

Code Compaction of an Operating System Kernel*

Haifeng He, John Trimble, Somu Perianayagam, Saumya Debray, Gregory Andrews
Department of Computer Science, The University of Arizona, Tucson, AZ 85721, USA
{hehf, trimble, somu, debray, greg}@cs.arizona.edu

Abstract

General-purpose operating systems, such as Linux, are increasingly being used in embedded systems. Computational resources are usually limited, and embedded processors often have a limited amount of memory. This makes code size especially important. This paper describes techniques for automatically reducing the memory footprint of general-purpose operating systems on embedded platforms. The problem is complicated by the fact that kernel code tends to be quite different from ordinary application code, including the presence of a significant amount of hand-written assembly code, multiple entry points, implicit control flow paths involving interrupt handlers, and a significant amount of indirect control flow via function pointers. We use a novel “approximate decompilation” technique to apply source-level program analysis to hand-written assembly code. A prototype implementation of our ideas, applied to a Linux kernel that has been configured to exclude unnecessary code, obtains a code size reduction of over 25%.

1 Introduction

Recent years have seen increasing use of general-purpose operating systems, such as Linux, deployed in embedded contexts such as cell phones, media players, and other consumer electronics [1]. This is due in great part to technological trends that make it uneconomical for vendors to develop custom in-house operating systems for devices with shorter and shorter life cycles. At the same time, however, these operating systems—precisely because they are general-purpose—contain features that are not needed in every application context, and which incur unnecessary overheads, e.g., in execution speed or memory footprint. Such overheads are especially undesirable in embedded processors and applications because they usually have resource constraints, such as a limited amount of memory. Thus, the memory footprint of the code is especially important, and a program that requires more memory than is available will not be able to run.¹

This paper focuses on automatic techniques for reducing the memory footprint of operating system kernels in embedded systems. Such systems tend to have relatively static configurations: at the hardware end, they are limited in the set of devices with which they interact (e.g., a cell phone or digital camera will typically not have a mouse interface); at the software end, they usually support a fixed set of applications (e.g., we do not routinely download or build new applications on a cell phone or digital camera). This implies that an embedded system will typically use only some of the functionality offered by a general-purpose operating system. The code corresponding to the unused functionality is unnecessary overhead, and should be removed. Some of this overhead can be removed simply by configuring the kernel carefully so as to exclude as much unnecessary code as possible. However, not all overheads can be removed in this manner. For example, a given set of applications running on an embedded platform will typically use only a subset of the system calls supported by the operating system; the code for the unused system calls is then potentially unnecessary. Such unnecessary code typically cannot be eliminated simply by tweaking the configuration files; additional analysis is required. This paper discusses how such analysis may be carried out in order to identify code that can be guaranteed to be unnecessary. Specifically, this paper makes the following contributions:

*This work was supported in part by NSF Grants EIA-0080123, CCR-0113633, and CNS-0410918.

¹Virtual memory is not always an option in embedded systems; even where it is available, the energy cost of paging out of secondary storage can be prohibitive.

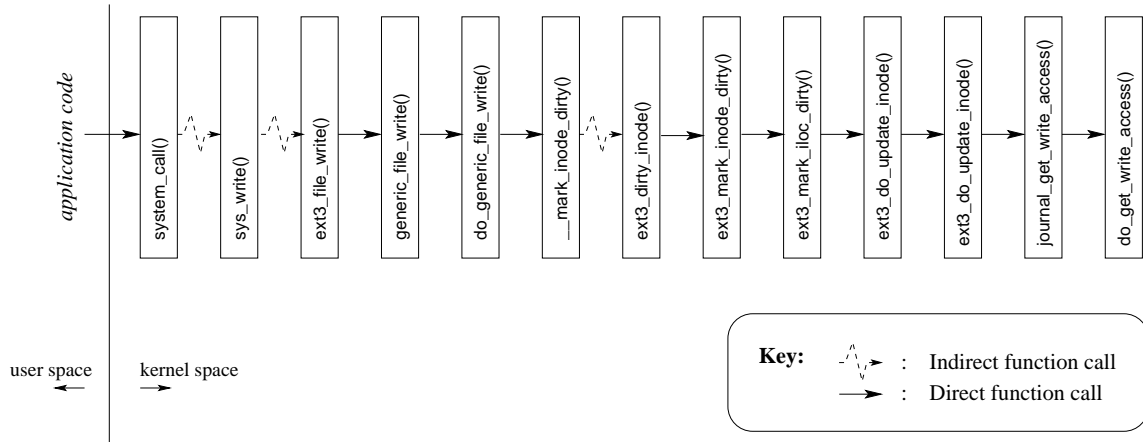


Figure 1: An example call chain in the Linux kernel

1. It discusses issues that arise in binary rewriting of operating system kernels and discusses how they may be handled.
2. It introduces a notion of “approximate decompilation” that allows us to apply source-level program analysis to hand-written assembly code. This improves the precision of our analysis considerably, and in a simple way, while ensuring that the safety of our transformations is not compromised.

