

Online Paging with Heterogeneous Cache Slots*

Marek Chrobak^{†1}, Samuel Haney^{‡2}, Mehraneh Liaee^{§3}, Debmalya Panigrahi^{¶4}, Rajmohan Rajaraman^{||3}, Ravi Sundaram^{**3}, and Neal E. Young^{††1}

¹University of California at Riverside; Email:

marek@cs.ucr.edu, neal.young@ucr.edu

²Tumult Labs; Email: sam.m.haney@gmail.com

³Northeastern University; Email: {mehraneh, rraj, koods}@ccs.neu.edu

⁴Duke University; Email: debmalya@cs.duke.edu

Abstract

It is natural to generalize the online k -Server problem by allowing each request to specify not only a point p , but also a subset S of servers that may serve it. To initiate a systematic study of this generalization, we focus on uniform and star metrics. For uniform metrics, the problem is equivalent to a generalization of Paging in which each request specifies not only a page p , but also a subset S of cache slots, and is satisfied by having a copy of p in some slot in S . We call this problem *Slot-Heterogeneous Paging*.

In realistic settings only certain subsets of cache slots or servers would appear in requests. Therefore we parameterize the problem by specifying a family $\mathcal{S} \subseteq 2^{[k]}$ of requestable slot sets, and we establish bounds on the competitive ratio as a function of the cache size k and family \mathcal{S} :

- If all request sets are allowed ($\mathcal{S} = 2^{[k]} \setminus \{\emptyset\}$), the optimal deterministic and randomized competitive ratios are exponentially worse than for standard Paging ($\mathcal{S} = \{[k]\}$).
- As a function of $|\mathcal{S}|$ and k , the optimal deterministic ratio is polynomial: at most $O(k^2|\mathcal{S}|)$ and at least $\Omega(\sqrt{|\mathcal{S}|})$.
- For any laminar family \mathcal{S} of height h , the optimal ratios are $O(hk)$ (deterministic) and $O(h^2 \log k)$ (randomized).
- The special case of laminar \mathcal{S} that we call *All-or-One Paging* extends standard Paging by allowing each request to specify a specific slot to put the requested page in. The optimal deterministic ratio for *weighted All-or-One Paging* is $\Theta(k)$. Offline All-or-One Paging is NP-hard .

Some results for the laminar case are shown via a reduction to the generalization of Paging in which each request specifies a set P of *pages*, and is satisfied by fetching any page from P into the cache. The optimal ratios for the latter problem (with laminar family of height h) are at most hk (deterministic) and hH_k (randomized).

*The conference version of this paper appears in STACS 2023 [25].

[†]Research partially supported by National Science Foundation grants CCF-1536026 and CCF-2153723.

[‡]Research partially supported by National Science Foundation grants CCF-1527084 and CCF-1535972.

[§]Research partially supported by National Science Foundation grants CCF-1535929 and CCF-1909363.

[¶]Research partially supported by National Science Foundation grants CCF-1527084, CCF-1535972, CCF-1750140, CCF-1955703, an Army Research Office grant W911NF2110230, and the Indo-US Joint Center for Algorithms under Uncertainty.

^{||}Research partially supported by National Science Foundation grants CCF-1535929 and CCF-1909363.

^{**}Research partially supported by National Science Foundation grants CCF-1535929 and IIS-2039945.

^{††}Research partially supported by National Science Foundation grant CCF-1619463.

1 Introduction

The standard k -Server and Paging models assume homogenous (interchangeable) servers and cache slots. They don't model applications where servers have different capabilities, nor the fact that modern cache systems often partition the slots, sometimes dynamically, with some parts exclusively accessible by specific processors, cores, processes, threads, or page sets (e.g., [30, 40, 47–49]).

This motivates us to generalize the online k -Server problem to allow each request to specify not only a point p , but also a subset S of servers that may serve it. We call this generalization *Heterogenous k -Server*. To date, only a few special cases of this problem have been studied [22, 44]. Here, following the strategy taken for other hard generalizations of k -Server [6, 7, 12, 23, 31, 39], we initiate a systematic study of this problem by focusing on its restriction to uniform and star metrics. For uniform metrics, the problem is equivalent to a variant of Paging in which each request specifies a page p and a subset S of k cache slots, to be satisfied by having a copy of p in some slot in S . We call this problem *Slot-Heterogenous Paging*. For star metrics the problem reduces to a weighted variant where the cost of retrieving a page is the weight of the page. For reasons discussed below, we parameterize these problems by allowing the requestable sets S to be restricted to an arbitrary but pre-specified family $\mathcal{S} \subseteq 2^{[k]}$. (Restricting to $\mathcal{S} = \{[k]\}$ gives standard Paging and k -Server.) Next is a summary of our results, followed by a summary of related work.

Slot-Heterogenous Paging (Section 3). As we point out, Slot-Heterogenous Paging easily reduces (preserving the competitive ratio) to the Generalized k -Server problem in uniform metrics, for which upper bounds of $k2^k$ and $O(k^2 \log k)$ on the deterministic and randomized ratios are known [7, 12].

- We show that the optimal deterministic and randomized competitive ratios for Slot-Heterogenous Paging are at least $\Omega(2^k/\sqrt{k})$ and $\Omega(k)$, respectively (Theorems 3.2 (i) and 3.3).

Hence, the optimal ratios for Slot-Heterogenous Paging are exponentially worse than for standard Paging. The proofs of Theorems 3.2 and 3.3 employ some novel ideas that may be useful for other problems: the lower bound in Theorems 3.2 (i) uses an adversary argument that requires the construction of a set family not yet studied in the literature, while the proof of Theorem 3.3 is carried out via a reduction *from* standard Paging with a cache of size $\exp(\Theta(k))$.

The large competitive ratios in these lower bounds occur only for instances that use exponentially many distinct request sets S . In realistic settings only certain subsets of cache slots or servers can appear in requests, namely those that represent capabilities or functionalities relevant in a given setting. This motivates us to study the optimal ratios as a function of the cache size k and the family \mathcal{S} of requestable slot sets, and to try to identify natural families that admit more reasonable ratios.

- We show that the optimal deterministic ratio is at most $k^2|\mathcal{S}|$ for any family \mathcal{S} (Theorem 3.1). Theorem 3.2 (ii) shows a complementary lower bound: for infinitely many families \mathcal{S} , every deterministic online algorithm has competitive ratio $\Omega(\sqrt{|\mathcal{S}|})$.

Together Theorems 3.1 and 3.2 (ii) imply that, as a function of $|\mathcal{S}|$ and k , the optimal deterministic ratio for Slot-Heterogenous Paging is polynomial.

Page-Laminar Paging (Section 4). We take a brief detour to consider *Page-Subset Paging*, a natural generalization of Paging in which each request is a set P of *pages* from an arbitrary but fixed family \mathcal{P} , and is satisfiable by fetching any page from P into any slot in the cache. We focus on its special case of *Page-Laminar Paging*, where this set family \mathcal{P} is laminar.

- We show that the optimal deterministic and randomized ratios for Page-Laminar Paging are at most hk and hH_k , where h is the height of the laminar family and $H_k = \sum_{i=1}^k 1/i = \ln k + O(1)$ (Theorem 4.1).

75 The proof is by a reduction that replaces each set request P by a request to one carefully chosen page in P ,
76 yielding an instance of Paging, while increasing the optimal cost by at most a factor of h .

77 **Slot-Laminar Paging (Section 5).** We then return to Slot-Heterogenous Paging, now considering the spe-
78 cific structure of \mathcal{S} , showing better bounds when \mathcal{S} is laminar. This case, which we call *Slot-Laminar Paging*,
79 models applications where slot (or server) capabilities are hierarchical. Laminarity implies that $|\mathcal{S}| < 2k$, so
80 (per Theorem 3.1 above) the optimal deterministic ratio is $O(k^3)$.

- 81 • We show that the optimal deterministic and randomized ratios for Slot-Laminar Paging are $O(h^2k)$
82 and $O(h^2 \log k)$, where $h \leq k$ is the height of \mathcal{S} (Theorem 5.1). We next tighten the deterministic
83 bound to $O(hk)$ (Theorem 5.2).

84 The proof of Theorem 5.1 is via a reduction to Page-Laminar Paging (discussed above), while the proof
85 of Theorem 5.2 refines the generic algorithm from Theorem 3.1. The dependence on k in these bounds is
86 asymptotically tight, as Slot-Laminar Paging generalizes standard Paging.

87 *Reducing Slot-Laminar Paging to Page-Laminar Paging.* The reduction of Slot-Laminar Paging to Page-
88 Laminar Paging in Theorem 5.1 is achieved via a relaxation of Slot-Laminar Paging that drops the constraint
89 that each slot holds at most one page, while still enforcing the cache-capacity constraint of k . This relaxed
90 instance is naturally equivalent to an instance of Page-Laminar Paging. The proof then shows how any
91 solution for the relaxed instance can be “rounded” back to a solution for the original Slot-Laminar Paging
92 instance, losing an $O(h)$ factor in the cost and competitive ratio.

93 **All-or-One Paging (Section 6).** *All-or-One Paging* is the restriction of Slot-Laminar Paging (with height
94 $h = 2$) to $\mathcal{S} = \{[k]\} \cup \{\{j\}\}_{j \in [k]}$. That is, only two types of requests are allowed: *general requests* (allowing
95 the requested page to be anywhere in the cache), and *specific requests* (requiring the page to be in a specified
96 slot). Specific requests don’t give the algorithm any choice, so may appear easy to handle, but in fact make
97 the problem substantially harder than standard Paging. Recent independent work on All-or-One Paging [22]
98 has shown that the optimal deterministic ratio is twice that of Paging, to within an additive constant.

- 99 • We show that the optimal randomized ratio of All-or-One Paging is also at least twice that for Paging
100 (Theorem 6.1), while Theorem 5.1 upper bounds the optimal randomized ratio to within a constant
101 factor of that for Paging. We also show that the offline problem is NIP-hard (Theorem 6.2), in sharp
102 contrast to even k -Server, which can be solved in polynomial time for arbitrary metrics.

103 **Weighted All-Or-One Paging (Section 7).** We initiate a study of Heterogenous k -Server in non-uniform
104 metrics through *Weighted All-Or-One Paging*, which extends All-or-One Paging so that each page has a
105 non-negative weight and the cost of each retrieval is the weight of the page instead of 1.

- 106 • We show that the optimal deterministic ratio for Weighted All-Or-One Paging is $O(k)$, matching the
107 ratio for standard Weighted Paging up to a small constant factor (Theorem 7.1).

108 The algorithm in the proof is implicitly a linear-programming primal-dual algorithm. With this approach
109 the crucial obstacle to overcome is that the standard linear program for standard Weighted Paging does
110 not force pages into specific slots. Indeed, doing so makes the natural integer linear program an NIP-hard
111 multicommodity-flow problem. (Section 7 has an example that illustrates the challenge.) We show how to
112 augment the linear program to partially model the slot constraints.

113 **Related work.** Paging and k -Server have played a central role in the theory of online computation since
114 their introduction in the 1980s [13, 41, 46]. For k -Server, the optimal deterministic ratio is between k and
115 $2k - 1$ [38]. Recent work [29] offers hope for closing this gap, and substantial progress towards resolving the
116 randomized case has been reported in [4, 18]. For Weighted Paging the optimal ratios are k (deterministic)
117 and $\Theta(\log k)$ (randomized) [1, 5, 32, 42, 46].

problem	set family \mathcal{S} (or \mathcal{P})	deterministic	randomized	where
Slot-Heterogenous Paging	$2^{[k]} \setminus \{\emptyset\}$	$\leq k2^k$	$\leq O(k^2 \log k)$	via [7, 12]
”	arbitrary \mathcal{S}	$\leq k \min(\mathcal{S}^* , \text{mass}(\mathcal{S}))$		Thm. 3.1
One-of- m Paging, $m \approx k/2$	$\binom{[k]}{m}$	$\geq \Omega(2^k/\sqrt{k})$	$\geq \Omega(k)$	Thms. 3.2(i), 3.3
One-of- m Paging, any m	$\binom{[k]}{m}$	$\gtrsim \Omega((4k/m)^{m/2}/\sqrt{m})$		Thm. 3.2(ii)
Slot-Laminar Paging	laminar \mathcal{S} , height h	$\leq (2h - 1)k$	$\leq 3h^2 H_k$	Thms. 5.1, 5.2
All-or-One Paging	$\{[k]\} \cup \{\{s\} : s \in [k]\}$	$\geq 2k - 1$	$\geq 2H_k - 1$	[22, 35], Thm. 6.1
”	”	$\leq 2k + 14$		[22]
Weighted All-Or-One Paging	$\{[k]\} \cup \{\{s\} : s \in [k]\}$	$\leq O(k)$		Thm. 7.1
Page-Subset Paging restricted to $\mathcal{P} = \binom{\text{all pages}}{m}$		$\geq \binom{k+m}{k} - 1$		[31]
”		$\leq k(\binom{k+m}{k} - 1)$	$\leq O(k^3 \log m)$	[23]
Page-Laminar Paging	\mathcal{P} laminar, height h	$\leq hk$	$\leq hH_k$	Thm. 4.1

Table 1: Summary of upper (\leq) and lower (\geq) bounds on optimal competitive ratios. Here $\text{mass}(\mathcal{S}) = \sum_{S \in \mathcal{S}} |S|$ and $\mathcal{S}^* = \bigcup_{S \in \mathcal{S}} S$. The lower bound for One-of- m Paging holds for some but not all m and k —see Theorem 3.2(ii). The upper bound for Slot-Laminar Paging in the deterministic case (Theorem 5.1) is in fact $2 \cdot \text{mass}(\mathcal{S}) - k$, which is at most $(2h - 1)k$. Also, offline All-or-One Paging and its generalizations are \mathbb{NP} -hard (Theorem 6.2), as is offline Page-Subset Paging ([23]).

118 *Restricted Caching* is one previously studied model with *heterogenous* cache slots. It is the restriction
119 of Slot-Heterogenous Paging in which each page p has one *fixed* set $S_p \subseteq [k]$ of slots, and each request
120 to p requires p to be in some slot in S_p . For this problem the optimal randomized ratio is $O(\log^2 k)$ [20].
121 Better bounds are possible given further restrictions on the sets, as in *Companion Caching*, which models a
122 cache partitioned into set-associative and fully associative parts [16, 17, 33, 43]. It is natural to ask whether
123 *Restricted k -Server*—the restriction of Heterogenous k -Server that requires each point p to be served by a
124 server in a *fixed* set S_p —is easier than Heterogenous k -Server. While the two problems are different for
125 many metric classes, they can be shown to be equivalent in metric spaces with no isolated points, such
126 as Euclidean spaces. The \mathbb{NP} -hardness result for Restricted Caching from [17] implies that offline Slot-
127 Heterogenous Paging with $\mathcal{S} = \{\{s, k\} : s \in [k - 1]\}$ is \mathbb{NP} -hard.

128 Other sophisticated online caching models include *Snoopy Caching*, in which multiple processors each
129 have their own cache and coordinate to maintain consistency across writes [37], *Multi-Level Caching*, where
130 the cost to access a slot depends on the slot [26], and *Writeback-Aware Caching*, where each page has
131 multiple copies, each with a distinct level and weight, and each request specifies a page and a level, and can
132 be satisfied by fetching a copy of this page at the given or a higher level [8, 9]. (This is a special case of
133 weighted Page-Laminar Paging where \mathcal{P} consists of pairwise-disjoint chains.) *Multi-Core Caching* models
134 the fact that faults can change the request sequence (e.g. [36]).

135 Patel’s master thesis [44] studies Heterogenous k -Server with just two types of requests—general re-
136 quests (to be served by any server) and specialized requests (to be served by any server in a fixed subset S'
137 of “specialized” servers)—and bounds the optimal ratios for uniform metrics and the line. Recent indepen-
138 dent work on deterministic algorithms for online All-or-One Paging establishes a $2k - 1$ lower bound and
139 a $2k + 14$ upper bound [22]. Earlier work in [35] presents a $2k - 1$ lower bound and a $3k$ upper bound on
140 deterministic algorithms.

141 Heterogenous k -Server reduces (see Section 3) to the *Generalized k -Server* problem, in which each
142 server moves in its own metric space, each request specifies one point in each space, and the request is
143 satisfied by moving any one server to the requested point in its space [39]. For uniform metrics, the op-
144 timal competitive ratios for this problem are between 2^k and $k2^k$ (deterministic) and between $\Omega(k)$ and

145 $O(k^2 \log k)$ (randomized) [7, 12]. These ratios are exponentially worse than the ratios for standard k -Server.
 146 Heterogenous k -Server, parameterized by \mathcal{S} , provides a spectrum of problems that bridges the two extremes.

147 *Weighted k -Server* is a restriction of Generalized k -Server in which servers move in the same metric
 148 space but have different weights, and the cost is the weighted distance [34]. For this problem (in non-
 149 uniform metrics) the deterministic and randomized ratios are at least (respectively) doubly exponential [6, 7]
 150 and exponential [3, 24].

151 For Page-Subset Paging restricted to m -element sets of pages, the optimal ratios are between $\binom{k+m}{k} - 1$
 152 and $k \binom{k+m}{k} - 1$ (deterministic) and between $\Omega(\log km)$ and $O(k^3 \log m)$ (randomized) [23, 31]. This
 153 problem has been studied as uniform *Metrical Service Systems with Multiple Servers (MSSMS)*. MSSMS is
 154 the generalization of k -Server where each request is a *set* of points, one of which needs to be covered by
 155 some server.

156 The *k -Chasing* problem extends k -Server by having each request P be a convex subset of \mathbb{R}^d , to be
 157 satisfied by moving any server to any point in P [19]. For k -Chasing, no online algorithm is competitive
 158 even for $d = k = 2$ [19], while for $k = 1$ the ratios grow with d [2, 45].

159 In the *k -Taxi problem* each request (p, q) requires any server to move to p then (for free) to q . For this
 160 problem the optimal ratios are exponentially worse than for standard k -Server [21, 28].

161 2 Formal Definitions

162 **Slot-Heterogenous Paging.** A problem instance consists of a set $[k] = \{1, 2, \dots, k\}$ of cache slots, a family
 163 $\mathcal{S} \subseteq 2^{[k]} \setminus \{\emptyset\}$ of requestable slot sets, and a request sequence $\sigma = \{\sigma_t\}_{t=1}^T$, where each request has the
 164 form $\sigma_t = \langle p_t, S_t \rangle$ for some *page* p_t and set $S_t \in \mathcal{S}$. A *cache configuration* C is a function that specifies
 165 the content of each slot $s \in [k]$; this content is either a single page (said to be assigned to the slot) or empty.
 166 Configuration C is said to *satisfy* a request $\langle p, S \rangle$ if it assigns page p to at least one slot in S . A *solution*
 167 for a given request sequence σ is a sequence $\{C_t\}_{t=1}^T$ of cache configurations such that, for each $t \in [T]$,
 168 C_t satisfies request σ_t . The objective is to minimize the number of *retrievals*, where a page p is retrieved in
 169 slot s at time t if C_t assigns p to s , but C_{t-1} does not (or $t = 1$). An event when C_{t-1} does not assign p_t to
 170 any slot in S_t is called a *fault*. Obviously a fault triggers a retrieval but, while this is not strictly necessary, it
 171 is convenient to also allow an algorithm to make spontaneous retrievals, either by fetching into the cache a
 172 non-requested page or by moving pages within the cache.

