

ON PROBABILISTIC ANALYSIS OF A COALESCED HASHING ALGORITHM

BY B. PITTEL

The Ohio State University

An allocation model [n balls, m ($\geq n$) cells, at most one ball in a cell] related to a hashing algorithm is studied. A ball x goes into the cell $h(x)$, where $h(\cdot): \{1, \dots, n\} \rightarrow \{1, \dots, m\}$ is random. In case the cell $h(x)$ is already occupied, the ball x is rejected and moved into the leftmost empty cell. This empty cell is found via the sequential search from left to right starting with the cell occupied by the last (before x) rejected ball. Denote $T_2(x)$ the number of the necessary probes. In the end, due to a resulting system of references, the n occupied cells form a disjoint union of ordered chains, and to locate a ball x it suffices to search only the cells of a subchain originating at the cell $h(x)$. Denote $T_1(x)$ the length of this subchain. The main result of the paper is: in probability,

$$\max T_1(x) = \log_b n - 2 \log_b \log n + O(1),$$

$$\max T_2(x) = \log_b n - \log_b \log n + O(1),$$

as $n \rightarrow \infty$, if n/m is bounded away from 0, $b = (1 - e^{-n/m})^{-1}$.

1. Introduction. Suppose we have n different objects (keys, in computer science terminology) labelled $1, 2, \dots, n$, which are to be allocated one after another in an m -long array of cells (a table), so that each cell would contain at most one key. (Naturally, $m \geq n$.) Assume also that there is *given* a function $h(\cdot): \{1, \dots, n\} \rightarrow \{1, \dots, m\}$, so that, for $1 \leq x \leq n$, $h(x)$ is the index of a cell the key x is assigned to, if this cell is still empty. [In applications, the keys may be numbers, records, i.e., elements of a given key space K . Thus, we should begin with a function $H(\cdot): K \rightarrow \{1, \dots, m\}$. But if the keys to be allocated in the table are k_1, \dots, k_n we may and shall restrict our attention to $h(i) = H(k_i)$, $1 \leq i \leq n$.] An allocation algorithm must include instructions about how to handle collisions, that is, the situations when $h(\cdot)$ has the same value for several keys. In addition, an algorithm must be able to locate a key if it is already stored in the table.

Out of the many existing (so called hashing) algorithms meeting these conditions, we consider in this paper a coalesced hashing algorithm [4], [5], [12], [16], [24] and [25].

This is how it works. Suppose that the keys $1, 2, \dots, x-1$ are already stored. Suppose also that the $x-1$ occupied cells form a union of disjoint subsets (chains), with an order \leq induced by the labels of the keys occupying them,

Received May 1985; revised February 1986.

AMS 1980 *subject classifications*. 60C05, 60F99, 68P10, 68P20, 68R05, 05C80.

Key words and phrases. Search algorithm, hashing, probabilistic analysis, limiting distributions, largest search time.

having the property: If a key x' is in a cell y , then the cells $h(x')$, y belong to the same chain and $h(x') \leq y$. (Technically, the chains are specified by references associated with the cells; see the example that follows.) Now, if the cell $h(x)$ is not occupied by a key $x' < x$, then the key x moves into this cell, and a new chain of length 1 is formed. Otherwise, the key x is rejected and moves into the leftmost empty cell; the latter joins, as a new last cell, the chain which passes through the cell $h(x)$, thus producing a longer chain. [The empty cell itself is found via the sequential (left-to-right) search starting from the right neighbor of a cell occupied by the key rejected last before the key x .] In the end, all the keys are stored, and the n occupied cells form several disjoint chains. [If needed, we can locate now any key x in the table by searching the cells of a subchain which starts at the cell $h(x)$.]

Efficiency of the algorithm is measured by a sequence of the search times $\{T_1(x), T_2(x)\}_{1 \leq x \leq n}$. Here $T_1(x)$ is the length of a subchain connecting the cell $h(x)$ and a cell which actually contains the key x . $T_2(x)$ is the number of extra cells searched sequentially to accommodate the key x , in case the cell $h(x)$ is occupied by a key $x' < x$.

EXAMPLE. $m = 10, n = 8$ and $\{h(x): 1 \leq x \leq 8\} = \{5, 9, 10, 5, 2, 9, 2, 1\}$. The keys are allocated as shown on Figure 1. The chains of cells are $\{10\}$, $\{2, 4\}$, $\{9, 3\}$ and $\{5, 1, 6\}$. Also, $\{T_1(x): 1 \leq x \leq 8\} = \{1, 1, 1, 2, 1, 2, 2, 2\}$ and $\{T_2(x): 1 \leq x \leq 8\} = \{0, 0, 0, 1, 0, 2, 1, 2\}$.

Assume that the hashing function $h(\cdot)$ is chosen at random, which means that $h(1), \dots, h(n)$ are independent and each $h(x)$ is uniformly distributed on $\{1, \dots, m\}$. Then the search times become random. Let also $n, m \rightarrow \infty$ so that $n = am$ and $0 < a \leq 1$ is fixed. (The parameter a is usually called the load factor.)

Exact and asymptotic formulas for $E(T_i(x))$ and $\sigma^2(T_1(x))$ were obtained in [16], and it was observed there that the *average* search times remain bounded even when $a = 1$, which corresponds to the completely filled table. (For the mean values analysis of other versions of the coalesced hashing, see also [4], [5], [15], [24] and [25].)

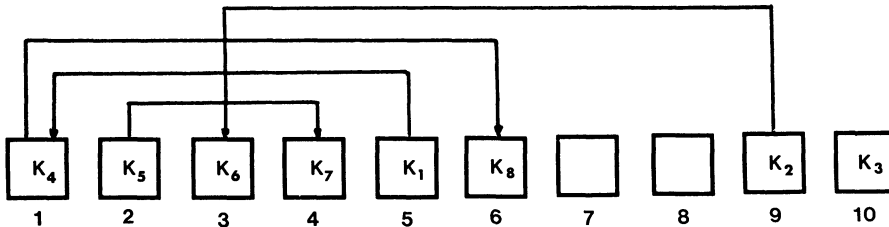


FIG. 1.

Our first result is

THEOREM 1. (a) $\lim P(T_1(n) = 1) = 1 - a$, and for $l \geq 0$

$$(1.1) \quad \begin{aligned} &\lim P(T_1(n) = l + 2) \\ &= a(1 + H_l) - \frac{a^2}{2} + \sum_{k=1}^l (-1)^k \binom{k+1}{k^2} \binom{l}{k} (1 - e^{-ka}), \end{aligned}$$

where $H_l = 1/1 + \dots + 1/l$. Also

$$(1.2) \quad P(T_1(n) = l) \sim e^a(1 - e^{-a})^{l-1}/l, \quad a = n/m,$$

if $l \rightarrow \infty$ and $l^2 = o(n)$.

(b) $T_2(n) \rightarrow \delta(1 + \tau)$ in distribution, where δ and τ are independent, $P(\delta = 1) = 1 - P(\delta = 0) = a$ and τ is geometrically distributed with parameter e^{-a} , i.e., $P(\tau = s) = e^{-a}(1 - e^{-a})^s, s \geq 0$.

Consider now $U_i(n) = \max\{T_i(x): 1 \leq x \leq n\}$, $i = 1, 2$; these random variables represent the largest search times. It turns out that $U_i(n)$ grow logarithmically with n . More precisely, we shall prove

THEOREM 2. In probability,

$$(1.3) \quad U_1(n) = \log_b n - 2 \log_b \log n + O(1),$$

$$(1.4) \quad U_2(n) = \log_b n - \log_b \log n + O(1),$$

where $b = (1 - e^{-a})^{-1}$.

NOTES. (1) It is quite surprising that $U_1(n)$ and $U_2(n)$ behave so similarly, since the corresponding searches are very different. (2) Consider $T(x) = T_1(x) + T_2(x)$ which is the total number of probes needed to accommodate the key x , $1 \leq x \leq n$, and denote $U(n) = \max\{T(x): 1 \leq x \leq n\}$. In view of Theorem 2,

$$\log_b n(1 + o(1)) \leq U(n) \leq 2 \log_b n(1 + o(1))$$

in probability, and we conjecture that, out of these two estimates, the upper case is sharp.

The reader interested in a description and probabilistic analysis of other hashing techniques (linear probing, double hashing, etc.) and related search algorithms can look at [1], [3], [6]–[16], [19]–[23] and [26]. (The list is by no means complete!) Undoubtedly, the invention of hashing algorithms and a need to analyze them have stimulated a new interest in the well known combinatorial schemes—random allocations [17] and mappings [18].