There has been a significant body of work on code size reduction, with different objectives (e.g., reducing wire transmission time vs. reducing memory footprint during execution), on different program representations (e.g., syntax trees, byte code, or native code), and making different assumptions regarding the requirement for runtime decompression and the availability of hardware support for any such runtime decompression. See Beszédes *et al.* [4] for a survey. Almost all of this work has been done in the context of application code. By contrast, the work described in this paper focuses on operating system kernels in native code representation, and it aims to reduce their runtime memory footprint. We do so using automatic program analysis and transformations, with a minimum of programmer intervention in the form of code annotations or other similar input. A prototype implementation of our ideas, applied to the Linux kernel configured minimally so as to exclude all unnecessary code, is able to achieve a code size reduction of over 25%, on average, on the MiBench suite [9].

2 Background

OS kernel code tends to be quite different from ordinary application code. A significant problem is the presence of considerable amounts of hand-written assembly code that often does not follow the familiar conventions of compiler-generated code, e.g., with regard to function prologues, epilogues, and argument passing. This makes it difficult to use standard compiler-based techniques for whole-system analysis and optimization of kernel code. To deal with this problem, we decided to use binary rewriting for kernel compaction. However, while processing a kernel at the binary level provides a uniform way to handle code heterogeneity arising from the combination of source code, assembly code and legacy code such as device drivers, it introduces its own set of problems. For example, it is necessary to deal with a significant amount of data embedded within executable sections, implicit addressing requirements that constrain code movement, and occasional unusual instruction sequences. As a result, binary rewriting techniques applicable to application code do not always carry over directly to kernel code.

An especially important issue for code compaction is that of control flow analysis, both intra-procedural and inter-procedural. This is because in practice, most of the code size reductions arising from compaction

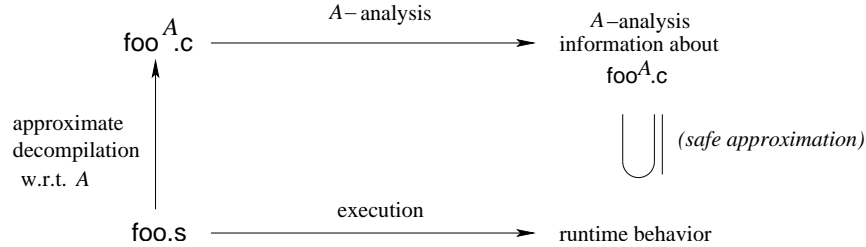


Figure 2: Using approximate decompilation for program analysis

comes from the detection and elimination of dead and unreachable code [6]. For soundness reasons, we have to ensure that we only eliminate code that can be guaranteed never to be needed during any future execution. This means that imprecision in control flow analysis directly affects the amount of code that can be eliminated.

Unfortunately, control flow analysis in operating system kernels is complicated by the interaction of two separate problems. First, there are significant amounts of hand-written assembly code, as mentioned above. Second, operating system kernels often make extensive use of indirect function calls in order to enhance maintainability and extensibility. The situation is illustrated by Figure 1, which shows a particular frequently executed chain of function calls within the Linux kernel corresponding to the *write()* system call. It can be seen that there are three indirect calls on the call path leading from the top-level routine *system_call()* to the deeper-level routine *do_get_write_access()*. This is a problem because static analyses are generally quite conservative in their treatment of indirect function calls.² Each of these problems—hand-written assembly and indirect function calls—is nontrivial in its own right, and the situation is exacerbated further by the fact that they interact: the hand-written assembly code in an operating system kernels may itself contain indirect function calls, and identifying those targets requires pointer alias analysis of the assembly code.

A final problem in dealing with control flow in operating system kernels is that not all entry points into the kernel, and control flow within the kernel, are explicit. There are implicit entry points such as system calls and interrupt handlers, as well as implicit control flow arising from interrupts, that have to be taken into account in order to guarantee soundness.

3 Pointer Analysis: Resolving Indirect Function Call Targets

As mentioned above, operating system kernels contain a significant amount of hand-written assembly code. This makes program analysis problematic. On the one hand, dealing with hand-written assembly code in a source-level or intermediate-code-level analysis is messy and awkward because of the need to inject architecture-specific knowledge into the analysis—such as aliasing between registers (e.g., in the Intel x86 architecture, the register `%al` is an alias for the low byte of the register `%eax`) and idiosyncrasies of various machine instructions. On the other hand, if the analysis is implemented at the assembly code or machine code level, much of the semantic information present at the source level is lost—in particular, information about types and pointer aliasing—resulting in overly conservative analysis that loses a great deal of precision.

One possible solution to this problem would be to decompile the hand-written assembly code back to equivalent C source code that could then be analyzed by source-level analysis. The problem with such an approach is that it is not obvious that all of the kernel assembly code can be reverse engineered back to equivalent C source code; for example, “system instructions” on the Intel x86 architecture, which include instructions such as “load interrupt descriptor table register” and “invalidate TLB entry”, may not have

²In general, identifying the possible targets of indirect function calls is equivalent to pointer alias analysis, which is a hard problem both theoretically and in practice.

obvious C-level counterparts. Moreover, even in situations where this is possible, it can be complicated and involve a great deal of engineering effort. Instead, we deal with this problem using an approach we call “approximate decompilation,” which maps hand-written assembly code back to C source files for analysis purposes. The idea, illustrated in Figure 2, is that given an assembly file `foo.s` and a program analysis A , we create a source file `fooA.c` that has the property that an A -analysis of `fooA.c` is a safe approximation of the behavior of `foo.s`, even though `fooA.c` is not semantically equivalent to `foo.s`. For example, if A focuses on control flow analysis, then `fooA.c` may elide those parts of `foo.s` that are irrelevant to control flow. We have applied this approach to use a source-level pointer alias analysis technique called FA-analysis to identify the possible targets of indirect function calls. The remainder of this section discusses how this is carried out.

