

# Graceful Interruption of Request-response Service Interactions\*

Mila Dalla Preda, Maurizio Gabbrielli, Ivan Lanese,  
Jacopo Mauro, and Gianluigi Zavattaro

Lab. Focus, Department of Computer Science/INRIA, University of Bologna, Italy.  
dallapre | gabbri | lanese | jmauro | zavattar@cs.unibo.it

**Abstract.** Bi-directional request-response interaction is a standard communication pattern in Service Oriented Computing (SOC). Such a pattern should be interrupted in case of faults. In the literature, different approaches have been considered: WS-BPEL discards the response, while Jolie waits for it in order to allow the fault handler to appropriately close the conversation with the remote service. We investigate an intermediate approach in which it is not necessary for the fault handler to wait for the response, but it is still possible on response arrival to gracefully close the conversation with the remote service.

## 1 Introduction

Service-oriented computing (SOC) is a programming paradigm based on the composition of services, computational entities available on the net. According to WSDL [9], the standard for describing web service interfaces, services can be invoked according to two main modalities: one-way and request-response. In one-way communication a message is sent to a remote service. In request-response communication a message is sent and an answer is waited for before continuing the computation.

Interaction with remote services may incur in errors of different kinds: the remote service may disconnect, messages may be lost, or a client may interrupt the interaction with a remote service exactly in between the request and the corresponding response. To avoid that such an error causes the failure of the whole application, error handling techniques have been developed. They are commonly based on the concept of fault handler and compensation. A fault handler is a piece of code devoted to take the application to a consistent state after a fault has been caught. A compensation is a piece of code devoted to undoing the effect of a previous activity because of a later error.

As an example, consider a hotel reservation service. A reservation can be canceled, but if it is not annulled the cost of one night will be charged in case of no show. If the trip has to be annulled, the compensation for the hotel reservation has to be executed, thus canceling the reservation and avoiding the cost of a no show.

Jolie [3] is a language for programming service-oriented applications. Jolie request-response invocation establishes a strong connection between the caller and callee, thus it should not be disrupted by faults. To this end, callee faults are notified to the caller that can thus manage them. Symmetrically, in case of caller faults the answer from the

\* Partly funded by the projects EU FP7-231620 HATS and ANR-2010-SEGI-013 AEOLUS.

callee is waited for and used during recovery. This allows, in particular, to compensate successful remote activities which are no more needed because of the local fault. This is the case of the hotel reservation above.

WS-BPEL [8], a main standard in the field, has a different approach: in case of caller faults execution can continue without waiting for the response, and the response is discarded upon arrival. In particular, it is not possible to write code that will be executed upon receipt of the response. The Jolie approach allows for programming safer applications. The fact that the request-response pattern is not disrupted by errors has been proved in [3], by relying on SOCK [4, 2], a calculus defining the formal semantics of Jolie. A nasty side effect of the Jolie approach is that the client has to wait for answers of request-response invocations before proceeding in its execution. This slows down the caller execution. For instance, referring to the hotel reservation example, the client cannot continue its operations before the answer from the hotel has been received and (s)he gets stuck whenever the answer is lost. This drawback is unacceptable for programming applications over the net. Such a kind of problem is normally solved using timeouts, but they are not available in Jolie. Also, they are not easy to mimic.

We propose here a new approach to error handling in Jolie, allowing on one side to compensate undesired remote side effects, and ensuring on the other side that local computation is not slowed down in case of late answers. In particular, this new approach allows to easily program timeouts. We also extend the approach to deal with concurrent invocations of multiple services, as needed for implementing speculative parallelism.

## 2 SOCK

We first introduce SOCK [4], the calculus that defines the semantics of Jolie [3] programs, and then we extend it to account for request-response and multiple request-response service invocations. SOCK is suitable for illustrating our approach since it has a formal SOS semantics, it provides request-response as a native operator, and it has a refined approach to error handling.

In the following we present the three layers in which SOCK is structured, omitting the aspects that are not central for us (see [4] for a detailed description).

**Service behavior layer.** The service behavior layer describes the actions performed by services. Actions can be operations on the state or communications. Services are identified by the name of their operations, and by their location.

SOCK error handling is based on the concepts of *scope*, *fault*, and *compensation*. A scope is a process container denoted by a unique name. A fault is a signal raised by a process when an error state is reached. A compensation is used either to smoothly stop a running activity in case of an external fault, or to compensate the activity after its successful termination (this encompasses both WS-BPEL termination and compensation mechanisms). Recovering mechanisms are implemented by exploiting processes called *handlers*. We use *fault handlers* and *compensation handlers*. They are executed to manage respectively internal faults and external faults/compensation requests.

