# Towards a Scalable and Robust DHT

Baruch Awerbuch
Department of Computer Science
Johns Hopkins University
3400 N. Charles Street
Baltimore, MD 21218, USA
baruch@cs.jhu.edu

Christian Scheideler
Institute for Computer Science
Technical University of Munich
Boltzmannstr. 3
85748 Garching, Germany
scheideler@in.tum.de

Regular Submission

### Abstract

The problem of scalable and robust distributed data storage has recently attracted a lot of attention. A common approach in the area of peer-to-peer systems has been to use a distributed hash table (or DHT). DHTs are based on the concept of virtual space. Peers and data items are mapped to points in that space, and local-control rules are used to decide, based on these virtual locations, how to interconnect the peers and how to map the data to the peers.

DHTs are known to be highly scalable and easy to update as peers enter and leave the system. It is relatively easy to extend the DHT concept so that a constant fraction of faulty peers can be handled without any problems, but no DHT is known that can even survive a small number of adversarial peers.

The biggest threats appear to be join-leave attacks (i.e., adaptive join-leave behavior by the adversarial peers to take over certain parts of the network or isolate honest peers) and attacks on the data management level (i.e., adaptive insert and lookup attacks by the adversarial peers to achieve a highly imbalanced request or data load) against which no provably robust mechanisms are known so far. We show, on a high level, that both of these threats can be handled in a scalable manner, even if a constant fraction of the peers in the system is adversarial, demonstrating that open systems for scalable distributed data storage that are robust against even massive adversarial behavior are feasible.

# 1 Introduction

In a distributed storage system, information is distributed among multiple storage devices, simply called *nodes* in the following. To provide a basic lookup service, the following operations have to be implemented:

- Insert($d$): this inserts data item $d$ into the system.

- Lookup(name): this returns any data item $d$ with $\text{Name}(d) = \text{name}$, if it exists.

Once a distributed storage system becomes large enough, one also has to deal with nodes leaving and joining the system, simply because storage devices may break down or new devices have to be added in order to maintain a desired service quality. Hence, two more operations are needed:

- Join($v$): node $v$ joins the system

- Leave($v$): node $v$ leaves the system

How can these four operations be implemented to obtain a robust and scalable distributed storage system? The most prominent approach studied in the research community is to implement a distributed hash table, or DHT. DHTs have been realized in various contexts including server-based systems such as Akamai and peer-to-peer systems such as Chord [35], CAN [25], Pastry [9], and Tapestry [38]. Most of the DHT-based systems are based on two influential papers: a paper by Plaxton, Rajaraman, and Richa on locality-preserving data management in distributed environments [23] and a paper by Karger, Lehman, Leighton et al. on consistent hashing and web caching [14]. The consistent hashing approach is a very simple and elegant approach that is based on the following rule:

Suppose that we have two random functions $f$ and $g$. The function $f$ maps the nodes randomly to real numbers in $[0, 1)$, and the function $g$ maps the data items randomly to real numbers in $[0, 1)$. Every data item $d$ is stored at the node $v$ with a minimum distance between $f(v)$ and $g(d)$ (viewing $[0, 1)$ here as a ring). It turns out that this rule has several nice features [14]:

- On expectation, every node has the same load,

- when integrating a node into or removing a node from a system of $n$ nodes, only a $1/n$-fraction of the data has to be replaced on expectation, and

- when storing $O(\log n)$ copies of each data item, the system is robust against a constant fraction of faulty nodes.

However, all of these properties only hold if *f and g are random* and the names selected for the nodes and the data items are *independent of f and g*. Unfortunately, in a DHT, $g$ must be a fixed hash function because otherwise it would not be possible to compute the location of the data items. Hence, it is easy for the adversary to generate many data items that all need to be stored at the same node, even if it cannot invert $g$. It just has to try sufficiently many names. In fact, millions of names can be quickly tested with hash functions like SHA-1.

But even a truly random mapping does not protect against adversarial attacks. Suppose, for example, that every new node is mapped uniformly at random to a point in $[0, 1)$. If the adversary wants to overpopulate a certain area of the $[0, 1)$ space, say some interval $I$, it just needs to execute sufficiently many join and leave operations in which it keeps all nodes that made it into $I$ in the system and removes all others for another join attempt.

DHT constructions are very vulnerable to a node and data imbalance in the virtual space since this will effect their correctness and scalability and can be used for denial-of-service attacks. For example, in systems like Chord join-leave attacks are possible with a logarithmic number of adversarial nodes that can isolate

any honest node from the rest of the system [3]. Can we design simple protocols for the operations join, leave, insert and lookup that are *provably* robust against these attacks without restricting the openness of the system?

In this paper we show that, on a high level, this is possible. More precisely, we will show that there are scalable join and leave protocols so that for a polynomial number of join and leave requests the nodes will be evenly distributed in the $[0, 1)$ space, with high probability[1], and the honest and adversarial nodes will be well-spread so that quorums of size $O(\log n)$ can be formed to wash out any adversarial behavior violating the protocols by simple majority decision. Moreover, we will show that there are scalable and robust insert and lookup protocols so that for a polynomial number of attempts the adversary will not manage to find data names so that it can create a high request or load imbalance in the system.

## 1.1 Our contributions

Next we give a detailed description of our contributions. For simplicity, we assume that the number of honest nodes in the system will only change by a constant factor over time. In this way, notions like "in a polynomial number of rounds" and "with high probability" are well-defined. Let $n$ be the maximum number of honest nodes in the system at any time and let $\epsilon n$ be the maximum number of nodes that the adversary can have in the system at any time. I.e., we assume that the adversary has bounded resources, but apart from that the adversary can do what it likes, such as choosing any names it likes for the data items and the nodes.

**Join-leave attacks**

We first focus on making a DHT robust against join-leave attacks. More precisely, we consider the following scenario. There are $n$ white (or honest) nodes and $\epsilon n$ black (or adversarial) nodes for some fixed constant $\epsilon < 1$. There is a rejoin request that, when applied to node $v$, lets $v$ first leave the system and then join it again from scratch. The leaving is done by simply removing $v$ from the system and the joining is done with the help of a join operation to be specified by the system. We assume that (after an initial joining of all the nodes) the sequence of rejoin requests is controlled by an adversary, which is a typical assumption in the analysis of online algorithms. However, whereas rejoin requests of white nodes have to be specified *in advance*, rejoin requests of black nodes can be specified in an *adaptive* manner. Thus, our adversarial model is a mix of an oblivious and an adaptive adversary. More precisely, the adversary has to fix a sequence $\sigma$ of rejoin requests of the white nodes in advance. Afterwards, we proceed in rounds. In each round, in which the first $i$ of the rejoin requests in $\sigma$ have already been executed, the adversary can decide either to initiate a rejoin request for any one of the black nodes or to initiate the $(i + 1)$st rejoin request in $\sigma$. Our goal is to find an *oblivious* join strategy, i.e., a strategy that cannot distinguish between the white and black nodes, so that for *any* adversarial strategy above the following two conditions can be preserved for every interval $I \subseteq [0, 1)$ of size at least $(c \log n)/n$ for a constant $c > 0$ and any polynomial number of rounds in $n$:

- *Balancing condition:* $I$ contains $\Theta(|I| \cdot n)$ nodes.

- *Majority condition:* the white nodes in $I$ are in the majority.

It is not difficult to see that the brute-force strategy of giving every node a new random place whenever a node rejoins will achieve the stated goal, with high probability, but this would be a very expensive strategy. The challenge is to find a join operation that needs as little randomness and as few rearrangements as possible to satisfy the two conditions. Fortunately, there is such a strategy, called the *cuckoo rule*. We first introduce some notation, and then we describe the strategy.

---

[1]By "with high probability" or "w.h.p." we mean a probability of at least $1 - 1/n$ where $n$ is the size of the system.

In the following, a *region* is an interval of size $1/2^r$ in $[0,1)$ for some positive integer $r$ that starts at an integer multiple of $1/2^r$. Hence, there are exactly $2^r$ regions of size $1/2^r$. A *k-region* is a region of size (closest from above to) $k/n$, and for any point $x \in [0,1)$, the $k$-region $R_k(x)$ is the unique $k$-region containing $x$.

**Cuckoo rule:** If a new node $v$ wants to join the system, pick a random $x \in [0,1)$. Place $v$ into $x$ and move all nodes in $R_k(x)$ to points in $[0,1)$ chosen uniformly and independently at random (without replacing any further nodes).

Our first main result is summarized in the following theorem.

**Theorem 1.1** *For any constants $\epsilon$ and $k$ with $\epsilon < 1 - 1/k$, the cuckoo rule with parameter $k$ satisfies the balancing and majority conditions for a polynomial number of rounds, with high probability, for any adversarial strategy within our model. The inequality $\epsilon < 1 - 1/k$ is sharp as counterexamples can be constructed otherwise.*

Hence, a constant $k > 1$ would be sufficient to prevent *adaptive* join-leave attacks of a *constant fraction* of adversarial peers. Thus, it is remarkably easy to defend an open distributed storage system against even massive join-leave attacks. The cuckoo rule allows us to use local quorum strategies in order to wash out adversarial behavior violating the protocols.

