

On-Line Paging Against Adversarially Biased Random Inputs

Neal E. Young¹

Dartmouth College, Hanover, New Hampshire 03755

Received June 24, 1998

In evaluating an algorithm, worst-case analysis can be overly pessimistic. Average-case analysis can be overly optimistic. An intermediate approach shows that an algorithm does well on a *broad class* of input distributions. E. Koutsoupias and C. H. Papadimitriou (1994, in “Proc. of the 35th IEEE Annual Symp. on Foundation of Computer Science,” pp. 394–400, IEEE Press, New York) recently analyzed the least-recently-used (LRU) paging strategy in this manner, analyzing its performance on an input sequence generated by a so-called *diffuse* adversary—one that must choose each request probabilistically so that no page is chosen with probability more than some fixed $\epsilon > 0$. They showed that LRU achieves the optimal competitive ratio (for deterministic on-line algorithms), but they did not determine the actual ratio. In this paper we estimate the optimal ratios within roughly a factor of two for both deterministic strategies and randomized strategies. Around the threshold $\epsilon \approx 1/k$ (where k is the cache size), the optimal ratios are both $\Theta(\ln k)$. Below the threshold the ratios tend rapidly to $O(1)$. Above the threshold the ratio is unchanged for randomized strategies but tends rapidly to $\Theta(k)$ for deterministic ones. We also show that the competitive ratios for First-in-first-out (FIFO) and Flush-when-full (FWF) are both k when $\epsilon \geq 1/k$. In contrast, the ratio for LRU is less than $2 \ln k + 4$ when $\epsilon = 1/k$. It is folklore that LRU outperforms FIFO in practice, but to date the only other variant of competitive analysis in which LRU has been shown to outperform FIFO is the access graph model. For completeness, we give an alternate proof of the optimality of LRU. © 2000 Academic Press

1. INTRODUCTION AND BACKGROUND

The *paging problem* was originally studied in the context of two-level virtual memory systems composed of a large, slow-access memory augmented with a cache (a small, fast-access memory, holding likely-to-be accessed pages in order to minimize access time).

¹E-mail: ney@cs.dartmouth.edu. Research partially funded by NSF CAREER Award CCR-9720664.



This paper concerns the following standard abstraction of this simple and common problem. The input is an integer k and a finite sequence $s = s_1 s_2 \dots s_n$ of requests. The parameter k is called *the cache size*. The output is a *schedule*—a sequence $S_1 S_2, \dots, S_n$ of sets, where each set is of size at most k and each S_t contains s_t . Each request s_t is said to occur *at time* t . The items in S_t are said to be *in the cache* after time t up to and including time $t + 1$. An item is said to be *evicted at time* t if the item is in S_{t-1} but not in S_t . The *cost* of the schedule is the number of evictions. A schedule for an input is *optimal* if it achieves the minimum possible cost.

Paging algorithms considered in this paper include the following. *Least-recently-used* (LRU) evicts the item whose most recent request is the least recent among all items in the cache. *First-in-first-out* (FIFO) evicts the item that has been in the cache the longest. *Flush-when-full* (FWF) evicts all items in the cache (when the cache is full and does not contain the requested item). The *randomized marking algorithm* (RMARK [6]) operates as follows. After an item is requested, it is marked. When an item must be evicted, a nonmarked item is chosen uniformly at random, with the caveat that if all items in the cache are marked, then all marks are first erased. By a *deterministic marking algorithm*, we mean any deterministic algorithm that maintains marks as RMARK does and evicts only unmarked items. LRU and FWF are examples, FIFO is not. By a *lazy deterministic marking algorithm* (DMARK), we mean a deterministic marking algorithm that evicts an item only when necessary, and then only one item. This additional requirement excludes FWF.

An algorithm for the problem is *on-line* if, for any request sequence and any request in that sequence, the items in the cache after the request are independent of later requests. In many contexts, on-line algorithms are necessary, but online algorithms are necessarily suboptimal on some request sequences. Hence, a natural question is how on-line algorithms can be effectively analyzed and compared.

This paper is concerned with a generalization of the standard *competitive analysis* [15] of on-line algorithms. The standard model measures the quality of an algorithm A by its *competitive ratio*: the minimum (to be precise, infimum) c such that, for some constant b , for all request sequences s ,

$$A(s) \leq c \cdot \text{OPT}(s) + b.$$

Here $A(s)$ denotes the cost of the schedule produced by A on input s ; $\text{OPT}(s)$ denotes the cost of an optimal schedule. If A is a randomized algorithm, then $A(s)$ denotes the expected cost of A on input s . Note that k is an implicit, and fixed, parameter in these definitions. Standard

competitive analysis is a worst-case type of analysis, in contrast to much of the earlier work on paging, which is concerned with average-case analysis.²

In the standard competitive-analysis framework the following results are known. Any deterministic marking algorithm, including LRU and FWF, has a competitive ratio of k ; the ratio k is also achieved by FIFO and is the best possible for any deterministic on-line strategy [3, 15]. The randomized marking algorithm RMARK has a competitive ratio of $2H_k - 1$ [1, 6], where $H_k \doteq \sum_1^k 1/i < 1 + \ln k$. PARTITION [13] and EQUITABLE [1], more complicated randomized algorithms, each have competitive ratio H_k . No randomized strategy can have a better ratio than H_k [6].

Largely due to the unrealistic magnitude of the optimal competitive ratios [17], many variations on the standard model have been considered (e.g. [4, 5, 7, 9, 10, 12, 17, 19]). For a survey on competitive analysis of paging, we refer the reader to the recent book by Borodin and El-Yaniv [3, Chaps. 3–5].

