

Comprehensive Explanation of SLA Violations at Runtime

C. Müller⁽¹⁾, M. Oriol⁽²⁾, X. Franch⁽²⁾, J. Marco⁽²⁾, M. Resinas⁽¹⁾, A. Ruiz-Cortés⁽¹⁾, M. Rodríguez⁽²⁾

⁽¹⁾University of Seville, LSI, Seville (Spain), ISA research group, <http://www.isa.us.es/>

{cmuller, resinas, aruiz}@us.es

⁽²⁾Universitat Politècnica de Catalunya, Barcelona (Spain), GESSI research group, <http://www.essi.upc.edu/~gessi/>

{moriol, jmarco}@lsi.upc.edu, {franch, marcr}@essi.upc.edu

Abstract—

Service Level Agreements (SLAs) establish the Quality of Service (QoS) agreed between service-based systems consumers and providers. Since the violation of such SLAs may involve penalties, quality assurance techniques have been developed to supervise the SLAs fulfillment at runtime. However, existing proposals present some drawbacks: (1) the SLAs they support are not expressive enough to model real-world scenarios, (2) they couple the monitoring configuration to a given SLA specification, (3) the explanations of the violations are difficult to understand and even potentially inaccurate, (4) some proposals either do not provide an architecture, or present low cohesion within their elements. In this paper, we propose a comprehensive solution, from a conceptual reference model to its design and implementation, that overcomes these drawbacks. The resulting platform, SALMonADA, receives the SLA agreed between the parties as input and reports the explanations of SLA violations in a timely and highly understandable way. SALMonADA performs an automated monitoring configuration and it analyses highly expressive SLAs by means of a constraint satisfaction problems based technique. We have evaluated the impact of SALMonADA over the resulting service consumption time performance. The results are satisfactory enough to consider SALMonADA for SLA supervision because of its low intrusiveness.

Keywords-service level agreement; SLA; monitoring; analysis; violation detection and explanation; QoS.

I. INTRODUCTION AND MOTIVATION

Service Level Agreements (SLAs) establish the Quality of Service (QoS) agreed between service-based systems consumers and providers, as well as penalty/reward clauses to be applied when the SLA is violated. Detecting such violations may be complicated, consider for instance the following case.

The Amazon Simple Storage Service (AmazonS3) guarantees a *monthly uptime percentage* equal to or greater than 99.9% in its SLA¹ including a clause to reward their consumers against a lack of service. However, Amazon requires the customer to proof this violation by sending an email within ten business days after the end of the billing cycle in which the errors occurred. Proving the violation demands the customer to compute the monthly uptime percentage by subtracting from 100% the average of the *error rates*².

This Amazon scenario illustrates the need of having techniques to supervise the fulfillment of SLAs [1]. In response to

this need, several quality assurance proposals have been formulated. They deal with aspects such as: *violation detection*, either at runtime [2]–[8] or testing time [9], [10]; *violation notification* either by *push* (i.e. notifications to subscribed clients as soon as violations are detected at runtime) [4], [5], [11] or *pull* strategies [2], [8], [12]; and *violation explanation* [13], [14]. The possible adoption of these proposals in real-world scenarios is greatly influenced by the following factors:

- 1) Which SLAs are supported. Usually, real-world SLAs (e.g. AmazonS3 SLA) describe the parties obligations in natural language and they may comprise complex elements such as: (1) conditional terms subject to a precondition (i.e. if the precondition holds, then the term applies), for instance, AmazonS3 SLA offers a data durability of 99.99% only if the client choose a cheaper reduced redundancy storage; (2) optional or alternative terms (i.e. a set of terms that can be chosen by the customers), for instance, AmazonS3 offers a set of alternative support plan terms.
- 2) How the interaction with the monitor is performed. In order to detect violations we need to monitor the QoS offered by the providers at runtime (i.e. specific values for monitorable service properties such as the service availability or response time). The interaction with the monitoring techniques [15] that gather such QoS monitoring results, requires both: (a) to configure the monitor with the location of the service, and its monitorable service properties; and (b) to decide how the QoS monitoring results are specified. For instance, in the case of AmazonS3, the SLA determines the need for monitoring and obtaining a specific value for the monthly uptime percentage for each billing cycle in order to detect a possible violation.
- 3) How the violations are detected and explained. The SLA violations are detected by checking the agreed QoS against the obtained QoS monitoring result. And, ideally, the explanation of violations needs to be both comprehensive and timely. By comprehensive, we mean to provide a user-friendly and accurate violation report including not only the violated terms, but also the violation causes and even possible effects, in a form that is easy to understand by humans. For instance, in the AmazonS3 scenario, a comprehensive explanation

¹<http://aws.amazon.com/s3-sla/>

²internal server errors divided by the requests during each five minute period

would require a report including that the violated term is "monthly uptime percentage equal to or greater than 99.9%" by a monthly uptime percentage of 92% measured in a specific monthly billing cycle. By timely, we mean to communicate the violation as soon as it happens, i.e. when the QoS monitoring result has just been retrieved at runtime and it is checked against the SLA to detect and explain the violations.

- 4) Which architecture is proposed. In order to ease the maintainability and adaptability of the system, the architecture should comprise loosely coupled and highly cohesive elements separating monitoring from analysis.

As shown in Sec. II, current proposals do not satisfactorily deal with the characteristics above. Some of them may be completely coupled to a particular SLA notation that are not able to deal with all the aspects required in real-world scenarios, others couple the monitoring configuration to a given SLA specification, and most (if not all) provide very basic information when a violation occurs. Furthermore, not all proposals provide architectures with a clear separation of concerns between monitoring and analysis.

The goal of this work is to design and implement SALMonADA, a service-based system to monitor and analyse SLAs in order to provide timely and comprehensive explanations of violations. SALMonADA main features are:

- 1) specification of a wide range of SLA structures and complex elements based on the use of WS-Agreement [16].
- 2) automated monitoring configuration through the analysis of the SLAs but without coupling to a given SLA specification.
- 3) powerful detection and explanation of the SLAs violations by means of a Constraint Satisfaction Problem (CSP)-based technique. A comprehensive and timely explanation of the SLA violations is notified either by push or pull strategies.
- 4) low coupling in its service-oriented architecture that supports the independent evolution of the core monitoring and analysis components, or even their independent substitution when moving from one problem domain to another.

The impact of SALMonADA over the resulting service consumption time performance and its scalability have been evaluated over real services. We analysed the results in alternative deployment scenarios and we state that SALMonADA has a low intrusiveness compared to the benefits it provides. Moreover, we have developed a web client as a front-end to try it online.

The paper is organised as follows. The related work is evaluated in Sec. II. The conceptual reference model of SALMonADA is detailed in Sec. III, while the details of its design and development are included in Sec. IV. Section V provides information about the CSP-based technique to analyse the SLA fulfillment. Section VI provides information about how SALMonADA checks the QoS monitoring result against the SLA to provide timely and comprehensive explanations. In turn, Section VII evaluates the impact of SALMonADA over the resulting service consumption time performance and

Table I
ANALYSIS OF THE RELATED WORK

PROPOSALS	Functionality				Architecture	
	Supported SLAs	Monitor config.	Monitoring results	Explanation of violations	Architecture elements	Architecture structure
WSLA [2,3]	General purpose. Not H-U. ★★★	Automatic. Coupled to the SLA. ★★★	Through API. ★★★	Detection and partial explanation. Not H-U. ★★★	Monitor and Analyser separated. ★★★	CBS. ★★★
Comuzzi Kotsokalis [4]	Particular. ★★★	Automatic. Decoupled from the SLA. ★★★	Through API. ★★★	Detection. ★★★	Monitor and Analyser separated. ★★★	SOA. ★★★
Michlmayr et al. [5]	Particular. ★★★	Automatic. Coupled to the SLA. ★★★	Through query lang. ★★★	Detection. ★★★	Monitor and Analyser separated. ★★★	SOA. ★★★
Raimondi et al. [20]	General purpose. Not H-U. ★★★	Automatic. Decoupled from the SLA. ★★★	Through log. ★★★	Detection. ★★★	Monitor and Analyser in 1 component. ★★★	CBS. ★★★
Sahai et al. [17]	Particular. ★★★	Automatic. Coupled to the SLA. ★★★	Through a formal model. ★★★	Detection and partial explanation. Not H-U. ★★★	Monitor and Analyser separated. ★★★	CBS. ★★★
Palacios et al. [9]	General purpose. Not H-U. ★★★	Automatic. Coupled to the SLA. ★★★	[N/A] ★★★	Detection. ★★★	[N/A] ★★★	[N/A] ★★★
Di Penta et al. [10]	Particular. ★★★	Automatic. Coupled to the SLA. ★★★	[N/A] ★★★	Detection. ★★★	[N/A] ★★★	[N/A] ★★★
SLA@SOI [11]	General purpose. Not H-U. ★★★	Automatic. Decoupled from the SLA. ★★★	Through API. ★★★	Detection and partial explanation. Not H-U. ★★★	Monitor and Analyser separated. ★★★	SOA. ★★★
TRUSTCOM [12]	General purpose. Not H-U. ★★★	Automatic. Coupled to the SLA. ★★★	Through API. ★★★	Detection and partial explanation. Not H-U. ★★★	Monitor and Analyser separated. ★★★	SOA. ★★★
Mahbub Spanoudakis [13,14]	General purpose. Not H-U. ★★★	Automatic. Coupled to the SLA. ★★★	[N/A] ★★★	Detection and explanation. Not H-U. ★★★	Monitor and Analyser separated. ★★★	CBS. ★★★
Comuzzi Spanoudakis [19]	General purpose. Not H-U. ★★★	Automatic. Decoupled from the SLA. ★★★	Through API. ★★★	Detection. ★★★	Monitor and Analyser separated. ★★★	SOA. ★★★
our proposal SALMonADA	General purpose. H-U. ★★★	Automatic. Decoupled from the SLA. ★★★	Through a formal document. ★★★	Detection and explanation. H-U. ★★★	Monitor and Analyser separated. ★★★	SOA. ★★★

H-U: human understandable

scalability. Finally, Sec. VIII concludes the paper with a discussion of our contributions.

