

The Relative Worst Order Ratio Applied to Paging*

Joan Boyar, Lene M. Favrholdt**, and Kim S. Larsen

Department of Mathematics and Computer Science, University of Southern Denmark, Odense, Denmark,
{joan, lenem, kslarsen}@imada.sdu.dk

Abstract. The relative worst order ratio, a new measure for the quality of on-line algorithms, was recently defined and applied to two bin packing problems. Here, we extend the definition and apply it to the paging problem. Work in progress by other researchers shows that the measure gives interesting results and new separations for bin coloring and scheduling problems as well.

Using the relative worst order ratio, we obtain the following results: We devise a new deterministic paging algorithm, Retrospective-LRU, and show that it performs better than LRU. This is supported by experimental results, but contrasts with the competitive ratio. All deterministic marking algorithms have the same competitive ratio, but here we find that LRU is better than FWF. Furthermore, LRU is better than PERM_π , which also has the same competitive ratio as LRU. No deterministic marking algorithm can be significantly better than LRU, but the randomized algorithm MARK is better than LRU. Finally, look-ahead is shown to be a significant advantage, in contrast to the competitive ratio, which does not reflect that look-ahead can be helpful.

1 Introduction

The standard measure for the quality of on-line algorithms is the competitive ratio [17, 28, 20], which is, roughly speaking, the worst-case ratio, over all possible input sequences, of the on-line performance to the optimal off-line performance. The definition of the competitive ratio is essentially identical to that of the approximation ratio. This seems natural in that on-line algorithms can be viewed as a special class of approximation algorithms. However, for approximation algorithms, the comparison to an optimal off-line algorithm, OPT, is natural, since the approximation algorithm is compared to another algorithm of the same general type, just with more computing power, while for on-line algorithms, the comparison to OPT is to a different type of algorithm.

Although the competitive ratio has been an extremely useful notion, in many cases, and particularly for the paging problem, it has appeared inadequate at differentiating between on-line algorithms. In a few cases (bin coloring [25] and dual bin packing [9]), one algorithm \mathbb{A} even has a better competitive ratio than another algorithm \mathbb{B} , though intuitively, \mathbb{B} is clearly better than \mathbb{A} .

When differentiating between on-line algorithms is the goal, performing a direct comparison between the algorithms, instead of involving an intermediate comparison to OPT, seems the obvious choice. A direct comparison on exactly the same sequences will produce the result that many algorithms are not comparable, because one algorithm does well on one type of ordering, while the other does well on another type. With the relative worst order ratio, on-line algorithms are compared directly to each other on their respective worst permutations of sequences. In this way, the relative worst order ratio [8] combines some of the desirable properties of the Max/Max ratio [5] and the random order ratio [22].

The Max/Max Ratio The Max/Max ratio [5] allows direct comparison of two on-line algorithms for an optimization problem, without the intermediate comparison to OPT. Rather than comparing two algorithms on the same sequence, they are compared on their respective worst-case sequences of

* Supported in part by the Danish Natural Science Research Council (SNF) and in part by the Future and Emerging Technologies program of the EU under contract number IST-1999-14186 (ALCOM-FT).

** Part of this work was done while working at the Department of Computer Science, University of Copenhagen, Denmark. Supported in part by the Villum Kann Rasmussen Fund.

the same length. The Max/Max Ratio applies only when the length of an input sequence yields a bound on the profit/cost of an optimal solution. Technically, it applies to the paging problem, but the Max/Max ratio of any paging algorithm (deterministic or randomized) approaches 1 as the size of the slow memory approaches infinity.

The Random Order Ratio The random order ratio [22] gives the possibility of considering some randomness of the request sequences without specifying a complete probability distribution. For an on-line algorithm \mathbb{A} , the random order ratio is the worst-case ratio, over all sequences of requests, of the expected performance of \mathbb{A} on a random permutation of the sequence, compared with an optimal solution. If, for all possible sequences of requests, any ordering of these requests is equally likely, this ratio gives a meaningful worst-case measure of how well an algorithm can do. Unfortunately, the random order ratio seems to be difficult to compute.

The Relative Worst Order Ratio With the relative worst order ratio, one considers the worst-case performance over all permutations instead of the average-case performance as with the random order ratio. Thus, when comparing two on-line algorithms, one considers a worst-case sequence and takes the ratio of how the two algorithms perform on their respective worst permutations of that sequence. Note that the two algorithms may have different “worst orderings” for the same sequence.

The relative worst order ratio can be viewed as a worst case version of Kenyon’s random order ratio, with the modification that on-line algorithms are compared directly, rather than indirectly through OPT. It can also be viewed as a modification of the Max/Max ratio, where a finer partition of the request sequences is used; instead of finding the worst sequence among those having the same length, one finds the worst sequence among those which are permutations of each other. This particular finer partition was inspired by the random order ratio.

The Paging Problem We consider the well studied paging problem. The input sequence consists of requests for pages in a slow memory, which contains N pages. There is a fast memory, called the cache, which has space for $k < N$ pages. A request for a page currently in cache is a *hit*, while a request for a page not in cache is a *page fault*. When a page fault occurs, the requested page must be brought into the cache. If the cache already contains k pages when this happens, at least one of them must be evicted. A paging algorithm decides which page to evict on a fault. This decision must usually be made on-line, i.e., without any knowledge about future requests. The goal is to minimize the number of faults.

Paging Algorithms Two major classes of deterministic algorithms for the paging problem are conservative algorithms [32] and marking algorithms [7].

A paging algorithm \mathbb{A} is called *conservative*, if no request sequence has a consecutive subsequence with requests to at most k distinct pages causing \mathbb{A} to fault more than k times. The algorithms, Least-Recently-Used (LRU) and First-In/First-Out (FIFO) are examples of conservative algorithms. On a page fault, LRU evicts the least recently used page in cache and FIFO evicts the page which has been in cache longest.

Marking algorithms work in phases. Each time a page is requested, this page is marked. When a page must be evicted, one of the unmarked pages is chosen, if one exists. Otherwise all marks are erased, and the requested page is marked. This request starts a new phase. Note that LRU is a marking algorithm, whereas FIFO is not. Another example of a marking algorithm is Flush-When-Full (FWF), the algorithm which evicts all pages in cache at the end of each phase. The randomized marking algorithm MARK chooses the unmarked page to be evicted uniformly at random.

We also study $\text{MARK}^{\text{LIFO}}$, LIFO, and PERM_π [6] defined in Sections 4 and 6.

Previous Results All conservative and marking algorithms have competitive ratio k [31, 29] and this is optimal among deterministic algorithms [28]. However, in practice, these algorithms do not all have the same performance: LRU is better than FIFO and much better than FWF [32]. Moreover, results from [15] suggest there may be algorithms that perform even better than LRU.

In [3] an alternative model, the Max-/Average-Model, for the paging problem capturing locality of reference was suggested. It was proven that, in this model, LRU is slightly better than FIFO, but LRU is still best possible among deterministic algorithms. Using access graphs, it has been proven that LRU is better than FIFO [12] and algorithms have been designed that are better than LRU [7]. Hence, these alternative ways of measuring the quality of paging algorithms give more satisfactory results. However, they are only defined for paging and paging-like problems.

In contrast to deterministic algorithms, MARK [14] has been shown to have a competitive ratio of $2H_k - 1$ [1], where H_k is the k th harmonic number, i.e., $H_k = \sum_{i=1}^k \frac{1}{i} \approx \ln k$. Other randomized algorithms have been shown to have the optimal competitive ratio for randomized algorithms of H_k [26, 1].

Look-Ahead. Look-ahead, where the algorithm deciding which page to evict is allowed to see the next ℓ page requests before making that decision, is a model which intuitively lies between on-line and off-line. It is well known that look-ahead cannot reduce the competitive ratio of any algorithm, but clearly it can be useful when it can be implemented.

Previously, alternative definitions of look-ahead have led to results showing that look-ahead helps. In each case, the algorithm is allowed to see a sequence of future requests satisfying some property. Young [31] proposed *resource-bounded look-ahead*, where the sequence is a maximal sequence of future requests for which it would incur ℓ page faults, Albers [2] proposed *strong look-ahead*, where the sequence of future requests contains ℓ distinct pages different from the current request, and Breslauer [10] proposed *natural look-ahead*, where the sequence of future requests contains ℓ pages not currently in cache. Here we retain the original definition, so the algorithm is only allowed to see the next ℓ pages, regardless of what they are.

The Max/Max Ratio [5] has been somewhat successfully applied to the standard definition of look-ahead, showing that a greedy strategy achieves a Max/Max ratio of $\frac{N-1}{\ell}$ for $N - k < \ell \leq N - 1$ (recall that N is the size of the slow memory). *Comparative analysis* [24] is more successful, showing that look-ahead gives a result which is a factor $\min\{k, \ell + 1\}$ better than without look-ahead. This is the same result we obtain with the relative worst order ratio.

Other Measures. Many alternatives to or variations on the competitive ratio have been proposed. We have already mentioned the Max/Max ratio, the random order ratio, access graphs, the Max-/Average-Model, and comparative analysis. Other alternatives are Markov paging [21], diffuse adversaries [24], extra resource analysis [19, 28], the accommodating function [9], and statistical adversaries [27]. Most of these techniques have been applied to only a few closely related problems. So far, the techniques which have been applied to a broader range of problems, extra resource analysis and the accommodating function, for instance, have given new separation results for only a limited number of different types of problems.

The Relative Worst Order Ratio. In contrast to previous measures, the relative worst order ratio has already been applied quite successfully to two very different problem types: bin packing [8] and now paging, here. For Classical Bin Packing, Worst-Fit is better than Next-Fit according to the relative worst order ratio, even though they have the same competitive ratio [18]. Thus, the advantage of keeping all bins open, instead of just one, is reflected by the relative worst order ratio. For Dual Bin Packing, the relative worst order ratio shows that First-Fit is better than Worst-Fit, while the competitive ratio indicates the opposite [9].

Other New Results on the Relative Worst Order Ratio The wide applicability of the relative worst order ratio has been confirmed by other new results. Recently, various researchers have applied the relative worst order ratio to other problems and obtained separations not given by the competitive ratio, but consistent with intuition and/or practice.

One simple example is given in [13]. For the problem of minimizing makespan on two related machines with speed ratio s , the optimal competitive ratio of $\frac{s+1}{s}$ for $s \geq \Phi \approx 1.618$ is obtained both by the post-greedy algorithm, which schedules each job on the machine where it will finish earliest, and by the algorithm which simply schedules all jobs on the fast machine. In contrast, the relative worst order ratio shows that the post-greedy algorithm is better. A similar result is obtained for the problem of minimizing makespan on $m \geq 2$ identical machines with preemption.

