

1 Journal of Circuits, Systems, and Computers
 Vol. 16, No. 4 (2007) 1–23
 3 © World Scientific Publishing Company



5 **LOW POWER SYSTEM DESIGN BY COMBINING SOFTWARE
 PREFETCHING AND DYNAMIC VOLTAGE SCALING**

7 SUMITKUMAR N. PAMNANI* and DEEPAK N. AGARWAL†

9 *Microprocessor Verification, AMD, M/S — 615,
 5204 East Ben White Blvd, Austin, TX, 78741, USA*

*sumitkumar.pamnani@amd.com

†deepak.agarwal@amd.com

11 GANG QU‡ and DONALD YEUNG§

13 *Electrical and Computer Engineering Department
 and Institute for Advanced Computer Study,
 University of Maryland, College Park, MD, 20742, USA*

‡gangqu@eng.umd.edu

§yeung@eng.umd.edu

17 Received 13 December 2004

Revised 28 October 2005

19 Accepted 21 September 2006

21 Performance-enhancement techniques improve CPU speed at the cost of other valu-
 23 able system resources such as power and energy. Software prefetching is one such tech-
 25 nique, tolerating memory latency for high performance. In this article, we quantitatively
 27 study this technique's impact on system performance and power/energy consumption.
 29 First, we demonstrate that software prefetching achieves an average of 36% performance
 31 improvement with 8% additional energy consumption and 69% higher power consump-
 33 tion on six memory-intensive benchmarks. Then we combine software prefetching with
 a (unrealistic) static voltage scaling technique to show that this performance gain can
 be converted to an average of 48% energy saving. This suggests that it is promising to
 build low power systems with techniques traditionally known for performance enhance-
 ment. We thus propose a practical on-line profiling based dynamic voltage scaling (DVS)
 algorithm. The algorithm monitors system's performance and adapts the voltage level
 accordingly to save energy while maintaining the observed system performance. Our
 proposed on-line profiling DVS algorithm achieves 38% energy saving without any sig-
 nificant performance loss.

35 *Keywords:* Dynamic voltage scaling; software prefetch; energy efficiency; performance.

37 **1. Introduction**

Low power and low energy design is important for battery-operated portable sys-
 tems such as personal computing devices, wireless communication, and imaging

‡Corresponding author.

2 *S. N. Pamnani et al.*

1 systems. Interestingly, many such systems are not hungry for performance since
2 they are already “fast enough” to satisfy end-user’s desire. For example, the phys-
3 ical limitation of human visual and auditory systems implies that a DVD player
4 does not necessarily need to be implemented using the fastest processor. One natural
5 question is *what will be the role, if any, of traditional high-performance techniques*
6 *such as pipelining, caches, prefetching, and branch prediction in low power system*
7 *design?* We believe that the performance gain achieved by these techniques can be
8 converted into power and energy reduction. We propose a general approach that
9 monitors the performance gain and transfers it into power and energy saving at run-
10 time by dynamic voltage scaling (DVS). In this article, we use software prefetching
11 as an example to demonstrate this idea of utilizing high-performance techniques for
12 low power system design.

13 **1.1. Prefetching**

14 Prefetching is a latency tolerance technique, which is generally used to overcome
15 the gap between the (short) processor cycle time and the (long) memory access
16 latency. In prefetching, data is brought to a level of the memory hierarchy that
17 is closer to the processor in advance rather than on demand. This will reduce or
18 eliminate the cache miss penalty and improve performance. More important, it has
19 been regarded as a promising approach for solving the memory wall problem.¹

20 Prefetches can be triggered either by a hardware mechanism, or by a software
21 instruction, or by a combination of both. While hardware prefetching requires extra
22 hardware resources to determine what to prefetch, *software prefetching* approaches
23 detect data access patterns by static program analysis and allow prefetching to be
24 done selectively and effectively. Several prefetching techniques have been proposed
25 in the past to improve performance.^{2–6} However, there has been very little reported
26 on the power and energy overhead associated with prefetching.

27 **1.2. Dynamic voltage scaling**

28 Dynamic voltage scaling is a technique that varies the supply voltage and clock
29 frequency based on the computation load to provide desired performance with the
30 minimal amount of energy consumption. The most practical and well-studied DVS
31 system is the multiple voltage DVS system where a set of predefined discrete volt-
32 ages are available simultaneously due to their simplicity of implementation and
33 effectiveness in power reduction.^{7–9} In fact, the International Technology Roadmap
34 for Semiconductors has predicted such systems as the trend of future low power
35 systems.¹⁰

36 DVS reduces the operating voltage and thus clock frequency. It saves power and
37 energy at the cost of longer execution time. Therefore, it is applicable only when
38 there exists slack time. That is, the system goes back and forth between “busy”
39 and “idle”. When the execution time information is available, there exist optimal
40 algorithms to achieve the maximal energy saving.^{8,9,11} Without knowing a task’s

1 real execution time, voltage can also be scaled based on system level information
such as current computation load and predicted future behavior.^{12–14} Normally,
3 voltage scaling is performed on the arrival and completion of a task, and periodically
during task execution.

5 **1.3. Our approach**

7 In this article, we investigate the power/energy vs. performance tradeoff of software
prefetching on a multiple-voltage system. We first make two observations from sim-
9 ulation: software prefetching is a very power- and energy-consuming technique, but
voltage scaling has the potential to balance the software prefetching’s performance
11 gain and power/energy overhead. Therefore, we propose a framework that auto-
matically transfers any “excess” performance enhancement (provided by prefetch-
13 ing or any other techniques) to energy reduction. At the heart of this approach is
an on-line profiling based algorithm that monitors the system’s performance and
15 adjusts the system voltage accordingly to save energy without causing noticeable
performance degradation. Although we illustrate this for software prefetching, our
17 approach is general and applicable to other high-performance techniques. In sum,
we claim that system designers can and should consider traditional performance-
enhancement techniques for low power system design.

19 The rest of this article is organized as follows: In Sec. 2, we review the previous
work on prefetching for high-performance and the low power DVS technique. In
21 Sec. 3, we elaborate on the proposed on-line profiling based DVS algorithm that
transfers the performance gain of prefetching automatically to energy saving via
23 careful voltage scaling. We report our experimental methodology in Sec. 4 and
evaluate the simulation results in Sec. 5. We conclude in Sec. 6 with discussion on
25 future research directions.

27 **2. Related Work**

In this section, we elaborate on the previous work in hardware and software prefetch-
ing as well as how slack time is utilized in DVS to reduce power and energy.

29 **2.1. Prefetching for high-performance**

The hardware prefetching approach detects accesses with regular patterns and
31 issues prefetches at run-time. The main advantages of the hardware-based approach
are that prefetches are handled dynamically without compiler intervention and
33 that code compatibility is preserved. However, it suffers from the requirement of
extra hardware and also that unwanted data could be prefetched. Jouppi³ intro-
35 duced stream buffers to improve direct mapped cache performance. Chen and Baer²
investigated a mechanism for prefetching data references characterized by regular
37 strides. Mowry and Gupta⁶ showed that prefetching can be effective for tolerating
large memory latencies in a shared memory multiprocessor environment.

4 *S. N. Pamnani et al.*

1 On the other hand, software-directed prefetching approaches rely on data access
2 patterns detected by static program analysis and allow the prefetching to be done
3 selectively and effectively. An intelligent compiler inserts data prefetch instructions
4 so that they execute several cycles before their corresponding memory instruc-
5 tions. These prefetch instructions are explicitly executed by the processor to initiate
6 prefetch requests. Porterfield^{15,16} presented a compiler algorithm which prefetched
7 all array references in inner loops one iteration ahead. Klaiber and Levy¹⁷ extended
8 this work by recognizing the need to prefetch more than a single iteration ahead.
9 Gornish *et al.*¹⁸ presented an algorithm for determining the earliest time when it is
10 safe to prefetch an entire subarray. The software scheme has some drawbacks. First,
11 some useful prefetching may not be uncovered so cache misses may not be elimi-
12 nated completely. Second, there is an execution overhead due to the extra prefetch
13 instructions.

