

The Case for Power-Agile Computing

Geoffrey Challen
MIT, SUNY Buffalo*

Mark Hempstead
Drexel University

1 Introduction

Battery-powered devices are trapped by trends. More powerful performance requires more power, battery technologies are improving slowly [8], and users want their devices to last as long as they do today—or longer [19]. A way to escape this trap leverages power-proportional hardware architectures [5] that scale performance and power consumption to perform when needed and draw little power when idle. Because most device components are tuned to operate efficiently within a narrow power-performance range, we expect future power-proportional architectures to be increasingly *heterogeneous*: featuring multiple different processors, memory banks, storage devices and radios, each component embodying a particular power-performance tradeoff. Heterogeneity produces single devices that can morph into many others: a phone that can sprint like a laptop and sleep like a sensor node.

Today many devices already incorporate multiple processors, storage devices and radios with different power-performance characteristics. Researchers have proposed operating system designs acknowledging this heterogeneity [6], performance- or power-driven component combinations [13, 4], approaches harnessing the efficiency of a particular set of components for a certain task [1, 18], and systems organized into multiple power-performance tiers [17]. Inspired by these efforts, we introduce the term *power agility* to describe the ability of a system to operate a heterogeneous power-proportional device balancing performance and power consumption.

Given increasingly heterogeneous devices, power agility requires not merely adjusting individual components but activating and deactivating them to react to changes in demand caused by variations in device usage. The idle phone in my pocket consumes less power than the one routing me to my destination, and while the mapping application wants the high-power radio, the game prefers a faster processor. So while power-

proportional *hardware* allows the device to sprint and sleep, power-agile *software* guides it nimbly between states in response to demand. We consider microarchitectural attempts to mask power fluctuations [16] from the OS a mistake. Only the operating system has enough information—from user priorities to a view of all running processes—to achieve power agility.

This paper outlines the principles of power-agile computing. First, we construct a heterogeneous power-proportional device illustrating the size and diversity of the state space inherent to these architectures. Next, we present a scenario demonstrating our device responding to changes in demand. Using this scenario, we develop a set of challenges inherent to power-agile operation and discuss approaches to overcoming them.

2 Example Architecture

To begin, we assemble a heterogeneous power-proportional device combining two general-purpose processors¹, two memory chips, two storage devices and two radios. Table 1 presents the components we selected.

The relationship between power and performance varies for each component. Processors may transition smoothly over a restricted power envelope using DVFS, but cannot scale to zero due to leakage current. DRAM memory chips have a constant refresh cost scaling roughly with capacity plus additional power draw corresponding to the rate of reads and writes. Storage devices differ based on whether or not they include spinning components. Flash drives do not and scale approximately with usage but are limited in size. Radios exhibit wide power-performance variation because their usage depends both on the hardware and the protocol. 802.11 clients can enter power-saving mode (PSM) which uses base station buffering to save power. Bluetooth has limited range but lower power consumption balanced between both sides of the link.

*Beginning August, 2011.

¹Distinguished from task-specific processors like GPUs or DSPs.

ID	Name	mW	Performance
P1	ARM Cortex-M4 ¹ [3]	0.9 ²	75 MHz
		15.6	300 MHz
P2	ARM Cortex-A9 [2]	23.5 ²	415 MHz
		400.	830 MHz
M1	32 MB ISSI SDRAM [11]	81.	Refresh only
		108.	166 MHz
M2	1 GB Micron “Slow” DDR2	322. ⁴	Refresh only
		482.	266 MHz
S1	2 GB MicroSD Card	20. ⁵	Idle
		100. ⁵	25 MBps ⁶
S2	64 GB OCZ SSD	200. ⁷	Idle
		1000. ⁷	5.5 MBps ⁷
R1	250 kbps TI CC2540 BLE	6.7 ³	10% duty cycle ⁸
		66.3 ⁹	Receive mode ⁹
R2	11 Mbps Marvell 802.11bg	30.9 ³	10% duty cycle ⁸
		309.3 ¹⁰	Idle mode ¹⁰

¹ Capable of running a subset of the full P2 ISA.