**Lookup and insert attacks**

Our basic strategy to handle attacks on the data layer is to use $2c - 1 = \Theta(\log n)$ one-way hash functions mapping each data item to $2c - 1$ points in the $[0,1)$ space. These hash functions are fixed but have certain expansion properties to make it hard for the adversary to create bad sets of insert or lookup requests. In order to achieve an even load balancing of the requests and the data items, we use the majority trick of Upfal and Wigderson [37]: for each insert request, store copies of the data item in at least $c$ of the $2c - 1$ locations, and for each lookup request, access at least $c$ locations of the data item. This indeed suffices for the correct implementation of these requests because if the copies are stored in a reliable way, then the lookup operation will always retrieve at least one copy of the data item.

Given this basic scheme, we present a scalable dynamic overlay network and robust protocols for the insert and lookup operations. In the following, $U$ may represent the space of all names or the names that the adversary can sample in a polynomial number of time steps.

**Theorem 1.2** *For any collection of lookup requests for data items out of a set $U$ of polynomial size with one request per node, the lookup protocol can serve all of these requests correctly and reliably in polylogarithmic time so that each node is passed by at most $O(\log^5 n)$ requests.*

Notice that the upper bound is *guaranteed* for *any adaptively* chosen set of lookup requests, including data items with multiple lookup requests. The hash functions just need to be selected so that they satisfy certain expansion properties. Also, notice that the adversary cannot modify or delete a request in transit since we use quorum strategies. Certainly, the bound is still too high for practical purposes, but the best bound previously known for adaptively chosen lookup requests in overlay networks is the trivial linear bound. In light of this, our result is an exponential improvement, and an interesting problem for future research will certainly be whether further improvements are possible. For insert requests we obtain a similar result.

**Theorem 1.3** *For any collection of insert requests for data items out of a set $U$ of polynomial size with one request per node, the insert protocol can serve all of these requests in polylogarithmic time so that each node is passed by at most $O(\log^5 n)$ requests. Moreover, the maximum amount of copies to be stored by any node to serve all of the requests is bounded by $O(\log^2 n)$.*

The $O(\log^2 n)$ bound is asymptotically optimal if $n$ insert requests are executed since we use quorums of size $O(\log n)$ and for $n$ requests every quorum is responsible for $O(\log^2 n)$ copies on average.

**Prerequisites**

An important prerequisite for our join and leave operations to work correctly is a distributed random number generator that can generate an unbiased random number even under the influence of a constant fraction of adversarial nodes, and an important prerequisite for our insert and lookup operations to work correctly is that one-way hash functions with certain expansion properties are available so that the adversary has no other way then sampling names in order to design malicious collections of names. It is commonly believed that one-way hash functions exist though no formal proof has been found yet. But we will at least prove that random hash functions will have expansion properties, w.h.p., that are good enough for our results to hold. A distributed random number generator sufficient for our purposes can be built on top of existing verifiable-secret-sharing (VSS) protocols. In fact, we have developed a distributed random number generator based on the $\lfloor \frac{n-1}{4} \rfloor$-VSS protocol in [12] that we include in the appendix for completeness. This generator may be used by the $\Theta(\log n)$ nodes in any quorum in our DHT to correctly and efficiently generate random numbers under the assumption that the honest nodes are in a sufficient majority in that quorum.

## 1.2 Previous work

In the area of peer-to-peer systems, work on robustness in the context of overlay network maintenance has mostly focused on how to handle a large fraction of faulty nodes (e.g., [36, 2, 27]) or churn, that is, peers frequently enter and leave the system (e.g., [26, 15]). However, none of these approaches can protect a DHT against the join-leave attacks considered in this paper because just assigning a random or pseudo-random point to each new node (by using some random number generator or cryptographic hash function) does not suffice to preserve the balancing and majority conditions [3]. People in the peer-to-peer community are aware of the danger of these attacks [8, 6] and various solutions have been proposed that may help thwart these attacks in practice [4, 5, 21, 31, 30, 32] but until recently no mechanism was known that can *provably* cope with these attacks without sacrificing the openness of the system.

One such mechanism, that can only cope with a linear number of adversarial join requests, was proposed in [11]. The first mechanism that was shown to preserve randomness in the system under adaptive adversarial behavior for a polynomial number of adversarial join-leave requests uses *random* node IDs and enforces a *limited* lifetime on every node in the system, i.e., every node *has* to reinject itself after a certain amount of time steps [3]. However, this leaves the system in a hyperactive mode that may unnecessarily consume resources that could be better used for other purposes. Ideally, one would like to use *competitive* strategies. That is, the resources consumed by the mixing mechanism should scale with the join-leave activity of the system. Recently, it was shown that for a pebble-shuffling game this can be achieved [28]. In this game, there are $n$ white pebbles and $\epsilon n$ black pebbles for some fixed constant $\epsilon < 1$. The pebbles are laid out on a ring and the black pebbles can join and leave the ring in an adaptive adversarial fashion. It was shown in [28] that a simple protocol called $k$-rotation exists that can preserve the majority condition with high probability. That is, for any sequence of $\Theta(\log n)$ pebbles along the ring, a majority of them is white. However, the result in [28] cannot be taken over to a virtual space setting as adversarial strategies exist for which the $k$-rotation rule cannot satisfy the balancing condition. Therefore, we had to design a new strategy, which we called the cuckoo rule.

Also attacks on the data management layer have been considered in the past. Most of the work considers the flash crowd scenario in which many peers want to access the same information at the same time. When using a pure DHT design, this can lead to severe bottlenecks. To remove these bottlenecks, various caching strategies have been proposed. Among them are CoopNet [22], Backslash [33], PROOFS [34] in the sys-

tems community and [20, 29] in the theory community. However, being able to handle flash crowds is not sufficient to survive the attacks considered in this paper because much worse than having many requests to the *same* data item is to have many requests to *different* data items residing at the *same* peer. Standard combining or caching strategies do not work here, so new strategies are needed. It turns out that, interestingly, work on deterministic simulations of CRCW PRAMs comes to the rescue here. This was pioneered by Mehlhorn and Vishkin [18] and further developed in a series of papers [1, 13, 16, 37]. The basic ideas behind our insert and lookup protocols are based on these results though we had to significantly extend their techniques since we need strategies that work in a dynamic well-structured overlay network whereas the PRAM results above only consider fixed complete networks or networks with expander-like properties.

## 1.3 Organization of the paper

In Section 2 we will prove Theorem 1.1 and in Section 3 we will prove Theorems 1.2 and 1.3. The paper ends with conclusions. Some of the proofs can be found in the appendix.

# 2 Analysis of the cuckoo rule

We assume that we first have $n$ rounds in which all the white nodes are placed in the system and then we have $\epsilon n$ rounds in which all the black nodes are placed in the system, using the cuckoo rule in each round. Afterwards, the adversary starts with its rounds.

Consider any fixed round $t_0$ after the initial $(1 + \epsilon)n$ rounds. Recall that a *region* is an interval of size $1/2^r$ in $[0, 1)$ for some positive integer $r$ that starts at an integer multiple of $1/2^r$. Let $\hat{R}$ be any fixed region of size $(c \log n) \cdot k/n$, for some constant $c$, for which we want to check the balancing and majority conditions over polynomial in $n$ many steps. Thus, $\hat{R}$ contains exactly $c \log n$ many $k$-regions. The *age* of a $k$-region is the difference between $t_0$ and the first round back in time when a new node was placed into it (and all old nodes got evicted) and the age of $\hat{R}$ is defined as the sum of the ages of its $k$-regions. A node in $\hat{R}$ is called *new* if it was placed in $\hat{R}$ when it joined the system, and otherwise it is called *old*. Before we start with our analysis, we state some technical lemmas.