SOCK *syntax* is based on the following (disjoint) sets:  $Var$ , ranged over by  $x, y$ , for variables,  $Val$ , ranged over by  $v$ , for values,  $\mathcal{O}$ , ranged over by  $o$ , for one-way operations,  $Faults$ , ranged over by  $f$ , for faults, and  $Scopes$ , ranged over by  $q$ , for

$P, Q ::= \bar{o}@l(\mathbf{y})$	output	$o(\mathbf{x})$	input
$x := e$	assignment	$P; Q$	sequence
$P Q$	parallel comp.	$\sum_{i \in I} o_i(\mathbf{x}_i); P_i$	external choice
$\text{if } \chi \text{ then } P \text{ else } Q$	det. choice	$\text{while } \chi \text{ do } (P)$	iteration
$\mathbf{0}$	null process	$\{P : \mathcal{H} : u\}_{q\perp}$	active scope
$\text{inst}(\mathcal{H})$	install handler	$\text{throw}(f)$	throw
$\text{comp}(q)$	compensate	$\langle P \rangle$	protection

**Table 1.** Service behavior syntax with faults

scope names.  $Loc$  is a subset of  $Val$  containing locations, ranged over by  $l$ . We denote as  $SC$  the set of service behavior processes, ranged over by  $P, Q, \dots$ . We use  $q\perp$  to range over  $Scopes \cup \{\perp\}$ , whereas  $u$  ranges over  $Faults \cup Scopes \cup \{\perp\}$ . Here  $\perp$  is used to specify that a handler is undefined.  $\mathcal{H}$  denotes a function from  $Faults$  and  $Scopes$  to processes (or  $\perp$ ). The function associating  $P_i$  to  $u_i$  for  $i \in \{1, \dots, n\}$  is  $[u_1 \mapsto P_1, \dots, u_n \mapsto P_n]$ . Finally, we use the notation  $\mathbf{k} = \langle k_0, k_1, \dots, k_i \rangle$  for vectors.

The syntax of service behavior processes is defined in Table 1. A one-way output  $\bar{o}@l(\mathbf{y})$  invokes the operation  $o$  of a service at location  $l$ , where  $\mathbf{y}$  are the variables that specify the values to be sent. Dually, in a one-way  $o(\mathbf{x})$ ,  $\mathbf{x}$  contains the variables that will receive the communicated values. Assignment, sequence, parallel composition, external and deterministic choice, iteration, and null process are standard.

We denote with  $\{P\}_q$  a scope named  $q$  executing process  $P$ . An active scope has instead the form  $\{P : \mathcal{H} : u\}_{q\perp}$ , where  $\mathcal{H}$  specifies the defined handlers. Term  $\{P\}_q$  is a shortcut for  $\{P : \mathcal{H}_0 : \perp\}_q$ , where  $\mathcal{H}_0$  evaluates to  $\perp$  for all fault names and to  $\mathbf{0}$  for all scope names. The argument  $u$  is the name of a handler waiting to be executed, or  $\perp$  if there is no such handler. When a scope has failed its execution, either because it has been killed from a parent scope, or because it has not been able to manage an internal fault, it reaches a zombie state. Zombie scopes have  $\perp$  as scope name. Primitives  $\text{throw}(f)$  and  $\text{comp}(q)$  respectively raises fault  $f$  and asks to compensate scope  $q$ .  $\langle P \rangle$  executes  $P$  in a protected way, i.e. not influenced by external faults. Handlers are installed into the nearest enclosing scope by  $\text{inst}(\mathcal{H})$ , where  $\mathcal{H}$  is the required update of the handler function. We assume that  $\text{comp}(q)$  occurs only within handlers, and  $q$  can only be a child of the enclosing scope. For each  $\text{inst}(\mathcal{H})$ ,  $\mathcal{H}$  is defined only on fault names and on the name of the nearest enclosing scope. Finally, scope names are unique.

The service behavior layer *semantics* generates all the transitions allowed by the process behavior, specifying the constraints on the state that have to be satisfied for them to be performed. The state is a substitution of values for variables. We use  $\sigma$  to range over substitutions, and write  $[v/x]$  for the substitution assigning values in  $v$  to variables in  $x$ . Given a substitution  $\sigma$ ,  $\text{Dom}(\sigma)$  is its domain.

Let  $Act$  be the set of labels of the semantics, ranged over by  $a$ . We use structured labels of the form  $\iota(\sigma : \theta)$  where  $\iota$  is the kind of action while  $\sigma$  and  $\theta$  are substitutions containing respectively the assumptions and the effects on the state. We also use the unstructured labels  $th(f), cm(q, P), inst(\mathcal{H})$ . We use operator  $\boxplus$  for updating the handler function:

$$(\mathcal{H} \boxplus \mathcal{H}')(u) = \begin{cases} \mathcal{H}'(u) & \text{if } u \in \text{Dom}(\mathcal{H}') \\ \mathcal{H}(u) & \text{otherwise} \end{cases}$$

Intuitively, handlers in  $\mathcal{H}'$  replace the corresponding ones in  $\mathcal{H}$ . We also use  $\text{cmp}(\mathcal{H})$  to denote the part of  $\mathcal{H}$  dealing with compensations.

