
An Analysis of Scalar Memory Accesses in Embedded and Multimedia Systems*

Osman S. Unsal¹, Zhenlin Wang², Israel Koren¹,
C. Mani Krishna¹, Csaba Andras Moritz¹

¹ Department of Electrical and Computer Engineering,

² Department of Computer Science,
University of Massachusetts, Amherst, MA 01003

Summary. In an earlier paper about the FlexCache project [24], we described our vision of a multipartitioned cache where memory accesses are separated based on their static predictability and memory footprint, and managed with various compiler controlled techniques supported by instruction set architecture extensions, or with traditional hardware control.

In line with that vision, this paper describes our work in progress related to the memory performance and memory management of scalars. Our focus in this paper is embedded and multimedia architectures, but the methodology described can be applied to other classes of applications.

In particular, we establish the minimum size of a memory partition that would allow us to map and manage all scalar accesses in a program statically, and describe compiler techniques to automate the extraction of this information. We evaluate the impact of register file size on the volume of scalar related memory accesses, and its impact on the applications' overall cache performance. We study the cache behavior of scalar accesses for embedded architectures, including reduction in cache misses due to separation of scalars from other types of memory accesses. Additionally, we develop an energy-efficient data caching strategy for multimedia processors, based on our scalar partitioning approach.

1 Introduction and Motivation

The recent proliferation of palmtops, MP3 players and internet-enabled wireless phones has ignited interest in embedded and multimedia systems. These systems have to be fast and energy efficient. As such, they have tight memory/processing requirements. Therefore, understanding their memory/caching behavior is of paramount importance.

In an earlier paper about the FlexCache project [24], we described our vision of a multipartitioned cache where memory accesses are separated based

* Supported in part by NSF grant EIA-0102696.

on their static predictability and memory footprint, and managed with various compiler controlled techniques. This paper addresses the cache behavior of scalar accesses and its memory footprint in order to enable a fully static memory management in a logical memory partition.

Although prior studies into memory behavior of arrays for embedded and multimedia systems have been conducted, the study of the memory footprint of scalars has lagged behind. Here we report our ongoing work in closing this gap. This paper presents techniques and results for scalar memory accesses in embedded and multimedia systems. Our preliminary results show promise and we hope that this work will heighten interest in this area.

This research spans compiler and architectural domains. Our particular contributions in this paper are threefold:

- First, we experimentally establish the memory size requirements of scalars for embedded and multimedia systems. We present a new compiler algorithm to automatically extract this information, as would be required in a multipartitioned cache.
- Second, by separating scalar accesses from array accesses, we expect decreased cache interference and improved static predictability. This aspect is especially important for hard deadline embedded systems.
- Third, we study the energy implications of partitioning the scalars from non-scalars in media processors. In particular, we compare the energy consumption of a regular data cache with a multipartioned one in which scalars are exclusively assigned to *scratchpad* memory.

The rest of this paper is organized as follows. In Section 2, we provide a brief literature survey and reiterate our motivation. Section 3 describes the experimental setup, we include baseline cases for both embedded and media processors. Section 4 provides the results, we consider separate case studies for embedded (Section 4.4) and media (Section 4.5) processors. In Section 5 we conclude with a brief summary and a synopsis of future work.

2 Previous Work

This work builds upon the framework in [23, 24]. Previous memory behavior research effort primarily targeted array structures [17, 26]. Delazuz et al. [12] discuss energy-directed compiler optimizations for array data structures on partitioned memory architectures; they use the SUIF compiler framework for their analysis. On the other hand, architectural support to improve memory behavior include split caches which were discussed in [22]. Albonesi [2] proposed selective cache ways, a vertical cache partitioning scheme. Benini et al. [4] discuss an optimal SRAM partitioning scheme for an embedded system-on-a-chip. Kin et al. [16] study a small L0 cache that saves energy while reducing performance by 21%. Lee and Tyson [19] use the mediabench benchmarks and have a coarse-granularity partitioning scheme: they opt for dividing the cache