The relative worst order ratio was also found by [23] to give the intuitively correct result for the bin coloring problem, where the competitive ratio gives the opposite result [25]: a trivial algorithm using only one open bin has a better competitive ratio than a natural greedy-type algorithm.

New Results on Paging First, we propose a new algorithm, Retrospective-LRU (RLRU), which is a variation on LRU that takes into account which pages would be in the cache of the optimal off-line algorithm, LFD, if it were given the subsequence of page requests seen so far. We show that, according to the relative worst order ratio, RLRU is better than LRU. This is interesting, since it contrasts with results on the competitive ratio and with results in [3] where a new model of locality of reference is studied.

It is easily shown that RLRU does not belong to either of the common classes of algorithms, conservative and marking algorithms which all have the optimal competitive ratio k . In fact, the competitive ratio of RLRU is $k + 1$ and thus slightly worse than that of LRU. Initial testing of RLRU indicates that it may perform better than LRU in practice.

Analyzing paging algorithms with the relative worst order ratio, we obtain far more detailed information than with competitive analysis. With the relative worst order ratio, LRU is better than FWF, so not all marking algorithms are equivalent, but no marking algorithm is significantly better than LRU. All conservative algorithms are equivalent, so LRU and FIFO have the same performance, but LRU is better than the k -competitive algorithm PERM_π . MARK is better than LRU, which is consistent with competitive analysis.

Look-ahead is shown to help significantly with respect to the relative worst order ratio. Compared to the competitive ratio which does not reflect that look-ahead can be of any use, this is a very nice property of the relative worst order ratio.

A new phenomenon with respect to the relative worst order ratio is observed: in [8], the pairs of algorithms investigated were either *comparable* or *incomparable*, but here some are found to be *weakly comparable*, i.e., while one algorithm performs marginally better than the other on some sequences and their permutations, the other algorithm performs significantly better on other sequences and their permutations. Furthermore, algorithms can be *asymptotically comparable*, which for the paging problem means that, for arbitrarily large cache sizes, the pair of algorithms are “arbitrarily close to being comparable”. This is defined more formally in Section 2.

The definition of the relative worst order ratio is extended to randomized algorithms, and MARK and LRU are compared, giving that MARK is the better algorithm, as with the competitive ratio.

2 The New Measure

In this section, we define the relative worst order ratio and the notion of two algorithms being comparable (Definition 2) as in [8]. This is the most important definition, but the new notions of being weakly comparable and asymptotically comparable (defined in Definitions 4 and 5) give the possibility of adding more detail to the description of the relation between two algorithms.

	minimization	maximization
A better than B	< 1	> 1
B better than A	> 1	< 1

Table 1. Ratio values for minimization and maximization problems

The Relative Worst Order Ratio. The definition of the relative worst order ratio uses $\mathbb{A}_W(I)$, the performance of an on-line algorithm \mathbb{A} on the “worst ordering” of the sequence I of requests, formally defined in the following way.

Definition 1. Consider an on-line optimization problem P and let I be any request sequence of length n . If σ is a permutation on n elements, then $\sigma(I)$ denotes I permuted by σ . Let \mathbb{A} be any algorithm for P .

If P is a maximization problem, $\mathbb{A}(I)$ is the value of running \mathbb{A} on I , and $\mathbb{A}_W(I) = \min_{\sigma} \mathbb{A}(\sigma(I))$.
If P is a minimization problem, $\mathbb{A}(I)$ is a cost, and $\mathbb{A}_W(I) = \max_{\sigma} \mathbb{A}(\sigma(I))$.

For many on-line problems, some algorithms perform well on particular types of orderings of the input, while other algorithms perform well on other types of orderings. The purpose of comparing on the worst permutation of sequences, rather than on each sequence independently, is to be able to differentiate between such pairs of algorithms, rather than just concluding that they are incomparable. Sequences with the same “content” are considered together, but the measure is worst case, so the algorithms are compared on their respective worst permutations. This was originally motivated by problems where all orderings are equally likely, but appears to be applicable to other problems as well.

Definition 2. Let $S_1(c)$ and $S_2(c)$ be statements about algorithms \mathbb{A} and \mathbb{B} defined in the following way.

$S_1(c)$: There exists a constant b such that $\mathbb{A}_W(I) \leq c \cdot \mathbb{B}_W(I) + b$ for all I .

$S_2(c)$: There exists a constant b such that $\mathbb{A}_W(I) \geq c \cdot \mathbb{B}_W(I) - b$ for all I .

The relative worst order ratio $WR_{\mathbb{A},\mathbb{B}}$ of on-line algorithm \mathbb{A} to algorithm \mathbb{B} is defined if $S_1(1)$ or $S_2(1)$ holds. In this case, \mathbb{A} and \mathbb{B} are said to be comparable.

If $S_1(1)$ holds, then $WR_{\mathbb{A},\mathbb{B}} = \sup \{r \mid S_2(r)\}$, and if $S_2(1)$ holds, then $WR_{\mathbb{A},\mathbb{B}} = \inf \{r \mid S_1(r)\}$.

The sentences $S_1(1)$ and $S_2(1)$ check that the one algorithm is always at least as good as the other on every sequence (on their respective worst permutations). When one of them holds, the relative worst order ratio is a bound on how much better the one algorithm can be. Note that if $S_1(1)$ holds, the supremum involves S_2 rather than S_1 , and vice versa. A ratio of 1 means that the two algorithms perform identically with respect to this quality measure; the further away from 1, the greater the difference in performance. The ratio may be greater than or less than one, depending on whether the problem is a minimization or a maximization problem and on which of the two algorithms is better. These possibilities are illustrated in Table 1.

It is easily shown [8] that the relative worst order ratio is a *transitive measure*, i.e., for any three algorithms \mathbb{A} , \mathbb{B} , and \mathbb{C} , $WR_{\mathbb{A},\mathbb{B}} \geq 1$ and $WR_{\mathbb{B},\mathbb{C}} \geq 1$ implies $WR_{\mathbb{A},\mathbb{C}} \geq 1$. The proof of transitivity shows that when $WR_{\mathbb{A},\mathbb{B}} \geq 1$, $WR_{\mathbb{B},\mathbb{C}} \geq 1$, and both are bounded above by some constant, then $\max(WR_{\mathbb{A},\mathbb{B}}, WR_{\mathbb{B},\mathbb{C}}) \leq WR_{\mathbb{A},\mathbb{C}} \leq WR_{\mathbb{A},\mathbb{B}} \cdot WR_{\mathbb{B},\mathbb{C}}$. Thus, when a new algorithm is analyzed, it need not be compared to all algorithms.

Although one of the goals in defining the relative worst order ratio was to avoid the intermediate comparison of any on-line algorithm, \mathbb{A} , to the optimal off-line algorithm, OPT , it is still possible to compare on-line algorithms to OPT . In this case, the measure is called the *worst order ratio* and denoted $WR_{\mathbb{A}}$.

Relatedness. Even if a pair of algorithms are not comparable, there may be something interesting to say about their relative performance. Therefore, we introduce the notion of relatedness that applies to most pairs of algorithms.

Definition 3. Let \mathbb{A} and \mathbb{B} be algorithms for an on-line optimization problem P , and let S_1 and S_2 be defined as in Definition 2.

Assume first that P is a minimization problem. If there exists a positive constant c such that $S_1(c)$ is true, let $c_{\mathbb{A},\mathbb{B}} = \inf \{r \mid S_1(r)\}$. Otherwise, $c_{\mathbb{A},\mathbb{B}}$ is undefined.

- If $c_{\mathbb{A},\mathbb{B}}$ and $c_{\mathbb{B},\mathbb{A}}$ are both defined, \mathbb{A} and \mathbb{B} are $(c_{\mathbb{A},\mathbb{B}}, c_{\mathbb{B},\mathbb{A}})$ -related.
- If $c_{\mathbb{A},\mathbb{B}}$ is defined and $c_{\mathbb{B},\mathbb{A}}$ is undefined, \mathbb{A} and \mathbb{B} are $(c_{\mathbb{A},\mathbb{B}}, \infty)$ -related.
- If $c_{\mathbb{A},\mathbb{B}}$ is undefined and $c_{\mathbb{B},\mathbb{A}}$ is defined, \mathbb{A} and \mathbb{B} are $(\infty, c_{\mathbb{B},\mathbb{A}})$ -related.

If P is a maximization problem, relatedness is defined similarly: If there exists a positive constant c such that $S_2(c)$ is true, let $c_{\mathbb{A},\mathbb{B}} = \sup \{r \mid S_2(r)\}$. Otherwise, $c_{\mathbb{A},\mathbb{B}}$ is undefined.

- If $c_{\mathbb{A},\mathbb{B}}$ and $c_{\mathbb{B},\mathbb{A}}$ are both defined, \mathbb{A} and \mathbb{B} are $(c_{\mathbb{A},\mathbb{B}}, c_{\mathbb{B},\mathbb{A}})$ -related.
- If $c_{\mathbb{A},\mathbb{B}}$ is defined and $c_{\mathbb{B},\mathbb{A}}$ is undefined, \mathbb{A} and \mathbb{B} are $(c_{\mathbb{A},\mathbb{B}}, 0)$ -related.
- If $c_{\mathbb{A},\mathbb{B}}$ is undefined and $c_{\mathbb{B},\mathbb{A}}$ is defined, \mathbb{A} and \mathbb{B} are $(0, c_{\mathbb{B},\mathbb{A}})$ -related.

This notation can also be used for algorithms which are comparable. In this case, one of the values is the relative worst order ratio and the other is typically 1 (unless one algorithm is strictly better than the other in all cases).

Weakly and Asymptotically Comparable. In Section 6, it is shown that LRU and Last-In/First-Out (LIFO) are $(\frac{k+1}{2}, \infty)$ -related. With this result it seems reasonable to prefer LRU to LIFO, even though they are not comparable by Definition 2. We say, therefore, that the pair of algorithms are *weakly comparable*.

Definition 4. Let \mathbb{A} and \mathbb{B} be algorithms for an on-line optimization problem P and let $c_{\mathbb{A},\mathbb{B}}$ be defined as in Definition 3. \mathbb{A} and \mathbb{B} are *weakly comparable* if \mathbb{A} and \mathbb{B} are comparable, if exactly one of $c_{\mathbb{A},\mathbb{B}}$ and $c_{\mathbb{B},\mathbb{A}}$ is defined, or if both are defined and $c_{\mathbb{A},\mathbb{B}} \notin \Theta(c_{\mathbb{B},\mathbb{A}})$.

More specifically, if P is a minimization problem and $c_{\mathbb{A},\mathbb{B}} \in o(c_{\mathbb{B},\mathbb{A}})$, or if P is a maximization problem and $c_{\mathbb{A},\mathbb{B}} \in \omega(c_{\mathbb{B},\mathbb{A}})$, \mathbb{A} and \mathbb{B} are *weakly comparable in \mathbb{A} 's favor*. Similarly, if $c_{\mathbb{A},\mathbb{B}}$ is defined and $c_{\mathbb{B},\mathbb{A}}$ is undefined, \mathbb{A} and \mathbb{B} are *weakly comparable in \mathbb{A} 's favor*.