$\frac{}{\bar{o} @ l(\mathbf{x}) \xrightarrow{\bar{o}(v) @ l(v/x:\emptyset)} \mathbf{0}} \quad \text{(ONE-WAYOUT)}$	$\frac{}{o(\mathbf{x}) \xrightarrow{o(v)(\emptyset:v/x)} \mathbf{0}} \quad \text{(ONE-WAYIN)}$
$\frac{\text{Dom}(\sigma) = \text{Var}(e) \quad \llbracket e\sigma \rrbracket = v}{x := e \xrightarrow{\tau(\sigma:v/x)} \mathbf{0}} \quad \text{(ASSIGN)}$	$\frac{\text{Dom}(\sigma) = \text{Var}(\chi) \quad \llbracket \chi\sigma \rrbracket = \text{true}}{\text{if } \chi \text{ then } P \text{ else } Q \xrightarrow{\tau(\sigma:\emptyset)} P} \quad \text{(IF-THEN)}$
$\frac{P \xrightarrow{a} P'}{P; Q \xrightarrow{a} P'; Q} \quad \text{(SEQUENCE)}$	$\frac{P \xrightarrow{a} P'}{P \mid Q \xrightarrow{a} P' \mid Q} \quad \text{(PARALLEL)}$
$\frac{o_i(\mathbf{x}_i) \xrightarrow{a} Q_i \quad i \in I}{\sum_{i \in I} o_i(\mathbf{x}_i); P_i \xrightarrow{a} Q_i; P_i} \quad \text{(CHOICE)}$	
<p>STRUCTURAL CONGRUENCE</p> $P \mid Q \equiv Q \mid P \quad P \mid \mathbf{0} \equiv P \quad P \mid (Q \mid R) \equiv (P \mid Q) \mid R \quad \mathbf{0}; P \equiv P \quad \langle \mathbf{0} \rangle \equiv \mathbf{0}$	

**Table 2.** Standard rules for service behavior layer ( $a \neq th(f)$ )

The SOCK semantics is defined as a relation  $\rightarrow \subseteq SC \times Act \times SC$ . The main rules for standard actions are in Table 2, while Table 3 defines the fault handling mechanism.

Rule ONE-WAYOUT defines the output operation, where  $v/x$  is the assumption on the state. Rule ONE-WAYIN corresponds to the input operation: it makes no assumption on the state, but it specifies a state update. The other rules in Table 2 are standard. The internal process  $P$  of a scope can execute thanks to rule SCOPE in Table 3. Handlers are installed in the nearest enclosing scope by rules ASKINST and INSTALL. According to rule SCOPE-SUCCESS, when a scope successfully ends, its compensation handlers are propagated to the parent scope. Compensation execution is required by rule COMPENSATE. The actual compensation code  $Q$  is guessed, and the guess is checked by rule COMPENSATION. Faults are raised by rule THROW. A fault is caught by rule CATCH-FAULT when a scope defining the corresponding handler is met. Activities involving the termination of a sub-scope and the termination of internal error recovery are managed by the rules for fault propagation THROW-SYNC, THROW-SEQ and RETHROW, and by the partial function *killable*. Function *killable* computes the activities that have to be completed before the handler is executed and it is applied to parallel components by rule THROW-SYNC. Moreover, function *killable* guarantees that when a fault is thrown there is no pending handler update. This is obtained by making *killable*( $P, f$ ) undefined (and thus rule THROW-SYNC not applicable) if some handler installation is pending in  $P$ . The  $\langle P \rangle$  operator (described by rule PROTECTION) guarantees that the enclosed activity will not be killed by external faults. Rule SCOPE-HANDLE-FAULT executes a handler for a fault. A scope that has been terminated from the outside is in zombie state. It can execute its compensation handler thanks to rule SCOPE-HANDLE-TERM, and then terminate with failure using rule SCOPE-FAIL. Similarly, a scope enters the zombie state when reached by a fault it cannot handle (rule RETHROW). The fault is propagated up along the scope hierarchy. Zombie scopes cannot throw faults any more, since rule IGNORE-FAULT has to be applied instead of RETHROW.

**Service engine layer.** The service engine layer manages the service state and instances. A service engine  $Y$  can be a session  $(P, \mathcal{S})$ , where  $P$  is a service behavior process and  $\mathcal{S}$  is a state, or a parallel composition  $Y|Y$  of them. The service engine layer allows to propagate only labels such that the condition  $\sigma$  (if available) is satisfied by the current state, and applies to the state the state update  $\rho$ .