173 **Slot-Laminar Paging.** This is the restriction of Slot-Heterogenous Paging to instances where \mathcal{S} is *laminar*:
 174 every pair $R, R' \in \mathcal{S}$ of sets is either disjoint or nested. (This implies $|\mathcal{S}| \leq 2k$.) A laminar family \mathcal{S} can
 175 be naturally represented by a forest (a collection of disjoint trees), with a set R being a descendant of R' if
 176 $R \subseteq R'$. When discussing Slot-Laminar Paging we will routinely use tree-related terminology; for example,
 177 we will refer to some sets in \mathcal{S} as leaves, roots, or internal nodes. The height h of a laminar family \mathcal{S} is one
 178 more than the maximum height of a tree in \mathcal{S} , that is the maximum h for which \mathcal{S} contains a sequence of h
 179 strictly nested sets: $R_1 \subsetneq R_2 \subsetneq \dots \subsetneq R_h$.

180 **All-or-One Paging.** This is the restriction of Slot-Laminar Paging to instances with $\mathcal{S} = \{[k]\} \cup \{\{j\}\}_{j \in [k]}$.
 181 That is, there are two types of requests: *general*, of the form $\langle p, [k] \rangle$, requiring page p to be in at least one
 182 slot of the cache, and *specific*, of the form $\langle p, \{j\} \rangle$, $j \in [k]$, requiring page p to be in slot j . For convenience,
 183 $\langle p, * \rangle$ is a synonym for $\langle p, [k] \rangle$, while $\langle p, j \rangle$ is a synonym for $\langle p, \{j\} \rangle$.

184 **Weighted All-Or-One Paging.** This is the natural extension of All-or-One Paging in which each page p is
 185 assigned a non-negative weight $\text{wt}(p)$, and the cost of retrieving p is $\text{wt}(p)$ instead of 1.

186 **One-of- m Paging.** This is the restriction of Slot-Heterogenous Paging to instances with $\mathcal{S} = \binom{[k]}{m} = \{S \subseteq$
 187 $[k] : |S| = m\}$, that is, every request specifies a set of m slots.

188 **Page-Subset Paging.** An instance consists of k cache slots, a collection \mathcal{P} of requestable sets of pages, and
 189 a request sequence $\pi = \{P_t\}_{t=1}^T$, where each P_t is drawn from \mathcal{P} . A solution is a sequence $\{C_t\}_{t=1}^T$ of
 190 cache configurations (as previously defined) such that, at each time $t \in [T]$, C_t assigns at least one page
 191 in P_t to at least one slot. The objective is to minimize the number of retrievals. (Slots are interchangeable
 192 here, so a cache configuration could be defined as a multiset of at most k pages, but using slot assignments
 193 is technically more convenient.)

194 **Page-Laminar Paging.** This is the restriction of Page-Subset Paging to instances where \mathcal{P} is *laminar*.

195 **Generalized k -Server.** In this variant of k -Server, each server moves in its own metric space; each request
 196 specifies one point in each space, and the request is satisfied by moving any one server to the requested point
 197 in its space [39].

198 **Approximation algorithms.** An algorithm \mathbb{A} for a given cost minimization problem is called a c -*approximation*
 199 *algorithm* if, for each instance σ , \mathbb{A} satisfies $\text{cost}_{\mathbb{A}}(\sigma) \leq c \cdot \text{opt}(\sigma) + b$, where $\text{cost}_{\mathbb{A}}(\sigma)$ is the cost of \mathbb{A} on
 200 σ , $\text{opt}(\sigma)$ is the optimum cost of σ , and b is a constant independent of σ . We follow the standard convention
 201 that when we are considering \mathbb{A} as an offline algorithm, the constant b must be 0.

202 **Online algorithms and competitive ratio.** In the online variants of the paging problems studied in this
 203 paper the requests arrive online, one per time step, and an online algorithm needs to satisfy each request
 204 before the next one is revealed. To simplify presentation we assume that the algorithm knows the underlying
 205 set family \mathcal{S} (or \mathcal{P}), but many of our algorithms work (or can be adapted to work) without knowing the set
 206 family in advance. An online algorithm \mathbb{A} is called c -*competitive* if \mathbb{A} is a c -approximation algorithm. As
 207 common in the literature, we will use the term “optimal deterministic (resp. randomized) competitive ratio”
 208 to refer to the optimal ratio of of a deterministic (resp. randomized) online algorithm for the given problem.

209 3 Slot-Heterogenous Paging

210 Any instance of Slot-Heterogenous Paging can be reduced to an instance of Generalized k -Server in uniform
 211 spaces, as follows. Represent each cache slot by a server in a uniform metric space whose points are the
 212 pages, then simulate each request $\langle p, S \rangle$ by a sufficiently long sequence of requests, each of which specifies
 213 point p for each server in S and alternates between two different points for the remaining servers, in $[k] \setminus S$.
 214 Composing this reduction with the upper bounds from [7] yields immediate upper bounds of $O(k2^k)$ and
 215 $O(k^3 \log k)$ on the deterministic and randomized ratios for unrestricted Slot-Heterogenous Paging (that is,
 216 with $\mathcal{S} = 2^{[k]} \setminus \{\emptyset\}$).

217 Theorem 3.2 (i) (Section 3.2) and Theorem 3.3 (Section 3.3) show that these bounds are tight within
 218 $\text{poly}(k)$ factors: the optimal ratios are at least $\Omega(2^k/\sqrt{k})$ and $\Omega(k)$, respectively. But restricting \mathcal{S} allows
 219 better ratios: Theorem 3.1 (Section 3.1) shows an upper bound of $k^2|\mathcal{S}|$ on the optimal deterministic ratio for
 220 any family \mathcal{S} . For One-of- m Paging, Theorem 3.1 and Theorem 3.2 (ii) imply that the optimal deterministic
 221 ratio is $O(k^{m+1})$ and $\Omega((4k/m)^{m/2}/\sqrt{m})$.

222 3.1 Upper bounds for deterministic Slot-Heterogenous Paging

223 This section gives upper bounds on the optimal deterministic competitive ratio for Slot-Heterogenous Paging
 224 with any slot-set family \mathcal{S} , as a function of $\text{mass}(\mathcal{S}) = \sum_{S \in \mathcal{S}} |S| \leq k|\mathcal{S}|$ and $|\mathcal{S}^*|$, where $\mathcal{S}^* = \bigcup_{S \in \mathcal{S}} 2^S$.
 225 The first bound follows from an easy counting argument. The second bound uses a refinement of the rank
 226 method of [7], which bounds the number of steps of a natural exhaustive-search algorithm by the rank of a
 227 certain upper-triangular matrix.

228 **Theorem 3.1.** Fix any $\mathcal{S} \subseteq 2^{[k]} \setminus \{\emptyset\}$. The competitive ratio of Algorithm EXHSEARCH in Figure 1 for
 229 Slot-Heterogenous Paging with requestable sets from \mathcal{S} is at most $k \cdot \min\{|\mathcal{S}^*|, \text{mass}(\mathcal{S})\}$.

input: Slot-Heterogenous Paging instance $(k, \mathcal{S}, \sigma = (\sigma_1, \dots, \sigma_T))$

1. let the initial cache configuration C_0 be arbitrary; let $\ell \leftarrow 1$ *— ℓ is the start of the current phase*
2. for each time $t \leftarrow 1, 2, \dots, T$:
 - 2.1. if current configuration C_{t-1} satisfies request σ_t : ignore the request (set $C_t = C_{t-1}$)
 - 2.2. else:
 - 2.2.1. if any configuration satisfies all requests $\sigma_\ell, \sigma_{\ell+1}, \dots, \sigma_t$: let C_t be any such configuration
 - 2.2.2. else: let $\ell \leftarrow t$; let C_t be any configuration satisfying σ_t *— start the next phase*

251 That is, $M = AB$; A and B (and M) have rank at most $|\mathcal{S}^*|$. And M has rank L , so $L \leq |\mathcal{S}^*|$. To
252 bound the optimum cost, consider any phase other than the last. Let t' and t'' be the start and end times.
253 Suppose for contradiction that the optimal solution incurs no cost (has no retrievals) during $[t' + 1, t'' + 1]$.
254 Then its configuration at time t' satisfies all requests in $[t', t'' + 1]$, contradicting the algorithm's condition
255 for terminating the phase. So the optimal solution pays at least 1 per phase (other than the last). In any phase
256 of length L the algorithm pays at most kL (at most k per step). This and the two upper bounds on L imply
257 Theorem 3.1. \square

258 3.2 Lower bounds for deterministic Slot-Heterogenous Paging

259 We establish our lower bounds for Slot-Heterogenous Paging and One-of- m Paging given in Table 1.

260 **Theorem 3.2.** (i) For all odd k , the optimal deterministic ratio for One-of- m Paging with $m = (k + 1)/2$
261 is at least $\binom{k}{m} = \Omega(2^k/\sqrt{k})$. For all k , the optimal ratio with $m = \lfloor (k + 1)/2 \rfloor$ is $\Omega(2^k/\sqrt{k})$. (ii) For any
262 even $m \geq 2$ and any $k > m$ that is an odd multiple of $m - 1$, the optimal deterministic ratio for One-of- m
263 Paging is at least $\binom{m-1}{m/2} \binom{k}{m-1}^{m/2} = \Theta((4k/m)^{m/2}/\sqrt{m}) = \Omega(\sqrt{|\mathcal{S}|})$, where $\mathcal{S} = \binom{[k]}{m}$.

264 Before proving Theorem 3.2, we prove Lemma 3.1. It states that for any \mathcal{S} the existence of a family
265 $\mathcal{Z} \subseteq 2^{[k]}$ with certain properties implies a lower bound of $|\mathcal{Z}|$ on the competitive ratio. The proof of the
266 theorem then constructs such families \mathcal{Z} for appropriate families \mathcal{S} of requestable sets. Throughout this
267 section \overline{X} denotes the complement of set $X \subseteq [k]$, that is $\overline{X} = [k] \setminus X$.

268 **Lemma 3.1.** For some $\mathcal{S} \subseteq 2^{[k]}$, suppose there are two set families $\mathcal{G} \subseteq \mathcal{S}$ and $\mathcal{Z} \subseteq 2^{[k]}$ such that

269 (gz0) For each $X \subseteq [k]$ there is $S \in \mathcal{G}$ such that $S \subseteq X$ or $S \subseteq \overline{X}$.

270 (gz1) If $Z \in \mathcal{Z}$ then $\overline{Z} \notin \mathcal{Z}$.

271 (gz2) For each $S \in \mathcal{G}$ there is $Y \in \mathcal{Z}$ such that $S \not\subseteq Y$ and $S \not\subseteq \overline{Y}$ for all $Z \in \mathcal{Z} \setminus \{Y\}$.

272 Then the optimal deterministic ratio for Slot-Heterogenous Paging with family \mathcal{S} is at least $|\mathcal{Z}|$.

273 *Proof.* The proof is an adversary argument based on the following idea. At each step, the adversary chooses
274 a request that forces the algorithm to fault but causes at most two faults total among a fixed set of $2|\mathcal{Z}|$
275 other solutions. At the end, the algorithm's total cost is at least $|\mathcal{Z}|$ times the average cost of these other
276 solutions, so its competitive ratio is at least $|\mathcal{Z}|$. This general approach is common for lower bounds on
277 deterministic online algorithms (see e.g. lower bounds on the optimal ratios for k -Server [41], for Metrical
278 Task Systems [15] and for Generalized k -Server on uniform metrics [39]).

279 Here are the details. Let \mathbb{A} be any deterministic online algorithm for Slot-Heterogenous Paging with
280 slot-set family \mathcal{S} . The adversary will request just two pages, p_0 and p_1 . For a set $X \subseteq [k]$, let C_X denote the
281 cache configuration where the slots in X contain p_0 and the slots in \overline{X} contain p_1 . Without loss of generality
282 assume that each slot of \mathbb{A} 's cache always holds p_0 or p_1 —its cache configuration is C_X for some X .

283 At each step, if the current configuration of \mathbb{A} is C_X , the adversary chooses $S \in \mathcal{G}$ such that either
284 $S \subseteq X$ or $S \subseteq \overline{X}$. (Such an S exists by Property (gz0).) If $S \subseteq X$, then all slots in S hold p_0 , and the
285 adversary requests $\langle p_1, S \rangle$, causing a fault. Otherwise, $S \subseteq \overline{X}$, so all slots in S hold p_1 . In this case the
286 adversary requests $\langle p_0, S \rangle$, causing a fault. The adversary repeats this K times, where K is arbitrarily large.
287 Since \mathbb{A} faults at each step, the overall cost of \mathbb{A} is at least K .

288 It remains to bound the optimal cost. Let $\tilde{\mathcal{Z}} = \{\overline{Z} : Z \in \mathcal{Z}\}$. By (gz1), we have $\tilde{\mathcal{Z}} \cap \mathcal{Z} = \emptyset$. For each
289 $Z \in \mathcal{Z} \cup \tilde{\mathcal{Z}}$ define a solution called the Z -strategy, as follows. The solution starts in configuration C_Z . It
290 stays in C_Z for the whole computation, except that on requests $\langle p_a, S \rangle$ that are not served by C_Z (that is,
291 when all slots of C_Z in S contain p_{1-a}), it retrieves p_a to any slot $j \in S$, serves the request, then retrieves
292 p_{1-a} back into slot j , restoring configuration C_Z . This costs 2.

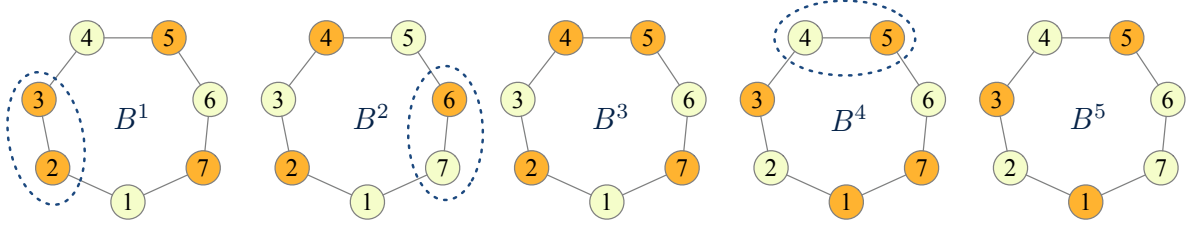


Figure 2: Illustration of the proof of Theorem 3.2 Part (ii) for $k = 35$, $m = 6$, and $\ell = 7$. The figure shows the partition of all slots into $m - 1 = 5$ sets B^1, \dots, B^5 , each represented by a cycle. To avoid clutter, each slot b_c^e is represented by its index c within B^e . The picture shows set $S = \{b_2^1, b_3^1, b_6^2, b_7^2, b_4^4, b_5^4\} \in \mathcal{G}$, marked by dashed ovals. It also shows $Z_{S'} \in \mathcal{Z}$, represented by orange/shaded circles, for $S' = \{b_2^1, b_3^1, b_4^3, b_5^3, b_7^4, b_1^4\}$.

293 We next observe that in each step at most one Z -strategy faults (and pays 2). Assume that the request
 294 at a given step is to p_0 (the case of a request to p_1 is symmetric). Let this request be $\langle p_0, S \rangle$, where $S \in \mathcal{G}$.
 295 Let $Y \subseteq [k]$ be the set from Property (gz2). For all $Z \in (\mathcal{Z} \cup \tilde{\mathcal{Z}}) \setminus \{Y, \bar{Y}\}$, then, $S \cap Z \neq \emptyset$, implying
 296 that configuration C_Z has a slot in S that contains p_0 —in other words, configuration C_Z satisfies S . Also,
 297 either $S \cap Y \neq \emptyset$ or $S \cap \bar{Y} \neq \emptyset$, so one of the configurations C_Y or $C_{\bar{Y}}$ also satisfies S . Therefore only one
 298 Z -strategy (Y or \bar{Y}) might not satisfy S . So, in each step, at most one Z -strategy faults (and pays 2).

299 Thus the combined total cost for all Z -strategies (not counting the cost of at most k for moving to Z at
 300 the beginning) is at most $2K$. There are $2|\mathcal{Z}|$ such strategies, so their average cost is at most $(2K + k)/2|\mathcal{Z}|$.
 301 The cost of \mathbb{A} is at least K , so the ratio is at least $\frac{K}{(2K+k)/2|\mathcal{Z}|} = \frac{|\mathcal{Z}|}{1+k/2K}$. Taking K arbitrarily large, the
 302 lemma follows. \square

303 *Proof of Theorem 3.2. Part (i).* Recall that $m = \lfloor (k+1)/2 \rfloor$. First consider the case when k is odd. Apply
 304 Lemma 3.1, taking both \mathcal{G} and \mathcal{Z} to be $\binom{[k]}{m}$. Properties (gz0) and (gz1) follow directly from k being odd
 305 and the definitions of \mathcal{G} and \mathcal{Z} . Property (gz2) also holds with $Y = S$. (For any $S \in \mathcal{G}$, every $Z \in \mathcal{Z}$
 306 satisfies $|Z| = |S| > |\bar{Z}|$, so $S \not\subseteq \bar{Z}$, while $S \subseteq Z$ implies $Z = S$.) Thus, by Lemma 3.1, the ratio is at least
 307 $|\mathcal{Z}| = \binom{k}{(k+1)/2} = \Omega(2^k/\sqrt{k})$. This proves Part (i) for odd k .

308 For even k , let $k' = k - 1$. Then apply Part (i) to k' using just cache slots in $[k']$, that is, using slot-set
 309 family $\mathcal{S}' = \binom{[k']}{m} \subseteq \binom{[k]}{m} = \mathcal{S}$, with slot k playing no role as it is never requested. This proves Part (i).

310 *Part (ii).* Fix such an m and k . Let $\ell = k/(m-1)$ so $\ell \geq 3$ is odd. Recall that $\mathcal{S} = \binom{[k]}{m}$
 311 is the family of requestable slot sets. Partition $[k]$ arbitrarily into $m-1$ disjoint subsets B^1, B^2, \dots, B^{m-1} ,
 312 each of cardinality ℓ . For each B^e , order its slots arbitrarily as $B^e = \{b_1^e, b_2^e, \dots, b_\ell^e\}$. For any index
 313 $c \in \{1, 2, \dots, \ell\}$ and an integer i , let $c \oplus i$ denote $((c+i-1) \bmod \ell) + 1$. In other words, we view each B^e
 314 as an odd-length cycle, and this cyclic structure is important in the proof. Any consecutive pair $\{b_c^e, b_{c \oplus 1}^e\}$ of
 315 slots on this cycle is called an *edge*. Thus each cycle B^e has ℓ edges.

316 First we define $\mathcal{G} \subseteq \mathcal{S}$ for Lemma 3.1. The sets S in \mathcal{G} are those obtainable as follows: choose any $m/2$
 317 edges, no two from the same cycle, then let S contain the m slots in those $m/2$ chosen edges. (The six slots
 318 inside the three dashed ovals in Figure 2 show one S in \mathcal{G} .) This set of $m/2$ edges uniquely determines S ,
 319 and vice versa.

320 We verify that \mathcal{G} has Property (gz0) from Lemma 3.1. Indeed, consider any $X \subseteq [k]$. Call the slots in X
 321 *white* and the slots in \bar{X} *black*. Each cycle B^e has odd length, so has an edge $\{b_c^e, b_{c \oplus 1}^e\}$ that is white (with
 322 two white slots) or black (with two black slots). So either (i) at least half the cycles have a white edge, or (ii)
 323 at least half have a black edge. Consider the first case (the other is symmetric). There are $m-1$ cycles, and
 324 m is even, so at least $m/2$ cycles have a white edge. So there are $m/2$ white edges with no two in the same
 325 cycle. The set S comprised of the m white slots from those edges is in \mathcal{G} , and is contained in X (because its
 326 slots are white). So \mathcal{G} has Property (gz0).