2. The proof of $T_1(n)$ and $U_1(n)$ asymptotic characteristics. Denote by $\tilde{h}(x)$ the index of a cell which actually contains the key x . By the definition of

the algorithm, $T_1(n) = 1$ iff $n(n) \neq \tilde{h}(x)$, $x < n$. Since $h(n)$ is independent of $\tilde{h}(x)$, $x < n$, and is uniformly distributed,

$$(2.1) \quad P(T_1(n) = 1) = 1 - (n - 1)/m \rightarrow 1 - a.$$

How can we evaluate $P(T_1(n) = l + 2)$, $l \geq 0$? Introduce $X(\nu, l)$ the number of the chains induced by the first ν keys, with the length *more* than l . Denote $E(\nu, l) = E(X(\nu, l))$. Clearly, $T_1(n) = l + 2$ iff the cell $h(n)$ belongs to a chain (induced by the first $n - 1$ keys) with length *at least* $l + 1$, and the length of the subchain connecting the cell $h(n)$ and the maximal cell of the chain is *exactly* $l + 1$. Hence, again by independence,

$$(2.2) \quad \begin{aligned} P(T_1(n) = l + 2) &= E(P(T_1(n) = l + 2 | \tilde{h}(x), x < n)) \\ &= E(X(n - 1, l)/m) = E(n - 1, l)/m. \end{aligned}$$

In view of (2.2), the relation (1.1) follows directly from

LEMMA 1. For $l + 1 \leq \nu \leq m$,

$$(2.3) \quad \begin{aligned} E(\nu, l) &= \nu(1 + H_l) - m^{-1} \binom{\nu}{2} \\ &+ m \sum_{k=1}^l (-1)^k \left(\frac{k+1}{k^2} \right) \binom{l}{k} \left[1 - \left(1 - \frac{k}{m} \right)^\nu \right], \\ H_l &= \sum_{k=1}^l \frac{1}{k}. \end{aligned}$$

PROOF. For $A = \{x_1 < \dots < x_{l+1}\}$, $x_j \leq \nu$, put $\delta(A) = 1$ if the cells $\tilde{h}(x_1), \dots, \tilde{h}(x_{l+1})$ form an initial segment of a chain; if not, put $\delta(A) = 0$. Observe that $\delta(A) = 1$ iff $h(x_1) \neq \tilde{h}(x)$ for $x < x_1$, $h(x) \neq \tilde{h}(x_1)$ ($= h(x_1)$) for $x_1 < x < x_2$, $h(x_2) = \tilde{h}(x_1), \dots, h(x_l) \in \{\tilde{h}(x_1), \dots, \tilde{h}(x_{l-1})\}$, $h(x) \notin \{\tilde{h}(x_1), \dots, \tilde{h}(x_l)\}$ for $x_l < x < x_{l+1}$, and $h(x_{l+1}) \in \{\tilde{h}(x_1), \dots, \tilde{h}(x_l)\}$. Obviously, $X(\nu, l) = \sum_A \delta(A)$, so

$$E(\nu, l) = \sum_A E(\delta(A)) = \sum_A P(\delta(A) = 1).$$

Here

$$(2.4) \quad \begin{aligned} P(\delta(A) = 1) &= [1 - (x_1 - 1)/m](1 - 1/m)^{x_2 - x_1 - 1} (1/m) \dots \\ &[(l - 1)/m](1 - l/m)^{x_{l+1} - x_l - 1} (l/m) \\ &= (1 - s_0/m)(l!/m^l) \prod_{j=1}^l (1 - j/m)^{s_j}, \end{aligned}$$

$s_j = x_{j+1} - x_j - 1, x_0 = 0.$

Therefore,

$$E(\nu, l) = (l!/m^l) \left[\sum_s (1 - s_0/m) \prod_{j=1}^l (1 - j/m)^{s_j} \right];$$

here the summation is taken over all the tuples $s = (s_0, \dots, s_{l+1})$, s_j being nonnegative integers and $\sum_{j=0}^{l+1} s_j = \nu - (l + 1)$. In other words, the sequence $\{E(u + (l + 1), l)(l!/m^l)^{-1}\}_{u \geq 0}$ is an $(l + 2)$ -fold convolution of $\{1 - u/m\}_{u \geq 0}$, $\{(1 - j/m)^u\}_{u \geq 0}$, $1 \leq j \leq l$, and $\{1\}_{u \geq 0}$. Hence

$$E(\nu, l) = (l!/m^l) \text{coeff}_{z^{\nu-(l+1)}} F(z, l),$$

$$F(z, l) = \left[(1 - z)^{-2} - z(1 - z)^{-3}/m \right] \prod_{j=1}^l [1 - z(1 - j/m)]^{-1}, \quad |z| < 1.$$

So, by the Cauchy formula,

$$E(\nu, l) = (l!/m^l)(2\pi i)^{-1} \int_C G(z) dz, \quad G(z) = F(z, l)/z^{\nu-l};$$

here $C = \{z = re^{i\phi}: -\pi \leq \phi < \pi\}$, $r < 1$. Outside C , $G(z)$ has the poles $z_j = (1 - j/m)^{-1}$, $0 \leq j \leq l$. Hence, by the residue theorem,

$$E(\nu, l) = -(l!/m^l) \left[\text{res}_{z=\infty} G(z) + \sum_{j=0}^l \text{res}_{z=z_j} G(z) \right].$$

Here $\text{res}_{z=\infty} G(z) = 0$ because $d \max\{|G(z)|: |z| = d\} \rightarrow 0$ as $d \rightarrow \infty$. A direct evaluation of the other residues leads, after simple but tedious manipulations, to the formula (2.3). \square

Our next lemma yields immediately [see (2.2)] the formula (1.2), and it will be also used later to prove the relation (1.3).

LEMMA 2. *Suppose that $l, \nu, m \rightarrow \infty$, $l^2/m = o(1)$ and ν/m is bounded away from 0. Then*

$$(2.5) \quad E(\nu, l) = me^{\nu} l^{-1} (1 - e^{-\nu})^{l+1} (1 + o(1)), \quad \nu = \nu/m.$$

PROOF. Introduce two auxiliary functions

$$(2.6) \quad \begin{aligned} E_1(t; \nu, l) &= t\nu + m \sum_{k=1}^l (-1)^k k^{-1} \binom{l}{k} \left[1 - \left(1 - t \frac{k}{m} \right)^\nu \right], \\ E_2(t; \nu, l) &= -t\nu H_l + t^2 m^{-1} \binom{\nu}{2} \\ &\quad - m \sum_{k=1}^l (-1)^k k^{-2} \binom{l}{k} \left[1 - \left(1 - t \frac{k}{m} \right)^\nu \right], \end{aligned}$$

where $t \in [0, 1]$. Clearly,

$$(2.7) \quad E(\nu, l) = E_1(1; \nu, l) - E_2(1; \nu, l).$$

Now, differentiating $E_1(t; \nu, l)$ and $E_2(t; \nu, l)$ with respect to t , we have

$$(2.8) \quad E_1'(t; \nu, l) = \nu \sum_{k=0}^l (-1)^k \binom{l}{k} \left(1 - t \frac{k}{m}\right)^{\nu-1},$$

$$(2.9) \quad \begin{aligned} E_2'(t; \nu, l) &= -\nu H_l + tm^{-1}\nu(\nu - 1) - \nu \sum_{k=1}^l (-1)^k k^{-1} \binom{l}{k} \left(1 - t \frac{k}{m}\right)^{\nu-1} \\ &= tm^{-1}\nu(\nu - 1) + \nu \sum_{k=1}^l (-1)^k k^{-1} \binom{l}{k} \left[1 - \left(1 - t \frac{k}{m}\right)^{\nu-1}\right] \\ &= m^{-1}\nu \left\{ t(\nu - 1) + m \sum_{k=1}^l (-1)^k k^{-1} \binom{l}{k} \left[1 - \left(1 - t \frac{k}{m}\right)^{\nu-1}\right] \right\} \\ &= m^{-1}\nu E_1(t; \nu - 1, l) \quad (!). \end{aligned}$$

(We have used here the identity $H_l = \sum_{k=1}^l (-1)^{k-1} k^{-1} \binom{l}{k}$, which can be easily proved by induction.)

The sum in (2.8) has a nice probabilistic interpretation. Namely, consider an allocation scheme with w balls and m cells, which differs from the classical scheme in that a ball is accepted by a cell, chosen by the ball, with probability t independently of all other balls and cells. Consider some fixed l cells, and denote $P(t; w, l)$ the probability that none of these l cells is empty. Since the probability that some k fixed cells are empty equals $(1 - t(k/m))^w$, the inclusion-exclusion principle shows that

$$(2.10) \quad P(t; w, l) = \sum_{k=0}^l (-1)^k \binom{l}{k} \left(1 - t \frac{k}{m}\right)^w.$$

Combining (2.8) and (2.9) with (2.10), and using $E_i(0; \nu, l) = 0$ [see (2.6)], we get

$$(2.11) \quad E_1(1; \nu, l) = \nu \int_0^1 P(t; \nu - 1, l) dt,$$

$$(2.12) \quad E_2(1; \nu, l) = m^{-1}\nu(\nu - 1) \int_0^1 (1 - t)P(t; \nu - 2, l) dt.$$

To proceed, we need a sharp estimate of the function $P(t; w, l)$, and it could hardly be done directly via (2.10). Fortunately, using the probabilistic interpretation of $P(t; w, l)$, we can show that

$$(2.13) \quad P(t; w, l) = w! \text{coeff}_{z^w} [e^z(1 - e^{-zt/m})^l].$$

Further, an application of the saddle-point method yields that, uniformly over $t \in (0, 1]$,

$$(2.14) \quad P(t; w, l) \sim (1 - e^{-gt})^l, \quad g = w/m,$$

provided that $l^2/m = o(1)$, $g \neq o(1)$. (See Appendix A.)