This paper concerns the following generalization of the standard model, recently proposed by Koutsoupias and Papadimitriou [11]. For any class Δ of distributions on the input sequences and any deterministic or randomized algorithm A , define $\mathcal{R}(\Delta, A)$, the *competitive ratio of A against the Δ -diffuse adversary*, to be the minimum (again, to be precise, infimum) c such that for each distribution D in Δ , there is a constant b such that

$$E_D[A(r)] \leq c \cdot E_D[\text{OPT}(r)] + b.$$

Here r is a random sequence chosen according to D . Define the *optimal ratio for deterministic on-line algorithms (against the Δ -diffuse adversary)* to be

$$\mathcal{R}(\Delta) \doteq \inf_A \mathcal{R}(\Delta, A),$$

where A ranges over all deterministic on-line algorithms. Analogously, define the *optimal ratio for randomized on-line algorithms (against the Δ -diffuse adversary)* to be

$$\bar{\mathcal{R}}(\Delta) \doteq \inf_{\bar{A}} \mathcal{R}(\Delta, \bar{A}),$$

where \bar{A} ranges over all randomized on-line algorithms.

²At least one work [8] preceding competitive analysis blends average-case and worst-case analysis. It considers input sequences where each request is chosen from a fixed but unknown distribution on the pages and compares known paging strategies to the optimal *on-line* strategy for that distribution.

The particular class of distributions considered by Koutsoupias and Papadimitriou is denoted Δ_ϵ and is defined as follows. Any distribution D specifies, for each item x and sequence of requests s , the probability $\Pr_D(x | s)$ that the next request of the random sequence r is x given that the sequence so far is s . Then Δ_ϵ contains those distributions D such that, for any request sequence s and item x , $\Pr_D(x | s) \leq \epsilon$. The parameter ϵ is a measure of the inherent uncertainty of each request. Koutsoupias and Papadimitriou show that LRU achieves the optimal ratio in this model (i.e., $\mathcal{R}(\Delta_\epsilon, \text{LRU}) = \mathcal{R}(\Delta_\epsilon)$), but they leave open the question of what the ratio is.

Here we estimate the optimal ratios within roughly a factor of two, for both deterministic and randomized algorithms, and we show that FIFO and FWF are not optimal, at least for $\epsilon = 1/k$. Here is our main theorem.

THEOREM 1. *Define $\Phi(\epsilon, k) \doteq 1 + \sum_{i=1}^{k-1} \max\{\epsilon^{-1} - i, 1\}^{-1}$. For any ϵ , let $\epsilon' = 1/\lceil \epsilon^{-1} \rceil$. The competitive ratios of deterministic (\mathcal{R}) and randomized ($\overline{\mathcal{R}}$) on-line algorithms against the Δ_ϵ -diffuse adversary are bounded as follows:*

range	deterministic - $\mathcal{R}(\Delta_\epsilon)$		randomized - $\overline{\mathcal{R}}(\Delta_\epsilon)$	
	lower bound	upper bound	lower bound	upper bound
$\epsilon \leq 1/(k + 1)$	$\Phi(\epsilon, k) - 1$	$2\Phi(\epsilon, k)$	$\Phi(\epsilon', k) - 1$	$2\Phi(\epsilon, k)$
$\epsilon \geq 1/(k + 1)$	$\Phi(\epsilon, k)$	$2\Phi(\epsilon, k)$	H_k	H_k

The upper bound 2Φ for deterministic algorithms holds for any lazy marking algorithm (e.g., LRU), but not for FIFO or FWF: in fact, $\mathcal{R}(\Delta_\epsilon, \text{FIFO}) = \mathcal{R}(\Delta_\epsilon, \text{FWF}) = k$ for $\epsilon \geq 1/k$. The upper bound H_k for randomized algorithms holds for PARTITION and EQUITABLE. The weaker upper bound $2H_k - 1$ holds for RMARK.

In all cases except one, the competitive ratios of (lazy) deterministic and randomized marking algorithms are at least $\Phi - 1$ and at most 2Φ . The exception is that for ϵ above the threshold $1/(k + 1)$, the randomized ratio is H_k (independent of ϵ). To understand the behavior of the function Φ , consider the case $\epsilon = 1/n$ for some integer n . Then

$$\Phi(1/n, k) = 1 + H_{n-1} + \begin{cases} -H_{n-k} & \text{when } n \geq k, \\ k - n & \text{when } n \leq k. \end{cases}$$

Recall that $H_k \doteq \sum_1^k 1/i \approx \ln(k + 1)$. The threshold of Φ around $\epsilon \approx 1/k$ is very sharp:

$$\Phi(\epsilon, k) \text{ is } \begin{cases} \leq 1 + \ln \frac{1}{\delta} & \text{when } \epsilon = (1 - \delta)/k, \\ \approx \ln k & \text{when } \epsilon = 1/k, \\ \geq k \frac{\delta}{1 + \delta} & \text{when } \epsilon = (1 + \delta)/k. \end{cases}$$

2. TECHNICAL OVERVIEW

We refine an existing *worst-case* competitive analysis for paging [3, 6, 15] to take into account the probabilistic restrictions on the adversary. We call this particular analysis the *factor-two-analysis* because for our purposes (and when used to analyze the randomized marking algorithm [6]) it (at best) can approximate OPT only within a factor of two.

2.1. Review of Factor-Two-Analysis in Standard Model

Let A be any paging algorithm and let $s = s_1 s_2 \dots s_n$ be any sequence of requests. The *phases* of s partition the times $\{1, 2, \dots, n\}$ into intervals as follows. Define $t(1) = 1$. For $l \in \mathbb{N}$ inductively define

$$t(l + 1) \doteq 1 + \max\{j \leq n : |\{s_{t(l)}, s_{t(l)+1}, s_{t(l)+2}, \dots, s_j\}| \leq k\}.$$

For each $l \in \mathbb{N}$ such that $t(l) \leq n$, the *lth phase* of s is defined to be the time interval $\{t(l), t(l) + 1, \dots, t(l + 1) - 1\}$. Thus, during each phase except the last, k distinct items are requested.

