# Near-Optimal Bounds for Online Caching with Machine Learned Advice

Dhruv Rohatgi
MIT
drohatgi@mit.edu

October 31, 2019

## Abstract

In the model of online caching with machine learned advice, introduced by Lykouris and Vassilvitskii, the goal is to solve the caching problem with an online algorithm that has access to next-arrival predictions: when each input element arrives, the algorithm is given a prediction of the next time when the element will reappear. The traditional model for online caching suffers from an $\Omega(\log k)$ competitive ratio lower bound (on a cache of size $k$). In contrast, the augmented model admits algorithms which beat this lower bound when the predictions have low error, and asymptotically match the lower bound when the predictions have high error, even if the algorithms are oblivious to the prediction error. In particular, Lykouris and Vassilvitskii showed that there is a prediction-augmented caching algorithm with a competitive ratio of $O(1 + \min(\sqrt{\eta/\text{OPT}}, \log k))$ when the overall $\ell_1$ prediction error is bounded by $\eta$, and OPT is the cost of the optimal offline algorithm.

The dependence on $k$ in the competitive ratio is optimal, but the dependence on $\eta/\text{OPT}$ may be far from optimal. In this work, we make progress towards closing this gap. Our contributions are twofold. First, we provide an improved algorithm with a competitive ratio of $O(1 + \min((\eta/\text{OPT})/k, 1) \log k)$. Second, we provide a lower bound of $\Omega(\log \min((\eta/\text{OPT})/(k \log k), k))$.

## 1 Introduction

In the *online caching* problem (also known as *paging*), we are given a sequence of elements which arrive one at a time, and we must maintain a cache of some fixed size $k$. The cost of a caching algorithm on some input is the number of cache misses. The standard goal is to design an online algorithm with minimal competitive ratio, relative to the optimal offline algorithm.

As a fundamental problem in the study of online algorithms, caching has been extensively studied [14, 9, 3, 1]; it is well-known that the optimal competitive ratio of any deterministic algorithm is $k$, and the optimal competitive ratio of any randomized algorithm is $2H(k) - 1 = \Theta(\log k)$, where $H(k)$ is the $k$-th harmonic number [1].

However, the traditional framework for analyzing online algorithms—namely, worst-case competitive ratios—is overly pessimistic, by virtue of requiring worst-case analyses. Real-world data often satisfies nice properties—it may be predictable, or simply random, or even just not adversarial—and for this reason, theoretically unsound algorithms can perform very well in practice. Numerous attempts have been made to theoretically ground this observation; some of the more prominent are average-case analyses [6, 2, 10] and smoothed analyses [15, 16], both within online algorithms and beyond.

One such attempt which has recently garnered significant attention is the framework of online algorithms with machine learned advice. In this model, the online algorithm is augmented with an oracle that makes certain predictions about future data. In practice, this oracle is likely to be a machine learned predictor. Since machine learning is imperfect (and can sometimes be wildly wrong), the algorithm must incorporate the oracle's advice judiciously, without being given any bound on the oracle's error. The goal is to develop an algorithm which is both *consistent*, in that it nearly matches the best offline algorithm when the predictor is nearly perfect, as well as *robust*, in that its worst-case performance is good even when the oracle is arbitrarily bad. The performance of the algorithm should be bounded as a function of some measure of the oracle error, even though the algorithm is oblivious to this error.

The ML advice model has in the past been applied to the ski rental problem [13, 4], job scheduling [13, 12] and online revenue maximization [11]; it has also been used to achieve theoretical and practical gains in streaming frequency estimation [5] and data structures [7]. Most relevant to this paper is prior work by [8] in which it was shown how the model can be applied to the online caching problem. In particular, they considered augmenting caching algorithms with an oracle that predicts the next arrival of each element. The oracle's $\ell_1$ error is defined as the sum over all elements of the absolute difference between the element's true and predicted next arrival. In this model, they developed a "predictive marker algorithm" with the following guarantee:

**Theorem 1.** *[8] The predictive marker algorithm achieves a competitive ratio of* $\min(2+4\sqrt{\eta/\text{OPT}}, 4H(k))$ *when the oracle has $\ell_1$ error of at most $\eta$, and the cost of the optimal offline algorithm is* OPT.

Note that the competitive ratio achieved is $O(\min(\sqrt{\eta/\text{OPT}}, \log k))$. This ratio is of course (asymptotically) optimal as a function of $k$. However, it is an open question how far the dependence on $\eta/\text{OPT}$ (a measure of relative error, in some sense) can be improved. In particular, it would be interesting to understand how accurate the predictions need to be in order to provide an improvement in the competitive ratio. From the previous work, it is only shown that $\eta/\text{OPT} = o(\log^2 k)$ suffices.

In this paper, we work with the same model, and make progress on this question. Building upon the techniques used in [8], we provide an algorithm with an improved competitive ratio:

**Theorem 2.** *There is an algorithm for caching with predictions that achieves a competitive ratio of*

$$O\left(1 + \min\left(1, \frac{\eta/\text{OPT}}{k}\right)\log k\right)$$

*when the oracle has $\ell_1$ error of at most $\eta$.*

Our bound matches the prior work when $\eta/\text{OPT} \geqslant k$ (in that neither algorithm improves upon the classical $O(\log k)$ competitive ratio) and is strictly better when $\eta/\text{OPT} < k$. For example, if $\eta/\text{OPT} = k/\log k$, then the prior algorithm had a competitive ratio of $O(\log k)$, whereas we show that $O(1)$ is possible.

Furthermore, we provide a lower bound, stated informally as follows:

**Theorem 3.** *Any randomized algorithm for caching with predictions must have a competitive ratio which is*

$$\Omega\left(\log\min\left(\frac{\eta/\text{OPT}}{k\log k}, k\right)\right)$$

*as a function of $\eta/\text{OPT}$ and $k$.*

To our knowledge, this is the only known lower bound. The upper bound and lower bound are asymptotically tight when $\eta/\text{OPT} \leqslant k/\log k$ or $\eta/\text{OPT} \geqslant k^{1+\epsilon}$. There is a still a significant gap; the two bounds are non-trivial on disjoint regimes, and in the regime $k \leqslant \eta/\text{OPT} \leqslant k \log k$, neither bound is non-trivial.

Nonetheless, these results make progress towards determining the largest possible error bound that still admits a non-trivial competitive ratio. Where prior work only showed that $\eta/\text{OPT} = o(\log^2 k)$ suffices to obtain a competitive ratio of $o(\log k)$, the above results imply that $\eta/\text{OPT} = o(k)$ suffices and $\eta/\text{OPT} \leqslant k^{1+o(1)}$ is necessary.

## 1.1 Roadmap

In Section 2, we formally describe the online caching model with machine learned advice, and define the predictor. We also review some facts from traditional caching algorithms (specifically, facts about marker-based algorithms) that we will rely on later in the paper.

Having defined the necessary terminology, we then outline in Section 3 the techniques used to achieve our results.

In Section 4, we provide the first of these results, a warm-up algorithm for caching with predictions. This algorithm improves upon the prior work, and is simpler than our final algorithm (and thus may have some practical advantages). Furthermore, as a marker-based algorithm, this first algorithm is somewhat simpler to analyze.

In Section 5, we describe our second and final algorithm, departing from the marker-based framework to achieve an improvement in the competitive ratio.

Finally, in Section 6, we prove the lower bound.

