

Runtime Prediction of Service Level Agreement Violations for Composite Services

Philipp Leitner¹, Branimir Wetzstein², Florian Rosenberg³, Anton Michlmayr¹, Schahram Dustdar¹, and Frank Leymann²

¹ Distributed Systems Group
Vienna University of Technology
Argentinierstrasse 8/184-1
A-1040, Vienna, Austria
`lastname@infosys.tuwien.ac.at`

² Institute of Architecture of Application Systems
University of Stuttgart
Stuttgart, Germany
`lastname@iaas.uni-stuttgart.de`

³ CSIRO ICT Centre
GPO Box 664
Canberra ACT 2601, Australia
`florian.rosenberg@csiro.au`

Abstract. SLAs are contractually binding agreements between service providers and consumers, mandating concrete numerical target values which the service needs to achieve. For service providers, it is essential to prevent SLA violations as much as possible to enhance customer satisfaction and avoid penalty payments. Therefore, it is desirable for providers to predict possible violations before they happen, while it is still possible to set counteractive measures. We propose an approach for predicting SLA violations at runtime, which uses measured and estimated facts (instance data of the composition or QoS of used services) as input for a prediction model. The prediction model is based on machine learning regression techniques, and trained using historical process instances. We present the basics of our approach, and briefly validate our ideas based on an illustrative example.

1 Introduction

In service-oriented computing [1], finer-grained basic functionality provided using Web services can be composed to more coarse-grained services. This model is often used by Software-as-a-Service providers to implement value-added applications, which are built upon existing internal and external Web services. Very important for providers and consumers of such services are Service Level Agreements (SLAs), which are legally binding agreements governing the quality that the composite service is expected to provide (Quality of Service, QoS). SLAs contain Service Level Objectives (SLOs), which are concrete numerical target values (e.g., “maximum response time is 45 seconds”). For the provider

it is essential to not violate these SLOs, since typically violations are coupled with penalty payments. Additionally, violations can negatively impact service consumer satisfaction. Therefore, it is vitally important for the service provider to be aware of SLA violations, in order to react to them accordingly.

Typically, SLA monitoring is done *ex post*, i.e., violated SLOs can only be identified after the violation happened. While this approach is useful in that it alerts the provider to potential quality problems, it clearly cannot directly help preventing them. In that regard an *ex ante* approach is preferable, which allows to predict possible SLA violations before they have actually occurred. The main contribution of this paper is the introduction of a general approach to prediction of SLA violations for composite services, taking into account both QoS and process instance data, and using estimates to approximate not yet available data. Additionally, we present a prototype implementation of the system and an evaluation based on an order processing example. The ideas presented here are most applicable for long-running processes, where human intervention into problematic instances is possible. Our system introduces the notions of checkpoints (points in the execution of the composition where prediction can be done), facts (data which is already known in a checkpoint, such as the response times of already used services) and estimates (data which is not yet available, but can be estimated). Facts and estimates can refer to both typical QoS data (e.g., response times, availability, system load) and process instance data (e.g., customer identifiers, ordered products). Our implementation uses regression classifiers, a technique from the area of machine learning, to predict concrete SLO values.

2 Illustrative Example

To illustrate the ideas presented in this paper we will use a simple purchase order scenario (see Figure 2 below). In this example there are a number of roles to consider (reseller, customer, banking service, shipping service, and two external suppliers). Whenever the reseller service receives an order from the customer, it first checks if all ordered items are available in the internal stock. If this is not the case, it checks if the missing item(s) can be ordered from Supplier 1, and, if this is not the case, from Supplier 2. If both cannot deliver the order has to be cancelled, otherwise the missing items are ordered from the respective supplier. When all ordered items are available she will (in parallel) proceed to charge the customer using the banking service and initialize shipment of the ordered goods (using the Shipping Service). Please refer to [2] for more details on this case.

In this case study, the reseller has an SLA with its customers, with an SLO specifying that the end-to-end response time of the composition cannot be more than a certain threshold of time units. For every time the SLO is violated the customer is contractually entitled a discount for the order. Note that even though our explanations in this paper will be based on just one single SLO, our approach can be generalized to multiple SLOs. Additionally, even though we present our approach based on a numerical SLO, our ideas can be also applied to estimation of nominal objectives.

