

# Paging for Multi-Core Shared Caches<sup>\*</sup>

Alejandro López-Ortiz  
David R. Cheriton School of Computer Science  
University of Waterloo  
200 University Ave. West  
Waterloo, Ontario, N2L 3G1, Canada  
alopez-o@uwaterloo.ca

Alejandro Salinger  
David R. Cheriton School of Computer Science  
University of Waterloo  
200 University Ave. West  
Waterloo, Ontario, N2L 3G1, Canada  
ajsalinger@uwaterloo.ca

## ABSTRACT

Paging for multi-core processors extends the classical paging problem to a setting in which several processes simultaneously share the cache. Recently, Hassidim proposed a model for multi-core paging [25], studying cache eviction policies for multi-cores under the traditional competitive analysis metric and showing that LRU is not competitive against an offline policy that has the power to arbitrarily delay request sequences to its advantage. While Hassidim brought attention to this problem, an effective and realistic model with accompanying competitive caching algorithms remains to be introduced.

In this paper we propose a more conventional model in which requests must be served as they arrive. We study the problem of minimizing the number of faults, deriving bounds on the competitive ratios of natural strategies to manage the cache. We show that traditional online paging algorithms are not competitive in our model. We then study the offline paging problem and show that the problem of deciding if a request can be served such that at a given time each sequence has faulted at most a given number of times is NP-complete and that its optimization version is APX-hard (for an unbounded number of sequences). We show as well that although offline algorithms can benefit from properly aligning future requests by means of faults, an algorithm that does so by forcing faults on pages that it has in its cache has no advantage over an honest algorithm that evicts pages only when faults occur. Lastly, we describe offline algorithms for the decision problem and for minimizing the total number of faults that run in polynomial time in the length of the sequences.

## Categories and Subject Descriptors

F.1.2 [Computation by Abstract Devices]: Modes of Computation—*Parallelism and concurrency, Online compu-*

<sup>\*</sup>A two page research announcement of this result appeared in the brief announcement session in SPAA'11 [28].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ITCS '12 Cambridge, Massachusetts USA

Copyright 2012 ACM 978-1-4503-1115-1/12/01 ...\$10.00.

*tation*; F.2.m [Analysis of Algorithms and Problem Complexity]: Miscellaneous

## General Terms

Algorithms, Theory

## Keywords

Cache, chip multiprocessor, multi-core, online algorithms, paging

## 1. INTRODUCTION

The paging problem models a two-level memory system consisting of a slow memory of infinite size and a fast memory of size  $K$  (the cache). The input is a sequence of page requests. Upon a request to a page  $\sigma$ , if  $\sigma$  is in the cache (hit) no action is required. Otherwise  $\sigma$  must be brought from slow to fast memory (fault), possibly requiring the eviction of another page in the cache. A paging algorithm must decide which page to evict upon each fault so as to minimize the number of faults.

In the last few years, multi-core processors have become the dominant processor architecture. While cache eviction policies have been widely studied both in theory and practice for sequential processors, the performance of even the most common eviction policies is not yet fully understood for the case in which various simultaneous processes share a common cache. In particular, there is almost no theoretical backing for the use of current eviction policies in multi-core processors. Recently, a work by Hassidim [25] initiated the theoretical study of paging strategies for shared caches in multi-cores or Chip Multiprocessors (CMPs). In a CMP system with  $p$  cores, a shared cache might receive up to  $p$  page requests simultaneously. Hassidim proposed a somewhat unconventional model in which the paging strategy can schedule the execution of threads. While in principle there is no reason why this cannot be so, historically the operating system has kept the scheduling of execution and the paging tasks separate. Within the operating system, the scheduler concentrates on fairness and throughput considerations to determine which task should be executed while the paging algorithm focuses on which of the pages currently in the cache should be evicted upon a fault.

The introduction of an accurate and effective model together with accompanying competitive paging algorithms for this natural extension of the classic paging problem remains an important open research problem. In this work we propose a more conservative and conventional model for

multi-core paging, in which cache algorithms are not allowed to make any scheduling decisions but must serve all active requests. In this model, a paging strategy serves a set  $\mathcal{R}$  of  $p$  request sequences using a shared cache of size  $K$ . Requests can be served in parallel, thus various pages can be read from cache or fetched from memory simultaneously, and a page fault delays the remaining requests of the corresponding sequence by  $\tau$  units of time. We define as FINAL-TOTAL-FAULTS (FTF) the problem of minimizing the total number of faults, and as PARTIAL-INDIVIDUAL-FAULTS (PIF) the problem of deciding, given a request sequence  $\mathcal{R}$ , a time  $t$ , and a bound vector  $\vec{b} \in \mathbb{N}^p$ , whether  $\mathcal{R}$  can be served such that at time  $t$  the number of faults on each of the sequences  $R_i$  is at most  $b_i$ .

Without loss of generality we define a paging strategy as a combination of a possible partition policy and an eviction policy, and compare the performance of natural strategies for FTF within this framework. We show that when restricted to static partition strategies, the choice of the partition has more impact than the choice of an eviction policy. We also show, however, that partition strategies cannot be competitive with respect to shared strategies if they do not update the partition often, even for disjoint request sequences. We finally show that shared strategies with traditional eviction policies (LRU, FIFO, CLOCK, and FWF) have competitive ratios that are arbitrarily large in the worst case.

We then study the offline cache problem and show properties of optimal offline algorithms. An algorithm that knows future requests can benefit from delaying one or more sequences with respect to others in order to balance the demands of different sequences. This could be achieved by evicting pages in order to force faults on pages that would otherwise result in hits. We show however that forcing faults in this manner is not beneficial. More specifically, we show that any offline algorithm for FTF that forces faults can be transformed into an *honest* algorithm that never evicts a page unless triggered by a fault, and which incurs in no more faults than the original algorithm. In particular, there exists an optimal algorithm that does not force faults. We show hardness results for PIF, showing that PIF is NP-complete and that a natural optimization version is APX-hard. Interestingly, PIF is NP-complete even in a simplified model with  $\tau = 0$ , i.e. when faults do not delay the remaining requests of a sequence. This result immediately implies that this offline problem is hard in the multiapplication caching model [6, 14]. Finally we present optimal offline algorithms for both FTF and PIF that run in polynomial time in the length of the sequences (and exponential in the number of processors<sup>1</sup>).

The rest of the paper is organized as follows. We review related work in Section 2. In Section 3 we describe the CMP cache model and formally define the problems we address in this paper. In Section 4 we derive bounds on the performance of natural strategies to minimize the number of faults. We study the offline problem in Section 5. We provide concluding remarks and future directions of research in Section 6. We include in most cases only proof sketches describing the main ideas, while full proofs are provided in the Appendix.

<sup>1</sup>Observe that for multi-cores the number of processors  $p$  is much smaller than the problem size  $n$ , in fact, we argued elsewhere that in practice  $p = O(\log n)$  [19].

## 2. RELATED WORK

The performance of the cache in the presence of multiple threads has been extensively studied, and research on the subject has increased markedly since the appearance of mainstream multi-core architectures. A variety of works have studied cache strategies in practice, developing heuristics to dynamically partition the cache [36, 38, 29, 15, 35] or to manage cache at the operating system level [20, 41, 31, 26].

From a theoretical perspective, much attention has been given recently to schedulers and algorithms for multithreaded computations with provable cache performance, either for private caches [11, 1, 4, 5, 2], shared caches [9], both private and shared [10, 16], or hierarchical cache configurations [8, 17, 40]. In these works the focus is on the design of either algorithms for specific problems or on general schedulers that yield cache performance guarantees. Cache replacement policies are not directly studied and usually an ideal cache is assumed. In this work, we focus directly on cache replacement policies independent of the application or specific thread schedule. In particular, we do not consider synchronization between threads.

More directly related to cache replacement policies, various models have been proposed to analyze the performance of paging algorithms in the presence of multiple request sequences, either modeling multiple applications or multiple threads. In what follows, we briefly review some of these models, which differ mainly in the assumptions they make with respect to the abilities of the paging algorithm to schedule page requests.

Fiat and Karlin [22] study paging algorithms in the access graph model, in which request sequences are restricted to paths in a graph [13]. They study the multi-pointer case, in which several paths through an access graph might be performed simultaneously, modeling both different applications (the graph could be disconnected) or multithreaded computations (having several paths in one same connected component). The order of requests is independent of the paging algorithm, which makes the paging problem substantially different from the one we study in this paper.

Barve *et al.* [6] study multiapplication caching in the competitive analysis framework [34]. In this model a set of page requests corresponding to different applications is served with a shared cache, and bounds are given with respect to a worst possible interleaving of the request sequences. As in Fiat and Karlin's model, however, the order of requests is the same for all algorithms.

Feuerstein and Strejilevich de Loma introduced Multi-Threaded Paging (MTP) [21]. In this problem, given a set of page requests, an algorithm must decide at each step which request to serve next, and how to serve it. They show that when no fairness restrictions are imposed there exist algorithms with bounded competitive ratio, but that no competitive algorithms exist when general fairness restrictions are imposed. Results in this model were further extended in [37] and [33]. In this model, a paging algorithm has the capability to schedule requests, and thus the order in which requests are served is algorithm dependent.

A recent paper by Hassidim [25] introduced a model for cache replacement policies specific to multi-core caches and studied the performance of algorithms in the competitive analysis framework with makespan as the performance measure. Hassidim's model includes the fetching time of pages

from memory. Thus if the algorithm incurs a fault on a sequence, it can continue serving other sequences while the faulting sequence’s page is fetched from memory. Hassidim shows that the competitive ratio of LRU with a cache of size  $K$  is  $\Omega(\tau/\alpha)$ , where  $\tau$  is the ratio between miss and hit times, and the offline optimal has a cache of size  $K/\alpha$ . As in the MTP model, Hassidim’s model assumes that the paging strategy can choose to serve requests of some sequences and delay others. In particular, the offline strategy is able to modify the schedule of requests, and hence is more powerful than a regular cache eviction algorithm. We discard this possibility, assuming that the order in which requests of different processors arrive to the cache is given by a scheduler over which the caching strategy has no influence. Given a request, the algorithm must serve the request either from cache or slow memory, and it cannot be delayed.

Hence, our model is different from previous models in that we assume no explicit scheduling capabilities of the paging strategy, while at the same time faults introduce delays in sequences, thus changing the order of requests.

### 3. THE CACHE MODEL

The model we use in this paper is broadly based on Hassidim’s model [25]. We have a multi-core processor with  $p$  cores  $\{1, \dots, p\}$ , and a shared cache of size  $K$  pages. The input is a multiset of request sequences  $\mathcal{R} = \{R_1, \dots, R_p\}$ , where  $R_j = \sigma_{s_1}^j, \dots, \sigma_{s_{n_j}}^j$  is the request sequence of core  $j$  of length  $n_j$ .  $\sigma_{s_i}^j$  is the identifier of the  $i$ -th page of the request, with  $1 \leq s_i \leq N$ , where  $N$  is the size of the universe of pages. The total number of page requests is  $n = \sum_{j=1}^p n_j$ . We assume  $K \gg p$  and  $n_j \gg K$ , for all  $1 \leq j \leq p$ . In particular, we assume that  $K \geq p^2$ , which can be regarded as a CMP variant of the tall cache assumption. We say that a request  $\mathcal{R}$  is *disjoint* if  $\forall i, j, i \neq j, R_i \cap R_j = \emptyset$  and *non-disjoint* otherwise. In practice, a single instruction of a core can involve more than one page. We treat each request as a request for one page, which models the case of separate data and instruction caches on a RISC architecture.

Page requests arrive at discrete timesteps. At any timestep, the cache might receive up to  $p$  page requests in parallel, and these are served in one parallel step. This assumes that requested pages from different cores can be read in parallel from the cache. We assume as well that fetching can be done in parallel, i.e. pages from memory corresponding to requests of different cores can be brought simultaneously from memory to the cache.

In our model, when a page request arrives, it must be serviced. The only choice the paging algorithm has is in which page to evict shall the request be a fault. A fault delays the remaining requests of the corresponding processor by an additive term  $\tau$ <sup>2</sup>. In other words, if a request  $\sigma_{s_{i^*}}^j$  is a fault, then for all  $i > i^*$ , the earliest time at which  $\sigma_{s_i}^j$  can be served increases by  $\tau$ .

To be consistent with [25], we adhere to the convention than when a page needs to be evicted to make space, first the page is evicted and the cache cell is unused until the fetching of the new page is finished. For non-disjoint requests we use the convention that when there is a request by processor  $j$ , of a page that is currently in the process of being fetched, then the sequence of processor  $j$  is only delayed until the page is

<sup>2</sup>Note that in [25]  $\tau$  is defined as the fetching time, which in this paper is  $\tau + 1$ .

fetched into the cache  $\tau$  units of time after the initial request. We also assume that cache coherency is provided at no cost to the algorithms. Finally, we adopt the convention that simultaneous requests are served logically in a fixed order (e.g. by increasing number of processor).

Under this model, various natural choices of objective functions may be considered. We define and address the following problems in this paper:

*Definition 1.* FINAL-TOTAL-FAULTS (FTF) Given a set of requests  $\mathcal{R} = \{R_1, \dots, R_p\}$ , a cache size  $K$ , and an integer  $\tau \geq 0$ , minimize the total number of faults when serving  $\mathcal{R}$  with a cache of size  $K$ .

*Definition 2.* PARTIAL-INDIVIDUAL-FAULTS (PIF) Given a set of requests  $\mathcal{R} = \{R_1, \dots, R_p\}$ , a cache size  $K$ , a time  $t$ , an integer  $\tau \geq 0$ , and  $\vec{b} \in \mathbb{N}^p$ , can  $\mathcal{R}$  be served with a cache of size  $K$  such that at time  $t$  the number of faults on each sequence  $R_i$  is at most  $b_i$ ?

Intuitively PIF is harder than FTF, since the former poses more restrictions on feasible solutions. Posing a bound on individual faults might be required to ensure fairness, and furthermore, doing so at arbitrary times can be used to ensure fairness throughout the execution of an algorithm.

### 4. BOUNDS OF ONLINE STRATEGIES FOR MINIMIZING FAULTS

Natural strategies to manage the cache in the multi-core cache model can be classified in two families: shared and partitioned. In the first one, the entire cache is shared by all processors, and a cache cell can hold a page corresponding to any processor. In the second one, the cache is partitioned in  $p$  parts, with each part destined exclusively to store pages of requests from one processor, throughout.