The rest is simpler. By (2.11) and (2.14),

$$(2.15) \quad E_1(1; \nu, l) \sim \nu \int_0^1 (1 - e^{-g_1 t})^l dt, \quad g_1 = (\nu - 1)/m.$$

Introducing $\psi(t) = \log(1 - e^{-g_1 t})$, we observe $\psi'(1) = g_1(e^{g_1} - 1)^{-1}$, $\psi''(t) \leq 0$, $t \in (0, 1]$ and is bounded for $t \in [\frac{1}{2}, 1]$. By Jensen's inequality,

$$(2.16) \quad \begin{aligned} \int_0^1 \exp[l\psi(t)] dt &\leq \int_0^1 \exp\{l[\psi(1) + \psi'(1)(t - 1)]\} dt \\ &= (1 - e^{-g_1})^l \int_0^1 \exp[l\psi'(1)(t - 1)] \\ &\sim (1 - e^{-g_1})^l (l\psi'(1))^{-1}, \quad l \rightarrow \infty, \\ &= e^{g_1} (lg_1)^{-1} (1 - e^{-g_1})^{l+1}. \end{aligned}$$

On the other hand, denoting $t_0 = 1 - l^{-2/3}$,

$$(2.17) \quad \begin{aligned} \int_0^1 (1 - e^{-g_1 t})^l dt &\geq \int_{t_0}^1 (1 - e^{-g_1 t})^l dt \\ &= \int_{t_0}^1 \exp\{l[\psi(1) + \psi'(1)(t - 1)] + O(l(t - 1)^2)\} dt \\ &\sim (1 - e^{-g_1})^l (l\psi'(1))^{-1} \int_0^{l^{1/3}\psi'(1)} e^{-u} du \\ &\sim e^{g_1} (lg_1)^{-1} (1 - e^{-g_1})^{l+1}, \quad l(t_0 - 1)^2 = l^{-1/3}. \end{aligned}$$

By (2.15)–(2.17),

$$(2.18) \quad E_1(1; \nu, l) \sim me^{\nu} l^{-1} (1 - e^{-\nu})^{l+1}, \quad \nu = \nu/m.$$

To estimate $E_2(1; \nu, l)$, observe first that, by the definition, $P(t; w, l)$ increases with w . Hence [see (2.11), (2.12) and (2.14)]

$$(2.19) \quad \begin{aligned} E_2(1; \nu, l) &\leq (\nu/m)\nu \int_0^1 (1 - t)P(t; \nu - 1, l) dt \\ &\sim (\nu/m)\nu \int_0^1 (1 - t)(1 - e^{-g_1 t})^l dt \\ &= (\nu/m) \int_0^1 \Psi(t) dt, \end{aligned}$$

where

$$(2.20) \quad \Psi(t) = \nu \int_0^t (1 - e^{-g_1 s})^l ds.$$

Now, the direct differentiation shows that $\Psi''(t) - lc\Psi'(t) \geq 0$, $t \in [0, 1]$, $c = (e - 1)^{-1}$; therefore, $\Psi'(t) - lc\Psi(t) \geq \Psi'(0) - lc\Psi(0) = 0$, and

$$\int_0^1 \Psi(t) dt \leq (lc)^{-1} \int_0^1 \Psi'(t) dt = (lc)^{-1} \Psi(1).$$

Together with (2.15), (2.19) and (2.20), it implies that

$$E_2(1; \nu, l) = O(l^{-1}E_1(1; \nu, l)) = o(E_1(1; \nu, l))$$

(remember, $l \rightarrow \infty$). This relation, (2.7) and (2.18) combined complete the proof of Lemma 2. \square

Let us turn to $U_1(n) = \max\{T_1(x) : 1 \leq x \leq n\}$. Introduce

$$(2.21) \quad \begin{aligned} l_1 &= \log_b n - 2 \log_b \log_b n + \omega(n), \\ l_2 &= \log_b n - 2 \log_b \log_b n - \omega(n), \end{aligned}$$

where $\omega(n) \rightarrow \infty$ and $\omega(n) = o(\log n)$.

PROOF OF (1.3). We need to show that

$$P(U_1(n) > l_1) \rightarrow 0, \quad P(U_1(n) > l_2) \rightarrow 1, \quad n \rightarrow \infty.$$

Fix an integer $0 \leq l \leq n - 1$. For a set $A = \{x_1 < \dots < x_{L+1}\}$, $l \leq L \leq n - 1$, put $\epsilon(A) = 1$ if the cells $\tilde{h}(x_1), \dots, \tilde{h}(x_{L+1})$ form an initial segment of a chain, and $T_1(x_1), \dots, T_1(x_L) \leq l$, but $T_1(x_{L+1}) \geq l + 1$, that is,

$$h(x_j) \in \begin{cases} \{\tilde{h}(x_1), \dots, \tilde{h}(x_{j-1})\}, & 2 \leq j \leq l, \\ \{\tilde{h}(x_{j-(l-1)}), \dots, \tilde{h}(x_{j-1})\}, & l + 1 \leq j \leq L, \\ \{\tilde{h}(x_1), \dots, \tilde{h}(x_{L-(l-1)})\}, & j = L + 1. \end{cases}$$

Otherwise, put $\epsilon(A) = 0$.

Introduce the random variables $Y(n, l)$ and $Z(n, l)$ which are the total number of such segments and the total number of segments of length *exactly* $l + 1$, respectively. Clearly,

$$Y(n, l) = \sum_{A: |A| \geq l+1} \epsilon(A), \quad Z(n, l) = \sum_{A: |A|=l+1} \epsilon(A).$$

Also

$$(2.22) \quad P(U_1(n) > l) = P(Y(n, l) > 0) \leq E(Y(n, l)),$$

and by Chebyshev's inequality,

$$(2.23) \quad P(U_1(n) > l) \geq P(Z(n, l) > 0) \geq E^2(Z(n, l))/E(Z^2(n, l)).$$

We will use (2.21) with $l = l_1$ and (2.22) with $l = l_2$. To estimate $E(Y(n, l))$, $E(Z(n, l))$ and $E(Z^2(n, l))$, we observe first that, by (2.4), and the definition of $\epsilon(A)$,

$$(2.24) \quad \begin{aligned} E(\epsilon(A)) &= P(\epsilon(A) = 1) \\ &= (1 - s_0/m) [(l - 1)!(l - 1)^{L-l}(L - l + 1)/m^L] \\ &\times \prod_{j=1}^L (1 - j/m)^{s_j}, \quad s_j = x_{j+1} - x_j - 1, \quad x_0 = 0, \\ &= E(\delta(A))(l - 1)^{L-l}(L - l + 1)/(L)_{L-l+1}. \end{aligned}$$

Let $l \sim \log_b n$, $b = (1 - e^{-a})^{-1}$, $a = n/m$. Denote $\mu = [3 \log_b n]$ and write

$$E(Y(n, l)) = \Sigma'_A E(\epsilon(A)) + \Sigma''_A E(\epsilon(A)) = \Sigma' + \Sigma'',$$

where $|A| \leq \mu$ in Σ' and $|A| > \mu$ in Σ'' . Since $\epsilon(A) \leq \delta(A)$ and $E(n, L)$ decreases as a function of L , we have by Lemma 2,

$$(2.25) \quad \begin{aligned} \Sigma'' &\leq \Sigma''_A E(\delta(A)) = \sum_{L \geq \mu} \sum_{A: |A|=L+1} E(\delta(A)) = \sum_{L \geq \mu} E(n, L) \\ &\leq nE(n, \mu) = O(n^2 b^{-\mu}) = O(n^{-1}). \end{aligned}$$

Further, by (2.24) (and Lemma 2 again),

$$\begin{aligned} \Sigma' &= \sum_{L=l}^{\mu-1} E(n, L)(l-1)^{L-l}(L-l+1)/(L)_{L-l+1} \\ &\sim me^a(l^2 b^{l+1})^{-1} f(n, l), \\ f(n, l) &= \sum_{L=l}^{\mu-1} (b^{L-l} L/l)^{-1} (l-1)^{L-l} (L-l+1)/(L)_{L-l}. \end{aligned}$$

Here

$$1 \leq f(n, l) \leq \sum_{L=l}^{\infty} (L-l+1)(b^{-1})^{L-l} = e^{2a} \leq e^2.$$

Therefore, if $l = l_1$ [see (2.21)], then

$$\Sigma' = O(n(\log_b n)^{-2} b^{-l}) = O(b^{-\omega(n)}),$$

and [see (2.22) and (2.25)]

$$P(U_1(n) > l) \leq \Sigma' + \Sigma'' = O(b^{-\omega(n)}) + O(n^{-1}) = o(1).$$

Let now $l = l_2$ [see (2.21)]. First of all, according to (2.24),

$$(2.26) \quad \begin{aligned} E(Z(n, l)) &= \sum_{A: |A|=l+1} E(\epsilon(A)) = l^{-1} \sum_{A: |A|=l+1} E(\delta(A)) = l^{-1} E(n, l) \\ &\sim (e^a - 1)m(l^2 b^l)^{-1} = a^{-1}(e^a - 1)b^{\omega(n)} \rightarrow \infty. \end{aligned}$$

Further, since $\epsilon(A_1)\epsilon(A_2) = 0$ for $A_1 \neq A_2$ and $A_1 \cap A_2 \neq \emptyset$,

$$(2.27) \quad E(Z^2(n, l)) = E(Z(n, l)) + \sum_{A_1, A_2} P(\epsilon(A_1) = 1, \epsilon(A_2) = 1),$$

where the sum is taken over all ordered pairs of disjoint subsets A_1, A_2 with $|A_1| = |A_2| = l + 1$.

