

An Optimal Memory Allocation Scheme for Scratch-Pad-Based Embedded Systems

OREN AVISSAR and RAJEEV BARUA

University of Maryland, College Park

and

DAVE STEWART

Embedded Research Solutions

This article presents a technique for the efficient compiler management of software-exposed heterogeneous memory. In many lower-end embedded chips, often used in microcontrollers and DSP processors, heterogeneous memory units such as scratch-pad SRAM, internal DRAM, external DRAM, and ROM are visible directly to the software, without automatic management by a hardware caching mechanism. Instead, the memory units are mapped to different portions of the address space. Caches are avoided due to their cost and power consumption, and because they make it difficult to guarantee real-time performance. For this important class of embedded chips, the allocation of data to different memory units to maximize performance is the responsibility of the software.

Current practice typically leaves it to the programmer to partition the data among different memory units. We present a compiler strategy that automatically partitions the data among the memory units. We show that this strategy is optimal, relative to the profile run, among all static partitions for global and stack data. For the first time, our allocation scheme for stacks distributes the stack among multiple memory units. For global and stack data, the scheme is provably equal to or better than any other compiler scheme or set of programmer annotations. Results from our benchmarks show a 44.2% reduction in runtime from using our distributed stack strategy vs. using a unified stack, and a further 11.8% reduction in runtime from using a linear optimization strategy for allocation vs. a simpler greedy strategy; both in the case of the SRAM size being 20% of the total data size. For some programs, less than 5% of data in SRAM achieves a similar speedup.

Categories and Subject Descriptors: Computing Systems [**Embedded Systems**]: Memory Management

General Terms: Design, Measurement, Performance

Additional Key Words and Phrases: Memory, heterogeneous, allocation, storage, embedded

1. INTRODUCTION

This article presents an automatic compiler method for allocating program data among different heterogeneous memory units in embedded systems. The kind

Authors' addresses: O. Avissar and R. Barua, Dept. of Electrical & Computing Engineering, University of Maryland, College Park, MD 20742; email: {oavissar,barua}@eng.umd.edu; D. Stewart, Embedded Research Solutions, LLC, 9687F, Gerwig Lane, Columbia, MD 21046; email: dstewart@embedded-zone.com.

Permission to make digital/hard copy of all or part of this material without fee for personal or classroom use provided that the copies are not made or distributed for profit or commercial advantage, the ACM copyright/server notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists requires prior specific permission and/or a fee.

© 2002 ACM 1539-9087/02/0011-0006 \$5.00

of embedded chips targeted are those without caches, and with at least two kinds of writable memory—only in such systems is intelligent memory allocation useful. Such embedded chips without caches and with multiple memory units are an important class of embedded chips, common in many lower-end embedded devices such as microcontrollers and some DSP chips. The writable memory units typically found in such systems include any two or more of internal SRAM, external SRAM, integrated DRAM, external DRAM, and even EEPROM that are writable, but with very high latency. In such chips, each of the memory units may have differing latencies and sizes, therefore the choice of memory allocation affects the overall performance.

The allocation strategy presented in this article models the problem using a 0/1 integer linear program, and solves it using commercially available software [Matlab 2001]. The formulation is provably optimal, relative to the profile run, for global and stack data. It is easily adapted to both nonpreemptive and preemptive (context-switching) scenarios. For the first time ever, the solution automatically distributes the program stack among multiple memory banks, effectively growing the stack simultaneously in several places. Stack distribution is unusual—presently programmer annotations and compiler methods place the entire stack in one memory unit. The resulting flexibility allows a more custom allocation for better performance. Finally, the optimality guarantee ensures that the solution is equal to or better than any programmer-specified allocation using annotations, or any existing or future compiler methods.

The method presented is motivated by a need to improve the quality of automatically compiled code. Compilers today, while better than before, still suffer from a large performance penalty compared to programs directly written in assembly language [Bhattacharyya et al. 2000; Paulin et al. 1997]. This forces many performance-critical kernels to be written directly in assembly. Assembly programming has well-known disadvantages: more tedious, expensive, and error-prone code development, difficulty in porting between different platforms, and a longer time-to-market between successive implementations [Bhattacharyya et al. 2000; Paulin et al. 1997]. Further, optimizations that benefit from whole-program analysis, such as memory allocation, cannot be captured by rewriting certain kernels.

One of the major remaining impediments to efficient compilation is the presence of multiple, possibly heterogeneous, memory units mapped to different portions of the address space. Many low-end embedded processors [Motorola MCore 1998; Texas Instruments TMS370Cx 1997; Motorola 68HC12 2000] have such memory units as on-chip scratch-pad SRAM, on-chip DRAM, off-chip DRAM, and ROM that are *not* members of a unified memory hierarchy. Caches are not used for reasons of real-time constraints, cost, and power dissipation. In contrast, the memory units in desktop processors are unified through caches. Directly addressed memory units in embedded processors require the software to allocate the data to different memories, a complex task for which good strategies are not available.

This work proposes a method for automatically allocating program data among the heterogeneous memory units in embedded processors without

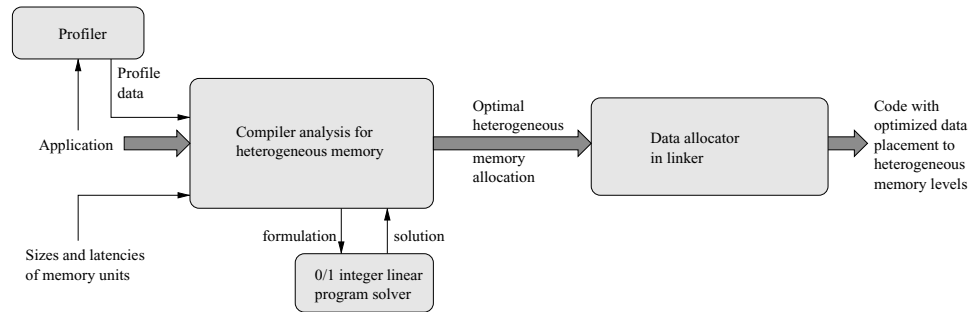


Fig. 1. Flow diagram for heterogeneous memory management on embedded systems. The thick arrows denote the primary flow of compilation and synthesis, while the thin arrows provide supporting data.

caches. The standard practice today is that the allocation is left to the programmer. An automated solution for a different kind of embedded processor has been proposed [Panda et al. 2000], namely, that for which the external memory has an internal hardware cache. We show, however, that the optimality criteria are very different without caches. Other related works are discussed later in the related work section.

In typical embedded processors without caches, a small, fast SRAM memory called *scratch pad* replaces the cache, but its allocation is under software control. An important recent article [Banakar et al. 2002] studied the trade-offs of a cache vs. scratch pad. Their results were startling: a scratch-pad memory has 34% smaller area and 40% lower power consumption than a cache memory of the same capacity. Even more surprising, the runtime measured in cycles was 18% better with a scratch pad when using a simple knapsack-based allocation algorithm [Banakar et al. 2002]. With the superior allocation schemes proposed here, the runtime improvement will be even larger. Thus, defying conventional wisdom, they found absolutely no advantage in using a cache, even in high-end systems in which performance is important. In conclusion, given the power, cost, performance, and real-time advantages of scratch pad, we expect that systems without caches will continue to dominate embedded systems in the future.