In the context of a particular time t , the *current request* refers to the request s_t . *This phase* or, synonymously, *the current phase* means the phase containing the time t . An item is *requested previously in this phase* if it is requested during this phase before time t . In the additional context of a particular schedule $S = S_1 S_2 \dots S_n$ for s , *the cache* refers to the set S_{t-1} of items in the cache before request t . Then at each time t , each item is classified with respect to its status before request s_t as follows:

new—not requested previously in this phase or in the last phase.

old—requested during the last phase, but not previously in this phase.

redundant—requested previously in this phase.

worrisome—requested in the last phase or previously in this phase, but not in the on-line algorithm's cache.

Each *request* is classified as well, according to the status of the requested item. For instance, a request s_i is *new* if the requested item was new after request s_{i-1} . Each phase (except possibly the last) has k nonredundant requests, each one of which is either new or old. Define

$\text{new}(s)$ —the total number of new requests in sequence s .

$\text{new_in_ph}(l)$ —(in the context of some sequence) the total number of new requests in the l th phase of the sequence. Here l is any positive integer. If $l = 0$ or there is no l th phase, define $\text{new_in_ph}(l)$ to be 0.

The relevance of the new requests is as follows.

LEMMA 1 [6, 16]. $\text{new}(s)/2 \leq \text{OPT}(s) \leq \text{new}(s)$.

Proof. Consider the $(l-1)$ st and l th phases of s for any l . The number of distinct items requested in the two phases is $k + \text{new_in_ph}(l)$. Thus, the number of evictions incurred by OPT during the two phases is at least $\text{new_in_ph}(l)$ and

$$\begin{aligned} \text{OPT}(s) &\geq \max \left\{ \sum_{l \text{ odd}} \text{new_in_ph}(l), \sum_{l \text{ even}} \text{new_in_ph}(l) \right\} \\ &\geq \sum_l \text{new_in_ph}(l)/2 = \text{new}(s)/2. \end{aligned}$$

On the other hand, the following schedule costs at most $\text{new}(s)$. At the beginning of each phase, evict those items that are not requested during the phase and bring in the items that are not in the cache but are requested during the phase. After each phase ends, the items requested during that phase are in the cache, so the number of evictions in the next phase is just the number of new requests in that phase. Thus, the cost of this schedule is $\text{new}(s)$. Since the schedule produced by OPT is at least as good, $\text{OPT}(s) \leq \text{new}(s)$. ■

By the *amortized* cost incurred by OPT during a phase, we mean half the number of new requests in that phase. By the lemma above, the total cost incurred by OPT is at least the total of these amortized costs and at most twice the total. To show bounds on the competitive ratio of A , we use the standard method of bounding the cost incurred by A during a phase divided by the amortized cost incurred by OPT during the phase. For instance, if this ratio is at most c for each phase of a sequence s , then it follows immediately that $A(s) \leq c \cdot \text{OPT}(s)$.

One intuition for understanding RMARK and other marking algorithms such as LRU and even FWF is that they are emulating the schedule described in the proof above that $\text{OPT}(s) \leq \text{new}(s)$. That is, during each phase, the goal (intuitively speaking) is to get the items that will be

requested during the phase into the cache. From this point of view, once an item is requested during a phase, it should be kept in the cache. This is the principle that defines a deterministic marking algorithm.

If this principle is followed, then only nonredundant requests can cause evictions. Since the phase ends after k nonredundant requests, any deterministic marking algorithm incurs a cost of at most k during the phase. This means that in the standard model, the competitive ratio is at most k (OPT also incurs at least one eviction per phase). Conversely, the adversary can force a ratio of k against a deterministic on-line algorithm by making one new request each phase and then making $k - 1$ requests, each to whichever old item is not currently in the cache.

2.2. Factor-Two-Analysis for the Diffuse Adversary

During each phase there are k nonredundant requests. New requests increase the cost to OPT. The only other requests that cost DMARK are worrisome requests. Each worrisome request is nonredundant. So, intuitively, a good adversary assigns probability to worrisome items most, then to marked items, then to unmarked items in the cache, and, as a last resort, to new items. (This is not quite accurate, because requesting a new item can benefit the adversary in that it increases the number of worrisome items, but this turns out to be a minor effect, accounting for the occasional “ -1 ” in the lower bound.)

We next calculate the upper bound glossing over an issue of probabilistic conditioning. In the subsequent section we give a formally correct treatment. The intuition for the lower bound is essentially described above.

Consider the l th phase for any l . There are k nonredundant requests in the phase (except possibly for the last phase, which may have fewer). Consider the state of any marking algorithm DMARK just before the $(i + 1)$ st nonredundant request, for $1 \leq i \leq k - 1$.

The i redundant items are marked and in the cache. Of the k items requested last phase, at most $\text{new_in_ph}(l)$ are worrisome (out of the cache). Thus, the adversary can assign at most $\epsilon \text{new_in_ph}(l)$ probability to worrisome items. Since there are only i redundant items, the adversary has to assign at least $1 - \epsilon i$ probability to nonredundant items. Therefore, the probability that the request will be worrisome, given that the request turns out to be nonredundant, is at most

$$\frac{\epsilon \text{new_in_ph}(l)}{1 - \epsilon i} = \frac{\text{new_in_ph}(l)}{\epsilon^{-1} - i}$$

(or 1 if this quantity is negative or more than 1). Summing over i , adding $\text{new_in_ph}(l)$ for the evictions due to new requests, and dividing by

$\text{new_in_ph}(l)/2$ (the amortized cost incurred by OPT for the phase) gives the desired upper bound 2Φ on the competitive ratio.

The issue we have glossed over in our intuitive explanation is that $\text{new_in_ph}(l)$ is a random variable, because the adversary generates the sequence s randomly. Further, conditioning on $\text{new_in_ph}(l)$ taking on a given value may “break” the probabilistic restriction on the adversary.

