
Distributed Network Algorithms

(Lecture Notes for GIAN Course)

Assoc.-Prof. Dr. Stefan Schmid

Assoc-Prof. Dr. Partha Sarathi Mandal

**Thanks to Prof. Dr. Roger Wattenhofer and
Prof. Dr. Christian Scheideler for basis of manuscript!**

Summer 2016

Contents

1	Topology & Routing	3
2	Vertex Coloring	9
2.1	Introduction	9
2.2	Coloring Trees	11
3	Leader Election	17
3.1	Distributed Algorithms and Complexity	17
3.2	Anonymous Leader Election	18
3.3	Asynchronous Ring	19
3.4	Lower Bounds	21
3.5	Synchronous Ring	23
4	Tree Algorithms	25
4.1	Broadcast	25
4.2	Convergecast	26
4.3	BFS Tree Construction	27
4.4	MST Construction	29
5	Distributed Sorting	33
5.1	Array & Mesh	33
5.2	Sorting Networks	36
5.3	Counting Networks	39
6	Shared Memory	45
6.1	Introduction	45
6.2	Mutual Exclusion	46
6.3	Store & Collect	49
6.3.1	Problem Definition	49
6.3.2	Splitters	50
6.3.3	Binary Splitter Tree	51
6.3.4	Splitter Matrix	53
7	Shared Objects	55
7.1	Introduction	55
7.2	Arrow and Friends	56
7.3	Ivy and Friends	61

8	Maximal Independent Set	65
8.1	MIS	65
8.2	Fast MIS from 1986	67
8.3	Fast MIS from 2009	70
8.4	Applications	74
9	Locality Lower Bounds	77
9.1	Locality	78
9.2	The Neighborhood Graph	80
10	Social Networks	85
10.1	Small-World Networks	86
10.2	Propagation Studies	92
11	Synchronization	95
11.1	Basics	95
11.2	Synchronizer α	96
11.3	Synchronizer β	97
11.4	Synchronizer γ	98
11.5	Network Partition	100
11.6	Clock Synchronization	102
12	Hard Problems	105
12.1	Diameter & APSP	105
12.2	Lower Bound Graphs	107
12.3	Communication Complexity	110
12.4	Distributed Complexity Theory	116
13	Stabilization	117
13.1	Self-Stabilization	117
13.2	Advanced Stabilization	122
14	Wireless Protocols	125
14.1	Basics	125
14.2	Initialization	127
14.2.1	Non-Uniform Initialization	127
14.2.2	Uniform Initialization with CD	127
14.2.3	Uniform Initialization without CD	128
14.3	Leader Election	129
14.3.1	With High Probability	129
14.3.2	Uniform Leader Election	129
14.3.3	Fast Leader Election with CD	130
14.3.4	Even Faster Leader Election with CD	131
14.3.5	Lower Bound	133
14.3.6	Uniform Asynchronous Wakeup without CD	133
14.4	Useful Formulas	134
15	All-to-All Communication	137

16 Consensus	143
16.1 Impossibility of Consensus	143
16.2 Randomized Consensus	148
17 Multi-Core Computing	153
17.1 Introduction	153
17.1.1 The Current State of Concurrent Programming	153
17.2 Transactional Memory	155
17.3 Contention Management	156
18 Dominating Set	163
18.1 Sequential Greedy Algorithm	164
18.2 Distributed Greedy Algorithm	165
19 Routing	171
19.1 Array	171
19.2 Mesh	172
19.3 Routing in the Mesh with Small Queues	173
19.4 Hot-Potato Routing	174
19.5 More Models	176

Welcome!

Welcome to the GIAN course on Distributed Network Algorithms!

Distributed network algorithms play a major role in many networked systems, ranging from computer networks (such as sensor networks, peer-to-peer networks, software-defined networks, datacenter networks, networks on chip) to social or even biological networks.

The goal of this course is to introduce the fundamental formal models and methods needed to reason about the correctness and performance of distributed network algorithms. In particular, we will teach essential algorithmic and analytic techniques which, after the course, are a useful toolbox and allow the students to develop and study their own distributed network algorithms. The course also discusses applications, e.g., in the context of wireless networks, peer-to-peer networks, datacenter networks, and virtual networks.

Due to time constraints, we will focus on some selected aspects, and we will not be able to cover all the material you find in these lecture notes. However, I will be very happy to talk to you about any other related topic during this week in person.

Enjoy the course!

Stefan Schmid (Aalborg University, Denmark & TU Berlin, Germany)

Partha Sarathi Mandal (IIT Guwahati, India)

Chapter 1

Topology & Routing

The most basic network topologies used in practice are trees, rings, grids or tori. Many other suggested networks are simply combinations or derivatives of these. The advantage of trees is that the routing is very easy: for every source-destination pair there is only one possible simple path. However, since the root of a tree is usually a severe bottleneck, so-called *fat trees* have been used. These trees have the property that every edge connecting a node v to its parent u has a capacity that is equal to all leaves of the subtree rooted at v . See Figure 1.1 for an example.

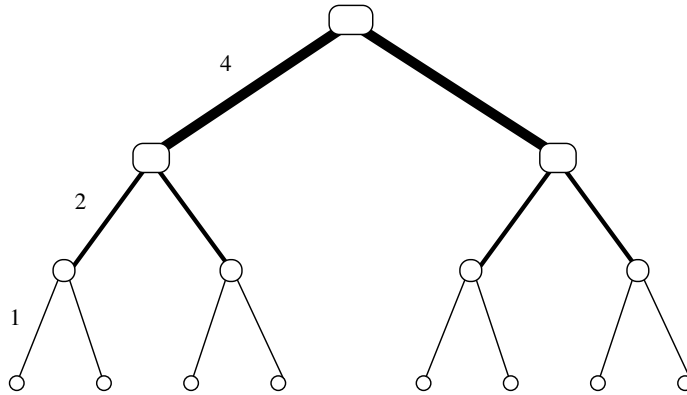


Figure 1.1: The structure of a fat tree.

Remarks:

- Fat trees belong to a family of networks that require edges of non-uniform capacity to be efficient. Easier to build are networks with edges of uniform capacity. This is usually the case for grids and tori. Unless explicitly mentioned, we will treat all edges in the following to be of capacity 1. In the following, $[x]$ means the set $\{0, \dots, x-1\}$.

Definition 1.1 (Torus, Mesh). Let $m, d \in \mathbb{N}$. The (m, d) -mesh $M(m, d)$ is a

graph with node set $V = [m]^d$ and edge set

$$E = \left\{ \{ (a_1, \dots, a_d), (b_1, \dots, b_d) \} \mid a_i, b_i \in [m], \sum_{i=1}^d |a_i - b_i| = 1 \right\}.$$

The (m, d) -torus $T(m, d)$ is a graph that consists of an (m, d) -mesh and additionally wrap-around edges from nodes $(a_1, \dots, a_{i-1}, m, a_{i+1}, \dots, a_d)$ to nodes $(a_1, \dots, a_{i-1}, 1, a_{i+1}, \dots, a_d)$ for all $i \in \{1, \dots, d\}$ and all $a_j \in [m]$ with $j \neq i$. In other words, we take the expression $a_i - b_i$ in the sum modulo m prior to computing the absolute value. $M(m, 1)$ is also called a line, $T(m, 1)$ a cycle, and $M(2, d) = T(2, d)$ a d -dimensional hypercube. Figure 1.2 presents a linear array, a torus, and a hypercube.

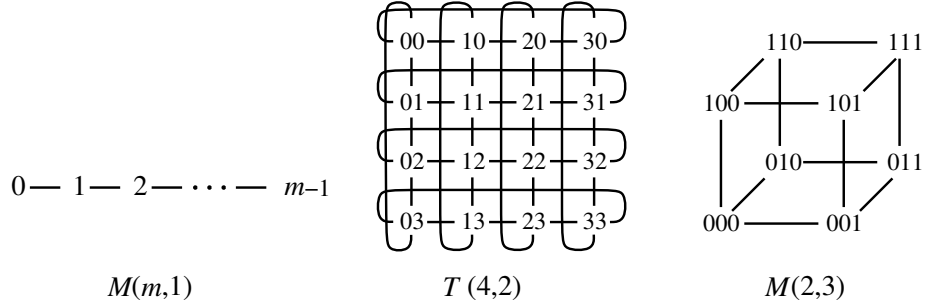


Figure 1.2: The structure of $M(m, 1)$, $T(4, 2)$, and $M(2, 3)$.

Remarks:

- Routing on mesh, torus, and hypercube is trivial. On a d -dimensional hypercube, to get from a source bitstring s to a target bitstring t one only needs to fix each “wrong” bit, one at a time; in other words, if the source and the target differ by k bits, there are $k!$ routes with k hops.
- The hypercube can directly be used for a structured Peer-to-Peer (P2P) architecture. It is trivial to construct a distributed hash table (DHT): We have n nodes, n for simplicity being a power of 2, i.e., $n = 2^d$. As in the hypercube, each node gets a unique d -bit ID, and each node connects to d other nodes, i.e., the nodes that have IDs differing in exactly one bit. Now we use a globally known hash function f , mapping file names to long bit strings; SHA-1 is popular in practice, providing 160 bits. Let f_d denote the first d bits (prefix) of the bitstring produced by f . If a node is searching for file name X , it routes a request message $f(X)$ to node $f_d(X)$. Clearly, node $f_d(X)$ can only answer this request if all files with hash prefix $f_d(X)$ have been previously registered at node $f_d(X)$.
- The hypercube has many derivatives, the so-called *hypercubic networks*. Among these are the butterfly, cube-connected-cycles, shuffle-exchange, and de Bruijn graph. We start with the butterfly, which is basically a “rolled out” hypercube (hence directly providing replication!).

Definition 1.2 (Butterfly). *Let $d \in \mathbb{N}$. The d -dimensional butterfly $BF(d)$ is a graph with node set $V = [d+1] \times [2]^d$ and an edge set $E = E_1 \cup E_2$ with*

$$E_1 = \{(i, \alpha), (i+1, \alpha)\} \mid i \in [d], \alpha \in [2]^d\}$$

and

$$E_2 = \{(i, \alpha), (i+1, \beta)\} \mid i \in [d], \alpha, \beta \in [2]^d, \alpha \text{ and } \beta \text{ differ only at the } i^{\text{th}} \text{ position}\}.$$

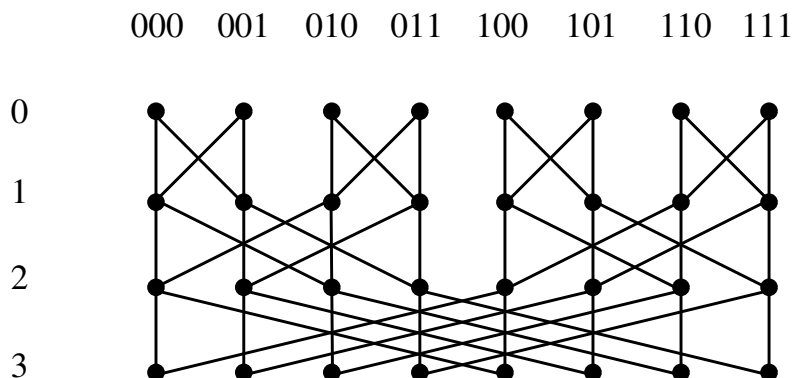
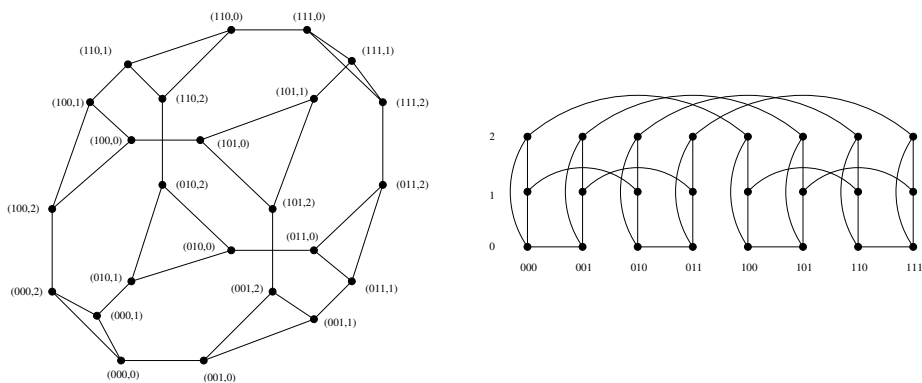
A node set $\{(i, \alpha) \mid \alpha \in [2]^d\}$ is said to form level i of the butterfly. The d -dimensional wrap-around butterfly $W-BF(d)$ is defined by taking the $BF(d)$ and identifying level d with level 0 .

Remarks:

- Figure 1.3 shows the 3-dimensional butterfly $BF(3)$. The $BF(d)$ has $(d+1)2^d$ nodes, $2d \cdot 2^d$ edges and degree 4. It is not difficult to check that combining the node sets $\{(i, \alpha) \mid i \in [d]\}$ into a single node results in the hypercube.
- Butterflies have the advantage of a constant node degree over hypercubes, whereas hypercubes feature more fault-tolerant routing.
- The structure of a butterfly might remind you of sorting networks from Chapter 5. Although butterflies are used in the P2P context (e.g. Viceroy), they have been used decades earlier for communication switches. The well-known Benes network is nothing but two back-to-back butterflies. And indeed, butterflies (and other hypercubic networks) are even older than that; students familiar with fast fourier transform (FFT) will recognize the structure without doubt. Every year there is a new application for which a hypercubic network is the perfect solution!
- Indeed, hypercubic networks are related. Since all structured P2P architectures are based on hypercubic networks, they in turn are all related.
- Next we define the cube-connected-cycles network. It only has a degree of 3 and it results from the hypercube by replacing the corners by cycles.

Definition 1.3 (Cube-Connected-Cycles). *Let $d \in \mathbb{N}$. The cube-connected-cycles network $CCC(d)$ is a graph with node set $V = \{(a, p) \mid a \in [2]^d, p \in [d]\}$ and edge set*

$$E = \{(a, p), (a, (p+1) \bmod d)\} \mid a \in [2]^d, p \in [d]\} \cup \{(a, p), (b, p)\} \mid a, b \in [2]^d, p \in [d], a = b \text{ except for } a_p\}.$$

Figure 1.3: The structure of $BF(3)$.Figure 1.4: The structure of $CCC(3)$.**Remarks:**

- Two possible representations of a CCC can be found in Figure 1.4.
- The shuffle-exchange is yet another way of transforming the hypercubic interconnection structure into a constant degree network.

Definition 1.4 (Shuffle-Exchange). *Let $d \in \mathbb{N}$. The d -dimensional shuffle-exchange $SE(d)$ is defined as an undirected graph with node set $V = [2]^d$ and an edge set $E = E_1 \cup E_2$ with*

$$E_1 = \{ \{ (a_1, \dots, a_d), (a_1, \dots, \bar{a}_d) \} \mid (a_1, \dots, a_d) \in [2]^d, \bar{a}_d = 1 - a_d \}$$

and

$$E_2 = \{ \{ (a_1, \dots, a_d), (a_d, a_1, \dots, a_{d-1}) \} \mid (a_1, \dots, a_d) \in [2]^d \} .$$

Figure 1.5 shows the 3- and 4-dimensional shuffle-exchange graph.

Definition 1.5 (DeBruijn). *The b -ary DeBruijn graph of dimension d $DB(b, d)$ is an undirected graph $G = (V, E)$ with node set $V = \{v \in [b]^d\}$ and edge set*

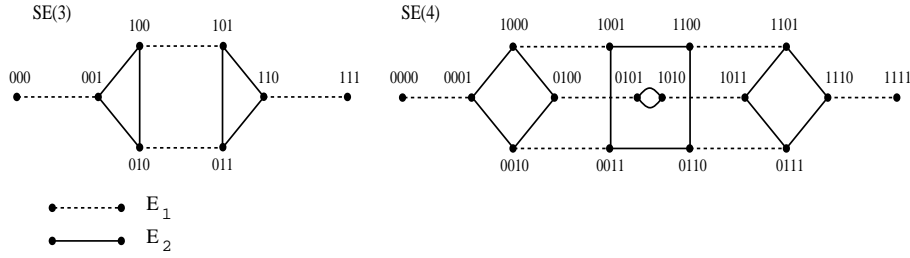


Figure 1.5: The structure of SE(3) and SE(4).

E that contains all edges $\{v, w\}$ with the property that $w \in \{(x, v_1, \dots, v_{d-1}) : x \in [b]\}$, where $v = (v_1, \dots, v_d)$.

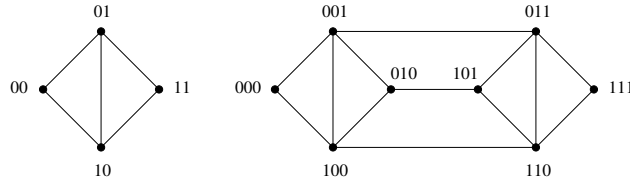


Figure 1.6: The structure of DB(2, 2) and DB(2, 3).

Remarks:

- Two examples of a DeBruijn graph can be found in Figure 1.6. The DeBruijn graph is the basis of the Koorde P2P architecture.
- There are some data structures which also qualify as hypercubic networks. An obvious example is the Chord P2P architecture, which uses a slightly different hypercubic topology. A less obvious (and therefore good) example is the skip list, the balanced binary search tree for the lazy programmer:

Definition 1.6 (Skip List). *The skip list is an ordinary ordered linked list of objects, augmented with additional forward links. The ordinary linked list is the level 0 of the skip list. In addition, every object is promoted to level 1 with probability $1/2$. As for level 0, all level 1 objects are connected by a linked list. In general, every object on level i is promoted to the next level with probability $1/2$. A special start-object points to the smallest/first object on each level.*

Remarks:

- Search, insert, and delete can be implemented in $O(\log n)$ expected time in a skip list, simply by jumping from higher levels to lower ones when overshooting the searched position. Also, the amortized memory cost of each object is constant, as on average an object only has two forward pointers.

- The randomization can easily be discarded, by deterministically promoting a constant fraction of objects of level i to level $i + 1$, for all i . When inserting or deleting, object o simply checks whether its left and right level i neighbors are being promoted to level $i + 1$. If none of them is, promote object o itself. Essentially we establish a MIS on each level, hence at least every third and at most every second object is promoted.
- There are obvious variants of the skip list, e.g., the skip graph. Instead of promoting only half of the nodes to the next level, we always promote all the nodes, similarly to a balanced binary tree: All nodes are part of the root level of the binary tree. Half the nodes are promoted left, and half the nodes are promoted right, on each level. Hence on level i we have 2^i lists (or, more symmetrically: rings) of about $n/2^i$ objects. This is pretty much what we need for a nice hypercubic P2P architecture.
- One important goal in choosing a topology for a network is that it has a small diameter. The following theorem presents a lower bound for this.

Theorem 1.7. *Every graph of maximum degree $d > 2$ and size n must have a diameter of at least $\lceil (\log n)/(\log(d-1)) \rceil - 2$.*

Proof. Suppose we have a graph $G = (V, E)$ of maximum degree d and size n . Start from any node $v \in V$. In a first step at most d other nodes can be reached. In two steps at most $d \cdot (d-1)$ additional nodes can be reached. Thus, in general, in at most k steps at most

$$1 + \sum_{i=0}^{k-1} d \cdot (d-1)^i = 1 + d \cdot \frac{(d-1)^k - 1}{(d-1) - 1} \leq \frac{d \cdot (d-1)^k}{d-2}$$

nodes (including v) can be reached. This has to be at least n to ensure that v can reach all other nodes in V within k steps. Hence,

$$(d-1)^k \geq \frac{(d-2) \cdot n}{d} \quad \Leftrightarrow \quad k \geq \log_{d-1}((d-2) \cdot n/d).$$

Since $\log_{d-1}((d-2)/d) > -2$ for all $d > 2$, this is true only if $k \geq \lceil (\log n)/(\log(d-1)) \rceil - 2$. \square

Remarks:

- In other words, constant-degree hypercubic networks feature an asymptotically optimal diameter.
- There are a few other interesting graph classes, e.g., expander graphs (an expander graph is a sparse graph which has high connectivity properties, that is, from every not too large subset of nodes you are connected to a larger set of nodes), or small-world graphs (popular representations of social networks). At first sight hypercubic networks seem to be related to expanders and small-world graphs, but they are not.

Chapter 2

Vertex Coloring

2.1 Introduction

Vertex coloring is an infamous graph theory problem. Vertex coloring does have quite a few practical applications, for example in the area of wireless networks where coloring is the foundation of so-called TDMA MAC protocols. Generally speaking, vertex coloring is used as a means to break symmetries, one of the main themes in distributed computing. In this chapter we will not really talk about vertex coloring applications but treat the problem abstractly. At the end of the class you probably learned the fastest (but not constant!) algorithm ever! Let us start with some simple definitions and observations.

Problem 2.1 (Vertex Coloring). *Given an undirected graph $G = (V, E)$, assign a color c_u to each vertex $u \in V$ such that the following holds: $e = (v, w) \in E \Rightarrow c_v \neq c_w$.*

Remarks:

- Throughout this course, we use the terms *vertex* and *node* interchangeably.
- The application often asks us to use few colors! In a TDMA MAC protocol, for example, less colors immediately imply higher throughput. However, in distributed computing we are often happy with a solution which is sub-optimal. There is a tradeoff between the optimality of a solution (efficacy), and the work/time needed to compute the solution (efficiency).

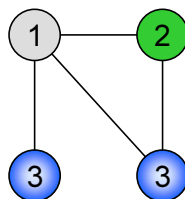


Figure 2.1: 3-colorable graph with a valid coloring.

Assumption 2.2 (Node Identifiers). *Each node has a unique identifier, e.g., its IP address. We usually assume that each identifier consists of only $\log n$ bits if the system has n nodes.*

Remarks:

- Sometimes we might even assume that the nodes exactly have identifiers $1, \dots, n$.
- It is easy to see that node identifiers (as defined in Assumption 3.11) solve the coloring problem 2.1, but not very well (essentially requiring n colors). How many colors are needed at least is a well-studied problem.

Definition 2.3 (Chromatic Number). *Given an undirected Graph $G = (V, E)$, the chromatic number $\chi(G)$ is the minimum number of colors to solve Problem 2.1.*

To get a better understanding of the vertex coloring problem, let us first look at a simple non-distributed (“centralized”) vertex coloring algorithm:

Algorithm 1 Greedy Sequential

```

1: while  $\exists$  uncolored vertex  $v$  do
2:   color  $v$  with the minimal color (number) that does not conflict with the
     already colored neighbors
3: end while

```

Definition 2.4 (Degree). *The number of neighbors of a vertex v , denoted by $\delta(v)$, is called the degree of v . The maximum degree vertex in a graph G defines the graph degree $\Delta(G) = \Delta$.*

Theorem 2.5 (Analysis of Algorithm 1). *The algorithm is correct and terminates in n “steps”. The algorithm uses $\Delta + 1$ colors.*

Proof: Correctness and termination are straightforward. Since each node has at most Δ neighbors, there is always at least one color free in the range $\{1, \dots, \Delta + 1\}$.

Remarks:

- For many graphs coloring can be done with much less than $\Delta + 1$ colors.
- This algorithm is not distributed at all; only one processor is active at a time. Still, maybe we can use the simple idea of Algorithm 1 to define a distributed coloring subroutine that may come in handy later.

Now we are ready to study distributed algorithms for this problem. The following procedure can be executed by every vertex v in a distributed coloring algorithm. The goal of this subroutine is to improve a given initial coloring.

Procedure 2 First Free

Require: Node Coloring {e.g., node IDs as defined in Assumption 3.11}Give v the smallest admissible color {i.e., the smallest node color not used by any neighbor}

Algorithm 3 Reduce

```

1: Assume that initially all nodes have ID's (Assumption 3.11)
2: Each node  $v$  executes the following code
3: node  $v$  sends its ID to all neighbors
4: node  $v$  receives IDs of neighbors
5: while node  $v$  has an uncolored neighbor with higher ID do
6:   node  $v$  sends "undecided" to all neighbors
7:   node  $v$  receives new decisions from neighbors
8: end while
9: node  $v$  chooses a free color using subroutine First Free (Procedure 2)
10: node  $v$  informs all its neighbors about its choice

```

Remarks:

- With this subroutine we have to make sure that two adjacent vertices are not colored at the same time. Otherwise, the neighbors may at the same time conclude that some small color c is still available in their neighborhood, and then at the same time decide to choose this color c .

Theorem 2.6 (Analysis of Algorithm 3). *Algorithm 3 is correct and has time complexity n . The algorithm uses $\Delta + 1$ colors.*

Remarks:

- Quite trivial, but also quite slow.
- However, it seems difficult to come up with a fast algorithm.
- Maybe it's better to first study a simple special case, a tree, and then go from there.

2.2 Coloring Trees

Lemma 2.7. $\chi(\text{Tree}) \leq 2$

Constructive Proof: If the distance of a node to the root is odd (even), color it 1 (0). An odd node has only even neighbors and vice versa. If we assume that each node knows its parent (root has no parent) and children in a tree, this constructive proof gives a very simple algorithm:

Remarks:

- With the proof of Lemma 2.7, Algorithm 4 is correct.
- How can we determine a root in a tree if it is not already given? We will figure that out later.

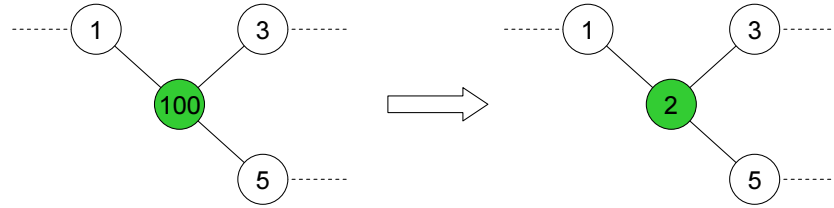


Figure 2.2: Vertex 100 receives the lowest possible color.

Algorithm 4 Slow Tree Coloring

-
- 1: Color the root 0, root sends 0 to its children
 - 2: **Each node** v concurrently executes the following code:
 - 3: **if** node v receives a message x (from parent) **then**
 - 4: node v chooses color $c_v = 1 - x$
 - 5: node v sends c_v to its children (all neighbors except parent)
 - 6: **end if**
-

- The time complexity of the algorithm is the height of the tree.
- If the root was chosen unfortunately, and the tree has a degenerated topology, the time complexity may be up to n , the number of nodes.
- Also, this algorithm does not need to be synchronous ...!

Theorem 2.8 (Analysis of Algorithm 4). *Algorithm 4 is correct. If each node knows its parent and its children, the (asynchronous) time complexity is the tree height which is bounded by the diameter of the tree; the message complexity is $n - 1$ in a tree with n nodes.*

Remarks:

- In this case the asynchronous time complexity is the same as the synchronous time complexity.
- Nice trees, e.g. balanced binary trees, have logarithmic height, that is we have a logarithmic time complexity.
- This algorithm is not very exciting. Can we do better than logarithmic?!?

The following algorithm terminates in $\log^* n$ time. Log-Star?! That's the number of logarithms (to the base 2) you need to take to get down to at least 2, starting with n :

Definition 2.9 (Log-Star).

$$\forall x \leq 2 : \log^* x := 1 \quad \forall x > 2 : \log^* x := 1 + \log^*(\log x)$$

Remarks:

- Log-star is an amazingly slowly growing function. Log-star of all the atoms in the observable universe (estimated to be 10^{80}) is 5! There are functions which grow even more slowly, such as the inverse Ackermann function, however, the inverse Ackermann function of all the atoms is 4. So log-star increases indeed very slowly!

Here is the idea of the algorithm: We start with color labels that have $\log n$ bits. In each synchronous round we compute a new label with exponentially smaller size than the previous label, still guaranteeing to have a valid vertex coloring! But how are we going to do that?

Algorithm 5 “6-Color”

```

1: Assume that initially the vertices are legally colored. Using Assumption
   3.11 each label only has  $\log n$  bits
2: The root assigns itself the label 0.
3: Each other node  $v$  executes the following code (synchronously in parallel)
4: send  $c_v$  to all children
5: repeat
6:   receive  $c_p$  from parent
7:   interpret  $c_v$  and  $c_p$  as little-endian bit-strings:  $c(k), \dots, c(1), c(0)$ 
8:   let  $i$  be the smallest index where  $c_v$  and  $c_p$  differ
9:   the new label is  $i$  (as bitstring) followed by the bit  $c_v(i)$  itself
10:  send  $c_v$  to all children
11: until  $c_w \in \{0, \dots, 5\}$  for all nodes  $w$ 

```

Example:

Algorithm 5 executed on the following part of a tree:

Grand-parent	0010110000	→	10010	→	...
Parent	1010010000	→	01010	→	111
Child	0110010000	→	10001	→	001

Theorem 2.10 (Analysis of Algorithm 5). *Algorithm 5 terminates in $\log^* n$ time.*

Proof: A detailed proof is, e.g., in [Peleg 7.3]. In class we do a sketch of the proof.

Remarks:

- Colors 11* (in binary notation, i.e., 6 or 7 in decimal notation) will not be chosen, because the node will then do another round. This gives a total of 6 colors (i.e., colors 0, ..., 5).
- Can one reduce the number of colors in only constant steps? Note that algorithm 3 does not work (since the degree of a node can be much higher than 6)! For fewer colors we need to have siblings monochromatic!
- Before we explore this problem we should probably have a second look at the end game of the algorithm, the UNTIL statement. Is this algorithm truly local?! Let's discuss!

Algorithm 6 Shift Down

-
- 1: Root chooses a new (different) color from $\{0, 1, 2\}$
 - 2: **Each** other **node** v concurrently executes the following code:
 - 3: Recolor v with the color of parent
-

Lemma 2.11 (Analysis of Algorithm 6). *Algorithm 6 preserves coloring legality; also siblings are monochromatic.*

Now Algorithm 3 (Reduce) can be used to reduce the number of used colors from six to three.

Algorithm 7 Six-2-Three

-
- 1: **Each node** v concurrently executes the following code:
 - 2: Run Algorithm 5 for $\log^* n$ rounds.
 - 3: **for** $x = 5, 4, 3$ **do**
 - 4: Perform subroutine Shift down (Algorithm 6)
 - 5: **if** $c_v = x$ **then**
 - 6: choose new color $c_v \in \{0, 1, 2\}$ using subroutine **First Free** (Algorithm 2)
 - 7: **end if**
 - 8: **end for**
-

Theorem 2.12 (Analysis of Algorithm 7). *Algorithm 7 colors a tree with three colors in time $O(\log^* n)$.*

Remarks:

- The term $O()$ used in Theorem 2.10 is called “big O” and is often used in distributed computing. Roughly speaking, $O(f)$ means “in the order of f , ignoring constant factors and smaller additive terms.” More formally, for two functions f and g , it holds that $f \in O(g)$ if there are constants x_0 and c so that $|f(x)| \leq c|g(x)|$ for all $x \geq x_0$. For an elaborate discussion on the big O notation we refer to other introductory math or computer science classes.
- As one can easily prove, a fast tree-coloring with only 2 colors is more than exponentially more expensive than coloring with 3 colors. In a tree degenerated to a list, nodes far away need to figure out whether they are an even or odd number of hops away from each other in order to get a 2-coloring. To do that one has to send a message to these nodes. This costs time linear in the number of nodes.
- Also other lower bounds have been proved, e.g., any algorithm for 2-coloring the d -regular tree of radius r which runs in time at most $2r/3$ requires at least $\Omega(\sqrt{d})$ colors.
- The idea of this algorithm can be generalized, e.g., to a ring topology. Also a general graph with constant degree Δ can be colored with $\Delta + 1$ colors in $O(\log^* n)$ time. The idea is as follows: In each step, a node compares its label to each of its neighbors, constructing a logarithmic difference-tag

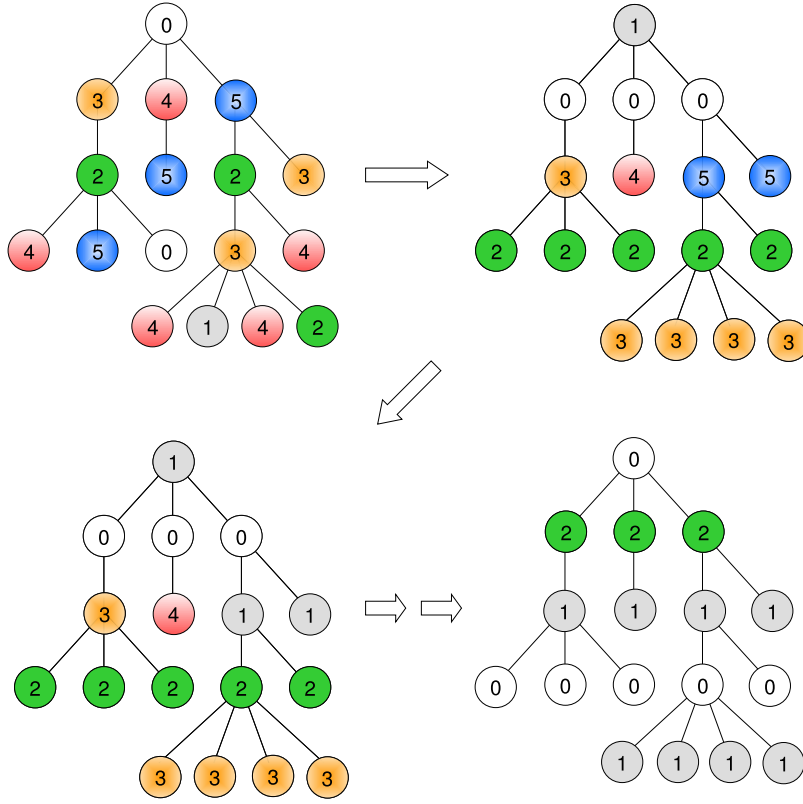


Figure 2.3: Possible execution of Algorithm 7.

as in 6-color (Algorithm 5). Then the new label is the concatenation of all the difference-tags. For constant degree Δ , this gives a 3Δ -label in $O(\log^* n)$ steps. Algorithm 3 then reduces the number of colors to $\Delta + 1$ in $2^{3\Delta}$ (this is still a constant for constant Δ) steps.

- Recently, researchers have proposed other methods to break down long ID's for log-star algorithms. With these new techniques, one is able to solve other problems, e.g., a maximal independent set in bounded growth graphs in $O(\log^* n)$ time. These techniques go beyond the scope of this course.
- Unfortunately, coloring a general graph is not yet possible with this technique. We will see another technique for that in Chapter 8. With this technique it is possible to color a general graph with $\Delta + 1$ colors in $O(\log n)$ time.
- A lower bound by Linial shows that many of these log-star algorithms are asymptotically (up to constant factors) optimal. This lower bound uses an interesting technique. However, because of the one-topic-per-class policy we cannot look at it today.

Chapter 3

Leader Election

3.1 Distributed Algorithms and Complexity

In the second part of this course we will often model the distributed system as a network or graph, and study protocols in which nodes (i.e., the processors) can only communicate with their neighbors to perform certain tasks. We are often interested in the following synchronous model or algorithm.

Definition 3.1 (Synchronous Distributed Algorithm). *In a synchronous algorithm, nodes operate in synchronous rounds. In each round, each processor executes the following steps:*

1. *Do some local computation (of reasonable “local complexity”).*
2. *Send messages to neighbors in graph (of reasonable size).*
3. *Receive messages (that were sent by neighbors in step 2 of the same round).*

Remarks:

- Any other step ordering is fine.

The other cornerstone model is the asynchronous algorithm.

Definition 3.2 (Asynchronous Distributed Algorithm). *In the asynchronous model, algorithms are event driven (“upon receiving message ..., do ...”). Processors cannot access a global clock. A message sent from one processor to another will arrive in finite but unbounded time.*

Remarks:

- The asynchronous model and the synchronous model (Definition 3.1) are the cornerstone models in distributed computing. As they do not necessarily reflect reality there are several models in between synchronous and asynchronous. However, from a theoretical point of view the synchronous and the asynchronous model are the most interesting ones (because every other model is in between these extremes).
- Note that in the asynchronous model, messages that take a longer path may arrive earlier.

In order to evaluate an algorithm, apart from the local complexity mentioned above, we consider the following metrics.

Definition 3.3 (Time Complexity). *For synchronous algorithms (as defined in 3.1) the time complexity is the number of rounds until the algorithm terminates.*

Definition 3.4 (Time Complexity). *For asynchronous algorithms (as defined in 3.1) the time complexity is the number of time units from the start of the execution to its completion in the worst case (every legal input, every execution scenario), assuming that each message has a delay of at most one time unit.*

Remarks:

- You cannot use the maximum delay in the algorithm design. In other words, the algorithm has to be correct even if there is no such delay upper bound.

Definition 3.5 (Message Complexity). *The message complexity of a synchronous and asynchronous algorithm is determined by the number of messages exchanged (again every legal input, every execution scenario).*

3.2 Anonymous Leader Election

Some algorithms (e.g., for medium access) ask for a special node, a so-called “leader”. Computing a leader is a most simple form of symmetry breaking. Algorithms based on leaders do generally not exhibit a high degree of parallelism, and therefore often suffer from poor (parallel) time complexity. However, sometimes it is still useful to have a leader to make critical decisions in an easy (though non-distributed!) way.

The process of choosing a leader is known as *leader election*. Although leader election is a simple form of symmetry breaking, there are some remarkable issues that allow us to introduce notable computational models.

In this chapter we concentrate on the ring topology. The ring is the “*dro-sophila*” of distributed computing as many interesting challenges already reveal the root of the problem in the special case of the ring. Paying special attention to the ring also makes sense from a practical point of view as some real world systems are based on a ring topology, e.g., the token ring standard for local area networks.

Problem 3.6 (Leader Election). *Each node eventually decides whether it is a leader or not, subject to the constraint that there is exactly one leader.*

Remarks:

- More formally, nodes are in one of three states: *undecided*, *leader*, *not leader*. Initially every node is in the undecided state. When leaving the undecided state, a node goes into a *final state* (leader or not leader).

Definition 3.7 (Anonymous). *A system is anonymous if nodes do not have unique identifiers.*

Definition 3.8 (Uniform). *An algorithm is called uniform if the number of nodes n is not known to the algorithm (to the nodes, if you wish). If n is known, the algorithm is called non-uniform.*

Whether a leader can be elected in an anonymous system depends on whether the network is symmetric (ring, complete graph, complete bipartite graph, etc.) or asymmetric (star, single node with highest degree, etc.). Simplifying slightly, in this context a symmetric graph is a graph in which the extended neighborhood of each node has the same structure. We will now show that non-uniform anonymous leader election for synchronous rings is impossible. The idea is that in a ring, symmetry can always be maintained.

Lemma 3.9. *After round k of any deterministic algorithm on an anonymous ring, each node is in the same state s_k .*

Proof by induction: All nodes start in the same state. A round in a synchronous algorithm consists of the three steps sending, receiving, local computation (see Definition 3.1). All nodes send the same message(s), receive the same message(s), do the same local computation, and therefore end up in the same state.

Theorem 3.10 (Anonymous Leader Election). *Deterministic leader election in an anonymous ring is impossible.*

Proof (with Lemma 3.9): If one node ever decides to become a leader (or a non-leader), then every other node does so as well, contradicting the problem specification 3.6 for $n > 1$. This holds for non-uniform algorithms, and therefore also for uniform algorithms. Furthermore, it holds for synchronous algorithms, and therefore also for asynchronous algorithms.

Remarks:

- Sense of direction is the ability of nodes to distinguish neighbor nodes in an anonymous setting. In a ring, for example, a node can distinguish the clockwise and the counterclockwise neighbor. Sense of direction does not help in anonymous leader election.
- Theorem 3.10 also holds for other symmetric network topologies (e.g., complete graphs, complete bipartite graphs, ...).
- Note that Theorem 3.10 does not hold for randomized algorithms; if nodes are allowed to toss a coin, symmetries can be broken.

3.3 Asynchronous Ring

We first concentrate on the asynchronous model from Definition 3.2. Throughout this section we assume non-anonymity; each node has a unique identifier as proposed in Assumption 3.11:

Assumption 3.11 (Node Identifiers). *Each node has a unique identifier, e.g., its IP address. We usually assume that each identifier consists of only $\log n$ bits if the system has n nodes.*

Having ID's seems to lead to a trivial leader election algorithm, as we can simply elect the node with, e.g., the highest ID.

Algorithm 8 Clockwise

-
- 1: **Each node** v executes the following code:
 - 2: v sends a message with its identifier (for simplicity also v) to its clockwise neighbor. {If node v already received a message w with $w > v$, then node v can skip this step; if node v receives its first message w with $w < v$, then node v will immediately send v .}
 - 3: **if** v receives a message w with $w > v$ **then**
 - 4: v forwards w to its clockwise neighbor
 - 5: v decides not to be the leader, if it has not done so already.
 - 6: **else if** v receives its own identifier v **then**
 - 7: v decides to be the leader
 - 8: **end if**
-

Theorem 3.12 (Analysis of Algorithm 8). *Algorithm 8 is correct. The time complexity is $O(n)$. The message complexity is $O(n^2)$.*

Proof: Let node z be the node with the maximum identifier. Node z sends its identifier in clockwise direction, and since no other node can swallow it, eventually a message will arrive at z containing it. Then z declares itself to be the leader. Every other node will declare non-leader at the latest when forwarding message z . Since there are n identifiers in the system, each node will at most forward n messages, giving a message complexity of at most n^2 . We start measuring the time when the first node that “wakes up” sends its identifier. For asynchronous time complexity (Definition 3.4) we assume that each message takes at most one time unit to arrive at its destination. After at most $n - 1$ time units the message therefore arrives at node z , waking z up. Routing the message z around the ring takes at most n time units. Therefore node z decides no later than at time $2n - 1$. Every other node decides before node z .

Remarks:

- Note that in Algorithm 8 nodes need to distinguish between clockwise and counterclockwise neighbors. In fact they do not: It is okay to simply send your own identifier to any neighbor, and forward a message m to the neighbor you did not receive the message m from. So nodes only need to be able to distinguish their two neighbors.
- Can we improve this algorithm?