14 Both hardware-based techniques and software-directed prefetching approaches
15 are successful in identifying prefetches of data with simple constant stride patterns.
16 These techniques have been studied solely for improving performance. Little is
17 known about the power and energy overhead of prefetching.

2.2. *Dynamic voltage scaling for low power*

19 We have already mentioned that DVS is applicable in the presence of slack time.
20 We now survey the existing efforts on DVS in two categories. The first set of work
21 assumes slack comes from the difference among a task's deadline, worst case execu-
22 tion time, and real execution time. The scheduler must know these information to
23 scale voltage accordingly. The second set of work does not require such application
24 information and scales voltage based on system level information such as current
25 computation load and predicted future behavior.

26 Yao *et al.*¹⁹ were among the first to study scheduling policies for energy reduc-
27 tion on variable speed CPUs. They built a model for the tasks to be executed,
28 where each task is associated with its arrival time, deadline, and worst case execu-
29 tion time. They discussed how to schedule the tasks and determine the CPU's speed,
30 controlled by varying voltage to complete all tasks with the minimum amount of
31 energy. There have been rich follow-up works addressing similar problems for peri-
32 odic tasks, aperiodic tasks, or mixed tasks; for priority tasks; for ideal voltage scaling
33 system where there is no overhead for voltage changing; for practical voltage scaling
34 system where there are physical constraints on maximal/minimal voltage and how
35 fast voltage can be changed; etc.^{7-9,11,20}

36 Weiser *et al.*¹⁴ proposed to divide time into intervals of 10-50 ms and adjust
37 the CPU clock speed by the task-level scheduler based on the processor utilization
38 over the preceding interval. Govil *et al.*¹² enhanced this by looking for recurring
39 patterns of processor utilization and setting processor speed accordingly. Pillai and
40 Shin²¹ presented a class of algorithms that modify the operating system's real-time
41 scheduler and task management service to save energy while maintaining real-time

1 deadline guarantee. Most recently, Hua *et al.*²² proposed the probabilistic design for
2 soft real-time systems, such as multimedia applications, where occasional execution
3 failure can be tolerated. They intentionally dropped tasks to reduce energy while
4 providing a statistical performance guarantee.

5 The novelty of our approach is that we utilize performance-enhancement tech-
6 niques (such as software prefetching) to reduce the actual execution time of the
7 task, and thus create slack that may not exist otherwise. Such slack will provide us
8 opportunities to apply DVS in order to reduce both power and energy consumption.

9 **2.3. Multiple-voltage DVS system**

10 Most DVS systems mentioned above assume on-chip DC–DC converters to provide
11 continuous and flexible operating voltage at run-time to achieve the maximal power
12 saving.^{9,13,19} This inevitably introduces time and energy overhead, most notably
13 the time for voltage to stabilize at the new level. Moreover, such overhead cannot
14 be treated as constant because the *transient response time* and consequently, the
15 energy consumption for voltage switch depend on the difference between the source
16 voltage and the desired voltage.

17 Unlike these systems, we consider DVS systems with multiple levels of supply
18 voltage and threshold voltage physically implemented on the chip. Such systems
19 are referred to as multiple-voltage DVS systems and can be implemented by using
20 multiple-voltage regulators, each of which is dedicated to generating one specific
21 voltage and clock frequency. This eliminates the use of DC–DC converters or other
22 auxiliary devices (such as buffers, delay line, and/or charge pump) to implement
23 DVS. By doing so, although we sacrifice the ability to generate voltage at any level,
24 we minimize the time and energy overhead during voltage switches.

25 First, the delay overhead for each voltage switch will be a constant. The oper-
26 ating system controls the clock frequency at run-time by writing to a register in
27 the system control state the same way as in Ref. 13. This takes one clock cycle
28 and there is a fixed transient response time (normally a few cycles as reported in
29 Refs. 24 and 25).

30 Second, the power dissipation on the voltage regulators will be the sum of one
31 regulator's dynamic power and the static power for other regulators. This is because
32 the system, when active, uses only one regulator at any time and can shut down
33 the rest for energy conservation. Note that recent advances in linear voltage regu-
34 lator design have led us to low dropout (LDO) regulators with high efficiency, low
35 power, and low transient response time. The LDO regulator's quiescent current is
36 in the order of μA 's²⁵ and its dynamic power can be well below the mW mark. For
37 example, Ref. 24 reports an adaptive voltage scaling controller that consists of a
38 voltage regulator, a charge pump, a resettable delay-line, level shifters, and a clock
39 generator among other devices. This controller consumes 2 mW at 3 V and 20 MHz,
with most energy dissipation spent on devices other than the regulator.

6 *S. N. Pamnani et al.*

1 Third, using multiple LDO regulators may not cause more area overhead than
2 a DC–DC converter. This is because we do not require the auxiliary devices to
3 produce flexible voltage and, as it has been reported,²⁶ two or three regulators will
4 be sufficient for energy efficiency for most application-specific embedded systems.

5 **3. Software Prefetching for Energy Reduction**

6 In this section, we first layout our framework of transferring processor’s perfor-
7 mance improvement to power/energy saving. Then we use software prefetching as
8 a performance-enhancing technique and apply the proposed framework to convert
9 the performance gain from prefetching to energy saving.

10 **3.1. Bridging performance enhancement and energy reduction**

11 The continuing push for high performance has brought us to a point where the
12 microprocessor is sufficiently fast to handle specific applications. Such improved
13 performance comes, in general, with larger amount of transistors and higher power
14 consumption. For example, our simulation shows that software prefetching is able
15 to improve the performance of our benchmarks by 36%. However, this takes 8%
16 more energy consumption and the average power consumption is 69% higher. (The
17 power is much higher because with prefetching, more energy is consumed in a
18 shorter period. See Sec. 5.1 for the detailed report.)

19 This experimental finding clearly suggests that it is important to study the
20 power and energy behavior when evaluating any performance-enhancement tech-
21 nique. A less obvious issue is whether such techniques should be considered at all in
22 the design of low power systems where only moderate performance is required. Con-
23 sider the same example as above. With the 36% performance gain from prefetching,
24 we are able to slow the system down without suffering any performance loss com-
25 pared to the system without prefetching. We enable the clock slow down by scaling
26 voltage down and this results in a 48% power and energy saving! (We have the same
27 average power and energy saving because we slow the clock such that we complete
28 the execution at the same time that the system without prefetching would. See
29 Sec. 5.2 for the detailed report.)

30 This suggests that it is possible to use performance-enhancement techniques for
31 low power system design. We propose a framework that automatically transfers the
32 performance gain to energy and power reduction. Such transfer requires:

- 33 • *Balanced performance degradation and energy reduction.* Ideally, we should main-
34 tain the user specified performance level and convert all the “excess” performance
35 beyond that level to power and energy saving.
- 36 • *Minimum re-design effort.* It should be easy to integrate this transfer in existing
37 systems. This implies that we should not assume information, which requires
38 extensive offline profiling, is known *a priori*.

```

A(N,N,N),B(N,N,N)
do j=2,N-1
  do i=2,N-1
    A(i,j) = 0.25 * (B(i-1,j) + B(i+1,j) + B(i,j-1) + B(i,j+1))
  
```

Fig. 1. Example affine array traversal code from the 2D Jacobi kernel.

1 In the remainder of this section, we elaborate on this proposed framework by
 2 presenting an algorithm that transfers “excess” performance, which is provided by
 3 software prefetching, to energy saving by DVS. Our algorithm does not require the
 4 unrealistic information such as a task’s real execution time, the success ratio and
 5 distribution of prefetches, and the performance gain by each successful prefetch.

3.2. Software prefetching as a performance-enhancing technique

7 Software prefetching is a promising technique for overcoming the speed gap between
 8 processor and memory. It relies on the programmer or compiler to insert explicit
 9 prefetch instructions into the application code for memory references that are likely
 10 to miss in the cache. At run-time, the inserted prefetch instructions bring the data
 11 into the processor’s cache in advance of its use, thus overlapping the cost of the
 12 memory access with useful work in the processor. This section discusses how appli-
 13 cation code is instrumented with the prefetch instructions.