II. RELATED WORK

Several service-based systems quality assurance proposals that aggregate monitoring and analysis facilities can be found in the literature. To conduct the search of the related work, we have revised the most relevant conferences and journals in the area, selecting those papers that were scoped in the field of SLA monitoring and analysis. Furthermore, we have increased the results by adding relevant papers obtained from experts in the field. Table I summarizes the results of this study.

We have examined the selected papers under the four factors described in the introduction. The first three factors fall into the functionality of the proposed solution, whereas the fourth factor falls into its architecture.

Functionality. Considering the three factors for functionality identified in the introduction, we focus on the following issues: (1) Which SLAs are supported, (2a) How the information to configure the monitor is specified, (2b) How the QoS monitoring result is specified and (3) How the violations are explained.

Architecture. The issues arising from this factor are: (4a) Which architectural elements are needed and (4b) How the architectural elements are structured.

We analyse these issues below:

Which SLAs are supported. The proposals fall into one of the following categories:

- 171 • Ad-hoc SLA notation [4], [5], [10], [17]. The supported
172 SLAs include ad-hoc information without considering
173 a general-purpose structure or notation. In addition the
174 proposals are not able to deal with all the aspects re-
175 quired on some real scenarios. For instance, they do
176 not support conditional terms subject to preconditions,
177 expressive Service Level Objectives (SLOs), and optional
178 or alternative terms to specify agreement variants [18]
- 179 • General-purpose SLA notation [2], [3], [9], [11]–[14],
180 [19], [20]. The supported SLAs consider a general-
181 purpose structure and/or notation. Specifically, in pro-
182 posals such as [9], [13], [14], [19] the SLAs sup-
183 port the general-purpose structure proposed in the
184 WS–Agreement specification [16], a highly flexible
185 and widespread SLA notation. An advantage of WS–
186 Agreement over other SLA proposals is that it supports
187 the aforementioned aspects that are necessary to model
188 agreements of real scenarios, namely: (1) optional or
189 alternative terms to specify agreement variants, and (2)
190 expressive SLOs that can be guarded by a qualifying
191 condition (QC) to specify conditional terms. However,
192 WS–Agreement just provides a general-purpose schema
193 that must be extended up to eight different points with
194 an internal sublanguage. More specifically, two of the
195 eight points need to be necessarily extended to become
196 a fully-fledged language, namely the service description
197 terms (SDT) that defines the service functionality, and
198 the SLOs. This causes that a system that can deal
199 with a particular WS–Agreement notation is not able
200 to deal with other WS–Agreement notations. Moreover,
201 the sublanguages used in [9], [13], [14], [19] are not
202 neither general-purpose to be easily mapped to each other,
203 nor human understandable. Other works support SLAs
204 specified with WSLA [2], [3] and SLA* [11] that include
205 general-purpose structures and notations. However, such
206 notations are based on XML schemas even for the SLOs
207 assertions (being recursive in the case of SLA*) and thus,
208 they are not as human understandable as they should be.

209 *How the information to configure the monitor is specified.* We
210 find the following situations:

- 211 • Automatic, coupled to the SLA. Some approaches include
212 a mechanism to automatically configure a monitor from
213 the SLA [2], [3], [5], [9], [10], [12]–[14], [17]. However,
214 in these solutions, the monitor can only be used for
215 a concrete SLA specification, and if this specification
216 changes, the monitor must be modified as well due to
217 its high coupling.
- 218 • Automatic, decoupled from the SLA. There are some
219 works [4], [11], [19], [20] which decouple the SLA
220 from the monitor by translating automatically the SLA
221 to another document which includes the information
222 required for monitoring.

223 *How the QoS monitoring result is specified.* Some approaches
224 do not describe how the monitoring results are reported [9],
225 [10], [13], [14]. From those that describe it, some provide
226 a log file [20] or an Application Programming Interface
227 (API) [2]–[4], [11], [12], [19] to access the monitoring results.

228 Since these logs and APIs are not standardized, they present
229 significant differences to each other. The lack of a standard
230 prevents the possibility to easily change the monitor and also
231 it requires the analyser to be compatible with the monitor’s
232 API or log. A more effective solution is found in [5] where the
233 authors propose a query language to access the measurements
234 of the monitor. However, in this solution, the monitor must
235 deal with that query language. A more advanced solution is
236 also found at [17], where the authors propose a model to store
237 the results, however such model is not explicitly described.
238 *How the violations are explained.* The proposals fall into one
239 of the following categories:

- 240 • Just detection. Some works detect SLA violations without
241 explaining the cause [4], [5], [9], [10], [19], [20].
- 242 • Partial explanation. Other approaches provide the SLOs
243 that have been violated [2], [3], [11], [12], [17], leading
244 to a *partial explanation* of the violation as an SLO could
245 be violated for several reasons that are not detailed.
- 246 • Precise explanation but not human-understandable.
247 In [13], [14] the authors provide an accurate violation
248 explanation that identifies the violated SLOs and the
249 monitoring results which caused the violation. However,
250 they use Event-Calculus to express both the conditions
251 and monitoring results, which as the authors state, is not
252 user-friendly.

253 *Which architectural elements are needed.* We find the follow-
254 ing situations:

- 255 • The proposal does not include an architecture [9], [10].
- 256 • The proposal includes an architecture where monitoring
257 and analysis are performed in the same component [20].
- 258 • The proposal implements the QoS monitoring and SLA
259 analysis in two separated components with different re-
260 sponsibilities, increasing the cohesion and reusability [2]–
261 [5], [11]–[14], [17], [19].

262 *How the architectural elements are structured.* We classify the
263 approaches into:

- 264 • Component-Based Systems (CBS) [2], [3], [13], [14],
265 [17], [20]. The components are specific of the system,
266 and no details regarding deployment or protocols used are
267 described, which makes them difficult to reuse or replace.
- 268 • Service-Oriented Architecture (SOA) [4], [5], [11], [12],
269 [19]. This architecture adds the capability of deploying
270 the different constituent services in a distributed manner,
271 adapting or replacing them, in a highly cohesive and
272 loosely coupled system.

273 As a summary of this state of the art (see Table I), we
274 can observe that in the functional part, the proposals cover
275 satisfactorily just one or even none of the four identified
276 issues. Improving this situation is the aim of our work. Our
277 solution takes all the features mentioned and either: (1) cover
278 the issues to the same degree of the best proposal of the state
279 of the art or improves the existing proposals by introducing a
280 new strategy, as it will be described in the following sections.
281 This improvement on functionality is accompanied by optimal
282 architectonic decisions (highly decoupled SOA solution).

283 Our solution includes the usage of platforms from our
284 previous works, SALMon [21] and ADA [22]. The main

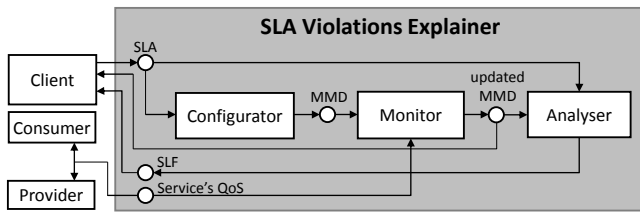


Figure 1. Conceptual Reference Model for SLA violations explainers

285 goals and initial features they had before this approach are as
 286 follows: ADA is a platform whose main goal is to check the
 287 consistency of an SLA and compatibility between SLA offers
 288 and SLA demands. In turn, SALMon is a monitoring platform
 289 whose main goal is to gather the QoS of web services and
 290 check simple conditions for several activities. None of these
 291 platforms successfully cover the aforementioned issues on their
 292 own. However, under the SALMonADA platform, both plat-
 293 forms have been extended to fulfill the required functionality.
 294 ADA has been extended to support detection and explanation
 295 of violations, and SALMon has been extended to support
 296 automatic monitoring configuration and to provide monitoring
 297 results in a formal document. A detailed description of these
 298 platforms, including their enhancements in this approach are
 299 described in Sec. IV-C for SALMon and Sec. IV-D for ADA.

300 III. THE SALMONADA CONCEPTUAL REFERENCE 301 MODEL

302 This section presents the conceptual reference model of the
 303 SALMonADA platform. This model introduces the relevant
 304 human and software agents that participate in the platform,
 305 and the data that they need to interchange. Its purpose is
 306 to provide a high-level view of the platform before going
 307 into the architectural and technological details. The conceptual
 308 reference model is shown in Fig. 1, using the SAP-TAM
 309 notation [23]. It includes the following agents:

310 *Client*: is the user of the platform. The client goal is
 311 to retrieve the explanations of SLA violations and/or the
 312 monitoring results. To fulfill such a goal the client has the
 313 responsibility of providing the SLA to monitor. The role of
 314 client may be played by the service consumer, the service
 315 provider or even a third party interested in monitoring the
 316 assessment of the SLA.

317 *Configurator*: is the agent that configures the monitor
 318 with the information included in the SLA. Thus, it decouples
 319 the SLA (a contractual specification understood by SLA-
 320 dependent agents) from the Monitor, by generating from the
 321 SLA a Monitoring Management Document (MMD), which
 322 is a specification of the monitoring directives to configure a
 323 monitor (see Sect. IV-B for more details about its structure).

324 *Monitor*: is the agent responsible of monitoring the
 325 interaction between the provider and the consumer according
 326 to the monitoring directives given in the MMD. The Monitor
 327 obtains the measured metrics from such an interaction and
 328 updates the monitoring results in the MMD right after each
 329 consumer request.

