

Run-time Energy Consumption Estimation Based On Workload in Server Systems

Adam Wade Lewis, *Member, IEEE*, Soumik Ghosh, *Member, IEEE*, and N.-F. Tzeng, *Fellow, IEEE*

Abstract—This paper proposes a system-wide energy consumption model for servers using hardware performance counters and experimental measurements. We develop a real-time energy prediction model that relates server energy consumption to its overall thermal envelope. While previous studies have attempted system-wide modeling of server power consumption through subsystem models, our approach is different in that it creates a model relating system energy input to subsystem energy consumption based on a small set of tightly correlated parameters. We develop a linear regression model that relates processor power, bus activity, and system ambient temperatures into real-time predictions of the power consumption of long running jobs (and as a result controlling their thermal impact). Using the HyperTransport and Front-Side Bus models as case studies and through electrical measurements on example server subsystems, we develop a statistical model for estimating run-time power consumption. Our model is accurate within an error of four percent (4%) for the HyperTransport bus and nine percent (9%) for the Front-Side Bus as verified using a set of common processor benchmarks.

Index Terms—Analysis of variance, data centers, energy consumption, HyperTransport bus, measurement techniques, modeling techniques, performance counters, thermal envelope.

1 INTRODUCTION

The upwardly spiraling operating costs of the infrastructure for enterprise-scale computing demand efficient power management in server environments. It is difficult in practice to achieve efficient power management as data centers usually over-provision their power capacity to address worst case scenarios. This results in either waste of considerable power budget or severe under-utilization of capacity. Thus, it is critical to quantitatively understand the relationship between power consumption and thermal load at the system level so as to optimize the use of deployed power capacity in the data center.

Power management techniques developed for mobile and desktop computers have been applied with some success to managing the power consumption of microprocessors used in server hardware. The current generation of Intel and AMD processors uses different techniques for processor-level power management including (1) per core clock gating, (2) power-gating functional blocks (for processors to turn off certain blocks that are not in use), (3) multiple clock domains, (4) multiple voltage domains for cores, caches, and memory, (5) dynamic voltage and frequency scaling per core and processor, and (6) hardware support for virtualization techniques. In general, these techniques take advantage of the fact that reducing switching activity within the processor reduces energy consump-

tion and that application performance can be adjusted to utilize idle time on the processor for energy savings [1].

This paper presents a statistical full-system model that provides run-time system-wide prediction of energy consumption on server blades. Our model considers a single server blade as a closed black-box system. The black-box system model lets us converge upon an upper bound of the thermal, energy, and power envelopes of the system. We develop our model by measuring the energy input into the system as a function of the work done by the system in executing its computational tasks and of residual thermal energy given off by the system in doing that work. A hardware performance counter (PeC) based relationship between server blade power consumption and the consequent thermal envelope is necessary to dynamically control the thermal footprint of large workloads. It is important to note that we aim to establish an energy relationship between the workload and the overall thermodynamics of the system.

We begin by measuring the total DC power input to the system at the power supply output. We partition total energy delivered to the system as a sum of the energy components consumed by different subsystems in the server blade. We measure the computational load in the system by observing bus transactions that occur in the system as reported through the PeCs together with a set of system metrics. These metrics are combined into a single model estimated through linear regression.

This work demonstrates that appropriate provision of additional PeCs beyond what are provided by a typical processor is required to obtain more accurate

• A. W. Lewis, S. Ghosh, and N.-F. Tzeng are with the Center for Advanced Computer Studies, University of Louisiana, Lafayette, LA 70504, E-mail: awlewis@cacs.louisiana.edu; sxg5317@cacs.louisiana.edu; tzeng@cacs.louisiana.edu

prediction of system-wide energy consumption. The model takes into account key thermal indicators and system parameters such as ambient temperatures, die temperatures, and hardware performance counters as metrics for system energy consumption estimation within a given power and thermal envelope.

We perform case studies of electrical measurements of server architecture based upon the HyperTransport [2] and Intel NetBurst [3] bus models to develop our model for estimating run-time power consumption. Scheduler-based mechanisms are being developed to take advantage of this estimation model when dispatching jobs to confine server power consumption within a given power budget and thermal envelope while minimizing impact upon server performance.

2 PRIOR WORK

Power models have been used to predict invocation of power management mechanisms in running server systems. These models can be classified into two broad categories: simulation-based models and detailed analytical power models. Although simulations can provide detailed analysis and breakdown of energy consumption, they are usually statically done off-line, are slow, do not scale well, and do not apply favorably to realistic applications and large data sets. In general, neither of these approaches have taken thermal effects of power dissipation into account. Management of system-critical thermal issues due to excessive power consumption, is further complicated by the existence of multiple cores per processor. Existing simulation-based models do not fit well in scenarios where dynamic power and thermal optimization for application performance is required [4].

Analytical models use detailed knowledge of the underlying hardware to directly measure energy consumption at the hardware level. Measurement-based models attempt to collate power measurements to the micro-architectural units on the processor via sampling hardware and software performance metrics. Two distinct classes of metrics have been used in these models: processor performance counters and operating system performance metrics. Processor performance counters are hardware registers that can be configured to count various micro-architectural events, such as branch mispredictions and cache misses. This type of model does not take into account the energy consumption of devices other than the processor and rely upon detailed knowledge of the micro-architecture of the processor. Attempts have been made to reconcile these approaches by attempting to map programs phases to events [5]. The most common technique used to associate PeCs and/or operating systems metrics to energy consumption uses linear regression models to map collected metrics to the energy consumed during the execution of a program [1] [4] [6] [7] [8].

In general, the number of recordable events exceeds the number of available registers. As a result, models that use these counters must time-multiplex different sets of events on the available registers. While this allows for more events to be monitored, it results in increased overhead and lower accuracy due to sampling issues [4] [9]. It is critical that a power model be based upon the smallest possible set of metrics (either a PeC or operating system metric) required to accurately model the system behavior in order to avoid the need for time-multiplexing. High level black-box models sacrifice some accuracy by avoiding extensive detailed knowledge of the underlying hardware. At the processor level, Contreras *et al.* [1], and Bellosa [10] created power models that linearly correlated power consumption with performance counters. Models have been built for the processor, storage devices, single systems, and groups of systems in data centers. These models have the advantage of being simple, fast and low-overhead but they do not model the full-system power consumption.

In server environments, it has been shown that full-system models using operating system CPU utilization can be highly accurate [11]. Others have used similar approaches [12] to develop linear models for energy-aware server consolidation in clusters. Full-system models such as MANTIS [4] [9] relate usage information to the power of the entire system rather than an individual component. Each of these cases requires one or more calibration phases that evaluates the contribution of each system component to overall power consumption. The accuracy and portability of full system power models is considered in Rivoire *et al.* [13]. This analysis indicated that to ensure reasonable accuracy across machines and workload required a model based on a combination of both PeCs and operating system metrics, which directly accounted for all components of a system's dynamic power and provided some insight into memory and disk power consumption.

Extensive study of the power profiles of the Intel Pentium architecture has occurred at the workstation [5] [6] [14] and servers [7] [15] [16]. However, very little consideration has been given to the power profiles of servers constructed using the NUMA-based architecture adopted by processors such as the AMD64 family of processors [17].

