

Fast, Optimized Sun RPC Using Automatic Program Specialization*

Gilles Muller, Renaud Marlet,
Eugen-Nicolae Volanschi and Charles Consel
IRISA, Campus Universitaire de Beaulieu,
35042 Rennes Cedex - France
{muller,marlet,volanski,consel}@irisa.fr
<http://www.irisa.fr/compose>

Calton Pu and Ashvin Goel
Dept. of Computer Science and Engineering,
Oregon Graduate Institute,
P.O. BOX 91000, Portland OR 97291-1000, USA
{calton,ashvin}@cse.ogi.edu

Abstract

Fast remote procedure call (RPC) is a major concern for distributed systems. Many studies aimed at efficient RPC consist of either new implementations of the RPC paradigm or manual optimization of critical sections of the code. This paper presents an experiment that achieves automatic optimization of an existing, commercial RPC implementation, namely the Sun RPC. The optimized Sun RPC is obtained by using an automatic program specializer. It runs up to 1.5 times faster than the original Sun RPC. Close examination of the specialized code does not reveal further optimization opportunities which would lead to significant improvements without major manual restructuring.

The contributions of this work are: (1) the optimized code is safely produced by an automatic tool and thus does not entail any additional maintenance; (2) to the best of our knowledge this is the first successful specialization of mature, commercial, representative system code; and (3) the optimized Sun RPC runs significantly faster than the original code.

Keyword: partial evaluation, RPC protocol, Sun RPC, distributed systems, automatic optimization.

1 Introduction

Specialization is a well-known technique for improving the performance of operating systems [3, 11, 20]. However, only recently have programming tools begun to be used to help system programmers perform specialization. To the best of our knowledge, this paper reports the first successful specialization of a significant OS component (the Sun RPC) using a partial evaluator. This work is significant for a combination of three main reasons: (1) automatic optimization of existing system code using a partial evaluator preserves the original source code (2) specialization applies to mature, commercial, representative system code, and (3) the specialized Sun RPC exhibits significant performance gains. We elaborate each reason in turn.

First, partial-evaluation-based specialization is qualitatively different from manual specialization done in the past

[3, 11, 20]. Manual specialization requires the system programmer to identify every occurrence of the invariants to be exploited and to write the specialized code exploiting these invariants. Although this approach may lead to significant performance gains, the manual specialization process is error-prone and results in code that is expensive to maintain. In contrast, a partial evaluator preserves the source code, and generates automatically the specialized code guided by the declarations of invariants specified by system programmers. Since we are specializing mature commercial code, the preservation of original code and semantics also preserves safety and maintainability. In our view, tools such as partial evaluators may help the industry to address the operating system code complexity concerns.

Second, we specialize mature, commercial code (Sun RPC) that we believe to be representative of production quality code. Sun RPC is one layer in the communication stack, and RPC itself is divided into micro-layers, each concerned with a reasonably small task, e.g., managing the underlying transport protocol such as TCP or UDP. The RPC code has been ported to a variety of software and hardware foundations, while preserving its layered structure.

Third, we obtain significant performance gains using partial-evaluation-based specialization. In our experiment, the optimized Sun RPC runs up to 1.5 times faster than the original Sun RPC. In addition, the specialized marshaling process runs up to 3.75 times faster than the original one. Close examination of the specialized code does not reveal further optimization opportunities which would lead to significant improvements without major manual restructuring.

Our partial-evaluation-based specialization experiment shows the promise of direct industrial relevance to commercial systems code.

The rest of the paper is organized as follows. Section 2 presents Sun RPC protocol and optimization issues. Section 3 examines opportunities for specialization in the Sun RPC. Section 4 gives an overview of the partial evaluator Tempo and shows its relevance for Sun RPC specialization. Section 5 describes the performance experiments. Section 6 discusses our experience with partial-evaluation based specialization of Sun RPC. Section 7 summarizes related work and Section 8 concludes the paper.

*This work has been partly supported by FRANCE TELECOM CTI-CNET 951W009.

2 The Sun RPC and Optimization Issues

The Sun RPC (Remote Procedure Call) protocol was introduced in 1984 to support the implementation of distributed services. This protocol has become a *de facto* standard in distributed service design and implementation, e.g., NFS [18] and NIS [22]. Since large networks are often heterogeneous, support for communicating machine-independent data involves encoding and decoding. Such environments (e.g., PVM [13] for a message passing model and Stardust [4] for a Distributed Shared Memory model) often rely on Sun XDR. The two main functionalities of the Sun RPC are:

1. A stub generator (`rpcgen`) that produces the client and server stub functions. The stub functions translate procedure call parameters into a machine-independent message format called XDR, and XDR messages back into procedure parameters. The translation of parameters into messages is known as *marshaling*.
2. The management of message exchange through the network.

In concrete terms, the Sun RPC code consists of a set of micro-layers, each one devoted to a small task. For example, there are micro-layers to write data during marshaling, to read data during unmarshaling and to manage specific transport protocols such as TCP or UDP. Each micro-layer has a generic function, but it may have several implementations. As such, the micro-layer organization of RPC code is fairly representative of modular production system software.

A Simple Example

We consider a simple example to illustrate the micro-layer organization of Sun RPC code: a function `rmin` that sends two integers to a remote server, which returns their minimum.

