

CS612 Algorithm Design and Analysis

Lecture 16. Paging problem ¹

Dongbo Bu

Institute of Computing Technology
Chinese Academy of Sciences, Beijing, China

¹The slides are made based on Algorithm Design, Randomized algorithm by R. Motwani and P. Raghavan, and a lecture by T. Chan.

Outline

- Introduction
- Greedy algorithm: Furthest-Future principle;
- Label on-line algorithm framework;
- The performance of LRU principle;
- A randomized on-line algorithm (Fiat et al '91);

PAGING problem I

INPUT:

Given a sequence of requests r_1, r_2, \dots, r_n , and a cache of size k ;

OUTPUT:

schedule the eviction decisions to reduce cache-missing as much as possible.

An example I

An eviction sequence: see a fig.

A dynamic-programming method I

- Subproblem: finding the optimal evictions for requests $r_i \dots r_n$ when cache contents are c_1, c_2, \dots, c_k , $c_i \in \Sigma$, $|\Sigma| = m$.
- Let $OPT(i, c_1, c_2, \dots, c_k)$ be the optimal solution value to the subproblem. We have the following recursion:

$$OPT(i, c_1, c_2, \dots, c_k) = \min \begin{cases} OPT(i+1, c'_1, c'_2, \dots, c'_k) + 1 \\ OPT(i+1, c_1, c_2, \dots, c_k) \end{cases} \quad \text{and}$$

$$OPT(n, c_1, c_2, \dots, c_k) = 0.$$

Here, c'_1, c'_2, \dots, c'_k differs from c_1, c_2, \dots, c_k at only one page.

Time-complexity: DP table size: $O(nC_m^k)$. Filling each entry takes $k(m-k) + 1$ time.

A greedy solution: Furthest-Future principle (L. Belady, '66) I

FF rule: evicts the farthest-future element;
Furthest-Future eviction sequence S_{FF} : see a fig.

Theorem

S_{FF} incurs no more missing than any other schedule S^ and hence is optimal.*

Proof:

- Exchange argument again!
- Basic idea: From an optimal schedule S^* , we generate a series of schedule S_1, S_2, \dots, S_n , such that:
 - 1 The first i evictions of S_i are the same to that of S_{FF} . Thus, $S_n = S_{FF}$.
 - 2 S_{i+1} incurs no more missing than S_i .

A greedy solution: Furthest-Future principle (L. Belady, '66) II

$S_0=S^*$	—	—	—	7	5	3	0	5	7
S_1	—	—	—	3	5	7	0	5	7
S_2	—	—	—	3	0	7	0	5	7
.....									
S_{FF}	—	—	—	3	0	9	—	—	7

Basic idea: interpolating a sequence of S_i to change S^* to S_{FF} .

Two requirements:

1. S_i : the first i evictions of S_i are same to that of S_{FF} .
2. S_{i+1} incurs no more evictions than S_i .

- Difficulty: how to construct S_{i+1} based on S_i ?

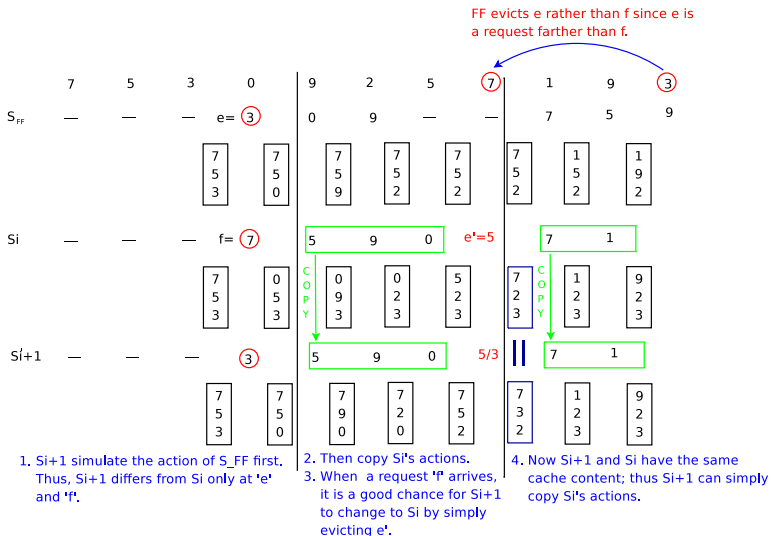
- Consider the $j + 1$ -th request d . S_i and S_{FF} have the same cache content till now.

A greedy solution: Furthest-Future principle (L. Belady, '66) III

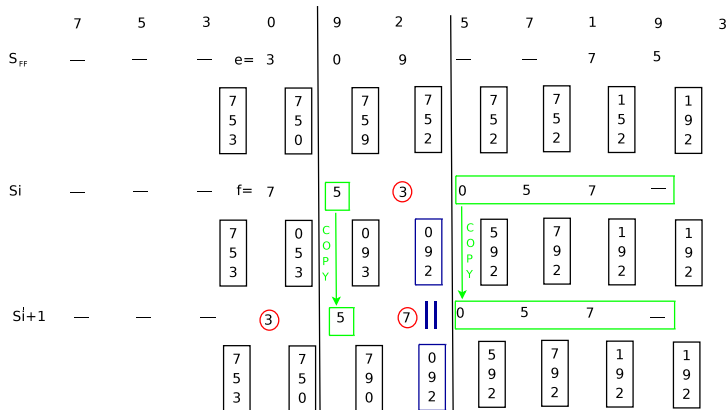
- 2 Suppose S_i evicts f but S_{FF} evicts $e \neq f$. We design S_{i+1} as follows:
- 3 S_{i+1} evicts e at the $j + 1$ -th step. Now, S_{i+1} and S_i has different cache content.
- 4 S_{i+1} simulates the actions of S_i from $j + 2$ step until the following two events occurs:
- 5 Case 1: request $g \neq e, f$, and S_i evict e .
We let S_{i+1} evict f . The cache are same now. Thus, we can copy the remaining of S_i to S_{i+1} .
- 6 Case 2: request f and S_i evicts e' .
We let S_{i+1} evict e' , too. And fill e if needed.