**Lemma 2.1 (Chernoff-Hoeffding Bounds [17])** *Suppose that $X_1, \ldots, X_n$ are independent random variables with $X_i \leq b$ for all $i$. Let $X = \sum_{i=1}^{n} X_i$ and $\mu = \mathrm{E}[X]$. Then it holds for all $\epsilon \geq 0$ that $\Pr[X \geq (1 + \epsilon)\mu] \leq e^{-\epsilon^2 \mu/(2b(1+\epsilon))}$ and for all $0 \leq \epsilon \leq 1$ that $\Pr[X \leq (1 - \epsilon)\mu] \leq e^{-\epsilon^2 \mu/(2b)}$.*

**Lemma 2.2** *Suppose that $X_1, \ldots, X_n$ are independent random variables on $\mathrm{I\!N}$ with $\Pr[X_i = t] = p(1 - p)^{t-1}$ for every $t \in \mathrm{I\!N}$ for some fixed $0 < p < 1$. Let $X = \sum_{i=1}^{n} X_i$ and $\mu = \mathrm{E}[X]$. Then it holds for every $\epsilon > 0$ that $\Pr[X \geq (1 + \epsilon)\mu] \leq e^{-\epsilon^2 n/2(1+\epsilon)}$ and for all $0 \leq \epsilon \leq 1$ that $\Pr[X \leq (1 - \epsilon)\mu] \leq e^{-\epsilon^2 n/2}$.*

**Proof.** Can be easily derived from Lemma 2.1. □

**Theorem 2.3 (Method of Bounded Martingale Differences [10])** *Let $X_1, \ldots, X_n$ be an arbitrary set of random variables and let $f$ be a function satisfying the property that for each $i \in \{1, \ldots, n\}$ there is a non-negative $c_i$ such that $|\mathrm{E}[f \mid X_1, \ldots, X_i] - \mathrm{E}[f \mid X_1, \ldots, X_{i-1}]| \leq c_i$. Then*

$$\Pr[f \geq \mathrm{E}[f] + t] \leq e^{-t^2/(2 \sum_{i=1}^{n} c_i^2)} \quad and \quad \Pr[f \leq \mathrm{E}[f] - t] \leq e^{-t^2/(2 \sum_{i=1}^{n} c_i^2)}.$$

Note that the $X_i$'s in this theorem are *not* required to be independent. Now we can start with the proof of Theorem 1.1. The next lemma follows directly from the cuckoo rule because every $k$-region can have at most one new node at any time.

**Lemma 2.4** *At any time, $\hat{R}$ contains at most $c \log n$ new nodes.*

In order to bound the number of old nodes in $\hat{R}$, we first have to bound the age of $\hat{R}$ (Lemma 2.5). Then we bound the maximum number of nodes in a $k$-region (Lemma 2.6) and use this to bound the number of evicted white and black nodes in a certain time interval (Lemma 2.7). After that, we can combine all lemmas to bound the number of old white and black nodes in $\hat{R}$ (Lemma 2.8). For simplicity, we assume in this extended abstract that the system is sufficiently old. It is not hard to extend our analysis so that the majority and balancing conditions can be shown to hold for any fixed time step, w.h.p., as long as $\epsilon < 1 - 1/k$.

**Lemma 2.5** *At any (sufficiently large) time, $\hat{R}$ has an age within $(1 \pm \delta)(c \log n)n/k$, with high probability, where $\delta > 0$ can be made arbitrarily small depending on $c$.*

**Proof.** See the appendix for a proof. □

**Lemma 2.6** *For any $k$-region $R$ in $\hat{R}$ it holds at any time that $R$ has at most $O(k \log n)$ nodes, with high probability.*

**Proof.** See the appendix for a proof. □

Next we bound the number of white and black nodes that are evicted in a certain time interval.

**Lemma 2.7** *For any time interval $I$ of size $T = (\gamma/\epsilon) \log^3 n$, the number of white nodes that are evicted in $I$ is within $(1 \pm \delta)T \cdot k$, with high probability, and the number of black nodes that are evicted in $I$ is within $(1 \pm \delta)T \cdot \epsilon k$, with high probability, where $\delta > 0$ can be made arbitrarily small depending on $\gamma$.*

**Proof.** We start with the proof for the white nodes. Consider any time interval $I$ consisting of $T$ rounds. For every $t \in \{1, \ldots, T\}$ let the random variable $X_t$ denote the number of white nodes evicted in the $t$th round of $I$, and let $X = \sum_{t=1}^{T} X_t$. Since there are $n$ white nodes in the system, it follows that $\mathrm{E}[X_t] = k$ for every $t$, no matter how the white nodes are distributed, and therefore $\mathrm{E}[X] = T \cdot k$. From Lemma 2.6 we know that every $X_t$ is at most $4 \log n$, w.h.p. Also, for any given $X_1, \ldots, X_{t-1}$, $\mathrm{E}[X_t \mid X_1, \ldots, X_{t-1}] = \mathrm{E}[X_t] = k$. Hence, it holds for the function $f(X_1, \ldots, X_n) = \sum_{i=1}^{n} X_i$ that

$$
\begin{aligned}
&|\mathrm{E}[f \mid X_1, \ldots, X_i] - \mathrm{E}[f \mid X_1, \ldots, X_{i-1}]| \\
&= \left| \mathrm{E}\left[ \sum_{j=i+1}^{n} (X_j \mid X_1, \ldots, X_i) \right] + \sum_{j=1}^{i} X_j - \left( \mathrm{E}\left[ \sum_{j=i}^{n} (X_j \mid X_1, \ldots, X_{i-1}) \right] + \sum_{j=1}^{i-1} X_j \right) \right| \\
&= |X_i - \mathrm{E}[X_i]| \leq 4 \log n
\end{aligned}
$$

Thus, the method of bounded martingale differences (Lemma 2.3) implies that, for any constant $\delta \geq 0$,

$$
\Pr[X \geq (1 + \delta)T \cdot k] \leq e^{-\delta^2 (T \cdot k)^2 / (2 \sum_{i=1}^{T} (4 \log n)^2)}
$$

which is polynomially small in $n$ if the constant in $T$ is sufficiently large. The same holds for $\Pr[X \leq (1 - \delta)T \cdot k]$, which proves the lemma for the white nodes. The proof for the black nodes is the same. □

Combining Lemmas 2.5 to 2.7, we obtain the following lemma.

**Lemma 2.8** *At any time, $\hat{R}$ has within $(1 \pm \delta)(c \log n) \cdot k$ old white nodes and within $(1 \pm \delta)(c \log n) \cdot \epsilon k$ old black nodes, with high probability.*

**Proof.** Consider any age distribution $t_1, \ldots, t_C$ for the $k$-regions $R_1, \ldots, R_C$ of $\hat{R}$, where $C = c \log n$. Let $T = (\gamma/\epsilon) \log^3 n$ be selected as in Lemma 2.7. Then it follows from Lemma 2.7 and the Chernoff bounds that $\hat{R}$ hat at least

$$\frac{(1+\delta)T \cdot k}{n/k} \sum_{i=1}^{C} \lfloor t_i/T \rfloor \geq \frac{(1+\delta)k^2}{n} \left( \sum_{i=1}^{C} t_i - C \cdot T \right)$$

white nodes and at most

$$\frac{(1-\delta)T \cdot k}{n/k} \sum_{i=1}^{C} \lceil t_i/T \rceil \leq \frac{(1-\delta)k^2}{n} \left( \sum_{i=1}^{C} t_i + C \cdot T \right)$$

white nodes, w.h.p. Since $\sum_{i=1}^{C} t_i$ is within $(1 \pm \delta')Cn/k$ according to Lemma 2.5, Lemma 2.8 follows for the white nodes.

The same calculations (with an additional $\epsilon$ factor) apply to the black nodes. $\qquad \square$

Combining Lemmas 2.4 and 2.8, we can now prove when the balancing and majority conditions are satisfied.

- **Balancing condition:** From Lemmas 2.4 and 2.8 it follows that every region $R$ of size $(c \log n)k/n$ has at least $(1 - \delta)(c \log n) \cdot k$ and at most $(1 + \delta)(c \log n + (c \log n)k + (c \log n)\epsilon k) = (1 + \delta)(c \log n)(1 + (1 + \epsilon)k)$ nodes, where the constant $\delta > 0$ can be made arbitrarily small. Hence, the regions are balanced within a factor of close to $(1 + \epsilon + 1/k)$.

- **Majority condition:** From Lemmas 2.4 and 2.8 it also follows that every region of size $(c \log n)k/n$ has at least $(1 - \delta)(c \log n) \cdot k$ white nodes and at most $(1 + \delta)(c \log n + (c \log n) \cdot \epsilon k)$ black nodes, w.h.p., where the constant $\delta > 0$ can be made arbitrarily small. These bounds are also tight in the worst case, which happens if the adversary focuses on a specific region $R$ of size $(c \log n)k/n$ and continuously rejoins with any black node outside of $R$. Hence, the adversary is not able to obtain the majority in any region of size $(c \log n)/n$ as long as $(c \log n)(\epsilon k + 1) < (c \log n) \cdot k$ which is true if and only if $\epsilon < 1 - 1/k$.

Hence, for $\epsilon < 1 - 1/k$ the balancing and majority conditions are satisfied, w.h.p., and this is sharp, which proves Theorem 1.1.

## 3 Robust insert and lookup protocols

In this section we present and analyze our robust insert and lookup protocols. These protocols are based on a dynamic de Bruijn graph and $2c - 1$ one-way hash functions with certain expansion properties. We first describe the dynamic de Bruijn graph and how to store data and how to route in it, and then we specify what kind of expansion properties the hash functions need to satisfy. Afterwards, we present the insert and lookup protocols.