330 *Analyser*: is the agent that checks whether the monitor-
 331 ing results of a service, available in the updated MMD, is

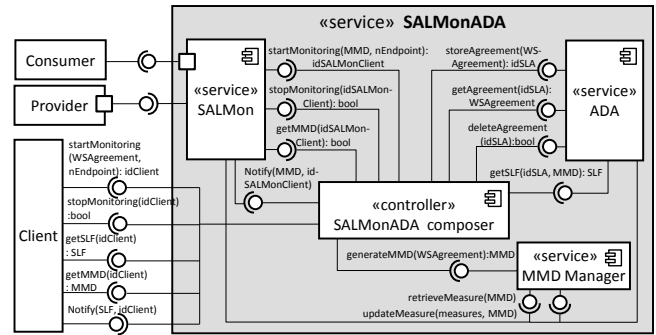


Figure 2. Architectural Model of SALMonADA

332 compliant or not with the agreed QoS included in the SLA. It
 333 ultimately produces the explanations of violations structured
 334 in a document designed for this purpose, the Service Level
 335 Fulfillment (SLF).

336 As a summary, the conceptual reference model shows a clear
 337 separation of concerns on the management of the SLAs, the
 338 MMDs and the SLFs, independent of the concrete technologies
 339 used, that are described in the next section.

340 IV. THE SALMONADA PLATFORM

341 In this section we present the details of SALMonADA
 342 platform which can be viewed as an instantiation of the
 343 conceptual reference model presented above. The platform is
 344 able to monitor and analyse expressive SLAs specified with
 345 WS-Agreement. SALMonADA has a decoupled architecture
 346 that integrates into a service-based system two previously
 347 existing systems which in turn, have been extended to re-
 348 alize this project: the SALMon monitor [21] and the ADA
 349 analyser [22]. Such an architecture is shown in Fig. 2. The
 350 core component of SALMonADA is its *composer* that provides
 351 the external interface and controls the execution flow of the
 352 system keeping *SALMon* and *ADA* decoupled from each other.
 353 Moreover, SALMonADA also comprises the *MMD Manager*
 354 service, which is used to generate and manipulate the MMDs
 355 independently of the underlying structure of such documents.

356 In the following subsections, we provide more details about
 357 these SALMonADA components. We focus on the internal
 358 architectures and responsibilities.

359 A. The SALMonADA composer

360 The SALMonADA composer is the component that presents
 361 the external interface to the client and controls the execution
 362 flow of the system. It also orchestrates the composition of
 363 ADA and SALMon to: (1) extract the monitoring information
 364 to be included in an MMD from an SLA, (2) monitor SLAs,
 365 and (3) analyse SLAs. The composer follows the low coupling
 366 design principle, so that it is possible to replace the monitor,
 367 the analyser, or the MMD manager without affecting the other
 368 elements of the platform. This implies that also the formats of
 369 SLAs and MMDs can be changed with minimal impact.

370 Thus, the SALMonADA composer controls the information
 371 flow needed for the different constituent services such as: (1)
 372 which client is interested in knowing which SLA violations;
 373 (2) the MMDs obtained from the SLAs through the MMD


```

1 <MonitoringManagementDocument>
2 <WebServiceInformation name="SALMonADA-complaint...">
3 <!--extracted from the service description term-->
4 <description>ADA is an Agreement Document Analysis
5   tool for WS-Agreement documents</...>
6 <domain>Analysis tool</...>
7 <wsdlURL>http://www.isa.us.es:8081/ADAService?wsdl</...>
8 <endpoint>http://www.isa.us.es:8081/ADAService</...>
9 <operation name="checkDocumentConsistency">
10 <soapAction>checkDocumentConsistency</...>
11 </...>
12 <operation name="explainNonCompliance">
13 <soapAction>explainNonCompliance</...>
14 </...>
15 <!--more operations were included-->
16 </WebServiceInformation>
17
18 <monitorConfiguration>
19 <globalPeriodInit>2013-05-18T18:02:38</...>
20 <!--starting monitoring time-->
21 <globalPeriodEnd>2014-01-01T00:00:00</...>
22 <!--extracted from the expiration time-->
23 </...>
24
25 <!--QoS attributes of the whole service:-->
26 <serviceMetric>
27 <metric>AverageAvailability</...>
28 <localPeriodInit>2013-05-18T18:02:38</...>
29 <localPeriodEnd>2014-01-01T00:00:00</...>
30 </...>
31 <serviceMetric>
32 <metric>GeneralResponseTime</...>
33 <localPeriodInit>2013-05-18T18:02:38</...>
34 <localPeriodEnd>2014-01-01T00:00:00</...>
35 </...>
36
37 <!--QoS attributes of specific operations:-->
38 <operationMetric opName="explainNonCompliance">
39 <metric>AverageResponseTime</...>
40 <localPeriodInit>2013-05-18T18:02:38</...>
41 <localPeriodEnd>2014-01-01T00:00:00</...>
42 </...>
43 <operationMetric opName="checkDocumentConsistency">
44 <metric>AverageResponseTime</...>
45 <localPeriodInit>2013-05-18T18:02:38</...>
46 <localPeriodEnd>2014-01-01T00:00:00</...>
47 </...>
48 <!--more operation metrics were included-->
49 </MonitoringManagementDocument>

```

Figure 5. Excerpt of the MMD obtained from document of Fig. 4.

are the service properties used within SLOs whose scope is one or more specific operations but not to all of them. For instance, `AverageResponseTime` is an operation metric for the `explainNonCompliance` and `checkDocumentConsistency` operations, among others; and `AverageAvailability` is a service metric obtained from the `GTgeneralAvailability` term. The information stored in the MMD is always accessed or modified through the MMD Manager. Therefore, the MMD structure is completely independent of the system. To achieve this, we have considered it as a variability point at design level and hence any other MMD Manager service could be used.

C. The SALMon service

SALMon is a service-based system aimed at monitoring the QoS of web services. The main features of SALMon are (see [21] for details): it may operate on any available technology (SOAP-based, RESTful, etc.) with minor and localised changes; it may interoperate easily with other frameworks (e.g., self-adaptive service-based systems [25] or cloud infrastructures [26]); it is easily extensible to monitor new metrics; it combines passive monitoring and on-line testing [15].

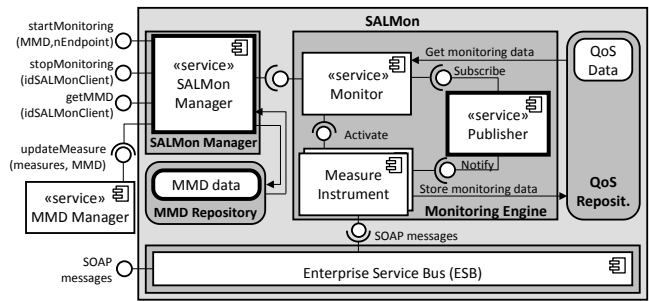


Figure 6. Architectural Model of the extended SALMon

In the context of SALMonADA, SALMon has been enhanced with new components and services (see elements with a thicker frame in Fig. 6): the SALMon Manager service, the MMD repository and the Publisher service. The existing components of SALMon have also been adapted to interact with the new services. We describe below each module:

SALMon Manager: is responsible to (1) store the MMDs in the repository, (2) configure the Monitoring Engine and (3) update the MMD when new measured metrics are retrieved. A measured metric is a metric with its value in a specific timestamp (see Fig. 7). The SALMon Manager uses the MMD Manager whenever it requires to get or store the monitoring results to the MMD. In such a way, it can be extended to support different MMD structures.

Monitoring Engine: is responsible to monitor the services. The Monitor service creates and manages the Measure Instruments. Each Measure Instrument is responsible to gather a specific metric such as response time or availability (metrics of the service from the client's perspective) and store the results in the QoS Repository. The module also includes the Publisher to notify when a new measured metric is obtained.

Enterprise Service Bus (ESB): all requests and responses are sent through this communication channel, which in turn, feeds the Measure Instruments with the intercepted messages.

QoS Repository: stores the measured metrics. To reduce storage consumption, it saves only the last measurements, the average values and the number of invocations.

The generated MMD with the monitoring results is shown in Fig. 7. Notice that the measured metrics related to time are defined in milliseconds, as a convenient unit to express the monitoring results. The platform performs afterwards the required unit conversions for the SLA analysis.

Finally, it is important to address performance and privacy issues during monitoring. The former is evaluated in depth in the experimental results in Sec. VII. The latter, which is strongly affected by the service policies and the deployment configuration, is discussed here considering the different factors involved. On secure services, the ESB does not interfere with the security level in the communications. However the issue arises for non-secure services dealing with sensitive data. For such a reason, the body of the message is never stored. Moreover, the ESB can be deployed on either the client or provider side under their management, so any sensitive data can be encrypted using WS-Security before forwarding the message to the Measure Instruments.

```

1 <MonitoringManagementDocument>
2 ...
3 <serviceMetric>
4 <metric>AverageAvailability</...>
5 <localPeriodInit>2013-05-18T18:02:38</...>
6 <localPeriodEnd>2014-01-01T00:00:00</...>
7 <measure>
8 <value>100</value>
9 <timeStamp>2013-05-18T18:02:38</timeStamp>
10 </measure>
11 </...>
12 ...
13 <operationMetric opName="explainNonCompliance">
14 <metric>AverageResponseTime</...>
15 <localPeriodInit>2013-05-18T18:02:38</...>
16 <localPeriodEnd>2014-01-01T00:00:00</...>
17 <measure>
18 <value>3421</value>
19 <timeStamp>2013-05-18T18:02:38</timeStamp>
20 </measure>
21 </...>
22 ...
23 </MonitoringManagementDocument>

```

Figure 7. Excerpt of the MMD of Fig. 5 with the monitoring results.

