

## 5 Distributed Data Management II – Caching

In this section we will study the approach of using caching for the management of data in distributed systems. Caching always tries to keep data at the place where it is needed. The problem with this approach is that it may not be known in advance which processor will be the next to access a certain data item. Is it nevertheless possible to find good caching strategies? When is a caching strategy “good”? Obviously, a model is needed that allows us to measure the performance of caching strategies.

### 5.1 The model

Let  $G = (V, E)$  be an arbitrary graph, and let  $U$  be the set of all shared objects. We assume that an adversary specifies a distributed application running on the nodes of the network, i.e., the adversary initiates read and write requests to the objects at the nodes of the network  $G$ . These requests must be served by an on-line caching strategy.

We restrict the class of allowed applications specified by the adversary to be *data-race free*, i.e., a write access to an object is not allowed to overlap with other accesses to the same object, and there is some order among the accesses to the same object such that for each read and write access, there is a unique least recent write. Note that this still allows arbitrary concurrent accesses to different objects and concurrent read accesses to the same object.

A caching strategy is called *consistent* if it ensures that a read request directed to an object returns the value of the most recent write access to that object. Write accesses are assumed to be object alterations rather than overwrites. That is, a writing node cannot just build a new copy of the object from scratch, but first has to search for an existing copy of the object.

We are interested in finding caching strategies that are efficient even though they have no prior knowledge of the actions of the adversary. These strategies are allowed to migrate, create, and invalidate copies of an object. We do not allow the use of encoding strategies.

Initially, each object has one copy somewhere in the network. Every message or migration of an object along an edge  $e$  increases its load (also called *congestion* in the following) by one. We want to keep the congestion as small as possible. Let  $\sigma$  be an arbitrary data-race free sequence of read and write requests. Furthermore, let  $C_e^A(\sigma)$  be the total congestion caused at  $e$  when using the caching algorithm  $A$  for  $\sigma$ , let  $C_A(\sigma) = \max_e C_e^A(\sigma)$ , and let  $C_{\text{OPT}}(\sigma) = \min_A C_A(\sigma)$ . Algorithm  $A$  is called *c-competitive* if there is a constant  $d$  so that for all sequences of requests  $\sigma$  we have

$$C_A(\sigma) \leq c \cdot C_{\text{OPT}}(\sigma) + d .$$

### 5.2 Caching in trees

Let  $T = (V, E)$  be an arbitrary tree. Our caching strategy for  $T$  works in the following way:

- $v$  wants to read  $x$ :  $v$  sends a request to the closest node  $u$  in  $T$  that has a copy of  $x$ . Upon receipt of the request,  $u$  sends a copy of  $x$  back to  $v$ . Every node visited by the copy will store a copy of  $x$ .
- $v$  wants to (re)write  $x$ :  $v$  sends a (re)write request for  $x$  to the closest node  $u$  in  $T$  that has a copy of  $x$ . Upon receipt of the message,  $u$  starts an invalidation broadcast to all other nodes that have

a copy of  $x$ . Afterwards,  $u$  stores the new copy of  $x$  and sends it back to  $v$ . Every node visited by the copy will store a copy of  $x$ .

The strategy maintains the following invariant.

**Fact 5.1** *For every object  $x$  and every time step  $t$ , the nodes that have a copy of  $x$  form a connected component.*

The invariant allows us to efficiently locate the nodes that store copies of an object  $x$ : every node has a sign post that points to the last node that issued a write request. (Initially, all sign posts point to the unique copy of  $x$ .) Additionally, we set markers at the border nodes of the connected component, so that the invalidation broadcast is confined to the nodes of the connected component. Although setting sign posts everywhere seems to be a wasteful strategy, observe that the only node at which sign posts have to point downwards are the nodes from the root of the tree down to the highest node in the tree storing a copy of  $x$ . Hence, when using the strategy of only storing information for sign posts that point downwards, we only have to store at most  $D$  sign posts for any object  $x$ , where  $D$  is the depth of the tree. Our simple caching strategy achieves the following remarkable result.

