



Energy trade-offs analysis using equal-energy maps



Maciej Drozdowski*, Jędrzej M. Marszałkowski, Jakub Marszałkowski

Institute of Computing Science, Poznań University of Technology, Piotrowo 2, 60-965 Poznań, Poland

HIGHLIGHTS

- The idea of isoenergy maps is introduced.
- Three models of energy consumption are proposed.
- Isoenergy maps for the proposed models are studied.
- Scalability of energy optimizations and their limitations are analyzed.

ARTICLE INFO

Article history:

Received 31 October 2012

Received in revised form

9 May 2013

Accepted 17 July 2013

Available online 26 July 2013

Keywords:

Energy use modeling

Performance evaluation

Visualization

Isolines

Parallel processing

Divisible load theory

ABSTRACT

Energy consumption has become a factor limiting further progress of supercomputing. Grasping the relationships determining the transformation of energy into computations is often difficult. Therefore, we propose a new method of visualizing the relationships as two-dimensional maps similar to isotherms or isobars in weather maps. Complex models of energy consumption are projected onto two-dimensional maps with isolines representing points of equal energy consumption. As an illustration of the concept, we present isoenergy maps for three models of parallel computations. The first two models are derived from Amdahl's and Gustafson's laws of parallel performance. The third model applies to divisible loads representing data-parallel computations in distributed systems. Equal-energy maps for the three models provide qualitative and quantitative insights into the interactions between certain energy-saving approaches, and their consequent limitations.

© 2013 Elsevier B.V. All rights reserved.

1. Introduction

Energy consumption is limiting further growth of supercomputers and datacenters [1–4]. For example, the top supercomputer in the June 2012 top500 list uses 7.9 MW to deliver 16.3E15 flops [5]. Scaling this technology to 1 exaflop would require roughly 480 MW. Such an installation would be prohibitively expensive and environmentally burdening. Thus, understanding the interactions between the determinants of energy use is essential for designing future computer systems.

Reduction of energy use can be achieved in many ways: by optimizing hardware, applications, and the methods of managing them. Possible hardware optimizations consist in efficient construction of CPU, memory, I/O and network devices, electric power distribution, and cooling systems. It is also necessary to open gates in the hardware to allow for power management. These can be,

for example, CPU dynamic voltage and frequency scaling (DVFS), slowdown and suspension modes of RAM, I/O devices under the ACPI framework, and similar methods for the network equipment. The management techniques are scheduling algorithms that coordinate the low-energy modes of the hardware with the state of the machine, network, and applications so that an acceptable trade-off between the performance and energy cost is found. Application enhancements consist in compiler optimizations, tuning application parameters, and designing energy-efficient algorithms right from the start. All these approaches interact with each other in the computing platform and the application.

As the techniques mentioned above are not completely independent, the savings achieved in one way may expose deficiencies elsewhere [6]. For example, parallelizing applications reduces the runtime and the cost of holding the computing equipment. But this hinges on communication, and hence on the efficiency of the network equipment. DVFS and I/O suspension methods reduce the computer power consumption, but this effect may be limited by the energy wasted in electric power supply and cooling systems. Furthermore, effective use of slowdown and suspension modes exposes the need for the specialized scheduling algorithms. Hence, it is necessary to have a global view of the factors and the optimization techniques, to understand their interrelated limitations and

* Corresponding author. Tel.: +48 616652981; fax: +48 618771525.

E-mail addresses: Maciej.Drozdowski@cs.put.poznan.pl (M. Drozdowski), Jedrzej.Marszalkowski@cs.put.poznan.pl (J.M. Marszałkowski), Jakub.Marszalkowski@cs.put.poznan.pl (J. Marszałkowski).

coordinated contribution to energy consumption. In this paper, we propose a new technique of visualizing relationships determining the performance of changing energy into computation.

The method is based on the idea of isolines widely used in science and engineering as, for example, cartography contour lines, isotherms, isobars, isogones, etc. This technique is building understanding of sensitivities and relationships in many systems and is guiding decision making. Here, maps with points of equal energy consumption (the isolines) will be presented. The relationships determining energy consumption can be perceived as a set of points of equal energy consumption in the multidimensional space of the system and application parameters. The proposed visualization method is a projection of such a multidimensional body onto a plane of just two parameters; this will be called an *isoenergy map*. Formally, the isoenergy map may be defined as follows. Let $E(x_1, \dots, x_n)$ be the energy usage depending on parameters x_1, \dots, x_n . Isoenergy map $IM(x_i, x_j)$ is a scalar field in domain $x_i \times x_j$ with values $E(x_1, \dots, x_n)$, where the parameters in set $\{x_1, \dots, x_n\} \setminus \{x_i, x_j\}$ are constant. An *isoenergy line* in $IM(x_i, x_j)$ is a set of points in $IM(x_i, x_j)$ with $E(x_1, \dots, x_n) = \text{const}$. For simplicity of perception, we will be referring to isoenergy maps as to collections of isoenergy lines, rather than fields of scalars. Before proceeding further, let us make a few observations on the patterns in isoenergy maps. A diagonal line in an isoenergy map means that changes in one parameter (e.g., x_i) must be accompanied by appropriate changes in the other parameter (x_j) to keep the energy constant. Lines perpendicular to one axis mean that the parameter on this axis dominates the energy usage. An isoline parallel to an axis means that changes of the parameter on that axis will not reduce the energy usage. A negative gradient of the scalar field shows the direction of possible energy savings. Thus, the relationships in the isoenergy maps guide decision making because they indicate the changes necessary to control energy use.