3 SYSTEM MODEL DESIGN

We start by considering the type of power supplied into the system. Most server blades operate on an AC input. The DC output from the power supply in our experimental platform is delivered in the domains of +/-12V, +/-5V, and +/-3.3V [18]. In the case of the Sun server used in this study, two 12 Vdc lines supply power to the processor's hard drive(s) and cooling fans in the system. The 5 Vdc and 3.3 Vdc lines are

dedicated to supplying power to the support chips and peripherals on the board. Most switched mode power supplies for servers have a power conversion efficiency from AC to DC of 72 - 80 % , depending on the load of the system. Typically for a server that is idling (running the operating system and no other jobs), the power consumption is between 40 to 42% of the rated power of the system (in our case 450W). The conversion efficiency increases to about 80% when the server is heavily loaded and the SMPS regulates the power supply to work at 75% conversion efficiency at loads over 50% of the rating.

Hence, at the very outset, we lose 20% of the power supplied into the system to conversion losses even for the best conversion factors. Studies have shown [19] that most DC systems perform better in terms of power efficiency than AC systems. A typical AC-based server system has a power supply efficiency of 73% as compared to a 92% for a DC system. Also, the overall system efficiency for a AC system is 61% as compared to 85% for a DC system.

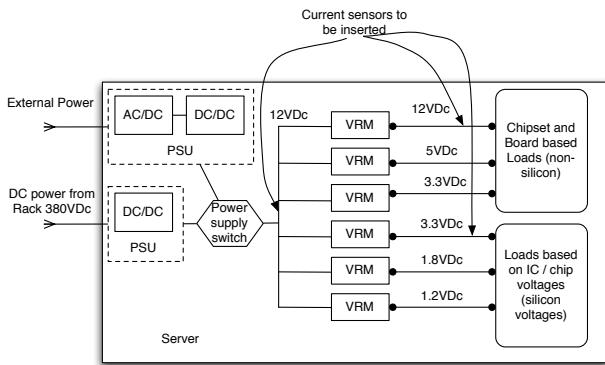


Fig. 1. Power distribution model for the server blade.

While a rack level DC power distribution system easily translates into large power savings at the server level, our model considers a part of this power conversion unit as it is not software tunable in its present state. On the other hand, current sensors, such as MAXIM's 4473 [20], placed at the outputs from the power supply would immensely aid in dynamically tracking DC power draw which varies according to the system load into the system. A proposed system diagram for a combined AC and DC power-supply based system with current sensors as measurable performance counters, along with power distribution, is shown in Fig. 1. In the interim, our model externally monitors the input power to the system and controls the power, and consequently, thermal envelope of the system based on the processing load.

In order to develop an energy consumption model based on computational load of the system, we begin by measuring the total DC power input to the system, at the output of the SMPS. As mentioned earlier, the DC power is delivered in +/-12 V, +/-5V, and +/-

3.3V power domains. Most SMPS limit the total power delivered through the 5V and 3.3V lines to about 20% of the rated power supply (P_R). Now assuming each of the voltage lines $v_k(t)$ draws current $i_k(t)$, then each line draws an instantaneous power $p_k(t) = v_k(t) \cdot i_k(t)$. If a voltage domain has M DC lines as output, then total power delivered for that voltage domain is:

$$p_{v1}(t) = \sum_{k=1}^M v_k(t) \cdot i_k(t)$$

If the board has N voltage domains: v_1, v_2, \dots, v_N , then the total DC power delivered into the system is:

$$p_{dc}(t) = \sum_{j=1}^N p_{vj}(t) = \sum_{j=1}^N \sum_{k=1}^{M_j} v_k(t) \cdot i_k(t).$$

So total energy delivered to the system between times t_2 to t_1 is:

$$E_{dc} = \int_{t_1}^{t_2} p_{dc}(t) dt = \int_{t_1}^{t_2} \sum_{j=1}^N \sum_{k=1}^{M_j} v_k(t) \cdot i_k(t) dt. \quad (1)$$

For the 3.3V and 5V lines, we have the following constraint:

$$\begin{aligned} E_{dclv} &= \int_{t_1}^{t_2} p_{dclv}(t) dt \\ &= \int_{t_1}^{t_2} \left(\sum_{k=1}^{M1} v_k(t) \cdot i_k(t) + \sum_{k=1}^{M2} v_k(t) \cdot i_k(t) \right) dt \\ &\leq 0.2 P_R \end{aligned} \quad (2)$$

where $M1$ and $M2$ are the total 3.3V and 5V lines, respectively. Thus in our 450W rated system the power delivered by the 3.3V and 5V lines is capped at 90W.

This energy delivered to the system $E_{dc} = E_{system}$ can now be expressed as a sum of energy consumed by the different sub-systems in the server blade. Broadly we define five sources of energy consumption within a system:

- E_{proc} : Energy consumed in the processor due to all computations,
- E_{mem} : Energy consumed in the DRAM chips,
- E_{hdd} : Energy consumed by the hard disk drive during the server's operation.
- E_{em} : Energy consumed by all electrical and electromechanical components in the server blade including fans and other server components which consume AC power,
- E_{board} : Energy consumed by peripherals that support the operation of the board. These include all devices in the multiple voltage domains across the board, including chipset chips, voltage regulation, bus control chips, connectors, interface devices, etc.,

The total energy consumed by the system for a given computational workload can be calculated as the sum of the five terms:

$$E_{system} = E_{proc} + E_{mem} + E_{hdd} + E_{board} + E_{em} \quad (3)$$

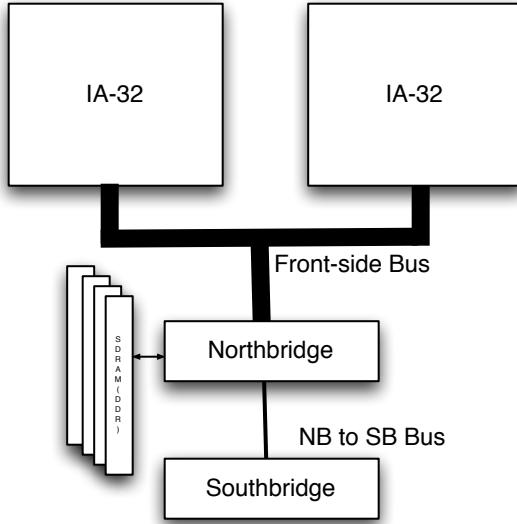


Fig. 2. Intel Pentium P6 server architecture.

We explore each of these terms in turn by following an energy conservation principle in the system. In order to get a true measure of the computational load on the system, our method looks to snoop on completed bus transactions per unit time in the system and measure the relative change in energy consumption (as indicated by change in temperature) as computation tasks are completed. Use of this performance counter metric as compared to other metrics fits well with the architecture of microprocessors used in NUMA-based processors in multi-core environments.

Consider the Intel Pentium P6 and AMD Opteron processor architectures connected in a dual core configuration shown in Figs. 2 and 3. The Pentium architecture (and its successors up to the recent Intel Nehalem processor architecture) employ a Front-Side-Bus (FSB) to connect individual cores to the Northbridge chip on the motherboard as depicted in Fig. 2. This chip provides the interface between the cores and memory. A coherent bus protocol is used to ensure consistency in memory access between the cores. For this architecture, the FSB becomes a performance bottleneck as processor cores have to moderate themselves to the slower speed of the bus.

