

A Self-Repairing Peer-to-Peer System Resilient to Dynamic Adversarial Churn*

Fabian Kuhn, Stefan Schmid, Roger Wattenhofer
{kuhn, schmiste, wattenhofer}@tik.ee.ethz.ch

Computer Engineering and Networks Laboratory (TIK), ETH Zurich, 8092 Zurich, Switzerland

Abstract

We present a dynamic distributed hash table where peers may join and leave at any time. Our system tolerates a powerful adversary which has complete visibility of the entire state of the system and can continuously add and remove peers. Our system provides worst-case fault-tolerance, maintaining desirable properties such as a low peer degree and a low network diameter.

1 Introduction

Storing and handling data in an efficient way lie at the heart of any data-driven computing system. Compared to a traditional client/server approach, decentralized peer-to-peer (P2P) systems have the advantage to be more reliable, available, and efficient. P2P systems are based on common desktop machines (“peers”), distributed over a large-scale network such as the Internet. These peers share data (as well as the management of the data) that is conventionally stored on a central server. Usually, peers are under control of individual users who turn their machines on or off at any time. Such peers join and leave the P2P system at high rates (“churn”), a problem that is not existent in orthodox distributed systems. In other words, a P2P system consists of unreliable components only. Nevertheless, the P2P system should provide a reliable and efficient service.

Most P2P systems in the literature are analyzed against an adversary who can crash a functionally bounded number of random peers. After crashing a few peers the system is given sufficient time to recover again. The scheme described in this paper significantly differs from this in two major aspects. First, we assume that joins and leaves occur in a worst-case manner. We think of an adversary which can remove and add a bounded number of peers. The adversary cannot be fooled

by any kind of randomness. It can choose which peers to crash and how peers join.¹ Note that we use the term “adversary” to model worst-case behavior. We do not consider Byzantine faults. Second, the adversary does not have to wait until the system is recovered before it crashes the next batch of peers. Instead, the adversary can constantly crash peers while the system is trying to stay alive. Indeed, our system is *never fully repaired* but *always fully functional*. In particular, our system is resilient against an adversary which continuously attacks the “weakest part” of the system. Such an adversary could for example insert a crawler into the P2P system, learn the topology of the system, and then repeatedly crash selected peers, in an attempt to partition the P2P network. Our system counters such an adversary by continuously moving the remaining or newly joining peers towards the sparse areas.

Clearly, we cannot allow our adversary to have unbounded capabilities. In particular, in any constant time interval, the adversary can at most add and/or remove $O(\log n)$ peers, n being the total number of peers currently in the system. This model covers an adversary which repeatedly takes down machines by a distributed denial of service attack, however only a logarithmic number of machines at each point in time. Our algorithm relies on messages being delivered timely, in at most constant time between any pair of operational peers. In distributed computing such a system is called *synchronous*. Note that if nodes are synchronized locally, our algorithm also runs in an asynchronous environment. In this case, the propagation delay of the slowest message defines the notion of time which is needed for the adversarial model.

The basic structure of our P2P system is a hypercube. Each peer is part of a distinct hypercube node; each hypercube node consists of $\Theta(\log n)$ peers. Peers have connections to other peers of their hypercube node and to peers of the neighboring hypercube nodes. In the case of joins or leaves, some of the peers have to change to another hy-

*Research (in part) supported by the Hasler Stiftung under grant number 1828 and the Swiss National Science Foundation. A short version of this paper has been published at the 4th International Workshop on Peer-To-Peer Systems (IPTPS), Ithaca, New York, USA, 2005.

¹We assume that a joining peer knows a peer which already belongs to the system. This is known as the *bootstrap* problem.

percube node such that up to constant factors, all hypercube nodes own the same number of peers at all times. If the total number of peers grows or shrinks above or below a certain threshold, the dimension of the hypercube is increased or decreased by one, respectively.

The balancing of peers among the hypercube nodes can be seen as a dynamic token distribution problem [13] on the hypercube. Each node of a graph (hypercube) has a certain number of tokens, the goal is to distribute the tokens along the edges of the graph such that all nodes end up with the same or almost the same number of tokens. While tokens are moved around, an adversary constantly inserts and deletes tokens. Our P2P system builds on two basic components: i) an algorithm which performs the described dynamic token distribution and ii) an information aggregation algorithm which is used to estimate the number of peers in the system and to adapt the dimension accordingly.

Based on the described structure, we get a fully scalable, efficient P2P system which tolerates $O(\log n)$ worst-case joins and/or crashes per constant time interval. As in other P2P systems, peers have $O(\log n)$ neighbors, and the usual operations (e.g. search) take time $O(\log n)$. In our view a main contribution of the paper, however, is to propose and study a model which allows for dynamic adversarial churn. We believe that our basic algorithms (dynamic token distribution and information aggregation) can be applied to other P2P topologies, such as butterflies, skip graphs, chordal rings, etc. It can even be used for P2P systems that go beyond distributed hash tables (DHT).

The paper is organized as follows. In Section 2 we discuss relevant related work. Section 3 gives a short description of the model. A detailed discussion of our P2P system is given in Sections 4 and 5. Finally, the paper is concluded in Section 6.

2 Related Work

A plethora of different overlay networks with various interesting technical properties have been proposed over the last years (e.g. [1, 3, 5, 6, 9, 10, 12, 15, 16, 20, 23]). Due to the nature of P2P systems, fault-tolerance has been a prime issue from the beginning. The systems usually tolerate a large number of random faults. However after crashing a few peers the systems are given sufficient time to recover again. From an experimental point of view, churn has been studied in [17], where practical design tradeoffs in the implementation of existing P2P networks are considered.

