

A Compilation Framework for Power and Energy Management on Mobile Computers

Ulrich Kremer*
Rutgers University
Department of Computer Science
New Brunswick, NJ 08854
uli@cs.rutgers.edu

Jamey Hicks and James M. Rehg
Compaq Computer Corporation
Cambridge Research Lab
Cambridge, MA 02142
{hicks,rehg}@crl.dec.com

Rutgers University
Department of Computer Science Technical Report
DCS-TR-446, June 2001

Abstract

This paper discusses the potential benefits of application-specific power management through remote task execution. Power management is crucial for mobile devices that have to rely on battery power for extended periods of time. Image processing and understanding is a class of applications that is important in mobile environments. Image processing can be used in autonomous robot navigation, target acquisition/classification, keyboard-less input, and aerial surveillance (Micro Air Vehicles), just to mention a few. Experimental results on an image processing application, namely a human face detection and recognition system, indicate the power savings that can be achieved for this important class of applications.

We discuss a compilation strategy that generates two versions of the initial application, one to be executed on the mobile device (client), and the other on a machine connected to the mobile device via a wireless network (server). The client and server codes have to be able to deal with disconnection events. Our proposed compilation strategy uses checkpointing techniques to allow the client to monitor program progress on the server, and to request checkpoint data in order to reduce the performance penalty in case of a possible server and/or network failure.

The reported results have been obtained by actual power measurements on three client systems, (1) the

StrongARM based low-power SKIFF system developed at Compaq's Cambridge Research Laboratory, (2) Compaq's commercially available StrongARM based iPAQ H3600, and (3) a PentiumII based laptop. Initial experiments show that energy consumption can be reduced significantly, in some cases up to one order of magnitude, depending on the selected characteristics of the mobile device, remote host, and wireless network. A prototype implementation of the discussed compilation framework is underway, and preliminary results are reported.

1 Introduction

Power dissipation has become one of the crucial design challenges of current and future computer systems. In a mobile environment, power savings are important to prolong battery life. For a desk-top "wall-powered" system, heat emission has become a severe design limitation with respect to transistor densities and clock frequencies. Power and energy management addresses both of these issues. However, in the context of this paper, prolonging battery life is the main objective¹.

Mobile devices come in many flavors, including laptop computers, webphones, pocket computers, Personal Digital Assistance (PDAs), and intelligent sensors. Many such devices already have wireless communication capabilities, and we expect most future systems to have such capabilities. There are two main differences between mobile

*This research was partially conducted while the first author was a visiting researcher at the Compaq's Cambridge Research Lab (CRL). Additional funding has been provided by NSF CAREER award No. 9985050.

¹In this paper, we consider a reduction in power and energy as the same optimization goal. This is a simplifying assumption, and not true in general.

and desk-top computing systems, namely the source of the power supply and the amount of available resources. Mobile systems operate entirely on battery power most or all the time. The resources available on a mobile system can be expected to be at least one order of magnitude less than those of a “wall-powered” desk-top system with similar technology. This fact is mostly due to space, weight, and power limitations placed on mobile platforms. Such resources include the amount and speed of the processor, memory, secondary storage, and I/O. With the development of new and even more power-hungry technology, we expect this gap to widen even more.

Image processing and image understanding will be one of the key applications for low-power mobile devices, in addition to speech recognition. Low power single-chip imagers, such as the Photobit PB-0100, are becoming available that will allow mobile devices to capture high-quality images. Image processing can be used in the context of robot control and navigation, wide-scale video surveillance, guidance systems for blind people to identify and classify objects, keyboard-less human-computer interfaces, and Micro Air Vehicles (MAVs) to support aerial surveillance, target acquisition and target classification. Face detection and recognition in video images is a key technology for several of these applications.

The compilation approach for power and energy management discussed in this paper is complementary to operating systems and hardware techniques. The latter techniques rely on observed past program behavior and resource requirements to predict future program characteristics. A compile-time whole program analysis is often able to determine future behavior and requirements since the entire program is available to the compiler and not just a limited window of recent events. In addition, a compiler may perform high-level transformations and thereby “reshape”, i.e., change program behavior and resource requirements, allowing further program optimizations. In cases where compile-time information is not available to perform a particular optimization task, the compiler may generate code that will make optimization decisions at run time based on run-time values of compile-time determined variables and conditions.

In this paper, we give a first assessment of the potential benefits of compiler-directed remote task mapping as an optimization to save energy on mobile devices. We use the tasks of face detection and recognition, in the context of a TourGuide system for visitors to Compaq’s CRL lab, as an experimental vehicle for our investigation. The main challenges in remote task mapping is the identification of suitable remote tasks where the communication overhead is more than compensated for by the expected

power savings. We discuss a compilation framework for remote task execution, and give experimental results for a hand-simulated compilation of our face detection and recognition program. The framework addresses the problem of network disconnection. Power measurements are reported for three platforms: a low-power single-board system called Skiff [8], which was developed at CRL, a recently introduced new handheld PC by Compaq (iPAQ H3650), and a PentiumII based laptop computer. Initial results with a prototype implementation of the compilation framework indicate the feasibility of our approach.