14 Several algorithms have been previously proposed to instrument software
 15 prefetching.^{4,5,23} The best algorithm to apply depends on the type of memory refer-
 16 ences performed by the application code. The most basic memory reference pattern
 17 is affine (linear) accesses to multidimensional arrays such as the 2D Jacobi kernel
 18 shown in Fig. 1. This benchmark performs a grid computation in which the array
 19 A elements are computed as the average of neighboring values in both dimensions
 20 from the array B elements. Prefetches for affine accesses like those in the 2D Jacobi
 21 kernel can be instrumented using Mowry’s algorithm.⁵ Figure 2 illustrates the 2D
 22 Jacobi kernel after Mowry’s algorithm has been applied.

23 Mowry’s prefetch algorithm involves three steps. First, the affine array refer-
 24 ences within inner-most loops are identified as prefetching candidates. To permit
 25 prefetching of only the missing dynamic instances of identified references, loop
 26 unrolling is performed to create multiple static instances of each memory reference
 27 within the loop body; the degree of loop unrolling is set to the number of back-to-
 28 back memory references co-located in the same cache block. By unrolling the loop
 29 this number of iterations, the dynamic instances that miss in the cache are isolated
 30 — the leading static memory reference in the unrolled loop always misses, while the
 31 remaining unrolled static memory references always hit. Hence, a single prefetch for
 32 the leading memory reference can be inserted into the unrolled loop, prefetching
 33 exactly the missing dynamic instances, thus avoiding unnecessary prefetches and
 34 reducing prefetch overhead.

8 *S. N. Pamnani et al.*

```

A(N,N,N),B(N,N,N)

do j=2,N-1
  do i=2,N-1                                // Prologue Loop
    prefetch(&B[i][j])
    prefetch(&B[i][j+1])
    prefetch(&B[i][j-1])
    prefetch(&A[i][j])

  do i=2,N-PD-1,step=4                       // Unrolled Loop
    prefetch(&B[i+PD][j])
    prefetch(&B[i+PD][j+1])
    prefetch(&B[i+PD][j-1])
    prefetch(&A[i+PD][j])

    A(i,j) = 0.25 * (B(i-1,j) + B(i+1,j) + B(i,j-1) + B(i,j+1))
    A(i+1,j) = 0.25 * (B(i,j) + B(i+2,j) + B(i+1,j-1) + B(i+1,j+1))
    A(i+2,j) = 0.25 * (B(i+1,j) + B(i+3,j) + B(i+2,j-1) + B(i+2,j+1))
    A(i+3,j) = 0.25 * (B(i+2,j) + B(i+4,j) + B(i+3,j-1) + B(i+3,j+1))

  do i=N-PD,N-1                              // Epilogue Loop
    A(i,j) = 0.25 * (B(i-1,j) + B(i+1,j) + B(i,j-1) + B(i,j+1))

```

Fig. 2. 2D Jacobi kernel from Fig. 1 instrumented with software prefetches using Mowry’s algorithm. The instrumented code contains a prologue loop, an unrolled “steady-state” loop, and an epilogue loop.

1 Figure 2 illustrates the loop unrolling and prefetch insertion transformations
 2 for the 2D Jacobi kernel. In this loop, there are five affine references, four for the
 3 B array, and one for the A array. Assuming each reference accesses 8 bytes and
 4 assuming a 32-byte cache block, the loop should be unrolled by a factor of 4. After
 5 loop unrolling, four prefetches are inserted for the five leading memory references:
 6 $A(i, j)$, $B(i - 1, j)$, $B(i + 1, j)$, $B(i, j - 1)$, and $B(i, j + 1)$. Notice the $B(i - 1, j)$
 7 and $B(i + 1, j)$ references lie on the same cache block, thus saving 1 prefetch.

8 The next step in Mowry’s algorithm is to perform *prefetch scheduling*. Given
 9 the high memory latency of most modern memory systems, a single loop iteration
 10 normally contains insufficient work under which the cost of memory accesses is
 11 hidden. To ensure that data arrive in time, the prefetches inserted in the first step
 12 of the algorithm must be initiated multiple iterations in advance. The minimum
 13 number of loop iterations needed to fully overlap a memory access is known as the
 14 *prefetch distance*. Prefetch scheduling determines the prefetch distance and applies
 15 it to the inserted prefetches. Assuming the memory latency is l cycles and the work
 16 per loop iteration is w cycles, the prefetch distance, PD, is simply $\lceil \frac{l}{w} \rceil$. Figure 2
 17 illustrates the indices of prefetched array elements contain a PD term, providing
 18 the early prefetch initiation required.

19 Finally, the last step is to “fix up” inefficiencies in prefetching created by
 prefetch scheduling. By modifying prefetches to initiate PD iterations in advance as

1 illustrated in Fig. 2, the first PD iterations do not get prefetched, and the last PD
 2 iterations prefetch past the end of each array. These inefficiencies can be addressed
 3 by performing loop peeling to handle the first and last PD iterations of the loop
 4 separately. The loop peeling transformation creates a *prologue loop* to execute PD
 5 prefetches for the first PD array elements, and an *epilogue loop* to execute the last
 6 PD iterations without prefetching. Figure 2 illustrates the prologue and epilogue
 7 loops created by loop peeling.

8 In addition to software prefetching for affine array accesses, in Sec. 5, we also
 9 consider software prefetching for indexed array and pointer-chasing accesses. We
 10 use the algorithm in Ref. 5 for prefetching indexed array references which is very
 11 similar to the algorithm described above except for a small extension to the prefetch
 12 scheduling step. For pointer-chasing accesses, we use the prefetch arrays technique.⁴
 13 This approach requires creating and managing *prefetch pointers*, in addition to
 14 instrumenting software prefetches, to overcome the serialization effects of pointer-
 15 chasing memory references. We omit the detailed explanation of these techniques
 and refer the reader to the cited papers for more information.

17 3.3. Transferring performance gain to energy saving via DVS

18 Let t' and t be the time to run an application with and without prefetching, respec-
 19 tively. When prefetching is successful, the system will complete the application $t - t'$
 20 earlier than required. That is, prefetching improves performance and creates slack.
 21 This slack allows the system to slow down (to conserve power and energy) while
 22 still managing to complete at time t . The ratio $\frac{t-t'}{t}$ measures the performance gain
 23 and now we discuss how to transfer such gain into energy saving.

24 The optimal way for voltage scaling to reduce energy is by switching the voltage
 25 level from v to v' such that the gate delay, which is proportional to $\frac{v_{dd}}{(v_{dd}-v_{th})^2}$, will
 26 be increased by a factor of $\frac{t}{t'}$. In this case, the application will finish at time t if
 27 the system applies prefetching and operates at the constant voltage v' . If, however,
 28 the system can only operate at certain pre-designed voltage levels (such a system
 29 is referred to as a *multiple-voltage system*) that do not include v' , then the most
 30 energy-efficient way is to run at the voltage v_1 which is slightly higher than v' for
 31 some time and then reduce to the next lower level v_2 (where $v_1 > v' > v_2$) such
 32 that the execution terminates at t .^{8,9} Note that power dissipation is approximately
 33 proportional to Cv_{dd}^2f , where C is the capacitance and f is the clock frequency. The
 34 energy consumption will be reduced because the application is operated at lower
 35 voltages and frequencies.

36 Unfortunately, the transfer from performance gain to energy saving is not so
 37 simple and straightforward. First, for most real-time applications, we do not have
 38 the luxury to know both t and t' and therefore will be unable to determine the opti-
 39 mal voltage scaling strategy. Off-line profiling will not provide a complete solution
 40 because the success of prefetching depends on the application's run-time behavior
 41 and the performance gain from each successful prefetch may be different and is,

10 *S. N. Pamnani et al.*

1 in general, unpredictable. Second, even if we know both t and t' , each prefetch is
 3 independent and successful prefetches may not be uniformly distributed throughout
 5 the application's execution. Consider that prefetching helps to improve the perfor-
 7 mance of a multimedia application (e.g., by reducing the frame decoding time).
 9 However, most of the prefetchings fail in the first half of the execution and we will
 11 constantly fall behind the execution at the reference voltage without prefetches.
 Although (successful) prefetchings will eventually help us to catch up (by complet-
 ing the multimedia application at the same time as the run at reference voltage), the
 slow down in the first half may cause unacceptable user-perceivable performance
 degradation.