Resilience to worst-case failures has been studied by Fiat, Saia et al. in [8, 18]. They propose a system where, w.h.p., $(1 - \varepsilon)$ -fractions of peers and data survive the adversarial deletion of up to half of all nodes. In contrast to our work

the failure model is static. Moreover, if the total number of peers changes by a constant factor, the whole structure has to be rebuilt from scratch.

Scalability and resilience to worst-case joins and leaves has been addressed by Abraham et al. in [2]. The focus lies on maintaining a balanced network rather than on fault-tolerance in the presence of concurrent faults. In contrast to our paper, whenever a join or leave happens, the network has some time to adapt.

The only paper which explicitly treats arbitrarily concurrent worst-case joins and leaves is by Li et al. [11]. In contrast to our work, Li et al. consider a completely asynchronous model where messages can be arbitrarily delayed. The stronger communication model is compensated by a weaker failure model. It is assumed that peers do not crash. Leaving peers execute an appropriate “exit” protocol and do not leave before the system allows this; crashes are not allowed.

3 Model

We consider the *synchronous message passing model*. In each round, each peer can send a message to all its neighbors. Additionally, we have an adversary $\mathcal{A}(J, L, \lambda)$ which may perform J arbitrary joins and L arbitrary leaves (crashes) in each interval of λ rounds.

We assume that a joining peer π_1 contacts an arbitrary peer π_2 which already belongs to the system; π_2 then triggers the necessary actions for π_1 's integration. A peer may be contacted by several joining peers simultaneously. In contrast to other systems where peers have to do some finalizing operations before leaving, we consider the more general case where peers depart or crash without notice.

4 Algorithm

In this section, we describe the maintenance algorithm which maintains the simulated hypercube in the presence of an adversary which constantly adds and removes peers. The goal of the maintenance algorithm is twofold. It guarantees that each node always contains at least one peer which stores the node's data. Further, it adapts the hypercube dimension to the total number of peers in the system.

This is achieved by two basic components. First, we present a dynamic token distribution algorithm for the hypercube. Second, we describe an information aggregation scheme which allows the nodes to simultaneously change the dimension of the hypercube.

4.1 Dynamic Token Distribution

The problem of distributing peers uniformly throughout a hypercube is a special instance of a *token distribution prob-*

lem, first introduced by Peleg and Upfal [13]. The problem has its origins in the area of load balancing, where the workload is modeled by a number of *tokens* or jobs of unit size; the main objective is to distribute the total load equally among the processors. Such load balancing problems arise in a number of parallel and distributed applications including job scheduling in operating systems, packet routing, large-scale differential equations and parallel finite element methods. More applications can be found in [19].

Formally, the goal of a token distribution algorithm is to minimize the maximum difference of tokens at any two nodes, denoted by the *discrepancy* ϕ . This problem has been studied intensively; however, most of the research is about the *static variant* of the problem, where given an arbitrary initial token distribution, the goal is to redistribute these tokens uniformly. In the *dynamic variant* on the other hand, the load is dynamic, that is, tokens may arrive and depart during the execution of the token distribution algorithm. In our case, peers may join and leave the simulated hypercube at arbitrary times, so the emphasis lies on the dynamic token distribution problem on a d -dimensional hypercube topology.

We use two variants of the token distribution problem: In the *fractional token distribution*, tokens are arbitrarily divisible, whereas in the *integer token distribution* tokens can only move as a whole. In our case, tokens represent peers and are inherently integer. However, it turns out that the study of the fractional model is useful for the analysis of the integer model.

We use a token distribution algorithm which is based on the *dimension exchange method* [7, 14]. Basically, the algorithm cycles continuously over the d dimensions of the hypercube. In step s , where $i = s \bmod d$, every node $u := \beta_0 \dots \beta_i \dots \beta_{d-1}$ having a tokens balances its tokens with its adjacent node in dimension i , $v := \beta_0 \dots \overline{\beta_i} \dots \beta_{d-1}$, having b tokens, such that both nodes end up with $\frac{a+b}{2}$ tokens in the fractional token distribution. On the other hand, if the tokens are integer, one node is assigned $\lceil \frac{a+b}{2} \rceil$ tokens and the other one gets $\lfloor \frac{a+b}{2} \rfloor$ tokens.

It has been pointed out in [7] that the described algorithm yields a perfect discrepancy $\phi = 0$ after d steps for the static fractional token distribution. In [14], it has been shown that in the worst case, $\phi = d$ after d steps in the static integer token distribution. We can show that if the decision to which node to assign $\lceil \frac{a+b}{2} \rceil$ and to which node to assign $\lfloor \frac{a+b}{2} \rfloor$ tokens is made randomly, the final discrepancy is constant in expectation. However, we do not make use of this because it has no influence on our asymptotic results.

In the following, the dynamic integer token distribution problem is studied, where a “token adversary” $\mathcal{A}(J, L, 1)$ adds at most J and removes at most L tokens at the beginning of each step. In particular, we will show that if the initial

distribution is perfect, i.e., $\phi = 0$, our algorithm maintains the invariant $\phi \leq 2J + 2L + d$ at every moment of time.

Lemma 4.1. *The dimension exchange algorithm yields a perfect discrepancy $\phi = 0$ after d steps for the static fractional token distribution.*

Proof. We prove by induction that after the i^{th} iteration, all nodes with the same postfix of length $d - 1 - i$ have the same number of tokens.

$i = 0$: Let v be node 0α having a tokens and u its neighbor 1α having b tokens, for an arbitrary $d - 1$ bit vector α . After balancing, both nodes have $\frac{a+b}{2}$ tokens.