Figure 1 outlines our method for heterogeneous memory management on embedded systems. The application program on the left is fed to our compiler analysis that derives the optimal static allocation for the data given the inputted profile data. The compiler analysis incorporates application-specific runtime information through the use of profile data. The analysis is provided with the sizes and latencies of the memory units available on the target chip, as shown. The compiler analysis models the problem for global and stack data as a 0/1 integer linear programming problem and solves it using Matlab [Matlab 6.1 2001]. The solution is always provably optimal, relative to the profiled run of the program, among static methods for handling global and stack data. As depicted, the derived allocation specification is output to the linking stage of compilation. The linker adds assembly directives to its output code to implement the desired allocation. The resulting code not only improves performance, but is likely to

reduce energy consumption—it is well-known that compiler optimizations that reduce runtime usually reduce energy consumption as well [Lee et al. 1997].

This article is organized as follows. Section 2 motivates the problem and outlines our approach. Section 3 describes related work; Section 4 shows a simple example to illustrate the trade-offs involved. Section 5 describes our method for global variables. Section 6 shows how the formulation can be extended for stack variables. Section 7 presents some preliminary results. Section 8 addresses certain real-world issues; while Section 9 concludes.

2. MOTIVATION AND APPROACH

Multiple heterogeneous memory units in many embedded systems are motivated by the varying performance, cost, and writability characteristics of different memory technologies. Typically, such chips contain some or all of the following: a small amount of on-chip SRAM, a moderate amount of on-chip DRAM, a moderate amount of off-chip SRAM, a large amount of off-chip DRAM, and some ROM to store programs and data constants. Each kind has its advantages: SRAM is fast but expensive; off-chip SRAM is somewhat slower; integrated on-chip DRAM is slower than SRAM but faster than external DRAM; whereas external DRAM is slow, but is the cheapest. ROM is cheap and nonvolatile, but it cannot be written to. Certain kinds of ROM, such as EEPROM, are writable with high latency.

The organization of the memories is different from that in desktop systems. Although desktops also contain many of these memories, they automatically manage them hierarchically using caches. Caches, however, consume area and power, and make it difficult to provide real-time guarantees. Consequently, many embedded chips, except for some at the high end, do not use caches. Examples include the Motorola 68HC12 [2000], Motorola MCore [1998] and Texas Instruments TMS370Cx [1997]. The lack of caches has meant that the different memories are mapped to different nonoverlapping portions of the address space.

A result of the lack of caches is that the allocation of data to memories must be software-managed—in most systems today, it is left to the programmer. This work presents a compiler method to manage the data. Compiler methods are preferable to programmer directives as they do not require programmer effort; are portable across different systems; and are likely to make better decisions, especially for large, complex programs. The dominance of chips without caches in embedded systems implies that good compiler methods for data allocation will have a large impact.

Our profile-guided compiler method is static: it fixes the allocation at compile-time; data is not relocated to another location during runtime. Alternate strategies could be dynamic—software-managed caches [Moritz et al. 2000; Hallnor and Reinhardt 2000] are one class of dynamic algorithms. There are, however, no software caching strategies available today that are optimized for the class of embedded processors we target. Dynamic strategies are not studied in this work—future work might study such schemes. Significant challenges will need to be overcome in designing a software-caching scheme for embedded

chips, including providing real-time guarantees in the face of unpredictable cache behavior; reducing software caching overhead; and restricting code size increase from the overhead instructions. Fortunately, profile data allows our static method to incorporate runtime information to some extent.

Even for static allocations, several factors make the problem a difficult one to solve optimally. Let us consider global, stack, and heap variables. We present a formulation that is optimal, relative to the profile run, for global variables among static partitions. Stack variables, however, unlike global variables, have limited lifetimes, allowing sharing of space between variables with disjoint lifetimes. This complicates the analysis, but we present a method that is able to retain the optimality guarantee for nonrecursive procedures. For heap data, the situation is worse: no static method can be optimal for all heap data, as the sizes and allocation frequencies are unknown for heap data. Allocation of heap data is not considered any further in this work. In our scheme, all heap data is allocated to external DRAM, which is the same as the default in today's compilers. Better schemes will be investigated in future work.

3. RELATED WORK

For embedded processors with heterogeneous memory units that are not cache-managed, there has been little related work that automatically allocates data to banks while increasing performance. The usual approach is to leave the task to the programmer. Compiler methods are preferable to programmer directives for three reasons: they do not require programmer effort; are portable across different systems; and are likely to make better decisions, especially for large, complex programs.

To our knowledge, Panda et al. [2000], Sjodin et al. [1998], and Sjodin and von Platen [2001] are the only published methods that aim to allocate data to on-chip and off-chip memories mapped to different portions of the address space. However, the architecture class targeted by Panda et al. [2000] is different from ours, they target embedded processors that, in addition to having scratch-pad SRAM, use hardware-managed caches on top of slower, external memory.¹ The presence of caches completely changes the goals of the allocation strategy. Instead of aiming to reduce the number of accesses to data in external memory, it becomes far more important to ensure that those accesses hit in cache. Consequently, the goal of the method in Panda [2000] is to map the variables that are likely to cause the most conflicts to scratch pad. It makes no attempt to maximize the number of accesses to scratch pad, and thus is unsuitable for our architectural model.

The method proposed by Sjodin et al. [1998] also differs from ours in several ways. Like our method, Sjodin et al. [1998] also utilize an allocation scheme that tries to keep variables with the highest number of accesses per byte in on-chip SRAM and allocate the less critical variables to slower, external RAM. However, unlike our formulation they only address two memory levels (on-chip SRAM and external RAM) and do not offer any methods for extending to allocate stack

¹Our architectural model without caches is employed by a variety of embedded processors, as caches consume area and power, and make it difficult to provide real-time guarantees.

variables. Our formulation automatically handles N levels of memory with varying latencies and sizes. Also, we have extended the formulation to account for stack (local) variables. Another difference is that Sjodin et al. primarily use a static profiling scheme, which can be inaccurate—especially for larger programs. Our scheme always uses dynamic profiling which accounts for every load and store throughout the program.

Recent work by Sjodin and von Platen [2001] uses a linear formulation such as ours. The key differences of our method from theirs are as follows. First, they do not incorporate a distributed stack—this accounts for 44% performance gain in our results. Also, they have no optimized analysis of stack variables. The goals of their work differ from ours in that they aim to optimize the use of native pointer types to improve performance and reduce code size.

As described earlier, our method yields a static allocation; dynamic strategies are possible. In a dynamic strategy, data may be moved from one location to another during program execution. One class of dynamic strategies are software-caching methods [Moritz et al. 2000; Hallnor and Reinhardt 2000]—these emulate a cache in fast memory using software. The tag, data, and valid bits are all managed by compiler-inserted software at each program data access. Software overhead is incurred to manage these fields, although the [Moritz et al. 2000] compiler optimizes away the overhead in some cases. Moritz et al. [2000] target the primary cache, while Hallnor and Reinhardt [2000] intend to manage the secondary cache and assume different costs—for this reason the method proposed by Moritz et al. [2000] is more applicable for the problem we consider.

Dynamic strategies are not studied in this work—future work might study such schemes. Significant challenges will need to be overcome in designing a software-caching scheme for embedded chips, including providing real-time guarantees in the face of unpredictable cache behavior; reducing software caching overhead; and restricting code size increase from the overhead instructions. Fortunately, profile data allows our static method to incorporate runtime information to some extent.