along OS regions for energy reduction. Chiou et al. [10] employ a software-controlled cache and use a cache way based partitioning scheme. A recent paper by Huang et al. [15] also uses a way-prediction scheme; their cache partitioning includes a specialized stack cache and compiler implementation concerns are addressed. Mueller [25] sketched some broad ideas on compiler support for cache partitioning. Combined compiler/architectural efforts toward increasing cache locality [21] have also exclusively focused on arrays. For multimedia systems, one previous work has considered reconfigurable caches [28], using the recently introduced Mediabench benchmark in the performance analysis, with comments on compiler controlled memory. Burlin [9] concentrates on optimizing stack frame layout in embedded systems. Cooper and Harvey [11] look at compiler-controlled memory. Their analysis includes spill memory requirements for some Spec '89 and Spec '95 applications. Engblom [13] and Lee et al. [18] discuss why Spec is not a suitable benchmark for embedded systems.

The above research, although preoccupied primarily with memory behavior of arrays, provided valuable pointers for our work. In this paper we consider scalar memory accesses, not only array or spill memory accesses, and we target embedded systems running a suite of media applications. We develop a compiler heuristic to calculate the memory requirements of scalars and discuss the impact of architectural design choices for embedded and multimedia systems on scalars.

3 Experimental Setup

We use the recently developed Mediabench benchmarks [18] in our experiments. Mediabench is a collection of popular embedded applications for communications and multimedia. We chose Mediabench, since other benchmarks such as SPECint, DSPstone or Dhrystone are not suitable for embedded or multimedia systems [6, 13].

We needed a detailed compiler framework that would give us sufficient feedback, is easy to understand, and allow us to change the source code for our modifications. With this in mind, we chose the SUIF/Machsuif suite as our compiler framework. SUIF [29] does high-level passes while Machsuif [20] makes machine specific optimizations. Our main focus is Machsuif's register allocator pass, Raga. Raga makes the transition from virtual registers into real registers and performs register allocation. The allocation uses a graph coloring heuristic to assign registers to temporaries. We have made modifications to Raga to annotate scalar memory accesses. The resulting annotated assembler code targets the Alpha processor. We have amended the assembler code by inserting NOP instructions around the scalar memory operations, thus *marking* them. The scalar memory accesses consist of spills and register promotion related memory accesses.

We used the SimpleScalar tool suite [7] to run the Alpha binaries and collect the results. We have modified SimpleScalar to recognize the scalar memory operations in the *marked* code. SimpleScalar was modified to squash the marker instructions on fetch, therefore the marker instructions do not impact the results in any way. Our baseline machine model is a single-issue in-order processor. Lee et. al. [19] use an identical SimpleScalar configuration in their power dissipation analysis of region-based caches for embedded processors. Most embedded processors employ an in-order microarchitecture. Using an out-of-order, non-blocking load type of microarchitecture would, to some degree, decrease the performance penalty of scalar/non-scalar conflict cache misses. We use the Wattch [8] tool suite to run the binaries and collect the energy results. Wattch is built on top of the SimpleScalar framework. We use the activity sensitive conditional clocking power model in Wattch, i.e., the cache consumes power when it is accessed.

For embedded processors, we did a survey of cache sizes to determine the baseline. As Table 1 indicates, embedded processor data cache sizes are usually small. Therefore, we have selected a data cache size of 2K for our experiments with embedded systems.

Table 1. Data cache sizes for typical embedded CPUs. SRAM scratchpad areas are available in the Samsung ARM7, Hitachi SH2 and Fujitsu SparcLite.