3. UPPER BOUND FOR MARKING ALGORITHMS

Next we prove the upper bounds on deterministic strategies in Theorem 1:

LEMMA 2. *For any lazy deterministic marking algorithm DMARK and $D \in \Delta_\epsilon$,*

$$E_D[\text{DMARK}(r)] \leq 2\Phi(\epsilon, k) \cdot E_D[\text{OPT}(r)] + O(1)$$

Proof. Without loss of generality, assume that D generates only sequences whose last phase has k nonredundant requests. (Otherwise we can easily modify the distribution so that the condition is satisfied, while increasing $E[\text{OPT}(r)]$ by at most the constant k .) In the context of the random sequence r , define the following random variables and events.

$R_{l,i}$ —the $(i + 1)$ st nonredundant request in the l th phase of r , if there is an l th phase.

$\text{prefix}(R)$ —the prefix of r up to but not including request R of r .

$\text{new_bef}(R)$ —the number of new requests before request R in the phase of r containing R .

$\text{new_in_ph}(l)$ —the total number of new requests made in the l th phase of r , if there is an l th phase, otherwise 0.

$\text{worrisome}(R)$ —the event that request R of r is worrisome.

In what follows, we abuse notation slightly as follows. By the event $\text{prefix}(R_{l,i}) = s$, we mean that there is an l th phase in r and the prefix of r preceding request $R_{l,i}$ is sequence s . Similarly, by the event $\text{worrisome}(R_{l,i})$, we mean that there is an l th phase in r and the request $R_{l,i}$ in that phase is worrisome.

We start by proving the following claim:

Claim 1. Fix any l and $i(1 \leq i \leq k - 1)$. Let s be any sequence such that the event $\text{prefix}(R_{l,i}) = s$ can happen. That is, s has l phases, and

the last phase of s has i nonredundant requests. Then

$$\begin{aligned} & \Pr[\text{worrisome}(R_{l,i}) \mid \text{prefix}(R_{l,i}) = s] \\ & \leq \mathbb{E} \left[\frac{\text{new_bef}(R_{l,i})}{\max\{1, \epsilon^{-1} - i\}} \mid \text{prefix}(R_{l,i}) = s \right]. \end{aligned}$$

Conditioning on $\text{prefix}(R_{l,i}) = s$ lets us use the restrictions on the adversary.

Here is the proof of Claim 1. In the event that s is a prefix of r , consider the random variable r_t where $t = |s| + 1$. (There must be such a request because $i < k$ and each phase of r , including the last, by the assumption at the beginning of the proof, has k nonredundant requests.)

The event $\text{prefix}(R_{l,i}) = s$ happens if and only if s is a prefix of r and r_t is nonredundant. If $\text{prefix}(R_{l,i}) = s$, then the event $\text{worrisome}(R_{l,i})$ happens if and only if r_t is worrisome. Thus,

$$\begin{aligned} & \Pr(\text{worrisome}(R_{l,i}) \mid \text{prefix}(R_{l,i}) = s) \\ & = \Pr(\text{worrisome}(r_t) \mid s \text{ is a prefix of } r \text{ and } r_t \text{ is nonredundant}) \\ & = \frac{\Pr(\text{worrisome}(r_t) \mid s \text{ is a prefix of } r)}{\Pr(r_t \text{ is nonredundant} \mid s \text{ is a prefix of } r)}. \end{aligned}$$

Assume that s is a prefix of r . After processing s , DMARK has all but $\text{new_bef}(r_t)$ of the items requested in the previous phase in the cache. Thus, the adversary can assign at most $\epsilon \text{new_bef}(r_t)$ probability to worrisome items. Thus, the numerator above is at most $\epsilon \text{new_bef}(r_t)$. Since there have been i nonredundant requests in this phase before r_t , there are only i redundant items, so the denominator above is at least $1 - \epsilon i$. To finish the proof of Claim 1, note that $\mathbb{E}[\text{new_bef}(R_{l,i}) \mid \text{prefix}(R_{l,i}) = s] = \text{new_bef}(r_t)$.

Now fix i and l . In the set of events $\{\text{prefix}(R_{l,i}) = s \mid s \text{ is a sequence}\}$, exactly one event happens. Thus, the bound in Claim 1 holds unconditionally:

$$\Pr[\text{worrisome}(R_{l,i})] \leq \mathbb{E} \left[\frac{\text{new_bef}(R_{l,i})}{\max\{1, \epsilon^{-1} - i\}} \right].$$

Since $\text{new_bef}(R_{l,i}) \leq \text{new_in_ph}(l)$, it follows that for all l and i ,

$$\Pr[\text{worrisome}(R_{l,i})] \leq \mathbb{E} \left[\frac{\text{new_in_ph}(l)}{\max\{1, \epsilon^{-1} - i\}} \right].$$

Since $\text{DMARK}(r)$ is the number of new or worrisome requests in r ,

$$\begin{aligned} \mathbb{E}[\text{DMARK}(r)] &\leq \mathbb{E}\left[\sum_l \text{new_in_ph}(l) + \sum_{l,i} \frac{\text{new_in_ph}(l)}{\max\{1, \epsilon^{-1} - i\}}\right] \\ &= \left(1 + \sum_i \max\{1, \epsilon^{-1} - i\}^{-1}\right) \cdot \mathbb{E}\left[\sum_l \text{new_in_ph}(l)\right] \\ &= \Phi(\epsilon, k) \mathbb{E}[\text{new}(r)] \\ &\leq \Phi(\epsilon, k) \mathbb{E}[\text{OPT}(r)/2] \quad (\text{by Lemma 1}). \end{aligned}$$

■

4. LOWER BOUND FOR DETERMINISTIC ALGORITHMS

Next we prove the lower bounds on deterministic strategies in Theorem 1:

LEMMA 3. *For any $\epsilon > 0$, any k , and any deterministic on-line algorithm A , there is a distribution $D \in \Delta_\epsilon$ such that*

$$\mathbb{E}_D[A(r)] \geq (\Phi(\epsilon, k) - 1 + 1/m) \cdot \mathbb{E}_D[\text{OPT}(r)],$$

where $m = \max\{1, \lceil \epsilon^{-1} \rceil - k\}$ and $\mathbb{E}_D[\text{OPT}(r)]$ is arbitrarily large.