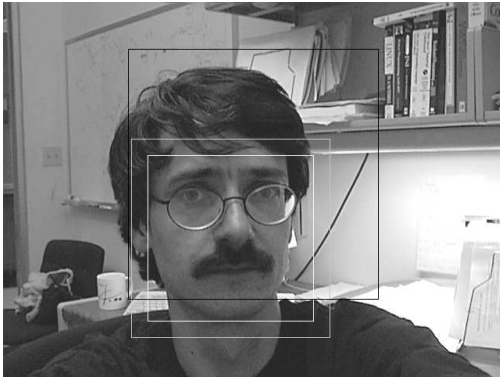
2 Remote Task Mapping

The input to the envisioned compilation system is a program written in C or Java that does not contain any user annotations, for instance, to specify tasks for remote execution. Compiler-directed remote task mapping for power and energy management identifies program tasks that (1) can be safely executed on a remote server, and that (2) will lead to power savings on the mobile client system. The profitability of remote task execution will depend on the amount of communication needed between the mobile client and remote server, the time it takes to complete the remote task, and the communication time itself. The latter two points are of particular importance in mobile systems that have a significant power consumption during the time the system is waiting for the response of the server.

Our proposed power management strategy considers disconnectivity. The system will be able to continue program execution while being disconnected from the server, possibly with a degradation of program performance such as an increase in execution time. This is an important feature for any mission critical program with real-time response requirements. Simply putting the mobile system into a hibernation state until the connection to the server has been reestablished is not acceptable. The goal of power management is to save as much energy as possible while running the program to completion within a set of real-time constraints required by the given application. For example, leaving a robot or a MAV without image processing capabilities for extended periods of time may lead to total system loss.

3 TourGuide : A Face Detection and Recognition System

We have developed a prototype application based on face detection and recognition, called TourGuide , which iden-



Ulrich Kremer: match distance 192

Figure 1: Example Image and TourGuide's output.

tifies researchers and staff members at Compaq's Cambridge Research Laboratory (CRL). The TourGuide system will consist of a mobile client and a network of embedded servers. The client can be given to visitors to the lab. It allows them to more easily identify and keep track of the members of the laboratory that they encounter during their visit. The system is designed to run on a portable low-power device with wireless communication capabilities that is equipped with a digital camera.

The face detection component of the TourGuide architecture utilizes a neural network-based detector for frontal, upright faces developed at Carnegie Mellon Univ. by Rowley et. al. [17]. The detector uses two networks which were trained to detect faces in windows of 20×20 and 30×30 pixels. These networks are scanned across an image pyramid in a coarse-to-fine fashion. The resulting system can detect any number of frontal, upright faces in a single image, regardless of their size or position. Faces which are tilted in the plane of the image or away from the camera will not be detected. However, if the detector is applied to a sequence of video frames taken from an interaction with a person, it is quite likely that there will be several correct detections.

The face recognition module in TourGuide is applied whenever a face is successfully detected. It is based on the ARENA [23, 22] system developed by Sim et. al. from CMU and JustResearch. ARENA uses a nearest neighbor classifier to identify faces by comparing their distance using an L1 norm to labeled example faces (prototypes) taken from a training set. Input faces are correctly identified when they fall within some predefined distance threshold of the correct prototype face. On relatively small populations the ARENA system has gener-

ated impressive results using fairly unconstrained input imagery.

An example input image and the resulting output is shown in Figure 1. A 320×320 central section of the 640×480 input image is passed through the face detector. If successful, the detected face is first enlarged by 20%, summarized, and then processed by the face recognizer. In the shown example, the recognizer reports a successful match and a confidence level for that match.

The overall code structure of the face detection and recognition system is shown in the appendix. In the envisioned application environment, the system may process several input images before successfully recognizing a person. The entire system has been implemented in C and runs under the Linux operating system. The source code is approx. 6000 lines, including comments. Since the current family of StrongARM processors does not have a floating-point unit, we have developed a fixed-point package for efficient floating-point emulation.

The TourGuide system is an interesting application for our benefit analysis of remote task execution for power management purposes since it exhibits important aspects of mobile applications, including

- *necessary mobility*: The system will be carried by the laboratory visitor who is meeting with researchers inside and outside the laboratory space.
- *potential disconnection*: Some part of the laboratory may not be covered by the wireless communication network. In addition, the communication is interrupted while the visitor is leaving the laboratory space.
- *continuous operation*: Even if the visitor is not within the range of a base station, the system should continue to work, although not at the same performance level as within the range of the network.

Dealing with all or just a subset of these aspects is a challenge, and we believe that compiler support will play a crucial role for effective power management of such applications.

4 Compilation Strategy

The basic idea of our compilation strategy is to identify program tasks that are safe to be executed remotely, and that will lead to power savings on the mobile architecture due to remote execution. During remote execution, the mobile machine will enter a hibernation state in order to save as much energy as possible. It is the compiler's

responsibility to insert code that will initiate power state transitions on the mobile machine. While waiting for the completion of a remote task, the mobile client and server may get disconnected. If after a predetermined waiting time the remote server has not finished executing the assigned task, the local client will reexecute the task locally.

The proposed compilation strategy is based on a client/server model and consists of several steps. In the first step, program tasks have to be identified that can safely be executed remotely, i.e., on the server. In the second step, the profitability of remote task execution is examined for remote task candidates. In the final step, two versions of application are generated, one that will be executed on the mobile machine, and one on the remote server accessible from the mobile machine via a wireless connection.