# 2 Preliminaries

## 2.1 Traditional caching

In the traditional online caching problem, the input is a sequence $\sigma = (z_1, z_2, \ldots, z_n)$ of elements which become available one by one. The cache has fixed size $k$ and is initially empty. As elements arrive, if an element is not present in the cache, then it counts as a "cache miss", and the algorithm must add it to the cache, and choose which cache element to evict. Otherwise nothing happens. The cost $\text{cost}_A(\sigma)$ of the algorithm $A$ on the input $\sigma$ is the number of cache misses. If the algorithm is randomized, this cost is the expected number of cache misses.

We define $\text{OPT}(\sigma)$ to be the minimum number of caches misses achievable by an "offline" algorithm—an algorithm which is given $\sigma$ in advance. Our online algorithm $A$ is $\alpha$-competitive if there is some constant $c$ such that for every input $\sigma$,

$$\text{cost}_A(\sigma) \leqslant \alpha \cdot \text{OPT}(\sigma) + c.$$

It is known that there is a $k$-competitive deterministic algorithm and an $O(\log k)$-competitive randomized algorithm. Furthermore, these ratios are optimal.

## 2.2 ML advice

In this paper we consider not the traditional model but rather an extension of it, in which our algorithm is also given some advice [8]. In particular, when input element $z_i$ arrives, an oracle gives the algorithm $h_i(z_i)$, which is an estimate of $y_i = \min_{j>i}\{j : z_j = z_i\}$, the next time when element $z_i$ will appear (if $z_i$ never appears again, and the input sequence has length $n$, then we set

3

$y_i = n + 1$). These estimates may not be correct, and we want to bound the performance of our algorithm as a function of the error. We define the error as the $\ell_1$ distance between the real and predicted next arrivals: for each input element $z_i$ we define $\mathrm{Err}_i(z_i) = |h_i(z_i) - y_i|$, and then define

$$\eta = \sum_{i=1}^{n} \mathrm{Err}_i(z_i).$$

For any time $i$ and input element $w$, we also define $L(w, i)$ to be the last time $j < i$ such that $z_j = w$.

When analyzing an algorithm $A$ in this model, the goal is to bound its competitive ratio as a function of $\eta/\mathrm{OPT}$: more precisely, to show for a desired function $\alpha$ and a constant $c$, that $\mathrm{cost}_A(\sigma) \leqslant \alpha(\eta/\mathrm{OPT}) \cdot \mathrm{OPT}(\sigma) + c$ for every input $\sigma$. This is the approach taken in prior work (see Theorem 1), and we will see how $\eta/\mathrm{OPT}$ arises naturally in our algorithms' analyses. It would not make sense for the competitive ratio to be a function of the absolute error $\eta$, since duplicating the input sequence would double $\eta$ but leave the ratio $(\mathrm{cost}_A(\sigma)/\mathrm{OPT}(\sigma))$ approximately unchanged.

## 2.3 Marker-based algorithms

Our first predictive caching algorithm will be a *marker-based algorithm* which judiciously incorporates the oracle's advice; our second algorithm will depart from but still rely heavily on the marker-based framework. Marker-based caching algorithms have the following structure. The execution of the algorithm comes in phases, and at the beginning of each phase, every cache element is said to be *unmarked*. When a cache hit occurs, the corresponding element is *marked*. When a cache miss occurs, some unmarked element is evicted from the cache, and the new element is inserted in its place, and immediately marked. If all elements of the cache are marked and another cache miss occurs, then the whole cache is unmarked and a new phase begins.

The decision that a marker-based algorithm has to make is which unmarked element to evict in the event of a cache miss. In the traditional online model, the randomized marker algorithm achieves an $O(\log k)$-competitive ratio by evicting a random unmarked element. Additionally, any marker-based algorithm is $k$-competitive.

For any marker-based algorithm, we can make the following definition.

**Definition 4.** *An input element is called* clean *for some phase $r$ of the algorithm execution if it appeared in phase $r$ but did not appear in phase $r - 1$. If it appeared in both phase $r$ and phase $r - 1$, it is called* stale.

The following lemma is known, relating the number of clean elements to the optimal offline cost.

**Lemma 5.** *[3] Let $L$ be the number of clean elements in an execution of a marker-based algorithm on some input $\sigma$. Then $L/2 \leqslant \mathrm{OPT}(\sigma) \leqslant L$.*

The phases (and consequently, the clean/stale elements) are in fact independent of the exact algorithm.

**Definition 6.** *An* arrival *in phase $r$ is an element $z_i$ which has not previously appeared in the same phase.*

Then the following fact can be readily derived for any marker-based algorithm:

**Claim 7.** *Every phase contains exactly $k$ arrivals.*

More specifically, each phase continues as long as possible without containing $k + 1$ distinct elements.

## 2.4 Eviction chains

To design and analyze marker-based algorithms, it is useful to decompose the set of cache misses into *eviction chains*, a concept perhaps first explicitly utilized in [8].

For any marker-based algorithm and any phase, each clean arrival in the phase causes a cache miss. This yields a chain of evictions in that phase which can be blamed on that clean arrival: the clean element's arrival evicts some element, whose next appearance evicts another element, and so forth until an element is evicted which never reappears in the phase. Each element in the chain must be an arrival, since elements which have previously appeared were marked and are therefore immune to eviction for the remainder of the phase. Thus, each element in the chain after the first clean element must be stale, since to be evicted it must have been present in the cache.

These clean-element chains account for all cache misses in the phase, since every stale element was in the cache at the start of the phase, so for it to cause a cache miss it must have been evicted by a previous element. So the total number of cache misses in a phase is simply the total length of the eviction chains.

## 3 Our techniques and related work

Given next-arrival predictions, an algorithmically naive approach is to trust the predictions completely. The optimal offline algorithm evicts, at each cache miss, the cache element with the latest next-arrival time. Thus, if the predictions are perfect then this approach will have a competitive ratio of 1. However, even small errors in the predictions can lead to an unbounded prediction. So it is necessary to balance trusting the predictions with making provably competitive decisions.

Our work builds on [8], which proposed a marker-based algorithm for caching with predictions. Their algorithm utilizes eviction chains. Since eviction chains partition the set of cache misses, and the number of eviction chains is equal to the number of clean arrivals, which is asymptotically equal to OPT, bounding the average chain length bounds the competitive ratio of the algorithm. Thus, algorithms which work with each chain independently can often be cleanly analyzed. In [8], this approach is carried out: for each eviction chain, the predictions are trusted (to choose which element to evict next) until the chain reaches length $\Omega(\log k)$, after which evictions are random. Each chain's length can be bounded by the prediction error of elements in that chain. Since the chains are disjoint, this implies a bound on the total cost of the algorithm by the total prediction error (see Theorem 1).

As a warm-up, we first show that a small modification to the algorithm from [8] improves the competitive ratio from $O(1 + \min(\sqrt{\eta/\text{OPT}}, \log k))$ to $O(1 + \min(\log(\eta/\text{OPT}), \log k))$. The modification is simple—trust the predictions only *once* in each chain—but the analysis requires more care. Unlike before, the length of each chain now depends on the next-arrival prediction errors of elements which may or may not appear in the chain. Thus, adding up the errors could double-count. To avoid this issue, we do not directly bound each chain's length by prediction error. Instead, we charge the length of each chain against a set of inversions in the order of element arrivals relative to the predictions. These sets of inversions are disjoint for different chains, so summing across chains does not double count. A combinatorial lemma then relates the total number of inversions to the total error.