In contrast, the AMD Opteron processor is based on the NUMA-based architecture (Fig. 3 where the Northbridge functionality is incorporated in the processor core and each core is responsible for local access to the memory connected to that Northbridge logic (shown in Fig. 3 as "Integrated Memory Controller"). Processor cores on a single die are connected via a crossbar to the Hypertransport bus (i.e., HT1) between processors. Again, a coherent bus protocol is used to ensure memory consistency between processor cores on each die. In addition, the master processor in the

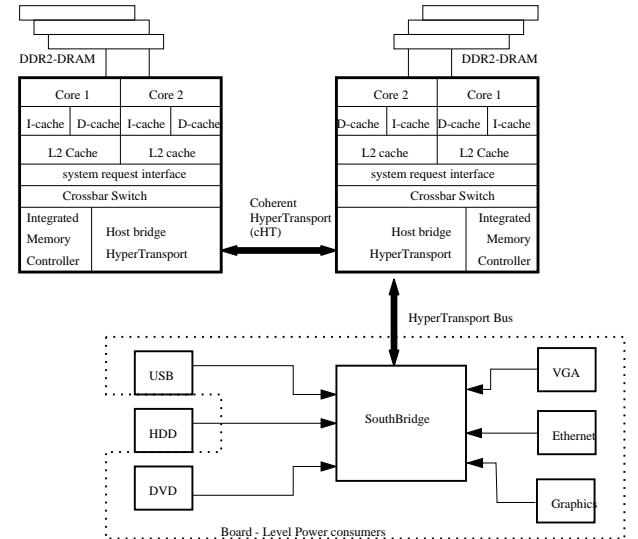


Fig. 3. AMD Opteron server architecture.

system is connected via a second Hypertransport bus (i.e., HT2) to the Southbridge device that manages connections to the outside world.

It is observed that the work done by any of these processors, as the heart of energy consumption in a server system, can be quantified in terms of bus transactions in and out of the processors. Traffic on the external buses provides a measure of how much data is being processed by the processor and what would be an upper limit of the work done by a single processor.

3.1 Processor energy consumption

Our processor model aims to treat each processor as a black box, whose energy consumption is a function of its work load as manifested by the core die-temperature and system ambient temperature (measured at a system level by ipmitool through sensors in the path of the outgoing airflow from the processor). A practical issue with trying to estimate processor power using a large number of performance counters (PeCs) is that there are only a limited number of PeCs that tools like cpustat can track simultaneously. In order to track the energy-thermal load relationship for a job, we had to develop a model with the least number of PeCs that would accurately reflect the energy consumption-thermal load relationship.

Given the AMD Opteron processor architecture connected in a dual core configuration shown in Fig. 3, we consider traffic on the HyperTransport buses as a representative of the processor work load and a reflection of the amount of data being processed by a processor or any of its cores. The HT2 bus is non-coherent and connects one of the two processors to the Southbridge (whereas the Northbridge is included

on the Opteron processor die). Thus, traffic on the HT2 bus reflects hard-disk and network traffic. The model therefore scales when considering the effect of network traffic and disk I/O operations. The HT1 bus is a coherent bus between the two SMP processors and, as such, PeCs on that bus provides an accurate reflection on the processing load of cores executing jobs. Per-core die temperature readings and, consequently, ambient temperature per processor are thus greatly affected by the number of transactions over the HT buses. A similar analysis can be applied to the Front-Side-Side bus used in the Intel NetBurst processor architecture. We also include L2 cache misses as one of our variables (to be explained in Section 3.2).

Thus the total processor power consumption to reflect the thermal change due to workload can be expressed as:

$$\mathbf{P}_{\text{proc}} = \mathbf{H} \cdot \mathbf{X} = [\beta_0 \dots \beta_{10}]^T \cdot [\text{Var}_0 \dots \text{Var}_{10}]^T$$

where \mathbf{X} vector contains the following variables: ambient temperatures and die temperatures for processors 0 and 1, HT1 and HT2 transactions, and last-level cache misses per core. From this equation, we can determine the total processor energy consumption between times t_1 and t_2 by integration:

$$E_{\text{proc}} = \int_{t_1}^{t_2} \{\mathbf{P}_{\text{proc}}\} dt \quad (4)$$

3.2 DRAM energy

Energy consumed by the DRAM banks can be computed by a combination of measuring the counts of the last level cache misses (CM_i) for each of the N cores in the system together with the DRAM Read/Write power (P_{DRAM}) along with the DRAM background power(activation power) (P_{activ}):

$$E_{\text{mem}} = \int_{t_1}^{t_2} \left(\sum_{i=1}^N CM_i \times P_{\text{DRAM}} + P_{\text{activ}} \right) dt.$$

As illustrated in [21], DRAM background power and activation power can be obtained from the DRAM documentation. For a single DRAM in our case, a total of 493mW would be consumed. However, given the number of L2 cache misses per second when a job is running on a certain core (over 22 Million per second at the peak of bzip2 SPEC2006 benchmark on the AMD Opteron processor), a significant amount of heat is generated from the DRAM chips. The thermal airflow proximity of the DRAM banks to their respective processors makes it possible for us to combine the energy consumption and the consequent thermal output of the memory banks with the processor ambient temperature. This value is reported by system management tools commonly found in server environments (such as IPMI) and is incorporated into our model as a contribution to the processor ambient temperature.

3.3 Hard disk energy

The energy consumed by the hard disk while operating can be approximated to its energy consumption upper bound using a combination of performance counters and drive ratings. In our server, Hitachi's 7200 RPM 250G, SATA hard disk is used. Table 1 lists the typical power consumption numbers for the hard disk used. Based on the physical, electrical and electromechanical parameters of the hard disk, very detailed power consumption models can be constructed. However, we can achieve a cruder but simpler model based on the typical power consumption data of the hard disk and performance counters.

TABLE 1
Hitachi HDT725025VLA360 disk power parameters

Parameter	Value
Interface	Serial ATA
Capacity	250 GB
Rotational speed	7200 rpm
Power	
Spin up	5.25 W (max)
Random read, write	9.4 W (typical)
Silent read, write	7 W (typical)
Idle	5 W (typical)
Low RPM idle	2.3 W (typical) for 4500 RPM
Standby	0.8 W (typical)
Sleep	0.6 W (typical)

The utility `iostat` can be used to measure the number of read and writes per second to the disk as well as the kilobytes read from and written to the disk. Thus, based on this performance counter, we can compute an approximate disk power consumption E_{hdd} value as:

$$E_{\text{hdd}} = P_{\text{spin-up}} \times T_{\text{su}} + P_{\text{read}} \sum N_r \times T_r + P_{\text{write}} \sum N_w \times T_w + \sum P_{\text{idle}} \times T_{\text{id}} \quad (5)$$

where $P_{\text{spin-up}}$ is the power required to spin-up the disk from 0 to full rotation. T_{su} is the time required to achieve spin up. T_{su} is typically about 10s. P_{read} is the power consumed per kilobyte of data read from the disk. N_r is the number of kilobytes of data read in time-slice T_r from the disk. The variables are analogous for the write energy consumption. P_{read} for our Hitachi disk can be computed as follows: The read operation at 1.5 Gbits/s consumes 530 mA current at +5V. Hence every kilobyte read, consumes approximately $13.3\mu\text{W}/\text{Kbyte}$. Similarly, every write operation consumes $6.67\mu\text{W}/\text{Kbyte}$. The numbers N_r and N_w can be obtained using `iostat` and our choice of time-slice.

The idle state has two conditions, idle and unloaded idle, where in the latter case the heads are unloaded. The time to go from unloaded to idle is usually less than 1 second, which is less than the resolution of `iostat`. Thus, a history match count in the `iostat`

statistics where the reads and writes have been zero, tells us the period in which the disk is idle, and the idle energy consumption can be computed accordingly. `iostat` reading based conditions for switching to different disk power states can be obtained with more in-depth analysis, but the net results fall into this equation's framework. Such an analysis is the topic of a future work.

