

Verifying Generic Erlang Client–Server Implementations

Thomas Arts¹ and Thomas Noll²

¹ Ericsson, Computer Science Lab,

Box 1505, S-12525 Älvsjö, Sweden, thomas@cslab.ericsson.se

² Department of Teleinformatics, Royal Institute of Technology (KTH),
S-16440 Kista, Sweden, noll@it.kth.se

Abstract. The Erlang Verification Tool is an interactive theorem prover tailored to verify properties of distributed systems implemented in Erlang. It is being developed by the Swedish Institute of Computer Science in collaboration with Ericsson.

In this paper we present an extension of this tool which allows to reason about the Erlang code on an architectural level. We present a verification method for client–server systems designed using the generic server implementation of the Open Telecom Platform. For this purpose, we specify a set of transition rules which characterize the abstract behaviour of the generic server functions. By this means we can reason in a partitioned way about any client–server application without having to consider the concrete implementation details of the generic part, which simplifies proofs dramatically.

The generic server architecture is just an example, and the technique extends to many other generic components. Moreover, the idea of considering standard components to reason on the architectural level of a concrete implementation can also be explored when using other verifications tools for Erlang or in the context of another language.

1 Introduction

The high quality demands on software for telecommunication applications may partly be ensured by the use of formal methods in the design and development. By the high degree of concurrency in those applications, *testing* is often not sufficient to guarantee correctness to a satisfactory degree. *Verification*, namely formally proving that a system has the desired properties, is therefore becoming a more and more widespread practice (see [CW96] for an overview).

Although a complete formal specification of an application would probably be one of the best ways to ensure its correctness, in practice the descriptions are rather informal, written in natural language in combination with some fragments of, for example, the Standard Description Language SDL [SDL93]. Reasons for the absence of a complete formal specification can be found in the fact that the specification changes several times during development, triggered by experiments with a release or by changed requirements. It is felt too time consuming to modify

both the formal specification and the code. Even the informal specification tends to run out of phase with the actual implementation and is often only updated after a release of the product.

Towards the end of a project, it is the running code that represents the best ‘specification’ of the implementation. Questions about its correctness are therefore best formulated in terms of this code: ‘is there a possibility that this finite-state machine implementation deadlocks?’, ‘does this server implementation correctly respond to all possible requests?’.

In order to find answers for these questions, one might abstract from the code (having the informal specification helping in this) and check the questions in an obtained model [CGL94,Huc99]. If one takes this realistically, the verification is used for finding errors, more than for proving correctness. If the model does *not* fulfill a certain property, then this might indicate an error in the code. It is common practice to analyze the given trace that leads to the detected error and to check whether this is also a valid trace in the actual code. The latter need not be, since a model (which is an abstraction) neglects some, potentially essential, details.

In general the constructed model depends on the property one wants to prove. Often one does not directly construct the final model, but many models are built, which are refinements of an initial rough model. The model is refined until either a detected error can be identified as a real error, or until one has enough confidence in the detailedness of the model to believe that the code is error-free. In this analysis method, finding a trace to an error can efficiently and automatically be performed by a model checker. The construction of a model and its refinements, including checking the trace in the real code, can often only be done by hand or with some minor computer assistance.

Given that the code is the only available formal description of the software, as an alternative to building the models and checking the traces one can use a theorem prover to reason about the code directly. An interactive proof assistant for the purpose of verifying properties of programs written in the functional programming language Erlang [AVWW96] has been developed by the Swedish Institute of Computer Science (SICS) in collaboration with Ericsson. This Erlang Verification Tool (EVT; [ADFG98]) can be regarded as a tableau-based prover with proof rules for first-order modal logic in Gentzen-style, extended by rules that reflect the semantics of Erlang, rules for decomposing proofs about compound systems to proofs about the components, and rules for induction and co-induction [DFG98].¹ The disadvantage that proofs have to be provided by hand should be put against the advantages of obtaining certainty that a property holds for the code and of the possibility to reason about unbounded data structures, unbounded message queues, and dynamic creation of processes. Moreover, bugs are detected by the fact that proofs cannot be provided, and the attempt to prove the property usually clearly indicates a trace in the code which can be used as a counterexample.

¹ EVT is available at <ftp://ftp.sics.se/pub/fdt/evt/index.html>.

In this paper we present an addition to this verification tool which allows to reason about the Erlang code on an architectural level. In this way we provide a general abstraction, obtained automatically from the code, that is detailed enough to prove many different properties of the program.