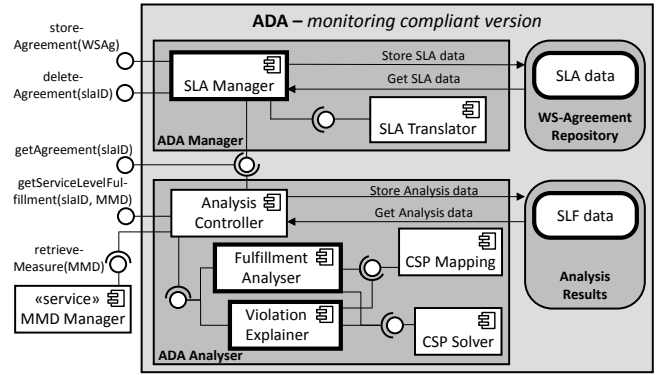


Figure 8. Architectural Model of the monitoring-compliant ADA.

```

1 Unfulfillment Explanation:
2   Violated terms: GTexplainNonComplianceRespTime
3   (AverageResponseTime < 2 s.)
4   explainedBy: AverageResponseTime = 3421 ms.

```

Figure 9. SLF for the violation of the SLA (Fig. 4) by the MMD (Fig. 7)

503 D. The ADA service

504 ADA is an *Agreement Document Analysis* framework
505 aimed at extracting useful information from agreement docu-
506 ments [22]. It has been developed based on our previous the-
507 oretical works on applying the *constraint satisfaction problem*
508 (CSP) [27] paradigm to the automated procurement of web
509 services [28]. The main features of ADA [29] are: (1) *ready-to-*
510 *use* by detecting and explaining conflicts within and between
511 WS–Agreement documents [29], [30]; (2) *functional suitabil-*
512 *ity* by supporting the analysis of expressive WS–Agreement
513 documents with conditional, optional, or alternative terms,
514 term scopes, arithmetic-logic expressions inside SLOs, et
515 cetera [18]; (3) *understandability* by supporting a plain-text
516 notation [18] that makes reading and writing WS–Agreement
517 documents easier for humans; (4) *interoperability* through a
518 triple distribution model (Java library, OSGi⁴ service, and web
519 service); and (5) *CSP solver independence*.

520 Similarly to SALMon, ADA has been enhanced with new
521 components (see elements with a thicker frame in Fig. 8): the
522 SLA Manager, several analysis facilities to detect and explain
523 violations at monitoring, and repositories for SLAs and SLFs.
524 The existing ADA components have also been adapted to
525 interact with other components of the SALMonADA frame-
526 work. ADA components are grouped into two modules with
527 the following responsibilities.

528 *ADA Manager*: is responsible for SLA storage and re-
529 trieval from the repository; as well as the translation between
530 several SLA models to a WS–Agreement-based normalised
531 one using XML that ADA is able to analyse.

532 *ADA Analyser*: is responsible for: (1) analysing the ful-
533 fillment of the WS–Agreement document given the monitoring
534 results stored in the corresponding MMD; (2) the creation of
535 violation explanations when a violation is detected; and (3) the
536 storage and retrieval of the SLF information. Such analysis is
537 performed by means of a CSP solver tool (see Sec. V for more
538 details) that supports the following assertion language:

$$\begin{aligned}
 P &::= P \text{ op}_L P \mid T, \text{ predicate, where } \text{op}_L \in \{\wedge \mid \vee \mid \neg \mid \Rightarrow \mid \Leftrightarrow\} \\
 T &::= E \text{ op}_C E, \text{ term, where } \text{op}_C \in \{= \mid \neq \mid > \mid \geq \mid < \mid \leq\} \\
 E &::= E \text{ op}_A E \mid \text{var} \mid \text{lit}, \text{ expression, where } \text{op}_A \text{ is an algebraic} \\
 &\quad \text{operator defined on the domain of variables and literals}
 \end{aligned}$$

539 As a result, we work with an assertion language inside
540 WS–Agreement documents that is not just expressive but
541 also easy to understand by non-technical users (cf. SLOs of
542 Fig. 4). Moreover, as Fig. 9 shows, the reported SLF is also
543 understandable because the violated terms are associated with
544 the violating monitored values, both expressed in a human
545 understandable way.

546 V. A CSP-BASED TECHNIQUE TO EXPLAIN VIOLATIONS

547 The detection and explanation of violations, which is the
548 ultimate goal of this paper, is necessarily based on the align-
549 ment of the WS–Agreement document, which expresses the
550 requirements on the service, and the MMD document, which
551 collects monitoring information at runtime. To implement this
552 alignment, we need to provide semantics to both documents
553 and then define the concept of violation and the procedure to
554 get explanations. In our technique, the semantics is defined
555 through semantic mappings. Under this view, the elements
556 of the documents, that are considered as source models are
557 mapped into a target domain whose semantics have been
558 formally defined [31]. The main advantage of semantic map-
559 pings is that they enable the usage of techniques, preferably
560 automated, which are specific to the target domain in order to
561 infer properties in the source models [32].

562 In our case, we have chosen constraint satisfaction problems
563 (CSP) [27] as the target domain. Solving problems by means
564 of CSPs has been a research topic in Artificial Intelligence
565 for years. In short [33], a CSP is a three–tuple of the form
566 (V, D, C) where V is a set of variables, D is a set of domains,
567 and C is a set of constraints. Each variable $V_i \in V$ has a
568 finite domain $D_i \in D$. Each constraint in C applies to a subset
569 of the variables, and restricts the combination of values that

⁴www.osgi.org

570 those variables can take at the same time. A solution to a
 571 CSP is an assignment of a value to each variable such that all
 572 the constraints are satisfied. Consider, for instance, the CSP:
 573 $(\{a, b\}, \{[0, 2], [0, 2]\}, \{a + b < 4\})$. Then, $(a = 2, b = 0)$
 574 is a possible solution since it verifies that $2 + 0 < 4$, whilst
 575 $(a = 2, b = 2)$ is not a valid solution.

576 The reason for selecting CSP is twofold. On the one hand,
 577 the MMD and the most significant part of an WS–Agreement
 578 document are sets of constraints over service properties and,
 579 therefore, CSPs can be used to model the detection and
 580 explanation problem in a very natural way, as we showed in
 581 previous works [28]–[30], [34]–[36]. On the other hand, there
 582 is a plethora of CSP solvers available that support a wide
 583 range of constraints and can be used to automatically analyse
 584 WS–Agreement documents in an efficient manner.

585 Our proposal to interpret the SLA fulfillment using CSPs
 586 is based on checking that the monitored metrics values are a
 587 possible solution for the CSP mapped from the SLA.

588 The proposed CSP-based technique to detect and explain
 589 SLA unfulfillments takes into consideration that those WS–
 590 Agreement terms whose scope is a specific service oper-
 591 ation are only affected by monitored metrics measured
 592 while such operation is being executed. For example, the
 593 `AverageResponseTime` attribute of a term whose scope is
 594 the `explainNonCompliance` operation, would be checked
 595 for fulfillment with the average of monitored response time
 596 of `explainNonCompliance` operation. To consider this
 597 in our system, the original WS–Agreement document Δ is
 598 separated into several views Δ_{op} by the service operation
 599 scope. Each view is a WS–Agreement document by itself
 600 restricted to both: the set of terms whose scope is a specific
 601 service operation; and such terms whose scope are all service
 602 operations (i.e. such terms that do not specify any scope
 603 such as `GTgeneralAvailability` of Fig. 4). In the
 604 case a term whose scope includes more than one service
 605 operation, but not all of them, it would be included in as many
 606 views as scoped service operations. For instance, the term
 607 `GTgeneralResponseTimeRelations` of Fig. 4 would
 608 be included in four views, one for each scoped operation.

609 In general, a WS–Agreement view Δ_{op} defines some guar-
 610 antees for a set of service properties whose domain has been
 611 previously defined. WS–Agreement views can be modelled
 612 in a CSP by means of the semantic mapping $map(\Delta_{op})$
 613 that is summarised in Table II, whereas Fig. 10 shows the
 614 CSP mapped from the `explainNonCompliance` view
 615 ($\Delta_{explainNonCompliance}$) of the ADA SLA included in Fig. 4.

616 Concerning MMDs, the original MMD M is also separated
 617 into several views M_{op} by the monitored service operation.
 618 Thus, each view comprises the QoS monitoring result that
 619 includes metric–measure pairs of a unique service operation⁵,
 620 and those affecting all service operations. Table III summarises
 621 the MMD semantic mapping $map(M_{op})$, whereas Fig. 11

⁵Note that the measure units must match for a given scale (e.g. time) in the service properties of the SLA and the MMD. If they do not match, a conversion formula should be applied for each CSP constraint, as proposed in some QoS ontologies for Web Services [37]. Fig. 15 includes different measure units for the `AverageResponseTime` property: seconds in the SLA and milliseconds in the MMD.

Table II
TERMS MAPPING FROM WS–AGREEMENT VIEW TO CSP.

WS–Agreement Element	CSP Mapping
1 Agreement ... 2 Context: ... 3 Responder/Initiator: ...	This information is not mapped into the CSP
1 name: Service Property 2 propertyName 3 – measured by metricDefinition	$V \leftarrow V \cup \text{propertyName}$ $D \leftarrow D \cup \text{domain}(\text{metricDefinition})$ $C \leftarrow C$
1 name: Service Description 2 SDTEExpr	This information is not mapped into the CSP
1 name: Guaranteed by x 2 Scope: ServiceOperation', 3 QC: QCExpr, 4 SLO: SLOExpr	$V \leftarrow V$ $D \leftarrow D$ $C \leftarrow C \cup (\text{QCExpr} \Rightarrow \text{SLOExpr})$
1 name: Guaranteed by x 2 Scope: ServiceOperation', 3 SLO: SLOExpr	$V \leftarrow V$ $D \leftarrow D$ $C \leftarrow C \cup (\text{SLOExpr})$

Table III
MONITORING QOS DATA MAPPING FROM MMD VIEW TO CSP.

Measure	CSP Mapping
1 metric = value 2 for ServiceOperation'	$V \leftarrow V \cup \text{metric}$ $D \leftarrow D \cup \text{domain}(\text{metric})$ $C \leftarrow C \cup (\text{metric} = \text{value})$

622 shows the CSP mapped from the monitoring results of the
 623 ADA SLA corresponding to the `explainNonCompliance`
 624 operation included in Fig. 7.