327 Next we define $\mathcal{Z} \subseteq 2^{[k]}$ for Lemma 3.1. The set \mathcal{Z} contains, for each set $S' \in \mathcal{G}$, one set $Z_{S'}$, defined
328 as follows. For each of the $m/2$ cycles B^e having an edge $\{b_c^e, b_{c \oplus 1}^e\}$ in S' , add to $Z_{S'}$ the two slots on that
329 edge, together with the $(\ell - 3)/2$ slots $b_{c \oplus 3}^e, b_{c \oplus 5}^e, \dots, b_{c \oplus (\ell - 2)}^e$. For each of the $m/2 - 1$ remaining cycles
330 B^e , add to $Z_{S'}$ the $(\ell - 1)/2$ slots $b_1^e, b_3^e, \dots, b_{\ell - 2}^e$. (The orange/shaded slots in Figure 2 show one set $Z_{S'}$
331 in \mathcal{Z} .) Then $Z_{S'}$ contains exactly $m/2$ edges (the ones in S') while its complement $\overline{Z}_{S'}$ contains exactly
332 $m/2 - 1$ edges (one from each cycle with no edge in S'). This implies Property (gz1). Note that $Z_{S'} \neq Z_{S''}$
333 for different sets $S', S'' \in \mathcal{G}$.

334 Next we show Property (gz2). Given any set $S \in \mathcal{G}$, let $Y = Z_S \in \mathcal{Z}$. Consider any $Z_{S'} \in \mathcal{Z}$ such
335 that $S \subseteq Z_{S'}$ or $S \subseteq \overline{Z}_{S'}$. We need to show $Z_{S'} = Z_S$, i.e., $S' = S$. It cannot be that $S \subseteq \overline{Z}_{S'}$, because
336 S contains $m/2$ edges, whereas $\overline{Z}_{S'}$ contains $m/2 - 1$ edges. So $S \subseteq Z_{S'}$. But S and $Z_{S'}$ each contain
337 exactly $m/2$ edges, which therefore must be the same. It follows from the definition of $Z_{S'}$ that $S' = S$. So
338 Property (gz2) holds.

339 So \mathcal{G} and \mathcal{Z} have Properties (gz0)-(gz2) from Lemma 3.1. Directly from definition we have $|\mathcal{Z}| = |\mathcal{G}|$,
340 while $|\mathcal{G}| = \binom{m-1}{m/2} \ell^{m/2}$ because there are $\binom{m-1}{m/2}$ ways to choose $m/2$ distinct cycles, and then for each
341 of these $m/2$ cycles there are ℓ ways to choose one edge. Lemma 3.1 and $\ell = k/(m - 1)$ imply that the
342 optimal deterministic ratio is at least $f(m, k) = \binom{m-1}{m/2} (k/(m - 1))^{m/2}$. To complete the proof of part (ii)
343 we lower-bound $f(m, k)$. We observe that

$$4^m = \Omega(\sqrt{m} (k/(k - m))^{k-m+1/2}). \quad (1)$$

This can be verified by considering two cases: If $k \geq m + 2$ then, using $1 + z \leq e^z$, we have $\sqrt{m} (k/(k - m))^{k-m+1/2} = \sqrt{m}(1 + m/(k - m))^{k-m+1/2} \leq \sqrt{m} \cdot e^{5m/4} \leq 2 \cdot 4^m$, for all $m \geq 1$. In the remaining case, for $k = m + 1$, we have $\sqrt{m}(k/(k - m))^{k-m+1/2} = \sqrt{m}(1 + m)^{3/2} \leq 2 \cdot 4^m$. Thus (1) indeed holds. Now, recalling that $f(m, k) = \binom{m-1}{m/2} (k/(m - 1))^{m/2}$, we derive

$$\begin{aligned} f(m, k) &= \Theta((2^m/\sqrt{m}) \cdot (k/(m - 1))^{m/2}) && \text{(Stirling's approximation)} \\ &= \Theta((4k/m)^{m/2} \cdot (1 + 1/(m - 1))^{m/2}/\sqrt{m}) && \text{(rewriting)} \\ &= \Theta((4k/m)^{m/2}/\sqrt{m}) && ((1 + 1/(m - 1))^{m/2} \leq e) \end{aligned} \quad (2)$$

This gives us one estimate on the competitive ratio in Theorem 3.2(ii). To obtain a second estimate, squaring both sides of (2), we obtain

$$\begin{aligned} f(m, k)^2 &= \Omega((4k/m)^m/m) = \Omega((k/m)^m \cdot 4^m/m) \\ &= \Omega((k/m)^m \cdot (k/(k - m))^{k-m+1/2}/\sqrt{m}) && \text{(using (1))} \\ &= \Omega(\binom{k}{m}) = \Omega(|\mathcal{S}|) && \text{(Stirling's approximation)} \end{aligned}$$

344 Therefore $f(m, k) = \Omega(\sqrt{|\mathcal{S}|})$, as claimed, completing the proof of Theorem 3.2(ii). \square

345 3.3 Lower bound for randomized Slot-Heterogenous Paging

346 Next we present a lower bound on the optimal competitive ratio for randomized algorithms:

347 **Theorem 3.3.** *The optimal randomized ratio for One-of- m Paging with $m = \lfloor k/2 \rfloor$ is $\Omega(k)$.*

348 The proof is by a reduction from standard Paging with some N pages and a cache of size $N - 1$. For
349 any N , this problem has optimal randomized competitive ratio $H_{N-1} = \Theta(\log N)$ [32]. This and the next
350 lemma imply the theorem.

351 **Lemma 3.2.** *Every $f(k)$ -competitive (randomized) online algorithm \mathbb{A} for One-of- m Paging with $m =$
352 $\lfloor k/2 \rfloor$ can be converted into an $O(f(k))$ -competitive (randomized) online algorithm \mathbb{B} for standard Paging
353 with N pages and a cache of size $N - 1$, where $N = 2^{\Theta(k)}$.*

354 *Proof.* Fix a sufficiently large k . Assume without loss of generality that k is even (otherwise apply the
355 construction below to slots in $[k - 1]$, ignoring slot k as it is never requested). Take $N = \lfloor e^{k/16} \rfloor$.

356 To ease exposition, view the Paging problem with N pages and a cache of size $N - 1$ as the following
357 equivalent online *Cat and Rat* game on any set \mathcal{H} of N holes (see e.g. [14, §11.3]). The input is a sequence
358 $\mu = (R_0, C_1, \dots, C_T)$ of holes (i.e., $R_0 \in \mathcal{H}$ and $C_t \in \mathcal{H}$ for all t). A solution is any sequence (R_1, \dots, R_T)
359 of holes such that $R_t \neq C_t$ for all $t \in [T]$. Informally, at each time $t \in [T]$, the cat inspects hole C_t , and
360 if the rat's hole R_{t-1} at time $t - 1$ was C_t , the rat is required to move to some other hole $R_t \in \mathcal{H} \setminus \{C_t\}$.
361 (For the solution to be online, R_t must be independent of $C_{t+1}, C_{t+2}, \dots, C_T$ for all $t \geq 1$.) The goal is to
362 minimize the number of times the rat moves, that is $|\{t \in [T] : R_t \neq R_{t-1}\}|$.

363 The claimed algorithm \mathbb{B} for Paging will work by reducing a given instance μ on a set \mathcal{H} of N holes to
364 an instance σ of One-of- $(k/2)$ Paging, simulating \mathbb{A} on σ , and converting its solution to a solution for μ .
365 This instance σ uses just two pages, p_0 and p_1 . To describe the reduction, we need a few more definitions
366 and observations.

367 For two disjoint sets $S_0, S_1 \subseteq [k]$, by $Q(S_0, S_1)$ we denote the cache configuration that assigns p_0 to slots
368 in S_0 and p_1 to slots in S_1 , with the remaining slots empty. A configuration $Q(S_0, S_1)$ is called *balanced*
369 if $|S_0| = |S_1| = k/2$. (This obviously implies that $\bar{S}_0 = S_1$, where $\bar{S}_0 = [k] \setminus S_0$.) The following easy
370 observation will be useful:

371 **Observation 3.3.** *Let $S \subseteq [k]$ with $|S| = k/2$. Any request $\langle p_0, S \rangle$ is satisfied by every balanced configura-
372 tion except $Q(\bar{S}, S)$, and any request $\langle p_1, S \rangle$ is satisfied by every balanced configuration except $Q(S, \bar{S})$.*

373 For any set \mathcal{C} of cache configurations, a *forcing sequence* for \mathcal{C} is a request sequence such that the cache
374 configurations that satisfy all requests in this sequence without cost are exactly those in \mathcal{C} .

375 **Claim 3.4.** *Let \mathcal{C} be a set of balanced cache configuration such that any two configuration in \mathcal{C} are at
376 distance at least 3. Then there is a request sequence $\psi(\mathcal{C})$ that is forcing for \mathcal{C} .*

377 To verify the claim, take $\psi(\mathcal{C})$ to be the sequence formed by (any ordering of) all those allowed requests
378 that are satisfied by all configurations in \mathcal{C} . Specifically, $\psi(\mathcal{C})$ consists of all requests $\langle p_i, S \rangle$ such that
379 $i \in \{0, 1\}$ and $|S| = k/2$, with $S \cap S_i \neq \emptyset$ for all $Q(S_0, S_1) \in \mathcal{C}$. By definition, each configuration in \mathcal{C}
380 satisfies all requests in $\psi(\mathcal{C})$.

381 It remains to show that for any configuration $Q(S'_0, S'_1) \notin \mathcal{C}$ there is a request in $\psi(\mathcal{C})$ not satisfied
382 by $Q(S'_0, S'_1)$. In the case that $Q(S'_0, S'_1)$ is balanced, we can take this request to be $\langle p_1, S'_0 \rangle$, because, by
383 Observation 3.3, it is included in $\psi(\mathcal{C})$ but it is not satisfied by $Q(S'_0, S'_1)$.

384 Next consider the case that $Q(S'_0, S'_1)$ is not balanced. Assume without loss of generality that $|S'_0| > k/2$.
385 Let B_0 and B'_0 be any two size- $k/2$ subsets of S'_0 such that $|B_0 \setminus B'_0| = |B'_0 \setminus B_0| = 1$. Then $Q(B_0, \bar{B}_0)$
386 and $Q(B'_0, \bar{B}'_0)$ are at Hamming distance 2 so both cannot be in \mathcal{C} . Assume without loss of generality
387 that $Q(B_0, \bar{B}_0)$ is not in \mathcal{C} . Then request $\langle p_1, B_0 \rangle$ is satisfied by every configuration in \mathcal{C} because, by
388 Observation 3.3, the only balanced configuration that doesn't satisfy this request is $Q(B_0, \bar{B}_0)$, which is
389 not in \mathcal{C} . So $\langle p_1, B_0 \rangle$ is in $\psi(\mathcal{C})$. On the other hand, since $B_0 \subseteq S'_0$, request $\langle p_1, B_0 \rangle$ is not satisfied by
390 $Q(S'_0, S'_1)$. This proves Claim 3.4.

391 We now describe how Algorithm \mathbb{B} computes its solution $R = (R_1, \dots, R_T)$ for a given request se-
392 quence μ . To streamline presentation, we present it first as an offline algorithm. At the beginning \mathbb{B} chooses
393 any collection \mathcal{C} of N balanced configurations such that the Hamming distance between every two distinct
394 configurations in \mathcal{C} is at least $k/16$. Such a collection \mathcal{C} can be constructed using a greedy method, as
395 in [12]. (Here we only need existence, which can be also established by a probabilistic proof: if one forms \mathcal{C}

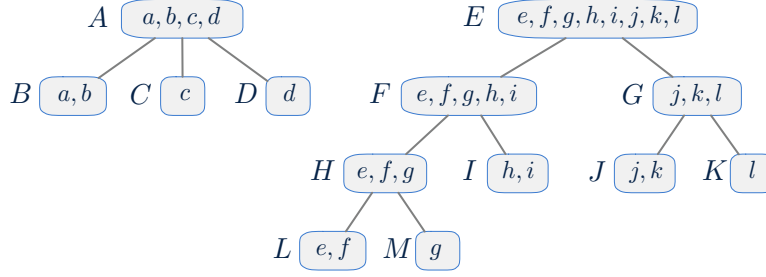


Figure 3: An example of a laminar family \mathcal{P} of height 4.

396 by randomly and uniformly sampling N times with replacement from the balanced configurations, then by
 397 a standard Chernoff bound and the naïve union bound \mathcal{C} has the required property with positive probability.)
 398 Algorithm \mathbb{B} lets $\mathcal{H} = \mathcal{C}$, identifying holes with cache configurations. Then, for each time $t \in [T]$, it replaces
 399 the request C_t in μ by $\psi(\mathcal{C} \setminus \{C_t\})^k$, that is, k repetitions of the forcing sequence for $\mathcal{C} \setminus \{C_t\}$. (This sequence
 400 exists if k is large enough, by Claim 3.4 and the choice of \mathcal{C} .) This will produce sequence σ , an instance
 401 of One-of- m Paging. Next, \mathbb{B} simulates \mathbb{A} on σ . Given a solution D for σ produced by \mathbb{A} , algorithm \mathbb{B}
 402 produces a solution R for μ as follows: For each $t \in [T]$, as D responds to $\psi(\mathcal{C} \setminus \{C_t\})^k$, it either incurs cost
 403 at least k , or uses at least one configuration $P \in \mathcal{C} \setminus \{C_t\}$. Assume without loss of generality that only the
 404 latter case occurs (otherwise modify D to move into any configuration in $\mathcal{C} \setminus \{C_t\}$ at the end of its response
 405 to $\psi(\mathcal{C} \setminus \{C_t\})^k$; these modifications at most double D 's cost), and let $R_t = P$. The produced sequence
 406 $R = (R_1, \dots, R_T)$ is a valid solution to μ , because $R_t \in \mathcal{C} \setminus \{C_t\}$ for $t \in [T]$.

407 To complete the description of \mathbb{B} , it remains to observe that R can be indeed produced in an online
 408 fashion, since R_t does not depend on any future requests in μ . If \mathbb{A} is deterministic then so is \mathbb{B} .

409 **Claim 3.5.** $\text{opt}(\sigma) \leq k \text{opt}(\mu)$.

410 To prove this claim, given an optimal solution (R_1^*, \dots, R_T^*) for $\mu = (R_0, C_1, \dots, C_T)$, consider the
 411 corresponding solution D^* for σ that starts in configuration $R_0^* = R_0$, then, for each $t \in [T]$, responds to
 412 $\psi(\mathcal{C} \setminus \{C_t\})^k$ by having its cache in configuration $R_t^* \in \mathcal{C} \setminus \{C_t\}$ for all requests in $\psi(\mathcal{C} \setminus \{C_t\})^k$. For each
 413 $t \in [T]$, its response to $\psi(\mathcal{C} \setminus \{C_t\})^k$ costs D^* 0 if $R_{t-1}^* = R_t^*$ (the rat didn't move) and otherwise at most
 414 k (to transition the cache from R_{t-1}^* to R_t^*). This proves Claim 3.5.

415 **Claim 3.6.** Algorithm \mathbb{B} is $O(f(k))$ -competitive.

416 Since \mathbb{A} is a $f(k)$ -competitive, $\text{cost}(D) \leq f(k)\text{opt}(\sigma) + O(1)$. Whenever the rat moves (i.e., $R_{t-1} \neq$
 417 R_t), by the definition of \mathcal{B} , the Hamming distance between R_{t-1} and R_t is $\Omega(k)$, so D paid $\Omega(k)$ to transition
 418 from R_{t-1} to R_t (possibly in multiple steps). Using Claim 3.5, we obtain that $\text{cost}(R) = O(\text{cost}(D)/k) =$
 419 $O(f(k)\text{opt}(\sigma)/k) = O(f(k)\text{opt}(\mu))$. That is, R is an $O(f(k))$ -competitive solution for μ . This proves
 420 Claim 3.6, completing the proof of the lemma. \square

421 4 Upper Bounds for Page-Laminar Paging

422 Recall that Page-Laminar Paging generalizes Paging by allowing each request to be a set P of pages. The
 423 request P is satisfiable by having any page $p \in P$ in the cache. We require $P \in \mathcal{P}$, where \mathcal{P} is a pre-specified
 424 laminar collection of sets of pages, whose height we denote by h . (See the example in Figure 3.) To our
 425 knowledge, this problem has not been yet studied in the literature. In particular, we do not know whether the
 426 optimum solution can be computed in polynomial time.

427 **Theorem 4.1.** *Page-Laminar Paging admits the following polynomial-time algorithms: an hk -competitive*
428 *deterministic online algorithm, an hH_k -competitive randomized online algorithm, and an offline h -approximation*
429 *algorithm.*

430 The proof is by reduction to standard Paging. Known polynomial-time algorithms for standard Paging
431 include an optimal offline algorithm [10], a deterministic k -competitive online algorithm [46] and a ran-
432 domized H_k -competitive online algorithm [1]. Theorem 4.1 follows directly from composing these known
433 results with the following lemma.

434 **Lemma 4.1.** *Every $f(k)$ -approximation algorithm \mathbb{A} for Paging can be converted into an $hf(k)$ -approximation*
435 *algorithm \mathbb{B} for Page-Laminar Paging, preserving the following properties: being polynomial-time, online,*
436 *and/or deterministic.*

437 *Proof.* Let \mathbb{A} be any (possibly online, possibly randomized) $f(k)$ -approximation algorithm for Paging. Let
438 Page-Laminar Paging instance π be the input to algorithm \mathbb{B} . For any time step t and set $P \in \mathcal{P}$, let $c_t(P)$
439 denote the child of P whose subtree contains π 's most recent request to a proper descendant of P . This is
440 the child c of P such that $P_{t'} \subseteq c$, where $t' = \max\{i \leq t : P_i \subset P\}$. If there is no such request (t' is
441 undefined or P is a leaf), then define $c_t(P) = P$. Define $p_t(P)$ inductively via $p_t(P) = p_t(c_t(P))$ when
442 $c_t(P) \neq P$, and otherwise $p_t(P)$ is an arbitrary (but fixed) page in P . Call $c_t(P)$ and $p_t(P)$ the *preferred*
443 *child* and *preferred page* of P at time t . At any time, P 's preferred page can be found by starting at P and
444 tracing the path down through preferred children.

445 Define a Paging instance σ from the given instance π by replacing each request P_t in π by its preferred
446 page $p_t(P_t)$ (so $\sigma_t = p_t(P_t)$). Algorithm \mathbb{B} just simulates Paging algorithm \mathbb{A} on input σ , and maintains
447 its cache exactly as \mathbb{A} does. (Note that σ can be computed online, deterministically, in polynomial time.)
448 Algorithm \mathbb{B} is correct because any solution to σ is also a solution to π (because $\sigma_t = p_t(P_t) \in P_t$). And
449 $\text{cost}(\mathbb{B}(\pi)) = \text{cost}(\mathbb{A}(\sigma))$. To finish proving the lemma, we show $\text{opt}(\sigma) \leq h \text{opt}(\pi)$.