Disconnection events are dealt with through a checkpointing mechanism that provides on-demand data exchange at compile-time determined program points. This checkpointing mechanism serves two purposes, namely (1) informing the client about the progress of the remote task on the server, and (2) allowing the client to receive partial results in order to avoid entire task reexecution in the event of a server or network failure. The version of the application running on the remote server can be considered a specialized server for that application.

The compilation strategy assumes that the entire source program is available during compilation. The compiler will classify I/O operations as performed on *replicated* data or *client* data.

Copies of replicated data are stored on the client and the server machines at the time the client and server programs are compiled on the two platforms, i.e., are made available to the client and server as part of the program installation process. Neither replicated data nor program code will need to be communicated during program execution. Client data is acquired by the client, and may lead to communication to the remote server. We will refer to I/O operations such as image acquisition as client I/O. Figure 2 shows the task graph of our TourGuide system, together with its object life ranges. The task graph nodes represent system initialization (SI:1), image acquisition (I/O:2), face detection preprocessing (PFD:3), face detection (FD:4), face recognition preprocessing (PFR:5), face recognition (FR:6), and final answer (I/O:7). Table 1 summarizes the data objects and their life ranges relative to the subtask nodes in the task graph.

A strategy for remote task execution may evaluate the profitability of remote execution for each individual task between a pair of client I/O operations. Our basic strategy is not as general, and is based on the assumption that the

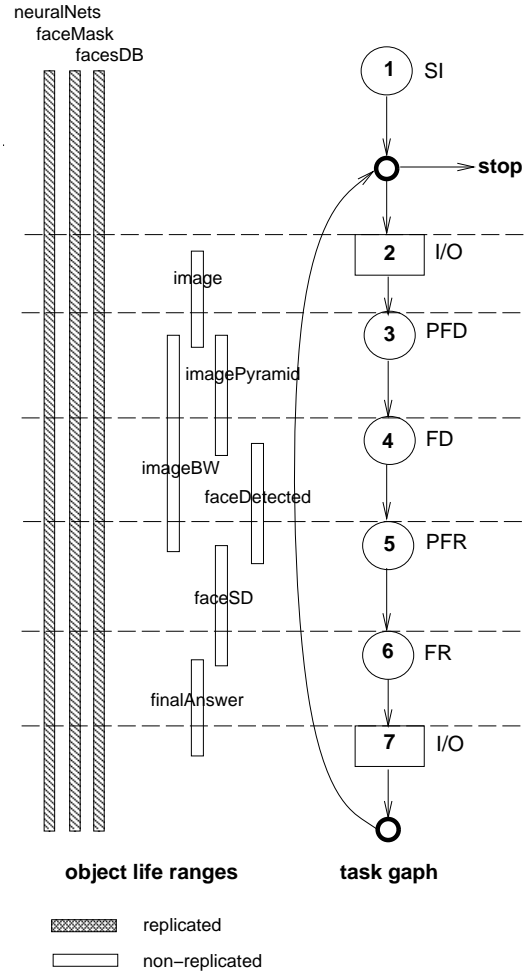


Figure 2: Task graph with object life ranges for the Tour-Guide system.

network and server is available most of the time. In such a scenario, the client may want to initiate remote execution as early as possible, and once initiated, may want to execute as many tasks as possible remotely, i.e., execute all subsequent tasks remotely until the next client I/O operation is reached. During remote execution, the client may enter a hibernation state. The client may wake up periodically to check the progress on the server, and may request checkpoint data from the server before entering the next hibernation interval. Suitable policies for wake-up and on-demand checkpointing are currently under investigation. For example, a client may request the execution of a checkpoint if the wireless signal strength is below a predefined threshold, indicating a potential disconnection

objects	life range ¹	size in bytes	comments
replicated			
neuralNets	2 - 7	300K	neural networks
faceMask	2 - 7	412	20 × 20 pnm image
facesDB	2 - 7	2.8M	faces data base of CRL staff members
non-replicated			
image	3	60K	640 × 480 jpeg image
imageBW	4 - 5	310K	640 × 480 pnm image
imagePyramid	4	153K	10 levels, based on cropped 320 × 320 “imageBW”
faceDetected	5	12	location & size of det. face
faceSD	6	269	16 × 16 scaled-down image
finalAnswer	7	200	output of recognizer

¹ life on entry to node(s).

Table 1: Data objects and their life ranges for the main routine of the TourGuide system program.

event in the near future. The wake-up and checkpointing policies will be implemented through a runtime library and code generated by the compiler as part of the client and server programs. No communication is necessary between the client and the server once the remote execution has been initiated, except for checkpointing purposes and the final answer computed by the server. The resulting execution model has strong similarities with the concept of program continuations in the sense that once a remote task execution is initiated, it will be able to execute the task to completion without any further communication. Due to space limitations, a discussion of our checkpointing mechanisms is beyond the scope of this paper.

A *candidate thread* starts with the initial client I/O operation and ends with the final client I/O operation. Once a remote execution has been initiated, the thread will execute all remaining tasks remotely, returning the final answer to the client just before the final I/O operation. The first task executed remotely is called the *remote entry point*. It is important to note that this remote task execution strategy serves as a base-line strategy, and may need to be refined further.