4. EXAMPLE

Figure 2 shows a simple example that illustrates how a programmer or compiler might make decisions about data allocation. Figure 2(a) is the application program for which we wish to allocate the data. Two byte arrays $A[100]$ and $B[1000]$ are passed as arguments to procedure foo where they are accessed using formal arguments $x[]$ and $y[]$. Assume that the program is compiled for an embedded chip that has 1 KB of fast scratch-pad SRAM, no on-chip DRAM and 8 KB of slower, external DRAM. The problem we are trying to solve is: which program variables should be allocated to which memory bank? It is clear that either array can individually fit in the 1 KB SRAM, but that both cannot simultaneously fit. For simplicity of illustration, assume that no other accesses to A and B occur, although there is no such requirement in our method.

For the code in Figure 2(a), the compiler needs to choose between the two possible data allocations shown in Figures 2(b) and (c). In Figure 2(b) $A[100]$

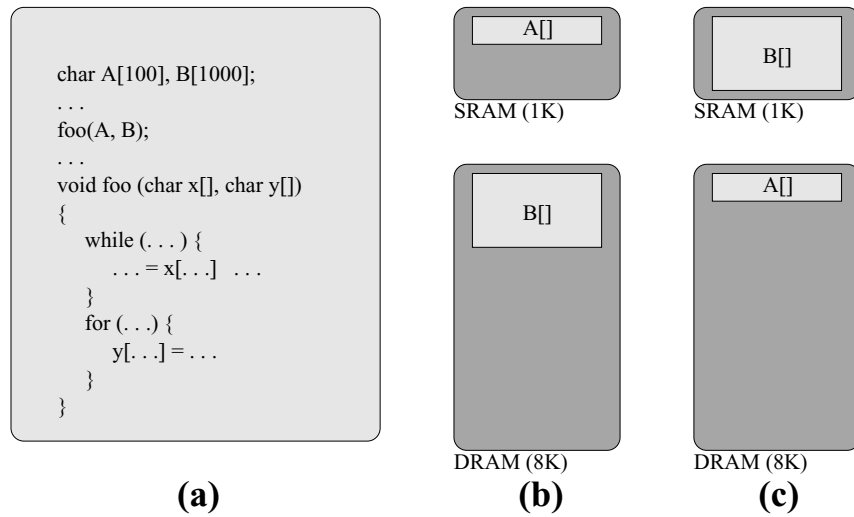


Fig. 2. Example showing choices in heterogeneous memory allocation: (a) the source code of the application; (b) $A[]$ allocated to SRAM (allocated portions are shaded lighter); (c) $B[]$ allocated to SRAM.

is in SRAM; in Figure 2(c) $B[1000]$ is in SRAM. The choice between the two depends on access frequencies. Two cases illustrate the choice. In the first case, suppose the *while* loop at runtime actually executes for more iterations than the *for* loop, then the allocation in Figure 2(b) is preferred, as it makes more accesses to faster SRAM. In the second case, suppose the *for* loop executes for more iterations instead. In this case, the allocation in Figure 2(c) is superior, as it makes more SRAM accesses. However, making estimates of relative access frequencies is difficult, since many loops have bounds that are unknown at compile; further, data-dependent control flow makes static prediction difficult. Fortunately, profile data gives good estimates provided the data set used is representative, and is far more accurate than static frequency prediction methods. Our allocation method uses profile data to find access frequencies of memory references.

The above example illustrates that for making good allocation decisions, the compiler must integrate at least three technologies. First, the general problem of optimal data allocation, which is NP-complete, must be solved either exactly or approximately using heuristics. We propose a method that returns an optimal static solution using an integer programming framework. There is some evidence that 0/1 integer programming has fast solution times (under a minute on modern computers) even for large programs with a few thousand variables [Appel and George 2001; New York City, Office of Budget and Management 1999]. Fortunately, the number of variables in our formulation is proportionate to the number of variables in the original program, which is usually no more than a few thousand for even large programs. Quick solution times are borne out by our results, where in all cases the solution was returned in under a minute. Second, the compiler must collect accurate frequency estimates using profiling,

and use them to guide allocation. Third, to collect frequency counts for variables, the profiler must correlate accesses with the variables they access. For example, in Figure 2(a), it must be known that $x[]$ is $A[]$ and $y[]$ is $B[]$. Knowing this statically requires computationally expensive interprocedural pointer analysis. Moreover, pointer analysis may return ambiguous data if, for example, there is more than one call to $foo()$ with different array arguments each time. For this reason, we avoid pointer analysis by taking a different approach—runtime address checks during profiling. During the profile run, each accessed address is checked against a table of address ranges for the different variables. A match indicates that the variable accessed has been found. This profile-based approach yields exact statistics, unlike the inexact information using pointer analysis.

5. FORMULATION FOR GLOBAL VARIABLES

Here we present the formulation for global variables; it is extended to handle stack variables later. The following symbols are used:²

- U = Number of heterogeneous memory units;
- T_{rj} = Time (latency) to read memory unit $j \in [1, U]$ in cycles;
- T_{wj} = Time (latency) to write memory unit $j \in [1, U]$ in cycles;
- M_j = Size of memory unit $j \in [1, U]$ in bytes;
- G = Number of global variables in application;
- v_i = i th global variable, $i \in [1, G]$;
- $N_r(v_i)$ = Number of times v_i is read (from profiling);
- $N_w(v_i)$ = Number of times v_i is written (from profiling);
- $S(v_i)$ = Size of variable v_i in bytes.

The optimization problem is formulated as a 0/1 integer linear program. The following set of 0/1 integer variables ($\forall j \in [1, U], \forall i \in [1, G]$) is defined:

$$I_j(v_i) = \begin{cases} 1 & \text{if variable } v_i \text{ is allocated on memory unit } j \\ 0 & \text{otherwise.} \end{cases}$$

The objective function to be minimized is the total access time of all the memory accesses in the application. For architectures allowing at most one memory access per cycle, the total time is

$$\sum_{j=1}^U \sum_{i=1}^G I_j(v_i) [T_{rj} N_r(v_i) + T_{wj} N_w(v_i)]. \quad (1)$$

It is easy to see how the above is the total time for all memory accesses. The term $T_{rj} N_r(v_i)$ is the time for all the read accesses to variable v_i , if it were allocated to memory unit j . A similar term is added for all the write accesses. When multiplied by the 0/1 variable $I_j(v_i)$ the result contributes the memory

²The T_{rj}, T_{wj} values used for DRAMs are averages. Some modern DRAMs have slightly lower latencies for sequential accesses compared to nonsequential accesses, by about 20%. The difference is not modeled.

access time if v_i were indeed allocated to unit j ; zero otherwise. Summing this term over all variables (the inner sigma) yields the total access time for a given memory unit. The outer sigma yields the total across all memory units.

For machines that allow at most one memory access per cycle, the formula in (1) is accurate. Most low-end embedded processors we target allow only one memory access per cycle, including some Very Long Instruction Word (VLIW) architectures,³ and hence the formula in (1) is accurate for most targeted chips. For higher-end VLIWs that allow more than one memory access per cycle, however, (1) does not take into account the overlap of memory latencies. To do so requires the formula to include the maximum of latencies of different memory accesses in the same cycle; unfortunately, the objective function does not remain linear, since the maximum function is not linear. Thus, heuristics instead of optimal 0/1 integer linear solvers must be used; these are not evaluated in this work.