To improve the competitive ratio further and achieve the bound stated in Theorem 2, we depart from the marker-based framework. In our previous algorithm, trusting the prediction once meant evicting the unmarked element $e$ with latest predicted next arrival. Intuitively, if the prediction error "per chain" is $O(\eta/\text{OPT})$, then if $e$ reappears in the same phase, it should on average be at most $O(\eta/\text{OPT})$ elements away from the end of the phase. Our first algorithm would then proceed

by random eviction of unmarked elements, so $e$'s eviction chain would have length $O(\log(\eta/\text{OPT}))$ in expectation (for the same reason that the traditional marker algorithm has competitive ratio $O(\log k)$).

However, if $\eta/\text{OPT} \ll k$, then a uniformly random cache element (potentially marked) will probably not appear after $e$, so evicting it might terminate the eviction chain at length $O(1)$ instead of $O(\log(\eta/\text{OPT}))$. This is the motivation for our final algorithm. However, it relies on evicting marked cache elements, which significantly complicates the analysis, since facts about marker-based algorithms no longer directly apply.

In particular, the number of eviction chains may no longer be $\Theta(\text{OPT})$; evicting a marked element may cause an "extra" eviction chain in the next phase. To deal with this issue, our key tool is a bijection from eviction chains to special elements that began the phase in the cache but never appeared. With this bijection, we show that every extra chain can be charged against either prediction error or a chain which *does not* exist but could have. An added complication is showing that these prediction errors are disjoint.

For the lower bound, the idea is to choose an input distribution and predictions such that the predictions give no information about the future input, but have reasonably small error. Generalizing the traditional lower bound against randomized caching algorithms—where each input element is uniformly distributed over $k+1$ pages, so each phase has one clean element—our input distribution consists of phases each with $t$ clean elements, where $t$ is a variable parameter of the distribution. Increasing the parameter $t$ allows smaller prediction error, at the cost of a smaller bound on the competitive ratio.

## 4 Marker-based predictive algorithm

Our first algorithm LMARKER is a modification of the algorithm proposed in [8], which is itself a balance between two paradigms: trust the oracle, or ignore the oracle. A marker-based algorithm which trusts the oracle would always evict the cache element with highest predicted next arrival; if the oracle had zero error, this algorithm would match the optimal offline algorithm. A marker-based algorithm which ignores the oracle is the random marker algorithm. To combine the gains of the former algorithm when the oracle has low error with the robustness of the latter algorithm when the oracle has high error, the algorithm from [8] uses the following strategy: for each eviction chain, trust the oracle until the chain becomes long, and subsequently ignore the oracle.

The modification we make is simple: trust the oracle less—only once at the beginning of each eviction chain. Intuitively, a total prediction error of $\eta$ translates to an error of $\Theta(\eta/\text{OPT})$ in each of the $\Theta(\text{OPT})$ chains. Evicting the unmarked element with highest predicted next arrival means that its true next arrival should be only $O(\eta/\text{OPT})$ from the end of the phase, resulting in $O(\log \eta/\text{OPT})$ more cache misses for that eviction chain.

We now describe LMARKER in more detail, and formalize the analysis. The chains cannot quite be analyzed independently (unlike in [8]), so care is needed.

**Algorithm description**   LMARKER is a marker-based algorithm with the following eviction strategy upon a cache miss: if the incoming element is clean, then evict the unmarked element with the highest predicted arrival time. If the incoming element is stale, then evict a random unmarked element.

**Algorithm analysis**   Fix a phase. Exactly $k$ distinct elements arrive during the phase; let $i_1, \ldots, i_k$ be the times of the first arrivals for the phase. These are the only times when cache

misses may occur, and the only times when an unmarked element becomes marked.

Consider a single clean element with arrival time $i_t$. It evicts the unmarked element with the highest predicted arrival time, which must be either (1) a stale element with arrival time in the set $\{i_{t+1}, \ldots, i_k\}$ or (2) an element which does not appear in this phase. Case (2) results in no more cache misses along this chain, so we analyze case (1): some stale element is evicted. Let $i_{e(t)}$ be the time at which it arrives.

For any $1 \leqslant a \leqslant b \leqslant k$, define $N_a(b)$ to be the number of stale elements which are unmarked and in the cache at time $i_a$, and have arrival times after $i_b$. Let $\mathcal{E}_t$ be the distribution over executions of the algorithm up to time $t$.

**Claim 8.** *If the clean element arriving at time $i_t$ evicts the stale element arriving at time $i_{e(t)}$, then the expected length of the eviction chain begun at $i_t$ is at most $O(\mathbb{E}[1 + \log(N_t(e(t)))])$.*

*Proof.* Conditioning on $\mathcal{E}_t$, the random variable $N_t(e(t))$ is determined. In the worst case, each stale element evicts another stale element until this is no longer possible. There are at most $N_t(e(t))$ unmarked stale cache elements at time $i_{e(t)}$. Say that there are $m$ such elements, and order them $1, \ldots, m$ by arrival time. If $z_{i_{e(t)}}$ evicts the $j$-th such element then there will be at most $m - j$ unmarked cache elements when that element arrives. But $j$ is uniformly distributed over $\{1, 2, \ldots, m\}$. Thus, the expected length of the chain, conditioned on $\mathcal{E}_t$, is bounded by $R_{N_t(e(t))}$ defined by the recurrence

$$R_m = 1 + \frac{1}{m} \sum_{j=1}^{m} R_{m-j},$$

which solves to $R_m = O(\log m)$. By the law of total expectation, the unconditional expected length of the chain is bounded by $O(\mathbb{E}[1 + \log(N_t(e(t)))])$. $\qquad\square$

Next, fix any execution of the algorithm on the entire input. We relate $N_t(e(t))$ to prediction error; more specifically, we relate it to the number of inversions in the predicted arrival order of stale elements in phase $r$. For each stale arrival time $i_s$, let $j_s$ be the most recent appearance of $z_{i_s}$ in phase $r - 1$. Let $J = \{j_s | z_{i_s} \text{ is stale in phase } r\}$ be the set of most recent appearances.

If $z_{i_{e(t)}}$ is evicted by $z_{i_t}$, then at time $i_t$, all unmarked stale elements $w$ in the cache with arrivals in the set $\{i_{e(t)+1}, \ldots, i_k\}$ had earlier predicted next arrivals than $z_{i_{e(t)}}$. Thus, the number of $u > e(t)$ for which $z_{i_u}$ is stale and $h_{j_{e(t)}}(z_{j_{e(t)}}) \geqslant h_{j_u}(z_{j_u})$ is at least $N_t(e(t))$.

Define $N = \sum N_t(e(t))$, summing over all clean arrival times $i_t$. Then the sequence of predicted stale arrivals $(h_j(z_j))_{j \in J}$ has at least $N$ inversions with respect to the strictly increasing integer sequence of actual stale arrivals. It follows from Lemma 11, which we state and prove in the next section, that

$$\sum_{j_s \in J} |h_{j_s}(z_{j_s}) - i_s| \geqslant N/2.$$

Let $\eta_{r-1}$ be the sum of prediction errors over all predictions made in phase $r - 1$. Since $J$ is a set of times in phase $r - 1$, we get from the above inequality that $\eta_{r-1} \geqslant N/2$. On the other hand, by Claim 8, the expected number of cache misses in phase $r$ is at most $\mathbb{E}\left[\sum_t O(1 + \log N_t(e(t)))\right]$, summing over all $t$ such that $i_t$ is a clean arrival time.

To simplify notation a bit, for each clean arrival time $c = i_t$, define $N_c = N_t(e(t))$. Then we have shown that $2\eta_{r-1} \geqslant \sum_{c \in r} N_c$, where the sum is over clean arrivals $c$ in phase $r$. Furthermore, we have shown that

$$\mathbb{E}[\# \text{ of cache misses in phase } r] \leqslant \sum_{c \in r} O(\mathbb{E}[1 + \log(N_c)]) \leqslant \sum_{c \in r} O(1 + \log(\mathbb{E}[N_c]))$$

7

by Jensen's inequality.