The communication necessary before the initiation of each remote execution corresponds to the set of non-replicated data objects that are life at the remote entry point. A candidate thread is safe if all data object instances that need to be communicated to the remote server are known. For example, an object instance that cannot be uniquely identified at compile time must not be life at a remote entry point. Such a life range will force the compiler to eliminate this particular candidate thread. Compile-time techniques such as *points-to* analysis (e.g.:

[3]), *escape analysis* (e.g.: [16]), and type analysis will be used to determine *what* objects and *how* objects need to be transmitted between client and server.

For example, Figure 3 shows the candidate thread PFD+FD+PFR+FR (entry point 3) of our TourGuide application with its three compile-time determined on-demand checkpoints. The compiler determines the data objects that need to be communicated at a checkpoint in order to resume program execution on the mobile client. The decision whether or not to execute a checkpoint will depend on the selected checkpointing policy and will be evaluated at runtime. The only life object at entry point 3 is `image`. The object is communicated to the server as part of initiating remote execution. The remote server will send object `finalAnswer` to the mobile client once the task `FR` has been successfully executed on the server.

5 Benefit Experiments

The TourGuide face detection and recognition system serves as a test case for the discussed compilation strategy. The four target configurations for this study are shown in Table 5. All machines run the Linux operating system and used wireless 802.11b PC cards as their wireless network connection. The experiments focus on the execution time reductions obtained through remote task execution for the different configurations.

In general, most performance oriented transformations will also improve the overall energy consumption of an application [25, 26]. However, Kandemir et al. showed that for loop tiling, the best tile sizes for performance and energy consumption are different, i.e., energy consumption and execution time improvements are different optimization goals [11].

In this paper, we focus on execution time reduction as the main tool to decrease energy consumption on a mobile computer, taking advantage of the significant resource gap between mobile computers and high-performance servers and workstation. In the presence of resource hibernation [24], remote execution is an *enabling* transformation, allowing significant energy savings even in cases where the overall program execution time is increased.

5.1 Skiff

Skiff has been designed as a low-power computing appliance and is based on a 233MHz StrongARM processor (SA-110). It has 32MB of SDRAM, 16KB I-cache, 16KB D-cache, 8MB flash memory, one 10Mbps Ethernet port, and two PCMCIA card slots, but no floating point unit.

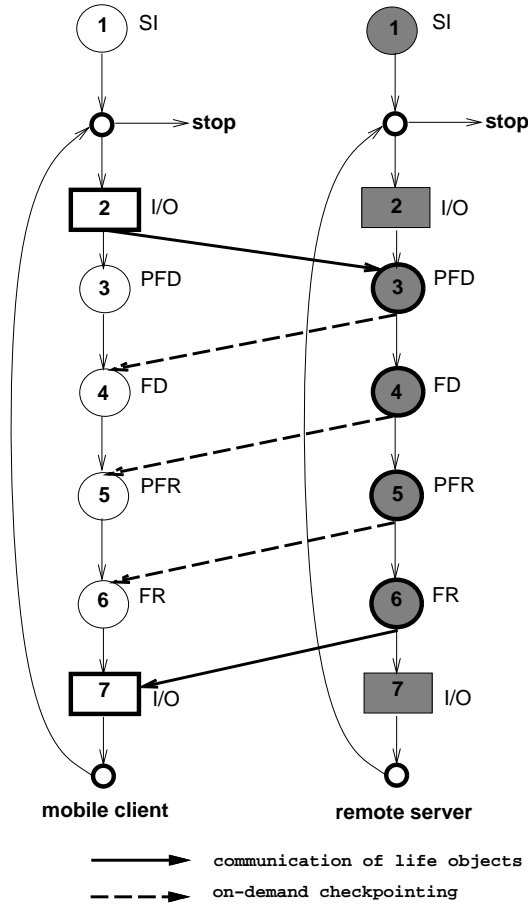


Figure 3: Client and server task graphs for remote thread PFD+FD+PFR+FR. The remote thread has three possible checkpoints.

A full system specification is given elsewhere [8]. The system has separate power planes (2V, 3.3V, and 5V) and allows power measurements for individual system components such as the processor, memory subsystem, and controllers. The power measurements were taken using a Tektronix oscilloscope at each of these power measurement points. In addition, the power supply reported the overall power consumption of the entire system, allowing power measurements for system components such as the power regulators.

Although Skiff has been designed to serve as a low-power computing appliance and not a mobile device, its system characteristics are similar to other handheld computers such as Compaq’s Itsy [5] and Compaq’s iPAQ H3600, which are also based on a StrongARM processor

mobile client	remote server
SA110 (SKIFF) 233 MHz	PentiumII 300 MHz
SA110 (SKIFF) 233 MHz	Dual Proc. PentiumIII 450 MHz
SA1100 (iPAQ) 206 MHz	Dual Proc. PentiumIII 450 MHz
PentiumII 300 MHz	Dual Proc. PentiumIII 450 MHz

Table 2: The target systems for the case study.