$i \rightarrow i + 1$: Consider two i -dimensional sub-cubes H_0 and H_1 consisting of all nodes with postfixes 0α and 1α respectively, where α is an arbitrary bit string of length $d - 1 - (i + 1)$. By the induction hypothesis, all nodes in H_0 have the same number of tokens, say a , and all nodes in H_1 have the same number of tokens, say b . After step $i + 1$, all nodes in $V(H_0) \cup V(H_1)$ — sharing the postfix α — will have $\frac{a+b}{2}$ tokens. To see this, consider the nodes $v_0 = \zeta 0\alpha$, $v_1 = \eta 0\alpha \in V(H_0)$, and $u_0 = \zeta 1\alpha$, $u_1 = \eta 1\alpha \in V(H_1)$, where ζ and η are arbitrary bit vectors of length $i + 1$. During the exchange of round $i + 1$, v_0 balances with u_0 and v_1 with u_1 . Obviously, all four nodes end up with $\frac{a+b}{2}$ tokens. \square

For the dynamic fractional token distribution, the tokens inserted and deleted at different times can be treated independently and be superposed. Therefore, the following lemma holds.

Lemma 4.2. *For the dynamic fractional token distribution, the number of tokens at a node depends only on the token insertions and deletions of the last d steps and on the total number of tokens in the system.*

Proof. Assume that a total amount of T tokens are distributed in two different ways on the d -dimensional hypercube. According to Lemma 4.1, each node has exactly $\frac{T}{2^d}$ tokens after d steps in the absence of an adversary. On the other hand, the token insertions and removals of the adversary that happen in-between can be treated as an independent superposition, as the corresponding operations are all linear. \square

We can now bound the discrepancy of the integer token distribution algorithm by comparing it with the fractional problem.

Lemma 4.3. *Let v be a node of the hypercube. Let $\tau_v(t)$ and $\tau_{v,f}(t)$ denote the number of tokens at v for the integer and fractional token distribution algorithms at time t , respectively. We have $\forall t : |\tau_v(t) - \tau_{v,f}(t)| \leq \frac{d}{2}$.*

Proof. For $t = 0$, we have $\tau_v(t) = \tau_{v,f}(t)$. For symmetry reasons, it is sufficient to show the upper bound $\tau_v(t) \leq \tau_{v,f}(t) + \frac{d}{2}$. We first prove by induction that $\tau_v(t) \leq \tau_{v,f}(t) + \frac{t}{2}$ at time t .

For the induction step, we consider two neighbors u and v which exchange tokens. We have

$$\begin{aligned} \tau_v(t+1) &\leq \left\lceil \frac{\tau_v(t) + \tau_u(t)}{2} \right\rceil \\ &\leq \left\lceil \frac{\lfloor \tau_{v,f}(t) + \frac{t}{2} \rfloor + \lfloor \tau_{u,f}(t) + \frac{t}{2} \rfloor}{2} \right\rceil \\ &\leq \frac{\lfloor \tau_{v,f}(t) + \frac{t}{2} \rfloor + \lfloor \tau_{u,f}(t) + \frac{t}{2} \rfloor}{2} + \frac{1}{2} \\ &\leq \tau_{v,f}(t+1) + \frac{t+1}{2}. \end{aligned}$$

The second inequality follows from the induction hypothesis and the fact that $\tau_v(t)$ and $\tau_u(t)$ are integers. Note that adding or removing tokens has no influence on the difference between τ_v and $\tau_{v,f}$ because it modifies τ_v and $\tau_{v,f}$ in the same way.

So far, we have seen that the number of integer tokens can deviate from the number of fractional tokens by at most $\frac{d}{2}$ after the first d steps. In order to show that this holds for all times t , we consider a fractional token distribution problem $\hat{\tau}_{v,f}$ for which $\hat{\tau}_{v,f}(t-d) = \tau_v(t-d)$. Using the above argument, we have $\tau_v(t-d) \leq \hat{\tau}_{v,f}(t) + \frac{d}{2}$ and by Lemma 4.2, we get $\hat{\tau}_{v,f}(t) = \tau_{v,f}(t)$. This concludes the proof. \square

Let ϕ_f and ϕ_i be the discrepancies of the fractional and the integer token distribution algorithms, respectively. From Lemma 4.3, we get Corollary 4.4.

Corollary 4.4. *At all times, $\phi_i - \phi_f \leq d$.*

Lemma 4.5. *In the presence of an adversary $\mathcal{A}(J, L, 1)$, it always holds that the integer discrepancy $\phi \leq 2J + 2L + d$.*

Proof. We show that the fractional discrepancy ϕ_f is bounded by $2J + 2L$. Together with Corollary 4.4, the claim follows. Let $J_t \leq J$ and $L_t \leq L$ be the insertions and deletions that happen at the beginning of step t . First, we consider the case of joins only, i.e., $L_t = 0$. Assume that all J_t tokens are inserted at node $v = \beta_0 \dots \beta_i \dots \beta_{d-1}$ where $i := t \bmod d$. In the upcoming paragraph, all indices are implicitly modulo d . In step t , according to the token distribution algorithm, v keeps $J_t/2$ tokens and sends $J_t/2$ to node $u = \beta_0 \dots \beta_i \dots \beta_{d-1}$. In step $t+1$, $J_t/4$ are sent to nodes $\beta_0 \dots \beta_i \beta_{i+1} \dots \beta_{d-1}$ and $\beta_0 \dots \beta_i \beta_{i+1} \dots \beta_{d-1}$, and so on. Thus, after step $t+d-1$, every node in the d -dimensional hypercube has the same share of $\frac{J_t}{2^d}$ tokens from that insertion. We conclude that a node can have at most all insertions of this

step, half of the insertions of the last step, a quarter of all insertions two steps ago and so on:

$$\begin{aligned} &\underbrace{J_t + \frac{J_{t-1}}{2} + \frac{J_{t-2}}{4} + \dots + \frac{J_{t-(d-1)}}{2^{d-1}}}_{< 2J} \\ &+ \underbrace{\frac{J_{t-d}}{2^d} + \frac{J_{t-(d+1)}}{2^d} + \frac{J_{t-(d+2)}}{2^d} + \dots}_{\text{shared by all nodes}} \end{aligned}$$