**Theorem 5.2** *The caching strategy is 3-competitive.*

**Proof.** Since the caching algorithm treats the objects independent of each other, it suffices to show that it is 3-competitive for any single object. So let us consider some fixed object  $x$ . Let  $e = (a, b)$  be any edge in the tree. Removing  $e$  from the tree breaks it down into two subtrees,  $T_a$  and  $T_b$ . We distinguish between three cases:

- $[A]$ : All copies are in  $T_a$ .
- $[B]$ : All copies are in  $T_b$ .
- $[AB]$ : Both subtrees contain copies. (Observe that this implies that  $a$  and  $b$  hold a copy.)

Consider any sequence  $\sigma$  of read and write requests, and let  $c_t$  be the configuration ( $[A]$ ,  $[B]$ , or  $[AB]$ ) after processing the  $t$ th request in  $\sigma$ . Then the sequence  $c_0, c_1, c_2, \dots$  is of the form

$$\dots [?]^+ [AB]^+ [?]^+ [AB]^+ [?]^+ [AB]^+ \dots$$

where  $[?]$  is a placeholder for  $[A]$  or  $[B]$  and  $[X]^+$  means any sequence of  $X$ 's of length at least 1. Without loss of generality, consider any period of the form  $[A]^+ [AB]^+$ . We investigate the online and offline cost of edge  $e$  during that period. For the online cost we get:

subphase	kind of request	prev. config.	online cost
$[A]^+$	starts with write from $T_a$ (*)	$[AB]^+$	1
	followed by requests from $T_a$	$[A]^+$	0
$[AB]^+$	starts with request from $T_b$ (**)	$[A]^+$	2
	followed by reads from $T_a$ or $T_b$	$[AB]^+$	0

What is the optimal offline cost in this period? Suppose that there exists an offline strategy with cost 0. Then request (\*) requires that there is no copy in  $T_b$ . Furthermore, no copy is migrated. Consequently, request (\*\*) has cost at least 1, contradicting our assumption. Therefore, the online cost in each period is 3 whereas the offline cost is at least 1.  $\square$

Notice that if we count the work for sending requests as 0 (which is a realistic assumption for large objects), then the competitiveness of the online strategy would even be 2.

### 5.3 Caching in meshes

Consider an arbitrary  $n \times n$ -mesh  $M$ . Our caching strategy for the mesh is based on a hierarchical decomposition of its nodes that is done in the following recursive way:

If  $M$  only consists of a single node, we are done. Otherwise, we partition  $M$  into two submeshes along the dimension with the most nodes (see Figure 1).