Both shared and partitioned strategies are accompanied by an eviction policy  $A$ . We use  $S_A$  to denote the algorithm that uses a shared cache with eviction policy  $A$ , and  $P_A^B$  to denote a partitioned strategy with partition function  $B$  and eviction policy  $A$  in each part. A partition function  $B$  is *static* if the size of all parts remain constant during an execution and is *dynamic* otherwise. For the latter, when the reduction of the size of a part involves page evictions, these are carried out according to the eviction policy. We make the restriction that all partitions must assign at least one unit of cache to all processors whose requests are active. For example,  $S_{LRU}$  evicts the least recently used page in the entire cache and  $P_{LRU}^{OPT}$  performs LRU on each part of the partition, which is determined offline so as to minimize the total number of faults. Table 1 shows an example of the execution of a shared strategy and a static partition strategy.

In the remainder of this section we compare the performance of partitioned and shared strategies for FTF. We denote the number of faults of a strategy  $Alg$  on a sequence  $\mathcal{R}$  as  $Alg(\mathcal{R})$ .

Partitioning the cache may be desirable to avoid costs of managing concurrency issues that arise when different threads access a shared page. In addition, a static partition allows for the execution of regular paging algorithms in each part, oblivious to the presence of other threads, and thus provides performance guarantees based on individual threads or processes. In fact, if we restrict paging strategies

$S_{LRU}$			$P_{LRU}^B$		
$t$	Remaining sequence	Cache	Remaining sequence	Cache	
0	$\sigma_1, \sigma_2, \sigma_4, \sigma_2, \sigma_3, \sigma_4$ <u><math>\sigma_5, \sigma_6, \sigma_2, \sigma_4, \sigma_5, \sigma_5</math></u>	$\sigma_1\sigma_2\sigma_3\sigma_5\sigma_6$	$\sigma_1, \sigma_2, \sigma_4, \sigma_2, \sigma_3, \sigma_4$ <u><math>\sigma_5, \sigma_6, \sigma_2, \sigma_4, \sigma_5, \sigma_5</math></u>	$\sigma_1\sigma_2\sigma_3$ $\sigma_5\sigma_6$	
1	<u><math>\sigma_2, \sigma_4, \sigma_2, \sigma_3, \sigma_4</math></u> $\sigma_6, \sigma_2, \sigma_4, \sigma_5, \sigma_5$	$\sigma_5\sigma_1\sigma_2\sigma_3\sigma_6$	$\sigma_2, \sigma_4, \sigma_2, \sigma_3, \sigma_4$ <u><math>\sigma_6, \sigma_2, \sigma_4, \sigma_5, \sigma_5</math></u>	$\sigma_1\sigma_2\sigma_3$ $\sigma_5\sigma_6$	
2	<u><math>\sigma_4, \sigma_2, \sigma_3, \sigma_4</math></u> <u><math>\sigma_2, \sigma_4, \sigma_5, \sigma_5</math></u>	$\sigma_6\sigma_2\sigma_5\sigma_1\sigma_3$	<u><math>\sigma_4, \sigma_2, \sigma_3, \sigma_4</math></u> <u><math>\sigma_2, \sigma_4, \sigma_5, \sigma_5</math></u>	$\sigma_2\sigma_1\sigma_3$ $\sigma_6\sigma_5$	
3	$\perp, \perp, \sigma_2, \sigma_3, \sigma_4$ <u><math>\sigma_4, \sigma_5, \sigma_5</math></u>	$\sigma_2\sigma_6\sigma_5\sigma_1^*$	$\perp, \perp, \sigma_2, \sigma_3, \sigma_4$ $\perp, \perp, \sigma_4, \sigma_5, \sigma_5$	$\sigma_2\sigma_1^*$ $\sigma_6^*$	
4	$\perp, \sigma_2, \sigma_3, \sigma_4$ $\perp, \sigma_5, \sigma_5$	$\sigma_2\sigma_6\sigma_5\sigma_1^*$	$\perp, \sigma_2, \sigma_3, \sigma_4$ $\perp, \sigma_4, \sigma_5, \sigma_5$	$\sigma_2\sigma_1^*$ $\sigma_6^*$	
5	$\sigma_2, \sigma_3, \sigma_4$ $\sigma_5, \sigma_5$	$\sigma_4\sigma_2\sigma_6\sigma_5\sigma_1$	$\sigma_2, \sigma_3, \sigma_4$ <u><math>\sigma_4, \sigma_5, \sigma_5</math></u>	$\sigma_4\sigma_2\sigma_1$ $\sigma_2\sigma_6$	
6	<u><math>\sigma_3, \sigma_4</math></u> $\sigma_5$	$\sigma_5\sigma_2\sigma_4\sigma_6\sigma_1$	<u><math>\sigma_3, \sigma_4</math></u> $\perp, \perp, \sigma_5, \sigma_5$	$\sigma_2\sigma_4\sigma_1$ $\sigma_2^*$	
7	$\perp, \perp, \sigma_4$	$\sigma_5\sigma_2\sigma_4\sigma_6^*$	$\perp, \perp, \sigma_4$ $\perp, \sigma_5, \sigma_5$	$\sigma_2\sigma_4^*$ $\sigma_2^*$	
8	$\perp, \sigma_4$	$\sigma_5\sigma_2\sigma_4\sigma_6^*$	$\perp, \sigma_4$ <u><math>\sigma_5, \sigma_5</math></u>	$\sigma_2\sigma_4^*$ $\sigma_4\sigma_2$	
9	$\sigma_4$	$\sigma_3\sigma_5\sigma_2\sigma_4\sigma_6$	$\sigma_4$ $\perp, \perp, \sigma_5$	$\sigma_3\sigma_2\sigma_4$ $\sigma_4^*$	
10		$\sigma_4\sigma_3\sigma_5\sigma_2\sigma_6$	$\perp, \sigma_5$	$\sigma_4\sigma_3\sigma_2$ $\sigma_4^*$	
11			$\sigma_5$	$\sigma_4\sigma_3\sigma_2$ $\sigma_5\sigma_4$	
12				$\sigma_4\sigma_3\sigma_2$ $\sigma_5\sigma_4$	

**Table 1: Example of execution of shared LRU and a static partition strategy with LRU and partition  $B = \{3, 2\}$  on the input  $\mathcal{R} = \{R_1, R_2\}$ , with  $R_1 = \sigma_1, \sigma_2, \sigma_4, \sigma_2, \sigma_3, \sigma_4$  and  $R_2 = \sigma_5, \sigma_6, \sigma_2, \sigma_4, \sigma_5, \sigma_5$ . The cache size is  $K = 5$  and  $\tau = 2$ . Underlined pages denote faults, and a ‘ $\perp$ ’ indicates a timestep in which a page is being fetched. Pages in the caches are shown from left to right in order of most recent use, and a ‘\*’ in the cache indicates that the cell will be used by a page currently being fetched. The number of faults incurred by  $S_{LRU}$  and  $P_{LRU}^B$  are 3 and 5, respectively.**

to a fixed static partition, any marking or conservative algorithm<sup>3</sup> (e.g. LRU, FIFO) has a competitive ratio of at most  $K$ , as in the sequential setting. Formally:<sup>4</sup>

**LEMMA 1. Online vs. offline eviction policies with a fixed static partition.** *Let  $A$  be any deterministic online eviction algorithm and let  $B = \{k_1, k_2, \dots, k_p\}$  be any online static partition. There exists a sequence  $\mathcal{R}$  such that  $P_A^B(\mathcal{R})/P_{OPT}^B(\mathcal{R}) = \Omega(\max_j \{k_j\})$ . When  $A$  is any marking or conservative algorithm (e.g. LRU), there is a matching upper bound, i.e.,  $\forall \mathcal{R}, P_A^B(\mathcal{R})/P_{OPT}^B(\mathcal{R}) \leq \max_j \{k_j\}$ .*

The choice of a good partition has more influence on the performance of static partition strategies than the eviction policy. In fact, if the partition function can be computed offline, then no online static partition strategy is competitive, even with an offline eviction policy.

**LEMMA 2. Online static partition strategies are not competitive.** *Let  $B = \{k_1, \dots, k_p\}$  be any online static partition. Let  $k^* = \min_j \{k_j | k_j \geq 2\}$ . Then there exists a sequence  $\mathcal{R}$  s.t. for all eviction policies  $A$ ,  $P_A^B(\mathcal{R})/P_{LRU}^{OPT}(\mathcal{R}) \geq \min\{k^*, p-1\} \frac{n}{K^2 p} = \Omega(n)$ .*

**Proof sketch<sup>4</sup>:** Consider  $A=LRU$ . All processors but one request repeatedly  $k_j + 1$  different pages, and the other one requests repeatedly the same page. Thus,  $P_A^B$  faults on every request of  $p - 1$  processors. An optimal partition allocates

<sup>3</sup>See [12, Ch. 3] for the definitions of marking and conservative algorithms.

<sup>4</sup>See the Appendix for full proofs.

enough cache for all processors and hence it faults a constant number of times. The result for any algorithm  $A$  (possibly offline) follows by Lemma 1.

Although strategies that partition the cache only once or that even allow a small number of changes during the execution might be simpler to manage as compared to general strategies, the performance of these strategies is not competitive when shared strategies are allowed. While this is perhaps to be expected for non-disjoint sequences, interestingly this holds even for disjoint sequences. Theorem 1 shows that shared strategies are preferable over static partitions—even if the partition is computed offline—as well as dynamic partitions that do not change often.

**THEOREM 1.** *Let  $A$  be any deterministic online cache eviction policy, let  $sOPT$  be an optimal static partition, and let  $D$  be any online dynamic partition strategy that changes the sizes of the parts  $o(n)$  times. The following statements hold:*

1. There exists a sequence  $\mathcal{R}$  s.t.  $\frac{P_{OPT}^{sOPT}(\mathcal{R})}{S_{LRU}(\mathcal{R})} = \Omega(n)$ .
2. For all sequences  $\mathcal{R}$ ,  $S_{LRU}(\mathcal{R})/P_{OPT}^{sOPT}(\mathcal{R}) \leq K$ .
3. There exists a sequence  $\mathcal{R}$  s.t.  $P_A^D(\mathcal{R})/S_{LRU}(\mathcal{R}) = \omega(1)$ . Furthermore, if  $D$  varies the partition a constant number of times,  $P_A^D(\mathcal{R})/S_{LRU}(\mathcal{R}) = \Omega(n)$ .

**Proof sketch<sup>4</sup>:**

1. Given any static partition (even the optimal offline), there is always a sequence that can demand more than

the part assigned to each processor at different mutually exclusive times, while a shared strategy can use the entire cache when serving the high demand period of each sequence.

2. Divide a sequence  $R_j$  in phases such that a new phase starts every time there is a request for the  $(k_j + 1)$ -th distinct page since the beginning of the previous phase, and the first phase begins at the first page of  $R_j$ . Define a shared phase of  $\mathcal{R}$  analogously for a cache of size  $K$ , with the order of page requests given by the execution of  $S_{LRU}$  on  $\mathcal{R}$ . A shared phase cannot start and end without the sequence of some processor changing phases, thus there are no more shared phases than phases in the sequences themselves.  $S_{LRU}$  faults at most  $K$  times per shared phase, and a static partition strategy must fault at least once in each phase of each processor's sequence.
3. Since the number of times the partition changes is  $o(n)$ , at least one period with constant partition has non-constant length. We can thus apply the same argument as in the proof in (1) for static partitions in this period. Each processor requests more pages than the size of its part in this period, but at different times, allowing a shared strategy to use all the cache to serve the requests of each processor.

Theorem 1 suggests that competitive strategies must either be shared or have a partition that changes often. In fact both types of strategies are equivalent for disjoint sequences. Although a dynamic partition strategy executes an eviction policy in each part separately, if the variation in the partition can be determined globally, then shared strategies can be simulated by a dynamic partition on disjoint sequences by reducing the part of the cache holding the page to be evicted. This implies that a dynamic partition strategy can be as effective as any other strategy. In fact, Hassidim [25] showed in his model that there exists an optimal dynamic partition strategy that upon a fault reduces the part of some processor, evicting the page that is furthest in the future in that processor's sequence. We show in Section 5 that the same result holds in our model.

Multi-core paging differs from sequential paging in that the actions of algorithms modify the order of future requests. Hence paging strategies must decide which page to evict not only with the goal of delaying further faults, but at the same time trying to properly align the demand periods of future requests. An online strategy, however, is oblivious to future requests and hence in general it cannot work toward the second goal. In this sense, in multi-core paging an optimal offline strategy has significantly more advantages over an online strategy than in sequential paging. While in the latter setting any online marking algorithm has a bounded competitive ratio of  $\frac{K}{K-h+1}$  when the offline has a cache of size  $h \leq K$  [27, 39], in multi-core paging the competitive ratio of traditional algorithms such as LRU, FIFO, CLOCK, and FWF can be arbitrarily large. Although the optimal offline strategy cannot explicitly schedule requests, it can increase the fault rate of a process thus effectively delaying that sequence in order to align periods of high demand with periods of low demands of other processors. The following theorem shows that the competitive ratio of  $S_{LRU}$  can grow proportionally to the square root of the length of the se-

quences, even when the offline strategy has a cache of about half the size. The theorem also applies to FIFO, CLOCK, and FWF.