3.1 FA Analysis

There is a large volume of literature on pointer alias analysis, with a variety of assumptions, goals, and trade-offs (see, for example, the discussion by Hind and Pioli [10]). In general, these analyses exhibit a trade-off between efficiency and precision: the greater the precision, the greater the analysis cost, i.e., the lower the efficiency. FA-analysis is a flow-insensitive context-insensitive pointer alias analysis, originally due to Zhang *et al.* [18, 19], that is at the low end of this efficiency/precision trade-off. Milanova *et al.* have observed, however, that even though this analysis is not always very precise for general-purpose pointer alias analysis, it turns out to be quite precise in practice for identifying the targets of indirect function calls [14]. The authors attribute this to the fact that programmers typically use function pointers in a few specific and relatively simple stylistic ways.³

3.2 Approximate Decompilation of Kernel Assembly Code for FA Analysis

As Figure 2 suggests, the way in which approximate decompilation is carried out depends in part on the source-level analysis that will be applied to the resulting source files. This section discusses approximate decompilation of assembly code in the Linux kernel code for FA analysis. For concreteness, we discuss kernel assembly code on the Intel x86 architecture.

The hand-written assembly instructions in the Linux kernel falls into two broad groups: (1) general-purpose instructions that perform basic data movement, arithmetic, logic, and program control flow operations, and (2) system instructions that provide support for operating systems and executives [11]. We process these instructions as follows:

- System instructions (the second group above) manipulate only the hardware (or data related to the hardware) and have no effect on pointer aliasing in the kernel code. For pointer alias analysis, therefore, we simply ignore these instructions.
- Since FA analysis is flow-insensitive and context-insensitive, instructions whose only effect is on intra-procedural control flow, such as conditional and unconditional branches, have no effect on the analysis. Inter-procedural control flow cannot be ignored, however, since it induces aliasing between the actual parameters at the call site and the formal parameters at the callee. Our decompiler therefore ignores conditional and unconditional control flow instructions whose targets are within the same function, but translates inter-procedural control transfers.
- The remaining instructions are those that move data and those that perform arithmetic and logic operations. These instructions are translated to the corresponding operations in C. For example, a register load instruction, ‘`mov $0, %eax,`’ is translated to an assignment ‘`eax = 0`’.

³A detailed discussion of FA-analysis is beyond the scope of this paper. The interested reader is referred to the original papers on this analysis [14, 18, 19].

Our decompiler maps registers in the assembly code to global variables of type *int* with 32-bit values. For example, the 32-bit register `%eax` is mapped to a variable `eax` of type *int*. Since we are only interested in capturing potential aliasing relationships between objects, and not necessarily the actual values computed, we map the 16-bit and 8-bit registers (which are aliases of parts of the 32-bit registers) into the appropriate 32-bit global. Thus, the 8-bit register `%al` and the 16-bit register `%ax`, which refer to the low 8 bits and the low 16 bits of the 32-bit register `%eax` respectively, are both mapped to the variable `eax` denoting the 32-bit register `%eax`.

Memory locations referenced in the assembly code are also treated as global variables. Since there is little type information available at the assembly level, we declare memory locations as having type `MEMOBJ`, which denotes a word in memory.⁴ An object spanning a series of memory locations in the assembly code is treated as an array of `MEMOBJ` in the generated C code. This is illustrated in Figure 3(a). The segment of assembly code shown on the left side in Figure 3(a), taken from the file `entry.S` in the Linux kernel, defines the system call table that contains the function addresses of all system call handlers. Since `sys_call_table` spans a series of (initialized) memory locations in the assembly code, we map it to an (initialized) array in the generated C code shown on the right side in Figure 3(a). Moreover, since the symbols for the system call handlers are not themselves defined in `entry.S`, they are declared as `extern` objects in the generated C code. Before we start the actual pointer analysis, we scan the entire kernel source code and match memory objects to functions so that the source-level FA analysis can deal properly with function pointers in the assembly code.

Functions in the assembly code are identified from symbol table information and mapped to functions in the generated C code. Memory locations accessed through the stack pointer register `%esp` are assumed to be on the stack; these are mapped to local variables in the corresponding C function. Since the value of the stack pointer may change from one point to another in the code, different references to the same stack location from different points in the code may use different displacements off the stack pointer. To handle this, we track the location of stack register `esp`. At the entry to a function, we assign a symbolic constant ℓ denoting the entry value for the stack pointer. Then, each time the value of the stack pointer is changed—either explicitly via an arithmetic operation or implicitly via a `push` or `pop` instruction—we update the distance between the current stack pointer and ℓ . This allows us to map each data reference off the stack pointer to an offset from the initial value ℓ , which can then be translated to a reference to a local variable.⁵ For example, in the assembly code in Figure 3(b), the first `push` instruction access the initial stack location at ℓ . We define a local variable, `LOCAL1` for this stack location and translate the instruction into assignment statement `LOCAL1 = 0` in C, then update the value of `%esp` to reflect the effect of the `push` instruction. The second `push` instruction, which writes to location $\ell - 4$, is then translated into an assignment statement to a different local variable `LOCAL2`.