Since $J_{t-i} \leq J$ for $i = 0, 1, 2, \dots$, we have $\phi_f \leq 2J$. For the case of only token deletions, the same argument can be applied, yielding a discrepancy of at most $2L$. Finally, if there are both insertions and deletions which do not cancel out each other, we have $\phi_f \leq 2J + 2L$. \square

4.2 Information Aggregation

When the total number of peers in the d -dimensional hypercube system exceeds a certain threshold, all nodes $\beta_0 \dots \beta_{d-1}$ have to split into two new nodes $\beta_0 \dots \beta_{d-1}0$ and $\beta_0 \dots \beta_{d-1}1$, yielding a $(d+1)$ -dimensional hypercube. Analogously, if the number of peers falls beyond a certain threshold, nodes $\beta_0 \dots \beta_{d-2}0$ and $\beta_0 \dots \beta_{d-2}1$ have to merge their peers into a single node $\beta_0 \dots \beta_{d-2}$, yielding a $(d-1)$ -dimensional hypercube. Based on ideas also used in [4, 21, 22], we present an algorithm which provides the same estimated number of peers in the system to all nodes in every step allowing all nodes to split or merge synchronously, that is, in the same step. The description is again made in terms of *tokens* rather than peers.

Assume that in order to compute the total number of tokens in a d -dimensional hypercube, each node $v = \beta_0 \dots \beta_{d-1}$ maintains an array $\Gamma_v[0 \dots d]$, where $\Gamma_v[i]$ for $i \in [0, d]$ stores the estimated number of tokens in the sub-cube consisting of the nodes sharing v 's prefix $\beta_0 \dots \beta_{d-1-i}$. Further, assume that at the beginning of each step, an adversary inserts and removes an arbitrary number of tokens at arbitrary nodes. Each node $v = \beta_0 \dots \beta_{d-1-i} \dots \beta_{d-1}$ then calculates the new array $\Gamma'_v[0 \dots d]$. For this, v sends $\Gamma_v[i]$ to its adjacent node $u = \beta_0 \dots \beta_{d-1-i} \dots \beta_{d-1}$, for $i \in [0, d-1]$. Then, $\Gamma'_v[0]$ is set to the new number of tokens at v which is the only node with prefix $\beta_0 \dots \beta_{d-1}$. For $i \in [1, d]$, the new estimated number of tokens in the prefix domain $\beta_0 \dots \beta_{d-1-(i+1)}$ is given by the total number of tokens in the domain $\beta_0 \dots \beta_{d-1-i}$ plus the total number of tokens in domain $\beta_0 \dots \beta_{d-1-i}$ provided by node u , that is, $\Gamma'_v[i+1] := \Gamma_v[i] + \Gamma_u[i]$.

Lemma 4.6. *Consider two arbitrary nodes v_1 and v_2 of the d -dimensional hypercube. Our algorithm guarantees that $\Gamma_{v_1}[d] = \Gamma_{v_2}[d]$ at all times t . Moreover, it holds that this*

value is the correct total number of tokens in the system at time $t - d$.

Proof. We prove by induction that at time $t + k$, all nodes sharing the prefix $\beta_0 \dots \beta_{d-1-k}$ for $k \in [0, d]$ store the same value $\Gamma_v[k]$ which represents the correct state of that sub-domain in step t .

$k = 0$: There is only one node having the prefix $\beta_0 \dots \beta_{d-1}$, so the claim trivially holds.

$k \rightarrow k + 1$: By the induction hypothesis, all nodes v with prefix $\beta_0 \dots \beta_{d-1-(k+1)} \beta_{d-1-k}$ share the same value $\Gamma_v[k]$ which corresponds to the state of the system k steps earlier, and the same holds for all nodes u with prefix $\beta_0 \dots \beta_{d-1-(k+1)} \beta_{d-1-k}$. In step $k+1$, all these nodes having the same prefix $\beta_0 \dots \beta_{d-1-(k+1)}$ obviously store the same value $\Gamma'_v[k+1] = \Gamma'_u[k+1] = \Gamma_v[k] + \Gamma_u[k]$. \square

5 Simulated Hypercube

Based on the components presented in the previous sections, both the topology and the maintenance algorithm are now described in detail. In particular, we show that, given an adversary $\mathcal{A}(d+1, d+1, 6)$ which inserts and removes at most $d+1$ peers in any time interval of 6 rounds, 1) the out-degree of every peer is bounded by $\Theta(\log^2 n)$ where n is the total number of peers in the system, 2) the network diameter is bounded by $\Theta(\log n)$, and 3) every node of the simulated hypercube has always at least one peer which stores its data items, so no data item will ever be lost.

5.1 Topology

We start with a description of the overlay topology. As already mentioned, the peers are organized to simulate a d -dimensional hypercube, where the hypercube's nodes are represented by a group of peers. A data item with identifier id is stored at the node whose identifier matches the first d bits of the hash-value of id .

The peers of each node v are divided into a *core* \mathcal{C}_v of at most $2d + 3$ peers and a *periphery* \mathcal{P}_v consisting of the remaining peers; all peers within the same node are completely connected (*intra-connections*). Moreover, every peer is connected to all *core* peers of the neighboring nodes (*inter-connections*). Figure 1 shows an example for $d = 2$.