**THEOREM 2. Lower bound on the competitive ratio of LRU.** *Assume  $p \geq 4$  and  $\tau > 0$ . There exists a sequence  $\mathcal{R}$  and an offline eviction policy  $A$  such that  $S_{LRU}(\mathcal{R})/S_A(\mathcal{R}) = \Omega(\sqrt{n\tau/K})$  when  $A$ 's cache size is  $h \geq K/2 + 3p/2$ .*

**PROOF.** The request sequence  $\mathcal{R}$  consists of two alternating phases which we call easy and hard, respectively (see Figure 1, top). In an easy phase, each sequence requests only 2 different pages, while in a hard phase, the total number of different pages requested is greater than  $K$ . More specifically, let

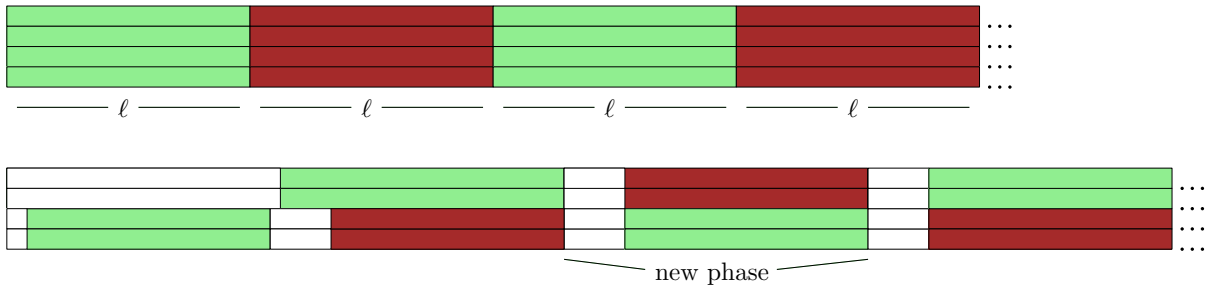
$$R_j = \left( (\sigma_a^j \sigma_b^j)^{\ell/2} (\sigma_1^j, \sigma_2^j, \dots, \sigma_{K/p+1}^j)^{\ell/(K/p+1)} \right)^\phi$$

The length of each phase is  $\ell$  pages and there are  $2\phi$  phases. The number of pages in a hard phase is  $K + p$ , and since every request is to the least recently used page,  $S_{LRU}$  faults on every request of a hard phase. Thus  $S_{LRU}(\mathcal{R}) \geq n/2$ , where  $n$  is the total number of pages.

An offline algorithm  $A$  can benefit from aligning hard phases of some sequences with easy phases of others, in order to keep the total number of requested pages below  $h$  in the aligned periods.  $A$  does this by initially assigning only one cell from the cache to a group of sequences  $\mathcal{R}_1 = \{R_1, \dots, R_{p/2}\}$ . Hence every request in these two sequences is a fault. Since  $A$  will fault only on the first two requests of the rest of the sequences  $\mathcal{R}_2 = \{R_{p/2+1}, \dots, R_p\}$ , sequences in  $\mathcal{R}_1$  will be delayed with respect to sequences in  $\mathcal{R}_2$ .  $A$  will delay sequences in  $\mathcal{R}_1$  so that the last page of their easy phase is aligned with the last page of the hard phase of sequences in  $\mathcal{R}_2$ . Consider a sequence  $R_i \in \mathcal{R}_2$ . The total number of pages requested during the first hard phase of  $R_i$  is  $(p/2)(K/p + 1) + p \leq h$ , and thus  $A$  faults only on the first  $K/p + 1$  requests of  $R_i$  in its first hard phase. Since  $A$  also faults twice in the first easy phase of  $R_i$ , the first hard phase of  $R_i$  is completed at time  $t = 2\tau + \ell + \tau(K/p + 1) + \ell = 2\ell + \tau(3 + K/p)$ . Therefore,  $A$  needs to incur in  $\ell/\tau + K/p + 3$  faults in each of the sequences in  $\mathcal{R}_1$  in order for their  $\ell$ -th request to be served at time  $t$ .

After the initial faults, easy phases of sequences in  $\mathcal{R}_1$  are aligned with hard phases of sequences in  $\mathcal{R}_2$  and vice versa (see Figure 1, bottom). Call each of these  $\ell$ -page phases a new phase. Since each new phase has at most  $h$  different pages,  $A$  faults only at the beginning of these phases.  $A$  keeps the alignment of sequences by partitioning the cache so that the first  $K/p + 1$  of all sequences are faults. Since any one sequence has at most  $K/p + 1$  distinct pages in a new phase,  $p(K/p + 1) \leq 2K$  faults are enough to maintain the alignment. These faults, plus the initial faults to arrange the alignment, add up to  $A(\mathcal{R}) \leq 4K\phi + (p/2)(\ell/\tau + K/p + 3)$ . The total number of pages is  $n = \phi\ell p$ , and thus  $\ell = n/(\phi p)$ . Substituting  $\ell$  in the number of faults and minimizing for  $\phi$  yields  $\phi = \sqrt{n/(8K\tau)}$  and thus  $A(\mathcal{R}) \leq 4\sqrt{nK/3\tau} + K/2 + 3p/2 = O(\sqrt{nK/\tau})$ . Since  $S_{LRU}(\mathcal{R}) \geq n/2$ , it follows that  $S_{LRU}(\mathcal{R})/S_A(\mathcal{R}) = \Omega(\sqrt{n\tau/K})$ .  $\square$

Note that if the offline algorithm has a cache of size  $K$  the lower bound in Theorem 2 can be made even larger: the offline algorithm delays only two sequences and the competitive ratio becomes  $\Omega(\sqrt{n\tau p/K})$ . The proof of Theorem 2



**Figure 1: Top: A request with 4 sequences that alternates phases of low and high demand. For each sequence, green (light) phases consist of alternating requests to two pages, while red (dark) phases are consecutive requests of  $K/p + 1$  pages. Bottom: sequence after being served by the offline strategy  $S_A$ . White periods denote faults. After the initial faults each new phase has at most  $h$  distinct pages, and  $K + p$  faults are necessary in each new phase to keep future alignments.**

can be used to show that not only Furthest-In-The-Future (FITF) [7] —an optimal algorithm in classical paging— is not optimal in multi-core paging, but that it performs badly compared to the true optimal.

**COROLLARY 1.** *Assume  $p \geq 4$  and  $\tau > 0$ . There exists a sequence  $\mathcal{R}$  and an offline eviction policy  $A$  such that  $S_{FITF}(\mathcal{R})/S_A(\mathcal{R}) = \Omega(\sqrt{n\tau}/K^{3/2})$  when  $A$ 's cache size is  $h \geq K/2 + 3p/2$ .*

**PROOF.** Let  $\mathcal{R}$  be the sequence given in the proof of Theorem 2. The number of faults of FITF is at least  $n/2K$ , and thus  $S_{FITF}(\mathcal{R})/S_{OPT}(\mathcal{R}) = \Omega(\sqrt{n\tau}/K^{3/2})$ .  $\square$

## 5. THE OFFLINE PROBLEM

In reality requests sequences are not known in advance, and thus paging is an online problem. In general, however, in any online problem setting deriving efficient optimal offline solutions is both of theoretical interest as well as useful in evaluating online algorithms in practice in the competitive analysis framework. Furthermore, an online solution can be designed based on properties of the offline solution. For example, in traditional paging, LRU approximates FITF using the past as the best approximation of the future. An example of an inherently online problem for which the offline problem has been extensively studied is the List Update problem (See, e.g., [32, 30, 3, 24]).

### 5.1 Hardness of Multi-Core Paging

In this section we show that even if the sequence of requests is known in advance multi-core paging is hard. More specifically, we show that PARTIAL-INDIVIDUAL-FAULTS is NP-complete and that there is no PTAS for its optimization version. This is in contrast to sequential paging: if there is only one processor, both FTF and PIF are solvable by FITF. As in the proof of Hassidim's makespan problem [25], the proof of NP-completeness of PIF uses a reduction from 3-PARTITION [23]. However, since our model disallows explicit scheduling, the reduction is quite different.

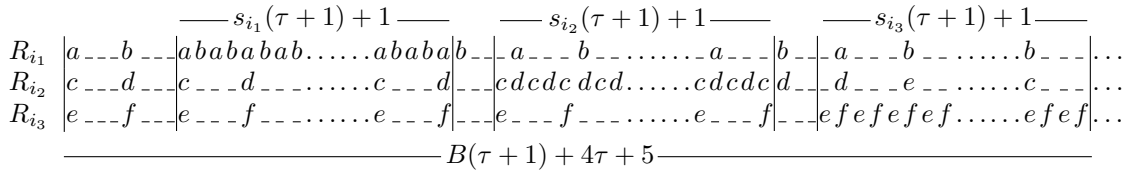
**THEOREM 3.** *PARTIAL-INDIVIDUAL-FAULTS is NP-complete.*

**Proof sketch**<sup>4</sup>: An instance of 3-PARTITION consists of a set of  $n$  integers  $S = \{s_1, \dots, s_n\}$ , and a bound  $B$ , such that

$B/4 < s_i < B/2$  for all  $1 \leq i \leq n$ . The problem is to determine if  $S$  can be partitioned into  $n/3$  sets  $A_1, \dots, A_{n/3}$  such that for all  $1 \leq j \leq n/3$ ,  $\sum_{i \in A_j} s_i = B$  and  $|A_j| = 3$ . Note that the restrictions of the problem imply that each subset  $A_j$  must have exactly 3 elements [23]. Given an instance  $\mathcal{J}$  of 3-PARTITION, we build an instance  $\mathcal{I}$  of PIF as follows. There are  $p = |S|$  sequences. Each sequence  $R_i$  consists of alternating requests to two distinct pages  $\alpha^i$  and  $\beta^i$ , i.e.  $R_i = \alpha^i \beta^i \alpha^i \beta^i \dots$ , and all sequences are disjoint. The length of  $R_i$  is  $|R_i| = B(\tau + 1) + 4\tau + 5$ , where  $\tau \geq 0$  is any integer. The size of the cache is  $K = (4/3)p$ , and the maximum allowed number of faults in each sequence  $R_i$  is  $b_i = B - s_i + 4$ , at time  $t = B(\tau + 1) + 4\tau + 5$ .

A solution to  $\mathcal{J}$  gives a solution to  $\mathcal{I}$  as follows. Let  $A_1, \dots, A_{n/3}$  be a solution to  $\mathcal{J}$ . Each group of three sequences corresponding to a group  $A_j$  will share a group of four cache cells. Each sequence in a group will have a dedicated cell at all times, while the extra fourth cell is assigned to each sequence at different times. It is not difficult to show that if each sequence  $R_i$  uses the extra cell continuously so that it incurs in exactly  $h_i = s_i(\tau + 1) + 1$  hits, then at time  $t = B(\tau + 1) + 4\tau + 5$  the number of faults of this sequence is exactly  $(t - h_i)/(\tau + 1) = B - s_i + 4$  (see Figure 2).

We argue that a solution to the instance  $\mathcal{I}$  gives a solution to  $\mathcal{J}$ . A sequence  $R_i$  must have at least  $h_i = s_i(\tau + 1) + 1$  hits by time  $t$  in order to satisfy the bound on the number of faults. Given the alternating pages in a sequence, a sequence can have a hit only if it has two cells for consecutive timesteps. Each sequence uses at least one cell until time  $t$ , and hence there are only  $p/3$  extra cells. Therefore there can be at most  $p/3$  hits in any timestep. In addition, it is not hard to see that any change in the sequence to which a cell is assigned implies at least  $\tau$  timesteps without hits for the two sequences involved. Considering the total number of possible hits until time  $t$  and the minimum total hits required, we can show that the maximum number of changes in the partition is  $2p/3$  and hence each sequence must have two cells for a continuous period of time. Furthermore, it can be shown that exactly three sequences must share one extra cell, otherwise at least one sequence will exceed the allowed faults before time  $t$ . Finally, it is not hard to show that the three sequences  $R_{i_1}, R_{i_2}, R_{i_3}$  that share one cell must be the ones whose corresponding  $s_i$ 's satisfy  $s_{i_1} + s_{i_2} + s_{i_3} = B$ . The groups of three sequences that share an extra cell define



**Figure 2:** Sequences  $\mathcal{R}_{i_1}, \mathcal{R}_{i_2}, \mathcal{R}_{i_3}$  share 4 cells of the cache. In the example,  $a, b, c, d, e, f$  are all different pages,  $\tau = 3$ , and the fetching period is represented by  $\tau$  consecutive underscore characters ( $\_$ ). Each sequence  $R_i$  holds the extra cell continuously so that it incurs in  $s_i(\tau + 1)$  hits. The number of faults of each  $R_i$  in the group is  $B - s_i + 4$  at time  $t = B(\tau + 1) + 4\tau + 5$  if and only if  $s_{i_1} + s_{i_2} + s_{i_3} = B$ .

a solution for the instance  $\mathcal{J}$  of 3-PARTITION.

Observe that PIF remains NP-complete even when  $\tau = 0$ , i.e. when sequences are not delayed due to faults. This means that this problem is hard also in the multiapplication caching model of Barve *et al.* [6]. Note that this is not the case for FINAL-TOTAL-FAULTS, for which FITF is optimal when  $\tau = 0$ . This provides evidence that achieving a fair distribution of faults is more difficult than merely minimizing the number of overall faults.

We show now that PIF is also hard to approximate. We define as MAX-PIF the problem of, given an instance of PIF, maximizing the number of sequences whose number of faults at a given time is within the given bound. We show that MAX-PIF is APX-hard, i.e. there is no polynomial time algorithm that can approximate this problem within a factor of  $(1 - \epsilon)$  for every  $\epsilon$ , unless  $P=NP$ . In order to show this, we describe a gap-preserving reduction from MAX-4-PARTITION (shown to be APX-hard in [18]). Therefore unless  $P=NP$ , there is no efficient way of serving the request sequences ensuring that an arbitrarily large part of them will fault within the allowed bounds.<sup>4</sup>

**THEOREM 4.** *MAX-PARTIAL-INDIVIDUAL-FAULTS is APX-hard.*

## 5.2 Properties of Offline Algorithms for FTF

The changes in the relative alignment of sequences can significantly affect the performance of an algorithm (see the proof of Theorem 2 for an example). Offline algorithms can benefit from properly aligning the demand periods of future requests by means of faults and their corresponding delays. For this purpose, an algorithm could evict a page voluntarily (i.e. not forced by a page fault) before it is requested in order to force a fault. We show, however, that forcing faults for the purpose of changing the alignments in this way is not beneficial for minimizing the number of faults. We say that an algorithm is *honest* if it does not evict a page unless there is a fault, and show that there exist an optimal algorithm that is honest.

**THEOREM 5.** *Let  $Alg$  be an offline optimal algorithm that forces faults. There exists an offline algorithm  $Alg'$  that is honest such that for all disjoint requests  $\mathcal{R}$ ,  $Alg'(\mathcal{R}) = Alg(\mathcal{R})$ .*

**Proof sketch<sup>4</sup>:** Let  $Alg$  be any offline algorithm. The proof is based on the following claim: for each timestep  $i$  we can build an algorithm  $Alg_i$  that behaves exactly like  $Alg$  until  $t = i - 1$  and that at time  $t = i$ , if  $Alg$  forces a fault, then  $Alg_i$  does not. Furthermore,  $Alg_i(\mathcal{R}) \leq Alg(\mathcal{R})$ . If this

claim is true, given any optimal algorithm, we can build an optimal algorithm that does not force faults by successively applying the claim for each  $i \geq 1$ .