² Optimistic estimate based on an optimistic estimate of DVFS providing 1:5 performance and 1:17 power scaling [9].

³ Estimated based on scaled full-power performance.

⁴ Estimated based on Micron leakage numbers.

⁵ Estimated due to lack of publicly-available datasheets.

⁶ Maximum achievable.

⁷ Measured by Tom’s Hardware [10].

⁸ Duty cycling shifts power usage from the receiver to the sender, which has to remain online (as in 802.11 PSM) or send longer packets (as in 802.15.4 Low-Power Listening [14]).

⁹ Receive-only in high-sensitivity mode. Transmit is similar.

¹⁰ Transmit and receive vary so usage is workload-dependent.

Table 1: Performance and power consumption of selected hardware components. We assume voltage gating can reduce the usage of disabled components to near zero [15]. The 10 notes reflect the challenge in obtaining these numbers. Most data sheets omit this information.

We define a *component ensemble* as the set of components currently active, constraining the set of valid ensembles to include only those that can support the device operating system. For our example, these include (a) one or both processors, (b) one or both memory chips², (c) zero, one or both storage devices and (d) zero, one or both radios. By switching between components our device can operate across a wide power range. In its lowest-power ensemble, the device has a 75 MHz CPU, 32 MB of RAM, and draws 82³ mW and is roughly-equivalent to a embedded sensor node. In its highest-power ensemble the device has multiple cores, over 1 GB of RAM, over 320 GB of storage, Wifi and Bluetooth. Consuming almost 2.5 W, it is similar to emerging smartphones.

²While many low-power processors come with small amounts of integrated memory, we have conservatively chosen to require 32 MB of RAM in order to run embedded versions of Linux. It is conceivable that our candidate device could enter an active sleep state with a micro-kernel capable of fitting in the processor’s onboard RAM.

³Actual power consumption would be higher due to system buses, memory controllers, and other components of a complete architecture.

This device can activate *144 valid component ensembles*⁴. Figure 1 shows the composition and power envelope of each, and motivates two observations. First, there are many valid ensembles and wide usage variation even in an architecture with only two components per class. Incorporating more components would produce even more options. Second, at any power level there are many diverse ensembles the device can use: a fast processor, small memory chip, and slow disk; a slow processor, large memory chip, and fast radio; etc. These differ not in their total power consumption but in how they perform and distribute power across components, and while some ensembles may seem too weird to be useful they may suit certain applications. Finally, while it may seem best to avoid inefficient ensembles—those achieving low utilization and a low active- to idle-power ratio—given the speed of temporal changes in demand and the overhead of ensemble transitions we expect devices to spend some time at the low end of ensemble power envelopes.

3 Challenges

To illustrate how a power-agile device might operate we imagine a phone performing a background task interrupted by an interactive session. Figure 2 shows how overall and per-component power allocations change to respond to the needs of the two applications. We refer to this scenario throughout the rest of this section as we examine the challenges inherent to power-agile computing. These are related to five roles that the operating system plays while operating power-agile hardware: measuring (3.1) and predicting (3.2) performance; and selecting (3.3), preparing (3.4) and executing (3.5) ensemble transitions. Throughout we demonstrate how traditional scheduling and resource-allocation problems are complicated by the flexible nature of the underlying hardware.

3.1 Measuring Efficiency

Determining performance differences between ensembles requires application metrics weighting both power and performance. We believe that variants of the energy-delay product—EDP = $E\Delta$, which the system tries to minimize—used in circuit design [12] may be appropriate for combining these concerns. E measures the energy consumed during some time quantum and Δ measures a performance characteristic of interest to the application: the time necessary to process a block of data or respond to user input. Controlling the strength of the performance component using an exponential, $EDP = E\Delta^n$, allows applications to weight their preference for performance v. efficiency. In our scenario, the

⁴3 processor choices \times 3 memory choices \times 4 storage choices \times 4 radio choices.