The client uses `rpcgen` (the RPC stub compiler) to compile a procedure interface specification for `rmin` into an assortment of source files. These files implement both the call on the client's side and the dispatch of procedures on the server's side. To emphasize the actual code executed, instead of including all the files generated by `rpcgen`, Figure 1 shows an abstract execution trace of a call to `rmin`.¹

Performance of RPC

Communication using the RPC paradigm is at the root of many distributed systems. As such, the performance of this component is critical. As a result, a lot of research has been carried out on the optimization of this paradigm [5, 8, 14, 19, 25, 29]. Many studies have been

¹For clarity, we omit some clutter in code listings: declarations, "uninteresting" arguments and statements, error handling, casts, and a level of function call.

carried out, but they often result in using new protocols that are incompatible with an existing standard such as the Sun RPC. The problem in reimplementing a protocol that is specified only by its implementation is that features (and even bugs) may be lost, resulting in incompatible implementation.

Optimizing the Existing Code

An alternative to reimplementing a system component for performance reasons is to directly derive an optimized version from the existing code. An advantage of starting with existing code is that the derived version remains compatible with existing standards. Another advantage is that the systematic derivation process can be repeated for different machines and systems.

The question that naturally arises at this point is: are there important opportunities for deriving significantly optimized versions of existing system components?

In fact, many existing system components are known to be generic and structured in layers and modules. This results in various forms of interpretation which are important sources of overhead as shown, for example, in the HP-UX file systems [20]. In case of the Sun RPC, this genericity takes the form of several layers of functions which interpret descriptors (i.e., data structures) to determine the parameters of the communication process: choice of protocol (TCP or UDP), whether to encode or decode, buffer management, ...

Importantly, most of these parameters are known for any given remote procedure call. This information can be exploited to generate specialized code where these interpretations are eliminated. The resulting code is tailored for specific situations.

Let us now examine forms of these interpretations in the Sun RPC code and how they can be optimized via specialization.

3 Opportunities for Specialization in the Sun RPC

The Sun RPC relies on various data structures (such as `CLIENT` or `XDR`). Some fields of those data structures have values that can be available before execution actually takes place; they do not depend on the run-time arguments of the RPC. The values of those fields are either repeatedly interpreted or propagated throughout the layers of the encoding/decoding process. Because these values can be available before execution, they may be the source of optimizations: the computations depending only on known (a.k.a *static*) values can be performed during a specialization phase. The specialized program only consists of the computations depending on the unknown (a.k.a *dynamic*) values.

We now describe typical specialization opportunities in the Sun RPC. We illustrate these opportunities with actual code excerpts, annotated to show static and dynamic computations. In the following figures, dynamic computations

```

arg.int1 = ... // Set first argument
arg.int2 = ... // Set second argument
rmin(&arg) // RPC User interface generated by rpcgen
  clnt_call(argsp) // Generic procedure call (macro)
    clntudp_call(argsp) // UDP generic procedure call
      // Write procedure identifier
      XDR_PUTLONG(&proc) // Generic marshaling to memory, stream... (macro)
        xdrmem_putlong(lp) // Write in output buffer and check overflow
          htonl(*lp) // Choice between big and little endian (macro)
      xdr_pair(argsp) // Stub function generated by rpcgen
        // Write first argument
        xdr_int(&argsp->int1) // Machine dependent switch on integer size
          xdr_long(intp) // Generic encoding or decoding
            XDR_PUTLONG(lp) // Generic marshaling to memory, stream... (macro)
              xdrmem_putlong(lp) // Write in output buffer and check overflow
                htonl(*lp) // Choice between big and little endian (macro)
            // Write second argument
            xdr_int(&argsp->int2) // Machine dependent switch on integer size
              xdr_long(intp) // Generic encoding or decoding
                XDR_PUTLONG(lp) // Generic marshaling to memory, stream... (macro)
                  xdrmem_putlong(lp) // Write in output buffer and check overflow
                    htonl(*lp) // Choice between big and little endian (macro)

```

Figure 1: Abstract trace of the encoding part of a remote call to `rmin`

correspond to code fragments printed in bold face; static computations are printed in Roman.

3.1 Eliminating Encoding/Decoding Dispatch

We examine an opportunity for specialization that illustrates a first form of interpretation. The function `xdr_long` (see Figure 2) is capable of both encoding and decoding long integers. It selects the appropriate operation to perform based on the field `x_op` of its argument `xdrs`. This form of interpretation is used in other similar functions for other data types.

In practice, the field `x_op` is known from the execution context (i.e., encoding or decoding process). This information, contained in the `xdrs` structure, can be propagated interprocedurally down to the function `xdr_long`. As a result, the dispatch on `xdrs->x_op` always yields a known result and can totally be eliminated; the specialized version of this function is reduced to only one of the branch, i.e., a return constructs. In this case, the specialized `xdr_long()`, being small enough, disappears after inlining.

3.2 Eliminating Buffer Overflow Checking

Another form of interpretation appears when buffers are checked for overflow. This situation applies to function `xdrmem_putlong` displayed in Figure 3. More specifically, as parameter marshaling proceeds, the remaining space in the buffer is maintained in the field `x_handy`. Similar to the first example, `xdrs->x_handy` is first initialized (i.e., given a static value), and then decremented by static values and tested several times (for each call to `xdrmem_putlong` and related functions). Since the entire process involves static values, the whole buffer overflow checking can be performed during a specialization

phase, before actually running the program. Only the buffer copy remains in the specialized version (unless a buffer overflow is discovered at specialization time).