Many applications that we consider consist of several servers that communicate with their clients. These servers are all implemented in a predefined, generic way. The generic implementation takes care of starting the server, receiving synchronous and asynchronous messages, and providing debug and log information for maintenance purposes. We added the specification of the functions of this generic server as proof rules to the verification tool. More precise: based on the transition-system semantics of Erlang as presented in [Fre], we provide rules which describe the possible transitions that any Erlang process evaluating the respective function can take, restricted by the shape of the environment if necessary. Thus, we abstract from the actual implementation of the server and concentrate on its specific behaviour instead, such that we can argue about any client-server application without having to consider the source code of the generic part. In this way we support a relativized style of reasoning which is based on the assumption that the concrete implementation of the generic module follows its specification. This abstraction is property-independent and can hence be used for all properties we are interested in, whereas we gain that we may skip many details, leading to much smaller proofs.

The remainder of this paper is organized as follows. Sect. 2 describes the class of systems addressed by our approach, namely, client-server systems implemented in Erlang using the `gen_server` module. In Sect. 3 we give the abstract representation of some of the `gen_server` functions (the complete specification can be found in [AN00]). Their implementation in EVT is discussed in Sect. 4. In Sect. 5 we address the correctness of our approach and conclude with some comparative remarks.

2 Generic Client-Server Implementations

Large software applications are built using a software architecture. Elements of such architectures are: databases, device drivers, finite-state machines, supervisors, monitors, servers, and many more. After putting the architecture together, the actual implementation of the components may start. Software engineering practice has taught that having all servers implemented in some general way is an advantage, both for development and for maintenance. Even better, when parts of the server software are already written and are used as the basis for all specific servers, it serves the correctness of the whole application, since the generic part is well developed and tested. Therefore, the Open Telecom Platform (OTP; [OTP]), the set of libraries and design principles that comes along with Erlang, supports a standard, generic implementation of a server by providing the `gen_server` module. This module implements several interface functions providing synchronous and asynchronous communication, debugging support, error and timeout handling, and other administrative tasks. In order to obtain

the required specific server functionality the programmer provides an instantiation for this generic server. This instantiation consists of a separate module, the so-called *callback module*, which contains (*callback*) functions that are invoked by the generic part of the server. Thanks to this software engineering practice we are able to easily abstract from the actual server implementation in the code.

The typical flow of control in a `gen_server`-based client-server application is as follows. When a client process wants to synchronously communicate with the server, it uses the standard `gen_server:call` function with a certain message as an argument. The generic part sends the message to the server process and blocks the client. In the server process, another function of the generic part receives the message and forwards it to the application-specific part by calling a function in the callback module. This callback function should return the response and the new server state. The new state is stored in the server process, and the reply is returned to the client by the generic part of the server, completing therewith the synchronous event.

In greater detail, the following single steps are taken.

- To start the server process the `gen_server:start` function is called. This function creates a new process, the *server process*, in which a function is started that implements the server. The first thing this function does is computing its initial state by calling the `init` function in the callback module. After that, the process waits for a request from a client process.
- The client uses the `gen_server:call` function to send a synchronous request to the server. The request is handled by the `handle_call` function in the callback module while the client is being suspended, waiting for the response. The current state of the server is passed as an argument to the callback function, which on its turn returns both a reply message and a new state. The reply message is sent to the suspended process.
- Alternatively, `gen_server:cast` can be used to send asynchronous requests to the server. Here only the internal state of the server is changed according to the result of the `handle_cast` function in the callback module.
- Both `handle_call` and `handle_cast` can return a value indicating that the server should terminate. In this case, `gen_server` invokes the `terminate` function in the callback module to clean up before the process terminates.

Clearly the flow of control as described above may look different when error situations occur, such as a server that cannot be started or a call that cannot be handled. In those cases some standard error handling is performed by the generic server. In addition several options can be provided to the standard calls in order to have them behave slightly different.

The following example of a simple locker server implements a scheduler that arbitrates the access to a single resource. It can be used by several clients at a time, communicating synchronously by executing function calls of the form `gen_server:call(Server, request)` and `gen_server:call(Server, release)` to request a lock and release it thereafter, respectively.

The example is classical and the properties of interest are likewise (formulated for a server with arbitrary many clients), such as: no deadlock, no starvation, mutual exclusion.

For details on the syntax and semantics of Erlang see [AVWW96]. It is a concurrent programming language with processes that execute functions. Erlang is an eager, dynamically typed language with only a few data types. In this paper we use atoms (constants) which are denoted by lowercase symbols, tuples, and lists. Variables start with an uppercase character, except for the special variable '_' which matches any value without getting bound to it (i.e., it is always a free variable).

