# Research Report

## Analysis of page replacement policies in the fluid limit

Ryo Hirade, Takayuki Osogami

IBM Research, Tokyo Research Laboratory
IBM Japan, Ltd.
1623-14 Shimotsuruma, Yamato
Kanagawa 242-8502, Japan

# Analysis of page replacement policies in the fluid limit

Ryo Hirade[*]        Takayuki Osogami

### Abstract

The performance of storage systems and database systems depends significantly on the page replacement policies. Although many page replacement policies have been discussed in the literature, their performances are not fully understood except for simple page replacement policies such as Least Recently Used. We introduce analytical techniques for evaluating the performances of page replacement policies including Two Queue (`2Q`), which manages two buffers to capture both the recency and frequency of requests. We derive an exact expression for the probability that a requested item is found (the hit probability) in a buffer managed by `2Q` in the fluid limit, where the number of items is scaled by $n$, the size of items is scaled by $1/n$, and $n$ approaches infinity. The hit probability in the fluid limit approximates the hit probability in the original system, and we find that the relative error in the approximation is typically within 1%. Our analysis also illuminates several fundamental properties of `2Q` useful for system designers.

## 1   Introduction

Caching data is prevalent in today's computer and communication systems. Operating systems and database management systems cache data in faster main memory to avoid accessing slower disks [21, 22]. Webpages are cached at intermediate servers to reduce network traffic, delays perceived by users, and loads at the original Web servers [20]. The effectiveness of caching is determined by what is cached. In database management systems, when a requested item is not found in main memory, the item must be copied from a disk to the main memory, and some item may need to be evicted from the main memory to make room for the requested item. A page replacement policy determines the item to be evicted, where the primary goal is to maximize the probability that the items requested in the future will be cached and found in the main memory. Below, we assume that items have a fixed size (*i.e.*, items are pages of data), and we refer to a cache to store the items as a buffer.

The most popular page replacement policy is Least Recently Used (`LRU`), which replaces the item that was requested least recently with a new item [22]. `LRU` is efficient in that replacement can be performed in $O(1)$ time. A well known drawback of `LRU` is that an item that is requested only infrequently is kept in a buffer until the item becomes least-recently requested and is evicted without ever being requested again [19]. Also well known is Least Frequently Used (`LFU`), which replaces the item that has been requested least frequently with a new item [22]. `LFU` only keeps items that are frequently requested in a buffer, but it requires $O(\log K)$ time for replacement, where $K$ is the size of the buffer. Also, `LFU` cannot quickly adapt to changes in the workload, since it ignores the recency of requests.

The complementary properties of `LRU` and `LFU` motivated researchers to investigate page replacement policies that take into account both the recency and frequency of requests. O'Neil *et al.* propose `LRU-`$k$ [19], which replaces the item whose $k^{\text{th}}$-to-last request is least recent with a new item. Although `LRU-`$k$ requires $O(\log K)$ time for replacement, it initiated a stream of research on efficient approximations of `LRU-`$k$. Johnson and Shasha propose Two Queue (`2Q`) [13], which mimics `LRU-2` by dividing a buffer into two parts, $B_0$ and $B_1$, and performs replacement in $O(1)$ time. Stored in $B_0$ are the items that are requested only once since the last time that the items are evicted from the buffer. Stored in $B_1$ are the items that are requested at least twice since the last eviction. When an item needs to be evicted from a part, the least recently requested item in the part is evicted. An intuition is that $B_0$ operates as a low pass filter that only allows frequently requested items to be stored in $B_1$. Adaptive Replacement Cache (`ARC`) is a variant of `2Q` and dynamically changes the sizes of $B_0$ and $B_1$ [17]. Versions of `2Q` and `ARC` are used in recent versions of database management systems[1] and file systems[2]. See [12, 24] for other page replacement policies that use multiple buffers to capture both the recency and frequency of requests.

---

[*]{rhirade,osogami}@jp.ibm.com
[1]www.postgresql.org/docs/8.0/static/release-8-0-2.html
[2]www.opensolaris.org/os/community/zfs/source/

Although numerous page replacement policies have been discussed in the literature, their relative performances are only partially understood. Page replacement policies are usually evaluated by measuring the performances against benchmarks or by trace-driven or discrete-event simulations, which are limited and time consuming. Our goal is to provide an analytical framework that not only allows us to quickly evaluate the performances of page replacement policies but also provides intuitions on their fundamental properties. Toward that end, this paper proposes analytical techniques for evaluating the performance of 2Q and studies its fundamental properties. The proposed analytical techniques may be useful for evaluating other page replacement polices, particularly those dividing a buffer into multiple parts such as [17, 24, 12].

Our primary contribution is an exact analysis of the probability that a requested item is found in a buffer managed by 2Q (the hit probability for 2Q) in the fluid limit, where the number of items is scaled by $n$, the size of items is scaled by $1/n$, and $n$ approaches infinity. We assume that requests are issued according to independent Poisson processes. An analysis in the fluid limit has been shown to be effective in understanding systems with many interactive objects, including communication networks and human systems (*e.g.*, see [1]). In our case, the hit probability in the fluid limit can be used to approximate the hit probability in the original system. In fact, the hit probability in a system with $N$ items and the hit probability in the fluid limit of the system with $N$ items usually converge as $N$ approaches infinity. Our numerical experiments suggest that the relative error in the approximation is small even for a small $N$ and within 1% for $N > 1000$. A key idea in our analysis is that $B_0$ and $B_1$ are analyzed as coupled buffers where items receive requests and invalidations that have particular correlation having partial insensitivity to the behavior of the buffers. Here, when an item is invalidated, the item is simply removed from a buffer.

Our secondary contribution is a characterization of the fundamental properties of 2Q using the analytical results in the fluid limit and simulations. In particular, we find that the hit probability for 2Q can in general be made higher than that in LRU by choosing the size of $B_0$ appropriately. We also find that the *stationary* hit probability for 2Q is higher when the size of $B_0$ is set smaller. However, simulations suggest that it takes longer for the buffer to reach the steady state when $B_0$ is smaller. As a result, 2Q may have poor *transient* hit probability when $B_0$ is set too small.

## 1.1 Prior work

Relatively little work has been done on stochastic analyses of the performances of page replacement policies. As we will review below, exact expressions for the hit probability or its fluid limit been derived only for LRU and other simpler page replacement policies, although various approximations have been proposed.

The hit probability for LRU can be derived by studying a corresponding move-to-front (MTF) list, where an item is moved to the head of the list when it is requested. The hit probability for LRU for a buffer of size $K$ coincides with the probability that the requested item is at the $K$-th position or closer to the head of the MTF list. McCabe [16] derives the first two moments of the stationary position of a requested item in an MTF list where requests are issued according to an "independent reference model," which is essentially equivalent to independent Poisson processes. The results of McCabe are extended to all moments by Gonnet *et al.* [7], to the distribution by Hendricks [9], and to the generating function by Frajolet *et al.* [6] and Fill and Holst [5].

Unfortunately, the distribution and the generating function in [9, 6, 5] are computationally hard to evaluate numerically and provide little intuition due to the complexity of their expressions. Fill [4] shows that the generating function of the stationary position, $C_N$, is simplified in the limit where the number, $N$, of items approaches infinity. Using the results of Fill, Jelenković [11] studies the fluid limit of the stationary position, $\lim_{n\to\infty} \frac{1}{n} C_{nN}$, which can be translated into the hit probability for LRU in the fluid limit.

Che *et al.* [3], Laoutaris *et al.* [14], and Hama and Hirade [8] study an approximation for the hit probability for LRU, which coincides with the hit probability in the fluid limit proved by Jelenković [11]. Further, Che *et al.* [3] and Laoutaris *et al.* [14] extend the approximation to hierarchical buffers, each of which is managed by LRU. Although these approximations are based on the idea of the fluid limit, it is unknown whether these approximations coincide with the fluid limits. Also the hierarchical buffers managed by LRU are essentially different from a buffer managed by 2Q, since an item may be stored at multiple positions in hierarchical buffers. In [14], a version of hierarchical buffers that stores an item exclusively at one position is also studied *by simulation*, not analytically. Our result is the first that derives and proves the fluid limit of the hit probability for a page replacement policy that is more sophisticated than LRU.

The rest of the paper is organized as follows. In Section 2, we start with an analysis of the hit probability for LRU in the fluid limit. We derive an expression that is essentially equivalent to Jelenković [11], but we find that our derivation is simpler. In Section 3, we extend the analysis of LRU in the fluid limit to the case where items are requested and invalidated, and the requests and the invalidations have a particular correlation. The analysis in Section 3 is extended to an analysis of the hit probability for 2Q in the fluid limit in Section 4. In
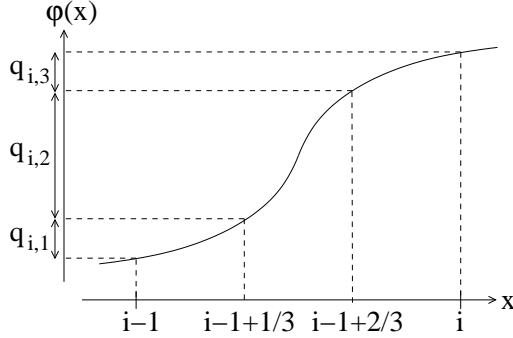
Figure 1: A $\varphi(\cdot)$ and $q_{i,j}$ when $n = 3$.

Section 5, we validate approximating the hit probability for 2Q by its fluid limit and study the fundamental properties of 2Q.

# 2 Analysis of LRU

We start by studying a buffer managed by LRU (an LRU buffer). Let $K$ be the size of the LRU buffer and $N$ be the number of items. The $N$ items, $e_i$ for $1 \le i \le N$, have size 1 and are requested independently of each other. The inter-request times of $e_i$ are independent and have a distribution function, $F_i(\cdot)$, for $1 \le i \le N$. When a requested item, $e_i$, is not in an LRU buffer, the least-recently requested item in the LRU buffer is replaced with the $e_i$ if the LRU buffer is full or the $e_i$ is simply added to the LRU buffer otherwise. When the $e_i$ is found in an LRU buffer, the $e_i$ becomes most-recently requested. Recall that the hit probability for LRU is equivalent to the probability that the position of the requested item is at most $K$ in a corresponding MTF list. To analyze the hit probability for LRU, we consider a system where requests are issued for an LRU buffer and for a corresponding MTF list at the same moments. In Section 2.1, we introduce the fluid limit of the system. In Section 2.2, we analyze the hit probability for LRU in the fluid limit.

## 2.1 Fluid limit

We consider a sequence of systems, $\mathcal{S}^{(n)}$ for $n = 1, 2, ...$, where each system is associated with an LRU buffer of size $K$ and with a corresponding MTF list. In $\mathcal{S}^{(n)}$, requests are generated independently for $n\,N$ items, $e_{i,j}$ for $1 \le i \le N$ and $1 \le j \le n$, of size $1/n$. For each $i$, the inter-request times of $e_{i,j}$ are independent and have the distribution function $F_i(\cdot)$ for $1 \le j \le n$. Note that, when an $e_{i,j}$ is requested in the MTF list of $\mathcal{S}^{(n)}$, the items ahead of the $e_{i,j}$ are moved backward by $1/n$ and the $e_{i,j}$ is moved to the head of the list. In particular, $\mathcal{S}^{(1)}$ is associated with the original LRU buffer, and $\mathcal{S}^{(\infty)} \equiv \lim_{n \to \infty} \mathcal{S}^{(n)}$ is associated with the fluid limit of the original LRU buffer.

The sequence of systems that we consider is similar to that in Jelenković [11]. Similarly to $\mathcal{S}^{(n)}$, the $n$-th system, $\mathcal{S}'^{(n)}$, considered in [11] has $n\,N$ items, $e'_{i,j}$ for $1 \le i \le N$ and $1 \le j \le n$, of size $1/n$. Unlike $\mathcal{S}^{(n)}$, however, the distribution of the inter-request times of $e'_{i,j}$ may differ for each $j$. Jelenković [11] assumes that an item is requested at each time step and that the probability of the requested item being $e'_{i,j}$ is $q_{i,j} = \varphi(i - 1 + \frac{j}{n}) - \varphi(i - 1 + \frac{j-1}{n})$ for $1 \le i \le N$ and $1 \le j \le n$, where $\varphi(\cdot)$ is some non-decreasing function such that $\varphi(0) = 0$ and $\varphi(N) = 1$ (see Figure 1). If $\varphi(\cdot)$ is linear on $[i - 1, i]$ for $1 \le i \le N$, then $\mathcal{S}'^{(n)}$ is equivalent to $\mathcal{S}^{(n)}$, since $q_{i,j}$ is independent of $j$ in this case. Also, for large $N$ and smooth $\varphi(\cdot)$, $\mathcal{S}'^{(n)}$ and $\mathcal{S}^{(n)}$ are approximately equivalent. It turns out that $\mathcal{S}^{(\infty)}$ and $\mathcal{S}'^{(\infty)}$ are essentially equivalent, but we find that our derivation is simpler. This simplicity allows us to extend the analysis to 2Q.