Processor	Cache Size	Processor	Cache Size
Samsung ARM7	2K	SparcLite	2K
PA-RISC HP	1K to 2K	Power PC 403GA	1K
Hitachi SH-II	4K unified	Coldfire 5102	1K
Embedded Pentium	8K	MIPS Jade	1 to 8K
Sandcraft SR-1-GX	8K		

On the other hand, as Table 2 indicates, the trend is towards larger caches for media processors. Therefore, for media processors we have selected a 64Kbyte 2-way cache as our baseline. The table also indicates that media processors do not typically have L2 data caches. Therefore, we only have L1 caches in our baseline architecture. However, our framework is applicable to media processors with L2 caches as well. In this case, one issue that must be addressed is the consistency between the L2 cache and the L1 data + scratchpad. Namely, the block fetched from L2 into the L1 caches could contain a mix of scalar/non-scalar data. We avoid this problem by keeping the block sizes the same across the caches. If the block sizes were different, then the issue could be addressed by clustering the scalar data to the beginning of the address space and padding them appropriately to the size of the L2 cache block size and boundary.

Table 2. Cache configurations for typical media processors.

Processor	L1 Cache	L2 Cache
ARM ARM10	32K	None
Transmeta Crusoe TM3200	32K	None
Transmeta Crusoe TM5400	64K	256K
Intel StrongARM SA-110	16K	None
Equator Map-CA	32K	None
Intel StrongArm 110	16K	None
Intel StrongARM 1100	8K	None

4 Results

4.1 Motivational Example

We start with a motivational example. Consider the sample program in Figure 1. The program consists of x scalar variables being written in a chain-dependent fashion, after which a single array element is written per loop iteration. We define the scalar miss ratio to be the ratio of scalar misses to total misses. Consider the scalar miss ratio for 32 scalars which is 34%. When we increase the number of scalars to 64 the ratio increases to 46%, although the memory footprint of 32 additional scalars is small. This points to the fact that interference between the scalar and array accesses are chiefly responsible for the increase in the scalar miss ratio. Therefore, if we can separate the array accesses from the scalar accesses, this ratio and the overall miss rate will decrease. Next, we present our results based on Mediabench applications for embedded systems.

<pre> main() { loop { scalar_accesses; array_access[]; } }</pre>		Number of scalar temporaries		
		32	64	96
	Scalar Miss Ratio	0.34	0.46	0.47
		Number of scalar temporaries		
		32	64	96
Scalar Footprint%		1.55	3.10	4.65

Fig. 1. Scalar misses for the synthetic example. Here the integer array is of size 2048, and the columns denote the number of scalar variables in the example. Scalar operations are of the form: $Variable_{n+1} = Variable_n \mp constant$. There is a single array access per loop iteration. The loop is iterated 100000 times. The cache is 2K-direct mapped.

4.2 Memory Size

We use two yardsticks for experimental evaluation of the scalar memory size requirements of media applications. The first of these is the static memory evaluation. It is static in the sense that the results were extracted by a compile-time analysis of assembler code. We isolated the scalar memory operations in every routine. We then determined the granularity of data by instruction analysis, i.e., the granularity is 8 if the move is a quadword instruction, 2 if it is a word instruction, and so on. We then identified the unique scalar accesses by counting multiple accesses into the same memory location only once and by taking the maximum of the pertaining granularities. The results given in the first column of Table 3 indicate that memory size requirements are modest.

Table 3. The Memory Size Requirements.

(In Bytes)	Static Dynamic		(In Bytes)	Static Dynamic	
ADPCM	0	0	EPIC	321	203
G721 Encode	48	32	GSM	202	146
JPEG Encode	502	83	MPEG Encode	2125	604
PEGWIT	98	16	RASTA	618	152
PGP	394	358	MESA	2191	770

However, the static estimate is pessimistic since not all of the data space is traversed during execution. We therefore, developed a second yardstick, a dynamic memory evaluation which provides a tighter, more robust bound. We recompiled the Mediabench benchmarks to record runtime routine use information. We executed each benchmark with its default input set and extracted the *dynamic* call-tree information by using the *gprof* profiling utility. Then, for every routine we noted the scalar memory requirements as in the static technique. Traversing the tree from the root to each leaf, adding up the unique scalar accesses from each routine, and finding the critical path, i.e., the path with the maximum size requirement, yields the result. We supply the *dynamic* call-tree for the EPIC application in Figure 2 as an example of this process. The memory requirements thus obtained are shown in the second column of Table 3. The results suggest that the memory footprint of scalars in media applications for embedded systems is quite small. These results will guide the choice of our architectural optimization schemes. We next present our compiler technique to automate the scalar memory size estimation.