The data items belonging to node v are replicated on all core peers, while the peripheral peers are used for the balancing between the nodes according to the peer distribution algorithm and do not store any data items. The partition into core and periphery has the advantage that the peers which move between nodes do not have to replace the data of the old node by the data of the new nodes in most cases.

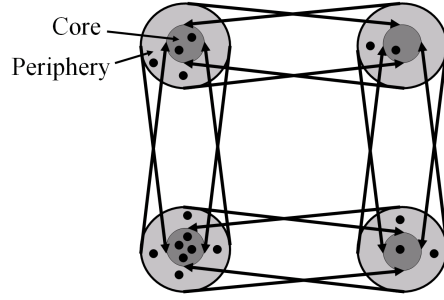


Figure 1: A simulated 2-dimensional hypercube with four nodes, each consisting of a core and a periphery. All peers within the same node are completely connected to each other, and additionally, all peers of a node are connected to all core peers of the neighboring nodes. Only the core peers store data items, while the peripheral peers may move between the nodes to balance biased adversarial changes.

5.2 6-Round (Maintenance) Algorithm

The *6-round (maintenance) algorithm* maintains the simulated hypercube topology described in the previous section given an adversary $\mathcal{A}(d+1, d+1, 6)$. In particular, it ensures that 1) every node has at least one core peer all the times and hence no data is lost; 2) each node always has between $3d + 10$ and $45d + 86$ peers; 3) only peripheral peers are moved between nodes, thus the unnecessary copying of data is avoided.

In the following, we refer to a complete execution of all six rounds (ROUND 1 – ROUND 6) of the maintenance algorithm as a *phase*. Basically, the 6-round algorithm balances the peers across one dimension in every phase according to the token distribution algorithm as described in Section 4.1; additionally, the total number of peers in the system is computed with respect to an earlier state of the system by the information aggregation algorithm of Section 4.2 to expand or shrink the hypercube if the total number of peers exceeds or falls below a certain threshold. In our system, we use the lower threshold $LT := 8d + 16$ and the upper threshold $UT := 40d + 80$ for the total number of peers *per node on average*.²

While peers may join and leave the system at arbitrary times, the 6-round algorithm considers the (accumulated) changes only once per phase. That is, a snapshot of the system is made in ROUND 1; ROUNDS 2 – 6 then ignore the changes that might have happened in the meantime and depend solely on the snapshot at the beginning of the phase.

²Note that since we consider the threshold *on average*, and since these values are provided with a delay of d phases in a d -dimensional hypercube (see Lemma 4.6), the number of peers at an individual node may lie outside $[LT, UT]$.

ROUND 1

Outline: Each node v makes a snapshot of the currently active peers, denoted by the ID set \mathcal{S}_v . The later rounds will only be based on these sets.

Sent Messages: Each peer of a node v sends a packet with its own ID and the (potentially empty) ID set of its joiners to all adjacent peers *within* v .

ROUND 2

Outline: Based on the snapshot of ROUND 1, the core peers of a node v know the total number of peers in the node, $size(v) := |\mathcal{S}_v|$. This information is needed for the peer distribution algorithm (see Section 4.1) and for the estimation of the total number of peers in the system (see Section 4.2).

Local Computations: The core peers compute $size(v) := |\mathcal{S}_v|$.

Sent Messages: Each peer informs its joiners about \mathcal{S}_v . The core peers \mathcal{C}_v additionally send the number $size(v)$ to their neighboring core \mathcal{C}_u , where node u is v 's neighbor in dimension i — the node with which v has to balance its peers in this phase. The core also exchanges the new estimated total number of peers in its domains with the corresponding adjacent cores (according to the algorithm presented in Section 4.2).

ROUND 3

Outline: At the beginning of this round, every peer within a node v knows \mathcal{S}_v , and the transfer for the peer distribution algorithm can be prepared. Let v again be an arbitrary node and u its adjacent node in dimension i . We assume that $size(v) > size(u)$; the case where $size(v) \leq size(u)$ is analogous and not described further here. The ID set \mathcal{T} of peers that have to move from node v to node u are the $\frac{size(v)-size(u)}{2}$ (arbitrarily rounded) peers in the periphery \mathcal{P}_v having the smallest identifiers.

Local Computations: The peers in each node v compute the new periphery $\mathcal{P}_v := \mathcal{S}_v \setminus \mathcal{C}_v$. The core remains the same.

Sent Messages: All cores forward the information about the new estimated total number of peers in the system to their peripheral peers. Moreover, the core of the larger node \mathcal{C}_v sends the identifiers of the to be transferred peers \mathcal{T} to \mathcal{C}_u , and the number $\frac{size(v)-size(u)}{2}$ to the new periphery \mathcal{P}_v .

ROUND 4

Outline: The transfer for the peer distribution algorithm is continued. Moreover, this round prepares the dimension reduction if necessary.

Sent Messages: The core \mathcal{C}_u informs the peers in \mathcal{T} about all neighboring cores \mathcal{C}_{u_j} , where u_j is the neighbor of u in dimension j for $j \in [0, d-1]$, about \mathcal{C}_u itself, about \mathcal{S}_u and about its peripheral peers \mathcal{P}_u . Moreover, \mathcal{C}_u informs its own periphery \mathcal{P}_u about the newcomers \mathcal{T} .

If the estimated total number of peers in the system is beyond the threshold, the core peers of a node \bar{v} which will be reduced send their data items plus the identifiers of all their peripheral peers (with respect to the situation *after* the transfer) to the core of their adjacent node \underline{v} .

ROUND 5

Outline: This round finishes the peer distribution, establishes the new peripheries, and prepares the building of a new core. If the hypercube has to grow in this phase, the nodes start to split, and vice versa if the hypercube is going to shrink.