## 2.2 Analysis of hit probability

We first analyze the stationary distribution of the position of a requested item in the MTF list of $\mathcal{S}^{(\infty)}$, which will be used to derive the hit probability for LRU in the fluid limit. For $1 \le i \le N$, let $f_i(\cdot)$ be the density function of the inter-request times, $R_i$, of $e_i$ and $G_i(\cdot)$ be the distribution function of the equilibrium distribution of $R_i$ (i.e., $f_i(t) = \frac{d}{dt}F_i(t)$ $\frac{d}{dt}G_i(t) = (1 - F_i(t))/\mathsf{E}\,[R_i]$). In particular, $f_i(t) = \lambda_i\,e^{-\lambda_i\,t}$ and $G_i(t) = 1 - e^{-\lambda_i\,t}$ if $e_i$ is requested according to a Poisson process with rate $\lambda_i$.

3

**Lemma 1** *Let $C_{i,j}^{(n)}$ be the stationary position of an $e_{i,j}$ in the MTF list of $\mathcal{S}^{(n)}$ when the $e_{i,j}$ is requested. As $n \to \infty$, $C_{i,j}^{(n)}$ converges in distribution to $C_i$ whose Laplace transform is given by*

$$\mathsf{E}\left[e^{-s\,C_i}\right] \quad = \quad \int_0^\infty e^{-s\sum_{k=1}^N G_k(t)}\, f_i(t)\, dt.$$

**Proof**: Let $t = 0$ be the stationary moment when an $e_{i,j}$ is requested in $\mathcal{S}^{(n)}$. Let $C_{i,j}^{(n)}(t)$ be the position of the $e_{i,j}$ in the MTF list of $\mathcal{S}^{(n)}$ at time $t$ given that the $e_{i,j}$ has not been requested by $t$. Since the time to the first request for the $e_{i,j}$ after time 0 has the density function $f_i(t)$, we have

$$\mathsf{E}\left[e^{-s\,C_{i,j}^{(n)}}\right] \quad = \quad \int_0^\infty \mathsf{E}\left[e^{-s\,C_{i,j}^{(n)}(t)}\right]\, f_i(t)\, dt.$$

Since $0 \le C_{i,j}^{(n)}(t) \le N$, the dominated convergence theorem can be used to show that

$$\lim_{n\to\infty} \mathsf{E}\left[e^{-s\,C_{i,j}^{(n)}}\right] \quad = \quad \int_0^\infty \lim_{n\to\infty}\mathsf{E}\left[e^{-s\,C_{i,j}^{(n)}(t)}\right]\, f_i(t)\, dt. \tag{1}$$

To derive $\mathsf{E}\left[e^{-s\,C_{i,j}^{(n)}(t)}\right]$, observe that $C_{i,j}^{(n)}(t)$ is incremented by $1/n$ when an $e_{k,\ell} \ne e_{i,j}$ is requested for the first time after time 0. Let $I_{k,\ell}(t)$ be the indicator function such that $I_{k,\ell}(t) = 1$ iff $e_{k,\ell}$ is requested at least once by time $t$. Let $U_{i,j}$ be the set of $(k,\ell)$ for $1 \le k \le N$ and $1 \le \ell \le n$, where $(k,\ell) \ne (i,j)$. Then

$$C_{i,j}^{(n)}(t) \quad = \quad \sum_{(k,\ell)\in U_{i,j}} \frac{1}{n} I_{k,\ell}(t). \tag{2}$$

Taking the Laplace transform of (2), we have

$$\mathsf{E}\left[e^{-s\,C_{i,j}^{(n)}(t)}\right] \quad = \quad \prod_{(k,\ell)\in U_{i,j}} \mathsf{E}\left[e^{-s\,I_{k,\ell}(t)/n}\right]. \tag{3}$$

Since the items are requested independently and the system under consideration is regenerative and at the steady state, the ASTA (Arrivals See Time Averages) principle [18] implies that the time to the first request for $e_{k,\ell} \ne e_{i,j}$ after time 0 has the distribution function $G_k(\cdot)$. Therefore,

$$\begin{aligned}
\mathsf{E}\left[e^{-s\,C_{i,j}^{(n)}(t)}\right] \quad &= \quad \prod_{(k,\ell)\in U_{i,j}} \mathsf{E}\left[e^{-s/n}\,G_k(t) + (1 - G_k(t))\right] \\
&= \quad \frac{\prod_{k=1}^N \left(e^{-s/n}\,G_k(t) + 1 - G_k(t)\right)^n}{G_i(t)\,e^{-s/n} + 1 - G_i(t)}.
\end{aligned} \tag{4}$$

Finally, we study the limit of (4) as $n \to \infty$. By Lemma 4 in Appendix A,

$$\lim_{n\to\infty}\left(G_k(t)\,e^{-s/n} + 1 - G_k(t)\right)^n \quad = \quad e^{-s\,G_k(t)} \tag{5}$$

for $1 \le k \le N$. Also, observe that

$$\lim_{n\to\infty}\left(G_i(t)\,e^{-s/n} + 1 - G_i(t)\right)^{-1} \quad = \quad 1. \tag{6}$$

Therefore, (4), (5), and (6) imply

$$\lim_{n\to\infty}\mathsf{E}\left[e^{-s\,C_{i,j}^{(n)}(t)}\right] \quad = \quad \prod_{k=1}^N e^{-s\,G_k(t)}. \tag{7}$$

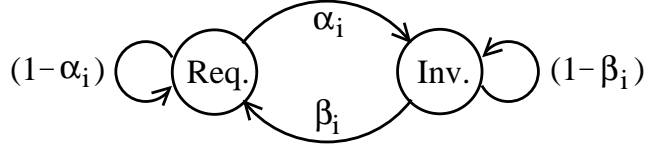Now, the lemma follows from (1) and (7). ∎

Figure 2: A discrete time Markov chain that determines the type (request (Req.) or invalidation (Inv.)) of an event based on the preceding type.

Observe that (7) implies that

$$\lim_{n \to \infty} C_{i,j}^{(n)}(t) \ = \ \sum_{k=1}^{N} G_k(t) \tag{8}$$

in probability for each $t$. Since the right hand side of (8) is a deterministic function of time, an item moves in the MTF list of $\mathcal{S}^{(\infty)}$ according to a deterministic process until the item is requested at a random time. Since an item moves toward the tail of the MTF list when other items are requested, the law of large numbers suggests that the movement is close to deterministic when there are many items. Also note that the right hand side of (8) does not depend on the particular item, so that every item moves at the same speed that depends only on the position of the item. This is because the effect of a single item is negligible when there are many items. These observations also suggest that $\mathcal{S}^{(\infty)}$ is a good approximation of $\mathcal{S}^{(1)}$ when $N$ is large.

Next, we will use Lemma 1 to derive the hit probability for `LRU` in the fluid limit. In the fluid limit, an item moves in the MTF list according to a deterministic process, so that there is a time $T$ such that the position of an item is at most $K$ iff the time since the last request of the item is at most $T$. Hence, a requested item is in an LRU buffer iff the time since the last request of the item is at most $T$. Formally,

**Theorem 1** *Let $p_{i,j}^{(n)}$ be the stationary probability that an $e_{i,j}$ is in the LRU buffer of $\mathcal{S}^{(n)}$ when the $e_{i,j}$ is requested. Let $T$ be the unique $t$ such that $\sum_{k=1}^{N} G_k(t) = K$. Then $p_{i,j}^{(n)} \to \int_0^T f_i(t) \, dt$ as $n \to \infty$.*

**Proof**: Recall that the hit probability of an $e_{i,j}$ in the LRU buffer of $\mathcal{S}^{(n)}$ coincides with the probability that, when the $e_{i,j}$ is requested, the position of the $e_{i,j}$ is at most $K$ in the MTF list of $\mathcal{S}^{(n)}$. Given that the $e_{i,j}$ was not requested, the position of the $e_{i,j}$ increases and reaches $K$ at time $T$ by (8). Therefore, $\Pr(C_{i,j}^{(n)} \le K) \to \int_0^T f_i(t) \, dt$ as $n \to \infty$, which proves the theorem. ∎

## 3 Analysis of LRU with invalidations

In this section, we study a buffer that is managed by `LRU` and where items are requested and invalidated. We refer to the buffer as the LRUI (LRU with Invalidations) buffer. When an item is requested, the LRUI buffer is updated in the same way as an LRU buffer. When an item in the LRUI buffer is invalidated, the item is removed from the LRUI buffer. When an item not in the LRUI buffer is invalidated, the LRUI buffer is not updated. In Section 3.1, we introduce a particular arrival process that generates correlated requests and invalidations. In Section 3.2, we analyze, in the fluid limit, the probability that a requested item is found in an LRUI buffer (the hit probability in an LRUI buffer) for the arrival process introduced in Section 3.1.

### 3.1 Arrival process of events

The arrival process for an $e_i$ generates events for the $e_i$, and an event for the $e_i$ is either a request or an invalidation for the $e_i$. We assume that the events for an $e_i$ are generated according to a Poisson process with rate $\lambda_i$ and that the probability that an event is a request or an invalidation depends on a past event. Specifically, when the preceding event for an $e_i$ is a request, the succeeding event for the $e_i$ is an invalidation with probability $\alpha_i$ and a request otherwise. When the preceding event for an $e_i$ is an invalidation, the succeeding event for the $e_i$ is a request with probability $\beta_i$ and an invalidation otherwise. Figure 2 shows a Markov chain that determines the type of an event for an $e_i$ based on the preceding event for the $e_i$. We assume that at least one of $\alpha_i$ and $\beta_i$ is nonzero.

The stationary probabilities that the event for an $e_i$ is a request and an invalidation are respectively given by

$$\pi_i^{\mathrm{R}} = \frac{\beta_i}{\alpha_i + \beta_i} \quad \text{and} \quad \pi_i^{\mathrm{I}} = \frac{\alpha_i}{\alpha_i + \beta_i}. \tag{9}$$

Therefore, for an $e_i$, requests are generated with average rate $\pi_i^{\mathrm{R}} \lambda_i$, and invalidations are generated with average rate $\pi_i^{\mathrm{I}} \lambda_i$. Note that the requests and the invalidations are correlated and do not follow Poisson processes unless $\alpha_i = 1 - \beta_i$. When $\alpha_i = 1 - \beta_i$, the requests and the invalidations of an $e_i$ are respectively generated according to independent Poisson processes with rate $(1 - \alpha_i) \lambda_i$ and $\alpha_i \lambda_i$.

## 3.2 Analysis of hit probability

We consider a sequence of systems, $\hat{\mathcal{S}}^{(n)}$ for $n = 1, 2, ...$, where each system is associated with an LRUI buffer of size $K$ and a corresponding list, which we refer to as a move-to-front-or-remove (MFR) list. When an $e_i$ in an MFR list is requested, the MFR list is updated in the same way as an MTF list. When an $e_i$ is invalidated in an MFR list, the $e_i$ is removed from the list. When an $e_i$ not in an MFR list is requested, the $e_i$ is inserted at the head of the MFR list. In $\hat{\mathcal{S}}^{(n)}$, events are generated independently for $n N$ items, $\hat{e}_{i,j}$ for $1 \leq i \leq N$ and $1 \leq j \leq n$, of size $1/n$. The events for an $\hat{e}_{i,j}$ are generated according to a Poisson process with rate $\lambda_i$, and the type of an event is determined by the Markov chain in Figure 2. Note that, when an invalidation is generated for an $\hat{e}_{i,j}$ in the MFR list of $\hat{\mathcal{S}}^{(n)}$, the items behind the $\hat{e}_{i,j}$ are moved forward by $1/n$ and the $\hat{e}_{i,j}$ is removed from the list. We will study the position of an item in an MFR list of $\hat{\mathcal{S}}^{(\infty)}$, which will be used to derive the hit probability in an LRUI buffer of $\hat{\mathcal{S}}^{(\infty)}$.