(SA-1100). Figure 4 shows the average power dissipation and distribution of Skiff for the four main tasks of TourGuide in a stand-alone mode. The TourGuide system program was compiled using `gcc -O3` and then executed and measured for several input images. For the Skiff experiments, we used Compaq’s WL100 11 Mbps wireless PC card. Overall, Skiff dissipates 3.77W with the wireless PC card, and 2.82W without such a card. Note that even though the wired Ethernet connection is not “in use”, the Ethernet controller consumes significant power. The same holds for the wireless LAN card.

5.2 iPAQ Handheld and PentiumII Laptop

We were only able to take overall power measurements for these two systems. Compaq’s iPAQ H3600 has a 206MHz Intel StrongARM SA-1100 processor. Since the basic unit does not have a PCMCIA slot, we used a PCMCIA expansion pack for the wireless PC card. For the image processing application, the iPAQ dissipates on average 2.20W with a wireless Orinoco (Lucent WaveLan) 11 Mbps card, and 1.25W without the card. Throughout the experiments, the back light of the color display was disabled.

The 300MHz PentiumII based laptop dissipates 19W on average without the display. Here, we used Compaq’s WL100 11 Mbps PC card as the wireless network connection.

5.3 Initial Benefit Analysis

In the initial experiment, we measured the execution times of the basic tasks in the task graph and the cost of communicating non-replicated life objects at each of the four remote entry points. The communication was performed by reading (receiving) and writing (sending) the life data objects to and from files residing on the remote host. In other words, communication was done through the network file system (NFS). The reported results are based on four input pictures. Three measurements were performed

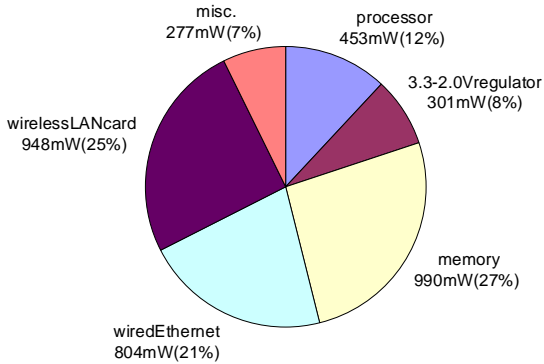


Figure 4: Average power measurements for TourGuide on Skiff.

for each picture, and the reported numbers are the medians over all runs and images. Timings are wall clock times in seconds.

Mainly due to the fact that the StrongARM (Skiff and iPAQ) does not have a floating point unit, floating point intensive computations such as in PFD and FD are executed on the order of $35\times$ faster on the dual processor PentiumIII system than on Skiff. For this case study, all four possible candidate threads, namely PFD+FD+PFR+FR (entry point 3), FD+PFR+FR (entry point 4), PFR+FR (entry point 5), and FR (entry point 6) are safe. Table 3 shows the expected overall execution times for the different candidate threads and target systems. The figures were computed as the sum of execution times for tasks and communications performed on the client and host. Communicating a data object involves writing the object into a file on the sender side, and reading it out of the file on the receiver side. The reported figures exclude system initialization times (SI), and the time needed to acquire the input jpeg image.

Execution time reductions correspond to the expected energy savings if the mobile client does not support hibernation states for its components. The wireless PC cards lead to additional power dissipation of approximately 950mW as shown in Figure 4. A Skiff board without the wireless LAN card dissipates 2.82W. Figure 5 shows the expected energy consumption for Skiff given the additional power costs of the wireless connection. The stand-alone version does not initiate any remote execution and does not use the wireless PC card, i.e., the PC card is

candidate thread	execution time (in seconds)
Skiff local only	8.691
iPAQ local only	9.510
PII local only	0.708

	Skiff &		iPAQ &	PII &
	PII	PIII	PIII	PIII
PFD+FD+PFR+FR	0.877	0.624	0.620	0.469
FD+PFR+FR	4.586	4.347	4.345	0.655
PFR+FR	7.716	7.543	8.089	0.655
FR	8.702	8.632	9.425	0.626

Table 3: Execution times of four candidate threads.

not inserted into one of Skiff's PCMCIA slots and therefore does not consume any power.

For the iPAQ, the power dissipation difference is even more significant, namely 2.20W with, and 1.25W without the card, i.e., over 40% of the overall power budget are due to the wireless connection.

For the PentiumII laptop, the additional power consumption of the wireless card is not significant (5% increase). For all target systems, we did not observe a significant difference in power consumption due to program initiated communication over the wireless Ethernet connection, i.e., the costs for sending and receiving were comparable.

The results show speedups of up to 13.9x with resulting energy savings of up to 10x on Skiff, and speedups of up to 15.3x with resulting energy savings of up to 8.6x on the iPAQ. However, the energy consumption can be significantly larger for some candidate threads as compared to the version of the code that is executed only on the mobile client due to the additional energy consumed by the plugged-in wireless LAN card. On the laptop, energy savings on the order of 30% are possible. In contrast to Skiff and iPAQ, all candidate threads lead to energy savings. The reported figures do not consider power savings due to transitions into low-power hibernation states on the mobile client.

6 Prototype Compiler

The implementation of a prototype system is currently underway and is based on the SUIF2 compilation infrastructure. The current prototype compiler is limited, and the final paper will contain updated results with respect to this compilation system.