This second example is important not only because of the immediate performance gain, but also because in contrast with a manual, unwarranted deletion of the buffer overflow checking, the elimination described here is strictly and systematically derived from the original program.

3.3 Propagating Exit Status

The third example relies on the previous examples. The return value of the procedure `xdr_pair` (shown in Figure 4) depends on the return value of `xdr_int`, which in turn depends on the return value of `xdr_putlong`. We have seen that `xdr_int` and `xdr_putlong` have a static return value. Thus the return value of `xdr_pair` is static as well. If we specialize the caller of `xdr_pair` (i.e., `clntudp_call`) as well to this return value, `xdr_pair` no longer needs to return a value: the type of the function can be turned in `void`. The specialized procedure, with the specialized calls to `xdr_int` and `xdr_putlong` inlined, is shown in Figure 5. The actual result value, which is always `TRUE` independently of dynamic `objp` argument (writing the two integers never overflows the buffer), is used to reduce an extra test in `clntudp_call` (not shown).

3.4 Assessment

The purpose of encoding is to copy a network-independent representation of the arguments into an output buffer. The minimal code that we can expect using the approach of a separate output buffer is basically what is shown in Figure 5. The same situation applies for de-

```

bool_t xdr_long(xdrs,lp)           // Encode or decode a long integer
XDR *xdrs;                         // XDR operation handle
long *lp;                           // pointer to data to be read or written
{
  if( xdrs->x_op == XDR_ENCODE )    // If in encoding mode
    return XDR_PUTLONG(xdrs,lp); // Write a long int into buffer
  if( xdrs->x_op == XDR_DECODE )    // If in decoding mode
    return XDR_GETLONG(xdrs,lp); // Read a long int from buffer
  if( xdrs->x_op == XDR_FREE )      // If in "free memory" mode
    return TRUE;                   // Nothing to be done for long int
  return FALSE;                    // Return failure if nothing matched
}

```

Figure 2: Reading or writing of a long integer: xdr_long()

```

bool_t xdrmem_putlong(xdrs,lp)     // Copy long int into output buffer
XDR *xdrs;                         // XDR operation handle
long *lp;                           // pointer to data to be written
{
  if((xdrs->x_handy -= sizeof(long)) < 0) // Decrement space left in buffer
    return FALSE;                   // Return failure on overflow
  *(xdrs->x_private) = htonl(*lp);   // Copy to buffer
  xdrs->x_private += sizeof(long);   // Point to next copy location in buffer
  return TRUE;                      // Return success
}

```

Figure 3: Writing a long integer: xdrmem_putlong()

```

bool_t xdr_pair(xdrs, objp) {      // Encode arguments of rmin
  if (!xdr_int(xdrs, &objp->int1)) // Encode first argument
    return (FALSE);               // Possibly propagate failure
  if (!xdr_int(xdrs, &objp->int2)) // Encode second argument
    return (FALSE);               // Possibly propagate failure
  return (TRUE);                  // Return success status
}

```

Figure 4: Encoding routine xdr_pair() used in rmin()

```

void xdr_pair(xdrs,objp) {         // Encode arguments of rmin
  // Overflow checking eliminated
  *(xdrs->x_private) = objp->int1; // Inlined specialized call
  xdrs->x_private += 4u;           // for writing the first argument
  *(xdrs->x_private) = objp->int2; // Inlined specialized call
  xdrs->x_private += 4u;           // for writing the second argument
  // Return code eliminated
}

```

Figure 5: Specialized encoding routine xdr_pair()

coding, except that additional dynamic tests must be performed to ensure the soundness and authenticity of the server reply.

We have seen that a systematic approach to specializing system code can achieve significant code simplifications. However, it is not feasible to manually specialize RPC for each remote function, as the process is long, tedious, and error-prone. Since the process is systematic, it can be automated. In fact, it could be as automatic as `rpcgen`. We now discuss how specialization of the previous examples is automated using a partial evaluator for C programs, named Tempo.

4 Automatic Specialization Using the Tempo Partial Evaluator

Partial evaluation [6] is a program transformation approach aimed at specializing programs. We have developed a partial evaluator for C programs, named Tempo. It takes a source program $P_{generic}$ written in C together with a known subset of its inputs, and produces a specialized C source program $P_{special}$, simplified with respect to the known inputs.

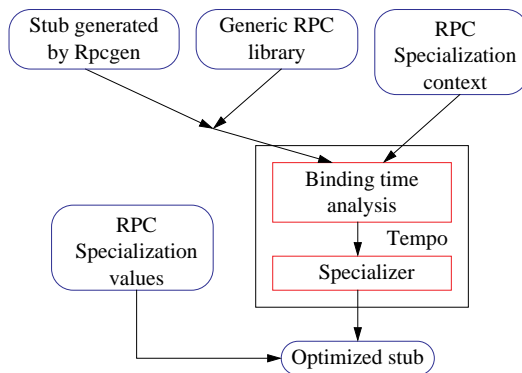


Figure 6: Compile-time specialization with Tempo

The heart of an off-line partial evaluator such as Tempo is the *Binding Time Analysis* (i.e., BTA). The BTA propagates a specialization context describing which inputs are static and which inputs are dynamic throughout the whole program (see Figure 6). After the BTA, $P_{generic}$ is divided into *static* and *dynamic* parts. The static part of $P_{generic}$ is evaluated using concrete specialization values for each known input, while the dynamic part is *residualized* (copied) into the output specialized program. The result $P_{special}$ is typically simpler than $P_{generic}$ since the static part has been pre-computed and only the dynamic part will be executed at runtime.