We conclude with a definition which is relevant for optimization problems with some limited resource, such as the size of the cache in the paging problem, the capacity of the knapsack in a knapsack problem, or the number of machines in a machine scheduling problem.

Definition 5. A resource dependent problem is an on-line optimization problem, where each problem instance, in addition to the input data given on-line, also has a parameter k , referred to as the amount of resources, such that for each input, the optimal solution depends monotonically on k .

Let \mathbb{A} and \mathbb{B} be algorithms for a resource dependent problem P and let $c_{\mathbb{A},\mathbb{B}}$ be defined as in Definition 3. \mathbb{A} and \mathbb{B} are *asymptotically comparable*, if

$$\left(\lim_{k \rightarrow \infty} \{c_{\mathbb{A},\mathbb{B}}\} \leq 1 \text{ and } \lim_{k \rightarrow \infty} \{c_{\mathbb{B},\mathbb{A}}\} \geq 1 \right) \text{ or } \left(\lim_{k \rightarrow \infty} \{c_{\mathbb{A},\mathbb{B}}\} \geq 1 \text{ and } \lim_{k \rightarrow \infty} \{c_{\mathbb{B},\mathbb{A}}\} \leq 1 \right)$$

Let \mathbb{A} and \mathbb{B} be asymptotically comparable algorithms. For a minimization problem, \mathbb{A} and \mathbb{B} are asymptotically comparable in \mathbb{A} 's favor if $\lim_{k \rightarrow \infty} \{c_{\mathbb{B},\mathbb{A}}\} > 1$; and for a maximization problem, if $\lim_{k \rightarrow \infty} \{c_{\mathbb{B},\mathbb{A}}\} < 1$.

If one algorithm is not clearly better than the other, as captured by the notions comparable, weakly comparable, and asymptotically comparable, we say the pair of algorithms are *incomparable*.

Definition 6. Let \mathbb{A} and \mathbb{B} be algorithms for an on-line optimization problem. If \mathbb{A} and \mathbb{B} are neither weakly nor asymptotically comparable, we say that they are *incomparable*.

3 A Better Algorithm than LRU

In this section, we introduce an algorithm which turns out to be better than LRU according to the relative worst order ratio. This is in contrast to the competitive ratio which says that LRU is best possible among deterministic algorithms. The algorithm, called Retrospective-LRU (RLRU), is defined in Figure 1. The name comes from the algorithm's marking policy. When evicting pages, RLRU uses the LRU policy, but it chooses only from the unmarked pages in cache, unless they are all marked. Marks are set according to what the optimal off-line algorithm, LFD [4], would have in cache, if given the part of the sequence seen so far. LFD is the algorithm that, on a fault, evicts the page that will be requested farthest in the future.

If RLRU has a fault and LFD does not, RLRU marks the page requested. If RLRU has a hit, the page p requested is marked if it is different from the page of the previous request. Requiring the page to be different from the previous page ensures that at least one other page has been requested since p was brought into the cache. A phase of the execution starts with the removal of all marks and this occurs whenever there would otherwise be a second fault on the same page within the current phase.

The first phase begins with the first request.

On request r to page p :

Update p 's timestamp

if p is not in cache **then**

if there is no unmarked page **then**

evict the least recently used page in cache

else

evict the least recently used unmarked page

if this is the second fault on p

since the start of the current phase **then**

unmark all pages

start a new phase with r

if p was in LFD's cache just before this request **then**

mark p

else

if p is different from the previous page **then**

mark p

Fig. 1. Retrospective-LRU (RLRU)

Lemma 1. *For any request sequence, each complete phase defined by RLRU contains requests to at least $k + 1$ distinct pages.*

Proof. Consider any phase P and the page p which starts the next phase. Page p was requested in phase P , and was later evicted, also within phase P . At that point, all other pages in the cache must either be marked or have been requested since the last request to p , so every page in cache at that point has been requested in phase P . The page requested when p is evicted must be different from the k pages in cache at that point. Thus, there must be at least $k + 1$ different pages requested in phase P . □

Lemma 2. *For any sequence I of page requests, $RLRU_W(I) \leq LRU_W(I)$.*

Proof. Consider a worst ordering I_{RLRU} of I with respect to RLRU. By definition, RLRU never faults twice on the same page within any single phase of I_{RLRU} .

Move the last, possibly incomplete, phase of I_{RLRU} to the beginning and call the resulting sequence I_{LRU} . Process the requests in this sequence phase by phase (the phases are the original RLRU phases), starting at the beginning. LRU faults on each distinct page in the first phase. Since there are at least $k + 1$ distinct pages in each of the later phases, all of the distinct pages in a phase can be ordered so that there will be a fault by LRU on each of them. Hence, in each phase, LRU faults at least as many times as RLRU, i.e., LRU has at least as many faults on I_{LRU} as RLRU on I_{RLRU} . \square

Theorem 1. $WR_{\text{LRU},\text{RLRU}} = \frac{k+1}{2}$.

Proof. Since the previous lemma shows that $WR_{\text{LRU},\text{RLRU}} \geq 1$, for the lower bound, it is sufficient to find a family of sequences I_n with $\lim_{n \rightarrow \infty} \text{LRU}(I_n) = \infty$, where there exists a constant b such that for all I_n ,

$$\text{LRU}_W(I_n) \geq \frac{k+1}{2} \text{RLRU}_W(I_n) - b.$$

Let I_n consist of n phases, where, in each phase, the first $k - 1$ requests are to the $k - 1$ pages p_1, p_2, \dots, p_{k-1} , always in that order, and the last two requests are to completely new pages. LRU will fault on every page, so it will fault $n(k + 1)$ times.

Regardless of the order this sequence is given in, LFD will never evict the pages p_1, p_2, \dots, p_{k-1} from cache, so RLRU will mark them the first time they are requested in each phase, if they have ever been requested before. Thus, for each of these pages p' , at most one phase is ended because of a fault on p' , and the requests to the pages which only occur once cannot end phases. This gives at most $k - 1$ phases, each containing at most one fault on each of the pages p_1, p_2, \dots, p_{k-1} , which limits the number of faults RLRU has on these $k - 1$ pages to a constant (dependent on k , but not n), so RLRU faults at most $2n + c$ times for some constant c . Asymptotically, the ratio is $\frac{k+1}{2}$.

For the upper bound, suppose there exists a sequence I , where LRU faults s times on its worst permutation, I_{LRU} , RLRU faults s' times on its worst ordering, I_{RLRU} , and $s > \frac{k+1}{2} \cdot s'$. Then, $s > \frac{k+1}{2} \cdot s''$, where s'' is the number of times RLRU faults on I_{LRU} . Assume by Lemma 3, that I_{LRU} is such that LRU faults on each request of a prefix I_1 of I_{LRU} and on no request after I_1 . Then there must exist a subsequence, $J = \langle r_1, r_2, \dots, r_{k+1} \rangle$, of consecutive requests in I_1 , where RLRU faults at most once. Since LRU faults on every request, they must be to $k + 1$ different pages. One may assume that r_1 is not the first request, since then RLRU would fault on all the requests in J . Let p be the page requested immediately before J . Clearly, p must be in RLRU's cache when it begins processing J . If r_{k+1} is not a request to p , then the fact that LRU faulted on every request in J means that J contains $k + 1$ pages different from p , but at most $k - 1$ of them could be in RLRU's cache when it begins processing J . Thus, LRU must fault at least twice on the requests in J . On the other hand, if r_{k+1} is a request to p , there are exactly k requests in J which are different from p . At least one of them must cause a fault, since at most $k - 1$ of them could have been in the cache when RLRU began processing J . If no others caused faults, then they must have all been marked. In this case RLRU evicts the least recently used page in cache, which cannot be a page requested in J before this fault, so it must be a later page in J , causing a second fault. This is a contradiction. \square

The proofs which establish that RLRU is better than LRU according to the relative worst order ratio rely on a few basic properties of RLRU. Modifications to the algorithm which do not change these basic properties will result in other algorithms which, according to the relative worst order ratio, are also better than LRU. One example of this is the test as to whether or not the current page is the same as the previous. This test could be removed and the page marked unconditionally and the proofs still hold. Another example is the decision when to end a phase. The important property is that each phase consists of requests to at least $k + 1$ distinct pages and there is at most one fault on each of them. This leaves room for experimentally testing a number of variations, and it could lead to algorithms which are even better in practice than the one we have found here.

Note that RLRU is neither a conservative nor a marking algorithm. This can be seen from the sequence $\langle p_1, p_2, p_3, p_4, p_1, p_2, p_3, p_4, p_3 \rangle$ for $k = 3$, where RLRU faults on every request.

In contrast to Theorem 1, the competitive ratio of RLRU is slightly worse than that of LRU:

Theorem 2. *The competitive ratio of RLRU is $k + 1$.*

Proof. The upper bound of $k + 1$ follows since each phase of the algorithm contains requests to at least $k + 1$ different pages, and RLRU faults at most once on each page within a phase. If there are $s > k$ different pages in a phase, OPT must fault at least $s - k$ times in that phase. The worst ratio is obtained when there are exactly $k + 1$ different pages in a phase, giving a ratio of $k + 1$.

The lower bound follows from a sequence with $k + 1$ distinct pages p_1, p_2, \dots, p_{k+1} , where each request is to the page not in RLRU's cache. This sequence is $\langle p_1, p_2, \dots, p_{k+1} \rangle^2 \langle p_k, p_1, p_2, \dots, p_{k-1}, p_{k+1} \mid p_{k-1}, p_k, p_1, p_2, \dots, p_{k-2}, p_{k+1} \mid p_{k-2}, p_{k-1}, p_k, p_1, p_2, \dots, p_{k-3}, p_{k+1} \mid \dots \mid p_1, p_2, \dots, p_{k+1} \rangle^n$, where \mid marks the beginning of a new phase. The part of the sequence which is repeated n times is called a superphase and consists of k phases, the i th phase consisting of the sequence $\langle p_{k+1-i}, \dots, p_k, p_1, \dots, p_{k-i}, p_{k+1} \rangle$, for $1 \leq i \leq k - 1$, and $\langle p_1, p_2, \dots, p_{k+1} \rangle$, for $i = k$. The optimal strategy is to evict page p_{k-1-i} in the i th phase of a superphase for $1 \leq i \leq k - 2$, p_k for $i = k - 1$, and p_{k-1} for $i = k$. Hence, an optimal off-line algorithm faults $k + 1$ times on the initial $2k + 1$ requests and then exactly once per phase, while RLRU faults on all $k + 1$ requests of each phase. \square

When LRU and RLRU are compared to OPT using the worst order ratio, instead of the competitive ratio, one finds that they have the same ratio k , confirming the intuition that information is lost in an indirect comparison to OPT.