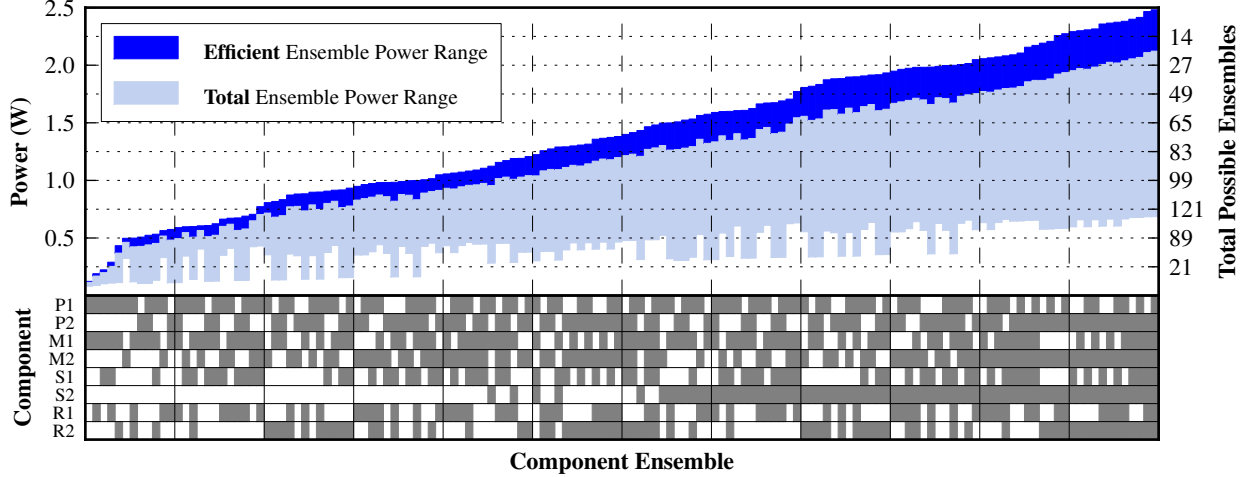


Figure 1: **Power envelopes of all 144 example device component ensembles.** Ensembles are sorted by increasing maximum power draw. For each ensemble, the bottom shows which components are active and the top displays the power envelope. The top 80% of the envelope—the most efficient operating range—is drawn in dark blue. The right axis counts the total number of ensembles that might draw that much power: e.g., there are 121 ensembles that could consume 0.75 W, depending on the workload.

interactive application uses $E\Delta^2$ —causing the system to activate high-performance ensembles—while the background task uses $E\sqrt{\Delta}$, causing the system to remain in lower-power states. Application metrics are likely to change over time as their needs change.

3.2 Predicting Ensemble Performance

Given the size of the ensemble state space, predicting ensemble performance is a key part of transitions. Assuming an application with power-performance preference $E\Delta^n$, both E and Δ will vary across ensembles: E with the cost and utilization of system components, and Δ with performance. The most direct way to determine power-performance is to run the application on many ensembles, but given the number of states and transition cost this is infeasible online. However, offline experimentation could produce binary annotations. Another approach is to have executables include hints about performance characteristics important to various stages. Before transmitting a large amount of information, a hint would alert the system to the imminent need for a high-bandwidth radio. While hints require programmer or compiler support, they may be portable across devices.