In order to prove that the claim is true, we consider the execution of  $Alg$  and  $Alg_i$  after  $Alg$  has forced a fault at  $t = i$ . Both executions can differ in the contents of the cache, the relative alignment of the sequences, and the number of current faults. Based on these, we define 5 possible states in which the executions can be, for which all satisfy  $Alg_i(\mathcal{R}, t) \leq Alg(\mathcal{R}, t)$  at the current timestep, where  $Alg(\mathcal{R}, t)$  is the number of faults of  $Alg$  at time  $t$ . We show that  $Alg_i$  manages to always keep the executions in one of these 5 states, possibly coming back to total synchronization with  $Alg$ , or reaching the end of the sequence with  $Alg_i(\mathcal{R}) \leq Alg(\mathcal{R})$ .

Hassidim shows that there is an optimal solution for minimizing the makespan that on each fault it evicts the page that is furthest in the future for some core. In other words, if the sequence whose page should be evicted is known, the page to be evicted is the furthest in the future in that sequence. We show that the same result holds in our model for minimizing the number of faults.<sup>4</sup>

**THEOREM 6.** *There exists an optimal offline algorithm for FTF on disjoint sequences that upon each fault evicts a page  $\sigma \in R_j$  whose next request time is maximal in  $R_j$ , for some  $j$ .*

## 5.3 Optimal Algorithms for FTF and PIF

Theorem 6 implies an  $O(p^n)$  time optimal algorithm for FTF that upon each fault chooses the sequence to evict from optimally by trying all possibilities. Using dynamic programming, however, we can obtain a faster algorithm that is exponential in the number of sequences, but polynomial in the length of the sequences (recall we assume  $n \gg p$ ). In particular, it has been observed that  $p$  can be effectively assumed to be  $O(\log n)$  [19]). This algorithm can be extended to solve PIF as well. Next we describe these algorithms, showing that if the number of sequences is constant, then both FTF and PIF admit polynomial time algorithms<sup>5</sup>.

### 5.3.1 Minimizing the number of faults

Let  $\vec{x} = (x_1, \dots, x_p)$ , where each  $x_i$ ,  $1 \leq x_i \leq n_i(\tau + 1) + 1$ , is an index of a place in sequence  $i$  when serving it, either

<sup>5</sup>Note that this does not contradict Theorem 3: we can think of PIF as having two input sizes, the length of the sequences  $n$ , and the number of sequences  $p$ , which is the number of cores. PIF is NP-hard when the number of sequences is unbounded, but can be solved in polynomial time in  $n$ .

at a page, or at a time when a page is being fetched. If  $x_i$  is of the form  $x_i = (j-1)(\tau+1) + 1$ , then  $x_i$  is the index of the  $j$ -th page in  $R_i$ . In this case, we say  $x_i$  points to a page (denoted as  $R_i(x_i)$ ). Otherwise,  $x_i$  points to the fetching period of page  $\lceil x_i/(\tau+1) \rceil$ , and  $R_i(x_i)$  denotes the page being fetched. Let  $\mathcal{R}(\vec{x})$  denote the set of pages  $p$  indexed by  $\vec{x}$ , including pages in the fetching period.

Given a request  $\mathcal{R}$  we want to compute, for each possible cache configuration  $C$  and possible vector of positions  $\vec{x}$ , the minimum number of faults required to serve  $\mathcal{R}$  up to  $\vec{x}$ , and arriving at a cache configuration  $C$ . If  $x_i$  is in a page, this cost does not include serving  $R_i(x_i)$ . If  $x_i$  points to a fetching period, this cost includes the fault on  $R_i(x_i)$ . We compute and store these values in a  $(p+1)$ -dimensional table storing the minimum cost for each configuration and position. A cell  $F[C, x_1, \dots, x_p]$  can contribute to a cell  $F[C', x'_1, \dots, x'_p]$ , where  $C'$  is any configuration that can be obtained by removing  $|\mathcal{R}(\vec{x}) \setminus C|$  pages from  $C$  and adding the ones in  $\mathcal{R}(\vec{x})$ , and  $x'_i > x_i$  is the next index on sequence  $i$ . If  $x_i$  points to a page and  $R_i(x_i)$  is a hit, then  $x'_i = x_i + \tau + 1$ , i.e. the index jumps to the next page. If  $R_i(x_i)$  is being fetched or is a miss, then  $x'_i = x_i + 1$ , unless the page  $R_i(x_i)$  finished fetching in this timestep for another sequence, in which case the index also jumps to the next page and thus  $x'_i = \lceil x_i/(\tau+1) \rceil(\tau+1) + 1$  (this case only applies to non-disjoint sequences). We fill the table in a bottom up fashion, updating from a cell  $c$  only the cells that  $c$  can contribute to with a lower number of faults. The total minimum number of faults is then the minimum among all cache configurations  $C$  of  $F[C, n_1(\tau+1)+1, \dots, n_p(\tau+1)+1]$ . Algorithm 1 shows this procedure in pseudocode.

Let  $w$  be the total number of different pages requested in an instance. The number of possible cache configurations is  $\sum_{i=0}^K \binom{w}{i} \leq (w+1)^K$ . Hence, the total number of cells in table  $F$  is  $O((w+1)^K (n(\tau+1)+1)^p)$ . Counting the faults in  $\mathcal{R}(\vec{x})$  for each  $\vec{x}$  takes time  $O(p^2)$ . Since  $|\mathcal{R}(\vec{x})| \leq p$ , at most  $\binom{K}{p} = O(K^p)$  cache configurations can be reached from any configuration  $C$ , thus the time to process each cell and update the cells that it can contribute to is  $O(p^2 + K^p)$ . Since  $w \leq n$ , when  $K$  and  $p$  are constants, the total running time is  $O(n^{K+p}(\tau+1)^p)$ .

**THEOREM 7.** *Given a set  $\mathcal{R}$  of  $p$  sequences of total length  $n$ , a cache size  $K$ , and  $\tau \geq 0$ , with  $p = O(1)$  and  $K = O(1)$ , the minimum number of faults to serve  $\mathcal{R}$  can be determined in  $O(n^{K+p}(\tau+1)^p)$  time.*

While the running time of Algorithm 1 is not practical for realistic values of  $K$ , this result is important in that it shows that the complexity of multi-core paging stems from the number of sequences and not their length, placing the problem in P when  $p$  is constant.

### 5.3.2 Deciding PARTIAL-INDIVIDUAL-FAULTS

The algorithm for FTF can be extended to solve PIF as follows. For each cache configuration  $C$  and positions  $(x_1, \dots, x_p)$  we store a set of pairs  $(\vec{f}, t)$ , where  $\vec{f} = (f_1, f_2, \dots, f_p)$  specifies the faults on each sequence when reaching configuration  $(C, x_1, \dots, x_p)$  at time  $t$ . Thus, each pair  $(\vec{f}, t)$  associated with  $F[C, x_1, \dots, x_p]$  represents the number of faults in each sequence for a possible way of serving sequence  $\mathcal{R}$  up to time  $t$ . Algorithm 2 in Appendix A.3 shows the algorithm in pseudocode. The number of entries of the table  $F$

---

### Algorithm 1 Minimum Final Total Faults( $\mathcal{R}, K$ )

---

```

for all configurations  $C$  do
   $F[C, 1, \dots, 1] = 0$ 
  for each  $(x_1, \dots, x_p) \in \{2, \dots, n_i(\tau+1)+1\}^p$  do
     $F[C, x_1, \dots, x_p] = \infty$ 
  for each  $(x_1, \dots, x_p) \in \{1, \dots, n_i(\tau+1)\}^p$  do
    for all configurations  $C$  do
      if  $F[C, x_1, \dots, x_p] \neq \infty$  then
         $f = 0$  {faults in  $\mathcal{R}(\vec{x})$ }
        for  $i = 1$  to  $p$  do
          if  $x_i = (j-1)(\tau+1)+1$  for some  $j$ , and  $R_i(x_i) \in C$  then { $R_i(x_i)$  is a hit}
             $x'_i = x_i + \tau + 1$ 
          else { $R_i(x_i)$  is being fetched or it is a fault}
            if  $R_i(x_i) \notin C$  then { $R_i(x_i)$  is a fault}
               $f = f + 1$ 
            if  $R_i(x_i) = R_\ell(x_\ell)$  and  $x_\ell = (j'-1)(\tau+1)$  for some  $j'$  and some  $\ell \neq i$  then { $R_i(x_i)$  was fetched for another sequence}
               $x'_i = \lceil x_i/(\tau+1) \rceil(\tau+1) + 1$ 
            else
               $x'_i = x_i + 1$ 
          for all  $C'$  s.t.  $\mathcal{R}(\vec{x}) \subseteq C'$  and  $(C' \setminus \mathcal{R}(\vec{x})) \subseteq C$  do
            if  $F[C, x_1, \dots, x_p] + f < F[C', x'_1, \dots, x'_p]$  then
               $F[C', x'_1, \dots, x'_p] = F[C, x_1, \dots, x_p] + f$ 
        return  $\min_C \{F[C, n_1(\tau+1)+1, \dots, n_p(\tau+1)+1]\}$ 

```

---

in Algorithm 2 is  $O(n^{K+p}(\tau+1)^p)$  as in the algorithm for FTF. However, now each entry stores a list of pairs of fault vectors and time. Since at any time the number of faults in a sequence is at most  $n$ , the total number of different fault vectors is  $O((n+1)^p)$ . The time component of each pair can have at most  $n(\tau+1)$  values, and hence each entry can have at most  $O(n^{p+1}(\tau+1))$  pairs. For each entry in  $F$  we have to go through the list of pairs and compute the new vectors, and hence processing an entry takes  $O(pn^{p+1}(\tau+1))$  time. Once the new vectors are computed, these might have to be added to at most at most  $O(K^p)$  other entries. Hence the total time to process one entry is  $O(pn^{p+1}(\tau+1) + K^p)$ , and therefore the total time is  $O(n^{K+2p+1}(\tau+1)^{p+1})$ .

**THEOREM 8.** *Given a set  $\mathcal{R}$  of  $p$  sequences of total length  $n$ , a cache size  $K$ ,  $\tau \geq 0$ , a checkpoint time  $t$ , and a vector  $\vec{b} \in \mathbb{N}^p$ , with  $p = O(1)$  and  $K = O(1)$ , it can be decided if  $\mathcal{R}$  can be served s.t. at time  $t$  each sequence  $R_i$  has incurred in at most  $b_i$  faults in  $O(n^{K+2p+1}(\tau+1)^{p+1})$  time.*

## 6. CONCLUSIONS

We have proposed a model for multi-core paging that extends classical paging to a setting in which various sequences must be served simultaneously with a shared cache. The presence of multiple sequences and the fact that faults delay future requests make this problem significantly more difficult than classical paging. Neither traditional online algorithms nor the optimal strategy for sequential paging are competitive for this problem. Moreover, we have shown that serving a set of requests while limiting the number of faults in each sequence is NP-complete for an unbounded number of sequences. We showed, on the other hand, that multi-core paging admits algorithms that run in polynomial time in the length of the sequences, thus the problem is in P when the number of sequences is constant. Given the unfeasible



running time of the proposed algorithms it would be desirable to obtain more efficient exact or approximate offline algorithms.

Other directions of research include determining the complexity of FINAL-TOTAL-FAULTS and obtaining competitive online algorithms. Given the apparent excessive advantage of an offline algorithm over an online strategy that cannot do anything about future alignments, perhaps comparing online strategies to an optimal offline algorithm that can align sequences to its advantage might not lead to interesting online strategies. Hence the definition of a good evaluation framework for online strategies is open for debate, and perhaps other measures such as fairness or relative progress of sequences should be considered over minimizing faults globally.

## 7. ACKNOWLEDGMENTS

The authors would like to thank Reza Dorrigiv, Robert Fraser, Patrick Nicholson, and Francisco Claude for helpful discussions on this topic.