625 In order to carry out the detection and explanation of SLA
 626 unfulfillment, we need to use a pair of analysis techniques that
 627 have been widely used in CSPs:

- 628 • $solve(V, D, C)$ tries to find all CSP solutions. To this end,
 629 many heuristics and techniques have been developed to
 630 solve CSPs in an efficient manner [38], [39].
- 631 • $explain(V, D, C)$ tries to provide an explanation when
 632 such solution is not possible [40]. This explanation is
 633 a minimal set of constraints $c \in C$ that makes impossible
 634 to find a solution that satisfies C . For instance, the CSP
 635 $(\{a, b, d\}, \{[0..2], [0..2], [0..2]\}, \{a + b < 1, a > 0, d >$
 636 $1\})$ is not satisfiable, and its possible explanations c are
 637 either $\{(a + b < 1)\}$ or $\{(a > 0)\}$.

638 On the basis of these operations and the previously de-
 639 scribed map functions for the SLA and MMD, we may
 640 trace back in order to infer the SLA fulfillment with the
 641 MMD measures, providing the precise semantics for detect-
 642 ing the unfulfillment and their explanation, as follows. Let
 643 $map(\Delta_{op}) = (V, D, C)$ and $map(M_{op}) = (V', D', C')$

$$unfulfillment(\Delta_{op}, M_{op}) \Leftrightarrow solve(V \cup V', D \cup D', C \cup C') = \emptyset$$

644 In the example above, the WS–Agreement view
 645 $\Delta_{explainNonCompliance}$ is not fulfilled by the MMD measures
 646 $M_{explainNonCompliance}$, due to the lack of solutions for the joint
 647 CSP generated.
 648

649 In this case, an explanation like such provided in Fig. 9
 650 would be reported using the following formula:

$$unfulfillment_{exp}(\Delta_{op}, M_{op}) = explain(V \cup V', D \cup D', C \cup C')$$

$V = \{ \text{AverageResponseTime}, \text{GeneralResponseTime}, \text{AverageAvailability} \}$
 $D = \{ [1..∞], [1..∞], [1..100] \}$
 $C = \{ \text{AverageAvailability} > 95, \text{GeneralResponseTime} \leq 1.5, \text{GeneralResponseTime} > 0.5 \Rightarrow \text{AverageResponseTime} < 2 \}$

Figure 10. CSP generated from WS–Agreement view: $\Delta_{\text{explainNonCompliance}}$.

$V = \{ \text{AverageResponseTime}, \text{GeneralResponseTime}, \text{AverageAvailability} \}$
 $D = \{ [1..∞], [1..∞], [1..100] \}$
 $C = \{ \text{AverageAvailability} = 100, \text{GeneralResponseTime} = 1.7, \text{AverageResponseTime} = 3.421 \}$

Figure 11. CSP generated from MMD view: $M_{\text{explainNonCompliance}}$.

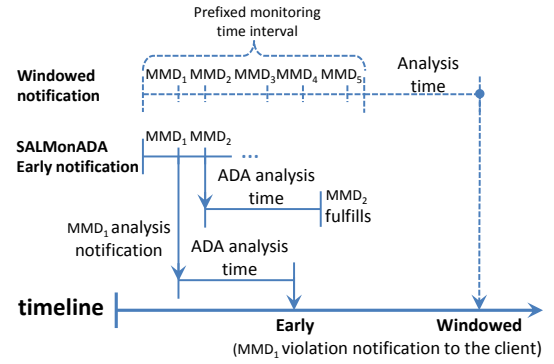


Figure 12. Early analysis notification vs. windowed analysis notification.

A. Push Interaction Approach

The *push approach* is the most convenient way to interact with SALMonADA due to the push nature of its service monitoring and analysing. In this sense, as Fig. 12 shows, the platform incorporates an *early analysis notification* that supports the analysis of the SLA fulfillment as soon as new measured metrics have been updated in the MMD. Thus, the SLF notification is sent to the client without any further delay than the analysis time, in contrast with windowed proposals [13] that get periodically the monitored values in a prefixed time interval to analyse them later. Assuming a similar analysis time, the difference between both approaches is higher when the violation affects to a measured metric at the beginning of the prefixed time interval (see MMD_1 in Fig. 12). In the case that the violation affects to a measured metric at the end of the prefixed time interval, the notification will be delivered at the same time in both approaches. As the sequence diagram of Fig. 14 depicts, once the client has started to monitor, the *provider service* included in the reported WS–Agreement document is monitored by the *SALMon* service. Next, the MMD created from the monitored WS–Agreement document is sent to the *MMD Manager* with the measured metrics to be updated. Finally, the new MMD is notified to the *SALMonADA composer* that sends it to the *ADA* service to analyse the service level fulfillment of the corresponding WS–Agreement document (cf. Sec. IV-D for more details). Then, the client is notified about the WS–Agreement document fulfillment by means of the SLF. If the WS–Agreement is not fulfilled, both the specific violated WS–Agreement terms and the violating metrics, are included as violation explanation. Note that SALMonADA supports the same endpoint acting as different clients, for instance, one of them to get the SLF, and another to store reputation analytics of the service consumer and provider, or even to perform self-adaptation strategies.

For instance, let us suppose that SALMonADA is monitoring the WS–Agreement document of ADA depicted in Fig. 4. That document specifies that some operations have a higher priority and are required to be faster than the average response time of the different methods of the service ($\text{AverageResponseTime} \leq \text{GeneralResponseTime}$). If there was a violation, SALMonADA would report an explanation right after detecting it. The explanation would identify if the

651

652 The algorithm is applied subsequently to the different views
 653 of both documents and trace the CSP constraints back to the
 654 corresponding WS–Agreement term or MMD measure. An
 655 example of the final result is provided in Fig. 9 and in Sec. VI.

656

VI. SALMONADA IN USE

657 SALMonADA is designed and developed to support push
 658 and pull interaction styles with its clients. Thus, a client, based
 659 on its own needs, may choose the approach that best fits
 660 them. Independently of the selected approach, a client shall
 661 start the use of SALMonADA by subscribing as such a client.
 662 Similarly, the client shall stop the SALMonADA monitoring
 663 at the end of the interaction, to be unsubscribed as client.

664 As the sequence diagram of Fig. 13 shows, the client
 665 starts the monitoring process by providing a WS–Agreement
 666 document to monitor its fulfillment. In case of push approach,
 667 the client shall also provide the endpoint where the notification
 668 is awaited. Independently of the approach used, ADA stores
 669 the WS–Agreement document. This document is also sent
 670 to the *MMD Manager* for it to generate the MMD. Such
 671 MMD will store the information to configure the monitor
 672 and all the measured metrics obtained while monitoring the
 673 service agreed in the WS–Agreement document. Afterwards,
 674 a monitoring session is started in *SALMon* providing the
 675 endpoint where the monitoring result must be notified to
 676 update the MMD after a service consumption. Finally, a client
 677 identifier (`clientID`) is generated to denote the specific
 678 monitoring session for this SLA. If several SALMonADA
 679 clients wanted to monitor a unique SLA, for instance the
 680 service consumer and provider to be informed about the SLA
 681 fulfillment, a different monitoring session would be started and
 682 thus, a different `clientID` would be returned.

683 When a client wants to be unsubscribed from
 684 SALMonADA, it provides its `clientID` to stop its
 685 monitoring process. The WS–Agreement document and the
 686 monitoring session is removed from SALMonADA by ADA
 687 and SALMon if there is no other client monitoring the same
 688 SLA, otherwise SALMonADA will keep them while in
 689 use by these other clients. In any case, the `clientID` is
 690 removed from the system.

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

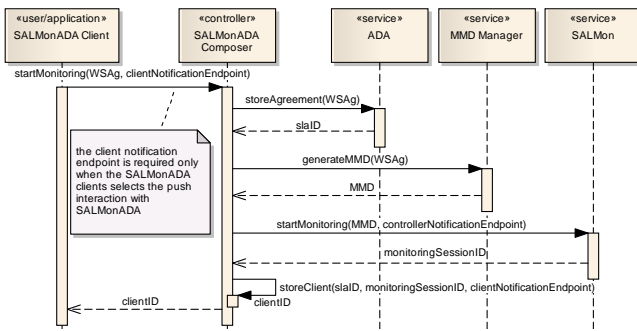


Figure 13. A SALMonADA client starts monitoring.

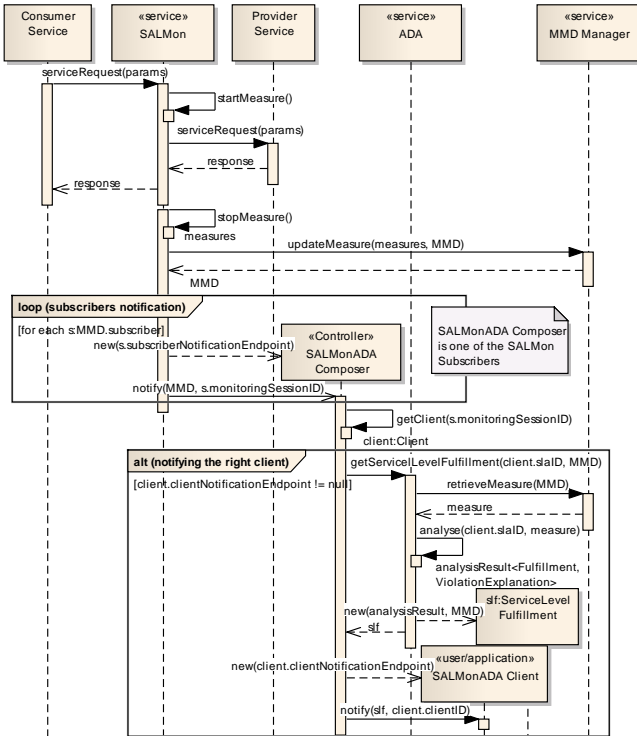


Figure 14. Push SALMonADA approach.