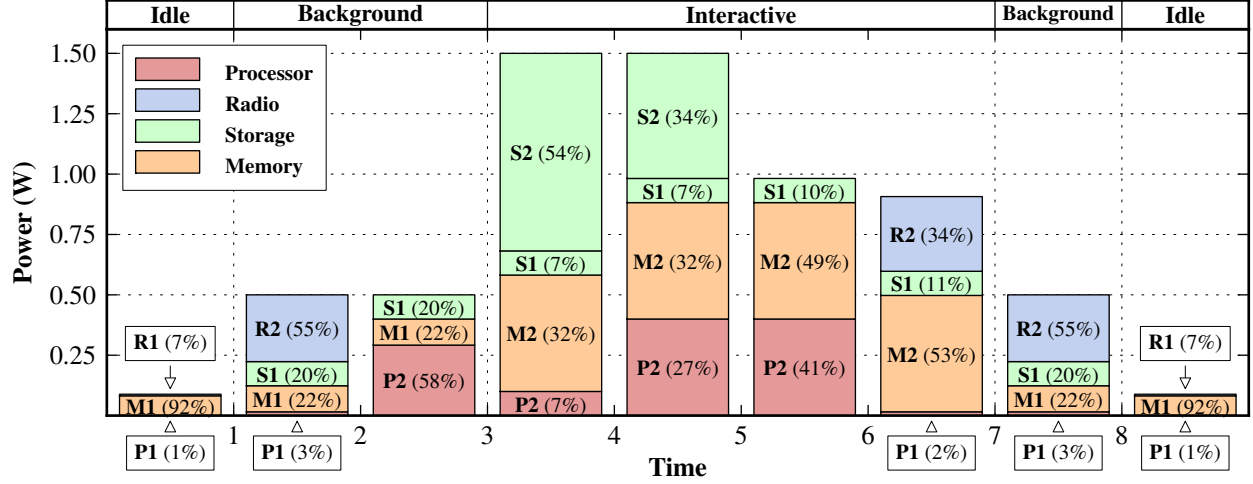
When running unannotated binaries, or mixtures of applications with complicated performance dependencies, the system may need to estimate the impact of ensemble changes before performing them. In some cases, the currently running ensemble can be *artificially* constrained to estimate how performance might change after a component is disabled or enabled. For example, when moving from **M2** to **M1** at $t = 7$ in the scenario, the system might be concerned about the impact of this

transition on the usage of **S1**. If disabling the large memory chip causes **S1** usage to increase dramatically, the system will fail to achieve the intended power reduction. To uncover a link between memory size and disk usage, the operating system can artificially limit the amount of memory in use by trimming pages from **M2**. It may do this in a smooth fashion until it is using only roughly the same amount of the larger chip as the smaller chip size, and then, assuming no serious component relationships have been uncovered, initiate the transition. This strategy is generally more applicable to transitions that attempt to trim power by disabling components, but this is also precisely where it is most useful, as it allows the operating system to discover relationships between component usage that might neuter power reductions.

3.3 Selecting Component Ensembles

Scheduling ensemble transitions relies on the capabilities already presented—metrics for evaluating performance and predicting performance across ensembles. When running only a single application the system can respond directly to its hints, annotations, or estimated performance, weighting possible improvements in efficiency against the cost of transitioning ensembles.

Running multiple applications creates new challenges. First, there is the question of how to assign performance metrics to applications. In our scenario the background task would complete faster if it were allowed to use the higher exponent used by the interactive application. The goal is to assign the most efficient metric to the application that produces acceptable performance, and doing so is likely to require user feedback.



- 0 When idle **P1** and **M1** are idled and **R1** operates at low duty cycle.
- 1 Receiving data over **R1** the phone initiates a background task. The device activates **R2** to rapidly receive data and **S1** to store it.
- 2 As the phone begins processing the task it activates **P2** and disables **R2**.
- 3 The user removes the phone from their pocket and begins interacting with an application, which activates **M2** and retrieves data from **S2**.
- 4 As the interactive application continues energy usage shifts from **S2** to **P2**.
- 5 When the interactive application is finished with **S2** it is disabled.
- 6 As the interactive session completes, the phone offloads data using **R2** driven by **P1**.
- 7 Background processing resumes in the same ensemble it was using previously.
- 8 The background task completes, idling the phone.

Figure 2: **Scenario.** The figure and table describe the scenario referred to throughout Section 3. Bars indicate the total energy consumed, broken down and labeled by component. The table describes what is happening at each time step.

Choosing the correct ensemble for both applications is the next challenge. If their performance requirements are aligned, then an ensemble may exist that works well for both. Applications differing in their performance requirements complicate the process. If the system has sufficient energy it may choose to operate a combination of both ideal ensembles, but this produces inefficiency as the set of distinct resources needed by one application is idled while the other runs.