## 8. REFERENCES

- [1] U. A. Acar, G. E. Blelloch, and R. D. Blumofe. The data locality of work stealing. *Theory Comput. Syst.*, 35(3):321–347, 2002.
- [2] D. Ajwani, N. Sitchinava, and N. Zeh. Geometric algorithms for private-cache chip multiprocessors. In M. de Berg and U. Meyer, editors, *Proceedings of the 18th Annual European Symposium (ESA 2010)*, volume 6347 of *LNCS*, pages 75–86. Springer, 2010.
- [3] C. Ambühl. Offline list update is NP-hard. In *Proceedings of the 8th Annual European Symposium (ESA 2000)*, volume 1879 of *LNCS*, pages 42–51. Springer, 2000.
- [4] L. Arge, M. T. Goodrich, M. J. Nelson, and N. Sitchinava. Fundamental parallel algorithms for private-cache chip multiprocessors. In *SPAA 2008: Proceedings of the 20th Annual ACM Symposium on Parallelism in Algorithms and Architectures*, pages 197–206. ACM, 2008.
- [5] L. Arge, M. T. Goodrich, and N. Sitchinava. Parallel external memory graph algorithms. In *Proceedings of the 24th IEEE International Symposium on Parallel and Distributed Processing*, pages 1–11. IEEE, 2010.
- [6] R. D. Barve, E. F. Grove, and J. S. Vitter. Application-controlled paging for a shared cache. *SIAM J. Comput.*, 29:1290–1303, February 2000.
- [7] L. A. Belady. A study of replacement algorithms for virtual-storage computer. *IBM Systems Journal*, 5(2):78–101, 1966.
- [8] G. E. Blelloch, R. A. Chowdhury, P. B. Gibbons, V. Ramachandran, S. Chen, and M. Kozuch. Provably good multicore cache performance for divide-and-conquer algorithms. In *Proceedings of the 2008 ACM-SIAM Symposium on Discrete Algorithms*, January 2008.
- [9] G. E. Blelloch and P. B. Gibbons. Effectively sharing a cache among threads. In *SPAA 2004: Proceedings of the 16th Annual ACM Symposium on Parallelism in Algorithms and Architectures*, pages 235–244, New York, NY, USA, 2004. ACM.
- [10] G. E. Blelloch, P. B. Gibbons, and H. V. Simhadri. Low depth cache-oblivious algorithms. In *SPAA 2010: Proceedings of the 22nd Annual ACM Symposium on Parallelism in Algorithms and Architectures*, pages 189–199. ACM, 2010.
- [11] R. D. Blumofe, M. Frigo, C. F. Joerg, C. E. Leiserson, and K. H. Randall. An analysis of dag-consistent distributed shared-memory algorithms. In *SPAA 1996: Proceedings of the 8th Annual ACM Symposium on Parallel Algorithms and Architectures*, pages 297–308, 1996.
- [12] A. Borodin and R. El-Yaniv. *Online computation and competitive analysis*. Cambridge University Press, New York, NY, USA, 1998.
- [13] A. Borodin, S. Irani, P. Raghavan, and B. Schieber. Competitive paging with locality of reference. *J. Comput. Syst. Sci.*, 50:244–258, April 1995.
- [14] P. Cao, E. W. Felten, and K. Li. Application-controlled file caching policies. In *Proceedings of the USENIX Summer 1994 Technical Conference on USENIX Summer 1994 Technical Conference - Volume 1, USTC'94*, pages 11–11, Berkeley, CA, USA, 1994. USENIX Association.
- [15] J. Chang and G. S. Sohi. Cooperative cache partitioning for chip multiprocessors. In *ICS '07: Proceedings of the 21st annual international conference on Supercomputing*, pages 242–252, New York, NY, USA, 2007. ACM.
- [16] R. A. Chowdhury and V. Ramachandran. Cache-efficient dynamic programming algorithms for multicores. In *SPAA 2008: Proceedings of the 20th Annual ACM Symposium on Parallelism in Algorithms and Architectures*, pages 207–216. ACM, 2008.
- [17] R. A. Chowdhury, F. Silvestri, B. Blakeley, and V. Ramachandran. Oblivious algorithms for multicores and network of processors. In *Proceedings of the 24th IEEE International Symposium on Parallel and Distributed Processing*, pages 1–12. IEEE, 2010.
- [18] M. Cieliebak, S. Eidenbenz, and G. Woeginger. Double digest revisited: Complexity and approximability in the presence of noisy data. In T. Warnow and B. Zhu, editors, *Computing and Combinatorics*, volume 2697 of *Lecture Notes in Computer Science*, pages 519–527. Springer Berlin / Heidelberg, 2003. 10.1007/3-540-45071-8\_52.
- [19] R. Dorrigiv, A. López-Ortiz, and A. Salinger. Optimal speedup on a low-degree multi-core parallel architecture (LoPRAM). In *SPAA 2008: Proceedings of the 20th Annual ACM Symposium on Parallelism in Algorithms and Architectures*, pages 185–187. ACM, 2008.
- [20] R. Fedorova, M. Seltzer, and M. D. Smith. Cache-fair thread scheduling for multicore processors. Technical Report TR-17-06, Harvard University, 2006.
- [21] E. Feuerstein and A. Strejilevich de Loma. On-line multi-threaded paging. *Algorithmica*, 32(1):36–60, 2002.
- [22] A. Fiat and A. R. Karlin. Randomized and multipointer paging with locality of reference. In *Proceedings of the 27th annual ACM symposium on Theory of computing, STOC '95*, pages 626–634, New York, NY, USA, 1995. ACM.

- [23] M. R. Garey and D. S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman, 1979.
- [24] T. Hagerup. Online and offline access to short lists. In *Proceedings of the 32nd International Symposium on Mathematical Foundations of Computer Science*, volume 4708 of *Lecture Notes in Computer Science*, pages 691–702. Springer, 2007.
- [25] A. Hassidim. Cache replacement policies for multicore processors. In *Proceedings of 1st Symposium on Innovations in Computer Science*. Tsinghua University Press, 2010.
- [26] A. Jaleel, W. Hasenplaugh, M. Qureshi, J. Sebot, S. Steely, Jr., and J. Emer. Adaptive insertion policies for managing shared caches. In *PACT '08: Proceedings of the 17th international conference on Parallel architectures and compilation techniques*, pages 208–219, New York, NY, USA, 2008. ACM.
- [27] A. R. Karlin, M. S. Manasse, L. Rudolph, and D. D. Sleator. Competitive snoopy caching. *Algorithmica*, 3:77–119, 1988.
- [28] A. López-Ortiz and A. Salinger. Brief announcement: paging for multicore processors. In *SPAA 2011: Proceedings of the 23rd Annual ACM Symposium on Parallelism in Algorithms and Architectures*, pages 137–138. ACM, 2011.
- [29] A. M. Molnos, S. D. Cotofana, M. J. M. Heijligers, and J. T. J. van Eijndhoven. Throughput optimization via cache partitioning for embedded multiprocessors. In G. Gaydadjiev, C. J. Glossner, J. Takala, and S. Vassiliadis, editors, *ICSAMOS*, pages 185–192. IEEE, 2006.
- [30] J. I. Munro. On the competitiveness of linear search. In *Proceedings of the 8th Annual European Symposium on Algorithms*, ESA '00, pages 338–345, London, UK, 2000. Springer-Verlag.
- [31] M. K. Qureshi, A. Jaleel, Y. N. Patt, S. C. Steely, and J. Emer. Adaptive insertion policies for high performance caching. In *ISCA '07: Proceedings of the 34th annual international symposium on Computer architecture*, pages 381–391, New York, NY, USA, 2007. ACM.
- [32] N. Reingold and J. Westbrook. Off-line algorithms for the list update problem. *Inf. Process. Lett.*, 60(2):75–80, 1996.
- [33] S. S. Seiden. Randomized online multi-threaded paging. *Nord. J. Comput.*, 6(2):148–161, 1999.
- [34] D. D. Sleator and R. E. Tarjan. Amortized efficiency of list update and paging rules. *Commun. ACM*, 28(2):202–208, 1985.
- [35] S. Srikantaiah, M. Kandemir, and Q. Wang. Sharp control: controlled shared cache management in chip multiprocessors. In *MICRO 42: Proceedings of the 42nd Annual IEEE/ACM International Symposium on Microarchitecture*, pages 517–528, New York, NY, USA, 2009. ACM.
- [36] H. Stone, J. Turek, and J. Wolf. Optimal partitioning of cache memory. *IEEE Transactions on Computers*, 41:1054–1068, 1992.
- [37] A. Streljevič de Loma. New results on fair multi-threaded paging. *Electronic Journal of SADIO*, 1(1):21–36, 1998.
- [38] G. E. Suh, L. Rudolph, and S. Devadas. Dynamic cache partitioning for simultaneous multithreading systems. In *Proceedings of the IASTED International Conference on Parallel and Distributed Computing and Systems*, pages 116–127, 2001.
- [39] E. Torng. A unified analysis of paging and caching. *Algorithmica*, 20(2):175–200, 1998.
- [40] L. G. Valiant. A bridging model for parallel computation. *Commun. ACM*, 33(8):103–111, 1990.
- [41] Y. Xie and G. H. Loh. Pipp: promotion/insertion pseudo-partitioning of multi-core shared caches. *SIGARCH Comput. Archit. News*, 37(3):174–183, 2009.

## APPENDIX

### A.

#### A.1 Proofs for Section 4

**LEMMA 1. Online vs. offline eviction policies with a fixed static partition.** *Let  $A$  be any deterministic online eviction algorithm and let  $B = \{k_1, k_2, \dots, k_p\}$  be any online static partition. There exists a sequence  $\mathcal{R}$  such that  $P_A^B(\mathcal{R})/P_{OPT}^B(\mathcal{R}) = \Omega(\max_j \{k_j\})$ . When  $A$  is any marking or conservative algorithm (e.g. LRU), there is a matching upper bound, i.e.,  $\forall \mathcal{R}, P_A^B(\mathcal{R})/P_{OPT}^B(\mathcal{R}) \leq \max_j \{k_j\}$ .*

**PROOF. Lower bound.** Let  $j^* = \operatorname{argmax}_j \{k_j\}$ . The sequence  $\mathcal{R}$  is such that for  $j \neq j^*$ ,  $R_j = (\sigma_1^j)^{n/p}$ , i.e., the same page is requested  $n/p$  times, while  $R_{j^*}$  consists of requesting, among pages  $\{\sigma_1, \sigma_2, \dots, \sigma_{k_{j^*}+1}\}$ , the page just evicted by  $A$ , where  $\sigma_{i_1} \neq \sigma_{i_2}$  for  $i_1 \neq i_2$ , and all sequences are disjoint.  $P_A^B(\mathcal{R}) = n/p + p - 1$ , since it faults on every request of  $R_{j^*}$  and once on each other sequence. On the other hand, since  $P_{OPT}^B$  only evicts a page of sequence  $R_{j^*}$  if it is not requested in the following  $k_{j^*}$  requests, we have  $P_{OPT}^B(\mathcal{R}) \leq (n/p)/k_{j^*} + p - 1$  and the lemma follows.

**Upper bound.** Divide sequence  $R_j$  in phases such that a new phase starts every time there is a request for the  $(k_j + 1)$ -th distinct page since the beginning of the previous phase, and the first phase begins at the first page of  $R_j$ .  $P_{LRU}^B$  faults at most  $k_j$  times in each phase of  $R_j$ , while any algorithm must fault at least once in each phase. Let  $\phi_j$  denote the number of phases of sequence  $R_j$ , then  $P_{LRU}^B(\mathcal{R}) \leq \sum_{j=1}^p \phi_j k_j \leq \max_j \{k_j\} \sum_{j=1}^p \phi_j$ . On the other hand,  $P_{OPT}^B(\mathcal{R}) \geq \sum_{j=1}^p \phi_j$ , and thus  $P_{LRU}^B(\mathcal{R})/P_{OPT}^B(\mathcal{R}) \leq \max_j \{k_j\}$ .  $\square$

**LEMMA 2. Online static partition strategies are not competitive.** *Let  $B = \{k_1, \dots, k_p\}$  be any online static partition. Let  $k^* = \min_j \{k_j | k_j \geq 2\}$ . Then there exists a sequence  $\mathcal{R}$  s.t. for all eviction policies  $A$ ,  $P_A^B(\mathcal{R})/P_{LRU}^{OPT}(\mathcal{R}) \geq \min\{k^*, p-1\} \frac{n}{K^2 p} = \Omega(n)$ .*

**PROOF.** Consider first  $A=LRU$ . Let  $j^* = \operatorname{argmin}_j \{k_j | k_j \geq 2\}$  (i.e.  $k_{j^*} = k^*$ ). Let  $P$  denote the set of the first  $k_{j^*}$  processors in decreasing order by size of part of the cache according to  $B$ . Note that if  $k_{j^*} \geq (p-1)$  then  $P$  is equal to the set of all processors. Let  $P' = P \setminus \{j^*\}$ . Let  $R_j = (\sigma_1^j \sigma_2^j \dots \sigma_{k_j+1}^j)^{x_j}$  with  $x_j$  such that  $x_j(k_j+1) = n/p$  for all  $j \in P'$ , and let  $R_{j^*} = (\sigma_1^{j^*} \sigma_2^{j^*} \dots \sigma_{k_{j^*}}^{j^*})^{x_{j^*}}$   $j \notin P'$  and  $j \neq j^*$ , where  $x_{j^*}$  is such that  $x_{j^*} k_{j^*} = n/p$ . Let  $R_{j^*} = (\sigma_1^{j^*})^{n/p}$ .

$P_{LRU}^B$  faults on every request of  $|P'|$  processors and faults only on the first request of processor  $j^*$ . Hence,  $P_{LRU}^B(\mathcal{R}) \geq \min\{k_{j^*}, (p-1)\}n/p$ .

On the other hand, an optimal partition for  $\mathcal{R}$  would be one such that all different pages of each request  $R_j$  fit in the cache. Intuitively, an optimal partition takes units of cache from  $j^*$  and assigns them to other processors. Let  $k_j^{OPT}$  denote the size of the cache for processor  $j$  according to the optimal partition, then  $k_j^{OPT} = k_j + 1$  if  $j \in P'$  and  $k_j^{OPT} = \min\{1, k_j - (p-1)\}$  for  $j = j^*$ . The number of faults of  $P_{LRU}^{OPT}$  on  $\mathcal{R}$  is  $K$ , since it only faults on the first request to each different page. Hence  $P_{LRU}^B(\mathcal{R})/P_{LRU}^{OPT}(\mathcal{R}) \geq \min\{k_{j^*}, (p-1)\}n/Kp = \Omega(n)$ . Now, by Lemma 1  $P_{OPT}^B(\mathcal{R}) \geq \frac{P_{LRU}^B(\mathcal{R})}{K}$ , and the lemma follows.  $\square$

**THEOREM 1.** *Let  $A$  be any deterministic online cache eviction policy, let  $sOPT$  be an optimal static partition, and let  $D$  be any online dynamic partition strategy that changes the sizes of the parts  $o(n)$  times. The following statements hold:*

1. *There exists a sequence  $\mathcal{R}$  s.t.  $\frac{P_{LRU}^{sOPT}(\mathcal{R})}{S_{LRU}(\mathcal{R})} = \Omega(n)$ .*
2. *For all sequences  $\mathcal{R}$ ,  $S_{LRU}(\mathcal{R})/P_{OPT}^{sOPT}(\mathcal{R}) \leq K$ .*
3. *There exists a sequence  $\mathcal{R}$  s.t.  $P_A^D(\mathcal{R})/S_{LRU}(\mathcal{R}) = \omega(1)$ . Furthermore, if  $D$  varies the partition a constant number of times,  $P_A^D(\mathcal{R})/S_{LRU}(\mathcal{R}) = \Omega(n)$ .*

**PROOF.** 1. Let  $K_1 = 0$  and  $K_j = \sum_{i=1}^{j-1} k_i + 1$ , for all  $2 \leq j \leq p$ . Consider a sequence of requests  $\mathcal{R}$ , in which processor  $j$  requests the following pages, for all  $j$  simultaneously:

$$(\sigma_1^j)^{\alpha_j} (\sigma_1^j \sigma_2^j \dots \sigma_{K/p+1}^j)^x (\sigma_1^j)^{\beta_j}$$

where  $\sigma_{i_1}^j \neq \sigma_{i_2}^j$  for all  $i_1 \neq i_2$ ,  $\alpha_j = (j-1)(K/p+1)(\tau+x)$ ,  $\beta_j = (K+p-j(K/p+1))(\tau+x)$ , and  $x$  is a parameter. In other words, processor  $j$  requests the same page for a while, then repeatedly requests  $K/p+1$  distinct pages (call this the *distinct period*), and then goes back to requesting the same page again. All processors do the same, taking turns to be the processor currently in the distinct period: when one processor is in the distinct period, all other processors request repeatedly the same page. Given the request sequence, an optimal partition assigns  $K/p+1$  units of cache to  $p-1$  processors, and the rest to one processor: assigning more than  $K/p+1$  units of cache to any processor does not result in fewer faults, and assigning less than  $K/p+1$  to more than one processor increases the number of faults. Let  $j^*$  be the processor whose partition is  $k_{j^*} = K/p - (p-1)$ . Consider the distinct period of this processor. Let  $A$  be any eviction policy. No matter what the eviction policy  $A$  is, even the optimal offline,  $P_A^{sOPT}$  will fault at least once every  $k_{j^*}$  requests. Hence  $P_A^{sOPT}(\mathcal{R}) \geq x(K/p+1)/k_{j^*}$ . On the other hand,  $S_{LRU}(\mathcal{R})$  faults only on the first  $K/p+1$  requests of the distinct period of each processor, for a total of  $K+p$  faults. Hence  $P_A^{sOPT}(\mathcal{R})/S_{LRU}(\mathcal{R}) \geq x/(pk_{j^*})$ .  $x$  can be made arbitrarily large, in fact  $n = \tau(K+p)(p-1) + xp(K+p)$  and thus  $x = n/(p(K+p)) + \tau(p-1)/p$ , and thus  $x/(pk_{j^*}) = \Omega(n)$ .