Actual parameters to a call are identified similarly: if the index of a data reference, computed as described above, indicates that it is deeper in the stack than the initial stack pointer value ℓ , then the location being referred to is an argument that has been passed to the current function. In this manner, by examining the references to actual parameters in the body of a function, we can determine the number of arguments it takes, and thereby generate a function prototype in the C code. Such prototypes are then used by the FA analysis to identify aliasing between actuals and formals. A control transfer to a symbol S is translated as a function call if either the instruction is a `call` instruction, or if the target S is a function. For example, in the segment of assembly code in Figure 3(c), there are two references to stack locations that are deeper than the initial stack pointer value. From this, we infer that there are two parameters needed by function `error_code`. Since the symbol `error_code` has been identified as a function, the instruction `jmp er-`

⁴Our implementation defines this type as `typedef MEMOBJ int;`

⁵This discussion refers to base-displacement addressing, which is typically used to access scalar variables on the stack. Similar reasoning can be used to process more complex addressing modes typically used for local arrays.

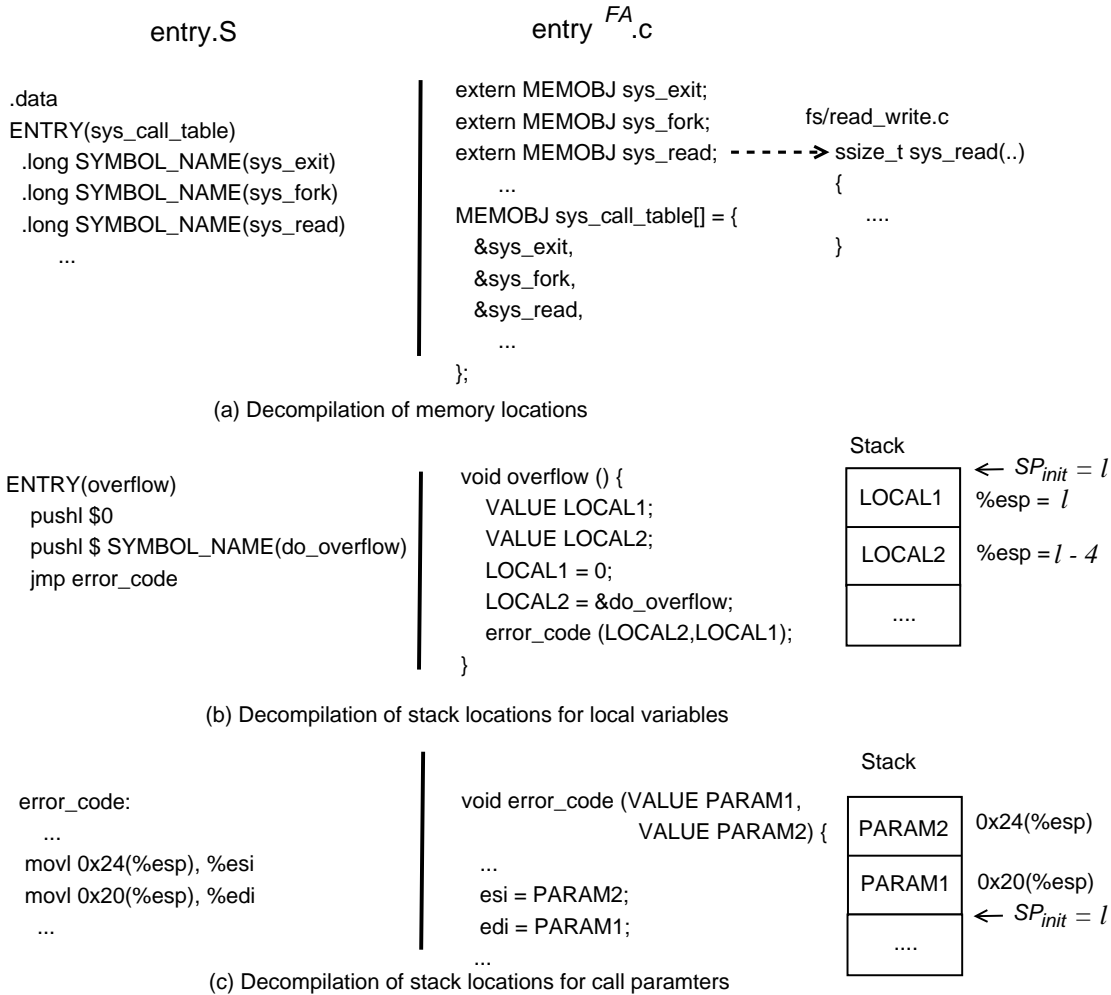


Figure 3: Examples of approximate decompilation of assembly code to C code for FA analysis.

ror_code’ in the assembly code in Figure 3(b) is translated into the function call ‘error_code (LOCAL2 , LOCAL1) ’ in the C code in Figure 3(b).⁶

4 Identifying Reachable Code

4.1 Reachability Analysis

For each indirect procedure call in the kernel, the source-level FA analysis produces a set of possible call targets for that indirect call. Our kernel binary rewriter takes this information as an input and constructs a program call graph for the entire kernel.