The hard disk power can also be measured in real-time if current sensors are provided at the output of the DC voltage lines delivering power to the hard disk drives. The +5V lines will draw a maximum current of 730mA and the +12V lines will draw a maximum current of 630mA. Thus E_{hdd} can also be formulated as:

$$E_{hdd} = \int_{t1}^{t2} \{v_1(t) \times i_1(t) + v_2(t) \times i_2(t)\} dt. \quad (6)$$

This approach can be applied in the presence of current sensors which in our experiments was measured with a current probe and logged through an oscilloscope.

3.4 Board components

The quantity E_{board} represents energy required by the support chipsets and usually falls in the 3.3V and 5V power domains. In our case, this value is obtained using current-probe based measurements. However, as in earlier cases, current sensors for the power lines going into the board can provide instantaneous energy draw from the power supply. For our server, at most 28 additional current sensors might be required for the entire blade [18]. The processor, disk, fan, and optical-drive power lines are excluded here which leads to:

$$E_{board} = \left(\sum V_{power-line} \times I_{power-line} \right) \times t_{time-slice}. \quad (7)$$

3.5 Electromechanical energy

There is a basic electrical cost related to running the computer. The quantity E_{em} in our model takes that into account. E_{em} is calculated as the summation of the DC and AC power consumption in the peripherals supporting the processor, particularly the electromechanical components. This is due to the power consumption in the power supply unit and the power consumption in the cooling fans. As discussed earlier, the AC to DC conversion process is a load based number, exhibiting only a best case conversion efficiency of 80% and a nominal efficiency of 75% .

Power drawn by the fans for cooling can be given by the following equation:

$$P_{fan} = P_{base} \cdot \left(\frac{RPM_{fan}}{RPM_{base}} \right)^3 \quad (8)$$

where P_{base} defines the base power consumption of the unloaded system. In our case, it refers to the power consumption of the system when running only the base operating system and no other jobs. That value is obtained experimentally by measuring the current drawn on the +12V and +5V lines, using a current probe and an oscilloscope. There is a current surge at system start, which is neglected. Under nominal conditions, the +12V line draws approximately 2.2A, which powers both blower fans in the system. The two peripheral fans running at +5V draw around 2.1A of current. Thus, the base power for the fans is known and it is possible to quantify the electrical power consumption in the system. From this, the electrical power consumption can be quantified as:

$$P_{elect} = V(t) \cdot I(t) + \sum_{i=1}^N P_i \quad (9)$$

where the first term in the equation is the instantaneous DC power output from the power supply and is the DC power consumed by the system. N is the number of fans in the server and P_i is the instantaneous power consumed by the $i - th$ fan according to Eq. 8. Total energy consumption during a given task period T_p due to electrical energy in the system can now be given by:

$$E_{em} = \int_0^{T_p} [V(t) \cdot I(t) + \sum_{i=1}^N P_i] dt. \quad (10)$$

4 APPLICATION OF MODEL TO PHYSICAL SYSTEM

Our model is based on application data driven activity and data flow across the server system. However, the current generation of server systems lack the complete set of measurement and monitoring capabilities and data flow state capture mechanisms required in order to formulate an exact analytical model. For example, the system board DC and AC power consumption cannot be easily split in terms of measurements or exact analytical models due to the presence of large numbers of voltage and current domains and components. However, the concept of AC and DC power consumption on the board can be captured as voltage and current domain based power summation in through external measurement. Thus, to use our model as a predictive tool requires that we introduce an element of statistical approximations to compensate for the cases where we cannot exactly determine components in our model.

4.1 Case study 1: HyperTransport bus

The first case study evaluates the behavior of our model on a AMD Opteron server using the HyperTransport bus (as detailed in Table 5). Physical predictors were selected for each term in the model

TABLE 2
Overall regression model for AMD Opteron Processor

Coeff.	Variable	Measurement
β_0		
β_1	T_{A_0}	Ambient Temp0
β_2	T_{A_1}	Ambient Temp1
β_3	T_{C_0}	CPU0 Die Temp
β_4	T_{C_1}	CPU1 Die Temp
β_5	HT_1	HT1 Bus X-Actions
β_6	HT_2	HT2 Bus X-Actions
β_7	CM_0	L1/L2 Cache Miss for Core0
β_8	CM_1	L1/L2 Cache Miss for Core1
β_9	CM_2	L1/L2 Cache Miss for Core2
β_{10}	CM_3	L1/L2 Cache Miss for Core3
β_{11}	D_r	Disk bytes read
β_{12}	D_w	Disk bytes written
β_{13}	F_C	CPU Cooling Fan Speed
β_{14}	F_M	Memory Cooling Fan Speed

shown in Eq. 3. Each predictor listed in Table 2 contributes to one of the terms in our combined model. Note the importance of the linear estimation coefficients in our statistical model. It is the units associated with each coefficient that is responsible for converting each quantity to map the units associated with each variable into energy units.

For example, we predict that the energy consumed in the process of computation E_{proc} as defined by Eq. 3.1 can be estimated by a linear combination of the CPU Die Temperatures and the amount of data transmitted across the HyperTransport bus connecting the physical cores:

$$E_{proc} \approx \beta_3 * T_{C_0} + \beta_4 * T_{C_1} + \beta_5 * HT_1. \quad (11)$$

These terms in this equation provide direct estimates for each of the quantities in Eq. 4.

In a similar fashion, the energy consumed in the memory system E_{mem} can be estimated by a combination of the amount of data transmitted across the HyperTransport bus between processor and the Southbridge and the number of last level cache misses for each physical core:

$$E_{mem} \approx \beta_6 * HT_2 + \beta_7 * CM_0 + \beta_8 * CM_1 + \beta_9 * CM_2 * \beta_{10} * CM_3.$$

By combining this expression for E_{mem} with the expression for E_{proc} from Eq.11, we can approximate the combined terms for processor and memory energy consumption in our model (Eq.3).

The energy consumed as result of disk activity (Eq. 5) is estimated by measuring the number of bytes read and written to the desk

$$E_{hdd} \approx \beta_{11} * D_r + \beta_{12} * D_w.$$

We estimate the energy consumed by the electrical components in the server blade as a linear combination of the ambient temperature as measured in the air flow moving over each physical processor:

$$E_{em} \approx \beta_1 * T_{A_0} + \beta_2 * T_{A_1}.$$

The energy consumed as result of the peripheral operation of the board E_{board} as defined in Equation 7 is estimated as a linear combination of the operating speeds of the fans cooling the CPU and memory plus the linear estimation constant term:

$$E_{board} \approx \beta_0 + \beta_{13} * F_C + \beta_{14} * F_M.$$

The predictors are combined together into a linear regression model:

$$\begin{aligned} E_{system} \approx & \beta_0 + \beta_1 * T_{A_0} + \beta_2 * T_{A_1} \\ & + \beta_3 * T_{C_0} + \beta_4 * T_{C_1} + \beta_5 * HT_1 \\ & + \beta_6 * HT_2 + \beta_7 * CM_0 + \beta_8 * CM_1 \\ & + \beta_9 * CM_2 + \beta_{10} * CM_3 + \beta_{11} * D_r \\ & + \beta_{12} * D_w + \beta_{13} * F_C + \beta_{14} * F_M. \end{aligned}$$

This model is fit to a collection of representative benchmarks from the SPEC CPU2006 benchmark suite [22] as listed in Table 3. The benchmarks were selected using two criteria: sufficient coverage of the functional units in the processor and reasonable applicability to the problem space. Components of the processor affect the thermal envelope in different ways [23]. This issue is addressed by balancing the benchmark selection between integer and floating point benchmarks in the SPEC CPU2006 benchmark suite. Second, the benchmarks were selected from the suite based upon fit into the problem space. Each benchmark represents an application typical of the problems solved on high-performance application servers.