### 3.1 The dynamic de Bruijn graph

In the classical $d$-dimensional de Bruijn graph, $\{0,1\}^d$ represents the set of nodes and two nodes $x, y \in \{0,1\}^d$ are connected by an edge if and only if there is a $b \in \{0,1\}$ so that $x = (x_1 \ldots x_d)$ and $y = (bx_1 \ldots x_{d-1})$ (i.e., $y$ is the result of a right shift of the bits in $x$ with the highest bit position taken by $b$) or $y = (x_2 \ldots x_d b)$. We aim to maintain the following dynamic variant of the de Bruijn graph (e.g., [20]):

Recall the definition of a region. For simplicity, we assume that the number of honest nodes in the system only deviates by a constant factor over time and that $n$ is the maximum number of honest nodes in

the system at any time. (We just need this so that we can focus on a fixed region size. Local-control update mechanisms for the region size (e.g., [3]) can be used if $n$ significantly changes over time.) We identify nodes by their points in $[0, 1)$. Given a node $v \in [0, 1)$, we define its *quorum region* $R_v$ as the unique region of size closest from above to $(\gamma \log n)/n$, $\gamma > 1$ fixed, that contains $v$. For any set of nodes $V \subset [0, 1)$, we require that every node $v \in V$ maintains connections to all nodes whose quorum regions contain a point in $\{v, v/2, (1 + v)/2, 2v \bmod 1, (2v - 1) \bmod 1\}$. Let us call the resulting graph $DB(V)$. The following lemma is easy to show.

**Lemma 3.1** *For any node set $V \subset [0, 1)$ that satisfies the balancing condition, $DB(V)$ has a diameter of $O(\log n)$ and a degree of $O(\log n)$ and every quorum region forms a clique of $\Theta(\log n)$ nodes.*

## 3.2 Reliable storage and routing

If also the majority condition holds, then the honest nodes are in the majority in each quorum region. This allows the honest nodes to wash out adversarial behavior violating the protocols by simple majority decision. In order to reliably store data, we therefore demand that copy $i$ of data item $x$ has to be stored in all nodes of the unique quorum region containing $h_i(x)$.

In order to route a message from a node $v \in [0, 1)$ to a point $w \in [0, 1)$ (representing a node or location of a copy) in a reliable way, we execute the following Route$(v, w)$ protocol:

Focus only on the first $\log n$ bits of the binary representation of $v$, denoted by $(v_1 v_2 \ldots v_{\log n})$ and forward the message from $R_v$ along the quorum regions containing the points $(v_2 v_3 \ldots v_{\log n} w_1)$, $(v_3 v_4 \ldots v_{\log n} w_1 w_2)$, and so on, until the quorum region containing the point $(w_1 w_2 \ldots w_{\log n})$ is reached, which also contains $w$. It is easy to check that this routing strategy can be performed along adjacent regions in $DB(V)$ and that the following lemma holds.

**Lemma 3.2** *If the balancing and majority conditions are satisfied, then Route$(v, w)$ combined with majority decision reliably routes a message in at most $\log n$ communication rounds between any pair of honest nodes $v$ and $w$ in the system.*

The only problem is the congestion caused by routing multiple messages. As a prerequisite for this we need suitable hash functions.

## 3.3 Robust hash functions

Next we specify the properties the $2c - 1$ hash functions $h_1, \ldots, h_{2c-1}$ have to satisfy for our protocols to work. Namely, Lemma 3.3 below will be used to keep the congestion low and Lemma 3.4 will be used to keep the data load evenly balanced.

For any $d \geq 1$, the $d$-dimensional Shuffle graph $SH(d)$ consists of $d + 1$ levels numbered from 0 to $d$. The node set of level $i$ is given as $V_i = \{0, 1\}^d$, and for every $0 \leq i < d$, every pair of nodes $v \in V_i$ and $w \in V_{i+1}$ is connected if and only if their binary representations satisfy $v = (v_1 v_2 \ldots v_d)$ and $w = (v_2 v_3 \ldots v_d b)$ for some bit $b \in \{0, 1\}$.

The Shuffle graph is related to the well-known Omega network. It is a leveled form of the de Bruijn graph and contains for every source $s \in V_0$ and destination $t \in V_d$ a unique path of length $d$ from $s$ to $t$. In fact, these paths represent the paths the packets will move along through $DB(V)$ when using the reliable routing strategy above. Notice that for every $1 \leq \ell \leq d$, $SH(d)$ contains $2^{d-\ell}$ disjoint $SH(\ell)$ graphs from level 0 to $\ell$. Let $\mathcal{SH}(\ell)$ denote their set.

Now, let $U$ be the universe of data items and consider the bipartite graph $G = (U, V_d, E)$ that contains an edge between any two nodes $u \in U$ and $v \in V_d$ if and only if there is an $i \in \{1, \ldots, 2c - 1\}$ so that $h_i(u) = v$ (or more precisely, the first $d$ bits in the binary representation of $h_i(u)$ are equal to $v$).

8

Let $V = \bigcup_{i=1}^{d} V_i$. For any edge set $F \subseteq U \times V$ let $U(F) = \{u \in U \mid \exists v \in V : (u, v) \in F\}$ and $V(F) = \{v \in V \mid \exists u \in U : (u, v) \in F\}$. We call $F$ a *k-bundle* if

1. $|F| = k|U(F)|$ and

2. there is an injective mapping $f : U(F) \to V_0$ and an injective mapping $g_u : F_u \to \{1, \ldots, 2c - 1\}$ for every node $u \in U(F)$ ($F_u$ is the edge set of $u$ in $F$) so that for every edge $e = \{u, v\} \in F_u$ the path from $f(u) \in V_0$ to $h_{g_u(e)}(u) \in V_d$ in $SH(d)$ passes through $v$.

$F$ is called a *k-bundle of* $U(F)$, and we define $\Gamma_F(U(F)) = V(F)$. A $k$-bundle is called $\sigma$-*sparse* if for any subset $V'_\ell \subseteq V_\ell$ representing all nodes in level $\ell$ of some graph in $\mathcal{SH}(\ell)$ it holds that $|V(F) \cap V'_\ell| \leq \sigma|V'_\ell|$.

The graph $G$ above is called a $(\sigma, k)$-*expander* if for every $1 \leq \ell \leq d$ it holds that for any $S \subseteq U$ with $|S| \leq n$ and any $\sigma$-sparse $k$-bundle $F$ of $S$, $|\Gamma_F(S)| \geq |S|$. In the following, let $|U| = m$.

**Lemma 3.3** *If $c \geq 6 \log m$ and $m \geq n^3$, $\sigma \leq 1/(16 \log n)$ and the functions $h_1, \ldots, h_{2c-1}$ are chosen uniformly and independently at random, then $G$ is a $(\sigma, c/2)$-expander with high probability.*

**Proof.** See the appendix for a proof. □

Although it seems that we only proved an expansion property for a static graph, notice that as long as $n$ is the maximum number of honest nodes in the system at any time, we just need to consider all $SH(d)$ with $d \leq \log n$ to capture all possible routing strategies for the nodes in $DB(V)$ at any time. Since the probability bound in Lemma 3.3 is polynomially small, the hash functions can also satisfy all of these $SH(d)$ together with high probability.

Next we prove another expansion result that is needed for an even load balancing of the data. We call $G$ a $(\lambda, k, \sigma)$-*expander* if for any $S \subseteq U$ with $|S| \leq \sigma n/c$ and any $k$-bundle $F$ of $S$ with $V(F) \subseteq V_d$ it holds that $|\Gamma_F(S)| \geq \lambda c|S|$.

**Lemma 3.4** *Let $0 < \lambda < 1/2$ be any constant. Then it holds that for any $c \geq 6 \log m$ and $\sigma \leq 1/(\lambda (16e^{2\lambda})^{1/(1-2\lambda)})$ that if the functions $h_1, \ldots, h_{2c-1}$ are chosen uniformly and independently at random, then $G$ is a $(\lambda, c/2, \sigma)$-expander with high probability.*

**Proof.** See the appendix for a proof. □

We are now ready to describe and analyze our lookup and update protocols.

## 3.4 The lookup protocol

Suppose, first, that we have $n$ lookup requests, one per peer, and when routing $2c - 1$ packets for each of these requests through the dynamic de Bruijn network, one for each hash function, every quorum region will be passed by at most $O(c \log n)$ packets for different data items in each level (i.e., multiple packets for the same data item are counted as one). This is the best we can expect, i.e., the bound holds if the paths of the packets are evenly distributed. Then a combination of the growing rank protocol [19] with Ranade's trick of combining packets to the same data item [24] can be used to route all packets to their destinations in $O(c \log^2 n)$ communication rounds, w.h.p., where in each communication round every quorum region can forward at most one packet according to the reliable routing strategy. (The reason why we have a probabilistic result here is that the growing rank protocol is a randomized scheduling protocol. Also a deterministic scheduling protocol could be used, so that we have a completely deterministic lookup protocol, but that would increase the runtime by another $\log n$ factor in the worst case.)