Unlike ordinary applications, an operating system kernel contains multiple entry points. These entry points are the starting points for our reachability analysis. We classify kernel entry points into four categories: (1) the entry point for initializing the kernel (for the Linux kernel, this is the function `startup_32`), (2) system calls invoked during the kernel boot process, (3) interrupt handlers, and (4) system calls invoked by user applications. Once the entry points into the kernel have been identified, our reachability analysis per-

⁶In the assembly code in Figure 3(b), since `jmp` is used instead of `call` for control transfer, no return address is stored on the stack.

Procedure Reachability-Analysis

```
worklist ← functions that are entry points into the kernel
while sizeof(worklist) > 0 do
  f ← Remove a function from worklist
  Mark f as reachable
  for every indirect/direct call target c of f do
    if c is static ∧ c is not executed based on profile then
      continue
    else if c is not marked as reachable then
      Add c into worklist
    end if
  end for
end while
```

Figure 4: The improved reachability analysis algorithm

forms a straightforward depth-first traversal of the program call graph to identify all the reachable functions in the kernel.

4.2 Improving the Reachability Analysis

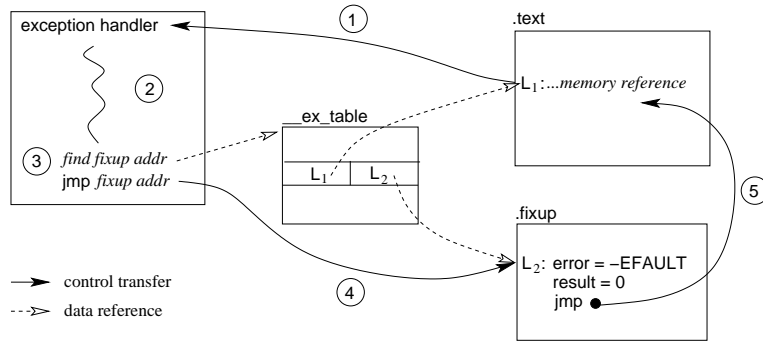
During the initialization phase of kernel bootup (e.g., before the `init` program in Linux begins execution), execution is deterministic because there is only one active thread and execution depends only on the hardware configuration and the configuration options passed through boot command line. In other words, the initialization code can be considered to be “static” in the partial evaluation sense [12]. This means that if the configuration is not changed, we can safely remove any initialization code that is not executed. We use this idea to further improve our reachability analysis.

Our goal is to identify the static functions in the kernel, i.e., functions whose execution is completely determined once the command-line configuration options and the hardware are fixed. To this end, we take advantage of a Linux kernel feature used to identify initialization code, most of which is not needed, and can be reclaimed, after initialization is complete. In particular, the Linux kernel simplifies this reclamation by segregating data and code used only for initialization into two sections in the ELF binary: `.text.init` and `.data.init`. Once the initialization of the kernel finishes during bootup, the kernel frees the memory pages occupied by these two sections to save physical kernel memory.

We use this knowledge to initialize the set of *static* functions to those appearing in the `.text.init` section. We then propagate this information as follows to find other functions that are not in the `.text.init` section but whose execution can be inferred to be completely determined given the command-line configuration options and hardware setup:

1. Mark all functions in `.text.init` section as *static*.
2. Based on the call graph of the Linux kernel, mark all functions that are not called by any other function as *static*.
3. If all the direct and indirect callers of a function F are *static*, then mark F as *static*. Repeat this process until there are no changes.

Once we have computed the set of functions that are considered to be *static* during kernel initialization, we use the results to improve our reachability analysis as shown in Figure 4. The improvement is that when a potentially reachable function is found, if the function is marked *static* and if, based on profile data, it was not called during kernel initialization, then we do not add it to the set of reachable functions.



Key:

- ① A memory exception at `L1` causes control to branch to the exception handler.
- ② Exception handling code.
- ③ Exception handler searches `_ex_table` with the address `L1`, where the exception occurred, to find the associated fixup code address `L2`.
- ④ Control branches from the exception handler to the fixup code.
- ⑤ Control branches from the fixup code to the instruction following the location `L1` where the exception had occurred.

Figure 5: Control flow during the handling of exceptions in the Linux kernel

4.3 Handling Exception Handlers

In order to identify all reachable code in the kernel, it is not enough to consider ordinary control transfers, which are explicit in the code: we also have to take into account control transfers that are implicit in the exception handling mechanisms of the kernel. For this, we examine the exception table in the kernel.

Locations in the kernel where an exception could be generated are known when the kernel is built. For example, the kernel code that copies data to/from user space is known as a potential source for a page fault exception. The Linux kernel contains an exception table, `_ex_table`, that specifies, for each such location, the code that is to be executed after handling an exception. Additionally, a special section, `.fixup`, contains snippets of code that carry out the actual control transfer from the exception handlers to the appropriate destination locations. The flow of control when handling an exception is shown in Figure 5: after the exception handler deals with an exception from an address `L1`, it searches `_ex_table` with `L1` as the key, finds the associated address `L2` of the corresponding fixup code, and jumps to `L2`. The key point to note here is that the control flow path from `L1` to `L2` is not explicit in the code, but is implicit in `_ex_table`. It is necessary to take such implicit execution paths into account for code compaction to ensure that we find all reachable code. We do this by examining the exception table and adding pseudo-control-flow edges to indicate such implicit control flow. For the example in Figure 5, we would add such an edge from `L1` to `L2`.