2. Divide a sequence  $R_j$  of processor  $j$  in phases such that in a sequential traversal of pages, a new phase

begins either on the first page, or at the  $(k_j+1)$ -th different page since the beginning of the current phase, where  $k_j$  is the size of the cache assigned by  $sOPT$  to processor  $j$ . Define a shared phase equivalently for the cache size  $K$  and the sequence  $\mathcal{R}'$  containing the pages of  $\mathcal{R}$  in the order in which they are requested during the execution of  $S_{LRU}$  (with simultaneous requests sorted by increasing number of processor). We claim that a shared phase cannot start and end without at least one sequence changing phase. In other words, the phase of at least one sequence must end before the end of a shared phase. If this was not the case, within the shared phase, the number of different pages in the sequence of each processor  $j$  would be at most  $k_j$ , and therefore the total number of different pages in the shared phase would be at most  $K$ , which is a contradiction. Let  $\phi$  denote the number of shared phases of sequence  $\mathcal{R}$  and  $\phi_j$  denote the number of phases of sequence  $R_j$ . The above claim implies that  $\phi \leq \sum_{j=1}^p \phi_j$ . Assuming that  $S_{LRU}$  timestamps each page at the moment of request, it is not difficult to see that the fact that  $S_{LRU}$  faults at most  $K$  times per phase in the sequential setting extends to the above definition of shared phases. Since any cache eviction algorithm must fault at least once per phase, it follows that  $S_{LRU}(\mathcal{R}) \leq K\phi \leq K \sum_{j=1}^p \phi_j \leq KP_{OPT}^{sOPT}(\mathcal{R})$ .

3. Let a stage of  $D$  denote a period in which the sizes of the partition are constant. If the number of stages of  $D$  is  $o(n)$ , then at least one stage has non-constant length  $\ell = \omega(1)$  (in number of parallel page requests). We then apply the same argument as in the proof of statement 1. Let  $\mathcal{R}$  in this stage consist of a sequence in the form of the sequence in that proof: each processor's sequence has three periods: (1) only one page  $\sigma_1^j$  is requested repeatedly, (2) the page requested is any page not in the cache of processor  $j$  (the *distinct period*), and (3) again only one page  $\sigma_1^j$  is requested. The length of period (2) is  $m$  pages, and each processor takes turns to be in the distinct period. Hence the total number of requests in the stage is  $mp^2 = \ell p$ . Let  $t$  be the time where the long stage begins. During the distinct period of processor  $j$ ,  $R_j$  consists of repeatedly requesting the page not in  $j$ 's cache, among the pages  $\{\sigma_1^j, \dots, \sigma_{k(j,t)+1}^j\}$ , where  $k(j,t)$  is the size of the cache of processor  $j$  at time  $t$ .  $P_A^D$  faults on every request of the distinct period of all processors, and hence in this stage  $P_A^D(\mathcal{R}) = pm = \ell$ . On the other hand, in this stage,  $S_{LRU}$  faults only on the first request to a distinct page in the distinct period of each processor, and thus in this stage  $S_{LRU}(\mathcal{R}) = K+p$  (recall we assume  $k_j \geq 1$  for all  $1 \leq j \leq p$  at all times). Let the rest of  $\mathcal{R}$  be such that neither algorithm faults. Then  $P_A^D(\mathcal{R})/S_{LRU}(\mathcal{R}) \geq \ell/(K+p) = \omega(1)$ . Note that if partitions are allowed only a constant number of stages, then  $P_A^D(\mathcal{R})/S_{LRU}(\mathcal{R}) = \Omega(n)$ .

$\square$

## A.2 Proofs for Section 5

**THEOREM 3.** *PARTIAL-INDIVIDUAL-FAULTS is NP-complete.*

**PROOF.** It is easy to see that the problem is in NP: given an instance  $\mathcal{I} = \{\mathcal{R}, K, t, \tau, \vec{b}\}$  with a "yes" answer, and a

certificate consisting of the pages to be evicted after each fault, it can be verified in time  $O(tp)$  that the number of faults in each sequence  $R_i$  is at most  $b_i$ . Note that  $t \leq \max_i \{|R_i|\}(\tau + 1)$ , and  $\tau = O(\max_i \{|R_i|\})$ , therefore the verifier runs in time polynomial in the size of the input.

In order to show that the problem is NP-complete, we build a reduction from 3-PARTITION. Recall that an instance of 3-PARTITION consists of a set of  $n$  integers  $S = \{s_1, \dots, s_n\}$ , and a bound  $B$ , such that  $B/4 < s_i < B/2$  for all  $1 \leq i \leq n$ . The problem is to determine if  $S$  can be partitioned into  $n/3$  sets  $A_1, \dots, A_{n/3}$  such that for all  $1 \leq j \leq n/3$ ,  $\sum_{i \in A_j} s_i = B$ . Note that the restrictions of the problem imply that each subset  $A_i$  must have exactly 3 elements [23].

Let  $\mathcal{J}$  be an instance of 3-PARTITION. We build an instance  $\mathcal{I}$  of PARTIAL-INDIVIDUAL-FAULTS as follows. There are  $p = |S|$  sequences. Each sequence  $R_i$  consists of alternating requests to 2 pages  $\alpha^i$  and  $\beta^i$ , where  $\alpha^i \neq \beta^i$ , and  $\alpha^i \neq \alpha^j$  and  $\beta^i \neq \beta^j$  for all  $i \neq j$ , and  $\alpha^i \neq \beta^j$  for all  $i, j$ . In other words,  $R_i = \alpha^i \beta^i \alpha^i \beta^i \dots$ , and all sequences are disjoint. The length of  $R_i$  is  $|R_i| = B(\tau + 1) + 4\tau + 5$ , where  $\tau \geq 0$  is any integer. The size of the cache is  $K = (4/3)p$ , and we want to know if the number of faults in each sequence  $R_i$  is at most  $b_i = B - s_i + 4$ , at time  $t = B(\tau + 1) + 4\tau + 5$ . Note that since 3-PARTITION is strongly NP-complete (it remains NP-complete if the input is encoded in unary), the reduction can be done from an instance encoded in unary, and hence it can be done in time polynomial in the size of  $\mathcal{J}$ .

We show now that there exists a solution for  $\mathcal{J}$  if and only if we can serve each  $R_i$  with at most  $b_i$  faults.

( $\Rightarrow$ ) We show first that if  $\mathcal{J}$  admits a solution then we can serve each  $R_i$  with at most  $b_i$  faults. Let  $A_1, \dots, A_{n/3}$  be the partition for  $\mathcal{J}$ . Divide the sequences in groups according to the partition, so that the sequences corresponding to  $A_j$  will share a group of 4 cells of the cache. Let  $R_{i_1}, R_{i_2}$ , and  $R_{i_3}$  be the sequences in group  $j$ . Each of these sequences will be assigned one cell for some time and two at other times. In other words, the three sequences will have one dedicated cell at least until time  $t$  and will share the extra cell of the group. Sequence  $R_i$  will use the extra cell continuously for enough time so it incurs in exactly  $h_i = s_i(\tau + 1) + 1$  hits (see Figure 2).

Say  $R_{i_1}, R_{i_2}$ , and  $R_{i_3}$  use the extra cell in that order. The first request to each sequence results in a fault and it is fetched to the dedicated cell of the corresponding sequence. The second request of  $R_{i_1}$  (also a fault) is fetched to the extra cell. Now both pages of  $R_{i_1}$  are in the cache, and they are kept there for the next  $h_{i_1}$  requests of  $R_{i_1}$ . Meanwhile, every request of  $R_{i_2}$  and  $R_{i_3}$  results in a fault and the eviction of the page in their corresponding dedicated cell. The last hit of  $R_{i_1}$  occurs at time  $(2 + s_{i_1})(\tau + 1) + 1$ , which coincides with a new request  $\sigma$  for  $R_{i_2}$ , since all pages have been faults for  $R_{i_2}$ . Instead of fetching  $\sigma$  to this sequence's dedicated cell,  $\sigma$  is fetched into the extra cell or  $R_{i_1}$ 's dedicated cell, depending on which page can be evicted at the time (if  $\sigma$  is fetched into  $R_{i_1}$ 's dedicated cell then this cell becomes the shared cell and the former shared cell becomes  $R_{i_1}$ 's dedicated cell). Now  $R_{i_2}$  has the extra cell and the remaining requests of  $R_{i_1}$  will result in faults. After  $h_{i_2}$  hits of  $R_{i_2}$ , the extra cell is now passed to  $R_{i_3}$  (again the last hit of  $R_{i_2}$  coincides with a request of  $R_{i_3}$ ), and this sequence keeps this cell until it completes  $h_{i_3}$  hits. At this point, the

time elapsed is the sum of the hits of each sequence, plus  $2\tau$  for the transitions of the extra cell from  $R_{i_1}$  to  $R_{i_2}$  and from  $R_{i_2}$  to  $R_{i_3}$ , plus the initial  $2(\tau + 1)$  time corresponding to the first 2 faults of the three sequences. Hence the time is  $t = h_{i_1} + h_{i_2} + h_{i_3} + 4\tau + 2 = (s_{i_1} + s_{i_2} + s_{i_3})(\tau + 1) + 4\tau + 5 = B(\tau + 1) + 4\tau + 5$ . The same strategy is used for each group in the partition, and the number of faults of sequence  $R_i$  is exactly  $(t - h_i)/(\tau + 1) = B - s_i + 4$ . Thus, if there is a solution to  $\mathcal{J}$ ,  $\mathcal{R}$  can be served such that at time  $t = B(\tau + 1) + 4\tau + 5$ , each sequence has incurred in  $b_i = B - s_i + 4$  faults.

( $\Leftarrow$ ) We show now that a solution to the instance  $\mathcal{I}$  of PIF gives a solution to the instance  $\mathcal{J}$  of 3-PARTITION. If  $\mathcal{R}$  can be served so that each sequence faults at most  $B - s_i + 4$  times by time  $t$ , then at least  $h_i = s_i(\tau + 1) + 1$  of  $R_i$ 's requests must be hits. Note that for all  $i$ ,  $|R_i| = t$ , and hence each sequence uses at least one cell until time at least  $t$ . A request can only be a hit if a sequence has two cells of the cache for consecutive timesteps. Since there are only  $(4/3)p$  cells and each sequence uses at least one cell, there are only  $p/3$  extra cells which can be used to store the second page of a sequence, and hence there can be at most  $p/3$  hits in one timestep.

In addition, any change in the partition that removes one cell from a sequence that had 2 cells and gives one more cell to another sequence implies at least  $\tau$  time without hits for the sequences involved. To see this, let  $R_i$  and  $R_j$  be two sequences such that  $k(i, t_1) = 2$  and  $k(j, t_1) = 1$ , respectively and  $k(i, t_1 + 1) = 1$  and  $k(j, t_1 + 1) = 2$  (where  $k(i, t_1)$  is the size of the cache of processor  $i$  at time  $t_1$ ), for some  $t_1 < t$ . Say page  $\alpha^i$  was a hit in  $R_i$ . Then at  $t_1 + 1$ ,  $\beta^i$  and  $\alpha^j$  are requested in  $R_i$  and  $R_j$ , respectively<sup>7</sup>.  $\alpha^i$  is evicted from the cache, and  $\alpha^j$  is fetched in to the cell previously used by  $\alpha^i$ . Sequence  $R_j$  must wait  $\tau$  more timesteps before having its first hit, while sequence  $R_i$ 's next request ( $\alpha^i$ ) will result in a fault. Hence, there are no hits in these two sequences in the period  $[t_1 + 2, t_1 + 1 + \tau]$ , i.e.  $\tau$  timesteps without hits.

Let  $C$  denote the number of all such changes in partition between a pair of sequences.  $C$  does not count the initial assignment of 2 cells to a sequence when starting to serve  $\mathcal{R}$ , but only when a sequence acquires an extra cell that held a page of another sequence. Note that no hits can happen until time  $2(\tau + 1) + 1$ , and hence the total number of possible hits before time  $t$  can be at most  $H_1 = (p/3)(t - 2(\tau + 1)) - C\tau = (p/3)((B - 2)(\tau + 1) + 4\tau + 5) - C\tau$ . On the other hand, the minimum number of required hits is  $H_2 = \sum_{i=1}^p h_i = \sum_{i=1}^p s_i(\tau + 1) + 1 = (p/3)B(\tau + 1) + p$ . We require  $H_1 \geq H_2$ , which implies that  $C \leq 2p/3$ , i.e. we can have at most  $2p/3$  changes in partitions. Note that initially at most  $p/3$  sequences can be assigned two cells, therefore every sequence must have 2 cells during a continuous period of time. We call this period the hit period of a sequence.