LEMMA 3. For all such pairs (A_1, A_2) ,

$$P(\epsilon(A_1) = 1, \epsilon(A_2) = 1) \leq c(l, m)P(\epsilon(A_1) = 1)P(\epsilon(A_2) = 1),$$

where

$$c(l, m) = [1 - 2(l + 1)/m]^{-2(l+1)}.$$

PROOF. For $B \subset \{1, \dots, n\}$, denote $I(B)$ the smallest interval which contains B , and $J(B) = I(B) \setminus B$. Introduce also $f_i(x)$, the total number of the keys

from A_i strictly less than x , $i = 1, 2$. Then [see (2.24)]

$$(2.28) \quad P(\varepsilon(A_i) = 1) = [(l - 1)!/m^l](1 - s_{0i}/m) \prod_{x \in J_i} (1 - f_i(x)/m),$$

$$J_i = J(A_i),$$

where $s_{0i} = x_{1i} - 1$, x_{1i} being the first key from A_i , $i = 1, 2$. It is not difficult to see that

$$(2.29) \quad P(\varepsilon(A_1) = 1, \varepsilon(A_2) = 1) = [(l - 1)!/m^l]^2(1 - s_{01}/m)(1 - s_{02}/m) \times \prod_{x \in J} (1 - f(x)/m), \quad J = J(A_1 \cup A_2),$$

where

$$f(x) = \begin{cases} f_1(x) + f_2(x), & \text{if } \max(f_1(x), f_2(x)) \leq l, \\ \min(f_1(x), f_2(x)), & \text{if } \max(f_1(x), f_2(x)) \geq l + 1. \end{cases}$$

By (2.28),

$$\begin{aligned} P(\varepsilon(A_1) = 1)P(\varepsilon(A_2) = 1) &= [(l - 1)!/m^l]^2(1 - s_{01}/m)(1 - s_{02}/m)\Pi_1\Pi_2\Pi_3, \\ \Pi_1 &= \prod_{x \in J_1 \cap J_2} (1 - f_1(x)/m)(1 - f_2(x)/m), \\ \Pi_2 &= \prod_{x \in J_1 \setminus J_2} (1 - f_1(x)/m), \\ \Pi_3 &= \prod_{x \in J_2 \setminus J_1} (1 - f_2(x)/m). \end{aligned}$$

But $f_i(x) \leq l$ on J_i , $i = 1, 2$, so, by the definition of $f(\cdot)$,

$$\begin{aligned} (1 - f_1(x)/m)(1 - f_2(x)/m) &\geq (1 - (f_1(x) + f_2(x))/m) \\ &= (1 - f(x)/m), \quad x \in J_1 \cap J_2, \\ 1 - f_i(x)/m &\geq 1 - f(x)/m, \quad x \in J_i, \quad i = 1, 2. \end{aligned}$$

Hence,

$$(2.30) \quad \begin{aligned} P(\varepsilon(A_1) = 1)P(\varepsilon(A_2) = 1) &\geq [(l - 1)!/m^l]^2(1 - s_{01}/m)(1 - s_{02}/m) \prod_{x \in J_1 \cup J_2} (1 - f(x)/m). \end{aligned}$$

Besides, it can be checked that

$$J_1 \cup J_2 \subset J \cup (I_1 \cap A_2) \cup (I_2 \cap A_1),$$

so

$$(2.31) \quad |(J_1 \cup J_2) \setminus J| \leq |I_1 \cap A_2| + |I_2 \cap A_1| \leq |A_2| + |A_1| = 2(l + 1).$$

Since $f(x) \leq 2(l + 1)$ as well, the combination of (2.29)–(2.31) leads to

$$P(\varepsilon(A_1) = 1)P(\varepsilon(A_2) = 1) \geq (1 - 2(l + 1)/m)^{2(l+1)}P(\varepsilon(A_1) = 1, \varepsilon(A_2) = 1).$$

□

By this lemma and (2.26) and (2.27),

$$E(Z^2(n, l)) \leq E(Z(n, l)) + c(l, m)E^2(Z(n, l)) \\ \leq E^2(Z(n, l)) [1 + O(l^2/m) + O(b^{-\omega(n)})].$$

Hence [see (2.23)],

$$P(U_1(n) > l) \geq [1 + O(l^2/m) + O(b^{-\omega(n)})]^{-1} = 1 + o(1).$$

The proof of (1.3) is now complete. \square

NOTE. Introduce $C(n)$, the length of the longest chain (“monster”; cf. [11]). It is obvious that $C(n) \geq T_1(n)$. Using $X(n, l)$ instead of $Y(n, l)$ and $Z(n, l)$, we could have proved in the same way that, in probability,

$$C(n) = \log_b n - \log_b \log n + O(1).$$

3. The proof of $T_2(n)$ and $U_2(n)$ asymptotic behavior. Recall that $T_2(x)$ is the number of extra cells searched sequentially to find an available cell for the key x , in case the cell $h(x)$ is occupied by a previous key, $1 \leq x \leq n$. [This search starts from the right neighbor of a cell occupied by the last (before x) rejected key, or from the cell 1 if no key $< x$ has been rejected.] In particular, $T_2(x) = 0$ iff $h(x) \notin \{\tilde{h}(x')\}_{x' < x}$. Denote $K(x) = \sum_{x' \leq x} T_2(x')$ [if $K(x) \neq 0$, then it is the index of a cell filled by a key rejected last among the keys $\leq x$]. We shall also need $L(x)$, the length of the maximal block of cells, beginning from the first cell, which are all occupied by the keys $\leq x$. [“Maximal” means that no key $\leq x$ occupies the cell $(L(x) + 1)$.] Put finally $M(x) = L(x) - K(x)$. (The expected values of $K(x)$, $L(x)$, $M(x)$ were obtained in [16].)

NOTE. Observe that $T_2(x) = \delta(x)[1 + M(x - 1)]$, where $P(\delta(x) = 1) = 1 - P(\delta(x) = 0) = (x - 1)/m$, and $\delta(x)$ is independent of $M(x - 1)$ (in fact of the whole $\{h(x') : x' < x\}$).

LEMMA 4. Let $k \leq l \leq x \leq n$. If $k \geq 1$, then

$$(3.1) \quad P(K(x) = k, L(x) \geq l) = m^{-x} \binom{m-l}{m-x} (x-1)^x P(x, x-1, x-k) \\ = m^{-x} \binom{m-l}{m-x} \sum_{j=0}^{x-k} (-1)^j \binom{x-k}{j} (x-1-j)^x,$$

where $P(u, v, w)$ is the probability that in the usual allocation model, with u balls and v cells, some w fixed cells are not empty. Also

$$(3.2) \quad P(K(x) = 0, L(x) \geq l) = m^{-x} \binom{m-l}{m-x} x!.$$

PROOF. The relation (3.2) is easy, since $K(x) = 0$ iff no cell is chosen by two or more keys $\leq x$. [We say a cell y is chosen by a key x' if $h(x') = y$.]

Let $K(x) = k \geq 1$ and $L(x) = l$. By the definition of the algorithm and $K(x)$, $L(x)$, we have then (a) no key $\leq x$ has chosen the cell k , (b) all the cells

$k + 1, \dots, l$ are chosen by the keys $\leq x$ and (c) exactly $x - l$ cells among those to the right from the cell l are chosen by the keys $\leq x$.

Conversely, suppose that the conditions (a), (b), (c) are met. If $K(x) < k$ then [by (a)] the total number of the cells occupied by the keys $\leq x$ is at most $x - 1$ —impossible. If $K(x) > k$, then [by (b)] $K(x) > l$ and [by (c)] the keys $\leq x$ occupy at least $x + 1$ cells—impossible again. Hence, $K(x) = k$ and [see (b)] $L(x) \geq l$.

Consequently,

$$\begin{aligned}
 P(K(x) = k, L(x) \geq l) &= P(\text{(a), (b), (c) are satisfied}) \\
 &= \binom{m-l}{x-l} \left(\frac{x-1}{m}\right)^x P(x, x-1, x-k) \\
 &= \binom{m-l}{x-l} \left(\frac{x-1}{m}\right)^x \sum_{j=0}^{x-k} (-1)^j \binom{x-k}{j} \left(1 - \frac{j}{x-1}\right)^x.
 \end{aligned}$$

□

COROLLARY 1 (straightforward).

$$(3.3) \quad P(K(x) = k) = \begin{cases} m^{-x} \binom{m-k}{m-x} (x-1)^x P(x, x-1, x-k), & k \geq 1, \\ m^{-x} \binom{m}{m-x} x!, & k = 0. \end{cases}$$

COROLLARY 2. If $m - x = O(1)$, then $M(x) \Rightarrow G_1$, where G_1 is geometrically distributed with parameter e^{-1} .

[In view of the note preceding Lemma 4, the last corollary implies Theorem 1(b) in case $a = 1$.]