**Lemma 2** Let $\hat{C}_{i,j}^{(n)}$ be the stationary position of an $\hat{e}_{i,j}$ in the MFR list of $\hat{\mathcal{S}}^{(n)}$ when the $\hat{e}_{i,j}$ is requested given that the preceding event for the $\hat{e}_{i,j}$ is a request. As $n \to \infty$, $\hat{C}_{i,j}^{(n)}$ converges in distribution to $\hat{C}_i$ whose Laplace transform is

$$\mathsf{E}\left[e^{-s \hat{C}_i}\right] = \int_0^\infty e^{-s \sum_{k=1}^N \hat{H}_k(t)} f_i(t) \, dt, \tag{10}$$

where $f_i(t) = \lambda_i e^{-\lambda_i t}$ is the density function of the inter-event times of $e_i$ for $1 \leq i \leq N$, and

$$\hat{H}_i(t) = \frac{\beta_i}{\alpha_i + \beta_i} \left(1 - e^{-\lambda_i t}\right). \tag{11}$$

**Proof**: Let $t = 0$ be the stationary moment when a request for an $\hat{e}_{i,j}$ is generated in $\hat{\mathcal{S}}^{(n)}$. Let $\hat{C}_{i,j}^{(n)}(t)$ be the position of the $\hat{e}_{i,j}$ in the MFR list of $\hat{\mathcal{S}}^{(n)}$ at time $t$ given that an event for the $\hat{e}_{i,j}$ has not been generated by $t$. By the memoryless property of the exponential distribution, the time to the first event for the $\hat{e}_{i,j}$ after time 0 given that the first event is a request has the density function $f_i(\cdot)$. Thus,

$$\mathsf{E}\left[e^{-s \hat{C}_{i,j}^{(n)}}\right] = \int_0^\infty \mathsf{E}\left[e^{-s \hat{C}_{i,j}^{(n)}(t)}\right] f_i(t) \, dt.$$

Since $0 \leq \hat{C}_{i,j}^{(n)}(t) \leq N$, the dominated convergence theorem can be used to show that

$$\lim_{n \to \infty} \mathsf{E}\left[e^{-s \hat{C}_{i,j}^{(n)}}\right] = \int_0^\infty \lim_{n \to \infty} \mathsf{E}\left[e^{-s \hat{C}_{i,j}^{(n)}(t)}\right] f_i(t) \, dt. \tag{12}$$

Observe that $\hat{C}_{i,j}^{(n)}(t)$ is incremented by $1/n$ when an $\hat{e}_{k,\ell} \neq \hat{e}_{i,j}$ is requested for the first time after time 0, is decremented by $1/n$ if the $\hat{e}_{k,\ell}$ is invalidated after the first request, and is incremented again by $1/n$ if the $\hat{e}_{k,\ell}$ is requested after the invalidation. Therefore, an $\hat{e}_{k,\ell} \neq \hat{e}_{i,j}$ contributes to an increment of $\hat{C}_{i,j}^{(n)}(t)$ by $1/n$ iff the $\hat{e}_{k,\ell}$ is requested at least once between time 0 and $t$ and the last event for the $\hat{e}_{k,\ell}$ at time $t$ is a request.

Thus, in the same way as (4) is proven, we can show that

$$\mathsf{E}\left[e^{-s \hat{C}_{i,j}^{(n)}(t)}\right] = \frac{\prod_{k=1}^N \left(e^{-s/n} \hat{H}_k(t) + 1 - \hat{H}_k(t)\right)^n}{\hat{H}_i(t) e^{-s/n} + 1 - \hat{H}_i(t)}, \tag{13}$$

6

Figure 3: Continuous time Markov chains used to derive $\hat{H}_k(t)$.

where $\hat{H}_k(t)$ is the probability that an $\hat{e}_{k,\ell}$ is requested at least once between time 0 and $t$ and the last event for the $\hat{e}_{k,\ell}$ at time $t$ is a request. In the same way as (7) is proven, it can be shown that

$$\lim_{n\to\infty} \mathsf{E}\left[ e^{-s\,\hat{C}_{i,j}^{(n)}(t)} \right] \;=\; \prod_{k=1}^{N} e^{-s\,\hat{H}_k(t)}, \tag{14}$$

which together with (12) implies (10).

What remains to be shown is that $\hat{H}_k(t)$ is given by (11). To derive $\hat{H}_k(t)$, we condition on the type of the last event for an $\hat{e}_{k,\ell}$ before time 0. When the last event for the $\hat{e}_{k,\ell}$ before time 0 is a request, $\hat{H}_k(t)$ is the probability that the Markov chain in Figure 3(a) is in State "Req." at time $t$ given that the Markov chain is in the state denoted by a double circle at time 0. Similarly, when the last event is an invalidation, $\hat{H}_k(t)$ is the probability that the Markov chain in Figure 3(b) is in State "Req." at time $t$. By the ASTA principle and the memoryless property of the exponential distribution, the last event for an $\hat{e}_{k,\ell} \neq \hat{e}_{i,j}$ before time 0 is a request with probability $\pi_k^{\mathrm{R}}$ and an invalidation with probability $\pi_k^{\mathrm{I}}$, as shown in (9). Hence,

$$\hat{H}_k(t) \;=\; \pi_k^{\mathrm{R}}\, \mathbf{v}_1^{\mathrm{t}}\, e^{-\mathbf{Q}_k^{\mathrm{R}}\, t}\, \mathbf{v}_2 + \pi_k^{\mathrm{I}}\, \mathbf{v}_1^{\mathrm{t}}\, e^{-\mathbf{Q}_k^{\mathrm{I}}\, t}\, \mathbf{v}_2, \tag{15}$$

where $\mathbf{v}_i$ is a unit column vector with three elements such that the $i$-th element is 1, $\mathbf{v}_i^{\mathrm{t}}$ is a corresponding unit row vector, and $\mathbf{Q}_k^{\mathrm{R}}$ and $\mathbf{Q}_k^{\mathrm{I}}$ are, respectively, the generator matrices of the Markov chains in Figure 3(a) and Figure 3(b):

$$\mathbf{Q}_k^{\mathrm{R}} \;=\; \begin{pmatrix} -\lambda_k & (1-\alpha_k)\,\lambda_k & \alpha_k\,\lambda_k \\ 0 & -\alpha_i\,\lambda_k & \alpha_k\,\lambda_k \\ 0 & \beta_k\,\lambda_k & -\beta_k\,\lambda_k \end{pmatrix}$$

$$\mathbf{Q}_k^{\mathrm{I}} \;=\; \begin{pmatrix} -\lambda_k & \beta_k\,\lambda_k & (1-\beta_k)\,\lambda_k \\ 0 & -\alpha_k\,\lambda_k & \alpha_k\,\lambda_k \\ 0 & \beta_k\,\lambda_k & -\beta_k\,\lambda_k \end{pmatrix}$$

Lemma 5 in Appendix A implies that (15) is equivalent to (11), which completes the proof of the lemma. ∎

Since the right hand side of (14) is a deterministic function of time, an item moves in the MFR list of $\mathcal{S}^{(\infty)}$ according to a deterministic process until the item is requested or invalidated at a random time. Note that this deterministic process is insensitive to the correlations in the arrival processes, since $\hat{H}_i(\cdot)$ depends only on the *marginal* probability, $\pi_i^{\mathrm{R}}$, that the event for an $e_i$ is a request and on the distribution function, $F_i(\cdot)$, of the inter-event times of an $e_i$.

In contrast to this insensitivity, the hit probability in an LRUI buffer in general depends on the correlations of the arrival processes. For example, when $\alpha_i = 1$, a request for an $e_i$ and an invalidation for the $e_i$ alternate, and the hit probability of the $e_i$ is 0. We will now use Lemma 2 to derive the hit probability in an LRUI buffer in the fluid limit.

**Corollary 1** *Let $\hat{p}_{i,j}^{(n)}$ be the stationary probability that an $\hat{e}_{i,j}$ is in the LRUI buffer of $\hat{\mathcal{S}}^{(n)}$ when the $\hat{e}_{i,j}$ is requested. Let $\hat{T}$ be the unique $t$ such that $\sum_{k=1}^{N} \hat{H}_k(t) = K$. Then $\hat{p}_{i,j}^{(n)} \to (1-\alpha_i) \int_0^{\hat{T}} f_i(t)\,dt$ as $n \to \infty$.*

**Proof**: A requested $\hat{e}_{i,j}$ is in an LRUI buffer of $\hat{\mathcal{S}}^{(n)}$ iff the preceding event for the $\hat{e}_{i,j}$ is a request and the position of the $\hat{e}_{i,j}$ in the corresponding MFR list of $\hat{\mathcal{S}}^{(n)}$ is at most $K$. Observe that the event for the $\hat{e}_{i,j}$
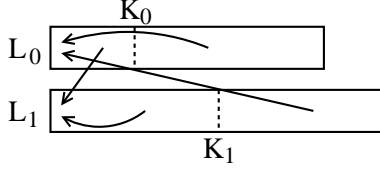
Figure 4: Rules of updating $L_0$ and $L_1$. When an $e_i$ is requested, the $e_i$ is moved to the head of $L_0$ or $L_1$ based on the position of the $e_i$ upon the request.

that precedes a request for the $\hat{e}_{i,j}$ is a request with probability $1 - \alpha_i$. Thus, $p_{i,j}^{(n)} = (1 - \alpha_i) \Pr(\hat{C}_{i,j} \leq K)$. Now, Lemma 2 can be used to show the corollary in the same way as Theorem 1. ∎

# 4    Analysis of 2Q

In this section, we analyze the hit probability for 2Q in the fluid limit. We study a simplest version of 2Q, but our analysis may be extended to other versions such as those introduced in [13]. 2Q divides a buffer of size $K$ into a part, $B_0$, of size $K_0 < K$ and a part, $B_1$, of size $K_1 = K - K_0 < K$. If a requested item, $e_i$, is neither in $B_0$ nor in $B_1$, the $e_i$ is added to $B_0$. If $B_0$ is full and the $e_i$ cannot be added, the least-recently requested item in $B_0$ is replaced with the $e_i$. If the $e_i$ is in $B_0$, the $e_i$ is removed from $B_0$ and added to $B_1$. If $B_1$ is full and the $e_i$ cannot be added, the least-recently requested item in $B_1$ is replaced with the $e_i$. If the $e_i$ is in $B_1$, the $e_i$ becomes the most-recently requested in $B_1$. Below, 2Q is referred to as 2Q$(\kappa)$ when $K_0/K = \kappa$, and a buffer managed by 2Q$(\kappa)$ is referred to as a 2Q$(\kappa)$ buffer. In Section 4.1, we introduce the fluid limit of a 2Q$(\kappa)$ buffer. In Section 4.2, we derive an analytical expression for the hit probability for 2Q$(\kappa)$ in the fluid limit. In Section 4.3, we show how we evaluate the analytical expression numerically. We assume that the requests for an $e_i$ are issued according to a Poisson process with rate $\lambda_i$ for $1 \leq i \leq N$, where $N$ is the number of items, each of which has size 1. Let $f_i(t) = \lambda_i e^{-\lambda_i t}$ be the density function of the inter-request times of $e_i$ for $1 \leq i \leq N$.

## 4.1    Fluid limit

To analyze the hit probability for 2Q$(\kappa)$, we consider a corresponding pair of MFR lists, $L_0$ and $L_1$, where each $e_i$ is either in $L_0$ or in $L_1$. When a request of an $e_i$ is generated in a 2Q$(\kappa)$ buffer, $L_0$ and $L_1$ are updated as follows (see also Figure 4). If the $e_i$ is in $L_0$ and its position is at most $K_0 = \kappa K$, the $e_i$ is removed from $L_0$ and inserted at the head of $L_1$. In other words, an invalidation of the $e_i$ is generated in $L_0$, and a request of the $e_i$ is generated in $L_1$. If the $e_i$ is in $L_0$ at a position greater than $K_0$, a request of the $e_i$ is generated in $L_0$, and an invalidation of the $e_i$ is generated in $L_1$. Note that $L_1$ is not updated, since the invalidated $e_i$ is not in $L_1$. If the $e_i$ is in $L_1$ at a position at most $K_1 = (1 - \kappa)K$, a request of the $e_i$ is generated in $L_1$, and an invalidation of the $e_i$ is generated in $L_0$ (again, $L_0$ is not updated). If the $e_i$ is in $L_1$ at a position greater than $K_1$, a request of the $e_i$ is generated in $L_0$, and an invalidation of the $e_i$ is generated in $L_1$. Observe that a requested $e_i$ is in a 2Q$(\kappa)$ buffer iff, in the corresponding pair of $L_0$ and $L_1$, the $e_i$ is either in $L_0$ at a position at most $K_0$ or in $L_1$ at a position at most $K_1$. In either case, an invalidation of the $e_i$ is generated in $L_0$ and a request of the $e_i$ is generated in $L_1$ upon the request of the $e_i$ in the 2Q$(\kappa)$ buffer.