450 For any requested set P , define a P -phase of π to be a maximal contiguous interval $[i, j] \subseteq [1, T]$ such
451 that $\pi_t \not\subseteq P$ for $t \in [i + 1, j]$. The P -phases for a given P partition $[1, T]$. Each P -phase $[i, j]$ (except
452 possibly the first, with $i = 1$) starts with a request to a proper descendant of P , but there are no such requests
453 during $[i + 1, j]$. It follows that $c_i(P)$ and $p_i(P)$ remain the preferred child and page of P throughout the
454 phase, that the preferred child $c_i(P)$ of P also has the same preferred page $p_i(P)$ throughout the phase,
455 and that P -phase $[i, j]$ is contained in some $c_i(P)$ -phase. By definition of σ , each request to P in the given
456 instance π in interval $[i, j]$ is replaced in σ by a request to P 's preferred page, $p_i(P)$.

457 **Example.** Consider the laminar family \mathcal{P} in Figure 3. Consider also a request sequence π whose requests
458 $\pi_{61}, \pi_{62}, \dots, \pi_{79}$ are shown in the table below:

time step:	...	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	...	
sequence π :	...	A	B	H	E	C	K	A	I	D	E	M	B	H	A	D	E	G	A	F	...	
E -phases:	...	-	-	H	-	-	K	-	I	-	-	M	-	H	-	-	-	G	-	F	...	
F -phases:	...	-	-	H	-	-	-	-	I	-	-	M	-	H	-	-	-	-	-	-	-	...
sequence σ :	...	a	a	e	e	c	l	c	h	d	h	g	a	g	a	d	g	l	d	g	...	

460 The third row in the table shows E -phases, marking the beginning of each phase with the request at that step.
461 The fourth row shows F -phases. Assume that at time 61 the preferred page of each set in \mathcal{P} is the leftmost
462 page in the leftmost leaf of its subtree. For example, $p_{61}(E) = e$. Then the fifth row shows the sequence of
463 preferred pages that forms the resulting sequence σ .

464 To prove that $\text{opt}(\sigma) \leq h \text{opt}(\pi)$, we will start with an optimal solution for π and convert it into a
465 solution of σ while increasing the cost at most by a factor of h . So let $C = (C_1, \dots, C_T)$ be an optimal
466 solution for π . The conversion of C into a solution for σ is given in Figure 4, and is described as an algorithm

1. Initialize the current instance π' and current solution C' to the given instance π and its solution C .
2. Incrementally modify π' and C' by *repairing* each phase, as follows.
3. While there is an unrepaired phase, choose any unrepaired P -phase $[i, j]$ such that all proper descendants of P have already been repaired, then repair the chosen phase as follows:
 - 3.1. Modify the current instance π' by replacing each request to P during $[i, j]$ in π' by a request to P 's preferred page $p = p_i(P)$. (So, after all phases are repaired, the current instance π' will equal σ .)
 - 3.2. Modify the current solution C' during $[i, j]$ accordingly, to ensure that C' continues to satisfy π' . To do that, we will establish a stronger property throughout $[i, j]$, namely: *whenever C' has at least one page in P cached, C' has p cached.*

Say that time $t \in [i, j]$ *needs repair* if, at time t , C' caches at least one page in P , but not p . For $t \leftarrow i, i + 1, \dots, j$, if time t needs repair, modify what C' caches at time t by replacing one of its currently cached pages $q_t \in P$ by p , where q_t is defined greedily as follows

$$q_t = \begin{cases} q_{t-1} & \text{if } q_{t-1} \text{ is defined and still cached at time } t \\ \text{any page in } P \text{ cached at time } t & \text{otherwise.} \end{cases}$$

This completes the repair of this P -phase $[i, j]$. The algorithm terminates after it has repaired all phases.

Figure 4: The algorithm that transforms (π, C) into (σ, C') , by repairing each phase.

467 that incrementally modifies, or “repairs”, both C and π , phase by phase. It maintains the invariant that the
468 current solution, denoted C' , is always correct for the current instance, denoted π' . (In π' , some requests will
469 be to a page, rather than a set. Any such request is satisfied only by having the requested page in the cache.)
470 At the end the modified instance π' will equal σ , so that the modified solution C' will be a correct solution
471 for σ .

472 Specifically, we will show the following claim (whose proof we postpone):

473 **Claim 4.2.** *The repair algorithm maintains the invariant that the current solution C' is correct for the current*
474 *instance π' , so at termination C' is a correct solution for σ .*

475 Next we bound the cost, as follows. Call a phase *costly* if its repair increases the cost of C' , and *free*
476 otherwise. We show that the number of costly phases is at most $(h - 1)\text{cost}(C)$, and that the repair of each
477 phase increases the cost of C' by at most 1. This implies that the final cost of C' is at most $\text{cost}(C) +$
478 $(h - 1)\text{cost}(C) = h \text{cost}(C)$, as desired. Specifically, we will show the following claims (whose proofs we
479 postpone):

480 **Claim 4.3.** *For any requested set P , the repair of any P -phase $[i, j]$ increases the cost of C' by at most 1,*
481 *and only if $j \neq T$.*

482 **Claim 4.4.** *For any non-leaf set P , the number of costly P -phases is at most the cost paid by C for pages in*
483 *P (that is, the number of retrievals of pages in P by C).*

484 By the definition of P -phases, each leaf set P has only one P -phase $[1, T]$, so by Claim 4.3 only non-leaf
485 sets have costly phases. Each page p is in at most $h - 1$ non-leaf sets P , so Claim 4.4 implies that the total
486 number of costly phases is at most $(h - 1)\text{cost}(C)$. This and Claim 4.3 imply that the final cost of C' is at
487 most $h \text{cost}(C) = h \text{opt}(\pi)$, proving Lemma 4.1.

488 It remains only to prove the three claims.

489 *Proof of Claim 4.2.* The invariant holds initially when $C' = C$ and $\pi' = \pi$, just because by definition C is
 490 an (optimal) solution for π . Suppose the invariant holds just before the repair of some P -phase $[i, j]$. We
 491 will show that it continues to hold after. The repair modifies π' by replacing each request to P during $[i, j]$
 492 by a request to its preferred page $p = p_i(P)$.

493 Consider any time $t \in [i, j]$. First consider the case that (before the repair) π' requested P at time t .
 494 In this case, C' cached some page in P , so (by Step 3.2 of the repair algorithm) after the repair C' has the
 495 preferred page p cached, and thus C' satisfies the modified request (for p).

496 The other case is when π' requested page at time t is not P . Then the repair doesn't modify the request
 497 in π' at time t . In this case, either the repair doesn't modify the cache at time t (in which case C' continues
 498 to satisfy π'), or the repair replaces some page q_t in the cache by p . If that happens, the page q_t is also in
 499 P . Also, the request in π at time t cannot be to a proper descendant of P (by definition of P -phase), so the
 500 request in π' at time t is either to an ancestor of P , a set disjoint from P , or an already repaired page not in
 501 P . (We use here that no proper ancestors of P have yet had their phases repaired.) In all three cases, after
 502 swapping p for q_t (with $p, q_t \in P$), the request must still be satisfied. So the invariant holds after the repair.
 503 At termination the invariant holds so C' is correct for σ . \square

504 *Proof of Claim 4.3.* Consider the repair of any P -phase $[i, j]$ for any requested set P . This repair modifies
 505 the cache only at times in $[i, j]$. Recall that the cost of C' at time t is the number of retrievals at time t , where
 506 a retrieval is a page that C' caches at time t but not at time $t - 1$.

507 At each time $t \in [i, j]$ that needed repair (as defined in Step 3.2 of the algorithm), the repair of the phase
 508 replaced some page q_t in the cache at time t by the preferred page $p = p_i(P)$. This can increase the cost
 509 only at times in $[i, j + 1]$, and by at most 1 at each such time.

510 We claim the cost cannot increase at time i . Indeed, in the case that $i = 1$, the cost at time i equals the
 511 number of pages cached at time i , which a repair at time 1 doesn't change. In the case that $i > 1$, at time i no
 512 repair is done, because π requests a proper descendant d of P , and the d -phase containing $[i, j]$ has already
 513 been repaired, so π' requests $p_i(d)$ at time i , and now our invariant implies that C' already caches $p_i(d)$ at
 514 time i . But by definition $p_i(P) = p_i(d)$, so C' already caches P 's preferred page p at time i . Thus the cost
 515 cannot increase at time i .

516 Next, consider any time $t \in [i + 1, j]$. To prove the claim, we show that the repair didn't increase the
 517 cost at time t . The repair can introduce up to two new retrievals at time t : a new retrieval of p , and/or a new
 518 retrieval of q_{t-1} . We show that, for each retrieval that the repair introduced, it removed another one.

519 Suppose it introduced a new retrieval of p at time t . That is, it replaced q_t by p in the cache at time t and
 520 (after modification) C' doesn't cache p at time $t - 1$. The latter property (by inspection of Step 3.2 of the
 521 algorithm) implies that C' caches no pages in P at time $t - 1$ (before or after modification). It follows that
 522 the repair removed one retrieval of q_t at time t .

523 Now suppose the repair introduced a new retrieval of q_{t-1} at time t . That is, it removed q_{t-1} from the
 524 cache at time $t - 1$ (replacing it by p) and (after modification) C' caches q_{t-1} at time t . The latter property
 525 implies that time t did not need repair (because if it did the repair would have taken $q_t = q_{t-1}$). So C' cached
 526 p at time t (before and after modification). Thus the repair removed a retrieval of p at time t .

527 Overall, for each retrieval introduced at times in $[i, j]$, another was removed, and therefore the cost at
 528 times in $[i, j]$ didn't increase. The cost increase at time $j + 1$, if any, can only be caused by the repair at time
 529 $t = j$ that replaced q_j by p , so this increase is at most 1. This completes the proof of the claim. \square

530 *Proof of Claim 4.4.* Fix any non-leaf set P . First consider the repair of any P -phase $[i, j]$ with $i > 1$. Let
 531 $c = c_i(P) \neq P$ and $p = p_i(P) = p_i(c)$ be the preferred child and page throughout the phase. When the
 532 repair starts, the c -phase containing $[i, j]$ has already been repaired. By inspection, that repair established
 533 the following property of C' throughout $[i, j]$: *whenever C' has at least one page in c cached, C' has c 's*
 534 *preferred page p cached.* None of the ancestors of c (including P) have had their phases repaired since then,
 535 so this property still holds just before the repair of P .

536 Since $i > 1$, the definition of P -phases implies that at time i the instance π requests a descendant of c ,
537 so C caches at least one page p' in c . Suppose that C doesn't evict p' during $[i + 1, j]$. Then, *at every time*
538 *during* $[i, j]$, C caches at least one page in c . Each repair for c or a descendant of c preserves this property
539 (because such a repair only replaces cached pages in c by other pages in c). Throughout $[i, j]$, then, C' also
540 caches at least one page in c and, by the previous paragraph, must cache P 's preferred page p . (Recall that P
541 and c have the same preferred page throughout $[i, j]$.) In this case, by inspection of the algorithm the repair
542 of this P -phase does not change C' , and the phase is not costly. We conclude that, for the phase to be costly,
543 C must evict p' during $[i + 1, j]$.

544 Note that $p' \in c \subset P$. Also, by Claim 4.3, this is not the final P -phase (that is, $j < T$). It follows that
545 the number of costly P -phases $[i, j]$ with $i > 1$ is at most the number of evictions of pages in P by C , before
546 the final P -phase (with $j = T$).

547 Regarding the P -phase $[i, j]$ with $i = 1$, if it is costly, then $j < T$ and, by the reasoning in the paragraph
548 before last, in the final P -phase, there is a page p' in P that C either evicts or leaves in the cache at time T .

549 By the above reasoning, the number of costly P -phases is at most the number of evictions by C of pages
550 in P , plus the number of pages left cached by C at time T . This sum is the number of retrievals of pages in
551 P by C , proving Claim 4.4. \square

552 As explained earlier, Claims 4.2, 4.3 and 4.4 imply the lemma, thus the proof is now complete. \square

553 5 Slot-Laminar Paging

554 In this section we prove upper bounds for Slot-Laminar Paging given in Table 1. Recall that in Slot-Laminar
555 Paging the family \mathcal{S} is assumed to be a laminar family of slot sets whose height we denote by h . Theorem 5.1
556 bounds the optimal ratios by $3h^2k$ (deterministic), $3h^2H_k$ (randomized) and $3h^2$ (offline polynomial-time
557 approximation). The proof of Theorem 5.1 (Section 5.1) is by a reduction of Slot-Laminar Paging to Page-
558 Laminar Paging, studied in Section 4. Theorem 5.2, presented in Section 5.2, tightens the deterministic
559 upper bound to $2hk$.

560 5.1 Upper bounds for randomized and offline Slot-Laminar Paging

561 **Theorem 5.1.** *Slot-Laminar Paging admits the following polynomial-time algorithms: a deterministic $3h^2k$ -*
562 *competitive online algorithm, a randomized $3h^2H_k$ -competitive online algorithm, and an offline $3h^2$ -approximation*
563 *algorithm.*

564 Our focus here is on uniform treatment of the three variants of Slot-Laminar Paging in the above theorem.
565 The ratios in this theorem have not been optimized. For example, in Section 5.2 we give a better deterministic
566 algorithm. For the special case when $h = 2$ the problem can be reduced to All-or-One Paging, for which the
567 ratio can be improved even further [22].

568 The proof of Theorem 5.1 is by a reduction of Slot-Laminar Paging to Page-Laminar Paging, in Lemma 5.1.
569 The reduction uses a relaxation of Slot-Laminar Paging that relaxes the constraint that each slot hold at most
570 one page (but still enforces the cache-capacity constraint), yielding an instance of Page-Laminar Paging.
571 The reduction simulates the given Page-Laminar Paging algorithm on multiple instances of Page-Laminar
572 Paging—one for each set $S \in \mathcal{S}$, obtained by relaxing the subsequence that contains just those requests con-
573 tained in S —then aggregates the resulting Page-Laminar Paging solutions to obtain the global Slot-Laminar
574 Paging solution. Lemma 5.1 and Theorem 4.1 (for Page-Laminar Paging) immediately imply Theorem 5.1.

575 **Lemma 5.1.** *Every $f_h(k)$ -approximation algorithm \mathbb{A} for Page-Laminar Paging can be converted into a*
576 *$3hf_h(k)$ -approximation algorithm \mathbb{B} for Slot-Laminar Paging, preserving the following properties: being*
577 *polynomial-time, online, and/or deterministic.*

578 *Proof.* We first define the Page-Laminar Paging *relaxation* of a given Slot-Laminar Paging instance. The
579 idea is to relax the constraint that each slot can hold at most one page, while keeping the cache-capacity
580 constraint. The relaxed problem is equivalent to a Page-Laminar Paging instance over “virtual” pages $v(p, s)$
581 corresponding to page/slot pairs (p, s) . This virtual page can be placed in any slot, although it represents
582 page p being in slot s .

583 Formally, this relaxation is defined as follows. Fix any k -slot Slot-Heterogenous Paging instance $\sigma =$
584 $(\sigma_1, \dots, \sigma_T)$ with requestable slot-set family \mathcal{S} . For any page p and $S \in \mathcal{S}$, define $V(p, S) = \{v(p, s) : s \in$
585 $S\}$, where $v(p, s)$ is a *virtual page* for the pair (p, s) . Define the *relaxation* of σ to be the k -slot Page-Subset
586 Paging instance $\pi = (P_1, \dots, P_T)$ defined by $P_t = V(p_t, S_t)$ (where $\sigma_t = \langle p_t, S_t \rangle$, for $t \in [T]$). The
587 requestable-set family for π is $\mathcal{P} = \{V(p, S) : p \text{ is any page and } S \in \mathcal{S}\}$. Crucially, if \mathcal{S} is slot-laminar
588 with height h , then \mathcal{P} is page-laminar with the same height h .

589 Instance π is a relaxation of σ in the sense that for any solution C for σ there is a solution D for π with
590 $\text{cost}(D) \leq \text{cost}(C)$. (Namely, have D keep in its cache the virtual pages $v(p, s)$ such that C has page p
591 cached in slot s .) It follows that $\text{opt}(\pi) \leq \text{opt}(\sigma)$.

592 Next we define the algorithm \mathbb{B} . Fix an $f_h(k)$ -approximation algorithm \mathbb{A} for Page-Laminar Paging. Fix
593 the input σ with $\sigma_t = \langle p_t, S_t \rangle$ (for $t \in [T]$) to Slot-Laminar Paging algorithm \mathbb{B} . We assume for ease of
594 presentation that Algorithm \mathbb{A} is an online algorithm, and present Algorithm \mathbb{B} as an online algorithm. If \mathbb{A}
595 is not online, \mathbb{B} can easily be executed as an offline algorithm instead.

596 Assume that the family \mathcal{S} has just one root R with $|R| \leq k$. (This is without loss of generality, as multiple
597 roots, being disjoint, naturally decouple any Slot-Laminar Paging instance into independent problems, one
598 for each root.)

599 For each $S \in \mathcal{S}$, define S 's Slot-Laminar Paging *subinstance* σ_S to be obtained from σ by deleting all
600 requests that are not subsets of S . Let π_S denote the (Page-Laminar Paging) relaxation of σ_S . Algorithm \mathbb{B}
601 on input σ executes, simultaneously, $\mathbb{A}(\pi_S)$ for every requestable set $S \in \mathcal{S}$, giving each execution $\mathbb{A}(\pi_S)$
602 its own independent cache of size $|S|$ composed of copies of the slots in S .

603 For each such S , Algorithm \mathbb{B} will build its own solution, denoted $\mathbb{B}(\sigma_S)$, for σ_S , also using its own
604 independent cache of size $|S|$ composed of copies of the slots in S . The desired solution to σ will then be
605 $\mathbb{B}(\sigma_R)$ (note that $\sigma = \sigma_R$).

606 For internal bookkeeping purposes only, in presenting Algorithm \mathbb{B} , we consider each virtual page $v(p, s)$
607 (as defined for Page-Laminar Paging) to be a *copy* of page p , and we have \mathbb{B} maintain cache configurations
608 that place these virtual pages in specific slots, with the understanding that the actual cache configurations are
609 obtained by replacing each virtual page $v(p, s)$ (in whatever slot it's in) by a copy of page p . This virtual
610 copy $v(p, s)$ is functionally equivalent to p ; for example, if placed in slot s' , it will satisfy any request
611 $\langle p, S' \rangle$ with $s' \in S'$. When we analyze the cost, we will consider two copies $v(p, s)$ and $v(p', s')$ to be
612 distinct unless $(p', s') = (p, s)$. In particular, if \mathbb{B} evicts $v(p, s)$ while retrieving $v(p, s')$ (with $s' \neq s$) in the
613 same slot, this contributes 1 to the cost of \mathbb{B} . We will upper bound \mathbb{B} 's cost overestimated in this way.