13 We propose an on-line profiling DVS algorithm that automatically selects volt-
 15 age at run-time based on the accumulated performance gain to reduce energy and
 17 provide performance guarantees (if prefetching does not have any negative impact
 on system performance. See Theorem 1 and its proof for detail). Figure 3 depicts
 our algorithm.

19 We periodically conduct real-time profiling to estimate the performance gain
 21 of prefetching. For each N instructions, we execute the first M instructions with
 23 prefetching and the next M instructions without prefetching (which can be easily
 controlled by executing prefetch instructions as NOPs). We then compute c_p and
 c_{np} , the cycles for these two sets of instructions as shown in steps 2-5. Next, we
 calculate the prefetching gain in step 6. If there is no gain due to prefetching, then
 we pull up to the reference voltage for execution. This is because prefetching has

```

1. repeat each cycle for each  $N$  instructions {
2.   for  $M$  instructions
3.     compute  $c_p$ , number of cycles to execute  $M$  instructions with prefetching;
4.   for  $M$  instructions
5.     compute  $c_{np}$ , number of cycles to execute  $M$  instructions without prefetching;
6.    $gain = \frac{c_{np} - c_p}{M}$ ;
7.   if ( $gain < 0$ )
8.     pull to the reference voltage and execute  $W$  instructions;
9.   goto step 2;
10.  repeat each cycle for  $N - 2M$  instructions {
11.    compute  $present_{saving}$ ;
12.    if ( $present_{saving} \geq c_{th}$ )    voltage_down();
13.    if ( $present_{saving} < 0$ ) {
14.      if ( $present_{saving} + cumulative_{saving} \geq 2c_{th}$ )    voltage_down();
15.      if ( $present_{saving} + cumulative_{saving} < c_{th}$ )    voltage_up();
16.    }
17.  }
18. }
```

Fig. 3. Pseudo-code for real-time profiling DVS algorithm.

1 not been helpful and executing at any voltage lower than reference would result in
 2 performance loss. After W instructions, we start profiling again (step 9).

3 If prefetching does provide gain, we repeat steps 10–17 for the next $N - 2M$
 4 instructions. We compute $present_{savings}$, which keeps track of how far ahead we
 5 are as compared to execution at reference voltage without prefetching. It can be
 6 quantitatively defined as the difference between prefetching gain and slow down
 7 of running at lower voltage. We then scale voltage up and down by instructions
 8 $voltage_up()$ and $voltage_down()$ based on both the current and accumulated per-
 9 formance gain. The algorithm is designed in a manner such that $voltage_up()$ is
 10 called only when the accumulated performance saving is below a pre-determined
 11 threshold c_{th} (step 15).

We use the following parameter values for this algorithm in our simulation:

- 13 (1) N : 100 K instructions (how often do we profile).
- 14 (2) M : 5 K instructions (profile period).
- 15 (3) c_{th} : $50\mu s$ (threshold against which the savings are compared).
- 16 (4) W : 5 K instructions (waiting time before re-profiling).

17 The setting of these parameters is application dependent and directly affects
 18 the performance of the algorithm. In order to guarantee no loss of performance,
 19 the values of N , M , and c_{th} have to be selected carefully. Although there is no
 20 systematic way to determine the optimal values of these parameters to achieve the
 21 most energy saving, they can be set rather precisely for application-specific systems
 22 by studying the behavior of the application. For general purpose systems, we can
 23 use the following basic guidelines and as we will show in the experimental evaluation
 24 section, these guidelines lead to the above parameter settings that give significant
 25 energy saving for a variety of data-intensive applications.

26 First, the profiling period M needs to be sufficiently long to distinguish the
 27 system performance with and without prefetching. On the other hand, our basic
 28 assumption is that prefetching will improve performance and we will use prefetching
 29 whenever possible. This implies that M should be relatively small compared to N
 30 so only a small portion of the instructions will be executed without prefetching (in
 31 step 3). Next, when we execute with prefetching for the $N - 2M$ instructions (steps
 32 10–17), we are conservative in converting the performance saving to energy saving
 33 with the use of threshold c_{th} . The selection of c_{th} is based on the setting of the
 34 voltages. The larger the gap is between two voltage levels, the larger c_{th} is.

35 Finally, we mention that the proposed method is only applicable when the
 36 prefetching technique is effective in giving performance gain. This includes many
 37 real life applications that have regular behavior as we will elaborate in the next sec-
 38 tion. If the application has irregular behavior, for example, when its computation
 39 loads in the first M instructions (when we do not prefetch), in the second M instruc-
 40 tions (when we do prefetch), and the next $N - 2M$ instructions are dramatically
 41 different, our estimated performance gain will not be accurate and our algorithm
 may fail to provide the performance guarantee. This is because when we re-evaluate

12 *S. N. Pamnani et al.*

1 the performance gain from executing the $N - 2M$ instructions with prefetching, we
 2 use the performance of the first M instructions without prefetching as the baseline.
 3 In such case, varying the setting of above parameters can help.

4 To summarize this section, we give two theorems about the low power and
 5 performance-guarantee features of the proposed algorithm.

7 **Theorem 1.** *Compared to the execution at the reference voltage without
 prefetching, the proposed algorithm guarantees that there is no performance degra-
 dation due to voltage scaling.*

9 **Proof.** We start with prefetching at the reference voltage (steps 1–3). If we assume
 10 that prefetching does not have any negative impact on system performance, we will
 11 not fall behind the non-prefetching execution for the first M instructions. The next
 12 M instructions will be executed without prefetching, so we maintain our perfor-
 13 mance gain, if any, from the successful prefetches in the first M instructions. For
 14 the rest of the execution, when we call *voltage_down()*, the system’s speed will slow
 15 down and we will lose the cycle savings that we have accumulated. However, when
 16 we set the parameters (profiling frequency N , period M , and threshold c_{th}) care-
 17 fully such that the *voltage_up()* calls will pull up voltage back to reference before
 we consume all the performance gain by prefetching. \square

19 **Theorem 2.** *On multiple-voltage system, the proposed algorithm converges to the
 optimal voltage setting.*

21 **Proof.** For multiple-voltage system, if the optimal v_{opt} is available, the optimal
 22 voltage scaling strategy is to run at v_{opt} ; otherwise, the system will run only at
 23 the two voltages v_1 and v_2 such that $v_1 > v_{opt} > v_2$ and there is no other voltage
 24 available between v_1 and v_2 . Our algorithm starts with the reference voltage and
 25 reduces voltage gradually via *voltage_down()* calls. If the prefetching performance
 26 gain is estimated accurately, the *present_saving* computed in step 11 will become
 27 negative when the current voltage is below v_{opt} . However, our algorithm may still
 28 lower voltage if the accumulated performance saving is high and will not raise
 29 voltage until the performance saving is below the threshold c_{th} . The *present_saving*
 30 starts to be positive when the voltage level passes v_{opt} and reaches v_1 . Now, our
 31 algorithm will not raise voltage any higher (one condition for *voltage_up()* call in
 32 step 13 is false). Once the accumulated performance saving exceeds c_{th} , the voltage
 33 goes down to v_2 (step 12). At this level, we are running slower than v_{opt} (i.e.,
 34 *present_saving* < 0), so the voltage will not go below v_2 (step 14). If the accumulated
 35 performance gain is below c_{th} , we will raise voltage back to v_1 (step 15). This implies
 36 that we have reached the optimal voltage setting: toggling between the two voltages
 37 v_1 and v_2 next to v_{opt} . This claim is also confirmed by a couple of applications in
 our simulation. \square

1 4. Experimental Methodology

3 We use software prefetching to improve application performance. All the bench-
 5 marks used are instrumented with software prefetching by hand, following the algo-
 7 rithms discussed in Sec. 3.2. The performance of these optimized codes was then
 9 measured on a detailed architectural simulator. The performance boost achieved
 11 was later traded for power savings by voltage scaling.