The client function is trivially implemented in a module called `client` by a function with the same name that takes the process identifier of the server as a parameter (to establish communication).

```
-module(client).  
  
client(Server) ->  
  gen_server:call(Server,request),  
  access_the_resource(),  
  gen_server:call(Server,release).
```

The state of the server consists of a list of pending clients. More exactly, the client that currently has access to the resource is stored in the head, and all waiting clients are kept in the tail of this list. We start the server by evaluating `gen_server:start([])` to initialize it with an empty list of pending processes.

```
-module(locker).  
-behaviour(gen_server).  
  
init(Requests) ->  
  {ok, Requests}.  
  
handle_call(request, From, Requests) ->  
  case Requests of  
    [] ->  
      {reply, ok, [From]};  
    _ ->  
      {noreply, Requests++[From]}  
  end;  
handle_call(release, From, [_|Waiting]) ->  
  case Waiting of  
    [] ->  
      {reply, done, Waiting};  
    _ ->  
      gen_server:reply(hd(Waiting), ok),  
      {reply, done, Waiting}  
  end;
```

```
handle_call(stop, From, Requests) ->
  {stop, normal, ok, Requests}.

terminate(Reason, Requests) ->
  ok.
```

The `gen_server:call` function in the client causes the `handle_call` function in the callback module to be executed. Note that the return value (either `ok` or `done`) is ignored by the client process, i.e., no check is performed to see if the value really is the expected one. The return of a value is a synchronization mechanism which in this case is independent of the actual value. In this particular example we could have restricted us to asynchronous communication for releasing by using `gen_server:cast` instead of `gen_server:call`, but for readability we aim to concentrate on only one communication primitive. Also note that the server expects the clients to stick to the locking protocol. In other words, in this simple version we left out any effort to program defensively. For example, misbehaving programmers can crash the locker by sending a `release` without a previous `request` or by sending a message that is not recognized by the locker at all.

On the level of the actual execution, Erlang supports only one way of communicating messages, which is asynchronous. However, the `gen_server` implementation ensures a synchronized behaviour: the `gen_server:start` function will not return before the `init` function has returned a state, and the `gen_server:call` function only returns if the `handle_call` function returns a reply to the callee. In the `gen_server` module this synchronous communication is implemented by using the asynchronous primitives: a message is sent and a `receive` statement directly succeeds this output operation, waiting for the response. The main goal of the technique we present in this paper is to be able to abstract from the implementation of the synchronous communication. The fact that the message is read from the message queue in a certain way, that a timeout primitive is supported and all that, need not be of our concern.

In this abstract setting, the given server implementing the locker can be described as a server that stores a list of clients which claim access to the resource. The number of clients is arbitrary and in properties or proofs we want to make no assumptions about an upper bound. Access is granted to the first client in this list, the other clients are suspended. Only after a release by the client that currently accesses the resource the client that is next in line gets access to it. The suspension of clients is implemented by not providing the reply immediately (returning `{noreply, ...}`), but sending it later (`gen_server:reply(..., ok)`).

3 The Verification Approach

Already without our addition, the Erlang Verification Tool can be used to verify the Erlang code of an application which makes use of the generic server imple-

mentation². When establishing such a proof, one has to follow the simulation of the synchronous communication by the underlying asynchronous implementation. By the nature of the Erlang semantics this means that one should also prove some properties about the message queues of the client(s) and of the server, which seems irrelevant given the knowledge how the server works. In particular, when many clients are involved, a lot of nondeterminism can be introduced in the proof by observationally equivalent traces. As such, the number of proof goals may be much larger than it seems strictly necessary.

Another disadvantage in such a proof is that one gets confronted with details such as debug features implemented in the `gen_server` module. Although the verification is performed in a context where debug facilities are assumed to be disabled, one still generates an extra proof goal for testing the debug flag. This test is not atomic, and since we work in a concurrent setting, even those few steps cause duplication of work, since another action may be chosen in another process at that same time.

In our approach we simplify the verification task by ignoring the concrete implementation of the `gen_server` module. We specify its abstract behaviour by making its syntactic constructs recognizable as keywords by EVT, and by adding appropriate transition rules to the proof system. Since these transition rules are of a general nature, they can also be used for implementing our approach in other tools supporting Erlang. The actual implementation of the rules in EVT is described in Sect. 4.

3.1 Extending the Erlang Verification Tool

An Erlang system is specified by a composition of *processes*, each represented as $\langle e, pid, q \rangle$, where e is the Erlang *expression* being evaluated, pid is the uniquely defined *process identifier*, and q is the *mailbox queue* in which incoming messages are stored. In order to have the tool recognize a generic server function call, we add these as special syntactic constructs to the set of expressions:

$$e ::= \dots \mid \text{gen_server:start}(e_1, e_2, e_3) \mid \text{gen_server:call}(e_1, e_2) \mid \\ \text{gen_server:reply}(e_1, e_2) \mid \text{gen_server:wait}(e) \mid \\ \text{gen_server:ready}(e) \mid \text{gen_server:busy}(e_1, e_2) \mid \\ \text{gen_server:down}(e_1, e_2, e_3) \mid \dots$$

In this way, those function calls can be treated in a different way than the other function calls. The standard method is to search for a definition of the function, to substitute the arguments, and to continue with evaluating the body of the definition. The way the special function calls are treated is defined by an extension of the operational semantics which is defined by labeled transition rules (in the style of [Fre]).