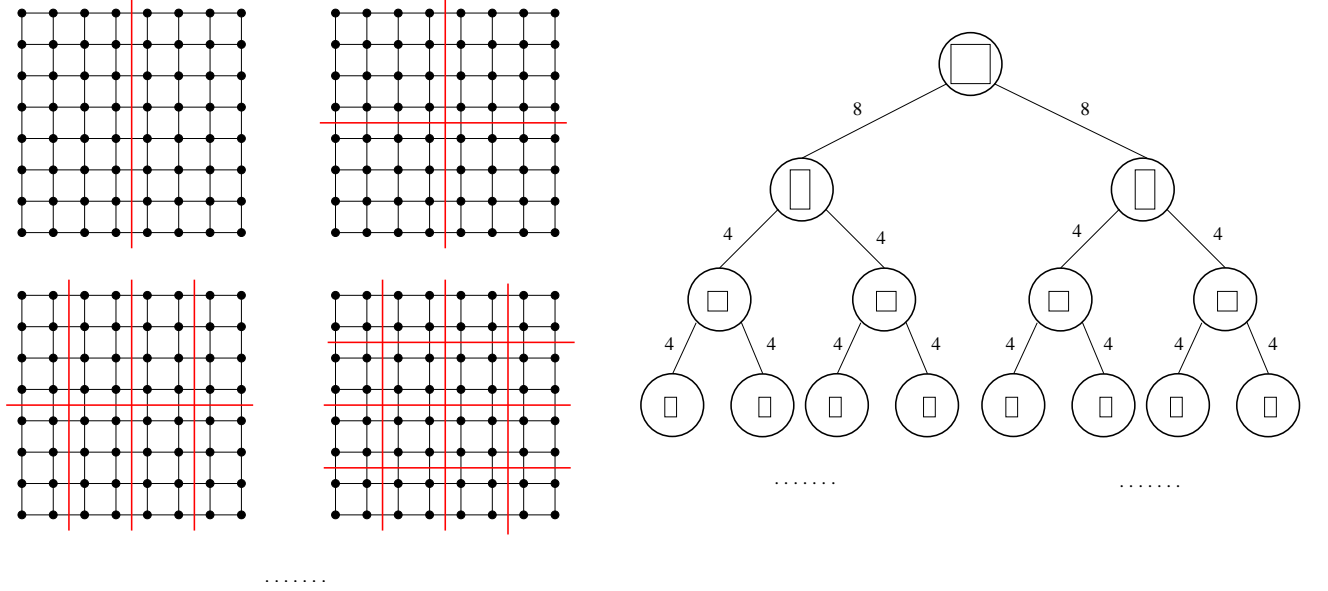


Figure 1: The decomposition of  $M$  into submeshes.

This hierarchical decomposition can be represented as a *decomposition tree*  $T(M)$ . Its root represents  $M$ , every internal node of  $T(M)$  represents a submesh of  $M$ , every leaf a node of  $M$ , and every edge a partition of a submesh into two submeshes. Thus,  $T(M)$  is a binary tree of depth  $O(\log n)$ . We view now  $T(M)$  as a virtual network that we want to simulate by  $M$ . To be able to compare the congestion in both networks, we define the *bandwidth*  $b(e)$  of an edge  $e = \{v, w\}$  in  $T(M)$  (where  $v$  is the father of  $w$ ) as the number of edges that leave the submesh  $M(w)$ .

For every object  $x$ , we define an access tree  $T_x(M)$  of  $x$  as a copy of the decomposition tree  $T(M)$ . We embed  $T_x(M)$  randomly into  $M$  with the help of a (pseudo-)random hash function  $h$  so that we can recover the positions of the nodes in  $T_x(M)$  at any time. Using  $h$ , every node  $v$  in  $T_x(M)$  is mapped to a random processor  $h(v)$  in  $M(v)$  and every edge  $\{v, w\}$  in  $T_x(M)$  is mapped to a path from  $h(v)$  to  $h(w)$  using the  $x - y$  routing strategy presented in Section 2.1. This allows us to simulate the movement of any message in  $T_x(M)$  by moving it along the corresponding paths in  $M$ .

We apply now to  $T_x(M)$  the same caching strategy that we used for the simple tree in Section 5.2. When comparing our caching strategy with an optimal *static* placement strategy (i.e. every object is stored at a single, fixed node in  $M$ ), we obtain the following theorem.

**Theorem 5.3 ([1])** *The caching strategy for the  $n \times n$ -mesh is  $O(\log n)$ -competitive, with high probability.*

**Proof.** In order to prove the above result we investigate two caching simulations. First, we consider the problem of simulating an optimal caching strategy in the mesh  $M$  by the tree  $T(M)$ , and then we consider the problem of simulating a caching strategy in  $T(M)$  by the mesh  $M$ .

In order to simulate an optimal caching strategy in the mesh by  $T(M)$ , we use the simple strategy that whenever a request is sent from node  $v$  to  $w$  in  $M$ , we send it along the unique path from the leaf representing  $v$  to the leaf representing  $w$  in  $T(M)$ .