The simplest approach is to transition between the ideal ensembles while increasing both application's time quanta sufficient to amortize the transition cost. In many cases, however, we expect that this will lead to unacceptable interactive performance. A second possible approach is to pick an ensemble that produces acceptable—but not ideal—performance for both applications, potentially weighted towards the application with higher priority. Another option is to select an ensemble optimized for one application while allocating resources within that ensemble in favor of the other. For example, given one application that requires a high-speed disk and another than needs a large memory chip, we can choose to use the large memory chip and a slower disk allocating a large portion of the memory to a buffer cache to improve performance for the I/O-bound application.

3.4 Preparing Ensemble Transitions

Because ensemble transitions are both important and costly, the operating system should prepare the ground to minimize their overhead. Preparation is particularly important in the memory and storage hierarchy, where the location of data has a significant impact on component transitions. Preparation also requires the system forecast future application demand and ensemble dwell times.

Consider an example transition that activates a larger memory chip with superior performance. If the system will be in that ensemble for a significant length of time, all applications will benefit from having data relocated from the smaller to the larger chip. This also allows the smaller chip to be shut off to save power. However, if and when the device wants to disable the larger memory chip in order to shift power toward some other necessary component, the amount of data stored in the larger memory bank creates a high overhead for this transition.

If the system predicts brief use of the larger memory bank, it may try several strategies to reduce the eventual transition overhead. First, if the transition is due to a particular application, it may continue to operate the smaller chip for other applications while allocating new pages on the larger component. Once the memory-hungry ap-

plication is finished with these pages, they can be discarded and the memory disabled without migrating data. Another approach is to copy accessed pages on demand but mirror writes to both memory banks to minimize the eventual transition cost. Assuming that the smaller chip is never shut off—possible if consumes little power—the physical address space may be configured to always mirror a portion to both chips when the larger bank is active. The operating system may try to allocate memory from the mirrored portion of the address space for pages that have long expected lifetimes, are used by applications that prefer more power-efficient states, or based on explicit application requests. These pages will benefit from better performance when the larger bank is active while never requiring migration.

3.5 Executing Ensemble Transitions

Ensemble transitions tailor the device to application demands but require potentially complex or energy-expensive component transitions. The ACPI specification [7] standardizes per-component and overall power states, but does not consider the overhead of moving from component to component. Below we outline for each component class, the complexity and cost of transitions and a brief description of how to perform one:

- **Processor:** Difficulty: *high*, Cost: *medium*. Transitioning between processors, even ones with highly-compatible instruction sets, requires migrating process state, correcting for processor differences, and potentially reloading new process executables enabling or disabling certain instructions.
- **Memory:** Difficulty: *medium*, Cost: *high*. Moving to a smaller chip requires migrating some pages to the new memory area while flushing others to the backing store, along with kernel adjustments to its own memory footprint. Transitioning to a larger chip requires migrating data.
- **Storage:** Difficulty: *low*, Cost: *low*. Disabling requires writing out dirty buffers. Enabling will cause a performance dip while caches fill.
- **Radio:** Difficulty: *medium*, Cost: *medium*. Disabling requires flushing any outstanding buffers, closing connections and potentially coordinating with the receiver to move together to a new radio technology. Enabling may require association—potentially costly, depending on the protocol—and a delay while link parameters necessary for efficient operation can be determined.

4 Summary

Power-proportional heterogeneous devices demand operating system support to achieve power-agility: the ability to operate these devices balancing performance and power efficiency. Given the constraints facing battery-powered devices, we consider power-agility critical to continued performance improvements, and have outlined the challenges on the road to power-agile computing.