Theorem 3.13 (Analysis of Algorithm 9). *Algorithm 9 is correct. The time complexity is $O(n)$. The message complexity is $O(n \log n)$.*

Proof: Correctness is as in Theorem 3.12. The time complexity is $O(n)$ since the node with maximum identifier z sends messages with round-trip times $2, 4, 8, 16, \dots, 2 \cdot 2^k$ with $k \leq \log(n + 1)$. (Even if we include the additional wake-up overhead, the time complexity stays linear.) Proving the message complexity is slightly harder: if a node v manages to survive round r , no other node in distance 2^r (or less) survives round r . That is, node v is the only node in its 2^r -neighborhood that remains active in round $r + 1$. Since this is the same for every node, less than $n/2^r$ nodes are active in round $r + 1$. Being active in round

Algorithm 9 Radius Growth (For readability we provide pseudo-code only; for a formal version please consult [Attiya/Welch Alg. 3.1])

- 1: **Each node** v does the following;
 - 2: Initially all nodes are *active*. {all nodes may still become leaders}
 - 3: Whenever a node v sees a message w with $w > v$, then v decides to not be a leader and becomes *passive*.
 - 4: Active nodes search in an exponentially growing neighborhood (clockwise and counterclockwise) for nodes with higher identifiers, by sending out *probe* messages. A probe message includes the ID of the original sender, a bit whether the sender can still become a leader, and a time-to-live number (*TTL*). The first probe message sent by node v includes a TTL of 1.
 - 5: Nodes (active or passive) receiving a probe message decrement the TTL and forward the message to the next neighbor; if their ID is larger than the one in the message, they set the leader bit to zero, as the probing node does not have the maximum ID. If the TTL is zero, probe messages are returned to the sender using a *reply* message. The reply message contains the ID of the receiver (the original sender of the probe message) and the leader-bit. Reply messages are forwarded by all nodes until they reach the receiver.
 - 6: Upon receiving the reply message: If there was no node with higher ID in the search area (indicated by the bit in the reply message), the TTL is doubled and two new probe messages are sent (again to the two neighbors). If there was a better candidate in the search area, then the node becomes passive.
 - 7: If a node v receives its own probe message (not a reply) v decides to be the leader.
-

r costs $2 \cdot 2 \cdot 2^r$ messages. Therefore, round r costs at most $2 \cdot 2 \cdot 2^r \cdot \frac{n}{2^{r-1}} = 8n$ messages. Since there are only logarithmic many possible rounds, the message complexity follows immediately.

Remarks:

- This algorithm is asynchronous and uniform as well.
- The question may arise whether one can design an algorithm with an even lower message complexity. We answer this question in the next section.

3.4 Lower Bounds

Lower bounds in distributed computing are often easier than in the standard centralized (random access machine, RAM) model because one can argue about messages that need to be exchanged. In this section we present a first lower bound. We show that Algorithm 9 is asymptotically optimal.

Definition 3.14 (Execution). *An execution of a distributed algorithm is a list of events, sorted by time. An event is a record (time, node, type, message), where type is “send” or “receive”.*

Remarks:

- We assume throughout this course that no two events happen at exactly the same time (or one can break ties arbitrarily).
- An execution of an asynchronous algorithm is generally not only determined by the algorithm but also by a “god-like” scheduler. If more than one message is in transit, the scheduler can choose which one arrives first.
- If two messages are transmitted over the same directed edge, then it is sometimes required that the message first transmitted will also be received first (“FIFO”).

For our lower bound, we assume the following model:

- We are given an asynchronous ring, where nodes may wake up at arbitrary times (but at the latest when receiving the first message).
- We only accept uniform algorithms where the node with the maximum identifier can be the leader. Additionally, every node that is not the leader must know the identity of the leader. These two requirements can be dropped when using a more complicated proof; however, this is beyond the scope of this course.
- During the proof we will “play god” and specify which message in transmission arrives next in the execution. We respect the FIFO conditions for links.

Definition 3.15 (Open Schedule). *A schedule is an execution chosen by the scheduler. A schedule for a ring is open if there is an open edge in the ring. An open (undirected) edge is an edge where no message traversing the edge has been received so far.*

The proof of the lower bound is by induction. First we show the base case:

Lemma 3.16. *Given a ring R with two nodes, we can construct an open schedule in which at least one message is received. The nodes cannot distinguish this schedule from one on a larger ring with all other nodes being where the open edge is.*

Proof: Let the two nodes be u and v with $u < v$. Node u must learn the identity of node v , thus receive at least one message. We stop the execution of the algorithm as soon as the first message is received. (If the first message is received by v , bad luck for the algorithm!) Then the other edge in the ring (on which the received message was not transmitted) is open. Since the algorithm needs to be uniform, maybe the open edge is not really an edge at all, nobody can tell. We could use this to glue two rings together, by breaking up this imaginary open edge and connect two rings by two edges.

Lemma 3.17. *By gluing together two rings of size $n/2$ for which we have open schedules, we can construct an open schedule on a ring of size n . If $M(n/2)$ denotes the number of messages already received in each of these schedules, at least $2M(n/2) + n/4$ messages have to be exchanged in order to solve leader election.*

Proof by induction: We divide the ring into two sub-rings R_1 and R_2 of size $n/2$. These subrings cannot be distinguished from rings with $n/2$ nodes if no messages are received from “outsiders”. We can ensure this by not scheduling such messages until we want to. Note that executing both given open schedules on R_1 and R_2 “in parallel” is possible because we control not only the scheduling of the messages, but also when nodes wake up. By doing so, we make sure that $2M(n/2)$ messages are sent before the nodes in R_1 and R_2 learn anything of each other!

Without loss of generality, R_1 contains the maximum identifier. Hence, each node in R_2 must learn the identity of the maximum identifier, thus at least $n/2$ additional messages must be received. The only problem is that we cannot connect the two sub-rings with both edges since the new ring needs to remain open. Thus, only messages over one of the edges can be received. We “play god” and look into the future: we check what happens when we close only one of these connecting edges. With the argument that $n/2$ new messages must be received, we know that there is at least one edge that will produce at least $n/4$ additional messages when being scheduled. (These messages may not be sent over the closed link, but they are *caused* by a message over this link. They cannot involve any message along the other (open) edge at distance $n/2$.) We schedule this edge and the resulting $n/4$ messages, and leave the other open.

Lemma 3.18. *Any uniform leader election algorithm for asynchronous rings has at least message complexity $M(n) \geq \frac{n}{4}(\log n + 1)$.*

Proof by induction: For simplicity we assume n being a power of 2. The base case $n = 2$ works because of Lemma 3.16 which implies that $M(2) \geq 1 = \frac{2}{4}(\log 2 + 1)$. For the induction step, using Lemma 3.17 and the induction hypothesis we have

$$\begin{aligned} M(n) &= 2 \cdot M\left(\frac{n}{2}\right) + \frac{n}{4} \\ &\geq 2 \cdot \left(\frac{n}{8} \left(\log \frac{n}{2} + 1\right)\right) + \frac{n}{4} \\ &= \frac{n}{4} \log n + \frac{n}{4} = \frac{n}{4} (\log n + 1). \end{aligned}$$

□

Remarks:

- To hide the ugly constants we use the “big Omega” notation, the lower bound equivalent of $O()$. A function f is in $\Omega(g)$ if there are constants x_0 and $c > 0$ such that $|f(x)| \geq c|g(x)|$ for all $x \geq x_0$. Again we refer to standard text books for a formal definition. Rewriting Lemma 3.18 we get:

Theorem 3.19 (Asynchronous Leader Election Lower Bound). *Any uniform leader election algorithm for asynchronous rings has $\Omega(n \log n)$ message complexity.*

3.5 Synchronous Ring

The lower bound relied on delaying messages for a very long time. Since this is impossible in the synchronous model, we might get a better message complexity

in this case. The basic idea is very simple: In the synchronous model, *not* receiving a message is information as well! First we make some additional assumptions:

- We assume that the algorithm is non-uniform (i.e., the ring size n is known).
- We assume that every node starts at the same time.
- The node with the minimum identifier becomes the leader; identifiers are integers.

Algorithm 10 Synchronous Leader Election

- 1: **Each node** v concurrently executes the following code:
 - 2: The algorithm operates in synchronous phases. Each phase consists of n time steps. Node v counts phases, starting with 0.
 - 3: **if** phase = v **and** v did not yet receive a message **then**
 - 4: v decides to be the leader
 - 5: v sends the message “ v is leader” around the ring
 - 6: **end if**
-

Remarks:

- Message complexity is indeed n .
- But the time complexity is huge! If m is the minimum identifier it is $m \cdot n$.
- The synchronous start and the non-uniformity assumptions can be dropped by using a wake-up technique (upon receiving a wake-up message, wake up your clockwise neighbors) and by letting messages travel slowly.
- There are several lower bounds for the synchronous model: comparison-based algorithms or algorithms where the time complexity cannot be a function of the identifiers have message complexity $\Omega(n \log n)$ as well.
- In general graphs efficient leader election may be tricky. While time-optimal leader election can be done by parallel flooding-echo (see next chapter), bounding the message complexity is generally more difficult.

Chapter 4

Tree Algorithms

In this chapter we learn a few basic algorithms on trees, and how to construct trees in the first place so that we can run these (and other) algorithms. The good news is that these algorithms have many applications, the bad news is that this chapter is a bit on the simple side. But maybe that's not really bad news?!

4.1 Broadcast

Definition 4.1 (Broadcast). *A broadcast operation is initiated by a single processor, the source. The source wants to send a message to all other nodes in the system.*

Definition 4.2 (Distance, Radius, Diameter). *The distance between two nodes u and v in an undirected graph G is the number of hops of a minimum path between u and v . The radius of a node u is the maximum distance between u and any other node in the graph. The radius of a graph is the minimum radius of any node in the graph. The diameter of a graph is the maximum distance between two arbitrary nodes.*

Remarks:

- Clearly there is a close relation between the radius R and the diameter D of a graph, such as $R \leq D \leq 2R$.
- The world is often fascinated by graphs with a small radius. For example, movie fanatics study the who-acted-with-whom-in-the-same-movie graph. For this graph it has long been believed that the actor Kevin Bacon has a particularly small radius. The number of hops from Bacon even got a name, the Bacon Number. In the meantime, however, it has been shown that there are “better” centers in the Hollywood universe, such as Sean Connery, Christopher Lee, Rod Steiger, Gene Hackman, or Michael Caine. The center of other social networks has also been explored, Paul Erdős for instance is well known in the math community.

Theorem 4.3 (Broadcast Lower Bound). *The message complexity of broadcast is at least $n - 1$. The source's radius is a lower bound for the time complexity.*

Proof: Every node must receive the message.

Remarks:

- You can use a pre-computed spanning tree to do broadcast with tight message complexity. If the spanning tree is a breadth-first search spanning tree (for a given source), then the time complexity is tight as well.

Definition 4.4 (Clean). *A graph (network) is clean if the nodes do not know the topology of the graph.*

Theorem 4.5 (Clean Broadcast Lower Bound). *For a clean network, the number of edges is a lower bound for the broadcast message complexity.*

Proof: If you do not try every edge, you might miss a whole part of the graph behind it.

Remarks:

- This lower bound proof directly brings us to the well known *flooding* algorithm.

Algorithm 11 Flooding

- 1: The source (root) sends the message to all neighbors.
 - 2: **Each other node** v upon receiving the message the first time forwards the message to all (other) neighbors.
 - 3: Upon later receiving the message again (over other edges), a node can discard the message.
-

Remarks:

- If node v receives the message first from node u , then node v calls node u *parent*. This parent relation defines a spanning tree T . If the flooding algorithm is executed in a synchronous system, then T is a breadth-first search spanning tree (with respect to the root).
- More interestingly, also in asynchronous systems the flooding algorithm terminates after R time units, R being the radius of the source. However, the constructed spanning tree may not be a breadth-first search spanning tree.

4.2 Convergecast

Convergecast is the same as broadcast, just reversed: Instead of a root sending a message to all other nodes, all other nodes send information to a root. The simplest convergecast algorithm is the echo algorithm:

Algorithm 12 Echo

Require: This algorithm is initiated at the leaves.

- 1: A leaf sends a message to its parent.
 - 2: If an inner node has received a message from each child, it sends a message to the parent.
-

Remarks:

- Usually the echo algorithm is paired with the flooding algorithm, which is used to let the leaves know that they should start the echo process; this is known as flooding/echo.
- One can use convergecast for termination detection, for example. If a root wants to know whether all nodes in the system have finished some task, it initiates a flooding/echo; the message in the echo algorithm then means “This subtree has finished the task.”
- Message complexity of the echo algorithm is $n - 1$, but together with flooding it is $O(m)$, where $m = |E|$ is the number of edges in the graph.
- The time complexity of the echo algorithm is determined by the depth of the spanning tree (i.e., the radius of the root within the tree) generated by the flooding algorithm.
- The flooding/echo algorithm can do much more than collecting acknowledgements from subtrees. One can for instance use it to compute the number of nodes in the system, or the maximum ID (for leader election), or the sum of all values stored in the system, or a route-disjoint matching.
- Moreover, by combining results one can compute even fancier aggregations, e.g., with the number of nodes and the sum one can compute the average. With the average one can compute the standard deviation. And so on ...

4.3 BFS Tree Construction

In synchronous systems the flooding algorithm is a simple yet efficient method to construct a breadth-first search (BFS) spanning tree. However, in asynchronous systems the spanning tree constructed by the flooding algorithm may be far from BFS. In this section, we implement two classic BFS constructions—Dijkstra and Bellman-Ford—as asynchronous algorithms.

We start with the Dijkstra algorithm. The basic idea is to always add the “closest” node to the existing part of the BFS tree. We need to parallelize this idea by developing the BFS tree layer by layer:

Theorem 4.6 (Analysis of Algorithm 13). *The time complexity of Algorithm 13 is $O(D^2)$, the message complexity is $O(m + nD)$, where D is the diameter of the graph, n the number of nodes, and m the number of edges.*

Proof: A broadcast/echo algorithm in T_p needs at most time $2D$. Finding new neighbors at the leaves costs 2 time units. Since the BFS tree height is bounded

Algorithm 13 Dijkstra BFS

-
- 1: The algorithm proceeds in phases. In phase p the nodes with distance p to the root are detected. Let T_p be the tree in phase p . We start with T_1 which is the root plus all direct neighbors of the root. We start with phase $p = 1$:
 - 2: **repeat**
 - 3: The root starts phase p by broadcasting “start p ” within T_p .
 - 4: When receiving “start p ” a leaf node u of T_p (that is, a node that was newly discovered in the last phase) sends a “join $p + 1$ ” message to all quiet neighbors. (A neighbor v is quiet if u has not yet “talked” to v .)
 - 5: A node v receiving the first “join $p+1$ ” message replies with “ACK” and becomes a leaf of the tree T_{p+1} .
 - 6: A node v receiving any further “join” message replies with “NACK”.
 - 7: The leaves of T_p collect all the answers of their neighbors; then the leaves start an echo algorithm back to the root.
 - 8: When the echo process terminates at the root, the root increments the phase
 - 9: **until** there was no new node detected
-

by the diameter, we have D phases, giving a total time complexity of $O(D^2)$. Each node participating in broadcast/echo only receives (broadcasts) at most 1 message and sends (echoes) at most once. Since there are D phases, the cost is bounded by $O(nD)$. On each edge there are at most 2 “join” messages. Replies to a “join” request are answered by 1 “ACK” or “NACK”, which means that we have at most 4 additional messages per edge. Therefore the message complexity is $O(m + nD)$.

Remarks:

- The time complexity is not very exciting, so let’s try Bellman-Ford!

The basic idea of Bellman-Ford is even simpler, and heavily used in the Internet, as it is a basic version of the omnipresent border gateway protocol (BGP). The idea is to simply keep the distance to the root accurate. If a neighbor has found a better route to the root, a node might also need to update its distance.

Algorithm 14 Bellman-Ford BFS

-
- 1: Each node u stores an integer d_u which corresponds to the distance from u to the root. Initially $d_{\text{root}} = 0$, and $d_u = \infty$ for every other node u .
 - 2: The root starts the algorithm by sending “1” to all neighbors.
 - 3: **if** a node u receives a message “ y ” with $y < d_u$ from a neighbor v **then**
 - 4: node u sets $d_u := y$
 - 5: node u sends “ $y + 1$ ” to all neighbors (except v)
 - 6: **end if**
-

Theorem 4.7 (Analysis of Algorithm 14). *The time complexity of Algorithm 14 is $O(D)$, the message complexity is $O(nm)$, where D, n, m are defined as in Theorem 4.6.*

Proof: We can prove the time complexity by induction. We claim that a node at distance d from the root has received a message “ d ” by time d . The root

knows by time 0 that it is the root. A node v at distance d has a neighbor u at distance $d - 1$. Node u by induction sends a message “ d ” to v at time $d - 1$ or before, which is then received by v at time d or before. Message complexity is easier: A node can reduce its distance at most $n - 1$ times; each of these times it sends a message to all its neighbors. If all nodes do this we have $O(nm)$ messages.

Remarks:

- Algorithm 13 has the better message complexity and Algorithm 14 has the better time complexity. The currently best algorithm (optimizing both) needs $O(m + n \log^3 n)$ messages and $O(D \log^3 n)$ time. This “trade-off” algorithm is beyond the scope of this course.

4.4 MST Construction

There are several types of spanning trees, each serving a different purpose. A particularly interesting spanning tree is the minimum spanning tree (MST). The MST only makes sense on weighted graphs, hence in this section we assume that each edge e is assigned a weight ω_e .

Definition 4.8 (MST). *Given a weighted graph $G = (V, E, \omega)$, the MST of G is a spanning tree T minimizing $\omega(T)$, where $\omega(G') = \sum_{e \in G'} \omega_e$ for any subgraph $G' \subseteq G$.*

Remarks:

- In the following we assume that no two edges of the graph have the same weight. This simplifies the problem as it makes the MST unique; however, this simplification is not essential as one can always break ties by adding the IDs of adjacent vertices to the weight.
- Obviously we are interested in computing the MST in a distributed way. For this we use a well-known lemma:

Definition 4.9 (Blue Edges). *Let T be a spanning tree of the weighted graph G and $T' \subseteq T$ a subgraph of T (also called a fragment). Edge $e = (u, v)$ is an outgoing edge of T' if $u \in T'$ and $v \notin T'$ (or vice versa). The minimum weight outgoing edge $b(T')$ is the so-called blue edge of T' .*

Lemma 4.10. *For a given weighted graph G (such that no two weights are the same), let T denote the MST, and T' be a fragment of T . Then the blue edge of T' is also part of T , i.e., $T' \cup b(T') \subseteq T$.*

Proof: For the sake of contradiction, suppose that in the MST T there is edge $e \neq b(T')$ connecting T' with the remainder of T . Adding the blue edge $b(T')$ to the MST T we get a cycle including both e and $b(T')$. If we remove e from this cycle we still have a spanning tree, and since by the definition of the blue edge $\omega_e > \omega_{b(T')}$, the weight of that new spanning tree is less than the weight of T . We have a contradiction.

Remarks:

- In other words, the blue edges seem to be the key to a distributed algorithm for the MST problem. Since every node itself is a fragment of the MST, every node directly has a blue edge! All we need to do is to grow these fragments! Essentially this is a distributed version of Kruskal's sequential algorithm.
- At any given time the nodes of the graph are partitioned into fragments (rooted subtrees of the MST). Each fragment has a root, the ID of the fragment is the ID of its root. Each node knows its parent and its children in the fragment. The algorithm operates in phases. At the beginning of a phase, nodes know the IDs of the fragments of their neighbor nodes.

Algorithm 15 GHS (Gallager–Humblet–Spira)

```

1: Initially each node is the root of its own fragment. We proceed in phases:
2: repeat
3:   All nodes learn the fragment IDs of their neighbors.
4:   The root of each fragment uses flooding/echo in its fragment to determine
     the blue edge  $b = (u, v)$  of the fragment.
5:   The root sends a message to node  $u$ ; while forwarding the message on the
     path from the root to node  $u$  all parent-child relations are inverted {such
     that  $u$  is the new temporary root of the fragment}
6:   node  $u$  sends a merge request over the blue edge  $b = (u, v)$ .
7:   if node  $v$  also sent a merge request over the same blue edge  $b = (v, u)$ 
     then
8:     either  $u$  or  $v$  (whichever has the smaller ID) is the new fragment root
9:     the blue edge  $b$  is directed accordingly
10:  else
11:    node  $v$  is the new parent of node  $u$ 
12:  end if
13:  the newly elected root node informs all nodes in its fragment (again using
     flooding/echo) about its identity
14: until all nodes are in the same fragment (i.e., there is no outgoing edge)

```

Remarks:

- Algorithm 15 was stated in pseudo-code, with a few details not really explained. For instance, it may be that some fragments are much larger than others, and because of that some nodes may need to wait for others, e.g., if node u needs to find out whether neighbor v also wants to merge over the blue edge $b = (u, v)$. The good news is that all these details can be solved. We can for instance bound the asynchronicity by guaranteeing that nodes only start the new phase after the last phase is done, similarly to the phase-technique of Algorithm 13.

Theorem 4.11 (Analysis of Algorithm 15). *The time complexity of Algorithm 15 is $O(n \log n)$, the message complexity is $O(m \log n)$.*

Proof: Each phase mainly consists of two flooding/echo processes. In general, the cost of flooding/echo on a tree is $O(D)$ time and $O(n)$ messages. However,

the diameter D of the fragments may turn out to be not related to the diameter of the graph because the MST may meander, hence it really is $O(n)$ time. In addition, in the first step of each phase, nodes need to learn the fragment ID of their neighbors; this can be done in 2 steps but costs $O(m)$ messages. There are a few more steps, but they are cheap. Altogether a phase costs $O(n)$ time and $O(m)$ messages. So we only have to figure out the number of phases: Initially all fragments are single nodes and hence have size 1. In a later phase, each fragment merges with at least one other fragment, that is, the size of the smallest fragment at least doubles. In other words, we have at most $\log n$ phases. The theorem follows directly.

Remarks:

- Algorithm 15 is called “GHS” after Gallager, Humblet, and Spira, three pioneers in distributed computing. Despite being quite simple the algorithm won the prestigious Edsger W. Dijkstra Prize in Distributed Computing in 2004, among other reasons because it was one of the first (1983) non-trivial asynchronous distributed algorithms. As such it can be seen as one of the seeds of this research area.
- We presented a simplified version of GHS. The original paper by Gallager et al. featured an improved message complexity of $O(m + n \log n)$.
- In 1987, Awerbuch managed to further improve the GHS algorithm to get $O(n)$ time and $O(m + n \log n)$ message complexity, both asymptotically optimal.
- The GHS algorithm can be applied in different ways. GHS for instance directly solves leader election in general graphs: The leader is simply the last surviving root!

Chapter 5

Distributed Sorting

“Indeed, I believe that virtually *every* important aspect of programming arises somewhere in the context of sorting [and searching]!”

– Donald E. Knuth, The Art of Computer Programming

In this chapter we study a classic problem in computer science—sorting—from a distributed computing perspective. In contrast to an orthodox single-processor sorting algorithm, no node has access to all data, instead the to-be-sorted values are *distributed*. Distributed sorting then boils down to:

Definition 5.1 (Sorting). *We choose a graph with n nodes v_1, \dots, v_n . Initially each node stores a value. After applying a sorting algorithm, node v_k stores the k^{th} smallest value.*

Remarks:

- What if we route all values to the same central node v , let v sort the values locally, and then route them to the correct destinations?! According to the message passing model studied in the first few chapters this is perfectly legal. With a star topology sorting finishes in $\mathcal{O}(1)$ time!

Definition 5.2 (Node Contention). *In each step of a synchronous algorithm, each node can only send and receive $\mathcal{O}(1)$ messages containing $\mathcal{O}(1)$ values, no matter how many neighbors the node has.*

Remarks:

- Using Definition 5.2 sorting on a star graph takes linear time.

5.1 Array & Mesh

To get a better intuitive understanding of distributed sorting, we start with two simple topologies, the array and the mesh. Let us begin with the array:

Algorithm 16 Odd/Even Sort

-
- 1: Given an array of n nodes (v_1, \dots, v_n) , each storing a value (not sorted).
 - 2: **repeat**
 - 3: Compare and exchange the values at nodes i and $i + 1$, i odd
 - 4: Compare and exchange the values at nodes i and $i + 1$, i even
 - 5: **until** done
-

Remarks:

- The compare and exchange primitive in Algorithm 16 is defined as follows:
Let the value stored at node i be v_i . After the compare and exchange node i stores value $\min(v_i, v_{i+1})$ and node $i + 1$ stores value $\max(v_i, v_{i+1})$.
- How fast is the algorithm, and how can we prove correctness/efficiency?
- The most interesting proof uses the so-called 0-1 Sorting Lemma. It allows us to restrict our attention to an input of 0's and 1's only, and works for any "oblivious comparison-exchange" algorithm. (Oblivious means: Whether you exchange two values must only depend on the relative order of the two values, and not on anything else.)

Lemma 5.3 (0-1 Sorting Lemma). *If an oblivious comparison-exchange algorithm sorts all inputs of 0's and 1's, then it sorts arbitrary inputs.*

Proof. We prove the opposite direction (does not sort arbitrary inputs \Rightarrow does not sort 0's and 1's). Assume that there is an input $x = x_1, \dots, x_n$ that is not sorted correctly. Then there is a smallest value k such that the value at node v_k after running the sorting algorithm is strictly larger than the k^{th} smallest value $x(k)$. Define an input $x_i^* = 0 \Leftrightarrow x_i \leq x(k)$, $x_i^* = 1$ else. Whenever the algorithm compares a pair of 1's or 0's, it is not important whether it exchanges the values or not, so we may simply assume that it does the same as on the input x . On the other hand, whenever the algorithm exchanges some values $x_i^* = 0$ and $x_j^* = 1$, this means that $x_i \leq x(k) < x_j$. Therefore, in this case the respective compare-exchange operation will do the same on both inputs. We conclude that the algorithm will order x^* the same way as x , i.e., the output with only 0's and 1's will also not be correct. \square

Theorem 5.4. *Algorithm 16 sorts correctly in n steps.*

Proof. Thanks to Lemma 5.3 we only need to consider an array with 0's and 1's. Let j_1 be the node with the rightmost (highest index) 1. If j_1 is odd (even) it will move in the first (second) step. In any case it will move right in every following step until it reaches the rightmost node v_n . Let j_k be the node with the k^{th} rightmost 1. We show by induction that j_k is not "blocked" anymore (constantly moves until it reaches destination!) after step k . We have already anchored the induction at $k = 1$. Since j_{k-1} moves after step $k - 1$, j_k gets a right 0-neighbor for each step after step k . (For matters of presentation we omitted a couple of simple details.) \square

Algorithm 17 Shearsort

-
- 1: We are given a mesh with m rows and m columns, m even, $n = m^2$.
 - 2: The sorting algorithm operates in phases, and uses the odd/even sort algorithm on rows or columns.
 - 3: **repeat**
 - 4: In the odd phases $1, 3, \dots$ we sort all the rows, in the even phases $2, 4, \dots$ we sort all the columns, such that:
 - 5: Columns are sorted such that the small values move up.
 - 6: Odd rows $(1, 3, \dots, m-1)$ are sorted such that small values move left.
 - 7: Even rows $(2, 4, \dots, m)$ are sorted such that small values move right.
 - 8: **until** done
-

Remarks:

- Linear time is not very exciting, maybe we can do better by using a different topology? Let's try a mesh (a.k.a. grid) topology first.

Theorem 5.5. *Algorithm 17 sorts n values in $\sqrt{n}(\log n + 1)$ time in snake-like order.*

Proof. Since the algorithm is oblivious, we can use Lemma 5.3. We show that after a row and a column phase, half of the previously unsorted rows will be sorted. More formally, let us call a row with only 0's (or only 1's) *clean*, a row with 0's and 1's is *dirty*. At any stage, the rows of the mesh can be divided into three regions. In the north we have a region of all-0 rows, in the south all-1 rows, in the middle a region of dirty rows. Initially all rows can be dirty. Since neither row nor column sort will touch already clean rows, we can concentrate on the dirty rows.

First we run an odd phase. Then, in the even phase, we run a peculiar column sorter: We group two consecutive dirty rows into pairs. Since odd and even rows are sorted in opposite directions, two consecutive dirty rows look as follows:

$$00000 \dots 11111$$

$$11111 \dots 00000$$

Such a pair can be in one of three states. Either we have more 0's than 1's, or more 1's than 0's, or an equal number of 0's and 1's. Column-sorting each pair will give us at least one clean row (and two clean rows if " $|0| = |1|$ "). Then move the cleaned rows north/south and we will be left with half the dirty rows.

At first glance it appears that we need such a peculiar column sorter. However, any column sorter sorts the columns in exactly the same way (we are very grateful to have Lemma 5.3!).

All in all we need $2 \log m = \log n$ phases to remain only with 1 dirty row in the middle which will be sorted (not cleaned) with the last row-sort. \square

Remarks:

- There are algorithms that sort in $3m + o(m)$ time on an m by m mesh (by dividing the mesh into smaller blocks). This is asymptotically optimal, since a value might need to move $2m$ times.
- Such a \sqrt{n} -sorter is cute, but we are more ambitious. There are non-distributed sorting algorithms such as quicksort, heapsort, or mergesort that sort n values in (expected) $\mathcal{O}(n \log n)$ time. Using our n -fold parallelism effectively we might therefore hope for a distributed sorting algorithm that sorts in time $\mathcal{O}(\log n)$!

5.2 Sorting Networks

In this section we construct a graph topology which is carefully manufactured for sorting. This is a deviation from previous chapters where we always had to work with the topology that was given to us. In many application areas (e.g. peer-to-peer networks, communication switches, systolic hardware) it is indeed possible (in fact, crucial!) that an engineer can build the topology best suited for her application.

Definition 5.6 (Sorting Networks). *A comparator is a device with two inputs x, y and two outputs x', y' such that $x' = \min(x, y)$ and $y' = \max(x, y)$. We construct so-called comparison networks that consist of wires that connect comparators (the output port of a comparator is sent to an input port of another comparator). Some wires are not connected to comparator outputs, and some are not connected to comparator inputs. The first are called input wires of the comparison network, the second output wires. Given n values on the input wires, a sorting network ensures that the values are sorted on the output wires. We will also use the term width to indicate the number of wires in the sorting network.*

Remarks:

- The odd/even sorter explained in Algorithm 16 can be described as a sorting network.
- Often we will draw all the wires on n horizontal lines (n being the “width” of the network). Comparators are then vertically connecting two of these lines.
- Note that a sorting network is an oblivious comparison-exchange network. Consequently we can apply Lemma 5.3 throughout this section. An example sorting network is depicted in Figure 5.1.

Definition 5.7 (Depth). *The depth of an input wire is 0. The depth of a comparator is the maximum depth of its input wires plus one. The depth of an output wire of a comparator is the depth of the comparator. The depth of a comparison network is the maximum depth (of an output wire).*

Definition 5.8 (Bitonic Sequence). *A bitonic sequence is a sequence of numbers that first monotonically increases, and then monotonically decreases, or vice versa.*

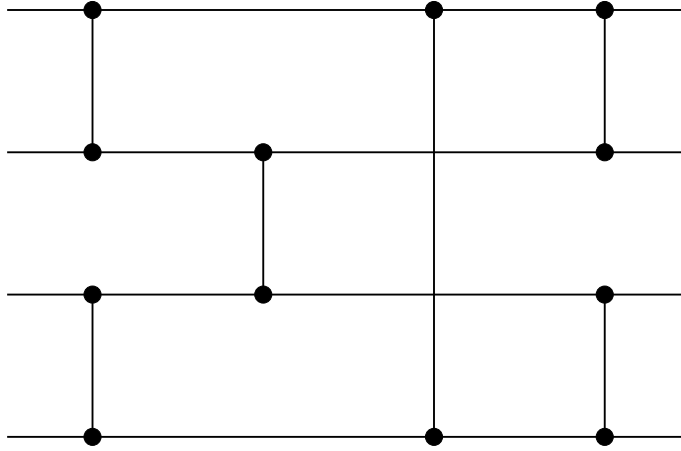


Figure 5.1: A sorting network.

Remarks:

- $\langle 1, 4, 6, 8, 3, 2 \rangle$ or $\langle 5, 3, 2, 1, 4, 8 \rangle$ are bitonic sequences.
- $\langle 9, 6, 2, 3, 5, 4 \rangle$ or $\langle 7, 4, 2, 5, 9, 8 \rangle$ are not bitonic.
- Since we restrict ourselves to 0's and 1's (Lemma 5.3), bitonic sequences have the form $0^i 1^j 0^k$ or $1^i 0^j 1^k$ for $i, j, k \geq 0$.

Algorithm 18 Half Cleaner

-
- 1: A half cleaner is a comparison network of depth 1, where we compare wire i with wire $i + n/2$ for $i = 1, \dots, n/2$ (we assume n to be even).
-

Lemma 5.9. *Feeding a bitonic sequence into a half cleaner (Algorithm 18), the half cleaner cleans (makes all-0 or all-1) either the upper or the lower half of the n wires. The other half is bitonic.*

Proof. Assume that the input is of the form $0^i 1^j 0^k$ for $i, j, k \geq 0$. If the midpoint falls into the 0's, the input is already clean/bitonic and will stay so. If the midpoint falls into the 1's the half cleaner acts as Shearsort with two adjacent rows, exactly as in the proof of Theorem 5.5. The case $1^i 0^j 1^k$ is symmetric. \square

Algorithm 19 Bitonic Sequence Sorter

-
- 1: A bitonic sequence sorter of width n (n being a power of 2) consists of a half cleaner of width n , and then two bitonic sequence sorters of width $n/2$ each.
 - 2: A bitonic sequence sorter of width 1 is empty.
-

Lemma 5.10. *A bitonic sequence sorter (Algorithm 19) of width n sorts bitonic sequences. It has depth $\log n$.*

Proof. The proof follows directly from the Algorithm 19 and Lemma 5.9. \square

Remarks:

- Clearly we want to sort arbitrary and not only bitonic sequences! To do this we need one more concept, merging networks.

Algorithm 20 Merging Network

- 1: A merging network of width n is a merger of width n followed by two bitonic sequence sorters of width $n/2$. A merger is a depth-one network where we compare wire i with wire $n - i + 1$, for $i = 1, \dots, n/2$.
-

Remarks:

- Note that a merging network is a bitonic sequence sorter where we replace the (first) half-cleaner by a merger.

Lemma 5.11. *A merging network of width n (Algorithm 20) merges two sorted input sequences of length $n/2$ each into one sorted sequence of length n .*

Proof. We have two sorted input sequences. Essentially, a merger does to two sorted sequences what a half cleaner does to a bitonic sequence, since the lower part of the input is reversed. In other words, we can use same argument as in Theorem 5.5 and Lemma 5.9: Again, after the merger step either the upper or the lower half is clean, the other is bitonic. The bitonic sequence sorters complete sorting. \square

Remarks:

- How do you sort n values when you are able to merge two sorted sequences of size $n/2$? Piece of cake, just apply the merger recursively.

Algorithm 21 Batcher's "Bitonic" Sorting Network

- 1: A batcher sorting network of width n consists of two batcher sorting networks of width $n/2$ followed by a merging network of width n . (See Figure 5.2.)
 - 2: A batcher sorting network of width 1 is empty.
-

Theorem 5.12. *A sorting network (Algorithm 21) sorts an arbitrary sequence of n values. It has depth $\mathcal{O}(\log^2 n)$.*

Proof. Correctness is immediate: at recursive stage k ($k = 1, 2, 3, \dots, \log n$) we merge 2^k sorted sequences into 2^{k-1} sorted sequences. The depth $d(n)$ of the sorting network of level n is the depth of a sorting network of level $n/2$ plus the depth $m(n)$ of a merging network with width n . The depth of a sorter of level 1 is 0 since the sorter is empty. Since a merging network of width n has the same depth as a bitonic sequence sorter of width n , we know by Lemma 5.10 that $m(n) = \log n$. This gives a recursive formula for $d(n)$ which solves to $d(n) = \frac{1}{2} \log^2 n + \frac{1}{2} \log n$. \square

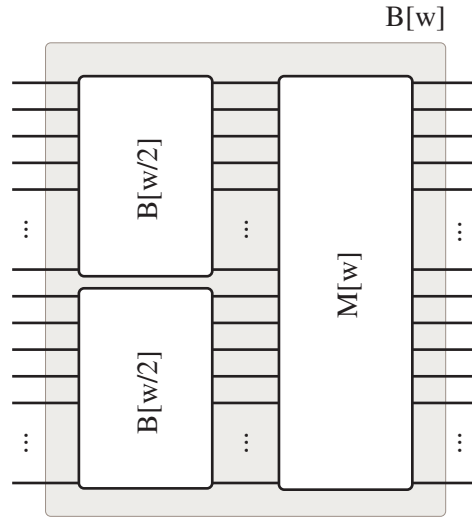


Figure 5.2: A batcher sorting network

Remarks:

- Simulating Batcher's sorting network on an ordinary sequential computer takes time $\mathcal{O}(n \log^2 n)$. As said, there are sequential sorting algorithms that sort in asymptotically optimal time $\mathcal{O}(n \log n)$. So a natural question is whether there is a sorting network with depth $\mathcal{O}(\log n)$. Such a network would have some remarkable advantages over sequential asymptotically optimal sorting algorithms such as heapsort. Apart from being highly parallel, it would be completely oblivious, and as such perfectly suited for a fast hardware solution. In 1983, Ajtai, Komlos, and Szemerédi presented a celebrated $\mathcal{O}(\log n)$ depth sorting network. (Unlike Batcher's sorting network the constant hidden in the big- \mathcal{O} of the "AKS" sorting network is too large to be practical, however.)
- It can be shown that Batcher's sorting network and similarly others can be simulated by a Butterfly network and other hypercubic networks, see next chapter.
- What if a sorting network is asynchronous?!? Clearly, using a synchronizer we can still sort, but it is also possible to use it for something else. Check out the next section!

5.3 Counting Networks

In this section we address distributed counting, a distributed service which can for instance be used for load balancing.

Definition 5.13 (Distributed Counting). *A distributed counter is a variable that is common to all processors in a system and that supports an atomic test-and-increment operation. The operation delivers the system's counter value to the requesting processor and increments it.*

Remarks:

- A naive distributed counter stores the system's counter value with a distinguished central node. When other nodes initiate the test-and-increment operation, they send a request message to the central node and in turn receive a reply message with the current counter value. However, with a large number of nodes operating on the distributed counter, the central processor will become a bottleneck. There will be a congestion of request messages at the central processor, in other words, the system will not scale.
- Is a scalable implementation (without any kind of bottleneck) of such a distributed counter possible, or is distributed counting a problem which is inherently centralized?!?
- Distributed counting could for instance be used to implement a load balancing infrastructure, i.e. by sending the job with counter value i (modulo n) to server i (out of n possible servers).

Definition 5.14 (Balancer). *A balancer is an asynchronous flip-flop which forwards messages that arrive on the left side to the wires on the right, the first to the upper, the second to the lower, the third to the upper, and so on.*

Algorithm 22 Bitonic Counting Network.

- 1: Take Batcher's bitonic sorting network of width w and replace all the comparators with balancers.
 - 2: When a node wants to count, it sends a message to an arbitrary input wire.
 - 3: The message is then routed through the network, following the rules of the asynchronous balancers.
 - 4: Each output wire is completed with a "mini-counter."
 - 5: The mini-counter of wire k replies the value " $k + i \cdot w$ " to the initiator of the i^{th} message it receives.
-

Definition 5.15 (Step Property). *A sequence y_0, y_1, \dots, y_{w-1} is said to have the step property, if $0 \leq y_i - y_j \leq 1$, for any $i < j$.*

Remarks:

- If the output wires have the step property, then with r requests, exactly the values $1, \dots, r$ will be assigned by the mini-counters. All we need to show is that the counting network has the step property. For that we need some additional facts...

Facts 5.16. *For a balancer, we denote the number of consumed messages on the i^{th} input wire with x_i , $i = 0, 1$. Similarly, we denote the number of sent messages on the i^{th} output wire with y_i , $i = 0, 1$. A balancer has these properties:*

- (1) *A balancer does not generate output-messages; that is, $x_0 + x_1 \geq y_0 + y_1$ in any state.*
- (2) *Every incoming message is eventually forwarded. In other words, if we are in a quiescent state (no message in transit), then $x_0 + x_1 = y_0 + y_1$.*

- (3) The number of messages sent to the upper output wire is at most one higher than the number of messages sent to the lower output wire: in any state $y_0 = \lceil (y_0 + y_1)/2 \rceil$ (thus $y_1 = \lfloor (y_0 + y_1)/2 \rfloor$).

Facts 5.17. If a sequence y_0, y_1, \dots, y_{w-1} has the step property,

- (1) then all its subsequences have the step property.
 (2) then its even and odd subsequences satisfy

$$\sum_{i=0}^{w/2-1} y_{2i} = \left\lfloor \frac{1}{2} \sum_{i=0}^{w-1} y_i \right\rfloor \text{ and } \sum_{i=0}^{w/2-1} y_{2i+1} = \left\lceil \frac{1}{2} \sum_{i=0}^{w-1} y_i \right\rceil.$$

Facts 5.18. If two sequences x_0, x_1, \dots, x_{w-1} and y_0, y_1, \dots, y_{w-1} have the step property,

- (1) and $\sum_{i=0}^{w-1} x_i = \sum_{i=0}^{w-1} y_i$, then $x_i = y_i$ for $i = 0, \dots, w-1$.
 (2) and $\sum_{i=0}^{w-1} x_i = \sum_{i=0}^{w-1} y_i + 1$, then there exists a unique j ($j = 0, 1, \dots, w-1$) such that $x_j = y_j + 1$, and $x_i = y_i$ for $i = 0, \dots, w-1, i \neq j$.

Remarks:

- An alternative representation of Batcher's network has been introduced by Aspnes et al. It is isomorphic to Batcher's network, and relies on a Merger Network $M[w]$ which is defined inductively: $M[w]$ consists of two $M[w/2]$ networks (an upper and a lower one) whose output is fed to $w/2$ balancers. The upper balancer merges the even subsequence x_0, x_2, \dots, x_{w-2} , while the lower balancer merges the odd subsequence x_1, x_3, \dots, x_{w-1} . Call the outputs of these two $M[w/2]$, z and z' respectively. The final stage of the network combines z and z' by sending each pair of wires z_i and z'_i into a balancer whose outputs yield y_{2i} and y_{2i+1} .
- It enough to prove that a merger network $M[w]$ preserves the step property.

Lemma 5.19. Let $M[w]$ be a merger network of width w . In a quiescent state (no message in transit), if the inputs $x_0, x_1, \dots, x_{w/2-1}$ resp. $x_{w/2}, x_{w/2+1}, \dots, x_{w-1}$ have the step property, then the output y_0, y_1, \dots, y_{w-1} has the step property.

Proof. By induction on the width w .

For $w = 2$: $M[2]$ is a balancer and a balancer's output has the step property (Fact 5.16.3).