3 Predicting SLA Violations

In this section we present the core ideas of our approach towards prediction of SLA violations. Generally, the approach is based on the idea of predicting concrete SLO values based on whatever information is already available at a concrete point in the execution of a composite service. We distinguish three different types of information. (1) **Facts** represent data which is already known at prediction time. Typical examples of facts are the QoS of already used services, such as the response time of a service which has already been invoked in this execution, or instance data which has either been passed as input or which has been generated earlier in the process execution. (2) **Unknowns** are the opposites of facts, in that they represent data which is entirely unknown at prediction time. Oftentimes, instance data which has not yet been produced falls into this category. If important factors are unknown at prediction time the prediction quality will be very bad, e.g., in our illustrative example a prediction cannot be accurate before it is known whether the order can be delivered from the reseller’s internal stock. (3) **Estimates** are a kind of middle ground between facts and unknowns, in that they represent data which is not yet available, but can be estimated. This is often the case for QoS data, since techniques such as QoS monitoring [3] can be used to get an idea of e.g., the response time of a service before it is actually invoked. Estimating instance data is more difficult, and generally domain-specific.

The overall architecture of our system is depicted in Figure 1. The most important concept used is that the user defines **checkpoints** in the service composition, which indicate points in the execution where a prediction should be carried out. The exact point in the execution model which triggers the checkpoint is called the **hook**. Every checkpoint is associated with one **checkpoint predictor**. Essentially, the predictor uses a function taking as input all facts which are already available in the checkpoint, and, if applicable, a number of estimates of not yet known facts, and produces a numerical estimation of the SLO

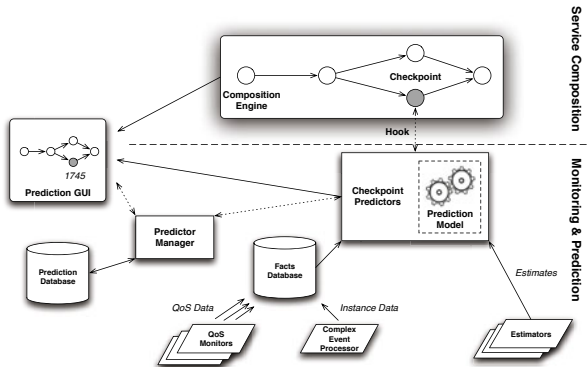


Fig. 1. Overall System Architecture

value(s). This function is generated using machine learning techniques. We refer to this function as the **prediction model** of a checkpoint predictor. Facts are retrieved from a **facts database**, which is filled using a number of **QoS monitors** (which provide QoS data) and a **Complex Event Processing (CEP)** engine (which extracts and correlates the instance data, as emitted by the process engine). A detailed discussion of our event-based approach to monitoring is out of scope of this paper, but can be reviewed in related work [2, 4]. **Estimators** are a generic framework for components which deliver estimates. Finally, the prediction result is transferred to a graphical user interface (**prediction GUI**), which visualizes the predicted value(s) for the checkpoint. A **predictor manager** component is responsible for the lifecycle management of predictors, i.e., for initializing, destroying and retraining them. Additionally, predictions are stored in a **prediction database** to be available for future analysis.

3.1 Checkpoint Definition

At design-time, the main issue is the definition of checkpoints in the composition model. For every checkpoint, the following input needs to be provided: (1) The hook, which defines the concrete point in the execution that triggers the prediction, (2) a list of available facts, (3) a list of estimates, and the estimator component as well as the parameters used to retrieve or calculate them, (4) the retraining strategy, which governs at which times a rebuilding of the prediction model should happen, and (5) as a last optional step, a parameterization of the machine learning technique used to build the prediction model. After all these inputs are defined the checkpoint is deployed using the predictor manager, and an initial model is built. For this a set of historical executions of the composite service need to be available, for which all facts (including those associated with estimates) have been monitored. If no or too little historical data is available the checkpoint is suspended by the predictor manager until enough training data has been collected. The amount of data necessary is case-specific, since it vastly depends on the complexity of the composition. We generally use the Training Data Correlation as a metric for evaluating the quality of a freshly trained model (see below for a definition), however, a detailed discussion of this is out of scope of this paper. After the initial model is built the continuous optimization of the predictor is governed by the predictor manager, according to the retraining strategy. Finally, the checkpoint can be terminated by the user via the prediction GUI. We will now discuss these concepts in more depth.