In order to simulate a caching strategy in  $T(M)$  by  $M$ , we use the strategy already explained above: For every object  $x$  we embed  $T_x(M)$  randomly into  $M$ , and whenever a request is sent along an edge  $\{v, w\}$  in  $T_x(M)$ , we send it along the  $x - y$  path from  $\pi(v)$  to  $\pi(w)$  in  $M$ .

Let  $C_{\text{OPT}}(M)$  be the best possible congestion achievable in the mesh  $M$  for the given application and let  $C_{\text{OPT}}(T(M))$  be the congestion caused in  $T(M)$  when simulating this best possible strategy by  $T(M)$ . Then we get the following result.

**Lemma 5.4**  $C_{\text{OPT}}(T(M)) \leq C_{\text{OPT}}(M)$ .

**Proof.** Let  $e$  denote an edge of  $T(M)$  with bandwidth  $b(e)$ , and let  $C(e)$  be the congestion caused by messages traversing  $e$  when using the simulation strategy above (i.e. the total number of messages traversing  $e$  is  $b(e) \cdot C(e)$ ). Any message that crosses  $e = \{v, w\}$  in  $T(M)$  corresponds to a message that either has to leave or enter the submesh  $M(w)$  in  $M$ . Since there are only  $b(e)$  edges leaving  $M(w)$ , there must be an edge in  $M(w)$  with congestion at least  $(C(e) \cdot b(e))/b(e) = C(e)$ . Thus, the maximum congestion over all edges in  $T(M)$  is at most the maximum congestion over all edges in  $M$ , which proves the lemma.  $\square$

The next lemma gives an upper bound on the expected congestion caused by simulating the access tree strategy on  $M$ . In the following, let  $C_T(e)$  denote the congestion caused at  $e$  for simulating  $T(M)$ .

**Lemma 5.5** For any edge  $e$  of  $M$ ,  $E[C_T(e)] = O(\log n \cdot C_{\text{OPT}}(T(M)))$ .

**Proof.** Let  $h$  denote the height of  $T(M)$  and let  $C_\ell(e)$  denote the congestion caused at  $e$  due to the simulation of edges on level  $\ell$  of  $T(M)$ ,  $1 \leq \ell \leq h$ . We show that  $E[C_\ell(e)] = O(C_{\text{OPT}}(T(M)))$  for all  $\ell \in \{1, \dots, h\}$ , which yields the lemma as  $h = O(\log n)$ .

Consider some fixed level  $\ell$ . Let  $v$  be a node of  $T(M)$  on level  $\ell - 1$  such that  $M(v)$  includes edge  $e$ . (If such a node does not exist then  $E[C_\ell(e)] = 0$ .) Let  $v'$  be one of the two children of  $v$ . We bound the expected congestion on  $e$  due to the simulation of the tree edge  $e_T = \{v, v'\}$  on level  $\ell$  of  $T(M)$ .

Notice that every two levels in  $T(M)$  the side length of the submeshes reduces by a factor of two. Hence,  $M(v)$  is a  $\Theta(n/2^{\ell/2}) \times \Theta(n/2^{\ell/2})$ -mesh. Suppose that the congestion at  $e_T$  is  $C(e_T)$ . Since  $e_T$  has a bandwidth of  $\Theta(n/2^{\ell/2})$ , this means that  $\Theta(C(e_T) \cdot n/2^\ell)$  requests pass through  $e_T$ . Since every object  $x$  chooses a random place for  $v$  and  $v'$  in  $M(v)$  and the  $x - y$  routing strategy only goes through  $e$  if  $v$  is in the row of  $e$  or  $v'$  is in the column of  $e$ , the probability that the path for object  $x$  simulating  $\{v, v'\}$  moves through  $e$  is equal to  $\Theta(2^{\ell/2}/n)$ . So for  $\Theta(C(e_T) \cdot n/2^{\ell/2})$  requests using  $e_T$  the expected congestion at  $e$  is

$$\Theta(C(e_T) \cdot n/2^{\ell/2}) \cdot \Theta(2^{\ell/2}/n) = \Theta(C(e_T)) .$$