Local Computations: Given the number $\frac{size(v)-size(u)}{2}$, the peripheral peers \mathcal{P}_v can compute the set \mathcal{T} selecting the $\frac{size(v)-size(u)}{2}$ smallest elements in \mathcal{P}_v . From this, the new periphery $\mathcal{P}_v := \mathcal{P}_v \setminus \mathcal{T}$ is computed. Analogously, the peers in node u (including \mathcal{T}) can compute the new periphery $\mathcal{P}_u := \mathcal{P}_u \cup \mathcal{T}$.

Then, all peers of each node v calculate the new core \mathcal{C}_v^{new} : It consists of the peers of the old core which have still been alive in ROUND 1, i.e., $\mathcal{C}_v^{old} := \mathcal{C}_v \cap \mathcal{S}_v$, plus the $2d+3 - |\mathcal{C}_v \cap \mathcal{S}_v|$ smallest IDs in the new periphery \mathcal{P}_v , denoted by \mathcal{C}_v^Δ . Hence, the new core is given by $\mathcal{C}_v^{new} := \mathcal{C}_v^{old} \cup \mathcal{C}_v^\Delta$, and the new periphery by $\mathcal{P}_v^{new} := \mathcal{P}_v \setminus \mathcal{C}_v^\Delta$.

If the hypercube has to grow in this phase, the smallest $2d+3$ peers in the new periphery \mathcal{P}_v^{new} become the new core of the expanded node, $\mathcal{C}_{\bar{v}}$. Half of the remaining peripheral peers, the ones with the smaller identifiers, build the new periphery $\mathcal{P}_{\bar{v}}$, and the other half becomes $\mathcal{P}_{\underline{v}}$. All these operations can be computed locally by every peer.

Sent Messages: The old core \mathcal{C}_v^{old} informs all its neighboring nodes (i.e., their old cores) about the new core \mathcal{C}_v^{new} . Moreover, \mathcal{C}_v^{old} sends its data items to the peers in \mathcal{C}_v^Δ .

If the hypercube is about to grow, \mathcal{C}_v^{old} sends the necessary data items to the core peers of the new node, $\mathcal{C}_{\bar{v}}$. Moreover, \mathcal{C}_v^{old} informs its neighboring (old) cores about the IDs of its expanded core $\mathcal{C}_{\bar{v}}$.

If the hypercube is about to shrink, all cores \mathcal{C}_v^{old} inform their periphery about the peers arriving from the expanded node and the peers in the expanded node about the new core $\mathcal{C}_{\underline{v}}^{new}$ and its periphery. \mathcal{C}_v^{old} copies also the data items of $\mathcal{C}_{\bar{v}}^{old}$ to the peers $\mathcal{C}_{\underline{v}}^\Delta$.

ROUND 6

Outline: Building the new cores and accomplishing the dimension change if necessary.

Local Computations: If the hypercube has been reduced, every peer can now compute the new periphery \mathcal{P}_v .

Sent Messages: The old core \mathcal{C}_v^{old} forwards the information about the new neighboring cores to the peers $\mathcal{C}_v^\Delta \cup \mathcal{P}_v$.

If the hypercube has grown, \mathcal{C}_v^{old} forwards the expanded cores of its neighboring nodes to *all* peers in its expanded node \bar{v} . Note that this requires that \mathcal{C}_v^{old} remembers the peripheral peers that have been transferred to \bar{v} in ROUND 5.

Theorem 5.1. *Given an adversary $\mathcal{A}(d+1, d+1, 6)$ which inserts and removes at most $d+1$ peers per phase, the described 6-round algorithm ensures that 1) every node always has at least one core peer and hence no data is lost; 2) each node has between $3d+10$ and $45d+86$ peers, yielding a logarithmic network diameter; 3) only peripheral peers are moved between nodes, thus the unnecessary copying of data is avoided.*

Proof. We first consider a simpler system without the separation into core and periphery, where the maintenance algorithm simply runs the peer distribution algorithm and the information aggregation algorithm to count the total number of peers in the system, and expands or reduces the hypercube with respect to the thresholds $LT = 8d+16$ and $UT = 40d+80$ presented above. Moreover, assume that these operations are performed in quiet phases, where the adversary may remove at most $d+1$ and add at most $d+1$ peers only in-between.

For this simpler system, it holds that every node in the simulated d -dimensional hypercube has at least $3d+10$ and at most $45d+86$ peers at every moment of time. Moreover, after the hypercube has changed its dimension from d_{old} to d_{new} , the dimension will remain stable for at least $2d_{new}+1$ phases. We will now prove these properties.

We consider the cases where the average number of peers per node μ falls beyond the lower threshold $8d_{old}+16$ or exceeds the upper threshold $40d_{old}+80$ in turn. Note that such an event will lead to a dimension change with a delay of d_{old} phases only, see Lemma 4.6. We prove that after the change, $\mu \in [8d_{new}+16, 40d_{new}+80]$ for at least $d_{new}+1$ phases. The dimension remains stable for at least $2d_{new}+1$ phases which implies — together with Lemma 4.5 — that the discrepancy before the next change is limited by $2(d_{new}+1) + 2(d_{new}+1) + d_{new} = 5d_{new}+4$.

Case $\mu < 8d+16$: At time $t - d_{old}$, it held that $\mu < 8d_{old}+16$ while at time $t - d_{old} - 1$ we had $\mu \geq 8d_{old}+16$. In $d_{old}+1$ phases, there are at most $(d_{old}+1)(d_{old}+1) = d_{old}^2 + 2d_{old} + 1$ leaves, so $\mu \geq 8d_{old}+16 - \frac{d_{old}^2 + 2d_{old} + 1}{2^{d_{old}}} >$

$8d_{old}+14$ before merging. Clearly, there must be a node with more than $8d_{old}+14$ peers, hence, given the discrepancy of $5d_{old}+4$ (see Lemma 4.5), every node has more than $3d_{old}+10$ peers before merging.