Hooks. Hooks can be inserted either before or after any WS-BPEL activity (for instance, an *Invoke* activity). Generally, there is a tradeoff to take into account here, since early predictions are usually more helpful (in that they rather allow for corrections if violations are predicted), but also less accurate since less facts are available and more estimates are necessary. Figure 2 depicts the (simplified) example from Section 2, and shows two possible checkpoints. In C_1 the only facts available are the ones given as input to the composition (such as a customer identifier, or the ordered products). Some other facts (mainly QoS metrics) can already

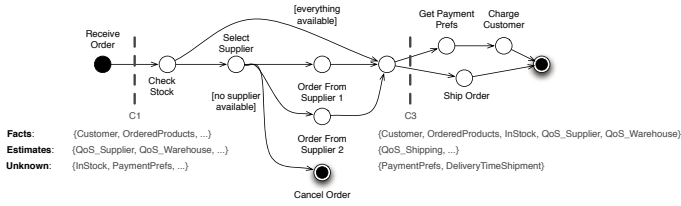


Fig. 2. Illustrative Example With Possible Checkpoints

be estimated, however, other important information, such as whether the order can be served directly from stock, is simply unavailable in C_1 , not even as an estimate. Therefore, the prediction cannot be very accurate. In checkpoint C_3 , on the other hand, most of the processes important raw data is already available as facts, allowing for good predictions. However, compared to C_1 , the possibilities to react to problems are limited, since only the payment and shipping steps are left to adapt (e.g., a user may still decide to use express shipping instead of the regular one if a SLA violation is predicted in C_3). Finding good checkpoints at which the prediction is reasonably accurate and still timely enough to react to problems demands for some domain knowledge about influential factors of composition performance. Dependency analysis as discussed in [4] can help providing this crucial information. Dependency analysis is the process of using historical business process instance data to find out about the main factors which dictate the performance of a process. When defining checkpoints, a user can assume that the majority of important factors of influence need to be available as either facts or at least as good estimates in order to achieve accurate predictions.

Facts and Estimates: Facts represent all important information which can already be measured in this checkpoint. This includes both QoS and instance data. Note that the relationship between facts and the final SLO values does not need to be known (e.g., a user can include instance data such as user identifiers or ordered items, even if she is not sure if this has any relevance for the SLO). However, dependency analysis can again be used to identify the most important facts for a checkpoint. Additionally, the user can also define estimates. In the example above, in C_1 the response time of the warehouse service is not yet known, however, it can e.g., be estimated using a QoS monitor. Since estimating instance data is inherently domain-specific, our system is extensible in that more specific estimators (which are implemented as simple Java classes) can be integrated seamlessly. Estimates are linked to facts, in the sense that they have to represent an estimation of a fact which will be monitorable at a later point.

Retraining Strategy: Generally, the prediction model needs to be rebuilt whenever enough new information is available to significantly improve the model. The retraining strategy is used to define when the system should check whether rebuilding the prediction model is necessary. We currently support the following retraining strategies: (1) periodic retraining (retraining is done in fixed intervals),

(2) instance-based retraining (retraining is done whenever a certain amount of new instances have been monitored), (3) on demand retraining (retraining only on explicit user demand), (4) on error retraining (retraining when the mean prediction error \bar{e} exceeds a given threshold), and (5) custom retraining (retraining based on user-defined conditions).

Prediction Model Parameterization: A user can also define the machine learning technique that should be used to build the prediction model. This is done by specifying an algorithm and the respective parameterization for the WEKA toolkit¹, an open source machine learning toolkit which we internally use in our prototype implementation. In this way the prediction quality can be tuned by a machine learning savvy user, however, we also provide a default configuration which can be used out of the box.