Proof. We describe D by describing an adversary that requests items probabilistically subject to the limitations of Δ_ϵ . Fix $\epsilon > 0$ and $k > 0$. Assume $\epsilon > 1/2k$ (otherwise the desired lower bound is trivially satisfied, because $\Phi(1/2k, k) - 1 + 1/m$ is less than 1).

The adversary requests the items in an on-line fashion, phase by phase. In the first part of each phase, the adversary makes m new requests by assigning probability only to items not previously requested.

For each remaining request, the adversary assigns a probability to each item as follows. First priority is given to worrisome items (those previously requested in this phase or in the last one but not in the cache of A). Second priority is given to redundant items (those requested previously in this phase and in the cache). Third priority is given to the remaining old items (the items not yet requested this phase, but in the cache).

Items are selected in order of priority and assigned as much probability as possible, subject to the constraint that no item is assigned probability more than ϵ and the total probability assigned is 1. By the choice of m , we have $(k + m)\epsilon \geq 1$, so all three kinds of items suffice for all probability to be assigned.

The adversary follows this strategy until k distinct items have been requested, at which point the adversary begins a new phase. The adversary continues for N phases, where N is arbitrarily large so that $\text{OPT}(r)$ is also arbitrarily large.

This defines the distribution $D \in \Delta_\epsilon$. Let r be a random request chosen from D . Next we prove that $E[A(r)] \geq Nm(\Phi(\epsilon, k) - 1 + 1/m)$. This proves the claimed bound, since $\text{OPT}(r) \leq Nm$ (by Lemma 1). Consider any l s.t. $1 \leq l \leq N$. For $i = m, \dots, k - 1$, define

$\text{worrisome}(R_{l,i})$ —the event that the i th nonredundant request of the l th phase is worrisome.

The expectation of $A(r)$ is $Nm + \sum_{l,i} \Pr[\text{worrisome}(R_{l,i})]$. For any l and i s.t. $m \leq i \leq k - 1$, consider the time just before the $(i + 1)$ st nonredundant request of the l th phase. There have been i nonredundant requests so far in the phase, so there are i redundant items. There have been m new requests so far, so there are $k + m$ items that were requested last phase or already this phase. Since the on-line algorithm has at least m of these items not in the cache, there are at least m worrisome items. Thus, the adversary assigns at least ϵm probability to worrisome items and at least ϵi probability to redundant items. (Unless $\epsilon m + \epsilon i > 1$, in which case $\Pr[\text{worrisome}(R_{l,i})] = 1$ —the adversary forces a worrisome request.) Thus, the probability that the request is worrisome, conditioned on it being nonredundant, is

$$\begin{aligned} & \Pr[\text{worrisome}(R_{l,i})] \\ & \geq \frac{\epsilon m}{\max\{1 - \epsilon i, \epsilon m\}} = \frac{m}{\max\{\epsilon^{-1} - i, m\}} = \frac{m}{\max\{\epsilon^{-1} - i, 1\}}. \end{aligned}$$

The rightmost equality holds because the choice of m implies that either $m = 1$ or $\epsilon^{-1} - i \geq m$. Adding the m new requests and summing over $i = m, \dots, k - 1$, the expected cost to A for each of the N phases is at least

$$m + \sum_{i=m}^{k-1} \frac{m}{\max\{\epsilon^{-1} - i, 1\}} \geq 1 + \sum_{i=1}^{k-1} \frac{m}{\max\{\epsilon^{-1} - i, 1\}}.$$

The rightmost expression is $m(\Phi(\epsilon, k) - 1 + 1/m)$. \blacksquare

The adversary can probably be made a little stronger to get a slightly better lower bound when $\epsilon \leq 1/(k + 1)$. In this case the issue of how the optimal adversary should fix m appears to be relatively subtle. This is why the lower bound loses the additive 1 with respect to the upper bound in this case. One small improvement to the above adversary would be, when

the adversary is requesting new items, to use the opportunity to also allocate probability to worrisome items.

5. FIFO AND FWF ARE NOT OPTIMAL

In this section we show that FIFO and FWF are not optimal.

LEMMA 4. *For $\epsilon \geq 1/k$, there exists a distribution $D \in \Delta_\epsilon$ such that*

$$E_D[\text{FIFO}(r)] = E_D[\text{FWF}(r)] = kE_D[\text{OPT}(r)].$$

Proof. We first consider FIFO. The adversary requests the items in an on-line fashion, phase by phase. Let S denote the set of items in FIFO's cache at the start of a phase. Let $x \in S$ be the item that was most recently brought into the cache. Let y be some item not in the cache. In the remainder of the phase, the adversary simply assigns probability $1/k$ to each page in $S \cup \{y\} - \{x\}$. The phase continues until x is evicted from the cache.

By the choice of x , each item in S is evicted at some point during the phase. Thus, FIFO incurs k evictions for the phase. On the other hand, the number of new requests in the phase is 1, so by Lemma 1 OPT incurs a cost of at most one per phase. This shows the lower bound of k for FIFO's ratio; equality follows because FIFO is known to be k -competitive in the standard model. The same adversary with an even simpler analysis shows the claim for FWF.

6. BOUNDS ON RANDOMIZED STRATEGIES

In this section we finish the proof of Theorem 1 by proving the upper and lower bounds for randomized strategies claimed there. By using what we already know, very little work is required to get the bounds.

We first consider lower bounds. Fix $\epsilon > 0$ and $k > 0$. We start with the case $\epsilon \leq 1/(k + 1)$. For simplicity we make the technical assumption that ϵ^{-1} is an integer. This assumption is not too restrictive and allows us to reuse the deterministic lower bound as follows.