References

- [1] ANDERSEN, D. G., FRANKLIN, J., KAMINSKY, M., PHANISHAYEE, A., TAN, L., AND VASUDEVAN, V. Fawn: a fast array of wimpy nodes. In *Proceedings of the ACM SIGOPS 22nd symposium on Operating systems principles* (New York, NY, USA, 2009), SOSP '09, ACM, pp. 1–14.
- [2] ARM. Cortex-A9 Processor. <http://www.arm.com/products/processors/cortex-a/cortex-a9.php>.
- [3] ARM. Cortex-M4 Processor. <http://www.arm.com/products/processors/cortex-m/cortex-m4-processor.php>.
- [4] BALASUBRAMANIAN, A., MAHAJAN, R., AND VENKATARAMANI, A. Augmenting mobile 3g using wifi. In *MobiSys '10: Proceedings of the 8th international conference on Mobile systems, applications, and services* (New York, NY, USA, 2010), ACM, pp. 209–222.
- [5] BARROSO, L. A., AND HÖLZLE, U. The case for energy-proportional computing. *Computer* 40 (December 2007), 33–37.
- [6] BAUMANN, A., BARHAM, P., DAGAND, P.-E., HARRIS, T., ISAACS, R., PETER, S., ROSCOE, T., SCHÜPBACH, A., AND SINGHANIA, A. The multikernel: a new os architecture for scalable multicore systems. In *Proceedings of the ACM SIGOPS 22nd symposium on Operating systems principles* (New York, NY, USA, 2009), SOSP '09, ACM, pp. 29–44.
- [7] CORPORATION, H.-P., CORPORATION, I., CORPORATION, M., LTD., P. T., AND CORPORATION, T. Advanced configuration and power interface. <http://www.acpi.info/>.
- [8] ECONOMIST, T. In search of the perfect battery. <http://www.economist.com/node/10789409>.
- [9] ET AL., K. J. N. A 32-bit powerpc system-on-a-chip with support for dynamic voltage scaling and dynamic frequency scaling. *IEEE Journal of Solid-State Circuits* 37, 11 (Nov 2002), 1441–1447.
- [10] HARDWARE, T. Flash SSD Update: More Results, Answers. <http://www.tomshardware.com/reviews/ssd-hard-drive,1968.html>.
- [11] ISSI. IS42VM32100C Advanced Information. www.issi.com/pdf/42VM32100C.pdf.
- [12] MARTIN, A. J., NYSTRÖM, M., AND PÉNZES, P. I. ET2: a metric for time and energy efficiency of computation. Kluwer Academic Publishers, Norwell, MA, USA, 2002, pp. 293–315.
- [13] MOGUL, J. C., ARGOLLO, E., SHAH, M., AND FARABOSCHI, P. Operating system support for nvm+dram hybrid main memory. In *Proceedings of the 12th conference on Hot topics in operating systems* (Berkeley, CA, USA, 2009), HotOS'09, USENIX Association, pp. 14–14.
- [14] MOSS, D., HUI, J., LEVIS, P., AND CHOI, J. I. Cc2420 radio stack. TinyOS Extension Proposal TEP-126, <http://www.tinyos.net/tinyos-2.x/doc/txt/tep126.txt>, 2007.
- [15] POWELL, M., YANG, S.-H., FALSAFI, B., ROY, K., AND VIJAYKUMAR, T. Gated-Vdd: A Circuit Technique to Reduce Leakage in Deep-Submicron Cache Memories. In *International Symposium on Low Power Electronics and Design (ISLPED)* (June 2000).
- [16] RANGAN, K., POWELL, M., WEI, G.-Y., AND BROOKS, D. Achieving Uniform Performance and Maximizing Throughput in the Presence of Heterogeneity. In *International Symposium on High Performance Microarchitecture (HPCA)* (January 2011).
- [17] SORBER, J., BANERJEE, N., CORNER, M. D., AND ROLLINS, S. Turducken: hierarchical power management for mobile devices. In *Proceedings of the 3rd international conference on Mobile systems, applications, and services* (New York, NY, USA, 2005), MobiSys '05, ACM, pp. 261–274.
- [18] SZALAY, A. S., BELL, G. C., HUANG, H. H., TERZIS, A., AND WHITE, A. Low-power amdahl-balanced blades for data intensive computing. *SIGOPS Oper. Syst. Rev.* 44 (March 2010), 71–75.
- [19] WALKO, J. What do cell phone users want? better batteries! <http://www.informationweek.com/news/showArticle.jhtml?articleID=171100095>.