7 4.1. Application overview

9 Our experimental evaluation employs six benchmarks, representing three classes of
 11 data-intensive applications. Table 1 lists the benchmarks along with their problem
 13 sizes and memory access patterns. The first two applications perform affine array
 15 accesses. MatMult multiplies two matrices and Jacobi performs a 3D Jacobi relax-
 17 ation, which is frequently found in multigrid PDE solvers such as MGRID from
 19 the SPEC/NAS benchmark suite. The next three applications perform indexed
 array accesses. Irreg is an iterative PDE solver for an irregular mesh. Moldyn is
 abstracted from the non-bonded force calculation in CHARMM, a key molecu-
 lar dynamics application used at NIH to model macromolecular systems. NBF is
 abstracted from the GROMOS molecular dynamics code. Finally, the last applica-
 tion, Health, performs pointer-chasing accesses. It simulates the Columbian health
 care system and is taken from the OLDEN benchmark suite.²⁷

4.2. Simulation environment

21 We use Wattch, an architecture level power analysis tool,²⁸ in our simulation.
 23 Wattch provides a framework for analyzing and optimizing microprocessor power
 25 dissipation and is an order of magnitude faster than the existing layout-level power
 27 tools, and at the same time maintains accuracy within 10%. Wattch models main
 processor units as one of the following four structures: array structures, fully asso-
 ciative content-addressable memories, combinational logic and wires and clocking.

27 The power consumption of the units modeled very much depends on their spec-
 29 ific implementation, particularly on the internal capacitances for the circuits that
 make up the processor. Wattch models the internal capacitance using assumptions
 similar to those made by Wilton and Jouppi²⁹ and Palacharla *et al.*³⁰ Current CPU

Table 1. Description of the benchmark applications.

Application	Problem Size	Access Pattern
MATMULT	200 × 200 matrices	Affine array
JACOBI	200 × 200 × 8 grid	Affine array
IRREG	14 K node mesh	Indexed array
MOLDYN	13 K molecules	Indexed array
NBF	14 K node mesh	Indexed array
HEALTH	5 levels, 500 iters	Pointer-chasing

14 *S. N. Pamnani et al.*

1 designs increasingly use conditional clocking to disable all or part of a hardware unit
 2 to reduce power consumption when it is not needed. We use Conditional clocking-3
 3 as described in Ref. 28 for our simulation. In this clock gating, power is scaled
 4 linearly with port or unit usage; unused units dissipate 10% of their maximum
 5 power.

6 The Wattch power model is interfaced with Simplescalar sim-outorder,³¹ which
 7 is used to keep track of different units accessed per cycle and hence record the total
 8 energy consumed for a given application. Sim-outorder is a detailed simulator sup-
 9 porting out-of-order issue and execution based on Tomasulo's approach for dynamic
 10 instruction scheduling. The processor's memory system is simple and does not accu-
 11 rately model DRAMs or the memory bus contention. We have replaced this by a
 12 cycle accurate memory controller and DRAM memory sub-system. Each L2 request
 13 to the memory controller simulates several actions: queuing of the request in the
 14 memory controller, Row Access Strobe (RAS) and Column Access Strobe (CAS)
 15 cycles, and data transfer across the memory system bus. We simulate concurrency
 16 between DRAM banks, but bank conflicts require back-to-back DRAM accesses to
 17 perform serially. Also memory requests are serialized at the memory bus, i.e, we
 18 model bus contention, but we assume infinite bandwidth between the L1 and L2
 19 caches. The bottom portion of Table 2 lists the parameters for our baseline memory
 20 sub-system model.

21 The reason for going to a detailed memory model is to get realistic performance
 22 numbers for prefetching. Prefetching increases memory activity by issuing memory
 23 requests frequently. Thus, modeling memory contention is crucial for accurately
 24 simulating prefetching performance. For our power study, we focus on the processor
 25 because in most systems, the processor consumes the most energy. This is true for
 small hand-held devices like PDAs³² which have very few components, but also

Table 2. Simulation parameters for the processor, cache, and memory sub-system models. Latencies are reported either in processor cycles or in nanoseconds. We assume a 1.25-ns processor cycle time. (^alatency is for floating operations.)

<i>Processor model: 600 MHz</i>			
Issue width	8	Integer latency	1 cycle
Instruction window size	64	Add/mult/Div latency ^a	2/4/12 cycles
Load-store queue size	32	Branch predictor	gshare
Fetch queue size	32	Branch predictor size	2048 entries
Integer/floating point units	4/4	BTB size	2048 entries
<i>Cache model: 1cycle = 1.25 ns</i>			
L1/L2 Cache size	16K-split/512K-unified	L1/L2 associativity	2/4 cycles
L1/L2 Cache block size	32/64 bytes	L1/L2 Latency	1/10 cycles
L1/L2 MSHRs	8/16	L1/L2 write buffers	8/16
<i>Memory sub-system model</i>			
DRAM banks	32	Row access strobe	22.5 ns
Memory system bus width	64 bytes	Column access strobe	22.5 ns
Address send	7.5 ns	Data transfer (per 8 bytes)	7.5 ns

1 on large laptop computers³³ that have many components, including large displays
with backlighting. As a result, the design problem generally boils down to a trade-off
3 between the computational power of the processor and the system’s battery life.

Our baseline simulations assume a SimpleScalar processor model running at
5 600 MHz and at 1.6 V. For voltage scaling, we use the same five different voltage
levels, 1.60 V, 1.50 V, 1.40 V, 1.25 V, and 1.10 V, as those in Transmeta’s Crusoe
7 processor. They provide frequencies of 600 MHz, 500 MHz, 400 MHz, 300 MHz, and
200 MHz, respectively.³⁴ As we have discussed in Sec. 2.3, we assume that the time
9 and power overhead for such multiple-voltage DVS system to switch voltage from
one level to another is negligible. In our proposed technique to convert performance
11 gain into power reduction, only a few counters are required for the real-time profil-
ing. Their power dissipation and the execution time of the profiling algorithm are
13 both considered in the simulation.

Finally, we mention that memory has become one of the major energy consumers
15 in modern systems. In the memory hierarchy, the memory’s power consumption
goes up as its speed goes up (or as it gets closer to the processor). Our simulator
17 models both L1 and L2 cache and hence their energy consumption is included in
the simulation. In addition, when the processor slows down due to voltage scaling,
19 the memory speed can also be reduced accordingly without affecting the system
performance. Therefore, in our simulation, we scale the voltage for both L1 and L2
21 caches whenever the processor speed changes. This reduces the energy consump-
tion in the memory hierarchy without affecting the system performance. However,
23 peak power may increase due to heavy memory activity during prefetching. In our
experimental results, we focus only on average power or the energy consumption in
25 a given period of time.

5. Experimental Evaluation

27 This section describes the experimental result and evaluates the performance of
software prefetching and power savings achievable through voltage scaling.

29 5.1. Prefetching performance

We begin by evaluating the performance of software prefetching over no prefetch-
31 ing. Figure 4 plots the normalized execution time along y -axis for the prefetching
and no-prefetching versions of different applications. The detailed execution time
33 information is also reported in the second and third columns of Table 3.

Each execution-time bar has been broken down into three components: use-
35 ful computation, prefetch-related software overhead, and memory stall, labeled
“Busy,” “Overhead,” and “Mem,” respectively. “Busy” is the execution time of
37 the “NP” version assuming a perfect memory system (e.g., all memory accesses
complete in one cycle). “Overhead” is the incremental increase in execution time
39 of the “P” version over “Busy,” again on a perfect memory system. “Mem” is the

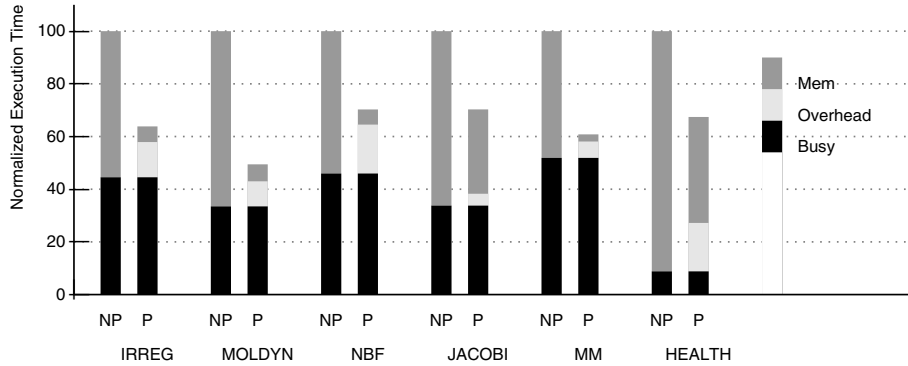
16 *S. N. Pamnani et al.*

Fig. 4. Execution time breakdown for annotated memory instructions. Individual bars show performance without prefetching labeled “NP” and with prefetching labeled “P”.