Now sum over all phases. The number of cache misses is $O(L + \sum_c \log(\mathbb{E}[N_c]))$, where the sum is over the $L$ clean arrivals in all phases. And we know that $2\eta \geqslant \sum_c N_c$ for all executions, implying that $2\eta \geqslant \sum_c \mathbb{E}[N_c]$. By another application of Jensen's inequality, the expected number of cache misses can be bounded by $O(L + L \log(2\eta/L))$.

Since we also know that $N_c \leqslant k$ for any clean arrival $c$, we can alternately bound the number of cache misses by $O(L + L \log k)$.

Combining the two bounds and using the fact that $L = \Theta(\text{OPT})$, we have proven the following theorem.

**Theorem 9.** *The algorithm* LMARKER *has competitive ratio* $O(1 + \min(\log(\eta/\text{OPT}), \log k))$, *where $\eta$ is the unknown $\ell_1$ prediction error.*

Let $H(m)$ be the $m$-th harmonic number. We can keep track of the exact constants in the proof, rather than using asymptotic notation:

**Theorem 10.** *The algorithm* LMARKER *has competitive ratio* $4 + 2H(\min(2\eta/\text{OPT}, k))$, *where $\eta$ is the unknown $\ell_1$ prediction error.*

## 4.1 Proof of combinatorial lemma

Let $M = (m_1, \ldots, m_n)$ be a strictly increasing integer sequence. For any integer sequence $A = (a_1, \ldots, a_n)$ let $\text{inv}(A)$ be the number of pairs of indices $i < j$ such that $a_i \geqslant a_j$. Let $\text{cost}(A) = \sum_{i=1}^n |a_i - m_i|$, and define $\Delta(A) = 2\text{cost}(A) - \text{inv}(A)$.

**Lemma 11.** *Let* $A = (a_1, \ldots, a_n)$ *be an arbitrary integer sequence. Then* $\text{inv}(A) \leqslant 2\text{cost}(A)$, *with* $\text{inv}(A)$ *and* $\text{cost}(A)$ *as defined above.*

*Proof.* Without loss of generality we can assume that all elements of $A$ are bounded between $m_1$ and $m_n$, since outliers can be thresholded without decreasing $\text{inv}(A)$ or increasing $\text{cost}(A)$. So the set of sequences is finite. Let $B = (b_1, \ldots, b_n)$ be a sequence which minimizes $\Delta$, and assume that $B$ has the maximum possible sum of elements out of all sequences minimizing $\Delta$.

Suppose that $B$ were to have a strict inversion. Then there is some $1 \leqslant i < n$ such that $b_i > b_{i+1}$. Define two sequences $B^l$ and $B^h$ so that $b_i^l = b_{i+1}$ and $b_{i+1}^h = b_i$, and in all other locations $B^l$ and $B^h$ agree with $B$. Then by construction,

$$\text{inv}(B) - \text{inv}(B^l) = \text{inv}(B^h) - \text{inv}(B).$$

Furthermore,

$$\begin{aligned}
\text{cost}(B) - \text{cost}(B^l) &= |b_i - m_i| - |b_{i+1} - m_i| \\
&\geqslant |b_i - m_{i+1}| - |b_{i+1} - m_{i+1}| \\
&= \text{cost}(B^h) - \text{cost}(B)
\end{aligned}$$

since $m_i < m_{i+1}$, and the function $|b_i - x| - |b_{i+1} - x|$ is non-increasing. It follows from the above inequalities and the optimality of $B$ that

$$\Delta(B^h) - \Delta(B) \leqslant \Delta(B) - \Delta(B^l) \leqslant 0.$$

Therefore $B^h$ minimizes $\Delta$ as well. But it has a larger sum of elements than $B$. Contradiction.

So we know that $B$ is non-decreasing, and thus $\text{inv}(B)$ is exactly the number of pairs of equal elements. But a constant sequence of length $l$ has cost at least $\binom{l}{2}/2$, and contributes only $\binom{l}{2}$ pairs of equal elements. Partitioning $B$ into constant sequences, we get $\text{inv}(B) \leqslant 2\text{cost}(B)$. □

# 5 Improved algorithm

In the previous section, we saw that using the predictions once at the beginning of each eviction chain takes the chain most of the way through the phase; random evictions of unmarked elements are then used until the chain ends. Suppose that the second element of an eviction chain does not reappear until $O(\eta/\text{OPT})$ steps away from the end of the phase. Then evicting a uniformly random element of the cache—marked or unmarked—would terminate the chain immediately with probability $1 - O((\eta/\text{OPT})/k)$.

In this section we present an improved algorithm LNONMARKER motivated by the above observation. As the name may suggest, it is not quite a marker-based algorithm, and we need to give some new names to familiar concepts:

**Definition 12.** *For any input sequence $z$ and cache size $k$, the* phases *of the input sequence are defined recursively as follows: phase $r$ begins right after the end of phase $r-1$, and extends as long as possible without containing $k + 1$ distinct elements.*

**Definition 13.** *Fix an algorithm. An input element is called* initial *for some phase $r$ if it appeared in phase $r$, and was present in the cache at the beginning of phase $r$. If it appeared in phase $r$ but was not present in the cache at the beginning of the phase, it is called* non-initial.

Note that the definition of phases given here coincides with the phases of any marker-based algorithm. This definition is algorithm-independent, and thus is also useful for non-marker-based algorithms. For any marker-based algorithm, the definitions of *clean* and *non-initial* coincide, as do the definitions of *stale* and *initial*. However, it can be seen that these may diverge in the execution of a non-marker-based algorithm. Some facts about clean and stale elements which we used in the previous section are now facts about initial and non-initial elements:

**Claim 14.** *Every phase contains exactly $k$ arrivals. Every non-initial arrival causes a cache miss. Every other cache miss in the phase was caused by some previous cache miss in the same phase.*

Thus, in analogy with clean arrivals, every non-initial arrival heads an eviction chain, and the eviction chains partition the set of cache misses. However, unlike before, not all cache misses are necessarily upon arrivals: an element might arrive, be evicted, and reappear all in one phase.

**Algorithm description** The new algorithm LNONMARKER still maintains markings on cache elements. At the beginning of each phase, all cache elements are unmarked. Whenever a cache hit occurs, the element is marked. Whenever a cache miss occurs, the algorithm evicts some element and marks the new element. In particular, the algorithm has the following eviction strategy upon a cache miss: if the incoming element was non-initial, then evict the unmarked element with the highest predicted arrival time. If the incoming element was evicted by a non-initial element, then evict a uniformly random element of the cache (not necessarily unmarked). In all other cases, evict a uniformly random unmarked element.

Before analyzing the algorithm, we must show that it is in fact well-defined.

**Claim 15.** *At any cache miss, there is at least one unmarked element. Thus, the algorithm is well-defined.*

*Proof.* Each marked element of the cache must have appeared in the current phase. The element which caused the cache miss is, of course, distinct from all elements of the cache. Since the phase contains at most $k$ distinct elements, the cache contains at most $k - 1$ marked elements. $\square$

**Algorithm analysis** Fix a phase $r$ in which the set of elements is $A$. Fix a single execution of the algorithm. Let $S$ be the cache at the beginning of phase $r$. Then $A \backslash S$ is the set of non-initial elements, and $S \backslash A$ is the set of cache elements which do not arrive in phase $r$. Two facts relate these sets:

- For each non-initial element $c$, the corresponding eviction chain evicts at most one element of $S \backslash A$.