PROOF OF COROLLARY 2. We may and shall assume that $m - x = \mu$ is fixed. Denote $R(x) = x - K(x)$. By Lemma 4, for fixed $r \geq s \geq 0$,

$$\begin{aligned}
 P(R(x) = r, M(x) \geq s) &= P(K(x) = x - r, L(x) \geq x + s - r) \\
 &= \binom{\mu + r - s}{\mu} (1 - x^{-1})^x (1 + \mu x^{-1})^{-x} \\
 &\quad \times \sum_{j=0}^r (-1)^j \binom{r}{j} (1 - j(x-1)^{-1})^x \\
 &\rightarrow \binom{\mu + r - s}{\mu} (e^{-1})^{\mu+1} (1 - e^{-1})^r
 \end{aligned}$$

or

$$\begin{aligned}
 P(R(x) = r, M(x) = s) \\
 \rightarrow \left[(e^{-1})^\mu (1 - e^{-1})^{r-s} \binom{\mu - 1 + r - s}{\mu - 1} \right] [e^{-1} (1 - e^{-1})^s].
 \end{aligned}$$

So

$$(R(x), M(x)) \Rightarrow (G_1^{(\mu)} + G_1, G_1),$$

where $G_1^{(\mu)}$ is negative binomially distributed with parameters e^{-1}, μ , and $G_1^{(\mu)}$ and G_1 are independent. [Just as easily, we could have proved that $M(x)$ and the lengths of all μ consecutive blocks of occupied cells to the right from the cell $L(x)$ are asymptotically independent and geometrically distributed with parameter e^{-1} .] \square

To handle the case $m - x \rightarrow \infty$, and later $-U_2(n)$, the following two lemmas are needed.

LEMMA 5. *There exists an absolute constant c_0 such that*

$$(3.4) \quad P(u, v, w) \leq c_0(u/rv)^{u+1/2} \exp(rv - u)(1 - e^{-r})^w, \quad \forall r > 0.$$

LEMMA 6. *Introduce*

$$(3.5) \quad \bar{k} = \bar{k}(x) = m - (m - x)e^\lambda, \quad \lambda = x/m.$$

For every $\varepsilon \in (0, 1)$, there exists $\mu_0 = \mu_0(\varepsilon)$ such that for $x \geq \varepsilon m$ and $\mu = m - x \geq \mu_0$,

$$(3.6) \quad P(|K(x) - \bar{k}| \geq (m - x)^{3/5}) \leq \exp[-(m - x)^{1/5}/30].$$

PROOF OF LEMMA 5. Arguing as in the case of the probability $P(t; w, l)$ in Section 2 (see Appendix A), we have

$$(3.7) \quad \begin{aligned} P(u, v, w) &= (u!/v^u) \text{coeff}_{z^u} \{ \exp[z(v - w)](e^z - 1)^w \} \\ &= (u!/v^u)(2\pi)^{-1} \int_{-\pi}^{\pi} \mathcal{H}(re^{i\phi}) d\phi, \end{aligned}$$

$$\mathcal{H}(z) = \exp[z(v - w)](e^z - 1)^w/z^u, \quad r > 0.$$

Here

$$(3.8) \quad \begin{aligned} |\mathcal{H}(re^{i\phi})| &\leq \mathcal{H}(r) \exp[(rv/2)(\cos \phi - 1)] \\ &\leq \mathcal{H}(r) \exp(-c_1rv\phi^2), \quad c_1 > 0, \end{aligned}$$

since

$$|\exp(re^{i\phi})| = e^r \exp[r(\cos \phi - 1)]$$

and it can be checked that

$$|\exp(re^{i\phi}) - 1| \leq |e^r - 1| \exp[r(\cos \phi - 1)/2]$$

(see [22, Appendix]).

Hence

$$P(u, v, w) \leq c_0(u!/v^u)(rv)^{-1/2} \mathcal{H}(r), \quad c_2 = 2^{-1}(\pi c_1)^{-1/2},$$

and, using an inequality $u! \leq c_3 u^{1/2}(u/e)^u$, we obtain (3.4) with $c_0 = c_2 c_3$. \square

PROOF OF LEMMA 6. First of all

$$(3.9) \quad \begin{aligned} |K(x) - \bar{k}| &= |x - R(x) - \bar{k}| \\ &= |R(x) - (m - x)(e^\lambda - 1)|, \quad \lambda = x/m. \end{aligned}$$

Further, by Lemmas 4 and 5 with $r = \lambda$, we have, for $0 \leq s \leq j < x$,

$$(3.10) \quad \begin{aligned} P(R(x) = j, M(x) \geq s) &= P(K(x) = x - k, L(x) \geq x + s - k) \\ &= m^{-x}(x - 1)^x \binom{\mu + j - s}{\mu} P(x, x - 1, j) \\ &\leq c_0 \binom{\mu + j - s}{\mu} m^{-x}(x - 1)^x [m/(x - 1)]^{x+1/2} (e^{-\lambda})^{\mu+1} (1 - e^{-\lambda})^j \\ &\leq c \binom{\mu + j - s}{\mu} (e^{-\lambda})^{\mu+1} (1 - e^{-\lambda})^j \\ &= cP(G_\lambda^{(\mu)} + G_\lambda = j, G_\lambda \geq s), \quad \mu = m - x, c = 2c_0. \end{aligned}$$

[G_λ and $G_\lambda^{(\mu)}$ are independent, $G_\lambda(G_\lambda^{(\mu)})$ is geometrically (negative binomially) distributed with parameter $e^{-\lambda}$ (with parameters $e^{-\lambda}, \mu$).] Consequently,

$$(3.11) \quad P(R(x) = j) \leq cP(G_\lambda^{(\mu')} = j), \quad j < x, \mu' = \mu + 1.$$

Also [see (3.3)],

$$(3.12) \quad P(R(x) = x) = P(K(x) = 0) = (m)_x/m^x \leq \exp[-x(x - 1)/2m].$$

Observe that

$$E(G_\lambda^{(\mu')}) = \mu' E(G_\lambda) = \mu'(e^\lambda - 1),$$

so [see (3.9), (3.11) and (3.12)]

$$\begin{aligned} P(|K(x) - \bar{k}| \geq \mu^{3/5}) &\leq \exp[-x(x - 1)/2m] \\ &\quad + cP(|G_\lambda^{(\mu')} - E(G_\lambda^{(\mu')})| \geq \chi(\mu')^{3/5}), \end{aligned}$$

where $\chi \rightarrow 1$ as $\mu \rightarrow \infty$.

The rest of the proof follows a general idea, owing to Chernoff [2]. Namely, we write

$$P(G_\lambda^{(\mu')} \geq (\leq) j) \leq g(z)/z^j, \quad z \geq (\leq) 1,$$

where

$$g(z) = E(z^{G_\lambda^{(\mu')}}) = [E(z^{G_\lambda})]^\mu = [p/(1 - qz)]^\mu, \quad p = 1 - q = e^{-\lambda}.$$

Choosing in each case z which minimizes $g(z)/z^j$, we obtain

$$(3.13) \quad \begin{aligned} P(G_\lambda^{(\mu')} \geq (\leq) j) &\leq \exp[\Phi_{mx}(j)], \quad j \geq (\leq) q\mu'/p = E(G_\lambda^{(\mu')}), \\ \Phi_{mx}(y) &= \mu' \log[p(\mu' + y)/\mu'] + y \log[q(\mu' + y)/y]. \end{aligned}$$

Now, since $\Phi_{mx}(q\mu'/p) = \Phi'_{mx}(q\mu'/p) = 0$ and

$$\Phi''_{mx}(y) = -\mu'/[y(\mu' + y)],$$

we easily get from (3.13) that ($p = e^{-\lambda}$)

$$(3.14) \quad P(|G_\lambda^{(\mu')} - E(G_\lambda^{(\mu')})| \geq \chi(\mu')^{3/5}) \leq 2 \exp[-p^2(m-x)^{6/5} \chi^2/4\mu'] \\ \leq \exp[-(m-x)^{1/5}/29.6],$$

provided that μ' is large enough.

Using (3.12) (remember, $x \geq \epsilon m$) and (3.14), we obtain (3.6). \square

This lemma enables us to complete the proof of Theorem 1(b) by proving

LEMMA 7. *If $\liminf_m x/m \in (0, 1]$, $m - x \rightarrow \infty$ and $s = o((m - x)^{1/5})$, then*

$$P(M(x) \geq s) \sim (1 - e^{-\lambda})^s, \quad \lambda = x/m.$$

PROOF. By (3.2),

$$P(M(x) \geq s) = m^{-x} \binom{m-s}{m-x} x! + P_1 + P_2,$$

where

$$P_1 = P(M(x) \geq s, |K(x) - \bar{k}| \leq (m-x)^{3/5}),$$

$$P_2 = P(M(x) \geq s, |K(x) - \bar{k}| > (m-x)^{3/5}), \quad \bar{k} = m - (m-x)e^\lambda.$$

It can be checked that, uniformly over k such that $|k - \bar{k}| \leq (m-x)^{3/5}$,

$$(3.15) \quad \binom{m-k-s}{m-x} = \binom{m-k}{m-x} [(x-\bar{k})/(m-\bar{k})]^s [1 + O(s/(m-x)^{2/5})] \\ = \binom{m-k}{m-x} (1 - e^{-\lambda})^s (1 + o(1)).$$

By (3.3), (3.15) and Lemma 6, we have then

$$P_1 = (1 - e^{-\lambda})^s P(|K(x) - \bar{k}| \leq (m-x)^{3/5}) (1 + o(1)) \\ \sim (1 - e^{-\lambda})^s.$$

It remains to observe that

$$P_2 \leq \exp[-(m-x)^{1/5}/30] = o[(1 - e^{-\lambda})^s],$$

[$s = o((m-x)^{1/5})$] and

$$m^{-x} \binom{m-s}{m-x} x! \leq m^{-x} \binom{m}{m-x} x! = O(\exp(-\epsilon^2 m)) \\ = o[(1 - e^{-\lambda})^s]. \quad \square$$