What about the maximum? At time $t - d_{old}$, $\mu < 8d_{old}+16$, and there have been at most $d_{old}(d_{old}+1)$ joins in d_{old} steps, so $\mu < 8d_{old}+16 + \frac{d_{old}(d_{old}+1)}{2^{d_{old}}} < 8d_{old}+18$ before merging, and $\mu < 16d_{old}+36$ afterwards. The maximum node has less than $21d_{new}+61$ peers.

Next, we show that $\mu \geq 8d_{new}+16$ for the next $d_{new}+1$ phases after a reduction. At time $t - d_{old} - 1$, $\mu \geq 8d_{old}+16 = 8d_{new}+24$. The reduction doubles the average number of peers per node, so $\mu \geq 16d_{new}+48$. Further, there are at most $(d_{old}+1)(d_{old}+1) + (d_{new}+1)(d_{new}+1) = 2d_{new}^2 + 6d_{new} + 5$ leaves in the meantime, so $\mu \geq 16d_{new}+48 - \frac{2d_{new}^2 + 6d_{new} + 5}{2^{d_{new}}} > 16d_{new}+41 > 8d_{new}+16$.

Finally, $\mu \leq 40d_{new}+80$ for $d_{new}+1$ phases. At time $t - d_{old}$, $\mu < 8d_{new}+24$, so $\mu < 16d_{new}+48$ after the reduction. There are at most $d_{old}(d_{old}+1) + (d_{new}+1)(d_{new}+1) = 2d_{new}^2 + 5d_{new} + 3$ joins, so $\mu < 16d_{new}+48 + \frac{2d_{new}^2 + 5d_{new} + 3}{2^{d_{new}}} < 16d_{new}+54 < 40d_{new}+80$.

Case $\mu > 40d+80$: At time $t - d_{old}$, $\mu > 40d_{old}+80 = 40d_{new}+40$, so $\mu > 20d_{new}+20$ after splitting; there are at most $d_{old}(d_{old}+1) = d_{new}^2 - d_{new}$ leaves in d_{old} steps, so $\mu > 20d_{new}+20 - \frac{d_{new}^2 - d_{new}}{2^{d_{new}}} > 20d_{new}+19$. According to Lemma 4.5, the minimum node has more than $15d_{new}+15$ peers after splitting. At time $t - d_{old} - 1$, $\mu \leq 40d_{old}+80$, and there are at most $(d_{old}+1)(d_{old}+1) = d_{old}^2 + 2d_{old} + 1$ joins. So before splitting, $\mu \leq 40d_{old}+80 + \frac{d_{old}^2 + 2d_{old} + 1}{2^{d_{old}}} < 40d_{old}+82$, and the maximum node has at most $45d_{old}+86$ peers.

Next, we show that $\mu \geq 8d_{new}+16$ for the next $d_{new}+1$ phases after the expansion. At time $t - d_{old}$, $\mu > 40d_{old}+80 = 40d_{new}+40$, so $\mu > 20d_{new}+20$ after the expansion. Moreover, there are at most $d_{old}(d_{old}+1) + (d_{new}+1)(d_{new}+1) = 2d_{new}^2 + d_{new} + 1$ leaves, and $\mu > 20d_{new}+20 - \frac{2d_{new}^2 + d_{new} + 1}{2^{d_{new}}} > 20d_{new}+17 \geq 8d_{new}+16$. Finally, $\mu \leq 40d_{new}+80$ for the next $d_{new}+1$ steps: At time $t - d_{old} - 1$, $\mu \leq 40d_{old}+80 = 40d_{new}+40$, so $\mu \leq 20d_{new}+20$ after the expansion; moreover, there are at most $(d_{old}+1)(d_{old}+1) + (d_{new}+1)(d_{new}+1) = 2d_{new}^2 + 2d_{new} + 1$ joins, so $\mu \leq 20d_{new}+20 + \frac{2d_{new}^2 + 2d_{new} + 1}{2^{d_{new}}} < 20d_{new}+24 < 40d_{new}+80$.

In our real system, repairing takes six rounds and runs *concurrently* to the adversary. However, as all operations in the whole phase are based upon the state of ROUND 1, a phase can be considered as running uninterruptedly, that is, as if the adversary inserted $d+1$ and removed $d+1$ peers only *between* the phases. Thus, the properties shown above

also hold in our system. However, we additionally have to postulate that there is always at least one *core peer*. We know that it is always possible to select $2d + 3$ core peers in ROUND 5 with respect to the state of ROUND 1. These peers have to survive until ROUND 6 of the next phase, so for twelve normal rounds in total; however, as the adversary $\mathcal{A}_{adv}(d + 1, d + 1, 6)$ may remove at most $2d + 2$ peers in twelve rounds, this clearly holds.

Finally, we show that there are indeed enough peripheral peers in ROUND 3 such that core peers do not have to change the node for the peer distribution, that is: In ROUND 3, it holds that $|\mathcal{P}_v| > \frac{size(v) - size(u)}{2}$. From the considerations made above, we know that $size(v) \geq 3d + 10$ and $size(u) \geq 3d + 10$. As v has at most $2d + 3$ core peers, we have $|\mathcal{P}_v| \geq size(v) - (2d + 3) \geq size(v) - size(u) > \frac{size(v) - size(u)}{2}$. \square

In order to enhance clarity, we described a scheme which is as simple as possible. Instead of a complete bipartite graph between adjacent hypercube nodes one could e.g. use a bipartite matching. This reduces the node degree from $O(\log^2 n)$ to $O(\log n)$. Apart from better node degrees, all our results still hold up to constant factors.