For $w > 2$: Let z resp. z' be the output of the upper respectively lower $M[w/2]$ subnetwork. Since $x_0, x_1, \dots, x_{w/2-1}$ and $x_{w/2}, x_{w/2+1}, \dots, x_{w-1}$ both have the step property by assumption, their even and odd subsequences also have the step property (Fact 5.17.1). By induction hypothesis, the output of both $M[w/2]$ subnetworks have the step property. Let $Z := \sum_{i=0}^{w/2-1} z_i$ and $Z' := \sum_{i=0}^{w/2-1} z'_i$. From Fact 5.17.2 we conclude that $Z = \lceil \frac{1}{2} \sum_{i=0}^{w-1} x_i \rceil + \lfloor \frac{1}{2} \sum_{i=w/2}^{w-1} x_i \rfloor$ and $Z' = \lfloor \frac{1}{2} \sum_{i=0}^{w/2-1} x_i \rfloor + \lceil \frac{1}{2} \sum_{i=w/2}^{w-1} x_i \rceil$. Since $\lceil a \rceil + \lfloor b \rfloor$ and $\lfloor a \rfloor + \lceil b \rceil$ differ by at most 1 we know that Z and Z' differ by at most 1.

If $Z = Z'$, Fact 5.18.1 implies that $z_i = z'_i$ for $i = 0, \dots, w/2 - 1$. Therefore, the output of $M[w]$ is $y_i = z_{\lfloor i/2 \rfloor}$ for $i = 0, \dots, w - 1$. Since $z_0, \dots, z_{w/2-1}$ has the step property, so does the output of $M[w]$ and the lemma follows.

If Z and Z' differ by 1, Fact 5.18.2 implies that $z_i = z'_i$ for $i = 0, \dots, w/2 - 1$, except a unique j such that z_j and z'_j differ by only 1, for $j = 0, \dots, w/2 - 1$. Let $l := \min(z_j, z'_j)$. Then, the output y_i (with $i < 2j$) is $l + 1$. The output y_i (with $i > 2j + 1$) is l . The output y_{2j} and y_{2j+1} are balanced by the final balancer resulting in $y_{2j} = l + 1$ and $y_{2j+1} = l$. Therefore $M[w]$ preserves the step property. \square

A bitonic counting network is constructed to fulfill Lemma 5.19, i.e., the final output comes from a Merger whose upper and lower inputs are recursively merged. Therefore, the following theorem follows immediately.

Theorem 5.20 (Correctness). *In a quiescent state, the w output wires of a bitonic counting network of width w have the step property.*

Remarks:

- Is every sorting network also a counting network? No. But surprisingly, the other direction is true!

Theorem 5.21 (Counting vs. Sorting). *If a network is a counting network then it is also a sorting network, but not vice versa.*

Proof. There are sorting networks that are not counting networks (e.g. odd/even sort, or insertion sort). For the other direction, let C be a counting network and $I(C)$ be the isomorphic network, where every balancer is replaced by a comparator. Let $I(C)$ have an arbitrary input of 0's and 1's; that is, some of the input wires have a 0, all others have a 1. There is a message at C 's i^{th} input wire if and only if $I(C)$'s i input wire is 0. Since C is a counting network, all messages are routed to the upper output wires. $I(C)$ is isomorphic to C , therefore a comparator in $I(C)$ will receive a 0 on its upper (lower) wire if and only if the corresponding balancer receives a message on its upper (lower) wire. Using an inductive argument, the 0's and 1's will be routed through $I(C)$ such that all 0's exit the network on the upper wires whereas all 1's exit the network on the lower wires. Applying Lemma 5.3 shows that $I(C)$ is a sorting network. \square

Remarks:

- We claimed that the counting network is correct. However, it is only correct in a quiescent state.

Definition 5.22 (Linearizable). *A system is linearizable if the order of the values assigned reflects the real-time order in which they were requested. More formally, if there is a pair of operations o_1, o_2 , where operation o_1 terminates before operation o_2 starts, and the logical order is " o_2 before o_1 ", then a distributed system is not linearizable.*

Lemma 5.23 (Linearizability). *The bitonic counting network is not linearizable.*

Proof. Consider the bitonic counting network with width 4 in Figure 5.3: Assume that two *inc* operations were initiated and the corresponding messages entered the network on wire 0 and 2 (both in light gray color). After having passed the second resp. the first balancer, these traversing messages “fall asleep”; In other words, both messages take unusually long time before they are received by the next balancer. Since we are in an asynchronous setting, this may be the case.

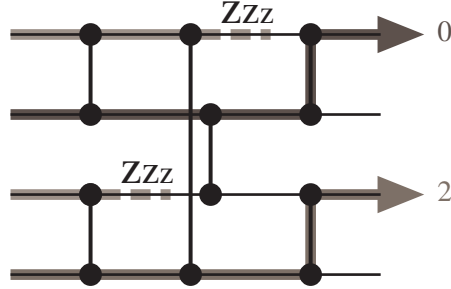


Figure 5.3: Linearizability Counter Example.

In the meantime, another *inc* operation (medium gray) is initiated and enters the network on the bottom wire. The message leaves the network on wire 2, and the *inc* operation is completed.

Strictly afterwards, another *inc* operation (dark gray) is initiated and enters the network on wire 1. After having passed all balancers, the message will leave the network wire 0. Finally (and not depicted in Figure 5.3), the two light gray messages reach the next balancer and will eventually leave the network on wires 1 resp. 3. Because the dark gray and the medium gray operation do conflict with Definition 5.22, the bitonic counting network is not linearizable. \square

Remarks:

- Note that the example in Figure 5.3 behaves correctly in the quiescent state: Finally, exactly the values 0, 1, 2, 3 are allotted.
- It was shown that linearizability comes at a high price (the depth grows linearly with the width).

Chapter 6

Shared Memory

6.1 Introduction

In distributed computing, various different models exist. So far, the focus of the course was on loosely-coupled distributed systems such as the Internet, where nodes asynchronously communicate by exchanging messages. The “opposite” model is a tightly-coupled parallel computer where nodes access a common memory totally synchronously—in distributed computing such a system is called a Parallel Random Access Machine (PRAM).

A third major model is somehow between these two extremes, the *shared memory* model. In a shared memory system, asynchronous processes (or processors) communicate via a common memory area of shared variables or registers:

Definition 6.1 (Shared Memory). *A shared memory system is a system that consists of asynchronous processes that access a common (shared) memory. A process can atomically access a register in the shared memory through a set of predefined operations. An atomic modification appears to the rest of the system instantaneously. Apart from this shared memory, processes can also have some local (private) memory.*

Remarks:

- Various shared memory systems exist. A main difference is how they allow processes to access the shared memory. All systems can atomically read or write a shared register R . Most systems do allow for advanced *atomic* read-modify-write (RMW) operations, for example:
 - test-and-set(R): $t := R$; $R := 1$; return t
 - fetch-and-add(R, x): $t := R$; $R := R + x$; return t
 - compare-and-swap(R, x, y): if $R = x$ then $R := y$; return **true**; else return **false**; endif;
 - load-link(R)/store-conditional(R, x): Load-link returns the current value of the specified register R . A subsequent store-conditional to the same register will store a new value x (and return **true**) only if no updates have occurred to that register since the load-link. If any updates have occurred, the store-conditional is guaranteed to fail

(and return **false**), even if the value read by the load-link has since been restored.

- The power of RMW operations can be measured with the so-called *consensus-number*: The consensus-number k of a RMW operation defines whether one can solve consensus for k processes. Test-and-set for instance has consensus-number 2 (one can solve consensus with 2 processes, but not 3), whereas the consensus-number of compare-and-swap is infinite. It can be shown that the power of a shared memory system is determined by the consensus-number (“universality of consensus”.) This insight has a remarkable theoretical and practical impact. In practice for instance, after this was known, hardware designers stopped developing shared memory systems supporting weak RMW operations.
- Many of the results derived in the message passing model have an equivalent in the shared memory model. Consensus for instance is traditionally studied in the shared memory model.
- Whereas programming a message passing system is rather tricky (in particular if fault-tolerance has to be integrated), programming a shared memory system is generally considered easier, as programmers are given access to global variables directly and do not need to worry about exchanging messages correctly. Because of this, even distributed systems which physically communicate by exchanging messages can often be programmed through a shared memory middleware, making the programmer’s life easier.
- We will most likely find the general spirit of shared memory systems in upcoming multi-core architectures. As for programming style, the multi-core community seems to favor an accelerated version of shared memory, *transactional memory*.
- From a message passing perspective, the shared memory model is like a bipartite graph: On one side you have the processes (the nodes) which pretty much behave like nodes in the message passing model (asynchronous, maybe failures). On the other side you have the shared registers, which just work perfectly (no failures, no delay).

6.2 Mutual Exclusion

A classic problem in shared memory systems is mutual exclusion. We are given a number of processes which occasionally need to access the same resource. The resource may be a shared variable, or a more general object such as a data structure or a shared printer. The catch is that only one process at the time is allowed to access the resource. More formally:

Definition 6.2 (Mutual Exclusion). *We are given a number of processes, each executing the following code sections:*

<Entry> → <Critical Section> → <Exit> → <Remaining Code>

A mutual exclusion algorithm consists of code for entry and exit sections, such that the following holds

- *Mutual Exclusion:* At all times at most one process is in the critical section.
- *No deadlock:* If some process manages to get to the entry section, later some (possibly different) process will get to the critical section.

Sometimes we in addition ask for

- *No lockout:* If some process manages to get to the entry section, later the same process will get to the critical section.
- *Unobstructed exit:* No process can get stuck in the exit section.

Using RMW primitives one can build mutual exclusion algorithms quite easily. Algorithm 23 shows an example with the test-and-set primitive.

Algorithm 23 Mutual Exclusion: Test-and-Set

Input: Shared register $R := 0$

<Entry>

1: **repeat**

2: $r := \text{test-and-set}(R)$

3: **until** $r = 0$

<Critical Section>

4: ...

<Exit>

5: $R := 0$

<Remainder Code>

6: ...

Theorem 6.3. *Algorithm 23 solves the mutual exclusion problem as in Definition 6.2.*

Proof. Mutual exclusion follows directly from the test-and-set definition: Initially R is 0. Let p_i be the i^{th} process to successfully execute the test-and-set, where successfully means that the result of the test-and-set is 0. This happens at time t_i . At time t'_i process p_i resets the shared register R to 0. Between t_i and t'_i no other process can successfully test-and-set, hence no other process can enter the critical section concurrently.

Proving no deadlock works similar: One of the processes loitering in the entry section will successfully test-and-set as soon as the process in the critical section exited.

Since the exit section only consists of a single instruction (no potential infinite loops) we have unobstructed exit. \square

Remarks:

- No lockout, on the other hand, is not given by this algorithm. Even with only two processes there are asynchronous executions where always the same process wins the test-and-set.
- Algorithm 23 can be adapted to guarantee fairness (no lockout), essentially by ordering the processes in the entry section in a queue.

- A natural question is whether one can achieve mutual exclusion with only reads and writes, that is without advanced RMW operations. The answer is yes!

Our read/write mutual exclusion algorithm is for two processes p_0 and p_1 only. In the remarks we discuss how it can be extended. The general idea is that process p_i has to mark its desire to enter the critical section in a “want” register W_i by setting $W_i := 1$. Only if the other process is not interested ($W_{1-i} = 0$) access is granted. This however is too simple since we may run into a deadlock. This deadlock (and at the same time also lockout) is resolved by adding a priority variable Π . See Algorithm 24.

Algorithm 24 Mutual Exclusion: Peterson’s Algorithm

Initialization: Shared registers W_0, W_1, Π , all initially 0.

Code for process p_i , $i = \{0, 1\}$

<Entry>

1: $W_i := 1$

2: $\Pi := 1 - i$

3: **repeat until** $\Pi = i$ or $W_{1-i} = 0$

<Critical Section>

4: ...

<Exit>

5: $W_i := 0$

<Remainder Code>

6: ...

Remarks:

- Note that line 3 in Algorithm 24 represents a “spinlock” or “busy-wait”, similarly to the lines 1-3 in Algorithm 23.

Theorem 6.4. *Algorithm 24 solves the mutual exclusion problem as in Definition 6.2.*

Proof. The shared variable Π elegantly grants priority to the process that passes line 2 first. If both processes are competing, only process p_Π can access the critical section because of Π . The other process $p_{1-\Pi}$ cannot access the critical section because $W_\Pi = 1$ (and $\Pi \neq 1 - \Pi$). The only other reason to access the critical section is because the other process is in the remainder code (that is, not interested). This proves mutual exclusion!

No deadlock comes directly with Π : Process p_Π gets direct access to the critical section, no matter what the other process does.

Since the exit section only consists of a single instruction (no potential infinite loops) we have unobstructed exit.

Thanks to the shared variable Π also no lockout (fairness) is achieved: If a process p_i loses against its competitor p_{1-i} in line 2, it will have to wait until the competitor resets $W_{1-i} := 0$ in the exit section. If process p_i is unlucky it will not check $W_{1-i} = 0$ early enough before process p_{1-i} sets $W_{1-i} := 1$ again in line 1. However, as soon as p_{1-i} hits line 2, process p_i gets the priority due to Π , and can enter the critical section. \square

Remarks:

- Extending Peterson's Algorithm to more than 2 processes can be done by a tournament tree, like in tennis. With n processes every process needs to win $\log n$ matches before it can enter the critical section. More precisely, each process starts at the bottom level of a binary tree, and proceeds to the parent level if winning. Once winning the root of the tree it can enter the critical section. Thanks to the priority variables Π at each node of the binary tree, we inherit all the properties of Definition 6.2.

6.3 Store & Collect

6.3.1 Problem Definition

In this section, we will look at a second shared memory problem that has an elegant solution. Informally, the problem can be stated as follows. There are n processes p_1, \dots, p_n . Every process p_i has a read/write register R_i in the shared memory where it can *store* some information that is destined for the other processes. Further, there is an operation by which a process can *collect* (i.e., read) the values of all the processes that stored some value in their register.

We say that an operation *op1* precedes an operation *op2* iff *op1* terminates before *op2* starts. An operation *op2* follows an operation *op1* iff *op1* precedes *op2*.

Definition 6.5 (Collect). *There are two operations: A STORE(val) by process p_i sets val to be the latest value of its register R_i . A COLLECT operation returns a view, a partial function V from the set of processes to a set of values, where $V(p_i)$ is the latest value stored by p_i , for each process p_i . For a COLLECT operation *cop*, the following validity properties must hold for every process p_i :*

- If $V(p_i) = \perp$, then no STORE operation by p_i precedes *cop*.
- If $V(p_i) = v \neq \perp$, then v is the value of a STORE operation *sop* of p_i that does not follow *cop*, and there is no STORE operation by p_i that follows *sop* and precedes *cop*.

Hence, a COLLECT operation *cop* should not read from the future or miss a preceding STORE operation *sop*.

We assume that the read/write register R_i of every process p_i is initialized to \perp . We define the step complexity of an operation *op* to be the number of accesses to registers in the shared memory. There is a trivial solution to the *collect* problem as shown by Algorithm 25.

Algorithm 25 Collect: Simple (Non-Adaptive) Solution

Operation STORE(val) (by process p_i) :

1: $R_i := val$

Operation COLLECT:

2: **for** $i := 1$ **to** n **do**

3: $V(p_i) := R_i$ // read register R_i

4: **end for**

Remarks:

- Algorithm 25 clearly works. The step complexity of every STORE operation is 1, the step complexity of a COLLECT operation is n .
- At first sight, the step complexities of Algorithm 25 seem optimal. Because there are n processes, there clearly are cases in which a COLLECT operation needs to read all n registers. However, there are also scenarios in which the step complexity of the COLLECT operation seems very costly. Assume that there are only two processes p_i and p_j that have stored a value in their registers R_i and R_j . In this case, a COLLECT in principle only needs to read the registers R_i and R_j and can ignore all the other registers.
- Assume that up to a certain time t , $k \leq n$ processes have finished or started at least one operation. We call an operation op at time t *adaptive* to contention if the step complexity of op only depends on k and is independent of n .
- In the following, we will see how to implement adaptive versions of STORE and COLLECT.

6.3.2 Splitters

Algorithm 26 Splitter Code

Shared Registers: $X : \{\perp\} \cup \{1, \dots, n\}$; $Y : \text{boolean}$ **Initialization:** $X := \perp$; $Y := \text{false}$ **Splitter access by process p_i :**

```

1:  $X := i$ ;
2: if  $Y$  then
3:   return right
4: else
5:    $Y := \text{true}$ 
6:   if  $X = i$  then
7:     return stop
8:   else
9:     return left
10:  end if
11: end if

```

To obtain adaptive collect algorithms, we need a synchronization primitive, called a *splitter*.

Definition 6.6 (Splitter). *A splitter is a synchronization primitive with the following characteristic. A process entering a splitter exits with either **stop**, **left**, or **right**. If k processes enter a splitter, at most one process exits with **stop** and at most $k - 1$ processes exit with **left** and **right**, respectively.*

Hence, it is guaranteed that if a single process enters the splitter, then it obtains **stop**, and if two or more processes enter the splitter, then there is at most one process obtaining **stop** and there are two processes that obtain

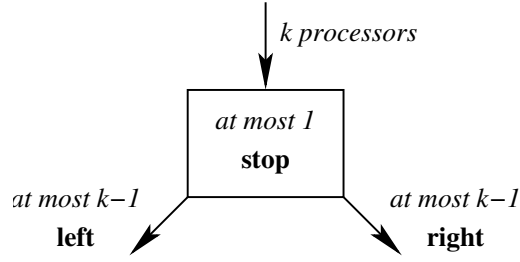


Figure 6.1: A Splitter

different values (i.e., either there is exactly one **stop** or there is at least one **left** and at least one **right**). For an illustration, see Figure 6.1. The code implementing a splitter is given by Algorithm 26.

Lemma 6.7. *Algorithm 26 correctly implements a splitter.*

Proof. Assume that k processes enter the splitter. Because the first process that checks whether $Y = \mathbf{true}$ in line 2 will find that $Y = \mathbf{false}$, not all processes return **right**. Next, assume that i is the last process that sets $X := i$. If i does not return **right**, it will find $X = i$ in line 6 and therefore return **stop**. Hence, there is always a process that does not return **left**. It remains to show that at most 1 process returns **stop**. For the sake of contradiction, assume p_i and p_j are two processes that return **stop** and assume that p_i sets $X := i$ before p_j sets $X := j$. Both processes need to check whether Y is **true** before one of them sets $Y := \mathbf{true}$. Hence, they both complete the assignment in line 1 before the first one of them checks the value of X in line 6. Hence, by the time p_i arrives at line 6, $X \neq i$ (p_j and maybe some other processes have overwritten X by then). Therefore, p_i does not return **stop** and we get a contradiction to the assumption that both p_i and p_j return **stop**. \square

6.3.3 Binary Splitter Tree

Assume that we are given $2^n - 1$ splitters and that for every splitter S , there is an additional shared variable $Z_S : \{\perp\} \cup \{1, \dots, n\}$ that is initialized to \perp and an additional shared variable $M_S : \mathbf{boolean}$ that is initialized to **false**. We call a splitter S marked if $M_S = \mathbf{true}$. The $2^n - 1$ splitters are arranged in a complete binary tree of height $n - 1$. Let $S(v)$ be the splitter associated with a node v of the binary tree. The STORE and COLLECT operations are given by Algorithm 27.

Theorem 6.8. *Algorithm 27 correctly implements STORE and COLLECT. Let k be the number of participating processes. The step complexity of the first STORE of a process p_i is $\mathcal{O}(k)$, the step complexity of every additional STORE of p_i is $\mathcal{O}(1)$, and the step complexity of COLLECT is $\mathcal{O}(k)$.*

Proof. Because at most one process can stop at a splitter, it is sufficient to show that every process stops at some splitter at depth at most $k - 1 \leq n - 1$ when invoking the first STORE operation to prove correctness. We prove that at most $k - i$ processes enter a subtree at depth i (i.e., a subtree where the root has distance i to the root of the whole tree). By definition of k , the number of

Algorithm 27 Adaptive Collect: Binary Tree Algorithm**Operation** STORE(*val*) (by process p_i) :

```

1:  $R_i := val$ 
2: if first STORE operation by  $p_i$  then
3:    $v :=$  root node of binary tree
4:    $\alpha :=$  result of entering splitter  $S(v)$ ;
5:    $M_{S(v)} := \mathbf{true}$ 
6:   while  $\alpha \neq \mathbf{stop}$  do
7:     if  $\alpha = \mathbf{left}$  then
8:        $v :=$  left child of  $v$ 
9:     else
10:       $v :=$  right child of  $v$ 
11:     end if
12:      $\alpha :=$  result of entering splitter  $S(v)$ ;
13:      $M_{S(v)} := \mathbf{true}$ 
14:   end while
15:    $Z_{S(v)} := i$ 
16: end if

```

Operation COLLECT:**Traverse marked part of binary tree:**

```

17: for all marked splitters  $S$  do
18:   if  $Z_S \neq \perp$  then
19:      $i := Z_S$ ;  $V(p_i) := R_i$  // read value of process  $p_i$ 
20:   end if
21: end for //  $V(p_i) = \perp$  for all other processes

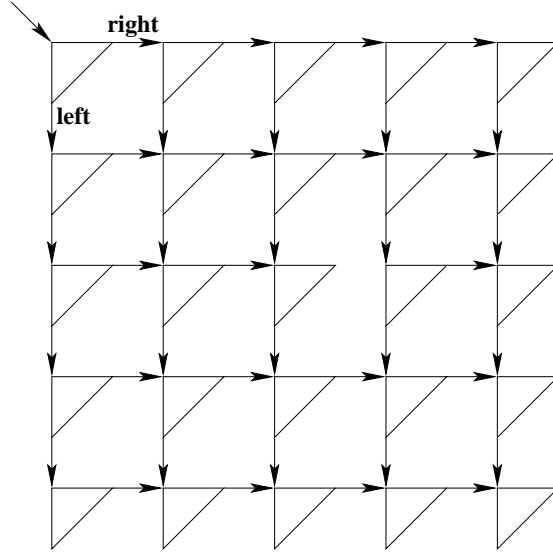
```

processes entering the splitter at depth 0 (i.e., at the root of the binary tree) is k . For $i > 1$, the claim follows by induction because of the at most $k - i$ processes entering the splitter at the root of a depth i subtree, at most $k - i - 1$ obtain **left** and **right**, respectively. Hence, at the latest when reaching depth $k - 1$, a process is the only process entering a splitter and thus obtains **stop**. It thus also follows that the step complexity of the first invocation of STORE is $\mathcal{O}(k)$.

To show that the step complexity of COLLECT is $\mathcal{O}(k)$, we first observe that the marked nodes of the binary tree are connected, and therefore can be traversed by only reading the variables M_S associated to them and their neighbors. Hence, showing that at most $2k - 1$ nodes of the binary tree are marked is sufficient. Let x_k be the maximum number of marked nodes in a tree, where k processes access the root. We claim that $x_k \leq 2k - 1$, which is true for $k = 1$ because a single process entering a splitter will always compute **stop**. Now assume the inequality holds for $1, \dots, k - 1$. Not all k processes may exit the splitter with **left** (or **right**), i.e., $k_l \leq k - 1$ processes will turn left and $k_r \leq \min\{k - k_l, k - 1\}$ turn right. The left and right children of the root are the roots of their subtrees, hence the induction hypothesis yields

$$x_k = x_{k_l} + x_{k_r} + 1 \leq (2k_l - 1) + (2k_r - 1) + 1 \leq 2k - 1,$$

concluding induction and proof. \square

Figure 6.2: 5×5 Splitter Matrix**Remarks:**

- The step complexities of Algorithm 27 are very good. Clearly, the step complexity of the COLLECT operation is asymptotically optimal. In order for the algorithm to work, we however need to allocate the memory for the complete binary tree of depth $n - 1$. The space complexity of Algorithm 27 therefore is exponential in n . We will next see how to obtain a polynomial space complexity at the cost of a worse COLLECT step complexity.

6.3.4 Splitter Matrix

Instead of arranging splitters in a binary tree, we arrange n^2 splitters in an $n \times n$ matrix as shown in Figure 6.2. The algorithm is analogous to Algorithm 27. The matrix is entered at the top left. If a process receives **left**, it next visits the splitter in the next row of the same column. If a process receives **right**, it next visits the splitter in the next column of the same row. Clearly, the space complexity of this algorithm is $\mathcal{O}(n^2)$. The following theorem gives bounds on the step complexities of STORE and COLLECT.

Theorem 6.9. *Let k be the number of participating processes. The step complexity of the first STORE of a process p_i is $\mathcal{O}(k)$, the step complexity of every additional STORE of p_i is $\mathcal{O}(1)$, and the step complexity of COLLECT is $\mathcal{O}(k^2)$.*

Proof. Let the top row be row 0 and the left-most column be column 0. Let x_i be the number of processes entering a splitter in row i . By induction on i , we show that $x_i \leq k - i$. Clearly, $x_0 \leq k$. Let us therefore consider the case $i > 0$. Let j be the largest column such that at least one process visits the splitter in row $i - 1$ and column j . By the properties of splitters, not all processes entering the splitter in row $i - 1$ and column j obtain **left**. Therefore, not all processes entering a splitter in row $i - 1$ move on to row i . Because at least one process

stays in every row, we get that $x_i \leq k - i$. Similarly, the number of processes entering column j is at most $k - j$. Hence, every process stops at the latest in row $k - 1$ and column $k - 1$ and the number of marked splitters is at most k^2 . Thus, the step complexity of COLLECT is at most $\mathcal{O}(k^2)$. Because the longest path in the splitter matrix is $2k$, the step complexity of STORE is $\mathcal{O}(k)$. \square

Remarks:

- With a slightly more complicated argument, it is possible to show that the number of processes entering the splitter in row i and column j is at most $k - i - j$. Hence, it suffices to only allocate the upper left half (including the diagonal) of the $n \times n$ matrix of splitters.
- The binary tree algorithm can be made space efficient by using a randomized version of a splitter. Whenever returning left or right, a randomized splitter returns left or right with probability $1/2$. With high probability, it then suffices to allocate a binary tree of depth $\mathcal{O}(\log n)$.
- Recently, it has been shown that with a considerably more complicated deterministic algorithm, it is possible to achieve $\mathcal{O}(k)$ step complexity and $\mathcal{O}(n^2)$ space complexity.

Chapter 7

Shared Objects

7.1 Introduction

Assume that there is a common resource (e.g. a common variable or data structure), which different nodes in a network need to access from time to time. If the nodes are allowed to change the common object when accessing it, we need to guarantee that no two nodes have access to the object at the same time. In order to achieve this mutual exclusion, we need protocols that allow the nodes of a network to store and manage access to such a shared object. A simple and obvious solution is to store the shared object at a central location (see Algorithm 28).

Algorithm 28 Shared Object: Centralized Solution

Initialization: Shared object stored at root node r of a spanning tree of the network graph (i.e., each node knows its parent in the spanning tree).

Accessing Object: (by node v)

- 1: v sends request up the tree
 - 2: request processed by root r (atomically)
 - 3: result sent down the tree to node v
-

Remarks:

- Instead of a spanning tree, one can use routing.
- Algorithm 28 works, but it is not very efficient. Assume that the object is accessed by a single node v repeatedly. Then we get a high message/time complexity. Instead v could store the object, or at least cache it. But then, in case another node w accesses the object, we might run into consistency problems.
- Alternative idea: The accessing node should become the new master of the object. The shared object then becomes mobile. There exist several variants of this idea. The simplest version is a home-based solution like in Mobile IP (see Algorithm 29).

Algorithm 29 Shared Object: Home-Based Solution

Initialization: An object has a home base (a node) that is known to every node. All requests (accesses to the shared object) are routed through the home base.

Accessing Object: (by node v)

- 1: v acquires a lock at the home base, receives object.
-

Remarks:

- Home-based solutions suffer from the triangular routing problem. If two close-by nodes access the object on a rotating basis, all the traffic is routed through the potentially far away home-base.

7.2 Arrow and Friends

We will now look at a protocol (called the Arrow algorithm) that always moves the shared object to the node currently accessing it without creating the triangular routing problem of home-based solutions. The protocol runs on a precomputed spanning tree. Assume that the spanning tree is rooted at the current position of the shared object. When a node u wants to access the shared object, it sends out a *find* request towards the current position of the object. While searching for the object, the edges of the spanning tree are redirected such that in the end, the spanning tree is rooted at u (i.e., the new holder of the object). The details of the algorithm are given by Algorithm 30. For simplicity, we assume that a node u only starts a find request if u is not currently the holder of the shared object and if u has finished all previous find requests (i.e., it is not currently waiting to receive the object).

Remarks:

- The parent pointers in Algorithm 30 are only needed for the find operation. Sending the variable to u in line 13 or to w .successor in line 23 is done using routing (on the spanning tree or on the underlying network).
- When we draw the parent pointers as arrows, in a quiescent moment (where no “find” is in motion), the arrows all point towards the node currently holding the variable (i.e., the tree is rooted at the node holding the variable)
- What is really great about the Arrow algorithm is that it works in a completely asynchronous and concurrent setting (i.e., there can be many find requests at the same time).

Theorem 7.1. (*Arrow, Analysis*) *In an asynchronous and concurrent setting, a “find” operation terminates with message and time complexity D , where D is the diameter of the spanning tree.*

Algorithm 30 Shared Object: Arrow Algorithm

Initialization: As for Algorithm 28, we are given a rooted spanning tree. Each node has a pointer to its parent, the root r is its own parent. The variable is initially stored at r . For all nodes v , $v.\text{successor} := \text{null}$, $v.\text{wait} := \text{false}$.

Start Find Request at Node u :

```

1: do atomically
2:    $u$  sends “find by  $u$ ” message to parent node
3:    $u.\text{parent} := u$ 
4:    $u.\text{wait} := \text{true}$ 
5: end do

```

Upon w Receiving “Find by u ” Message from Node v :

```

6: do atomically
7:   if  $w.\text{parent} \neq w$  then
8:      $w$  sends “find by  $u$ ” message to parent
9:      $w.\text{parent} := v$ 
10:  else
11:     $w.\text{parent} := v$ 
12:    if not  $w.\text{wait}$  then
13:      send variable to  $u$       //  $w$  holds var. but does not need it any more
14:    else
15:       $w.\text{successor} := u$       //  $w$  will send variable to  $u$  a.s.a.p.
16:    end if
17:  end if
18: end do

```

Upon w Receiving Shared Object:

```

19: perform operation on shared object
20: do atomically
21:    $w.\text{wait} := \text{false}$ 
22:   if  $w.\text{successor} \neq \text{null}$  then
23:     send variable to  $w.\text{successor}$ 
24:      $w.\text{successor} := \text{null}$ 
25:   end if
26: end do

```

Before proving Theorem 7.1, we prove the following lemma.

Lemma 7.2. *An edge $\{u, v\}$ of the spanning tree is in one of four states:*

- 1.) *Pointer from u to v (no message on the edge, no pointer from v to u)*
- 2.) *Message on the move from u to v (no pointer along the edge)*
- 3.) *Pointer from v to u (no message on the edge, no pointer from u to v)*
- 4.) *Message on the move from v to u (no pointer along the edge)*

Proof. W.l.o.g., assume that initially the edge $\{u, v\}$ is in state 1. With a message arrival at u (or if u starts a “find by u ” request, the edge goes to state 2. When the message is received at v , v directs its pointer to u and we are therefore in state 3. A new message at v (or a new request initiated by v) then brings the edge back to state 1. \square

Proof of Theorem 7.1. Since the “find” message will only travel on a static tree, it suffices to show that it will not traverse an edge twice. Suppose for the sake of contradiction that there is a first “find” message f that traverses an edge $e = \{u, v\}$ for the second time and assume that e is the first edge that is traversed twice by f . The first time, f traverses e . Assume that e is first traversed from u to v . Since we are on a tree, the second time, e must be traversed from v to u . Because e is the first edge to be traversed twice, f must re-visit e before visiting any other edges. Right before f reaches v , the edge e is in state 2 (f is on the move) and in state 3 (it will immediately return with the pointer from v to u). This is a contradiction to Lemma 7.2. \square

Remarks:

- Finding a good tree is an interesting problem. We would like to have a tree with low stretch, low diameter, low degree, etc.
- It seems that the Arrow algorithm works especially well when lots of “find” operations are initiated concurrently. Most of them will find a “close-by” node, thus having low message/time complexity. For the sake of simplicity we analyze a synchronous system.

Theorem 7.3. *(Arrow, Concurrent Analysis) Let the system be synchronous. Initially, the system is in a quiescent state. At time 0, a set S of nodes initiates a “find” operation. The message complexity of all “find” operations is $\mathcal{O}(\log |S| \cdot m^*)$ where m^* is the message complexity of an optimal (with global knowledge) algorithm on the tree.*

Proof Sketch. Let d be the minimum distance of any node in S to the root. There will be a node u_1 at distance d from the root that reaches the root in d time steps, turning all the arrows on the path to the root towards u_1 . A node u_2 that finds (is queued behind) u_1 cannot distinguish the system from a system where there was no request u_1 , and instead the root was initially located at u_1 . The message cost of u_2 is consequentially the distance between u_1 and u_2 on the spanning tree. By induction the total message complexity is exactly as if a collector starts at the root and then “greedily” collects tokens located at the nodes in S (greedily in the sense that the collector always goes towards the closest token). Greedy collecting the tokens is not a good strategy in general because it will traverse the same edge more than twice in the worst

case. An asymptotically optimal algorithm can also be translated into a depth-first-search collecting paradigm, traversing each edge at most twice. In another area of computer science, we would call the Arrow algorithm a nearest-neighbor TSP heuristic (without returning to the start/root though), and the optimal algorithm TSP-optimal. It was shown that nearest-neighbor has a logarithmic overhead, which concludes the proof. \square

Remarks:

- An average request set S on a not-too-bad tree gives usually a much better bound. However, there is an almost tight $\log |S| / \log \log |S|$ worst-case example.
- It was recently shown that Arrow can do as good in a dynamic setting (where nodes are allowed to initiate requests at any time). In particular the message complexity of the dynamic analysis can be shown to have a $\log D$ overhead only, where D is the diameter of the spanning tree (note that for logarithmic trees, the overhead becomes $\log \log n$).
- What if the spanning tree is a star? Then with Theorem 7.1, each find will terminate in 2 steps! Since also an optimal algorithm has message cost 1, the algorithm is 2-competitive. . . ? Yes, but because of its high degree the star center experiences contention. . . It can be shown that the contention overhead is at most proportional to the largest degree Δ of the spanning tree.
- Thought experiment: Assume a balanced binary spanning tree—by Theorem 7.1, the message complexity per operation is $\log n$. Because a binary tree has maximum degree 3, the time per operation therefore is at most $3 \log n$.
- There are better and worse choices for the spanning tree. The stretch of an edge $\{u, v\}$ is defined as distance between u and v in a spanning tree. The maximum stretch of a spanning tree is the maximum stretch over all edges. A few years ago, it was shown how to construct spanning trees that are $\mathcal{O}(\log n)$ -stretch-competitive.

What if most nodes just want to read the shared object? Then it does not make sense to acquire a lock every time. Instead we can use caching (see Algorithm 31).

Theorem 7.4. *Algorithm 31 is correct. More surprisingly, the message complexity is 3-competitive (at most a factor 3 worse than the optimum).*

Proof. Since the accesses do not overlap by definition, it suffices to show that between two writes, we are 3-competitive. The sequence of accessing nodes is $w_0, r_1, r_2, \dots, r_k, w_1$. After w_0 , the object is stored at w_0 and not cached anywhere else. All reads cost twice the smallest subtree T spanning the write w_0 and all the reads since each read only goes to the first copy. The write w_1 costs T plus the path P from w_1 to T . Since any data management scheme must use an edge in T and P at least once, and our algorithm uses edges in T at most 3 times (and in P at most once), the theorem follows. \square

Algorithm 31 Shared Object: Read/Write Caching

- Nodes can either read or write the shared object. For simplicity we first assume that reads or writes do not overlap in time (access to the object is sequential).
 - Nodes store three items: a parent pointer pointing to one of the neighbors (as with Arrow), and a cache bit for each edge, plus (potentially) a copy of the object.
 - Initially the object is stored at a single node u ; all the parent pointers point towards u , all the cache bits are false.
 - When initiating a read, a message follows the arrows (this time: without inverting them!) until it reaches a cached version of the object. Then a copy of the object is cached along the path back to the initiating node, and the cache bits on the visited edges are set to true.
 - A write at u writes the new value locally (at node u), then searches (follow the parent pointers and reverse them towards u) a first node with a copy. Delete the copy and follow (in parallel, by flooding) all edge that have the cache flag set. Point the parent pointer towards u , and remove the cache flags.
-

Remarks:

- Concurrent reads are not a problem, also multiple concurrent reads and one write work just fine.
- What about concurrent writes? To achieve consistency writes need to invalidate the caches before writing their value. It is claimed that the strategy then becomes 4-competitive.
- Is the algorithm also time competitive? Well, not really: The optimal algorithm that we compare to is usually offline. This means it knows the whole access sequence in advance. It can then cache the object before the request even appears!
- Algorithms on trees are often simpler, but have the disadvantage that they introduce the extra stretch factor. In a ring, for example, any tree has stretch $n - 1$; so there is always a bad request pattern.

Algorithm 32 Shared Object: Pointer Forwarding

Initialization: Object is stored at root r of a precomputed spanning tree T (as in the Arrow algorithm, each node has a parent pointer pointing towards the object).

Accessing Object: (by node u)

- 1: follow parent pointers to current root r of T
 - 2: send object from r to u
 - 3: $r.\text{parent} := u$; $u.\text{parent} := u$; *// u is the new root*
-

Algorithm 33 Shared Object: Ivy

Initialization: Object is stored at root r of a precomputed spanning tree T (as before, each node has a parent pointer pointing towards the object). For simplicity, we assume that accesses to the object are sequential.

Start Find Request at Node u :

- 1: u sends “find by u ” message to parent node
- 2: $u.\text{parent} := u$

Upon v receiving “Find by u ” Message:

- 3: **if** $v.\text{parent} = v$ **then**
 - 4: send object to u
 - 5: **else**
 - 6: send “find by u ” message to $v.\text{parent}$
 - 7: **end if**
 - 8: $v.\text{parent} := u$ *// u will become the new root*
-

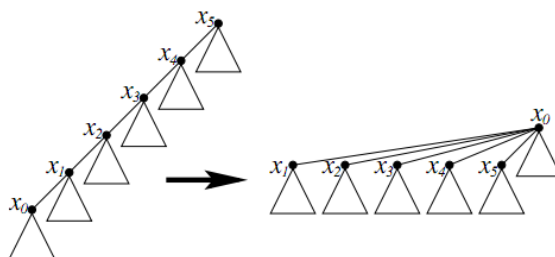
7.3 Ivy and Friends

In the following we study algorithms that do not restrict communication to a tree. Of particular interest is the special case of a complete graph (clique). A simple solution for this case is given by Algorithm 32.

Remarks:

- If the graph is not complete, routing can be used to find the root.
- Assume that the nodes line up in a linked list. If we always choose the first node of the linked list to acquire the object, we have message/time complexity n . The new topology is again a linear linked list. Pointer forwarding is therefore bad in a worst-case.
- If edges are not FIFO, it can even happen that the number of steps is unbounded for a node having bad luck. An algorithm with such a property is named “not fair,” or “not wait-free.” (Example: Initially we have the list $4 \rightarrow 3 \rightarrow 2 \rightarrow 1$; 4 starts a find; when the message of 4 passes 3, 3 itself starts a find. The message of 3 may arrive at 2 and then 1 earlier, thus the new end of the list is $2 \rightarrow 1 \rightarrow 3$; once the message of 4 passes 2, the game re-starts.)

There seems to be a natural improvement of the pointer forwarding idea. Instead of simply redirecting the parent pointer from the old root to the new root, we can redirect all the parent pointers of the nodes on the path visited

Figure 7.1: Reversal of the path $x_0, x_1, x_2, x_3, x_4, x_5$.

during a find message to the new root. The details are given by Algorithm 33. Figure 7.1 shows how the pointer redirecting affects a given tree (the right tree results from a find request started at node x_0 on the left tree).

Remarks:

- Also with Algorithm 33, we might have a bad linked list situation. However, if the start of the list acquires the object, the linked list turns into a star. As the following theorem shows, the search paths are not long on average. Since paths sometimes can be bad, we will need amortized analysis.

Theorem 7.5. *If the initial tree is a star, a find request of Algorithm 33 needs at most $\log n$ steps on average, where n is the number of processors.*

Proof. All logarithms in the following proof are to base 2. We assume that accesses to the shared object are sequential. We use a potential function argument. Let $s(u)$ be the size of the subtree rooted at node u (the number of nodes in the subtree including u itself). We define the potential Φ of the whole tree T as (V is the set of all nodes)

$$\Phi(T) = \sum_{u \in V} \frac{\log s(u)}{2}.$$

Assume that the path traversed by the i^{th} operation has length k_i , i.e., the i^{th} operation redirects k_i pointers to the new root. Clearly, the number of steps of the i^{th} operation is proportional to k_i . We are interested in the cost of m consecutive operations, $\sum_{i=1}^m k_i$.