5 Kernel Compaction

Once all the potentially reachable code in the kernel has been identified, code compaction can be carried out by deleting unreachable code [6]. Further code size reductions can be obtained by applying code factoring to the reachable parts of the kernel, which entails identifying repeated code fragments and abstracting them into procedures. However, applying these transformations to the kernel code involves some subtleties. For example, some kernel functions are required to be at specific fixed addresses, and cannot be moved. Another issue is that some forms of procedural abstraction require a memory location be allocated to save the return

Benchmark set	Programs	No. of unique system calls	No. of non-bootup system calls
Auto./Industrial	basicmath, bitcount, qsort, susan	33	9
Consumer	jpeg, mad, lame, tiff2bw, tiff2rgba, tiffdither, tiffmedian, typeset	46	11
Network	dijkstra, patricia (blowfish, CRC32, sha)	43	12
Office	ghostscript, ispell, rsynth, stringsearch	57	15
Security	blowfish, pgp, rijndael, sha	49	10
Telecomm	adpcm, CRC32, FFT, gsm	39	11
Entertainment	jpeg, lame, mad	43	10
Cellphone	blowfish, sha, CRC32, FFT, gsm, typeset	45	12

Table 1: Characteristics of the benchmarks used (from the MiBench suite [9])

address of the procedure. We currently do not consider such code fragments for procedural abstraction within the kernel for two reasons. First, if a page fault occurs when accessing this location to store a return address and the page tables have not yet been initialized, the kernel will crash. Second, since the kernel is multi-threaded in general, using a single global location can lead to incorrect results if one thread overwrites the return address stored there by another thread; this means that the memory allocation has to be done on a per-thread basis, which complicates the implementation and reduces its benefits.

6 Experimental Results

To get an accurate evaluation of the efficacy of our ideas, we begin with a minimally configured kernel where as much unnecessary code as possible has been eliminated by configuring the kernel carefully. For our experiments, therefore, we configured the Linux 2.4.31 kernel to remove modules, such as the sound card, video support, and network driver, that are not required to run our benchmarks. The size of the resulting Linux kernel image, compiled with *gcc* version 3.4.4 using the default makefile flags of ‘-O2 -fomit-frame-pointer’, is 1,152,869 bytes. In order to simplify the booting process of the Linux kernel, we modified the kernel boot up file *inittab* so that the Linux kernel will run in single user mode (level 1). Based on the profile data, there are 81 different system calls called during the booting process.

We used programs from the MiBench suite [9], a widely used and freely available collection of benchmark programs for embedded systems, to evaluate our approach. The MiBench suite is organized into six sets of benchmarks, corresponding to different kinds of embedded environments: Automotive and industrial control, Consumer devices, Networking, Office automation, Security, and Telecommunications; each of these sets contains several different application programs. We augmented this with two additional sets: Entertainment, representing a multi-media consumer appliance for music and digital pictures; and Cellphone, representing a cell phone with security features. Characteristics of these sets are shown in Table 1.

Before we can carry out kernel compaction for any given benchmark set, we have to identify the system calls that can arise from programs in that set. It is not enough to examine their executions using tools such as *strace*, since this may not cover all the execution paths in the programs; nor is it enough to simply examine the source code of the benchmarks for system calls, since these actually call library routines that may contain additional system calls not visible in the source code. We therefore analyze statically linked binaries of the programs to ensure that we find all the system calls that may be invoked. This, however, causes the entire C library to be linked in. We address this problem by first carrying out a reachability analysis on the application program binaries to identify and eliminate unreachable library routines (using a conservative approximation to deal with indirect function calls) and then traversing the resulting whole-program control flow graph to

Kernel	.text.init section		.text section		Total	
	size (bytes)	ratio	size (bytes)	ratio	size (bytes)	ratio
Original	75,638	1.000	735,139	1.000	810,777	1.000
Auto./Industrial	52899	0.699	548,146	0.746	601,045	0.741
Consumer	52915	0.700	550,607	0.749	603,522	0.744
Network	52899	0.699	548,209	0.746	601,108	0.741
Office	52915	0.700	549,245	0.747	602,160	0.743
Security	52899	0.699	548,868	0.747	601,767	0.742
Telecomm	52899	0.699	548,338	0.746	601,237	0.742
Entertainment	52899	0.699	549,885	0.748	602,784	0.743
Cellphone	52899	0.699	549,060	0.747	601,959	0.742
GEOM, MEAN:		0.699		0.747		0.742

Table 2: compaction result

determine the set of possible system calls. These data are shown in Table 1: the third column of this table gives the number of different system calls across all of the programs in each set of benchmarks, while the fourth column gives, for each benchmark set, the number of system calls not occurring in the set of system calls invoked during the kernel bootup process. Once we have the system calls that may be invoked by a set of programs, we use them to identify and eliminate unreachable code in the kernel.

The results of code compaction are shown in Table 2. For each benchmark set, we present three sets of numbers; these give the amount of compaction achieved for the `.text.init` section (the code used for kernel bootup), the `.text` section (the kernel code used during steady-state execution), and the total amount of code (`.text.init` and `.text` together). We are able to reduce the size of the kernel bootup code by 30%, and the steady-state kernel code by just over 25%. Overall, the savings in code size come to just under 26% on average.