The compiler takes C programs as input. Candidates for remote execution are procedure calls in the main rou-

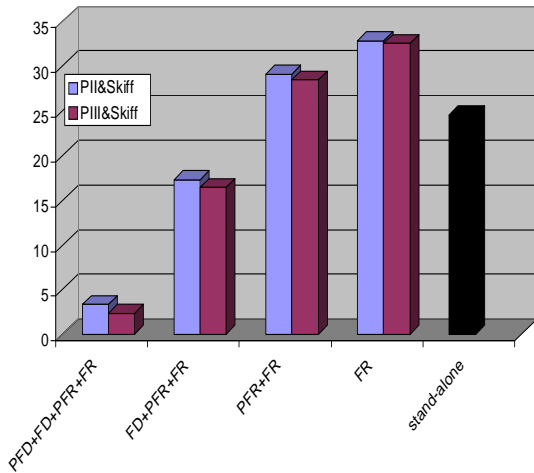


Figure 5: Expected energy consumption on Skiff in Joules for the four remote candidate threads of the TourGuide system and two server target architectures. Reported figures are based on candidate threads execution times and Skiff’s average power dissipation (3.77W with wireless LAN card, 2.82W without).

tine. All calls are assumed to be safe, i.e., procedure calls do not have side effects other than modifying data objects passed as reference parameters. This condition is currently not verified by the compiler. The compiler generates separate client and server programs, where client and server communicate with each other through Unix sockets. A parameterized, static performance model is used to determine the profitability of remote execution of a candidate procedure call. The compiler generates a call to a policy routine for each candidate procedure call. The policy routines evaluate the parameterized performance model for the sizes of the actual parameters at each particular candidate call site. The policy routines also check whether a communication link is available or not. In the latter case, the candidate procedure call will be executed locally. The basic decision strategy is as follows:

```

if local_computation_cost(candidate_call) >
    (remote_computation_cost(candidate_call) +
     communication_cost(parameter_list_with_sizes))
then
    remote_call(candidate_call)
else
    local_call(candidate_call)

```

The used cost functions have been derived by a micro-benchmarking approach [2, 20], and reflect only expected

execution times as an approximation to energy consumption. Executing a remote call requires the marshaling and unmarshaling of the actual parameters and final function value at the client and server sides. The compiler inserts the appropriate code based on the types of the formal parameters. Once a remote call is initiated, the client blocks until the remote call terminates. The preliminary prototype does not support on-demand checkpointing, and the client program will block in case of a disconnection event.

Preliminary experiments were performed for two simple programs, a selection sort and accumulation code, and a program that performs private RSA encryption [4]. The first code takes as input an integer array of size n and returns its sum, while the encryption code takes as input a message of size n , and returns the encrypted message. The experiments used a 233MHz Skiff board as its client machine, and a 440MHz SUN UltraSparc-10 workstation as the server. Figure 6 shows the execution times of the local and remote versions of the two example programs. A simple performance model is able to detect the cross-over points, allowing the correct selection between local and remote procedure execution. For large problem sizes, the remote version is up to 2.3x and 2.7x faster for the sort/accumulation and encryption programs, respectively, resulting in corresponding energy savings of up to 41% and 50%.

7 Related Work

Code generation for reducing energy consumption of an application has been discussed by researchers at Princeton and Fujitsu [25]. Their work showed that most classical optimizing for execution speed alone will also improve the overall power efficiency of an application. For instance, improving the locality of a computation will avoid cache misses which may require off-chip data to be fetched, an operation that consumes several orders of magnitude more energy than an access to an on-chip cache or register. In the context of wireless communications, most work has concentrated on low-power communication protocols [12] or using the wireless network as a secondary storage device [13]. The latter work mainly deals with problems arising from limited network bandwidth.

The Odyssey from CMU was designed by Satyanarayanan et al. to support a variety of mobile information access applications [14]. The applications run on the mobile client and access data objects on the remote server. The client tries to adapt to changes in the network bandwidth or other resources such as remaining battery life [7] by requesting data at different levels of fidelity or qual-

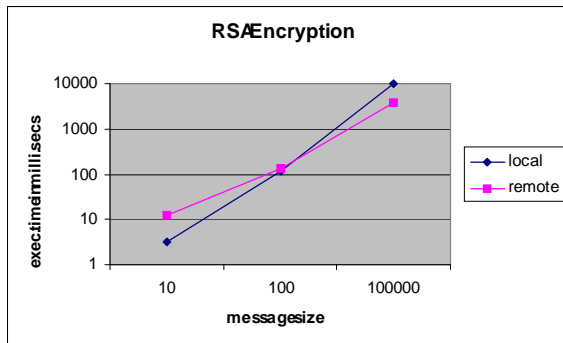
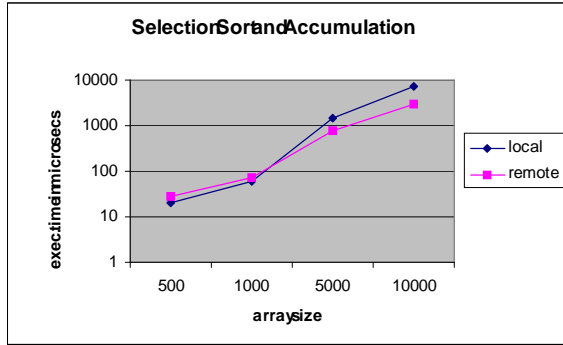