Table 3. Simulation results on static voltage scaling. The numbers in bold indicate the voltage setting such that there is no significant increase in the execution time.

Application	NP(1.6 V)		P(1.6 V)		P(1.5 V)		P(1.4 V)		P(1.25 V)		P(1.1 V)	
	ms	Energy	ms	Energy	ms	Energy	ms	Energy	ms	Energy	ms	Energy
IRREG	22.38	8.45	14.29	9.09	17.20	6.60	21.65	4.59	28.58	2.80	43.30	1.49
MOLDYN	22.23	7.36	10.99	6.50	13.24	4.70	16.65	3.29	21.98	2.01	33.30	1.07
NBF	13.21	5.46	09.28	6.45	11.10	4.69	14.07	3.25	18.57	1.99	28.14	1.05
JACOBI	14.04	4.99	09.87	4.96	11.89	3.60	14.96	2.49	19.75	1.51	29.90	0.80
MATMULT	86.74	40.5	52.75	39.4	63.55	28.6	79.93	19.8	105.5	12.1	159.8	6.40
HEALTH	10.45	2.32	07.05	3.13	08.49	2.27	10.68	1.57	14.10	9.60	21.36	5.09

1 incremental increase in execution time over “Busy”+“Overhead” assuming a real
2 memory system. All times are normalized against the “NP” bars.

3 Comparing NP and P bars we find prefetching enhances the performance of all
4 the applications significantly. For example, in MOLDYN, prefetching reduces the
5 execution time by 50.56%. On an average, software prefetching is able to boost the
6 performance by 36.04%. We now see the system’s power and energy consumption
7 with and without prefetching.

8 In Table 3, the second and the third columns give the execution time and energy
9 consumption for the runs at the same reference 1.6 V voltage without (“NP” col-
10 umn) and with (“P” column) prefetching, respectively. We see that the run with
11 prefetching does not necessarily save energy consumption. We mention that we have
12 used clock gating to turn off hardware units whenever they are idle. There are three
13 sources for prefetching’s energy overhead. First, it utilizes the parallelism among
14 different hardware units and therefore keeps them active, dissipating energy. Sec-
15 ond, the execution of prefetching instructions consumes energy. Third, unsuccessful
16 prefetches happen, further increasing energy consumption. Only when the perfor-
17 mance gain from prefetching is sufficiently high (such as in the case of Moldyn and
18 MatMult), can we observe energy saving. On an average, prefetching requires 7.67%
19 more energy to achieve the above 36.04% speed-up.

1 Although prefetching may still result in some small energy saving, its average
2 power consumption will be much higher than that of the system running at reference
3 voltage without prefetching. Intuitively, prefetching reduces the execution time but
4 consumes slightly more energy; therefore, the average power (which is the ratio of
5 energy consumption over execution time) will be high. Based on the values reported
6 in Table 3, prefetching consumes 69% more power on average. The reason for this
7 drastic increase is that, as we have mentioned above, more hardware units will be
8 active at the same time due to the prefetching instructions. In the next two sections,
9 we will show how to transfer the performance improvement from prefetching into
10 power and energy saving.

11 5.2. Static voltage scaling

12 In order to transfer the prefetching gains into power and energy savings, we look
13 into both static and DVS. For static voltage scaling, the no prefetching version
14 (NP) was run at the highest voltage, i.e., 1.6 V, with its corresponding frequency
15 of 600 MHz. The prefetching version (P) was run at each different voltage level
16 with the corresponding frequency while we kept the voltage fixed for each run.
17 Table 3 reports the time (in ms) and energy (in $10e + 7 W \cdot cycle$) to complete each
18 application at different voltage levels.

19 From the last five columns, we see that as we reduce the voltage from 1.6 V to
20 1.1 V, the system with prefetching requires more time to complete the application,
21 but consumes less energy (and power). For each application, the column in bold
22 identifies the voltage setting for static voltage scaling that achieves the same (or
23 close) performance as NP, the run at reference 1.6 V voltage without prefetching.
24 These are the closest that we can get with the five different settings of (voltage,
25 frequency) pair on the Crusoe processor. Consider the Moldyn application, which
26 has the maximum prefetching gain of 50.56% (see Fig. 4). We are able to reduce the
27 voltage to 1.25 V and still have an execution time 21.98 ms shorter than the reference
28 (22.23 ms). At this low voltage level, the energy consumption drops to $2.01 \times 10^7 W \cdot$
29 $cycle$ from $7.36 \times 10^7 W \cdot cycle$, a 72.69% saving. The average power reduction
30 is about the same because we have kept the execution time close (21.98 ms and
31 22.23 ms). For the six applications, on an average, prefetching with static voltage
32 scaling converts the 36.04% performance gain to 48.72% energy and average power
33 saving. The power saving and performance change for each application are reported
34 in the last two columns of Table 4.

35 The significance of this static voltage scaling experiment is that it shows the
36 potential of combining voltage scaling and prefetching techniques to reduce system's
37 power and energy consumption. In fact, due to the quadratic relationship between
38 energy consumption and voltage, one can predict that the energy saving converted
39 from "excess" performance by slowing the system down is more than linear. For
40 instance, on the above six benchmarks, 36.04% speed-up becomes 48.72% energy
41 saving. This promising energy and power reduction without performance loss is the
incentive behind our study of the on-line DVS algorithm.

Table 4. Simulation results on on-line profiling based DVS. The third column indicates the performance loss (if negative) or gain in terms of the execution time change. The last two columns give the comparison with the power saving by static voltage scaling as a reference.

Application	DVS-Power Savings %	DVS-Gain in Performance %	SVS-Power Savings %	SVS-Gain in Performance %
IRREG	39.2	+1.30	45.68	+3.26
MOLDYN	65.0	-1.00	72.69	+1.12
NBF	8.00	+13.0	40.00	-6.51
JACOBI	38.5	-1.57	50.01	-6.55
MATMULT	52.2	+3.00	53.11	+7.85
HEALTH	25.8	-2.70	32.32	-2.20

1 5.3. *Dynamic voltage scaling*

2 Although static voltage scaling can give excellent results, it requires the hardware
 3 to be designed with the right voltage to run the particular application. This require-
 4 ment of profiling means that the application is to be run at different voltage levels
 5 and the optimum voltage level determined and stored somehow for each applica-
 6 tion. For a general purpose environment, this is impractical. We propose the on-line
 7 profiling DVS algorithm (Fig. 3) to solve this problem.

8 Table 4 gives the power savings and performance change for different appli-
 9 cations using the proposed on-line profiling DVS algorithm. Both are compared
 10 with the system running at reference voltage without prefetching. On an average
 11 across the six applications, DVS has a 38.11% power savings and is 2.00% faster.
 12 Ideally, we want to convert all the performance gains to energy and power saving.
 13 This may not be possible due to the discrete nature of the five available voltages
 14 and frequencies of the baseline Crusoe processor. In addition, in three of the appli-
 15 cations (MOLDYN, JACOBI, and HEALTH), we have a slight performance loss
 16 (1.0%–2.7%). This is caused by the irregular behavior of these applications as we
 17 have discussed in Sec. 3.3 before Theorem 1. There is also a 13.0% performance
 18 gain in NBF that our algorithm fails to convert to further energy reduction. This
 19 is mainly because of the voltage setting. A detailed explanation of these results
 20 is given next when we analyze the voltage profile of each application. Finally, we
 21 mention that compared with the power saving potential provided by the impracti-
 22 cal static voltage scaling (the fourth column in Table 4), we see that the proposed
 23 on-line profiling based DVS algorithm is effective. In fact, the power saving by DVS
 24 is within 11% on average to the saving by static voltage scaling. This saving is
 25 impressive particularly when one considers the fact that the energy and execution
 26 time overhead for on-line profiling and voltage adjustment has been included in our
 27 simulation.