733 violating metric is either *AverageResponseTime* or *General-*
 734 *ResponseTime* because a simple identification of the violated
 735 term is not enough to grasp the violation cause. Similarly,
 736 SALMonADA supports the explanation of violations of more
 737 expressive SLOs. For instance, the provider may guarantee a
 738 different average response time limit for the slower service
 739 operations, depending on the general response time of the
 740 service, as follows: $((GeneralResponseTime \geq 0 \text{ AND } GeneralResponseTime < 2) \text{ IMPLIES } (AverageResponseTime < 3))$
 741 $\text{AND } ((GeneralResponseTime \geq 2 \text{ AND } GeneralResponseTime$
 742 $\leq 4) \text{ IMPLIES } (AverageResponseTime < 5))$.
 743

744 B. Pull Interaction Approach

745 The *pull approach* allows the client to actively request the
 746 results of SALMonADA for either: the current MMD with
 747 the most recent monitoring results obtained by *SALMon*; or
 748 the current SLF of the WS–Agreement document analysed
 749 by *ADA*. The former document is obtained by invoking the

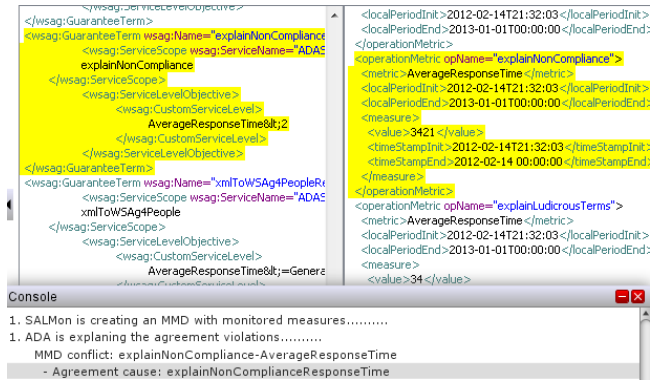


Figure 15. Reporting a violation with the SALMonADA client

750 getMMD method, whereas the latter is through the get-
 751 ServiceLevelFulfillment method. In both cases, the client is
 752 required to provide the `clientID` as input.

753 As usual in pull approaches, if the consumer(s) do not
 754 invoke the service, there is not new monitoring information
 755 and thus, it is possible that the client gets the same monitoring
 756 information in consecutive MMD requests.

757 For demonstration purposes, we have implemented a web
 758 application⁶ as a SALMonADA client in order to specify or
 759 upload the WS–Agreement documents to monitor, execute
 760 SALMonADA and receive the results. In this web applica-
 761 tion, we have introduced the WS–Agreements of ADA
 762 and SALMon themselves. By monitoring the SLAs of these
 763 services, we assess on the one hand, the functionality of
 764 SALMonADA, and on the other, the non-functional aspects of
 765 its main components. Such a SALMonADA client uses the pull
 766 interaction because the user press the corresponding interface
 767 controls to get the MMD and the analysis results. Moreover,
 768 as part of the demonstration and to assure that the service
 769 subject of the SLA is being requested, we have simulated the
 770 consumers that execute ADA and SALMon services.

771 With the client identifier, the SALMonADA client can get,
 772 at any time, the MMD, check if a violation has occurred,
 773 and in such a case, receive an explanation of the violation.

774 As Fig. 15 depicts, the web application highlights as viola-
 775 tion explanation that the *AverageResponseTime* of *explain-*
 776 *NonCompliance* operation is the violating metric because it
 777 was measured as 3421 milliseconds, while the guarantee term
 778 obligates the provider to respond in less than 2 seconds. As
 779 stated in Sec. V, SALMonADA handles different measurement
 780 units and the required unit transformation is performed to
 781 evaluate the conditions.

782 VII. PERFORMANCE AND SCALABILITY EVALUATION

783 In this section we evaluate both the performance and
 784 scalability of SALMonADA. Particularly, we focus on the
 785 overhead introduced by including SALMonADA within the
 786 consumer-provider interaction, and the maximum number of
 787 service request it is able to handle without incrementing such
 788 overhead. To do so, we first introduce how the components of
 789 its architectural model (see Fig. 2) affects the performance:

⁶The SALMonADA web application can be tried at www.isa.us.es/ada.source/SLAnalyzer/. A screencast is available at gessi.lsi.upc.edu/salmon/ada/

790 *SALMonADA composer, MMD Manager and ADA*: We
 791 have implemented two strategies to avoid them introducing
 792 overhead: (1) they are executed concurrently, without inter-
 793 fering the service consumption, because the response of the
 794 provider is returned to the consumer before analysing the
 795 SLA fulfillment (see Fig. 14), and (2) they are deployed in
 796 a different location from the monitored service, and hence
 797 they do not share the same resources (e.g. CPU, RAM,
 798 etc.). Consequently, although these components might be time-
 799 consuming (e.g. the analysis of the SLA), they do not introduce
 800 an overhead over the response time of the monitored service.

801 *SALMon*: The ESB Apache Synapse included in
 802 SALMon adds a low overhead while handling the HTTP
 803 messages. The ESB has a non-blocking HTTP transport and
 804 multi-threaded mediation, which as we measured, results in a
 805 negligible 1 - 3 ms overhead. Nevertheless, in our approach,
 806 there are three possible locations where SALMon can be
 807 deployed: at the server side, at the client side, or in an
 808 intermediate server (i.e. in the middle). Depending on the
 809 location, the overhead experienced by the consumer varies.

810 If SALMon is placed at the server or client side, there
 811 is an overhead on the resources due to the execution of the
 812 monitoring components. However, this overhead can be easily
 813 compensated by adding more resources.

814 If SALMon is placed in the middle, it does not produce an
 815 overhead on the resources of the client or server side. However,
 816 the deployment of SALMon in an intermediate server adds
 817 a network delay from Internet Service Providers due to the
 818 redirection of the messages. In this scenario, SALMonADA is
 819 not responsible of the overhead introduced by the network,
 820 but under heavy usage the components of SALMon might
 821 experience a bottleneck. To quantify it, we evaluate by means
 822 of an adequate benchmark, (1) the response time overhead
 823 under normal operating conditions (i.e. one invocation at a
 824 time), and (2) the maximum throughput SALMonADA is able
 825 to handle without incrementing such overhead.

826 *a) Setting up the experiment*: To perform the evaluation,
 827 we invoke a set of real services, and compare the response time
 828 by invoking the services both directly and through SALMon.
 829 The agents involved in the experiment are the monitored
 830 services, the client and SALMon.

831 To obtain a set of representative services, we started from
 832 a list of 393 services available in a public repository⁷. Then
 833 we applied the following criteria: (1) We first considered the
 834 most recently submitted services under the assumption that
 835 recent services are more likely to be available and running
 836 than older services. Considering the length of the list, we
 837 established as threshold the 1/3 of the complete list. (2)
 838 From the resulting 131 services, we removed those ones
 839 falling into any of the following situations: were not available,
 840 were payment services, required registration or didn't have
 841 stateless operations, resulting in 23 services. (3) We tested
 842 these 23 services and removed those ones that had errors
 843 in their descriptions (WSDL), or that gave faulty results in
 844 their functionality when invoked, resulting in a final list of 11
 845 services from 8 different service providers, deployed on their

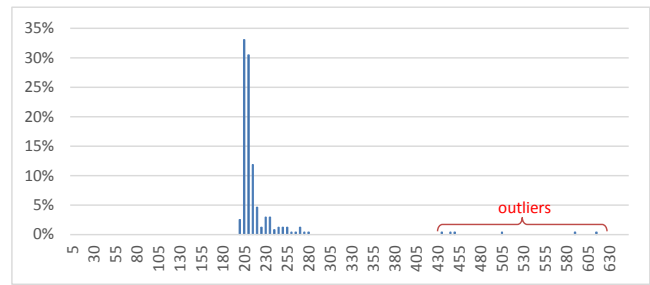


Figure 16. Response time distribution of ISBN service

846 respective servers (i.e. out of the control of the experiment),
 847 see Appendix A for the full list.

848 The client that invokes those services has been deployed in
 849 Seville in an Intel i7 of 2,20 GHz, 8 GB of RAM, a download
 850 speed of 13,92 Mbps and upload speed of 8,15 Mbps.

851 SALMon has been deployed in Barcelona in a dedicated
 852 server Intel 2,6Ghz, 6GB RAM, a download speed of 56,57
 853 Mbps and upload speed of 15,71 Mbps.

854 *b) Conducting the experiment*: We first conducted the
 855 experiments under normal conditions executing 100 service
 856 calls per each service in both direct and redirected forms by us-
 857 ing synchronous calls (i.e. one invocation at a time). Then, we
 858 conducted asynchronous calls to test the maximum throughput
 859 of the system starting from a throughput of 1 invocation
 860 per second to a throughput of 100 invocations per second.
 861 We performed 100 invocations per each throughput in both
 862 direct and redirected forms. Some services have limitations
 863 to be tested under these circumstances (e.g. restrictions on
 864 the number of concurrent invocations) and the scalability have
 865 been analysed on the services that didn't have these limitations
 866 (6 of the 11 services).