614 **Correctness.** Algorithm \mathbb{B} will somehow maintain the following invariant over time:

615 *For each requestable set S , for each virtual page $v(p, s)$ currently cached by $\mathbb{A}(\pi_S)$:*

- 616 1. *the solution $\mathbb{B}(\sigma_S)$ caches $v(p, s)$ in some slot in S , and*
- 617 2. *if S has a child c with $s \in c$, and $\mathbb{B}(\sigma_c)$ has $v(p, s)$ in its cache c , then in $\mathbb{B}(\sigma_S)$ copy $v(p, s)$ is*
618 *in the same slot as in $\mathbb{B}(\sigma_c)$.*

619 The invariant suffices to guarantee correctness of the solution $\mathbb{B}(\sigma_S)$ for each instance σ_S . Indeed, when
620 $\mathbb{B}(\sigma_S)$ receives a request $\langle p_t, S_t \rangle$, its relaxation $\mathbb{A}(\pi_S)$ has just received the request $\{v(p_t, s) : s \in S_t\}$,
621 so $\mathbb{A}(\pi_S)$ is caching a virtual page $v(p_t, s)$ (for some $s \in S_t$) in S . By Condition 1, then, $\mathbb{B}(\sigma_S)$ also has
622 $v(p_t, s)$ in some slot in S . In the case $S = S_t$, this suffices for $\mathbb{B}(\sigma_S)$ to satisfy the request. In the remaining

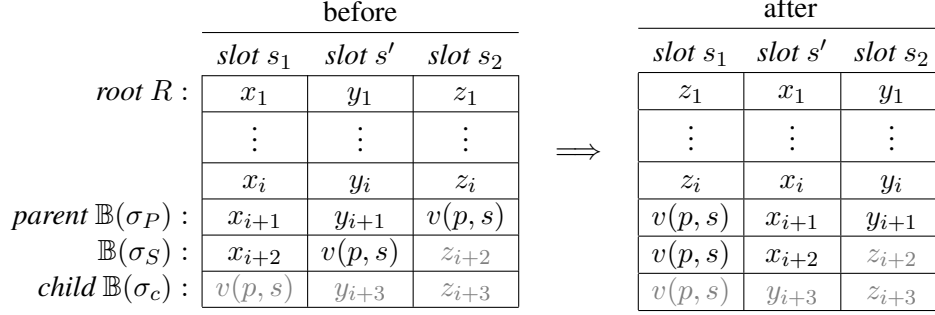


Figure 5: “Rotating” slots in $\mathbb{B}(\sigma_S)$ and ancestors to preserve the invariant. Pages in grey are not moved.

623 case S has a child c with $S_t \subseteq c$, and $\mathbb{B}(\sigma_C)$ just received the same request, so (assuming inductively that
624 $\mathbb{B}(\sigma_c)$ is correct for σ_c) $\mathbb{B}(\sigma_c)$ has $v(p_t, s)$ in some slot s' in S_t , so by Condition 2 of the invariant $\mathbb{B}(\sigma_S)$ has
625 $v(p_t, s)$ in the same slot s' in S_t , as required. In particular, $\mathbb{B}(\sigma_R)$ will be correct for σ_R .

626 To maintain the invariant \mathbb{B} does the following for each requestable set S . Whenever the relaxed solution
627 $\mathbb{A}(\pi_S)$ evicts a page $v(p, s)$, the solution $\mathbb{B}(\sigma_S)$ also evicts $v(p, s)$. After this eviction both Conditions 1
628 and 2 will be preserved. Whenever $\mathbb{A}(\pi_S)$ retrieves a page $v(p, s)$, the solution $\mathbb{B}(\sigma_S)$ also retrieves $v(p, s)$,
629 into any vacant slot in S (there must be one, because $\mathbb{A}(\pi_S)$ caches at most $|S|$ pages). This retrieval can
630 cause up to two violations of Condition 2 of the invariant: one at $\mathbb{B}(\sigma_S)$, because $v(p, s)$ is already cached
631 by a child $\mathbb{B}(\sigma_c)$ but in some slot $s_1 \neq s'$; the other at the parent $\mathbb{B}(\sigma_P)$ of $\mathbb{B}(\sigma_S)$ (if any), because $v(p, s)$ is
632 already cached by the parent, but in some slot $s_2 \neq s'$. In the case that the retrieval does create two violations
633 (and $s_1 \neq s_2$), \mathbb{B} restores the invariant by “rotating” the contents of the slots s_1 , s' , and s_2 in $\mathbb{B}(\sigma_S)$ and in
634 each ancestor, as shown in Figure 5. Note that y_{i+3} and z_{i+2} cannot be $v(p, s)$, so moving $v(p, s)$ out of slots
635 s' and s_2 doesn’t introduce a violation there. Thus this rotation indeed restores the invariant, at the expense
636 of three retrievals at the root. (The retrievals at other nodes only modify the internal state of \mathbb{B} .) There are
637 three other cases: two violations with $s_1 = s_2$, one violation at $\mathbb{B}(\sigma_S)$, or one violation at its parent, but all
638 these three cases can be handled similarly, also with at most three retrievals (in fact at most two) at the root.

Total cost. Each retrieval by $\mathbb{A}(\pi_S)$ causes at most 3 retrievals in $\mathbb{B}(\sigma_R)$, so $\text{cost}(\mathbb{B}(\sigma_R))$ is at most

$$\leq \sum_{S \in \mathcal{S}} 3 \text{cost}(\mathbb{A}(\pi_S)) \leq \sum_{S \in \mathcal{S}} 3 f_h(|S|) \text{opt}(\pi_S) \leq 3 f_h(k) \sum_{S \in \mathcal{S}} \text{opt}(\sigma_S) \leq 3h f_h(k) \text{opt}(\sigma_R).$$

639 The second step uses that $\mathbb{A}(\pi_S)$ is $f_h(|S|)$ -competitive for π_S . The third step uses that π_S is a relaxation of
640 σ_S so $\text{opt}(\pi_S) \leq \text{opt}(\sigma_S)$, and that $|S| \leq k$ so $f_h(|S|) \leq f_h(k)$.¹ The last step uses that the sets within any
641 given level $i \in \{1, 2, \dots, h\}$ of the laminar family are disjoint, so $\text{opt}(\sigma_R)$ is at least the sum, over the sets
642 S within level i , of $\text{opt}(\sigma_S)$. This shows that \mathbb{B} is a $3h f_h(k)$ -approximation algorithm. To finish, we observe
643 that \mathbb{B} is polynomial-time, online, and/or deterministic if \mathbb{A} is. \square

644 5.2 Improved upper bound for deterministic Slot-Laminar Paging

645 For Slot-Laminar Paging, this section presents a deterministic algorithm with competitive ratio $O(hk)$, im-
646 proving upon the bound of $O(h^2k)$ from Theorem 5.1. The algorithm, REFSEARCH, refines EXHSEARCH.
647 Like EXHSEARCH, it is phase-based and maintains a configuration that can satisfy all requests in a phase;
648 however, in order to satisfy the next request in the current phase, the particular configuration is chosen by
649 judiciously moving pages in certain slots that are serving requests along a path in the laminar hierarchy.

¹We assume here that $f_h(k') \leq f_h(k)$ for $k' \leq k$, which is without loss of generality as one can simulate a cache of size k' using a cache of size k by introducing artificial requests that force $k - k'$ slots to be continuously occupied.

<p>input: Slot-Laminar Paging instance $(k, \mathcal{S}, \sigma = (\sigma_1, \dots, \sigma_T))$</p> <ol style="list-style-type: none"> 1. for $t \leftarrow 1, 2, \dots, T$, respond to the current request $\sigma_t = \langle p, S \rangle$ as follows: <ol style="list-style-type: none"> 1.1. if $t = 1$ or $R_{t-1} \cup \{\sigma_t\}$ is not satisfiable: let $R_{t-1} = \emptyset$ and empty the cache — start new phase 1.2. let $R_t = R_{t-1} \cup \{\sigma_t\}$ 1.3. if C_{t-1} satisfies $\sigma_t = \langle p, S \rangle$: let $C_t = C_{t-1}$ — redundant request 1.4. else: — non-redundant request <ol style="list-style-type: none"> 1.4.1. find sequences $\langle s_1, \dots, s_m \rangle$, $\langle S_0 = S, S_1, \dots, S_{m-1} \rangle$, and $\langle p_0 = p, p_1, \dots, p_{m-1} \rangle$ s.t. <ol style="list-style-type: none"> (i) $S_{i-1} \subsetneq S_i$ and slot $s_i \in S_{i-1}$ of C_{t-1} satisfies $\langle p_i, S_i \rangle \in \text{rep}(R_{t-1})$, for $1 \leq i < m$, and (ii) slot $s_m \in S_{m-1}$ of C_{t-1} either <ol style="list-style-type: none"> (ii.1) does not satisfy any requests in $\text{rep}(R)$, or (ii.2) satisfies a request $\langle p, S' \rangle \in \text{rep}(R_{t-1})$ such that $S' \supsetneq S_{m-2}$ 1.4.2. to obtain C_t and satisfy $\langle p_{i-1}, S_{i-1} \rangle$, place p_{i-1} in slot s_i, for $1 \leq i \leq m$
--

Figure 6: Deterministic online Slot-Laminar Paging algorithm REFSEARCH. Note that in Step 1.4.1 we have $m \leq k + 1 - |S|$, and that in (ii), if s_m satisfies $\langle p, S' \rangle \in \text{rep}(R_{t-1})$ then $m \geq 2$ (because C_{t-1} does not satisfy σ_t); thus S_{m-2} is well-defined.

650 **Theorem 5.2.** For Slot-Laminar Paging, Algorithm REFSEARCH (Fig. 6) has competitive ratio at most
651 $2 \cdot \text{mass}(\mathcal{S}) - k \leq (2h - 1)k$.

652 We begin by defining the terminology used in the algorithm and the proof, and establish some useful
653 properties. Recall that a configuration D satisfies a request $r = \langle p, S \rangle$ if there exists a slot s in S such that
654 s holds p in D ; in this case, we also say that slot s satisfies r in D . A configuration D is said to satisfy
655 a set R of requests if it satisfies every request in R . A set R of requests will be called *satisfiable* if there
656 exists a configuration that satisfies R . To determine if a set R of requests is satisfied by a configuration, it
657 is sufficient (and necessary) to examine the maximal subset of “deepest” requests in the laminar hierarchy.
658 Formally, a request $\langle p, S \rangle$ is an *ancestor* (resp., *descendant*) of $\langle p, S' \rangle$ if $S \supseteq S'$ (resp., $S \subseteq S'$). For any set
659 R of requests, define $\text{rep}(R)$ as the set of requests in R that do not have any proper descendants in R . That
660 is, $\text{rep}(R) = \{ \langle p, S \rangle \in R : \forall S' \subsetneq S, \langle p, S' \rangle \notin R \}$. For $r = \langle p, S \rangle$, define $\text{anc}(r, R) = \{ \langle p, S' \rangle \in R : S \subseteq S' \}$.
661 Lemma 5.2 establishes some basic properties of $\text{rep}(R)$.

662 **Lemma 5.2.** Let R be a set of requests. Then,

663 (i) In any configuration, each slot can satisfy at most one request in $\text{rep}(R)$.

664 (ii) A configuration satisfies R if and only if it satisfies $\text{rep}(R)$.

665 (iii) R is satisfiable iff for any requestable set S , $\text{rep}(R)$ has at most $|S|$ requests to subsets of S .

666 *Proof.* (i) This part holds because any two requests in $\text{rep}(R)$ request either different pages or disjoint slot
667 sets. (ii) Since $\text{rep}(R) \subseteq R$, if R is satisfiable, so is $\text{rep}(R)$. On the other hand, if a configuration D satisfies
668 $\text{rep}(R)$ then D satisfies R , because every r in R is an ancestor of some r' in $\text{rep}(R)$ and can be satisfied by
669 the slot satisfying r' .

670 (iii) Suppose that R is satisfiable. If D is a configuration that satisfies R then it also satisfies $\text{rep}(R)$,
671 by (ii). By (i), for any requestable set S , all requests in $\text{rep}(R)$ to subsets of S must be satisfied in D by
672 different slots of S , so there can be at most $|S|$ such requests. To prove the reverse implication, assume that
673 for any requestable set S there are at most $|S|$ requests in $\text{rep}(R)$ to subsets of S . We construct D top-down.
674 Let T be the root of the laminar hierarchy \mathcal{S} . (We could assume that $T = [k]$, but it’s not necessary.) By
675 our assumption, there are at most $|T|$ requests in $\text{rep}(R)$. The children of T in \mathcal{S} are disjoint, so we can
676 distribute these requests to the children of T in such a way that each child Q is assigned at most $|Q|$ requests
677 from $\text{rep}(R)$, and each request assigned to Q is to a subset of Q . Continuing this recursively down the tree,

678 we will end up with requests assigned to leaves. Then, for any leaf L we can satisfy its assigned requests by
 679 different slots in L . \square

680 Algorithm REFSEARCH is given in Figure 6. It consists of phases. The first phase starts in time step 1,
 681 and each phase ends when adding the current request to the request set from this phase makes it unsatisfiable.
 682 Within a phase, redundant requests, that is those satisfied by the current configuration, are ignored (Step 1.3).
 683 To serve a non-redundant request $\sigma_t = \langle p, S \rangle$, the cache content is rearranged to free a slot in S . This
 684 rearrangement involves shifting the content of some slots that serve requests in $\text{rep}(R)$ along the path from
 685 S to the root, to find a slot that is either unused or holds p (Step 1.4.2).

686 For technical reasons, in the analysis of Algorithm REFSEARCH it will be useful to introduce a slightly
 687 refined concept of configurations. Given a request set R , an R -configuration is a configuration D in which
 688 each request in $\text{rep}(R)$ is served by exactly one slot. (By Lemma 5.2(i), each slot can serve only one
 689 request in $\text{rep}(R)$), but in general in a configuration serving R there may be multiple slots that serve the same
 690 request in $\text{rep}(R)$.) Slots in D that do not serve requests in $\text{rep}(R)$ are called *free in D* . Observe that each
 691 configuration C_t of Algorithm REFSEARCH implicitly is an R_t -configuration – due to the assignment of
 692 slots in Step 1.4.2. Also, if the slot s_m chosen by the algorithm in Step 1.4.1 satisfies condition (ii.1) then
 693 s_m is a free slot of D , according to our definition.

694 The following helper claim, which characterizes when a particular request is not satisfied by a given
 695 configuration, follows directly from Lemma 5.2(iii).

696 **Claim 5.3.** *Let R be a set of requests and D be an R -configuration. Let also $r = \langle p, S \rangle$ be a request such
 697 that D does not satisfy r , yet $R \cup \{r\}$ is satisfiable. Then D has a slot s in S that is either free or satisfies a
 698 request $\langle p', S' \rangle \in \text{rep}(R)$ where $S \subsetneq S'$.*

699 The following lemma establishes the validity of Steps 1.4.1 and 1.4.2 of Algorithm REFSEARCH.

700 **Lemma 5.4.** *Let R be a set of requests and D be an R -configuration. Let $r = \langle p_0, S_0 \rangle$ be a request
 701 such that r is not satisfied by D and $R \cup \{r\}$ is satisfiable. Then there exist sequences $\langle s_1, \dots, s_m \rangle$,
 702 $\langle S_0, S_1, \dots, S_{m-1} \rangle$, and $\langle p_0, p_1, \dots, p_{m-1} \rangle$ such that (i) $S_{i-1} \subsetneq S_i$ and $s_i \in S_{i-1}$ is currently satisfying
 703 request $\langle p_i, S_i \rangle \in \text{rep}(R)$, for $1 \leq i < m$, and (ii) $s_m \in S_{m-1}$ is either a free slot or is currently satisfying
 704 $\langle p_0, S' \rangle \in \text{rep}(R)$ for some $S' \supsetneq S_{m-2}$. Furthermore, transforming D by moving page p_{i-1} to slot s_i (and
 705 modifying the slot assignment in D accordingly), for $1 \leq i \leq m$, yields an $(R \cup \{r\})$ -configuration.*

706 *Proof.* The proof is by induction on the depth of S_0 in the laminar hierarchy. For the induction base, consider
 707 $S_0 = [k]$. Since r is not satisfied by D , $R \cup \{r\}$ is satisfiable, and every requestable slot set is subset of $[k]$,
 708 we obtain from Claim 5.3 that there is a free slot $s_1 \in S_0$. The desired claim of the lemma holds with $m = 1$
 709 and sequences $\langle s_1 \rangle$, $\langle S_0 \rangle$ and $\langle p_0 \rangle$ which satisfy (i). Since s_1 is free, bringing page p_0 to slot s_1 yields a
 710 $(R \cup \{r\})$ -configuration.

711 We now establish the induction step. Let R , D , and $r = \langle p_0, S_0 \rangle$ be as given. By Claim 5.3 there are
 712 two cases. In the first case, there is a free slot $s_1 \in S_0$ in D . Then the desired claim holds with $m = 1$, and
 713 sequences $\langle s_1 \rangle$, $\langle S_0 \rangle$ and $\langle p_0 \rangle$. Furthermore, as in the base case, since s_1 is free, bringing page p_0 to slot s_1
 714 yields an $(R \cup \{r\})$ -configuration.

715 The remainder of this proof concerns the second case, in which there is a slot $s_1 \in S_0$ currently satisfying
 716 a request $r' = \langle p_1, S_1 \rangle$ in $\text{rep}(R)$ with $S_0 \subsetneq S_1$. Let D' denote the configuration that is identical to D except
 717 that D has p_0 in slot s_1 . Since D is an R -configuration, no other slot satisfies r' in D ; the same holds in
 718 D' . Hence, D' does not satisfy r' . Furthermore, D' satisfies every request in $\text{rep}(R)$ other than r' . Let
 719 $R' = R \cup \{r\} \setminus \text{anc}(r', R)$. In D' , s_1 satisfies r . Consider any request x in $R \setminus \text{anc}(r', R)$. By definition of
 720 $\text{rep}(R)$, there exists a request x' in $\text{rep}(R)$ that is a descendant of x . Since R' does not include any ancestors
 721 of r' , x' is not r' and hence is satisfied by some slot in D' . We thus obtain that D' satisfies R' and, in fact D' is
 722 an R' -configuration. In D' slot s_1 is assigned to r , and if there is a request $\langle p, S' \rangle$ in $\text{rep}(R)$ then its assigned

723 slot is designated as free in D' . At the same time, D' does not satisfy r' . Further, since $R' \cup \{r'\}$ is a subset
724 of $R \cup \{r\}$, which is satisfiable, $R' \cup \{r'\}$ is also satisfiable. Since $S_1 \supseteq S_0$, by the induction hypothesis,
725 there are sequences $\langle s_2, \dots, s_m \rangle$, $\langle S_1, S_2, \dots, S_{m-1} \rangle$ and $\langle p_1, p_2, \dots, p_{m-1} \rangle$ such that (i) $S_{i-1} \subsetneq S_i$ and
726 $s_i \in S_{i-1}$ is currently satisfying $(p_i, S_i) \in \text{rep}(R')$, for $2 \leq i < m$; and either (ii.1) s_m is a free slot in D' or
727 (ii.2) is currently satisfying a request $(p_1, S') \in \text{rep}(R')$ for some $S' \supseteq S_1$. Note, however, that s_m has to be
728 a free slot in D' since (ii.2) above cannot hold: any request (p_1, S') is in $\text{anc}(r', R)$, all requests of which are
729 excluded from R' . Furthermore, transforming D' to D'' by moving page p_{i-1} to s_i for $2 \leq i \leq m$, satisfies
730 $R' \cup \{r'\}$.