Let T_0 be the initial tree and let T_i be the tree after the i^{th} operation. Further, let $a_i = k_i - \Phi(T_{i-1}) + \Phi(T_i)$ be the *amortized cost* of the i^{th} operation. We have

$$\sum_{i=1}^m a_i = \sum_{i=1}^m (k_i - \Phi(T_{i-1}) + \Phi(T_i)) = \sum_{i=1}^m k_i - \Phi(T_0) + \Phi(T_m).$$

For any tree T , we have $\Phi(T) \geq \log(n)/2$. Because we assume that T_0 is a star, we also have $\Phi(T_0) = \log(n)/2$. We therefore get that

$$\sum_{i=1}^m a_i \geq \sum_{i=1}^m k_i.$$

Hence, it suffices to upper bound the amortized cost of every operation. We thus analyze the amortized cost a_i of the i^{th} operation. Let $x_0, x_1, x_2, \dots, x_{k_i}$ be the path that is reversed by the operation. Further for $0 \leq j \leq k_i$, let s_j be the size of the subtree rooted at x_j before the reversal. The size of the subtree rooted at x_0 after the reversal is s_{k_i} and the size of the one rooted at x_j after the reversal, for $1 \leq j \leq k_i$, is $s_j - s_{j-1}$ (see Figure 7.1). For all other nodes, the sizes of their subtrees are the same, therefore the corresponding terms cancel out in the amortized cost a_i . We can thus write a_i as

$$\begin{aligned} a_i &= k_i - \left(\sum_{j=0}^{k_i} \frac{1}{2} \log s_j \right) + \left(\frac{1}{2} \log s_{k_i} + \sum_{j=1}^{k_i} \frac{1}{2} \log(s_j - s_{j-1}) \right) \\ &= k_i + \frac{1}{2} \cdot \sum_{j=0}^{k_i-1} (\log(s_{j+1} - s_j) - \log s_j) \\ &= k_i + \frac{1}{2} \cdot \sum_{j=0}^{k_i-1} \log \left(\frac{s_{j+1} - s_j}{s_j} \right). \end{aligned}$$

For $0 \leq j \leq k_i - 1$, let $\alpha_j = s_{j+1}/s_j$. Note that $s_{j+1} > s_j$ and thus that $\alpha_j > 1$. Further note, that $(s_{j+1} - s_j)/s_j = \alpha_j - 1$. We therefore have that

$$\begin{aligned} a_i &= k_i + \frac{1}{2} \cdot \sum_{j=0}^{k_i-1} \log(\alpha_j - 1) \\ &= \sum_{j=0}^{k_i-1} \left(1 + \frac{1}{2} \log(\alpha_j - 1) \right). \end{aligned}$$

For $\alpha > 1$, it can be shown that $1 + \log(\alpha - 1)/2 \leq \log \alpha$ (see Lemma 7.6). From this inequality, we obtain

$$\begin{aligned} a_i &\leq \sum_{j=0}^{k_i-1} \log \alpha_j = \sum_{j=0}^{k_i-1} \log \frac{s_{j+1}}{s_j} = \sum_{j=0}^{k_i-1} (\log s_{j+1} - \log s_j) \\ &= \log s_{k_i} - \log s_0 \leq \log n, \end{aligned}$$

because $s_{k_i} = n$ and $s_0 \geq 1$. This concludes the proof. \square

Lemma 7.6. For $\alpha > 1$, $1 + \log(\alpha - 1)/2 \leq \log \alpha$.

Proof. The claim can be verified by the following chain of reasoning:

$$\begin{aligned} 0 &\leq (\alpha - 2)^2 \\ 0 &\leq \alpha^2 - 4\alpha + 4 \\ 4(\alpha - 1) &\leq \alpha^2 \\ \log_2(4(\alpha - 1)) &\leq \log_2(\alpha^2) \\ 2 + \log_2(\alpha - 1) &\leq 2 \log_2 \alpha \\ 1 + \frac{1}{2} \log_2(\alpha - 1) &\leq \log_2 \alpha. \end{aligned}$$

\square

Remarks:

- Systems guys (the algorithm is called Ivy because it was used in a system with the same name) have some fancy heuristics to improve performance even more: For example, the root every now and then broadcasts its name such that paths will be shortened.
- What about concurrent requests? It works with the same argument as in Arrow. Also for Ivy an argument including congestion is missing (and more pressing, since the dynamic topology of a tree cannot be chosen to have low degree and thus low congestion as in Arrow).
- Sometimes the type of accesses allows that several accesses can be combined into one to reduce congestion higher up the tree. Let the tree in Algorithm 28 be a balanced binary tree. If the access to a shared variable for example is “add value x to the shared variable”, two or more accesses that accidentally meet at a node can be combined into one. Clearly accidental meeting is rare in an asynchronous model. We might be able to use synchronizers (or maybe some other timing tricks) to help meeting a little bit.

Chapter 8

Maximal Independent Set

In this chapter we present a highlight of this course, a fast maximal independent set (MIS) algorithm. The algorithm is the first randomized algorithm that we study in this class. In distributed computing, randomization is a powerful and therefore omnipresent concept, as it allows for relatively simple yet efficient algorithms. As such the studied algorithm is archetypal.

A MIS is a basic building block in distributed computing, some other problems pretty much follow directly from the MIS problem. At the end of this chapter, we will give two examples: matching and vertex coloring (see Chapter 2).

8.1 MIS

Definition 8.1 (Independent Set). *Given an undirected Graph $G = (V, E)$ an independent set is a subset of nodes $U \subseteq V$, such that no two nodes in U are adjacent. An independent set is maximal if no node can be added without violating independence. An independent set of maximum cardinality is called maximum.*

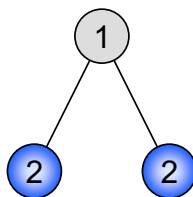


Figure 8.1: Example graph with 1) a maximal independent set (MIS) and 2) a maximum independent set (MaxIS).

Remarks:

- Computing a maximum independent set (MaxIS) is a notoriously difficult problem. It is equivalent to maximum clique on the complementary graph. Both problems are NP-hard, in fact not approximable within $n^{\frac{1}{2}-\epsilon}$.
- In this course we concentrate on the maximal independent set (MIS) problem. Please note that MIS and MaxIS can be quite different, indeed e.g. on a star graph the MIS is $\Theta(n)$ smaller than the MaxIS (cf. Figure 8.1).
- Computing a MIS sequentially is trivial: Scan the nodes in arbitrary order. If a node u does not violate independence, add u to the MIS. If u violates independence, discard u . So the only question is how to compute a MIS in a distributed way.

Algorithm 34 Slow MIS

Require: Node IDs**Every node** v executes the following code:

- 1: **if** all neighbors of v with larger identifiers have decided not to join the MIS **then**
 - 2: v decides to join the MIS
 - 3: **end if**
-

Remarks:

- Not surprisingly the slow algorithm is not better than the sequential algorithm in the worst case, because there might be one single point of activity at any time. Formally:

Theorem 8.2 (Analysis of Algorithm 34). *Algorithm 34 features a time complexity of $O(n)$ and a message complexity of $O(m)$.*

Remarks:

- This is not very exciting.
- There is a relation between independent sets and node coloring (Chapter 2), since each color class is an independent set, however, not necessarily a MIS. Still, starting with a coloring, one can easily derive a MIS algorithm: We first choose all nodes of the first color. Then, for each additional color we add “in parallel” (without conflict) as many nodes as possible. Thus the following corollary holds:

Corollary 8.3. *Given a coloring algorithm that needs C colors and runs in time T , we can construct a MIS in time $C + T$.*

Remarks:

- Using Theorem 2.12 and Corollary 8.3 we get a distributed deterministic MIS algorithm for trees (and for bounded degree graphs) with time complexity $O(\log^* n)$.

- With a lower bound argument one can show that this deterministic MIS algorithm for rings is asymptotically optimal.
- There have been attempts to extend Algorithm 5 to more general graphs, however, so far without much success. Below we present a radically different approach that uses randomization. Please note that the algorithm and the analysis below is not identical with the algorithm in Peleg's book.

8.2 Fast MIS from 1986

Algorithm 35 Fast MIS

The algorithm operates in synchronous rounds, grouped into phases.

A single phase is as follows:

- 1) Each node v marks itself with probability $\frac{1}{2d(v)}$, where $d(v)$ is the current degree of v .
 - 2) If no higher degree neighbor of v is also marked, node v joins the MIS. If a higher degree neighbor of v is marked, node v unmarks itself again. (If the neighbors have the same degree, ties are broken arbitrarily, e.g., by identifier).
 - 3) Delete all nodes that joined the MIS and their neighbors, as they cannot join the MIS anymore.
-

Remarks:

- Correctness in the sense that the algorithm produces an independent set is relatively simple: Steps 1 and 2 make sure that if a node v joins the MIS, then v 's neighbors do not join the MIS at the same time. Step 3 makes sure that v 's neighbors will never join the MIS.
- Likewise the algorithm eventually produces a MIS, because the node with the highest degree will mark itself at some point in Step 1.
- So the only remaining question is how fast the algorithm terminates. To understand this, we need to dig a bit deeper.

Lemma 8.4 (Joining MIS). *A node v joins the MIS in step 2 with probability $p \geq \frac{1}{4d(v)}$.*

Proof: Let M be the set of marked nodes in step 1. Let $H(v)$ be the set of neighbors of v with higher degree, or same degree and higher identifier. Using independence of the random choices of v and nodes in $H(v)$ in Step 1 we get

$$\begin{aligned}
 P[v \notin \text{MIS} | v \in M] &= P[\exists w \in H(v), w \in M | v \in M] \\
 &= P[\exists w \in H(v), w \in M] \\
 &\leq \sum_{w \in H(v)} P[w \in M] = \sum_{w \in H(v)} \frac{1}{2d(w)} \\
 &\leq \sum_{w \in H(v)} \frac{1}{2d(v)} \leq \frac{d(v)}{2d(v)} = \frac{1}{2}.
 \end{aligned}$$

Then

$$P[v \in \text{MIS}] = P[v \in \text{MIS} | v \in M] \cdot P[v \in M] \geq \frac{1}{2} \cdot \frac{1}{2d(v)}.$$

□

Lemma 8.5 (Good Nodes). *A node v is called good if*

$$\sum_{w \in N(v)} \frac{1}{2d(w)} \geq \frac{1}{6}.$$

Otherwise we call v a bad node. A good node will be removed in Step 3 with probability $p \geq \frac{1}{36}$.

Proof: Let node v be good. Intuitively, good nodes have lots of low-degree neighbors, thus chances are high that one of them goes into the independent set, in which case v will be removed in step 3 of the algorithm.

If there is a neighbor $w \in N(v)$ with degree at most 2 we are done: With Lemma 8.4 the probability that node w joins the MIS is at least $\frac{1}{8}$, and our good node will be removed in Step 3.

So all we need to worry about is that all neighbors have at least degree 3: For any neighbor w of v we have $\frac{1}{2d(w)} \leq \frac{1}{6}$. Since $\sum_{w \in N(v)} \frac{1}{2d(w)} \geq \frac{1}{6}$ there is a

subset of neighbors $S \subseteq N(v)$ such that $\frac{1}{6} \leq \sum_{w \in S} \frac{1}{2d(w)} \leq \frac{1}{3}$

We can now bound the probability that node v will be removed. Let therefore R be the event of v being removed. Again, if a neighbor of v joins the MIS in Step 2, node v will be removed in Step 3. We have

$$\begin{aligned} P[R] &\geq P[\exists u \in S, u \in \text{MIS}] \\ &\geq \sum_{u \in S} P[u \in \text{MIS}] - \sum_{u, w \in S; u \neq w} P[u \in \text{MIS} \text{ and } w \in \text{MIS}]. \end{aligned}$$

For the last inequality we used the inclusion-exclusion principle truncated after the second order terms. Let M again be the set of marked nodes after Step 1. Using $P[u \in M] \geq P[u \in \text{MIS}]$ we get

$$\begin{aligned} P[R] &\geq \sum_{u \in S} P[u \in \text{MIS}] - \sum_{u, w \in S; u \neq w} P[u \in M \text{ and } w \in M] \\ &\geq \sum_{u \in S} P[u \in \text{MIS}] - \sum_{u \in S} \sum_{w \in S} P[u \in M] \cdot P[w \in M] \\ &\geq \sum_{u \in S} \frac{1}{4d(u)} - \sum_{u \in S} \sum_{w \in S} \frac{1}{2d(u)} \frac{1}{2d(w)} \\ &\geq \sum_{u \in S} \frac{1}{2d(u)} \left(\frac{1}{2} - \sum_{w \in S} \frac{1}{2d(w)} \right) \geq \frac{1}{6} \left(\frac{1}{2} - \frac{1}{3} \right) = \frac{1}{36}. \end{aligned}$$

□

Remarks:

- We would be almost finished if we could prove that many nodes are good in each phase. Unfortunately this is not the case: In a star-graph, for instance, only a single node is good! We need to find a work-around.

Lemma 8.6 (Good Edges). *An edge $e = (u, v)$ is called bad if both u and v are bad; else the edge is called good. The following holds: At any time at least half of the edges are good.*

Proof: For the proof we construct a directed auxiliary graph: Direct each edge towards the higher degree node (if both nodes have the same degree direct it towards the higher identifier). Now we need a little helper lemma before we can continue with the proof.

Lemma 8.7. *A bad node has outdegree at least twice its indegree.*

Proof: For the sake of contradiction, assume that a bad node v does not have outdegree at least twice its indegree. In other words, at least one third of the neighbor nodes (let's call them S) have degree at most $d(v)$. But then

$$\sum_{w \in N(v)} \frac{1}{2d(w)} \geq \sum_{w \in S} \frac{1}{2d(w)} \geq \sum_{w \in S} \frac{1}{2d(v)} \geq \frac{d(v)}{3} \frac{1}{2d(v)} = \frac{1}{6}$$

which means v is good, a contradiction. \square

Continuing the proof of Lemma 8.6: According to Lemma 8.7 the number of edges directed into bad nodes is at most half the number of edges directed out of bad nodes. Thus, the number of edges directed into bad nodes is at most half the number of edges. Thus, at least half of the edges are directed into good nodes. Since these edges are not bad, they must be good.

Theorem 8.8 (Analysis of Algorithm 35). *Algorithm 35 terminates in expected time $O(\log n)$.*

Proof: With Lemma 8.5 a good node (and therefore a good edge!) will be deleted with constant probability. Since at least half of the edges are good (Lemma 8.6) a constant fraction of edges will be deleted in each phase.

More formally: With Lemmas 8.5 and 8.6 we know that at least half of the edges will be removed with probability at least $1/36$. Let R be the number of edges to be removed. Using linearity of expectation we know that $\mathbb{E}[R] \geq m/72$, m being the total number of edges at the start of a phase. Now let $p := P[R \leq \mathbb{E}[R]/2]$. Bounding the expectation yields

$$\mathbb{E}[R] = \sum_r P[R = r] \cdot r \leq p \cdot \mathbb{E}[R]/2 + (1 - p) \cdot m.$$

Solving for p we get

$$p \leq \frac{m - \mathbb{E}[R]}{m - \mathbb{E}[R]/2} < \frac{m - \mathbb{E}[R]/2}{m} \leq 1 - 1/144.$$

In other words, with probability at least $1/144$ at least $m/144$ edges are removed in a phase. After expected $O(\log m)$ phases all edges are deleted. Since $m \leq n^2$ and thus $O(\log m) = O(\log n)$ the Theorem follows. \square

Remarks:

- With a bit of more math one can even show that Algorithm 35 terminates in time $O(\log n)$ “with high probability”.
- The presented algorithm is a simplified version of an algorithm by Michael Luby, published 1986 in the SIAM Journal of Computing. Around the same time there have been a number of other papers dealing with the same or related problems, for instance by Alon, Babai, and Itai, or by Israeli and Itai. The analysis presented here takes elements of all these papers, and from other papers on distributed weighted matching. The analysis in the book by David Peleg is different, and only achieves $O(\log^2 n)$ time.
- Though not as incredibly fast as the \log^* -coloring algorithm for trees, this algorithm is very general. It works on any graph, needs no identifiers, and can easily be made asynchronous.
- Surprisingly, much later, there have been half a dozen more papers published, with much worse results!! In 2002, for instance, there was a paper with linear running time, improving on a 1994 paper with cubic running time, restricted to trees!
- In 2009, Métivier, Robson, Saheb-Djahromi and Zemmari found a slightly different (and simpler) way to compute a MIS in the same logarithmic time:

8.3 Fast MIS from 2009

Algorithm 36 Fast MIS 2

The algorithm operates in synchronous rounds, grouped into phases.

A single phase is as follows:

- 1) Each node v chooses a random value $r(v) \in [0, 1]$ and sends it to its neighbors.
 - 2) If $r(v) < r(w)$ for all neighbors $w \in N(v)$, node v enters the MIS and informs its neighbors.
 - 3) If v or a neighbor of v entered the MIS, v terminates (v and all edges adjacent to v are removed from the graph), otherwise v enters the next phase.
-

Remarks:

- Correctness in the sense that the algorithm produces an independent set is simple: Steps 1 and 2 make sure that if a node v joins the MIS, then v 's neighbors do not join the MIS at the same time. Step 3 makes sure that v 's neighbors will never join the MIS.
- Likewise the algorithm eventually produces a MIS, because the node with the globally smallest value will always join the MIS, hence there is progress.
- So the only remaining question is how fast the algorithm terminates. To understand this, we need to dig a bit deeper.

- Our proof will rest on a simple, yet powerful observation about expected values of random variables that *may not be independent*:

Theorem 8.9 (Linearity of Expectation). *Let X_i , $i = 1, \dots, k$ denote random variables, then*

$$\mathbb{E} \left[\sum_i X_i \right] = \sum_i \mathbb{E}[X_i].$$

Proof. It is sufficient to prove $\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$ for two random variables X and Y , because then the statement follows by induction. Since

$$\begin{aligned} P[(X, Y) = (x, y)] &= P[X = x] \cdot P[Y = y | X = x] \\ &= P[Y = y] \cdot P[X = x | Y = y] \end{aligned}$$

we get that

$$\begin{aligned} \mathbb{E}[X + Y] &= \sum_{(X, Y) = (x, y)} P[(X, Y) = (x, y)] \cdot (x + y) \\ &= \sum_{X=x} \sum_{Y=y} P[X = x] \cdot P[Y = y | X = x] \cdot x \\ &\quad + \sum_{Y=y} \sum_{X=x} P[Y = y] \cdot P[X = x | Y = y] \cdot y \\ &= \sum_{X=x} P[X = x] \cdot x + \sum_{Y=y} P[Y = y] \cdot y \\ &= \mathbb{E}[X] + \mathbb{E}[Y]. \end{aligned}$$

Remarks:

- How can we prove that the algorithm only needs $O(\log n)$ phases in expectation? It would be great if this algorithm managed to remove a constant fraction of nodes in each phase. Unfortunately, it does not.
- Instead we will prove that the number of *edges* decreases quickly. Again, it would be great if any single edge was removed with constant probability in Step 3. But again, unfortunately, this is not the case.
- Maybe we can argue about the expected number of edges to be removed in one single phase? Let's see: A node v enters the MIS with probability $1/(d(v) + 1)$, where $d(v)$ is the degree of node v . By doing so, not only are v 's edges removed, but indeed all the edges of v 's neighbors as well – generally these are much more than $d(v)$ edges. So there is hope, but we need to be careful: If we do this the most naive way, we will count the same edge many times.
- How can we fix this? The nice observation is that it is enough to count just some of the removed edges. Given a new MIS node v and a neighbor $w \in N(v)$, we count the edges only if $r(v) < r(x)$ for all $x \in N(w)$. This looks promising. In a star graph, for instance, only the smallest random value can be accounted for removing all the edges of the star.

Lemma 8.10 (Edge Removal). *In a single phase, we remove at least half of the edges in expectation.*

Proof: To simplify the notation, at the start of our phase, the graph is simply $G = (V, E)$. Suppose that a node v joins the MIS in this phase, i.e., $r(v) < r(w)$ for all neighbors $w \in N(v)$. If in addition we have $r(v) < r(x)$ for all neighbors x of a neighbor w of v , we call this event $(v \rightarrow w)$. The probability of event $(v \rightarrow w)$ is at least $1/(d(v) + d(w))$, since $d(v) + d(w)$ is the maximum number of nodes adjacent to v or w (or both). As v joins the MIS, all edges (w, x) will be removed; there are $d(w)$ of these edges.

In order to count the removed edges, we need to weigh events properly.

Whether we remove the edges adjacent to w because of event $(v \rightarrow w)$ is a random variable $X_{(v \rightarrow w)}$. If event $(v \rightarrow w)$ occurs, $X_{(v \rightarrow w)}$ has the value $d(w)$, if not it has the value 0. For each edge $\{v, w\}$ we have two such variables, the event $X_{(v \rightarrow w)}$ and $X_{(w \rightarrow v)}$. Due to Theorem 8.9, the expected value of the sum X of all these random variables is at least

$$\begin{aligned} \mathbb{E}[X] &= \sum_{\{v,w\} \in E} \mathbb{E}[X_{(v \rightarrow w)}] + \mathbb{E}[X_{(w \rightarrow v)}] \\ &= \sum_{\{v,w\} \in E} P[\text{Event } (v \rightarrow w)] \cdot d(w) + P[\text{Event } (w \rightarrow v)] \cdot d(v) \\ &\geq \sum_{\{v,w\} \in E} \frac{d(w)}{d(v) + d(w)} + \frac{d(v)}{d(w) + d(v)} \\ &= \sum_{\{v,w\} \in E} 1 = |E|. \end{aligned}$$

In other words, in expectation all edges are removed in a single phase?!? Probably not. This means that we still counted some edges more than once. Indeed, for an edge $\{v, w\} \in E$ our random variable X includes the edge if the event $(u \rightarrow v)$ happens, but X also includes the edge if the event $(x \rightarrow w)$ happens. So we may have counted the edge $\{v, w\}$ twice. Fortunately however, not more than twice, because at most one event $(\cdot \rightarrow v)$ and at most one event $(\cdot \rightarrow w)$ can happen. If $(u \rightarrow v)$ happens, we know that $r(u) < r(w)$ for all $w \in N(v)$; hence another $(u' \rightarrow v)$ cannot happen because $r(u') > r(u) \in N(v)$. Therefore the random variable X must be divided by 2. In other words, in expectation at least half of the edges are removed.

Remarks:

- This enables us to follow a bound on the expected running time of Algorithm 36 quite easily.

Theorem 8.11 (Expected running time of Algorithm 36). *Algorithm 36 terminates after at most $3 \log_{4/3} m + 1 \in O(\log n)$ phases in expectation.*

Proof: The probability that in a single phase at least a quarter of all edges are removed is at least $1/3$. For the sake of contradiction, assume not. Then with probability less than $1/3$ we may be lucky and many (potentially all) edges are removed. With probability more than $2/3$ less than $1/4$ of the edges are removed. Hence the expected fraction of removed edges is strictly less than $1/3 \cdot 1 + 2/3 \cdot 1/4 = 1/2$. This contradicts Lemma 8.10.

Hence, at least every third phase is “good” and removes at least a quarter of the edges. To get rid of all but two edges we need $\log_{4/3} m$ good phases in expectation. The last two edges will certainly be removed in the next phase. Hence a total of $3 \log_{4/3} m + 1$ phases are enough in expectation.

Remarks:

- Sometimes one expects a bit more of an algorithm: Not only should the expected time to terminate be good, but the algorithm should *always* terminate quickly. As this is impossible in randomized algorithms (after all, the random choices may be “unlucky” all the time!), researchers often settle for a compromise, and just demand that the probability that the algorithm does not terminate in the specified time can be made absurdly small. For our algorithm, this can be deduced from Lemma 8.10 and another standard tool, namely Chernoff’s Bound.

Definition 8.12 (W.h.p.). *We say that an algorithm terminates w.h.p. (with high probability) within $O(t)$ time if it does so with probability at least $1 - 1/n^c$ for any choice of $c \geq 1$. Here c may affect the constants in the Big- O notation because it is considered a “tunable constant” and usually kept small.*

Definition 8.13 (Chernoff’s Bound). *Let $X = \sum_{i=1}^k X_i$ be the sum of k independent $0 - 1$ random variables. Then Chernoff’s bound states that w.h.p.*

$$|X - \mathbb{E}[X]| \in O\left(\log n + \sqrt{\mathbb{E}[X] \log n}\right).$$

Corollary 8.14 (Running Time of Algorithm 36). *Algorithm 36 terminates w.h.p. in $O(\log n)$ time.*

Proof: In Theorem 8.11 we used that *independently* of everything that happened before, in each phase we have a constant probability p that a quarter of the edges are removed. Call such a phase *good*. For some constants C_1 and C_2 , let us check after $C_1 \log n + C_2 \in O(\log n)$ phases, in how many phases at least a quarter of the edges have been removed. In expectation, these are at least $p(C_1 \log n + C_2)$ many. Now we look at the random variable $X = \sum_{i=1}^{C_1 \log n + C_2} X_i$, where the X_i are independent $0 - 1$ variables being one with exactly probability p . Certainly, if X is at least x with some probability, then the probability that we have x good phases can only be larger (if no edges are left, certainly “all” of the remaining edges are removed). To X we can apply Chernoff’s bound. If C_1 and C_2 are chosen large enough, they will overcome the constants in the Big- O from Chernoff’s bound, i.e., w.h.p. it holds that $|X - \mathbb{E}[X]| \leq \mathbb{E}[X]/2$, implying $X \geq \mathbb{E}[X]/2$. Choosing C_1 large enough, we will have w.h.p. sufficiently many good phases, i.e., the algorithm terminates w.h.p. in $O(\log n)$ phases.

Remarks:

- The algorithm can be improved a bit more even. Drawing random real numbers in each phase for instance is not necessary. One can achieve the same by sending only a total of $O(\log n)$ random (and as many non-random) bits over each edge.
- One of the main open problems in distributed computing is whether one can beat this logarithmic time, or at least achieve it with a deterministic algorithm.

- Let's turn our attention to applications of MIS next.

8.4 Applications

Definition 8.15 (Matching). *Given a graph $G = (V, E)$ a matching is a subset of edges $M \subseteq E$, such that no two edges in M are adjacent (i.e., where no node is adjacent to two edges in the matching). A matching is maximal if no edge can be added without violating the above constraint. A matching of maximum cardinality is called maximum. A matching is called perfect if each node is adjacent to an edge in the matching.*

Remarks:

- In contrast to MaxIS, a maximum matching can be found in polynomial time (Blossom algorithm by Jack Edmonds), and is also easy to approximate (in fact, already any maximal matching is a 2-approximation).
- An independent set algorithm is also a matching algorithm: Let $G = (V, E)$ be the graph for which we want to construct the matching. The auxiliary graph G' is defined as follows: for every edge in G there is a node in G' ; two nodes in G' are connected by an edge if their respective edges in G are adjacent. A (maximal) independent set in G' is a (maximal) matching in G , and vice versa. Using Algorithm 36 directly produces a $O(\log n)$ bound for maximal matching.
- More importantly, our MIS algorithm can also be used for vertex coloring (Problem 2.1):

Algorithm 37 General Graph Coloring

- 1: Given a graph $G = (V, E)$ we virtually build a graph $G' = (V', E')$ as follows:
 - 2: Every node $v \in V$ clones itself $d(v) + 1$ times ($v_0, \dots, v_{d(v)} \in V'$), $d(v)$ being the degree of v in G .
 - 3: The edge set E' of G' is as follows:
 - 4: First all clones are in a clique: $(v_i, v_j) \in E'$, for all $v \in V$ and all $0 \leq i < j \leq d(v)$
 - 5: Second all i^{th} clones of neighbors in the original graph G are connected: $(u_i, v_i) \in E'$, for all $(u, v) \in E$ and all $0 \leq i \leq \min(d(u), d(v))$.
 - 6: Now we simply run (simulate) the fast MIS Algorithm 36 on G' .
 - 7: If node v_i is in the MIS in G' , then node v gets color i .
-

Theorem 8.16 (Analysis of Algorithm 37). *Algorithm 37 $(\Delta + 1)$ -colors an arbitrary graph in $O(\log n)$ time, with high probability, Δ being the largest degree in the graph.*

Proof: Thanks to the clique among the clones at most one clone is in the MIS. And because of the $d(v) + 1$ clones of node v every node will get a free color! The running time remains logarithmic since G' has $O(n^2)$ nodes and the exponent becomes a constant factor when applying the logarithm.

Remarks:

- This solves our open problem from Chapter 2.1!
- Together with Corollary 8.3 we get quite close ties between $(\Delta+1)$ -coloring and the MIS problem.
- However, in general Algorithm 37 is not the best distributed algorithm for $O(\Delta)$ -coloring. For fast distributed vertex coloring please check Kothapalli, Onus, Scheideler, Schindelhauer, IPDPS 2006. This algorithm is based on a $O(\log \log n)$ time *edge* coloring algorithm by Grable and Panconesi, 1997.
- Computing a MIS also solves another graph problem on graphs of bounded independence.

Definition 8.17 (Bounded Independence). $G = (V, E)$ is of bounded independence, if each neighborhood contains at most a constant number of independent (i.e., mutually non-adjacent) nodes.

Definition 8.18 ((Minimum) Dominating Sets). A dominating set is a subset of the nodes such that each node is in the set or adjacent to a node in the set. A minimum dominating set is a dominating set containing the least possible number of nodes.

Remarks:

- In general, finding a dominating set less than factor $\log n$ larger than a minimum dominating set is NP-hard.
- Any MIS is a dominating set: if a node was not covered, it could join the independent set.
- In general a MIS and a minimum dominating sets have not much in common (think of a star). For graphs of bounded independence, this is different.

Corollary 8.19. On graphs of bounded independence, a constant-factor approximation to a minimum dominating set can be found in time $O(\log n)$ w.h.p.

Proof: Denote by M a minimum dominating set and by I a MIS. Since M is a dominating set, each node from I is in M or adjacent to a node in M . Since the graph is of bounded independence, no node in M is adjacent to more than constantly many nodes from I . Thus, $|I| \in O(|M|)$. Therefore, we can compute a MIS with Algorithm 36 and output it as the dominating set, which takes $O(\log n)$ rounds w.h.p.

Chapter 9

Locality Lower Bounds

In Chapter 2, we looked at distributed algorithms for coloring. In particular, we saw that rings and rooted trees can be colored with 3 colors in $\log^* n + O(1)$ rounds. In this chapter, we will reconsider the distributed coloring problem. We will look at a classic lower bound by Nathan Linial that shows that the result of Chapter 2 is tight: Coloring rings (and rooted trees) indeed requires $\Omega(\log^* n)$ rounds. In particular, we will prove a lower bound for coloring in the following setting:

- We consider deterministic, synchronous algorithms.
- Message size and local computations are unbounded.
- We assume that the network is a directed ring with n nodes.
- Nodes have unique labels (identifiers) from 1 to n .

Remarks:

- A generalization of the lower bound to randomized algorithms is possible. Unfortunately, we will however not have time to discuss this.
- Except for restricting to deterministic algorithms, all the conditions above make a lower bound stronger. Any lower bound for synchronous algorithms certainly also holds for asynchronous ones. A lower bound that is true if message size and local computations are not restricted is clearly also valid if we require a bound on the maximal message size or the amount of local computations. Similarly also assuming that the ring is directed and that node labels are from 1 to n (instead of choosing IDs from a more general domain) strengthen the lower bound.
- Instead of directly proving that 3-coloring a ring needs $\Omega(\log^* n)$ rounds, we will prove a slightly more general statement. We will consider deterministic algorithms with time complexity r (for arbitrary r) and derive a lower bound on the number of colors that are needed if we want to properly color an n -node ring with an r -round algorithm. A 3-coloring lower bound can then be derived by taking the smallest r for which an r -round algorithm needs 3 or fewer colors.

Algorithm 38 Synchronous Algorithm: Canonical Form

-
- 1: In r rounds: **send** complete initial state to nodes at distance at most r
 - 2: *// do all the communication first*
 - 3: Compute output based on complete information about r -neighborhood
 - 4: *// do all the computation in the end*
-

9.1 Locality

Let us for a moment look at distributed algorithms more generally (i.e., not only at coloring and not only at rings). Assume that initially, all nodes only know their own label (identifier) and potentially some additional input. As information needs at least r rounds to travel r hops, after r rounds, a node v can only learn about other nodes at distance at most r . If message size and local computations are not restricted, it is in fact not hard to see, that in r rounds, a node v can exactly learn all the node labels and inputs up to distance r . As shown by the following lemma, this allows to transform every deterministic r -round synchronous algorithm into a simple canonical form.

Lemma 9.1. *If message size and local computations are not bounded, every deterministic, synchronous r -round algorithm can be transformed into an algorithm of the form given by Algorithm 38 (i.e., it is possible to first communicate for r rounds and then do all the computations in the end).*

Proof. Consider some r -round algorithm \mathcal{A} . We want to show that \mathcal{A} can be brought to the canonical form given by Algorithm 38. First, we let the nodes communicate for r rounds. Assume that in every round, every node sends its complete state to all of its neighbors (remember that there is no restriction on the maximal message size). By induction, after r rounds, every node knows the initial state of all other nodes at distance at most i . Hence, after r rounds, a node v has the combined initial knowledge of all the nodes in its r -neighborhood. We want to show that this suffices to locally (at node v) simulate enough of Algorithm \mathcal{A} to compute all the messages that v receives in the r communication rounds of a regular execution of Algorithm \mathcal{A} .

Concretely, we prove the following statement by induction on i . For all nodes at distance at most $r - i + 1$ from v , node v can compute all messages of the first i rounds of a regular execution of \mathcal{A} . Note that this implies that v can compute all the messages it receives from its neighbors during all r rounds. Because v knows the initial state of all nodes in the r -neighborhood, v can clearly compute all messages of the first round (i.e., the statement is true for $i = 1$). Let us now consider the induction step from i to $i + 1$. By the induction hypothesis, v can compute the messages of the first i rounds of all nodes in its $(r - i + 1)$ -neighborhood. It can therefore compute all messages that are received by nodes in the $(r - i)$ -neighborhood in the first i rounds. This is of course exactly what is needed to compute the messages of round $i + 1$ of nodes in the $(r - i)$ -neighborhood. \square

Remarks:

- It is straightforward to generalize the canonical form to randomized algorithms: Every node first computes all the random bits it needs throughout the algorithm. The random bits are then part of the initial state of a node.

Definition 9.2 (*r*-hop view). *We call the collection of the initial states of all nodes in the *r*-neighborhood of a node *v*, the *r*-hop view of *v*.*

Remarks:

- Assume that initially, every node knows its degree, its label (identifier) and potentially some additional input. The *r*-hop view of a node *v* then includes the complete topology of the *r*-neighborhood (excluding edges between nodes at distance *r*) and the labels and additional inputs of all nodes in the *r*-neighborhood.

Based on the definition of an *r*-hop view, we can state the following corollary of Lemma 9.1.

Corollary 9.3. *A deterministic *r*-round algorithm \mathcal{A} is a function that maps every possible *r*-hop view to the set of possible outputs.*

Proof. By Lemma 9.1, we know that we can transform Algorithm \mathcal{A} to the canonical form given by Algorithm 38. After *r* communication rounds, every node *v* knows exactly its *r*-hop view. This information suffices to compute the output of node *v*. \square

Remarks:

- Note that the above corollary implies that two nodes with equal *r*-hop views have to compute the same output in every *r*-round algorithm.
- For coloring algorithms, the only input of a node *v* is its label. The *r*-hop view of a node therefore is its labeled *r*-neighborhood.
- Since we only consider rings, *r*-hop neighborhoods are particularly simple. The labeled *r*-neighborhood of a node *v* (and hence its *r*-hop view) in a directed ring is simply a $(2r + 1)$ -tuple $(\ell_{-r}, \ell_{-r+1}, \dots, \ell_0, \dots, \ell_r)$ of distinct node labels where ℓ_0 is the label of *v*. Assume that for $i > 0$, ℓ_i is the label of the i^{th} clockwise neighbor of *v* and ℓ_{-i} is the label of the i^{th} counterclockwise neighbor of *v*. A deterministic coloring algorithm for directed rings therefore is a function that maps $(2r + 1)$ -tuples of node labels to colors.
- Consider two *r*-hop views $\mathcal{V}_r = (\ell_{-r}, \dots, \ell_r)$ and $\mathcal{V}'_r = (\ell'_{-r}, \dots, \ell'_r)$. If $\ell'_i = \ell_{i+1}$ for $-r \leq i \leq r-1$ and if $\ell'_r \neq \ell_r$ for $-r \leq i \leq r$, the *r*-hop view \mathcal{V}'_r can be the *r*-hop view of a clockwise neighbor of a node with *r*-hop view \mathcal{V}_r . Therefore, every algorithm \mathcal{A} that computes a valid coloring needs to assign different colors to \mathcal{V}_r and \mathcal{V}'_r . Otherwise, there is a ring labeling for which \mathcal{A} assigns the same color to two adjacent nodes.

9.2 The Neighborhood Graph

We will now make the above observations concerning colorings of rings a bit more formal. Instead of thinking of an r -round coloring algorithm as a function from all possible r -hop views to colors, we will use a slightly different perspective. Interestingly, the problem of understanding distributed coloring algorithms can itself be seen as a classical graph coloring problem.

Definition 9.4 (Neighborhood Graph). *For a given family of network graphs \mathcal{G} , the r -neighborhood graph $\mathcal{N}_r(\mathcal{G})$ is defined as follows. The node set of $\mathcal{N}_r(\mathcal{G})$ is the set of all possible labeled r -neighborhoods (i.e., all possible r -hop views). There is an edge between two labeled r -neighborhoods \mathcal{V}_r and \mathcal{V}'_r if \mathcal{V}_r and \mathcal{V}'_r can be the r -hop views of two adjacent nodes.*

Lemma 9.5. *For a given family of network graphs \mathcal{G} , there is an r -round algorithm that colors graphs of \mathcal{G} with c colors iff the chromatic number of the neighborhood graph is $\chi(\mathcal{N}_r(\mathcal{G})) \leq c$.*

Proof. We have seen that a coloring algorithm is a function that maps every possible r -hop view to a color. Hence, a coloring algorithm assigns a color to every node of the neighborhood graph $\mathcal{N}_r(\mathcal{G})$. If two r -hop views \mathcal{V}_r and \mathcal{V}'_r can be the r -hop views of two adjacent nodes u and v (for some labeled graph in \mathcal{G}), every correct coloring algorithm must assign different colors to \mathcal{V}_r and \mathcal{V}'_r . Thus, specifying an r -round coloring algorithm for a family of network graphs \mathcal{G} is equivalent to coloring the respective neighborhood graph $\mathcal{N}_r(\mathcal{G})$. \square

Remarks:

- If an algorithm is non-uniform, i.e., the nodes know n , we can see this as having different neighborhood graphs for different values of n (as opposed to a disconnected neighborhood graph).
- This does not make much of a difference for coloring algorithms on the ring, as we are interested in neighborhoods that are much smaller than n .

Instead of directly defining the neighborhood graph for directed rings, we define directed graphs $\mathcal{B}_{k,n}$ that are closely related to the neighborhood graph. Let k and n be two positive integers and assume that $n \geq k$. The node set of $\mathcal{B}_{k,n}$ contains all k -tuples of increasing node labels ($[n] = \{1, \dots, n\}$):

$$V[\mathcal{B}_{k,n}] = \{(\alpha_1, \dots, \alpha_k) : \alpha_i \in [n], i < j \rightarrow \alpha_i < \alpha_j\} \quad (9.1)$$

For $\underline{\alpha} = (\alpha_1, \dots, \alpha_k)$ and $\underline{\beta} = (\beta_1, \dots, \beta_k)$ there is a directed edge from $\underline{\alpha}$ to $\underline{\beta}$ iff

$$\forall i \in \{1, \dots, k-1\} : \beta_i = \alpha_{i+1}. \quad (9.2)$$

Lemma 9.6. *Viewed as an undirected graph, the graph $\mathcal{B}_{2r+1,n}$ is a subgraph of the r -neighborhood graph of directed n -node rings with node labels from $[n]$.*

Proof. The claim follows directly from the observations regarding r -hop views of nodes in a directed ring from Section 9.1. The set of k -tuples of increasing node labels is a subset of the set of k -tuples of distinct node labels. Two nodes of $\mathcal{B}_{2r+1,n}$ are connected by a directed edge iff the two corresponding r -hop views are connected by a directed edge in the neighborhood graph. Note that if there

is an edge between $\underline{\alpha}$ and $\underline{\beta}$ in $\mathcal{B}_{k,n}$, $\alpha_1 \neq \beta_k$ because the node labels in $\underline{\alpha}$ and $\underline{\beta}$ are increasing. \square

To determine a lower bound on the number of colors an r -round algorithm needs for directed n -node rings, it therefore suffices to determine a lower bound on the chromatic number of $\mathcal{B}_{2r+1,n}$. To obtain such a lower bound, we need the following definition.

Definition 9.7 (Diline Graph). *The directed line graph (diline graph) $\mathcal{DL}(G)$ of a directed graph $G = (V, E)$ is defined as follows. The node set of $\mathcal{DL}(G)$ is $V[\mathcal{DL}(G)] = E$. There is a directed edge $((w, x), (y, z))$ between $(w, x) \in E$ and $(y, z) \in E$ iff $x = y$, i.e., if the first edge ends where the second one starts.*

Lemma 9.8. *If $n > k$, the graph $\mathcal{B}_{k+1,n}$ can be defined recursively as follows:*

$$\mathcal{B}_{k+1,n} = \mathcal{DL}(\mathcal{B}_{k,n}).$$

Proof. The edges of $\mathcal{B}_{k,n}$ are pairs of k -tuples $\underline{\alpha} = (\alpha_1, \dots, \alpha_k)$ and $\underline{\beta} = (\beta_1, \dots, \beta_k)$ that satisfy Conditions (9.1) and (9.2). Because the last $k-1$ labels in $\underline{\alpha}$ are equal to the first $k-1$ labels in $\underline{\beta}$, the pair $(\underline{\alpha}, \underline{\beta})$ can be represented by a $(k+1)$ -tuple $\underline{\gamma} = (\gamma_1, \dots, \gamma_{k+1})$ with $\gamma_1 = \alpha_1$, $\gamma_i = \beta_{i-1} = \alpha_i$ for $2 \leq i \leq k$, and $\gamma_{k+1} = \beta_k$. Because the labels in $\underline{\alpha}$ and the labels in $\underline{\beta}$ are increasing, the labels in $\underline{\gamma}$ are increasing as well. The two graphs $\mathcal{B}_{k+1,n}$ and $\mathcal{DL}(\mathcal{B}_{k,n})$ therefore have the same node sets. There is an edge between two nodes $(\underline{\alpha}_1, \underline{\beta}_1)$ and $(\underline{\alpha}_2, \underline{\beta}_2)$ of $\mathcal{DL}(\mathcal{B}_{k,n})$ if $\underline{\beta}_1 = \underline{\alpha}_2$. This is equivalent to requiring that the two corresponding $(k+1)$ -tuples $\underline{\gamma}_1$ and $\underline{\gamma}_2$ are neighbors in $\mathcal{B}_{k+1,n}$, i.e., that the last k labels of $\underline{\gamma}_1$ are equal to the first k labels of $\underline{\gamma}_2$. \square

The following lemma establishes a useful connection between the chromatic numbers of a directed graph G and its diline graph $\mathcal{DL}(G)$.

Lemma 9.9. *For the chromatic numbers $\chi(G)$ and $\chi(\mathcal{DL}(G))$ of a directed graph G and its diline graph, it holds that*

$$\chi(\mathcal{DL}(G)) \geq \log_2(\chi(G)).$$

Proof. Given a c -coloring of $\mathcal{DL}(G)$, we show how to construct a 2^c coloring of G . The claim of the lemma then follows because this implies that $\chi(G) \leq 2^{\chi(\mathcal{DL}(G))}$.