Theorem 3. $WR_{LRU} = WR_{RLRU} = k$.

Proof. Consider any sequence I . Since no algorithm is better than OPT, on OPT's worst ordering of I , LRU will fault at least as many times as OPT, so it also will on its own worst ordering. The sequence consisting of n copies of $k + 1$ pages repeated cyclicly is a worst ordering of the underlying multi-set for both LRU and OPT. LRU faults k times for every time that OPT faults. Since the worst order ratio cannot be larger than the competitive ratio, and LRU's competitive ratio is k , $WR_{LRU} = k$.

Consider any sequence I . As above, on OPT's worst ordering of I , RLRU will fault at least as many times as OPT, so it also will on its own worst ordering. By Lemma 2, for any sequence I , $RLRU_W(I) \leq LRU_W(I)$. Thus, since $WR_{LRU} = k$, $WR_{RLRU} \leq k$. The sequence $\langle p_1, p_2, \dots, p_{k+1} \rangle^2 \langle p_k, p_1, p_2, \dots, p_{k-1}, p_{k+1} \mid p_{k-1}, p_k, p_1, p_2, \dots, p_{k-2}, p_{k+1} \mid p_{k-2}, p_{k-1}, p_k, p_1, p_2, \dots, p_{k-3}, p_{k+1} \mid \dots \mid p_1, p_2, \dots, p_{k+1} \rangle^n$, where \mid marks the beginning of a new phase, will cause RLRU to fault every time. A worst ordering for OPT will repeat the $k + 1$ pages in a cyclic manner and OPT will fault once on every k pages, giving the ratio k . \square

Implementation RLRU decides whether or not to mark pages based on whether or not they are in LFD's cache. At any given point in time, it is of course impossible to compute the entire contents of LFD's cache, since this depends on future requests. It is, however, possible, given a request and the request sequence up to that point, to compute whether or not LFD would have that particular page in cache.

RLRU can be implemented to run in time $O(\log N)$ and space $O(N)$, where N is the number of different pages requested. The question is whether or not these time and space bounds are good enough in practice. We believe there are at least two interesting scenarios to consider. One is the interaction between two high speed storage media, the speed of which differ by only a small multiplicative constant, such as primary versus secondary cache. Here, a paging algorithm must be very efficient, which also implies that it cannot be allowed much working space. In such a scenario, even LRU is most often too time and space consuming. Another scenario is the interaction of storage media, the speed of which

differ by orders of magnitude. This could be the buffer pool versus the disk in database systems or local file caching of Internet files. In those situations, we can use substantial space, and time logarithmic in either the number of different pages or just in the cache size would both be insignificant compared with almost any small improvement in cache behavior. A similar point is made in [15]. If, in some special application, space is a problem, then it could possibly be reduced to a function of k using the techniques of [1]. In summary, a comparison between LRU and RLRU is interesting because the circumstances under which they can reasonably be applied are close to identical.

Empirical Analysis To see if the positive theoretical results are also reflected in practice, we have investigated the behavior of LRU and RLRU on traces¹ collected from very different applications, including key words searches in text files, selections and joins in the Postgres database system, external sorting, and various kernel operations. We have used all ten data files from the site.

In Table 2, we list the results for each data file, and for cache sizes of 8, 16, 32, and 64. Each entry in the table shows first the number of page faults of LRU, then the number of faults of RLRU, and finally the percent-wise improvement of RLRU over LRU. Thus, if ℓ and r denote the number of faults by LRU and RLRU, respectively, then the improvement is computed as $100\frac{\ell-r}{\ell}$. This number is negative if LRU performs best.

Cache Size	File names and lengths									
	bigsort	j1	j2	j3	j4	j5	j6	pjoin	pq7	xds
	40167	18533	25881	38112	59744	95723	20709	41558	32989	88558
8	14632	494	8233	4262	7278	25412	5100	8014	9371	10768
	13204	491	8197	4276	7278	25385	4545	7200	9419	10724
	9.76	0.61	0.44	-0.33	0.00	0.11	10.88	10.16	-0.51	0.41
16	12619	470	8177	4243	7201	25332	4596	7718	9277	10762
	10736	468	8134	4255	7221	25326	4525	7003	9259	10709
	14.92	0.43	0.53	-0.28	-0.28	0.02	1.54	9.26	0.19	0.49
32	10744	463	8138	4239	7180	25307	4516	7401	9216	10756
	10561	425	8078	4248	7186	25303	4513	6888	9170	10697
	1.70	8.21	0.74	-0.21	-0.08	0.02	0.07	6.93	0.50	0.55
64	10587	136	8120	4230	7135	25276	4505	6879	9185	10754
	10402	137	8057	4239	7140	25278	4506	6838	9103	10695
	1.75	-0.74	0.78	-0.21	-0.07	-0.01	-0.02	0.60	0.89	0.55

Table 2. Empirical comparison of LRU and RLRU

Out of the 40 tests, 11 are negative. The largest negative result of -0.74% is from a short sequence and is due to a difference of only one page fault. The remaining negative results lie between zero and approximately half a per cent. There are 28 positive results. RLRU beats LRU with more than half a per cent in 18 cases, more than 1% in 10 cases, and more than 5% in 7 cases. This is illustrated in Figure 2.

¹ www.cs.wisc.edu/~cao/traces/

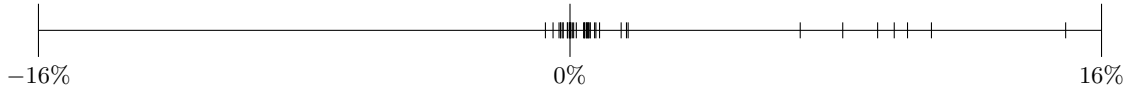


Fig. 2. Percentages with which RLRU is better than LRU

4 Conservative and Marking Algorithms

It is easy to see that both LRU and FIFO are conservative algorithms [32]: between any two faults on the same page there must be requests to at least k other pages. According to the relative worst order ratio, all conservative algorithms have the same performance:

For conservative algorithms, we first prove a technical lemma for LRU.

Lemma 3. *For any sequence I of page requests, there exists a worst ordering of I for LRU with all faults appearing before all hits.*

Proof. We describe how any permutation I' of I can be transformed, step by step, to a permutation I_{LRU} with all hits appearing at the end of the sequence, without decreasing the number of faults LRU will have on the sequence. Let I' consist of the requests r_1, r_2, \dots, r_n , in that order.

Consider the first hit r_i in I' with respect to LRU. We construct a new sequence I'' by moving r_i later in I' . Let p denote the page requested by r_i .

First, we remove r_i from the sequence. If p is evicted at some point after r_{i-1} in this shorter sequence, and is not evicted at the same point in I' , r_i is placed just after the first request r_j , $j > i$, causing p to be evicted (see Figure 3). Otherwise, r_i is inserted after r_n . In this case, let $j = n$.

$$\begin{aligned} I' &: r_1, \dots, r_{i-1}, r_i, r_{i+1}, \dots, r_j, r_{j+1}, \dots, r_n \\ I'' &: r_1, \dots, r_{i-1}, r_{i+1}, \dots, r_j, r_i, r_{j+1}, \dots, r_n \end{aligned}$$

Fig. 3. The two sequences I' and I'' in the case where p is evicted at r_j .

LRU maintains a queue of the pages in cache, and, on a fault, evicts the first page in the queue. Moving r_i within the sequence affects the position of p in the queue, but the mutual order of the other pages stays the same. Just before r_{i+1} , the cache contents are the same for both sequences. Therefore, for I'' , the behavior of LRU is the same as for I' until p is evicted. Just after this eviction, p is requested by r_i in I'' . Thus, just before r_{j+1} , the cache contents are again the same for both sequences, but for I'' , p is at the end of queue. This means that all pages that are in cache just before r_{j+1} , except p , are evicted no later for I'' than for I' . The first request to p after the j th request may be a fault in I' and a hit in I'' . On the other hand, r_i is a hit in I' and a fault in I'' .

Let r_ℓ be the first request after r_i in I'' , where p is either requested or evicted. After r_ℓ , the state of LRU is the same for both sequences.

By moving r_i , the number of faults among the first j requests is increased by at least one, and the total number of faults is not decreased. Thus, continuing in this way, we obtain I_{LRU} in a finite number of steps. \square

LRU is best possible among conservative algorithms.

Lemma 4. $WR_{\mathbb{C}, \text{LRU}} \geq 1$, for any conservative paging algorithm \mathbb{C} .

Proof. By Lemma 3, we can consider a sequence I where all faults by LRU occur before all hits. Let I_1 denote the subsequence consisting of the faults. We prove by induction on the lengths of prefixes of

I_1 that, on any request in I_1 , any conservative algorithm \mathbb{C} evicts the same page as LRU, and hence has as many faults on I as LRU.

For the base case, consider the first $k + 1$ requests in the sequence. Since LRU faults on each request, these $k + 1$ requests are all to different pages (ignoring the trivial case with at most k pages in I). Hence, on the $(k + 1)$ st request, any algorithm must evict a page. Since \mathbb{C} is conservative it evicts p_1 (if it evicted some page $p_i \neq p_1$, requesting p_i after p_{k+1} would yield a sequence with a subsequence $\langle p_2, \dots, p_{k+1}, p_i \rangle$ with requests to only k distinct pages, but with $k + 1$ faults).

The induction step is similar to the base case. By the induction hypothesis, \mathbb{C} has the same pages in cache as LRU. For each request r to some page p , the previous k requests were all to different pages different from p . Hence, \mathbb{C} must evict the first of these k pages, as LRU does. \square

In addition, LRU is a worst possible conservative algorithm.

Lemma 5. $WR_{LRU, \mathbb{C}} \geq 1$, for any conservative paging algorithm \mathbb{C} .

Proof. Consider any conservative algorithm \mathbb{C} and any request sequence I . Divide I into phases, so that \mathbb{C} faults exactly $k + 1$ times per phase, starting the next phase with a fault (the last phase may have fewer than $k + 1$ faults). Since \mathbb{C} is conservative, each phase, except possibly the last, contains at least $k + 1$ distinct pages. These can be reordered, phase by phase, so that LRU faults once on each distinct page within a phase. Thus, with this new ordering LRU faults at least as many times as \mathbb{C} on I , except possibly in the last phase. Since \mathbb{C} faults at least once in the last phase, LRU faults as many times on the new ordering of I as \mathbb{C} on I , except for at most k faults. \square

Thus, all conservative algorithms are equivalent under the relative worst order ratio:

Theorem 4. For any pair of conservative algorithms \mathbb{C}_1 and \mathbb{C}_2 , $WR_{\mathbb{C}_1, \mathbb{C}_2} = 1$.