As in any linear optimization problem, a set of constraints is also defined. The first is an exclusion constraint that enforces that, for every application variable v_i , it is allocated on only one memory unit:

$$\sum_{j=1}^U I_j(v_i) = 1 \quad (\forall i \in [1, G]). \quad (2)$$

Another constraint is that the sum of the sizes of all variables allocated to a particular memory unit must not exceed the size of that memory unit:

$$\sum_{i=1}^G I_j(v_i) * S(v_i) \leq M_j \quad (\forall j \in [1, U]). \quad (3)$$

The objective function (1) combined with the constraints (2) and (3) define the optimization problem. The function and constraints are then solved with an available mathematical solver; we use Matlab [2001]. The resulting values of $I_j(v_i)$ are the optimal static allocation.

6. EXTENSION TO STACK VARIABLES

For good performance, stack variables—procedure parameters, local variables, and return variables—must be distributed among the different heterogeneous memory units to achieve a more custom allocation. Stack distribution, however, is a complex objective, since the stack is normally a sequentially allocated abstraction. Normally, the stack grows in units of stack frames, one per procedure, where a stack frame is a block of contiguous memory locations containing all the variables and parameters of the procedure. The stack grows and shrinks sequentially in units of frames for every nested procedure call with a procedure. Consequent to this sequential abstraction, the entire stack is placed in one memory unit in all programmer-annotated strategies and automatic allocation methods existing today of which we are aware. Custom strategies in which

³VLIWs are architectures that allow the compiler to schedule multiple instructions per cycle, though usually not all of them may be memory instructions.

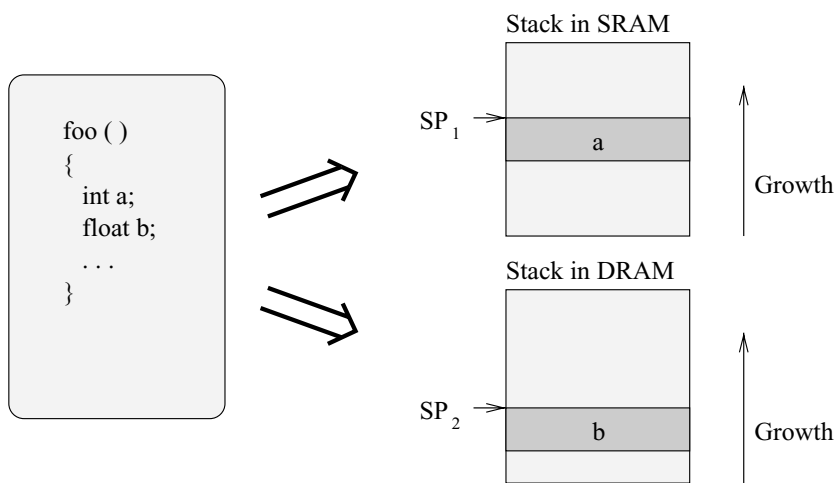


Fig. 3. Example of stack split into two separate memory units. Variables a and b are placed on SRAM and DRAM, respectively. A call to $foo()$ requires the stack pointers in both memories to be incremented.

frequently used stack variables are allocated to fast memory, and others to slow memory, have not been investigated.

This article presents a strategy for distributed stacks applied to heterogeneous memory allocation. Distributed stacks were proposed for the first time in an earlier work by one of the authors of Barua et al. [2001] for a different purpose for *homogeneous* memory units. For the first time, this article adapts distributed stacks for heterogeneous memory units. Our approach is best explained through an example. On its left, Figure 3 shows an example code fragment containing a procedure $foo()$ with two local variables a and b . On the right is shown how the stack is distributed into two memory units: variable a is allocated to SRAM, and b to DRAM. This results in two stack pointers, SP_1 and SP_2 , in SRAM and DRAM, respectively. Both stack pointers must be incremented upon procedure entry and decremented upon exit, instead of one. For ease of use, the implementation of distributed stacks is done automatically by the compiler; the sequential view of an undivided stack is retained for the programmer.

Distributed stacks as in Figure 3 incur some software overhead: multiple stack pointers must be updated upon procedure entry and exit, instead of one. Two solutions to overcome the overhead are presented below in Alternatives 1 and 2. Alternative 1 is to eliminate the overhead by forcing all the variables within one procedure to be allocated to the same memory level—this way only one stack pointer needs to be updated per procedure. The stack is still divided, however, as different stack frames can still be allocated to different memory units, and multiple stack pointers exist. However, eliminating the overhead that has a price, grouping together variables, results in loss of allocation flexibility. A second solution, presented in Alternative 2 below, is to tolerate the overhead, and distribute each individual variable within a stack frame to potentially different banks, as in Figure 3. The best solution

is to use a hybrid of Alternatives 1 and 2: selectively tolerate the overhead for long-running procedures. For long-running procedures, identified by profiling, the impact of a few overhead instructions will be insignificant, so for such procedures the overhead is tolerated. For short-running procedures, each stack frame is allocated to one memory unit, and the overhead is eliminated.

Modification of global formulation. To see how to modify the global variable formulation for stack variables, consider that the fundamental difference between the two is the limited lifetimes of stack variables. Although stack variables for nonrecursive procedures can be treated just as global variables by allocating them for all time, the resulting allocation would not be optimal for stack variables. The reason is that performance with stack variables can be further improved by taking advantage of their limited lifetimes. Stack variables are allocated upon procedure entry and freed upon exit. Thus, constraint (3) is no longer valid: the total size of variables allocated to a bank may exceed its size provided not all of them are live at the same time. Variables with nonoverlapping lifetimes may share the same space in memory.

One way to incorporate stack variables into our formulation is to treat each stack variable just as a global variable, with the modification that the maximum size constraint (3) is relaxed somewhat. Instead of requiring that all variables fit simultaneously in memory, the call graph of the program is analyzed to construct a new set of constraints that require that only variables that can be live simultaneously fit in memory. The intuition is that one new constraint is introduced for each unique path in the call graph from *main()* to each leaf node in the call graph.⁴ The details follow directly from the intuition and are presented below for the two alternative methods 1 and 2, both with their merits.

6.1 Alternative 1: Distribution Granularity = Stack Frames

In this first alternative, the stack for each procedure is combined into a single aggregate variable in the formulation. This ensures that the full stack frame for a procedure is allocated to a single memory unit, leading to simplicity in formulation and implementation, and no software overhead for updating multiple stack pointers. The stack is still distributed, since different frames could be allocated to different memory units. To describe this formulation, the following symbols are introduced in addition to the ones before for globals:

F = Number of aggregate stack variables in the application program
(Number of functions);

f_i = i th function, $i \in [1, F]$;

$NP(f_i)$ = Total number of unique paths to the function f_i in the call graph;

$P_j(f_i)$ = The j th unique path in the call graph to f_i , $j \in [1, NP(f_i)]$;

\mathcal{L} = The set of all leaf nodes in the call graph.

⁴Cycles in the call graph (recursion) cannot be handled this way. Instead they are collapsed to a single aggregate variable in the formulation before this step, and assigned a maximum allowable size based on recursive depth.

In addition, $N_r(f_i)$, $N_w(f_i)$, $S(f_i)$, $I_j(f_i)$ are defined as the number of reads to, number of writes to, size, and 0/1 variable for the stack frame for f_i , in an analogous manner to $N_r(v_i)$, $N_w(v_i)$, $S(v_i)$, $I_j(v_i)$. The solution for the $I_j(f_i)$ values yields the desired allocation.