We consider a sequence of systems, $\bar{S}^{(n)}$ for $n = 1, 2, ...$, where each system is associated with a 2Q$(\kappa)$ buffer of size $K$ and a corresponding pair of MFR lists, $L_0^{(n)}$ and $L_1^{(n)}$. In a 2Q$(\kappa)$ buffer of $\bar{S}^{(n)}$, requests are generated independently for $n N$ items, $\bar{e}_{i,j}$ for $1 \leq i \leq N$ and $1 \leq j \leq n$, and the items have size $1/n$. The requests for an $\bar{e}_{i,j}$ are generated according to a Poisson process with rate $\lambda_i$.

## 4.2    Analysis of hit probability

Now, we analyze the hit probability for 2Q$(\kappa)$ in the fluid limit, $\bar{S}^{(\infty)}$. We will find that the position of an $e_i$ in $L_0^{(\infty)}$ increases according to a deterministic process, $\bar{C}(\cdot)$, (and the position of an $e_i$ in $L_1^{(\infty)}$ increases according to a deterministic process, $\bar{C}'(\cdot)$) until the $e_i$ is requested or invalidated at a random time. The $\bar{C}(\cdot)$ and $\bar{C}'(\cdot)$ depend on the probability, $a_i$, that an $e_i$ is requested in the 2Q$(\kappa)$ buffer before the position of the $e_i$ reaches $K_0$ in $L_0^{(\infty)}$ and on the probability, $b_i$, that an $e_i$ is *not* requested in the 2Q$(\kappa)$ buffer before the

position of the $e_i$ reaches $K_1$ in $L_1^{(\infty)}$. In turn, the $a_i$ and $b_i$ for $1 \le i \le N$ depend on $\bar{C}(\cdot)$ and $\bar{C}'(\cdot)$. Hence, $\bar{C}(\cdot)$, $\bar{C}'(\cdot)$, $a_i$ for $1 \le i \le N$, and $b_i$ for $1 \le i \le N$ will be derived by solving a system of equations.

**Theorem 2** *Let $\bar{p}_{i,j}^{(n)}$ be the stationary probability that an $\bar{e}_{i,j}$ is in a $2Q(\kappa)$ buffer of $\bar{S}^{(n)}$ when the $\bar{e}_{i,j}$ is requested. Let $K_0 = \kappa K$ and $K_1 = (1-\kappa)K$. Then*

$$\lim_{n \to \infty} \bar{p}_{i,j}^{(n)} \;=\; \frac{a_i}{a_i + b_i}, \tag{16}$$

*where $a_i$ and $b_i$ for $1 \le i \le N$ are the unique constants that satisfy the following system of equations:*

$$\bar{C}(t) = \sum_{k=1}^{N} \frac{b_k}{a_k + b_k} \left( 1 - e^{-\lambda_k t} \right) \tag{17}$$

$$\bar{C}'(t) = \sum_{k=1}^{N} \frac{a_k}{a_k + b_k} \left( 1 - e^{-\lambda_k t} \right) \tag{18}$$

$$a_i = \int_{t=0}^{\infty} \Pr\left( \bar{C}(t) \le K_0 \right) f_i(t)\, dt \quad for \;\; 1 \le i \le N \tag{19}$$

$$b_i = \int_{t=0}^{\infty} \Pr\left( \bar{C}'(t) > K_1 \right) f_i(t)\, dt \quad for \;\; 1 \le i \le N. \tag{20}$$

**Proof**: We will analyze $\bar{p}_{i,j}^{(n)}$ by studying a corresponding pair of MFR lists, $L_0^{(n)}$ and $L_1^{(n)}$. We start by studying the arrival process of events (requests and invalidations) in $L_0^{(n)}$. Since an event for an $\bar{e}_{i,j}$ is generated in $L_0^{(n)}$ when a request for the $\bar{e}_{i,j}$ is generated in the corresponding $2Q(\kappa)$ buffer, the events for an $\bar{e}_{i,j}$ are generated in $L_0^{(n)}$ according to a Poisson process with rate $\lambda_i$. Given that the preceding event for an $\bar{e}_{i,j}$ is a request in $L_0^{(n)}$, the succeeding event for the $\bar{e}_{i,j}$ is an invalidation in $L_0^{(n)}$ iff the position of the $\bar{e}_{i,j}$ in $L_0^{(n)}$ is at most $K_0$ when the succeeding event is generated. Similarly, given that the preceding event for an $\bar{e}_{i,j}$ is an invalidation in $L_0^{(n)}$, the succeeding event for the $\bar{e}_{i,j}$ is a request in $L_0^{(n)}$ iff the position of the $\bar{e}_{i,j}$ in $L_1^{(n)}$ is greater than $K_1$ when the succeeding event is generated.

By the memoryless property of the exponential distribution, the system regenerates respectively when a request for an $\bar{e}_{i,j}$ is generated in $L_0^{(n)}$ and when an invalidation for an $\bar{e}_{i,j}$ is generated in $L_0^{(n)}$. Let $a_{i,j}^{(n)}$ be the probability that the position of an $\bar{e}_{i,j}$ is at most $K_0$ in $L_0^{(n)}$ when a request for the $\bar{e}_{i,j}$ is generated in $L_0^{(n)}$ given that the preceding event for the $\bar{e}_{i,j}$ is a request in $L_0^{(n)}$. Similarly, let $b_{i,j}^{(n)}$ be the probability that the position of an $\bar{e}_{i,j}$ is greater than $K_1$ in $L_1^{(n)}$ when the $\bar{e}_{i,j}$ is requested in $L_0^{(n)}$ given that the preceding event for the $\bar{e}_{i,j}$ is an invalidation in $L_0^{(n)}$.

Then the events for an $\bar{e}_{i,j}$ are generated in $L_0^{(n)}$ according to the arrival process introduced in Section 3.1, where $\alpha_i = a_{i,j}^{(n)}$ and $\beta_i = b_{i,j}^{(n)}$. Recall that a requested $\bar{e}_{i,j}$ is in a $2Q(\kappa)$ buffer iff an invalidation is generated for the $\bar{e}_{i,j}$ in $L_0^{(n)}$ upon the request of the $e_{i,j}$ in the $2Q(\kappa)$ buffer. Thus, $\bar{p}_{i,j}^{(n)}$ is equivalent to the marginal probability that an event for an $\bar{e}_{i,j}$ in $L_0^{(n)}$ is an invalidation. Therefore, (9) implies that

$$\bar{p}_{i,j}^{(n)} \;=\; \frac{a_{i,j}^{(n)}}{a_{i,j}^{(n)} + b_{i,j}^{(n)}}. \tag{21}$$

By the memoryless property of the exponential distribution, the time to the first event for an $\bar{e}_{i,j}$ after a request is generated for the $\bar{e}_{i,j}$ in $L_0^{(n)}$, given that the first event for the $\bar{e}_{i,j}$ is a request in $L_0^{(n)}$, has the density function $f_i(\cdot)$. Hence,

$$a_{i,j}^{(n)} \;=\; \int_{t=0}^{\infty} \Pr\left( \bar{C}_{i,j}^{(n)}(t) \le K_0 \right) f_i(t)\, dt,$$

where $\bar{C}_{i,j}^{(n)}(t)$ is the position of the $\bar{e}_{i,j}$ in $L_0^{(n)}$ at time $t$ given that a request for the $\bar{e}_{i,j}$ is generated in $L_0^{(n)}$ at time 0 and that no event is generated for the $\bar{e}_{i,j}$ between time 0 and $t$. Similarly,

$$b_{i,j}^{(n)} \;=\; \int_{t=0}^{\infty} \Pr\left( \bar{C}'_{i,j}^{(n)}(t) > K_1 \right) f_i(t)\, dt,$$

9

where $\bar{C'}_{i,j}^{(n)}(t)$ is the position of the $\bar{e}_{i,j}$ in $L_1^{(n)}$ at time $t$ given that a request for the $\bar{e}_{i,j}$ is generated in $L_1$ at time $0$ and that no event is generated for the $\bar{e}_{i,j}$ between time $0$ and $t$.

Next, we study $\bar{C}_{i,j}^{(n)}(t)$ and $\bar{C'}_{i,j}^{(n)}(t)$. In the same way as (13) is proven, it can be shown that

$$\mathsf{E}\left[e^{-s\,\bar{C}_{i,j}^{(n)}(t)}\right] \;\;=\;\; \frac{\prod_{k=1}^{N}\left(e^{-s/n}\,\bar{H}_k^{(n)}(t)+1-\bar{H}_k^{(n)}(t)\right)^n}{\bar{H}_i^{(n)}(t)\,e^{-s/n}+1-\bar{H}_i^{(n)}(t)}$$

where $\bar{H}_i^{(n)}(t)$ is the probability that an $\bar{e}_{i,j}$ is requested in $L_0^{(n)}$ at least once between time $0$ and $t$ and the last event for the $\bar{e}_{i,j}$ before time $t$ is a request in $L_0^{(n)}$. Thus, by (11),

$$\bar{H}_i^{(n)}(t) \;\;=\;\; \frac{b_i^{(n)}}{a_i^{(n)}+b_i^{(n)}}\left(1-e^{-\lambda_i t}\right).$$

To derive a similar expression for $\bar{C'}_{i,j}^{(n)}(t)$, we need to study the arrival process of the events in $L_1^{(n)}$. Recall that an event for an $\bar{e}_{i,j}$ is generated in $L_1^{(n)}$ at the same moment as an event for the $\bar{e}_{i,j}$ is generated in $L_0^{(n)}$, and that the event for the $\bar{e}_{i,j}$ is a request in $L_1^{(n)}$ iff the event for the $\bar{e}_{i,j}$ is an invalidation in $L_0^{(n)}$. Thus, the events for an $\bar{e}_{i,j}$ in $L_1^{(n)}$ are also generated according to the arrival process introduced in Section 3.1, but now $\alpha_i = b_{i,j}^{(n)}$ and $\beta_i = a_{i,j}^{(n)}$. Therefore, the Laplace transform of $\bar{C'}_{i,j}^{(n)}(t)$ is given by

$$\mathsf{E}\left[e^{-s\,\bar{C'}_{i,j}^{(n)}(t)}\right] \;\;=\;\; \frac{\prod_{k=1}^{N}\left(e^{-s/n}\,\bar{H'}_k^{(n)}(t)+1-\bar{H'}_k^{(n)}(t)\right)^n}{\bar{H'}_i^{(n)}(t)\,e^{-s/n}+1-\bar{H'}_i^{(n)}(t)},$$

where

$$\bar{H'}_i^{(n)}(t) \;\;=\;\; \frac{a_i^{(n)}}{a_i^{(n)}+b_i^{(n)}}\left(1-e^{-\lambda_i t}\right).$$

Finally, we study $\bar{C}_{i,j}^{(n)}(t)$ and $\bar{C'}_{i,j}^{(n)}(t)$ in the limit of $n \to \infty$. In the same way as (14) is proven, it can be shown that $\bar{C}_{i,j}^{(n)}(t)$ and $\bar{C'}_{i,j}^{(n)}(t)$, respectively, converge to deterministic processes as $n \to \infty$. Formally, for any $\epsilon > 0$, there exists $M$ such that, for all $n \geq M$,

$$\left|\mathsf{E}\left[e^{-s\,\bar{C}_{i,j}^{(n)}(t)}\right] - e^{-s\sum_{k=1}^{N}\bar{H}_k^{(n)}(t)}\right| < \epsilon$$

$$\left|\mathsf{E}\left[e^{-s\,\bar{C'}_{i,j}^{(n)}(t)}\right] - e^{-s\sum_{k=1}^{N}\bar{H'}_k^{(n)}(t)}\right| < \epsilon$$

for any $t$. Note that if $|\phi(a_n) - \psi(b_n)| \to 0$ as $n \to \infty$ and there exists a *unique* pair $(a,b)$ such that $\phi(a) = \psi(b)$, then $(a_n, b_n) \to (a,b)$ as $n \to \infty$. As we will see in Lemma 3, there exist *unique* constants, $a_i$ and $b_i$ for $1 \leq i \leq N$, and *unique* deterministic processes, $\bar{C}_i(\cdot)$ and $\bar{C'}_i(\cdot)$ for $1 \leq i \leq N$, that satisfy (17)-(20). Therefore, as $n \to \infty$, $a_{i,j}^{(n)} \to a_i$, $b_{i,j}^{(n)} \to b_i$, $\bar{C}_{i,j}^{(n)}(t) \to \bar{C}(t)$, and $\bar{C'}_{i,j}^{(n)}(t) \to \bar{C'}(t)$ for $1 \leq j \leq n$ and $1 \leq i \leq N$, where $a_i$, $b_i$, $\mathsf{E}\left[e^{-s\,\bar{C}(t)}\right]$, and $\mathsf{E}\left[e^{-s\,\bar{C'}(t)}\right]$ are defined by the unique solutions of (17)-(20). Now, the theorem follows from (21). ∎

**Lemma 3** *The system of equations (17)-(20) has a unique solution, $(a_i, b_i)$ for $1 \leq i \leq N$.*

**Proof**: First, consider the case where $K_1 \geq N$. In this case, it is expected that $\bar{p}_i \equiv \lim_{n \to \infty} \bar{p}_{i,j}^{(n)} = 1$ for all $i$ since all of the items are in $K_1$, so that the unique solution is that $a_i = 1$ and $b_i = 0$ for all $i$. In fact, (20) implies that $b_i = 0$ for all $i$, which in turn implies $\bar{C}(t) \equiv 0$ by (17). Then (19) implies $a_i = 1$ for all $i$. Therefore, the system of equations (17)-(20) have the unique solution.

