

Optimal Caching and Greedy Exchange

In the optimal caching problem, our computer has a main memory of size n , a cache of size k , and we are presented with a sequence of data $D = d_1, d_2, \dots, d_m$ that we must process. When an item is not in the cache, we have a cache miss, and must bring the item into the cache and evict something else if the cache is full. Our goal is to give an algorithm that minimizes the number of misses.

Example 1: a, b, c, b, c, a, b

$k = 2$ cache = $\{a, b\}$

Example 2: $a, b, c, d, a, d, e, a, d, b, c$

$k = 3$ cache = $\{a, b, c\}$

Definition 1. A schedule is *reduced* if it does the minimal amount of work necessary in a given step.

Lemma 1. *For every non-reduced schedule, there is an equally good reduced schedule (that brings in at most as many items as the original schedule).*

Prove this by construction.

Hint: You might *charge* a miss from one schedule to a miss in another schedule to show that it doesn't have any extra misses.

Proof. For any schedule S , define its reduced schedule \bar{S} : whenever S brings in some d that wasn't requested, \bar{S} leaves d in main memory and only brings it in the next time that it's requested. Then the *miss* of \bar{S} in that next step j when d is requested can be charged to step i in S . \square

Observation 1. *For any reduced schedule, the number of items that are brought in is exactly the number of misses.*

A Greedy Algorithm for Optimal Caching

Determine a cache maintenance algorithm by coming up with an eviction schedule.

A cache maintenance algorithm with an optimal greedy eviction schedule is the *Farthest-in-Future* Algorithm. When a new item needs to be brought into the cache, it greedily evicts the item that is needed the farthest into the future.

Goal: Prove that the Farthest-in-Future Algorithm is optimal. We'll call this schedule S_{FF} .

Proof Brainstorming

Given S_O as the reduced schedule from some unspecified (possibly optimal) algorithm, and S_A as an algorithm that we want to prove optimal, how can we prove $cost(S_A) \leq cost(S_O)$?

Here are a couple steps to build out a strategy that we will call Greedy Exchange.

a. Transform S_O into S_A in a way that provably does not increase cost. How?

b. Make a sequence of transformations (exchanges).

$$S_O = S'_0 \rightarrow S'_1 \rightarrow S'_2 \rightarrow \cdots \rightarrow S'_k = S_A$$

such that $cost(S'_i) \geq cost(S'_{i+1})$. That implies $cost(S'_0) \geq cost(S'_k)$ and $cost(S_O) \geq cost(S_A)$.

Now, we just need to design the transformation (sequence of S'_i s) so that it starts with $S'_0 = S_O$, ends with $S'_k = S_A$, and we can prove $cost(S'_i) \geq cost(S'_{i+1})$. If we can do this for any O , then A is optimal. However, this will not be possible for all A .

c. Same as previous, but add a constraint that S'_i shares the first i eviction decisions with A .

This new constraint forces $S'_k = S_A$ and constrains our transformation choices.

To go much further, we will need to reason about what S_A looks like, so we will now narrow our focus to $S_A = S_{FF}$ and construct the sequence of S'_i s.

Proving the Greedy Exchange

Let S_{FF} be the schedule created by the farthest-in-future algorithm and let S be an arbitrary reduced schedule.

Lemma 2. *Suppose S is a reduced schedule that makes the same eviction decisions as S_{FF} through the first j items in the sequence for some j . Then there exists a reduced schedule S' that makes the same eviction decisions as S_{FF} through the first $j + 1$ items and incurs no more misses in total than S does in total.*

Prove this by constructing S' . This is an *exchange* argument.

- a. What happens if the $j + 1^{st}$ item is in cache?

Neither evicts (recall that S is reduced). Let $S' = S$.

- b. What happens if the $j + 1^{st}$ item isn't in cache, but S evicts the same item as S_{FF} ?

Let $S' = S$.

- c. What happens if the $j + 1^{st}$ item isn't in cache, and S evicts a different item as S_{FF} ? What should S' do?

Suppose S evicts some f while S_{FF} evicts some e . Let S' match S_{FF} up to $j + 1$, and now S' and S 's caches differ by one. Suppose S_{FF} evicts some item e and S evicts some item f .

- d. How can you get S' 's cache back to the same as S 's without incurring more total misses?

From now on, S' does the same as S unless

- there's a request to $g \neq e$, f not in cache and S would evict e . Then S' evicts f .
- there's a request to f and S evicts e' . If $e' \neq e$, then S' evicts e' and brings in e .

How do we know there's not a request to e ? Well, e was evicted because it's seen farthest in the future, so certainly f or $g \neq e$, f is requested before e .

Now the caches and misses are equal. We need to transform S' to reduced.

- e. How do we know that S' is a reduced schedule?

It may not be, but by Lemma 1, we can transform it to be. The schedule construction used to prove Lemma 1 only adjusts the schedule after the first non-reduced step, so it won't affect the initial $j + 1$ steps. (Also S_{FF} will always be reduced so the reduction never conflicts there either.)

- f. Sanity check: Are all parts of the lemma true?

Matches S_{FF} through step $j + 1$, same caches and misses as S , is reduced. Yep!

By induction, for any m and algorithm S , there's a reduced schedule that incurs no more misses but makes the same evictions as S_{FF} .

Theorem 2. S_{FF} incurs no more misses than any other schedule S^* and hence is optimal.

Recap: Proof by Greedy Exchange

Step 1: Label. Label your algorithm's solution ($A = \{a_1, a_2, \dots, a_k\}$), and a general solution ($O = \{o_1, o_2, \dots, o_m\}$).

Step 2: Compare. Compare greedy with the other solution. Assume that they're not the same and isolate some difference.

Suppose they are the same through the first j items. Then we show in the following lemma, by an exchange argument, that we can modify them to be the same through the first $j + 1$ items.

Step 3: Exchange. Swap the elements in in O without making the solution worse. Argue that swapping a finite number of times will result in A . Hence, greedy is just as good as *any* optimal or arbitrary solution.

After k exchanges, O has been transformed into A without increasing the cost.