TABLE 3
SPEC CPU2006 benchmarks used for model calibration

Integer Benchmarks		
perlbench	C	PERL Programming Language
bzip2	C	Compression
mcf	C	Combinatorial Optimization
omnetpp	C++	Discrete Event Simulation
FP Benchmarks		
gromacs	C/F90	Biochemistry/Molecular Dynamics
cactusADM	C/F90	Physics/General Relativity
leslie3d	F90	Fluid Dynamics
lbm	C	Fluid Dynamics

Five classes of metrics are sampled at 5-second intervals during the experiment: (1) CPU temperature for all processors in the system, (2) ambient temperature in the computer case measured in one or more locations using the sensors provided by server manufacturer, (3) the number of completed transactions processed through the system bus, and (4) the number of misses that occur in the L2 cache associated with each CPU core in the system.

Two methods were considered for consolidation: arithmetic mean (average) and geometric mean. Trial models were constructed using each method and

a statistical analysis of variance was performed to determine which model generated the best fit to the collected data. Fig. 5 provides a visual comparison of the system power as measured on the SUT device versus the predicted power consumption from each device. For a practical usage scenario the statistical coefficients need to be computed only once using the SPEC benchmarks for a given server architecture. They can be used as embedded constants available either through the system firmware or the operating system kernel. A linear regression model was created from the consolidated data set to generate the parameters described in our theoretical model:

$$E_{system} \approx 32.71 + 1.31 * T_{A_0} + 0.54 * T_{A_1} + 0.54 * T_{C_0} \\ + 0.61 * T_{C_1} + 0.01 * HT_1 + 0.01 * HT_2 + 0.01 * CM_0 \\ + 0.01 * CM_1 + 0.01 * CM_2 + 0.01 * CM_3 + 0.50 * D_r \\ + 0.50 * D_w + 0.01 * F_C + 0.01 * F_M.$$

Fig. 5 compares the predicted power consumption from the model versus the actual power consumption for each benchmark. Table 8 shows the descriptive statistics and percentage errors for benchmark predictions.

4.2 Case study 2: Front-Side Bus

Our second case study considers the application of the model to the Intel Xeon Woodcrest processor. This processor is a dual-core processor based on the Intel Core architecture with linked L1-cache and a shared last-level L2-cache. The bus counters used for constructing our statistical model for the HyperTransport bus were replaced with the Intel performance counters that measure the transactions that occur on the Front-Side Bus on Xeon Woodcrest processors: BUS_TRANS_ANY, BUS_TRANS_MEM, and BUS_TRANS_BURST. Coefficients for the revised regression estimator are shown in Table 9. Our general linear regression model is adapted to this architecture by combining the beta coefficients and variables together as follows:

$$E_{system} \approx \beta_0 + \beta_1 T_{C_0} + \beta_2 T_{C_1} + \beta_3 T_{A_0} \\ + \beta_4 T_{A_1} + \beta_5 T_{A_2} + \beta_6 FSB_1 + \beta_7 CM_0 \\ + \beta_8 D_r + \beta_9 D_w + \beta_{10} F_C + \beta_{11} F_{M2a} \\ + \beta_{12} F_{M2b} + \beta_{13} F_{M3a} + \beta_{14} F_{M3b} \\ + \beta_{16} F_{M4a} + \beta_{17} F_{M4b} + \beta_{18} F_{M5a} + \beta_{19} F_{M5b}.$$

The process described in Section 4 was repeated and a new model was created from the consolidated data set to generate the parameters described in our theoretical

TABLE 4
Regression model as revised for the Intel architecture

Coeff.	Variable	Measurement
β_0		
β_1	T_{C_0}	CPU0 Die Temp
β_2	T_{C_1}	CPU1 Die Temp
β_3	T_{A_0}	Ambient Temp0
β_4	T_{A_1}	Ambient Temp1
β_5	T_{A_2}	Ambient Temp2
β_6	FSB_1	Front Side Bus Transactions
β_7	CM_0	L1/L2 Cache Miss for Core0
β_8	D_r	Disk bytes read
β_9	D_w	Disk bytes written
β_{10}	F_C	Memory Cooling Fan Speed
β_{11}	F_{M2a}	Memory Cooling Fan Speed 2a
β_{12}	F_{M2b}	Memory Cooling Fan Speed 2a
β_{13}	F_{M3a}	Memory Cooling Fan Speed 3a
β_{14}	F_{M3b}	Memory Cooling Fan Speed 3b
β_{16}	F_{M4a}	Memory Cooling Fan Speed 4a
β_{17}	F_{M4b}	Memory Cooling Fan Speed 4b
β_{18}	F_{M5a}	Memory Cooling Fan Speed 5a
β_{19}	F_{M5b}	Memory Cooling Fan Speed 5b

model:

$$E_{system} \approx 2.53 + 2.29 * T_{C_0} + 0.03 * T_{C_1} \\ + 0.03 * T_{A_0} + 0.01 * T_{A_1} + 0.01 * T_{A_2} \\ + 0.52 * FSB_1 + 0.35 * CM_0 + 0.01 * D_r \\ + 0.01 * D_w + 4.85 * F_C + 6.61 * F_{M2a} \\ + 3.92 * F_{M2b} + 0.28 * F_{M3a} + 0.52 * F_{M3b} \\ + 0.01 * F_{M4a} + 0.78 * F_{M5a} + 0.61 * F_{M5b}.$$

The fit of this model to the consolidated data was done through an Analysis of Variance (ANOVA) and related statistical tests of the parameters and overall model fit (shown in Table 10). An adjusted R-Square value of 0.98 gives an projected best-case error for the model of 11.0 percent with a 95% confidence level.

TABLE 5
Test hardware configuration

	Sun Fire 2200	Dell PowerEdge SC1425
CPU	2 AMD Opteron	2 Intel Xeon (NetBurst) 5110
CPU L2 cache	2x2MB	4MB
Memory	8GB	9GM
Internal disk	2060GB	500GM
Network	2x1000Mbps	1x1000Mbps
Video	On-board	NVIDIA Quadro FX4600
Height	1 rack unit	Mini-tower

5 EVALUATION

Case studies of the use of our power model to evaluate energy consumption on a test system were created for the hardware described in Table 5. The first case study evaluates the behavior of our model on a AMD Opteron server using the HyperTransport bus while the second evaluates model behavior on a Intel Xeon

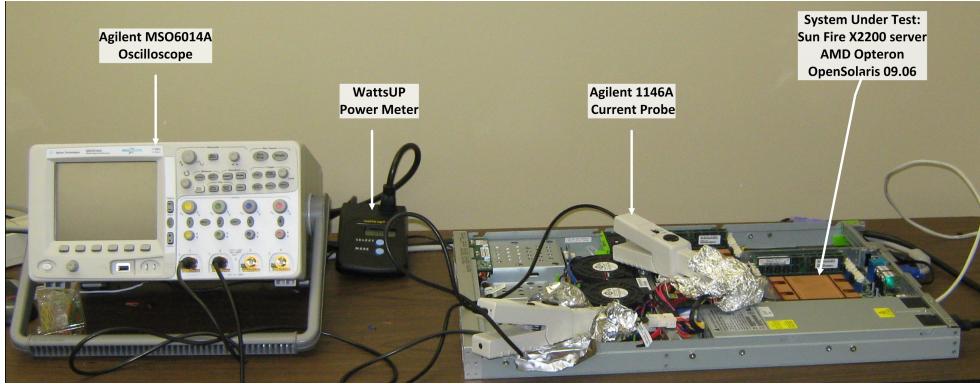


Fig. 4. Hardware test setup.

processor that uses the NetBurst micro-architecture and a Front-Side Bus. Four additional benchmarks from the SPEC CPU2006 benchmark suite (shown in Table 6) were executed to evaluate the predictive performance of the model.