² Although the EVT tool lacks support for modules at the moment, one can combine the callback, the `gen_server` and, if present, some client module into one bigger module by little effort.

A *reduction context* is an Erlang expression $r[\cdot]$ with a ‘hole’ \cdot in it, which identifies the position of r where the next evaluation step takes place. In this way, the rules for the actual expression evaluation have to be given only for exceptional cases, namely, when all parameters of an expression construct are *values*, i.e., have been fully evaluated.

For example, process creation is formally described by the following rule³:

$$\frac{\langle r[\mathbf{spawn}(f, [v_1, \dots, v_n])], pid, q \rangle}{\xrightarrow{\tau} \langle r[pid'], pid, q \rangle \parallel \langle f(v_1, \dots, v_n), pid', \varepsilon \rangle} \quad (1)$$

Here, f is a (function) atom, v_1, \dots, v_n are values, q is an arbitrary mailbox and ε denotes the empty mailbox. Thus, a process evaluating a `spawn` function call has a transition to a system of two processes (\parallel denotes parallel composition) which have to evaluate the expressions $r[pid']$ (pid' is the return value of `spawn`) and $f(v_1, \dots, v_n)$, respectively. For the process identifiers pid and pid' we require $pid' \neq pid$.

Assuming now that the set of reduction contexts has been extended accordingly to cope with the new syntactic constructs, we can formalize the intuitive meaning of the `gen_server` functions as described in the previous section as follows. Here we describe the starting of a server process and the handling of a server call; the complete specification is given in [AN00].

Starting a server is similar to spawning a process, but the continuation of the process depends on the evaluation of the `init` callback function. This is the reason for adding the special `gen_server:wait` construct:

$$\frac{\langle r[\mathbf{gen_server:start}(mod, arg, opt)], pid, q \rangle}{\xrightarrow{\tau} \langle r[\mathbf{gen_server:wait}(spid)], pid, q \rangle \parallel \langle \mathbf{init}(arg), spid, \varepsilon \rangle} \quad (2)$$

Here, $spid$ denotes a fresh (server) pid, and the server is created with an empty mailbox. The term `gen_server:wait(spid)` should not be treated as a normal form, since then reductions in the context $r[\cdot]$ would be allowed, but rather as a construct from which currently no transitions are possible (similar to a `receive` statement with an empty mailbox).

According to the generic server description, the result of evaluating the `init` function should be a tuple with the initial server state as its second component. This state should be kept as part of the looping server (looping over: receiving a request, computing the answer and the next state, and responding).

$$\frac{\langle r[\mathbf{gen_server:wait}(spid)], pid, q \rangle \parallel \langle \{\mathbf{ok}, state\}, spid, sq \rangle}{\xrightarrow{\tau} \langle r[\{\mathbf{ok}, spid\}], pid, q \rangle \parallel \langle \mathbf{gen_server:ready}(state), spid, sq \rangle} \quad (3)$$

Note that the identifier of the process that started the server is not known to the server in this specification. Since pids of newly created processes are unique,

³ Actually, in the definition of the semantics a two-layer scheme is employed which separates expression-level from process-level steps. We will consider this distinction in Sect. 4 in greater detail.

this causes no problem in our setting. The starting process ‘remembers’ which server it has started (by the obtained process identifier $spid$).

A call by a client can be handled by the server if it is in an idle state, denoted by the `gen_server:ready` construct. In this case, the server process invokes the `handle_call` callback function, and the client process is put into a waiting state until the request has been answered. Now, however, the server needs to store the pid of the calling client in order to be able to distinguish clients if several of them are waiting for the same server:

$$\begin{aligned} & \langle r[\text{gen_server:call}(spid, req)], pid, q \rangle \parallel \\ & \langle \text{gen_server:ready}(state), spid, sq \rangle \\ \xrightarrow{\tau} & \langle r[\text{gen_server:wait}(spid)], pid, q \rangle \parallel \\ & \langle \text{gen_server:busy}(\text{handle_call}(req, pid, state), pid), spid, sq \rangle \end{aligned} \quad (4)$$

If the `handle_call` function yields a triple of the form $\{\text{reply}, \text{answer}, \text{state}\}$, then this answer is immediately returned to the waiting client, and the server changes into the idle state again:

$$\begin{aligned} & \langle r[\text{gen_server:wait}(spid)], pid, q \rangle \parallel \\ & \langle \text{gen_server:busy}(\{\text{reply}, \text{answer}, \text{state}\}, pid), spid, sq \rangle \\ \xrightarrow{\tau} & \langle r[\text{answer}], pid, q \rangle \parallel \\ & \langle \text{gen_server:ready}(state), spid, sq \rangle \end{aligned} \quad (5)$$

The fact that the process identifier is stored in the second argument in the server expression guarantees that the reply is received by the right waiting client.