Intuitively, the upper bound of the size of the scalar buffer is the maximum of the distinct scalars along all program execution paths. An algorithm that accurately calculates this bound needs inter-procedural analysis and a complex data-flow analysis. Here we present a good approximation. Our algorithm conservatively assumes that the scalars along all paths are distinct. It simply adds the number of bytes needed for each scalar. Of concern here are loops in

call node, *resolved*, *scalarBound*, and *localScalarSize*. The *localScalarSize* of a basic block is the total number of bytes for all scalars in the basic block. The *localScalarSize* of a routine is its scalar buffer upper bound without taking routine calls into account. The *scalarBound* of a basic block is the scalar bound along all simple paths from the entry block to the current block in the control flow graph. The *scalarBound* of a routine is the scalar bound along all simple paths from the *main* routine to the current routine in the call graph. We say that a basic block or a routine is *resolved* when its *scalarBound* is known. Assuming there are N routines in a program and the maximal number of basic blocks of a routine is M , then the complexity of the algorithms is $O(NM^2 + N^2)$.

Algorithm 1 Find Routine Scalar Memory Requirement from CFG

Require: localScalarSize of block

```

/* Phase 1 */
/* For each routine, traverse its control flow graph */
for each routine do
  calculate scalar bound for each basic block;
  mark back edges in CFG;
  /* Add the entry back block to workList */
  E = entry basic block;
  E.scalarBound = E.localScalarSize;
  E.resolved = true;
  workList = successors of E;
  while !empty(workList) do
    B = next element in workList;
    allResolved = true;
    maxBound = 0;
    /* check if all B's predecessors are resolved */
    for each predecessor P of B do
      if the edge (P,B) is not marked and P is not resolved then
        allResolved = false;
        break;
      else
        maxBound = max(maxBound, P.scalarBound);
      end if
    end for
    if allResolved then
      remove B from workList;
      B.resolved = true;
      B.scalarBound = maxBound + B.localScalarSize;
      add all unresolved successors of B to workList;
    end if
  end while
  set localScalarSize of the current routine as scalarBound of its exit block.
end for

```

4.3 Register File Size

For memory analysis of arrays, optimizing the cache is more important than the register file architecture, since array accesses seldom use registers. However, for scalars the situation is different. The register file size can have a direct impact on spills and thus impact performance. We therefore, analyze the impact of register file sizes on scalars. We take a two-step approach: first, we do a survey of the register file sizes for *current* embedded CPUs and use these results to drive our experiment. Second, we gauge the impact of expanded register file sizes in *future* embedded processors.

Table 4 shows the register file sizes on some typical embedded CPUs: the size ranges between 16 and 32 except for the embedded Pentium which has 8 general-purpose registers. We therefore, varied the register file size from 16 to 32 in our experiments. We modified Machsuf passes and architectural definitions to output binaries for different register file sizes. Then, we noted the static number of scalars inserted into the instruction stream. The results in Figure 3 show that there is a considerable number of scalar memory accesses for 16 registers. Another point is that for some particular benchmarks (e.g., JPEG Encode), the number of scalar memory accesses is more dramatically decreased than others as more registers become available. This is because the register pressure is more unevenly distributed in those benchmarks, i.e., only a few routines exhibit intense register pressure. Once those are relieved through additional registers, the decrease in scalar register spills is more steep.

Table 4. Integer Register File Sizes in Current Embedded CPU's.