Next:

LEMMA 8. (a) *If $\liminf_m x/m \in (0, 1]$ and $s = o(m)$, then*

$$(3.16) \quad P(M(x) \geq s) = O[(1 - e^{-\lambda})^s], \quad \lambda = x/m.$$

(b) Denote $\mathcal{M}(x) = \max\{M(x'): x' \leq x\}$. Then, uniformly over $x \leq m$, and $s \geq 0$,

$$(3.17) \quad P(\mathcal{M}(x) \geq s) = O[m(1 - e^{-\lambda})^s].$$

PROOF. (a) According to (3.11) and (3.12),

$$\begin{aligned} P(M(x) \geq s) &\leq \exp[-x(x - 1)/2m] + c \sum_{j=s}^{\infty} P(G_\lambda^{(j)} + G_\lambda = j, G_\lambda \geq s) \\ &= \exp[-x(x - 1)/2m] + cP(G_\lambda \geq s) = O[(1 - e^{-\lambda})^s]. \end{aligned}$$

(b) By the definition, $M(x')$ is the length of the block of cells chosen by the keys $\leq x'$, which begins from the right neighbor of the cell occupied by a key rejected last among those keys. Hence, $\mathcal{M}(x) \leq B(x)$, where $B(x)$ is the length of the longest block of cells chosen by the keys $\leq x$. So,

$$(3.18) \quad P(\mathcal{M}(x) \geq s) \leq P(B(x) \geq s) \leq mP(x, m, s)$$

[see Lemma 4 for the definition of $P(u, v, w)$]. Applying Lemma 5 ($u = x$, $v = m$, $w = s$, $r = \lambda = x/m$), we get

$$P(\mathcal{M}(x) \geq s) \leq c_0 m(1 - e^{-\lambda})^s. \quad \square$$

With the last lemma at our disposal, we can now prove

LEMMA 9 [on an upper bound of $U_2(n)$]. In probability,

$$U_2(n) \leq \log_b n - \log_b \log_b n + O(1).$$

PROOF. It suffices to estimate from above $\mathcal{M}(n)$, since $1 + M(x) \geq T_2(x)$, $1 \leq x \leq n$. Denote $n_1 = [n/2]$, $\mathcal{M}(n_1, n) = \max\{M(x): n_1 \leq x \leq n\}$. Let also $s = \log_b n - \log_b \log_b n + \omega(n)$, where $\omega(n) \rightarrow \infty$, $\omega(n) = o(\log n)$ and is otherwise arbitrary. Since

$$\mathcal{M}(n) = \max\{\mathcal{M}(n_1), \mathcal{M}(n_1, n)\},$$

we have

$$\begin{aligned} (3.19) \quad P(\mathcal{M}(n) \geq s) &\leq P(\mathcal{M}(n_1) \geq s) + P(\mathcal{M}(n_1, n) \geq s) \\ &\leq P(\mathcal{M}(n_1) \geq s) + \sum_{x=n_1}^n P(M(x) \geq s) \\ &= O\left\{n[1 - \exp(-n/2m)]^s + \sum_{x=n_1}^n [1 - \exp(-x/m)]^s\right\}. \end{aligned}$$

Here $[b = (1 - e^{-a})^{-1}, a = n/m]$

$$\begin{aligned} n[1 - \exp(-n/2m)]^s &= nb^{-s}[(1 - \exp(-n/2m))/b]^s \\ &= O[(\log_b n)\rho^{\log_b n}] \quad (\rho = [(1 - \exp(-n/2m))/b]^{1/2}) \\ &= O(n^{-c}), \quad c > 0, \end{aligned}$$

since $\rho < 1$ and is bounded away from 1. Furthermore, estimating the sum in

(3.19) by an integral and arguing as in (2.20)–(2.23), we also have

$$(3.20) \quad \sum_{x=n_1}^n [1 - \exp(-x/m)]^s \sim n \int_{1/2}^1 [1 - \exp(-at)]^s dt \sim c(a)ns^{-1}b^{-s},$$

$c(a) = e^a(ab)^{-1}$. By the choice of s ,

$$ns^{-1}b^{-s} = b^{-\omega(n)}.$$

In view of (3.19) and (3.20), and $\omega(n) = o(\log n)$,

$$P(\mathcal{M}(n) \geq s) = O(b^{-\omega(n)}) \rightarrow 0, \quad n \rightarrow \infty,$$

and the lemma is proven. \square

It remains to establish a similar bound from below. Let us describe our plan of action. Denote $s = \log_b n - \log_b \log_b n - \omega(n)$ [$\omega(n) \rightarrow \infty$, $\omega(n) = o(\log n)$], $n_1 = \lfloor n/2 \rfloor$, $n_2 = n$ if $a = n/m < 1$, and $n_2 = \lfloor n - n^{1/2} \rfloor$ if $a = n/m = 1$. Define $\Delta(x) = 1$ [$\Delta(x) = 0$] if $M(x) \geq s$ [$M(x) < s$], and introduce $Y(n) = \sum_{x=n_1}^{n_2} \Delta(x)$.

We shall (a) estimate $E(Y(n))$ and $E(Y^2(n))$, (b) show—via Chebyshev’s inequality—that $Y(n)$ is unbounded in probability and (c) deduce from this that $U_2(n) \geq s$ with probability approaching 1 as $n \rightarrow \infty$.

LEMMA 10.

$$(3.21) \quad E(Y(n)) \sim e^a(ab)^{-1}b^{\omega(n)}, \quad n \rightarrow \infty \text{ (and } m \rightarrow \infty \text{)}.$$

PROOF. For m large enough, $m - x \geq n^{1/2}$, $s = o((m - x)^{1/5})$ if $n_1 \leq x \leq n_2$. Since $E(\Delta(x)) = P(M(x) \geq s)$, summing over $n_1 \leq x \leq n_2$ and using Lemma 7, we have [see (3.20)]

$$\begin{aligned} E(Y(n)) &= \sum_x P(M(x) \geq s) \sim \sum_x (1 - e^{-x/m})^s \\ &\sim n \int_{n_1/n}^{n_2/n} [1 - \exp(-at)]^s dt \sim n \int_{1/2}^1 [1 - \exp(-at)]^s dt \sim c(a)ns^{-1}b^s \\ &= e^a(ab)^{-1}b^{\omega(n)}. \end{aligned} \quad \square$$

LEMMA 11.

$$(3.22) \quad E(Y^2(n)) \sim E^2(Y(n)).$$

(See Appendix B for the proof.)

COROLLARY. $Y(n) \rightarrow \infty$ in probability.

PROOF. (For the sake of completeness.)

$$\begin{aligned} P(Y(n) \leq E(Y(n))/2) &\leq P(|Y(n) - E(Y(n))| \geq E(Y(n))/2) \\ &\leq 4\sigma^2(Y(n))/E^2(Y(n)) \\ &= 4[E(Y^2(n))/E^2(Y(n)) - 1] \rightarrow 0. \end{aligned} \quad \square$$

Finally,

LEMMA 12. *In probability, $U_2(n) \geq s$.*

PROOF. Introduce

$$t(n) = \begin{cases} \min\{1 \leq x \leq n: M(x) \geq s\}, & \text{if } \mathcal{M}(n) \geq s, \\ n, & \text{if } \mathcal{M}(n) < s. \end{cases}$$

[Recall that $\mathcal{M}(n) = \max\{M(x): 1 \leq x \leq n\}$.] Since $Y(n)$ is the number of keys x (between n_1 and n_2) such that $M(x) \geq s$, the previous corollary yields that $n - t(n) \rightarrow \infty$ in probability. Also, $t(n)$ is a stopping time adapted to the sequence $\{h(x): 1 \leq x \leq n\}$; in other words, for each j , $\{t(n) = j\}$ belongs to the σ -field generated by $\{h(x): 1 \leq x \leq j\}$. Hence, $\{t(n) = j\}$ is independent of $\{h(x): j + 1 \leq x \leq n\}$.

Consequently,

$$\begin{aligned} P(U_2(n) < s) &= \sum_{j=1}^n P(t(n) = j, \forall x > j \text{ chooses an empty cell}) \\ (3.23) \qquad &= \sum_{j=1}^n P(t(n) = j) \prod_{k=j}^{n-1} (1 - k/m). \end{aligned}$$

Fix an integer $d > 0$. Then, it follows from (3.23) that

$$\begin{aligned} P(U_2(n) < s) &\leq P(t(n) > n - d) + \sum_{j=1}^{n-d} P(t(n) = j) \prod_{k=j}^{n-1} (1 - k/m) \\ &\leq P(t(n) > n - d) + \prod_{k=n-d}^{n-1} (1 - k/m) \\ &\leq P(t(n) > n - d) + \exp[-d(n - d)/m]. \end{aligned}$$

Since $P(t(n) > n - d) \rightarrow 0$, as $n \rightarrow \infty$, we have

$$\limsup_n P(U_2(n) < s) \leq \exp(-da), \quad a = n/m,$$

and, letting $d \uparrow \infty$,

$$\limsup_n P(U_2(n) < s) = 0,$$

because $a > 0$. The lemma is proven. \square

This completes the proof of Theorem 1(b) and (1.4) of Theorem 2.