3.2 Run-Time Prediction

At runtime, the prediction process is triggered by lifecycle events from the WS-BPEL engine. These are events emitted by some engines, which contain lifecycle information about the service composition. Our approach is based on these events, therefore, a WS-BPEL engine which is capable of emitting these events is a preliminary of our approach. When checkpoints are deployed we use the hook information to register respective event listeners. For instance, for a checkpoint with the hook “After invoke CheckStock” we generate a listener for `ActivityExecEndEvents` which consider the invoke activity “CheckStock”. We show the sequence of actions which is triggered as soon as such an event is received in Figure 3.

After being triggered by a lifecycle event the checkpoint predictor first extracts some necessary correlation information from the event received. This includes the process instance ID as assigned by the composition engine, the instance start time (i.e., the time when the instance was created) and the timestamp of

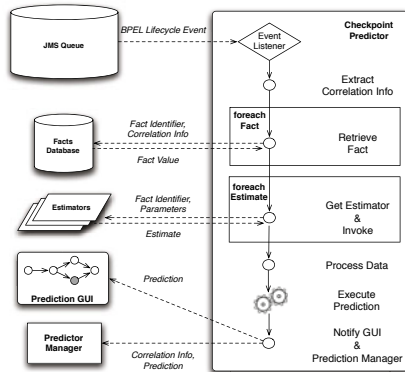


Fig. 3. Runtime View On Checkpoint Predictors

¹ <http://www.cs.waikato.ac.nz/ml/weka/>

the event. This information is necessary to be able to retrieve the correct facts from the facts database, which is done for every fact in the next step (e.g., in order to find the correct fact “CustomerNumber” for the current execution the process instance ID needs to be known). When all facts have been gathered, the predictor also collects the still missing estimates. For this, for every estimate the predictor instantiates the respective estimator component (if no instance of this estimator was available before), and invokes it (passing all necessary parameters as specified in the checkpoint definition). The gathered facts and estimates are then converted into the format expected by the prediction model (in the case of our prototype, this is the WEKA Attribute-Relation File Format ARFF²), and, if necessary, some data cleaning is done. Afterwards, the actual prediction is carried out by passing the gathered input to the prediction model producing a numerical estimation of the SLO value. This prediction is then passed to the prediction GUI (for visualization) and the prediction manager.

Note that the “intelligence” that actually implements the prediction of the SLO values is encapsulated in the prediction model. Since we (usually) want to predict numerical SLO values the prediction model needs to be a regression model. We consider the regression model to be a black-box function which takes a list of numeric and nominal values as input, and produces a numeric output. Generally, our approach is agnostic of how this is actually implemented. In our work we use multilayer perceptrons (a powerful variant of neural networks) to implement the regression model. Multilayer perceptrons are trained iteratively using a back-propagation technique (maximization of the correlation between the actual outcome of training instances and the outcome that the network would predict on those instances), and can (approximately) represent any relationship between input data and outcome.

3.3 Evaluation of Predictors

Another important task of the prediction manager is quality management of predictors, i.e., continually supervising how predictions compare to the actual SLO values once the instance is finished. Generally, we use three different quality metrics to measure the quality of predictions in checkpoints. The first metric, *Training Data Correlation* ($corr = \frac{cov(P,M)}{\sigma_p \sigma_m}$), is a standard machine learning approach to evaluating regression models. We use it mainly to evaluate freshly generated models, when no actual predictions have yet been carried out. This metric is defined as the statistical correlation between all training instance outcomes and the predictions that the model would deliver for these training instances. The definition given is the standard statistical definition of the correlation coefficient between a set of predicted values P and a set of measured values M . However, note that this metric is inherently overconfident in our case, since during training all estimates are replaced for the facts that they estimate (i.e., the training is done as if all estimates were perfect). Therefore, we generally measure the

² <http://www.cs.waikato.ac.nz/~ml/weka/arff.html>

prediction error later on, when actual estimates are being used. However, a low training data correlation is an indication that important facts are still unknown in the checkpoint, i.e., that the checkpoint may be too early. The actual quality of the prediction is measured using the *Mean Prediction Error* $\bar{e} = \frac{\sum_{i=0}^n |m_i - p_i|}{n}$, which is the average difference between predicted and monitored values. n is the total number of predictions, p_i is a predicted value, and m_i is the measured value to prediction p_i . Finally, we use the *Prediction Error Standard Deviation* ($\sigma = \sqrt{\frac{\sum_{i=0}^n (e_i - \bar{e})^2}{n}}$) to describe the variability of the prediction error (i.e., high σ essentially means that the actual error for an instance can be much lower or higher than \bar{e}). e_i is the actual prediction error for a process instance ($m_i - p_i$). These metrics are mainly used to give the user an estimation of how trustworthy a given prediction is. Additionally, the **on error** retraining strategy triggers on \bar{e} exceeding a certain threshold.