As can be seen in our locker example, the `handle_call` may also return a tuple of the form $\{\text{noreply}, \text{state}\}$. In this case the client process remains in the waiting state (and does not have to be considered therefore), whereas the server becomes idle again:

$$\begin{aligned} & \langle \text{gen_server:busy}(\{\text{noreply}, \text{state}\}, pid), spid, sq \rangle \\ \xrightarrow{\tau} & \langle \text{gen_server:ready}(state), spid, sq \rangle \end{aligned} \quad (6)$$

Note that when we have two clients, one suspended and one calling the server, then both clients are in the waiting state, but only one client can be activated by the return of the `handle_call`, viz. the process which called. The other process must be activated by explicitly using the `gen_server:reply` function:

$$\begin{aligned} & \langle r[\text{gen_server:wait}(spid)], pid, q \rangle \parallel \\ & \langle r'[\text{gen_server:reply}(pid, \text{answer})], spid, sq \rangle \\ \xrightarrow{\tau} & \langle r[\text{answer}], pid, q \rangle \parallel \langle r'[\text{true}], spid, sq \rangle \end{aligned} \quad (7)$$

In this way semantical rules can be used to accurately describe the given example of the locker server. As can be seen, the asynchronous communication actions that are used in the `gen_server` module to implement synchronous message passing are abstractly represented by simple handshaking operations which do not consider the message queues of the client nor of the server.

Asynchronous communication, not on the level of Erlang, but on the level of `gen_server` is also supported. This is implemented via the `gen_server:cast` and `handle_cast` functions. The `gen_server:cast` mechanism is formalized similar to the `gen_server:call` mechanism. The only difference is that the client immediately proceeds without waiting for a server response. Having evaluated the `handle_cast` function (which is indicated by a `noreply` result tuple), the server just changes into the `gen_server:ready` state, modifying its local data according to the result.

Apart from the `reply` and `noreply` values, `handle_call` (and `handle_cast`) can instead return a result that indicates that the server has to terminate. If so, the `terminate` function in the callback module is invoked. In this situation, the response is stored on the server side until `terminate` has finished:

$$\begin{aligned} & \langle \text{gen_server:busy}(\{\text{stop}, \text{reason}, \text{answer}, \text{state}\}, \text{pid}), \text{spid}, \text{sq} \rangle \\ \xrightarrow{\tau} & \langle \text{gen_server:down}(\text{terminate}(\text{reason}, \text{state}), \text{pid}, \text{answer}), \text{spid}, \text{sq} \rangle \end{aligned} \quad (8)$$

The `terminate` function is supposed to return the value `ok` and after that, the client is released and the server process is removed:

$$\begin{aligned} & \langle r[\text{gen_server:wait}(\text{spid})], \text{pid}, q \rangle \parallel \\ & \langle \text{gen_server:down}(\text{ok}, \text{pid}, \text{answer}), \text{spid}, \text{sq} \rangle \\ \xrightarrow{\tau} & \langle r[\text{answer}], \text{pid}, q \rangle \end{aligned} \quad (9)$$

Note that the callback functions such as `init` and `handle_call` are specified in the callback module. We use the rules provided by the system to reason about their behaviour. Thus, abstraction by means of the semantical rules is only provided for the generic part of the server.

4 Implementation

As mentioned in the introduction, our `gen_server` verification approach has been implemented using the Erlang Verification Tool. Some minor additions had to be made to the tool itself, basically the recognition of the special constructs whose handling is left to the user. Thus, whenever a special construct occurs, the tool is aware of the fact that this is not a normal form, but that transitions may arise from the respective term. It leaves it to the user to check which transitions are possible. We, as a user, provided tactics to analyze the term and to apply the corresponding transition rule. These tactics are combinations of proof rules that are applied to the proof goal. Hence, we specified the transition rules given in Sect. 3 as logical formulae.

Specifying the transition rules as logical formulae rather than integrating them as an extension of the original Erlang semantics into the EVT source code has two advantages: first of all, the reasoning within the tool remains sound with respect to the abstract `gen_server` semantics. By using an abstract model of a server one introduces a potential unsoundness. If the property that one wants to prove depends on the actual implementation of the server, one might

be able to falsely prove it for the abstraction, whereas it does not hold for the real program (this point is discussed in the conclusions). By using logical rules for the transitions, one explicitly states in the assumptions how one expects the server to behave. Since the reasoning that involves these assumptions is based on the (sound) EVT proof system, soundness is guaranteed under the premise that the (low-level) implementation of the `gen_server` module behaves as described by the (high-level) specification.