To perform a mapping of the multidimensional body of equal-energy points onto an isoenergy map, the body must be known. Equivalently, this means that a method of energy calculation for the given system and application parameters must be known. Energy consumption can be measured in existing systems, but not for systems yet to be built. It is economically unacceptable to build, for example, a big datacenter, just to check what would be the effect on performance of some system parameter tuning. Therefore, we are bound to use models of energy consumption to give directions for the changes needed. In this paper, we propose analytical models of energy consumption in parallel processing. The first two models are derived from Amdahl's [7] and Gustafson's speedup laws [8]. The third model is more detailed, and applies to divisible computations which represent data-parallel computations in distributed systems [9]. The contributions of this paper can be summarized as follows.

- The idea of isoenergy maps is introduced.
- Models of energy consumption are proposed.
- Isoenergy maps for the proposed models are studied for grasping the relationships guiding energy conservation optimizations and their limitations.

The rest of the paper is organized as follows. In the next section, we give an account on the work related to our study. In Section 3, energy use models and isoenergy maps for multicore systems are proposed. Similarly, the energy use model for divisible computations and the method of constructing the isolines are given in Section 4. The corresponding isoenergy maps are discussed in Section 5. The last section is dedicated to the conclusions. Notation used in the paper is summarized in Table 1.

Table 1
Summary of notation.

Symbol	Meaning
A	Computing rate (s/byte); see Section 4
C	Communication rate (s/byte); see Section 4
E	Energy (e.g. J)
f	Size of parallel part of the computation (Section 3)
k	Power reduction factor for the idle state
m	Number of processors
P_C	Processor power in active state (W); see Section 4
P_N	Network power in active state (W); see Section 4
S	Communication startup time (s); see Section 4
V	Size of load (bytes); see Section 4

2. Related work

Optimization of system and application energy use has been approached from different directions: for example, as a measurement and modeling issue, a scheduling and management issue, or a software and hardware construction problem.

Publications [10–14] serve as examples of the measurement and modeling direction. The energy cost and performance loss of parallel applications executed at different energy gears are empirically studied in [10]. An energy gear is a voltage–frequency combination of a CPU. In [11], multi-variable linear regression is used to model the execution time and energy consumption of the high-performance Linpack benchmark. In [12], a problem of constructing the shortest schedule for multi-phase parallel computation, meeting the energy limit, is considered. The energy use model distinguishing energy used in communication and in computation is experimentally validated. An index of iso-energy-efficiency is introduced in [13], as the ratio of the energy consumed in sequential computation to the energy consumed in parallel computation. Let us observe that, despite the similarity of the name the iso-energy-efficiency of [13] is conceptually different than the isoefficiency in [15–17] and the isoenergy maps in this paper. Analogously to [15–17], we consider isolines as relations (in the mathematical sense) linking system and application parameters such that the energy needed for the computation is constant. In [14], Amdahl's law is used to construct a general analytical model of energy consumption in multicore processors. The above papers introduce models of energy consumption. Some of the models are very detailed and tailored to a single platform and algorithm. Also, in this paper, we use energy consumption models. However, models with broader applicability will be used.

Energy optimization as a management and scheduling issue has been tackled in many publications. Scheduling parallel applications using various DVFS, low-power modes, virtualization, and load-shifting techniques under timing constraints is considered, for example, in [12,18–20]. Insofar as every scheduling model can be turned into a performance evaluation model, in this paper we use a scheduling model to evaluate performance.

The problem of effective algorithms, hardware, and their co-design has been studied, for example, in [21–23,6,24,25]. The issue of energy-efficient algorithms can be illustrated with the example of data compression. It is believed that compression may provide performance benefits and energy savings when transferring data between remote computers or different levels of memory hierarchy. However, it is demonstrated in [22,25] that the real picture is much more complicated. Only some compression algorithms, for some types of data, give any gain in energy. In [23], techniques of energy-efficient hardware and software design are reviewed. The authors distinguish three phases of system design: modeling and conceptualization, design and implementation, and runtime operation. Construction of efficient interconnections for big datacenters is studied in [21,24]. It appears that not only do computers consume considerable power, but also networking does. It is observed

Download English Version:

<https://daneshyari.com/en/article/425265>

Download Persian Version:

<https://daneshyari.com/article/425265>

[Daneshyari.com](https://daneshyari.com)