$$\begin{array}{c}
\text{(SCOPE)} \quad \frac{P \xrightarrow{a} P' \quad a \neq \text{inst}(\mathcal{H}), \text{cm}(q', \mathcal{H}')}{\{P : \mathcal{H} : u\}_{q_{\perp}} \xrightarrow{a} \{P' : \mathcal{H} : u\}_{q_{\perp}}} \quad \text{(INSTALL)} \quad \frac{P \xrightarrow{\text{inst}(\mathcal{H})} P'}{\{P : \mathcal{H}' : u\}_{q_{\perp}} \xrightarrow{\tau(\emptyset; \emptyset)} \{P' : \mathcal{H}' \boxplus \mathcal{H} : u\}_{q_{\perp}}} \\
\text{(ASKINST)} \quad \frac{\text{inst}(\mathcal{H}) \xrightarrow{\text{inst}(\mathcal{H})} \mathbf{0}}{\{P : \mathcal{H} : u\}_{q_{\perp}} \xrightarrow{a} \{P' : \mathcal{H} : u\}_{q_{\perp}}} \quad \text{(THROW)} \quad \frac{\text{throw}(f) \xrightarrow{th(f)} \mathbf{0}}{\{P : \mathcal{H}' : u\}_{q_{\perp}} \xrightarrow{\tau(\emptyset; \emptyset)} \{P' : \mathcal{H}' \boxplus \mathcal{H} : u\}_{q_{\perp}}} \quad \text{(COMPENSATE)} \quad \frac{\text{comp}(q) \xrightarrow{\text{cm}(q, Q)} Q}{\{P : \mathcal{H}' : u\}_{q_{\perp}} \xrightarrow{\tau(\emptyset; \emptyset)} \{P' : \mathcal{H}' \boxplus \mathcal{H} : u\}_{q_{\perp}}} \\
\text{(SCOPE-SUCCESS)} \quad \frac{\text{inst}(\text{cmp}(\mathcal{H})) \xrightarrow{\text{inst}(\text{cmp}(\mathcal{H}))} \mathbf{0}}{\{\mathbf{0} : \mathcal{H} : \perp\}_q \xrightarrow{\text{inst}(\text{cmp}(\mathcal{H}))} \mathbf{0}} \quad \text{(SCOPE-HANDLE-FAULT)} \quad \frac{\text{throw}(f) \xrightarrow{th(f)} \mathbf{0}}{\{\mathbf{0} : \mathcal{H} : f\}_{q_{\perp}} \xrightarrow{\tau(\emptyset; \emptyset)} \{\mathcal{H}(f) : \mathcal{H} \boxplus [f \mapsto \perp] : \perp\}_{q_{\perp}}} \\
\text{(COMPENSATION)} \quad \frac{P \xrightarrow{\text{cm}(q, Q)} P', \mathcal{H}(q) = Q}{\{P : \mathcal{H} : u\}_{q_{\perp}} \xrightarrow{\tau(\emptyset; \emptyset)} \{P' : \mathcal{H} \boxplus [q \mapsto \mathbf{0}] : u\}_{q_{\perp}}} \\
\text{(SCOPE-HANDLE-TERM)} \quad \frac{\text{throw}(f) \xrightarrow{th(f)} \mathbf{0}}{\{\mathbf{0} : \mathcal{H} : q\}_{\perp} \xrightarrow{\tau(\emptyset; \emptyset)} \{\mathcal{H}(q) : \mathcal{H} \boxplus [q \mapsto \mathbf{0}] : \perp\}_{\perp}} \quad \text{(SCOPE-FAIL)} \quad \frac{\text{throw}(f) \xrightarrow{th(f)} \mathbf{0}}{\{\mathbf{0} : \mathcal{H} : \perp\}_{\perp} \xrightarrow{\tau(\emptyset; \emptyset)} \mathbf{0}} \\
\text{(PROTECTION)} \quad \frac{P \xrightarrow{a} P'}{\langle P \rangle \xrightarrow{a} \langle P' \rangle} \quad \text{(THROW-SYNC)} \quad \frac{P \xrightarrow{th(f)} P', \text{killable}(Q, f) = Q'}{P|Q \xrightarrow{th(f)} P'|Q'} \quad \text{(THROW-SEQ)} \quad \frac{P \xrightarrow{th(f)} P'}{P; Q \xrightarrow{th(f)} P'} \\
\text{(CATCH-FAULT)} \quad \frac{P \xrightarrow{th(f)} P', \mathcal{H}(f) \neq \perp}{\{P : \mathcal{H} : u\}_{q_{\perp}} \xrightarrow{\tau(\emptyset; \emptyset)} \{P' : \mathcal{H} : f\}_{q_{\perp}}} \quad \text{(IGNORE-FAULT)} \quad \frac{P \xrightarrow{th(f)} P', \mathcal{H}(f) = \perp}{\{P : \mathcal{H} : u\}_{\perp} \xrightarrow{\tau(\emptyset; \emptyset)} \{P' : \mathcal{H} : u\}_{\perp}} \\
\text{(RETHROW)} \quad \frac{P \xrightarrow{th(f)} P', \mathcal{H}(f) = \perp}{\{P : \mathcal{H} : u\}_q \xrightarrow{th(f)} \{\{P' : \mathcal{H} : \perp\}_{\perp}\}}
\end{array}$$

where

$$\begin{aligned}
& \text{killable}(\{P : \mathcal{H} : u\}_q, f) = \{\{\text{killable}(P, f) : \mathcal{H} : q\}_{\perp}\} \text{ if } P \neq \mathbf{0} \\
& \text{killable}(P | Q, f) = \text{killable}(P, f) | \text{killable}(Q, f) \\
& \text{killable}(P; Q, f) = \text{killable}(P, f) \text{ if } P \neq \mathbf{0} \\
& \text{killable}(\langle P \rangle, f) = \langle P \rangle \text{ if } \text{killable}(P, f) \text{ is defined} \\
& \text{killable}(P, f) = \mathbf{0} \text{ if } P \in \{\mathbf{0}, o(\mathbf{x}), \bar{o} @ l(\mathbf{x}), x := e, \text{if } \chi \text{ then } P \text{ else } Q, \text{while } \chi \text{ do } (P) \\
& \quad \sum_{i \in W} o_i(\mathbf{x}_i); P_i, \text{throw}(f), \text{comp}(q)\}
\end{aligned}$$

**Table 3.** Faults-related rules for service behavior layer ( $a \neq th(f)$ )

**Services system layer.** The service system layer allows the interaction between different engines. A service system  $E$  can be a located service engine  $Y @ l$  or a parallel composition  $E \parallel E$  of them. The services system layer just allows complementary communication actions to interact, transforming them into internal steps  $\tau$ , and propagates the other actions.