28 Figure 5 gives the operating voltage level throughout a complete run of each
 29 application following the on-line profiling DVS algorithm. It reveals some interesting
 insights of the proposed algorithm and the application's run-time behavior.

1 As claimed earlier, we expect the DVS algorithm to reach the optimal voltage
 3 level such that the modified code with prefetching will not run slower than the
 5 original code without prefetching. This is possible only when there exists a volt-
 7 age level at which the system can slow down to offset completely the prefetching
 9 gain. Moldyn experiences this steady-state behavior. Table 3 indicates that the
 completion time of Moldyn at 1.25 V (21.98 ms) is very close to that in the non-
 prefetching version (22.23 ms). Therefore, after some initial toggling, DVS algorithm
 finds 1.25 V as the optimal voltage and stays there. The small performance gain is
 canceled by the prefetching overhead.

11 However, none of the other applications have the same “steady-state” behavior.
 13 The main reason is that their required optimal voltage levels lie between the avail-
 15 able “discrete voltage levels.” This leads to the toggling behavior as one can see
 17 in Irreg and MatMult (see Fig. 5). For example, Table 3 suggests that the optimal
 19 voltages for Irreg and MatMult should lie between 1.4 and 1.25 V.

21 Another reason for the toggling behavior is the dynamic behavior of an appli-
 cation, like that of Health. Compared to the earlier regular applications (Irreg,
 Moldyn, and MatMult) that execute the same loop over and over again, Health
 has irregular behavior, spending time in lots of function calls and different loops
 with different prefetching gains. Different prefetching gains lead to a more dynamic
 behavior from our DVS algorithm’s standpoint. From Fig. 5, it is clear that in the
 first half of the execution of the Health application, our algorithm tries to stabilize
 at a voltage level between 1.5 V and 1.4 V, while in the second half, it goes down to

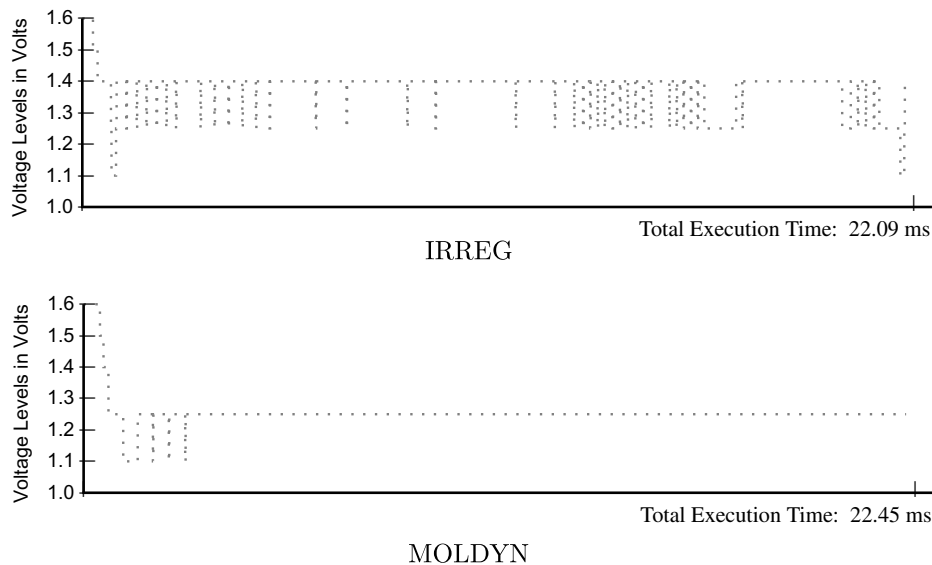


Fig. 5. DVS for different applications, with voltage levels in volts on x -axis and time in ms along the y -axis.

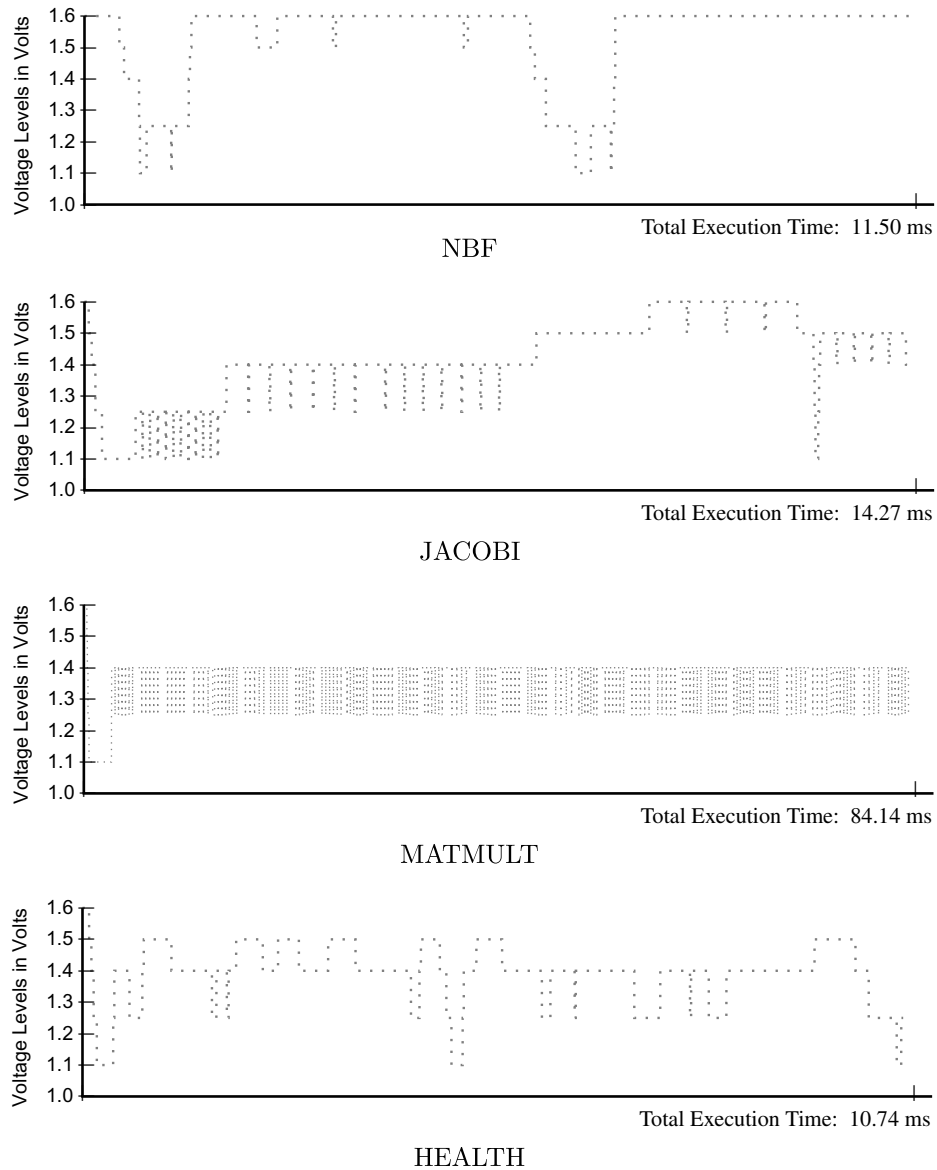
20 *S. N. Pamnani et al.*

Fig. 5. (Continued)

1 between 1.4 V and 1.25 V. We find the main reason is that prefetching has a higher
 2 performance gain in the second half.