- Every $x \in S \backslash A$ is evicted at most once in the phase: $x$ does not appear in the phase, so after it is evicted it will not return to the cache.

Together with the observation that $|A| = |S| = k$, these facts imply that there is a bijection $f : A \backslash S \to S \backslash A$ such that if $c$'s chain evicts $x \in S \backslash A$ then $f(c) = x$. For any $c \in A \backslash S$ such that $c$'s chain does not evict any element of $S \backslash A$, we set $f(c)$ arbitrarily, subject to the condition that $f$ is a bijection.

By Claim 14, the number of cache misses in phase $r$ is the total length of the eviction chains headed by the non-initial elements of phase $r$. Let $c$ be one such non-initial element, with eviction chain of length $length(c)$, which arrives at time $t_c$ and evicts some unmarked element $e$. Suppose $e$ reappears in the phase. Let $N(c)$ be the number of first arrivals after $e$ in the same phase, and let $N^*(c)$ be the number of distinct elements after $e$ (not necessarily first arrivals).

The above definitions were made for a single execution, but the algorithm defines a distribution $\mathcal{E}$ over executions. Let $\mathcal{E}_r$ be the distribution of executions of the first $r - 1$ phases, and let $\mathcal{E}_{t_c}$ be the distribution of executions through time $t_c$.

Conditioned on $\mathcal{E}_r$, the non-initial elements $c$ are determined, but $N^*(c)$ and $N(c)$ are random variables (over the randomness of the algorithm). Defining $length(c) = 0$ and $N^*(c) = 0$ if $c$ is not a non-initial element or if $e$ does not reappear in the phase, we get the following result.

**Claim 16.** *For any element $c$ in phase $r$,*

$$\mathbb{E}_{\mathcal{E}}[length(c)|\mathcal{E}_r] \leqslant \alpha \cdot \mathbb{E}_{\mathcal{E}}\left[1 + \frac{N^*(c)}{k}\log N(c)\Big|\mathcal{E}_r\right]$$

*for an absolute constant $\alpha$.*

*Proof.* Condition on $\mathcal{E}_{t_c}$. Then $e$ is determined, and $N^*(c)$ and $N(c)$ are determined. If $e$ does not reappear in the phase, then $\mathbb{E}[length(c)|\mathcal{E}_{t_c}] = 0$. Suppose otherwise. Since $e$ evicts a random element $g$ from the cache, the probability that $g$ subsequently appears during the same phase is at most $p = N^*(c)/k$. If $g$ does appear, the eviction chain continues by random eviction of unmarked elements, until an evicted element is not in $A$. The number of unmarked elements at the time of $g$'s cache miss which are present in $A$ is at most $N(c)$, since any element which arrived earlier in the phase is either marked or no longer in the cache. Thus, by the same argument as in Claim 8, the expected length of the chain is $O(\log N(c))$, conditioned on $g$'s cache miss. Since the length is $O(1)$ if $g$ does not appear, it follows that $\mathbb{E}[length(c)|\mathcal{E}_{t_c}]$ is at most $O(1 + p\log N(c))$. Taking the expectation over $\mathcal{E}_{t_c}|\mathcal{E}_r$ yields the claimed result. $\square$

Now we want to bound $N^*(c)$ in terms of the predictor error. We condition on the entire execution of the algorithm: that is, the following lemmas hold deterministically for all executions.

**Lemma 17.** *For any chain $(c, e, \dots)$ in which $e$ reappears after eviction, let $t$ be $c$'s arrival time. Let $r'$ be the next phase in which $f(c)$ arrives. Then*

$$\mathrm{Err}_{L(f(c),t)}(f(c)) + \mathrm{Err}_{L(e,t)}(e) \geqslant N^*(c) + k(r' - r - 1).$$

*Proof.* Note that $f(c)$ is an unmarked element of the cache at time $t$. It does not appear in phase $r$, so if phase $r+1$ begins at time $t_{r+1}$ then $f(c)$ does not appear until at least time $t_{r+1}+k(r'-r-1)$, but its predicted next appearance $h_{L(f(c),t)}(f(c))$ satisfies $h_{L(f(c),t)}(f(c)) \leqslant h_{L(e,t)}(e)$, since $e$ was the unmarked cache element which maximized predicted arrival at time $t$. The next appearance of $e$ is no later than time $t_{r+1} - N^*(c)$, since at least $N^*(c)$ elements lie between $e$'s appearance and the end of the phase. So $\text{Err}_{L(f(c),t)}(f(c)) + \text{Err}_{L(e,t)}(e) \geqslant N^*(c) + k(r'-r-1)$. □

*Remark* 1. There is an edge case in Lemma 17 if $f(c)$ never appears after phase $r$ and the last phase has less than $k$ elements. In this case the inequality in Lemma 17 can be off by as much as $k$. However, when the lemma is applied in Lemmas 19 and 20, the applications almost always implicitly weaken the inequality by at least $k$ anyway. Each application where the inequality is not weakened corresponds to a chain in one of the last two phases. There are only $O(k)$ such chains, each causing a discrepancy of at most $k$. To make the lemmas hold, it therefore suffices to replace $\eta$ by $\eta + O(k^2)$. However, $O(k^2)$ is bounded as the input length and OPT grow, so the competitive ratio is unaffected.

Next we would like to sum the inequality given by Lemma 17 across all chains in all phases, to bound $\sum N^*(c)$ in terms of the total prediction error. Within a single phase, $f(c)$ and $e$ each uniquely determine $c$, so the error of each prediction is counted at most once. However, the error of a single prediction may be counted in multiple chains in successive phases. The following lemmas do a more careful summation, taking into consideration the double-counting.

**Lemma 18.** *For distinct times $t_1, t_2$ suppose $c_1 = z_{t_1}$ and $c_2 = z_{t_2}$ are non-initial elements of phases $r_1$ and $r_2$, which evict $e_1$ and $e_2$ respectively. Suppose that $e_1$ and $e_2$ reappear after eviction. Then $L(e_1, t_1) \neq L(e_2, t_2)$.*

*Proof.* If equality were to hold, then $e_1 = e_2 = w$ for some $w$. Note that $w$ is evicted at $c_1$'s arrival, and $w$ is evicted at $c_2$'s arrival, but $w$ is unmarked both times. So the evictions occurred in different phases. Without loss of generality suppose that $r_1 < r_2$. Then $w$ appeared in phase $r_1$ after eviction by $c_1$, so $L(w, t_1) < L(w, t_2)$. □

**Lemma 19.** *Summing over all chains $(c, e, \dots)$ in which $e$ reappears after eviction,*

$$\sum_c N^*(c_i) \leqslant 3\eta$$

*where $\eta$ is the total prediction error.*

*Proof.* Fix a time $t$ and consider the set of chains $(c_i, e_i, \dots)$ such that $L(f(c_i), t_i) = t$, where $t_i$ is the arrival time of $c_i$: that is, $f(c_i) = z_t$ and the most recent appearance of $f(c_i)$ prior to time $t_i$ was at time $t$. Each chain is in a different phase, since $f$ is injective for any one phase. If there are $m_t$ such chains in phases $r_1 < \dots < r_{m_t}$, then $z_t$ does not appear after time $t$ until phase $r_{m_t} + 1$ or later. Applying Lemma 17 to each of the $m_t$ chains and summing, we get

$$m_t \text{Err}_t(z_t) + \sum_{i=1}^{m_t} \text{Err}_{L(e_i, t_i)}(e_i) \geqslant k\frac{m_t(m_t - 1)}{2} + \sum_{i=1}^{m_t} N^*(c_i) \geqslant \sum_{i=1}^{m_t} \frac{m_t + 1}{2} N^*(c_i).$$

