguments, an individual two-dimensional array is defined for each transaction type independently.

implementing the *Soft state* property. This module is instantiated only once and deals with all the data bases and their replications.

In addition to the above modules we have to prepare several functions to take the model executable and verify it. These functions are written by a C-like UPPAAL unique language while the model structure is common among application domains. These functions are application unique and must be built for each application domain.

Figure 1 through Figure 5 show an example of the above UPPAAL modules. As stated above the structure of their respective modules can be commonly used among different application domains including function names and channels associated with edges and locations in the model. However, the implementation of these functions and other supplemental functions are differently built among different application domains.

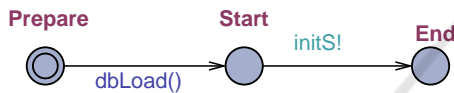


Figure 1 Initialization module

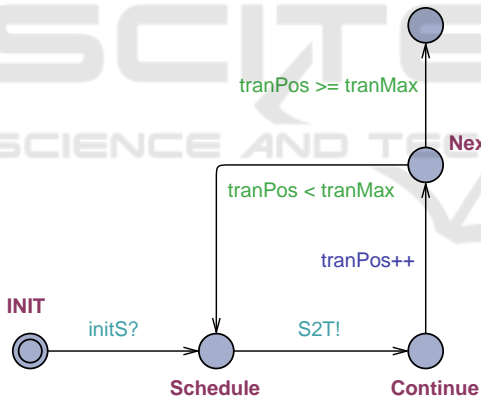


Figure 2 Scheduler module

For example, the function `dbLoad` in Figure 1 represents a function that initialize all the data bases in the system and the name is common for all applications. However, its implementation usually different among the depending on the structure and usage of the data bases. Figure 6 shows a sample implementation of the `dbLoad` function for a simplified library application.

When executing the model, these modules are instantiated through the system definition as shown in Figure 7. In this example, three concurrent threads and three data base replications are assumed.

These modules operate as follows:

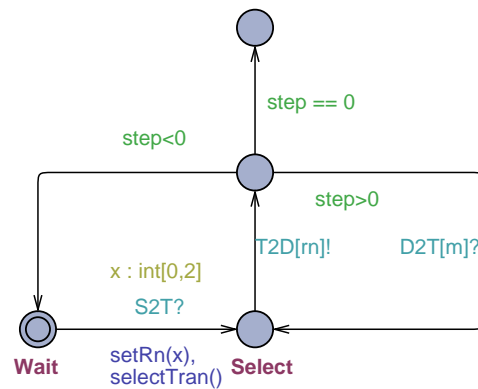


Figure 3 Thread module

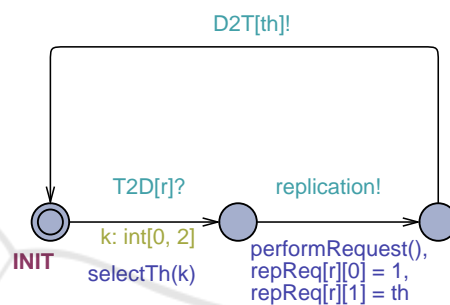


Figure 4 Database module

1. Firstly, the `dbLoad` function of the Initialization module is invoked to prepare all the data bases. At this time, only the associated edge is eligible for transition since other modules are waiting for signals through the UPPAAL channels.
2. After the completion of the `dbLoad` function, the Scheduler module is activated through the `initS` channel.
3. The scheduler module sends a signal to the Thread module through the channel `S2T`.
4. The Thread module selects a transaction from the predefined transaction list by the `selectTran` function and sends a signal to the Database module through the channel `T2D` where `th` is a replication number.
5. The Database module accesses and updates the data bases.

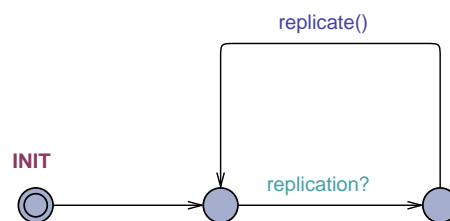


Figure 5 Replication module

```

void d Load {
  for i int 0 2 {
    for int 0 14 {
      int search oo i oo ey
      if > 0 {
        for int 0 3 oo i oo
          i
        }
      }
    }
  }

  for i int 0 2 {
    for int 0 8 {
      int searchAccount i account ey
      if > 0 {
        for int 0 3 account i
          account i
        }
      }
    }
  }

  for i int 0 2 {
    for int 0 4 {
      int searchLoan i loan ey 0 loan ey
      if > 0 {
        for int 0 3 loan i loan
          i
        }
      }
    }
  }
}
    
```

Figure 6 d Load Function

```

// Place template instantiations here.
In Initiator
Sc Scheduler
T1 Thread 0
T2 Thread 1
T3 Thread 2
D 1 D 0
D 2 D 1
D 3 D 2
EP eplication
// List one or more processes to be composed into
a system.
system In Sc T1 T2 T3 D 1 D 2 D 3 EP
    
```

Figure 7 System Definition for module Instantiation

6 The step 2 to 5 are repeated until all the predefined transactions are processed

The version control and commit a port processes are embedded in the Database module as functions

4 DATA INTEGRITY EVALUATION USING THE UPPAAL

The UPPAAL model checker provides us with three major functionalities. The first is a graphical model editor with programming capability that we have used in the previous section. The second is a model simulator that executes the model we build to show an instance of its behavior. The third is a model verifier that examines all the possible behavior whether the model satisfies the given properties written in the form of CTL (Computational Logic Tree for Ulae).