867 *c) Analysing the results*: One key issue regarding the
 868 analysis of the results is dealing with outliers (e.g. network
 869 failures that increase the response time of an invocation).
 870 Commonly used methods to deal with outliers require that
 871 the data follow a Gaussian distribution [41]. However, from
 872 the experiment results we have observed that response times
 873 do not follow a Gaussian distribution, but an exponentially
 874 modified Gaussian or inverse Gaussian distribution. For in-
 875 stance, the distribution of the response time of one of the
 876 monitored services is depicted in Fig. 16. As shown, the
 877 population grows rapidly on the left-hand side and decreases
 878 slowly on the right-hand side in the form of a tail. Those
 879 elements that are far away from the mean are considered
 880 outliers. To deal with these outliers, we followed the methods
 881 described and evaluated by Ratcliff for dealing with response
 882 time outliers [42]. Although Ratcliff studied response time of
 883 people in the field of psychology, the results can be applied
 884 to any model that follows the inverse Gaussian distribution.
 885 According to Ratcliff, we will not compute directly the average
 886 response time (which is not a robust estimator in front of
 887 outliers), but we will use two other robust estimators, namely,
 888 the inverse transformation and removing outliers at a standard
 889 deviation distance. The first estimator consists on applying
 890 the inverse response time (1/R) on each individual invocation,
 891 calculate the average, and then invert the result. The second

⁷<http://www.xmethods.net/ve2/Directory.po>

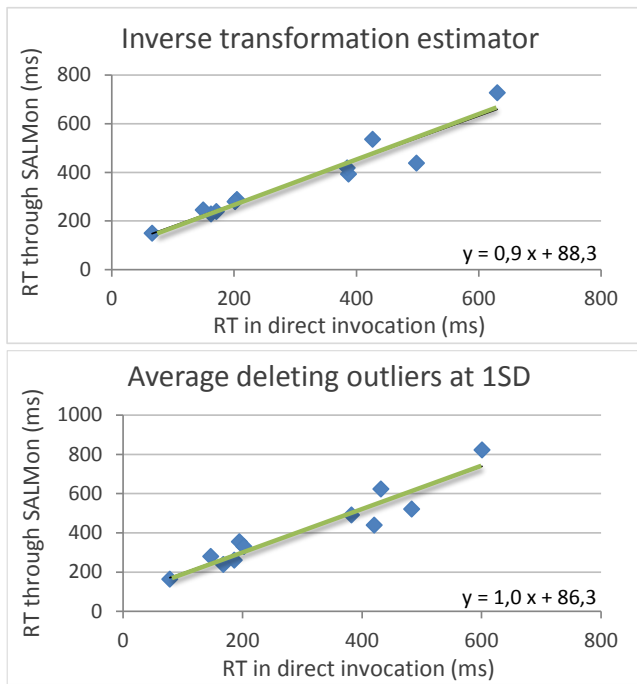


Figure 17. Response time of the service with SALMon with respect to direct invocation

estimator consists on calculating the average response time after removing the outliers at a standard deviation distance. We computed these methods over the invocations on each service for both directed and redirected invocations. As a result, we got two robust estimators per each service. We applied these estimators to the response time of direct and redirected invocations in order to calculate the response time overhead introduced in the service interaction by the deployment of SALMon in the middle. We decided to relate the two parameters with a linear interpolation curve fitting method with the aim of obtaining mathematical functions approximating the response time overhead. Fig. 17 shows the obtained functions for each of the two applied robust estimators methods, which are: $y = 0.9x + 88.3ms$ and $y = 1.0x + 86.3ms$.

Then, using the same estimators, we calculated what is the maximum throughput without incrementing such overhead. As shown in the results of Fig. 18, the maximum throughput of redirected invocations is 41 invocations per second. Above this number, the overhead response time raises significantly. Notice that some monitored services present a lower maximum throughput, but as shown, it is because they are unable to process such amount of invocations (either direct or redirected). We conclude that the response time overhead of the deployment of SALMon in the middle has a constant value between 86 and 89 ms, and its maximum throughput is 41 invocations per second.

We must remark that this overhead does not interfere with the SLA analysis, because the monitored response time corresponds to the real response time of the service. Nevertheless, as the consumer of the service experiences this delay, it is worthy to mention some deployment strategies to mitigate any concern. On the one hand, deploying SALMon at the

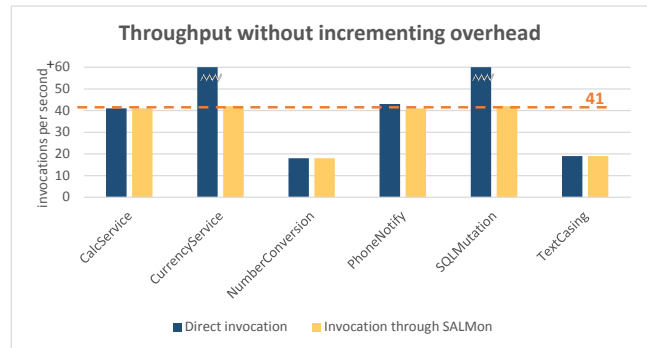


Figure 18. Throughput without incrementing response time

client or server side is suitable for service consumers that require extremely fast response times and need to avoid the 86-89 ms overhead caused by the network delay. Moreover, as discussed before, this deployment is also suitable for services dealing with sensitive data. On the other hand, for other types of services, we argue that a deployment in the middle is preferred, since this solution is less intrusive to both the client and the provider server, as it does not require the installation of the monitor in their infrastructures. Moreover, if a throughput higher than 41 requests per second is required, multiple instances of the ESB can be replicated in different servers under the same SALMon platform.

VIII. CONCLUSIONS AND DISCUSSION

In this paper we have presented a solution to monitor and analyse SLAs in order to provide timely detection and comprehensive explanations of their violations. Such information is really appealing for decision-making activities performed at runtime. For instance, the consumers of AmazonS3 scenario mentioned in Sec. I would be benefited from using SALMonADA because they are able to decide if they have to claim for a reward or not due to the service level fulfillment information provided. Other scenarios that may benefit from our proposal are the renegotiation of SLAs or the adaptation of SBSs. Our solution addresses satisfactorily the different issues identified in Sec. II:

- SALMonADA's supported SLAs are general-purpose because they follow the WS-Agreement [16] structure, completed with general-purpose sublanguages [18]. Moreover, the used notation is more human understandable than such proposed in [2], [3], [9], [11], [13], [14], [19] making it easier to be managed by human users.
- The SALMonADA platform is able to extract automatically from the SLA the information needed to configure the monitor in order to detect violations. This is done through a document, the Monitoring Management Document (MMD), which supports decoupling the SLA structure from the monitor service.
- SALMonADA uses the MMD itself to store the monitoring results coming from the monitor. The advantage of our approach is to have the MMD as the unique document that centralise all monitoring-related information without coupling to a specific API or query language.

- We introduce the Service Level Fulfillment (SLF) as a document to explain accurately the violations by identifying explicitly the violated terms and the violating monitored results. The explanations are computed through the application of a powerful CSP-based mechanism. SALMonADA supports reporting such SLF information either as soon as a violation is detected (push), or when the client requests it (pull).
- SALMonADA's architecture keeps the monitor and analyser services decoupled from each other allowing thus independent evolution and eventually selective substitution. The MMD and SLF documents support this decoupling.
- The organization of SALMonADA as a SOA makes the interoperability of the platform with other tools easier.

Moreover, we have performed an evaluation of the impact that SALMonADA has over the performance of the service consumption. The evaluation has been performed over real services using suitable estimators for response time to evaluate both its overhead and scalability. Although a low overhead is added either in the resources consumption or in the service response time, we expose how to mitigate it either by adding more resources, or deploying SALMonADA in alternative locations. Despite such an overhead, the client and the service provider are benefited from using our proposal. On the one hand, the client gets the real QoS of the service, which is required for several activities, such as self healing, claiming rewards due to SLA penalties, et cetera. On the other hand, the service provider can manage its services with the knowledge of the QoS, which is helpful to take appropriate decisions such as adding more resources to a specific service, renegotiating the SLA, et cetera.

As Future work we plan to study how to integrate SALMonADA on different self healing systems as such proposed in [13], [19]. In doing so, we will demonstrate that SALMonADA can be easily deployed with other monitors and analysers, as well as supporting other MMDs and SLAs.

ACKNOWLEDGMENT

This work has been partially supported by: S-Cube, the European Network of Excellence in Software Services and Systems; the European Commission (FEDER); the Spanish Government under the CICYT projects SETI (TIN2009-07366), TAPAS (TIN2012-32273) and ProS-Req (TIN2010-19130-C02-01); and by the Andalusian Government under the projects THEOS (TIC-5906) and ISABEL (P07-TIC-2533).