In the RPC experiment (see Figure 6), $P_{generic}$ is formed from the stub generated by `Rpcgen` and Sun’s generic library. The specialization context is written by a Tempo user, the *expert*, who has an intimate knowledge of the software. Once the context is written, Tempo is used in a fully transparent manner by ordinary RPC programmers.

Issues about writing the specialization context are detailed in Section 6.1.

Accuracy of Specialization. The more precise the BTA, the more it exhibits static constructions (that will be eliminated during specialization). We have more specifically targeted Tempo towards system software. The main refinements introduced in Tempo for this purpose include:

- *partially-static structures*: Figure 3 shows that some fields of the `xdrs` structure are static while others are dynamic. Effective specialization requires that we be able to access the static fields at specialization time. Without such a functionality the whole structure must conservatively be considered dynamic and the repeated buffer overflow checking cannot be eliminated.
- *flow sensitivity*: Possible runtime errors occurring in the decoding of input buffer introduces dynamic conditions after which static information is lost; however each branch of corresponding conditionals can still exploit static information. To this end, binding time of variables (i.e., static or dynamic) must not be a global property; it must depend on the program point considered.
- *static returns*: As seen in the example in section 3.3, the computation at specialization time of exit status tests relies on the ability to statically know the return value of a function call even though its arguments and its actions on input/output buffers are dynamic. More generally, the return value of a function may be static even though its arguments and side-effects are dynamic. Thus we can use the return value of a function call even when the call must be residualized.

Tempo also relies on several other analyses, such as alias (pointer) and dependency analysis. It goes much beyond conventional constant propagation and folding, in the sense that it is not limited to exploiting *scalar* values *intra-procedurally*. It propagates aliases and partially-static data structures *inter-procedurally*. These features are critical when tackling system code. In concrete terms, Tempo is able to achieve all the specializations described in Section 3.

5 Performance Experiments

Having explained the forms of specialization that Tempo performs on the RPC code, we now turn to the assessment of the resulting optimized RPC.

The test program. We have specialized both the client and the server code of the 1984 copyrighted version of Sun RPC. The unspecialized RPC code is about 1500 lines long (without comments) on the client side and 1700 on the server side. The test program, which utilizes remote procedure calls, emulates the behavior of parallel programs that

exchange large chunks of structured data. This is a benchmark representative of applications that use a network of workstations as large scale multiprocessors.

Platforms for measurements. Measurements have been made on two kinds of platforms:

- Two Sun IPX 4/50 workstations running SunOS 4.1.4 with 32 MB of memory connected with a 100 Mbits/s ATM link. The cache is write-through with a size of 64KB. ATM cards are model ESA-200 from Fore Systems. This platform is several years old and quite inefficient compared to up to date products, both in term of CPU, network latency and bandwidth. With faster, more recent ATM cards (i.e., 155 Mbit/s - 622 Mbit/s), we may expect better results due to higher throughput.
- Two 166 MHz Pentium PC machines running Linux with 96 MB of memory and a 100 Mbits/s Fast-Ethernet network connection. The cache size is 512KB. There were no other machines on this network during experiments.

Our specialization is tested on different environments in order to check that the results we obtain are not specific to a particular platform. All programs have been compiled using gcc version 2.7.2, with the option `-O2`.

Benchmarks. To evaluate the efficiency of specialization, we have made two kinds of measurements: (i) a micro-benchmark of the client marshaling process, and (ii) an application level benchmark which measures the elapsed total time of a complete RPC call (round-trip). The client test program loops on a simple RPC which sends and receives an array of integers. The intent of this second experiment is to take into account architectural features such as cache, memory and network bandwidth that affect global performance significantly. Performance comparisons for the two platforms and the two experiments are shown Figure 7. The marshaling and round-trip benchmark numbers result from the mean of 10000 iterations.

Not surprisingly, the PC/Linux platform is always faster than the IPX/SunOS's one. This is partly due to a faster CPU, but also to the fact that the Fast-Ethernet cards have a higher bandwidth and a smaller latency than our ATM cards. A consequence of instruction elimination by the specialization process is that the gap between platforms is lowered on the specialized code (see marshaling comparisons in Figure 7-1 and 7-2).

Marshaling. The specialized client stub code runs up to 3.75 faster than the non-specialized one on the IPX/SunOS, and 3.3 on the PC/Linux. Since the number of instructions eliminated by specialization is linear with array size and the cost header marshaling is fixed, one would expect the speedup to increase with the array size up to a constant asymptote. However, on the Sun IPX the speedup

Client code	Array size				
	20	100	500	1000	2000
generic	20004				
specialized	24340	27540	33540	63540	111348

Table 1: Size of the SunOS binaries (in bytes)

decreases with the size of the array of integers (see Figure 7-5). The explanation is that program execution time is dominated by memory accesses. When the array size grows, most of the marshaling time is spent in copying the integer array argument into the output buffer. Even though specialization decreases the number of instructions used to encode an integer, the number of memory moves remains constant between the specialized and non-specialized code. Therefore, the instruction savings becomes comparatively smaller as the array size grows. On the PC, which has a 512KB cache this behavior does not appear; as expected the speedup curve reaches a constant asymptote.