4 Experimentation

In order to provide a first validation of the ideas presented we have implemented the illustrative example as discussed in Section 2, and ran some experiments using our prototype tool. All experiments have been conducted on a single test machine with 3.0 GHz and 32 GByte RAM, running under Windows Server 2007 SP1.

For reasons of brevity we focus on prediction accuracy in this paper. To measure accuracy, we have implemented five checkpoints in the illustrative example (see top of Figure 4): C1 is located directly after the order is received, C2 after the internal warehouse is checked, C3 after eventual orders from external suppliers have been carried out, C4 during the payment and shipment process, and finally C5 when the execution is already finished. In each of those checkpoints we have trained a prediction model using 1000 historical process instances, and

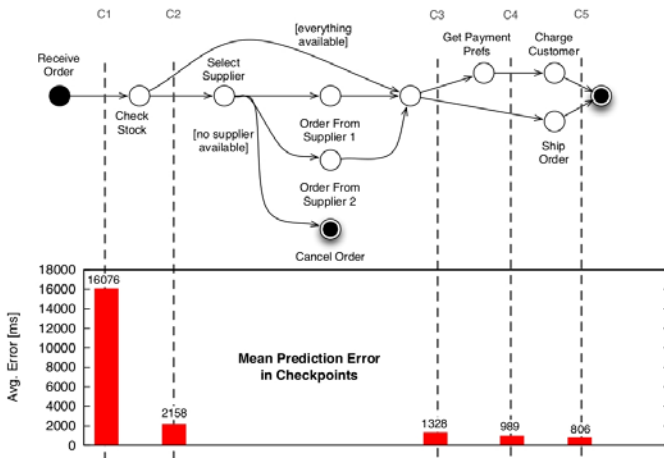


Fig. 4. Prediction Error in Checkpoints

have specified all available data as facts. For not yet available QoS metrics we have used the average of all previous invocations as estimate. Missing instance data has been treated as unknown. We have used each of those checkpoints to predict the outcome of 100 random executions, and calculated \bar{e} . As expected, \bar{e} is decreasing with the amount of factual data available. In C1, the prediction is mostly useless, since no real data except the user input is available. However, in C2 the prediction is already rather good. This is mostly due to the fact that in C2 the information whether the order can be delivered directly from stock is already available. In C3, C4 and C5 the prediction is continually improving, since more actual QoS facts are available, and less estimates are necessary. Speaking in absolute values, \bar{e} in e.g., C3 is 1328 ms. Since the average SLO value in our illustrative example was about 16000 ms, the error represents only about 8% of the actual SLO value, which seems satisfactory. We have also carried out experiments to determine the time necessary to build prediction models and to do predictions at runtime. These results are not discussed here because of space restrictions.

5 Related Work

The work presented in this paper is complementary to the more established concept of SLA management. SLA management incorporates the definition and monitoring of SLAs, as well as the matching of consumer and provider templates [5, 6]. In our work we add another facet to the more general area of SLA management, namely the prediction of SLA violations before they have actually occurred. Inherently, this prediction demands for some insight into the internal factors impacting composite service performance. In [7], the MoDe4SLA approach has been introduced to model dependencies of composite services on the used base services, and to analyze the impact that these dependencies have. Similarly, the work we have presented in [2] allows for an analysis of the impact that certain factors have on the performance of service compositions. SLA prediction as discussed in this paper has first been discussed in [8], which is based on some early work of HP Laboratories on SLA monitoring for Web services [9]. In [8], the authors introduced some concepts which are also present in our solution, such as the basic idea of using prediction models based on machine learning techniques, or the trade-off between early prediction and prediction accuracy. However, the authors do not discuss important issues such as the integration of instance and QoS data, or strategies for updating prediction models. Additionally, this work does not take estimates into account, and relatively little technical information about their implementation is publicly available. A second related approach to QoS prediction has been presented recently in [10]. In this paper the focus is on KPI prediction using analysis of event data. Generally, this work exhibits similar limitations as the work described in [8], however, the authors discuss the

influence of seasonal cycles on KPIs. This facet has not been examined in our work, even though seasons can arguably be integrated easily in our approach as additional facts.