Below, we assume that $K_1 < N$. We will first show that $\bar{C}(\cdot)$ and $\bar{C'}(\cdot)$ must be increasing functions. By (19) and (20), this will allow us to express the $2N$ variables, $(a_i, b_i)$ for $1 \leq i \leq N$, by two variables, $(T_0, T_1)$, such that $a_i = F_i(T_0)$ and $b_i = 1 - F_i(T_1)$, where $F_i(t) = 1 - e^{-\lambda_i t}$. Then our proof will be reduced to show

the existence of the unique pair $(T_0, T_1)$ so that $a_i = F_i(T_0)$ and $b_i = 1 - F_i(T_1)$ for $1 \le i \le N$ satisfy the system of equations.

Observe in (17) that $\bar{C}(\cdot)$ is an increasing function unless $b_i = 0$ for all $i$. Suppose that $b_i = 0$ for all $i$. Then we have seen above that $a_i = 1$ for all $i$. Then (18) implies $\bar{C}'(t) = \sum_{k=1}^{N}(1 - e^{-\lambda_k t})$, which converges to $N$ as $t \to \infty$. Therefore, (20) implies $b_i > 0$ for all $i$ when $K_1 < N$. This contradicts our assumption that $b_i = 0$ for all $i$. Therefore, $b_i > 0$ for at least one $i$, and $\bar{C}(\cdot)$ is an increasing function.

Observe in (18) that $\bar{C}'(\cdot)$ is an increasing function unless $a_i = 0$ for all $i$. Suppose that $a_i = 0$ for all $i$. Then (18) implies $\bar{C}'(t) \equiv 0$, which in turn implies $b_i = 0$ for all $i$ by (20). However, we have seen above that $b_i > 0$ for at least one $i$. Therefore, it must be that $a_i > 0$ for at least one $i$ and that $\bar{C}'(\cdot)$ is an increasing function.

Since $\bar{C}(\cdot)$ and $\bar{C}'(\cdot)$ are increasing functions, (19) and (20) imply that $a_i$ and $b_i$ can be expressed as

$$a_i = F_i(T_0) \quad \text{and} \quad b_i = 1 - F_i(T_1), \tag{22}$$

where $F_i(t) = 1 - e^{-\lambda_i t}$ and

$$T_0 = \begin{cases} \bar{C}^{-1}(K_0) & \text{if } \bar{C}(\infty) > K_0 \\ \infty & \text{otherwise} \end{cases} \tag{23}$$

$$T_1 = \begin{cases} \bar{C}'^{-1}(K_1) & \text{if } \bar{C}'(\infty) > K_1 \\ \infty & \text{otherwise}, \end{cases} \tag{24}$$

where $\bar{C}^{-1}(\cdot)$ and $\bar{C}'^{-1}(\cdot)$ are inverse functions of $\bar{C}(\cdot)$ and $\bar{C}'(\cdot)$. Roughly speaking, $T_0$ is the time it takes for an item to reach $K_0$ in $L_0$ from the head of $L_0$ given that no event is generated for the item. Similarly, $T_1$ is the time to reach $K_1$ in $L_1$ from the head of $L_1$ under the same conditions.

Observe that there exists a unique solution, $(a_i, b_i)$ for $1 \le i \le N$, that satisfies (17)-(20) iff there exists a unique pair, $(T_0, T_1)$, that satisfies (17)-(18) and (22)-(24). Therefore, it suffices to prove the existence of the unique pair $(T_0, T_1)$. It will turn out that $T_1$ is always finite, so that we will prove the existence of the unique pair for two cases, where $T_0$ is finite and where $T_0$ is infinite.

Notice that $T_1 < \infty$ follows immediately from (20), since $b_i > 0$ for at least one $i$, as we have seen above. When $T_1 < \infty$, we have by (18), (22), and (24) that

$$K_1 = \sum_{k=1}^{N} \frac{F_k(T_0)\, F_k(T_1)}{F_k(T_0) + 1 - F_k(T_1)}. \tag{25}$$

We now discuss the existence of the unique pair $(T_0, T_1)$ for the case where $T_0 < \infty$. In this case, the following relation must hold by (17), (22), and (23):

$$K_0 = \sum_{k=1}^{N} \frac{(1 - F_k(T_1))\, F_k(T_0)}{F_k(T_0) + 1 - F_k(T_1)}. \tag{26}$$

Let $\xi_0(T_0, T_1)$ be the right hand side of (26) and $\xi_1(T_0, T_1)$ be the right hand side of (25). Observe that $\xi_0(T_0, T_1)$ is increasing in $T_0$ and decreasing in $T_1$, and that $\xi_1(T_0, T_1)$ is increasing in $T_0$ and $T_1$ (see Figure 5 for contour curves of a $\xi_0(T_0, T_1)$ and a $\xi_1(T_0, T_1)$). Therefore, a unique pair, $T_0 < \infty$ and $T_1 < \infty$, that satisfies (25) and (26) exists iff $T_1^{(0)} > T_1^{(1)}$, where $T_1^{(0)}$ and $T_1^{(1)}$ are, respectively, unique $t_1^{(0)}$ and $t_1^{(1)}$ such that $\xi_0(\infty, t_1^{(0)}) = K_0$ and $\xi_1(\infty, t_1^{(1)}) = K_1$. In summary, if $T_1^{(0)} > T_1^{(1)}$, then there exists a unique pair, $T_0 < \infty$ and $T_1 < \infty$, that satisfies the system of equations, (17)-(18) and (22)-(24).

Finally, we will prove that if $T_1^{(0)} \le T_1^{(1)}$, then there exists a unique pair, $T_0 = \infty$ and $T_1 < \infty$, that satisfies the system of equations, (17)-(18) and (22)-(24). Note that if (26) is satisfied by a $T_0$ and a $T_1$, the $T_0$ must be finite, since

$$\bar{C}(\infty) = \sum_{k=1}^{N} \frac{1 - F_k(T_1)}{F_k(T_0) + 1 - F_k(T_1)} > K_0$$

follows from (17), (22), and (26), and $\bar{C}(\infty) > K_0$ iff $T_0 < \infty$. By the contrapositive, $\bar{C}(\infty) \le K_0$ iff $T_1^{(0)} \le T_1^{(1)}$. When $\bar{C}(\infty) \le K_0$, we have $T_0 = \infty$, which implies that $a_i = 1$ for all $i$ by (19). Thus, (17)-(18) and (22) imply

$$\bar{C}(t) = \sum_{k=1}^{N} \frac{(1 - F_k(T_1))F_k(t)}{2 - F_k(T_1)} \text{ and } \bar{C}'(t) = \sum_{k=1}^{N} \frac{F_k(t)}{2 - F_k(T_1)},$$
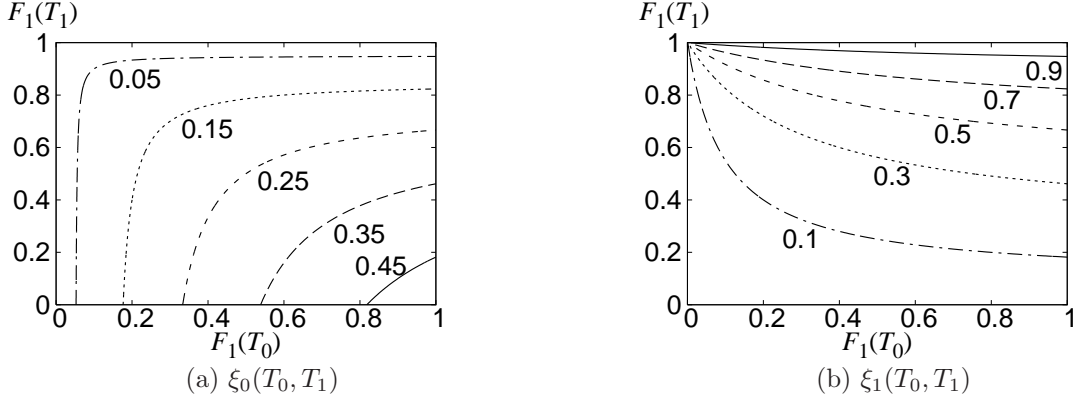
11

Figure 5: Contour curves of a $\xi_0(T_0, T_1)$ and a $\xi_1(T_0, T_1)$, where $N = 1$ and $\lambda_1 = 1$. Note that the horizontal and vertical axes are $F_1(T_0)$ and $F_1(T_1)$, which are increasing in $T_0$ and $T_1$, respectively.

where $T_1$ is a unique solution of the following equation:

$$K_1 = \sum_{k=1}^{N} \frac{F_k(T_1)}{2 - F_k(T_1)}.$$

Therefore, $T_0 = \infty$ and $T_1 = T_1^{(1)}$ is the unique pair that satisfies (17)-(18) and (22)-(24). ■

## 4.3   Numerical evaluation of hit probability

The arguments in the proof of Lemma 3 lead to the following algorithm for calculating $\bar{p}_i \equiv \lim_{n \to \infty} \bar{p}_{i,j}^{(n)}$. If $K_1 \geq N$, then $\bar{p}_i = 1$ for $1 \leq i \leq N$. If $K_1 < N$, find a unique pair $(T_1^{(0)}, T_1^{(1)})$ that satisfy

$$\xi_0(T_1^{(0)}) \equiv \sum_{k=1}^{N} \frac{1 - F_k(T_1^{(0)})}{2 - F_k(T_1^{(0)})} = K_0 \tag{27}$$

$$\xi_1(T_1^{(1)}) \equiv \sum_{k=1}^{N} \frac{F_k(T_1^{(1)})}{2 - F_k(T_1^{(1)})} = K_1. \tag{28}$$

Since $\xi_0(T_1^{(0)})$ is decreasing in $T_1^{(0)}$ and $\xi_1(T_1^{(1)})$ is increasing in $T_1^{(1)}$, the $T_1^{(0)}$ and the $T_1^{(1)}$ may be found by binary search. If $T_1^{(0)} \leq T_1^{(1)}$, then $\bar{p}_i = (2 - F_i(T_1^{(1)}))^{-1}$ for $1 \leq i \leq N$. If $T_1^{(0)} > T_1^{(1)}$, then

$$\bar{p}_i = \frac{F_i(T_0)}{F_i(T_0) + 1 - F_i(T_1)} \tag{29}$$

for $1 \leq i \leq N$, where $(T_0, T_1)$ is a unique pair that satisfies (25) and (26). The left hand side of (26) is increasing in $T_0$ and decreasing in $T_1$, and the left hand side of (25) is increasing in $T_0$ and $T_1$, respectively. Hence the $(T_0, T_1)$ may be found for example by Newton's method.