Similar to the formulation for global variables only, the objective function is the total time for all memory accesses (in this case global and stack variables). The objective function for the stack extended formulation is

$$\begin{aligned} & \sum_{j=1}^U \sum_{i=1}^G I_j(v_i) [T_{rj} N_r(v_i) + T_{wj} N_w(v_i)] \\ & + \sum_{j=1}^U \sum_{i=1}^F I_j(f_i) [T_{rj} N_r(f_i) + T_{wj} N_w(f_i)]. \end{aligned} \quad (4)$$

The first term in the above is the original objective function (1) for the global variables, which represents the total time needed to access the global variables. The second term is the total time needed to access the stack variables.

Regarding constraints, the exclusion constraint for global variables presented earlier in (2) is still needed unchanged. A similar constraint is added for stack variables:

$$\sum_{j=1}^U I_j(f_i) = 1 \quad (\forall i \in [1, F]). \quad (5)$$

As previously mentioned, changes are needed to account for the fact that stack variables can have disjoint lifetimes. Substantial changes are made to the memory size constraint (3) to accommodate the limited lifetimes of stack variables. The new memory size constraint is

$$\begin{aligned} & \forall j \in [1, U], \forall f_l \in \mathcal{L}, \forall t \in [1, NP(f_l)], : \\ & \sum_{i=1}^G I_j(v_i) S(v_i) + \sum_{\forall f_p \in P_t(f_l)} I_j(f_p) S(f_p) \leq M_j. \end{aligned} \quad (6)$$

The first line of the above states that the second line (the constraint) is replicated for all combinations of memory banks (j), leaf nodes (f_l) in the call graph, and paths to that leaf node (t). The constraint in the second line states that the global variables plus all the stack variables in the given path to the given leaf node must fit into memory. The first term represents the global variable size; the second term represents the size of the stack variables for every call-graph path to a leaf function. The stack is of a maximal size when a call-graph leaf is reached; consequently, ensuring that all paths to leaf nodes fit in memory ensures that the program allocation will fit in memory at all times.

The set of constraints in (6) is large, as it is replicated across all j , f_l , and t . Fortunately, however, this is not expected to adversely impact the runtime of the 0/1 solver by much, as the runtime for such solvers depends more on the number of 0/1 variables and less on the number of constraints. Indeed, more constraints may decrease the runtime by decreasing the space of feasible solutions.

Table I. Information About the Benchmark Programs

Benchmark	Source	Total Data Size	Runtime (cycles)	Description
BMM	Trimaran	120204 bytes	4.29 MB	Block matrix multiplication
BTOA	Rutter et al.	475 bytes	62 KB	Changes 8 bit byte streams into ASCII
CRC32	MiBench 1.0	2478 bytes	64 KB	32 BIT ANSI X3.66 CRC checksum files
DIJKSTRA	MiBench 1.0	1908 bytes	1.04 MB	Finds shortest paths using Dijkstra's alg
FFT	UTDSP	4196 bytes	471 KB	256-point complex FFT
FIR	Trimaran	1088 bytes	262 KB	Finite impulse response algorithm
IIR	UTDSP	1140 bytes	102 KB	4-cascaded IIR biquad filter (64 pts)
LATNRM	UTDSP	1116 bytes	432 KB	32nd-order Norm Lattice filter (64 pts)

6.2 Alternative 2: Distribution Granularity = Stack Variables

In this alternative, stack variables from the same procedure are allowed to be allocated to different memory units. The formulation is modified as follows. The stack variables are treated just as global variables in the formulation, leading to an objective function similar to (1). The exclusion constraint is similar to (2). The memory size constraint is similar to (6) with the second \sum function converted to a $\sum \sum$, the outer summation remaining the same as the second summation in (6), and the inner summation summing across all the individual variables in the procedure f_p .

7. RESULTS

7.1 Global and Stack Variables

7.1.1 Benchmark and Simulation Environment. Our formulation for global and stack variables has been implemented in the public-domain GCC cross-compiler set to target the Motorola M-Core [1998] embedded processor. A collection of small programs, BMM, BTOA, CRC32, DIJKSTRA, FFT, FIR, IIR, and LATNRM have been compiled and evaluated. Their characteristics are shown in Table I. The benchmarks FIR and BMM were obtained from the trimaran [Consortium 1999] benchmark suit; BTOA was obtained from Rutter et al.; CRC32 and DIJKSTRA were obtained from Guthaus et al. [2001]; and FFT, IIR, and LATNRM were obtained from [UTDSP 1992]. These eight benchmarks represent code that would be used in typical applications. The first benchmark, BMM, which has 7 functions and 6 global variables, creates and multiplies two matrices and sums up all the elements of the resulting matrices. The second benchmark, BTOA, which has 4 functions and 7 global variables, is a stream filter to change 8 bit bytes into printable ASCII characters. The third benchmark, CRC32, which has 2 functions and 3 global variables, creates the ANSI standard cyclic redundancy check checksum for a data set. The fourth benchmark, DIJKSTRA, which has 5 functions and 11 global variables, finds the

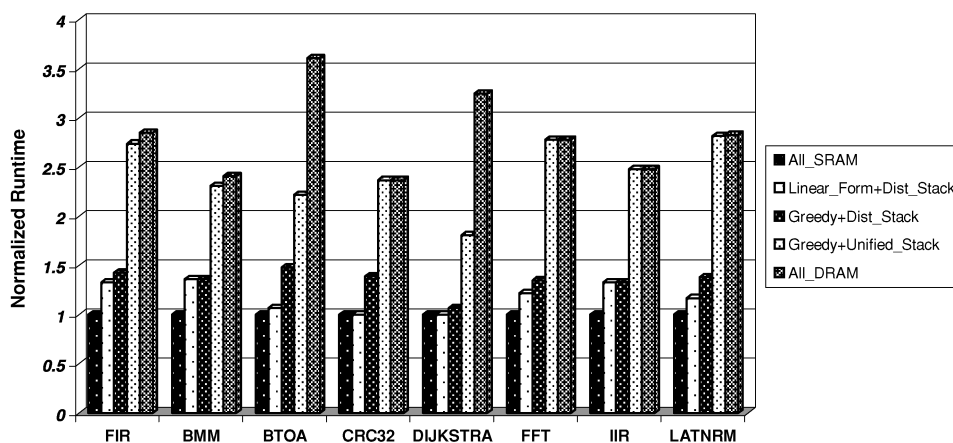


Fig. 4. Normalized simulated runtimes for all benchmarks with varying memory configurations. Each group of bars displays the cases, for all benchmarks, when the SRAM size is equal to the program data size, 20% of the program data size (for all three allocation schemes), and the case when all program data is placed in DRAM.

shortest path between two elements in a matrix using Dijkstra’s algorithm. The fifth benchmark, FFT, which has 2 functions and 5 global variables, performs a fast Fourier transform on a 256 point data set. The sixth benchmark, FIR, which has 2 functions and 4 global variables, is an implementation of the finite impulse filter algorithm. The seventh benchmark, IIR, which has 2 functions and 5 global variables, is an implementation of an infinite impulse response filter algorithm. The eighth benchmark, LATNRM, which has 2 functions and 5 global variables, is a 32nd-order normalized lattice filter that processes a 64 point data set.

Runtimes of the benchmarks are obtained using the M-Core simulator available from Motorola. The M-Core chip simulated has three levels of memory: a 256 KB EEPROM with 1-cycle read and 500-cycle write latency; a 256 KB external DRAM with 10-cycle read/write latency; and an internal SRAM with 1-cycle read/write latency. In the experiments below we analyze the effect of varying the SRAM size while providing ample DRAM and EEPROM memory. In these experiments only stack data (local variables) and global data were accounted for. The Alternative 1 strategy from Section 6.1 is used for stack variables. The benchmarks selected did not use the heap; our current scheme does not optimize for heaps. Relatively small benchmarks, which do not make calls to C standard libraries, were selected to avoid the complexity of compiling/optimizing many additional libraries. The runtime of the 0/1 optimizer was never a problem: it rarely took more than a few seconds to run.