Figure 6: Measured execution times of local and remote versions of two sample programs.

ity. For example, fewer image frames may be requested in response to a reduction in communication bandwidth. The discussed approach relies on the semantic properties of the data processing done by the client. Data loss or skipping data, for instance by dropping a frame in the context of an image processing application, is not acceptable for all classes of applications. This application dependent aspect is important and effective, but orthogonal to the compiler approach discussed in this paper. In addition, their work does not deal with disconnection, i.e., total loss of network connectivity. However, Satyanarayanan et al. have developed operating system techniques that can deal with disconnectivity in the context of CODA, a distributed file system [21]. The CODA project addressed the important problem of establishing a consistent state of the file system across the client and server after network reconnection. More recently, Flinn et al. have developed Spectra, an environment that support remote execution as part of the Aura architecture, which includes the CODA and

Odyssey systems [6]. Spectra uses monitoring technology to balance different optimization goals of an application that contains Spectra operations and system calls.

Executing task on a remote server for the purpose of power savings is not a new idea. Researchers at UCLA and the University College London have investigated the profitability of remote task execution for operating system level tasks [19, 15]. In contrast, our compilation framework identifies within a single program candidate computations that are profitable to be executed remotely. Rudenko et al. have developed a remote processing framework that supports remote task mapping [19] at the OS level. Files are the objects that are communicated between a client and server. For a set of compilation tasks, the framework achieved power reductions on the order of 3 to 6 times as compared to the compilation done on the client machine. On a task set that includes text processing and Gaussian Elimination tasks in addition to compilation tasks, they report power consumption savings of up to 50% [18]. Battery lifetime extensions of up to 21% were reported by Othman and Hailes [15]. These results were obtained by simulation, and different assumptions about the available network bandwidth and CPU processing power of the client vs. the server machines. To the best of our knowledge, no compiler work has been done to support remote task execution for the purposes of energy savings.

Our proposed compile-time techniques identify remote subtask within a single program and determine when their execution may be profitable in terms of overall energy savings. For many applications, compilers are able to analyze entire programs and predict their future behavior and resource requirements. Operating systems and hardware approaches rely on the past observed program behavior to predict future program characteristics. Preliminary results of compiler support for dynamic voltage and frequency scaling for power and energy management purposes has been very encouraging [9, 10]. In addition, our work addresses the important issue of network disconnection and failure recovery. Finally, our power measurements are more precise than previously reported figures since we are able to measure instantaneous power directly at the hardware level, in particular for the Skiff system.

8 Conclusions and Future Work

This paper presents a compilation framework that identifies program regions that may be candidates for remote execution. To the best of our knowledge, we are the first to discuss image processing application in the context of mo-

mobile environments, and report actual power measurements for such an application on existing low-power platforms, instead of reporting power predictions based on simulation.

Remote task execution can be an effective compiler optimization for image processing applications. Preliminary experiments on StrongARM based systems show potential energy savings of one order of magnitude, mainly due to the reduction in overall execution times. On a PentiumII based laptop computer, up to 30% of the energy was saved. However, the experiments also show that remote tasks have to be selected carefully since a poor choice can lead to a significant increase in overall energy consumption. An initial prototype system has been implemented based on the SUIF2 compiler infrastructure [1]. Experimental results show the effectiveness of our compilation framework for two sample programs. The compiler was able to correctly choose between local and remote execution of procedure calls, resulting in up to 41% overall energy savings for the first, and 50% overall energy savings for the second sample code.

The initial compiler prototype is currently extended to perform interprocedural safety and more advanced profitability analyses for remote execution. In addition, compiler techniques are investigated that will initiate transitions between hibernation states of system components, allowing additional energy savings.

Acknowledgement

We would like to thank Christoph Peery, Steve Sanbeg and Kiran Nagaraja for their work on the preliminary prototype implementation of the compilation framework.