Processor	Register File Size	Processor	Register File Size
Samsung ARM7	15	SparcLite	32
PA-RISC HP	16+16	Power PC 403GA	32
Hitachi SH-II	16	Coldfire 5102	16
Embedded Pentium	8	MIPS Jade	32
Sandcraft SR-1-GX	32		

Usually, the current methods and techniques used in general microprocessors migrate to embedded systems with a couple of years time lag. We believe that the integer register file sizes will follow the same trend. Therefore, we project the embedded CPU integer register file size to grow to 64, 128 and maybe 256. We extended our analysis by modifying Machsuf to output code for larger register file sizes. The results are shown in Figure 4. Note that for large register file sizes all the register spills are eliminated, the only remaining scalar memory operations are register promotion related; this is the reason for the flattening out of the scalar memory accesses. The implication is that, as far as scalars in media applications are concerned, increasing the register file size will not bring any additional benefits. This is especially true for register

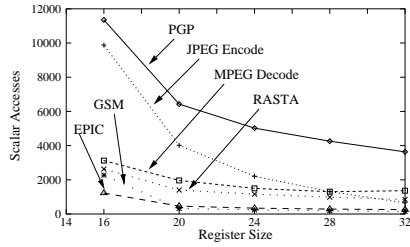


Fig. 3. Number of Scalar Memory Accesses With Register File Size

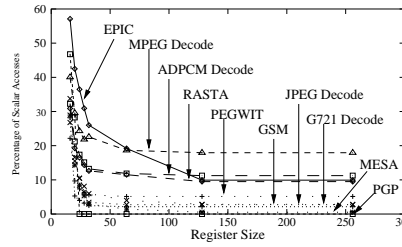


Fig. 4. Percentage of Scalar Memory Accesses for Extended Register File Size

file sizes larger than 64. Thus, we experimentally establish what has been an industry insight with general-purpose CPUs [5]. This provides a guidance to the designers of future embedded/multimedia CPUs: the number of integer general-purpose registers should be at most 64, the additional chip real estate could be devoted to other functional units (e.g., caches) that offer better incremental performance.

Our analysis also includes a cross-architectural comparison as seen in Figure 5. We used an Intel X86-family targeted version of SimpleScalar for this analysis. The 8-register X86 has significantly more scalar memory accesses than the 32-register Alpha. This is due to Machsuif’s register allocator, Raga. Raga is based on a graph coloring heuristic and as argued in [1], register allocation based on graph coloring is sensitive to the number of registers, in particular when the number of available registers is low. Here, we experimentally verify this argument. *Therefore, compiler designers for embedded CPUs, which typically have fewer registers, should develop new register allocation heuristics.* Work in this direction has already started [1]. We also comment on an important property leading to a dual conclusion. Sometimes, increasing the register file size can increase scalar memory accesses. This may seem counterintuitive at first. However, consider Figure 6 for the MPEG benchmark. As the register size is increased from 28 to 32, the number of scalar memory accesses actually increases. This is due to the graph-coloring heuristic used in Raga to assign registers. The use of this heuristic creates a phenomenon similar to the Belady anomaly in paging [3]. The conclusion that can be drawn: *embedded CPU designers should be aware of the characteristics of their target compiler in choosing their design point.* In summary, the above experiments show that the compiler/architecture coupling in embedded systems is stronger than previously assumed and should be considered at the design phase.

4.4 Case Study: Scalar Data Remapping for Embedded Processors

We assume that the reorganization and separation of scalar and array accesses are compiler level tasks. As mentioned in Section 2, there are several cache