Since the bulk of our code compaction is achieved through the identification of unreachable code, the amount of compaction achieved for each set of applications shown above depends on the particular set of system calls made by the set of application programs. To get an idea of the extent to which our results might generalize to other sets of embedded applications, we evaluated the “popularity” of different system calls across the MiBench suite. The popularity of a given system call s in a set of programs P is given by the fraction of programs in P that use s . Intuitively, if different kinds of applications use very different sets of system calls, i.e., many system calls have low popularity, then our results may not generalize well; on the other hand, if different kinds of applications tend to have mostly-similar sets of system calls, then we can expect these results to generalize. The results are shown in Figure 6. It can be seen that out of some 226 different possible system calls in our system,⁷ there is a small core of 32 system calls that are used by every program. Popularity drops off sharply outside this core set. All of our benchmark programs, taken together, refer to only 76 system calls, i.e., about a third of the total set of system calls.

This relative uniformity in system call usage across a wide variety of applications helps explain the surprising uniformity of our code compression results across all of the benchmark sets. While the applications themselves are very different in terms of their nature and code size, the popularity data shown in Figure 6 show that they do not differ from each other hugely in terms of their interactions with the operating system kernel: for example, they typically read data from some files, process that data, and write out the results.

⁷There were 259 syscall entries in our version of Linux; of these, 33 were not implemented (“`sys_ni_syscall`”), leaving a total of 226. Other sources put the number of Linux system calls much higher: e.g., Wikipedia mentions “almost 300” system calls for Linux, while according to `syscalls(2)` in the Linux programmer’s manual, there are some 1,100 system calls listed in `/usr/src/linux/include/asm-*/unistd.h`.

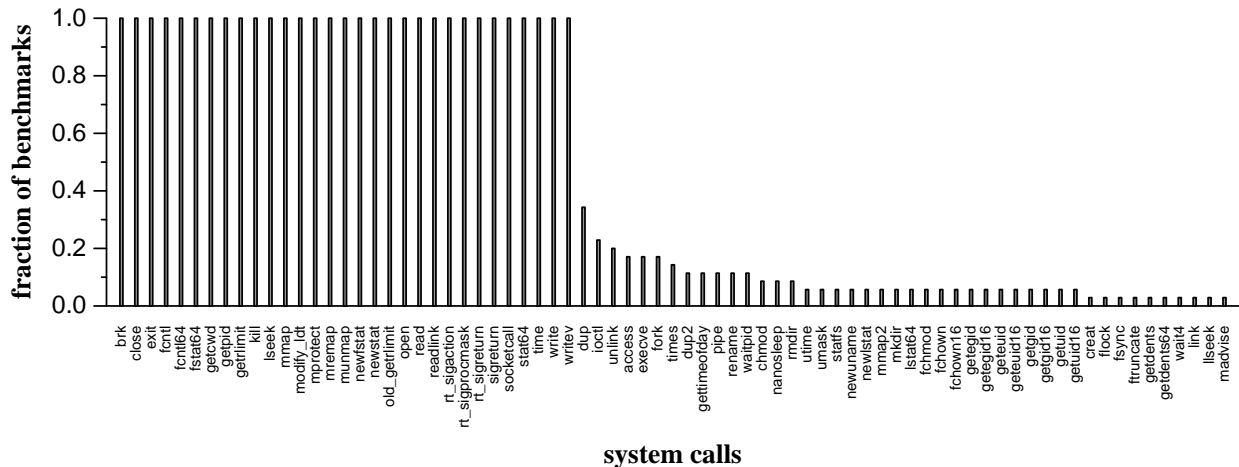


Figure 6: System call “popularity” in embedded applications

Moreover, in addition to the system calls made by application code, a significant number of system calls are made from within the kernel itself. For example, the kernel makes 81 different system calls during the bootup process. The overall result is that the set of “non-bootup” system calls arising in the application code is relatively small, and does not vary greatly from one benchmark set to another (Table 1, col. 4). Because of this, the compaction results for the different benchmark sets tend to be similar.

7 Related Work

The work that is closest to ours is that of Chanet *et al.*, who describe a system for code compaction of the Linux kernel to reduce its memory footprint [5]. Our approach is more general, in that the Diablo system described by Chanet *et al.* relies on a modified compiler tool chain and requires special annotations (currently manually applied) to deal with hand-coded assembly. Furthermore, the techniques we use, in particular the use of approximate decompilation for program analysis, are quite different from theirs. Finally, the code size reductions we obtain, of around 25% on average, are significantly higher than those of Chanet *et al.*, who report an overall code size reduction of about 19%; while this is not entirely an apples-to-apples comparison, because the set of application programs—and, therefore, the set of system calls—considered by Chanet *et al.* were different from those we considered, it gives an approximate idea of the relative compaction power of the two approaches.

We are not aware of a great deal of work on binary rewriting of operating systems kernels. Flowers *et al.* describe the use of Spike, a binary optimizer for the Compaq Alpha, to optimize the Unix kernel, focusing in particular on profile-guided code layout [7]. A number of researchers have looked into specializing operating system kernel code [13, 16, 15]. These generally focused on improving execution speed rather than reducing code size and therefore used techniques very different from ours.

There is a considerable body of work on code compaction; Beszédes *et al.* give a comprehensive survey [4]. Almost all of this work is in the context of application programs in high-level languages and does not consider the issues that arise when dealing with an OS kernel.