Dividing both sides by $(m_t + 1)/2$, we get that

$$2\text{Err}_t(z_t) + \sum_{i=1}^{m} \text{Err}_{L(e_i, t_i)}(e_i) \geqslant \sum_{i=1}^{m_t} N^*(c_i).$$

Now sum over all times $t$. By Lemma 18, all $L(e_i, t_i)$ are distinct, so the corresponding errors sum to at most $\eta$. Thus, we have

$$3\eta \geqslant \sum_c N^*(c_i)$$

summing over all chains $(c, e, \dots)$ in which $e$ reappears after eviction. $\qquad\square$

Finally, we must bound the number of eviction chains $C$ across all phases. In a marker-based algorithm we would have $C = L$, with one chain per clean element. But here there is one chain per non-initial element, and not all non-initial elements are clean (and vice versa). Nonetheless, we can still bound the discrepancy $C - L$ against the prediction error. The intuition for the following proof is that any "extra" non-initial element must have been caused by eviction of a marked element in the previous phase. But for each chain which ends by evicting a marked element, one less unmarked element is evicted. So whereas in a marker-based algorithm, every cache element which did not appear in a given phase would be evicted by the end of the phase, in this algorithm some absent elements might remain in the cache. These elements have the potential to be clean in the next phase and yet not start eviction chains. Hence, in one case an extra chain can be charged against a non-existent chain. In the other case—when the absent element does not appear in the next phase either—the extra chain can be charged against prediction error.

**Lemma 20.** *If $C$ is the number of non-initial elements across all phases, $L$ is the number of clean elements, and $\eta$ is the total prediction error, then $\eta \geqslant k(C - L)/2$.*

*Proof.* Suppose some non-initial element $g$ in phase $r + 1$ is not clean. Then it appeared in phase $r$, but was not in the cache at the end of phase $r$. In the last eviction of $g$ during phase $r$, $g$ was already marked. So by the algorithm design, that eviction must be third in some eviction chain $(c, e, g)$, where $c$ is some non-initial element in phase $r$. Then $f(c)$ does not appear in phase $r$. Nonetheless, $f(c)$ is not in $c$'s eviction chain, since $g$ is the unique element in $c$'s eviction chain which does not reappear in phase $r$. Furthermore, by definition of $f$, no other chain in phase $r$ evicts $f(c)$, so it is in the cache at the end of phase $r$. There are two cases:

1. $f(c)$ appears in phase $r + 1$. Then $f(c)$ is clean and initial in phase $r + 1$.

2. $f(c)$ does not appear in phase $r + 1$. Let $r'$ be the next phase in which $f(c)$ does appear. Then if $t$ is the arrival time of $c$, by Lemma 17

$$\mathrm{Err}_{L(f(c),t)}(f(c)) + \mathrm{Err}_{L(e,t)}(e) \geqslant k(r' - r - 1).$$

The composed map $g \mapsto c \mapsto f(c)$ is injective for fixed phase $r$, since $g$ is determined by the eviction chain of $c$, and $f$ is injective for a fixed phase. If case (2) never occurred, then every non-initial, non-clean element in phase $r + 1$ would be injectively mapped to an initial, clean element in phase $r + 1$, so we would have $C \leqslant L$. More generally, case (2) occurs for at least $C - L$ chains. Each case (2) chain provides an inequality bounding two prediction errors by at least $k$, and ideally we would simply add up the inequalities to bound $\eta$. However, some predictions may be counted in several of the inequalities.

By Lemma 18, all terms $\mathrm{Err}_{L(e,t)}(e)$ in the inequalities are contributed by different predictions—i.e. adding up those error terms does not double-count.

Consider a fixed time $t$. Consider the set of case-(2) chains $(c_i, e_i, g_i)$ such $L(f(c_i), t_i) = t$, where $t_i$ is the time of $c_i$'s arrival. Each chain is in a different phase, since $f$ is injective for any one

phase. If there are $m_t$ such chains in phases $r_1 < \cdots < r_{m_t}$, then $z_t$ does not appear after time $t$ until phase $r_{m_t} + 2$ or later. So the case-(2) inequality applied to the earliest chain $(c_1, e_1, g_1)$ gives

$$\mathrm{Err}_t(z_t) + \mathrm{Err}_{L(e_1, t_1)}(e_1) \geqslant k(r_{m_t} + 1 - r_1) \geqslant km_t.$$

Summing the above inequality over all times $t$, each prediction error is counted at most twice—once as the first term and once as the second—whereas $\sum_t m_t \geqslant C - L$ as previously shown. Hence, $2\eta \geqslant k(C - L)$, as desired. $\qquad\square$

With the above error/performance bounds in hand, we can prove a bound on the competitive ratio of the algorithm.

**Theorem 21.** *The algorithm* LNONMARKER *achieves a competitive ratio of*

$$O\left(1 + \frac{\eta/\mathrm{OPT}}{k} \log k\right)$$

*when the prediction error is no more than $\eta$.*

*Proof.* Fix any phase $r$. Then for any execution $E_r$ of the first $r - 1$ phases,

$$\mathbb{E}[\# \text{ cache misses in phase } r | \mathcal{E}_r = E_r] = \sum_{(c,e,\dots) \in r} \mathbb{E}[\mathrm{length}(c) | \mathcal{E}_r = E_r],$$

summing over chains in phase $r$.

By Claim 16, the right hand side is bounded by

$$\alpha C_r + \alpha \frac{\log k}{k} \mathbb{E}\left[\sum_{(c,e,\dots) \in r} N^*(c) \,\middle|\, \mathcal{E}_r = E_r\right],$$

where $C_r$ is the number of chains in phase $r$, and $\alpha$ is a constant. Applying law of total expectation and then summing over all phases,

$$\mathbb{E}[\# \text{ cache misses}] \leqslant \alpha \mathbb{E}[C] + \alpha \frac{\log k}{k} \mathbb{E}\left[\sum_{(c,e,\dots)} N^*(c)\right].$$

From Lemma 20, the inequality $C \leqslant 2\eta/k + L$ holds for all executions, and therefore in expectation. From Lemma 19, the inequality $\sum_c N^*(c) \leqslant 3\eta$ holds for all executions, where $c$ ranges over chains $(c, e, \dots)$ where $e$ reappears after eviction. Recall that $N^*(c)$ was defined to be 0 for all other chains. So we get that

$$\mathbb{E}[\# \text{ cache misses}] \leqslant 2\alpha \frac{\eta}{k} + \alpha L + 3\alpha \frac{\eta}{k} \log k.$$

Recalling from Lemma 5 that $L/2 \leqslant \mathrm{OPT}$, the competitive ratio follows. $\qquad\square$

The above algorithm has one flaw: it does not satisfy any robustness guarantee (at least, that has been proven), since as $\eta/\mathrm{OPT} \to \infty$ the bound on the number of chains disappears, and the competitive ratio becomes potentially unbounded. This can be resolved (albeit in a somewhat unsatisfactory manner) by appealing to a black-box simulation theorem. For concreteness, we recall the following theorem:

**Theorem 22.** *[8] Let $A, B$ be algorithms for the caching problem with competitive ratios of $\alpha$ and $\beta$ respectively. Then there is a black box algorithm $ALG$ with a competitive ratio of $9 \min(\alpha, \beta)$.*

The black box algorithm in the above theorem proceeds by simulating $A$ and $B$ and switching between them whenever one starts to heavily outperform the other. The proof generalizes without change to the learned caching problem. Since there is an $O(\log k)$-competitive algorithm for learned caching which simply ignores the predictions, Theorem 22 implies that we can obtain an $O(\log k)$ worst-case guarantee for our predictive caching algorithm with only an extra constant factor loss:

**Corollary 23.** *There is an algorithm for caching with predictions that achieves a competitive ratio of*

$$O \left( 1 + \min \left( 1, \frac{\eta/\text{OPT}}{k} \right) \log k \right).$$

Tracing through the proof, it turns out that the exact bound is as follows:

**Corollary 24.** *There is an algorithm for caching with predictions that achieves a competitive ratio of*

$$9 \min \left( 4 + 7\frac{\eta/\text{OPT}}{k} + 3\frac{\eta/\text{OPT}}{k} H(k), 2H(k) \right).$$

For example, as $\eta/\text{OPT} \to 0$, the competitive ratio of LNONMARKER approaches 4, and so the competitive ratio of this black box algorithm approaches 36.