### 3 Request-response interaction pattern

A request-response pattern is a bi-directional interaction where a client sends a message to a server and waits for an answer. When a server receives such a message, it elaborates

the answer and sends it back to the client. In the literature there are two proposals to deal with a client that fails during a request-response interaction. The WS-BPEL approach kills the receive activity and, when the message arrives, it is silently discarded. In Jolie instead, clients always wait for the answer and exploit it for error recovery.

Here we present an intermediate approach: in case of failure we wait for the answer, but without blocking the computation. Moreover, when the answer is received we allow for the execution of a compensation activity. Let  $\mathcal{O}_r$  be the set of request-response operations, ranged over by  $o_r$ . We define the request-response pattern in terms of the output primitive  $\overline{o_r} @ l(\mathbf{y}, \mathbf{x}, P)$ , also called solicit, and of the input primitive  $o_r(\mathbf{x}_1, \mathbf{y}_1, Q)$ . When interacting, the client sends the values from variables  $\mathbf{y}$  to the server, that stores them in variables  $\mathbf{x}_1$ . Then, the server executes process  $Q$  and, when  $Q$  terminates, the values in variables  $\mathbf{y}_1$  are sent back to the client who stores them in variables  $\mathbf{x}$ . Only at this point the execution of the client can restart. If a fault occurs on the client-side after the remote service has been invoked, but before the answer is received, we allow the client to handle the fault regardless of the reply, so that recovery can start immediately. However, we create a receiver for the missing message in a fresh session so that, if later on the message is received, the operation can be compensated. The compensation is specified by the parameter  $P$  of the solicit operation. If instead a fault is raised on the server-side during the computation of the answer, the fault is propagated to the client where it raises a local fault. In this case there is no need to compensate the remote invocation, since we assume that this is dealt with by local recovery of the server.

**Service behavior calculus - extension.** We extend here the behavioral layer with the request-response and with few auxiliary operators used to define its semantics.

$\overline{o_r} @ l(\mathbf{y}, \mathbf{x}, P)$	Solicit	$o_r(\mathbf{x}_1, \mathbf{y}_1, Q)$	Request-Response
$Exec(l, o_r, \mathbf{y}, P)$	Req.-Resp. execution	$Wait(o_r, \mathbf{y}, P)$	Wait
$\overline{o_r} ! f @ l$	Fault output	$Bubble(P)$	Bubble

$Exec(l, o_r, \mathbf{y}, P)$  is a server-side running request-response:  $P$  is the process computing the answer,  $o_r$  the name of the operation,  $\mathbf{y}$  the vector of variables to be used for the answer, and  $l$  the client location. Symmetrically,  $Wait(o_r, \mathbf{y}, P)$  is the process waiting for the response on client-side:  $o_r$  is request-response operation,  $\mathbf{y}$  is the vector of variables for storing the answer and  $P$  is the compensation code to run in case the client fails before the answer is received. When a fault is triggered on the server-side, an error notification has to be sent to the client: this is done by  $\overline{o_r} ! f @ l$ , where  $o_r$  is the operation,  $f$  the fault and  $l$  the client location. If a fault occurs on client-side, we have to move the receipt operation to a fresh, parallel session, so that error recovery can start immediately. This is done by the primitive  $Bubble(P)$ , which allows to create a new session (a “bubble”) executing code  $P$ . This primitive is the key element that allows a failed solicit to wait for a response outside its scope and potentially allowing its termination regardless of the arrival of the answer.

The semantics of the behavior layer is extended with the rules presented in Table 4 (the last rule refers to the engine). Function *killable* is also extended, as follows:

- $killable(Exec(l, o_r, \mathbf{y}, P), f) = killable(P, f) | \langle \overline{o_r} ! f @ l \rangle$
- $killable(Wait(o_r, \mathbf{x}, P), f) = Bubble(Wait(o_r, \mathbf{x}, \mathbf{0}); P)$
- $killable(\overline{o_r} ! f @ l, f) = \overline{o_r} ! f @ l$
- $killable(Bubble(P), f) = Bubble(P)$