Finally, we remark that a conservation law holds for $\bar{p}_i = \lim_{n \to \infty} \bar{p}_{i,j}^{(n)}$. Specifically, $\sum_{i=1}^{N} \bar{p}_i = K$ for any $\kappa$ as long as the 2Q($\kappa$) buffer is fully utilized in $\bar{\mathcal{S}}^{(\infty)}$. The conservation law implies that a $\bar{p}_i$ cannot be increased without decreasing another $\bar{p}_j$. Note, however, that this does not mean that the overall hit probability, $\sum_{k=1}^{N} r_k \bar{p}_k$, is insensitive to $\kappa$, where $r_k = \lambda_k / \sum_{j=1}^{N} \lambda_j$. The overall hit probability can be made higher by increasing $p_i$ having a large $\lambda_i$ and decreasing $p_i$ having a small $\lambda_i$. Notice that the 2Q($\kappa$) buffer is underutilized in $\bar{\mathcal{S}}^{(\infty)}$ iff $\bar{C}(\infty) \leq K_0$, since $\bar{C}(\infty) \leq K_0$ suggests that no item reaches $K_0$ in a finite time. When the 2Q($\kappa$) buffer is underutilized, $\sum_{i=1}^{N} \bar{p}_i < K$. Formally,

**Corollary 2 (Conservation law for 2Q)** *If $K < N$ and $T_1^{(0)} > T_1^{(1)}$, then $\sum_{k=1}^{N} \bar{p}_k = K$. If $K < N$ and $T_1^{(0)} \leq T_1^{(1)}$, or if $K \geq N$, then $\sum_{k=1}^{N} \bar{p}_k \leq K$.*

**Proof**: Consider the case where $K < N$ and $T_1^{(0)} > T_1^{(1)}$. Summing both sides of (26) and (25), we obtain $K = \sum_{k=1}^{N} \bar{p}_k$ by (29). When $K < N$ and $T_1^{(0)} \leq T_1^{(1)}$, we obtain $K = \xi_0(T_1^{(0)}) + \xi_1(T_1^{(1)})$ by summing both sides of (27) and (28). Since $\xi_0(\cdot)$ is a decreasing function and $T_1^{(0)} \leq T_1^{(1)}$, $K \geq \xi_0(T_1^{(1)}) + \xi_1(T_1^{(1)})$, which implies $K \geq \sum_{k=1}^{N} \bar{p}_k$ by $\bar{p}_k = (2 - F_k(T_1^{(1)}))^{-1}$. When $K \geq N$, $\sum_{k=1}^{N} \bar{p}_k = N \leq K$ follows immediately from $\bar{p}_k = 1$. $\blacksquare$

## 5 Results

In this section, we study the fundamental properties of $2\mathbb{Q}(\kappa)$. In Section 5.1, we start by a validation of approximating the hit probabilities for $2\mathbb{Q}(\kappa)$ and LRU by those in the fluid limit. In Section 5.2, we will study the hit probability for $2\mathbb{Q}(\kappa)$, comparing it against that for LRU and against a theoretical upper bound, which is calculated as the hit probability when the $K$ items having the largest $\lambda_i$'s are always stored in the buffer (optimal static arrangement). We refer to the upper bound as the hit probability for OPT. In particular, we will find that

- the relative error in approximating the hit probabilities for $2\mathbb{Q}(\kappa)$ and LRU by those in the fluid limit is within 1% for $N > 1000$;
- the (stationary) hit probability for $2\mathbb{Q}(\kappa)$ can in general be made higher than that for LRU;
- the (stationary) hit probability for $2\mathbb{Q}(\kappa)$ is in general maximized when $K_0 = \kappa K = 1$;
- when $K_0 = 1$, the (stationary) hit probability for $2\mathbb{Q}(\kappa)$ is close to that for OPT;
- when $\kappa$ is smaller, however, a longer time is required to reach the stationary hit probability, so that a larger $\kappa$ may be preferred to a smaller $\kappa$.

### 5.1 Validation

We evaluate the accuracy of approximating the overall hit probability, $H = \sum_{i=1}^{N} r_i p_i$, for $2\mathbb{Q}(\kappa)$ and for LRU by those in the fluid limit, where $p_i$ is the hit probability of an $e_i$ and $r_i = \lambda_i / \sum_{j=1}^{N} \lambda_j$ is the stationary fraction of the requests for the $e_i$. Let $H_{\text{flu}}$ be the overall hit probability in the fluid limit and $H_{\text{sim}}$ be the overall hit probability estimated by simulation. The relative error (%) in $H_{\text{flu}}$ is defined by $100 \, |H_{\text{flu}} - H_{\text{sim}}| / H_{\text{sim}}$. Below, we omit the discussion on LRU, but the relative error in approximating the overall hit probability for LRU is smaller than that for $2\mathbb{Q}(\kappa)$ by a factor of 2 to 200 for all of the cases studied. Also, although we show only a limited set of plots, our discussion is based on experiments with a wider range of parameter sets.

For each evaluation, simulation is run at least 20 times. In each run, $10^7$ requests are generated after a warm-up period of $10^5$ requests. When the 20 runs do not suffice to provide the confidence that the estimated value is within 0.0001 of the true value with probability at least 0.95, the simulation is repeated until this accuracy is achieved (see [15]).

First, we study the case where the distribution of $\lambda_i$ for $1 \leq i \leq n$ follows Zipf's law (Breslau *et al.* find that the distributions of the rates that webpages are requested follow Zipf's law approximately [2]). Specifically, we choose $\lambda_i = 1/i$ for $1 \leq i \leq N$. Figure 6(a) shows $H_{\text{flu}}$ and $H_{\text{sim}}$ against varying values of $\kappa$. The number of items, $N$, is as labeled in each row. The solid lines represent $H_{\text{flu}}$ when $K = N/4$, the dashed lines when $K = 3N/8$, and the dotted lines when $K = N/2$. The '$*$' marks represent $H_{\text{sim}}$. The value of $K$ for each $H_{\text{sim}}$ is understood by the value of $K$ of the nearest line. Observe that every '$*$' is on or very close to the corresponding line. Thus, $H_{\text{flu}}$ closely agrees with $H_{\text{sim}}$ for a wide range of conditions.

To take a closer look, Figure 6(b) shows the relative error in $H_{\text{flu}}$ under the same conditions as Column (a). We find that the relative error is in general smaller for a larger $N$ and that the relative error becomes less than 0.1% for $N = 2^{10}$. This makes intuitive sense, since the system approaches the fluid limit as $N \to \infty$. However, we find that the relative error is surprisingly small even for a small $N$, in particular within 2% for $N = 2^6$.

The relative error in $H_{\text{flu}}$ is also sensitive to $\kappa$. In general, we find that the relative error is an increasing function of $\kappa$ and that $H_{\text{flu}} > H_{\text{sim}}$ for a large $\kappa$. This may be explained by examining the utilization of $B_0$. Recall that $B_0$ may have less than $K_0$ items, since a requested item in $B_0$ moves to $B_1$. However, under the conditions of Figure 6, $B_0$ in the fluid limit is fully utilized, which in turn makes $H_{\text{flu}}$ higher than $H_{\text{sim}}$.

Next, we study the effect of the distribution of $\lambda_i$ on the relative error in $H_{\text{flu}}$. Above, we have assumed that the distribution of $\lambda_i$ follows Zipf's law. We now consider the case where $\lambda_i$ is geometrically distributed and where $\lambda_i$ is linearly distributed. Specifically, for $1 \leq i \leq N$, we choose $\lambda_i = 1/N^{\frac{i-1}{N-1}}$ when $\lambda_i$ is geometrically
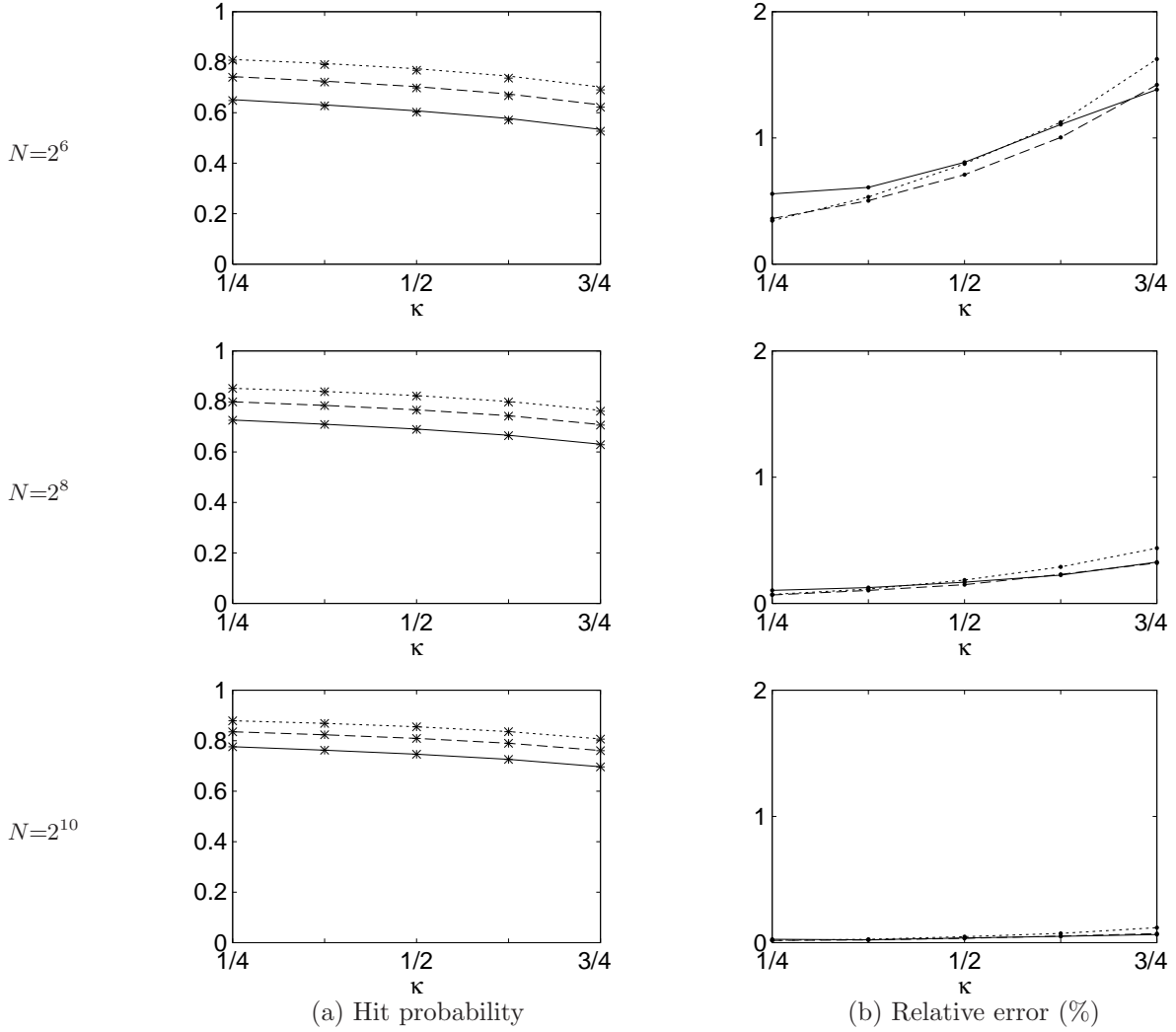
Figure 6: Accuracy of approximating the hit probability for $2\mathsf{Q}(\kappa)$ by $H_{\mathrm{flu}}$ when $\lambda_i$ follows Zipf's law. For each row, $N$ is varied as labeled. Column (a) shows $H_{\mathrm{flu}}$ by lines and $H_{\mathrm{sim}}$ by '$*$' marks, and Column (b) shows the relative error (%) in $H_{\mathrm{flu}}$, where $K = N/4$ for solid lines, $K = 3N/8$ for dashed lines, and $K = N/2$ for dotted lines.

distributed, and $\lambda_i = (N+1-i)/N$ when $\lambda_i$ is linearly distributed, so that $\lambda_1 = 1$ and $\lambda_N = \frac{1}{N}$ stay unchanged for all distributions under consideration. Figure 7 shows the values of $\lambda_i$ in log scale for the three distributions, where $N = 2^8$. Observe that a very small number of items have high $\lambda_i$ in Zipf's law (solid line), many items have high $\lambda_i$ in the linear distribution (dotted line), and the geometric distribution (dashed line) falls between the two distributions.