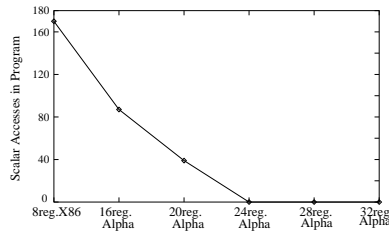


Fig. 5. Scalar accesses for Intel-X86 and Alpha

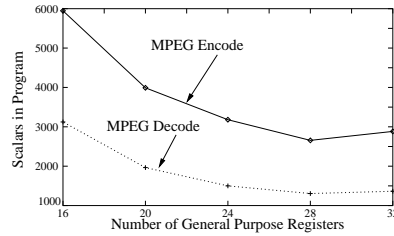


Fig. 6. Effect of register coloring heuristic

reorganization options for scalars: vertical or horizontal partitioning [2], which partition the cache along cache ways and lines, respectively. Another option is to use a scratchpad SRAM area and direct the scalar memory accesses to this partition. Here, an appropriate partitioning option must be selected. Vertical partitioning schemes are wasteful: the existing cache has to be divided into a power of two, and our results indicate that the memory footprint of the scalars in embedded media applications is small. Instead, we advocate the use of a scratchpad SRAM area. Separate SRAMs are widely used in DSP's: they are typically used to hold frequently used data such as floating-point constants. A scratchpad SRAM guarantees single cycle access time to scalars since there are no cache misses. Moreover, the on-chip scratchpad SRAMs have small sizes, making this scheme ideal for data with small memory footprint such as the scalars in embedded media applications. This is also beneficial for a software-directed approach, since as shown in [24] every hardware partition can be logically partitioned and the scalar buffer area can be implemented as a logical partition. We assume the SRAM area to be sufficient to hold all the scalar data. No architectural modifications are necessary since many embedded processors have a scratchpad buffer area, see Table 1.

Therefore, if the embedded processor is equipped with a scratchpad SRAM area, the scalar memory accesses can be annotated by the compiler and remapped to the scratchpad. If not, then the Instruction Set Architecture (ISA) can be augmented by special load-store instructions which would channel the scalar data to a separate cache area. The modifications to the compiler are minimal and consist of statically determining the application memory size and mapping the scalar accesses to the special load-store instructions.

We ran the Mediabench benchmarks with the baseline cache settings; we compared this with the same cache settings but with the scalar accesses being redirected to the SRAM buffer area by the compiler. We stress that capacity misses are not an issue here: Fritts et al. [14] have shown that data working set sizes of the considered benchmarks are very small. The results for selected

benchmarks are presented in Table 5. The improvement depends on the particular benchmark and ranges between 0.6 to 9.5 percent. This improvement is more pronounced for the benchmarks which have a significant percentage of scalars in their memory accesses. Our results also affirm that scratchpad warmup costs are extremely small compared to the number of cache misses.

Table 5. The number of misses for the baseline and for a design with a scalar SRAM buffer are shown in the first and third columns, respectively. The second column shows the baseline miss rate. The percent drop in miss rate for remapping scalars to scratchpad is given in the fourth column. The fifth column is the percentage of scalar accesses to total memory accesses. The last column shows the scratchpad buffer warmup costs, i.e., the cost associated with promoting scalars from main memory to the scratchpad SRAM.

	Baseline Miss (%)	Partitioned Improvement(%)	Scalars(%)	Warmup		
EPIC	1753939	13.6	1589065	9.5	32.0	15598
G721	1377675	2.0	1369395	0.6	4.5	9
GSM	239914	0.5	230549	3.9	2.3	19
JPEG	228644	9.3	224185	1.9	1.1	21
RASTA	216173	6.9	203373	5.9	16.0	59

We also replicated our experiments for a 2 Kbyte 2-way cache organization. Table 6 shows the results. Note that the percentage improvements due to remapping of scalars to scratchpad are similar to the direct mapped cache results.

Table 6. The results for the 2-way associative cache. The first column shows the baseline miss rate. The percentage reduction in miss rate due to remapping is shown in the second column.

	Miss Rate(%)	Improvement(%)
EPIC	12.7	6.0
G721	1.2	0.1
GSM	0.5	5.5
JPEG	5.9	1.6
RASTA	5.3	9.2

4.5 Case Study: Scratchpad Energy Savings for Media Processors

Unless otherwise stated, all the results in this section are with a scratchpad of size 1024 bytes, and the baseline cache is 64Kbyte 2-way associative. We ran the benchmarks using the modified Wattach/SimpleScalar and collected the

data cache energy results. Figure 7 shows the percentage energy savings for our 32 general-purpose register media processor model. We save 10.7% energy on average by using our scheme.

Many media processors such as the ARM have a smaller number of registers, usually 16. Therefore, we have repeated our energy analysis for a 16-register version of our media processor. For 16 registers we have significantly more scalar memory accesses due to register pressure. The results are also shown in Figure 7. Our technique saves in this case an average of 38.2% in energy.