Assume that we are given a c -coloring of $\mathcal{DL}(G)$. A c -coloring of the diline graph $\mathcal{DL}(G)$ can be seen as a coloring of the edges of G such that no two adjacent edges have the same color. For a node v of G , let S_v be the set of colors of its outgoing edges. Let u and v be two nodes such that G contains a directed edge (u, v) from u to v and let x be the color of (u, v) . Clearly, $x \in S_u$ because (u, v) is an outgoing edge of u . Because adjacent edges have different colors, no outgoing edge (v, w) of v can have color x . Therefore $x \notin S_v$. This implies that $S_u \neq S_v$. We can therefore use these color sets to obtain a vertex coloring of G , i.e., the color of u is S_u and the color of v is S_v . Because the number of possible subsets of $[c]$ is 2^c , this yields a 2^c -coloring of G . \square

Let $\log^{(i)} x$ be the i -fold application of the base-2 logarithm to x :

$$\log^{(1)} x = \log_2 x, \quad \log^{(i+1)} x = \log_2(\log^{(i)} x).$$

Remember from Chapter 2 that

$$\log^* x = 1 \text{ if } x \leq 2, \quad \log^* x = 1 + \min\{i : \log^{(i)} x \leq 2\}.$$

For the chromatic number of $\mathcal{B}_{k,n}$, we obtain

Lemma 9.10. *For all $n \geq 1$, $\chi(\mathcal{B}_{1,n}) = n$. Further, for $n \geq k \geq 2$, $\chi(\mathcal{B}_{k,n}) \geq \log^{(k-1)} n$.*

Proof. For $k = 1$, $\mathcal{B}_{k,n}$ is the complete graph on n nodes with a directed edge from node i to node j iff $i < j$. Therefore, $\chi(\mathcal{B}_{1,n}) = n$. For $k > 2$, the claim follows by induction and Lemmas 9.8 and 9.9. \square

This finally allows us to state a lower bound on the number of rounds needed to color a directed ring with 3 colors.

Theorem 9.11. *Every deterministic, distributed algorithm to color a directed ring with 3 or less colors needs at least $(\log^* n)/2 - 1$ rounds.*

Proof. Using the connection between $\mathcal{B}_{k,n}$ and the neighborhood graph for directed rings, it suffices to show that $\chi(\mathcal{B}_{2r+1,n}) > 3$ for all $r < (\log^* n)/2 - 1$. From Lemma 9.10, we know that $\chi(\mathcal{B}_{2r+1,n}) \geq \log^{(2r)} n$. To obtain $\log^{(2r)} n \leq 2$, we need $r \geq (\log^* n)/2 - 1$. Because $\log_2 3 < 2$, we therefore have $\log^{(2r)} n > 3$ if $r < \log^* n/2 - 1$. \square

Corollary 9.12. *Every deterministic, distributed algorithm to compute an MIS of a directed ring needs at least $\log^* n/2 - O(1)$ rounds.*

Remarks:

- It is straightforward to see that also for a constant $c > 3$, the number of rounds needed to color a ring with c or less colors is $\log^* n/2 - O(1)$.
- There basically (up to additive constants) is a gap of a factor of 2 between the $\log^* n + O(1)$ upper bound of Chapter 2 and the $\log^* n/2 - O(1)$ lower bound of this chapter. It is possible to show that the lower bound is tight, even for undirected rings (for directed rings, this will be part of the exercises).
- The presented lower bound is due to Nathan Linial. The lower bound is also true for randomized algorithms. The generalization for randomized algorithms was done by Moni Naor.
- Alternatively, the lower bound can also be presented as an application of Ramsey's theory. Ramsey's theory is best introduced with an example: Assume you host a party, and you want to invite people such that there are no three people who mutually know each other, and no three people which are mutual strangers. How many people can you invite? This is an example of Ramsey's theorem, which says that for any given integer c , and any given integers n_1, \dots, n_c , there is a Ramsey number $R(n_1, \dots, n_c)$, such that if the edges of a complete graph with $R(n_1, \dots, n_c)$ nodes are colored with c different colors, then for some color i the graph contains some complete subgraph of color i of size n_i . The special case in the party example is looking for $R(3, 3)$.

- Ramsey theory is more general, as it deals with hyperedges. A normal edge is essentially a subset of two nodes; a hyperedge is a subset of k nodes. The party example can be explained in this context: We have (hyper)edges of the form $\{i, j\}$, with $1 \leq i, j \leq n$. Choosing n sufficiently large, coloring the edges with two colors must exhibit a set S of 3 edges $\{i, j\} \subset \{v_1, v_2, v_3\}$, such that all edges in S have the same color. To prove our coloring lower bound using Ramsey theory, we form all hyperedges of size $k = 2r + 1$, and color them with 3 colors. Choosing n sufficiently large, there must be a set $S = \{v_1, \dots, v_{k+1}\}$ of $k + 1$ identifiers, such that all $k + 1$ hyperedges consisting of k nodes from S have the same color. Note that both $\{v_1, \dots, v_k\}$ and $\{v_2, \dots, v_{k+1}\}$ are in the set S , hence there will be two neighboring views with the same color. Ramsey theory shows that in this case n will grow as a power tower (tetration) in k . Thus, if n is so large that k is smaller than some function growing like $\log^* n$, the coloring algorithm cannot be correct.
- The neighborhood graph concept can be used more generally to study distributed graph coloring. It can for instance be used to show that with a single round (every node sends its identifier to all neighbors) it is possible to color a graph with $(1 + o(1))\Delta^2 \ln n$ colors, and that every one-round algorithm needs at least $\Omega(\Delta^2 / \log^2 \Delta + \log \log n)$ colors.
- One may also extend the proof to other problems, for instance one may show that a constant approximation of the minimum dominating set problem on unit disk graphs costs at least log-star time.
- Using r -hop views and the fact that nodes with equal r -hop views have to make the same decisions is the basic principle behind almost all locality lower bounds (in fact, we are not aware of a locality lower bound that does not use this principle). Using this basic technique (but a completely different proof otherwise), it is for instance possible to show that computing an MIS (and many other problems) in a general graph requires at least $\Omega(\sqrt{\log n / \log \log n})$ and $\Omega(\log \Delta / \log \log \Delta)$ rounds.

Chapter 10

Social Networks

Distributed computing is applicable in various contexts. This lecture exemplarily studies one of these contexts, social networks, an area of study whose origins date back a century. To give you a first impression, consider Figure 10.1.

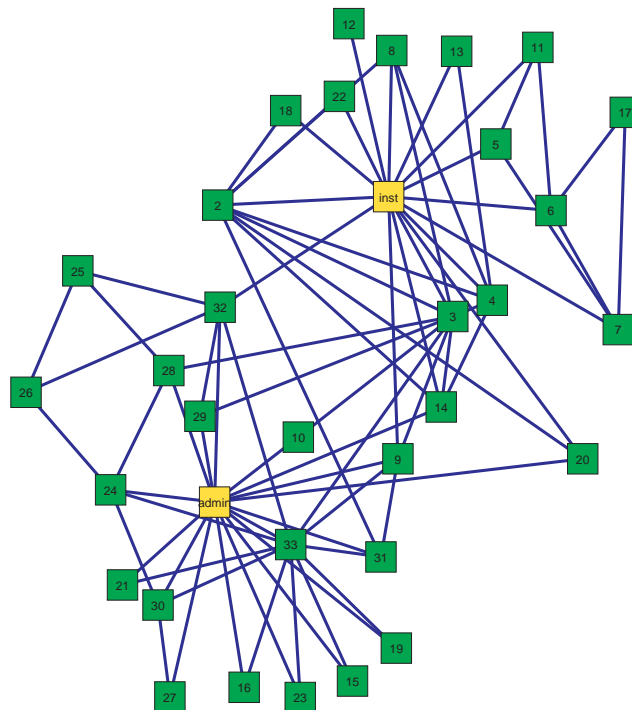


Figure 10.1: This graph shows the social relations between the members of a karate club, studied by anthropologist Wayne Zachary in the 1970s. Two people (nodes) stand out, the instructor and the administrator of the club, both happen to have many friends among club members. At some point, a dispute caused the club to split into two. Can you predict how the club partitioned? (If not, just search the Internet for Zachary and Karate.)

10.1 Small-World Networks

Back in 1929, Frigyes Karinthy published a volume of short stories that postulated that the world was “shrinking” because human beings were connected more and more. Some claim that he was inspired by radio network pioneer Guglielmo Marconi’s 1909 Nobel Prize speech. Despite physical distance, the growing density of human “networks” renders the actual social distance smaller and smaller. As a result, it is believed that any two individuals can be connected through at most five (or so) acquaintances, i.e., within six hops.

- The topic was hot in the 1960s. For instance, in 1964, Marshall McLuhan coined the metaphor “Global Village”. He wrote: “As electrically contracted, the globe is no more than a village”. He argues that due to the almost instantaneous reaction times of new (“electric”) technologies, each individual inevitably feels the consequences of his actions and thus automatically deeply participates in the global society. McLuhan understood what we now can directly observe – real and virtual world are moving together. He realized that the transmission medium, rather than the transmitted information is at the core of change, as expressed by his famous phrase “the medium is the message”.
- This idea has been followed ardently in the 1960s by several sociologists, first by Michael Gurevich, later by Stanley Milgram. Milgram wanted to know the average path length between two “random” humans, by using various experiments, generally using randomly chosen individuals from the US Midwest as starting points, and a stockbroker living in a suburb of Boston as target. The starting points were given name, address, occupation, plus some personal information about the target. They were asked to send a letter to the target. However, they were not allowed to *directly* send the letter, rather, they had to pass it to somebody they knew on first-name basis and that they thought to have a higher probability to know the target person. This process was repeated, until somebody knew the target person, and could deliver the letter. Shortly after starting the experiment, letters have been received. Most letters were lost during the process, but if they arrived, the average path length was about 5.5. The observation that the entire population is connected by short acquaintance chains got later popularized by the terms “six degrees of separation” and “small world”.
- Statisticians tried to explain Milgram’s experiments, by essentially giving network models that allowed for short diameters, i.e., each node is connected to each other node by only a few hops. Until today there is a thriving research community in statistical physics that tries to understand network properties that allow for “small world” effects.
- One of the keywords in this area are power-law graphs, networks where node degrees are distributed according to a power-law distribution, i.e. the number of nodes with degree δ is proportional to $\delta^{-\alpha}$, for some $\alpha > 1$. Such power-law graphs have been witnessed in many application areas, apart from social networks also in the web, or in Biology or Physics.

- Obviously, two power-law graphs might look and behave completely differently, even if α and the number of edges is exactly the same.

One well-known model towards this end is the Watts-Strogatz model. Watts and Strogatz argued that social networks should be modeled by a combination of two networks: As the basis we take a network that has a large cluster coefficient ...

Definition 10.1. *The cluster coefficient of a network is defined by the probability that two friends of a node are likely to be friends as well, summing up over all the nodes.*

..., then we augment such a graph with random links, every node for instance points to a constant number of other nodes, chosen uniformly at random. This augmentation represents acquaintances that connect nodes to parts of the network that would otherwise be far away.

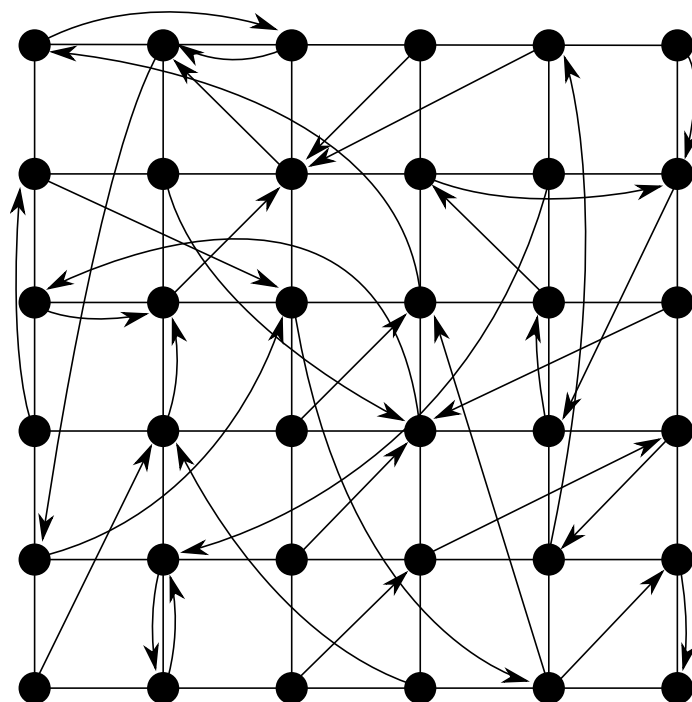
Remarks:

- Without further information, knowing the cluster coefficient is of questionable value: Assume we arrange the nodes in a grid. Technically, if we connect each node to its four closest neighbors, the graph has cluster coefficient 0, since there are no triangles; if we instead connect each node with its eight closest neighbors, the cluster coefficient is $3/7$. The cluster coefficient is quite different, even though both networks have similar characteristics.

This is interesting, but not enough to really understand what is going on. For Milgram's experiments to work, it is not sufficient to connect the nodes in a certain way. In addition, the nodes *themselves* need to know how to forward a message to one of their neighbors, even though they cannot know whether that neighbor is really closer to the target. In other words, nodes are not just following physical laws, but they make decisions themselves. In contrast to those mathematicians that worked on the problem earlier, Jon Kleinberg understood that Milgram's experiment essentially shows that social networks are "navigable", and that one can only explain it in terms of a *greedy routing*.

In particular, Kleinberg set up an artificial network with nodes on a grid topology, plus some additional random links per node. In a quantitative study he showed that the random links need a specific distance distribution to allow for efficient greedy routing. This distribution marks the sweet spot for any navigable network.

Definition 10.2 (Augmented Grid). *We take $n = m^2$ nodes $(i, j) \in V = \{1, \dots, m\}^2$ that are identified with the lattice points on an $m \times m$ grid. We define the distance between two nodes (i, j) and (k, ℓ) as $d((i, j), (k, \ell)) = |k - i| + |\ell - j|$ as the distance between them on the $m \times m$ lattice. The network is modeled using a parameter $\alpha \geq 0$. Each node u has a directed edge to every lattice neighbor. These are the local contacts of a node. In addition, each node also has an additional random link (the long-range contact). For all u and v , the long-range contact of u points to node v with probability proportional to $d(u, v)^{-\alpha}$, i.e., with probability $d(u, v)^{-\alpha} / \sum_{w \in V \setminus \{u\}} d(u, w)^{-\alpha}$. Figure 10.2 illustrates the model.*

Figure 10.2: Augmented grid with $m = 6$ **Remarks:**

- The network model has the following geographic interpretation: nodes (individuals) live on a grid and know their neighbors on the grid. Further, each node has some additional acquaintances throughout the network.
- The parameter α controls how the additional neighbors are distributed across the grid. If $\alpha = 0$, long-range contacts are chosen uniformly at random (as in the Watts-Strogatz model). As α increases, long-range contacts become shorter on average. In the extreme case, if $\alpha \rightarrow \infty$, all long-range contacts are to immediate neighbors on the grid.
- It can be shown that as long as $\alpha \leq 2$, the diameter of the resulting graph is polylogarithmic in n (polynomial in $\log n$) with high probability. In particular, if the long-range contacts are chosen uniformly at random ($\alpha = 0$), the diameter is $O(\log n)$.

Since the augmented grid contains random links, we do not know anything for sure about how the random links are distributed. In theory, all links could point to the same node! However, this is almost certainly not the case. Formally this is captured by the term *with high probability*.

Definition 10.3 (With High Probability). *Some probabilistic event is said to occur with high probability (w.h.p.), if it happens with a probability $p \geq 1 - 1/n^c$, where c is a constant. The constant c may be chosen arbitrarily, but it is considered constant with respect to Big-O notation.*

Remarks:

- For instance, a running time bound of $c \log n$ or $e^{cl} \log n + 5000c$ with probability at least $1 - 1/n^c$ would be $O(\log n)$ w.h.p., but a running time of n^c would not be $O(n)$ w.h.p. since c might also be 50.
- This definition is very powerful, as any polynomial (in n) number of statements that hold w.h.p. also holds w.h.p. at the same time, regardless of any dependencies between random variables!

Theorem 10.4. *The diameter of the augmented grid with $\alpha = 0$ is $O(\log n)$ with high probability.*

Proof Sketch. For simplicity, we will only show that we can reach a node w starting from some node v . However, it can be shown that (essentially) each of the intermediate claims holds with high probability, which then by means of the union bound yields that *all* of the claims hold simultaneously with high probability for *all* pairs of nodes.

Let N_g be the $\lceil \log n \rceil$ -hop neighborhood of v on the grid, containing $\Omega(\log^2 n)$ nodes. Each of the nodes in N_g has a random link, probably leading to distant parts of the graph. As long as we have reached only $o(n)$ nodes, any new random link will with probability $1 - o(1)$ lead to a node for which none of its grid neighbors has been visited yet. Thus, in expectation we find almost $|N_g|$ new nodes whose neighbors are “fresh”. Using their grid links, we will reach $(4 - o(1))|N_g|$ more nodes within one more hop. If bad luck strikes, it could still happen that many of these links lead to a few nodes, already visited nodes, or nodes that are very close to each other. But that is very unlikely, as we have lots of random choices! Indeed, it can be shown that not only in expectation, but with high probability $(5 - o(1))|N_g|$ many nodes are reached this way.

Because all these shiny new nodes have (so far unused) random links, we can repeat this reasoning inductively, implying that the number of nodes grows by (at least) a constant factor for every two hops. Thus, after $O(\log n)$ hops, we will have reached $n/\log n$ nodes (which is still small compared to n). Finally, consider the expected number of links from these nodes that enter the $(\log n)$ -neighborhood of some target node w with respect to the grid. Since this neighborhood consists of $\Omega(\log^2 n)$ nodes, in expectation $\Omega(\log n)$ links come close enough to w . This is large enough to almost guarantee that this happens. Summing everything up, we still used merely $O(\log n)$ hops in total to get from v to w .

□

This shows that for $\alpha = 0$ (and in fact for all $\alpha \leq 2$), the resulting network has a small diameter. Recall however that we also wanted the network to be navigable. For this, we consider a simple greedy routing strategy (Algorithm 39).

Algorithm 39 Greedy Routing

```

1: while not at destination do
2:   go to a neighbor which is closest to destination (considering grid distance
   only)
3: end while

```

Lemma 10.5. *In the augmented grid, Algorithm 39 finds a routing path of length at most $2(m-1) \in O(\sqrt{n})$.*

Proof. Because of the grid links, there is always a neighbor which is closer to the destination. Since with each hop we reduce the distance to the target at least by one in one of the two grid dimensions, we will reach the destination within $2(m-1)$ steps. \square

This is not really what Milgram's experiment promises. We want to know how much the additional random links speed up the process. To this end, we first need to understand how likely it is that two nodes u and v are connected by a random link in terms of n and their distance $d(u, v)$.

Lemma 10.6. *Node u 's random link leads to a node v with probability*

- $\Theta(1/(d(u, v)^\alpha m^{2-\alpha}))$ if $\alpha < 2$.
- $\Theta(1/(d(u, v)^2 \log n))$ if $\alpha = 2$,
- $\Theta(1/d(u, v)^\alpha)$ if $\alpha > 2$.

Moreover, if $\alpha > 2$, the probability to see a link of length at least d is in $\Theta(1/d^{\alpha-2})$.

Proof. For $\alpha \neq 2$, we have that

$$\sum_{w \in V \setminus \{u\}} \frac{1}{d(u, w)^\alpha} \in \sum_{r=1}^m \frac{\Theta(r)}{r^\alpha} = \Theta \left(\int_{r=1}^m \frac{1}{r^{\alpha-1}} dr \right) = \Theta \left(\left[\frac{r^{2-\alpha}}{2-\alpha} \right]_1^m \right).$$

If $\alpha < 2$, this gives $\Theta(m^{2-\alpha})$, if $\alpha > 2$, it is in $\Theta(1)$. If $\alpha = 2$, we get

$$\sum_{w \in V \setminus \{u\}} \frac{1}{d(u, w)^\alpha} \in \sum_{r=1}^m \frac{\Theta(r)}{r^2} = \Theta(1) \cdot \sum_{r=1}^m \frac{1}{r} = \Theta(\log m) = \Theta(\log n).$$

Multiplying with $d(u, v)^\alpha$ yields the first three bounds.

For the last statement, compute

$$\sum_{\substack{w \in V \\ d(u, w) \geq d}} \Theta(1/d(u, w)^\alpha) = \Theta \left(\int_{r=d}^m \frac{r}{r^\alpha} dr \right) = \Theta \left(\left[\frac{r^{2-\alpha}}{2-\alpha} \right]_d^m \right) = \Theta(1/d^{\alpha-2}).$$

\square

Remarks:

- For $\alpha \neq 2$, this is bad news for the greedy routing algorithm, as it will take $n^{\Omega(1)} = m^{\Omega(1)}$ expected steps to reach the destination. This is disappointing, we were hoping for something polylogarithmic.
- If $\alpha < 2$, in distance $m^{(2-\alpha)/3}$ to the target are $m^{2(2-\alpha)/3}$ many nodes. Thus it takes $\Theta(m^{(2-\alpha)/3})$ links in expectation to find a link that comes that close to the destination. Without finding such a link, we have to go at least this far using grid links only.

- If $\alpha > 2$, it takes $\Theta(m^{(\alpha-2)/(\alpha-1)})$ steps until we see a link of length at least $m^{1/(\alpha-1)}$ in expectation. Without such links, it takes at least $m/m^{1/(\alpha-1)} = m^{(\alpha-2)/(\alpha-1)}$ steps to travel a distance of m .
- Any algorithm that uses only the information on long-range contacts that it can collect at the so far visited nodes cannot be faster.
- However, the case $\alpha = 2$ looks more promising.

Definition 10.7 (Phase). *Consider routing from a node u to a node v and assume that we are at some intermediate node w . We say that we are in phase j at node w if the lattice distance $d(w, v)$ to the target node v is between $2^j < d(w, v) \leq 2^{j+1}$.*

Remarks:

- Enumerating the phases in decreasing order is useful, as notation becomes less cumbersome.
- There are $\lceil \log m \rceil \in O(\log n)$ phases.

Lemma 10.8. *Assume that we are in phase j at node w when routing from u to v . The probability for getting to phase $j-1$ in one step is at least $\Omega(1/\log n)$.*

Proof. Let B_j be the set of nodes x with $d(x, v) \leq 2^j$. We get from phase j to phase $j-1$ if the long-range contact of node w points to some node in B_j . Note that we always make progress while following the greedy routing path. Therefore, we have not seen node w before and the long-range contact of w points to a random node that is independent of anything seen on the path from u to w .

For all nodes $x \in B_j$, we have $d(w, x) \leq d(w, v) + d(x, v) \leq 2^{j+1} + 2^j < 2^{j+2}$. Hence, for each node $x \in B_j$, the probability that the long-range contact of w points to x is $\Omega(1/2^{2j+4} \log n)$. Further, the number of nodes in B_j is at least $(2^j)^2/2 = 2^{2j-1}$. Hence, the probability that some node in B_j is the long range contact of w is at least

$$\Omega\left(|B_j| \cdot \frac{1}{2^{2j+4} \log n}\right) = \Omega\left(\frac{2^{2j-1}}{2^{2j+4} \log n}\right) = \Omega\left(\frac{1}{\log n}\right). \quad \square$$

Theorem 10.9. *Consider the greedy routing path from a node u to a node v on an augmented grid with parameter $\alpha = 2$. The expected length of the path is $O(\log^2 n)$.*

Proof. We already observed that the total number of phases is $O(\log n)$ (the distance to the target is halved when we go from phase j to phase $j-1$). At each point during the routing process, the probability of proceeding to the next phase is at least $\Omega(1/\log n)$. Let X_j be the number of steps in phase j . Because the probability for ending the phase is $\Omega(1/\log n)$ in each step, in expectation we need $O(\log n)$ steps to proceed to the next phase, i.e., $\mathbb{E}[X_j] \in O(\log n)$. Let $X = \sum_j X_j$ be the total number of steps of the routing process. By linearity of expectation, we have

$$\mathbb{E}[X] = \sum_j \mathbb{E}[X_j] \in O(\log^2 n). \quad \square$$

10.2 Propagation Studies

In networks, nodes may influence each other's behavior and decisions. There are many applications where nodes influence their neighbors, e.g. they may impact their opinions, or they may bias what products they buy, or they may pass on a disease.

On a beach (modeled as a line segment), it is best to place an ice cream stand right in the middle of the segment, because you will be able to “control” the beach most easily. What about the second stand, where should it settle? The answer generally depends on the model, but assuming that people will buy ice cream from the stand that is closer, it should go right next to the first stand.

Rumors can spread astoundingly fast through social networks. Traditionally this happens by word of mouth, but with the emergence of the Internet and its possibilities new ways of rumor propagation are available. People write email, use instant messengers or publish their thoughts in a blog. Many factors influence the dissemination of rumors. It is especially important where in a network a rumor is initiated and how convincing it is. Furthermore the underlying network structure decides how fast the information can spread and how many people are reached. More generally, we can speak of diffusion of information in networks. The analysis of these diffusion processes can be useful for viral marketing, e.g. to target a few influential people to initiate marketing campaigns. A company may wish to distribute the rumor of a new product via the most influential individuals in popular social networks such as Facebook. A second company might want to introduce a competing product and has hence to select where to seed the information to be disseminated. Rumor spreading is quite similar to our ice cream stand problem.

More formally, we may study propagation problems in graphs. Given a graph, and two players. Let the first player choose a seed node u_1 ; afterwards let the second player choose a seed node u_2 , with $u_2 \neq u_1$. The goal of the game is to maximize the number of nodes that are closer to one's own seed node.

In many graphs it is an advantage to choose first. In a star graph for instance the first player can choose the center node of the star, controlling all but one node. In some other graphs, the second player can at least score even. But is there a graph where the second player has an advantage?

Theorem 10.10. *In a two player rumor game where both players select one node to initiate their rumor in the graph, the first player does not always win.*

Proof. See Figure 10.3 for an example where the second player will always win, regardless of the decision the first player. If the first player chooses the node x_0 in the center, the second player can select x_1 . Choice x_1 will be outwitted by x_2 , and x_2 itself can be answered by z_1 . All other strategies are either symmetric, or even less promising for the first player. \square

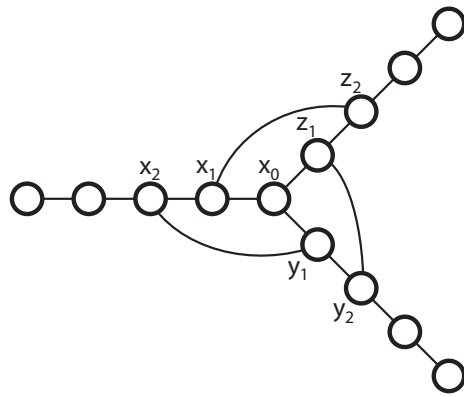


Figure 10.3: Counter example.

Chapter 11

Synchronization

So far, we have mainly studied synchronous algorithms. Generally, asynchronous algorithms are more difficult to obtain. Also it is substantially harder to reason about asynchronous algorithms than about synchronous ones. For instance, computing a BFS tree (Chapter 4) efficiently requires much more work in an asynchronous system. However, many real systems are not synchronous, and we therefore have to design asynchronous algorithms. In this chapter, we will look at general simulation techniques, called *synchronizers*, that allow running synchronous algorithms in asynchronous environments.

11.1 Basics

A synchronizer generates sequences of *clock pulses* at each node of the network satisfying the condition given by the following definition.

Definition 11.1 (valid clock pulse). *We call a clock pulse generated at a node v valid if it is generated after v received all the messages of the synchronous algorithm sent to v by its neighbors in the previous pulses.*

Given a mechanism that generates the clock pulses, a synchronous algorithm is turned into an asynchronous algorithm in an obvious way: As soon as the i^{th} clock pulse is generated at node v , v performs all the actions (local computations and sending of messages) of round i of the synchronous algorithm.

Theorem 11.2. *If all generated clock pulses are valid according to Definition 11.1, the above method provides an asynchronous algorithm that behaves exactly the same way as the given synchronous algorithm.*

Proof. When the i^{th} pulse is generated at a node v , v has sent and received exactly the same messages and performed the same local computations as in the first $i - 1$ rounds of the synchronous algorithm. \square

The main problem when generating the clock pulses at a node v is that v cannot know what messages its neighbors are sending to it in a given synchronous round. Because there are no bounds on link delays, v cannot simply wait “long enough” before generating the next pulse. In order to satisfy Definition 11.1, nodes have to send additional messages for the purpose of synchronization. The total

complexity of the resulting asynchronous algorithm depends on the overhead introduced by the synchronizer. For a synchronizer \mathcal{S} , let $T(\mathcal{S})$ and $M(\mathcal{S})$ be the time and message complexities of \mathcal{S} for each generated clock pulse. As we will see, some of the synchronizers need an initialization phase. We denote the time and message complexities of the initialization by $T_{\text{init}}(\mathcal{S})$ and $M_{\text{init}}(\mathcal{S})$, respectively. If $T(\mathcal{A})$ and $M(\mathcal{A})$ are the time and message complexities of the given synchronous algorithm \mathcal{A} , the total time and message complexities T_{tot} and M_{tot} of the resulting asynchronous algorithm then become

$$T_{\text{tot}} = T_{\text{init}}(\mathcal{S}) + T(\mathcal{A}) \cdot (1 + T(\mathcal{S})) \text{ and } M_{\text{tot}} = M_{\text{init}}(\mathcal{S}) + M(\mathcal{A}) + T(\mathcal{A}) \cdot M(\mathcal{S}),$$

respectively.

Remarks:

- Because the initialization only needs to be done once for each network, we will mostly be interested in the overheads $T(\mathcal{S})$ and $M(\mathcal{S})$ per round of the synchronous algorithm.

Definition 11.3 (Safe Node). *A node v is safe with respect to a certain clock pulse if all messages of the synchronous algorithm sent by v in that pulse have already arrived at their destinations.*

Lemma 11.4. *If all neighbors of a node v are safe with respect to the current clock pulse of v , the next pulse can be generated for v .*

Proof. If all neighbors of v are safe with respect to a certain pulse, v has received all messages of the given pulse. Node v therefore satisfies the condition of Definition 11.1 for generating a valid next pulse. \square

Remarks:

- In order to detect safety, we require that all algorithms send acknowledgements for all received messages. As soon as a node v has received an acknowledgement for each message that it has sent in a certain pulse, it knows that it is safe with respect to that pulse. Note that sending acknowledgements does not increase the asymptotic time and message complexities.

11.2 The Local Synchronizer α

Algorithm 40 Synchronizer α (at node v)

- 1: **wait** until v is safe
 - 2: **send** SAFE to all neighbors
 - 3: **wait** until v receives SAFE messages from all neighbors
 - 4: start new pulse
-

Synchronizer α is very simple. It does not need an initialization. Using acknowledgements, each node eventually detects that it is safe. It then reports this fact directly to all its neighbors. Whenever a node learns that all its neighbors are safe, a new pulse is generated. Algorithm 40 formally describes the synchronizer α .

Theorem 11.5. *The time and message complexities of synchronizer α per synchronous round are*

$$T(\alpha) = O(1) \quad \text{and} \quad M(\alpha) = O(m).$$

Proof. Communication is only between neighbors. As soon as all neighbors of a node v become safe, v knows of this fact after one additional time unit. For every clock pulse, synchronizer α sends at most four additional messages over every edge: Each of the nodes may have to acknowledge a message and reports safety. \square

Remarks:

- Synchronizer α was presented in a framework, mostly set up to have a common standard to discuss different synchronizers. Without the framework, synchronizer α can be explained more easily:
 1. Send message to all neighbors, include round information i and actual data of round i (if any).
 2. Wait for message of round i from all neighbors, and go to next round.
- Although synchronizer α allows for simple and fast synchronization, it produces awfully many messages. Can we do better? Yes.

11.3 The Global Synchronizer β

Algorithm 41 Synchronizer β (at node v)

```

1: wait until  $v$  is safe
2: wait until  $v$  receives SAFE messages from all its children in  $T$ 
3: if  $v \neq \ell$  then
4:   send SAFE message to parent in  $T$ 
5:   wait until PULSE message received from parent in  $T$ 
6: end if
7: send PULSE message to children in  $T$ 
8: start new pulse

```

Synchronizer β needs an initialization that computes a leader node ℓ and a spanning tree T rooted at ℓ . As soon as all nodes are safe, this information is propagated to ℓ by a convergecast. The leader then broadcasts this information to all nodes. The details of synchronizer β are given in Algorithm 41.

Theorem 11.6. *The time and message complexities of synchronizer β per synchronous round are*

$$T(\beta) = O(\text{diameter}(T)) \leq O(n) \quad \text{and} \quad M(\beta) = O(n).$$

The time and message complexities for the initialization are

$$T_{\text{init}}(\beta) = O(n) \quad \text{and} \quad M_{\text{init}}(\beta) = O(m + n \log n).$$

Proof. Because the diameter of T is at most $n - 1$, the convergecast and the broadcast together take at most $2n - 2$ time units. Per clock pulse, the synchronizer sends at most $2n - 2$ synchronization messages (one in each direction over each edge of T).

With an improvement (due to Awerbuch) of the GHS algorithm (Algorithm 15) you saw in Chapter 4, it is possible to construct an MST in time $\mathcal{O}(n)$ with $\mathcal{O}(m + n \log n)$ messages in an asynchronous environment. Once the tree is computed, the tree can be made rooted in time $\mathcal{O}(n)$ with $\mathcal{O}(n)$ messages. \square

Remarks:

- We now got a time-efficient synchronizer (α) and a message-efficient synchronizer (β), it is only natural to ask whether we can have the best of both worlds. And, indeed, we can. How is that synchronizer called? Quite obviously: γ .

11.4 The Hybrid Synchronizer γ

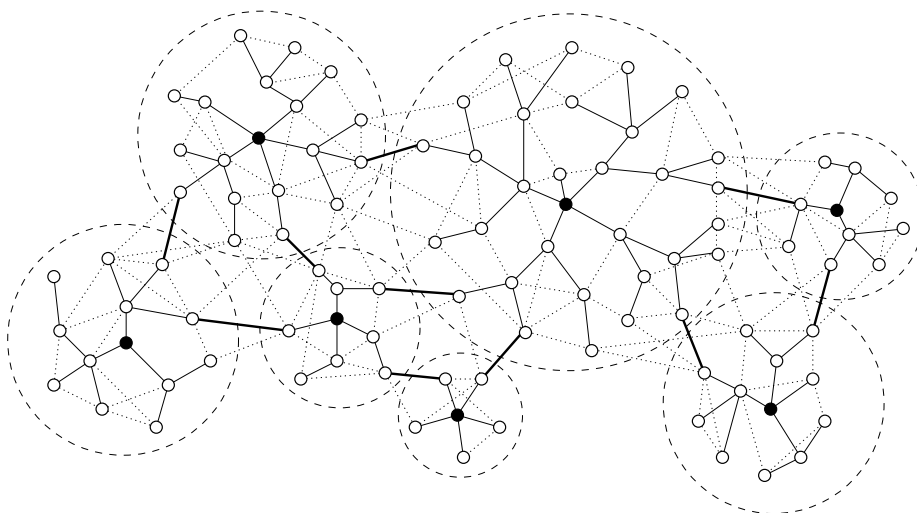


Figure 11.1: A cluster partition of a network: The dashed cycles specify the clusters, cluster leaders are black, the solid edges are the edges of the intracenter trees, and the bold solid edges are the intercluster edges

Synchronizer γ can be seen as a combination of synchronizers α and β . In the initialization phase, the network is partitioned into clusters of small diameter. In each cluster, a leader node is chosen and a BFS tree rooted at this leader node is computed. These trees are called the *intracenter trees*. Two clusters C_1 and C_2 are called neighboring if there are nodes $u \in C_1$ and $v \in C_2$ for which $(u, v) \in E$. For every two neighboring clusters, an *intercluster edge* is chosen, which will serve for communication between these clusters. Figure 11.1 illustrates this partitioning into clusters. We will discuss the details of how to construct such a partition in the next section. We say that a cluster is safe if all its nodes are safe.

Synchronizer γ works in two phases. In a first phase, synchronizer β is applied separately in each cluster by using the intracluster trees. Whenever the leader of a cluster learns that its cluster is safe, it reports this fact to all the nodes in the clusters as well as to the leaders of the neighboring clusters. Now, the nodes of the cluster enter the second phase where they wait until all the neighboring clusters are known to be safe and then generate the next pulse. Hence, we essentially apply synchronizer α between clusters. A detailed description is given by Algorithm 42.

Algorithm 42 Synchronizer γ (at node v)

```

1: wait until  $v$  is safe
2: wait until  $v$  receives SAFE messages from all children in intracluster tree
3: if  $v$  is not cluster leader then
4:   send SAFE message to parent in intracluster tree
5:   wait until CLUSTERSAFE message received from parent
6: end if
7: send CLUSTERSAFE message to all children in intracluster tree
8: send NEIGHBORSAFE message over all intercluster edges of  $v$ 
9: wait until  $v$  receives NEIGHBORSAFE messages from all adjacent inter-
   cluster edges and all children in intracluster tree
10: if  $v$  is not cluster leader then
11:   send NEIGHBORSAFE message to parent in intracluster tree
12:   wait until PULSE message received from parent
13: end if
14: send PULSE message to children in intracluster tree
15: start new pulse

```

Theorem 11.7. *Let m_C be the number of intercluster edges and let k be the maximum cluster radius (i.e., the maximum distance of a leaf to its cluster leader). The time and message complexities of synchronizer γ are*

$$T(\gamma) = O(k) \quad \text{and} \quad M(\gamma) = O(n + m_C).$$

Proof. We ignore acknowledgements, as they do not affect the asymptotic complexities. Let us first look at the number of messages. Over every intracluster tree edge, exactly one SAFE message, one CLUSTERSAFE message, one NEIGHBORSAFE message, and one PULSE message is sent. Further, one NEIGHBORSAFE message is sent over every intercluster edge. Because there are less than n intracluster tree edges, the total message complexity therefore is at most $4n + 2m_C = O(n + m_C)$.

For the time complexity, note that the depth of each intracluster tree is at most k . On each intracluster tree, two convergecasts (the SAFE and NEIGHBORSAFE messages) and two broadcasts (the CLUSTERSAFE and PULSE messages) are performed. The time complexity for this is at most $4k$. There is one more time unit needed to send the NEIGHBORSAFE messages over the intercluster edges. The total time complexity therefore is at most $4k + 1 = O(k)$. \square

11.5 Network Partition

We will now look at the initialization phase of synchronizer γ . Algorithm 43 describes how to construct a partition into clusters that can be used for synchronizer γ . In Algorithm 43, $B(v, r)$ denotes the ball of radius r around v , i.e., $B(v, r) = \{u \in V : d(u, v) \leq r\}$ where $d(u, v)$ is the hop distance between u and v . The algorithm has a parameter $\rho > 1$. The clusters are constructed sequentially. Each cluster is started at an arbitrary node that has not been included in a cluster. Then the cluster radius is grown as long as the cluster grows by a factor more than ρ .

Algorithm 43 Cluster construction

```

1: while unprocessed nodes do
2:   select an arbitrary unprocessed node  $v$ ;
3:    $r := 0$ ;
4:   while  $|B(v, r + 1)| > \rho |B(v, r)|$  do
5:      $r := r + 1$ 
6:   end while
7:   makeCluster( $B(v, r)$ )           // all nodes in  $B(v, r)$  are now processed
8: end while

```

Remarks:

- The algorithm allows a trade-off between the cluster diameter k (and thus the time complexity) and the number of intercluster edges m_C (and thus the message complexity). We will quantify the possibilities in the next section.
- Two very simple partitions would be to make a cluster out of every single node or to make one big cluster that contains the whole graph. We then get synchronizers α and β as special cases of synchronizer γ .

Theorem 11.8. *Algorithm 43 computes a partition of the network graph into clusters of radius at most $\log_\rho n$. The number of intercluster edges is at most $(\rho - 1) \cdot n$.*

Proof. The radius of a cluster is initially 0 and does only grow as long as it grows by a factor larger than ρ . Since there are only n nodes in the graph, this can happen at most $\log_\rho n$ times.

To count the number of intercluster edges, observe that an edge can only become an intercluster edge if it connects a node at the boundary of a cluster with a node outside a cluster. Consider a cluster C of size $|C|$. We know that $C = B(v, r)$ for some $v \in V$ and $r \geq 0$. Further, we know that $|B(v, r + 1)| \leq \rho \cdot |B(v, r)|$. The number of nodes adjacent to cluster C is therefore at most $|B(v, r + 1) \setminus B(v, r)| \leq \rho \cdot |C| - |C|$. Because there is only one intercluster edge connecting two clusters by definition, the number of intercluster edges adjacent to C is at most $(\rho - 1) \cdot |C|$. Summing over all clusters, we get that the total number of intercluster edges is at most $(\rho - 1) \cdot n$. \square

Corollary 11.9. *Using $\rho = 2$, Algorithm 43 computes a clustering with cluster radius at most $\log_2 n$ and with at most n intercluster edges.*

Corollary 11.10. *Using $\rho = n^{1/k}$, Algorithm 43 computes a clustering with cluster radius at most k and at most $\mathcal{O}(n^{1+1/k})$ intercluster edges.*

Remarks:

- Algorithm 43 describes a centralized construction of the partitioning of the graph. For $\rho \geq 2$, the clustering can be computed by an asynchronous distributed algorithm in time $\mathcal{O}(n)$ with $\mathcal{O}(m + n \log n)$ (reasonably sized) messages (showing this will be part of the exercises).
- It can be shown that the trade-off between cluster radius and number of intercluster edges of Algorithm 43 is asymptotically optimal. There are graphs for which every clustering into clusters of radius at most k requires $n^{1+c/k}$ intercluster edges for some constant c .

The above remarks lead to a complete characterization of the complexity of synchronizer γ .