References

- [1] National Compiler Infrastructure (NCI) project. Overview available online at <http://www-suif.stanford.edu/suif/nci/index.html>, Co-funded by NSF/DARPA, 1998.
- [2] V. Balasundaram, G. Fox, K. Kennedy, and U. Kremer. A static performance estimator to guide data partitioning decisions. In *ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 213–223, Williamsburg, VA, April 1991.
- [3] R. Chatterjee, B. Ryder, and W. Landi. Relevant context inference. In *ACM SIGPLAN Symposium on the Principles of Programming Languages (POPL)*, pages 133–146, January 1999.
- [4] T. H. Cormen, C. E. Leiserson, and R. L. Rivest. *Introduction to Algorithms*. The MIT Press, Cambridge, MA, 1990.
- [5] K. Farkas, J. Flinn, G. Back, D. Grunwald, and J. Anderson. Quantifying the energy consumption of a pocket computer and a Java virtual machine. In *ACM International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS)*, Santa Clara, CA, June 2000.
- [6] J. Flinn, D. Narayanan, and M. Satyanarayanan. Self-tuned remote execution for pervasive computing. In *Hot Topics on Operating Systems (HotOS-VIII)*, Schloss Elmau, Germany, May 2001.
- [7] J. Flinn and M. Satyanarayanan. Energy-aware adaptation for mobile applications. In *ACM Symposium on Operating Systems Principles (SOSP)*, December 1999.
- [8] J. Hicks and U. Kremer. Skiff - a platform for pervasive computing. Technical Report 2000-3, Compaq Cambridge Research Laboratory (CRL), 2000.
- [9] C.-H. Hsu, U. Kremer, and M. Hsiao. Compiler-directed dynamic frequency and voltage scheduling. In *Workshop on Power-Aware Computer Systems (PACS'00)*, Cambridge, MA, November 2000.
- [10] C.-H. Hsu, U. Kremer, and M. Hsiao. Compiler-directed dynamic voltage/frequency scheduling for energy reduction in microprocessors. In *Proceedings of the International Symposium on Low-Power Electronics and Design (ISLPED'01)*, August 2001.
- [11] M. Kandemir, N. Vijaykrishnan, M.J. Irwin, and H.S. Kim. Experimental evaluation of energy behavior of iteration space tiling. In *International Workshop on Languages and Compilers for Parallel Computing (LCPC)*, August 2000.
- [12] R. Kravets and P. Krishnan. Power management techniques for mobile communication. In *Proceedings of the Annual International Conference on Mobile Computing and Networking (MOBI-COM)*, Dallas, TX, 1999.
- [13] J. Lorch and A. Smith. Software strategies for portable computer energy management. *IEEE Personal Communications Magazine*, 5(3), June 1998.
- [14] B. Noble. System support for mobile, adaptive applications. *IEEE Personal Communications*, 7(1):44–49, February 2000.
- [15] M. Othman and S. Hailes. Power conservation strategy for mobile computers using load sharing. *Mobile Computing and Communications Review*, 2(1):44–50, 1998.
- [16] Y. Park and B. Goldberg. Escape analysis on lists. In *ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI'92)*, pages 116–127, June 1992.
- [17] H.A. Rowley, S. Baluja, and T. Kanade. Neural network-based face detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(1):23–38, 1998.
- [18] A. Rudenko, P. Reiher, G. Popek, and G. Kuenning. Saving portable computer battery power through remote process execution. *Mobile Computing and Communications Review*, 2(1):19–26, 1998.
- [19] A. Rudenko, P. Reiher, G. Popek, and G. Kuenning. The remote processing framework for portable computer power saving. In *Proceedings of the ACM Symposium on Applied Computing (SAC99)*, San Antonio, TX, February 1999.
- [20] R. Saavedra-Barrera. *CPU Performance Evaluation and Execution Time Prediction Using Narrow Spectrum Benchmarking*. PhD thesis, U.C. Berkeley, February 1992. UCB/CSD-92-684.
- [21] M. Satyanarayanan. Mobile information access. *IEEE Personal Communications*, 3(1):26–33, February 1996.
- [22] T. Sim, R. Sukthankar, M. Mullin, and S. Baluja. High-performance memory-based face recognition for visitor identification. Technical report, Just Research, 1999.
- [23] T. Sim, R. Sukthankar, M. Mullin, and S. Baluja. Memory-based face recognition for visitor identification. In *Proc. 4th Intl. Conf. on Face and Gesture Recognition*, pages 214–220, Grenoble, France, March 2000.

- [24] T. Simunic, L. Benini, P. Glynn, and G. De Micheli. Dynamic power management for portable systems. In *Proceedings of the Sixth Annual International Conference on Mobile Computing and Networking (MobiCom)*, Boston, MA, August 2000.
- [25] V. Tiwari, S. Malik, A. Wolfe, and M. Lee. Instruction level power analysis and optimization of software. *Journal of VLSI Signal Processing*, 13(2/3):1–18, 1996.
- [26] M. Valluri and L. John. Is compiling for performance == compiling for power? In *The 5th Annual Workshop on Interaction between Compilers and Computer Architectures (INTERACT-5)*, January 2001.

A Basic Structure of TourGuide

System Initialization (SI)

- create in-memory representation of neural networks (for face detection)
- read-in face mask
- create in-memory representation of faces data base (for face recognition)

while active **do**

Image Acquisition (IA)

- obtain 640×480 pixel jpeg image

Preprocessing for Face Detection (PFD)

- decompress jpeg image into black&white pnm image
- crop image to 320×320 pixels (center section)
- normalize image for contrast
- build image pyramid with scale factor 1.1, where last five levels are pruned, resulting in 10 remaining levels

Face Detection (FD)

- for** $i = \text{highest_pyramid_level}$ **downto** $\text{lowest_pyramid_level}$ **do**
 - compute set of 30×30 pixels target regions that may contain a face (pre-filter)
 - foreach** target region **do**
 - detect 20×20 pixels face in target region
 - if** detection successful **then** exit both surrounding loops
 - endforeach**
- endfor**

Preprocessing for Face Recognition (PFR)

- Increase size of region with detected face by 20%
- Normalize detected face for contrast
- Scale-down detected face to 16×16 pixels

Face Recognition (FR)

- determine best matching score with 16×16 pixels faces in data base
- if** best match is above predefined quality threshold **then**
 - report successful recognition and return persons identity
- else**
 - report unsuccessful recognition
- endif**

enddo