Round-trip RPC. The application level benchmark results are presented in Figure 7-3,7-4, and 7-6. The specialized code runs up to 1.55 faster than the non-specialized one on the IPX/SunOS, and 1.35 on the PC/Linux. On both platforms, the speedup curve reaches an asymptote. On the IPX/SunOS, this is due to memory behavior explained previously. On the PC/Linux, the reason of the asymptote is due to the fact that several Ethernet packets have to be sent that increases the network latency. In addition to these memory accesses, the Sun RPC includes a call to `bzero` to initialize the input buffer on both the client and server sides. These initializations further increase memory access overhead as the data size grows.

Code size. As shown in Table 1, the specialized code is always larger than the original one. The reason is that the default specialized code unrolls the array encoding/decoding loops completely. It should be noticed that the specialized code is also larger for small array size. This is due to the fact that the specialized code also contains some unspecialized generic functions because of runtime error handling.

While loop unrolling increases code sizes, it also affects cache locality. An additional experiment was conducted on the PC to measure this effect. Since completely unrolling large loops may exceed the instruction cache capacity, we only partially unrolled the loop to adjust its body to the cache size. As shown in Table 2, the resulting code exhibits a lower deterioration of performance as the number of elements grows. Also, this transformation allows code explosion to be limited to the 250 array size case. Currently, this transformation is done manually. However, in the future, we plan to introduce this strategy in Tempo so as to better control loop unrolling and code explosion.

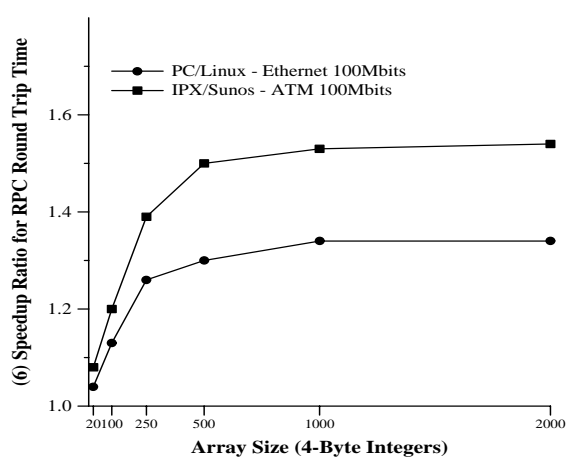
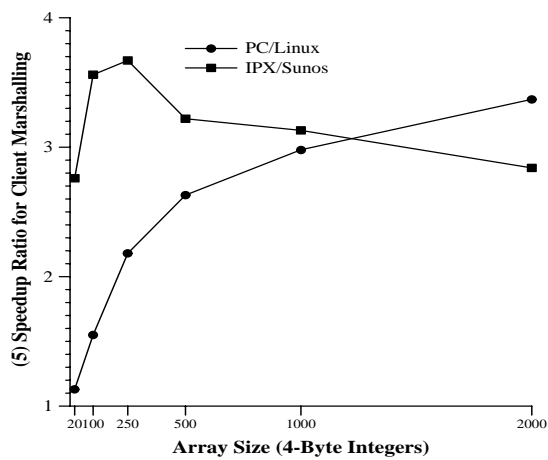
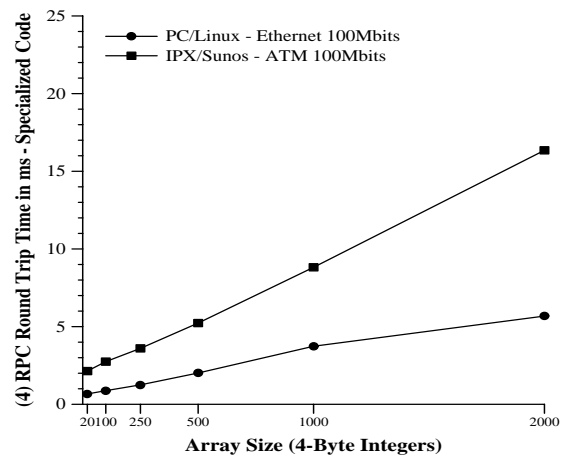
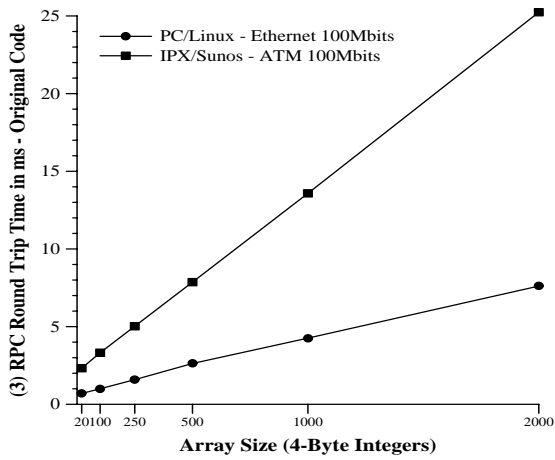
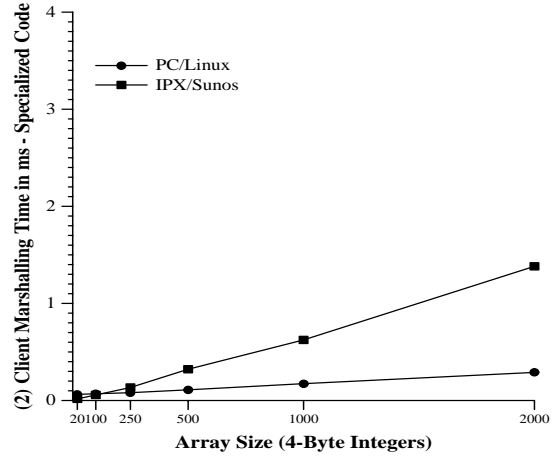
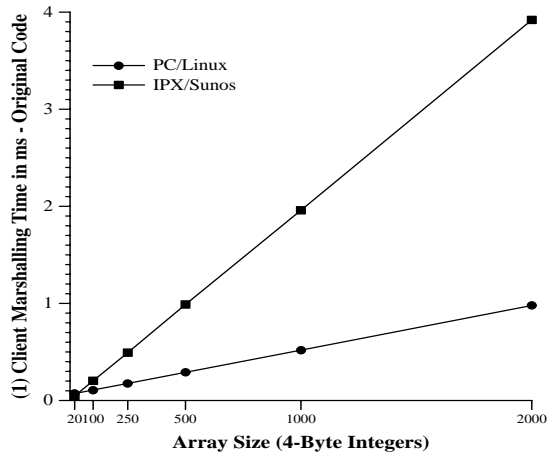