Corollary 11.11. *The time and message complexities of synchronizer γ per synchronous round are*

$$T(\gamma) = \mathcal{O}(k) \quad \text{and} \quad M(\gamma) = \mathcal{O}(n^{1+1/k}).$$

The time and message complexities for the initialization are

$$T_{\text{init}}(\gamma) = \mathcal{O}(n) \quad \text{and} \quad M_{\text{init}}(\gamma) = \mathcal{O}(m + n \log n).$$

Remarks:

- The synchronizer idea and the synchronizers discussed in this chapter are due to Baruch Awerbuch.
- In Chapter 4, you have seen that by using flooding, there is a very simple synchronous algorithm to compute a BFS tree in time $\mathcal{O}(D)$ with message complexity $\mathcal{O}(m)$. If we use synchronizer γ to make this algorithm asynchronous, we get an algorithm with time complexity $\mathcal{O}(n + D \log n)$ and message complexity $\mathcal{O}(m + n \log n + D \cdot n)$ (including initialization).
- The synchronizers α , β , and γ achieve global synchronization, i.e. every node generates every clock pulse. The disadvantage of this is that nodes that do not participate in a computation also have to participate in the synchronization. In many computations (e.g. in a BFS construction), many nodes only participate for a few synchronous rounds. An improved synchronizer due to Awerbuch and Peleg can exploit such a scenario and achieves time and message complexity $\mathcal{O}(\log^3 n)$ per synchronous round (without initialization).
- It can be shown that if all nodes in the network need to generate all pulses, the trade-off of synchronizer γ is asymptotically optimal.
- Partitions of networks into clusters of small diameter and coverings of networks with clusters of small diameters come in many variations and have various applications in distributed computations. In particular, apart from synchronizers, algorithms for routing, the construction of sparse spanning subgraphs, distributed data structures, and even computations of local structures such as a MIS or a dominating set are based on some kind of network partitions or covers.

11.6 Clock Synchronization

“A man with one clock knows what time it is – a man with two is never sure.”

Synchronizers can directly be used to give nodes in an asynchronous network a common notion of time. In wireless networks, for instance, many basic protocols need an accurate time. Sometimes a common time in the whole network is needed, often it is enough to synchronize neighbors. The purpose of the time division multiple access (TDMA) protocol is to use the common wireless channel as efficiently as possible, i.e., interfering nodes should never transmit at the same time (on the same frequency). If we use synchronizer β to give the nodes a common notion of time, every single clock cycle costs D time units!

Often, each (wireless) node is equipped with an internal clock. Using this clock, it should be possible to divide time into slots, and make each node send (or listen, or sleep, respectively) in the appropriate slots according to the media access control (MAC) layer protocol used.

However, as it turns out, synchronizing clocks in a network is not trivial. As nodes' internal clocks are not perfect, they will run at speeds that are time-dependent. For instance, variations in temperature or supply voltage will affect this *clock drift*. For standard clocks, the drift is in the order of parts per million, i.e., within a second, it will accumulate to a couple of microseconds. Wireless TDMA protocols account for this by introducing *guard times*. Whenever a node knows that it is about to receive a message from a neighbor, it powers up its radio a little bit earlier to make sure that it does not miss the message even when clocks are not perfectly synchronized. If nodes are badly synchronized, messages of different slots might collide.

In the *clock synchronization* problem, we are given a network (graph) with n nodes. The goal for each node is to have a logical clock such that the logical clock values are well synchronized, and close to real time. Each node is equipped with a hardware clock, that ticks more or less in real time, i.e., the time between two pulses is arbitrary between $[1 - \epsilon, 1 + \epsilon]$, for a constant $\epsilon \ll 1$. Similarly as in our asynchronous model, we assume that messages sent over the edges of the graph have a delivery time between $[0, 1]$. In other words, we have a bounded but variable drift on the hardware clocks and an arbitrary jitter in the delivery times. The goal is to design a message-passing algorithm that ensures that the logical clock skew of adjacent nodes is as small as possible at all times.

Theorem 11.12. *The global clock skew (the logical clock difference between any two nodes in the graph) is $\Omega(D)$, where D is the diameter of the graph.*

Proof. For a node u , let t_u be the logical time of u and let $(u \rightarrow v)$ denote a message sent from u to a node v . Let $t(m)$ be the time delay of a message m and let u and v be neighboring nodes. First consider a case where the message delays between u and v are $1/2$. Then all the messages sent by u and v at time i according to the clock of the sender arrive at time $i + 1/2$ according to the clock of the receiver.

Then consider the following cases

- $t_u = t_v + 1/2, t(u \rightarrow v) = 1, t(v \rightarrow u) = 0$
- $t_u = t_v - 1/2, t(u \rightarrow v) = 0, t(v \rightarrow u) = 1,$

where the message delivery time is always fast for one node and slow for the other and the logical clocks are off by $1/2$. In both scenarios, the messages sent at time i according to the clock of the sender arrive at time $i + 1/2$ according to the logical clock of the receiver. Therefore, for nodes u and v , both cases with clock drift seem the same as the case with perfectly synchronized clocks. Furthermore, in a linked list of D nodes, the left- and rightmost nodes l, r cannot distinguish $t_l = t_r + D/2$ from $t_l = t_r - D/2$. \square

Remarks:

- From Theorem 11.12, it directly follows that all the clock synchronization algorithms we studied have a global skew of $\Omega(D)$.
- Many natural algorithms manage to achieve a global clock skew of $\mathcal{O}(D)$.

As both the message jitter and hardware clock drift are bounded by constants, it feels like we should be able to get a constant drift between neighboring nodes. As synchronizer α pays most attention to the local synchronization, we take a look at a protocol inspired by the synchronizer α . A pseudo-code representation for the clock synchronization protocol α is given in Algorithm 44.

Algorithm 44 Clock synchronization α (at node v)

```

1: repeat
2:   send logical time  $t_v$  to all neighbors
3:   if Receive logical time  $t_u$ , where  $t_u > t_v$ , from any neighbor  $u$  then
4:      $t_v := t_u$ 
5:   end if
6: until done

```

Lemma 11.13. *The clock synchronization protocol α has a local skew of $\Omega(n)$.*

Proof. Let the graph be a linked list of D nodes. We denote the nodes by v_1, v_2, \dots, v_D from left to right and the logical clock of node v_i by t_i . Apart from the left-most node v_1 all hardware clocks run with speed 1 (real time). Node v_1 runs at maximum speed, i.e. the time between two pulses is not 1 but $1 - \epsilon$. Assume that initially all message delays are 1. After some time, node v_1 will start to speed up v_2 , and after some more time v_2 will speed up v_3 , and so on. At some point of time, we will have a clock skew of 1 between any two neighbors. In particular $t_1 = t_D + D - 1$.

Now we start playing around with the message delays. Let $t_1 = T$. First we set the delay between the v_1 and v_2 to 0. Now node v_2 immediately adjusts its logical clock to T . After this event (which is instantaneous in our model) we set the delay between v_2 and v_3 to 0, which results in v_3 setting its logical clock to T as well. We perform this successively to all pairs of nodes until v_{D-2} and v_{D-1} . Now node v_{D-1} sets its logical clock to T , which indicates that the difference between the logical clocks of v_{D-1} and v_D is $T - (T - (D - 1)) = D - 1$. \square

Remarks:

- The introduced examples may seem cooked-up, but examples like this exist in all networks, and for all algorithms. Indeed, it was shown that any natural clock synchronization algorithm must have a bad local skew. In particular, a protocol that averages between all neighbors is even worse than the introduced α algorithm. This algorithm has a clock skew of $\Omega(D^2)$ in the linked list, at all times.
- Recently, there was a lot of progress in this area, and it was shown that the local clock skew is $\Theta(\log D)$, i.e., there is a protocol that achieves this bound, and there proof that no algorithm can be better than this bound!
- Note that these are worst-case bounds. In practice, clock drift and message delays may not be the worst possible, typically the speed of hardware clocks changes at a comparatively slow pace and the message transmission times follow a benign probability distribution. If we assume this, better protocols do exist.

Chapter 12

Hard Problems

This chapter is on “hard” problems in distributed computing. In sequential computing, there are NP-hard problems which are conjectured to take exponential time. Is there something similar in distributed computing? Using flooding/echo (Algorithms 11,12) from Chapter 4, everything so far was solvable basically in $\mathcal{O}(D)$ time, where D is the diameter of the network.

12.1 Diameter & APSP

But how do we compute the diameter itself?! With flooding/echo, of course!

Algorithm 45 Naive Diameter Construction

- 1: all nodes compute their radius by synchronous flooding/echo
 - 2: all nodes flood their radius on the constructed BFS tree
 - 3: the maximum radius a node sees is the diameter
-

Remarks:

- Since all these phases only take $\mathcal{O}(D)$ time, nodes know the diameter in $\mathcal{O}(D)$ time, which is asymptotically optimal.
- However, there is a problem! Nodes are now involved in n parallel flooding/echo operations, thus a node may have to handle many and big messages in one single time step. Although this is not strictly illegal in the message passing model, it still feels like cheating! A natural question is whether we can do the same by just sending short messages in each round.
- In Definition 3.1 of Chapter 2 we postulated that nodes should send only messages of “reasonable” size. In this chapter we strengthen the definition a bit, and require that each message should have at most $\mathcal{O}(\log n)$ bits. This is generally enough to communicate a constant number of ID’s or values to neighbors, but not enough to communicate everything a node knows!
- A simple way to avoid large messages is to split them into small messages that are sent using several rounds. This can cause that messages are

getting delayed in some nodes but not in others. The flooding might not use edges of a BFS tree anymore! These floodings might not compute correct distances anymore! On the other hand we know that the maximal message size in Algorithm 45 is $\mathcal{O}(n \log n)$. So we could just simulate each of these “big message” rounds by n “small message” rounds using small messages. This yields a runtime of $\mathcal{O}(nD)$ which is not desirable. A third possible approach is “starting each flooding/echo one after each other” and results in $\mathcal{O}(nD)$ in the worst case as well.

- So let us fix above algorithm! The key idea is to arrange the flooding-echo processes in a more organized way: Start the flooding processes in a certain order and prove that at any time, each node is only involved in one flooding. This is realized in Algorithm 46.

Definition 12.1. (BFS_v) *Performing a breadth first search at node v produces spanning tree BFS_v (see Chapter 4). This takes time $\mathcal{O}(D)$ using small messages.*

Remarks:

- A spanning tree of a graph G can be traversed in time $\mathcal{O}(n)$ by sending a pebble over an edge in each time slot.
- This can be done using e.g. a depth first search (DFS): Start at the root of a tree, recursively visit all nodes in the following way. If the current node still has an unvisited child, then the pebble always visits that child first. Return to the parent only when all children have been visited.
- Algorithm 46 works as follows: Given a graph G , first a leader l computes its BFS tree BFS_l . Then we send a pebble P to traverse tree BFS_l . Each time pebble P enters a node v for the first time, P waits one time slot, and then starts a breadth first search (BFS) – using edges in G – from v with the aim of computing the distances from v to all other nodes. Since we start a BFS_v from every node v , each node u learns its distance to all these nodes v during the according execution of BFS_v . There is no need for a echo-process at the end of BFS_u .

Algorithm 46 Computes APSP on G .

```

1: Assume we have a leader node  $l$  (if not, compute one first)
2: compute  $BFS_l$  of leader  $l$ 
3: send a pebble  $P$  to traverse  $BFS_l$  in a DFS way;
4: while  $P$  traverses  $BFS_l$  do
5:   if  $P$  visits a new node  $v$  then
6:     wait one time slot; // avoid congestion
7:     start  $BFS_v$  from node  $v$ ; // compute all distances to  $v$ 
8:     // the depth of node  $u$  in  $BFS_v$  is  $d(u, v)$ 
9:   end if
10: end while

```

Remarks:

- Having all distances is nice, but how do we get the diameter? Well, as before, each node could just floods its radius (its maximum distance) into the network. However, messages are small now and we need to modify this slightly. In each round a node only sends the maximal distance that it is aware of to its neighbors. After D rounds each node will know the maximum distance among all nodes.

Lemma 12.2. *In Algorithm 46, at no time a node w is simultaneously active for both BFS_u and BFS_v .*

Proof. Assume a BFS_u is started at time t_u at node u . Then node w will be involved in BFS_u at time $t_u + d(u, w)$. Now, consider a node v whose BFS_v is started at time $t_v > t_u$. According to the algorithm this implies that the pebble visits v after u and took some time to travel from u to v . In particular, the time to get from u to v is at least $d(u, v)$, in addition at least node v is visited for the first time (which involves waiting at least one time slot), and we have $t_v \geq t_u + d(u, v) + 1$. Using this and the triangle inequality, we get that node w is involved in BFS_v strictly after being involved in BFS_u since $t_v + d(v, w) \geq (t_u + d(u, v) + 1) + d(v, w) \geq t_u + d(u, w) + 1 > t_u + d(u, w)$. \square

Theorem 12.3. *Algorithm 46 computes APSP (all pairs shortest path) in time $\mathcal{O}(n)$.*

Proof. Since the previous lemma holds for any pair of vertices, no two BFS “interfere” with each other, i.e. all messages can be sent on time without congestion. Hence, all BFS stop at most D time slots after they were started. We conclude that the runtime of the algorithm is determined by the time $\mathcal{O}(D)$ we need to build tree BFS_l , plus the time $\mathcal{O}(n)$ that P needs to traverse BFS_l , plus the time $\mathcal{O}(D)$ needed by the last BFS that P initiated. Since $D \leq n$, this is all in $\mathcal{O}(n)$. \square

Remarks:

- All of a sudden our algorithm needs $\mathcal{O}(n)$ time, and possibly $n \gg D$. We should be able to do better, right?!
- Unfortunately not! One can show that computing the diameter of a network needs $\Omega(n/\log n)$ time.
- On the other hand we can check fast whether a graph has diameter 1 or not: each node just checks whether its degree is $n - 1$ and tells the result to its neighbors.

12.2 Lower Bound Graphs

We define a family \mathcal{G} of graphs that we use to prove a lower bound on the rounds needed to compute the diameter. To simplify our analysis, we assume that $(n - 2)$ can be divided by 8. We start by defining four sets of nodes, each consisting of $q = q(n) := (n - 2)/4$ nodes. Throughout this chapter we write $[q]$ as a short version of $\{1, \dots, q\}$ and define:

$$\begin{aligned}
\mathbf{L}_0 &:= \{l_i \mid i \in [q]\} && // \text{ upper left in Figure 12.1} \\
\mathbf{L}_1 &:= \{l'_i \mid i \in [q]\} && // \text{ lower left} \\
\mathbf{R}_0 &:= \{r_i \mid i \in [q]\} && // \text{ upper right} \\
\mathbf{R}_1 &:= \{r'_i \mid i \in [q]\} && // \text{ lower right}
\end{aligned}$$

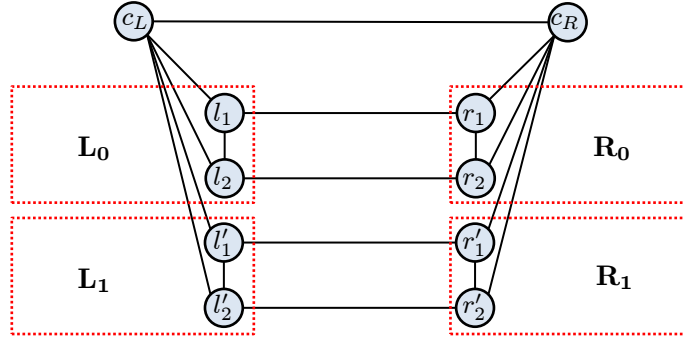


Figure 12.1: The above skeleton G' contains $n = 10$ nodes, such that $q = 2$.

We add node c_L and connect it to all nodes in \mathbf{L}_0 and \mathbf{L}_1 . Then we add node c_R , connected to all nodes in \mathbf{R}_0 and \mathbf{R}_1 . Furthermore, nodes c_L and c_R are connected by an edge. For $i \in [q]$ we connect l_i to r_i and l'_i to r'_i . Also we add edges such that nodes in \mathbf{L}_0 are a clique, nodes in \mathbf{L}_1 are a clique, nodes in \mathbf{R}_0 are a clique, and nodes in \mathbf{R}_1 are a clique. The resulting graph is called G' . Graph G' is the skeleton of any graph in family \mathcal{G} .

More formally skeleton $G' = (V', E')$ is:

$$V' := \mathbf{L}_0 \cup \mathbf{L}_1 \cup \mathbf{R}_0 \cup \mathbf{R}_1 \cup \{c_L, c_R\}$$

$$\begin{aligned}
E' &:= \bigcup_{v \in \mathbf{L}_0 \cup \mathbf{L}_1} \{(v, c_L)\} && // \text{ connections to } c_L \\
&\cup \bigcup_{v \in \mathbf{R}_0 \cup \mathbf{R}_1} \{(v, c_R)\} && // \text{ connections to } c_R \\
&\cup \bigcup_{i \in [q]} \{(l_i, r_i), (l'_i, r'_i)\} \cup \{(c_L, c_R)\} && // \text{ connects left to right} \\
&\cup \bigcup_{S \in \{\mathbf{L}_0, \mathbf{L}_1, \mathbf{R}_0, \mathbf{R}_1\}} \bigcup_{u \neq v \in S} \{(u, v)\} && // \text{ clique edges}
\end{aligned}$$

To simplify our arguments, we partition G' into two parts: **Part L** is the subgraph induced by nodes $\mathbf{L}_0 \cup \mathbf{L}_1 \cup \{c_L\}$. **Part R** is the subgraph induced by nodes $\mathbf{R}_0 \cup \mathbf{R}_1 \cup \{c_R\}$.

Family \mathcal{G} contains any graph G that is derived from G' by adding any combination of edges of the form (l_i, l'_j) resp. (r_i, r'_j) with $l_i \in \mathbf{L}_0$, $l'_j \in \mathbf{L}_1$, $r_i \in \mathbf{R}_0$, and $r'_j \in \mathbf{R}_1$.

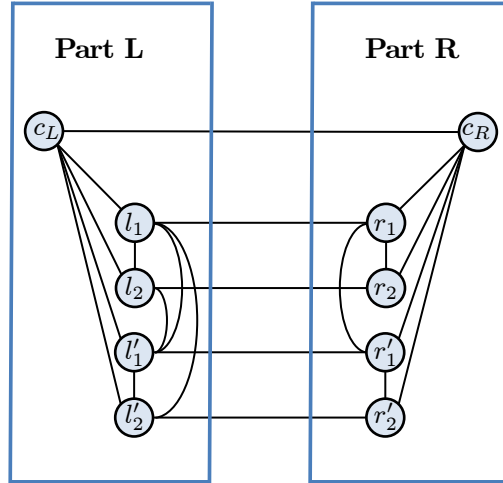


Figure 12.2: The above graph G has $n = 10$ and a member of family \mathcal{G} . What is the diameter of G ?

Lemma 12.4. *The diameter of a graph $G = (V, E) \in \mathcal{G}$ is 2 if and only if: For each tuple (i, j) with $i, j \in [q]$, there is either edge (l_i, l'_j) or edge (r_i, r'_j) (or both edges) in E .*

Proof. Note that the distance between most pairs of nodes is at most 2. In particular, the radius of c_L resp. c_R is 2. Thanks to c_L resp. c_R the distance between, any two nodes within **Part L** resp. within **Part R** is at most 2. Because of the cliques $\mathbf{L}_0, \mathbf{L}_1, \mathbf{R}_0, \mathbf{R}_1$, distances between l_i and r_j resp. l'_i and r'_j is at most 2.

The only interesting case is between a node $l_i \in \mathbf{L}_0$ and node $r'_j \in \mathbf{R}_1$ (or, symmetrically, between $l'_j \in \mathbf{L}_1$ and node $r_i \in \mathbf{R}_0$). If either edge (l_i, l'_j) or edge (r_i, r'_j) is present, then this distance is 2, since the path (l_i, l'_j, r'_j) or the path (l_i, r_i, r'_j) exists. If neither of the two edges exist, then the neighborhood of l_i consists of $\{c_L, r_i\}$, all nodes in \mathbf{L}_0 , and some nodes in $\mathbf{L}_1 \setminus \{l'_j\}$, and the neighborhood of r'_j consists of $\{c_R, l'_j\}$, all nodes in \mathbf{R}_1 , and some nodes in $\mathbf{R}_0 \setminus \{r_i\}$ (see for example Figure 12.3 with $i = 1$ and $j = 2$.) Since the two neighborhoods do not share a common node, the distance between l_i and r'_j is (at least) 3. \square

Remarks:

- Each part contains up to $q^2 \in \Theta(n^2)$ edges.
- There are $2q + 1 \in \Theta(n)$ edges connecting the left and the right part. Since in each round we can transmit $\mathcal{O}(\log n)$ bits over each edge (in each direction), the bandwidth between **Part L** and **Part R** is $\mathcal{O}(n \log n)$.

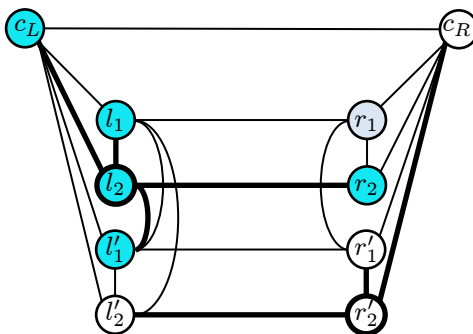


Figure 12.3: Nodes in the neighborhood of l_2 are cyan, the neighborhood of r'_2 is white. Since these neighborhoods do not intersect, the distance of these two nodes is $d(l_2, r'_2) > 2$. If e.g. edge (l_2, l'_2) was included, their distance was 2.

- If we transmit the information of the $\Theta(n^2)$ edges in a naive way with a bandwidth of $\mathcal{O}(n \log n)$, we need $\Omega(n/\log n)$ time. But maybe we can do better?!? Can an algorithm be smarter and only send the information that is really necessary to tell whether the diameter is 2?
- It turns out that any algorithm needs $\Omega(n/\log n)$ rounds, since the information that is really necessary to tell that the diameter is larger than 2 contains basically $\Theta(n^2)$ bits.

12.3 Communication Complexity

To prove the last remark formally, we can use arguments from two-party communication complexity. This area essentially deals with a basic version of distributed computation: two parties are given some input each and want to solve a task on this input.

We consider two students (Alice and Bob) at two different universities connected by a communication channel (e.g. via email) and we assume this channel to be reliable. Now Alice and Bob want to check whether they received the same problem set for homework (we assume their professors are lazy and wrote it on the black board instead of putting a nicely prepared document online.) Do Alice and Bob really need to type the whole problem set into their emails? In a more formal way: Alice receives an n -bit string x and Bob another n -bit string y , and the goal is for both of them to compute the equality function.

Definition 12.5. (*Equality.*) We define the equality function EQ to be:

$$\text{EQ}(x, y) := \begin{cases} 1 & : x = y \\ 0 & : x \neq y \end{cases}.$$

Remarks:

- In a more general setting, Alice and Bob are interested in computing a certain function $f : \{0,1\}^k \times \{0,1\}^k \rightarrow \{0,1\}$ with the least amount of communication between them. Of course they can always succeed by having Alice send her whole n -bit string to Bob, who then computes the function, but the idea here is to find clever ways of calculating f with less than n bits of communication. We measure how clever they can be as follows:

Definition 12.6. (Communication complexity CC .) The communication complexity of protocol A for function f is $CC(A, f) :=$ minimum number of bits exchanged between Alice and Bob in the worst case when using A . The communication complexity of f is $CC(f) := \min\{CC(A, f) \mid A \text{ solves } f\}$. That is the minimal number of bits that the best protocol needs to send in the worst case.

Definition 12.7. For a given function f , we define a $2^k \times 2^k$ matrix M^f representing f . That is $M_{x,y}^f := f(x, y)$.

Example 12.8. For EQ, in case $k = 3$, matrix M^{EQ} looks like this:

EQ	000	001	010	011	100	101	110	111	$\leftarrow x$
000	1	0	0	0	0	0	0	0	
001	0	1	0	0	0	0	0	0	
010	0	0	1	0	0	0	0	0	
011	0	0	0	1	0	0	0	0	
100	0	0	0	0	1	0	0	0	
101	0	0	0	0	0	1	0	0	
110	0	0	0	0	0	0	1	0	
111	0	0	0	0	0	0	0	1	
$\uparrow y$									

As a next step we define a (combinatorial) monochromatic rectangle. These are “submatrices” of M^f which contain the same entry.

Definition 12.9. (monochromatic rectangle.) A set $R \subseteq \{0,1\}^k \times \{0,1\}^k$ is called a monochromatic rectangle, if

- whenever $(x_1, y_1) \in R$ and $(x_2, y_2) \in R$ then $(x_1, y_2) \in R$.
- there is a fixed z such that $f(x, y) = z$ for all $(x, y) \in R$.

Example 12.10. The first three of the following rectangles are monochromatic, the last one is not:

R_1	$= \{011\} \times \{011\}$	Example 12.8: light gray
R_2	$= \{011, 100, 101, 110\} \times \{000, 001\}$	Example 12.8: gray
R_3	$= \{000, 001, 101\} \times \{011, 100, 110, 111\}$	Example 12.8: dark gray
R_4	$= \{000, 001\} \times \{000, 001\}$	Example 12.8: boxed

Each time Alice and Bob exchange a bit, they can eliminate columns/rows of the matrix M^f and a combinatorial rectangle is left. They can stop communicating when this remaining rectangle is monochromatic. Informally speaking, a fooling set can be used to fool a protocol that wants to be lazy: if the fooling

set is large, there will be many maximal monochromatic rectangles (maximal in the sense that they cannot be extended while staying monochromatic). Since by communicating one bit the set of possible monochromatic rectangles does not shrink too much, we can expect that it takes long time until a monochromatic rectangle is found in the worst case.

Definition 12.11. (*fooling set.*) A set $S \subset \{0, 1\}^k \times \{0, 1\}^k$ fools f if there is a fixed z such that

- $f(x, y) = z$ for each $(x, y) \in S$
- For any $(x_1, y_1) \neq (x_2, y_2) \in S$, the rectangle $\{x_1, x_2\} \times \{y_1, y_2\}$ is not monochromatic: Either $f(x_1, y_2) \neq z$, $f(x_2, y_1) \neq z$ or both $\neq z$.

Example 12.12. Consider $S = \{(000, 000), (001, 001)\}$. Take a look at the non-monochromatic rectangle R_4 in Example 12.10. Verify that S is indeed a fooling set for EQ!

Remarks:

- Can you find a larger fooling set for EQ?
- We assume that Alice and Bob take turns in sending a bit. This results in 2 possible actions (send 0/1) per round and in 2^t action patterns during a sequence of t rounds.

Lemma 12.13. If S is a fooling set for f , then $CC(f) = \Omega(\log |S|)$.

Proof. For simplicity we assume that $|S|$ is a power of 2. We prove the statement via contradiction: fix a protocol A and assume that it needs $t := \log_2(|S|/2)$ rounds in the worst case. Then there are $2^t = 2^{\log_2(|S|/2)} = |S|/2$ possible action patterns. On the other hand there are $|S| = 2^{\log_2 |S|}$ elements in S and we conclude that at least two elements (let's call them $(x_1, y_1), (x_2, y_2)$) in S cause the same action pattern P . Naturally, the action pattern on the alternative inputs $(x_1, y_2), (x_2, y_1)$ will be P as well: in the first round Alice and Bob have no information on the other party's string and send the same bit that was sent in P . Based on this, they determine the second bit to be exchanged, which will be the same as the second one in P for a similar reason. This continues for all t rounds. We conclude that after t rounds, Alice does not know whether Bob's input is y_1 or y_2 and Bob does not know whether Alice's input is x_1 or x_2 . By the definition of fooling set, either

- $f(x_1, y_2) \neq f(x_1, y_1)$ in which case Alice (with input x_1) does not know the solution yet,

or

- $f(x_2, y_1) \neq f(x_1, y_1)$ in which case Bob (with input y_1) does not know the solution yet.

This contradicts the assumption that A leads to a correct decision for all inputs after t rounds. Therefore at least $t + 1$ rounds are necessary, which is

$$t + 1 = \log_2(|S|/2) + 1 = \frac{\log |S|}{\log 2} \in \Omega(\log |S|).$$

□

Theorem 12.14. $CC(EQ) = \Omega(k)$.

Proof. The set $S := \{(x, x) \mid x \in \{0, 1\}^k\}$ fools EQ and has size 2^k . Now apply Lemma 12.13. \square

Definition 12.15. Denote the negation of a string z by \bar{z} and by $x \circ y$ the concatenation of strings x and y .

Lemma 12.16. Let x, y be k -bit strings. Then $x \neq y$ if and only if there is an index $i \in [2k]$ such that the i^{th} bit of $x \circ \bar{x}$ and the i^{th} bit of $\bar{y} \circ y$ are both 0.

Proof. If $x \neq y$, there is an $j \in [k]$ such that x and y differ in the j^{th} bit. Therefore either the j^{th} bit of both x and \bar{y} is 0, or the j^{th} bit of \bar{x} and y is 0. For this reason, there is an $i \in [2k]$ such that $x \circ \bar{x}$ and $\bar{y} \circ y$ are both 0 at position i .

If $x = y$, then for any $i \in [2k]$ it is always the case that either the i^{th} bit of $x \circ \bar{x}$ is 1 or the i^{th} bit of $\bar{y} \circ y$ (which is the negation of $x \circ \bar{x}$ in this case) is 1. \square

Remarks:

- With these insights we get back to the problem of computing the diameter of a graph and relate this problem to EQ .

Definition 12.17. Using the parameter q defined before, we define a bijective map between all pairs x, y of q^2 -bit strings and the graphs in \mathcal{G} : each pair of strings x, y is mapped to graph $G_{x,y} \in \mathcal{G}$ that is derived from skeleton G' by adding

- edge (l_i, l'_j) to **Part L** if and only if the $(j + q \cdot i)^{th}$ bit of x is 1.
- edge (r_i, r'_j) to **Part R** if and only if the $(j + q \cdot i)^{th}$ bit of y is 1.

Remarks:

- Clearly, **Part L** of $G_{x,y}$ depends on x only and **Part R** depends on y only.

Lemma 12.18. Let x and y be $\frac{q^2}{2}$ -bit strings given to Alice and Bob¹. Then graph $G := G_{x \circ \bar{x}, \bar{y} \circ y} \in \mathcal{G}$ has diameter 2 if and only if $x = y$.

Proof. By Lemma 12.16 and the construction of G , there is neither edge (l_i, l'_j) nor edge (r_i, r'_j) in E if and only if $x \neq y$. Applying Lemma 12.4 yields: G has diameter 2 if and only if $x = y$. \square

Theorem 12.19. Any distributed algorithm A that decides whether a graph G has diameter 2 might need $\Omega\left(\frac{n}{\log n} + D\right)$ time.

¹Thats why we need that $n - 2$ can be divided by 8.

Proof. Computing D for sure needs time $\Omega(D)$. It remains to prove $\Omega\left(\frac{n}{\log n}\right)$. Assume there is a distributed algorithm A that decides whether the diameter of a graph is 2 in time $o(n/\log n)$. When Alice and Bob are given $\frac{q^2}{2}$ -bit inputs x and y , they can simulate A to decide whether $x = y$ as follows: Alice constructs **Part L** of $G_{x \circ \bar{x}, \bar{y} \circ y}$ and Bob constructs **Part R**. As we remarked, both parts are independent of each other such that **Part L** can be constructed by Alice without knowing y and **Part R** can be constructed by Bob without knowing x . Furthermore, $G_{x \circ \bar{x}, \bar{y} \circ y}$ has diameter 2 if and only if $x = y$ (Lemma 12.18.)

Now Alice and Bob simulate the distributed algorithm A round by round: In the first round, they determine which messages the nodes in their part of G would send. Then they use their communication channel to exchange all $2(2q+1) \in \Theta(n)$ messages that would be sent over edges between **Part L** and **Part R** in this round while executing A on G . Based on this Alice and Bob determine which messages would be sent in round two and so on. For each round simulated by Alice and Bob, they need to communicate $\Theta(n \log n)$ bits: possibly $\Theta(\log n)$ bits for each of $\Theta(n)$ messages. Since A makes a decision after $o(n/\log n)$ rounds, this yields a total communication of $o(n^2)$ bits. On the other hand, Lemma 12.14 states that to decide whether x equals y , Alice and Bob need to communicate at least $\Omega\left(\frac{q^2}{2}\right) = \Omega(n^2)$ bits. A contradiction. \square

Remarks:

- Until now we only considered deterministic algorithms. Can one do better using randomness?

Algorithm 47 Randomized evaluation of EQ .

- 1: Alice and Bob use public randomness. That is they both have access to the same random bit string $z \in \{0, 1\}^k$
 - 2: Alice sends bit $a := \sum_{i \in [k]} x_i \cdot z_i \pmod 2$ to Bob
 - 3: Bob sends bit $b := \sum_{i \in [k]} y_i \cdot z_i \pmod 2$ to Alice
 - 4: **if** $a \neq b$ **then**
 - 5: we know $x \neq y$
 - 6: **end if**
-

Lemma 12.20. *If $x \neq y$, Algorithm 47 discovers $x \neq y$ with probability at least $1/2$.*

Proof. Note that if $x = y$ we have $a = b$ for sure.

If $x \neq y$, Algorithm 47 may not reveal inequality. For instance, for $k = 2$, if $x = 01$, $y = 10$ and $z = 11$ we get $a = b = 1$. In general, let I be the set of indices where $x_i \neq y_i$, i.e. $I := \{i \in [k] \mid x_i \neq y_i\}$. Since $x \neq y$, we know that $|I| > 0$. We have

$$|a - b| \equiv \sum_{i \in I} z_i \pmod 2,$$

and since all z_i with $i \in I$ are random, we get that $a \neq b$ with probability at least $1/2$. \square

Remarks:

- By excluding the vector $z = 0^k$ we can even get a discovery probability strictly larger than $1/2$.
- Repeating the Algorithm 47 with different random strings z , the error probability can be reduced arbitrarily.
- Does this imply that there is a fast randomized algorithm to determine the diameter? Unfortunately not!
- Sometimes public randomness is not available, but private randomness is. Here Alice has her own random string and Bob has his own random string. A modified version of Algorithm 47 also works with private randomness at the cost of the runtime.
- One can prove an $\Omega(n/\log n)$ lower bound for any randomized distributed algorithm that computes the diameter. To do so one considers the disjointness function *DISJ* instead of equality. Here, Alice is given a subset $X \subseteq [k]$ and Bob is given a subset $Y \subseteq [k]$ and they need to determine whether $Y \cap X = \emptyset$. (X and Y can be represented by k -bit strings x, y .) The reduction is similar as the one presented above but uses graph $G_{\bar{x}, \bar{y}}$ instead of $G_{x \circ \bar{x}, \bar{y} \circ y}$. However, the lower bound for the randomized communication complexity of *DISJ* is more involved than the lower bound for $CC(EQ)$.
- Since one can compute the diameter given a solution for APSP, an $\Omega(n/\log n)$ lower bound for APSP is implied. As such, our simple Algorithm 46 is almost optimal!
- Many prominent functions allow for a low communication complexity. For instance, $CC(PARITY) = 2$. What is the Hamming distance (number of different entries) of two strings? It is known that $CC(HAM \geq d) = \Omega(d)$. Also, $CC(\text{decide whether “}HAM \geq k/2 + \sqrt{k}\text{” or “}HAM \leq k/2 - \sqrt{k}\text{”}) = \Omega(k)$, even when using randomness. This problem is known as the Gap-Hamming-Distance.
- Lower bounds in communication complexity have many applications. Apart from getting lower bounds in distributed computing, one can also get lower bounds regarding circuit depth or query times for static data structures.
- In the distributed setting with limited bandwidth we showed that computing the diameter has about the same complexity as computing all pairs shortest paths. In contrast, in sequential computing, it is a major open problem whether the diameter can be computed faster than all pairs shortest paths. No nontrivial lower bounds are known, only that $\Omega(n^2)$ steps are needed – partly due to the fact that there can be n^2 edges/distances in a graph. On the other hand the currently best algorithm uses fast matrix multiplication and terminates after $\mathcal{O}(n^{2.3727})$ steps.

12.4 Distributed Complexity Theory

We conclude this chapter with a short overview on the main complexity classes of distributed message passing algorithms. Given a network with n nodes and diameter D , we managed to establish a rich selection of upper and lower bounds regarding how much time it takes to solve or approximate a problem. Currently we know five main distributed complexity classes:

- Strictly *local* problems can be solved in constant $\mathcal{O}(1)$ time, e.g. a constant approximation of a dominating set in a planar graph.
- Just a little bit slower are problems that can be solved in *log-star* $\mathcal{O}(\log^* n)$ time, e.g. many combinatorial optimization problems in special graph classes such as growth bounded graphs. 3-coloring a ring takes $\mathcal{O}(\log^* n)$.
- A large body of problems is *polylogarithmic* (or *pseudo-local*), in the sense that they seem to be strictly local but are not, as they need $\mathcal{O}(\text{polylog } n)$ time, e.g. the maximal independent set problem.
- There are problems which are *global* and need $\mathcal{O}(D)$ time, e.g. to count the number of nodes in the network.
- Finally there are problems which need *polynomial* $\mathcal{O}(\text{poly } n)$ time, even if the diameter D is a constant, e.g. computing the diameter of the network.

Chapter 13

Stabilization

A large branch of research in distributed computing deals with fault-tolerance. Being able to tolerate a considerable fraction of failing or even maliciously behaving (“Byzantine”) nodes while trying to reach *consensus* (on e.g. the output of a function) among the nodes that work properly is crucial for building reliable systems. However, consensus protocols require that a majority of the nodes remains non-faulty all the time.

Can we design a distributed system that survives transient (short-lived) failures, even if *all* nodes are temporarily failing? In other words, can we build a distributed system that *repairs itself*?

13.1 Self-Stabilization

Definition 13.1 (Self-Stabilization). *A distributed system is self-stabilizing if, starting from an arbitrary state, it is guaranteed to converge to a legitimate state. If the system is in a legitimate state, it is guaranteed to remain there, provided that no further faults happen. A state is legitimate if the state satisfies the specifications of the distributed system.*

Remarks:

- What kind of transient failures can we tolerate? An adversary can crash nodes, or make nodes behave Byzantine. Indeed, temporarily an adversary can do harm in even worse ways, e.g. by corrupting the volatile memory of a node (without the node noticing – not unlike the movie *Memento*), or by corrupting messages on the fly (without anybody noticing). However, as all failures are transient, eventually all nodes must work correctly again, that is, crashed nodes get resurrected, Byzantine nodes stop being malicious, messages are being delivered reliably, and the memory of the nodes is secure.
- Clearly, the read only memory (ROM) must be taboo at all times for the adversary. No system can repair itself if the program code itself or constants are corrupted. The adversary can only corrupt the variables in the volatile random access memory (RAM).

Definition 13.2 (Time Complexity). *The time complexity of a self-stabilizing system is the time that passed after the last (transient) failure until the system has converged to a legitimate state again, staying legitimate.*

Remarks:

- Self-stabilization enables a distributed system to recover from a transient fault regardless of its nature. A self-stabilizing system does not have to be initialized as it eventually (after convergence) will behave correctly.
- One of the first self-stabilizing algorithms was Dijkstra's token ring network. A token ring is an early form of a local area network where nodes are arranged in a ring, communicating by a token. The system is correct if there is exactly one token in the ring. Let's have a look at a simple solution. Given an oriented ring, we simply call the clockwise neighbor parent (p), and the counterclockwise neighbor child (c). Also, there is a leader node v_0 . Every node v is in a state $S(v) \in \{0, 1, \dots, n\}$, perpetually informing its child about its state. The token is implicitly passed on by nodes switching state. Upon noticing a change of the parent state $S(p)$, node v executes the following code:

Algorithm 48 Self-stabilizing Token Ring

```

1: if  $v = v_0$  then
2:   if  $S(v) = S(p)$  then
3:      $S(v) := S(v) + 1 \pmod{n}$ 
4:   end if
5: else
6:    $S(v) := S(p)$ 
7: end if

```

Theorem 13.3. *Algorithm 48 stabilizes correctly.*

Proof: As long as some nodes or edges are faulty, anything can happen. In self-stabilization, we only consider the system after it is correct (at time t_0 , however starting in an arbitrary state).

Every node apart from leader v_0 will always attain the state of its parent. It may happen that one node after the other will learn the current state of the leader. In this case the system stabilizes after the leader increases its state at most n time units after time t_0 . It may however be that the leader increases its state even if the system is not stable, e.g. because its parent or parent's parent accidentally had the same state at time t_0 .

The leader will increase its state possibly multiple times without reaching stability, however, at some point the leader will reach state s , a state that no other node had at time t_0 . (Since there are n nodes and n states, this will eventually happen.) At this point the system must stabilize because the leader cannot push for $s + 1 \pmod{n}$ until every node (including its parent) has s .

After stabilization, there will always be only one node changing its state, i.e., the system remains in a legitimate state. □

Remarks:

- Although one might think the time complexity of the algorithm is quite bad, it is asymptotically optimal.
- It can be a lot of fun designing self-stabilizing algorithms. Let us try to build a system, where the nodes organize themselves as a maximal independent set (MIS, Chapter 8):

Algorithm 49 Self-stabilizing MIS

Require: Node IDs**Every node** v executes the following code:

```

1: do atomically
2:   Leave MIS if a neighbor with a larger ID is in the MIS
3:   Join MIS if no neighbor with larger ID joins MIS
4:   Send (node ID, MIS or not MIS) to all neighbors
5: end do

```

Remarks:

- Note that the main idea of Algorithm 49 is from Algorithm 34, Chapter 8.
- As long as some nodes are faulty, anything can happen: Faulty nodes may for instance decide to join the MIS, but report to their neighbors that they did not join the MIS. Similarly messages may be corrupted during transport. As soon as the system (nodes, messages) is correct, however, the system will converge to a MIS. (The arguments are the same as in Chapter 8).
- Self-stabilizing algorithms always run in an infinite loop, because transient failures can hit the system at any time. Without the infinite loop, an adversary can always corrupt the solution “after” the algorithm terminated.
- The problem of Algorithm 49 is its time complexity, which may be linear in the number of nodes. This is not very exciting. We need something better! Since Algorithm 49 was just the self-stabilizing variant of the slow MIS Algorithm 34, maybe we can hope to “self-stabilize” some of our fast algorithms from Chapter 8?
- Yes, we can! Indeed there is a general transformation that takes any local algorithm (efficient but not fault-tolerant) and turns it into a self-stabilizing algorithm, keeping the same level of efficiency and efficacy. We present the general transformation below.