731 We now establish the desired claim for D , R , and r . Consider sequences $\langle s_1, \dots, s_m \rangle$, $\langle S_0, S_1, \dots, S_{m-1} \rangle$
732 and $\langle p_0, \dots, p_{m-1} \rangle$. The desired condition (i) follows from (i) of the induction hypothesis above and the fact
733 that in D , $s_1 \in S_0$ is currently satisfying a request (p_1, S_1) in $\text{rep}(R)$ with $S_0 \subsetneq S_1$. For (ii), note that
734 since s_m is a free slot in D' , either s_m is a free slot in D or (p_0, S') is in $\text{rep}(R)$ for some $S' \supseteq S_{m-2}$, thus
735 establishing (ii). Finally, transforming D to D'' by moving p_{i-1} to s_i for $1 \leq i \leq m$, satisfies $R' \cup \{r'\}$.
736 Since any request satisfying r' also satisfies all ancestors of r' , we have $\text{rep}(R \cup \{r\}) = \text{rep}(R' \cup \{r'\})$,
737 implying that D'' also satisfies $R \cup \{r\}$. This completes the induction step and the proof of the lemma. \square

738 *Proof of Theorem 5.2.* We first argue that at any time t , configuration C_t of REFSEARCH satisfies the set R_t
739 of requests from the current phase of the algorithm. The proof is by induction on the number of steps within
740 a phase. When the phase is about to start at time t then R_{t-1} is set to \emptyset , so the claim holds. For the induction
741 step, consider a step t within a phase and assume that C_{t-1} satisfies R_{t-1} . If C_{t-1} satisfies new request σ_t ,
742 then by Step 3.3, C_t satisfies R_t . Otherwise, $R_{t-1} \cup \{\sigma_t\}$ is satisfiable but C_{t-1} does not satisfy σ_t . Then,
743 by Lemma 5.4, Steps 1.4.1 and 1.4.2 derive a configuration C_t satisfying R_t , completing the induction step
744 and the argument that at any time t , C_t satisfies R_t .

745 We next analyze the competitive ratio. We first show that the number of page retrievals during a phase
746 of REFSEARCH is at most $2 \cdot \text{mass}(\mathcal{S})$. Let R denote the set of requests in the current phase. We charge
747 the cost in this phase to the depths of the requests in $\text{rep}(R)$. The cost of Step 1.4.2 is m . If s_m satisfies
748 condition (ii.1), then $\text{rep}(R \cup \{\sigma_t\}) = \text{rep}(R) \cup \{\sigma_t\}$ and the depth of S is at least m , so the charge per unit
749 depth is at most 1. Otherwise, condition (ii.2) holds and $\text{rep}(R \cup \{\sigma_t\}) = \text{rep}(R) \cup \{\sigma_t\} \setminus \{(p, S')\}$. In this
750 case we have σ_t inherit the charges to $\langle p, S' \rangle$, and we charge the cost of m to the difference in depths of S and
751 S' , which is at least $m - 1$ (because $S_{m-2} \subsetneq S'$), so the charge per unit of depth is at most $m/(m - 1) \leq 2$.
752 (Note that in this case $m \geq 2$.) When the phase ends, a request at depth d was charged at most d times, and
753 these charges include at least one unit charge, so its total charge is most $2d - 1$. Thus the algorithm's cost per
754 phase is at most $2 \cdot \text{mass}(\mathcal{S}) - k \leq (2h - 1)k$. The optimal cost in a phase is at least 1 as no configuration
755 satisfies all requests in the phase and the request that starts the next phase. The theorem follows. \square

756 6 All-or-One Paging

757 Recall that All-or-One Paging is the extension of standard Paging that allows two types of requests: A general
758 request for a page p , denoted $\langle p, * \rangle$, can be served by having p in any cache slot. A specific request $\langle p, j \rangle$,
759 where $j \in [k]$, must be served by having p in slot j of the cache. (Section 2 gives a formal definition.) It is a
760 restriction of Slot-Laminar Paging with $h = 2$.

761 For All-or-One Paging, this section first shows that the optimal randomized ratio is at least $2H_k - O(1)$.
762 It then shows that the offline problem is NP-hard .

763 6.1 Lower bound for randomized All-or-One Paging

764 **Theorem 6.1.** *Every online randomized algorithm \mathbb{A} for the All-or-One Paging problem has competitive*
765 *ratio at least $2H_k - 1$.*

766 *Proof.* We establish our lower bound by giving a probability distribution on the input sequences for which
767 any deterministic algorithm \mathbb{A} has expected cost at least $2H_k - 1$ times the optimum cost. Without loss
768 of generality we can assume that \mathbb{A} is *lazy*, in the sense that it retrieves a page only when it is necessary to
769 satisfy a request. We use some fixed $k + 1$ pages p_1, p_2, \dots, p_{k+1} and the random input sequence will consist
770 of L phases, where L is some large integer.

Consider any phase. To ease notation, by symmetry, assume without loss of generality that when the
phase starts the adversary has pages p_1, p_2, \dots, p_k in the cache, with each page p_i in slot i , for $i = 1, 2, \dots, k$.
To start the phase, the adversary chooses a random permutation $p_{i_1}, p_{i_2}, \dots, p_{i_k}$ of these k pages, replaces
 p_{i_k} in its cache by p_{k+1} , at cost 1, then makes request $\langle p_{k+1}, * \rangle$, followed by $k - 1$ stages. Each stage
 $s = 1, 2, \dots, k - 1$ consists of $L \cdot (2H_k - 1)$ repetitions of the request sequence

$$\langle p_{i_1}, i_1 \rangle, \langle p_{i_2}, i_2 \rangle, \dots, \langle p_{i_s}, i_s \rangle, \langle p_{k+1}, * \rangle,$$

771 which costs the adversary nothing.

772 It remains to bound the expected cost of \mathbb{A} . Let E denote the event that for every phase and every stage s
773 in the phase, the configuration of \mathbb{A} at the end of the stage has each page p_{i_r} , for $r = 1, 2, \dots, s - 1$, in slot
774 i_r and one of the slots in $[k] \setminus \{i_1, \dots, i_{s-1}\}$ contains p_{k+1} . We will separately bound the expected cost of
775 \mathbb{A} conditioned on E , and the expected cost of \mathbb{A} conditioned on \overline{E} .

776 We first analyze the expected cost of \mathbb{A} conditioned on E . Consider any stage s of any phase. If this
777 is the first stage of the first phase, then \mathbb{A} and the adversary start with the same configuration. Otherwise,
778 since event E holds, the configuration of \mathbb{A} has each page p_{i_r} , for $r = 1, 2, \dots, s - 1$, in slot i_r , and
779 one of the slots in $[k] \setminus \{i_1, \dots, i_{s-1}\}$ contains p_{k+1} . Since the probability distribution of i_s is uniform
780 in $[k] \setminus \{i_1, i_2, \dots, i_{s-1}\}$, the probability that \mathbb{A} has p_{k+1} in slot i_s equals $1/(k - s + 1)$. If it does, the
781 cost of \mathbb{A} is at least 2 in stage s , because p_{i_s} will need to be fetched into slot i_s and p_{k+1} will need to
782 be moved to a different slot. So the expected cost of \mathbb{A} in this stage is at least $2/(k - s + 1)$. Summing
783 over all stages $s = 1, 2, \dots, k - 1$ and adding 1 for the first request, the expected cost of \mathbb{A} for a phase,
784 conditioned on event E , will be at least $2(H_k - 1) + 1 = 2H_k - 1$. On the other hand, the adversary pays
785 1 for each phase. Therefore, the expected total cost of \mathbb{A} over L phases, conditioned on event E , is at least
786 $L \cdot (2(H_k - 1) + 1) = L \cdot (2H_k - 1)$, which grows with L , and is at least $2H_k - 1$ times the adversary's
787 total cost.

788 We next analyze the expected cost of \mathbb{A} conditioned on \overline{E} . The event \overline{E} implies that there is a stage s of
789 a phase in which \mathbb{A} does not end with a configuration in which page p_{i_r} , for $r = 1, 2, \dots, s - 1$, is in slot
790 i_r , and one of the slots in $[k] \setminus \{i_1, \dots, i_{s-1}\}$ contains p_{k+1} . Since such a configuration satisfies all requests
791 in the stage and \mathbb{A} is lazy, this implies that \mathbb{A} never reaches such a configuration in the stage. Therefore, the
792 cost of \mathbb{A} in this stage alone is at least $L \cdot (2H_k - 1)$. Since the adversary pays 1 for each phase, the total
793 cost of the adversary is L . Therefore, the expected total cost of \mathbb{A} conditioned on \overline{E} is at least $2H_k - 1$ times
794 the adversary's total cost.

795 We thus obtain the expected total cost of \mathbb{A} grows with L and is at least $2H_k - 1$ times the adversary's
796 total cost. Therefore, the competitive ratio of \mathbb{A} is at least $2H_k - 1$. \square

797 6.2 NP-completeness of offline All-or-One Paging

798 The off-line version of Paging, where the request sequence is given upfront, can be solved in time $O(n \log n)$
799 using the classical algorithm by Belady [11]. All-or-One Paging differs from standard Paging only by inclu-
800 sion of specific requests, which appear easy to handle because they don't give the algorithm any choice. In
801 this section we show that this intuition is not valid:

802 **Theorem 6.2.** *Offline All-or-One Paging is NP-complete.*

803 *Proof of Theorem 6.2.* Let $G = (V, E)$ be a graph with vertex set $V = \{0, 1, \dots, n-1\}$. Given an integer
804 k , $1 \leq k \leq n$, we compute in polynomial time a request sequence σ and an integer F such that the following
805 equivalence holds: G has a vertex cover of size k if and only if there is a solution for σ whose cost with a
806 cache size $k+2$ is at most F .

807 At a fundamental level our proof resembles the argument in [27], where \mathbb{NP} -completeness of an interval-
808 packing problem was proved. The basic idea of the proof is to represent the vertices by a collection of
809 intervals with specified endpoints that are to be packed into a strip of width k . These intervals will be
810 represented by pairs of requests, one at the beginning and one at the end of the interval, and the strip to be
811 packed is the cache. Since the strip's capacity is bounded by k , only a subset of intervals can be packed, and
812 the intervals that are packed correspond to a vertex cover.

813 There will actually be many “bundles” of such intervals, with each bundle containing n intervals cor-
814 responding to the n vertices. If we had $|E|$ bundles and if we forced each bundle's packing (that is, its
815 corresponding set of vertices) to be the same, we could add an edge-gadget to each bundle that will verify
816 that all edges are covered. While it does not seem possible to design these bundles to force all bundles'
817 packings to be equal, there is a way to design them to ensure that the packing of each bundle is dominated
818 (in the sense to be defined shortly) by the next one, and this dominance relation has polynomial depth. So
819 with polynomially many bundles we can ensure that there will be $|E|$ consecutive equally packed bundles,
820 allowing us to verify whether the vertex set corresponding to this packing is indeed a correct vertex cover.

821 **Set dominance.** We consider the family of all k -element subsets of V . For any two k -element sets $X, Y \subseteq$
822 V , we say that Y *dominates* X , and denote it $X \preceq Y$, if there is a 1-to-1 function $\psi : X \rightarrow Y$ such that
823 $x \leq \psi(x)$ for all $x \in X$. We write $X \prec Y$ iff $X \preceq Y$ and $X \neq Y$. The dominance relation is a partial
824 order. The following lemma from [27] will be useful:

825 **Lemma 6.1.** *Let $X_1, X_2, \dots, X_p \subseteq V$ be sets of cardinality k such that $X_1 \prec X_2 \prec \dots \prec X_p$. Then*
826 $p \leq k(n-k)$.

827 **Cover chooser.** We start by specifying the “cover chooser” sequence σ' of requests. In this sequence some
828 time slots will not have assigned requests. Some of these unassigned slots will be used later to insert requests
829 representing edge gadgets.

830 Let $m = |E| + 1$. (For notation-related reasons, it is convenient to have m be one larger than the number
831 of edges.) Let also $P = k(n-k) + 1$ and $B = mP$. In σ' we will use the following pages and requests:

- 832 • We have nB pages $x_{b,j}$, for $b = 0, 1, \dots, B-1$ and $j = 0, 1, \dots, n-1$. For each page $x_{b,j}$ there
833 are two general requests $\langle x_{b,j}, * \rangle$ in σ' at time steps $\tau_{b,j} = 9(bn+j)$ and $\tau'_{b,j} = 9(bn+j) + 9n - 6$.
834 These requests are called *vertex requests*. They are grouped into *bundles* of requests, where bundle b
835 consists of all $2n$ general requests to pages $x_{b,0}, x_{b,1}, \dots, x_{b,n-1}$. See Figures 7 and 8 for illustration.
- 836 • We have B pages y_b , for $b = 0, 1, \dots, B-1$. For each page y_b we have two specific requests
837 $\langle y_b, k+1 \rangle$ and $\langle y_b, k+2 \rangle$ in σ' , at times $\theta_b = \tau'_{b,0} - 2 = \tau_{b,n-1} + 1$ and $\theta'_b = \tau'_{b,0} - 1 = \tau_{b,n-1} + 2$,
838 respectively. For each b , these requests are called *b-blocking requests*, because for each page $x_{b,j}$ in
839 bundle b we have $\theta_b, \theta'_b \in [\tau_{b,j}, \tau'_{b,j}]$, so these two requests make it impossible to have both requests
840 $\langle x_{b,j}, * \rangle$ served in cache slots $k+1$ or $k+2$ with only one fault.

841 Slots $1, 2, \dots, k$ in the cache will be referred to as *vertex slots*. Slot $k+1$ is called the *edge-gadget slot*, and
842 slot $k+2$ is called the *junkyard slot*.

843 Let $F' = (2n-k+2)B$. For any solution S of σ' and any bundle b , denote by $V_{S,b}$ the set of vertices
844 $j \in V$ for which S does not fault in σ' on request $\langle x_{b,j}, * \rangle$ at time $\tau'_{b,j}$. In other words, S keeps $x_{b,j}$ in the
845 cache throughout the time interval $[\tau_{b,j}, \tau'_{b,j}]$.

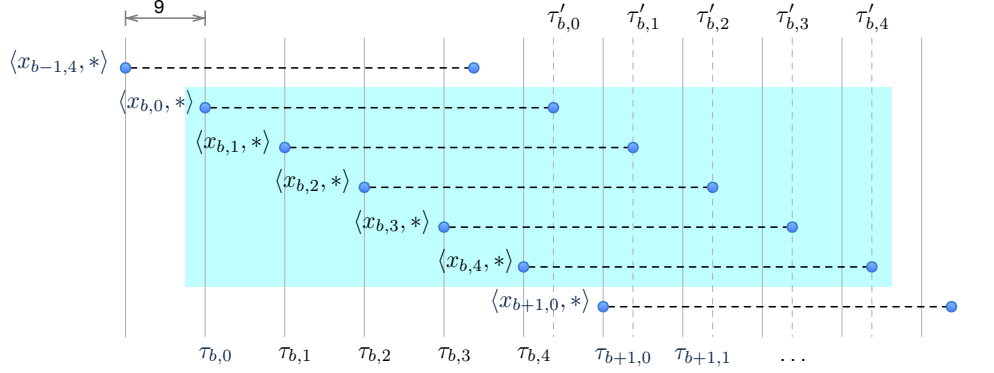


Figure 7: The sequence of vertex requests, for $n = 5$. The shaded region contains requests from bundle b .

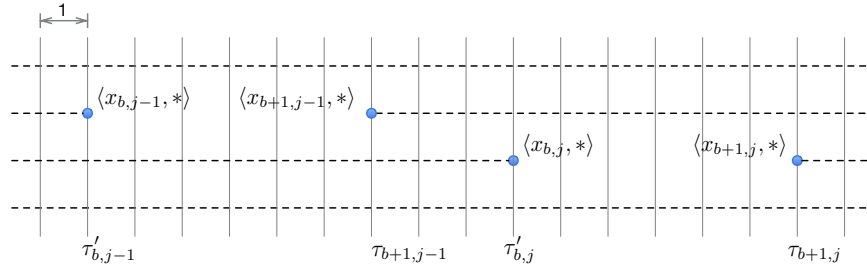


Figure 8: A more detailed picture showing relations between general requests to pages $x_{b,j-1}$, $x_{b,j}$, $x_{b+1,j-1}$, and $x_{b+1,j}$, where $1 \leq j \leq n - 1$.

846 **Lemma 6.2.** (a) The minimum number of faults on σ' in a cache of size $k + 2$ is F' . (b) If S is a solution for
847 σ' with at most F' faults, then for any $b = 0, 1, \dots, B - 2$ we have $V_{S,b} \preceq V_{S,b+1}$.

848 *Proof.* (a) There are $2B$ specific b -blocking requests $\langle y_b, k + 1 \rangle$ and $\langle y_b, k + 2 \rangle$ and all of these are faults.
849 Consider a bundle b . For this bundle, for each j , the two requests $\langle x_{b,j}, * \rangle$ at times $\tau_{b,j}$ and $\tau'_{b,j}$ are separated
850 by requests $\langle y_b, k + 1 \rangle$ and $\langle y_b, k + 2 \rangle$. Thus if S does not fault at time $\tau'_{b,j}$ then page $x_{b,j}$ must have been
851 stored in one of the vertex slots $1, 2, \dots, k$ throughout the time interval $[\tau_{b,j}, \tau'_{b,j}]$. As there are k vertex slots,
852 S can avoid faulting on at most k requests in bundle b . So, including the faults at $\langle y_b, k + 1 \rangle$ and $\langle y_b, k + 2 \rangle$,
853 the number of faults in S associated with this bundle b will be at least $2 + k + 2(n - k) = 2n - k + 2$. We
854 thus conclude that the total number of faults is at least F' .

855 It is also possible to achieve only F' faults on σ' , as follows: for each b , and for each vertex $j =$
856 $0, 1, \dots, k - 1$, at time $\tau_{b,j}$ load $x_{b,j}$ into cache slot $j + 1$ and keep it there until time $\tau'_{b,j}$. For $j = k, \dots, n - 1$,
857 load each request to $x_{b,j}$ into slot $k + 2$. This will give us exactly F' faults.

858 (b) If S makes at most F' faults, since there are $2B$ faults on the blocking requests and for each bundle
859 S makes at least $2n - k$ faults on vertex requests, S must make *exactly* $2n - k$ faults on vertex requests from
860 each bundle, including one request for each vertex $j \in V_{S,b}$ and two requests for each vertex $j \notin V_{S,b}$. If
861 $u \in V_{S,b}$ and $x_{b,u}$ is stored by S in slot ℓ of the cache throughout its interval $[\tau_{b,u}, \tau'_{b,u}]$, and if some $x_{b+1,v}$,
862 for $v \in V_{S,b+1}$, is stored by S in slot ℓ throughout its interval $[\tau_{b+1,v}, \tau'_{b+1,v}]$, then we must have $v \geq u$. This
863 is because otherwise we would have $\tau_{b+1,v} < \tau'_{b,u}$, that is the intervals of $x_{b,u}$ and $x_{b+1,v}$ would overlap, so
864 we would fault at least three times on the requests to these two pages. This implies part (b). \square

865 We partition all bundles into *phases*, where phase $p = 0, 1, \dots, P - 1$ consists of m bundles $b =$

866 $pm, pm + 1, \dots, pm + m - 1$. (Recall that $m = |E| + 1$.) The corollary below states that there is a phase p
867 in which all sets $V_{S,b}$ must be equal. It follows directly from Lemmas 6.1 and 6.2, by applying the pigeonhole
868 principle.

869 **Corollary 6.3.** *If S is a solution for σ' with F' faults, then there is index p , $0 \leq p \leq P - 1$, for which*
870 $V_{S,pm} = V_{S,pm+1} = \dots = V_{S,pm+m-1}$.