Figure 7: Performance Comparison between IPX/SunOS and PC/Linux

Array Size	PC/Linux				
	Original	Spec.	Speedup	250-unrolled	Speedup
500	0.29	0.11	2.65	0.108	2.70
1000	0.51	0.17	3.00	0.15	3.40
2000	0.97	0.29	3.35	0.25	3.90

Table 2: Specialization with loops of 250-unrolled integers (times in ms)

6 Discussion

In this section we discuss our experience in using Tempo for specialization, the lessons learned from working with existing commercial code, and the relevance of this kind of specialization for general system code.

6.1 Experience with Tempo

In order to treat large programs, Tempo allows the software expert (who writes the specialization context) to visualize the results of the analysis before specialization. Different colors are used to display the static and dynamic parts of a program, thus helping the expert to follow the propagation of the inputs declared as known and assess the degree of specialization to be obtained. After specialization, the expert can compare the original program with the specialized program, and decide whether appropriate reduction and residualization have been carried out.

6.2 Working with Existing Code

An important lesson learned in this experiment is that existing code is a challenge for an optimization technique such as partial evaluation. Indeed, like any other optimization technique, partial evaluation is sensitive to various program features such as program structure and data organization. As a result, specializing an existing program requires an intimate knowledge of its structure and algorithms. It also requires the programmer to estimate what parts of the program should be evaluated away. This is in contrast with a situation where the same programmer both writes and specializes some code: he can structure it with specialization in mind.

Careful inspection of the resulting specialized code shows few opportunities for further optimization without major restructuring of the RPC code. However, Tempo is not a panacea and we (in the expert role) occasionally had to slightly modify the original source code in order to obtain suitable specializations (in practice, to make values available for specialization). Still, the modifications of the original code are performed only once by the software expert. They are fully transparent to the final RPC user. Therefore, this does not contradict our claim of automatically treating existing system code.

6.3 General Applicability

We consider the Sun RPC to be representative of existing system code, not only because it is mature, commercial, and standard code, but also because its structure reflects production quality concerns as well as unrestrained use of the C programming language.

The examples that we highlighted in Section 3 (i.e., dispatching, buffer overflow checking, handling of exit status)

are typical instances of general constructions found in system code and protocols such as TCP/UDP/IP. The fact that Tempo is able to automatically specialize them reinforces our conviction that automatic optimization tools like partial evaluators are relevant for system code production.

Finally, it should be said that analysing an existing software and optimizing it with Tempo takes the software expert only few days; this has to be compared to the several weeks or months required to develop an optimizing compiler.

7 Related Work

The specialization techniques presented in this paper relate to many studies in various research domains such as specific RPC optimizations, kernel level optimizations, operating system structuring, and automatic program transformation. Let us outline the salient aspects of these research directions.

General RPC optimizations. A considerable amount of work has been dedicated to optimizing RPC (see [14, 25, 29]). In most of these studies, a fast path in the RPC is identified, corresponding to a performance-critical, frequently used case. The fast path is then optimized using a wide range of techniques. The optimizations address different layers of the protocol stack, and are performed either manually (by rewriting a layer), or by a domain-specific optimizer.

Marshaling layer optimizations. Clark and Tennenhouse [5] were the first to identify the presentation layer as an important bottleneck in protocol software. They attribute to it up to 97% of the total protocol stack overhead, in some practical applications. Rather than optimizing an existing implementation, they propose some design principles to build new efficient implementations.

O'Malley *et al.* [19] present a universal stub compiler, called USC. As opposed to XDR, which converts between a fixed host format and another fixed external representation, USC converts data between two user-specified formats. USC integrates several domain-specific optimizations, resulting in much faster code than the one produced by XDR. However, in order to perform these aggressive optimizations, USC imposes some restrictions over the marshaled data types: types such as floating point numbers or pointers are not allowed. In fact, USC is not designed for general argument marshaling, but rather for header conversions and interfacing to memory-mapped devices.