Second, one obtains a greater flexibility for experimenting. It is easier to change a logical expression and a tactic than to modify the implementation of EVT itself. Moreover, several different specifications of the server may exist at the same time, all using the same tool.

In the following we give the logical representation of the transition rules which describe the starting of a server process. As mentioned earlier, the formal semantics of Erlang is given by a two-layer scheme [Fre], which is also used in the EVT implementation. First the Erlang expressions are provided with a semantics on the expression level. The actions here are a functional computation step, an output, a receiving of a message, and a call of a builtin function (like `spawn` for process creation) with side effects on the process-level state. Second, the transition behaviour of Erlang systems (that is, concurrent processes evaluating expressions in the context of a unique process identifier and a mailbox of incoming messages) is captured through a set of transition rules which lift the expression actions to the process level, and which describe the interleaving of concurrent actions. Here, possible process actions are computation steps, input, and output actions.

In this setting, (1) is decomposed into an expression-level rule which indicates the spawning action, and a process-level rule which models the actual process creation:

$$\frac{\text{spawn}(f, [v_1, \dots, v_n]) \xrightarrow{\text{spawn}(f, [v_1, \dots, v_n]) \rightarrow pid'} pid' \quad e \xrightarrow{\text{spawn}(f, [v_1, \dots, v_n]) \rightarrow pid'} e' \quad pid' \neq pid}{\langle e, pid, q \rangle \xrightarrow{\tau} \langle e', pid, q \rangle \parallel \langle f(v_1, \dots, v_n), pid', \varepsilon \rangle}$$

These transition rules in the tool automatically generate possible next states of an Erlang system when we want to prove that something holds in some successor state or in all successor states (diamond or box modality, respectively). Thus, given a `spawn` call in the program, we can reason about the state in which a new process is created, among the other possible next states that the tool computes for us. For the generic server behaviour we want to obtain a similar level of comfort. However, now we specify the transitions as logical formulae. For example, (2) gives rise to the following logical formula which describes the

starting of a server.

$$\begin{aligned} & \forall mod : \text{Atom}. \forall arg : \text{Value}. \forall opt : \text{List}. \\ & \left(\forall spid : \text{Pid}. \right. \\ & \quad \left(\text{gen_server:start}(mod, arg, opt) \xrightarrow{\text{spawn}(\text{init}, [arg]) \rightarrow spid} \text{gen_server:wait}(spid) \right) \\ & \quad \wedge \\ & \quad \left(\forall a : \text{Action}. \forall e : \text{Expr}. \text{gen_server:start}(mod, arg, opt) \xrightarrow{a} e \Rightarrow \right. \\ & \quad \left. \exists spid : \text{Pid}. a = \text{spawn}(\text{init}, [arg]) \rightarrow spid \wedge e = \text{gen_server:wait}(spid) \right) \\ & \left. \right) \end{aligned}$$

By associating a `spawn` action with the `gen_server:start` function call, we employ the process-level rule for `spawn` as given above for computing the next state. The server starts evaluating the `init` function⁴, which may involve some standard reasoning with the tool.

However, since we now specify the transition as a logical rule, we also need to state that no other transition is possible from this point, as expressed by the second subformula of the conjunction. The implementation of our generic server primitives in the tool is such that it presents all known actions at every point where one of these primitives may enable a transition. We manually have to prove that a certain transition is possible and that the others are not. The information needed for this proof is provided in the assumptions by means of a logical formula as presented above.

Since we try to achieve a high degree of automation in our proofs, we provided tactics for automatically showing that a certain action can take place and that the others are not enabled. For example, when proving a diamond property for an expression where a server is started, the tactics scan in the assumptions of the goal for a property named `start_dia`, and this property is used to automatically prove that after the `gen_server:start` call a new process is created evaluating the `init` function.

Rule (3) expresses that, after the successful initialization, the server pid is returned to the spawning process. Since it makes a provision on the syntactic structure of the two expressions computed by the processes, we have to give four formulae in the EVT specification. The first two specify the local effects of sending/receiving a message by/from the server, respectively.