LEMMA 5. *If ϵ^{-1} is an integer greater than k , then the distribution D described in the proof of Lemma 3 is independent of the algorithm A .*

Proof. Consider that distribution. Within each phase, the random sequence r has requests to m new items, followed by requests restricted to a set of $k + m$ items, where $m = \max\{1, \epsilon^{-1} - k\}$, until k distinct items have been requested. The condition on ϵ and the choice of m imply that

$m = \epsilon^{-1} - k$, so that $\epsilon = 1/(k + m)$. In this case, each phase simply consists of requests to m new items, followed by a sequence of requests to the $k + m$ items, where each request is chosen *uniformly at random* from those $k + m$ items, until a total of k distinct items have been requested, after which the next phase begins. ■

This distribution generalizes a distribution defined in a previous lower bound on the competitive ratio of randomized on-line strategies against the standard adversary [3, Theorem 8.7], [14, Theorem 13.2]. (That lower bound is equivalent to our case $m = 1$.) There and here, Yao's principle implies that for a random input r from any input distribution D , any randomized on-line algorithm A_R satisfies

$$E[A_R(r)] \geq \inf_A E[A(r)],$$

where A ranges over all deterministic on-line algorithms.

(Briefly, this is because A_R may be viewed as probabilistically picking some deterministic algorithm A and then running A on the input r . Thus, $E_D[A_R(r)] = \sum_A \Pr[A_R \text{ chooses } A] \cdot E_D[A(r)] \geq \inf_A E_D[A(r)]$. Here D can be any distribution, but we take it to be the one defined in Lemma 3. The input r is randomly chosen from D . We refer the reader to [14, Theorem 13.2] or [3, Theorem 8.7] for a full explanation of Yao's principle in this context.)

By Lemma 5, in the special case when ϵ^{-1} is an integer greater than k , the distribution defined in the previous section is independent of the on-line algorithm A . Thus, by Yao's principle, the lower bounds proved there extend to randomized algorithms. This proves:

LEMMA 6. *Suppose $\epsilon \leq 1/(k + 1)$ and ϵ^{-1} is an integer. Then the lower bound established in Lemma 3 also applies to randomized on-line algorithms.*

Decreasing ϵ only weakens the adversary. Thus, when ϵ is not an integer, letting $\epsilon' = 1/\lceil \epsilon^{-1} \rceil < \epsilon$, the lower bounds hold with ϵ' replacing ϵ .

Also, when $\epsilon = 1/(k + 1)$ one can verify that the above lemma implies that the ratios are at least $H_k \doteq \sum_1^k 1/i$. This proves:

LEMMA 7. *Suppose $\epsilon \geq 1/(k + 1)$. Then $\bar{\mathcal{R}}(\Delta_\epsilon) \geq H_k$.*

So the above two lemmas prove the lower bounds for randomized strategies claimed in Theorem 1. What about the upper bounds? Because the diffuse adversary is no stronger than the standard adversary, we get immediately from previous results that:

LEMMA 8. *For $\epsilon \leq 1/(k + 1)$, $\bar{\mathcal{R}}(\Delta_\epsilon, \text{RMARK}) \leq 2\Phi(\epsilon, k)$.*

For $\epsilon \geq 1/(k + 1)$, $\bar{\mathcal{R}}(\Delta_\epsilon, \text{RMARK}) \leq 2H_k - 1$, while $\bar{\mathcal{R}}(\Delta_\epsilon, \text{PARTITION}) \leq H_k$, and $\bar{\mathcal{R}}(\Delta_\epsilon, \text{EQUITABLE}) \leq H_k$.

The first upper bound follows from the fact that Lemma 2 also applies to RMARK (since the upper bound applies to any deterministic marking algorithm, i.e., any conditioning of RMARK on a particular outcome of its random choices). The remaining upper bounds follow from known upper bounds on the competitive ratios of the various algorithms against the (stronger) standard adversary [1, 3, 6, 13]. Lemma 8 proves the upper bounds on randomized strategies in Theorem 1. This completes the proof of that theorem.

7. ALTERNATE PROOF THAT LRU IS OPTIMAL

For the record, we include here a distillation of Koutsoupias and Papadimitriou's proof that LRU is optimal against the diffuse adversary Δ_ϵ . This version of the proof is shorter and self-contained, but does not give the intermediate results about work functions in the original proof.

Given a request sequence s of items from a universe U and an (arbitrary) initial ordering π of the items, define the *rank of an item* $x \in U$ in s to be the rank of x in the following order: items that are requested in s are first, in order of last request; items that are not requested in s are next, ordered by π .

In analyzing an on-line paging algorithm, if s is the sequence of requests seen so far, then the most recently requested item currently has rank 1, the next most recently requested item currently has rank 2, etc. Without loss of generality, when specifying a request or the contents of the cache, we can specify each item by its current rank; this uniquely identifies the item. Except in the proof of Lemma 9 where we use both representations, *items in this section are assumed to be specified by their current rank*.

LEMMA 9. *Let r and r' be two equal-length request sequences. Let \mathbf{r} and \mathbf{r}' , respectively, be the same sequences but with each request specified by rank (w.r.t. the same initial ordering and universe). If \mathbf{r} dominates \mathbf{r}' in the sense that $r_t \geq r'_t$ for all t , then $\text{OPT}(\mathbf{r}) \geq \text{OPT}(\mathbf{r}')$.*

Proof. It suffices to prove the case when there is a single d such that $r'_d = r_d - 1$ but $r_t = r'_t$ for all $t \neq d$. The general case then follows by induction. Assume such a d .

How do r and r' differ? Consider the two sequences simultaneously for $t = 1, 2, \dots, |r|$ in an on-line fashion. At each t focus on the ranks the items in the two subsequences $s = r_1 r_2 \dots r_t$ and $s' = r'_1 r'_2 r' \dots r'_t$.

At each time $t < d$, for each item, the rank in s equals the rank in s' . Let x and x' be the items requested, respectively, in r and r' at time d . By assumption, just before time d , the respective ranks of x and x' are r_d and $r_d - 1$. What about just after time d ? In sequence s , the rank of x changes

to 1, while the rank of x' changes to r_d . In sequence s' , the rank of x stays r_d , while the rank of x' changes to 1. For each item other than x or x' , the rank of the item is equal in both sequences.