REFERENCES

- [1] M. P. Papazoglou, K. Pohl, M. Parkin, and A. Metzger, eds., *Service Research Challenges and Solutions for the Future Internet - S-Cube*, vol. 6500 of LNCS, Springer, 2010.
- [2] A. Keller and H. Ludwig, "Defining and monitoring service level agreements for dynamic e-business," in *Proc. of the 16th USENIX System Administration Conference*, no. November, 2002.
- [3] A. Keller and H. Ludwig, "The WSLA Framework : Specifying and Monitoring Service Level Agreements for Web Services," *Network*, vol. 11, no. 1, pp. 57-81, 2003.
- [4] M. Comuzzi and C. Kotsokalis, "Establishing and monitoring SLAs in complex service based systems," *Web Services*, 2009.
- [5] A. Michlmayr, F. Rosenberg, P. Leitner, and S. Dustdar, "Comprehensive qos monitoring of web services and event-based sla violation detection," in *4th Int. Workshop on Middleware for Service Oriented Comp.*, 2009.
- [6] B. Pernici, "Adaptation of web services based on QoS satisfaction," *Service-Oriented Computing*, pp. 65-75, 2011.
- [7] R. Kazhamiak, B. Wetzstein, D. Karastoyanova, M. Pistore, and F. Leymann, "Adaptation of service-based applications based on process quality factor analysis," in *ICSOC Workshops*, vol. 6275 of LNCS, 2009.
- [8] O. Moser and F. Rosenberg, "Non-intrusive monitoring and service adaptation for WS-BPEL," in *Proc. of the 17th Int. Conf. in WWW*, pp. 815-824, 2008.
- [9] M. Palacios, J. Garcia-Fanjul, J. Tuya, and C. De La Riva, "A Proactive Approach to Test Service Level Agreements," *5th Int. Conf. on Software Engineering Advances*, pp. 453-458, 2010.
- [10] M. Di Penta, G. Canfora, G. Esposito, V. Mazza, and M. Bruno, "Search-based testing of service level agreements," in *Proc. of the 9th Conf. on Genetic and evolutionary computation*, pp. 1090-1097, ACM, 2007.
- [11] SLA@SOI, "Reference architecture for an sla management framework (whitepaper of deliverable d.a1 a framework architecture)," 2011.
- [12] M. Wilson et al., "Trustcom framework v4," 2007.
- [13] K. Mahbub and G. Spanoudakis, "Monitoring ws-agreement s: An event calculus-based approach," in *Test and Analysis of Web Services*, pp. 265-306, Springer, 2007.
- [14] G. Spanoudakis and K. Mahbub, "Non-intrusive monitoring of service-based systems," *International Journal of Cooperative Information Systems*, vol. 15, no. 3, pp. 325-358, 2006.
- [15] S. Benbernou, L. Cavallaro, M. Sahid-Hacid, R. Kazhamiak, G. Kecskemeti, J.-L. Poizat, F. Silvestri, M. Uhlig, B. Wetzstein, "State of the Art Report, Gap Analysis of Knowledge on Principles, Techniques and Methodologies for Monitoring and Adaptation of SBAs (S-Cube deliverable PO-JRA-1.2.1)," 2008.
- [16] A. Andrieux et al., "Web Services Agreement Specification (WS-Agreement) (v. gfd-r.192)," 2011. OGF - Grid Resource Allocation Agreement Protocol WG.
- [17] A. Sahai, V. Machiraju, M. Sayal, A. P. A. v. Moorsel, and F. Casati, "Automated sla monitoring for web services," in *13th IEEE International Workshop on Distributed Systems: Operations and Management*, 2002.
- [18] C. Müller, A. Durán, M. Resinas, A. Ruiz-Cortés, and O. Martín-Díaz, "Experiences from building a ws-agreement document analyzer tool," *Tech. Rep. ISA-10-TR-03*, <http://www.isa.us.es/>, Jul 2010.
- [19] M. Comuzzi and G. Spanoudakis, "Dynamic set-up of monitoring infrastructures for service based systems," pp. 2414-2421, ACM, 2010.
- [20] F. Raimondi, J. Skene, and W. Emmerich, "Efficient online monitoring of web-service slas," in *16th ACM SIGSOFT Int. Symposium on Foundations of software engineering*, (New York, NY, USA), ACM, 2008.
- [21] M. Oriol, X. Franch, J. Marco, and D. Ameller, "Monitoring adaptable soa-systems using salmon," in *Workshop on Service Monitoring, Adaptation and Beyond (Mona+)*, pp. 19-28, 2008.
- [22] C. Müller, M. Resinas, and A. Ruiz-Cortés, "A Framework to Analyse WS-Agreement Documents," in *4th Workshop on Non-Functional Properties and SLA Management in SOC (NFPSLAM-SOC'10)*, 2010.
- [23] S. AG, *Standardized Technical Architecture Modeling: Conceptual and Design Level. Version 1.0*. SAP, 2007.
- [24] T. Erl, *SOA Design Patterns*. Prentice Hall PTR, 1st ed., 2009.
- [25] O. Sammodi, A. Metzger, X. Franch, M. Oriol, J. Marco, and K. Pohl, "Usage-based online testing for proactive adaptation of service-based applications," in *COMPASAC*, pp. 582-587, IEEE Comp. Soc., 2011.
- [26] A. Kertesz, G. Kecskemeti, M. Oriol, P. Kotcauer, S. Acs, M. Rodríguez, O. Mercè, A. Marosi, J. Marco, and X. Franch, "Enhancing federated cloud management with an integrated service monitoring approach," *Journal of Grid Computing*, pp. 1-22, 2013.
- [27] E. Tsang, *Foundations of Constraint Satisfaction*. Academic Press, 1995.
- [28] A. Ruiz-Cortés, O. Martín-Díaz, A. Durán, and M. Toro, "Improving the Automatic Procurement of Web Services using Constraint Programming," *I. J. on Cooperative Information Systems*, vol. 14, no. 4, 2005.
- [29] C. Müller, M. Resinas, and A. Ruiz-Cortés, "Automated Analysis of Conflicts in WS-Agreement Documents," *in press Transactions on Services Computing*, 2012.
- [30] C. Müller, M. Resinas, and A. Ruiz-Cortés, "Explaining the Non-Compliance between Templates and Agreement Offers in WS-Agreement*," in *Proc. of the 7th Int. Conf. on SOC*, pp. 237-252, 2009.
- [31] D. Harel and B. Rumpe, "Meaningful modeling: What's the semantics of "semantics"?," *IEEE Computer*, vol. 37, no. 10, pp. 64-72, 2004.
- [32] J. Rivera, E. Guerra, J. de Lara, and A. Vallecillo, "Analyzing rule-based behavioral semantics of visual modeling languages with Maude," in *Software Language Engineering*, LNCS, pp. 54-73, 2009.
- [33] A. Swearngin, B. Choueiry, and E. Freuder, "A reformulation strategy for multi-dimensional cpsps: The case study of the set game," 2011.

- 1100 [34] O. Martín-Díaz, A. Ruiz-Cortés, A. Durán, D. Benavides, and M. Toro,
1101 “Automating the Procurement of Web Services,” in *Proc. of The 1st Int.*
1102 *Conf. on Service-Oriented Computing*, LNCS, pp. 91–103, 2003.
- 1103 [35] O. Martín-Díaz, A. Ruiz-Cortés, A. Durán, and C. Müller, “An approach
1104 to temporal-aware procurement of web services,” in *Proc. of The 3rd Int.*
1105 *Conf. on Service-Oriented Computing*, pp. 170–184, 2005.
- 1106 [36] C. Müller, A. Ruiz-Cortés, and M. Resinas, “An Initial Approach to
1107 Explaining SLA Inconsistencies,” in *Proc. of the 6th Int. Conf. on*
1108 *Service-Oriented Computing (ICSOC)*, pp. 394–406, 2008.
- 1109 [37] X. Wang, T. Vitvar, M. Kerrigan, and I. Toma, “A qos-aware selection
1110 model for semantic web services,” in *Proc. of The 4th Int. Conf.*
1111 *in Service-Oriented Computing* (A. Dan and W. Lamersdorf, eds.),
1112 vol. 4294 of *LNCS*, pp. 390–401, Springer Berlin / Heidelberg, 2006.
- 1113 [38] F. Laburthe, N. Jussien, G. Rochart, H. Cambazard, C. Prud’homme,
1114 A. Malapert, and J. Menana, “Choco constraint solver web site,” 2010.
1115 <http://www.emn.fr/z-info/choco-solver/index.html>.
- 1116 [39] IBM, “IBM ILOG web site,” 2013. <http://www-01.ibm.com/software/websphere/ilog/>.
- 1117 [40] U. Junker, “Quickxplain: preferred explanations and relaxations for
1118 over-constrained problems,” in *Proc. of the 19th Int. Conf. on Artificial*
1119 *intelligence*, pp. 167–172, AAAI Press, 2004.
- 1120 [41] D. Hawkins, *Identification of outliers*. Chapman and Hall London, 1980.
- 1121 [42] R. Ratcliff, “Methods for dealing with reaction time outliers.,” *Psycho-*
1122 *logical bulletin*, vol. 114, pp. 510–532, Nov. 1993.
- 1123



1148 **Jordi Marco** is Associate professor and member
1149 of the GESSI research group at the Universitat
1150 Politècnica de Catalunya (UPC), Spain. He obtained
1151 his PhD and Msc in Computing from this Uni-
1152 versity. His current research lines include Service-
1153 Oriented Computing, conceptual modelling, con-
1154 tainer libraries, and computer graphics.

1155

1156 **Manuel Resinas** is a Lecturer at the University of
1157 Sevilla, Spain. He obtained his PhD in Computer
1158 Science from this University. His current research
1159 lines include analysis and management of service
1160 level agreements, business process compliance, and
1161 process performance management. Previously, he
1162 worked on automated negotiation of service level
1163 agreements.

1164

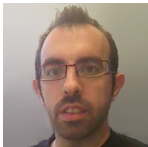
1124 **Carlos Müller** is a Research Assistant and member
1125 of the Applied Software Engineering Group (ISA,
1126 www.isa.us.es) at University of Sevilla, Spain. His
1127 current research line includes service oriented com-
1128 puting, specifically the automated analysis of service
1129 level agreements and the application of such analysis
1130 at SLA design and monitoring.



1124
1125
1126
1127
1128
1129
1130

1131

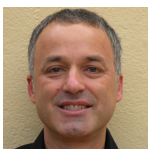
1132 **Marc Oriol** is a PhD student in Computer Science
1133 and member of the GESSI research group at the at
1134 Universitat Politècnica de Catalunya (UPC), Spain.
1135 He obtained a Msc in Computing from this Uni-
1136 versity. His current research lines include Service-
1137 Oriented Computing, Quality-of-Service and moni-
1138 toring.



1132
1133
1134
1135
1136
1137
1138

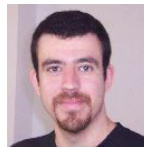
1139

1140 **Xavier Franch** is Associate professor and Head
1141 of the GESSI research group at the Universitat
1142 Politècnica de Catalunya (UPC), Spain. He obtained
1143 his PhD and Msc in Informatics from this University.
1144 His current research lines include Service-Oriented
1145 Computing, Requirements Engineering, Software
1146 Quality and Software Architecture, among others.



1140
1141
1142
1143
1144
1145
1146

1147



1165 **Antonio Ruiz-Cortés** is Associate Professor and
1166 Head of the Applied Software Engineering Group
1167 (ISA, www.isa.us.es) at University of Sevilla, Spain.
1168 He obtained his PhD in Computer Science from
1169 this University. His current research lines include
1170 service oriented computing, software product lines,
1171 and business process management.

1165
1166
1167
1168
1169
1170
1171

1172

1173 **Marc Rodríguez** is an undergraduate student and
1174 PAS researcher of the GESSI research group at the
1175 Universitat Politècnica de Catalunya (UPC), Spain.
1176 His current research lines include Service Oriented
1177 Computing, Quality-of-Service and monitoring.

1173
1174
1175
1176
1177

1178