One of the technical issues of considerable importance in code compaction of an OS kernel is that of resolving the possible targets of indirect function calls. This problem resembles the problem of identifying the possible targets of virtual method invocations in object-oriented programs, which has received considerable attention in the object-oriented community [2, 3, 8, 17]. Much of this work relies on using type information to resolve the possible targets of a virtual method invocation. We applied this idea of using type information to determine possible targets of indirect calls in the Linux kernel but the results were disappointing,

because the available type information is not discriminating enough to give an acceptable level of precision (especially when taking into account the possibility of type casts on pointers). The FA analysis proved to be far more effective: a type-based signature matching analysis we implemented gave, on average, about 96 targets per indirect function call in the Linux kernel; with FA analysis, by contrast, the average number of targets per indirect call is 10.

8 Conclusions

Because of the limited amount of memory typically available in embedded devices, it is important to reduce the memory footprint of the code running on such devices. This paper describes an approach to code compaction of operating system kernels. We begin with the observation that embedded systems typically run a small fixed set of applications. This knowledge can be used to identify the minimal functionality required of the kernel code to support those applications and then to discard unnecessary code. We discuss a number of technical challenges that have to be addressed in order to make this work; in particular, we describe a novel approach of “approximate decompilation” that allows us to apply source-level program analyses to hand-written assembly code. Our ideas have been implemented in a prototype binary rewriting tool that is able to achieve a code size reduction of over 25% on an already minimally-configured Linux kernel.

References

- [1] The Linux mobile phones showcase, February 2006.
<http://www.linuxdevices.com/articles/AT9423084269.html>.
- [2] G. Aigner and U. Hölzle. Eliminating virtual function calls in C++ programs. In *Proc. European Conference on Object-Oriented Programming (ECOOP)*, volume 1098 of *Lecture Notes in Computer Science*, pages 142–166. Springer, 1996.
- [3] D. F. Bacon and P. F. Sweeney. Fast static analysis of C++ virtual function calls. In *Conference Proceedings: Object-Oriented Programming Systems, Languages, and Applications (OOPSLA)*, pages 324–341, 1996.
- [4] Á. Beszédes, R. Ferenc, T. Gyimóthy, A. Dolenc, and K. Karsisto. Survey of code-size reduction methods. *ACM Computing Surveys*, 35(3):223–267, 2003.
- [5] D. Chanut, B. De Sutter, B. De Bus, L. Van Put, and K. De Bosschere. System-wide compaction and specialization of the Linux kernel. In *Proc. 2005 ACM SIGPLAN/SIGBED Conference on Languages, Compilers, and Tools for Embedded Systems (LCTES’05)*, pages 95–104, June 2005.
- [6] S. Debray, W. Evans, R. Muth, and B. De Sutter. Compiler techniques for code compaction. *ACM Transactions on Programming Languages and Systems*, 22(2):378–415, March 2000.
- [7] R. Flower, C.-K. Luk, R. Muth, H. Patil, J. Shakshober, R. Cohn, and P. G. Lowney. Kernel optimizations and prefetch with the Spike executable optimizer. In *Proc. 4th Workshop on Feedback-Directed and Dynamic Optimization (FDDO-4)*, December 2001.
- [8] D. Grove, G. DeFouw, J. Dean, and C. Chambers. Call graph construction in object-oriented languages. In *Proc. Conference on Object-Oriented Programming Systems, Languages & Applications (OOPSLA ’97)*, pages 108–124, 1997.
- [9] M. Guthaus, J. Ringenberg, D. Ernst, T. Austin, T. Mudge, and T. Brown. MiBench: A free, commercially representative embedded benchmark suite. pages 3–14, December 2001.

- [10] M. Hind and A. Pioli. Which pointer analysis should I use? In *ISSTA '00: Proceedings of the 2000 ACM SIGSOFT International Symposium on Software Testing and Analysis*, pages 113–123, 2000.
- [11] Intel Corp. *IA-32 Intel Architecture Software Developer's Manual, Volume 1: Basic Architecture*.
- [12] N. D. Jones, C. K. Gomard, and P. Sestoft. *Partial Evaluation and Automatic Program Generation*. Prentice Hall, 1993.
- [13] D. McNamee, J. Walpole, C. Pu, C. Cowan, C. Krasic, A. Goel, , and P. Wagle. Specialization tools and techniques for systematic optimization of system software. *ACM Trans. on Computer Systems*, 19(2):217–251, May 2001.
- [14] A. Milanova, A. Rountev, and B. G. Ryder. Precise call graphs for C programs with function pointers. *Automated Software Engineering*, 11(1):7–26, 2004.
- [15] C. Pu, H. Massalin, and J. Ioannidis. The Synthesis kernel. *Computing Systems*, 1(1):11–32, 1988.
- [16] C. Pu *et al.* Optimistic incremental specialization: Streamlining a commercial operating system. In *Proc. 15th ACM Symposium on Operating Systems Principles (SOSP'95)*, pages 314–324, Dec 1995.
- [17] F. Tip and J. Palsberg. Scalable propagation-based call graph construction algorithms. In *Proc. Conference on Object-Oriented Programming, Systems, Languages and Application (OOPSLA-00)*, pages 281–293, October 2000.
- [18] S. Zhang. *Practical Pointer Aliasing Analyses for C*. PhD thesis, 1998.
- [19] S. Zhang, B. G. Ryder, and W. Landi. Program decomposition for pointer aliasing: A step toward practical analyses. In *Proc. Fourth ACM SIGSOFT Symposium on the Foundations of Software Engineering*, pages 81–92, October 1996.