Proof. By Lemmas 4 and 5, $WR_{\mathbb{C}_1, LRU} \geq 1$ and $WR_{LRU, \mathbb{C}_2} \geq 1$. Since the relative worst order ratio is easily seen to be a transitive measure [8], this means that $WR_{\mathbb{C}_1, \mathbb{C}_2} \geq 1$. Similarly, $WR_{\mathbb{C}_1, LRU} \leq 1$ and $WR_{LRU, \mathbb{C}_2} \leq 1$, which implies $WR_{\mathbb{C}_1, \mathbb{C}_2} \leq 1$. \square

In particular, since LRU and FIFO are both conservative, they are equivalent.

Corollary 1. $WR_{LRU, FIFO} = 1$.

By Theorems 1 and 4 and the transitivity of the relative worst order ratio, we have the following:

Corollary 2. For any conservative algorithm \mathbb{C} , $WR_{\mathbb{C}, RLRU} = \frac{k+1}{2}$.

Similarly, the results below, showing that LRU is at least as good as or better than some algorithm \mathbb{A} , also show that RLRU is better than \mathbb{A} . In addition, any result holding for LRU also holds for any conservative algorithm. These results are not all stated explicitly.

Although all pairs of conservative algorithms are indistinguishable using the relative worst order ratio, not all pairs of marking algorithms are. In particular, LRU is better than FWF, as it is in practice, despite the fact that they have the same competitive ratio.

Lemma 6. $WR_{FWF, LRU} \geq \frac{2k}{k+1}$.

Proof. Consider any sequence I . If LRU faults on request r in I to page p , then p was not among the last k different pages that were requested. Thus, p could not be in FWF's cache when request r occurs and FWF will also fault. Hence, on any sequence, FWF will fault at least as many times on its worst ordering as LRU will on its. This shows that $WR_{FWF, LRU} \geq 1$.

It is now sufficient to find a family of sequences I_n with $\lim_{n \rightarrow \infty} \text{LRU}(I_n) = \infty$, where there exists a constant b such that for all I_n ,

$$\text{FWF}_W(I_n) \geq \frac{2k}{k+1} \text{LRU}_W(I_n) - b.$$

Let $I_n = \langle p_1, p_2, \dots, p_k, p_{k+1}, p_k, \dots, p_2 \rangle^n$. FWF will fault on every page, so it will fault $n(2k)$ times.

By Lemma 3, a worst ordering for LRU consists of some permutation of the $k+1$ pages, with p_1 and p_{k+1} at the end, repeated n times, followed by the first $k-1$ pages in the permutation repeated n times, i.e., $\langle p_2, p_3, \dots, p_{k+1}, p_1 \rangle^n \langle p_2, p_3, \dots, p_k \rangle^n$ is a worst ordering of I_n with respect to LRU. LRU will fault $n(k+1) + k - 1$ times. Asymptotically, the ratio is $\frac{2k}{k+1}$. \square

FWF is worst possible among marking algorithms:

Lemma 7. *For any marking algorithm \mathbb{M} and any request sequence I , $\mathbb{M}_W(I) \leq \frac{2k}{k+1} \text{LRU}_W(I) + k$.*

Proof. For any sequence with n complete phases, \mathbb{M} faults at most kn times. Since for any request sequence, every pair of consecutive phases contains requests to at least $k+1$ distinct pages, for the first $2 \cdot \lfloor \frac{n}{2} \rfloor$ complete phases, there is a permutation such that LRU faults at least $(k+1) \lfloor \frac{n}{2} \rfloor$ times. If the remaining requests consist of a complete phase, plus a partial phase, then LRU will also fault on $k+1$ of those requests if given in the correct order. Thus, the additive constant is bounded by the number of faults \mathbb{M} makes on the last partial or complete phase and is thus at most k . \square

Combining the above two lemmas, Theorem 4, and the the transitivity of the relative worst order ratio gives:

Theorem 5. *For any conservative algorithm \mathbb{C} , $WR_{\text{FWF}, \mathbb{C}} = \frac{2k}{k+1}$.*

Furthermore, LRU is close to being a best possible marking algorithm:

Lemma 8. *For any marking algorithm \mathbb{M} and any sequence I of page requests,*

$$\text{LRU}_W(I) \leq \frac{k+1}{k} \mathbb{M}_W(I).$$

Proof. Consider any sequence I of requests. By Lemma 3, there is a worst ordering I_{LRU} such that LRU faults on each request of a prefix I_1 of I_{LRU} and on no request after I_1 . Partition I_1 into consecutive subsequences, each consisting of exactly $k+1$ requests (the last subsequence may contain fewer). Since LRU faults on all requests in I_1 , each subsequence, except possibly the last, must contain $k+1$ distinct pages. Hence, for each subsequence with pages p_1, p_2, \dots, p_{k+1} , an adversary can create a marking phase, by choosing k of the pages p_1, p_2, \dots, p_{k+1} , such that the marking algorithm faults on all k pages. This is easily seen in the following way. Pages requested within a phase stay in cache throughout the phase. Therefore, when x of the pages p_1, p_2, \dots, p_{k+1} have been requested, the remaining $k+1-x$ pages cannot all be in the cache. \square

This immediately gives the following:

Corollary 3. *For any marking algorithm \mathbb{M} with $WR_{\text{LRU}, \mathbb{M}}$ defined, $WR_{\text{LRU}, \mathbb{M}} \leq \frac{k+1}{k}$.*

Let $\text{MARK}^{\text{LIFO}}$ denote the marking algorithm which, on a fault, evicts the unmarked page that was most recently brought into cache. On some sequences, $\text{MARK}^{\text{LIFO}}$ is slightly better than LRU, but on others, LRU is about twice as good as $\text{MARK}^{\text{LIFO}}$:

The following lemma shows that Lemma 8 is tight. It does not settle, however, whether Corollary 3 is tight.

Lemma 9. *There exists a family of sequences I_n of page requests such that*

$$LRU_W(I_n) = \frac{k+1}{k} MARK_W^{LIFO}(I_n),$$

and $\lim_{n \rightarrow \infty} LRU_W(I_n) = \infty$.

Proof. First note that, for any sequence I of page requests, there is a worst ordering of I with respect to $MARK^{LIFO}$ such that all faults precede all hits. This is because the eviction strategy considers only the order in which the pages were brought into the cache. Therefore, moving a hit to the end of the sequence does not affect which of the other requests will be faults, and hence can only increase the number of faults.

Consider the sequence $\langle p_1, p_2, \dots, p_{k+1} \rangle^n$ and a permutation for which $MARK^{LIFO}$ faults on a longest possible prefix and then has no more faults. Since there are only $k+1$ pages, once the first k requests are given, the remaining part of the prefix is fixed. Hence, there is essentially only one such permutation, namely $\langle p_1, p_2, \dots, p_k, p_{k+1}, p_k, \dots, p_2 \rangle^{n/2} \langle p_1, p_{k+1} \rangle^{n/2}$. $MARK^{LIFO}$ does not fault on the last $n-1$ requests, whereas LRU will fault on all requests in the permutation $\langle p_1, p_2, \dots, p_{k+1} \rangle^n$. \square

Despite Lemma 9, $MARK^{LIFO}$ is not better than LRU:

Lemma 10. *There exists a family of sequences I_n of page requests and a constant b such that*

$$MARK_W^{LIFO}(I_n) = \frac{2k}{k+1} \cdot LRU_W(I_n) - b,$$

and $\lim_{n \rightarrow \infty} MARK_W^{LIFO}(I_n) = \infty$.

Proof. On the sequence $I_n = \langle p_1, p_2, \dots, p_k, p_{k+1}, p_k, \dots, p_2 \rangle^n$, which was also used in the proof of Lemma 6, $MARK^{LIFO}$ faults on every request. As explained in the proof of Lemma 6, the sequence $\langle p_2, p_3, \dots, p_{k+1}, p_1 \rangle^n \langle p_2, p_3, \dots, p_k \rangle^n$ is a worst ordering of I_n with respect to LRU. Hence, LRU faults at most $(k+1)n + k - 1$ times, i.e.,

$$MARK_W^{LIFO}(I_n) = 2kn = \frac{2k}{k+1}(k+1)n = \frac{2k}{k+1}LRU_W(I_n) - \frac{2k(k-1)}{k+1}.$$

\square

This gives the following:

Theorem 6. *LRU and $MARK^{LIFO}$ are $(1 + \frac{1}{k}, c)$ -related with $c \geq 2 - \frac{1}{k}$, i.e., they are asymptotically comparable in LRU's favor.*

5 Look-Ahead

In the standard on-line model, requests arrive one by one. A model in which the algorithm is informed of the next $\ell \geq 1$ page requests before servicing the current one, is a look-ahead model. This model is in-between the standard on-line model and the off-line model.

It is well known that using standard competitive analysis one cannot show that knowing the next ℓ requests is any advantage for any fixed ℓ ; for any input sequence, an adversary can “fill up” the look-ahead by using $\ell+1$ consecutive copies of each request, adding no cost to the optimal off-line solution. In contrast, results on the relative worst order ratio, indicate that look-ahead helps significantly. Here we only look at a modification of LRU, using look-ahead, though the technique can be applied to other algorithms as well.

Define $\text{LRU}(\ell)$ to be the algorithm which on a fault evicts the least recently used page in cache which is not among the next ℓ requests. If $\ell \geq k$, all pages in cache may be among the next ℓ requests. In this case, the page whose next request is farthest in the future is evicted.

One can see that $\text{LRU}(\ell)$ is at least as good as LRU on any sequence by noting that $\text{LRU}(\ell)$ is conservative.

Lemma 11. *$\text{LRU}(\ell)$ is a conservative algorithm.*

Proof. Let I be a request sequence, and assume that there is an interval, I' , in I , containing only k distinct pages, on which $\text{LRU}(\ell)$ faults at least $k + 1$ times. Then it must fault on some page, p , twice in I' . Between these two faults, say at request r , page p must be evicted. First assume that $\ell < k$. At this point, p is the least recently used page which is not among the next ℓ . Clearly the second request causing a fault on p must be beyond these next ℓ . So the other $k - 1$ pages in cache, when request r occurs, must all have been requested between these two faults on p . In addition, the request r cannot be for p or any of the other pages in cache at that time. Thus, there must be at least $k + 1$ distinct pages in I' , giving a contradiction. Now assume that $\ell \geq k$. If p is not among the next ℓ requests when r occurs, then the previous argument holds, so assume that it is. In this case p must have been the page in cache which was requested furthest in the future, so the other $k - 1$ pages are requested between request r and the second fault on p . Again, counting the request r and p , there must be at least $k + 1$ distinct pages in I' , which is a contradiction. Thus, $\text{LRU}(\ell)$ is conservative. \square

Observe that Lemma 5 holds for algorithms using look-ahead, though Lemma 4 does not.