# 6 Lower bound

In this section we provide a lower bound against predictive caching algorithms. The basic strategy is to construct a distribution of inputs and predictions such that the relative prediction error $\eta/\text{OPT}$ is not too high, but any deterministic algorithm which has access to those predictions suffers a large number of cache misses in expectation, relative to OPT. Yao's minimax principle then implies a lower bound against the competitive ratios of randomized algorithms at that level of relative prediction error.

More specifically, the input distribution and predictions will be chosen such that any prefix of the input completely determines the state of the algorithm. Furthermore, each prefix is sufficiently independent of future inputs that the algorithm can essentially do no better than (a) keeping previously-seen elements in the cache (for the duration of the phase), and (b) guessing arbitrarily for the remaining unseen elements.

Fix $k$ and $n$. Let $1 \leqslant t \leqslant k$ be picked later; it is a free parameter which will determine the relative prediction error $\eta/\text{OPT}$. Let $\Omega$ be the set of sequences that can be constructed in the following manner. Let $C_1 = \{1, \ldots, t\}$ and let $S_1 = \{t + 1, \ldots, k\}$. For all $2 \leqslant r \leqslant n$, let $C_r = [k + t] \backslash (C_{r-1} \cup S_{r-1})$ and let $S_r$ be an arbitrary subset of $C_{r-1} \cup S_{r-1}$ of size $k - t$. Then each sequence of $\Omega$ is constructed as the concatenation of $n$ *phases*, where phase $r$ consists of $3k \log k$ elements of $C_r \cup S_r$ (possibly with some omissions and necessarily with some repetitions), followed by a single copy of $C_r \cup S_r$ in increasing order, without repetitions.

That is, each phase has length $m = 3k \log k + k$. Each phase $r$ has $t$ clean elements $C_r$ (which did not appear in the previous phase) and $k - t$ stale elements $S_r$ (which did appear). For any fixed $r > 1$, a uniformly random sequence of $\Omega$, conditioned on $C_r$, has a uniformly random set of stale elements $S_r$. Furthermore, conditioned on $C_r$ and $S_r$, each of the first $3k \log k$ elements of phase $r$ is independent and uniformly distributed over $C_r \cup S_r$.

We must also define the predictions for a fixed sequence. For the first $3k \log k$ elements of phase $r$, each prediction is the subsequent timestep (i.e. $h(z_i) = i + 1$). For the final copy of $C_r \cup S_r$, each prediction is the end of phase $r + 1$.

**Claim 25.** *The overall prediction error is $O(nk^2 \log k)$, where $n$ is the number of phases in $\Omega$.*

*Proof.* Fix a phase $r$. Among the first $3k \log k$ elements of the phase, there are at most $k$ distinct elements. For each, the prediction error telescopes to at most $3k \log k + k$. For each of the $k$ elements in the final copy of $C_r \cup S_r$, the true next arrival is either in phase $r + 1$ or phase $r + 2$, so the error is at most $O(k \log k)$. Thus, the error in phase $r$ is $O(k^2 \log k)$. $\square$

Now we would like to lower bound the average cache misses of any algorithm on $\Omega$. We must first give some simplifying notation and a probabilistic lemma:

For any time $T$ in phase $r$, let $\sim_T$ be the equivalence relation on $\Omega$ where sequences $z$ and $z'$ are equivalent if $z_j = z'_j$ for $j \leqslant T$. Observe that $z \sim_T z'$ implies that $h(z_j) = h(z'_j)$ for all $j \leqslant T$ by how the predictions were constructed. Thus, for any equivalence class $[z]_T$ of $\sim_T$, the algorithm has identical executions on sequences in $[z]_T$ up to time $T$.

For $z \in \Omega$ and time $T$, let $\mathcal{C}_{z,T}$ be the cache at time $T$ on input $[z]_T$.

**Lemma 26.** *Fix $k, l$ with $2 \leqslant l \leqslant k$. Let $X_1, \ldots, X_{3k \log k}$ be independent random variables uniformly distributed over $[k]$. Let $Y_i$ be the number of distinct elements in $\{X_1, \ldots, X_i\} \cap [l]$ and for $0 \leqslant j < l$ let $T_j$ be the number of $i$ such that $Y_i = j$. Then $\mathbb{E}[T_j] \geqslant k/(l - j) - 1/k$.*

*Proof.* Extend the sequence $X$ to an infinite sequence. Extend $Y$ accordingly, and define $\hat{T}_j$ as the number of steps in the infinite sequence at which $Y_i = j$. For $0 \leqslant j < l$, let $S_j = T_0 + \cdots + T_j$ and $\hat{S}_j = \hat{T}_0 + \cdots + \hat{T}_j$. Then $S_j = \min(\hat{S}_j, 3k \log k)$, so

$$\mathbb{E}[T_j] \geqslant \mathbb{E}[\min(\hat{S}_j, 3k \log k)] - \mathbb{E}[\hat{S}_{j-1}]$$
$$= \mathbb{E}[\hat{T}_j] - \mathbb{E}[(\hat{S}_j - 3k \log k)\mathbb{1}_{\hat{S}_j \geqslant 3k \log k}].$$

For the first term, $\hat{T}_j$ is a geometric random variable and $\mathbb{E}[\hat{T}_j] = k/(l - j)$. For the second,

$$\mathbb{E}[(\hat{S}_j - 3k \log k)\mathbb{1}_{\hat{S}_j \geqslant 3k \log k}] \leqslant \sum_{c=3k \log k}^{\infty} \Pr[\hat{S}_j \geqslant c].$$

Observe that for any $c$, if $\hat{S}_j \geqslant c$ then $Y_c \leqslant j < l$, which occurs with probability at most $l(1 - 1/k)^c$. It follows that

$$\mathbb{E}[(\hat{S}_j - 3k \log k)\mathbb{1}_{\hat{S}_j \geqslant 3k \log k}] \leqslant \sum_{c=3k \log k}^{\infty} l\left(1 - \frac{1}{k}\right)^c \leqslant \frac{1}{k}.$$

We conclude that $\mathbb{E}[T_j] \geqslant k/(l - j) - 1/k$. $\square$

**Theorem 27.** *For any deterministic algorithm with access to the predictions $h$, the expected number of cache misses achieved on an input sampled uniformly at random from $\Omega$ is at least*

$$O\left(nt \log \frac{k}{t}\right).$$

*Proof.* There are $n$ phases in $\Omega$. Fix a phase $r$. We claim that the expected number of cache misses in phase $r$ is at least $O(t \log k/t)$.