Key: the Furthest-Future principle ensures that before an request of e , there should be a request of f (Case 1).

Case 1:



Case 2:



- S_{i+1} simulate the action of S_{FF} first. Thus, S_{i+1} differs from S_i only at 'e' and 'f'.
- Then copy S_i 's actions.
- When S_i evicts 'e', S_{i+1} evicts 'f' to reduce the difference.
- Now S_{i+1} and S_i have the same cache content; thus S_{i+1} can simply copy S_i 's actions.

From FF to LRU

- LRU: least recently used.
- Intuition: “longest in the past rather than the farthest in the future” since we have no idea of the future requests.
- Reason: locality of reference, i.e., a program will generally keep accessing the things it has just been accessing.

Theoretical analysis of LRU principle (Sleator and Tarjan)

Key idea: divide the requests into “phases”. Each phase consists of a set of evictions.

A Label-based algorithm framework:

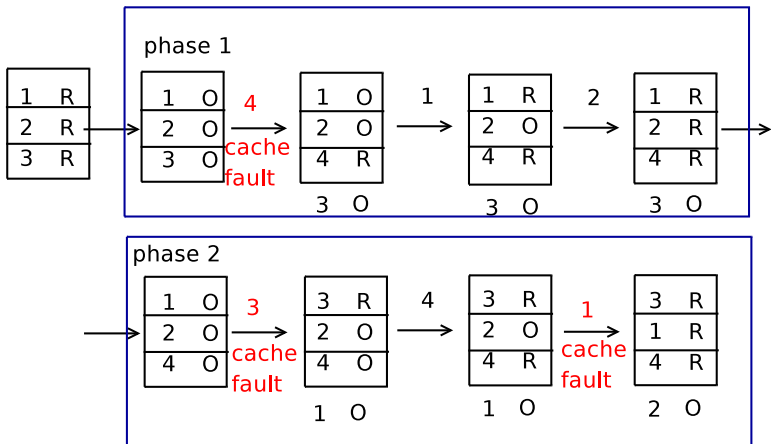
- 1 initially make all pages in cache as “old”;
- 2 when a request of page P arrives,
- 3 mark P “recent”;
- 4 if P is not in cache
- 5 if all cache pages are “recent”,
- 6 remark all pages “old”, and begin a new phase;
- 7 choose an old page q to evict;

Old: the pages have already been loaded before this phase.

Recent: the pages are loaded in this phase.

Theoretical analysis of LRU principle (Sleator and Tarjan)

II



Theoretical analysis of LRU principle (Sleator and Tarjan)

III

Analysis:

Suppose there are r phases.

- Fact 1: Every phase contains k distinct requests. (Reason: when a page changes from “Old” to “Recent”, it will stay in cache till the phase ends.)
- Fact 2: At each phase, there are at most k evictions. Thus, there are at most rk evictions. (Intuition: evicting a page cause a page remarking from “Old” to “Recent” .)
- Fact 3: An optimal solution incurs at least $r - 1$ missing. (Reason: the first request of a phase i cause a remarking of a page from “Old” to “Recent” .)
- Therefore, the ratio of any Label-based algorithm is k .

Worst-case: repeating a cycle of requests $1, 2, \dots, k + 1$ when cache size is k .

Note: LRU is a Label-based method.

A randomize algo I

Algo: choose a “random” old page to evict.

Theorem

Let denote the minimal eviction number as F^* . The expected number of evictions of RandomEviction algorithm is at most $2 \ln k F^*$.

Proof:

Consider phase i .

- Let A be the cache content at beginning. Sort A according to the request order in this phase, say $A = \{a_1, a_2, \dots, a_k\}$
- Let b_i be the requests that are not in A .
- When a_j is requested, and a_j is marked “Old”; (Reason: the case that a_j is “Recent” is omitted since it causes no cache fault.)
- $\#OldPages = k - (j - 1)$; (Reason: a_1, a_2, \dots, a_{j-1} are marked as “Recent”.)

A randomize algo II

- $\#OldPagesInCache = k - (j - 1) - |b_i|$; (Reason: a request in b_i evicts an “Old” page out of cache.)

Pages in cache are in random. Thus, we have:

- $\Pr(a_j \text{ is in cache}) \geq \frac{k - (j - 1) - |b_i|}{k - (j - 1)}$;
- $\Pr(a_j \text{ is NOT in cache}) \leq 1 - \frac{k - (j - 1) - |b_i|}{k - (j - 1)}$ (cache missing);
-

$$E(\#missing) \leq b_i + \sum_{j=1}^k \frac{|b_i|}{k - (j - 1)} \quad (1)$$

$$= |b_i| \log(k - (j - 1)) \quad (2)$$

$$= O(\log(k)) \quad (3)$$

In phase i and $i + 1$, there are $k + b_i$ distinct pages requested. Thus we can bound the number of faults as follows:

- $\#missing \geq b_i$

A randomize algo III

- $\#total - missing \geq \frac{1}{2} \sum_i |b_i|$
- ratio: $\leq 2 \log(k)$.