6 Conclusions

In this paper we have presented an approach to runtime prediction of SLA violations. Central to our approach are checkpoints, which define concrete points in the execution of a composite service at which prediction has to be carried out, facts, which define the input of the prediction, and estimates, which represent predictions about data which is not yet available in the checkpoint. We use techniques from the area of machine learning to construct regression models from recorded historical data to implement predictions in checkpoints. Retraining strategies govern at which times these regression models should be refreshed. Our Java-based implementation uses the WEKA Machine Learning framework to build regression models. Using an illustrative example we have shown that our approach is able to predict SLO values accurately.

Acknowledgements

The research leading to these results has received funding from the European Community's Seventh Framework Programme [FP7/2007-2013] under grant agreement 215483 (S-Cube).

References

1. Papazoglou, M.P., Traverso, P., Dustdar, S., Leymann, F.: Service-Oriented Computing: State of the Art and Research Challenges. *IEEE Computer* 11 (2007)
2. Wetzstein, B., Leitner, P., Rosenberg, F., Brandic, I., Leymann, F., Dustdar, S.: Monitoring and Analyzing Influential Factors of Business Process Performance. In: *EDOC 2009: Proceedings of the 13th IEEE International Enterprise Distributed Object Computing Conference* (2009)
3. Michlmayr, A., Rosenberg, F., Leitner, P., Dustdar, S.: Comprehensive QoS Monitoring of Web Services and Event-Based SLA Violation Detection. In: *MW4SOC 2009: Proceedings of the 4th International Workshop on Middleware for Service Oriented Computing* (2009)
4. Wetzstein, B., Strauch, S., Leymann, F.: Measuring Performance Metrics of WS-BPEL Service Compositions. In: *ICNS 2009: Proceedings of the Fifth International Conference on Networking and Services*. IEEE Computer Society, Los Alamitos (2009)
5. Dan, A., Davis, D., Kearney, R., Keller, A., King, R., Kuebler, D., Ludwig, H., Polan, M., Spreitzer, M., Youssef, A.: Web Services on Demand: WSLA-Driven Automated Management. *IBM Systems Journal* 43(1), 136–158 (2004)
6. Tomic, V., Pagurek, B., Patel, K., Esfandiari, B., Ma, W.: Management applications of the web service offerings language (wsol). *Information Systems* 30(7), 564–586 (2005)

7. Bodenstaff, L., Wombacher, A., Reichert, M., Jaeger, M.C.: Monitoring Dependencies for SLAs: The MoDe4SLA Approach. In: SCC 2008: Proceedings of the 2008 IEEE International Conference on Services Computing, pp. 21–29 (2008)
8. Castellanos, M., Casati, F., Dayal, U., Shan, M.C.: Intelligent Management of SLAs for Composite Web Services. In: Bianchi-Berthouze, N. (ed.) DNIS 2003. LNCS, vol. 2822, pp. 28–41. Springer, Heidelberg (2003)
9. Sahai, A., Machiraju, V., Sayal, M., Moorsel, A.P.A.v., Casati, F.: Automated SLA Monitoring for Web Services. In: Feridun, M., Kropf, P.G., Babin, G. (eds.) DSOM 2002. LNCS, vol. 2506, pp. 28–41. Springer, Heidelberg (2002)
10. Zeng, L., Lingenfelder, C., Lei, H., Chang, H.: Event-Driven Quality of Service Prediction. In: Bouguettaya, A., Krueger, I., Margaria, T. (eds.) ICSOC 2008. LNCS, vol. 5364, pp. 147–161. Springer, Heidelberg (2008)