Theorem 13.4 (Transformation). *We are given a deterministic local algorithm A that computes a solution of a given problem in k synchronous communication rounds. Using our transformation, we get a self-stabilizing system with time complexity k . In other words, if the adversary does not corrupt the system for k time units, the solution is stable. In addition, if the adversary does not corrupt any node or message closer than distance k from a node u , node u will be stable.*

Proof: In the proof, we present the transformation. First, however, we need to be more formal about the deterministic local algorithm \mathcal{A} . In \mathcal{A} , each node of the network computes its decision in k phases. In phase i , node u computes its local variables according to its local variables and received messages of the earlier phases. Then node u sends its messages of phase i to its neighbors. Finally node u receives the messages of phase i from its neighbors. The set of local variables of node u in phase i is given by L_u^i . (In the very first phase, node u initializes its local variables with L_u^1 .) The message sent from node u to node v in phase i is denoted by $m_{u,v}^i$. Since the algorithm \mathcal{A} is deterministic, node u can compute its local variables L_u^i and messages $m_{u,*}^i$ of phase i from its state of earlier phases, by simply applying functions f_L and f_m . In particular,

$$L_u^i = f_L(u, L_u^{i-1}, m_{*,u}^{i-1}), \text{ for } i > 1, \text{ and} \quad (13.1)$$

$$m_{u,v}^i = f_m(u, v, L_u^i), \text{ for } i \geq 1. \quad (13.2)$$

The self-stabilizing algorithm needs to simulate all the k phases of the local algorithm \mathcal{A} in parallel. Each node u stores its local variables L_u^1, \dots, L_u^k as well as all messages received $m_{*,u}^1, \dots, m_{*,u}^k$ in two tables in RAM. For simplicity, each node u also stores all the sent messages $m_{u,*}^1, \dots, m_{u,*}^k$ in a third table. If a message or a local variable for a particular phase is unknown, the entry in the table will be marked with a special value \perp (“unknown”). Initially, all entries in the table are \perp .

Clearly, in the self-stabilizing model, an adversary can choose to change table values at all times, and even reset these values to \perp . Our self-stabilizing algorithm needs to constantly work against this adversary. In particular, each node u runs these two procedures constantly:

- For all neighbors: Send each neighbor v a message containing the complete row of messages of algorithm \mathcal{A} , that is, send the vector $(m_{u,v}^1, \dots, m_{u,v}^k)$ to neighbor v . Similarly, if neighbor u receives such a vector from neighbor v , then neighbor u replaces neighbor v 's row in the table of incoming messages by the received vector $(m_{v,u}^1, \dots, m_{v,u}^k)$.
- Because of the adversary, node u must constantly recompute its local variables (including the initialization) and outgoing message vectors using Functions (13.1) and (13.2) respectively.

The proof is by induction. Let $N^i(u)$ be the i -neighborhood of node u (that is, all nodes within distance i of node u). We assume that the adversary has not corrupted any node in $N^k(u)$ since time t_0 . At time t_0 all nodes in $N^k(u)$ will check and correct their initialization. Following Equation (13.2), at time t_0 all nodes in $N^k(u)$ will send the correct message entry for the first round ($m_{*,*}^1$) to all neighbors. Asynchronous messages take at most 1 time unit to be received at a destination. Hence, using the induction with Equations (13.1) and (13.2) it follows that at time $t_0 + i$, all nodes in $N^{k-i}(u)$ have received the correct messages $m_{*,*}^1, \dots, m_{*,*}^i$. Consequently, at time $t_0 + k$ node u has received all messages of local algorithm \mathcal{A} correctly, and will compute the same result value as in \mathcal{A} . \square

Remarks:

- Using our transformation (also known as “local checking”), designing self-stabilizing algorithms just turned from art to craft.
- As we have seen, many local algorithms are randomized. This brings two additional problems. Firstly, one may not exactly know how long the algorithm will take. This is not really a problem since we can simply send around all the messages needed, until the algorithm is finished. The transformation of Theorem 13.4 works also if nodes just send all messages that are not \perp . Secondly, we must be careful about the adversary. In particular we need to restrict the adversary such that a node can produce a reproducible sufficiently long string of random bits. This can be achieved by storing the sufficiently long string along with the program code in the read only memory (ROM). Alternatively, the algorithm might not store the random bit string in its ROM, but only the seed for a random bit generator. We need this in order to keep the adversary from reshuffling random bits until the bits become “bad”, and the expected (or with high probability) efficacy or efficiency guarantees of the original local algorithm \mathcal{A} cannot be guaranteed anymore.
- Since most local algorithms have only a few communication rounds, and only exchange small messages, the memory overhead of the transformation is usually bearable. In addition, information can often be compressed in a suitable way so that for many algorithms message size will remain polylogarithmic. For example, the information of the fast MIS algorithm (Algorithm 36) consists of a series of random values (one for each round), plus two boolean values per round. These boolean values represent whether the node joins the MIS, or whether a neighbor of the node joins the MIS. The order of the values tells in which round a decision is made. Indeed, the series of random bits can even be compressed just into the random seed value, and the neighbors can compute the random values of each round themselves.
- There is hope that our transformation as well gives good algorithms for mobile networks, that is for networks where the topology of the network may change. Indeed, for deterministic local approximation algorithms, this is true: If the adversary does not change the topology of a node’s k -neighborhood in time k , the solution will locally be stable again.
- For randomized local approximation algorithms however, this is not that simple. Assume for example, that we have a randomized local algorithm for the dominating set problem. An adversary can constantly switch the topology of the network, until it finds a topology for which the random bits (which are not really random because these random bits are in ROM) give a solution with a bad approximation ratio. By defining a weaker adversarial model, we can fix this problem. Essentially, the adversary needs to be oblivious, in the sense that it cannot see the solution. Then it will not be possible for the adversary to restart the random computation if the solution is “too good”.

- Self-stabilization is the original approach, and self-organization may be the general theme, but new buzzwords pop up every now and then, e.g. self-configuration, self-management, self-regulation, self-repairing, self-healing, self-optimization, self-adaptivity, or self-protection. Generally all these are summarized as “self-*”. One computing giant coined the term “autonomic computing” to reflect the trend of self-managing distributed systems.

13.2 Advanced Stabilization

We finish the chapter with a non-trivial example beyond self-stabilization, showing the beauty and potential of the area: In a small town, every evening each citizen calls all his (or her) friends, asking them whether they will vote for the Democratic or the Republican party at the next election.¹ In our town citizens listen to their friends, and everybody re-chooses his or her affiliation according to the majority of friends.² Is this process going to “stabilize” (in one way or another)?

Remarks:

- Is eventually everybody voting for the same party? No.
- Will each citizen eventually stay with the same party? No.
- Will citizens that stayed with the same party for some time, stay with that party forever? No.
- And if their friends also constantly root for the same party? No.
- Will this beast stabilize at all?!? Yes!

Theorem 13.5 (Dems & Reps). *Eventually every citizen is rooting for the same party every other day.*

Proof: To prove that the opinions eventually become fixed or cycle every other day, think of each friendship between citizens as a pair of (directed) edges, one in each direction. Let us say an edge is currently “bad” if the party of the *advising* friend differs from the next-day’s party of the *advised* friend. In other words, the edge is bad if the advised friend did not follow the advisor’s opinion (which means that the advisor was in the minority). An edge that is not bad, is “good”.

Consider the out-edges of citizen c on day t , during which (say) c roots for the Democrats. Assume that during day t , g out-edges of c are good, and b out-edges are bad. Note that $g + b$ is the degree of c . Since g out-edges were good, g friends of c root for the Democrats on day $t + 1$. Likewise, b friends of c root for the Republicans on day $t + 1$. In other words, on the evening of day $t + 1$ citizen c will receive g recommendations for Democrats, and b for Republicans. We distinguish two cases:

¹We are in the US, and as we know from The Simpsons, you “throw your vote away” if you vote for somebody else. As a consequence our example has two parties only.

²Assume for the sake of simplicity that everybody has an odd number of friends.

- $g > b$: In this case, citizen c will still (or again) root for the Democrats on day $t + 2$. Note that in this case, on day $t + 1$, exactly g in-edges of c are good, and exactly b in-edges are bad. In other words, the number of bad out-edges on day t is exactly the number of bad in-edges on day $t + 1$.
- $g < b$: In this case, citizen c will root for the Republicans on day $t + 2$. Note that in this case, on day $t + 1$, exactly b in-edges of c are good, and exactly g in-edges are bad. In other words, the number of bad out-edges on day t was exactly the number of good in-edges on day $t + 1$ (and vice versa). Since citizen c is rooting for the Republicans, the number of bad out-edges on day t was strictly larger than the number of bad in-edges on day $t + 1$.

We account for every edge as out-edge on day t , and as in-edge on day $t + 1$. Since in both of the above cases the number of bad edges does not increase, the total number of bad edges B cannot increase. In fact, if any node switches its party from day t to $t + 2$, we know that the total number of bad edges strictly decreases. But B cannot decrease forever. Once B hits its minimum, the system stabilizes in the sense that every citizen will either stick with his or her party forever or flip-flop every day – the system “stabilizes”. \square

Remarks:

- The model can be generalized considerably by, for example, adding weights to vertices (meaning some citizens’ opinions are more important than others), allowing loops (citizens who consider their own current opinions as well), allowing tie-breaking mechanisms, and even allowing different thresholds for party changes.
- How long does it take until the system stabilizes?
- Some of you may be reminded of Conway’s Game of Life: We are given an infinite two-dimensional grid of cells, each of which is in one of two possible states, *dead* or *alive*. Every cell interacts with its eight neighbors. In each round, the following transitions occur: Any live cell with fewer than two live neighbors dies, as if caused by loneliness. Any live cell with more than three live neighbors dies, as if by overcrowding. Any live cell with two or three live neighbors lives on to the next generation. Any dead cell with exactly three live neighbors is “born” and becomes a live cell. The initial pattern constitutes the “seed” of the system. The first generation is created by applying the above rules simultaneously to every cell in the seed, births and deaths happen simultaneously, and the discrete moment at which this happens is sometimes called a tick. (In other words, each generation is a pure function of the one before.) The rules continue to be applied repeatedly to create further generations. John Conway figured that these rules were enough to generate interesting situations, including “breeders” which create “guns” which in turn create “gliders”. As such Life in some sense answers an old question by John von Neumann, whether there can be a simple machine that can build copies of itself. In fact Life is Turing complete, that is, as powerful as any computer.

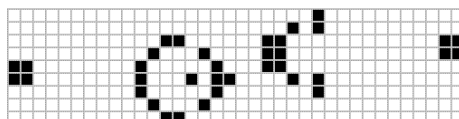


Figure 13.1: A “glider gun”...



Figure 13.2: ...in action.

Chapter 14

Wireless Protocols

Wireless communication was one of the major success stories of the last decades. Today, different wireless standards such as wireless local area networks (WLAN) are omnipresent. In some sense, from a distributed computing viewpoint wireless networks are quite simple, as they cannot form arbitrary network topologies. Simplistic models of wireless networks include geometric graph models such as the so-called unit disk graph. Modern models are more robust: The network graph is restricted, e.g., the total number of neighbors of a node which are not adjacent is likely to be small. This observation is hard to capture with purely geometric models, and motivates more advanced network connectivity models such as bounded growth or bounded independence.

However, on the other hand, wireless communication is also more difficult than standard message passing, as for instance nodes are not able to transmit a different message to each neighbor at the same time. And if two neighbors are transmitting at the same time, they interfere, and a node may not be able to decipher anything.

In this chapter we deal with the distributed computing principles of wireless communication: We make the simplifying assumption that all n nodes are in the communication range of each other, i.e., the network graph is a clique. Nodes share a synchronous time, in each time slot a node can decide to either transmit or receive (or sleep). However, two or more nodes transmitting in a time slot will cause interference. Transmitting nodes are never aware if there is interference because they cannot simultaneously transmit and receive.

14.1 Basics

The basic communication protocol in wireless networks is the medium access control (MAC) protocol. Unfortunately it is difficult to claim that one MAC protocol is better than another, because it all depends on the parameters, such as the network topology, the channel characteristics, or the traffic pattern. When it comes to the principles of wireless protocols, we usually want to achieve much simpler goals. One basic and important question is the following: How long does it take until one node can transmit successfully, without interference? This question is often called the wireless leader election problem (Chapter 2), with the node transmitting alone being the leader.

Clearly, we can use node IDs to solve leader election, e.g., a node with ID i transmits in time slot i . However, this may be incredibly slow. There are better deterministic solutions, but by and large the best and simplest algorithms are randomized.

Throughout this chapter, we use a random variable X to denote the number of nodes transmitting in a given slot.

Algorithm 50 Slotted Aloha

- 1: **Every node** v executes the following code:
 - 2: **repeat**
 - 3: transmit with probability $1/n$
 - 4: **until** one node has transmitted alone
-

Theorem 14.1. *Using Algorithm 50 allows one node to transmit alone (become a leader) after expected time e .*

Proof. The probability for success, i.e., only one node transmitting is

$$Pr[X = 1] = n \cdot \frac{1}{n} \cdot \left(1 - \frac{1}{n}\right)^{n-1} \approx \frac{1}{e},$$

where the last approximation is a result from Theorem 14.23 for sufficiently large n . Hence, if we repeat this process e times, we can expect one success. \square

Remarks:

- The origin of the name is the ALOHAnet which was developed at the University of Hawaii.
- How does the leader know that it is the leader? One simple solution is a “distributed acknowledgment”. The nodes just continue Algorithm 50, including the ID of the the leader in their transmission. So the leader learns that is the leader.
- One more problem?! Indeed, node v which managed to transmit the acknowledgment (alone) is the only remaining node which does not know that the leader knows that it is the leader. We can fix this by having the leader acknowledge v ’s successful acknowledgment.
- One can also imagine an unslotted time model. In this model two messages which overlap partially will interfere and no message is received. As everything in this chapter, Algorithm 50 also works in an unslotted time model, with a factor 2 penalty, i.e., the probability for a successful transmission will drop from $\frac{1}{e}$ to $\frac{1}{2e}$. Essentially, each slot is divided into t small time slots with $t \rightarrow \infty$ and the nodes start a new t -slot long transmission with probability $\frac{1}{2nt}$.

14.2 Initialization

Sometimes we want the n nodes to have the IDs $\{1, 2, \dots, n\}$. This process is called initialization. Initialization can for instance be used to allow the nodes to transmit one by one without any interference.

14.2.1 Non-Uniform Initialization

Theorem 14.2. *If the nodes know n , we can initialize them in $\mathcal{O}(n)$ time slots.*

Proof. We repeatedly elect a leader using e.g., Algorithm 50. The leader gets the next free number and afterwards leaves the process. We know that this works with probability $1/e$. The expected time to finish is hence $e \cdot n$. \square

Remarks:

- But this algorithm requires that the nodes know n in order to give them IDs from $1, \dots, n$. For a more realistic scenario we need a uniform algorithm, i.e., the nodes do not know n .

14.2.2 Uniform Initialization with CD

Definition 14.3 (Collision Detection, CD). *Two or more nodes transmitting concurrently is called interference. In a system with collision detection, a receiver can distinguish interference from nobody transmitting. In a system without collision detection, a receiver cannot distinguish the two cases.*

Let us first present a high-level idea. The set of nodes is recursively partitioned into two non-empty sets, similarly to a binary tree. This is repeated recursively until a set contains only one node which gets the next free ID. Afterwards, the algorithm continues with the next set.

Algorithm 51 RandomizedSplit(b)

```

1: Every node  $v$  executes the following code:
2: repeat
3:   if  $b_v = b$  then
4:     choose  $r$  uniformly at random from  $\{0, 1\}$ 
5:     in the next two time slots:
6:       transmit in slot  $r$ , and listen in other slot
7:   end if
8: until there was at least 1 transmission in both slots
9: if  $b_v = b$  then
10:   $b_v := b_v + r$  {append bit  $r$  to bitstring  $b_v$ }
11: end if
12: if some node  $u$  transmitted alone in slot  $r \in \{0, 1\}$  then
13:  node  $u$  gets ID  $m$  {and becomes passive}
14:   $m := m + 1$ 
15: else
16:  RandomizedSplit( $b + r$ )
17: end if

```

Remarks:

- In line 8 the transmitting nodes need to know if they were the only one transmitting. Since we have enough time, we can do a leader election first and use a similar trick as before to ensure this.
- In line 12 we check separately for $r = 0$ and $r = 1$

Algorithm 52 Initialization with Collision Detection

-
- 1: **Every node** v executes the following code:
 - 2: global variable $m := 0$ {number of already identified nodes}
 - 3: local variable $b_v := ''$ {current bitstring of node v , initially empty}
 - 4: RandomizedSplit('')
-

Theorem 14.4. *Algorithm 52 correctly initializes the set of nodes in $\mathcal{O}(n)$.*

Proof. A successful split is defined as a split in which both subsets are non-empty. We know that there are exactly $n - 1$ successful splits because we have a binary tree with n leaves and $n - 1$ inner nodes. Let us now calculate the probability for creating two non-empty sets from a set of size $k \geq 2$ as

$$Pr[1 \leq X \leq k - 1] = 1 - Pr[X = 0] - Pr[X = k] = 1 - \frac{1}{2^k} - \frac{1}{2^k} \geq \frac{1}{2}.$$

Thus, in expectation we need $\mathcal{O}(n)$ splits. \square

Remarks:

- What if we do not have collision detection?

14.2.3 Uniform Initialization without CD

Let us assume that we have a special node ℓ (leader) and let S denote the set of nodes which want to transmit. We now split every time slot from before into two time slots and use the leader to help us distinguish between silence and noise. In the first slot every node from the set S transmits, in the second slot the nodes in $S \cup \{\ell\}$ transmit. This gives the nodes sufficient information to distinguish the different cases (see Table 14.1).

	nodes in S transmit	nodes in $S \cup \{\ell\}$ transmit
$ S = 0$	no	yes
$ S = 1, S = \{\ell\}$	yes	yes
$ S = 1, S \neq \{\ell\}$	yes	no
$ S \geq 2$	no	no

Table 14.1: Using a leader to distinguish between noise and silence: *no* represents noise/silence, *yes* represents a successful transmission.

Remarks:

- As such, Algorithm 52 works also without CD, with only a factor 2 overhead.
- More generally, a leader immediately brings CD to any protocol.
- This protocol has an important real life application, for instance when checking out a shopping cart with items which have RFID tags.
- But how do we determine such a leader? And how long does it take until we are “sure” that we have one? Let us repeat the notion of *with high probability*.

14.3 Leader Election

14.3.1 With High Probability

Definition 14.5 (With High Probability). *Some probabilistic event is said to occur with high probability (w.h.p.), if it happens with a probability $p \geq 1 - 1/n^c$, where c is a constant. The constant c may be chosen arbitrarily, but it is considered constant with respect to Big-O notation.*

Theorem 14.6. *Algorithm 50 elects a leader w.h.p. in $\mathcal{O}(\log n)$ time slots.*

Proof. The probability for not electing a leader after $c \cdot \log n$ time slots, i.e., $c \log n$ slots without a successful transmission is

$$\left(1 - \frac{1}{e}\right)^{c \ln n} = \left(1 - \frac{1}{e}\right)^{e \cdot c' \ln n} \leq \frac{1}{e^{\ln n \cdot c'}} = \frac{1}{n^{c'}}.$$

□

Remarks:

- What about uniform algorithms, i.e. the number of nodes n is not known?

14.3.2 Uniform Leader Election

Algorithm 53 Uniform leader election

```

1: Every node  $v$  executes the following code:
2: for  $k = 1, 2, 3, \dots$  do
3:   for  $i = 1$  to  $ck$  do
4:     transmit with probability  $p := 1/2^k$ 
5:     if node  $v$  was the only node which transmitted then
6:        $v$  becomes the leader
7:       break
8:     end if
9:   end for
10: end for

```

Theorem 14.7. *By using Algorithm 53 it is possible to elect a leader w.h.p. in $\mathcal{O}(\log^2 n)$ time slots if n is not known.*

Proof. Let us briefly describe the algorithm. The nodes transmit with probability $p = 2^{-k}$ for ck time slots for $k = 1, 2, \dots$. At first p will be too high and hence there will be a lot of interference. But after $\log n$ phases, we have $k \approx \log n$ and thus the nodes transmit with probability $\approx \frac{1}{n}$. For simplicity's sake, let us assume that n is a power of 2. Using the approach outlined above, we know that after $\log n$ iterations, we have $p = \frac{1}{n}$. Theorem 14.6 yields that we can elect a leader w.h.p. in $\mathcal{O}(\log n)$ slots. Since we have to try $\log n$ estimates until $k \approx n$, the total runtime is $\mathcal{O}(\log^2 n)$. \square

Remarks:

- Note that our proposed algorithm has not used collision detection. Can we solve leader election faster in a uniform setting with collision detection?

14.3.3 Fast Leader Election with CD

Algorithm 54 Uniform leader election with CD

```

1: Every node  $v$  executes the following code:
2: repeat
3:   transmit with probability  $\frac{1}{2}$ 
4:   if at least one node transmitted then
5:     all nodes that did not transmit quit the protocol
6:   end if
7: until one node transmits alone

```

Theorem 14.8. *With collision detection we can elect a leader using Algorithm 54 w.h.p. in $\mathcal{O}(\log n)$ time slots.*

Proof. The number of active nodes k is monotonically decreasing and always greater than 1 which yields the correctness. A slot is called successful if at most half the active nodes transmit. We can assume that $k \geq 2$ since otherwise we would have already elected a leader. We can calculate the probability that a time slot is successful as

$$\Pr[1 \leq X \leq \lceil \frac{k}{2} \rceil] \geq \frac{1}{2} - \Pr[X = 0] = \frac{1}{2} - \frac{1}{2^k} \geq \frac{1}{4}.$$

Since the number of active nodes at least halves in every successful time slot, $\log n$ successful time slots are sufficient to elect a leader. Now let Y be a random variable which counts the number of successful time slots after $8 \cdot c \cdot \log n$ time slots. The expected value is $E[Y] \geq 8 \cdot c \cdot \log n \cdot \frac{1}{4} \geq 2 \cdot \log n$. Since all those time slots are independent from each other, we can apply a Chernoff bound (see Theorem 14.22) with $\delta = \frac{1}{2}$ which states

$$\Pr[Y < (1 - \delta)E[Y]] \leq e^{-\frac{\delta^2}{2}E[Y]} = e^{-\frac{1}{8} \cdot 2c \log n} \leq n^{-\alpha}$$

for any constant α . \square

Remarks:

- Can we be even faster?

14.3.4 Even Faster Leader Election with CD

Let us first briefly describe an algorithm for this. In the first phase the nodes transmit with probability $1/2^{2^0}, 1/2^{2^1}, 1/2^{2^2}, \dots$ until no node transmits. This yields a first approximation on the number of nodes. Afterwards, a binary search is performed to determine an even better approximation of n . Finally, the third phase finds a constant approximation of n using a biased random walk. The algorithm stops in any case as soon as only one node is transmitting which will become the leader.

Algorithm 55 Fast uniform leader election

```

1:  $i := 1$ 
2: repeat
3:    $i := 2 \cdot i$ 
4:   transmit with probability  $1/2^i$ 
5: until no node transmitted
   {End of Phase 1}
6:  $l := 2^{i-2}$ 
7:  $u := 2^i$ 
8: while  $l + 1 < u$  do
9:    $j := \lceil \frac{l+u}{2} \rceil$ 
10:  transmit with probability  $1/2^j$ 
11:  if no node transmitted then
12:     $u := j$ 
13:  else
14:     $l := j$ 
15:  end if
16: end while
   {End of Phase 2}
17:  $k := u$ 
18: repeat
19:  transmit with probability  $1/2^k$ 
20:  if no node transmitted then
21:     $k := k - 1$ 
22:  else
23:     $k := k + 1$ 
24:  end if
25: until exactly one node transmitted

```

Lemma 14.9. *If $j > \log n + \log \log n$, then $\Pr[X > 1] \leq \frac{1}{\log n}$.*

Proof. The nodes transmit with probability $1/2^j < 1/2^{\log n + \log \log n} = \frac{1}{n \log n}$. The expected number of nodes transmitting is $E[X] = \frac{n}{n \log n}$. Using Markov's inequality (see Theorem 14.21) yields $\Pr[X > 1] \leq \Pr[X > E[X] \cdot \log n] \leq \frac{1}{\log n}$. \square

Lemma 14.10. *If $j < \log n - \log \log n$, then $P[X = 0] \leq \frac{1}{n}$.*

Proof. The nodes transmit with probability $1/2^j < 1/2^{\log n - \log \log n} = \frac{\log n}{n}$. Hence, the probability for a silent time slot is $(1 - \frac{\log n}{n})^n = e^{-\log n} = \frac{1}{n}$. \square

Corollary 14.11. *If $i > 2 \log n$, then $Pr[X > 1] \leq \frac{1}{\log n}$.*

Proof. This follows from Lemma 14.9 since the deviation in this corollary is even larger. \square

Corollary 14.12. *If $i < \frac{1}{2} \log n$, then $P[X = 0] \leq \frac{1}{n}$.*

Proof. This follows from Lemma 14.10 since the deviation in this corollary is even larger. \square

Lemma 14.13. *Let v be such that $2^{v-1} < n \leq 2^v$, i.e., $v \approx \log n$. If $k > v + 2$, then $Pr[X > 1] \leq \frac{1}{4}$.*

Proof. Markov's inequality yields

$$Pr[X > 1] = Pr\left[X > \frac{2^k}{n} E[X]\right] < Pr[X > \frac{2^k}{2^v} E[X]] < Pr[X > 4E[X]] < \frac{1}{4}.$$

\square

Lemma 14.14. *If $k < v - 2$, then $P[X = 0] \leq \frac{1}{4}$.*

Proof. A similar analysis is possible to upper bound the probability that a transmission fails if our estimate is too small. We know that $k \leq v - 2$ and thus

$$Pr[X = 0] = \left(1 - \frac{1}{2^k}\right)^n < e^{-\frac{n}{2^k}} < e^{-\frac{2^{v-1}}{2^k}} < e^{-2} < \frac{1}{4}.$$

\square

Lemma 14.15. *If $v - 2 \leq k \leq v + 2$, then the probability that exactly one node transmits is constant.*

Proof. The transmission probability is $p = \frac{1}{2^{v \pm \Theta(1)}} = \Theta(1/n)$, and the lemma follows with a slightly adapted version of Theorem 14.1. \square

Lemma 14.16. *With probability $1 - \frac{1}{\log n}$ we find a leader in phase 3 in $\mathcal{O}(\log \log n)$ time.*

Proof. For any k , because of Lemmas 14.13 and 14.14, the random walk of the third phase is biased towards the good area. One can show that in $\mathcal{O}(\log \log n)$ steps one gets $\Omega(\log \log n)$ good transmissions. Let Y denote the number of times exactly one node transmitted. With Lemma 14.15 we obtain $E[Y] = \Omega(\log \log n)$. Now a direct application of a Chernoff bound (see Theorem 14.22) yields that these transmissions elect a leader with probability $1 - \frac{1}{\log n}$. \square

Theorem 14.17. *The Algorithm 55 elects a leader with probability of at least $1 - \frac{\log \log n}{\log n}$ in time $\mathcal{O}(\log \log n)$.*

Proof. From Corollary 14.11 we know that after $\mathcal{O}(\log \log n)$ time slots, the first phase terminates. Since we perform a binary search on an interval of size $\mathcal{O}(\log n)$, the second phase also takes at most $\mathcal{O}(\log \log n)$ time slots. For the third phase we know that $\mathcal{O}(\log \log n)$ slots are sufficient to elect a leader with probability $1 - \frac{1}{\log n}$ by Lemma 14.16. Thus, the total runtime is $\mathcal{O}(\log \log n)$.

Now we can combine the results. We know that the error probability for every time slot in the first two phases is at most $\frac{1}{\log n}$. Using a union bound (see Theorem 14.20), we can upper bound the probability that no error occurred by $\frac{\log \log n}{\log n}$. Thus, we know that after phase 2 our estimate is at most $\log \log n$ away from $\log n$ with probability of at least $1 - \frac{\log \log n}{\log n}$. Hence, we can apply Lemma 14.16 and thus successfully elect a leader with probability of at least $1 - \frac{\log \log n}{\log n}$ (again using a union bound) in time $\mathcal{O}(\log \log n)$. \square

Remarks:

- Tightening this analysis a bit more, one can elect a leader with probability $1 - \frac{1}{\log n}$ in time $\log \log n + o(\log \log n)$.
- Can we be even faster?

14.3.5 Lower Bound

Theorem 14.18. *Any uniform protocol that elects a leader with probability of at least $1 - \frac{1}{\log n}$ must run for at least $\log \log n$ time slots.*

Proof. The probability that exactly one node transmits is

$$\Pr[X = 1] = n \cdot p \cdot (1 - p)^{n-1}.$$

Consider now a system with only 2 nodes. The probability that exactly one transmits is at most

$$\Pr[X = 1] = p \cdot (1 - p) \leq \frac{1}{2}.$$

Thus, after $\log \log n$ time slots the probability that a leader was elected is at most $1 - \frac{1}{2}^{\log \log n} = 1 - \frac{1}{\log n}$. \square

14.3.6 Uniform Asynchronous Wakeup without CD

Until now we have assumed that all nodes start the algorithm in the same time slot. But what happens if this is not the case? How long does it take to elect a leader if we want a uniform and anonymous (nodes do not have an identifier and thus cannot base their decision on it) algorithm?

Theorem 14.19. *If nodes wake up in an arbitrary (worst-case) way, any algorithm may take $\Omega(n/\log n)$ time slots until a single node can successfully transmit.*

Proof. Nodes must transmit at some point, or they will surely never successfully transmit. With a uniform protocol, every node executes the same code. We focus on the first slot where nodes may transmit. No matter what the protocol is, this happens with probability p . Since the protocol is uniform, p must be a constant, independent of n .

The adversary wakes up $w = \frac{c}{p} \ln n$ nodes in each time slot with some constant c . All nodes woken up in the first time slot will transmit with probability p . We study the event E_1 that exactly one of them transmits in that first time slot. Using the inequality $(1 + t/n)^n \leq e^t$ from Lemma 14.23 we get

$$\begin{aligned} \Pr[E_1] &= w \cdot p \cdot (1 - p)^{w-1} \\ &= c \ln n (1 - p)^{\frac{1}{p}(c \ln n - p)} \\ &\leq c \ln n \cdot e^{-c \ln n + p} \\ &= c \ln n \cdot n^{-c} e^p \\ &= n^{-c} \cdot \mathcal{O}(\log n) \\ &< \frac{1}{n^{c-1}} = \frac{1}{n^{c'}}. \end{aligned}$$

In other words, w.h.p. that time slot will not be successful. Since the nodes cannot distinguish noise from silence, the same argument applies to every set of nodes which wakes up. Let E_α be the event that all n/w time slots will not be successful. Using the inequality $1 - p \leq (1 - p/k)^k$ from Lemma 14.24 we get

$$\Pr[E_\alpha] = (1 - \Pr[E_1])^{n/w} > \left(1 - \frac{1}{n^{c'}}\right)^{\Theta(n/\log n)} > 1 - \frac{1}{n^{c''}}.$$

In other words, w.h.p. it takes more than n/w time slots until some node can transmit alone. □

14.4 Useful Formulas

In this chapter we have used several inequalities in our proofs. For simplicity's sake we list all of them in this section.

Theorem 14.20. *Boole's inequality or union bound: For a countable set of events E_1, E_2, E_3, \dots , we have*

$$\Pr\left[\bigcup_i E_i\right] \leq \sum_i \Pr[E_i].$$

Theorem 14.21. *Markov's inequality: If X is any random variable and $a > 0$, then*

$$\Pr[|X| \geq a] \leq \frac{E[X]}{a}.$$

Theorem 14.22. *Chernoff bound: Let Y_1, \dots, Y_n be a independent Bernoulli random variables let $Y := \sum_i Y_i$. For any $0 \leq \delta \leq 1$ it holds*

$$\Pr[Y < (1 - \delta)E[Y]] \leq e^{-\frac{\delta^2}{2} E[Y]}$$

and for $\delta > 0$

$$\Pr[Y \geq (1 + \delta) \cdot E[Y]] \leq e^{-\frac{\min\{\delta, \delta^2\}}{3} \cdot E[Y]}$$

Theorem 14.23. *We have*

$$e^t \left(1 - \frac{t^2}{n}\right) \leq \left(1 + \frac{t}{n}\right)^n \leq e^t$$

for all $n \in \mathbb{N}$, $|t| \leq n$. Note that

$$\lim_{n \rightarrow \infty} \left(1 + \frac{t}{n}\right)^n = e^t.$$

Theorem 14.24. *For all p, k such that $0 < p < 1$ and $k \geq 1$ we have*

$$1 - p \leq (1 - p/k)^k.$$

Chapter 15

All-to-All Communication

In the previous chapters, we have mostly considered communication on a particular graph $G = (V, E)$, where any two nodes u and v can only communicate directly if $\{u, v\} \in E$. This is however not always the best way to model a network. In the Internet, for example, every machine (node) is able to “directly” communicate with every other machine via a series of routers. If every node in a network can communicate directly with all other nodes, many problems can be solved easily. For example, assume we have n servers, each hosting an arbitrary number of (numeric) elements. If all servers are interested in obtaining the maximum of all elements, all servers can simultaneously, i.e., in one communication round, send their local maximum element to all other servers. Once these maxima are received, each server knows the global maximum.

Note that we can again use graph theory to model this *all-to-all* communication scenario: The communication graph is simply the complete graph $K_n := (V, \binom{V}{2})$. If each node can send its entire local state in a single message, then all problems could be solved in 1 communication round in this model! Since allowing unbounded messages is not realistic in most practical scenarios, we restrict the message size: Assuming that all node identifiers and all other variables in the system (such as the numeric elements in the example above) can be described using $\mathcal{O}(\log n)$ bits, each node can only send a message of size $\mathcal{O}(\log n)$ bits to all other nodes. In other words, only a constant number of identifiers (and elements) can be packed into a single message. Thus, in this model, the limiting factor is the amount of information that can be transmitted in a fixed amount of time. This is fundamentally different from the model we studied before where nodes are restricted to local information about the network graph.

In this chapter, we study one particular problem in this model, the computation of a minimum spanning tree (MST), i.e., we will again look at the construction of a basic network structure. Let us first review the definition of a minimum spanning tree from Chapter 4. We assume that each edge e is assigned a weight ω_e .

Definition 15.1 (MST). *Given a weighted graph $G = (V, E, \omega)$. The MST of G is a spanning tree T minimizing $\omega(T)$, where $\omega(H) = \sum_{e \in H} \omega_e$ for any subgraph $H \subseteq G$.*

Remarks:

- Since we have a complete communication graph, the graph has $\binom{n}{2}$ edges in the beginning.
- As in Chapter 4, we assume that no two edges of the graph have the same weight. Recall that assumption ensures that the MST is unique. Recall also that this simplification is not essential as one can always break ties by using the IDs of adjacent vertices.

For simplicity, we assume that we have a synchronous model (as we are only interested in the time complexity, our algorithm can be made asynchronous using synchronizer α at no additional cost (cf. Chapter 11). As usual, in every round, every node can send a (potentially different) message to each of its neighbors. In particular, note that the message delay is 1 for every edge e independent of the weight ω_e . As mentioned before, every message can contain a constant number of node IDs and edge weights (and $\mathcal{O}(\log n)$ additional bits).

There is a considerable amount of work on distributed MST construction. Table 15.1 lists the most important results for various network diameters D . As we have a complete communication network in our model, we focus only on $D = 1$.

Upper Bounds		
Graph Class	Time Complexity	Authors
General Graphs	$\mathcal{O}(D + \sqrt{n} \cdot \log^* n)$	Kutten, Peleg
Diameter 2	$\mathcal{O}(\log n)$	Lotker, Patt-Shamir, Peleg
Diameter 1	$\mathcal{O}(\log \log n)$	Lotker, Patt-Shamir, Pavlov, Peleg

Lower Bounds		
Graph Class	Time Complexity	Authors
Diameter $\Omega(\log n)$	$\Omega(D + \sqrt{n} / \log^2 n)$	Peleg, Rubinovich
Diameter 4	$\Omega(n^{1/3} / \sqrt{\log n})$	Lotker, Patt-Shamir, Peleg
Diameter 3	$\Omega(n^{1/4} / \sqrt{\log n})$	Lotker, Patt-Shamir, Peleg

Table 15.1: Time complexity of distributed MST construction

Remarks:

- Note that for graphs of arbitrary diameter D , if there are no bounds on the number of messages sent, on the message size, and on the amount of local computations, there is a straightforward generic algorithm to compute an MST in time D : In every round, every node sends its complete state to all its neighbors. After D rounds, every node knows the whole graph and can compute any graph structure locally without any further communication.

- In general, the diameter D is also an obvious lower bound for the time needed to compute an MST. In a weighted ring, e.g., it takes time D to find the heaviest edge. In fact, on the ring, time D is required to compute any spanning tree.

In this chapter, we are not concerned with lower bounds, we want to give an algorithm that computes the MST as quickly as possible instead! We again use the following lemma that is proven in Chapter 4.

Lemma 15.2. *For a given graph G let T be an MST, and let $T' \subseteq T$ be a subgraph (also known as a fragment) of the MST. Edge $e = (u, v)$ is an outgoing edge of T' if $u \in T'$ and $v \notin T'$ (or vice versa). Let the minimum weight outgoing edge of the fragment T' be the so-called blue edge $b(T')$. Then $T' \cup b(T') \subseteq T$.*

Lemma 15.2 leads to a straightforward distributed MST algorithm. We start with an empty graph, i.e., every node is a fragment of the MST. The algorithm consists of phases. In every phase, we add the blue edge $b(T')$ of every existing fragment T' to the MST. Algorithm 56 shows how the described simple MST construction can be carried out in a network of diameter 1.

Algorithm 56 Simple MST Construction (at node v)

```

1: // all nodes always know all current MST edges and thus all MST fragments
2: while  $v$  has neighbor  $u$  in different fragment do
3:   find lowest-weight edge  $e$  between  $v$  and a node  $u$  in a different fragment
4:   send  $e$  to all nodes
5:   determine blue edges of all fragments
6:   add blue edges of all fragments to MST, update fragments
7: end while

```

Theorem 15.3. *On a complete graph, Algorithm 56 computes an MST in time $\mathcal{O}(\log n)$.*

Proof. The algorithm is correct because of Lemma 15.2. Every node only needs to send a single message to all its neighbors in every phase (line 4). All other computations can be done locally without sending other messages. In particular, the blue edge of a given fragment is the lightest edge sent by any node of that fragment. Because every node always knows the current MST (and all current fragments), lines 5 and 6 can be performed locally.

In every phase, every fragment connects to at least one other fragment. The minimum fragment size therefore at least doubles in every phase. Thus, the number of phases is at most $\log_2 n$. \square

Remarks:

- Algorithm 56 does essentially the same thing as the GHS algorithm (Algorithm 15) discussed in Chapter 4. Because we now have a complete graph and thus every node can communicate with every other node, things become simpler (and also much faster).
- Algorithm 56 does not make use of the fact that a node can send different messages to different nodes. Making use of this possibility will allow us to significantly reduce the running time of the algorithm.

Our goal is now to improve Algorithm 56. We assume that every node has a unique identifier. By sending its own identifier to all other nodes, every node knows the identifiers of all other nodes after one round. Let $\ell(F)$ be the node with the smallest identifier in fragment F . We call $\ell(F)$ the leader of fragment F . In order to improve the running time of Algorithm 56, we need to be able to connect every fragment to more than one other fragment in a single phase. Algorithm 57 shows how the nodes can learn about the $k = |F|$ lightest outgoing edges of each fragment F (in constant time!).

Algorithm 57 Fast MST construction (at node v)

```

1: // all nodes always know all current MST edges and thus all MST fragments
2: repeat
3:    $F :=$  fragment of  $v$ ;
4:    $\forall F' \neq F$ , compute min-weight edge  $e_{F'}$  connecting  $v$  to  $F'$ 
5:    $\forall F' \neq F$ , send  $e_{F'}$  to  $\ell(F')$ 
6:   if  $v = \ell(F)$  then
7:      $\forall F' \neq F$ , determine min-weight edge  $e_{F,F'}$  between  $F$  and  $F'$ 
8:      $k := |F|$ 
9:      $E(F) := k$  lightest edges among  $e_{F,F'}$  for  $F' \neq F$ 
10:    send send each edge in  $E(F)$  to a different node in  $F$ 
        // for simplicity assume that  $v$  also sends an edge to itself
11:  end if
12:  send edge received from  $\ell(F)$  to all nodes
13:  // the following operations are performed locally by each node
14:   $E' :=$  edges received by other nodes
15:  AddEdges( $E'$ )
16: until all nodes are in the same fragment

```

Given this set E' of edges, each node can locally decide which edges can safely be added to the constructed tree by calling the subroutine AddEdges (Algorithm 58). Note that the set of received edges E' in line 14 is the same for all nodes. Since all nodes know all current fragments, all nodes add the same set of edges!

Algorithm 58 uses the lightest outgoing edge that connects two fragments (to a larger super-fragment) as long as it is safe to add this edge, i.e., as long as it is clear that this edge is a blue edge. A (super-)fragment that has outgoing edges in E' that are surely blue edges is called *safe*. As we will see, a super-fragment \mathcal{F} is safe if all the original fragments that make up \mathcal{F} are still incident to at least one edge in E' that has not yet been considered. In order to determine whether all lightest outgoing edges in E' that are incident to a certain fragment F have been processed, a counter $c(F)$ is maintained (see line 2). If an edge incident to two (distinct) fragments F_i and F_j is processed, both $c(F_i)$ and $c(F_j)$ are decremented by 1 (see Line 8).

An edge connecting two distinct super-fragments \mathcal{F}' and \mathcal{F}'' is added if at least one of the two super-fragments is safe. In this case, the two super-fragments are merged into one (new) super-fragment. The new super-fragment is safe if and only if both original super-fragments are safe and the processed edge e is not the last edge in E' incident to any of the two fragments F_i and F_j that are incident to e , i.e., both counters $c(F_i)$ and $c(F_j)$ are still positive (see line 12).

The considered edge e may not be added for one of two reasons. It is possible that both \mathcal{F}' and \mathcal{F}'' are not safe. Since a super-fragment cannot become safe again, nothing has to be done in this case. The second reason is that $\mathcal{F}' = \mathcal{F}''$. In this case, this single fragment may become unsafe if e reduced either $c(F_i)$ or $c(F_j)$ to zero (see line 18).