Unfortunately, the adversary can generate lookup requests so that some quorum regions get congested. To handle this problem, we allow several attempts for each lookup request and use a simple threshold mechanism to control the congestion in each attempt. More precisely, we repeat the following procedure for sufficiently many rounds:

For each of the remaining lookup requests, $2c - 1$ packets are sent out, one for each of its $2c - 1$ destinations. The packets are routed as above with the difference that the quorum regions store the set of data items for which packets have already passed through them. If this counter exceeds a threshold of $\gamma c \log^2 n$ for some level $\ell$, then all further packets reaching that region in level $\ell$ are discarded. A lookup request is successful in that round if at least $c$ of its packets reach their destinations.

Each of these attempts will take at most $O(c \log^3 n)$ communication rounds, w.h.p., but the question is, how many attempts are needed until all of the lookup requests are successful, no matter which data items have been selected for them. The following theorem gives an answer to this, which implies Theorem 1.2.

**Theorem 3.5** *For any set of $n$ lookup requests out of a set $U$ of polynomial size with one request per node, the lookup protocol can serve all requests in a guaranteed number of $O(\log n)$ attempts and needs at most $O(c \log^4 n)$ communication rounds altogether, w.h.p.*

**Proof.** See the appendix for a proof. □

## 3.5 The insert protocol

The insert protocol is similar to the lookup protocol, with the difference that we play a collision game at the destinations (first used in [7] in the context of PRAM simulations) to keep the number of copies each node has to store as small as possible. The protocol proceeds in rounds. In each round, every remaining insert request generates a packet for each of its $2c - 1$ destinations. Besides the normal packets, also special counter packets are sent through the network whose sole purpose is to count how many packets are destined for each of the quorum regions. These counter packets can be combined during the routing so that they can be handled in $O(c \log^2 n)$ communication rounds, w.h.p. [24]. Once each quorum region knows its count, it will send acknowledgement packets backwards that will split appropriately in order to inform all the sources. The sources then discard all packets whose destination regions want to be reached by more than $\gamma c \log n$ packets in that round, where $\gamma > 1$ is a sufficiently large constant. Also, a normal packet may be discarded in transit if it reaches a region for some level $\ell$ with a congestion of more than $\gamma c \log^2 n$ in that level.

This insert protocol has the following performance, which implies Theorem 1.3.

**Theorem 3.6** *For any set of $n$ insert requests out of a set $U$ of polynomial size with one request per node, the insert protocol can serve all requests in a guaranteed number of $O(\log n)$ attempts and needs at most $O(c \log^4 n)$ communication rounds altogether, w.h.p. Furthermore, every node has to store at most $O(c \log n)$ copies.*

**Proof.** See the appendix for a proof. □

# 4 Conclusions

In this paper we showed that, on a high level, a scalable DHT can be designed that is provably robust against adaptive adversarial join-leave attacks as well as insert and lookup attacks. Certainly, low-level protocols still have to be designed for our operations that work well and correctly in an asynchronous environment. We believe that designing such protocols is possible though their design and formal correctness proofs may require a significant effort. Especially coping with low-level attacks such as network-layer denial-of-service attacks will be a major challenge.

# References

[1] H. Alt, T. Hagerup, K. Mehlhorn, and F.P. Preparata. Deterministic simulation of idealized parallel computers on more realistic ones. *SIAM Journal on Computing*, 16:808–835, 1987.

[2] J. Aspnes and G. Shah. Skip graphs. In *Proc. of the 14th ACM Symp. on Discrete Algorithms (SODA)*, pages 384–393, 2003. See also `http://cs-www.cs.yale.edu/homes/aspnes/skip-graphs-full.pdf`.

[3] B. Awerbuch and C. Scheideler. Group Spreading: A protocol for provably secure distributed name service. In *Proc. of the 31st International Colloquium on Automata, Languages and Programming (ICALP)*, 2004.

[4] M. Castro, P. Druschel, A. Ganesh, A. Rowstron, and D. Wallach. Security for structured peer-to-peer overlay networks. In *Proc. of the 5th Usenix Symp. on Operating Systems Design and Implementation (OSDI)*, 2002.

[5] M. Castro and B. Liskov. Practical Byzantine fault tolerance. In *Proc. of the 2nd Usenix Symp. on Operating Systems Design and Implementation (OSDI)*, 1999.

[6] S. Crosby and D. Wallach. Denial of service via algorithmic complexity attacks. In *Usenix Security*, 2003.

[7] M. Dietzfelbinger and F. Meyer auf der Heide. Simple, efficient shared memory simulations. In *Proc. of the 5 ACM Symp. on Parallel Algorithms and Architectures (SPAA)*, pages 110–119, 1993.

[8] J. R. Douceur. The sybil attack. In *Proc. of the 1st International Workshop on Peer-to-Peer Systems (IPTPS)*, 2002.

[9] P. Druschel and A. Rowstron. Pastry: Scalable, distributed object location and routing for large-scale peer-to-peer systems. In *Proc. of the 18th IFIP/ACM International Conference on Distributed Systems Platforms (Middleware 2001)*, 2001. See also `http://research.microsoft.com/~antr/Pastry`.

[10] D. Dubhashi and A. Panconesi. Concentration of measure for the analysis of randomized algorithms. Unpublished manuscript, accessible via http://www.cs.unibo.it/~pancones/papers.html, October 20 1998.

[11] A. Fiat, J. Saia, and M. Young. Making Chord robust to Byzantine attacks. In *Proc. of the European Symposium on Algorithms (ESA)*, 2005.

[12] R. Gennaro, Y. Ishai, E. Kushilevitz, and T. Rabin. The round complexity of verifiable secret sharing and secure multicast. In *Proc. of the 33rd ACM Symp. on Theory of Computing (STOC)*, pages 580–589, 2001.

[13] K.T. Herley and G. Bilardi. Deterministic simulations of PRAMs on bounded degree networks. *SIAM Journal on Computing*, 23:276–292, 1994.

[14] D. Karger, E. Lehman, T. Leighton, M. Levine, D. Lewin, and R. Panigrahi. Consistent hashing and random trees: Distributed caching protocols for relieving hot spots on the World Wide Web. In *29th ACM Symp. on Theory of Computing (STOC)*, pages 654–663, 1997.

[15] F. Kuhn, S. Schmid, and R. Wattenhofer. A self-repairing peer-to-peer system resilient to dynamic adversarial churn. In *Proc. of the 4th International Workshop on Peer-to-Peer Systems (IPTPS)*, 2005.

[16] F. Luccio, A. Pietracaprina, and G. Pucci. A new scheme for the deterministic simulation of PRAMs in VLSI. *Algorithmica*, 5:529–544, 1990.

[17] McDiarmid. Concentration. In M. Habib, C. McDiarmid, J. Ramirez-Alfonsin, and B. Reed, editors, *Probabilistic Methods for Algorithmic Discrete Mathematics*, pages 195–247. Springer Verlag, Berlin, 1998.

[18] K. Mehlhorn and U. Vishkin. Randomized and deterministic simulations of PRAMs by parallel machines with restricted granularity of parallel mamories. *Acta Informatica*, 21:339–374, 1984.

[19] F. Meyer auf der Heide and B. Vöcking. Shortest paths routing in arbitrary networks. *Journal of Algorithms*, 31(1):105–131, 1999.

[20] M. Naor and U. Wieder. Novel architectures for P2P applications: the continuous-discrete approach. In *Proc. of the 15th ACM Symp. on Parallel Algorithms and Architectures (SPAA)*, 2003.

[21] S. Nielson, S. Crosby, and D. Wallach. Kill the messenger: A taxonomy of rational attacks. In *Proc. of the 4th International Workshop on Peer-to-Peer Systems (IPTPS)*, 2005.

[22] V.N. Padmanabhan and K. Sripanidkulchai. The case for cooperative networking. In *Proc. of the 1st International Workshop on Peer-to-Peer Systems (IPTPS)*, 2002.

[23] G. Plaxton, R. Rajaraman, and A.W. Richa. Accessing nearby copies of replicated objects in a distributed environment. In *Proc. of the 9th ACM Symp. on Parallel Algorithms and Architectures (SPAA)*, pages 311–320, 1997.

[24] A. G. Ranade. How to emulate shared memory. *Journal of Computer and System Sciences*, 42:307–326, 1991.

[25] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker. A scalable content-addressable network. In *Proc. of the ACM SIGCOMM '01*, 2001.

[26] S. Rhea, D. Geels, T. Roscoe, and J. Kubiatowicz. Handling churn in a DHT. In *USENIX Annual Technical Conference*, 2004.

[27] J. Saia, A. Fiat, S. Gribble, A. Karlin, and S. Saroiu. Dynamically fault-tolerant content addressable networks. In *Proc. of the 1st International Workshop on Peer-to-Peer Systems (IPTPS)*, 2002.

[28] C. Scheideler. How to spread adversarial nodes? Rotate! In *Proc. of the 37th ACM Symp. on Theory of Computing (STOC)*, pages 704–713, 2005.