Flick [8] is a flexible optimizing IDL (Interface Description Language) compiler recently designed by Eide *et al.*. Flick supports several IDL such as Sun, CORBA and MIG IDLs. It can be used either in traditional distributed applications or for interfacing operating system layers [12]. Due to its flexible internal architecture, flick can match the characteristics of the target machines and implement aggressive optimizations that goes beyond the scope of partial evaluation.

All these studies require one to build a special-purpose code generator, with a complexity ranging from an ad-hoc template assembler to a full, domain-specific, optimizing compiler. In contrast, we take the stubs generated by an existing stub compiler, and derive the specialized stubs with Tempo, a general program specialization tool.

Kernel-level optimizations. It is well recognized that physical memory copy is an important cause of overhead in protocol implementation. Finding solutions to avoid or optimize copies is a constant concern of operating system designers. For instance, copy-on-write was the technique which made message passing efficient enough to allow operating systems to be designed based on a micro-kernel architecture [23, 24]. Buffers are needed when different modules or layers written independently for modularity reasons have to cooperate together at run time. This cause of overhead has been clearly demonstrated by Thekkath and Levy in their performance analysis of RPC implementations [29]. Recent proposals in the networking area explore solutions to improve network throughput and to reduce latency. Maeda and Bershad propose to restructure network layers and to move some functions into user space [16].

Manual specialization. In a first step, operating systems specialization has been performed manually in experiments such as Synthesis [17, 21], and Synthetix [20]. Manual specialization, however, tends to compromise other system properties such as maintainability and portability. Furthermore, manual specialization is typically uniquely tailored to each situation and therefore requires a high degree of programmer skill and system knowledge. While tool-based specialization may not fit the traditional kernel development process, we see it as a natural next step for operating system development the same way compilers became useful programming technology decades ago.

Recently, a semi-automatic approach to program transformation has been developed; it extends the C language to include syntactic constructs aimed at building code at run time [9]. It has been used for a realistic system application, namely the packet filter [10]. This work demonstrates that exploiting invariants can produce significant speedups for demultiplexing messages. This result has been obtained at the cost of manually rewriting a new algorithm that adapts to the specific usage context.

Automatic program transformation. Program transformation has been used successfully for specializing programs in domains such as computer graphics [15]. The key point of program transformation is that it preserves the semantics of the program. Therefore, if the transformation process can be automated, the final code has the same level of safety than the initial program. Tempo relies on partial evaluation [6], a form of program transformation which is now reaching a level of maturity that makes it possible to

develop specializers for real-sized languages like C [1, 7] and apply these specializers to real-sized problems.

C-Mix is the only other partial evaluator for C reported in the literature. Unfortunately, the accuracy of its analyses does not allow it to deal with partially-static structures and pointers to these objects interprocedurally [2].

8 Conclusion

We have described the specialization of Sun RPC using the Tempo partial evaluator. This is the first successful partial-evaluation-based specialization of a significant OS component. The experiment consists of declaring the known inputs of the Sun RPC code and allowing Tempo to automatically evaluate the static parts of code at specialization time. Examples of known information include the number and type of RPC parameters.

There are three reasons why partial-evaluation based specialization is a significant innovation, in comparison to manual specialization [11, 20]. First, Tempo preserves the source code at the programming level, thus preserving the safety and maintainability of a mature commercial code. Second, Tempo achieves speedup up to 3.75 in micro-benchmarks, and 1.5 in test programs. Close inspection of the specialized Sun RPC did not reveal obvious opportunities for further significant improvements without major code restructuring. Third, the successful specialization of commercial code is automatic and shows the promise of industrial application of Tempo.

We carried out experiments on two very different platforms, namely Sun 4/50 workstations running SunOS connected with a 100 Mbits/s ATM link, and 166 MHz Pentium machines running Linux, connected to a 100 Mbits/s Ethernet network. The differences of these platforms and the consistency of performance gains show a robust applicability of Tempo as a tool for partial-evaluation based specialization of layered operating systems code such as Sun RPC.

Acknowledgments

The authors would like to thank the other designers and implementors of Tempo: Jacques Noyé, Luke Hornof, Julia Lawall, Scott Thibault, François Noël and others.

References

- [1] L.O. Andersen. Self-applicable C program specialization. In *Partial Evaluation and Semantics-Based Program Manipulation*, pages 54–61, San Francisco, CA, USA, June 1992. Yale University, New Haven, CT, USA. Technical Report YALEU/DCS/RR-909.
- [2] L.O. Andersen. *Program Analysis and Specialization for the C Programming Language*. PhD thesis, Computer Science Department, University of Copenhagen, May 1994. DIKU Technical Report 94/19.
- [3] B.N. Bershad, S. Savage, P. Pardyak, E. Gün Sirer, M.E. Fiuczynski, D. Becker, C. Chambers, and S. Eggers. Extensibility, safety and performance in the SPIN operating system. In *SOSP95* [28], pages 267–283.