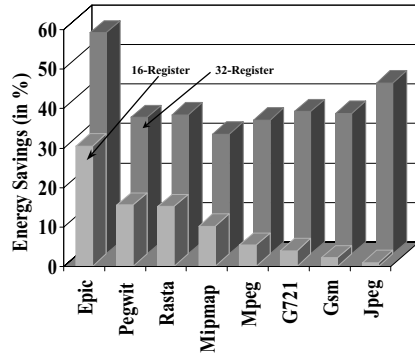


Fig. 7. Scratchpad Energy Savings.

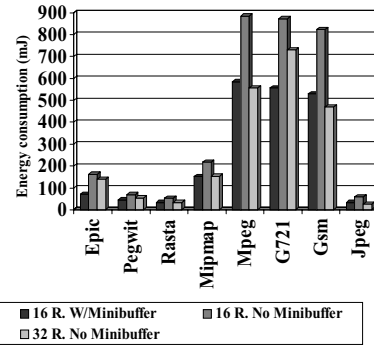


Fig. 8. 16-register architecture with scratchpad can be more energy efficient than 32-register architecture without scratchpad.

In fact, we show that we can be just as energy-efficient with a 16-register media processor with a scratchpad SRAM as a 32-register processor with no scratchpad, see Figure 8. Actually, the overall energy savings are even greater since we just concentrate on the data cache energy consumption: a 16-register file consumes substantially less power than a 32-register file.

5 Conclusion and Future Work

We have performed an analysis of scalars in embedded systems. We established the memory requirements of scalars in embedded applications and presented a compiler algorithm to extract this information. We then discussed several architectural issues pertaining to scalars in embedded systems.

This is ongoing work in line with our vision of creating memory systems with logical partitions where accesses are being mapped based on their static properties [24]. In particular, we are re-integrating our technique with other compiler/architectural techniques that handle diverse types of memory accesses.

References

1. Appel AW, George L (2001) Optimal Spilling for CISC Machines with Few Registers, In: Proceedings of the ACM Sigplan Conference on Programming Language Design and Implementation, pp. 243–253
2. Albonesei DH (1999) Selective Cache Ways: On-Demand Cache Resource Allocation, In: Proceedings of the 32nd International Symposium on Microarchitecture, MICRO32, pp. 248–258
3. Belady LA (1966) A Study of Replacement Algorithms for a Virtual-Storage Computer, IBM Systems Journal, 5(2):78–101
4. Benini L, Macii A, Poncino M (2000) A Recursive Algorithm for Low-Power Memory Partitioning, In: Proceedings of the International Symposium on Low Power Electronics and Design, ISLPED'00, Rapallo, Italy, pp. 78–83
5. Bhandarkar DP (1996) Alpha Implementations and Architecture, Complete Reference Guide, Digital Press, pp. 42–43
6. Bishop B, Kelliher T, Irwin N (1999) A Detailed Analysis of MediaBench, In: Proceedings of the IEEE Workshop on Signal Processing Systems, Taipei, Taiwan
7. Burger D, Austin TD (1997) The SimpleScalar Tool Set, Version 2.0, University of Wisconsin-Madison Computer-Sciences Department Technical Report #1342
8. Brooks D, Tiwari V, Martonosi M (2000) Wattch: A Framework for Architectural-Level Power Analysis and Optimizations, In: Proceedings of the 27th International Symposium on Computer Architecture, ISCA'00, Vancouver, Canada, pp. 83–94
9. Burlin J (2000) Optimizing Stack Frame Layout for Embedded Systems, Masters Thesis, Computing Science Department, Uppsala University, Uppsala, Sweden
10. Chiou D, Jain P, Rudolph L, Devadas S (2000) Application-Specific Memory Management for Embedded Systems Using Software-Controlled Caches, In: Proceedings of the 37th Design Automation Conference, DAC'00, Los Angeles, CA, pp. 416–419
11. Cooper KD, Harvey TJ (1998) Compiler-Controlled Memory, In: Proceedings of the Eighth International Conference on Architectural Support for Programming Languages and Systems (ASPLOS-VIII), pp. 2–11
12. Delaluz V, Kandemir M, Vijaykrishnan N, Irwin MJ (2000) Energy-Oriented Compiler Optimizations for Partitioned Memory Architectures, In: Proceedings of the International Conference on Compilers, Architectures, and Synthesis for Embedded Systems CASES00, San Jose, CA, pp. 138–147
13. Engblom J (1999) Why SpecInt95 Should Not Be Used to Benchmark Embedded Systems Tools, In: Proceedings of the ACM Sigplan Workshop on Languages, Compilers and Tools for Embedded Systems (LCTES'99), pp. 96–103
14. Fritts J, Wolf W, Liu B (1999) Understanding Multimedia Application Characteristics for Designing Programmable Media Processors, In: Proceedings of SPIE, Multimedia Hardware Architectures, San Jose, CA, pp.2–13
15. Huang M, Renau J, Torrellas J (2001) L1 Cache Decomposition for Energy Efficient Processors, In: Proceedings of the International Symposium on Low-Power Electronics and Design, ISLPED'01, Huntington Beach, CA, pp. 10–15
16. Kin J, Gupta M, Mangione-Smith WH (1997) The Filter Cache: An Energy Efficient Memory Structure, In: Proceedings of the 30th Annual Symposium on Microarchitecture, MICRO30, pp. 184–193