[29] C. Scheideler and W. Wang. A load-balanced peer-to-peer registration service. Unpublished manuscript. See http://www14.informatik.tu-muenchen.de/personen/scheideler, June 2005.

[30] A. Singh, M. Castro, A. Rowstron, and P. Druschel. Defending against Eclipse attacks on overlay networks. In *Proc. of the 11th ACM SIGOPS European Workshop*, 2004.

[31] E. Sit and R. Morris. Security considerations for peer-to-peer distributed hash tables. In *Proc. of 1st International Workshop on Peer-to-Peer Systems (IPTPS)*, 2002.

[32] M. Srivatsa and L. Liu. Vulnerabilities and security threats in structured overlay networks: A quantitative analysis. In *Proc. of the 20th IEEE Computer Security Applications Conference (ACSAC)*, 2004.

[33] T. Stading, P. Maniatis, and M. Baker. Peer-to-peer caching schemes to address flash crowds. In *Proc. of the 1st International Workshop on Peer-to-Peer Systems (IPTPS)*, 2002.

[34] A. Stavron, D. Rubenstein, and S. Sahn. A lightweight robust P2P system to handle flash crowds. In *Proc. of the IEEE Intl. Conf. on Network Protocols (ICNP)*, 2002.

[35] I. Stoica, R. Morris, D. Karger, M.F. Kaashoek, and H. Balakrishnan. Chord: A scalable peer-to-peer lookup service for Internet applications. In *Proc. of the ACM SIGCOMM '01*, 2001. See also http://www.pdos.lcs.mit.edu/chord/.

[36] I. Stoica, R. Morris, D. Liben-Nowell, D. Karger, M.F. Kaashoek, F. Dabek, and H. Balakrishnan. Chord: A scalable peer-to-peer lookup service for internet applications. Technical report, MIT, January 2002. See http://www.pdos.lcs.mit.edu/chord/.

[37] E. Upfal and A. Wigderson. How to share memory in a distributed system. *Journal of the ACM*, 34:116–127, 1987.

[38] B.Y. Zhao, J. Kubiatowicz, and A. Joseph. Tapestry: An infrastructure for fault-tolerant wide-area location and routing. Technical report, UCB/CSD-01-1141, University of California at Berkeley, 2001. See also http://www.cs.berkeley.edu/~ravenben/tapestry.

# A  Proofs in Section 2

## A.1  Proof of Lemma 2.5

Let $R_1, \ldots, R_C$ be the $k$-regions of $\hat{R}$, where $C = c \log n$. For every $k$-region $R_i$, let the random variable $X_i$ denote the age of $R_i$ at the beginning of the given round, and let $X = \sum_{i=1}^{C} X_i$. For all $i$ and $t \geq 1$ it holds that $\Pr[X_i = t] = (k/n)(1 - (k/n))^{t-1}$. Hence, $X_i$ is geometrically distributed with probability $p = k/n$. Thus, $\mathrm{E}[X_i] = 1/p = n/k$, and therefore, $\mathrm{E}[X] = \sum_{i=1}^{C} \mathrm{E}[X_i] = (n/k) \cdot C$. It remains to show that $X$ is concentrated around $\mathrm{E}[X]$.

Unfortunately, the ages of the $k$-regions are not independent as two $k$-regions cannot have the same age. However, there is an easy solution to this problem. Let $Y_1, \ldots, Y_C$ be independent random variables with the same probability distributions as $X_1, \ldots, X_C$ and let $Y = \sum_{i=1}^{C} Y_i$. Then it holds for all ages $0 < t_1 < t_2 \ldots < t_C$ that

$$\Pr[X_1 = t_1 \wedge \ldots \wedge X_C = t_C] = \prod_{i=1}^{C} \Pr[X_i = t_i \mid X_1 = t_1 \wedge \ldots \wedge X_{i-1} = t_{i-1}]$$

$$= \prod_{i=1}^{C} \left(1 - \frac{k}{n}\right)^{-(i-1)} \Pr[Y_i = t_i]$$

$$= \left(1 - \frac{k}{n}\right)^{-\binom{C}{2}} \Pr[Y_1 = t_1 \wedge \ldots \wedge Y_C = t_C]$$

Hence, because all $X_i$'s have the same probability distribution, it holds for any $T \geq 1$ that

$$\Pr[X \geq T] = \sum_{\{t_1,\ldots,t_C\} \subseteq \mathbb{N}, \sum_i t_i \geq T} \Pr[X_1 = t_1 \wedge \ldots \wedge X_C = t_C]$$

$$= \sum_{\{t_1,\ldots,t_C\} \subseteq \mathbb{N}, \sum_i t_i \geq T} \left(1 - \frac{k}{n}\right)^{-\binom{C}{2}} \Pr[Y_1 = t_1 \wedge \ldots \wedge Y_C = t_C]$$

$$\leq \left(1 - \frac{k}{n}\right)^{-\binom{C}{2}} \Pr[Y \geq T] \leq e^{C^2 k/n} \Pr[Y \geq T] \,.$$

Combining this with Lemma 2.2, it follows that $\Pr[X \geq (1+\delta)C \cdot n/k]$ is polynomially small for any constant $\delta > 0$ depending on the constant in $C$.

For a lower bound on the age of $\hat{R}$, we use the mapping $f : \mathbb{N}^C \to \mathbb{N}^C$ with

$$f(t_1, t_2, \ldots, t_C) = (t_1 + d_1, t_2 + d_2, \ldots, t_C + d_C)$$

where

$$d_j = |\{k \in \{1, \ldots, C\} \mid t_k < t_j \vee (t_k = t_j \wedge k < j)\}|$$

for all $j \in \{1, \ldots, C\}$. It is easy to check that this mapping is injective. Furthermore, for all $(t_1, \ldots, t_C) \in \mathbb{N}$ it holds for $(t'_1, \ldots, t'_C) = f(t_1, \ldots, t_C)$ that $t'_1, \ldots, t'_C$ are pairwise disjoint and

$$\Pr[Y_1 = t_1 \wedge \ldots \wedge Y_C = t_C] = \Pr[X_1 = t'_1 \wedge \ldots \wedge X_C = t'_C]$$

Hence, for any $T \geq 0$,

$$\Pr[Y \geq T] = \sum_{(t_1,\ldots,t_C) \in \mathbb{N}^C, \sum_i t_i \geq T} \Pr[(Y_1, \ldots, Y_C) = (t_1, \ldots, t_C)]$$

$$= \sum_{(t_1,\ldots,t_C) \in \mathbb{N}^C, \sum_i t_i \geq T} \Pr[(X_1, \ldots, X_C) = f(t_1, \ldots, t_C)]$$

$$\leq \sum_{\{t'_1,\ldots,t'_C\} \subseteq \mathbb{N}, \sum_i t'_i \geq T} \Pr[(X_1, \ldots, X_C) = (t'_1, \ldots, t'_C)]$$

$$= \Pr[X \geq T]$$

13

and therefore $\Pr[X \le T] \le \Pr[Y \le T]$. Combining this with Lemma 2.2, it follows that $\Pr[X \le (1 - \delta)C \cdot n/k]$ is polynomially small for any constant $\delta > 0$ depending on the constant in $C$.

## A.2 Proof of Lemma 2.6

We need two claims to prove the lemma. Let $T = \gamma(n/k) \ln n$.

**Claim A.1** *For any (black or white) node $v$ and any $c \ge 1$, $v$ gets replaced at most $(1 + \delta)\gamma \ln n$ times within $T$ rounds, w.h.p., where the constant $\delta > 0$ can be made arbitrarily small depending on the constant $\gamma$.*

**Proof.** For any $t \in \{1, \ldots, T\}$ let the binary random variable $X_t$ be 1 if and only if $v$ gets replaced in the $t$th round. Let $X = \sum_{t=1}^{T} X_t$. Because a $k$-region is chosen uniformly at random for eviction, $\Pr[X_t = 1] = k/n$ for every $t$. Hence, $\mathrm{E}[X] = (k/n) \cdot T = c \ln n$. Since the $X_t$'s are independent, the bound on $X$ follows from the Chernoff bounds. $\square$

**Claim A.2** *For any $k$-region $R$ in $\hat{R}$ it holds at any time that $R$ has an age of at most $T$, w.h.p.*

**Proof.** The probability that $R$ has been evicted at some round $t' \in \{t_0 - t - T, \ldots, t_0 - t - 1\}$ is $k/n$, and this probability is independent of other rounds. Hence, the probability that $R$ has an age of at least $T$ is equal to $(1 - \frac{k}{n})^T \le e^{(k/n) \cdot T} = n^{-c}$. $\square$

The two claims imply that there are at most $(1 + \epsilon)n \cdot (1 + \delta)\gamma \ln n$ node replacements during the lifetime of a $k$-region, w.h.p. Hence, the expected number of nodes in a $k$-region can be at most

$$\frac{k}{n} \cdot ((1 + \epsilon)n \cdot (1 + \delta)c \ln n + 1) = (1 + \epsilon)(1 + \delta)k \cdot \gamma \ln n + 1 .$$

Since the locations of the node replacements are independent of each other, it follows from Lemma 2.1 that the number of nodes in a $k$-region is at most $O(k \log n)$ at any time, w.h.p.