$$\begin{aligned} & \forall state : \text{Value}. \\ & \left(\forall spid : \text{Pid}. \forall pid : \text{Pid}. \right. \\ & \quad \left(\{ok, state\} \xrightarrow{\text{server}(spid, pid) \rightarrow \{ok, spid\}} \text{gen_server:ready}(state) \right) \wedge \\ & \quad \left(\forall a : \text{Action}. \forall e : \text{Expr}. \{ok, state\} \xrightarrow{a} e \Rightarrow \exists spid : \text{Pid}. \exists pid : \text{Pid}. \right. \\ & \quad \left. a = \text{server}(spid, pid) \rightarrow \{ok, spid\} \wedge e = \text{gen_server:ready}(state) \right) \\ & \quad \forall spid : \text{Pid}. \\ & \quad \left(\forall v : \text{Value}. \text{gen_server:wait}(spid) \xrightarrow{\text{client}(spid) \rightarrow v} v \right) \wedge \\ & \quad \left(\forall a : \text{Action}. \forall e : \text{Expr}. \text{gen_server:wait}(spid) \xrightarrow{a} e \Rightarrow \right. \\ & \quad \left. \exists v : \text{Value}. a = \text{client}(spid) \rightarrow v \wedge e = v \right) \\ & \left. \right) \end{aligned}$$

⁴ To be precise the function `mod:init` is called, but the tool, at the moment, lacks support for function calls in remote modules, such that we have to localize all calls.

Furthermore the expression-level synchronization actions have to be lifted to the process level. We only give the formula for the server side; the dual one (defining corresponding input actions of the form $spid?sync(pid, v)$) is of a similar shape.

$$\forall e : \text{Expr. } \forall e' : \text{Expr. } \forall spid : \text{Pid. } \forall pid : \text{Pid. } \forall v : \text{Value. } \forall sq : \text{Queue.}$$

$$\left(e \xrightarrow{\text{server}(spid, pid) \rightarrow v} e' \Rightarrow \langle e, spid, sq \rangle \xrightarrow{spid!sync(pid, v)} \langle e', spid, sq \rangle \right) \wedge$$

$$\left(\langle e, spid, sq \rangle \xrightarrow{spid!sync(pid, v)} \langle e', spid, sq \rangle \Rightarrow e \xrightarrow{\text{server}(spid, pid) \rightarrow v} e' \right)$$

Now the standard synchronization mechanism implemented in EVT is used to model the actual communication between the two `sync` events. Note that both the client and the server pid is used to match the synchronization actions.

Similar formulae are obtained for (4) to (9) describing the synchronous `gen_server:call` mechanism and the termination of a server process, and for the remaining `gen_server` constructs. All single properties are collected in a big conjunction named `gen_server`, which completely specifies the abstract behaviour of `gen_server` systems, and which can be found in [AN00].

By implementing the transition rules as logical formulae we gain a flexible and easily adaptable extension of the tool. However, by the fact that we also need the ‘negative’ information that a certain transition is the only possible, we have some overhead in disproving extra generated subgoals for all other types of transitions. These subgoals are, however, relatively easy to reject and we implemented tactics to deal with them automatically.

For example it is possible to establish the mutual exclusion property of the locker protocol as defined in Sect. 2. Let us assume that the `access_the_resource` function, which implements the client’s activity in the critical section, is given by

```
access_the_resource() ->
  self()!access.
```

Hence the fact that within a given process system `S` a client with pid `Client` has entered its critical section can be expressed by the formula

```
in_cs: erlangPid -> erlang_system -> prop =
  \Client: erlangPid.
  \S: erlang_system.
  (S: <Client!message(access)>tt).
```

In the simplest case of a system with two clients using one locker, the mutual exclusion property can be characterized by the following safety formula, asserting for every reachable state of the system that not both clients are in their critical sections at the same time:

```

mutex: erlangPid -> erlangPid -> erlang_system -> prop =>
  \Client1: erlangPid.
  \Client2: erlangPid.
  \S: erlang_system.
  ( not (S: (in_cs Client1) /\ (in_cs Client2))
    /\
    (S: [tau](mutex Client1 Client2))
  ).

```

We have experimented with our implementation of the generic server specifications applying it to some small examples like the above. Indeed the proofs are easier and shorter than without using this extension. On larger examples this benefit can only become more significant.

5 Conclusions and Future Work

In this article we presented an addition to the Erlang Verification Tool that enables us to reason on a higher level about the code implementing a client-server architecture using the generic server paradigm. The addition drastically reduces the amount of details one needs to consider when proving a property of an Erlang application that uses this client-server architecture, therefore resulting in shorter proofs.

We formalized the operational semantics of the generic server behaviour similar to the way in which the operational semantics of more primitive Erlang functions is specified. For implementing this formalization we defined it in terms of logical formulae such that the only change in the tool that we required was the extension of the list of recognizable keywords by the `gen_server`-specific function names. It turned out that, in comparison to proofs based on the concrete implementation, the logical formulae support a more effective reasoning about client-server systems. The efficiency of this reasoning has been increased by adding tactics that automatically prove subgoals about impossible alternative transitions.