Theorem 7. $WR_{\text{LRU}, \text{LRU}(\ell)} = \min\{k, \ell + 1\}$.

Proof. Since the previous lemma combined with Lemma 5 show that $WR_{\text{LRU}, \text{LRU}(\ell)} \geq 1$, to prove the lower bound, it is sufficient to find a family of sequences I_n with $\lim_{n \rightarrow \infty} WR_{\text{LRU}, \text{LRU}(\ell)}(I_n) = \infty$, where there exists a constant b such that for all I_n ,

$$\text{LRU}_W(I_n) \geq \min\{k, \ell + 1\} \text{LRU}(\ell)_W(I_n) - b.$$

Let I_n consist of n phases, each containing the pages $p_1, p_2, \dots, p_k, p_{k+1}$, in that order. LRU will fault $n(k + 1)$ times. However, if $\ell \leq k - 1$, after the first k faults, $\text{LRU}(\ell)$ never faults on any of the next ℓ pages after a fault. Thus, regardless of the order, $\text{LRU}(\ell)$ faults on at most $k + \lfloor \frac{n(k+1)-k}{\ell+1} \rfloor$ pages. Asymptotically, this gives a ratio of $\ell + 1$. If $\ell \geq k$, then $\text{LRU}(\ell)$ faults on at most one out of every k pages.

For the upper bound, suppose there exists a sequence I , where LRU faults s times on its worst permutation, I_{LRU} , $\text{LRU}(\ell)$ faults s' times on its worst ordering, $I_{\text{LRU}(\ell)}$, and $s > \min(k, \ell + 1) \cdot s'$. Then, $s > \min(k, \ell + 1) \cdot s''$, where s'' is the number of times $\text{LRU}(\ell)$ faults on I_{LRU} . One cannot have $\ell \geq k$, since then $\text{LRU}(\ell)$ faults fewer times than OPT. So suppose $\ell < k$, and assume by Lemma 3, that I_{LRU} is such that LRU faults on each request of a prefix I_1 of I_{LRU} and on no request after I_1 . Then there must exist a request r in I_{LRU} where $\text{LRU}(\ell)$ faults, but it does not fault on any of the next $\ell + 1$ requests, all of which are in I_1 . The last of these $\ell + 1$ requests caused LRU to fault, so it was not among the last k distinct requests at that point. Since $\ell < k$, it was not in any of the requests in the look-ahead when $\text{LRU}(\ell)$ processed request r , and all of the pages in the look-ahead were in cache then since $\text{LRU}(\ell)$ did not fault on any of them. Hence, this $\ell + 1$ st page was evicted by $\text{LRU}(\ell)$ when r was requested, and there must have been a fault the next time it was requested after that, giving a contradiction. \square

6 Other Algorithms

The algorithm LIFO, which evicts the page most recently brought into cache, is clearly much worse than any of the conservative or marking algorithms. This behavior is also reflected in the relative worst order ratio, since there exists a family of sequences where LIFO does unboundedly worse than LRU. However, there also exists a family of sequences, where LIFO does better than LRU by a factor of $\frac{k+1}{2}$. This factor is tight, though, so LRU can be unboundedly better than LIFO, while LIFO can be at most a factor $\frac{k+1}{2}$ better than LRU.

Theorem 8. *LRU and LIFO are $(\frac{k+1}{2}, \infty)$ -related, i.e., they are weakly comparable in LRU's favor.*

Proof. Let I_n be the sequence $\langle p_1, p_2, \dots, p_{k-1} \rangle \langle p_k, p_{k+1} \rangle^n$. LIFO faults on every request of this sequence. LRU's worst ordering is with $\langle p_k, p_{k+1} \rangle$ first and then the pages $\langle p_1, p_2, \dots, p_k, p_{k+1} \rangle$, since then it faults $k + 3$ times. Hence, for this sequence, the ratio is unbounded in LRU's favor.

On the other hand, with the sequence $J_n = \langle p_1, p_2, \dots, p_k, p_{k+1} \rangle^n$, LRU faults $n(k + 1)$ times, so this is clearly its worst order of J_n . Regardless of the order of J_n , LIFO faults exactly $2n + k - 1$ times, since it holds $k - 1$ pages in memory, never changing them. Thus, there is a ratio of $\frac{k+1}{2}$ in LIFO's favor.

Since LRU has a competitive ratio of k it cannot be more than a factor of k worse than LIFO. We can show, however, that the lower bound of $\frac{k+1}{2}$ is tight. Let I be any request sequence. By Lemma 3, there is a worst ordering I_{LRU} such that LRU faults on each request of a prefix I_1 of I_{LRU} and on no request after I_1 . Consider the prefix I_1 , divided into phases with $k + 1$ pages per phase, except possibly fewer in the last phase. We reorder the requests within each complete phase, starting at the beginning. The contents of LIFO's cache at the beginning of a phase will be with respect to the sequence I_1 with all modifications up to that point. Since LRU faults on every page in a phase in the original I_1 , the $k + 1$ pages must all be different. LIFO has at most k of the $k + 1$ pages from the current phase in its cache at the start of the phase. If there were only $k - 1$, instead of k , we would be done, since then it must fault at least twice in this phase. Suppose p_i is the page requested in the current phase, which was not in LIFO's cache. Move the request to p_i to the beginning of the phase. This will evict a page requested later in the phase, so LIFO will fault at least twice in this phase. \square

The algorithm MIXPERM [6] is a mixed, randomized, memory-bounded algorithm which obtains a competitive ratio of $\frac{k+1}{2}$ when $N = k + 1$. It is the uniform mixture of all the permutation algorithms $PERM_\pi$ [11]², which have competitive ratio k . The parameter π in the definition of $PERM_\pi$ is a cyclic permutation of the N pages in slow memory. Let $\pi^{(m)}$ denote the m -fold composition of the cyclic permutation π , and let $m(i)$ be the minimum m with $\pi^{(m)}(i)$ currently in cache. The algorithm $PERM_\pi$ is defined so that on a page fault on page i , the page $\pi^{(m(i))}(i)$ is evicted. $PERM_\pi$ and LRU are equivalent according to the competitive ratio, and when $N = k + 1$, they are equivalent according to the relative worst order ratio. However, for $N \geq k + 2$, $PERM_\pi$ performs worse than LRU according to the relative worst order ratio.

Theorem 9. *If $N = k + 1$, $WR_{PERM_\pi, LRU} = 1$. If $N \geq k + 2$, $WR_{PERM_\pi, LRU} > 2 - \frac{k-1}{N} + \frac{2}{k}$.*

Proof. We first show that $WR_{PERM_\pi, LRU} \geq 1$. Consider any request sequence I . By Lemma 3, there is a worst ordering I_{LRU} such that LRU faults on each request of a prefix I_1 of I_{LRU} and on no request after I_1 . Consider the prefix I_1 , divided into phases with $k + 1$ requests per phase, except possibly fewer in the last phase. We reorder the requests within each complete phase, starting at the beginning, such that $PERM_\pi$ will also fault on all of these requests. The contents of $PERM_\pi$'s cache at the beginning of a phase will be with respect to the sequence I_1 with all modifications up to that point. Since LRU faults on every page in a phase in the original I_1 , the $k + 1$ pages must all be different.

² Chrobak et al. called this algorithm ROTATE.

Consider the current phase being processed, phase P_i . Arrange the pages in P_i in a cycle C , in the same order in which they occur in π . We will find a starting point in C so that requesting the pages in the order given by C , starting at this point, will cause PERM_π to fault on every page.

Arrange the pages in PERM_π 's cache at the beginning of this phase in a cycle, C' , also in the order in which they occur in π . If the cycles C and C' have no pages in common, none of the pages in phase P_i are currently in cache, so PERM_π will fault on every request. Thus, we assume that they have some pages in common. Since both cycles are ordered according to π the pages they have in common occur in the same order in both cycles. Match up the identical pages on the two cycles, and count on each cycle the number of pages strictly between consecutive matched pages. Create a cycle C'' consisting of integers representing the difference between the number of pages in C' and the number of pages in C between consecutive matched pages. Since C has one more page than C' , the sum of the numbers in C'' is -1 .

Claim. There exists a starting point on C'' such that, in the string S defined by starting at this point on C'' and including the remainder of C'' in order, the sums of all prefixes are negative.

Proof. Suppose there is no such starting point. Take any starting point and follow a prefix until a nonnegative sum occurs. Continue like this, using the first value not in the previous prefix as the next starting point and stopping when a nonnegative sum has been found. This process can be repeated indefinitely. Since C'' is finite, eventually two starting points will be the same. Between these two starting points, one has been around C'' an integer number of times with the sum of all values adding to some nonnegative integer value. This is a contradiction since the sum of the numbers in C'' is -1 . \square

Choose a starting point on C'' where all prefix sums are negative. This corresponds to a subsequence of at least one page request in C which is not in C' and thus not in cache. This will be the starting point in C . Each request to a page in C , but not in C' , will cause the eviction of a page in C' and these evictions will occur in the order given by C' . When a page p in both C and C' is requested, the number of pages before p in C is greater than the number before p in C' , so p will have been evicted and will also cause a fault. Thus PERM_π will also fault $k + 1$ times on this set of $k + 1$ pages.

This occurs for every complete phase. Thus, LRU faults at most an additive constant k times more on its worst ordering than what PERM_π does on its worst ordering, on the pages in the incomplete phase, if it exists.

We now prove that for $N = k + 1$, $\text{WR}_{\text{PERM}_\pi, \text{LRU}} \leq 1$. This completes the proof that, for $N = k + 1$, $\text{WR}_{\text{PERM}_\pi, \text{LRU}} = 1$. Consider any sequence I with requests to at most $k + 1$ pages and let I_{PERM_π} be a worst ordering of I with respect to PERM_π . Whenever PERM_π faults on a page p , it goes through π starting at p until it finds a page which is in cache and evicts this page. Since $N = k + 1$, this will be the page p' immediately succeeding p in π , and the next fault will be on p' . Thus, if I_{PERM_π} is reordered such that the pages that PERM_π faults on occur at the beginning of the sequence, in the same mutual order as in I_{PERM_π} , then LRU will also fault on each of these requests.