# B Proofs in Section 3

## B.1 Proof of Lemma 3.3

Let $d = \log n$. Suppose that, for randomly chosen functions $h_1, \ldots, h_{2c-1}$, $G$ is not a $(\sigma, c/2)$-expander. Then there exists a set $S \subset U$ with $|S| \le n$ and a $\sigma$-sparse $c/2$-bundle $F$ of $S$ with $|\Gamma_F(S)| < |S|$. We claim that the probability $p_s$ that such a set $S$ of size $s$ exists is at most

$$\binom{m}{s}\binom{n}{s}\binom{(2c-1)s}{cs/2}d^{cs/2}\binom{dn}{s} \cdot \sigma^{cs/2}$$

This holds because there are $\binom{m}{s}$ ways of choosing a subset $S \subset U$, $\binom{n}{s}$ ways of assigning the nodes in $S$ to nodes in $V_0$, $\binom{(2c-1)s}{cs/2}$ ways of selecting $cs/2$ edges for $F$ and $d^{cs/2}$ ways of specifying a level for them. Furthermore, there are at most $\binom{dn}{s}$ ways of choosing a $\sigma$-sparse $W \subseteq V$ witnessing a bad expansion of $F$. Since $W$ is $\sigma$-sparse, the probability that any particular edge selected for $F$ indeed points to a node in $W$ is at most $\sigma$, and since the hash functions are chosen uniformly and independently at random, we obtain a total probability of at most $\sigma^{cs/2}$.

Next we simplify $p_s$. Using the conditions on $c$ and $\sigma$ in the lemma it holds that

$$\binom{m}{s}\binom{n}{s}\binom{(2c-1)s}{cs/2}d^{cs/2}\binom{dn}{s} \cdot \sigma^{cs/2}$$
$$\le \left(\frac{em}{s}\right)^s \left(\frac{en}{s}\right)^s 2^{2cs} d^{cs/2} \left(\frac{edn}{s}\right)^s \cdot \sigma^{cs/2}$$
$$\le \left(\frac{em}{s} \cdot \frac{en}{s} \cdot \frac{edn}{s} \cdot (8d\sigma)^{c/2}\right)^s$$
$$\le \left(\frac{m^2}{s^3} \cdot \left(\frac{1}{2}\right)^{3\log m}\right)^s \le \frac{1}{m^s}$$

14

Hence, summing up over all possible values of $s$, we obtain a probability of having a bad $\sigma$-sparse $c/2$-bundle of at most $2/m$, which proves the lemma.

## B.2   Proof of Lemma 3.4

The proof is similar to, e.g., the proof of Theorem 1 in [13]. Suppose that, for randomly chosen functions $h_1, \ldots, h_{2c-1}$, $G$ is not a $(\lambda, c/2, \sigma)$-expander. Then there exists a set $S \subset U$ with $|S| \leq \sigma n/c$ and a $c/2$-bundle $F$ of $S$ with $V(F) \subseteq V_d$ and $|\Gamma_F(S)| < \lambda c|S|$. We claim that the probability $p_s$ that such a set $S$ of size $s$ exists is at most

$$\binom{m}{s}\binom{(2c-1)s}{cs/2}\binom{n}{\lceil \lambda cs - 1 \rceil} \cdot \left( \frac{\lceil \lambda cs - 1 \rceil}{n} \right)^{cs/2}$$

This holds because there are $\binom{m}{s}$ ways of choosing a subset $S \subset U$. Furthermore, there are $\binom{(2c-1)s}{cs/2}$ ways of a choosing $cs/2$ edges for $F$ and at most $\binom{n}{\lceil \lambda cs-1 \rceil}$ ways of choosing a set $W \subseteq V_d$ witnessing a bad expansion of the edges in $F$. The fraction of graphs $G$ for which the selected edges indeed point to nodes in $W$ is at most $(\frac{\lceil \lambda cs-1 \rceil}{n})^{cs}$ because the hash functions $h_1, \ldots, h_{2c-1}$ are chosen independently and uniformly at random.

Next we simplify $p_s$. Using the conditions on $c$ and $\sigma$ in the lemma it holds that

$$
\begin{aligned}
&\binom{m}{s}\binom{(2c-1)s}{cs/2}\binom{n}{\lceil \lambda cs - 1 \rceil} \cdot \left( \frac{\lceil \lambda cs - 1 \rceil}{n} \right)^{cs/2} \\
&\leq \left( \frac{em}{s} \right)^s 2^{2cs} \left( \frac{en}{\lambda cs} \right)^{\lambda cs} \left( \frac{\lambda cs}{n} \right)^{cs/2} \\
&\leq \left[ \frac{em}{s} \cdot \left( 8e^{2\lambda} \lambda^{1-2\lambda} \cdot \left( \frac{cs}{n} \right)^{1-2\lambda} \right)^{c/2} \right]^s \\
&\leq \left[ m \cdot \left( 8e^{2\lambda} \cdot \lambda^{1-2\lambda} \cdot \sigma^{1-2\lambda} \right)^{c/2} \right]^s \\
&\leq \left[ m \cdot \left( \frac{1}{2} \right)^{c/2} \right]^s \leq \frac{1}{m^s}
\end{aligned}
$$

Hence, summing up over all possible values of $s$, we obtain a probability of having a bad $c/2$-bundle of at most $2/m$, which proves the lemma.

## B.3   Proof of Theorem 3.5

Let us consider any fixed attempt and let $s$ be the number of remaining lookup requests. For simplicity, let us view the routing of packets in the dynamic de Bruijn graph as routing them along the $d$-dimensional Shuffle graph $SH(d)$ with $d = \log n$.

Let the threshold set by the lookup protocol be $2c/\sigma$ with $\sigma = 1/(16 \log n)$ so that Lemma 3.3 can be applied. A node in the Shuffle graph is called *congested* if more than $2c/\sigma$ packets for different data items pass it, and let $W$ be the set of all congested nodes. Note that $W$ is $\sigma$-sparse because in every subgraph $SH(\ell)$ in $SH(d)$ there can be a total of at most $(2c-1)2^\ell$ packets and therefore at most $(2c-1)2^\ell/(2c/\sigma) \leq \sigma 2^\ell$ nodes at level $\ell$ with at least $2c/\sigma$ packets. Furthermore, since the total number of packets is $(2c-1)s$, $|W| \leq d \cdot (2c-1)s/(2c/\sigma) = \sigma ds$.

Since the graph $G$ is a $(\sigma, c/2)$-expander, there can be at most $|W|$ lookup requests with at least $c/2$ requests passing through nodes in $W$ because otherwise we would have a $\sigma$-sparse $c/2$-bundle $F$ with $\Gamma_F(U(F)) < |U(F)|$. Hence, at most $|W| = \sigma ds = s/16$ of the lookup requests fail. Therefore, at most $O(\log n)$ attempts are necessary until all requests have been served.

## B.4   Proof of Theorem 3.6

Consider any fixed attempt and let $s$ be the number of remaining insert requests. For simplicity, we again view the routing of packets in the dynamic de Bruijn graph as routing them along the $d$-dimensional Shuffle graph $SH(d)$ with $d = \log n$.

Let the threshold for intermediate levels be $2c/\sigma$ with $\sigma = 1/(16 \log n)$ and the threshold for the final level be $2\gamma c$. A node in an intermediate level is *congested* if more than $2c/\sigma$ packets for different data items pass it, and a node in the final level is *congested* if more than $2\gamma c$ packets for different data items reach it. Let $W_1$ be the set of all intermediate congested nodes and $W_2$ be the set of all final congested nodes. From the previous proof we know that $W_1$ is $\sigma$-sparse and $|W_1| \leq \sigma d s$, and it is easy to see that $|W_2| \leq s/\gamma$.

Let $S$ be the set of failed lookup requests. Then either $S$ and $W_1$ form a $\sigma$-sparse $c/2$-bundle or $S$ and $W_2$ form a $c/2$-bundle. In the first case, we know from the previous proof that $|S| \leq s/16$. In the second case, $|S| \leq |W_2|/(\lambda c)$ for some constant $0 < \lambda < 1/2$ because $G$ is a $(\lambda, c/2, 1/\gamma)$-expander if $\gamma$ is sufficiently large. In both cases, $|S| \leq s/16$. Hence, at most $O(\log n)$ attempts are necessary until all requests have been served.

Since each node only starts to accept packets once the total number of different data items in it with packets is below $2\gamma c$, every node only has to store $O(c)$ copies at the end. Moving from nodes to regions, this means that every region has to store at most $O(c \log n)$ copies.