Thus,  $E[C_\ell(e)] = O(C(e_T)) = O(C_{\text{OPT}}(T(M)))$ , which proves the lemma.  $\square$

Combining the two lemmata, we get  $E[C_T(e)] = O(\log n \cdot C_{\text{OPT}}(M))$  for all edges  $e$  in  $M$ . One can also show that this bound holds with high probability, using the fact that for an object  $x$  with  $k$  write accesses the congestion in OPT caused for object  $x$  must be at least  $k/4$  because all of these requests must update  $x$  at the same, static location.  $\square$

One can also extend the theorem to the case that OPT is allowed to use dynamic data management (i.e. the location of object  $x$  can change) by choosing a new random location of a node  $v$  in  $T_x(M)$  each time a write request passes  $v$  (and updating the pointers of its neighbors in  $T_x(M)$  correspondingly).

In addition to the upper bound, one can show the following result (see also Section 4).

**Theorem 5.6** *Any online caching strategy on the  $n \times n$ -mesh has a competitive ratio of  $\Omega(\log n)$ .*

Hence, the online caching strategy presented above is asymptotically optimal.

## 5.4 Caching in arbitrary networks

One can observe from the proof for the mesh that there is a general way of coming up with good online caching strategy: Given some network  $G$ , try to come up with a decomposition tree  $T(G)$  for  $G$  so that

- $C_{\text{OPT}}(T(G)) \leq C_{\text{OPT}}(G)$ , i.e. an optimal strategy in  $G$  can be simulated by  $T(G)$  with at most the congestion the optimal strategy needs, and
- for any edge  $e$  of  $G$ ,  $E[C(e)] = O(\gamma \cdot C_{\text{OPT}}(T(G)))$ , i.e. there is a way to simulate a caching strategy in  $T(G)$  by  $G$  so that the congestion increases by a factor of at most  $\gamma$ .

The first item is easy to achieve, because no matter how  $G$  is decomposed to form  $T(G)$ , Lemma 5.4 is always true. The difficult part is the second item. Here, it is important to come up with a tree  $T(G)$  so that

- the height  $h$  of  $T(G)$  is low, and
- for every edge  $e$  in  $G$  and every level  $\ell$  in  $T(G)$  it holds that  $E[C_\ell(e)] = O(\beta \cdot C(e_T))$  for some low  $\beta$ , where  $e_T$  is an edge from level  $\ell - 1$  to  $\ell$  in  $T(G)$  whose  $(\ell - 1)$ -level node contains  $e$ .

In this case we would get  $\gamma = O(h \cdot \beta)$ . For further information on this subject see [3, 4].

One may ask whether it is always possible to find decompositions so that  $h$  and  $\beta$  are polylogarithmic. Surprisingly, Räcke recently showed that this is possible:

**Theorem 5.7 ([2])** *For every network  $G$  with non-negative capacities there is an online caching strategy that achieves for every application a congestion of  $O(C_{\text{OPT}}(G) \cdot \text{polylog}(n))$ .*

Hence, even if no knowledge is available about the requests issued by the processors, efficient online caching is possible for arbitrary networks.

## References

- [1] B. Maggs, F. M. auf der Heide, B. Vöcking, and M. Westermann. Exploiting locality for networks of limited bandwidth. In *Proc. of the 38th IEEE Symp. on Foundations of Computer Science (FOCS)*, pages 284–293, 1997.
- [2] H. Räcke. Minimizing congestion in general networks. In *Proc. of the 43rd IEEE Symp. on Foundations of Computer Science (FOCS)*, 2002.
- [3] B. Vöcking. *Static and Dynamic Data Management in Networks*. PhD thesis, Dept. of Mathematics and Computer Science, University of Paderborn, October 1998.
- [4] M. Westermann. *Caching in Networks: Non-Uniform Algorithms and Memory Capacity Constraints*. PhD thesis, Dept. of Mathematics and Computer Science, University of Paderborn, November 2000.