Now, consider the case $N \geq k + 2$. To see that $\text{WR}_{\text{PERM}_\pi, \text{LRU}} > 2 - \frac{k-1}{N} + \frac{2}{k}$, assume without loss of generality that $\pi = (1, 2, \dots, N)$ and consider the family of request sequences

$$I_n = \langle p_1, p_2, \dots, p_{k+1}, p_1, p_{k+2}, p_1, p_{k+3}, \dots, p_1, p_N \rangle^n,$$

where each subsequence $\langle p_1, p_2, \dots, p_{k+1}, p_1, p_{k+2}, p_1, p_{k+3}, \dots, p_1, p_N \rangle$ is called a phase. We show that PERM_π faults on all requests in I_n . Since each phase has one request to each of the N pages and one extra request to p_1 for each of the $N - (k + 1)$ pages $p_{k+2}, p_{k+3}, \dots, p_N$, this implies

$$(\text{PERM}_\pi)_W(I_n) = (N + N - (k + 1))n = (2N - (k + 1))n.$$

Consider the first phase of I_n . We prove by induction that just before the i th request to p_1 , $2 \leq i \leq N - k$, the cache contains p_i, \dots, p_{i+k-1} . For $i = 2$, this is clearly the case. Therefore, consider the i th request r_i to p_1 and assume that the cache contains the pages p_i, \dots, p_{i+k-1} . Since p_1 is not in the cache, r_i will cause p_i to be evicted. After r_i , p_{k+i} is requested, and this causes p_1 to be evicted. Thus, the induction hypothesis is maintained. At the end of the first phase, the pages p_{N-k+1}, \dots, p_N are in cache. Therefore, for the following k requests, the request to p_i will cause p_{N-k+i} to be evicted. Thus, after the first k requests of the second phase of I_n , the cache contents are the same as after the first k requests of the first phase. Since the algorithm is memoryless, this is sufficient to prove that it will behave the same during all n phases of I_n .

Now consider LRU. Between any pair of faults on p_1 , at least k of the $N - 1$ other pages are requested. Thus, regardless of the ordering, LRU faults on at most $\lfloor \frac{N-1}{k} n \rfloor + 1$ requests to p_1 , i.e.

$$\text{LRU}_W(I_n) \leq \left\lfloor \frac{N-1}{k} n \right\rfloor + 1 + (N-1)n \leq \left(\frac{N-1}{k} + \frac{1}{n} + N-1 \right) n \leq \left(\frac{N}{k} + N-1 \right) n, \text{ for } n \geq k.$$

This gives a ratio of

$$\begin{aligned} \frac{(\text{PERM}_\pi)_W(I_n)}{\text{LRU}_W(I_n)} &\geq \frac{2N - (k+1)}{N-1 + \frac{N}{k}} = \frac{2N - 2 + \frac{2N}{k} + 2 - \frac{2N}{k} - (k+1)}{N-1 + \frac{N}{k}} \geq 2 - \frac{k-1 + \frac{2N}{k}}{N-1 + \frac{N}{k}} \\ &\geq 2 - \frac{k-1}{N} + \frac{2}{k}, \text{ since } N > k \end{aligned}$$

□

7 Randomized Algorithms

The relative worst order ratio can also be applied to randomized algorithms. The only change to the definitions is that an algorithm's expected value/cost on a worst ordering of a sequence is used in place of the value/cost obtained by a deterministic algorithm.

Definition 7. Consider an on-line optimization problem P and let I be any request sequence of length n . If σ is a permutation on n elements, then $\sigma(I)$ denotes I permuted by σ . Let \mathbb{A} be any randomized algorithm for P . If P is a maximization problem, $E[\mathbb{A}(I)]$ is the expected value of running \mathbb{A} on I , and $\mathbb{A}_W(I) = \min_\sigma E[\mathbb{A}(\sigma(I))]$. If P is a minimization problem, $\mathbb{A}(I)$ is a cost, and $\mathbb{A}_W(I) = \max_\sigma E[\mathbb{A}(\sigma(I))]$.

Using the above definition, the relative worst order ratio, related, weakly comparable, asymptotically comparable, and incomparable are now defined as in the deterministic case.

Consider the randomized paging algorithm MARK [14]. On a fault, MARK chooses the unmarked page to be evicted uniformly at random. In this section, we show that MARK and LRU are comparable and $\text{WR}_{\text{LRU}, \text{MARK}} \geq k/H_k$. This is consistent with the results one obtains with the competitive ratio where MARK has ratio $2H_k - 1$ [1], while LRU has ratio k .

In each phase of MARK's processing (except possibly the last), exactly k distinct pages are requested, and the first page requested within a phase was not requested in the previous phase. Thus, the subsequence processed within a phase (except possibly the last) is a maximal subsequence containing requests to exactly k distinct pages. A subsequence processed within one marking phase is called a k -phase. Note that the partitioning of a sequence into k -phases is independent of the particular marking algorithm.

For Lemma 12 and Theorem 10, we need the fact that MARK's expected number of faults in the i th k -phase is $m_i(H_k - H_{m_i} + 1)$ [14], where m_i is the number of *new* pages in the i th phase, i.e., the number of pages that are requested in the i th phase and not in the $(i-1)$ st phase.

Lemma 12. *There exists a sequence I , which is a worst ordering for both MARK and LRU, where MARK's expected number of faults is H_k per k -phase, while LRU faults k times per k -phase.*

Proof. Consider the sequence with $k + 1$ pages repeated cyclicly. \square

Lemma 13. *For any sequence I of page requests, there exists a worst ordering I_{MARK} of I with respect to MARK, such that all k -phases, except possibly the last, have the following properties:*

1. *The first page requested in a k -phase did not appear in the previous k -phase.*
2. *There are exactly k requests, all to distinct pages.*

Proof. Consider a worst ordering I_{MARK} of I for MARK and consider its k -phase partition. The first property follows by the definition of a k -phase partition. Within a phase, the first occurrence of each page is the only one MARK has any chance of faulting on. Thus, moving extra occurrences of a page within a phase to the end of the sequence will never decrease the probability of MARK faulting on any page. After this has been completed, each phase (except possibly the last) consists of exactly k requests, all to distinct pages. \square

Lemma 14. *For any request sequence I , $\text{LRU}(I_{\text{LRU}}) \geq \text{MARK}(I_{\text{MARK}}) - b$, where b is a constant independent of the sequence I .*

Proof. For any request sequence I , consider the worst permutation I_{MARK} of I with respect to MARK, satisfying the conditions of the previous lemma. Partition the sequence I_{MARK} in blocks, each consisting of variable numbers of consecutive k -phases satisfying the following property:

Definition 8. *Call pages requested in every phase of a block S frequent pages, and the others rare pages. A block S , containing s k -phases, has the sizing property if each k -phase in S contains $s - 1$ rare pages, and there is no $r < s$ such that the first block R consisting of the first r pages of S contains $r - 1$ rare pages.*

Each block in the partition will be analyzed separately, and it will be shown that each block can be reordered so that LRU faults at least as many times as the expected number of faults by MARK on that block.

Consider a block, S , containing s rare pages and thus $s + 1$ k -phases. Note that $s \geq 1$, since a new k -phase is started by a new page, which was not among the last k distinct pages.

If the first k -phase, P , in the block has at most s new pages, then MARK's expected number of faults is at most $s(s + 1)(1 + H_k - H_s)$. Since the rare pages occur at most s times in this block, one can reorder the block into s groups of $k + 1$ distinct pages, plus $(s + 1)k - s(k + 1) = k - s$ extra pages. Thus, LRU can be forced to fault $s(k + 1)$ times. Since MARK's expected number of faults on this block is $s(s + 1)(1 + H_k - H_s) \leq s(s + 1)(1 + \frac{k-s}{s+1}) = s(k + 1)$, the result holds if P has at most s new pages.

Now assume that the first k -phase, P , in the block, S , has at most $s + i$ new pages, where $i > 0$. Then, some frequent page in P is also a new page. MARK's expected number of faults is at most $s^2(1 + H_k - H_s) + (s + i)(1 + H_k - H_{s+i})$.

Let S' be the previous block. Since S' satisfies the sizing property, there are $|S'| - 1$ rare pages in each k -phase of S' . Consider any frequent page, p , in P which is also new. It is clearly not a frequent page in S' . If it occurs in all but the last k -phase of S' , then the block consisting of the first $|S'| - 1$ k -phases of S' has at least one more frequent page than all of S' does. Since S' has the sizing property, this cannot occur. Thus, p occurs at most $|S'| - 2$ times in S' . Hence, one can choose i frequent, new pages in P which can be moved back into the previous block, S' , ordering S' such that LRU faults on these i pages, in addition to the $(|S'| - 1)(k + 1)$ pages originally in S' which it faults on. Removing

these i pages from block S will not decrease the number of pages within block S which LRU faults on. Thus, one can count $s(k+1) + i$ pages from block S which LRU will fault on. The lemma follows by the following claim:

Claim. For all integers s, i, k , such that $1 \leq s \leq k$, $i \geq 0$, and $s + i \leq k$,

$$s^2(1 + H_k - H_s) + (s + i)(1 + H_k - H_{s+i}) \leq s(k + 1) + i.$$

Proof. To show this, one can check that it holds for all $k \leq 5$ by hand.

Assume first that $s > k/2$. Then, $s^2(1 + H_k - H_s) + (s + i)(1 + H_k - H_{s+i}) = s(s+1)(1 + \sum_{j=s+1}^k \frac{1}{j}) + i(1 + \sum_{j=s+i+1}^k \frac{1}{j}) - s(\sum_{j=s+1}^{s+i} \frac{1}{j}) \leq s(s+1)(1 + \frac{k-s}{s+1}) + i(1 + \frac{k-s-i}{s+i+1}) - s \frac{i}{s+i} = s(k+1) + i + (\frac{i(k-s-i)}{s+i+1} - \frac{si}{s+i})$. The last part of this is negative when $s > k/2$, so the result holds for $s > k/2$.

Note that approximating by integrals, one gets that $H_x - H_y \leq \ln(x) - \ln(y)$ for $x > y$. Thus, $s^2(1 + H_k - H_s) + (s + i)(1 + H_k - H_{s+i}) \leq s^2(1 + \ln(k) - \ln(s)) + (s + i)(1 + \ln(k) - \ln(s + i))$, so it is sufficient to prove that $f(s, i, k) = s^2(1 + \ln(k) - \ln(s)) + (s + i)(1 + \ln(k) - \ln(s + i)) - s(k + 1) - i \leq 0$ for $k \geq 6$ and $s \leq k/2$.

Taking the derivative of $f(s, i, k)$ with respect to i gives $f'_i(s, i, k) = \ln(k) - \ln(s + i) - 1$, which is zero at $i = \frac{k-se}{e}$, positive for smaller i and negative for larger. Thus, for fixed s and k , $f(s, i, k)$ has its maximum at $i = \frac{k-se}{e}$. For $s > \frac{k}{e}$, this value of i is negative, and hence outside the specified range. Since $f(s, i, k)$ is decreasing for $i > \frac{k-se}{e}$, for $s \geq \frac{k}{e}$, $f(s, i, k)$ is maximum at its smallest allowable value, $i = 0$.