7.1.2 Comparison with Other Methods. The first experiment performed is to compare the runtimes of the benchmarks for five significant cases. The results of this experiment can be seen in Figure 4. First, *ALL_SRAM* simulates the runtimes for all the benchmarks for the case when all the program data is allocated to internal SRAM—the baseline case (normalized runtime = 1.0).

Second, *Linear_Form + Dist_Stack* simulates the runtimes for the case when the linear formulation is used for allocation, along with distributed stacks, with an SRAM size that is 20% of the total program data obtained. Third, *Greedy + Dist_Stack* is the case when a greedy allocation with a distributed stack is used, when the SRAM size is 20% of the total program data size. Fourth, *Greedy + Unified_Stack* is the case when a greedy allocation, with a stack allocated entirely to external DRAM, is used when the SRAM size is 20% of the total program data size. Fifth, *ALL_DRAM* simulates the runtimes of the benchmarks when all the program data is allocated to external DRAM.

We now explain the three nontrivial allocation schemes above—*Linear_Form + Dist_Stack*, *Greedy + Dist_Stack*, and *Greedy + Unified_Stack*. To achieve performance improvement, *Linear_Form + Dist_Stack*, our best method, uses profile data to place frequently accessed data words in fast memory, and other data in slower memory. In most benchmarks, a large share of the memory accesses goes to a small fraction of the data [Hennessy and Patterson 1996]. These accesses can then be placed in fast memory. For all programs the *Linear_Form + Dist_Stack* is optimal among all possible static methods for globals and stacks—poor results only stem from program characteristics, not any deficiencies in the formulation. The *Greedy + Dist_Stack* optimization orders all the global and stack variables by size and then tries to allocate the smallest variables to the the fastest memory. This method does not use profiling and represents a strategy that would typically be used in the absence of sophisticated methods. By allocating the smallest variables to the fastest memory banks, scalars and other small variables are placed in fast memory. This generally increases the probability that fast memory will be used. In the *Greedy + Unified_Stack* optimization the greedy method described above is only applied to global variables and the entire stack is placed in external DRAM. This method represents a situation when a unified stack is used, as is typically the case.

There are three salient conclusions that can be derived from Figure 4. First, a distributed stack is significantly better than a single unified stack. This can be seen by comparing the bars for *Greedy + Dist_Stack* and *Greedy + Unified_Stack*. By comparing these bars we see that for every benchmark there is a distinct performance improvement when distributing the stack. On average there is a 44.2% reduction in normalized runtime with the distributed stack. In some cases the normalized runtime is more than halved! Second, the linear formulation is better than the greedy allocation. This is seen by comparing the bars for *Linear_Form + Dist_Stack* and *Greedy + Dist_Stack*. When comparing these cases we see that for many of the benchmarks there is a significant improvement when using the linear formulation. On average the *Linear_Form + Dist_Stack* produced a 11.8% improvement in the normalized runtime over the *Greedy + Dist_Stack*. Third, when keeping just 20% of the data in SRAM, the optimizations are able to deliver performance that is closer to that of the *ALL_SRAM* case than the *ALL_DRAM* case. This demonstrates the overall success of the method in reducing the SRAM size required for good performance to well below the *ALL_SRAM* case.

7.1.3 Variance with SRAM Size. The second experiment performed is to plot the runtimes of the benchmarks when the SRAM size is varied while the EEPROM and DRAM sizes are fixed to an ample size. The results of this experiment can be seen in Figure 5. The three curves in each graph are *Linear_Form + Dist_Stack*, *Greedy + Dist_Stack*, and *Greedy + Unified_Stack* as described above. All the runtimes in Figure 5 are normalized to the 100% SRAM case. The SRAM size is varied from 0% to 100% of the total data size of the program. The data sizes of the BMM, BTOA, CRC32, DIJKSTRA, FFT, FIR, IIR, and LATNRM benchmarks are 120204, 475, 2478, 1908, 4196, 1088, 1140, and 1116 bytes, respectively. The CRC32, DIJKSTRA, FIR, IIR, and LATNRM benchmarks are composed of a roughly equal number of scalars and arrays; in the BMM and FFT benchmarks much of the memory is used by a few large arrays; the BTOA benchmark is mostly scalars with one array.

In the plots of Figure 5 we see that the runtimes increase as the size of the SRAM decreases, as one would expect. The effectiveness of the *Linear_Form + Dist_Stack* can clearly be seen by the large jumps in runtime that occur as the SRAM size gets close to 0%. These jumps happen when the most-often-used variables, that are kept in SRAM as much as possible, are forced into slower memory banks. In Figure 5 we see that *Linear_Form + Dist_Stack* curve is always below the other curves. The *Linear_Form + Dist_Stack* guarantees that the runtimes are statically optimal for each memory size. Figure 5 also demonstrates the drastic improvements achieved by distributing the stack. For every benchmark we see that the performance is remarkably worse for *Greedy + Unified_Stack* than *Linear_Form + Dist_Stack*, and *Greedy + Dist_Stack*. The most noticeable difference is the starting point of the *Greedy + Unified_Stack* curve. Because the entire stack is placed in slow memory, even when there is ample SRAM, all the stack data accesses are slow.

Closer examination of Figure 5 reveals some more detailed observations. Three significant trends can be seen—the BMM numbers are used as an example. First, utilizing SRAM can lead to a large gain compared to using only DRAM (runtime 1 vs. 2.41), showing the importance of carefully allocating data to SRAM. Second, the formulation is able to get a fairly good runtime (1.36) for even a very small memory size of 32 bytes, compared to the 0 byte case (2.33), showing that the profile-guided optimization is successfully able to put heavily accessed data in SRAM. In the 32 byte case, just 0.03% of the total data is allocated to SRAM, yet it is able to achieve a 41.6% improvement (2.33 vs. 1.36) in runtime! Third, the formulation is successfully able to utilize EEPROM even for data that is written, when SRAM is not available. This is seen for the SRAM = 0 case—the runtime reduces to 2.33 (DRAM + EEPROM) compared to 2.41 (only DRAM). This 3.32% benefit is because the high cost of EEPROM writes may be recovered by frequent reads for data with infrequent writes—in the CRC32 benchmark there is a 41.3% improvement between the SRAM = 0 and *ALL_DRAM* case (2.37 vs. 1.39).