Consider a group of sequences whose hit periods are consecutive, i.e. a group  $I = \{i_1, i_2, \dots, i_\ell\}$  of the sequences such that a request in sequence  $i_{j+1}$  evicts a page from sequence  $i_j$  to start its hit period. These sequences can be served using a minimum of  $\ell + 1$  cells: one cell is dedicated for each sequence, while the other extra cell is used to assign 2 cells to some sequence during its hit period (this extra cell need not to be the same one during the entire execution). We claim that  $\ell \leq 3$ . To see this, assume  $\ell > 3$ , then since hit

<sup>7</sup>Note that  $R_j$  might be fetching a page instead, but the case when a new request comes is the one that minimizes the time that elapses from  $t_1$  until  $R_j$ 's next hit.

periods have no interruptions, the total time to serve these sequences is at least the total number of hits, plus the time for each change in partition, plus 2 initial faults, for a total of  $T = (\sum_{i \in I} h_i) + (\ell - 1)\tau + 2(\tau + 1) = (\sum_{i \in I} s_i(\tau + 1) + 1) + (\ell + 1)\tau + 2 > (\ell/4)B(\tau + 1) + \ell + (\ell + 1)\tau + 2$ , since  $s_i > B/4$  for all  $1 \leq i \leq p$ . Taking  $\ell = 4$ ,  $T > B(\tau + 1) + 5\tau + 6$ , which is strictly greater than  $t$ , and thus one sequence will not have all its required hits before  $t$  and hence it will exceed the allowed number of faults. Hence  $\ell \leq 3$ . Furthermore we argue that  $\ell$  is exactly 3. Assume, otherwise, that sequences are served in groups of 1, 2, and 3 sequences. Let  $n_\ell$  be the number of groups of  $\ell$  sequences, for  $\ell \in \{1, 2, 3\}$ . In order to satisfy the minimum number of hits for each sequence, it must be the case that a group of  $\ell$  sequences must use at least  $\ell + 1$  cells. Hence, the following must be satisfied:  $n_1 + 2n_2 + 3n_3 = p$ , and  $2n_1 + 3n_2 + 4n_3 \leq (4/3)p$ , with  $n_1, n_2, n_3 \geq 0$ . It is not difficult to see that a feasible solution must have  $n_1 = n_2 = 0$ , and  $R$  must be served only with groups of 3 sequences.

Finally, it must be the case that every group of sequences  $I = i_1, i_2, i_3$  must satisfy  $s_{i_1} + s_{i_2} + s_{i_3} = B$ . Suppose that a group  $I_1$  is such that  $\sum_{i \in I_1} s_i < B$ . Then since  $B = (3/p) \sum_{i=1}^p s_i$  there must exist another group  $I_2$  such that  $\sum_{i \in I_2} s_i > B$ . Then, since the minimum number of hits per sequence is  $h_i = s_i(\tau + 1)$ , the total time to serve the group would be  $T = (\sum_{i \in I_2} h_i) + 2\tau + 2(\tau + 1) > B(\tau + 1) + 4\tau + 5 = t$ , and thus at least one sequence in the group would have to fault more than its maximum number of allowed faults by time  $t$ . Therefore, the only possible way to serve the requests satisfying the faults requirement for each sequence is to divide them in groups of 3 sequences  $I_1, \dots, I_{p/3}$  such that  $\sum_{i \in I_j} s_i = B$  for all  $1 \leq j \leq p/3$ , which is a solution for the instance of 3-PARTITION.  $\square$

**THEOREM 4. MAX-PARTIAL-INDIVIDUAL-FAULTS is APX-hard.**

**PROOF.** We describe a gap preserving reduction from MAX-4-PARTITION to MAX-PIF. 4-PARTITION is an analog of 3-PARTITION in which the goal is to partition a set  $S = \{s_1, \dots, s_n\}$  into subsets  $A_1, \dots, A_{n/4}$  such that for all  $1 \leq j \leq n/4$ ,  $\sum_{i \in A_j} s_i = B$ , where  $B = (4/n) \sum_{i=1}^n s_i$  [23]. Each element  $s_i$  satisfies  $B/5 < s_i < B/3$  and thus each subset must have 4 elements. 4-PARTITION is also NP-complete [23], and a reduction to PIF can be built by modifying the proof of Theorem 3 in a straightforward way: the cache size is now  $K = (5/4)p$ , the length of each sequence is  $B(\tau + 1) + 5\tau + 6$  and the goal is to serve the sequences such that at time  $t = B(\tau + 1) + 5\tau + 6$  sequence  $i$  has incurred in at most  $b_i = B - s_i + 5$ . It is not hard to see that the same arguments in the proof of Theorem 3 apply to argue that an instance of 4-PARTITION admits a solution if and only if the instance of PIF admits a solution.

The MAX-4-PARTITION problem (as defined in [18]) is: given a set  $S$  and  $B$  as in the 4-PARTITION problem, find a maximum number of disjoint subsets whose elements add up to  $B$ . This problem is APX-hard, i.e. it does not admit a PTAS (assuming  $P \neq NP$ ) [18]. Given an instance  $\mathcal{J}$  ( $\mathcal{I}$ ) of MAX-4-PARTITION (MAX-PIF), let  $OPT_{4PART}(\mathcal{J})$  ( $OPT_{PIF}(\mathcal{I})$ ) denote the value of the optimal solution to  $\mathcal{J}$  ( $\mathcal{I}$ ). Let  $n = |S|$  in  $\mathcal{J}$ . In order to show that MAX-PIF is APX-hard, we build a reduction to an instance  $\mathcal{I}$  of MAX-PIF and show:

1.  $OPT_{4PART}(\mathcal{J}) \geq n/4 \Rightarrow OPT_{PIF}(\mathcal{I}) \geq n$
2.  $OPT_{4PART}(\mathcal{J}) < (1 - \epsilon)n/4 \Rightarrow OPT_{PIF}(\mathcal{I}) < (1 - \epsilon/4)n$

The reduction from an instance  $\mathcal{J}$  of MAX-4-PARTITION to an instance  $\mathcal{I}$  of MAX-PIF is exactly the same as the reduction from 4-PARTITION described above. Since a solution to  $\mathcal{J}$  gives a solution to  $\mathcal{I}$ , if  $OPT_{4PART}(\mathcal{J}) \geq n/4$  (and thus equal to  $n/4$ ) then all sequences of  $\mathcal{I}$  can be served with a number of faults within the given bounds, proving statement (1).

For statement (2), note that in the reduction from 4-PARTITION to PIF (adapted from the proof of Theorem 3), the only way of serving all sequences within the fault bounds is by partitioning them in groups of 4 that shared 5 cells of cache. Furthermore, the 4 sequences in a group can be served within the faults bounds if and only if the corresponding elements in  $S$  add up to  $B$ . Otherwise, at least one of the sequences will have to incur in more faults than the allowed bound. Therefore,  $OPT_{PIF}(\mathcal{I}) \leq 4OPT_{4PART}(\mathcal{J}) + 3(n/4 - OPT_{4PART}(\mathcal{J})) = OPT_{4PART}(\mathcal{J}) + 3n/4$ . Since  $OPT_{4PART}(\mathcal{J}) < (1 - \epsilon)n/4$ , we have  $OPT_{PIF}(\mathcal{I}) < (1 - \epsilon)n/4 + 3n/4 = (1 - \epsilon/4)n$ . Thus the reduction is gap-preserving, proving the theorem.  $\square$

**THEOREM 5.** *Let Alg be an offline optimal algorithm that forces faults. There exists an offline algorithm Alg' that is honest such that for all disjoint requests  $\mathcal{R}$ ,  $Alg'(\mathcal{R}) = Alg(\mathcal{R})$ .*

**PROOF.** We follow an inductive argument similar to the proof of optimality of Furthest-In-The-Future in the sequential setting [12]. The proof relies on the following claim: let  $Alg$  be any paging algorithm and  $R$  be any request sequence. Then, for all timesteps  $i$  it is possible to construct an algorithm  $Alg_i$  such that (i) for all  $t = 1, \dots, i - 1$ , it behaves exactly like  $Alg$ , (ii) if  $Alg$  forces a fault on  $t = i$ ,  $Alg_i$  does not, and (iii)  $Alg_i(R) \leq Alg(R)$ .

If the claim is correct, this implies that it is possible to obtain an optimal algorithm that does not force faults: for a given sequence  $R$ , start from any optimal algorithm  $OPT$ , and apply the claim with  $i = 1$  to obtain  $OPT_1$ , then apply the claim with  $i = 2$  to  $OPT_1$  to obtain  $OPT_2$ , and so on and so forth.  $OPT_{t'}$  is an optimal algorithm that does not force faults, where  $t'$  is the maximum timestep of the execution of the algorithm on  $R$ .

Let us prove the claim. Both algorithms start with an empty cache and hence  $Alg_i$  can do exactly as  $Alg$  does up to step  $i - 1$ . If at timestep  $i$   $Alg$  does not force a fault, then  $Alg_i$  continues behaving like  $Alg$  until the end of the request, and the number of faults of both algorithms is the same. Now, assume that  $Alg_i$  forces a fault on  $t = i$  on a page  $p_1$  on sequence  $R_s$ . Let  $C_{Alg}(t)$  and  $C_{Alg_i}(t)$  denote the caches of  $Alg$  and  $Alg_i$  right before serving  $R(t)$ . Since both algorithms behaved exactly the same up to  $t = i - 1$ , we have  $C_{Alg}(i) = C_{Alg_i}(i)$ . Also, let  $A(R, t)$  denote the number of faults of algorithm  $A$  right before serving request  $R(t)$ .

We argue that from that point on  $Alg_i$  can be such that both algorithms will fault on exactly the same pages in the rest of the sequences (and hence keeping the same alignment with respect to both algorithms), and that their caches will differ by at most one page. Furthermore,  $Alg_i(R) \leq Alg(R)$  at all times. We show this by defining a set of states that

describe the differences between the algorithms sequences, caches, and number of faults for each subsequent request in  $R_s$  for  $Alg$ . We define the following states at time  $t$ :

- A. All sequences in both algorithms have the same alignment,  $C_{Alg}(t) = C_{Alg_i}(t)$ , and  $Alg(R, t) = Alg_i(R, t)$ .
- B. All sequences other than  $R_s$  have the same alignment in both algorithms,  $C_{Alg}(t) = C_{Alg_i}(t)$ , and  $Alg(R, t) = Alg_i(R, t) + 1$ .
- C. All sequences other than  $R_s$  have the same alignment in both algorithms,  $C_{Alg}(t) = C_{Alg_i}(t) \cup \{p\}$ , and  $Alg(R, t) = Alg_i(R, t)$ , where  $p$  is a page previously requested in  $R_s$ , and the remaining cell of  $Alg_i$  cache is fetching a page  $p' \neq p$ .
- D. All sequences other than  $R_s$  have the same alignment in both algorithms,  $C_{Alg}(t) = (C_{Alg_i}(t) \cup \{p\}) \setminus \{\alpha\}$ , and  $Alg(R, t) = Alg_i(R, t) + 1$ , where  $p$  is a page previously requested in  $R_s$ , and  $\alpha$  is a page from a sequence other than  $R_s$ .
- E. All sequences other than  $R_s$  have the same alignment in both algorithms,  $C_{Alg}(t) = C_{Alg_i}(t) \cup \{p\}$ , and  $Alg(R, t) = Alg_i(R, t)$ , where  $p$  is a page previously requested in  $R_s$ , and the remaining cell of  $Alg_i$  cache is fetching  $p$ .

Let  $p_2, p_3, \dots$  be the pages in  $R_s$  after  $p_1$ . We define the request period of each page  $p_j \in R_s$  as the timesteps that include its request and possible fetching for algorithm  $Alg$ , i.e. if  $p_j$  is request at time  $t_j$  by  $Alg$ , then its request period is  $[t_j, t_j + \tau]$  if  $p_j$  is a fault, and just  $t_j$  if it is a hit. We will now show that the above are all the possible states that describe the algorithms in each period. We will prove this by induction on the request periods of pages in  $R_s$ .

Before  $p_1$ , the algorithms are in state (A). Consider the request period for  $p_1$ . Recall that  $C_{Alg}(t_1) = C_{Alg_i}(t_1)$ .  $Alg$  forces a fault on  $p_1$  and hence the cell corresponding to  $p_1$  in the cache is being used to fetch this page until the end of the period. Upon any request  $\sigma$  during the period, if  $Alg$  evicts  $\alpha$ ,  $Alg_i$  evicts  $\alpha$  as well. In this period up to  $\tau$  pages  $p_2, \dots, p_{1+\tau}$  after  $p_1$  might be requested for  $Alg_i$  in  $R_s$ . If none of these are faults, then at the end of the period  $Alg(R, t_2) = Alg_i(R, t_2) + 1$  and the caches of both algorithms are equal, hence the algorithms are in state (B)<sup>8</sup>. On the other hand, if any of the pages  $p_2, \dots, p_{1+\tau}$  is a fault for  $Alg_i$ , then  $Alg_i$  evicts any of the previous hit pages in  $R_s$  (there is a least one,  $p_1$ ). Hence at the end of the period the number of faults of both algorithms is the same, and both caches have the same pages, but for the evicted page by  $Alg_i$ , thus arriving at state (C).

For the request period of a page  $p_j$ ,  $j > 1$ , let  $s_j$  be the corresponding state of  $Alg$  and  $Alg_i$ . Assume that  $s_j$  is one of the states in  $S = \{A, B, C, D, E\}$ . We show now that  $Alg_i$  is able to reach a state  $s_{j+1} \in S$  in the next period when  $s_j = s$ , for each  $s \in S$ .

$[s_j = A]$  Suppose we have reached state (A). Since  $Alg_i$  cannot force faults only on request  $i$ , but it can do so on

later requests, it just behaves exactly like  $Alg$  for the rest of the sequence, maintaining state (A).

$[s_j = B]$  We can only arrive to state (B) if  $Alg_i$  had no faults for requests in  $R_s$  in the previous period. Since  $p_j$  was requested for  $Alg_i$  in the previous period (sequences  $R_s$  in both algorithms are misaligned by  $\tau$  at most),  $p_j$  is a hit for  $Alg$ . If there is any fault on a request of another sequence,  $Alg_i$  evicts whatever  $Alg$  evicts. At time  $t_j$  the request for  $Alg_i$  is  $p_{j+\tau}$ . If this page is a hit, then we stay in state (B). If this page is a fault,  $Alg_i$  evicts some page  $p_{j'}$  with  $j' < j + \tau$ . At least one page of  $R_s$  is in  $Alg_i$ 's cache since they were all hits in the previous period (and we assume they are not evicted during this period by  $Alg$ ). Hence, at the end of the period  $Alg(R, t_{j+1}) = Alg_i(R, t_{j+1})$  and  $C_{Alg}(t_{j+1}) = C_{Alg_i}(t_{j+1}) + \{p\}$ , for some  $p \in R_s$ , arriving at state (C).