This means that the sequence of items requested by r is the same as the sequence of items requested by r' , except that from time d to the end, the roles of x and x' are reversed: if r requests x (resp. x'), then r' requests x' (resp. x).

Let i and i' , respectively, be the times of the most recent requests to x and x' before time d . (If either item is being requested for the first time, then let i or i' equal 1, as appropriate.) By assumption $r'_d = r_d - 1$, so $i \leq i'$.

Consider any schedule S for r . For any j with $i' < j \leq d$, consider obtaining S' from S by reversing the roles of x and x' from time j onward (i.e., swapping the two in S_j, S_{j+1}, \dots). By the established relation between r and r' , S' will be a valid schedule for r' . To finish, we need only choose j so that S' costs no more than S . In particular, at time j , S' should evict no more of the two pages $\{x, x'\}$ than S does. If for some j , $|\{x', x\} \cap S_j| \in \{0, 2\}$ or $|\{x', x\} \cap S_{j-1}| \in \{0, 2\}$, then this j clearly suffices. Otherwise there is a j such that $\{x', x\} \cap S_{j-1} = \{x'\}$ and $\{x', x\} \cap S_j = \{x\}$. Using this j , S' is cheaper than S . ■

THEOREM 2 [11]. *Let D be any distribution $D \in \Delta_\epsilon$. Let A be any deterministic on-line algorithm. Then there is a distribution $D' \in \Delta_\epsilon$ such that*

$$E_D[\text{LRU}(r)] \leq E_{D'}[A(r')] \text{ and } E_D[\text{OPT}(r)] \geq E_{D'}[\text{OPT}(r')],$$

where r and r' are randomly chosen according to D and D' , respectively.

Thus, $\mathcal{R}(\Delta_\epsilon) = \mathcal{R}(\Delta_\epsilon, \text{LRU})$.

Proof. In what follows, we assume all items are specified not by name but by rank (with respect to some sequence implicit in context, the universe U of the items requested by D , and an arbitrary initial ordering).

The following random experiment defines the distribution D' by describing how to choose a random sequence r' according to that distribution. Choose a random sequence r according to D . Reveal r in an on-line fashion, one request at a time, producing each corresponding request of r' as follows.

Let L denote the cache of LRU (specified by rank with respect to \mathbf{s}) after processing $\mathbf{s} = r_1 \dots r_{t-1}$. Similarly, let A denote the cache (specified by rank with respect to \mathbf{s}') of A after processing $\mathbf{s}' = r'_1 \dots r'_{t-1}$. Let f be any 1-1 mapping from $A - L$ into $L - A$ (note $|A| \leq |L|$) and define (in the context of \mathbf{s} and \mathbf{s}')

$$\mathcal{X} \doteq \{x \in A - L \mid p(f(x)) < p(x)\}, \text{ where}$$

$$p(x) \doteq \Pr_D(x \mid \mathbf{s}).$$

Finally, determine r_t as follows. First set $r'_t = r_t$, but if $r_t \in \mathcal{X}$, change r'_t to $f(r_t)$ with probability $p(f(r_t))/p(r_t)$.

This completes the random experiment that gives r' and so defines D' . Each outcome of this experiment determines a *pair* of random variables (r, r') .

We use $\Pr_{D'}(X \mid \mathbf{s}, \mathbf{s}')$ to denote the probability of event X conditioned on \mathbf{s} and \mathbf{s}' being prefixes of r and r' . The following claim characterizes the distribution of r'_t conditioned on this event.

Claim 2. Fix any two sequences \mathbf{s} and \mathbf{s}' with length $t - 1$. Condition the random experiment above on \mathbf{s} and \mathbf{s}' being prefixes of the random variables r and r' , respectively. In the context of \mathbf{s} and \mathbf{s}' , define $p'(x) \doteq \Pr_{D'}(r'_t = x \mid \mathbf{s}, \mathbf{s}')$.

Then for each x , $p'(x) = p(\pi(x))$, where π is the permutation defined by $\pi(x) = f(x)$ and $\pi(f(x)) = x$ for $x \in \mathcal{X}$, and otherwise $\pi(x) = x$.

The claim follows by direct calculation based on the last line of the experiment.

Claim 3. Let r, r', \mathbf{s} , and \mathbf{s}' be as in Claim 2. Then

$$\Pr_{D'}[\text{LRU faults on } r_t \mid \mathbf{s}, \mathbf{s}'] \leq \Pr_{D'}[A \text{ faults on } r'_t \mid \mathbf{s}, \mathbf{s}'].$$

Why? It suffices to show that for every item x in \mathbf{A} , there is a unique item y in \mathbf{L} such that $p'(x) \leq p(y)$. But by Claim 2 and the choice of \mathcal{X} , this is the case: take $y = x$ unless $x \in \mathbf{A} - \mathbf{L}$, in which case take $y = f(x) \in \mathbf{L} - \mathbf{A}$.

Claim 4. The first part of the theorem is true: $E_D[\text{LRU}(r)] \leq E_{D'}[A(r')]$.

This follows directly from Claim 3. To see it formally, letting \mathbf{s} and \mathbf{s}' range over all equal-length pairs of sequences, we have

$$\begin{aligned} E_D[\text{LRU}(r)] &= \sum_{\mathbf{s}, \mathbf{s}'} \Pr_{D'}(\mathbf{s}, \mathbf{s}') \Pr_{D'}[\text{LRU faults on } r_{|\mathbf{s}|+1} \mid \mathbf{s}, \mathbf{s}'] \\ &\leq \sum_{\mathbf{s}, \mathbf{s}'} \Pr_{D'}(\mathbf{s}, \mathbf{s}') \Pr_{D'}[A \text{ faults on } r'_{|\mathbf{s}'|+1} \mid \mathbf{s}, \mathbf{s}'] \\ &= E_{D'}[A(r')]. \end{aligned}$$

Above $\Pr_{D'}(\mathbf{s}, \mathbf{s}')$ denotes the probability that \mathbf{s} is a prefix of r and \mathbf{s}' is a prefix of r' in the random experiment.

