



Contents lists available at ScienceDirect

Journal of Systems Architecture

journal homepage: www.elsevier.com/locate/sysarc

Energy efficient task allocation for hybrid main memory architecture

Xiaojun Cai*, Lei Ju, Xin Li, Zhiyong Zhang, Zhiping Jia*

School of Computer Science and Technology, Shandong University, China

ARTICLE INFO

Article history:

Received 26 November 2015

Revised 11 April 2016

Accepted 1 June 2016

Available online xxx

Keywords:

DRAM

PRAM

Hybrid main memory architecture

Task allocation

ABSTRACT

Compared with the conventional dynamic random access memory (DRAM), emerging non-volatile memory technologies provide better density and energy efficiency. However, current NVM devices typically suffer from high write power, long write latency and low write endurance. In this paper, we study the task allocation problem for the hybrid main memory architecture with both DRAM and PRAM, in order to leverage system performance and the energy consumption of the memory subsystem via assigning different memory devices for each individual task. For an embedded system with a static set of periodical tasks, we design an integer linear programming (ILP) based offline adaptive space allocation (offline-ASA) algorithm to obtain the optimal task allocation. Furthermore, we propose an online adaptive space allocation (online-ASA) algorithm for dynamic task set where arrivals of tasks are not known in advance. Experimental results show that our proposed schemes achieve 27.01% energy saving on average, with additional performance cost of 13.6%.

© 2016 Elsevier B.V. All rights reserved.

1. Introduction

Memory subsystem has significant impact on the performance and energy efficiency of contemporary computer systems. In addition to the power consumption of read/write operations, traditional DRAM constantly requires background and refresh power in order to retain data integrity. Researches show that the DRAM-based main memory accounts for 40% of the total energy consumption of modern computer systems [1,2]. For energy-constrained embedded systems with low-power embedded processors, memory management becomes a key consideration in energy-efficient system design.

Emerging non-volatile memory (NVM) technologies such as phase change random access memory (PRAM) attract extensive attention of both industry and research communities. PRAM can be used as the main memory since it has similar performance metrics as DRAM. Compared to DRAM, NVM typically has better power efficiency due to low background power and absence of refresh energy consumption. However, current NVM technologies such as PRAM usually suffer from low write endurance, as well as high write latency and power consumption. For instance, a detailed comparisons between DRAM and PRAM are shown in Table 1 as presented in [3].

Given the characteristics of DRAM and PRAM, many researchers have put forward the idea of design main memory architecture

with both DRAM and PRAM. Some studies suggest to use a small DRAM as an upper-level buffer for the main memory of only PRAM [4]. On the other hand, there also have been some work on having parallel PRAM and DRAM which constitute the unified main address space [5,6]. In this work, we focus on the latter approach where the entire main memory space with both DRAM and PRAM is managed by the operating system.

Most of the existing works on hybrid DRAM-PRAM memory focus on the architecture-level design methodologies. On the other hand, for application-specific embedded systems, it is of paramount importance to study the characteristics of the applications in order to achieve optimal system-level design choices. In particular, given the different memory access patterns of various applications, it is possible to determine an optimal ratio between DRAM and PRAM in the hybrid main memory architecture, as well as the corresponding task allocation schemes for high performance and energy efficiency design. For a given total main memory size, increasing the use of PRAM over DRAM saves static energy, but it may lead to excessive writes to PRAM which shortens PRAM's lifetime and degrades the system performance. Although there have been several studies on reducing DRAM static energy consumption with task allocation and scheduling on digital signal processing (DSP) systems [7], the proportion of DRAM and PRAM in the hybrid architectures has not been considered. In this paper, we study the optimal ratio between DRAM and PRAM and task allocation problem in a hybrid main memory architecture for embedded systems. We explore the trade-offs and propose task allocation schemes for hybrid main memory for utmost energy saving, while prolonging

* Corresponding authors.

E-mail addresses: xj_cai@sdu.edu.cn (X. Cai), jzp@sdu.edu.cn (Z. Jia).

Table 1
Performance comparison.

| | PRAM | DRAM |
|---------------------|--|--|
| Volatility | No | Yes |
| Access latency | Read: 60 ns Write: 100–1000 ns | Read: 20–60 ns Write: 20–60 ns |
| Lifetime limitation | Sustain $10^8 - 10^9$ writes | No |
| Power consumption | No refresh and activation power, less leakage and access power | Leakage, refresh, activation, and access power |

the write endurance of PRAM and minimizing the performance overheads. If the original task set is schedulable with the traditional DRAM-based system, our proposed task allocation schemes guarantee the schedulability of the task set with the hybrid main memory architecture. The main contributions of this paper are as follows.