Several benchmark-specific characteristics can also be seen in Figure 5. In the BMM, FFT, LATNRM, and IIR curves the first jumps in runtime occur in large SRAM intervals because these programs have many large arrays. The points in these plots indicate where each of the large arrays get pushed out of

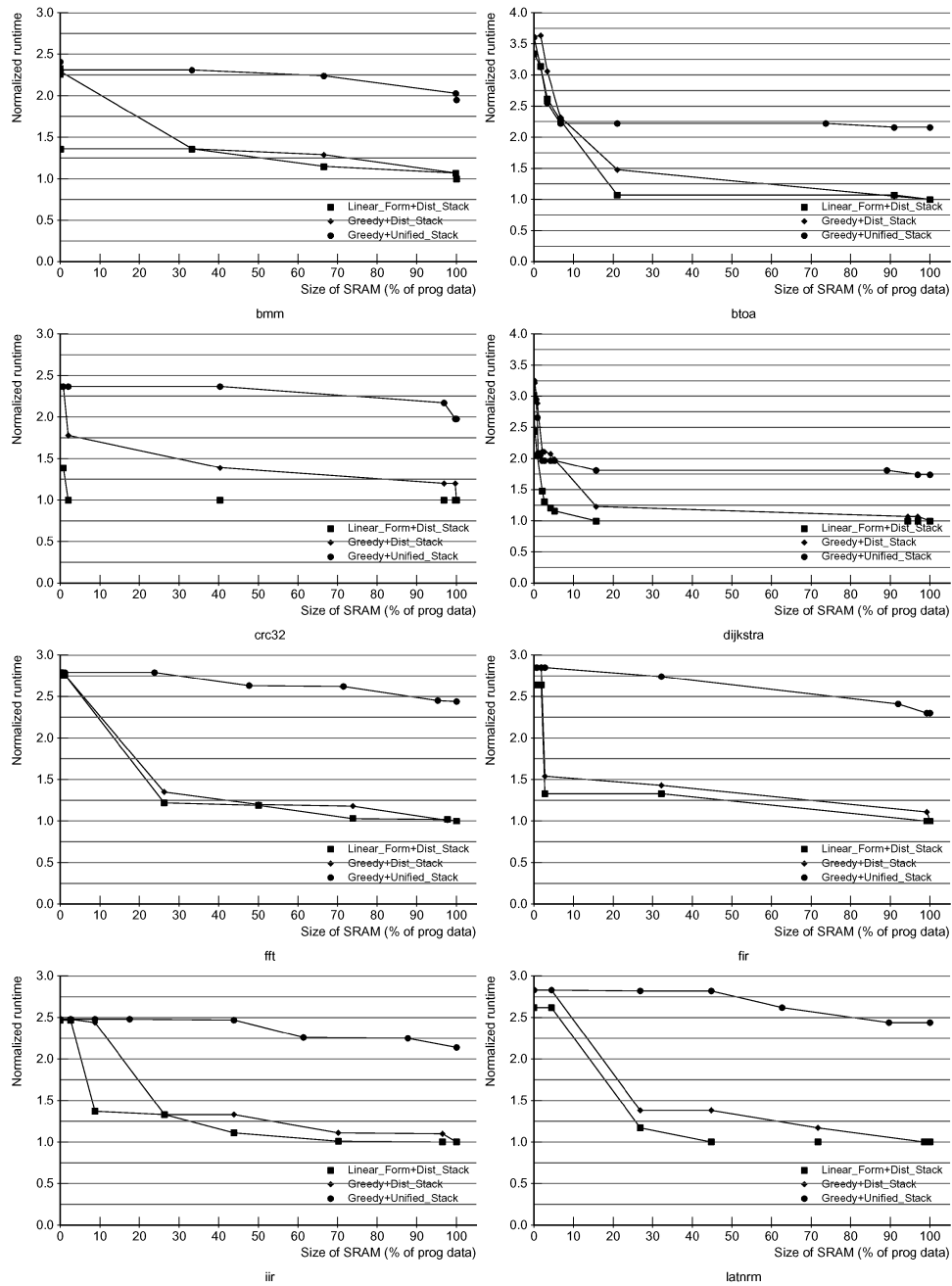


Fig. 5. Normalized runtimes for benchmarks with varying SRAM size. DRAM and EEPROM sizes are fixed. X-axis = SRAM size, as a percentage of the total data size for that program. Y-axis = runtime, normalized to 1.0 for SRAM size = 100% of data size. Note the steep jump in runtime as the SRAM size approaches zero.

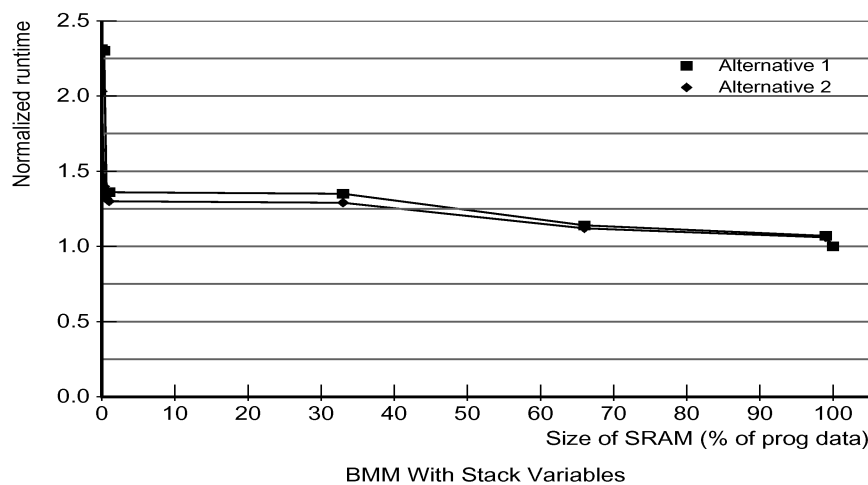


Fig. 6. Normalized runtimes of Alternative 1 and Alternative 2 stack formulations for the BMM benchmark with varying SRAM size.

SRAM—and allocated to slower DRAM. The final runtime jumps then occur at very small SRAM sizes—almost 0% SRAM—as the heavily used scalar variables are pushed out of SRAM. This effect can also be seen to a lesser degree in the FIR and CRC32 plots. The DIJKSTRA and BTOA benchmarks contain mostly scalar and small array variables with only a few infrequently accessed large arrays. As a result, there is not a substantial increase in runtime until the SRAM size gets close to 0%. As the SRAM size is reduced, the lesser-used variables are first allocated to slower memory banks, while the often-used variables are kept in fast memory as long as possible.

7.2 Comparison of Alternatives 1 and 2

To see the impact of using Alternative 1 vs. Alternative 2 we measured the runtime of the BMM benchmark, using both alternatives, when the SRAM size is varied while the EEPROM and DRAM sizes are fixed to an ample size. The data for Alternative 1 was generated using the formulation described in Section 6.1. The data for Alternative 2 was simulated by changing all stack variables to global variables. The side effect of using this approach is that the overhead incurred by the additional stack manipulation is not accounted for.

Figure 6 summarizes the results of the above experiment. The BMM benchmark was chosen because it contains numerous functions along with several global variables. From the figure we see that Alternative 1 performed slightly worse than Alternative 2, which was expected. However, the figure indicates that Alternative 1 may be a feasible strategy. For the larger SRAM sizes the runtimes are virtually the same while for smaller sizes the performance separation becomes more apparent. The average performance difference between the two alternatives, for SRAM sizes greater than 1% of the program data size, is 2.5%. In theory, Alternative 2 will always perform better than Alternative 1

because it provides a finer granularity in the data objects. The finer granularity creates more flexibility for the allocation algorithm, which increases its effectiveness.

8. REAL-WORLD ISSUES

8.1 Application to Synthesis

Our compiler method can be used as a synthesis tool to determine the minimum SRAM size needed by the chip during its assembly by a particular user. Using a smaller SRAM can result in significant cost savings. To obtain the smallest SRAM size that still meets users' performance requirements, the application domain is recompiled for ever-larger SRAM sizes, and the simulated runtime vs. SRAM size is plotted, as in Figure 5. The smallest size that delivers the performance needed is chosen.