<p>(SOLICIT)</p> $\overline{o_r} @ l(\mathbf{y}, \mathbf{x}, P) \xrightarrow{\overline{o_r}(\mathbf{v}) @ l(\emptyset: \mathbf{v} / \mathbf{x})} Wait(o_r, \mathbf{x}, P)$ <p>(REQUEST-EXEC)</p> $\frac{P \xrightarrow{a} P'}{Exec(l, o_r, \mathbf{y}, P) \xrightarrow{a} Exec(l, o_r, \mathbf{y}, P')}$ <p>(REQUEST-RESPONSE)</p> $Exec(l, o_r, \mathbf{y}, \mathbf{0}) \xrightarrow{\overline{o_r}(\mathbf{v}) @ l(\mathbf{v} / \mathbf{y} : \emptyset)} \mathbf{0}$ <p>(SEND-FAULT)</p> $\overline{o_r} ! f @ l \xrightarrow{\overline{o_r}(f) @ l(\emptyset: \emptyset)} \mathbf{0}$ <p>(CREATE BUBBLE)</p> $Bubble(P) \xrightarrow{\tau(\emptyset: \emptyset) [P]} \mathbf{0}$	<p>(REQUEST)</p> $o_r(\mathbf{x}, \mathbf{y}, P) \xrightarrow{o_r(\mathbf{v}) :: l(\emptyset: \mathbf{v} / \mathbf{x})} Exec(l, o_r, \mathbf{y}, P)$ <p>(THROW-REEXEC)</p> $\frac{P \xrightarrow{th(f)} P'}{Exec(l, o_r, \mathbf{y}, P) \xrightarrow{th(f)} P'   \langle \overline{o_r} ! f @ l \rangle}$ <p>(SOLICIT-RESPONSE)</p> $Wait(o_r, \mathbf{x}, P) \xrightarrow{o_r(\mathbf{v}) (\emptyset: \mathbf{v} / \mathbf{x})} \mathbf{0}$ <p>(RECEIVE FAULT)</p> $Wait(o_r, \mathbf{x}, P) \xrightarrow{o_r(f) (\emptyset: \emptyset)} throw(f)$ <p>(ENGINE-BUBBLE)</p> $\frac{P \xrightarrow{\tau(\emptyset: \emptyset) [Q]} P' \quad Q \neq \mathbf{0}}{(P, S) \xrightarrow{\tau} (P', S)   (Q, S)}$
---	---

**Table 4.** Request-response pattern and engine rules

Rules SOLICIT and REQUEST start a solicit-response operation on client and server side respectively. Upon invocation, the request-response becomes an active construct executing process  $P$ , and storing all the information needed to send back the answer. The execution of  $P$  is managed by rule REQUEST-EXEC. When  $P$  terminates, rule REQUEST-RESPONSE sends back an answer. This synchronizes with rule SOLICIT-RESPONSE on the client side.

A running request-response reached by a fault is transformed into a fault notification (see rule THROW-REEXEC and the definition of function *killable*) on server side. Fault notification is executed by rule SEND-FAULT, and it interacts with the waiting receive thanks to rule RECEIVE-FAULT. When received, the fault is re-thrown at the client side.

A fault on client side instead gives rise to a bubble, creating the process that will wait for the answer in a separate session. The bubble is created by rule CREATE BUBBLE, and will be installed at the service engine level by rule ENGINE-BUBBLE. The label for bubble creation has the form  $\tau(\emptyset : \emptyset) [P]$ , where  $P$  is the process to be run inside the new session. We will write  $\tau(\emptyset : \emptyset)$  for  $\tau(\emptyset : \emptyset) [\mathbf{0}]$ . The new receive operation inside the bubble has no handler update, since it will be executed out of any scope, and its compensating code  $P$  has been promoted as a continuation. In this way,  $P$  will be executed only in case of successful answer. In case of faulty answer, the generated fault will have no effect since it is in a session on its own.

**Service engine calculus - extension.** We have to add to the service engine layer a rule for installing bubbles: when a bubble reaches the service engine layer, a new session is started executing the code inside the bubble (rule ENGINE-BUBBLE in Table 4).

**Service system calculus - extension.** The service system calculus semantics is extended by allowing the labels for request-response communication to be matched.

**Example.** We present now an example of usage of the request-response primitive. A first solution for the hotel reservation example described in the introduction is:

```
CLIENT ::=  $\overline{book_r}$ @hotel_Imperial(\langle CC, dates \rangle, \langle res\_num \rangle,
            $\text{annul}@hotel_Imperial(\langle res\_num \rangle) ) ;$ 
```

P

The  $book_r$  operation transmits the credit card number  $CC$  and the dates of the reservation and waits for the reservation number. In case the user wants to cancel the reservation before receiving an answer from the hotel a fault can be used to kill this operation. In such a case the  $annul$  operation is invoked when the answer is received to compensate the  $book_r$  operation. The  $annul$  operation will be executed in a new session by using our mechanism based on bubbles.

As a more concrete instance, we could consider the case where the user is willing to wait a limited amount of time for the answer from the hotel, after which (s)he will cancel the reservation. This case could be programmed by assuming a service  $timeout$  that offers a request-response operation that sends back an answer after  $n$  seconds<sup>1</sup>:

```
CLIENT ::=
  res_num := 0;
  { inst(f ↦ if res_num == 0 then throw(tm));
    (timeout_r @ timeout(⟨60⟩, ⟨⟩, 0); throw(f)
    | book_r @ hotel_Imperial(⟨CC, dates⟩, ⟨res_num⟩,
      annul @ hotel_Imperial(⟨res_num⟩)); throw(f) )
  }_q ; P
```

In this scenario the timeout operation is in parallel with the booking. The first operation that finishes raises the fault  $f$  that is caught by the handler of the scope  $q$ . The fault will kill the remaining operation and if the hotel response has not arrived yet (i.e. the value of  $res\_num$  is still 0) then the fault  $tm$  is raised.  $P$  is executed otherwise.