- (1) For a given set of periodic tasks, we propose an Integer Linear Programming (ILP) based offline Allocation algorithm to explore the optimal ratio between DRAM and PRAM, which reduces the write operations on PRAM and leverages between the energy consumption and performance overhead.
- (2) In order to enable fast design space exploration, we design a heuristic algorithm, i.e. offline Adaptive Space Allocation algorithm (offline-ASA) for task allocation on the hybrid memory architecture, which achieves near-optimal results in polynomial time.
- (3) For a task set where the arrival time of individual task is unknown, we propose an online Adaptive Space Allocation algorithm (online-ASA) to balance the size of DRAM and PRAM while obtaining the minimum energy consumption and wear leveling of PRAM.
- (4) We have performed experiments to evaluate the proposed algorithms. Results show that the proposed ILP and offline-ASA achieve 42.8% and 35.1% energy saving, at a cost of 35.4% and 17.1% performance degradation, respectively. On the other hand, the online-ASA leads to 27.01% energy saving with a 13.6% performance overhead.

The rest of this paper is organized as follows. Section 2 describes the previous work and the overview of the proposed strategy. Section 3 presents the background of the hybrid main memory architecture, describes the energy and calculation model and provides the problem description. Section 4 gives the ILP formulations and offline-ASA algorithm description. The discussion about online-ASA algorithm is explained in Section 5. Section 6 presents the experimental results, and also gives out the detailed results of the analysis. Finally, this paper is concluded in Section 7.

2. Related work

In the past decades, phase change RAM (PRAM), one of the emerging non-volatile memory technologies, has been comprehensively studied as a promising main memory candidate. Compared with the traditional DRAM, PRAM owes the advantages such as non-volatility, higher density, higher throughput, less leakage power consumption, etc. However, PRAM has deficiencies of limited write endurance and longer access latency. Therefore, the hybrid main memory architecture, which is consisting of both DRAM and PRAM, has been proposed to benefit from both memory technologies. In this architecture the high-speed DRAM and the energy-saving PRAM are carefully managed by allocating write-intensive pages into DRAM, such that the write endurance of PRAM could be prolonged and long PRAM write latency is mitigated.

For wear leveling, N-chance cache replacement policy was proposed in [8], and the system architecture of PRAM translation

layer was also presented in [9]. [10] proposed a new dynamic wear-leveling method that reduces unnecessary data migrations by adopting a hot/cold swapping based dynamic method. Meanwhile, [11–13] presented the scheduling techniques to reduce the write activity. Page migration between DRAM and PRAM was put forward to prolong the lifetime of PRAM in [5,14]. An intra-line flipping based scheme is proposed as a fine-grained bit-level wear leveling scheme to balance writes within memory lines [15]. A wear leveling method coming with a security policy was presented in [16].

For performance optimization, a novelty architecture, where DRAM serves as a buffer of PRAM, was proposed in [4,17] to reduce the writes on PRAM and improve the performance of main memory. [18] presented a PRAM model, called PCRAMsim, to bridge the gap between the device-level and system-level research on PRAM technology. The results in [19] showed that the PRAM/DRAM hybrid main memory with a modest DRAM size could give comparable performance.

In modern embedded systems, an increasing amount of energy is consumed by the main memory. In [20] an RDRAM module was detailedly described, and the devices in this module could be in different power states, i.e. attention, standby, nap and power-down. The purpose was to save memory energy consumption by switching the state of the devices. [21] studied how to utilize PRAM for energy optimization in embedded systems. For energy saving, [22] presented two methods, DRAM bypass and dirty data keeping, for further reduction in refresh energy and memory access latency, respectively. And the power-aware methods for DSPs with hybrid PRAM and DRAM main memory were proposed in [7]. In [6], a new DRAM + PCM memory system design, which comprised a sophisticated memory controller and a page placement policy called Rank-based Page Placement, was proposed to consider the performance, power consumption and wear leveling. For task allocation, ILP formulations and heuristic algorithms were proposed to maximize the energy saving in [23].

In the hybrid main memory architecture, various factors and design goals including the performance, energy efficiency, write endurance are overlapping. Therefore, it is critical for an optimal system design to consider the balance between these factors. In this work, we consider the management of hybrid main memory at the level of task allocation, and focus on the decision of ratio between DRAM and PRAM for a given total size constraint. We study the trade-offs between various design choices and propose the task allocation scheme for the hybrid main memory to utmost save energy consumption while improving the write endurance of PRAM and minimizing the performance overheads.

3. Problem analysis

In this section, we first introduce the background of the hybrid main memory architecture to be studied in this paper, followed by the performance and energy model for the architecture. Finally, we describe our targeting problem and objectives in the design of hybrid memory subsystem.

Download English Version:

<https://daneshyari.com/en/article/4956209>

Download Persian Version:

<https://daneshyari.com/article/4956209>

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