$[s_j = C]$  At the beginning of this period the next page in the request ordering of  $Alg$  is  $p_j$  and  $Alg_i$  is fetching some page  $p_{j'}$  that resulted in a fault in the previous period. This is the case if we arrive from states (A), (B), and we will see that it holds if we stay in (C), or arrive from (D) or (E) as well. For all faults in sequences other than  $R_s$ ,  $Alg_i$  evicts the same page that  $Alg$  evicts, unless  $Alg$  evicts  $p$ , in which case  $Alg_i$  evicts another page in  $R_s$ . If  $p_j$  is a hit for  $Alg$  then nothing changes and we stay in (C) as well. If  $p_j$  is a fault, then  $Alg$  evicts a page  $\alpha$ . Recall that  $C_{Alg}(t_{j+1}) = C_{Alg_i}(t_{j+1}) + \{p\}$ . Assume first that  $\alpha \neq p$ . Then, if all the subsequent requests of  $R_s$  in the request ordering of  $Alg_i$  are hits and  $\alpha$  is not requested during this period, then we have one more fault for  $Alg$  and at the end of the period  $Alg$ 's cache still has  $p$  but not  $\alpha$ , while  $Alg_i$  does not have  $p$  but has  $\alpha$ , hence reaching state (D). If  $\alpha$  is requested during the period, then  $Alg_i$  forces a fault on this page and hence we remain in state (C). Now, if one of the requests from  $R_s$  for  $Alg_i$  results in a fault, say  $p_{j'}$ , then, again if  $\alpha$  was not requested before  $p_{j'}$ , then  $Alg_i$  evicts  $\alpha$  for this request. If  $\alpha$  is requested before  $p_{j'}$ ,  $Alg_i$  forces a fault on  $\alpha$  and evicts for  $p_{j'}$  whatever  $Alg$  evicted for  $\alpha$  (or some page  $p' \in R_s$  if  $Alg$  evicts  $p$ ). At the end of the period the difference in number of faults remains the same. Now, if  $p_{j'} \neq p$ , then we stay in state (C). However, if  $p_{j'} = p$ , i.e.  $Alg_i$  faulted in the page that it did not have but that  $Alg$  had, then we reach state (E).

Now, we analyze the case  $\alpha = p$ . The next page request in  $R_s$  in the ordering of  $Alg_i$  is  $p_{j+1}$ . If this page is a fault, then  $Alg_i$  evicts another page  $p' \in R_s$ , for example  $p_j$ . In this case the number of faults increases by 1 for both algorithms and we remain in state (C). If on the other hand,  $p_{j+1}$  is in  $Alg_i$ 's cache (and hence in  $Alg$ 's cache),  $Alg_i$  forces a fault on this page, reaching state (E).

$[s_j = D]$  We could only reach this state if  $p_j$  is a hit. Let  $p_{j'}$  be the next page of  $R_s$  in the request ordering of  $Alg_i$ . Suppose  $\alpha$  is not requested in  $R(t_j)$ . If  $p_{j'}$  is a hit, then nothing changes, and we remain in state (D). If  $p_{j'}$  is a fault,  $Alg_i$  evicts  $\alpha$  and we reach state (C). If  $\alpha$  is requested, since  $\alpha \notin C_{Alg}(t_j)$ ,  $Alg$  evicts a page  $\alpha'$ .  $Alg_i$  forces a fault on  $\alpha$ . If  $p_{j'}$  is a hit, then we remain in state (D). If  $p_{j'}$  is a fault,  $Alg_i$  evicts  $\alpha'$  if  $\alpha' \neq p$ , or another page  $p' \in R_s$  if  $\alpha' = p$ , reaching state (C).

$[s_j = E]$  This state is reached when the number of faults of both algorithms is the same and  $Alg_i$  is fetching the page  $p$  that is missing with respect to  $Alg$ 's cache. Suppose first that  $p_j \neq p$  (i.e. this is not the timestep in which  $Alg_i$

<sup>8</sup>We assume here and for the rest of the analysis that  $Alg_i$  lets  $Alg$  run ahead at least until its next fault in  $R_s$ , so that if  $Alg$  evicts a page  $p_j$  after the time a request for this page was served by  $Alg_i$ ,  $Alg_i$  forces the fault on this page as well. In other words, if  $p_j$  is a hit for  $Alg_i$  at time  $t$ , it will be a hit for  $Alg$  at time  $t + \tau$ .

finishes fetching  $p^9$ )  $p_j$  is a hit for  $Alg$  (this is the case in the two cases that we can arrive to this state from (C), and the one from (E)). Upon any fault on another sequence  $Alg_i$  evicts whatever  $Alg$  evicts (with the assumption that  $Alg$  will not evict a page of  $R_s$  that was a hit in the previous period for  $Alg_i$ ). If one of these evictions is for page  $p$ , then  $Alg_i$  evicts another page  $p' \in R_s$  and we reach state (C). If  $p$  was not evicted, then we remain in state (E). Now, if  $p_j = p$ , this page is in  $Alg$ 's cache. Again,  $Alg_i$  evicts what  $Alg$  evicts for other requests. If, however, in any of these evictions the page evicted is  $p^{10}$ ,  $Alg_i$  evicts another page  $p' \in R_s$ . In this case  $Alg$  faults in  $p_j$  and we are back in state (C). On the other hand, if  $p_j = p$  is a hit,  $Alg_i$  finishes at the same time to fetch  $p$ , therefore both caches are equal, sequences  $R_s$  in both algorithms are aligned, and the number of faults of both algorithms is the same, reaching state (A).

Hence,  $Alg_i$  can keep the execution in these states until the end. Since  $Alg_i(R) \leq Alg(R)$  in each of these states, this holds until the end of the sequence, proving the claim.  $\square$

**THEOREM 6.** *There exists an optimal offline algorithm for FTF on disjoint sequences that upon each fault evicts a page  $\sigma \in R_j$  whose next request time is maximal in  $R_j$ , for some  $j$ .*

**PROOF.** As in the proof of Theorem 5, we claim that for any algorithm  $Alg$  there exists an algorithm  $Alg_i$  that behaves exactly like  $Alg$  until time  $i - 1$ . If at time  $i$   $Alg$  evicts a page from sequence  $R_j$ ,  $Alg_i$  evicts the page from the same sequence that is furthest in the future, and  $Alg_i(\mathcal{R}) \leq Alg(\mathcal{R})$ . Applying this claim on an optimal algorithm successively at each timestep gives an optimal algorithm.

We now prove the claim. Suppose that at time  $t = i$   $Alg$  evicts a page  $\sigma_1 \in R_j$ .  $Alg_i$  behaves exactly like  $Alg$  until  $t = i - 1$  but at  $t = i$  it evicts a page  $\sigma_2 \in R_j$ , which is the page furthest in the future in this sequence. Assume  $\sigma_1 \neq \sigma_2$  as otherwise the claim is trivially true. Let  $C_{Alg}(t)$  and  $Alg(\mathcal{R}, t)$  denote the contents of  $Alg$ 's cache and the number of faults before serving  $R(t)$ . We have  $C_{Alg}(t) = C_{Alg_i}(t)$  and  $C_{Alg}(t + \tau) = (C_{Alg_i}(t + \tau) \cup \{\sigma_2\}) \setminus \{\sigma_1\}$ , and the sequences have the same alignment in both algorithms.

Before the request for  $\sigma_1$ , if upon a fault  $Alg$  evicts  $\sigma_2$ ,  $Alg_i$  evicts  $\sigma_1$ , then the caches are equal and from then on  $Alg_i$  behaves exactly like  $Alg$ , thus  $Alg_i(\mathcal{R}) \leq Alg(\mathcal{R})$ . If  $Alg$  instead evicts a page  $\alpha \neq \sigma_2$ ,  $Alg_i$  evicts the same page.

If  $\sigma_2$  was not evicted, when  $\sigma_1$  is requested  $Alg$  faults and evicts a page  $\alpha$  from some sequence.  $\sigma_1$  is a hit for  $Alg_i$  and thus  $Alg(\mathcal{R}, t) = Alg_i(\mathcal{R}, t) + 1$ ,  $C_{Alg}(t) = (C_{Alg_i}(t) \cup \{\sigma_2\}) \setminus \{\alpha\}$ , and  $Alg_i$  is ahead in  $R_j$  by  $\tau$  timesteps. Assume first that  $\sigma_2$  is not evicted by  $Alg$  before its request. Suppose  $Alg_i$  gets to  $\sigma_2$  with this configuration (for some  $\alpha$ ).  $\sigma_2$  is a fault for  $Alg_i$  and a hit for  $Alg$ .  $Alg_i$  evicts  $\alpha$  and now  $Alg(\mathcal{R}, t) = Alg_i(\mathcal{R}, t)$ ,  $C_{Alg}(t) = C_{Alg_i}(t)$  and the sequences are aligned equally. From then on both algorithms are equivalent and the claim is true. After the request for  $\sigma_1$ ,  $Alg_i$  evicts whatever  $Alg$  evicts (assume  $\sigma_2$  is not evicted by  $Alg$

<sup>9</sup>Since sequences  $R_s$  in the execution of  $Alg$  and  $Alg_i$  are either aligned or  $Alg_i$ 's sequence is behind by exactly  $\tau$ , if the current request to  $Alg$  is to the same page  $p_j$  that is being fetched by  $Alg_i$ , then this is the timestep in which  $Alg_i$  finishes fetching  $p_j$ .

<sup>10</sup> $Alg$  could evict  $p$  prior to its request according to the logical order of simultaneous requests.

during this period). If  $\alpha$  is requested (a fault for  $Alg$  but a hit for  $Alg_i$ ),  $Alg_i$  forces a fault and hence  $Alg(\mathcal{R}, t) = Alg_i(\mathcal{R}, t) + 1$  and  $C_{Alg}(t) = (C_{Alg_i}(t) \cup \{\sigma_2\}) \setminus \{\alpha\}$  still holds, with  $\alpha'$  being the page evicted by  $Alg$ . By Theorem 5, we can build another algorithm  $Alg'_i$  with the same number of faults that does not force faults, and hence the claim is true in this case. Now, if  $Alg$  evicts  $\sigma_2$  before getting to its request  $Alg_i$  evicts  $\alpha$ , and both caches are the same. When  $\sigma_2$  is requested both algorithms will fault and  $Alg(\mathcal{R}, t) = Alg_i(\mathcal{R}, t) + 1$  holds. However, sequences  $R_j$  in both algorithms do not have the same alignment with respect to the rest of the sequences, with  $Alg_i$ 's sequence being ahead by  $\tau$  timesteps. This setting is the same as the one in the proof of Theorem 5 after  $Alg$  has forced a fault. Applying that proof we can show that  $Alg_i$  can keep the execution within the 5 states defined in the proof, and hence  $Alg_i(\mathcal{R}) \leq Alg(\mathcal{R})$  at the end of the execution. Again, if at any point  $Alg_i$  forces a fault, then by the same Theorem 5 we can obtain an honest algorithm that does not exceed  $Alg_i$ 's faults. Since in all cases the claim is true, this proves the theorem.  $\square$

### A.3 Pseudocode of the Algorithm for PIF.

**Algorithm 2** Partial Individual Faults( $\mathcal{R}, K, time, \tau, \vec{b}$ )

---

```

for all configurations  $C$  do
   $F[C, 1, \dots, 1] = \{(\{0\}^p, 0)\}$ 
  for each  $(x_1, \dots, x_p) \in \{2, \dots, n_i(\tau + 1) + 1\}^p$  do
     $F[C, x_1, \dots, x_p] = \emptyset$ 
  for each  $(x_1, \dots, x_p) \in \{1, \dots, n_i(\tau + 1)\}^p$  do
    for all configurations  $C$  do
      if  $F[C, x_1, \dots, x_p] \neq \emptyset$  then
        for  $i = 1$  to  $p$  do
          if  $x_i = (j - 1)(\tau + 1) + 1$  for some  $j$ , and  $R_i(x_i) \in C$  then  $\{R_i(x_i)$  is a hit $\}$ 
             $x'_i = x_i + \tau + 1$ 
          else if  $R_i(x_i) = R_\ell(x_\ell)$  and  $x_\ell = (j' - 1)(\tau + 1)$  for some  $j'$  and some  $\ell \neq i$  then  $\{R_i(x_i)$  was fetched for another sequence $\}$ 
             $x'_i = \lceil x_i / (\tau + 1) \rceil (\tau + 1) + 1$ 
          else
             $x'_i = x_i + 1$ 
           $P = \emptyset$   $\{P$  is the list of updated fault vectors $\}$ 
        for each  $(\vec{f}, t)$  in  $F[C, x_1, \dots, x_p]$  do
          validVector = true
          for  $i = 1$  to  $p$  do
             $f'_i = f_i$   $\{f_i$  is the  $i$ -th element of  $\vec{f}\}$ 
            if  $R_i(x_i) \notin C$  then  $\{R_i(x_i)$  is a fault $\}$ 
               $f'_i = f'_i + 1$ 
            if  $f'_i > b_i$  then  $\{$ this path exceeded the maximum faults for sequence  $i$  $\}$ 
              validVector = false
          if validVector and  $t + 1 \leq time$  then
             $P = P \cup (f', t + 1)$ 
            if  $t + 1 = time$  or  $x'_i = n_i(\tau + 1) + 1$  for all  $i = 1..p$  then  $\{$ we reached the checkpoint time or the end of all sequences $\}$ 
              return TRUE
        for all  $C'$  s.t.  $\mathcal{R}(\vec{x}) \subseteq C'$  and  $(C' \setminus \mathcal{R}(\vec{x})) \subseteq C$  do
           $F[C', x'_1, \dots, x'_p] = F[C', x'_1, \dots, x'_p] \cup P$ 
           $\{$ we reached the end without finding a feasible solution $\}$ 
        return FALSE

```

---