Therefore two alternative ways are available to evaluate the data integrity of transaction processing. The first is to execute the model to obtain the values of the variables for the database records at each state transition. As discussed in the previous section the data integrity is expressed as a set of predicate logic for Ulae in the form of PCF. In the UPPAAL model these logic for Ulae refer to the variables associated with the database records and attributes. Therefore we can determine whether the data integrity is maintained in the transaction processing by examining the above variables using a function implementing each constraint logic for Ulae. Since this method can evaluate only one instance of the system behavior selected by the simulation we have to perform the simulation for every possible behavior. However this possible behavior could be uncountable. Therefore this method would be a sampling based evaluation.

On the other hand the UPPAAL verifier provides us with a capability of full state space search against a set of CTL for Ulae. In order to evaluate the data integrity in this way we have to transform a set of predicate logic for Ulae into a set of CTL for Ulae. Unlike the predicate logic for Ulae CTL for Ulae can include the path operator A and E which deal with state transition paths of a system and temporal operator \square and \diamond which define the validation points of the for Ulae. In addition there are no quantifiers \forall and \exists in CTL. Therefore several considerations should be taken into account in the above transformation from predicate logic for Ulae into CTL for Ulae. These considerations include

- 1 If a property P must always holds in a predicate logic for Ulae the CTL for Ulae is $A\square P$
- 2 If a property P always implies a property Q then the CTL for Ulae is $A\square(P \rightarrow Q)$
- 3 If a property P eventually implies a property Q the CTL for Ulae is $A\square(P \rightarrow \diamond Q)$

- 4 If a property P must hold at specific point we introduce a boolean variable to express the point and set it *true* at the point in the model. In this case we need to modify the model.
- 5 If the original predicate logic for formula includes the quantifiers \forall and \exists we introduce a boolean function into the model to evaluate whether all of or some of the variables in the model satisfy the formula. A model modification is required in this case again.

After the above transformation is completed we can evaluate the data integrity by running the verifier that the UPPAAL provides.

This CTL based evaluation seems simpler than the simulator based one however it performs full state space search and consumes huge computing resources. As a result it takes long time to obtain the result. In such cases we need to reduce the model by decreasing the number of variables or values to be assigned.

5 CONCLUSIONS

In cloud environments the behavior of transaction processing is considerably different from the traditional ones. One of the major reasons is that the cloud introduces a new principle for the data integrity called ASE instead of the traditional ACID. In order to analyze the transaction processing state in the cloud we need to reveal the behavior of it clearly and evaluate the data integrity rigorously.

This paper proposed a model based data integrity evaluation using the UPPAAL model checker. In order to analyze the model easily understandably and reusable we composed it using five functional modules namely Initialization, Scheduling, Thread Data base and replication following the ASE principle. While the model structure can be reused among different application domains we need to build application unique functions for the model.

The UPPAAL provides us with two different ways to evaluate the data integrity. One is a simulation based evaluation which evaluates only one instance of the behavior of transaction processing. The other is a verifier based evaluation which evaluates full state space search to determine whether the given constraints are satisfied while the latter way can evaluate the integrity more precisely we need to transform the original predicate logic for formulae into the CTL formulae. In addition it consumes huge computing resources for full state space search and takes long time to obtain the evaluation results.

ACKNOWLEDGEMENTS

This work was supported by SPS A E I grant number 25330094.

REFERENCES

- David A. Larsen and Legay A. 2015. UPPAAL S.C tutorial. In *International Journal on Software Tools for Technology Transfer Volume 17 Issue 4* pages 397-415. Springer.
- Fitzgerald, Larsen P., Uher, P., Plat and Erhoef. 2004. *Validated Designs for Object-oriented Systems*. Springer.
- Gray and Reuter A. 1993. *Transaction Processing: Concepts and Techniques*. Morgan Kaufmann.
- IBM. 2015. *Bluemix*. <http://www.ibm.com/cloud-computing/blue>.
- Jensen and Kristensen L. 2009. *Coloured Petri Nets: Modeling and Validation of Concurrent Systems*. Springer-Verlag.
- Shimada S. and Shinawa. 2014. Data Integrity in Cloud Transactions. In *Proc. 4th International Conference on Cloud Computing and Services Science* pages 457-462.
- Pritchett D. 2008. ASE: An ACID alternative. In *ACM QUEUE Volume 6 Issue 3* pages 48-55. ACM.
- Eisig. 1985. *Petri Nets: An Introduction*. Springer.
- Sanderson D. 2009. *Programming Google App Engine*. O'Reilly Associates Inc.
- Schoening U. 2008. *Logic for Computer Scientists (Modern Birkhauser Classics)*. Birkhauser Boston.
- Shinawa. 2012. CP based Data Integrity Evaluation for Cloud Transactions. In *Proc. 6th International Conference on Software Paradigm Trends* pages 267-272.
- Spivey. 2008. *Understanding Z: A Specification Language and its Formal Semantics*. Cambridge University Press.
- Thiel A. 2001. *Systems Engineering with SDL: Developing Performance-Critical Communication Systems*. Wiley.
- van Hilt and Paganelli F. 2011. *Programming Amazon EC2*. O'Reilly Associates Inc.