

Online algs for caching (the paging problem)

C blocks of cache
 a_1, a_2, \dots, a_n access sequence of block numbers
 cache control alg. \rightarrow cost $T := \#$ cache misses
 \downarrow
 offline \rightarrow OPT $:=$ cost of optimum offline alg.
 online

we want: find online alg. s.t. $T \leq k \cdot \text{OPT}$
 a constant \uparrow } alg. is k -competitive

Unfortunately not possible

\rightarrow we can force $T \geq \sim C \cdot \text{OPT}$.

Cheat: Give LRU a larger cache than OPT has $\rightarrow C_{LRU} > C_{OPT}$.

Theorem (Sleator & Tarjan; weakened):

For every $C_{LRU} > C_{OPT} \geq 1$ and every access sequence:

$$T_{LRU} \leq \frac{C_{LRU}}{C_{LRU} - C_{OPT}} \cdot T_{OPT} \rightarrow C_{OPT}$$

for $C_{LRU} = 2C_{OPT}$
 $T_{LRU} \leq 2 \cdot T_{OPT} + C_{OPT}$

Proof: Split access seq. to epochs:

Consider epoch E_i for $i > 0$:

Case 1: All LRU's misses on different blocks

\rightarrow accesses to C_{LRU} diff. blocks

\rightarrow opt. alg. misses at least $C_{LRU} - C_{OPT}$ times.



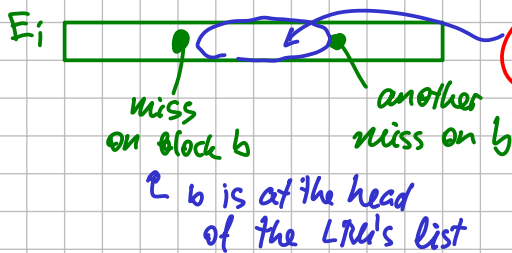
For E_i : # misses of LRU = C_{LRU} .

Except E_0 : # misses of LRU $\leq C_{LRU}$

in this epoch:

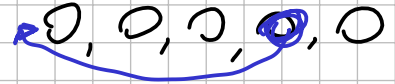
$$\frac{T_{LRU}}{T_{OPT}} \leq \frac{C_{LRU}}{C_{LRU} - C_{OPT}} \oplus$$

Case 2: LRU misses twice on the same block.



accesses to at least C_{LRU} different blocks

LRU keeps a list of cached blocks sorted by age



If we have $\frac{T_{LRU}}{T_{OPT}} \leq \oplus$ in every E_i , it also holds over the whole seq. (except for E_0).

Return to E_0 :

Case 1: The cache is empty at the start.

we have at most C_{LRU} diff. blocks accessed.

$$\Rightarrow T_{OPT} \geq T_{LRU} \quad \text{ratio } \frac{T_{LRU}}{T_{OPT}} \leq 1 \leq \oplus$$

Case 2: OPT starts with a non-empty cache...

it can save up to C_{OPT} misses.

Conclusion: For algs we studied: #reads $(N, B, M) \in O(\#reads(N, B, M/2))$

So: #reads with LRU, cache $M/2 \in O(\#reads with LRU for cache size M) \in O(\#reads with OPT, cache size M/2)$

LRU strategy

\uparrow least-recently used

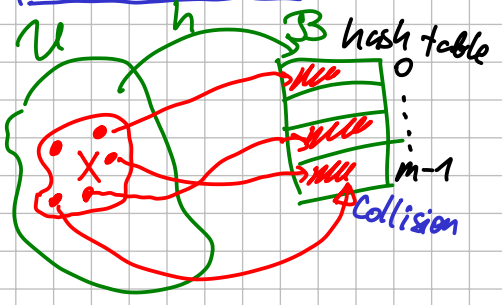
We discard the block not accessed for the longest time.

So the assumption of opt. strategy affects only constants.

HASHING

Notation: U universe, $U := |U|$
 in many cases:
 $U = \{0 - U-1\}$

Complexity: $\left\{ \begin{array}{l} \text{worst-case} \\ \text{amortized} \\ \text{average-case} \end{array} \right.$
 cost of hashing \uparrow
 \rightarrow no order-based operations (Min, Succ, ...)



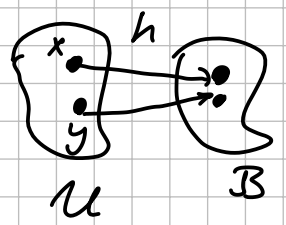
B set of buckets = $[m]$
 h hash function $h: U \rightarrow B$ evaluated in const. time
 $X \subseteq U$ a finite set stored in the D.S.
 $n := |X|$

In every bucket, we have a linked list of items (a chain) \rightarrow hashing with chaining

\rightarrow time per op. is $O(\text{chain length})$
 \hookrightarrow worst-case $O(n)$
 \rightarrow we hope for short chains on average $\frac{n}{m}$
 We choose $m \sim n$

Def: A family \mathcal{H} of functions from U to B is c-universal for $c > 0$ iff
 $\forall x, y \in U, x \neq y: \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq c/m$.

We usually want:
 ① $h \in \mathcal{H}$ can be picked randomly in $O(1)$ time
 ② $h(x)$ can be evaluated in $O(1)$ time



\hookrightarrow solution: $\mathcal{H} := \{h_a \mid a \in A\}$
 pick $a \in A$ parameter
 it can happen that $h_a = h_{a'}$
 \rightarrow consider \mathcal{H} a multi-set.

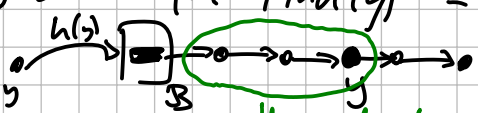
Def: \mathcal{H} is universal iff $\exists c > 0: \mathcal{H}$ is c-universal.

Lemma: Given $x_1 \dots x_n, y \in U$ all distinct and \mathcal{H} is c-universal:
 $\mathbb{E}_{h \in \mathcal{H}} [\sum_{i=1}^n \mathbb{1}_{h(x_i) = h(y)}] \leq \frac{cn}{m}$.

Proof: Use indicators $A_1 \dots A_n$:
 $A_i = \begin{cases} 0 & \\ 1 & h(x_i) = h(y) \end{cases}$
 $\mathbb{E}[A_i] = 0 \cdot \Pr[A_i=0] + 1 \cdot \Pr[A_i=1] = \Pr[h(x_i) = h(y)] \leq c/m$
 $\mathbb{E}[A_i] \leq c/m$
 $A = \sum_i A_i \xrightarrow{\text{linearity of } \mathbb{E}} \mathbb{E}[A] = \sum_i \mathbb{E}[A_i] \leq \frac{cn}{m}$
 Always: $\mathbb{E}[X+Y] = \mathbb{E}[X] + \mathbb{E}[Y]$

Use the lemma to analyze expected complexity of hashing with chaining with hash func. chosen u.a.r. from c-universal family.

- ① unsuccessful Find(y) \rightarrow by lemma $\mathbb{E}[\text{length of chain in } h(y)] \in O(n/m)$
- ② successful Insert(y) \rightarrow
- ③ successful Find(y) \leq succ. Insert



- ④ Unsucc. Ins. \rightarrow Succ. Find (these already were there at the time of Insert)
- ⑤ Delete \rightarrow Find

if we maintain $m \in \Omega(n)$, this is $O(1)$.
 use flexible arrays & rehash