Figure 8 shows the relative error in $H_{\mathrm{flu}}$ when $\lambda_i$ is geometrically distributed (Column (a)) and when $\lambda_i$ is linearly distributed (Column (b)). We show only the case where $N = 2^8$, but we find that the relative error is smaller for a larger $N$ as observed in Figure 6. We find that $K$ and the distribution of $\lambda_i$ have a rather complex impact on the relative error in $H_{\mathrm{flu}}$. When $\lambda_i$ is geometrically distributed, the relative error is larger for a smaller $K$. When $\lambda_i$ is linearly distributed, the relative error is larger for a larger $K$. In Figure 6(b), we have seen that the relative error is rather insensitive to $K$ when $\lambda_i$ follows Zipf's law. Our interpretation is that the relative error in $H_{\mathrm{flu}}$ is mainly due to the fact that the underutilization of $B_0$ becomes negligible in the fluid limit. When more items have high $\lambda_i$ (e.g., when linearly distributed), more items move from $B_0$ to $B_1$, which in turn makes $B_0$ more underutilized. However, $H_{\mathrm{flu}}$ overestimates the hit probability because requests are generated for items that are in $B_0$ but would not be in $B_0$ if the underutilization of $B_0$ was not negligible. Therefore, the the magnitude of the overestimation depends on $K$ and the distribution of $\lambda_i$ in a rather complex manner.
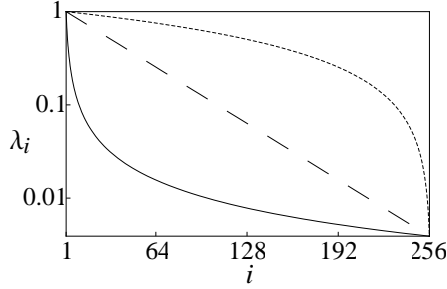
Figure 7: The values of $\lambda_i$ for Zipf's law (solid line), a geometric distribution (dashed line), and a linear distribution (dotted line), where $N = 2^8$.
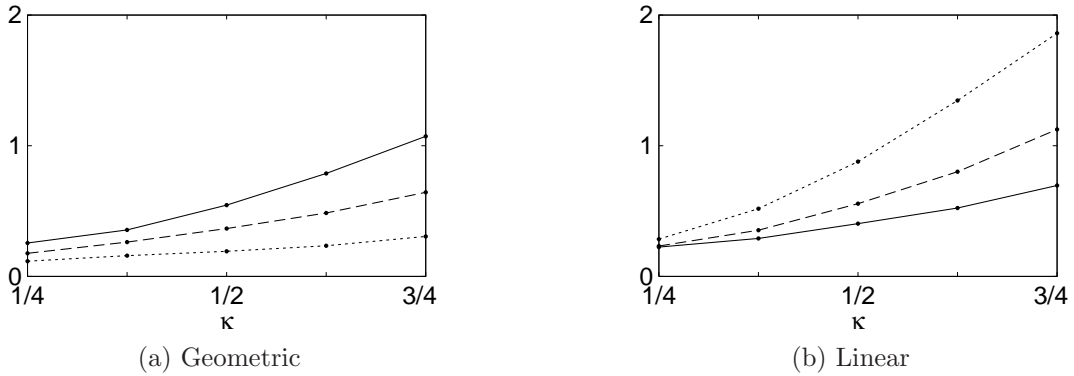


(a) Geometric

(b) Linear

Figure 8: The relative error in $H_{\text{flu}}$ when $\lambda_i$ is (a) distributed geometrically and (b) linearly, where $K = N/4$ for solid lines, $K = 3N/8$ for dashed lines, and $K = N/2$ for dotted lines, where $N = 2^8$.

## 5.2 Hit probabilities for 2Q and LRU

We will now evaluate the hit probabilities for $\texttt{2Q}(\kappa)$, $\texttt{LRU}$, and $\texttt{OPT}$. In Section 5.2.1, we will study the stationary hit probabilities in the fluid limit, using the analytical expressions derived in Sections 2 and 4. Above we have seen that the stationary hit probabilities are very well approximated by those in the fluid limit when $N > 1000$. In Section 5.2.2, we will study the transient hit probabilities via simulations.

### 5.2.1 Stationary hit probability

Figure 9 shows the stationary hit probabilities for $\texttt{2Q}(\kappa)$, $\texttt{LRU}$, and $\texttt{OPT}$. The stationary hit probabilities for $\texttt{2Q}(\kappa)$ are plotted by solid lines as functions of $\kappa$. The stationary probabilities of $\texttt{LRU}$ and $\texttt{OPT}$ are, respectively, plotted by straight dashed lines and by straight dotted lines. Four graphs correspond to four cases with different values of $N$ and $K$ as labeled. We assume that $\lambda_i$ follows Zipf's law, specifically $\lambda_i = 1/i$ for $1 \leq i \leq N$, but similar observations hold for other distributions as well.

Observe that the stationary hit probability for $\texttt{2Q}(\kappa)$ can be made higher than that for $\texttt{LRU}$ and very close to that for $\texttt{OPT}$ by choosing a sufficiently small $\kappa$ for any $N$ and $K$. Specifically, when $K_0 = \kappa K = 1$, the hit probability for $\texttt{2Q}(\kappa)$ is 3-17% higher than that for $\texttt{LRU}$ and within 1% of that for $\texttt{OPT}$ for the four cases in Figure 9. In fact, our analysis in Section 4 can be used to show that, in the fluid limit, the hit probability for $\texttt{2Q}(\kappa)$ can always be made equal to that for $\texttt{LRU}$ by choosing a particular value of $K_0$:

**Corollary 3** *Let $\kappa^\star \equiv \sum_{i=1}^{N} F_i(T)\,(1 - F_i(T))/K$, where $T$ is a unique $t$ such that $\sum_{i=1}^{N} F_i(t) = K$, where $F_i(t) = 1 - e^{-\lambda_i t}$. Then, in the fluid limit, the hit probability for $\texttt{2Q}(\kappa^\star)$ is equal to that for $\texttt{LRU}$.*

**Proof**: If $\kappa = \kappa^\star$, then $T_0 = T_1 = T$ satisfies (26) and (25). The arguments in the proof of Lemma 3 imply that, when (26) and (25) have a solution, the solution must be unique, and the hit probability in the fluid limit is given by (29). Hence, the hit probability of $e_i$ for $\texttt{2Q}(\kappa^\star)$ in the fluid limit is $F_i(T)/(F_i(T) + 1 - F_i(T))$, which is equal that for $\texttt{LRU}$ by Theorem 1, since $F_i(\cdot) \equiv G_i(\cdot)$ for a Poisson process. ∎

In general, $\kappa^\star K$ is fractional. We expect, however, that the hit probability for $\texttt{2Q}(\kappa)$ can be made close to that for $\texttt{LRU}$ by choosing $\kappa \approx \kappa^\star$ such that $\kappa K$ is an integer.

Figure 9: Stationary hit probabilities for $2\mathtt{Q}(\kappa)$ (solid lines), $\mathtt{LRU}$ (dashed lines), and $\mathtt{OPT}$ (dotted lines) where $\lambda_i = 1/i$ for $1 \leq i \leq N$.

Also, observe that, for fixed $N$ and $K$, the stationary hit probability for $2\mathtt{Q}(\kappa)$ is higher when $\kappa$ is smaller and maximized when $K_0 = \kappa K = 1$ (and $K_1 = K - 1$). Although we have not been able to provide an analytical proof, this is an observation that generally holds for all cases studied, including those not shown in the paper, unless $\lambda_i$ is a constant for all $i$. When $\lambda_i$ is a constant for all $i$, Corollary 2 and Corollary 3 imply that the hit probability for $2\mathtt{Q}(\kappa)$ cannot be made higher than that for $\mathtt{LRU}$ in the fluid limit.

Overall, the above observations suggest that we should choose $\kappa < \kappa^\star$ for $2\mathtt{Q}(\kappa)$ to achieve a hit probability higher than $\mathtt{LRU}$. In the next section, we will see that too small an $\kappa$ is not necessarily a good choice. Although a smaller $\kappa$ implies a higher *stationary* hit probability, a smaller $\kappa$ requires longer time to reach the stationary state. In particular, if we start from an empty buffer, $2\mathtt{Q}(\kappa)$ with a small $\kappa$ would suffer from a long period of low transient hit probabilities.

### 5.2.2 Transient hit probability

Figure 10 shows transient hit probabilities for $2\mathtt{Q}(\kappa)$ and $\mathtt{LRU}$. Specifically, starting with an empty buffer at time 0, the buffer is simulated until $10^8$ requests are generated. For each $10^4$ requests, the fraction of the requests that find the requested items in the buffer (fraction of hit) is recorded. The simulation with $10^8$ requests are repeated for 20 times, and the average fraction of hit over the 20 runs for every interval of $10^4$ requests is plotted for $\mathtt{LRU}$ and $2\mathtt{Q}(\kappa)$ with varying values of $\kappa$ (specifically, $\kappa = 1/2, 1/2^4, 1/2^7$). The four graphs correspond to the four cases studied in Figure 9.

We find that the hit probability for $\mathtt{LRU}$ (solid lines) quickly reaches the stationary hit probability, while the hit probability for $2\mathtt{Q}(\kappa)$ may need longer time particularly when $\kappa$ is small (dotted lines). When $K = N/8$ (top row), the stationary hit probability for $2\mathtt{Q}(\kappa)$ is higher than that for $\mathtt{LRU}$ for all of the three values of $\kappa$. However, when $\kappa = 1/2^4$ or $\kappa = 1/2^7$, the transient hit probability for $2\mathtt{Q}(\kappa)$ is lower than that for $\mathtt{LRU}$ until $10^4$ to $10^5$ requests are generated. When $\kappa = 1/2$, the hit probability for $2\mathtt{Q}(\kappa)$ quickly reaches the stationary hit probability and is higher than that for $\mathtt{LRU}$ after $10^4$ requests are generated. When $K = N/2$ (bottom row), the hit probability for $2\mathtt{Q}(\kappa)$ with $\kappa = 1/2$ also reaches the stationary hit probability quickly, but never exceeds that for $\mathtt{LRU}$, since the stationary hit probability for $2\mathtt{Q}(\kappa)$ is lower than that for $\mathtt{LRU}$ in this setting.

Although the time for the hit probability for $2\mathtt{Q}(\kappa)$ to reach the steady state is highly sensitive to $\kappa$, we find that it is relatively insensitive to $K$ (compare the top row and the bottom row of Figure 10). This may be explained by examining the time needed to fill $B_1$. When $K_1$ is larger, more items in $B_0$ need to be requested before $B_1$ is filled, so that a larger $K_0$ is needed to fill $B_1$ in a given time.
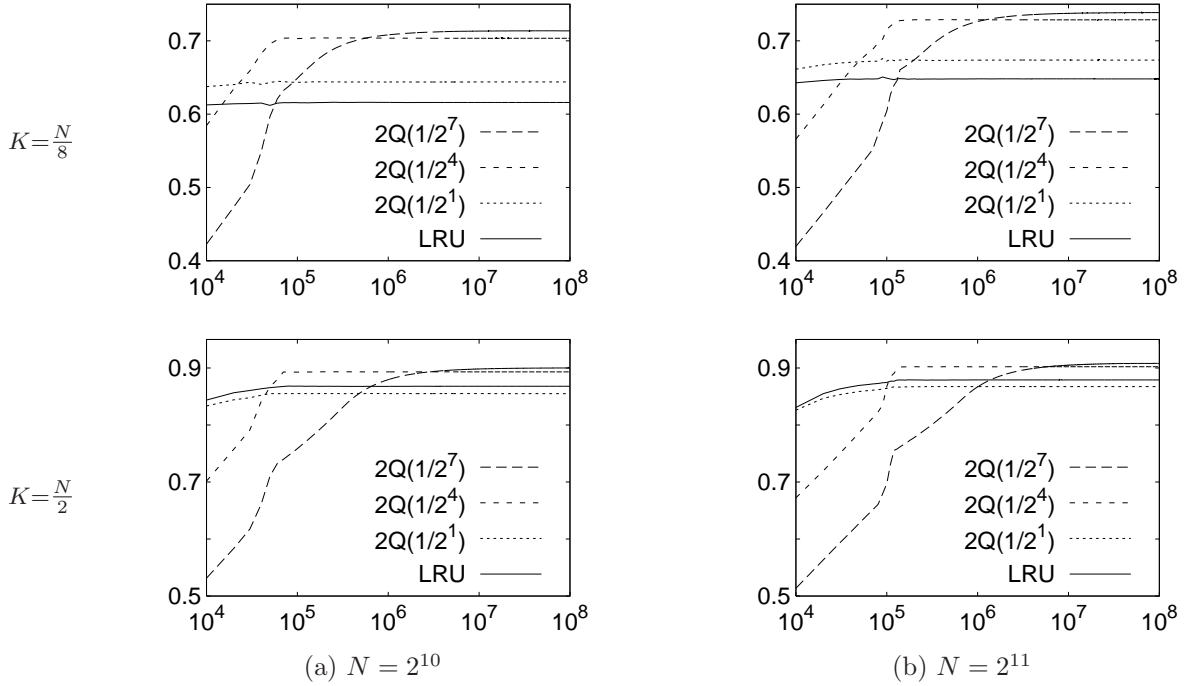
Figure 10: Transient hit probabilities for 2Q($\kappa$) and LRU against the number of requests generated, where $\lambda_i = 1/i$ for $1 \le i \le N$.