3 Similarly, the application Jacobi has a few different “for loops” in the source
 4 code. Our algorithm is able to track the optimum voltage for each “for loop” (by
 5 toggling about it) as the application behavior changes. Finally, NBF has anomalous
 behavior which we are currently investigating.

6. Conclusions and Future Work

We propose a low power system design method that couples DVS and performance-enhancement techniques to simultaneously reduce energy consumption and provide performance guarantees. Specifically, our approach exploits the quadratic power savings of voltage scaling, using the “excess” performance provided by software prefetching to absorb the increased gate delays resulting from voltage scaling. We present an on-line profiling based DVS algorithm that periodically measures the performance gain, and automatically adapts the voltage level to convert such gain to energy reduction. The proposed algorithm is general and applicable to other performance-enhancement techniques.

Our simulation results show that this approach is promising in building low power systems. On six memory-intensive applications, software prefetching provides a 36% performance boost on average. Using the impractical static voltage scaling, this gain has the potential to be converted to 48% average power reduction while still maintaining the same level of performance as the system without software prefetching. Furthermore, we show that our practical on-line DVS algorithm achieves similar energy reduction as the static scheme (within 11% on average) while maintaining approximately the same performance level as the original system (within 6%).

Based on this encouraging result, we believe that the idea of coupling DVS with performance-enhancement techniques is attractive for simultaneously reducing energy consumption and minimizing performance impact in low power systems, and deserves further investigation. Important future work includes coupling DVS with other performance-enhancement techniques, improving our on-line profiling based DVS algorithm to provide further power/energy savings and better performance guarantees, and studying in-depth the cache memory’s energy consumption when the processor is slowed down.

References

1. W. Wulf and S. McKee, Hitting the memory wall: Implications of the obvious. *Computer Archit. News*, **23**(1) (1995) 20–24.
2. T. Chen and J. Baer, Effective hardware-based data prefetching for high-performance processors, *Trans. Comput.* **44**(5) (1995) 609–623.
3. N. P. Jouppi, Improving direct-mapped cache performance by the addition of a small fully-associative cache and prefetch buffers, in *Proc. 17th Ann. Int. Symp. Comput. Archit.* pp. 364–373, Seattle, WA (ACM, 1990).
4. M. Karlsson, F. Dahlgren and P. Stenstrom, A prefetching technique for irregular accesses to linked data structures, in *Proc. 6th Int. Conf. High Performance Comput. Archit.* Toulouse, France, January 2000.
5. T. Mowry, Tolerating latency in multiprocessors through compiler-inserted prefetching, *Trans. Comput. Syst.* **16**(1) (1998) 55–92.
6. T. Mowry and A. Gupta, Tolerating latency through software-controlled prefetching in shared-memory multiprocessors, *J. Parallel Distributed Comput.* **12**(2) (1991) 87–106.

22 S. N. Pamnani et al.

- 1 7. J. Chang and M. Pedram, Energy minimization using multiple supply voltages, *IEEE Trans. Very Large Scale Integration Syst.* **5**(4) (1997) 436–443.
- 3 8. T. Ishihara and H. Yasuura, Voltage scheduling problem for dynamically variable voltage processors, *ISLPED'98: Int. Symp. Low Power Electron. Des.* August 1998, pp. 197–202.
- 5 9. G. Qu, What is the limit of energy saving by dynamic voltage scaling? *IEEE/ACM Int. Conf. Computer-Aided Des.* November 2001, pp. 560–563.
- 7 10. <http://public.itrs.net>, 2001.
- 9 11. C. Chen and M. Sarrafzadeh, Probably good algorithm for low power consumption with dual supply voltages, *IEEE/ACM Int. Conf. Computer Aided Des.* November 1999, pp. 76–79.
- 11 12. K. Govil, E. Chan and H. Wasserman, Comparing algorithms for dynamic speed-setting of a low-power CPU, *ACM Int. Conf. Mobile Comput. Network.* November 1995, pp. 13–25.
- 13 13. T. Pering, T. D. Burd and R. W. Brodersen, Voltage scheduling in the IpARM microprocessor system, *ISLPED'00: Int. Symp. Low Power Electron. Des.* July 2000, pp. 96–101.
- 15 14. M. Weiser, B. Welch, A. Demers and S. Shenker, Scheduling for reduced CPU energy, *USENIX Symp. Operating Syst. Des. Implemen.* November 1994, pp. 13–23.
- 17 15. D. Callahan, K. Kennedy and A. Porterfield Software methods for improvement of cache performance on supercomputer applications, *PhD thesis*, Department of Computer Science, Rice University (May 1989).
- 19 16. D. Callahan, K. Kennedy and A. Porterfield, Software prefetching, in *Proc. 4th Int. Conf. Archit. Support Program. Languages Operating Syst.* April 1991, pp. 40–52.
- 21 17. A. C. Klaiber and H. M. Levy, An architecture for software-controlled data prefetching, in *Proc. 18th Int. Symp. Comput. Archit.* Toronto, Canada (ACM, 1991), pp. 43–53.
- 23 18. E. Gornish, E. Granston and A. Veidenbaum, Compiler-directed data prefetching in multiprocessors with memory hierarchies, in *Proc. Int. Conf. Supercomputing* (ACM, 1990).
- 25 19. F. Yao, A. Demers and S. Shenker, A scheduling model for reduced CPU energy, *FOCS'95: IEEE Ann. Foundations Comput. Sci.* 1995, pp. 374–382.
- 27 20. G. Quan and X. Hu, Energy efficient fixed-priority scheduling for real-time systems on variable voltage processors, *Des. Automat. Conf.* 2001, pp. 828–833.
- 29 21. P. Pillai and G. Shin, Real-time dynamic voltage scaling for low-power embedded operating systems, *Proc. 18th ACM Symp. Operating Syst. Principles* 2001, pp. 89–102.
- 31 22. S. Hua, G. Qu and S. S. Bhattacharyya, Energy reduction techniques for multimedia applications with tolerance to deadline misses, *Design Automation Conf.* 2003, pp. 131–136.
- 33 23. C.-K. Luk and T. C. Mowry, Compiler-based prefetching for recursive data structures, in *Proc. Seventh Int. Conf. Archit. Support for Programming Languages Operating Syst.* Cambridge, MA (ACM, 1996), pp. 222–233.
- 35 24. S. Dhar, D. Maksimovic and B. Kranzen, Closed-loop adaptive voltage scaling controller for standardcell ASICs, *Int. Symp. Low Power Electron. Des.* 2002, pp. 103–107.
- 37 25. C. Simpson, Linear and switching voltage regulator fundamentals, National Semiconductor, www.national.com/appinfo/power/files/f4.pdf.
- 39 26. S. Hua and G. Qu, Voltage set-up problem for embedded systems with multiple voltages, in *IEEE Trans. Very Large Scale Integration (VLSI) Syst.* **13**(7) (2005) 869–872.

- 1 27. M. C. Carlisle and A. Rogers, Software caching and computation migration in olden,
3 in *Proc. 5th ACM SIGPLAN Symp. Principles and Practice of Parallel Programming*,
July 1995, pp. 29–38.
- 5 28. D. Brooks, V. Tiwari and M. Martonosi, Wattch: a framework for architectural-level
power analysis and optimizations, in *ISCA*, 2000, pp. 83–94.
- 7 29. S. J. E. Wilton and N. P. Jouppi, An enhanced access and cycle time model for on-chip
caches, *WRL Research Report*, June 1994.
- 9 30. S. Palacharla, N. P. Jouppi and J. E. Smith, Complexity-effective superscalar proces-
sors, in *ISCA*, 1997, pp. 206–218.
- 11 31. D. Burger and T. M. Austin, The SimpleScalar tool Set, Version 2.0. CS TR 1342,
University of Wisconsin-Madison, June 1997.
- 13 32. C. S. Ellis, The case for higher-level power management, in *Proc. 7th IEEE Workshop
Hot Topics Operating Syst.* 1999, pp. 162–167.
- 15 33. J. Lorch and A. J. Smith, Apple Macintosh’s energy consumption, *Micro* **18**(6) (1998)
54–63.
34. Transmeta-corporation. Tm5400 processor specifications.