TABLE 6
SPEC CPU2006 benchmarks used for evaluation

Integer Benchmark		
astar	C++	Path Finding
gobmk	C	Artificial Intelligence: Go
FP Benchmarks		
calculix	C++/F90	Structural Mechanics
zeusmp	F90	Computational Fluid Dynamics

5.1 Evaluation environment

The operating system used in our setup is OpenSolaris (Solaris 11). System data is collected from the system baseboard controller using the IPMI interface via the OpenSolaris ipmitool utility. Processor performance counters are collected on a system-wide basis using the OpenSolaris cpustat utility.

In terms of measuring performance counters, we have used the OpenSolaris cpustat, iostat, and ipmitool utilities. Of these, iostat and ipmitool are available across all UNIX-based operating systems commonly used in data centers. cpustat is an OpenSolaris specific utility but is already being ported to Linux. In future work, it is planned to use tools like dtrace and oprofile for more controllable and tunable performance parameters which have major impacts on system-wide and processor wide power consumption.

The power consumed is measured with a WattsUP power meter [24] connected between the AC Main and the system under test (SUT). The power meter measures the total and average wattage, voltage, and amperage over the run of a workload. The internal memory of the power meter is cleared at the start of

the run and the measures collected during the run are downloaded after the run completes from the meter's internal memory into a spreadsheet. Current flow on the different voltage domains in the server is measured using an Agilent MSO6014A oscilloscope with one Agilent 1146A current probe per system power domain (12v, 5v, and 3.3v). This data is collected from the oscilloscope at the end of the execution of a benchmark and stored in a spreadsheet on the test host.

5.2 Case study 1: HyperTransport bus

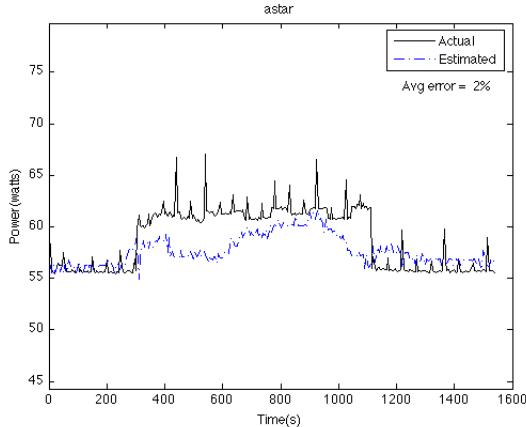
The first case study evaluates the behavior of our model on a AMD Opteron server using the HyperTransport bus (as detailed in Table 5). A linear regression model was created from the consolidated data set to generate the parameters described in our theoretical model:

$$\begin{aligned} E_{system} \approx & 32.71 + 1.31 * T_{A_0} + 0.54 * T_{A_1} + 0.54 * T_{C_0} \\ & + 0.61 * T_{C_1} + 0.01 * HT_1 + 0.01 * HT_2 + 0.01 * CM_0 \\ & + 0.01 * CM_1 + 0.01 * CM_2 + 0.01 * CM_3 \\ & 0.50 * D_r + 0.50 * D_w + 0.01 * F_C + 0.01 * F_M. \end{aligned}$$

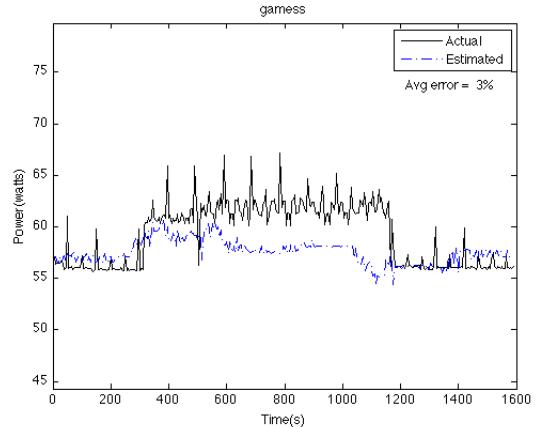
Fig. 5 compares the predicted power consumption from the model versus the actual power consumption for each benchmark. Table 8 shows the descriptive statistics and percentage errors for benchmark predictions. The fit of this model to the consolidated data was done through an Analysis of Variance (ANOVA) and related statistical tests of the parameters and overall model fit (shown in Table 7). An adjusted R-Square value of 0.965 gives an projected error for the model of 3.5 percent with a 95% confidence level.

5.3 Case study 2: Front-Side Bus

Our second case study considers the application of model to the Intel Xeon Woodcrest processor. This processor is a dual-core processor based on the Intel Core architecture with linked L1-cache and a shared last-level L2-cache. The bus counters used



(a) Astar.



(b) Gamess.

Fig. 5. Actual versus predicted energy consumption for AMD Opteron Processors.

TABLE 7
ANOVA for AMD Opteron processor

Source	df	SS	MS	F	P
Regr	10.00	2261.17	226.12	1150.23	0.00
Resid	399.00	78.44	0.20		
Total	409.00	2339.60			

R-sq	Adj. R-sq	0.97
------	-----------	------

TABLE 8
Model error for each benchmark

Benchmark	Mean	Median	Std	Var	Error
astar	1.12	0.98	0.95	0.90	1.9%
gamess	1.38	1.19	1.24	1.53	2.3%
gobmk	2.30	2.20	0.82	0.67	3.9%
zeusmp	1.84	1.76	1.13	1.28	3.1%

for constructing our statistical model for the Hyper-Transport bus were replaced with the Intel performance counters that measure the transactions that occur on the Front-Side Bus on Xeon Woodcrest processors: BUS_TRANS_ANY, BUS_TRANS_MEM, and BUS_TRANS_BURST. Coefficients for the revised regression estimator are shown in Table 9. Our general linear regression model is adapted to this architecture by combining the beta coefficients and variables together as follows:

$$\begin{aligned}
 E_{system} \approx & \beta_0 + \beta_1 T_{C_0} + \beta_2 T_{C_0} + \beta_3 T_{A_0} \\
 & + \beta_4 T_{A_1} + \beta_5 T_{A_2} + \beta_6 FSB_1 + \beta_7 CM_0 \\
 & + \beta_8 D_r + \beta_9 D_w + \beta_{10} F_C + \beta_{11} F_{M2a} \\
 & + \beta_{12} F_{M2b} + \beta_{13} F_{M3a} + \beta_{14} F_{M3b} \\
 & + \beta_{16} F_{M4a} + \beta_{17} F_{M4b} + \beta_{18} F_{M5a} + \beta_{19} F_{M5b}.
 \end{aligned}$$

The process described in Section 4 was repeated and a new model was created from the consolidated data set to generate the parameters described in our theoretical