By specifying only the ‘essential’ behaviour of any reasonable `gen_server` implementation, we introduce a certain unsoundness in the proofs. That is, if employing our approach we succeed to prove a certain property of an Erlang program that contains servers implemented by the `gen_server` module, then it depends on the property whether it really holds for the actual code. This, however, is a consequence of abstraction and not at all a problem in practice. First of all we are not so much interested in proving correctness but rather in finding errors in the program. If we find an error with respect to this abstraction, it is most likely an error in the real code as well. Second, if the property is independent of the specification of the generic server, then the property should hold for the actual code. Since we abstract several server steps into one handshake operation in our semantics, the property should be at least τ -insensitive, i.e., its validity should not depend on the number of internal actions the system evaluates. The τ -insensitivity is a minimal requirement of the property, but insufficient,

since there are several other issues involved as well that make it very hard to formalize the exact independence criteria. For example, the property should be independent of: the number of messages in the mailbox of a server, the priority used to read messages from a mailbox, the debug and fault tolerance additions not specified by us, etc.

Pragmatically, our concern is to provide a framework in which we can prove properties of the code in an abstract setting, where we use one abstraction for all possible properties. This abstraction is very close to the real implementation, but there will always exist properties for which it turns out to be too general. However, if we can prove a certain property about the abstraction, then we increased the level of confidence in the code; if we find that a certain property does not hold by reasoning in this abstracted setting, then, most likely, this corresponds to an error in the real program. For that part, our technique is therefore rather close to the model-checking approach. Here, however, we only need one abstraction for arbitrary properties and we do not have to abstract over unbounded data structures, dynamic process spawning, or dynamic network creation. Our approach obtains the abstraction automatically, but needs human assistance in non-trivial proofs, whereas the latter can often automatically be handled when using a model checker.

In order to compare efforts, we experiment with using the same formalization of the behaviour of the generic server and its callback module with model-checking tools (such as TRUTH/SLC [LLNT99]). Since one possible source of infinite state spaces, the unbounded message queue of an Erlang process, is abstracted away in our model, this approach should potentially be more successful than for arbitrary Erlang programs. However, we can still apply those tools only to examples where the state space is finite; in particular, the number of processes must be bounded.

The Open Telecom Platform contains several additional generic architectures, such as generic finite-state machines, generic event handlers, generic supervision trees, etc. Those concepts can be formalized and added to EVT along the same lines. A major part of the tactics that we have already written will directly be usable for those other generic concepts. In this way, with only little extra effort, the verification of even more realistic large applications can be simplified.

Acknowledgement

We thank Lars-åke Fredlund for his support in making the necessary changes in the tool and explaining and helping us with the implementation of the tactics.

References

- [ADFG98] T. Arts, M. Dam, L.-å. Fredlund, and D. Gurov. System description: Verification of distributed Erlang programs. In *Proc. CADE'98*, volume 1421, pages 38–41. Springer-Verlag, 1998.

- [AN00] T. Arts and T. Noll. Verifying generic Erlang client-server implementations. Technical Report 00-08, Aachen University of Technology, Aachen, Germany, 2000. <ftp://ftp.informatik.rwth-aachen.de/pub/reports/2000/00-08.ps.gz>.
- [AVWW96] J.L. Armstrong, S.R. Virding, M.C. Williams, and C. Wikström. *Concurrent Programming in Erlang*. Prentice Hall International, 2nd edition, 1996.
- [CGL94] E.M. Clarke, O. Grumberg, and D.E. Long. Model checking and abstraction. *ACM Transactions on Programming Languages and Systems*, 16(5):1512–1542, 1994.
- [CW96] E.M. Clarke and J.M. Wing. Formal methods: State of the art and future directions. Technical Report CMU-CS-96-178, Carnegie Mellon University, Pittsburg, USA, 1996.
- [DFG98] M. Dam, L.-å. Fredlund, and D. Gurov. Toward parametric verification of open distributed systems. In *Compositionality: the Significant Difference*, volume 1536 of *Lecture Notes in Computer Science*, pages 150–185. Springer-Verlag, 1998.
- [Fre] L.-å. Fredlund. Towards a semantics for Erlang. Unpublished manuscript. Swedish Institute of Computer Science.
- [Huc99] F. Huch. Verification of Erlang programs using abstract interpretation and model checking. *ACM SIGPLAN Notices*, 34(9):261–272, 1999. Proceedings of the ACM SIGPLAN International Conference on Functional Programming (ICFP '99).
- [LLNT99] M. Lange, M. Leucker, T. Noll, and S. Tobies. Truth – a verification platform for concurrent systems. In *Tool Support for System Specification, Development, and Verification*, Advances in Computing Science. Springer-Verlag Wien New York, 1999.
- [OTP] Open Telecom Platform (OTP). Ericsson Utvecklings AB, <http://www.erlang.org/documentation/doc/index.html>.
- [SDL93] CCITT Specification and Description Language (SDL). Technical Report 03/93, International Telecommunication Union, 1993. <http://www.itu.int/>.