APPENDIX A

On the asymptotic formula for $P(t; w, l)$. Observe first that by the definition of $P(t; w, l)$, rather than (2.10),

$$(A.1) \quad P(t; w, l) = \sum_s (w!/m^w) \prod_{j=1}^m (s_j!)^{-1} \prod_{k=1}^l [1 - (1 - t)^{s_k}],$$

where $s_j \geq 0$, $1 \leq j \leq m$ and $\sum_{j=1}^m s_j = w$. Indeed, $(w!/m^w) \prod_{j=1}^m (s_j!)^{-1}$ is the probability that s_j balls choose the j th cell, $1 \leq j \leq m$, and $\prod_{k=1}^l [1 - (1-t)^{s_k}]$ is the probability that each of the first l cells accepts at least one of the balls which have chosen it. (Of course, the latter probability equals 0 whenever $s_k = 0$ for some $1 \leq k \leq l$.)

Using (A.1), we can derive a surprisingly simple formula for the exponential generating function of the sequence $\{P(t; w, l)\}_{w \geq 0}$:

$$\begin{aligned} \sum_{w \geq 0} P(t; w, l) z^w / w! &= \sum_{s_1, \dots, s_m \geq 0} \prod_{j=1}^m (z/m)^{s_j} / s_j! \prod_{k=1}^l [1 - (1-t)^{s_k}] \\ &= \left\{ \sum_{s \geq 0} [1 - (1-t)^s] (z/m)^s / s! \right\}^l \left[\sum_{s \geq 0} (z/m)^s / s! \right]^{m-l} \\ &= \{ \exp(z/m) - \exp[z(1-t)/m] \}^l \exp[z(m-l)/m] \\ &= e^z [1 - \exp(-zt/m)]^l. \end{aligned}$$

So, choosing a contour $C = \{a = re^{i\phi}: -\pi \leq \phi < \pi\}$, we have

$$\begin{aligned} (A.2) \quad P(t; w, l) &= w! (2\pi i)^{-1} \int_C e^z (1 - e^{-zt/m})^l z^{-w-1} dz \\ &= w! (2\pi)^{-1} \int_{-\pi}^{\pi} \exp(F(re^{i\phi})) d\phi, \end{aligned}$$

where

$$F(z) = z + l \log(1 - e^{-zt/m}) - w \log z.$$

To estimate the last integral, we use the saddle-point method. Select r which minimizes $F(u)$, $u \in (0, \infty)$, so that

$$(A.3) \quad F'(r) = 1 + (e^{rt/m} - 1)^{-1} (lt/m) - w/r = 0.$$

[Such a point exists since $\lim_{u \rightarrow 0^+} F'(u) = \lim_{u \rightarrow 0^+} (l - w)/u = -\infty$ ($l < w$) and $\lim_{u \rightarrow \infty} F'(u) = 1$.] By (A.3),

$$r = w - l(rt/m)/(e^{rt/m} - 1) = w + O(l)$$

uniformly over $t \in (0, 1]$, since $\eta/(e^\eta - 1) < 1$ for all $\eta > 0$. (All the related estimates below are also uniform.) “Bootstrapping,” we get a sharper estimate,

$$r = w - lgt/(e^{gt} - 1) + O(l^2/m), \quad g = w/m.$$

Since

$$(A.4) \quad F''(u) = (w/u^2) \left[1 - (l/w) \eta^2 e^\eta / (e^\eta - 1)^2 \right], \quad \eta = ut/m,$$

and $F'(r) = 0$, we have

$$\begin{aligned} (A.5) \quad F(r) &= F(w) - (F''(\tilde{u})/2)(w-r)^2 \quad (r \leq \tilde{u} \leq w) \\ &= F(w) + O(l^2/w) = w - w \log w + l \log(1 - e^{-gt}) + O(l^2/w). \end{aligned}$$

In addition, a simple computation shows that

$$|r^3 F^{(3)}(re^{i\phi})| = 2w(1 + O(l/w)) = O(w)$$

for all small enough $|\phi|$. Then, for $|\phi| \leq \phi_0 = w^{-5/12}$,

$$\begin{aligned} F(re^{i\phi}) &= F(r) + (r^2F''(r)/2)(e^{i\phi} - 1)^2 + O(w|e^{i\phi} - 1|^3) \quad [\text{see (A.4)}] \\ &= F(r) - (r^2F''(r)/2)\phi^2 + O(w\phi_0^3) \end{aligned}$$

and $O(w\phi_0^3) = O(w^{-1/4})$. So $[\lambda^2 = r^2F''(r) = w(1 + O(l/w))]$

$$\begin{aligned} \text{(A.6)} \quad \int_{-\phi_0}^{\phi_0} \exp(F(re^{i\phi})) d\phi &\sim \exp(F(r))\lambda^{-1/2} \int_{-\lambda\phi_0}^{\lambda\phi_0} \exp(-u^2/2) du \\ &\sim (2\pi/w)^{1/2} \exp(F(r)) \end{aligned}$$

because $\lambda\phi_0 \sim w^{1/12}$.

On the other hand, since $|e^z - 1| \leq |e^{|z|} - 1|$, we have for $|\phi| \geq |\phi_0|$,

$$\begin{aligned} |\exp(F(re^{i\phi}))| &\leq \exp(F(r) + (1 - tl/m)r(\cos \phi - 1)) \\ &\leq \exp(F(r) - cw\phi^2) \leq \exp(F(r) - cw^{1/6}), \quad c > 0, \end{aligned}$$

whence

$$\text{(A.7)} \quad \left| \int_{|\phi| > \phi_0} \exp(F(re^{i\phi})) d\phi \right| = O(\exp(F(r) - cw^{1/6})).$$

Putting together (A.2), (A.5)–(A.7) and applying the Stirling formula to $w!$, we have

$$\begin{aligned} P(t; w, l) &\sim w!(2\pi)^{-1}(2\pi/w)^{1/2} \exp(w - w \log w + l \log(1 - e^{-gt})) \\ &\sim (1 - e^{-gt})^l. \end{aligned}$$

APPENDIX B

Proof of Lemma 11. By the definition of $Y(n)$,

$$\begin{aligned} E(Y^2(n)) &= \sum_{x_1, x_2} P(M(x_i) \geq s, i = 1, 2) \\ &= \sum_x P(M(x) \geq s) + \sum_{x_1 \neq x_2} P(M(x_i) \geq s, i = 1, 2) \\ &= \Sigma_1 + \Sigma_2, \quad n_1 \leq x, x_1, x_2 \leq n_2. \end{aligned}$$

Here $\Sigma_1 = E(Y(n))$, $\Sigma_2 = \Sigma_{21} + \Sigma_{22}$ and

$$\begin{aligned} \Sigma_{21} &= \sum_{x_1 \neq x_2} P(M(x_i) \geq s, i = 1, 2; K(x_1) \neq K(x_2)), \\ \Sigma_{22} &= \sum_{x_1 \neq x_2} P(M(x_i) \geq s, i = 1, 2; K(x_1) = K(x_2)). \end{aligned}$$

Notice that, for $x_1 < x_2$, the condition $K(x_1) = K(x_2)$ means that no key among the keys $x_1 + 1, \dots, x_2$ is rejected; so, in particular, $M(x_2) \geq M(x_1)$. Therefore,

$$\begin{aligned} P(M(x_i) \geq s, i = 1, 2; K(x_1) = K(x_2)) &= P(M(x_1) \geq s) \prod_{x_1 \leq j < x_2} (1 - j/m) \\ &\leq P(M(x_1) \geq s) \beta^{x_2 - x_1}, \\ &\quad \beta = 1 - [n/2]/m, \end{aligned}$$

and

$$(B.1) \quad \begin{aligned} \Sigma_{22} &\leq \Sigma_x P(M(x) \geq s) (\Sigma_{j \geq x} \beta^{j-x}) = ([n/2]/m)^{-1} E(Y(n)) \\ &= O(E(Y(n))). \end{aligned}$$

To estimate Σ_{21} , first we write

$$(B.2) \quad \begin{aligned} \Sigma_{21} &\leq 2 \Sigma_{x_1 < x_2} P(M(x_i) \geq s, |K(x_i) - \bar{k}(x_i)| \leq (m - x_i)^{3/5}, \\ &\quad i = 1, 2; K(x_1) \neq K(x_2)) \\ &\quad + 2 \Sigma_{x_1 < x_2} P(|K(x_1) - \bar{k}(x_1)| > (m - x_1)^{3/5}) \\ &= \Sigma'_{21} + \Sigma''_{21}. \end{aligned}$$

Here (see Lemma 6)

$$(B.3) \quad \Sigma''_{21} \leq 2 \Sigma_{x_1 < x_2} \exp[-(m - x_1)^{1/5}/30] = O[n^2 \exp(-n^{1/10}/30)] = o(1)$$

and

$$(B.4) \quad \Sigma'_{21} = 2 \Sigma_{x_1 < x_2} \Sigma_{k_1, k_2} P(M(x_i) \geq s, K(x_i) = k_i, i = 1, 2),$$

where k_1, k_2 satisfy the conditions

$$(B.5) \quad \begin{aligned} |k_i - \bar{k}(x_i)| &\leq (m - x_i)^{3/5}, \quad i = 1, 2, \\ k_1 < k_1 + s &\leq x_1, \quad k_1 + s < k_2 < k_2 + s \leq x_2. \end{aligned}$$