Consider $s \geq \frac{k}{e}$, so $f(s, 0, k) = s(s+1)(1 + \ln(k) - \ln(s)) - s(k+1)$ is an upper bound on $f(s, i, k)$. The derivative of this with respect to s is $-s - k - 2 + (1 + \ln(k) - \ln(s))(2s + 1)$ which is zero where $(s+k+2)/(2s+1) = (1 + \ln(k) - \ln(s))$. At this point, $f(s, 0, k) = s(s+1)(s+k+2)/(2s+1) - sk - s = \frac{s^3+s^2+s-s^2k}{2s+1}$. This is equal to zero, where $s^2 + s + 1 - sk = 0$, which is not in the range $\frac{k}{e} \leq s \leq \frac{k}{2}$ for $k \geq 4$. At the endpoints of this range $s = k/2$ and $s = k/e$, where $f(k/2, 0, k)$ is negative for $k \geq 6$, and $f(k/e, 0, k)$ is negative for $k \geq 4$. Thus, $f(s, i, k)$ is negative for $s \geq k/e$ and $k \geq 6$.

For $s \leq \frac{k}{e}$, $f(s, i, k)$ has its maximum at $f(s, \frac{k-se}{e}, k) = s^2(1 + \ln(k) - \ln(s)) + 2s + \frac{k-se}{e} - s(k+1)$. Taking the derivative with respect to s gives $s - k + 2s(\ln(k) - \ln(s))$, which is negative for s small enough and then positive, so the function has a local minimum where the derivative is zero. Thus, the maximum values are at the endpoints. For $s = 1$, one gets that $f(1, \frac{k}{e} - 1, k) = 1 + \ln(k) + \frac{k}{e} - k$, which is negative for $k \geq 4$. For $s = \frac{k}{e}$, one gets that $f(\frac{k}{e}, 0, k) = (\frac{2}{e} - 1)sk + s$, which is negative for $k \geq 4$.

Thus, $s^2(1 + H_k - H_s) + (s + i)(1 + H_k - H_{s+i}) \leq s(k + 1) + i$ for $1 \leq s \leq k$, $i \geq 0$, and $s + i \leq k$. \square

This proves the lower bound for the following:

Theorem 10. $WR_{LRU, MARK} = \frac{k}{H_k}$.

Proof. To see the the ratio cannot be higher, consider any k -phase in LRU's worst ordering. LRU never faults more than k times on any k -phase, and MARK never has an expected number of faults less than H_k on any complete k -phase. Thus, the result is tight. \square

8 Conclusion and Open Problems

This second problem area, paging, studied using the relative worst order ratio, gives even more convincing evidence than the first, bin packing, that this new performance measure could become an important tool for analyzing on-line algorithms. Comparing algorithms directly to each other, rather than doing it indirectly through a comparison to OPT, appears to give more meaningful results, both for paging and for bin packing. The results mentioned in the introduction for the bin coloring problem and two scheduling problems show more types of problems where the relative worst order ratio gives

“better” results than the competitive ratio. Previous measures and models, proposed as replacements or supplements to the competitive ratio, have been more limited as to applicability, usually to very few problems. Further study is needed to determine how widely applicable the relative worst order ratio is, but paging and bin packing are very different problems. Together with the results on the bin coloring and scheduling problems, this gives a convincing basis for further investigation.

For paging, the competitive ratio of many algorithms with widely varying performance in practice is k . The relative worst order ratio is able to distinguish between some of these algorithms. Most notably, LRU is found to be better than FWF, and look-ahead is shown to help. It is also promising that this new performance measure is leading to the discovery of new algorithms. Further testing is needed to determine if RLRU, or some variant of it, is better than LRU in practice.

We expect that for many other problems, the relative worst order ratio will distinguish between known algorithms which have the same competitive ratio, give the correct relationship between algorithms where the competitive ratio fails to, and lead to new algorithms which are better than algorithms which have an optimal competitive ratio.

Theorem 3 shows that no marking algorithm can be much better than LRU. It would be interesting to know if LRU is in fact the best marking algorithm according to the relative worst order ratio.

References

1. D. Achlioptas, M. Chrobak, and J. Noga. Competitive analysis of randomized paging algorithms. *Theoretical Computer Science*, 234(1–2):203–218, 2000.
2. S. Albers. On the influence of lookahead in competitive paging algorithms. *Algorithmica*, 18:283–305, 1997.
3. S. Albers, L. M. Favrholdt, and O. Giel. On paging with locality of reference. In *34th Annual ACM Symposium on the Theory of Computing*, pages 258–267, 2002.
4. L. A. Belady. A study of replacement algorithms for virtual storage computers. *IBM Systems Journal*, 5:78–101, 1966.
5. S. Ben-David and A. Borodin. A new measure for the study of on-line algorithms. *Algorithmica*, 11(1):73–91, 1994.
6. A. Borodin and R. El-Yaniv. On randomization in online computation. In *Proceedings of the IEEE Conference on Computational Complexity*, pages 226–238, 1997.
7. A. Borodin, S. Irani, P. Raghavan, and B. Schieber. Competitive paging with locality of reference. *Journal of Computer and System Sciences*, 50(2):244–258, 1995.
8. J. Boyar and L. M. Favrholdt. The relative worst order ratio for on-line algorithms. In *5th Italian Conference on Algorithms and Complexity*, volume 2653 of *LNCS*, pages 58–69, 2003.
9. J. Boyar, K. S. Larsen, and M. N. Nielsen. The accommodating function — a generalization of the competitive ratio. *SIAM Journal of Computation*, 31(1):233–258, 2001.
10. D. Breslauer. On competitive on-line paging with lookahead. *Theoretical Computer Science*, 209(1–2):365–375, 1998.
11. M. Chrobak, H. Karloff, T. Payne, and S. Vishwanathan. New results on server problems. *SIAM Journal on Discrete Mathematics*, 4(2):172–181, 1991.
12. M. Chrobak and J. Noga. LRU Is Better than FIFO. *Algorithmica*, 23:180–185, 1999.
13. L. Epstein, L. M. Favrholdt, and J. S. Kohrt. The relative worst order ratio applied to scheduling problems. Work in progress, 2003.
14. A. Fiat, R. M. Karp, M. Luby, L. A. McGeoch, D. D. Sleator, and N. E. Young. Competitive paging algorithms. *Journal of Algorithms*, 12:685–699, 1991.
15. A. Fiat and Z. Rosen. Experimental studies of access graph based heuristics: Beating the LRU standard? In *8th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 63 – 72, 1997.
16. J. S. Frederiksen and K. S. Larsen. On-line seat reservations via off-line seating arrangements. In *Eighth International Workshop on Algorithms and Data Structures*, volume 2748 of *LNCS*, pages 174–185. Springer-Verlag, 2003.
17. R. L. Graham. Bounds for certain multiprocessing anomalies. *Bell Systems Technical Journal*, 45:1563–1581, 1966.
18. D. S. Johnson. Fast algorithms for bin packing. *Journal of Computer and System Sciences*, 8:272–314, 1974.
19. B. Kalyanasundaram and K. Pruhs. Speed is as powerful as clairvoyance. *Journal of the ACM*, 47:617–643, 2000.
20. A. R. Karlin, M. S. Manasse, L. Rudolph, and D. D. Sleator. Competitive snoopy caching. *Algorithmica*, 3(1):79–119, 1988.
21. A. R. Karlin, S. J. Phillips, and P. Raghavan. Markov paging. *SIAM Journal of Computation*, 30:906–922, 2000.
22. C. Kenyon. Best-fit bin-packing with random order. In *7th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 359–364, 1996.
23. J. S. Kohrt. The relative worst order ratio applied to bin coloring. Work in progress, 2003.

24. E. Koutsoupias and C. H. Papadimitriou. Beyond competitive analysis. In *35th Annual Symposium on Foundations of Computer Science*, pages 394–400, 1994.
25. S. O. Krumke, W. E. de Paepe, J. Rambau, and L. Stougie. Online bin coloring. In *9th Annual European Symposium on Algorithms*, volume 2161 of *LNCS*, pages 74–85, 2001.
26. L. A. McGeoch and D. D. Sleator. A strongly competitive randomized paging algorithm. *Algorithmica*, 6:816–825, 1991.
27. P. Raghavan. A statistical adversary for on-line algorithms. In D. D. Sleator L. A. McGeoch, editor, *On-Line Algorithms*, volume 7 of *Series in Discrete Mathematics and Theoretical Computer Science*, pages 79–83. American Mathematical Society, 1992.
28. D. D. Sleator and R. E. Tarjan. Amortized efficiency of list update and paging rules. *Communications of the ACM*, 28(2):202–208, 1985.
29. E. Torng. A unified analysis of paging and caching. *Algorithmica*, 20:175–200, 1998.
30. M. J. van Kreveld and M. H. Overmars. Union-copy structures and dynamic segment trees. *Journal of the ACM*, 40(3):635–652, 1993.
31. N. Young. Competitive paging and dual-guided algorithms for weighted caching and matching (thesis). Technical Report CS-TR-348-91, Computer Science Department, Princeton University, 1991.
32. N. Young. The k -server dual and loose competitiveness for paging. *Algorithmica*, 11:525–541, 1994.

Appendix

Additional Comments on an Implementation of RLRU

In addition to the administration required to evict least recently used pages, which is similar to the administration necessary for LRU, RLRU needs to be able to perform the following operations:

1. Check if it faults on a page for the second time in a phase.
2. Mark a page, and unmark all pages.
3. Find least recently used page, possibly just among unmarked pages.
4. Check for a page in LFD's cache.

These can all be implemented in time $O(\log N)$ and space $O(N)$, where N is the number of different pages requested.

1. We use a balanced binary search tree over all the different pages on which RLRU has faulted during the phase.
2. Using a balanced binary search tree over all the different pages which have been requested, we mark a page by associating the current phase number with the page. Thus, by incrementing the phase number, we can unmark all pages in constant time.
3. Using a balanced binary search tree ordered on timestamp, the least recently used page can be found in logarithmic time. If the timestamp is also associated with pages in cache, then old timestamp entries can be found and updated when a page is requested. By adding information to the nodes in the tree regarding the last phase in which the page stored in the node was marked and information regarding the least recent phase of any node in the subtree, it is also possible in logarithmic time to find the least recently used page among those which are unmarked, i.e., not marked in the current phase. In an actual implementation, points 1, 2, and 3 can be combined.
4. At any given point in time, it is of course impossible to compute the entire contents of LFD's cache, since this depends on future requests. It is, however, possible, given a request and the request sequence up to that point, to compute whether or not LFD would have that particular page in cache. Using techniques [16] inspired by geometric algorithms [30], this can be done by registering the known time intervals of pages in LFD's cache in a balanced binary search tree. Also here, time $O(\log N)$ and space $O(N)$ can be obtained.