871 **Edge gadget.** For each fixed phase p , we create $m - 1$ edge gadgets, one for each edge. Ordering the edges
872 arbitrarily, the gadget for the e th edge, where $0 \leq e \leq m - 2$, will be denoted $\omega_{p,e}$, and it will consist of
873 8 requests between times θ'_{pm+e} and θ_{pm+e+1} , that is in the region where bundles $pm + b$ and $pm + b + 1$
874 overlap.

875 Let the e th edge be (u, v) , where $u < v$. Edge gadget $\omega_{p,e}$ uses six new pages $z_{p,u}, z_{p,v}, g_{p,u}, g_{p,v}, h_{p,u}$
876 and $h_{p,v}$, and consists of the following requests:

- 877 • Two specific requests $\langle z_{p,u}, k + 2 \rangle$ at times $\tau'_{pm+e,u} + 2$ and $\tau'_{pm+e,u} + 4$, and two specific requests
878 $\langle z_{p,v}, k + 2 \rangle$ at times $\tau'_{pm+e,v} + 2$ and $\tau'_{pm+e,v} + 4$.
- 879 • General requests $\langle g_{p,u}, * \rangle, \langle g_{p,v}, * \rangle$, at times $\tau'_{pm+e,u} + 3$ and $\tau'_{pm+e,v} + 3$.
- 880 • A pair of requests $\langle h_{p,u}, k + 1 \rangle, \langle h_{p,u}, * \rangle$, the first one specific and the second one general, at times
881 $\tau'_{pm+e,u} + 1$ and $\tau'_{pm+e,v} + 1$, respectively.
- 882 • A pair of requests $\langle h_{p,v}, * \rangle, \langle h_{p,v}, k + 1 \rangle$, the first one general and the second one specific, at times
883 $\tau'_{pm+e,u} + 5$ and $\tau'_{pm+e,v} + 5$, respectively.

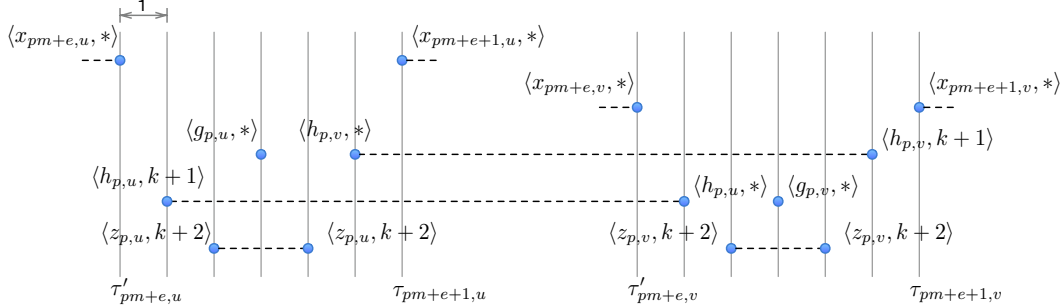


Figure 9: Gadget $\omega_{p,e}$. Requests $\langle x_{pm+e,u}, * \rangle, \langle x_{pm+e+1,u}, * \rangle, \langle x_{pm+e,v}, * \rangle$, and $\langle x_{pm+e+1,v}, * \rangle$ are not part of this gadget; they are shown only to illustrate how gadget $\omega_{p,e}$ fits into the overall request sequence.

884 Consider now possible solutions of gadget $\omega_{p,e}$. Notice that this gadget will require 7 faults regardless
885 of all other requests, since we need to make two faults on requests $\langle g_{p,u}, * \rangle, \langle g_{p,v}, * \rangle$, at least two faults
886 on requests to pages $\langle z_{p,u}, k + 2 \rangle, \langle z_{p,v}, k + 2 \rangle$, and at least three faults on requests $\langle h_{p,u}, k + 1 \rangle, \langle h_{p,u}, * \rangle,$
887 $\langle h_{p,v}, * \rangle$, and $\langle h_{p,v}, k + 1 \rangle$. (This is because if we retain page $h_{p,u}$ in slot $k + 1$ until time $\tau'_{pm+e,v} + 1$, so
888 that we do not fault on $\langle h_{p,u}, * \rangle$, then we will fault on both requests $\langle h_{p,v}, * \rangle$, and $\langle h_{p,v}, k + 1 \rangle$.) Another
889 important observation is that if we fault only 7 times on $\omega_{p,e}$ then one of requests $\langle g_{p,u}, * \rangle, \langle g_{p,v}, * \rangle$ must be
890 put in a vertex cache slot (that is, one of slots $1, 2, \dots, k$). A solution that puts $\langle g_{p,u}, * \rangle$ in a vertex slot is
891 called a u -solution of $\omega_{p,e}$ and a solution that puts $\langle g_{p,v}, * \rangle$ in a vertex slot is called a v -solution of $\omega_{p,e}$. (A
892 solution of $\omega_{p,e}$ can be both a u -solution and a v -solution.)

893 **Complete reduction.** Let $F = F' + 7P(m - 1)$. Our request sequence σ constructed for G consists of σ'
894 and of all $P(m - 1)$ edge gadgets $\omega_{p,e}$ defined above inserted into σ at their specified time steps. (At some

895 time steps there will not be any requests.) To complete the proof it is now sufficient to show the following
 896 claim.

897 **Claim 6.4.** *G has a vertex cover of size k if and only if σ has a solution with at most F faults in a cache of*
 898 *size $k + 2$.*

899 (\Rightarrow) Suppose that G has a vertex cover U of size k . We construct a solution for σ as follows. Each vertex
 900 $j \in U$ is assigned to some unique vertex cache slot and all $2B$ requests $\langle x_{b,j}, * \rangle$ associated with vertex j are
 901 served in this slot. This will create kB faults. For $j \notin U$, all requests $\langle x_{b,j}, * \rangle$ are served in the junkyard slot
 902 $k + 2$ at cost $2(n - k)B$. Together with the $2B$ blocking requests $\langle y_b, k + 1 \rangle$ and $\langle y_b, k + 2 \rangle$, this will give
 903 us $F' = (2n - k + 2)B$ faults. For each $p = 0, 1, \dots, P - 1$, and for each $e = 0, 1, \dots, m - 2$, we do this:
 904 Let the e th edge of G be (u, v) . Since U is a vertex cover, we either have $u \in U$ or $v \in U$. If $u \in U$, we
 905 use the u -solution for gadget $\omega_{p,e}$, with $g_{p,u}^*$ served in the cache slot associated with u . If $v \in U$, we use the
 906 v -solution for gadget $\omega_{p,e}$, with $\langle g_{p,v}, * \rangle$ served in the cache slot associated with v . This will give us 7 faults
 907 for this gadget, adding up to $7P(m - 1)$ faults on all edge gadgets. Then the total number of faults on σ will
 908 be $F' + 7P(m - 1) = F$.

909 (\Leftarrow) Now suppose that there is a solution S for σ with at most F faults. By the earlier observations,
 910 we know that S must have exactly F faults, including exactly F' faults on the request in σ' and exactly 7
 911 faults per each edge gadget. As there are F' faults on σ' , we can find some p , $0 \leq p \leq P - 1$, such that
 912 $V_{S,pm} = V_{S,pm+1} = \dots = V_{S,pm+m-1}$, per Corollary 6.3. Let $U = V_{S,pm}$. For each $j \in U$, all requests
 913 $\langle x_{b,j}, * \rangle$, for $b = pm, pm + 1, \dots, pm + m - 1$, must be in the same slot, that we refer to as the slot associated
 914 with vertex j . The size of U is k , and we claim that U must be a vertex cover. To show this, let (u, v) be an
 915 edge, and let e be its index. Solution S makes 7 faults on gadget $\omega_{p,e}$, so for this gadget it must be either a
 916 u -solution or a v -solution. If it is a u -solution then $\langle g_{p,u}, * \rangle$ is served in some vertex cache slot. But the only
 917 vertex slot available in that time step is the slot associated with vertex u . This means that u must be in U .
 918 The case of a v -solution is symmetric. Thus we obtain that either $u \in U$ or $v \in U$. This holds for each edge,
 919 implying that U is a vertex cover. This proves the claim, and completes the proof of the theorem. \square

920 7 Weighted All-Or-One Paging

921 This section initiates the study of Heterogenous k -Server in non-uniform metrics. Weighted All-Or-One
 922 Paging is the natural weighted extension of All-or-One Paging (allowing general and specific requests) in
 923 which the pages have weights and the cost of retrieving a page is its weight. This is equivalent to Heteroge-
 924 nous k -Server in star metrics with requestable-set family $\mathcal{S} = \{[k]\} \cup \{\{s\} : s \in [k]\}$. This section proves
 925 the following theorem:

926 **Theorem 7.1.** *Weighted All-Or-One Paging has a deterministic $O(k)$ -competitive online algorithm.*

927 The bound is optimal up to a small constant factor, as the optimal ratio for standard Weighted Paging
 928 is k . Figure 10 shows the algorithm. It is implicitly a linear-programming primal-dual algorithm. Note
 929 that the standard linear program for standard Weighted Paging doesn't have constraints that force pages into
 930 specific slots—indeed, those constraints make even the unweighted problem an NP -hard special case of
 931 Multicommodity Flow. As a small example that illustrates the challenge, consider a cache of size $k = 2$,
 932 and repeatedly make three requests: a general request to a weight-1 page, and specific requests to different
 933 weight-zero pages in slots 1 and 2. The weight-zero requests force the weight-1 page to be evicted with each
 934 round, so the optimal cost is the number of rounds. But the solution of the classical linear-program relaxation
 935 will have value 1. Thus this linear program cannot be used to bound the competitive ratio.

936 Here is a sketch of the proof of Theorem 7.1, then the detailed proof. Fix an optimal solution C , that
 937 is $\text{opt}(\sigma) = \text{cost}(C)$. For each $t \in [T]$, let $x_t \in \{0, 1\}$ be an indicator variable for the event that C evicts

input: Weighted All-Or-One Paging instance (k, σ) , where $\sigma_t = \langle p_t, s_t \rangle$ for $t \in [T]$

1. initialize $\text{cap}[t] \leftarrow \text{credit}[t] \leftarrow 0$ for each $t \in [T]$
2. assume that $\langle p_t, s_t \rangle = \langle 0, t \rangle$ for $t \in [k]$ — k specific requests to artificial weight-0 page in each slot
3. for $t \leftarrow k + 1, k + 2, \dots, T$:
 - 3.1. if $\langle p_t, s_t \rangle$ is a specific request with no equivalent request t' (s.t. $\langle p_{t'}, s_{t'} \rangle = \langle p_t, s_t \rangle$) in the cache:
 - 3.1.1. evict any cached general request to page p_t , and any cached request in slot s_t
 - 3.1.2. put t in slot s_t — note $\text{cap}[t] = \text{credit}[t] = 0$
 - 3.2. else if $\langle p_t, s_t \rangle$ is a general request not satisfied by any cached request t' (s.t. $p_{t'} = p_t$):
 - 3.2.1. define $\begin{cases} \ell_t(s) := \max\{t' \leq t : s_{t'} = s\} \text{ for } s \in [k] & \text{— most recent specific request to slot } s \\ A := \{s \in [k] : \text{cap}[\ell_t(s)] \geq \frac{1}{2} \text{wt}(p_t) \text{ and } s \text{ does not hold a specific request}\} \\ B := \{s \in [k] : \text{slot } s \text{ holds a general request of weight at least } \frac{1}{2} \text{wt}(p_t)\} \end{cases}$
 - 3.2.2. while $|A| \leq |B|$:
 - 3.2.2.1. continuously raise $\text{cap}[\ell_t(s)]$ for $s \in [k]$ and $\text{credit}[t']$ for each cached request t' , at unit rate,
 - 3.2.2.2. evicting each request t' such that $\text{credit}[t'] = \text{wt}(p_{t'})$, and updating A and B continuously
 - 3.2.3. choose a slot $s \in A \setminus B$; evict the request t' currently in slot s (if any)
 - 3.2.4. put t in slot s — note $\text{credit}[t] = 0$
 - 3.3. else: classify the (already satisfied) request as *redundant* and ignore it

Figure 10: An $O(k)$ -competitive online algorithm for Weighted All-Or-One Paging. For technical convenience, we present the algorithm as caching request times rather than pages, with the understanding that request t represents page p_t .

938 request t before satisfying another request $t' > t$ with the same page/slot pair that satisfied t . Let $R \subseteq [T]$
 939 be the set of all specific requests, and for each $t \in R$, let y_t be the amount C pays to retrieve pages into
 940 slot s_t before the next specific request to slot s_t (if any). Define the *pseudo-cost* of the optimal solution to
 941 be $\sum_{t=1}^T \text{wt}(p_t)x_t + \sum_{t \in R} y_t$. The pseudo-cost is at most $2 \text{opt}(\sigma)$. As the algorithm proceeds, define the
 942 *residual cost* to be $\sum_{t=1}^T \max(0, \text{wt}(p_t)x_t - \text{credit}[t]) + \sum_{t \in R} \max(0, y_t - \text{cap}[t])$. The residual cost is
 943 initially the pseudo-cost (at most $2 \text{opt}(\sigma)$), and remains non-negative throughout, so the total decrease in the
 944 residual cost is at most $2 \text{opt}(\sigma)$. One can show (Lemma 7.1) that whenever the algorithm is raising credits
 945 and capacities at time t , there is either a cached request t' with $x_{t'} = 1$ and $\text{credit}[t'] < \text{wt}(p_{t'})$, or there is a
 946 slot s with $y_{t'} > \text{cap}[t']$, where $t' = \ell_t(s) \in R$. It follows that the residual cost is decreasing at least at unit
 947 rate in Step 3.2.2.1.

948 On the other hand, the algorithm is raising k capacities and at most k credits, so the value of $\phi =$
 949 $\sum_{t=1}^T \text{credit}[t] + \sum_{t \in R} \text{cap}[t]$ is increasing at rate at most $2k$. So, the final value of ϕ is at most $4k \text{opt}(\sigma)$. To
 950 finish, we show by a charging argument that the algorithm's cost is at most $6\phi + 3 \text{opt}(\sigma) \leq (24k + 3) \text{opt}(\sigma)$.

951 Here is the detailed proof. Consider any execution of the algorithm on a k -slot instance σ . To ease
 952 notation and streamline the analysis, without loss of generality we will make the following assumptions:

- 953 • The first k requests are specific requests for an artificial weight-zero page in each of the k slots.
- 954 • Each request is not redundant (per Step 3.3).
- 955 • The last k requests are specific requests for an artificial weight-zero page in each of the k slots.

956 These assumptions can indeed be made without loss of generality, as the zero-weight requests do not have

957 any cost, the algorithm ignores redundant requests, and removing redundant requests doesn't increase the
 958 optimum cost. For technical convenience, we think of the algorithm and the optimal solution as caching
 959 request *times* rather than pages, with the understanding that request t represents page p_t . We first prove a key
 960 lemma used in the proof of the theorem.

961 **Lemma 7.1.** *Suppose that, while responding to a general request t , the algorithm is executing Step 3.2.2.1*
 962 *(that is, the loop condition in Step 3.2.2 is satisfied). Then, in any solution C , just after C has responded to*
 963 *request t , either*

- 964 (i) C has evicted some request t' currently cached by the algorithm, or
 965 (ii) for some slot $s \in [k]$, after the most recent specific request $\ell_t(s)$ to slot s solution C has incurred cost
 966 more than $\text{cap}[\ell_t(s)]$ for retrievals into s .

967 *Proof.* If C satisfies property (i), we are done. So assume that it doesn't, and we will show that then
 968 property (ii) holds. If (i) doesn't hold then, just after responding to request t , in addition to the current general
 969 request p_t , solution C caches every request t' that is cached by the algorithm. This, together with the loop
 970 condition, implies that C has at least $|B| + 1 \geq |A| + 1$ generally requested pages of weight at least $\frac{1}{2} \text{wt}(p_t)$
 971 in its cache. Thus one of these pages, say $p_{t'}$, is in a slot $s \notin A$. The choice of $p_{t'}$ and the definition of A
 972 imply then that the cost of C for retrievals into s after time $\ell_t(s)$ is at least $\text{wt}(p_{t'}) \geq \frac{1}{2} \text{wt}(p_t) > \text{cap}[\ell_t(s)]$,
 973 so property (ii) holds. \square

974 *Proof of Theorem 7.1.* Fix an optimal solution C , that is $\text{opt}(\sigma) = \text{cost}(C)$. For each $t \in [T]$, let $x_t \in \{0, 1\}$
 975 be an indicator variable for the event that C evicts request t before satisfying another request $t' > t$ with
 976 the same page/slot pair that satisfied t . Let $R \subseteq [T]$ be the set of all specific requests, and for each $t \in R$,
 977 let y_t be the amount C pays to retrieve pages into slot s_t before the next specific request to slot s_t (if any).
 978 Define the *pseudo-cost* of the optimal solution to be $\sum_{t=1}^T \text{wt}(p_t)x_t + \sum_{t \in R} y_t$. The pseudo-cost is at most
 979 $2 \text{opt}(\sigma)$. As the algorithm proceeds, define the *residual cost* to be $\sum_{t=1}^T \max(0, \text{wt}(p_t)x_t - \text{credit}[t]) +$
 980 $\sum_{t \in R} \max(0, y_t - \text{cap}[t])$. The residual cost is initially the pseudo-cost (at most $2 \text{opt}(\sigma)$), and remains
 981 non-negative throughout, so the total decrease in the residual cost is at most $2 \text{opt}(\sigma)$. By Lemma 7.1,²
 982 whenever the algorithm is raising credits and capacities at time t , there is either a cached request t' with
 983 $x_{t'} = 1$ and $\text{credit}[t'] < \text{wt}(p_{t'})$, or there is a slot s with $y_{t'} > \text{cap}[t']$, where $t' = \ell_t(s) \in R$. It follows that
 984 the residual cost is decreasing at least at unit rate in Step 3.2.2.1.

985 On the other hand, the algorithm is raising k capacities and at most k credits, so the value of $\phi =$
 986 $\sum_{t=1}^T \text{credit}[t] + \sum_{t \in R} \text{cap}[t]$ is increasing at rate at most $2k$. So, the final value of ϕ is at most $4k \text{opt}(\sigma)$.
 987 To finish, we show that the algorithm's cost is at most $6\phi + 3 \text{opt}(\sigma) \leq (24k + 3) \text{opt}(\sigma)$.³ Count the costs
 988 that the algorithm pays as follows:

- 989 1. *Requests remaining in the cache at the end (time T).* By the assumption on the last k requests, these
 990 cost nothing to bring in. All other requests are evicted.
- 991 2. *Requests evicted in Line 3.2.2.2.* Each such request t' is evicted only after $\text{credit}[t']$ reaches $\text{wt}(p_{t'})$.
 992 So these have total weight at most $\sum_{t'=1}^T \text{credit}[t']$.
- 993 3. *Specific requests t' evicted from slot s_t in Line 3.1.1.* Throughout the time interval $[t', t - 1]$, the
 994 algorithm has $p_{t'}$ in slot $s_{t'} = s_t$, and σ has neither an equivalent specific request nor a general request
 995 to p_t (by our non-redundancy assumption). The optimal solution C has $p_{t'}$ in slot $s_{t'}$ at time t' , but not
 996 at time t , so evicts it during $[t' + 1, t]$. So the total cost of such requests is at most the total weight of
 997 specific requests evicted by C , and thus at most $\text{opt}(\sigma)$.