Claim 5. The second part of the theorem is true: $E_D[\text{OPT}(r)] \geq E_{D'}[\text{OPT}(r')]$.

Since the random experiment described above produces the same distribution on r as D does, it suffices to prove the inequality assuming that the

pair (r, r') is generated by that experiment. Since LRU keeps the most recently requested items in its cache, and $f: (A - L) \rightarrow (L - A)$, we have $x \leq f(x)$. Thus, in any outcome, r dominates r' (in the sense of Lemma 9) and so $\text{OPT}(r) \geq \text{OPT}(r')$. This proves the claim.

Claim 6. The distribution D' defined by the random experiment is in Δ_ϵ .

This also follows directly from Claim 2. To prove it in detail, we need to show that for any s' and x , $\Pr_{D'}(x | s') \leq \epsilon$. But

$$\Pr_{D'}(x | s') = \sum_{\mathbf{s}} \Pr_D(\mathbf{s}) \Pr_{D'}(r'_t = x | \mathbf{s}, s') \leq \sum_{\mathbf{s}} \Pr_D(\mathbf{s}) \epsilon = \epsilon.$$

Above $\Pr_D(\mathbf{s})$ denotes the probability that \mathbf{s} is a prefix of r , and \mathbf{s} ranges over all sequences of length $|s'| = t - 1$. The second-to-last inequality follows because by Claim 2 each $\Pr_{D'}(r'_t = x | \mathbf{s}, s')$ equals $\Pr_D(y | \mathbf{s})$ for some y , and by the assumption that $D \in \Delta_\epsilon$, $\Pr_D(y | \mathbf{s}) \leq \epsilon$. This proves the claim (and the theorem!). ■

ACKNOWLEDGMENTS

Thanks to Elias Koutsoupias, Lenny Ng, and Jeff Westbrook for helpful discussions. Thanks to Prudence Wong who pointed out that, in contrast to a claim in an earlier version of this paper [18], FIFO is not a marking algorithm; that tip led to the discovery of Lemma 4. Thanks to an anonymous referee for helping to greatly clarify the presentation of the results.

REFERENCES

1. D. Achlioptas, M. Chrobak, and J. Noga, Competitive analysis of randomized paging algorithms, in "Algorithms—ESA '96, Fourth Annual European Symposium, Barcelona, Spain, 25–27 September 1996" (J. Díaz and M. Serna, Eds.), Lecture Notes in Computer Science, Vol. 1136, pp. 419–430, Springer-Verlag, Berlin/New York.
2. "Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms, San Francisco, 25–27 January 1998."
3. A. Borodin and R. El-Yaniv, "Online Computation and Competitive Analysis," Cambridge Univ. Press, Cambridge, UK, 1998.
4. A. Borodin, S. Irani, P. Raghavan, and B. Schieber, Competitive paging with locality of reference, *J. Comput. System Sci.* **50** (1995), 244–258.
5. A. Fiat and A. R. Karlin, Randomized and multipointer paging with locality of reference, in "Proceedings of the Twenty-Seventh Annual ACM Symposium on Theory of Computing, Las Vegas, 29 May–1 June 1995," pp. 626–634, Assoc. Comput. Mach., New York.
6. A. Fiat, R. M. Karp, M. Luby, L. A. McGeoch, D. D. Sleator, and N. E. Young, Competitive paging algorithms, *J. Algorithms* **12** (1991), 685–699.
7. A. Fiat and Z. Rosen, Experimental studies of access graph based heuristics: Beating the LRU standard? in "Proceedings of the Eighth Annual ACM-SIAM Symposium on

- Discrete Algorithms," New Orleans, 5–7 January 1997, pp. 63–72, Assoc. Comput. Mach., New York.
8. P. A. Franaszek and T. J. Wagner, Some distribution-free aspects of paging algorithm performance, *J. Assoc. Comput. Math.* **21** (1974), 31–39.
 9. S. Irani, A. R. Karlin, and S. Phillips, Strongly competitive algorithms for paging with locality of reference, *SIAM J. Comput.* **25** (1996), 477–497.
 10. A. R. Karlin, S. J. Phillips, and P. Raghavan, Markov paging (extended abstract), in "33rd Annual Symposium on Foundations of Computer Science, Pittsburgh, 24–27 October 1992," pp. 208–217, IEEE Press, New York.
 11. E. Koutsoupias and C. H. Papadimitriou, Beyond competitive analysis, in "Proc. of the 35th IEEE Annual Symp. on Foundation of Computer Science," pp. 394–400, IEEE Press, New York, 1994.
 12. C. Lund, S. Phillips, and N. Reingold, IP over connection-oriented networks and distributional paging, in "35th Annual Symposium on Foundations of Computer Science, Santa Fe, NM, 20–22 November 1994," pp. 424–434, IEEE Press, New York.
 13. L. A. McGeoch and D. D. Sleator, A strongly competitive randomized paging algorithm, *Algorithmica* **6** (1991), 816–825.
 14. R. Motwani and P. Raghavan, "Randomized Algorithms," Cambridge Univ. Press, Cambridge, UK, 1995.
 15. D. D. Sleator and R. E. Tarjan, Amortized efficiency of list update and paging rules, *Comm. Assoc. Comput. Math.* **28** (1985), 202–208.
 16. N. E. Young, "Competitive Paging and Dual-Guided Algorithms for Weighted Caching and Matching," (Thesis) Tech. Rep. CS-TR-348-91, Computer Science Department, Princeton University, October 1991.
 17. N. E. Young, The k -server dual and loose competitiveness for paging, *Algorithmica* **11** (1994), 525–541.
 18. N. E. Young, Bounding the diffuse adversary, in "Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms, San Francisco, 25–27 January 1998," pp. 420–425.
 19. N. E. Young, On-line file caching, in "Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms, San Francisco, 25–27 January 1998," pp. 82–86.