TABLE 9
Regression model as revised for the Intel architecture

Coeff.	Variable	Measurement
β_0		
β_1	T_{C_0}	CPU0 Die Temp
β_2	T_{C_1}	CPU1 Die Temp
β_3	T_{A_0}	Ambient Temp0
β_4	T_{A_1}	Ambient Temp1
β_5	T_{A_2}	Ambient Temp2
β_6	FSB_1	Front Side Bus Transactions
β_7	CM_0	L1/L2 Cache Miss for Core0
β_8	D_r	Disk bytes read
β_9	D_w	Disk bytes written
β_{10}	F_C	Memory Cooling Fan Speed
β_{11}	F_{M2a}	Memory Cooling Fan Speed 2a
β_{12}	F_{M2b}	Memory Cooling Fan Speed 2a
β_{13}	F_{M3a}	Memory Cooling Fan Speed 3a
β_{14}	F_{M3b}	Memory Cooling Fan Speed 3b
β_{16}	F_{M4a}	Memory Cooling Fan Speed 4a
β_{17}	F_{M4b}	Memory Cooling Fan Speed 4b
β_{18}	F_{M5a}	Memory Cooling Fan Speed 5a
β_{19}	F_{M5b}	Memory Cooling Fan Speed 5b

model:

$$\begin{aligned}
 E_{system} \approx & 2.53 + 2.29 * T_{C_0} + 0.03 * T_{C_0} \\
 & + 0.03 * T_{A_0} + 0.01 * T_{A_1} + 0.01 * T_{A_2} \\
 & + 0.52 * FSB_1 + 0.35 * CM_0 + 0.01 * D_r \\
 & + 0.01 * D_w + 4.85 * F_C + 6.61 * F_{M2a} \\
 & + 3.92 * F_{M2b} + 0.28 * F_{M3a} + 0.52 * F_{M3b} \\
 & + 0.01 * F_{M4a} + 0.78 * F_{M5a} + 0.61 * F_{M5b}.
 \end{aligned}$$

The fit of this model to the consolidated data was done through an Analysis of Variance (ANOVA) and related statistical tests of the parameters and overall model fit (shown in Table 10). An adjusted R-Square value of 0.98 gives an projected best-case error for the model of 11.0 percent with a 95% confidence level.

The four benchmarks described in Table 6 were executed to evaluate the predictive performance of the model. We compare the predicted power consumption

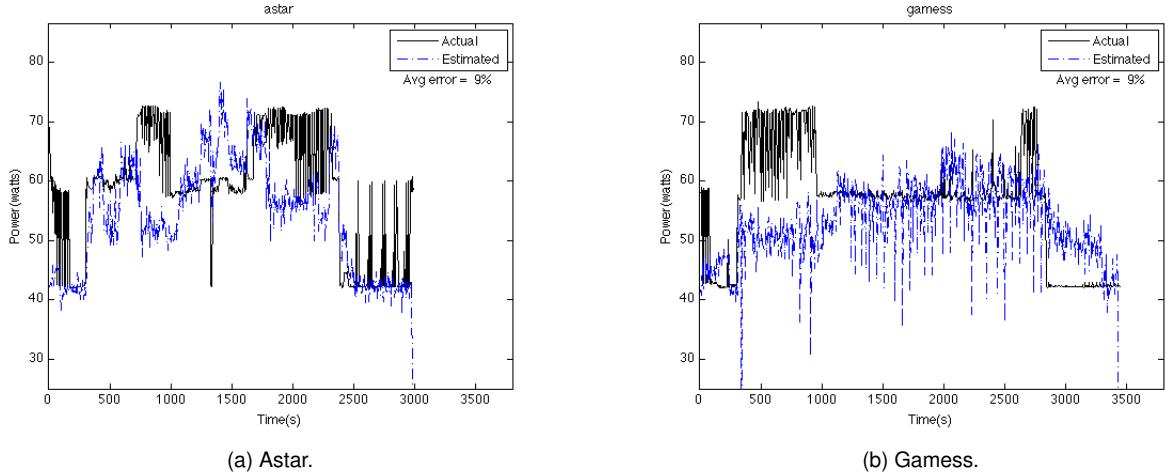


Fig. 6. Actual versus predicted energy consumption for Intel NetBurst Processors.

TABLE 10
ANOVA for Intel Xeon model

Source	df	SS	MS	F	P
Regr	19.00	2455.33	1292.82	143.0500	0.02
Resid	652.00	588.98	9.03		
Total	671.00	2514.23			
R-sq	0.97	Adj. R-sq	0.98		

from the model versus the actual power consumption for each benchmark as shown in Fig. 6. Table 11 shows the descriptive statistics and percentage errors for benchmark predictions.

TABLE 11
Model error for each benchmark for Intel processors

Benchmark	Mean	Median	Std	Var	Error
astar	7.57	5.45	7.41	54.97	9.0%
gamess	7.59	4.91	8.66	75.08	9.2%
gobmk	6.35	5.24	5.65	31.90	9.3%
zeusmp	7.87	6.62	6.90	47.56	9.4%

5.4 Analysis

The consolidated model is attempting to predict for all benchmarks. Given the large volume of data generated thorough the different logging mechanisms, it is nearly impossible to identify and discard bad data. Using the geometric mean as discussed in the previous section helps to smooth out some of the errors introduced in the cases. However, the diversity of the benchmarks used means that some discrepancies arise within variables where we expect to see tight correlations. Thus, the model predicts well in some cases and not in others. Also, the asymmetry of the β -coefficients for tightly correlated variables (the variables for the HyperTransport bus HT_1 and HT_2 , for

example) leads us to believe non-linear relationships may exist among these variables. Therefore, future work needs to consider the impact of use of non-linear regression models together with hardware performance counters for managing energy consumption and thermal envelopes.

Another observation from the model pertains to the placement of the temperature sensors in the server. The ambient temperature regression variable T_{A_0} reflects more of the hot air flow due to the server design. This illustrates that the factors controlling the thermal envelope for different server designs will be accurately reflected in the model. Thus, we would expect to see a more symmetric set of coefficients for ambient temperature variables T_{A_0} and T_{A_1} had the placement of the sensors been more balanced in the server. The benchmarks used to calibrate and

TABLE 12
Model error for a model built from only Integer benchmarks

Benchmark	Mean	Median	Std	Var	Error
astar	18.57	17.99	2.19	4.82	3.2%
gamess	17.80	16.34	6.87	47.22	3.0%
gobmk	15.01	13.87	2.44	5.94	2.5%
zeusmp	14.15	13.68	2.98	8.91	2.4%

evaluate the model were chosen based on two factors: instruction mix and workload type. Benchmarks were selected in an attempt to balance between integer and floating-point instructions. Follow-up tests were performed to evaluate the impact of instruction mix on the predictive ability of the model. Predictive models were created using only the integer benchmarks and only the floating-point benchmarks listed in Table 3. In Table 12 and Table 13, we use each of these models to predict the energy consumption for each of benchmarks in our study. We see that the overall model is a better predictor for the general

TABLE 13
Model error for a model built from only floating point benchmarks

Benchmark	Mean	Median	Std	Var	Error
astar	5.24	5.10	1.72	2.96	8.9%
gamess	3.99	4.21	2.15	4.61	6.7%
gobmk	4.85	3.71	2.44	5.94	8.2%
zeusmp	3.56	3.27	1.76	3.11	6.0%

case as opposed to using the integer-only or floating point-only models. In future work, we will consider benchmarks focused on disk-use or memory utilization. The model developed in this paper is valid for any AMD Opteron dual-core/dual-processor system using the HyperTransport system bus. However, it is scalable to any quad-core dual processors Opteron system using HyperTransport. One would expect to see a slight difference or variation in the predicted power due to the greater or diminished affect of the die temperatures on the other parameters and the model would have to be adjusted accordingly. For a dual-core quad-processor system, the additional term HyperTransport bus regression variable HT_0 would be introduced into the CPU power consumption term and the β coefficients would have to be recalculated and the CPU power equation will have more terms. For a quad-core dual-processor system, similar recalculations would be required.