- [4] G. Cabillic and I. Puaut. Stardust: an environment for parallel programming on networks of heterogeneous workstations. *Journal of Parallel and Distributed Computing*, 40:65–80, February 1997.
- [5] D.D. Clark and D.L. Tennenhouse. Architectural considerations for a new generation of protocols. In *SIGCOMM Symposium on Communications Architectures and Protocols*, pages 200–208, Philadelphia, PA, September 1990. ACM Press.
- [6] C. Consel and O. Danvy. Tutorial notes on partial evaluation. In *Conference Record of the Twentieth Annual ACM SIGPLAN-SIGACT Symposium on Principles Of Programming Languages*, pages 493–501, Charleston, SC, USA, January 1993. ACM Press.
- [7] C. Consel, L. Hornof, F. Noël, J. Noyé, and E.N. Volanschi. A uniform approach for compile-time and run-time specialization. In O. Danvy, R. Glück, and P. Thiemann, editors, *Partial Evaluation, International Seminar, Dagstuhl Castle*, number 1110 in Lecture Notes in Computer Science, pages 54–72, February 1996.
- [8] Eric Eide, Kevin Frei, Bryan Ford, Jay Lepreau, and Gary Lindstrom. Flick: A flexible, optimizing IDL compiler. In *Proceedings of the ACM SIGPLAN '97 Conference on Programming Language Design and Implementation*, pages 44–56, Las Vegas, Nevada, June 15–18, 1997.
- [9] D.R. Engler, W.C. Hsieh, and M.F. Kaashoek. 'C: A language for high-level, efficient, and machine-independent dynamic code generation. In *Conference Record of the 23rd Annual ACM SIGPLAN-SIGACT Symposium on Principles Of Programming Languages*, pages 131–144, St. Petersburg Beach, FL, USA, January 1996. ACM Press.
- [10] D.R. Engler and M.F. Kaashoek. DPF: Fast, flexible message demultiplexing using dynamic code generation. In *SIGCOMM'96* [26], pages 26–30.
- [11] D.R. Engler, M.F. Kaashoek, and J.W. O'Toole. Exokernel: An operating system architecture for application-level resource management. In *SOSP95* [28], pages 251–266.
- [12] B. Ford, G. Back, G. Benson, J. Lepreau, A. Lin, and O. Shivers. The Flux OSKit: A substrate for kernel and language research. In *Proceedings of the 1997 ACM Symposium on Operating Systems Principles*, pages 38–51, St-Malo, France, October 1997.
- [13] A. Geist, A. Beguelin, J. Dongarra, W. Jiang, R. Manchek, and V. Sunde. *PVM: Parallel Virtual Machine - A Users' Guide and Tutorial for Networked Parallel Computing*. MIT Press, 1994.
- [14] D.B. Johnson and W. Zwaenepoel. The Peregrine high-performance RPC system. *Software - Practice And Experience*, 23(2):201–221, February 1993.
- [15] B.N. Locanthi. Fast bitblt() with asm() and cpp. In *European UNIX Systems User Group Conference Proceedings*, pages 243–259, AT&T Bell Laboratories, Murray Hill, September 1987. EUUG.
- [16] C. Maeda and B.N. Bershad. Protocol service decomposition for high-performance networking. In *SOSP93* [27], pages 244–255.
- [17] H. Massalin and C. Pu. Threads and input/output in the Synthesis kernel. In *Proceedings of the Twelfth Symposium on Operating Systems Principles*, pages 191–201, Arizona, December 1989.
- [18] Sun Microsystem. NFS: Network file system protocol specification. RFC 1094, Sun Microsystem, March 1989. <ftp://ds.internic.net/rfc/1094.txt>.
- [19] S. O'Malley, T. Proebsting, and A.B. Montz. USC: A universal stub compiler. In *Proceedings of Conference on Communication Architectures, Protocols and Applications*, London (UK), September 1994.
- [20] C. Pu, T. Autrey, A. Black, C. Consel, C. Cowan, J. Inouye, L. Kethana, J. Walpole, and K. Zhang. Optimistic incremental specialization: Streamlining a commercial operating system. In *SOSP95* [28], pages 314–324.
- [21] C. Pu, H. Massalin, and J. Ioannidis. The Synthesis kernel. *Computing Systems*, 1(1):11–32, Winter 1988.
- [22] R. Ramsey. *All about administering NIS+*. SunSoft, 1993.
- [23] R.F. Rashid, A. Tevanian Jr., M.W. Young, D.B. Golub, R.V. Baron, D. Black, Bolosky W.J., and J. Chew. Machine-independent virtual memory management for paged uniprocessor and multiprocessor architectures. *IEEE Transactions on Computers*, 37(8):896–908, August 1988.
- [24] V. Rozier, V. Abrossimov, F. Armand, I. Boule, M. Gien, M. Guillemont, F. Herrmann, C. Kaiser, S. Langlois, P. Léonard, and W. Neuhauser. Overview of the Chorus distributed operating system. In *USENIX - Workshop Proceedings - Micro-kernels and Other Kernel Architectures*, pages 39–70, Seattle, WA, USA, April 1992.
- [25] M.D. Schroeder and M. Burrows. Performance of Firefly RPC. *ACM Transactions on Computer Systems*, 8(1):1–17, February 1990.
- [26] *SIGCOMM Symposium on Communications Architectures and Protocols*, Stanford University, CA, August 1996. ACM Press.
- [27] *Proceedings of the 1993 ACM Symposium on Operating Systems Principles*, Asheville, NC, USA, December 1993. ACM Operating Systems Reviews, 27(5), ACM Press.
- [28] *Proceedings of the 1995 ACM Symposium on Operating Systems Principles*, Copper Mountain Resort, CO, USA, December 1995. ACM Operating Systems Reviews, 29(5), ACM Press.
- [29] C.A. Thekkath and H.M. Levy. Limits to low-latency communication on high-speed networks. *ACM Transactions on Computer Systems*, 11(2):179–203, May 1993.