Recall that each phase has length $m = 3k \log k + k$. Let $T = (r-1)m$: the final index of phase $r-1$. Fix $z \in \Omega$. Fix $1 \leqslant i \leqslant k \log k$ and consider any class $[w]_{T+i} \subseteq [z]_T$. Picking a sequence $W \in [w]_{T+i}$ uniformly at random, the clean elements $C_r$ and the cache $\mathcal{C}_{W,T+i}$ are determined, but $S_r$ is a random variable. Let $S_{\text{seen}}$ be the set of elements which have already been seen in phase $r$ by time $T + i$, excluding $C_r$. Then $S_r$ must contain $S_{\text{seen}}$, but $S_r \backslash S_{\text{seen}}$ is a uniformly random subset of $[k+t] \backslash (C_r \cup S_{\text{seen}})$ of size $k - t - |S_{\text{seen}}|$. Hence,

$$\Pr_W[\text{cache hit on } W_{T+i}] = \sum_{\hat{S}} \Pr_W[S_r = \hat{S}] \Pr_W[\text{cache hit on } W_{T+i}|S_r = \hat{S}]$$

$$= \sum_{\hat{S}} \Pr_W[S_r = \hat{S}] \frac{|\mathcal{C}_{W,T+i} \cap (C_r \cup \hat{S})|}{k}$$

$$= \frac{|\mathcal{C}_{W,T+i} \cap (C_r \cup S_{\text{seen}})| + \mathbb{E}_W |\mathcal{C}_{W,T+i} \cap (S_r \backslash S_{\text{seen}})|}{k}$$

If $|\mathcal{C}_{W,T+i} \cap (C_r \cup S_{\text{seen}})| = a$, then $\mathcal{C}_{W,T+i}$ contains $k - a$ elements in $[k+t] \backslash (C_r \cup S_{\text{seen}})$, each of which is contained in $S_r \backslash S_{\text{seen}}$ with probability $(k-t-|S_{\text{seen}}|)/(k-|S_{\text{seen}}|)$. So the above expression is maximized when $a = |C_r \cup S_{\text{seen}}| = k + t$. Thus,

$$\Pr_W[\text{cache hit on } W_{T+i}] \leqslant \frac{t + |S_{\text{seen}}| + (k - t - |S_{\text{seen}}|) \cdot \frac{k-t-|S_{\text{seen}}|}{k-|S_{\text{seen}}|}}{k}.$$

Simplifying, it follows that

$$\Pr_W[\text{cache miss on } W_{T+i}] \geqslant \frac{t}{k} \frac{k - t - |S_{\text{seen}}|}{k - |S_{\text{seen}}|}.$$

If $Z \in [z]_T$ is chosen uniformly at random, then for any $0 \leqslant N < k - t$, the expected number of times at which $|S_{\text{seen}}| = N$ is at least $k/(k-t-N) - 1/k \geqslant k/(2(k-t-N))$ by Lemma 26. Summing over $N$, the expected number of cache misses in phase $r$ on input $Z$ is at least

$$\sum_{N=0}^{k-t-1} \frac{t}{k} \frac{k-t-N}{k-N} \frac{k}{2(k-t-N)} = \sum_{N=0}^{k-t-1} \frac{t}{2(k-N)} \geqslant \frac{t}{2}\left(\frac{1}{t+1} + \cdots + \frac{1}{k}\right) = \Omega(t \log k/t).$$

Since the equivalence classes of $\sim_T$ partition $\Omega$, and the above bound holds for each class, it holds for $\Omega$ as desired. $\qquad\square$

Now for any fixed "relative prediction error" $\epsilon$ with $k \log k \leqslant \epsilon \leqslant k^2 \log k$, we can pick $t = (k^2 \log k)/\epsilon$. Then every sequence in $\Omega$ has $nt$ clean elements, and thus $\text{OPT} = \Theta(nt)$. Furthermore, by Claim 25, we have $\eta \leqslant nk^2 \log k$. Thus, $\eta/\text{OPT} \leqslant \epsilon$.

But by Theorem 27, any deterministic algorithm requires $nt \log k/t$ cache misses in expectation on $\Omega$. By Yao's minimax principle, for any randomized algorithm there is some input $z \in \Omega$ for which the algorithm incurs $nt \log k/t$ cache misses in expectation. Hence, we have the following result:

**Theorem 28.** *Let $A$ be a randomized online algorithm for caching, which has access to next-arrival predictions. For any $\epsilon$, the algorithm achieves competitive ratio no better than*

$$\Omega\left(\log \min\left(\frac{\epsilon}{k \log k}, k\right)\right)$$

*when restricted to inputs with $\eta/\text{OPT} \leqslant \epsilon$.*

# References

[1] Dimitris Achlioptas, Marek Chrobak, and John Noga. Competitive analysis of randomized paging algorithms. *Theoretical Computer Science*, 234(1):203–218, 2000.

[2] Deepak Ajwani and Tobias Friedrich. Average-case analysis of online topological ordering. In *Proceedings of the 18th International Conference on Algorithms and Computation*, ISAAC'07, pages 464–475, Berlin, Heidelberg, 2007.

[3] Amos Fiat, Richard M. Karp, Michael Luby, Lyle A. McGeoch, Daniel D. Sleator, and Neal E. Young. Competitive paging algorithms. *Journal of Algorithms*, 12(4):685–699, 1991.

[4] Sreenivas Gollapudi and Debmalya Panigrahi. Online algorithms for rent-or-buy with expert advice. In *International Conference on Machine Learning*, pages 2319–2327, 2019.

[5] Chen-Yu Hsu, Piotr Indyk, Dina Katabi, and Ali Vakilian. Learning-based frequency estimation algorithms. In *International Conference on Learning Representations*, 2019.

[6] Elias Koutsoupias and Christos Papadimitriou. Beyond competitive analysis. *SIAM Journal on Computing*, 30(1):300–317, 2000.

[7] Tim Kraska, Alex Beutel, Ed H. Chi, Jeffrey Dean, and Neoklis Polyzotis. The case for learned index structures. In *Proceedings of the 2018 International Conference on Management of Data*, SIGMOD '18, pages 489–504, New York, NY, USA, 2018.

[8] Thodoris Lykouris and Sergei Vassilvitskii. Competitive caching with machine learned advice. In *International Conference on Machine Learning*, pages 3302–3311, 2018.

[9] Lyle A. McGeoch and Daniel D. Sleator. A strongly competitive randomized paging algorithm. *Algorithmica*, 6(1):816–825, Jun 1991.

[10] Andrew McGregor. Graph stream algorithms: A survey. *SIGMOD Record*, 43(1):9–20, May 2014.

[11] Andrés Muñoz Medina and Sergei Vassilvitskii. Revenue optimization with approximate bid predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, pages 1856–1864, USA, 2017.

[12] Michael Mitzenmacher. Scheduling with predictions and the price of misprediction. *arXiv preprint arXiv:1902.00732*, 2019.

[13] Manish Purohit, Zoya Svitkina, and Ravi Kumar. Improving online algorithms via ML predictions. In *Advances in Neural Information Processing Systems 31*, pages 9661–9670. 2018.

[14] Daniel D. Sleator and Robert E. Tarjan. Amortized efficiency of list update and paging rules. *Communications of the ACM*, 28(2):202–208, February 1985.

[15] Daniel A. Spielman and Shang-Hua Teng. Smoothed analysis of algorithms: Why the simplex algorithm usually takes polynomial time. *Journal of the ACM*, 51(3):385–463, May 2004.

[16] Daniel A. Spielman and Shang-Hua Teng. Smoothed analysis: an attempt to explain the behavior of algorithms in practice. *Communications of the ACM*, 52(10):76–84, 2009.