8.2 Adaption to Preemptive (Context Switching) Environments

The formulation for automatic data allocation is easily extended to preemptive systems that context-switch between multiple programs in a given workload. In context-switching, data from all the programs must be simultaneously present somewhere in memory. If only one program uses all of SRAM, however, the performance of the other programs will suffer. The only way to get better performance is to partition the SRAM among the different programs.⁵ In context-switching environments, our framework can partition among different programs in an elegant manner by combining their variables and solving together, after weighing (multiplying) the frequencies $N_r(v_i)$ and $N_w(v_i)$ of each variable by the relative frequency of that context. The solution is statically optimal, relative to the profile run, for the group of programs when run together.

8.3 Initial Values and Volatile Variables

For most embedded systems, globals and static local variables with initial values must be specially handled in the formulation. When a global or static local variable has an initial value, and in addition, is possibly written to in the program for some input data, then it must be allocated to RAM, but its initial value is stored in ROM. Therefore memory must be allocated in both RAM and ROM for such a variable. This is easily handled in the formulation by preallocating the variables in ROM and running the rest of the formulation unchanged with a reduced (remaining) size of ROM.

Another special consideration must be made for variables declared as volatile. Such variables may be changed by entities outside of the program, such as I/O devices or another processor. Under normal conditions, the formulation presented in this article can occasionally make use of ROM for variables that are seldom, or never, written to. If a variable is volatile, the variable might

⁵The alternative solution is to allocate the SRAM to one program only at a time, and save the entire SRAM contents to DRAM on a context-switch just as register files. This is infeasible as SRAMs are usually much larger than register files, making the context-switch overhead unacceptable.

actually be written to often by some outside source. Hence, volatile variables must never be placed in ROM. This is handled in the formulation by presetting the 0/1 variable corresponding to the volatile variable being allocated to ROM, to zero.

9. CONCLUSION

This article presents a compiler method to distribute program data among the different memory units of embedded processors without caching hardware. Without caching, the task of allocating data to the different banks falls to the software. The compiler method derives a static partition, that is, one in which the allocation of data to memory units is fixed and unchanging throughout program execution. Among static methods, the method presented is optimal, relative to the profile run, for global and stack data. For the first time, our method distributes stacks among various heterogeneous memory units, resulting in a more custom allocation. The optimality guarantee ensures that the method presented will be as good or better than any programmer-derived annotations or competing static compiler technique. The method models the problem as a 0/1 integer linear programming problem and solves it using available commercial packages.

We are encouraged by the results obtained. They show the benefits of optimal allocation: with just 20% of the data in SRAM, our method is able to decrease the runtime by 56% on average for our benchmarks vs. allocating all data to slow memory, without any programmer involvement. For some programs, less than 5% of data in SRAM achieves a similar speedup. Further, results from our benchmarks show a 44.2% reduction in runtime from using our distributed stack strategy vs. using a unified stack, and an additional 11.8% reduction in runtime from using a linear optimization strategy for allocation vs. a simpler greedy strategy; both in the case of the SRAM size being 20% of the total data size. The variance in results is due to application characteristics. Applications that access some data more frequently than other data see large decreases in runtime when those data are allocated to fast memory—the amount of runtime reduction depends on the relative frequencies of access.

REFERENCES

- APPEL, A. AND GEORGE, L. 2001. Optimal spilling for CISC machines with few registers. In *Proceedings of the SIGPLAN '01 Conference on Program Language Design and Implementation* (Snowbird, UT, June).
- BANAKAR, R., STEINKE, S., LEE, B.-S., BALAKRISHNAN, M., AND MARWEDEL, P. 2002. Scratchpad memory: A design alternative for cache on-chip memory in embedded systems. In *10th International Symposium on Hardware/Software Codesign (CODES)* (Estes Park, CO, May 6–8). ACM, New York.
- BARUA, R., LEE, W., AMARASINGHE, S., AND AGARWAL, A. 2001. Compiler support for scalable and efficient memory systems. *IEEE Trans. Comput., Special Issue on Advances in High Performance Memory Systems* (Nov.).
- BHATTACHARYYA, S. S., LEUPERS, R., AND MARWEDEL, P. 2000. Software synthesis and code generation for signal processing systems. *IEEE Trans. Circuits Syst.* 47, 9 (Sept.).
- CONSORTIUM, T. T. 1999. The Trimaran benchmark suite. Available at <http://www.trimaran.org/>.

- GUTHAUS, M. R., RINGENBERG, J. S., ERNST, D., AUSTIN, T. M., MUDGE, T., AND BROWN, R. B. 2001. MiBench: A free, commercially representative embedded benchmark suite. Available at <http://www.eecs.umich.edu/jringenb/mibench/>.
- HALLNOR, G., AND REINHARDT, S. K. 2000. A fully associative software-managed cache design. In *Proceedings of the 27th International Symposium on Computer Architecture (ISCA)* (Vancouver, B.C., Canada, June).
- HENNESSY, J. AND PATTERSON, D. 1996. *Computer Architecture A Quantitative Approach*. 2nd ed. Morgan-Kaufmann, Palo Alto, Calif.
- LEE, T.-C., TIWARI, V., MALIK, S., AND FUJITA, M. 1997. Power analysis and minimization techniques for embedded DSP software. *IEEE Trans. VLSI Syst.* (Mar.).
- Matlab 6.1*. The Math Works, Inc., 2001. <http://www.mathworks.com/products/matlab/>.
- MORITZ, C. A., FRANK, M., AND AMARASINGHE, S. 2001. FlexCache: A framework for flexible compiler generated data caching. In *The 2d Workshop on Intelligent Memory Systems* (Boston, MA, Nov. 12).
- CPU12 Reference Manual*. Motorola Corporation, 2000. http://e-www.motorola.com/brdata/PDFDB/MICROCONTROLLERS/16_BIT/68HC12_FAMILY/REF_MAT/CPU12RM.pdf.
- M-CORE—MMC2001 Reference Manual*. Motorola Corporation, 1998. <http://www.motorola.com/SPS/MCORE/info-documentation.htm>.
- New York City, Office of Budget and Management. 1999. Website on frequently asked questions on linear programming. <http://www.eden.rutgers.edu/~pil/FAQ.html>. New York, NY.
- University of Toronto Digital Signal Processing (UTDSP). 1992. University of Toronto Digital Signal Processing (UTDSP) Benchmark Suite. Available at <http://www.eecg.toronto.edu/>.
- PANDA, P. R., DUTT, N. D., AND NICOLAU, A. 2000. On-chip vs. off-chip memory: The data partitioning problem in embedded processor-based systems. *ACM Trans. Des. Autom. Electron. Syst.* 5, 3 (July).
- PAULIN, P., LIEM, C., CORNERO, M., NACABAL, F., AND GOOSSENS, G. 1997. Embedded software in real-time signal processing systems: Application and architecture trends. *Invited paper, Proc. IEEE* 85, 3 (Mar.).
- RUTTER, P., OROST, J., AND GLOISTEIN, D. BTOA: Binary to printable ASCII converter source code. Available at <http://www.bookcase.com/library/software/msdos.devel.lang.c.html>.
- SJODIN, J., FRODERBERG, B., AND LINDGREN, T. 1998. Allocation of global data objects in on-chip RAM. *Compiler and Architecture Support for Embedded Computing Systems*. Dec.
- SJODIN, J. AND VON PLATEN, C. 2001. Storage allocation for embedded processors. 2001 *Compiler and Architecture Support for Embedded Computing Systems*. Nov.
- TMS370Cx7x 8-bit microcontroller*. Texas Instruments, Revised Feb. 1997. <http://www-s.ti.com/sc/psheets/spns034c/spns034c.pdf>.

Received January 2002; accepted July 2002