# 6 Conclusion

This paper provides an exact analysis of the hit probability for 2Q($\kappa$) in the fluid limit when items are requested according to independent Poisson processes. The analysis of 2Q relies on an analysis of LRU when events (requests and invalidations) for each item are generated according to a Poisson process and the type of an event is determined by a particular Markov chain. We remark that the analysis of LRU under this model of correlated arrivals is of interest in its own right. In a Web system for online shopping, for example, an object (an item) is created and cached when a user logs in to identify the user and to record the status of her shopping cart. The object is used (the item is requested) while the user is shopping. When the user logs out, the object is removed from the cache (the item is invalidated). These correlated requests and invalidations may be well represented by our model of correlated arrivals.

The hit probability for LRU has been found to be closely approximated by that in the fluid limit [11], but the analysis in the fluid limit has not been applied to other page replacement policies. We find not only that the hit probability for 2Q is also well approximated by that in the fluid limit but also that an analysis of 2Q in the fluid limit illuminates several fundamental properties of 2Q. We expect that analytical techniques introduced in this paper will be useful for an analysis of other page replacement policies in the fluid limit, particularly those that divide a buffer into multiple parts [17, 24, 12].

Another future direction is an extension of the analysis of 2Q to the case of a non-Poisson arrival process. The hit probability for LRU has been shown to be insensitive to some types of dependency in an arrival process [10, 23]. It would be of interest to examine whether similar properties hold for 2Q and other page replacement policies.

# Acknowledgements

# References

[1] J. Le Boudec, D. McDonald, and J. Mundinger. A generic mean field convergence result for systems of interacting objects. In *Proceedings of the 4th International Conference on the Quantitative Evaluation of Systems*, September 2007.

[2] Lee Breslau, Pei Cao, Li Fan, Graham Phillips, and Scott Shenker. Web caching and Zipf-like distributions: Evidence and implications. In *Proceedings of the IEEE INFOCOM*, pages 126–134, March 1999.

[3] Hao Che, Ye Tung, and Zhijun Wang. Hierarchical Web caching systems: Modeling, design and experimental results. *IEEE Journal on Selected Areas in Communications*, 20(7):1305–1314, 2002.

[4] James Allen Fill. Limits and rates of convergence for the distribution of search cost under the move-to-front rule. *Theoretical Computer Scicience*, 164(1-2):185–206, 1996.

[5] James Allen Fill and Lars Holst. On the distribution of search cost for the move-to-front rule. *Random Structures & Algorithms*, 8(3):179–186, 1996.

[6] Philippe Flajolet, D. Gardy, and L. Thimonier. Birthday paradox, coupon collectors, caching algorithms and self-organizing search. *Discrete Applied Mathematics*, 39(3):207–229, 1992.

[7] G. H. Gonnet, J. I. Munro, and H. Suwanda. Exegesis of self-organizing linear search. *SIAM Journal on Computing*, 10(3):613–637, 1981.

[8] T. Hama and R. Hirade. Cache policy optimization for response cache of Web application server. Technical Report RT0572, IBM Research, 2004.

[9] W. J. Hendricks. The stationary distribution of an interesting Markov chain. *Journal of Applied Probability*, 9(1):231–233, 1972.

[10] P. R. Jelenković and A. Radovanović. Least-recently-used caching with dependent requests. *Theoretical Computer Science*, 326(1-3):293–327, 2004.

[11] Predrag R. Jelenković. Asymptotic approximation of the move-to-front search cost distribution and least-recently-used caching fault probabilities. *Annals of Applied Probability*, 9(2):430–464, 1999.

[12] Song Jiang and Xiaodong Zhang. LIRS: An efficient low inter-reference recency set replacement policy to improve buffer cache performance. *ACM SIGMETRICS Performance Evaluation Review*, 30(1):31–42, 2002.

[13] T. Johnson and D. Shasha. 2Q: A low overhead high performance buffer management replacement algorithm. In *Proceedings of the 20th International Conference on Very Large Data Bases*, pages 439–450, September 1994.

[14] Nikolaos Laoutaris, Hao Che, and Ioannis Stavrakakis. The LCD interconnection of LRU caches and its analysis. *Performance Evaluation*, 63(7):609–634, 2006.

[15] A. M. Law and W. D. Kelton. *Simulation Modeling and Analysis*. McGraw-Hill, third edition, 2000.

[16] J. McCabe. On serial files with relocatable records. *Operations Research*, 13(4):609–618, 1965.

[17] Nimrod Megiddo and Dharmendra S. Modha. ARC: A self-tuning, low overhead replacement cache. In *Proceedings of the 2nd USENIX Conference on File and Storage Technologies*, pages 115–130, March 2003.

[18] B. Melamed and D. D. Yao. The ASTA property. In J. H. Dshalalow, editor, *Advances in Queueing: Theory, Methods, and Open Problems*, chapter 7. CRC Press, 1995.

[19] Elizabeth J. O'Neil, Patrick E. O'Neil, and Gerhard Weikum. The LRU-K page replacement algorithm for database disk buffering. *ACM SIGMOD Record*, 22(2):297–306, 1993.

[20] Michael Rabinovich and Oliver Spatscheck. *Web Caching and Replication*. Addison-Wesley, 1st edition, 2002.

[21] Raghu Ramakrishnan and Johannes Gehrke. *Database Management Systems*. McGraw-Hill, 3rd edition, 2002.

[22] Abraham Silberschatz, Peter Baer Galvin, and Greg Gagne. *Operating System Concepts*. John Wiley and Sons, 7th edition, 2004.

[23] Toyoaki Sugimoto and Naoto Miyoshi. On the asymptotics of fault probability in least-recently-used caching with Zipf-type request distribution. *Random Structures & Algorithms*, 29(3):296–323, 2006.

[24] Yuanyuan Zhou, James Philbin, and Kai Li. The multi-queue replacement algorithm for second level buffer caches. In *Proceedings of the USENIX Annual Technical Conference*, pages 91–104, June 2001.

# A    Technical lemmas

**Lemma 4** *Let $c$ and $s$ be nonnegative constants. Then $\left(c\,e^{-s/n} + 1 - c\right)^n \to e^{-c\,s}$ as $n \to \infty$.*

**Proof**: Observe that

$$\left(c\,e^{-s/n} + 1 - c\right)^n = \left(1 - \frac{c\,s}{n} + c\sum_{\ell=2}^{\infty}\left(\frac{-s}{n}\right)^\ell \frac{1}{\ell!}\right)^n = \sum_{k=0}^{n}\sigma_k,$$

where

$$\sigma_k = \binom{n}{k}\left(1 - \frac{c\,s}{n}\right)^{n-k}\left(c\sum_{\ell=2}^{\infty}\left(\frac{-s}{n}\right)^\ell \frac{1}{\ell!}\right)^k$$

for $0 \le k \le n$. Since $\sigma_0 \to e^{-c\,s}$ as $n \to \infty$, it suffices to show that $\sum_{k=1}^{n}\sigma_k \to 0$ as $n \to \infty$. Since $c, s \ge 0$ and $1/n^{\ell\,k} \le 1/n^{2k}$ for $\ell \ge 2$, we have

$$\left|\sum_{k=1}^{n}\sigma_k\right| \le \sum_{k=1}^{n}\binom{n}{k}\left(1 + \frac{c\,s}{n}\right)^{n-k}\frac{1}{n^{2k}}\left|\left(c\sum_{\ell=2}^{\infty}\frac{(-s)^\ell}{\ell!}\right)^k\right|.$$

Now, since $\binom{n}{k} \le n^k/k!$, we have

$$\left|\sum_{k=1}^{n}\sigma_k\right| \le \sum_{k=1}^{n}\frac{1}{k!\,n^k}\left(1 + \frac{c\,s}{n}\right)^{n-k}\left(c\,|e^{-s} - 1 + s|\right)^k$$

$$= \left(1 + \frac{c\,s}{n}\right)^n \sum_{k=1}^{n}\frac{1}{k!}\left(\frac{c\,|e^{-s} - 1 + s|}{n + c\,s}\right)^k.$$

Since the summands are nonnegative, we have

$$\left|\sum_{k=1}^{n}\sigma_k\right| \le \left(1 + \frac{c\,s}{n}\right)^n \sum_{k=1}^{\infty}\frac{1}{k!}\left(\frac{c\,|e^{-s} - 1 + s|}{n + c\,s}\right)^k$$

$$= \left(1 + \frac{c\,s}{n}\right)^n \left(e^{\frac{c\,|e^{-s}-1+s|}{n+c\,s}} - 1\right).$$

Since the right hand size approaches 0 as $n \to \infty$, the left hand size also approaches 0 as $n \to \infty$.    ∎


**Lemma 5** *Let $\mathbf{P}$ and $\mathbf{Q}$ be generator matrices of Markov chains such that*

$$\mathbf{P} = \begin{pmatrix} -\lambda & (1-\alpha)\,\lambda & \alpha\,\lambda \\ 0 & -\alpha\,\lambda & \alpha\,\lambda \\ 0 & \beta\,\lambda & -\beta\,\lambda \end{pmatrix}$$

$$\mathbf{Q} = \begin{pmatrix} -\lambda & \beta\,\lambda & (1-\beta)\,\lambda \\ 0 & -\alpha\,\lambda & \alpha\,\lambda \\ 0 & \beta\,\lambda & -\beta\,\lambda \end{pmatrix}.$$

*Let $\mathbf{v}_i$ be a unit column vector with three elements such that the $i$-th element is 1 and $\mathbf{v}_i^{\mathrm{t}}$ be the corresponding unit row vector. Then, for $\gamma = \beta/(\alpha + \beta)$,*

$$\mathbf{v}_1^{\mathrm{t}}\left(\gamma\,e^{-\mathbf{P}\,t} + (1-\gamma)\,e^{-\mathbf{Q}\,t}\right)\mathbf{v}_2 = \gamma\left(1 - e^{-\lambda\,t}\right).$$

**Proof**: Since $\mathbf{v}_1^t \mathbf{A} \mathbf{v}_2$ is the $(1,2)$ element of a matrix, $\mathbf{A}$, of size $3 \times 3$, it suffices to obtain the $(1,2)$ element of $e^{-\mathbf{P} t}$ and $e^{-\mathbf{Q} t}$, respectively. Recall that $e^{\mathbf{A}} \equiv \sum_{n=0}^{\infty} \mathbf{A}^n / n!$. Observe that, for $n \geq 1$,

$$\frac{\mathbf{P}^n}{(-\lambda)^n} = \begin{pmatrix} 1 & -1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} + (\alpha + \beta)^{n-1} \begin{pmatrix} 0 & \alpha & -\alpha \\ 0 & \alpha & -\alpha \\ 0 & -\beta & \beta \end{pmatrix}$$

$$\frac{\mathbf{Q}^n}{(-\lambda)^n} = \begin{pmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} + (\alpha + \beta)^{n-1} \begin{pmatrix} 0 & -\beta & \beta \\ 0 & \alpha & -\alpha \\ 0 & -\beta & \beta \end{pmatrix},$$

which can be verified by induction. Note that the above expressions are invalid for $n = 0$ and that $\mathbf{P}^0 = \mathbf{Q}^0 = \mathbf{I}$, where $\mathbf{I}$ is the identity matrix of size $3 \times 3$. Therefore,

$$\mathbf{v}_1^t \, e^{-\mathbf{P} t} \, \mathbf{v}_2 \;\; = \;\; -e^{-\lambda t} + \frac{\alpha}{\alpha + \beta} e^{-(\alpha+\beta)\lambda t} + \frac{\beta}{\alpha + \beta}$$

$$\mathbf{v}_1^t \, e^{-\mathbf{Q} t} \, \mathbf{v}_2 \;\; = \;\; -\frac{\beta}{\alpha + \beta} e^{-(\alpha+\beta)\lambda t} + \frac{\beta}{\alpha + \beta}.$$

Now the lemma follows from the definition of $\gamma$. ∎