6 Conclusions

We presented a first distributed hash table which provably tolerates dynamic worst-case joins and leaves. We believe that our approach opens several exciting P2P research challenges. For example: How well perform classic P2P proposals when studied with a dynamic failure model or what is the adversary/efficiency tradeoff when studying dynamic models?

References

- [1] K. Aberer. P-Grid: A Self-Organizing Access Structure for P2P Information Systems. In *Proc. 9th Int. Conference on Cooperative Information Systems (CoopIS)*, pages 179–194, 2001.
- [2] I. Abraham, B. Awerbuch, Y. Azar, Y. Bartal, D. Malkhi, and E. Pavlov. A Generic Scheme for Building Overlay Networks in Adversarial Scenarios. In *Proc. 17th Int. Symp. on Parallel and Distributed Processing (IPDPS)*, page 40.2, 2003.
- [3] I. Abraham, D. Malkhi, and O. Dobzinski. LAND: Stretch $(1 + \epsilon)$ Locality-Aware Networks for DHTs. In *Proc. 15th Ann. ACM-SIAM Symp. on Discrete Algorithms (SODA)*, pages 550–559, 2004.
- [4] K. Albrecht, R. Arnold, M. Gähwiler, and R. Wattenhofer. Aggregating Information in Peer-to-Peer Systems for Improved Join and Leave. In *4th IEEE Int. Conference on Peer-to-Peer Computing (P2P)*, 2004.
- [5] J. Aspnes and G. Shah. Skip Graphs. In *Proc. 14th Ann. ACM-SIAM Symp. on Discrete Algorithms (SODA)*, pages 384–393, 2003.
- [6] B. Awerbuch and C. Scheideler. The Hyperring: A Low-Congestion Deterministic Data Structure for Distributed Environments. In *Proc. 15th Ann. ACM-SIAM Symp. on Discrete Algorithms (SODA)*, pages 318–327, 2004.
- [7] G. Cybenko. Dynamic Load Balancing for Distributed Memory Multiprocessors. *Journal on Parallel Distributed Computing*, 7:279–301, 1989.
- [8] A. Fiat and J. Saia. Censorship Resistant Peer-to-Peer Content Addressable Networks. In *Proc. 13th Symp. on Discrete Algorithms (SODA)*, 2002.
- [9] N. J. A. Harvey, M. B. Jones, S. Saroiu, M. Theimer, and A. Wolman. SkipNet: A Scalable Overlay Network with Practical Locality Properties. In *Proc. 4th USENIX Symp. on Internet Technologies and Systems (USITS)*, 2003.
- [10] J. Kubiatiowicz, D. Bindel, Y. Chen, P. Eaton, D. Geels, R. Gummadi, S. Rhea, H. Weatherspoon, W. Weimer, C. Wells, and B. Zhao. OceanStore: An Architecture for Global-scale Persistent Storage. In *Proc. of ACM ASPLOS*, November 2000.
- [11] X. Li, J. Misra, and C. G. Plaxton. Active and Concurrent Topology Maintenance. In *Proc. 18th Ann. Conference on Distributed Computing (DISC)*, 2004.
- [12] D. Malkhi, M. Naor, and D. Ratajczak. Viceroy: A Scalable and Dynamic Emulation of the Butterfly. In *Proc. 21st Ann. Symp. on Principles of Distributed Computing (PODC)*, pages 183–192, 2002.
- [13] D. Peleg and E. Upfal. The Token Distribution Problem. *SIAM Journal on Computing*, 18(2):229–243, 1989.
- [14] C. G. Plaxton. Load Balancing, Selection and Sorting on the Hypercube. In *Proc. 1st Ann. ACM Symp. on Parallel Algorithms and Architectures (SPAA)*, pages 64–73, 1989.
- [15] C. G. Plaxton, R. Rajaraman, and A. W. Richa. Accessing Nearby Copies of Replicated Objects in a Distributed Environment. In *Proc. 9th Ann. ACM Symp. on Parallel Algorithms and Architectures (SPAA)*, pages 311–320, 1997.
- [16] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker. A Scalable Content Addressable Network. In *Proc. of ACM SIGCOMM 2001*, 2001.
- [17] S. Rhea, D. Geels, T. Roscoe, and J. Kubiatiowicz. Handling Churn in a DHT. In *Proc. USENIX Ann. Technical Conference*, 2004.
- [18] J. Saia, A. Fiat, S. Gribble, A. Karlin, and S. Saroiu. Dynamically Fault-Tolerant Content Addressable Networks. In *Proc. 1st Int. Workshop on Peer-to-Peer Systems (IPTPS)*, 2002.
- [19] B. A. Shirazi, K. M. Kavi, and A. R. Hurson. *Scheduling and Load Balancing in Parallel and Distributed Systems*. IEEE Computer Science Press, 1995.
- [20] 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. ACM SIGCOMM Conference*, 2001.
- [21] R. Van Renesse, K. P. Birman, and W. Vogels. Astrolabe: A Robust and Scalable Technology for Distributed System Monitoring, Management, and Data Mining. *ACM Transactions on Computing Systems*, 21(2):164–206, 2003.
- [22] R. van Renesse and A. Bozdog. Willow: DHT, Aggregation, and Publish/Subscribe in One Protocol. In *Proc. 3rd Int. Workshop on Peer-To-Peer Systems (IPTPS)*, 2004.
- [23] B. Y. Zhao, L. Huang, J. Stribling, A. D. Joseph, and J. D. Kubiatiowicz. Tapestry: A Resilient Global-scale Overlay for Service Deployment. *IEEE Journal on Selected Areas in Communications*, 22(1), 2004.