A similar solution is not viable in BPEL: in case of timeout, the booking invocation is killed, and if an answer arrives, it is discarded. Thus one does not know whether the invocation succeeded or not, neither the reservation number in case of success.

In Jolie, the answer is used for error recovery. However, in case no answer is received from the booking service, the whole service engine gets stuck. In our approach instead the main session can continue its execution without delays.

It is difficult to apply the proposed solution when two or more solicits install handlers or require compensation. One may try to exploit the handler update primitive, but in this way compensations are executed inside the scope, thus they have to be terminated before execution can proceed. This problem, and other technical difficulties, justify the multiple solicit response primitive introduced in the next section.

## 4 Multiple Request-Response communication pattern

The request-response pattern allows one invocation to be sent and one answer to be received. For optimization reasons, it may be important to invoke many services in parallel, and only consider the first received answer (speculative parallelism).

We model this communication pattern using a dedicated primitive that we call multiple solicit-response (MSR for short). A MSR consists of a list of solicit-responses, each one equipped with its own continuation. Formally, we define the syntax of the MSR primitive as  $MSR\{z_1, \dots, z_n\}$  where each  $z_i$  is a *solicit-response with continuation* written  $z_i = \overline{\sigma_{r_i}} @ l_i(\mathbf{y}_i, \mathbf{x}_i, P_i) \mapsto Q_i$ . Intuitively, the continuation  $Q_i$  is executed only when  $\overline{\sigma_{r_i}} @ l_i(\mathbf{y}_i, \mathbf{x}_i, P_i)$  is the first to receive a successful answer.

<sup>1</sup> Clearly, because of network delay the answer may be received later than expected.



$$\begin{array}{c}
\text{(MSR-SOLICIT)} \\
z_1 = \overline{o_r} @ l(\mathbf{y}, \mathbf{x}, P) \mapsto Q \quad w_{m+1} = \text{Wait}(o_r, \mathbf{y}, P) \mapsto Q \\
\hline
\text{Wait}^+(z_1, \dots, z_n \triangleright w_1, \dots, w_m) \xrightarrow{\overline{o_r}(\mathbf{v}) @ l(\emptyset: \mathbf{v} / \mathbf{x})} \text{Wait}^+(z_2, \dots, z_n \triangleright w_1, \dots, w_m, w_{m+1}) \\
\text{(MSR-RESPONSE)} \\
\forall k \in \{1, \dots, n\} : w_k = \text{Wait}(o_{r_k}, \mathbf{y}_k, P_k) \mapsto Q_k \quad i \in \{1, \dots, n\} \quad J = \{1, \dots, n\} \setminus \{i\} \\
\hline
\text{Wait}^+(\triangleright w_1, \dots, w_n) \xrightarrow{o_{r_i}(\mathbf{v}) @ l(\emptyset: \mathbf{v} / \mathbf{y}_i)} Q_i | \prod_{j \in J} \text{Bubble}(\text{Wait}(o_{r_j}, \mathbf{y}_j, \mathbf{0}); P_j) \\
\text{(MSR-IGNORE FAULT)} \\
n > 1 \quad w_i = \text{Wait}(o_{r_i}, \mathbf{y}_i, P_i) \mapsto Q_i \quad i \in \{1, \dots, n\} \\
\hline
\text{Wait}^+(\triangleright w_1, \dots, w_n) \xrightarrow{o_{r_i}(f) @ l(\emptyset: \emptyset)} \text{Wait}^+(\triangleright w_1, \dots, w_{i-1}, w_{i+1}, \dots, w_n) \\
\hline
\text{Wait}^+(\triangleright \text{Wait}(o_r, \mathbf{y}, P) \mapsto Q) \equiv \text{Wait}(o_r, \mathbf{y}, P); Q
\end{array}$$

**Table 5.** Multiple request-response pattern rules

**Service behavior calculus - extension.** We extend the service behavior calculus with the MSR primitive and with some auxiliary operators:

$$\begin{array}{ll}
P, Q ::= \dots & \\
MSR\{z_1, \dots, z_n\} & \text{multiple solicit-response} \\
\text{Wait}^+(z_1, \dots, z_n \triangleright w_1, \dots, w_m) & \text{multiple wait} \\
z ::= \overline{o_r} @ l(\mathbf{y}, \mathbf{x}, P) \mapsto Q & \text{solicit with continuation} \\
w ::= \text{Wait}(o_r, \mathbf{y}, P) \mapsto Q & \text{wait with continuation}
\end{array}$$

In a MSR the solicits are sent one after the other, and only when all the requests have been sent the MSR can receive a response. For this reason we introduce the multiple wait  $\text{Wait}^+(z_1, \dots, z_n \triangleright w_1, \dots, w_m)$  that specifies the solicits that still have to be sent  $z_1, \dots, z_n$ , and the ones that will wait for an answer  $w_1, \dots, w_m$ . Thus, the MSR primitive  $MSR\{z_1, \dots, z_n\}$  above is a shortcut for  $\text{Wait}^+(z_1, \dots, z_n \triangleright)$ . Moreover, we have that a multiple wait with only one waiting process is structurally equivalent to a standard wait. We formally define the behavior of the MSR primitive by extending the service behavior semantics with the rules presented in Table 5.