The experimental validation of our model reveals opportunities for further investigation. The model has been validated for NUMA-based systems such as the AMD Opteron processor and FSB-based systems such as those from the Intel Xeon family constructed using the Intel NetBurst and Core architectures. The model needs to be validated on more recent processors built from the Intel Nehalem architecture. However, the Intel's QuickPath interconnect technology used in the Nehalem architecture [25], like HyperTransport, assumes that each core has integrated memory controllers and NUMA-style memory access. Thus, it is our proposition that further investigation will show that we can much more closely model the power and thermal envelope of these processors using the model proposed in this work versus the results from our case study of the older Intel architectures.

The model requires validation on other architectures such as the IBM Cell BE processors and NVIDIA GPU processors. Further study of the power and thermal envelope of non-CPU components (memory and disk) is required to better understand their contributions to the system thermal envelope.

6 CONCLUDING REMARKS

In this paper, we have introduced a comprehensive model which uses statistical methods to predict system-wide energy consumption on server blades.

The model measures energy input to the system as a function of the work done for completing tasks being gauged and the residual thermal energy given off by the system as a result. Traffic on the system bus, misses in the L2 cache, CPU temperatures, and ambient temperatures are combined together using linear regression techniques to create a predictive model which can be employed to manage the processor thermal envelope.

The experimental validation of our model reveals opportunities for further investigation. The model has been validated for NUMA-based systems such as the AMD Opteron processor and FSB-based systems such as the Intel Xeon processors based on the Intel NetBurst and Core architectures; it requires validation on other architectures such as the Intel Nehalem architecture, IBM Cell BE processors and NVIDIA GPU processors. Further study of the power and thermal envelope of non-CPU components (memory and disk) is required to better understand their contributions to the system thermal envelope.

A fast, accurate, and robust model for the power and thermal envelope for a single server blade is critical to understanding and solving the power management challenges unique in dense servers. The model presented in this work is the first step towards building solutions that bridge multi-core, multiple blade, and full data center power and thermal management.

7 ACKNOWLEDGEMENTS

This work is supported in part by the U.S. Department of Energy (DOE) under Award Number DE-FG02-04ER46136 and by the Board of Regents, State of Louisiana, under Contract Number DOE/LEQSF(2004-07)-ULL.

REFERENCES

- [1] G. Contreras and M. Martonosi, "Power Prediction for Intel XScale®Processors Using Performance Monitoring Unit Events," in *Proc. of the 2005 Int'l Symp. on Low Power Electronics and Design*. New York, NY, USA: ACM, 2005, pp. 221–226.
- [2] H. T. Consortium, "HyperTransport I/O Link Specification," HyperTransport Technology Consortium, Specification 3.0c, September 2007.
- [3] *Intel® 64 and IA-32 Architectures Optimization Reference Manual*, Intel Corporation, P.O. Box 5937; Denver, CO, November 2009.
- [4] D. Economou, S. Rivoire, C. Kozyrakis, and P. Ranganathan, "Full-System Power Analysis and Modeling for Server Environments," in *Proc. of the 2008 Workshop on Modeling Benchmarking and Simulation*, 2006.
- [5] C. Isci and M. Martonosi, "Phase Characterization for Power: Evaluating Control-flow-based and Event-counter-based Techniques," *Proc. of the 12th Int'l Symp. on High-Performance Computer Architecture*, pp. 121–132, 11-15 Feb. 2006.
- [6] ———, "Runtime Power Monitoring in High-end Processors: Methodology and Empirical Data," *Proc. of the 36th IEEE/ACM Int'l Symp. on Microarchitecture*, pp. 93–104, 3-5 Dec. 2003.
- [7] W. Bircher and L. John, "Complete System Power Estimation: A Trickle-Down Approach Based on Performance Events," *Ispass*, vol. 0, pp. 158–168, 2007.
- [8] A. Lewis, S. Ghosh, and N.-F. Tzeng, "Run-time energy consumption estimation based on workload in server systems," in *Proc. of the 2008 Workshop on Power Aware Computing and Systems (Hotpower'08)*, 2008.

- [9] S. Rivoire, "Models and Metrics for Energy-efficient Computer Systems," Ph.D. dissertation, Stanford University, 2008.
- [10] F. Bellosa, A. Weissel, M. Waitz, and S. Kellner, "Event-driven energy accounting for dynamic thermal management," *Proc. of the 2003 Workshop on Compilers and Operating Systems for Low Power*, 2003.
- [11] X. Fan, W.-D. Weber, and L. A. Barroso, "Power provisioning for a warehouse-sized computer," in *Proc. of the 34th Int'l Symp. on Computer Architecture*. New York, NY, USA: ACM, 2007, pp. 13–23.
- [12] T. Heath, B. Diniz, E. V. Carrera, W. M. Jr., and R. Bianchini, "Energy Conservation in Heterogeneous Server Clusters," in *Proc. of the 10th ACM SIGPLAN Symp. on Principles and Practice of Parallel Programming*. New York, NY, USA: ACM, 2005, pp. 186–195.
- [13] S. Rivoire, P. Ranganathan, and C. Kozyrakis, "A Comparison of High Level Full-System Power Models," in *Proc. of 2008 USENIX Workshop on Power Aware Computing and Systems*, 2008.
- [14] C. Isci and M. Martonosi, "Identifying Program Power Phase Behavior Using Power Vectors," *Proc. of the IEEE 2003 Int'l Workshop on Workload Characterization*, pp. 108–118, 27 Oct. 2003.
- [15] W. Bircher, J. Law, W. Valluri, and L. John, "Effective Use of Performance Monitoring Counters for Run-Time Prediction of Power," The University of Texas at Austin, Tech. Rep. TR-041104-01, November 2004.
- [16] K.-J. Lee and K. Skadron, "Using performance counters for runtime temperature sensing in high-performance processors," *Proc. of the 19th IEEE Int'l Symp. Parallel and Distributed Processing*, pp. 8 pp.–, 4-8 April 2005.
- [17] *AMD Opteron Processor Data Sheet*, 3rd ed., AMD, March 2007.
- [18] "EPS12v Power Supply Design Guide, V2.92," Server System Infrastructure Consortium, Spec. 2.92, 2004.
- [19] M. Ton, B. Fortenberry, and W. Tschudi. DC Power for Improved Data Center Efficiency.
- [20] Maxim. Practical Considerations for Advanced Current Sensing in High-Reliability Systems.
- [21] "Calculating Memory System Power for DDR3," Micron, Inc, Tech. Note TN41_01DDR3 Rev.B, August 2007.
- [22] J. L. Henning, "Spec cpu2006 benchmark descriptions," *Computer Architecture News*, vol. 34, no. 4, Sept 2006.
- [23] A. Kumar, L. Shang, L.-S. Peh, and N. Jha, "System-Level Dynamic Thermal Management for High-Performance Microprocessors," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 27, no. 1, pp. 96–108, Jan. 2008.
- [24] Electronic Educational Devices, Inc., "WattsUp Power Meter," December 2006.
- [25] Intel, "An Introduction to the Intel QuickPath Interconnect."