Now, according to an argument in the proof of Lemma 4, $M(x_i) \geq s, K(x_i) = k_i, i = 1, 2$, iff (a) no key $\leq x_i$ has chosen the cell $k_i, i = 1, 2$, and (b) all the cells $k_i + 1, \dots, l_i = k_i + s$ and exactly $x_i - l_i$ cells to the right from the cell l_i are chosen by the keys $\leq x_i, i = 1, 2$. To evaluate the probability $P(M(x_i) \geq s, K(x_i) = k_i, i = 1, 2)$, observe that each allocation of the keys $\leq x_2$ satisfying the conditions (a) and (b) is achieved as follows. Consider the sets of cells $A_1 = \{1, \dots, k_1 - 1\}, |A_1| = k_1 - 1, A_2 = \{k_1 + 1, \dots, l_1\}, |A_2| = l_1 - k_1, A_3 = \{l_1 + 1, \dots, k_2 - 1\}, |A_3| = k_2 - 1 - l_1$, and $A_4 = \{k_2 + 1, \dots, l_2\}, |A_4| = l_2 - k_2$. First, we choose a set $A_5, |A_5| = x_2 - l_2$, among the cells $\geq l_2 + 1$ (the number of ways is $\binom{m - l_2}{x_2 - l_2}$). Second, we choose the set $A_6, |A_6| = x_1 - l_1$, out of $A_3 \cup A_4 \cup A_5$. If $|A_6 \cap (A_4 \cup A_5)| = j$, then the number of ways is

$$\binom{|A_4 \cup A_5|}{j} \binom{|A_3|}{x_1 - l_1 - j} = \binom{x_2 - k_2}{j} \binom{k_2 - 1 - l_1}{x_1 - l_1 - j}.$$

Third, we allocate the keys $\leq x_1$ on the set $A_7 = A_1 \cup A_2 \cup A_6, |A_7| = x_1 - 1$, so that the cells of $A_8 = A_2 \cup A_6, |A_8| = x_1 - k_1$, are nonempty. Finally, fourth, we allocate the keys $x_1 + 1, \dots, x_2$ on the set $A_9 = \bigcup_{i=1}^5 A_i \cup \{k_1\}, |A_9| = x_2 - 1$, so that the cells of $A_{10} = (A_4 \cup A_5) \setminus A_6, |A_{10}| = x_2 - k_2 - j$, are nonempty.

Hence,

$$\begin{aligned}
 &P(M(x_i) \geq s, K(x_i) = k_i, i = 1, 2) \\
 &= m^{-x_2} \binom{m - l_2}{x_2 - l_2} \left[\sum_j \binom{x_2 - k_2}{j} \binom{k_2 - 1 - l_1}{x_1 - l_1 - j} \right. \\
 &\quad \left. \times N(x_2 - x_1, x_2 - 1, x_2 - k_2 - j) \right] \\
 &\quad \times N(x_1, x_1 - 1, x_1 - k_1),
 \end{aligned}
 \tag{B.6}$$

where $N(u, v, w)$ is the number of ways to allocate u balls among v cells, so that some w fixed cells are nonempty.

The rest is simple. Since $|k_2 - \bar{k}(x_2)| \leq (m - x_2)^{3/5}$, we can use (3.15) to obtain from (B.6) that, for k_1, k_2 satisfying the conditions (B.4),

$$\begin{aligned}
 &P(M(x_i) \geq s, K(x_i) = k_i, i = 1, 2) \\
 &\sim (1 - e^{-\lambda_2})^s m^{-x_2} \binom{m - k_2}{x_2 - k_2} \left[\sum_j \binom{x_2 - k_2}{j} \binom{k_2 - 1 - l_1}{x_1 - l_1 - j} \right. \\
 &\quad \left. \times N(x_2 - x_1, x_2 - 1, x_2 - k_2 - j) \right] \\
 &\quad \times N(x_1, x_1 - 1, x_1 - k_1) \\
 &= (1 - e^{-\lambda_2})^s P(M(x_1) \geq s; K(x_i) = k_i, i = 1, 2),
 \end{aligned}
 \tag{B.7}$$

$$\lambda_2 = x_2/m.$$

Besides, by Lemma 7,

$$P(M(x) \geq x) \sim (1 - e^{-\lambda})^s, \quad \lambda = x/m,$$

uniformly over $n_1 \leq x \leq n_2$. Combining (B.4) and (B.7) we obtain then

$$\begin{aligned}
 &\Sigma_{21} \leq (1 + o(1)) 2 \Sigma_{x_1 < x_2} (1 - e^{-\lambda_2})^s \Sigma_{k_1, k_2} P(M(x_1) \geq s; \\
 &\quad K(x_i) = k_i, i = 1, 2) \\
 &\leq (1 + o(1)) 2 \Sigma_{x_1 < x_2} (1 - e^{-\lambda_2})^s P(M(x_1) \geq s) \\
 &\leq (1 + o(1)) [\Sigma_x (1 - e^{-\lambda})^s]^2 \sim E^2(Y(n)).
 \end{aligned}
 \tag{B.8}$$

Since $E(Y(n)) \rightarrow \infty$, the relations (B.1)–(B.3) and (B.8) imply that

$$E(Y^2(n)) \leq [1 + o(1)] E^2(Y(n)).$$

It remains to observe that $E^2(Y(n)) \leq E(Y^2(n))$. \square

Acknowledgment. I am grateful to the referee for many helpful suggestions regarding the presentation of the paper.

REFERENCES

- [1] AJTAI, M., FREDMAN, M. and KOMLÓS, J. (1978). There is no fast single hashing algorithm. *Inform. Process. Lett.* **7** 270–273.
- [2] BILLINGSLEY, P. (1979). *Probability and Measure*. Wiley, New York.
- [3] BLAKE, I. F. and KONHEIM, A. G. (1977). Big buckets are (are not) better! *J. Assoc. Comput. Mach.* **24** 591–606.
- [4] CHEN, W.-C. and VITTER, J. S. (1983). Analysis of early-insertion standard coalesced hashing. *SIAM J. Comput.* **12** 667–676.
- [5] CHEN, W.-C. and VITTER, J. S. (1984). Analysis of new variants of coalesced hashing. *ACM Trans. Database Systems* **9** 616–645.
- [6] DEVROYE, L. (1982). A note on the average depth of tries. *Computing* **28** 367–371.
- [7] DEVROYE, L. (1984). A probabilistic analysis of the height of tries and of complexity of triesort. *Acta Inform.* **21** 229–237.
- [8] DEVROYE, L. (1985). The expected length of the longest probe sequence for bucket searching when the distribution is not uniform. *J. Algorithms* **6** 1–19.
- [9] FAGIN, R., NIEVERGELT, J., PIPPENGER, N. and STRONG, H. R. (1979). Extendible hashing—a fast access method for dynamic files. *ACM Trans. Database Systems* **4** 315–344.
- [10] FLAJOLET, PH. and STEYAERT, J. M. (1982). A branching process arising in dynamic hashing, trie searching and polynomial factorization. *Proc. Ninth ICALP Colloquium. Lecture Notes in Computer Sci.* **140** 239–251. Springer, Berlin.
- [11] GONNET, G. H. (1981). Expected length of the longest probe sequence in hash code searching. *J. Assoc. Comput. Mach.* **28** 289–304.
- [12] GREENE, D. H. and KNUTH, D. E. (1982). *Mathematics for the Analysis of Algorithms*, 2nd ed. Birkhäuser, Boston.
- [13] GUIBAS, L. J. (1978). The analysis of hashing techniques that exhibit k -ary clustering. *J. Assoc. Comput. Mach.* **25** 544–555.
- [14] GUIBAS, L. J. and SZEMEREDI, E. (1978). The analysis of double hashing. *J. Comput. Sci.* **16** 226–274.
- [15] KNOTT, G. D. (1984). Direct chaining with coalescing lists. *J. Algorithms* **5** 7–21.
- [16] KNUTH, D. E. (1973). *The Art of Computer Programming* **3**. Addison-Wesley, Reading, Mass. (Section 6.4).
- [17] KOLCHIN, V. F. (1982). *Random Mappings*. Nauka, Moscow. (In Russian.)
- [18] KOLCHIN, V. F., SEVAST'YANOV, B. A. and CHISTYAKOV, V. P. (1978). *Random Allocations*. Wiley, New York.
- [19] KONHEIM, A. G. and WEISS, B. (1966). An occupancy discipline and applications. *SIAM J. Appl. Math.* **14** 1266–1274.
- [20] MENDELSON, H. (1982). Analysis of extendible hashing. *IEEE Trans. Software Engrg.* **8** 611–619.
- [21] PITTEL, B. (1985). Asymptotical growth of a class of random trees. *Ann. Probab.* **13** 414–427.
- [22] PITTEL, B. (1986). Paths in a random digital tree: limiting distributions. *Adv. in Appl. Probab.* **18** 139–155.
- [23] PITTEL, B. (1987). Linear probing: The probable largest search time grows logarithmically with the number of records. *J. Algorithms* **8** 1–14.
- [24] VITTER, J. S. (1981). A shared-memory scheme for coalesced hashing. *Inform. Process. Lett.* **13** 77–79.
- [25] VITTER, J. S. (1983). Analysis of the search performance of coalesced hashing. *J. Assoc. Comput. Mach.* **30** 231–258.
- [26] YAO, A. C. (1980). A note on the analysis of extendible hashing. *Inform. Process. Lett.* **11** 84–86.

DEPARTMENT OF MATHEMATICS
THE OHIO STATE UNIVERSITY
COLUMBUS, OHIO 43210