²The lemma gives linear constraints on the vectors (x, y) . Minimizing the pseudo-cost subject to these constraints is a linear-program relaxation of the problem. Our argument implicitly defines a dual solution whose cost is our lower bound on $\text{opt}(\sigma)$.

³This constant can be reduced with more careful analysis.

- 998 4. *General requests evicted from slot s_t in Line 3.1.1.* By Line 3.2.3, any general request in slot s_t at time
999 t has weight at most $2 \text{cap}[\ell_{t-1}(s_t)]$. So the total weight of such requests is at most $2 \sum_{t' \in R} \text{cap}[t']$.
- 1000 5. *General requests to page p_t evicted in Line 3.1.1.* The algorithm replaces each such general request t'
1001 by a specific request t (which it later evicts, unless the weight is zero) to the same page. Have general
1002 request t' charge its cost $\text{wt}(p_{t'}) = \text{wt}(p_t)$, and any amount charged to t' (in Item 6 below), to specific
1003 request t . (We analyze the charging scheme for Items 5 and 6 below.)
- 1004 6. *General requests t' evicted in Line 3.2.3.* Have request t' charge the cost of its eviction, and any
1005 amount charged to t' to request t . Since the slot holding $p_{t'}$ is not in B , $\text{wt}(p_{t'}) < \frac{1}{2} \text{wt}(p_t)$.

1006 Each general request t receives at most one charge in Item 6, from a request t' of at most half the weight
1007 of t ; this general request t' may also receive such charges, forming a chain of charges, but since the weights
1008 of the requests in this chain decrease geometrically, t is charged at most its weight. In Item 5, each specific
1009 request t is charged by at most one general request t' of the same weight, that may also carry the chain charge
1010 not exceeding its weight. So this specific request is charged at most twice its weight. Overall, the charge of
1011 each request from Items 5 and 6 is at most twice its weight.

1012 The total weight of evictions considered in Items 1, 2, 3, and 4 is at most $2\phi + \text{opt}(\sigma)$. Adding also the
1013 charges to these items by evictions considered in Items 5 and 6, we obtain that the total cost of the algorithm
1014 is bounded by $3(2\phi + \text{opt}(\sigma)) = 6\phi + 3\text{opt}(\sigma)$. \square

1015 8 Open Problems

1016 The results here suggest many open problems and avenues for further research. Closing or tightening gaps
1017 left by our upper and lower bounds would be of interest. In particular:

- 1018 • For Slot-Heterogenous Paging, is the upper bound in Theorem 3.1 tight for every $\mathcal{S} \subseteq 2^{[k]} \setminus \{\emptyset\}$,
1019 within $\text{poly}(k)$ factors?
- 1020 • For Page-Laminar Paging it is easy to show a lower bound of $\Omega(h)$, even for $k = 1$ and for randomized
1021 algorithms. But it still may be possible to eliminate or reduce the multiplicative dependence on h . For
1022 example, is it possible to achieve ratio $O(h + k)$ with a deterministic algorithm and $O(h + H_k)$
1023 with a randomized algorithm? Similarly, does Slot-Laminar Paging (where $h \leq k$) admit an $O(k)$
1024 deterministic ratio and $O(\log k)$ randomized ratio?
- 1025 • For deterministic All-or-One Paging, we conjecture that the optimal ratio is $2k - 1$. (For $k = 2$ we
1026 can show an upper bound of 3.) In the randomized case, can ratio $2H_k - 1$ be achieved?
- 1027 • For Weighted All-Or-One Paging, is the optimal randomized ratio $O(\text{polylog}(k))$?
- 1028 • The status of Heterogenous k -Server in arbitrary metric spaces is wide open. Can ratio dependent only
1029 on k be achieved? This question, while challenging, could still be easier to resolve for Heterogenous
1030 k -Server than for Generalized k -Server.

1031 References

- 1032 [1] Dimitris Achlioptas, Marek Chrobak, and John Noga. Competitive analysis of randomized paging
1033 algorithms. *Theor. Comput. Sci.*, 234(1-2):203–218, 2000. doi:10.1016/S0304-3975(98)
1034 00116-9.

- 1035 [2] C. J. Argue, Anupam Gupta, Ziyi Tang, and Guru Guruganesh. Chasing convex bodies with linear
1036 competitive ratio. *Journal of the ACM*, 68(5):32:1–32:10, August 2021. doi:10.1145/3450349.
- 1037 [3] Nikhil Ayyadevara and Ashish Chiplunkar. The randomized competitive ratio of weighted k -server
1038 is at least exponential. *CoRR*, abs/2102.11119, 2021. URL: [https://arxiv.org/abs/2102.](https://arxiv.org/abs/2102.11119)
1039 [11119](https://arxiv.org/abs/2102.11119), arXiv:2102.11119.
- 1040 [4] Nikhil Bansal, Niv Buchbinder, Aleksander Mądry, and Joseph Naor. A polylogarithmic-competitive
1041 algorithm for the k -server problem. In Rafail Ostrovsky, editor, *IEEE 52nd Annual Symposium*
1042 *on Foundations of Computer Science, FOCS 2011, Palm Springs, CA, USA, October 22-25, 2011*,
1043 pages 267–276. IEEE Computer Society, 2011. URL: [http://ieeexplore.ieee.org/xpl/](http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=6108120)
1044 [mostRecentIssue.jsp?punumber=6108120](http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=6108120), doi:10.1109/FOCS.2011.63.
- 1045 [5] Nikhil Bansal, Niv Buchbinder, and Joseph Naor. A primal-dual randomized algorithm for
1046 weighted paging. In *48th Annual IEEE Symposium on Foundations of Computer Science (FOCS*
1047 *2007), October 20-23, 2007, Providence, RI, USA, Proceedings*, pages 507–517. IEEE Computer
1048 Society, 2007. URL: [http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?](http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=4389466)
1049 [punumber=4389466](http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=4389466), doi:10.1109/FOCS.2007.7.
- 1050 [6] Nikhil Bansal, Marek Eliás, and Grigorios Koumoutsos. Weighted k -server bounds via combinatorial
1051 dichotomies. In Chris Umans, editor, *58th IEEE Annual Symposium on Foundations of Computer Sci-*
1052 *ence, FOCS 2017, Berkeley, CA, USA, October 15-17, 2017*, pages 493–504. IEEE Computer Society,
1053 2017. doi:10.1109/FOCS.2017.52.
- 1054 [7] Nikhil Bansal, Marek Eliáš, Grigorios Koumoutsos, and Jesper Nederlof. Competitive algorithms
1055 for generalized k -server in uniform metrics. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM*
1056 *Symposium on Discrete Algorithms, SODA 2018, New Orleans, LA, USA, January 7-10, 2018*, pages
1057 992–1001, 2018. doi:10.1137/1.9781611975031.64.
- 1058 [8] Nikhil Bansal, Joseph (Seffi) Naor, and Ohad Talmon. Efficient online weighted multi-level paging. In
1059 Kunal Agrawal and Yossi Azar, editors, *SPAA '21: 33rd ACM Symposium on Parallelism in Algorithms*
1060 *and Architectures, Virtual Event, USA, 6-8 July, 2021*, pages 94–104. ACM, 2021. doi:10.1145/
1061 3409964.3461801.
- 1062 [9] Nathan Beckmann, Phillip B. Gibbons, Bernhard Haeupler, and Charles McGuffey. Writeback-aware
1063 caching. In Bruce M. Maggs, editor, *1st Symposium on Algorithmic Principles of Computer Systems,*
1064 *APOCS 2020, Salt Lake City, UT, USA, January 8, 2020*, pages 1–15. SIAM, 2020. doi:10.1137/
1065 1.9781611976021.1.
- 1066 [10] L. A. Belady. A study of replacement algorithms for a virtual-storage computer. *IBM Systems Journal*,
1067 5(2):78–101, 1966. doi:10.1147/sj.52.0078.
- 1068 [11] Laszlo A. Belady. A study of replacement algorithms for virtual-storage computer. *IBM Syst. J.*,
1069 5(2):78–101, 1966. doi:10.1147/sj.52.0078.
- 1070 [12] Marcin Bienkowski, Łukasz Jeż, and Pawel Schmidt. Slaying Hydrae: Improved bounds for general-
1071 ized k -server in uniform metrics. In Pinyan Lu and Guochuan Zhang, editors, *30th International Sympo-*
1072 *sium on Algorithms and Computation, ISAAC 2019*, volume 149 of *LIPICs*, pages 14:1–14:14. Schloss
1073 Dagstuhl - Leibniz-Zentrum für Informatik, 2019. doi:10.4230/LIPICs.ISAAC.2019.14.
- 1074 [13] Allan Borodin and Ran El-Yaniv. *Online computation and competitive analysis*. Cambridge University
1075 Press, 1998.

- 1076 [14] Allan Borodin and Ran El-Yaniv. On randomization in on-line computation. *Inf. Comput.*, 150(2):244–
1077 267, 1999. doi:10.1006/inco.1998.2775.
- 1078 [15] Allan Borodin, Nathan Linial, and Michael E. Saks. An optimal on-line algorithm for metrical task
1079 system. *J. ACM*, 39(4):745–763, 1992. doi:10.1145/146585.146588.
- 1080 [16] Mark Brehob, Richard J. Enbody, Eric Torng, and Stephen Wagner. On-line restricted caching. *J.*
1081 *Sched.*, 6(2):149–166, 2003. doi:10.1023/A:1022989909868.
- 1082 [17] Mark Brehob, Stephen Wagner, Eric Torng, and Richard J. Enbody. Optimal replacement is NP-hard
1083 for nonstandard caches. *IEEE Trans. Computers*, 53(1):73–76, 2004. doi:10.1109/TC.2004.
1084 1255792.
- 1085 [18] Sébastien Bubeck, Michael B. Cohen, Yin Tat Lee, James R. Lee, and Aleksander Mądry. K-server via
1086 multiscale entropic regularization. In Ilias Diakonikolas, David Kempe, and Monika Henzinger, edi-
1087 tors, *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018,*
1088 *Los Angeles, CA, USA, June 25-29, 2018*, pages 3–16. ACM, 2018. doi:10.1145/3188745.
1089 3188798.
- 1090 [19] Sébastien Bubeck, Yuval Rabani, and Mark Sellke. Online multiserver convex chasing and optimiza-
1091 tion. In *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 2093–
1092 2104. SIAM, 2021.
- 1093 [20] Niv Buchbinder, Shahar Chen, and Joseph Naor. Competitive algorithms for restricted caching and
1094 matroid caching. In Andreas S. Schulz and Dorothea Wagner, editors, *Algorithms - ESA 2014 -*
1095 *22th Annual European Symposium, Wroclaw, Poland, September 8-10, 2014. Proceedings*, volume
1096 8737 of *Lecture Notes in Computer Science*, pages 209–221. Springer, 2014. doi:10.1007/
1097 978-3-662-44777-2_18.
- 1098 [21] Niv Buchbinder, Christian Coester, and Joseph (Seffi) Naor. Online k-taxi via double coverage and
1099 time-reverse primal-dual. In Mohit Singh and David P. Williamson, editors, *Integer Programming and*
1100 *Combinatorial Optimization - 22nd International Conference, IPCO 2021, Atlanta, GA, USA, May 19-*
1101 *21, 2021, Proceedings*, volume 12707 of *Lecture Notes in Computer Science*, pages 15–29. Springer,
1102 2021. doi:10.1007/978-3-030-73879-2_2.
- 1103 [22] Jannik Castenow, Björn Feldkord, Till Knollmann, Manuel Malatyali, and Friedhelm Meyer auf der
1104 Heide. The k-server with preferences problem. In *SPAA '22: 34rd ACM Symposium on Parallelism in*
1105 *Algorithms and Architectures, 2022*. To appear. URL: <https://arxiv.org/abs/2205.11102>.
- 1106 [23] Ashish Chiplunkar and Sundar Vishwanathan. Metrical service systems with multiple servers. *Algo-*
1107 *rithmica*, 71(1):219–231, 2015. doi:10.1007/s00453-014-9903-7.
- 1108 [24] Ashish Chiplunkar and Sundar Vishwanathan. Randomized memoryless algorithms for the weighted
1109 and the generalized k-server problems. *ACM Trans. Algorithms*, 16(1), December 2019. doi:10.
1110 1145/3365002.
- 1111 [25] Marek Chrobak, Samuel Haney, Mehraneh Liaee, Debmalya Panigrahi, Rajmohan Rajaraman, Ravi
1112 Sundaram, and Neal E. Young. Online paging with heterogeneous cache slots. In *40th International*
1113 *Symposium on Theoretical Aspects of Computer Science (STACS 2023)*, 2023.
- 1114 [26] Marek Chrobak and John Noga. Competitive algorithms for relaxed list update and multilevel caching.
1115 *J. Algorithms*, 34(2):282–308, 2000. doi:10.1006/jagm.1999.1061.

- 1116 [27] Marek Chrobak, Gerhard J. Woeginger, Kazuhisa Makino, and Haifeng Xu. Caching is hard — even
1117 in the fault model. *Algorithmica*, 63(4):781–794, 2012.
- 1118 [28] Christian Coester and Elias Koutsoupias. The online k -taxi problem. In Moses Charikar and Edith
1119 Cohen, editors, *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing,*
1120 *STOC 2019, Phoenix, AZ, USA, June 23-26, 2019*, pages 1136–1147. ACM, 2019. doi:10.1145/
1121 3313276.3316370.
- 1122 [29] Christian Coester and Elias Koutsoupias. Towards the k -server conjecture: A unifying potential,
1123 pushing the frontier to the circle. In *48th International Colloquium on Automata, Languages,*
1124 *and Programming, ICALP 2021, July 12-16, 2021, Glasgow, Scotland (Virtual Conference)*, vol-
1125 ume 198 of *LIPICs*, pages 57:1–57:20. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2021.
1126 doi:10.4230/LIPICs.ICALP.2021.57.
- 1127 [30] Leonid Domnitser, Aamer Jaleel, Jason Loew, Nael Abu-Ghazaleh, and Dmitry Ponomarev. Non-
1128 monopolizable caches: Low-complexity mitigation of cache side channel attacks. *ACM Trans. Archit.*
1129 *Code Optim.*, 8(4), January 2012.
- 1130 [31] Esteban Feuerstein. Uniform service systems with k servers. In Gerhard Goos, Juris Hartmanis, Jan
1131 van Leeuwen, Cláudio L. Lucchesi, and Arnaldo V. Moura, editors, *LATIN'98: Theoretical Informatics*,
1132 volume 1380, pages 23–32, Berlin, Heidelberg, 1998. Springer Berlin Heidelberg. doi:10.1007/
1133 Bf0054307.
- 1134 [32] Amos Fiat, Richard M. Karp, Michael Luby, Lyle A. McGeoch, Daniel Dominic Sleator, and Neal E.
1135 Young. Competitive paging algorithms. *J. Algorithms*, 12(4):685–699, 1991. doi:10.1016/
1136 0196-6774(91)90041-v.
- 1137 [33] Amos Fiat, Manor Mendel, and Steven S. Seiden. Online companion caching. In Rolf Möhring and
1138 Rajeev Raman, editors, *Algorithms — ESA 2002*, pages 499–511, Berlin, Heidelberg, 2002. Springer
1139 Berlin Heidelberg.
- 1140 [34] Amos Fiat and Moty Ricklin. Competitive algorithms for the weighted server problem. *Theor. Comput.*
1141 *Sci.*, 130(1):85–99, 1994. doi:10.1016/0304-3975(94)90154-6.
- 1142 [35] Samuel Haney. *Algorithms for Networks With Uncertainty*. PhD thesis, Duke University, 2019. URL:
1143 <https://dukespace.lib.duke.edu/dspace/handle/10161/18661>.
- 1144 [36] Shahin Kamali and Helen Xu. Multicore paging algorithms cannot be competitive. In *Proceedings of*
1145 *the 32nd ACM Symposium on Parallelism in Algorithms and Architectures*, pages 547–549, 2020.
- 1146 [37] Anna R. Karlin, Mark S. Manasse, Larry Rudolph, and Daniel Dominic Sleator. Competitive snoopy
1147 caching. *Algorithmica*, 3:77–119, 1988. doi:10.1007/BF01762111.
- 1148 [38] Elias Koutsoupias and Christos H. Papadimitriou. On the k -server conjecture. *J. ACM*, 42(5):971–983,
1149 1995. doi:10.1145/210118.210128.
- 1150 [39] Elias Koutsoupias and David Scot Taylor. The CNN problem and other k -server variants. *Theor.*
1151 *Comput. Sci.*, 324(2-3):347–359, 2004. doi:10.1016/j.tcs.2004.06.002.
- 1152 [40] Fangfei Liu, Qian Ge, Yuval Yarom, Frank Mckeen, Carlos Rozas, Gernot Heiser, and Ruby B. Lee.
1153 CATalyst: Defeating last-level cache side channel attacks in cloud computing. In *2016 IEEE Inter-*
1154 *national Symposium on High Performance Computer Architecture (HPCA)*, pages 406–418, 2016.
1155 doi:10.1109/HPCA.2016.7446082.

- 1156 [41] Mark S. Manasse, Lyle A. McGeoch, and Daniel Dominic Sleator. Competitive algorithms for server
1157 problems. *J. Algorithms*, 11(2):208–230, 1990. doi:10.1016/0196-6774(90)90003-w.
- 1158 [42] Lyle A. McGeoch and Daniel Dominic Sleator. A strongly competitive randomized paging algorithm.
1159 *Algorithmica*, 6(6):816–825, 1991. doi:10.1007/BF01759073.
- 1160 [43] M. Mendel and Steven S. Seiden. Online companion caching. *Theoretical Computer Science*,
1161 324(2):183–200, 2004. Online Algorithms: In Memoriam, Steve Seiden. doi:https://doi.
1162 org/10.1016/j.tcs.2004.05.015.
- 1163 [44] Jignesh Patel. Restricted k -server problem. Master’s thesis, Michigan State University, 2004. URL:
1164 https://d.lib.msu.edu/etd/32678.
- 1165 [45] Mark Sellke. Chasing convex bodies optimally. In *Proceedings of the 2020 ACM-SIAM Symposium*
1166 *on Discrete Algorithms (SODA)*, Proceedings, pages 1509–1518. Society for Industrial and Applied
1167 Mathematics, 2020. doi:10.1137/1.9781611975994.92.
- 1168 [46] Daniel Dominic Sleator and Robert Endre Tarjan. Amortized efficiency of list update and paging rules.
1169 *Commun. ACM*, 28(2):202–208, 1985. doi:10.1145/2786.2793.
- 1170 [47] Zhenghong Wang and Ruby B. Lee. New cache designs for thwarting software cache-based side chan-
1171 nel attacks. In *Proceedings of the 34th Annual International Symposium on Computer Architecture*,
1172 ISCA ’07, page 494–505, New York, NY, USA, 2007. doi:10.1145/1250662.1250723.
- 1173 [48] Ying Ye, Richard West, Zhuoqun Cheng, and Ye Li. COLORIS: A dynamic cache partitioning system
1174 using page coloring. In *2014 23rd International Conference on Parallel Architecture and Compilation*
1175 *Techniques (PACT)*, pages 381–392, 2014. doi:10.1145/2628071.2628104.
- 1176 [49] Wei Zang and Ann Gordon-Ross. CaPPS: cache partitioning with partial sharing for multi-
1177 core embedded systems. *Des. Autom. Embed. Syst.*, 20(1):65–92, 2016. doi:10.1007/
1178 s10617-015-9168-7.