The multiple wait executes all the solicit-responses through rule MSR-SOLICIT. Once all the solicits have been sent, the multiple wait receives a successful answer through rule MSR-RESPONSE. It continues the execution with the corresponding continuation code, and kills all the other solicits by creating a bubble for each remaining waiting process. If a fault notification arrives as an answer, it is discarded by rule MSR-IGNORE FAULT if there is at least another available wait. If instead there is no other solicit waiting for an answer, the last fault received is raised (rule RECEIVE FAULT described in Table 4). When an external fault arrives a bubble containing a dead solicit response is created for every solicit that has been sent, as specified by the function *killable* that is extended in the following way:

$$\begin{aligned}
\text{killable}(\text{Wait}^+(z_1, \dots, z_n \triangleright w_1, \dots, w_m), f) = \\
\prod_{\text{Wait}(o_{r_j}, \mathbf{y}_j, P_j) \mapsto Q_j \in \{w_1, \dots, w_m\}} \text{Bubble}(\text{Wait}(o_{r_j}, \mathbf{y}_j, \mathbf{0}); P_j)
\end{aligned}$$

The MSR primitive is perfectly suited to capture speculative parallelism scenarios. Consider for instance the hotel reservation problem defined in the introduction. Suppose to use two booking services for making the hotel reservation, and that you would like to get the acknowledgment in 1 minute. If the booking services are located at A and B and if we use the *timeout* service introduced before, this service could be defined as:

```
CLIENT ::= msr {
   $\overline{timeout}_r @ \text{timeout} (\langle 60 \rangle, \langle \rangle, \mathbf{0}) \mapsto \text{throw}(tm)$ 
   $\overline{book}_r @ H_1 (\langle CC, \text{dates} \rangle, \langle \text{res\_num} \rangle, \text{annul}@H_1 (\langle \text{res\_num} \rangle)) \mapsto \mathbf{0}$ 
   $\overline{book}_r @ H_2 (\langle CC, \text{dates} \rangle, \langle \text{res\_num} \rangle, \text{annul}@H_2 (\langle \text{res\_num} \rangle)) \mapsto \mathbf{0}$ 
}
```

## 5 Related and future work

Among the most related approaches, Web- $\pi$  [6] has no request-response pattern and its treatment of faults is rather different from ours. Orc [5] has a pruning primitive similar to our MSR. However, since Orc has no notion of fault, all the difficulties coming from error management do not emerge. Finally, the service oriented calculi CaSPiS [1] and COWS [7] include low-level mechanisms allowing the programmer to close ongoing conversations. However, our approach is different, since we aim at providing primitives which free the programmer from this burden.

As a future work, we plan to incorporate the primitives we propose in Jolie. Also we would like to study their expressive power: the implementation of MSR in terms of the existing primitives is not easy and we believe that a separation result could be proved.

## References

1. M. Boreale, R. Bruni, R. De Nicola, and M. Loreti. Sessions and pipelines for structured service programming. In *FMOODS'08*, volume 5051 of *LNCS*, pages 19–38. Springer, 2008.
2. C. Guidi, I. Lanese, F. Montesi, and G. Zavattaro. On the interplay between fault handling and request-response service invocations. In *ACSD'08*, pages 190–198. IEEE Press, 2008.
3. C. Guidi, I. Lanese, F. Montesi, and G. Zavattaro. Dynamic error handling in service oriented applications. *Fundamentae Informaticae*, 95(1):73–102, 2009.
4. C. Guidi, R. Lucchi, R. Gorrieri, N. Busi, and G. Zavattaro. SOCK: A calculus for service oriented computing. In *ICSOC'06*, volume 4294 of *LNCS*, pages 327–338. Springer, 2006.
5. D. Kitchin, A. Quark, W. R. Cook, and J. Misra. The Orc programming language. In *FMOODS/FORTE'09*, volume 5522 of *LNCS*, pages 1–25. Springer, 2009.
6. C. Laneve and G. Zavattaro. Foundations of web transactions. In *FoSSaCS'05*, volume 3441 of *LNCS*, pages 282–298. Springer, 2005.
7. A. Lapadula, R. Pugliese, and F. Tiezzi. A calculus for orchestration of web services. In *ESOP'07*, volume 4421 of *LNCS*, pages 33–47. Springer, 2007.
8. OASIS. *Web Services Business Process Execution Language Version 2.0*. <http://docs.oasis-open.org/wsbpel/2.0/wsbpel-v2.0.html>.
9. World Wide Web Consortium. *Web Services Description Language (WSDL) 1.1*. <http://www.w3.org/TR/wSDL>.