## C   Robust random ID generation

In order to generate a random ID for a node, we use a verifiable secret sharing (VSS) scheme. In VSS, a dealer $D$ tries to store a secret $s$ in $n$ nodes so that it can be reliably recovered. More precisely, a protocol on $n$ nodes is called a $(n, k)$-VSS scheme if, for any adversary owning $k$ nodes, the following requirements hold:

- **Privacy**: If $D$ is honest, then the adversary's view during the sharing phase reveals no information about $s$.

- **Correctness:** If $D$ is honest, then the reconstructed value is *always* equal to the secret $s$.

- **Commitment:** Even if $D$ is dishonest, any successful execution of the sharing phase determines a unique value $s^*$ which will be reconstructed at the reconstruction phase.

A protocol fulfilling all these properties is, for example, the $\lfloor \frac{n-1}{4} \rfloor$-VSS in [12]. For completeness, we present it here:

- Sharing phase:

  1. $D$ chooses a random bivariate polynomial $F \in K[x, y]$ of degree (at most) $k$ in each variable s.t. $F(0, 0) = s$. It sends to each player $P_i$ the (univariate) polynomials $f_i(x) = F(x, i)$ and $g_i(y) = F(i, y)$.

  2. Player $P_i$ sends to each player $P_j$ the value $g_i(j)$.

  3. Player $P_i$ broadcasts a list $L_i$ of players $P_j$ for whom it holds that $f_i(j) \neq g_j(i)$.

- Local computation (by each player):

  1. Add edge $(i, j)$ to the consistency graph $G$ on $n$ nodes if $P_i$ is not in $L_j$ and $P_j$ is not in $L_i$.

  2. Find a maximal matching in $\bar{G}$.

  3. Define a vertex set $C$ to include all vertices not in the matching. ($C$ is a clique in $G$.)

  4. Define $ADD$ to be the set of vertices $i$ s.t. $i \notin C$ and there exist $2k + 1$ nodes $j \in C$ such that $(i, j) \in G$.

  5. If $|C| + |ADD| \geq 3k + 1$ then accept the sharing; otherwise, disqualify the dealer.

- Reconstruction phase: Each $P_i \in C \cup ADD$ provides $f_i(0)$. Use error correction on $\{f_i(0)\}_{i \in C \cup ADD}$ to recover the polynomial $g_0(y) = F(0, y)$. Compute $g_0(0)$.

The problem with applying it to our setting is that it uses a broadcast operation to disseminate information and it assumes the set of nodes to be static. However, it can be adapted so that it can be used to generate a random ID in a dynamic, asynchronous environment without a broadcast channel. This works as follows (the number of adversarial nodes is assumed to be at most $k/6$):

Suppose that node $u$ wants to generate a new ID and let $G$ be the group of nodes $u$ knows in its quorum region $R_u$ (which includes all honest node but may not include all adversarial nodes). Then $u$ asks all nodes $v \in G$ to execute the VSS scheme above on $G$ for a secret $s_v$ picked at random from $[0, 1)$ by $v$. If $v$ knows more than $4k$ nodes in $G$, then $v$ executes step 1 of the sharing phase on $s_v$ and attaches $G$ to its messages; otherwise, it aborts. Every node $w \in G$ that receives a message from $v$ with the same $G$ it received from $u$ executes step 2 of the sharing phase. Every node $w' \in G$ that receives step 2-messages from at least $|G| - k/6$ nodes $w \in G$ for some $v$, computes $L_{w'}^v$ and sends

$L_{w'}^v$ to $u$. Once $u$ has received from at least $|G| - k/6$ nodes $w \in G$ $L_w^v$'s for *all* nodes $v$ in some set $S \subseteq G$ with $|S| = |G| - k/2$, $u$ starts the local computation phase to determine the set $P$ of all $v$ with $|C_v| + |ADD_v| \geq 3k + 1$. If $|P| \geq |G| - k$, then $u$ sends to each $v \in G$ an ID reconstruction request together with $P$ and $\{(C_v, ADD_v)\}_{v \in P}$. Each node $w \in G$ receiving an ID reconstruction request from $u$ forwards this request to all other nodes in $G$. Each node $w \in G$ receiving the same ID reconstruction request from at least $|G| - k/6$ nodes in $G$ waits until it has sent $L_w^p$ to $u$ for all $p \in P$ and afterwards initiates the reconstruction phase by sending $\{(f_w^p(0), L_w^p)\}_{p \in P}$ to all $v \in G$. Each node $w \in G$ receiving at least $|G| - k/3$ reconstruction messages for all $v \in P$ first checks if there is a $v \in P$ s.t. no change of $L_i^p$ for $\leq k/6$ messages it received together with suitable $L_i^p$'s for the $\leq k/3$ missing messages would fulfill the conditions on $C_p$ and $ADD_p$. If so, it aborts. Otherwise, it recovers $g_0^p(y)$ for every $p$, computes $x = \bigoplus_{p \in P} g_0(y)$, and sends $x$ to $u$.

We show the correctness with two lemmata, using the assumption that all honest nodes in $G$ know each other and at most $k/6$ nodes in $G$ are adversarial.

**Lemma C.1** *In any case in which there is an honest node $v$ that computes some ID $x$ in the ID generation scheme initiated by some (honest or adversarial) node $u$, $x$ must be random.*

**Proof.** First of all, an honest node $v$ only participates in the ID generation stage if $|G| > 4k$, i.e. $G$ is sufficiently large for the VSS-protocol in [12] with threshold $k$ to work. Second, notice that every honest node $v \in G$ will only reveal any of its private information about keys $x_w$ if it received at least $|G| - k/6$ votes concerning the ID reconstruction message from $u$ matching the reconstruction message it got from $u$. Thus, $P$ and $(C_p, ADD_p)_{p \in P}$ are fixed for $v$ at that stage. Furthermore, due to at most $k/6$ adversarial nodes, no honest node can have a different view of $P$ and $(C_p, ADD_p)_{p \in P}$ when revealing its private information. Hence, at that point where the first honest node reveals private information about some keys, no adversarial node can influence $P$ or $(C_p, ADD_p)_{p \in P}$ any more.

Also, at that point where the first honest node $v$ reveals private information, it must have sent $L_v^p$'s for all $p \in P$ to $u$. Thus, it must have received $g_i^p(j)$'s from at least $|G| - k/6$ nodes in $G$ for all $p \in P$. This, in turn, means that for each $p \in P$ at least $|G| - k/3$ honest nodes in $G$ must have received a pair $(f_i^p(x), g_i^p(y))$ from $p$ before any private information is revealed by any honest node. Let us call these nodes *p-safe*.

Now, any honest node $v$ that is convinced that the sharing for some $p \in P$ is successful must have used at least $|G| - 5t/6$ $f_i^p(0)$'s from $p$-safe honest nodes, because among the at least $|G| - k/3$ $f_i^p(0)$'s it receives from nodes in $G$, at most $k/6$ can come from adversarial nodes and at most $k/3$ can come from non-$p$-safe honest nodes. If $v$ is convinced of the correct sharing for $p$, then because it may have changed $L_i^p$'s from at most $k/6$ of the $p$-safe honest nodes to justify this. Hence, at least $|G| - k$ $f_i^p(0)$'s from $p$-safe nodes will be considered when revealing $p$'s secret which, according to [12], will recover a unique, unbiased value $x_p$.

Because $|P| \geq |G| - k$ and $|G| > 4k$ in order for an honest node $v$ to participate, at least one $x_p$ must have been generated by an honest node and therefore must be random. Hence, $v$ is convinced that the recovery phase succeeded, it computes a random ID $x$. □

**Lemma C.2** *Any honest node $u$ initiating the random ID generation scheme will get the same value $x$ back from at least $|G| - k/3$ nodes in $G$.*

**Proof.** If $u$ is honest, then it will wait until it has at least $|G| - k/6$ nodes $w \in G$ that sent $L_w^v$'s to $u$ for all nodes $v$ of some set $S \subseteq G$ with $|S| \geq |G| - k/2$. In this case, at least $|G| - k/3$ honest nodes $w$ must have sent $L_w^v$'s for all nodes in $S$, and therefore at least $|G| - k/3$ honest nodes will initiate the reconstruction phase. This makes sure that every honest node in $G$ receives reconstruction messages from at least $|G| - k/3$ honest nodes in $G$, which allows them to recover the keys $x_p$. Notice that no honest node will abort, because if $u$ makes sure that the conditions for $ADD_p$ and $C_p$ are fulfilled for every $p \in P$, then every honest node can find corrections for the at most $k/6$ $L_i^P$'s it received from adversarial nodes and can come up with suitable $L_i^p$'s for the at most $k/3$ missing nodes so that the conditions on $ADD_p$ and $C_p$ are met. If these conditions can be met, it follows from [12] that unique keys can be recovered from the received parts, and from the lemma above it follows that these keys must be unbiased, and at least one of them must be random. Hence, all honest nodes in $G$ that participate in the recovery, which are at least $|G| - k/3$, will agree on the same, random value for $x$. □