17. Kulkarni C, Catthoor F, De Man H (2000) Advanced Data Layout Organization for Multi-media Applications, In: Workshop on Parallel and Distributed Computing in Image Processing, Video Processing, and Multimedia (PDIVM 2000), Cancun, Mexico
18. Lee C, Potkonjak M, Mangione-Smith WH (1997) Mediabench: A Tool for Evaluating and Synthesizing Multimedia and Communications Systems, In: Proceedings of the 30th Annual International Symposium on Microarchitecture, MICRO30, pp. 330–335
19. Lee HS, Tyson GS (2000) Region-Based Caching: An Energy Delay Efficient Memory Architecture for Embedded Processors, In: Proceedings of PACM (CASES'00), San Jose, CA, pp. 120–127
20. <http://www.eecs.harvard.edu/hube/software/software.html>
21. Memik G, Kandemir M, Haldar M, Choudhary A (1999) A Selective Hardware/Compiler Approach for Improving Cache Locality, Northwestern University Technical Report CPDC-TR-9909-016
22. Milutinovich V, Tomasevic M, Markovic B, Tremblay M (1996) The Split Temporal / Spatial Cache: Initial Performance Analysis, In: Proceedings of SCIZZL-5, Santa Clara, CA, pp. 63–69
23. Moritz CA, Frank M, Amarasinghe S (2000) FlexCache: A Framework for Compiler Generated Data Caching, In: Proceedings of the Second Workshop on Intelligent Memory Systems, IRAM00, Held in Conjunction with ASPLOS-IX, Cambridge, MA
24. Moritz CA, Frank M, Amarasinghe S (2001) FlexCache: A Framework for Compiler Generated Data Caching, Lecture Notes in Computer Science, Springer-Verlag
25. Mueller F (1995) Compiler Support for Software-Based Cache Partitioning, In: Proceedings of the ACM SIGPLAN Workshop on Languages, Compilers and Tools for Real-Time Systems, La Jolla, CA, pp. 125–133
26. O'Boyle M, Knijnenburg P (1996) Non-Singular Data Transformations: Definition, Validity, Applications, In: Proceedings of the 6th Workshop on Compilers for Parallel Computers (CPC'96), Aachen, Germany, pp. 287–297
27. Panda PR, Dutt ND, Nicolau A (1997) Efficient Utilization of Scratch-Pad Memory in Embedded Processor Applications, In: Proceedings of the European Design and Test Conference, Paris, France, pp. 7–11
28. Ranganathan P, Adve S, Jouppi NP (2000) Reconfigurable Caches and Their Application to Media Processing, In: Proceedings of the 27th International Symposium on Computer Architecture (ISCA-27), pp. 214–224
29. <http://suif.stanford.edu/>