Algorithm 58 AddEdges(E'): Given the set of edges E' , determine which edges are added to the MST

```

1: Let  $F_1, \dots, F_r$  be the initial fragments
2:  $\forall F_i \in \{F_1, \dots, F_r\}, c(F_i) := \#$  incident edges in  $E'$ 
3: Let  $\mathcal{F}_1 := F_1, \dots, \mathcal{F}_r := F_r$  be the initial super-fragments
4:  $\forall \mathcal{F}_i \in \{\mathcal{F}_1, \dots, \mathcal{F}_r\}, safe(\mathcal{F}_i) := true$ 
5: while  $E' \neq \emptyset$  do
6:    $e :=$  lightest edge in  $E'$  between the original fragments  $F_i$  and  $F_j$ 
7:    $E' := E' \setminus \{e\}$ 
8:    $c(F_i) := c(F_i) - 1, c(F_j) := c(F_j) - 1$ 
9:   if  $e$  connects super-fragments  $\mathcal{F}' \neq \mathcal{F}''$  and  $(safe(\mathcal{F}') \text{ or } safe(\mathcal{F}''))$  then
10:     add  $e$  to MST
11:     merge  $\mathcal{F}'$  and  $\mathcal{F}''$  into one super-fragment  $\mathcal{F}_{new}$ 
12:     if  $safe(\mathcal{F}') \text{ and } safe(\mathcal{F}'') \text{ and } c(F_i) > 0 \text{ and } c(F_j) > 0$  then
13:        $safe(\mathcal{F}_{new}) := true$ 
14:     else
15:        $safe(\mathcal{F}_{new}) := false$ 
16:     end if
17:   else if  $\mathcal{F}' = \mathcal{F}''$  and  $(c(F_i) = 0 \text{ or } c(F_j) = 0)$  then
18:      $safe(\mathcal{F}') := false$ 
19:   end if
20: end while

```

Lemma 15.4. *The algorithm only adds MST edges.*

Proof. We have to prove that at the time we add an edge e in line 9 of Algorithm 58, e is the blue edge of some (super-)fragment. By definition, e is the lightest edge that has not been considered and that connects two distinct super-fragments \mathcal{F}' and \mathcal{F}'' . Since e is added, we know that either $safe(\mathcal{F}')$ or $safe(\mathcal{F}'')$ is true. Without loss of generality, assume that \mathcal{F}' is safe. According to the definition of *safe*, this means that from each fragment F in the super-fragment \mathcal{F}' we know at least the lightest outgoing edge, which implies that we also know the lightest outgoing edge, i.e., the blue edge, of \mathcal{F}' . Since e is the lightest edge that connects *any* two super-fragments, it must hold that e is exactly the blue edge of \mathcal{F}' . Thus, whenever an edge is added, it is an MST edge. \square

Theorem 15.5. *Algorithm 57 computes an MST in time $\mathcal{O}(\log \log n)$.*

Proof. Let β_k denote the size of the smallest fragment after phase k of Algorithm 57. We first show that every fragment merges with at least β_k other fragments in each phase. Since the size of each fragment after phase k is at least β_k by definition, we get that the size of each fragment after phase $k + 1$ is at least $\beta_k(\beta_k + 1)$. Assume that a fragment F , consisting of at least β_k nodes,

does not merge with β_k other fragments in phase $k + 1$ for any $k \geq 0$. Note that F cannot be safe because being safe implies that there is at least one edge in E' that has not been considered yet and that is the blue edge of F . Hence, the phase cannot be completed in this case. On the other hand, if F is not safe, then at least one of its sub-fragments has used up all its β_k edges to other fragments. However, such an edge is either used to merge two fragments or it must have been dropped because the two fragments already belong to the same fragment because another edge connected them (in the same phase). In either case, we get that any fragment, and in particular F , must merge with at least β_k other fragments.

Given that the minimum fragment size grows (quickly) in each phase and that only edges belonging to the MST are added according to Lemma 15.4, we conclude that the algorithm correctly computes the MST. The fact that

$$\beta_{k+1} \geq \beta_k(\beta_k + 1)$$

implies that $\beta_k \geq 2^{2^{k-1}}$ for any $k \geq 1$. Therefore after $1 + \log_2 \log_2 n$ phases, the minimum fragment size is n and thus all nodes are in the same fragment. \square

Remarks:

- It is not known whether the $\mathcal{O}(\log \log n)$ time complexity of Algorithm 57 is optimal. In fact, no lower bounds are known for the MST construction on graphs of diameter 1 and 2.
- Algorithm 57 makes use of the fact that it is possible to send different messages to different nodes. If we assume that every node always has to send the same message to all other nodes, Algorithm 56 is the best that is known. Also for this simpler case, no lower bound is known.

Chapter 16

Consensus

This chapter is the first to deal with fault tolerance, one of the most fundamental aspects of distributed computing. Indeed, in contrast to a system with a single processor, having a distributed system may permit getting away with failures and malfunctions of parts of the system. This line of research was motivated by the basic question whether, e.g., putting two (or three?) computers into the cockpit of a plane will make the plane more reliable. Clearly fault-tolerance often comes at a price, as having more than one decision-maker often complicates decision-making.

16.1 Impossibility of Consensus

Imagine two cautious generals who want to attack a common enemy.¹ Their only means of communication are messengers. Unfortunately, the route of these messengers leads through hostile enemy territory, so there is a chance that a messenger does not make it. Only if both generals attack at the very same time the enemy can be defeated. Can we devise a protocol such that the two generals can agree on an attack time? Clearly general *A* can send a message to general *B* asking to e.g. “attack at 6am”. However, general *A* cannot be sure that this message will make it, so she asks for a confirmation. The problem is that general *B* getting the message cannot be sure that her confirmation will reach general *A*. If the confirmation message indeed is destroyed, general *A* cannot distinguish this case from the case where general *B* did not even get the attack information. So, to be safe, general *B* herself will ask for a confirmation of her confirmation. Taking again the position of general *A* we can similarly derive that she cannot be sure unless she also gets a confirmation of the confirmation of the confirmation...

To make things worse, also different approaches do not seem to work. In fact it can be shown that this two generals problem cannot be solved, in other words, there is no finite protocol which lets the two generals find consensus! To show this, we need to be a bit more formal:

¹If you don't fancy the martial tone of this classic example, feel free to think about something else, for instance two friends trying to make plans for dinner over instant messaging software, or two lecturers sharing the teaching load of a course trying to figure out who is in charge of the next lecture.

Definition 16.1 (Consensus). *Consider a distributed system with n nodes. Each node i has an input x_i . A solution of the consensus problem must guarantee the following:*

- *Termination: Every non-faulty node eventually decides.*
- *Agreement: All non-faulty nodes decide on the same value.*
- *Validity: The decided value must be the input of at least one node.*

Remarks:

- The validity condition infers that if all nodes have the same input x , then the nodes need to decide on x . Please note that consensus is not democratic, it may well be that the nodes decide on an input value promoted by a small minority.
- Whether consensus is possible depends on many parameters of the distributed system, in particular whether the system is synchronous or asynchronous, or what “faulty” means. In the following we study some simple variants to get a feeling for the problem.
- Consensus is a powerful primitive. With established consensus almost everything can be computed in a distributed system, e.g. a leader.

Given a distributed asynchronous message passing system with $n \geq 2$ nodes. All nodes can communicate directly with all other nodes, simply by sending a message. In other words, the communication graph is the complete graph. Can the consensus problem be solved? Yes!

Algorithm 59 Trivial Consensus

- 1: Each node has an input
 - 2: We have a leader, e.g. the node with the highest ID
 - 3: **if** node v is the leader **then**
 - 4: the leader shall simply decide on its own input
 - 5: **else**
 - 6: send message to the leader asking for its input
 - 7: wait for answer message by leader, and decide on that
 - 8: **end if**
-

Remarks:

- This algorithm is quite simple, and at first sight seems to work perfectly, as all three consensus conditions of Definition 16.1 are fulfilled.
- However, the algorithm is not fault-tolerant at all. If the leader crashes before being able to answer all requests, there are nodes which will never terminate, and hence violate the termination condition. Is there a deterministic protocol that can achieve consensus in an asynchronous system, even in the presence of failures? Let’s first try something slightly different.

Definition 16.2 (Reliable Broadcast). *Consider an asynchronous distributed system with n nodes that may crash. Any two nodes can exchange messages, i.e., the communication graph is complete. We want node v to send a reliable broadcast to the $n - 1$ other nodes. Reliable means that either nobody receives the message, or everybody receives the message.*

Remarks:

- This seems to be quite similar to consensus, right?
- The main problem is that the sender may crash while sending the message to the $n - 1$ other nodes such that some of them get the message, and the others not. We need a technique that deals with this case:

Algorithm 60 Reliable Broadcast

```

1: if node  $v$  is the source of message  $m$  then
2:   send message  $m$  to each of the  $n - 1$  other nodes
3:   upon receiving  $m$  from any other node: broadcast succeeded!
4: else
5:   upon receiving message  $m$  for the first time:
6:     send message  $m$  to each of the  $n - 1$  other nodes
7: end if

```

Theorem 16.3. *Algorithm 60 solves reliable broadcast as in Definition 16.2.*

Proof. First we should note that we do not care about nodes that crash during the execution: whether or not they receive the message is irrelevant since they crashed anyway. If a single non-faulty node u received the message (no matter how, it may be that it received it through a path of crashed nodes) all non-faulty nodes will receive the message through u . If no non-faulty node receives the message, we are fine as well! \square

Remarks:

- While it is clear that we could also solve reliable broadcast by means of a consensus protocol (first send message, then agree on having received it), the opposite seems more tricky!
- No wonder, it cannot be done!! For the presentation of this impossibility result we use the read/write shared memory model introduced in Chapter 6. Not only was the proof originally conceived in the shared memory model, it is also cleaner.

Definition 16.4 (Univalent, Bivalent). *A distributed system is called x -valent if the outcome of a computation will be x . An x -valent system is also called univalent. If, depending on the execution, still more than one possible outcome is feasible, the system is called multivalent. If exactly two outcomes are still possible, the system is called bivalent.*

Theorem 16.5. *In an asynchronous shared memory system with $n > 1$ nodes, and node crash failures (but no memory failures!) consensus as in Definition 16.1 cannot be achieved by a deterministic algorithm.*

Proof. Let us simplify the proof by setting $n = 2$. We have processes u and v , with input values x_u and x_v . Further let the input values be binary, either 0 or 1.

First we have to make sure that there are input values such that initially the system is bivalent. If $x_u = 0$ and $x_v = 0$ the system is 0-valent, because of the validity condition (Definition 16.1). Even in the case where process v immediately crashes the system remains 0-valent. Similarly if both input values are 1 and process u immediately crashes the system is 1-valent. If $x_u = 0$ and $x_v = 1$ and v immediately crashes, process u cannot distinguish from both having input 0, equivalently if u immediately crashes, process v cannot distinguish from both having 1, hence the system is bivalent!

In order to solve consensus an algorithm needs to terminate. All non-faulty processes need to decide on the same value x (agreement condition of Definition 16.1), in other words, at some instant this value x must be known to the system as a whole, meaning that no matter what the execution is, the system will be x -valent. In other words, the system needs to change from bivalent to univalent. We may ask ourselves what can cause this change in a deterministic asynchronous shared memory algorithm? We need an element of non-determinism; if everything happens deterministically the system would have been x -valent even after initialization which we proved to be impossible already.

The only nondeterministic elements in our model are the asynchrony of accessing the memory and crashing processes. Initially and after every memory access, each process decides what to do next: Read or write a memory cell or terminate with a decision. We take control of the scheduling, either choosing which request is served next or making a process crash. Now we hope for a *critical* bivalent state with more than one memory request, and depending which memory request is served next the system is going to switch from bivalent to univalent. More concretely, if process u is being served next the system is going x -valent, if process v (with $v \neq u$) is served next the system is going y -valent (with $y \neq x$). We have several cases:

- If the operations of processes u and v target different memory cells, processes cannot distinguish which memory request was executed first. Hence the local states of the processes are identical after serving both operations and the state cannot be critical.
- The same argument holds if both processes want to read the same register. Nobody can distinguish which read was first, and the state cannot be critical.
- If process u reads memory cell c , and process v writes memory cell c , the scheduler first executes u 's read. Now process v cannot distinguish whether that read of u did or did not happen before its write. If it did happen, v should decide on x , if it did not happen, v should decide y . But since v does not know which one is true, it needs to be informed about it by u . We prevent this by making u crash. Thus the state can only be univalent if v never decides, violating the termination condition!
- Also if both processes write the same memory cell we have the same issue, since the second writer will immediately overwrite the first writer, and hence the second writer cannot know whether the first write happened at all. Again, the state cannot be critical.

Hence, if we are unlucky (and in a worst case, we are!) there is no critical state. In other words, the system will remain bivalent forever, and consensus is impossible. \square

Remarks:

- The proof presented is a variant of a proof by Michael Fischer, Nancy Lynch and Michael Paterson, a classic result in distributed computing. The proof was motivated by the problem of committing transactions in distributed database systems, but is sufficiently general that it directly implies the impossibility of a number of related problems, including consensus. The proof also is pretty robust with regard to different communication models.
- The FLP (Fischer, Lynch, Paterson) paper won the 2001 PODC Influential Paper Award, which later was renamed Dijkstra Prize.
- One might argue that FLP destroys all the fun in distributed computing, as it makes so many things impossible! For instance, it seems impossible to have a distributed database where the nodes can reach consensus whether to commit a transaction or not.
- So are two-phase-commit (2PC), three-phase-commit (3PC) et al. wrong?! No, not really, but sometimes they just do not commit!
- What about turning some other knobs of the model? Can we have consensus in a message passing system? No. Can we have consensus in synchronous systems? Yes, even if all but one node fails!
- Can we have consensus in synchronous systems even if some nodes are mischievous, and behave much worse than simply crashing, and send for example contradicting information to different nodes? This is known as *Byzantine* behavior. Yes, this is also possible, as long as the Byzantine nodes are strictly less than a third of all the nodes. This was shown by Marshall Pease, Robert Shostak, and Leslie Lamport in 1980. Their work won the 2005 Dijkstra Prize, and is one of the cornerstones not only in distributed computing but also information security. Indeed this work was motivated by the “fault-tolerance in planes” example. Pease, Shostak, and Lamport noticed that the computers they were given to implement a fault-tolerant fighter plane at times behaved strangely. Before crashing, these computers would start behaving quite randomly, sending out weird messages. At some point Pease et al. decided that a malicious behavior model would be the most appropriate to be on the safe side. Being able to allow strictly less than a third Byzantine nodes is quite counterintuitive; even today many systems are built with three copies. In light of the result of Pease et al. this is a serious mistake! If you want to be tolerant against a single Byzantine machine, you need four copies, not three!
- Finally, FLP only prohibits deterministic algorithms! So can we solve consensus if we use randomization? The answer again is yes! We will study this in the remainder of this chapter.

16.2 Randomized Consensus

Can we solve consensus if we allow randomization? Yes. The following algorithm solves Consensus even in face of Byzantine errors, i.e., malicious behavior of some of the nodes. To simplify arguments we assume that at most f nodes will fail (crash) with $n > 9f$, and that we only solve binary consensus, that is, the input values are 0 and 1. The general idea is that nodes try to push their input value; if other nodes do not follow they will try to push one of the suggested values randomly. The full algorithm is in Algorithm 61.

Algorithm 61 Randomized Consensus

```

1: node  $u$  starts with input bit  $x_u \in \{0, 1\}$ , round:=1.
2: broadcast BID( $x_u$ , round)
3: repeat
4:   wait for  $n - f$  BID messages of current round
5:   if at least  $n - f$  messages have value  $x$  then
6:      $x_u := x$ ; decide on  $x$ 
7:   else if at least  $n - 2f$  messages have value  $x$  then
8:      $x_u := x$ 
9:   else
10:    choose  $x_u$  randomly, with  $Pr[x_u = 0] = Pr[x_u = 1] = 1/2$ 
11:   end if
12:   round := round + 1
13:   broadcast BID( $x_u$ , round)
14: until decided

```

Theorem 16.6. *Algorithm 61 solves consensus as in Definition 16.1 even if up to $f < n/9$ nodes exhibit Byzantine failures.*

Proof. First note that it is not a problem to wait for $n - f$ BID messages in line 4 since at most f nodes are corrupt. If all nodes have the same input value x , then all (except the f Byzantine nodes) will bid for the same value x . Thus, every node receives at least $n - 2f$ BID messages containing x , deciding on x in the first round already. We have consensus!

If the nodes have different (binary) input values the validity condition becomes trivial as any result is fine. What about agreement? Let u be one of the first nodes to decide on value x (in line 6). It may happen that due to asynchronicity another node v received messages from a different subset of the nodes, however, at most f senders may be different. Taking into account that Byzantine nodes may lie, i.e., send different BIDs to different nodes, f additional BID messages received by v may differ from those received by u . Since node u had at least $n - 2f$ BID messages with value x , node v has at least $n - 4f$ BID messages with x . Hence every correct node will bid for x in the next round, and then decide on x .

So we only need to worry about termination! We already have seen that as soon as one correct node terminates (in line 6) everybody terminates in the next round. So what are the chances that some node u terminates in line 6? Well, if push comes to shove we can still hope that all correct nodes randomly propose the same value (in line 10). Maybe there are some nodes not choosing

at random (i.e., entering line 8), but they unanimously propose either 0 or 1: For the sake of contradiction, assume that both 0 and 1 are proposed in line 8. This means that both 0 and 1 had been proposed by at least $n - 5f$ *correct* nodes. In other words, we have a total of $2(n - 5f) + f = n + (n - 9f) > n$ nodes. Contradiction!

Thus, at worst all $n - f$ correct nodes need to randomly choose the same bit, which happens with probability $2^{-(n-f)}$. If so, all will send the same BID, and the algorithm terminates. So the expected running time is smaller than 2^n . \square

Remarks:

- The presentation of Algorithm 61 is a simplification of the typical presentation in text books.
- What about an algorithm that allows for crashes only, but can manage more failures? Good news! Slightly changing the presented algorithm will do that for $f < n/4$! See exercises.
- Unfortunately Algorithm 61 is still impractical as termination is awfully slow. In expectation about the same number of nodes choose 1 or 0 in line 10. Termination would be much more efficient if all nodes chose the same random value in line 10! So why not simply replacing line 10 with “choose $x_u := 1$ ”?!! The problem is that a majority of nodes may see a majority of 0 bids, hence proposing 0 in the next round. Without randomization it is impossible to get out of this equilibrium. (Moreover, this approach is deterministic, contradicting Theorem 16.5.)
- The idea is to replace line 10 with a subroutine where all nodes compute a so-called *shared* (or common, or global) coin. A shared coin is a random variable that is 0 with constant probability and 1 with constant probability. Sounds like magic, but it isn't! We assume at most $f < n/3$ nodes may crash:

Algorithm 62 Shared Coin (code for node u)

```

1: set local coin  $x_u := 0$  with probability  $1/n$ , else  $x_u := 1$ 
2: use reliable broadcast to tell everybody about your local coin  $x_u$ 
3: memorize all coins you get from others in the set  $c_u$ 
4: wait for exactly  $n - f$  coins
5: copy these coins into your local set  $s_u$  (but keep learning coins)
6: use reliable broadcast to tell everybody about your set  $s_u$ 
7: wait for exactly  $n - f$  sets  $s_v$  (which satisfy  $s_v \subseteq c_u$ )
8: if seen at least a single coin 0 then
9:   return 0
10: else
11:   return 1
12: end if

```

Theorem 16.7. *If $f < n/3$ nodes crash, Algorithm 62 implements a shared coin.*

Proof. Since only f nodes may crash, each node sees at least $n - f$ coins and sets in lines 4 and 7, respectively. Thanks to the reliable broadcast protocol each node eventually sees all the coins in the other sets. In other words, the algorithm terminates in $\mathcal{O}(1)$ time.

The general idea is that a third of the coins are being seen by everybody. If there is a 0 among these coins, everybody will see that 0. If not, chances are high that there is no 0 at all! Here are the details:

Let u be the first node to terminate (satisfy line 7). For u we draw a matrix of all the seen sets s_v (columns) and all coins c_u seen by node u (rows). Here is an example with $n = 7, f = 2, n - f = 5$:

	s_1	s_3	s_5	s_6	s_7
c_1	X	X	X	X	X
c_2			X	X	X
c_3	X	X	X	X	X
c_5	X	X	X		X
c_6	X	X	X	X	
c_7	X	X		X	X

Note that there are exactly $(n - f)^2$ X's in this matrix as node u has seen exactly $n - f$ sets (line 7) each having exactly $n - f$ coins (lines 4 to 6). We need two little helper lemmas:

Lemma 16.8. *There are at least $f + 1$ rows that have at least $f + 1$ X's*

Proof. Assume (for the sake of contradiction) that this is not the case. Then at most f rows have all $n - f$ X's, and all other rows (at most $n - f$) have at most f X's. In other words, the number of total X's is bounded by

$$|X| \leq f \cdot (n - f) + (n - f) \cdot f = 2f(n - f).$$

Using $n > 3f$ we get $n - f > 2f$, and hence $|X| \leq 2f(n - f) < (n - f)^2$. This is a contradiction to having exactly $(n - f)^2$ X's! \square

Lemma 16.9. *Let W be the set of local coins for which the corresponding matrix row has more than f X's. All local coins in the set W are seen by all nodes that terminate.*

Proof. Let $w \in W$ be such a local coin. By definition of W we know that w is in at least $f + 1$ seen sets. Since each node must see at least $n - f$ seen sets before terminating, each node has seen at least one of these sets, and hence w is seen by everybody terminating. \square

Continuing the proof of Theorem 16.7: With probability $(1 - 1/n)^n \approx 1/e \approx .37$ all nodes chose their local coin equal to 1, and 1 is decided. With probability $1 - (1 - 1/n)^{|W|}$ there is at least one 0 in W . With Lemma 16.8 we know that $|W| \approx n/3$, hence the probability is about $1 - (1 - 1/n)^{n/3} \approx 1 - (1/e)^{1/3} \approx .28$. With Lemma 16.9 this 0 is seen by all, and hence everybody will decide 0. So indeed we have a shared coin. \square

Theorem 16.10. *Plugging Algorithm 62 into Algorithm 61 we get a randomized consensus algorithm which finishes in a constant expected number of rounds.*

Remarks:

- If some nodes go into line 8 of Algorithm 61 the others still have a constant probability to guess the same shared coin.
- For crash failures there exists an improved constant expected time algorithm which tolerates f failures with $2f < n$.
- For Byzantine failures there exists a constant expected time algorithm which tolerates f failures with $3f < n$.
- Similar algorithms have been proposed for the shared memory model.

Chapter 17

Multi-Core Computing

This chapter is based on the article “Distributed Computing and the Multicore Revolution” by Maurice Herlihy and Victor Luchangco. Thanks!

17.1 Introduction

In the near future, nearly all computers, ranging from supercomputers to cell phones, will be multiprocessors. It is harder and harder to increase processor clock speed (the chips overheat), but easier and easier to cram more processor cores onto a chip (thanks to Moore’s Law). As a result, uniprocessors are giving way to dual-cores, dual-cores to quad-cores, and so on.

However, there is a problem: Except for “embarrassingly parallel” applications, no one really knows how to exploit lots of cores.

17.1.1 The Current State of Concurrent Programming

In today’s programming practice, programmers typically rely on combinations of locks and conditions, such as monitors, to prevent concurrent access by different threads to the same shared data. While this approach allows programmers to treat sections of code as “atomic”, and thus simplifies reasoning about interactions, it suffers from a number of severe shortcomings.

- Programmers must decide between *coarse-grained* locking, in which a large data structure is protected by a single lock (usually implemented using operations such as test-and-set or compare and swap(CAS)), and *fine-grained* locking, in which a lock is associated with each component of the data structure. Coarse-grained locking is simple, but permits little or no concurrency, thereby preventing the program from exploiting multiple processing cores. By contrast, fine-grained locking is substantially more complicated because of the need to ensure that threads acquire all necessary locks (and only those, for good performance), and because of the need to avoid deadlocks, when acquiring multiple locks. The decision is further complicated by the fact that the best engineering solution may be

Algorithm **Move**(Element **e**, **Table** **from**, **Table** **to**)

```

1: if from.find(e) then
2:   to.insert(e)
3:   from.delete(e)
4: end if

```

platform-dependent, varying with different machine sizes, workloads, and so on, making it difficult to write code that is both scalable and portable.

- Conventional locking provides poor support for code composition and reuse. For example, consider a lock-based hash table that provides atomic **insert** and **delete** methods. Ideally, it should be easy to *move* an element atomically from one table to another, but this kind of composition simply does not work. If the table methods synchronize internally, then there is no way to acquire and hold both locks simultaneously. If the tables export their locks, then modularity and safety are compromised. For a concrete example, assume we have two hash tables T_1 and T_2 storing integers and using internal locks only. Every number is only inserted into a table, if it is not already present, i.e., multiple occurrences are not permitted. We want to atomically move elements using two threads between the tables using Algorithm Move. If we have external locks, we must pay attention to avoid deadlocks etc.

	Table T1 is contains 1 and T2 is empty	
Time	Thread 1	Thread 2
	Move(1,T1,T2)	Move(1,T2,T1)
1	T1.find(1)	delayed
2	T2.insert(1)	
3	delayed	T2.find(1)
4		T1.insert(1)
5	T1.delete(1)	T2.delete(1)
	both T1 and T2 are empty	

- Such basic issues as the mapping from locks to data, that is, which locks protect which data, and the order in which locks must be acquired and released, are all based on convention, and violations are notoriously difficult to detect and debug. For these and other reasons, today's software practices make lock-based concurrent programs (too) difficult to develop, debug, understand, and maintain.

The research community has addressed this issue for more than fifteen years by developing nonblocking algorithms for stacks, queues and other data structures. These algorithms are subtle and difficult. For example, the pseudo code of a delete operation for a (non-blocking) linked list, recently presented at a conference, contains more than 30 lines of code, whereas a delete procedure for a (non-concurrent, used only by one thread) linked list can be written with 5 lines of code.

17.2 Transactional Memory

Recently the *transactional memory* programming paradigm has gained momentum as an alternative to locks in concurrent programming. Rather than using locks to give the illusion of atomicity by preventing concurrent access to shared data with transactional memory, programmers designate regions of code as transactions, and the system guarantees that such code appears to execute atomically. A transaction that cannot complete is aborted—its effects are discarded—and may be retried. Transactions have been used to build large, complex and reliable database systems for over thirty years; with transactional memory, researchers hope to translate that success to multiprocessor systems. The underlying system may use locks or nonblocking algorithms to implement transactions, but the complexity is hidden from the application programmer. Proposals exist for implementing transactional memory in hardware, in software, and in schemes that mix hardware and software. This area is growing at a fast pace.

More formally, a transaction is defined as follows:

Definition 17.1. *A transaction in transactional memory is characterized by three properties (ACI):*

- *Atomicity: Either a transaction finishes all its operations or no operation has an effect on the system.*
- *Consistency: All objects are in a valid state before and after the transaction.*
- *Isolation: Other transactions cannot access or see data in an intermediate (possibly invalid) state of any parallel running transaction.*

Remarks:

- For database transactions there exists a fourth property called durability: If a transaction has completed, its changes are permanent, i.e., even if the system crashes, the changes can be recovered. In principle, it would be feasible to demand the same thing for transactional memory, however this would mean that we had to use slow hard discs instead of fast DRAM chips...
- Although transactional memory is a promising approach for concurrent programming, it is not a panacea, and in any case, transactional programs will need to interact with other (legacy) code, which may use locks or other means to control concurrency.
- One major challenge for the adoption of transactional memory is that it has no universally accepted specification. It is not clear yet how to interact with I/O and system calls should be dealt with. For instance, imagine you print a news article. The printer job is part of a transaction. After printing half the page, the transaction gets aborted. Thus the work (printing) is lost. Clearly, this behavior is not acceptable.
- From a theory perspective we also face a number of open problems. For example:

- System model: An abstract model for a (shared-memory) multiprocessor is needed that properly accounts for performance. In the 80s, the PRAM model became a standard model for parallel computation, and the research community developed many elegant parallel algorithms for this model. Unfortunately, PRAM assume that processors are synchronous, and that memory can be accessed only by read and write operations. Modern computer architectures are asynchronous and they provide additional operations such as test-and-set. Also, PRAM did not model the effects of contention nor the performance implications of multilevel caching, assuming instead a flat memory with uniform-cost access. More realistic models have been proposed to account for the costs of interprocess communication, but these models still assume synchronous processors with only read and write access to memory.
- How to resolve conflicts? Many transactional memory implementations “optimistically” execute transactions in parallel. Conflicts between two transactions intending to modify the same memory at the same time are resolved by a contention manager. A contention manager decides whether a transaction continues, waits or is aborted. The contention management policy of a transactional memory implementation can have a profound effect on its performance, and even its progress guarantees.

17.3 Contention Management

After the previous introduction of transactional memory, we look at different aspects of contention management from a theoretical perspective. We start with a description of the model.

We are given a set of *transactions* $S := \{T_1, \dots, T_n\}$ sharing up to s *resources* (such as memory cells) that are executed on n *threads*. Each thread runs on a separate processor/core P_1, \dots, P_n . For simplicity, each transaction T consists of a sequence of t_T operations. An operation requires one time unit and can be a write access of a resource R or some arbitrary computation.¹ To perform a write, the written resource must be acquired exclusively (i.e., locked) before the access. Additionally, a transaction must store the original value of a written resource. Only one transaction can lock a resource at a time. If a transaction A attempts to acquire a resource, locked by B , then A and B face a conflict. If multiple transactions concurrently attempt to acquire an unlocked resource, an arbitrary transaction A will get the resource and the others face a conflict with A . A *contention manager* decides how to resolve a conflict. Contention managers operate in a distributed fashion, that is to say, a separate instance of a contention manager is available for every thread and they operate independently. Contention managers can make a transaction wait (arbitrarily long) or abort. An aborted transaction undoes all its changes to resources and frees all locks before restarting. Freeing locks and undoing the changes can be done with one operation. A successful transaction finishes with a commit and simply frees

¹Reads are of course also possible, but are not critical because they do not attempt to modify data.

all locks. A contention manager is unaware of (potential) future conflicts of a transaction. The required resources might also change at any time.

The quality of a contention manager is characterized by different properties:

- **Throughput:** How long does it take until all transactions have committed? How good is our algorithm compared to an optimal?

Definition 17.2. *The makespan of the set S of transactions is the time interval from the start of the first transaction until all transactions have committed.*

Definition 17.3. *The competitive ratio is the ratio of the makespans of the algorithm to analyze and an optimal algorithm.*

- **Progress guarantees:** Is the system deadlock-free? Does every transaction commit in finite time?

Definition 17.4. *We look at three levels of progress guarantees:*

- *wait freedom (strongest guarantee): all threads make progress in a finite number of steps*
- *lock freedom: one thread makes progress in a finite number of steps*
- *obstruction freedom (weakest): one thread makes progress in a finite number of steps in absence of contention (no other threads compete for the same resources)*

Remarks:

- For the analysis we assume an *oblivious* adversary. It knows the algorithm to analyze and chooses/modifies the operations of transactions arbitrarily. However, the adversary does not know the random choices (of a randomized algorithm). The optimal algorithm knows all decisions of the adversary, i.e. first the adversary must say how transactions look like and then the optimal algorithm, having full knowledge of all transaction, computes an (optimal) schedule.
- Wait freedom implies lock freedom. Lock freedom implies obstruction freedom.
- Here is an example to illustrate how needed resources change over time: Consider a dynamic data structure such as a balanced tree. If a transaction attempts to insert an element, it must modify a (parent) node and maybe it also has to do some rotations to rebalance the tree. Depending on the elements of the tree, which change over time, it might modify different objects. For a concrete example, assume that the root node of a binary tree has value 4 and the root has a (left) child of value 2. If a transaction A inserts value 5, it must modify the pointer to the right child of the root node with value 4. Thus it locks the root node. If A gets aborted by a transaction B , which deletes the node with value 4 and commits, it will attempt to lock the new root node with value 2 after its restart.

- There are also systems, where resources are not locked exclusively. All we need is a correct serialization (analogous to transactions in database systems). Thus a transaction might speculatively use the current value of a resource, modified by an uncommitted transaction. However, these systems must track dependencies to ensure the ACI properties of a transaction (see Definition 17.1). For instance, assume a transaction T_1 increments variable x from 1 to 2. Then transaction T_2 might access x and assume its correct value is 2. If T_1 commits everything is fine and the ACI properties are ensured, but if T_1 aborts, T_2 must abort too, since otherwise the atomicity property was violated.
- In practice, the number of concurrent transactions might be much larger than the number of processors. However, performance may decrease with an increasing number of threads since there is time wasted to switch between threads. Thus, in practice, load adaption schemes have been suggested to limit the number of concurrent transactions close to (or even below) the number of cores.
- In the analysis, we will assume that the number of operations is fixed for each transaction. However, the execution time of a transaction (in the absence of contention) might also change, e.g., if data structures shrink, less elements have to be considered. Nevertheless, often the changes are not substantial, i.e., only involve a constant factor. Furthermore, if an adversary can modify the duration of a transaction arbitrarily during the execution of a transaction, then any algorithm must make the exact same choices as an optimal algorithm: Assume two transactions T_0 and T_1 face a conflict and an algorithm Alg decides to let T_0 wait (or abort). The adversary could make the opposite decision and let T_0 proceed such that it commits at time t_0 . Then it sets the execution time T_0 to infinity, i.e., $t_{T_0} = \infty$ after t_0 . Thus, the makespan of the schedule for algorithm Alg is unbounded though there exists a schedule with bounded makespan. Thus the competitive ratio is unbounded.

Problem complexity

In graph theory, coloring a graph with as few colors as possible is known to be hard problem. A (vertex) coloring assigns a color to each vertex of a graph such that no two adjacent vertices share the same color. It was shown that computing an optimal coloring given complete knowledge of the graph is NP-hard. Even worse, computing an approximation within a factor of $\chi(G)^{\log \chi(G)/25}$, where $\chi(G)$ is the minimal number of colors needed to color the graph, is NP-hard as well.

To keep things simple, we assume for the following theorem that resource acquisition takes no time, i.e., as long as there are no conflicts, transactions get all locks they wish for at once. In this case, there is an immediate connection to graph coloring, showing that even an *offline* version of contention management, where all potential conflicts are known and do not change over time, is extremely hard to solve.

Theorem 17.5. *If the optimal schedule has makespan k and resource acquisition takes zero time, it is NP-hard to compute a schedule of makespan less than*

$k^{\log k/25}$, even if all conflicts are known and transactions do not change their resource requirements.

Proof. We will prove the claim by showing that any algorithm finding a schedule taking $k' < k^{(\log k)/25}$ can be utilized to approximate the chromatic number of any graph better than $\chi(G)^{\frac{\log \chi(G)}{25}}$.

Given the graph $G = (V, E)$, define that V is the set of transactions and E is the set of resources. Each transaction (node) $v \in V$ needs to acquire a lock on all its resources (edges) $\{v, w\} \in E$, and then computes something for exactly one round. Obviously, this “translation” of a graph into our scheduling problem does not require any computation at all.

Now, if we knew a $\chi(G)$ -coloring of G , we could simply use the fact that the nodes sharing one color form an independent set and execute all transactions of a single color in parallel and the colors sequentially. Since no two neighbors are in an independent set and resources are edges, all conflicts are resolved. Consequently, the makespan k is at most $\chi(G)$.

On the other hand, the makespan k must be at least $\chi(G)$: Since each transaction (i.e., node) locks all required resources (i.e., adjacent edges) for at least one round, no schedule could do better than serve a (maximum) independent set in parallel while all other transactions wait. However, by definition of the chromatic number $\chi(G)$, V cannot be split into less than $\chi(G)$ independent sets, meaning that $k \geq \chi(G)$. Therefore $k = \chi(G)$.

In other words, if we could compute a schedule using $k' < k^{(\log k)/25}$ rounds in polynomial time, we knew that

$$\chi(G) = k \leq k' < k^{(\log k)/25} = \chi(G)^{(\log \chi(G))/25}.$$

□

Remarks:

- The theorem holds for a central contention manager, knowing all transactions and all potential conflicts. Clearly, the *online* problem, where conflicts remain unknown until they occur, is even harder. Furthermore, the distributed nature of contention managers also makes the problem even more difficult.
- If resource acquisition does not take zero time, the connection between the problems is not a direct equivalence. However, the same proof technique shows that it is NP-hard to compute a polynomial approximation, i.e., $k' \leq k^c$ for some constant $c \geq 1$.

Deterministic contention managers

Theorem 17.5 showed that even if all conflicts are known, one cannot produce schedules which makespan get close to the optimal without a lot of computation. However, we target to construct contention managers that make their decisions quickly without knowing conflicts in advance. Let us look at a couple of contention managers and investigate their throughput and progress guarantees.

- A first naive contention manager: Be aggressive! Always abort the transaction having locked the resource. Analysis: The throughput might be zero, since a livelock is possible. But the system is still obstruction free. Consider two transactions consisting of three operations. The first operation of both is a write to the same resource R . If they start concurrently, they will abort each other infinitely often.
- A smarter contention manager: Approximate the work done. Assume before a start (also before a restart after an abort) a transaction gets a unique timestamp. The older transaction, which is believed to have already performed more work, should win the conflict.

Analysis: Clearly, the oldest transaction will always run until commit without interruption. Thus we have lock-freedom, since at least one transaction makes progress at any time. In other words, at least one processor is always busy executing a transaction until its commit. Thus, the bound says that all transactions are executed sequentially. How about the competitive ratio? We have s resources and n transactions starting at the same time. For simplicity, assume every transaction T_i needs to lock at least one resource for a constant fraction c of its execution time t_{T_i} . Thus, at most s transactions can run concurrently from start until commit without (possibly) facing a conflict (if $s + 1$ transactions run at the same time, at least two of them lock the same resource). Thus, the makespan of an optimal contention manager is at least: $\sum_{i=0}^n \frac{c \cdot t_{T_i}}{s}$. The makespan of our timestamping algorithm is at most the duration of a sequential execution, i.e. the sum of the lengths of all transactions: $\sum_{i=0}^n t_{T_i}$. The competitive ratio is:

$$\frac{\sum_{i=0}^n t_{T_i}}{\sum_{i=0}^n \frac{c \cdot t_{T_i}}{s}} = \frac{s}{c} = O(s).$$

Remarks:

- Unfortunately, in most relevant cases the number of resources is larger than the number of cores, i.e., $s > n$. Thus, our timestamping algorithm only guarantees sequential execution, whereas the optimal might execute all transactions in parallel.

Are there contention managers that guarantee more than sequential execution, if a lot of parallelism is possible? If we have a powerful adversary, that can change the required resources after an abort, the analysis is tight. Though we restrict to deterministic algorithms here, the theorem also holds for randomized contention managers.

Theorem 17.6. *Suppose n transactions start at the same time and the adversary is allowed to alter the resource requirement of any transaction (only) after an abort, then the competitive ratio of any deterministic contention manager is $\Omega(n)$.*

Proof. Assume we have n resources. Suppose all transactions consist of two operations, such that conflicts arise, which force the contention manager to

abort one of the two transactions T_{2i-1}, T_{2i} for every $i < n/2$. More precisely, transaction T_{2i-1} writes to resource R_{2i-1} and to R_{2i} afterwards. Transaction T_{2i} writes to resource R_{2i} and to R_{2i-1} afterwards. Clearly, any contention manager has to abort $n/2$ transactions. Now the adversary tells each transaction which did not finish to adjust its resource requirements and write to resource R_0 as their first operation. Thus, for any deterministic contention manager the $n/2$ aborted transactions must execute sequentially and the makespan of the algorithm becomes $\Omega(n)$.

The optimal strategy first schedules all transactions that were aborted and in turn aborts the others. Since the now aborted transactions do not change their resource requirements, they can be scheduled in parallel. Hence the optimal makespan is 4, yielding a competitive ratio of $\Omega(n)$. \square

Remarks:

- The prove can be generalized to show that the ratio is $\Omega(s)$ if s resources are present, matching the previous upper bound.
- But what if the adversary is not so powerful, i.e., a transaction has a fixed set of needed resources?

The analysis of algorithm timestamp is still tight. Consider the dining philosophers problem: Suppose all transactions have length n and transaction i requires its first resource R_i at time 1 and its second R_{i+1} (except T_n , which only needs R_n) at time $n - i$. Thus, each transaction T_i potentially conflicts with transaction T_{i-1} and transaction T_{i+1} . Let transaction i have the i^{th} oldest timestamp. At time $n - i$ transaction $i + 1$ with $i \geq 1$ will get aborted by transaction i and only transaction 1 will commit at time n . After every abort transaction i restarts 1 time unit before transaction $i - 1$. Since transaction $i - 1$ acquires its second resource $i - 1$ time units before its termination, transaction $i - 1$ will abort transaction i at least $i - 1$ times. After $i - 1$ aborts transaction i may commit. The total time until the algorithm is done is bounded by the time transaction n stays in the system, i.e., at least $\sum_{i=1}^n (n - i) = \Omega(n^2)$. An optimal schedule requires only $\mathcal{O}(n)$ time: First schedule all transactions with even indices, then the ones with odd indices.

- Let us try to approximate the work done differently. The transaction, which has performed more work should win the conflict. A transaction counts the number of accessed resources, starting from 0 after every restart. The transaction which has acquired more resources, wins the conflict. In case both have accessed the same number of resources, the transaction having locked the resource may proceed and the other has to wait.

Analysis: Deadlock possible: Transaction A and B start concurrently. Transaction A writes to R_1 as its first operation and to R_2 as its second operation. Transaction B writes to the resources in opposite order.

Randomized contention managers

Though the lower bound of the previous section (Theorem 17.6) is valid for both deterministic and randomized schemes, let us look at a randomized approach:

Each transaction chooses a random priority in $[1, n]$. In case of a conflict, the transaction with lower priority gets aborted. (If both conflicting transactions have the same priority, both abort.)

Additionally, if a transaction A was aborted by transaction B , it waits until transaction B committed or aborted, then transaction A restarts and draws a new priority.

Analysis: Assume the adversary cannot change the resource requirements, otherwise we cannot show more than a competitive ratio of n , see Theorem 17.6. Assume if two transactions A and B (potentially) conflict (i.e., write to the same resource), then they require the resource for at least a fraction c of their running time. We assume a transaction T potentially conflicts with d_T other transactions. Therefore, if a transaction has highest priority among these d_T transactions, it will abort all others and commit successfully. The chance that for a transaction T a conflicting transaction chooses the same random number is $(1 - 1/n)^{d_T} > (1 - 1/n)^n \approx 1/e$. The chance that a transaction chooses the largest random number and no other transaction chose this number is thus at least $1/d_T \cdot 1/e$. Thus, for any constant $c \geq 1$, after choosing $e \cdot d_T \cdot c \cdot \ln n$ random numbers the chance that transaction T has committed successfully is

$$1 - \left(1 - \frac{1}{e \cdot d_T}\right)^{e \cdot d_T \cdot c \cdot \ln n} \approx 1 - e^{-c \ln n} = 1 - \frac{1}{n^c}.$$

Assuming that the longest transaction takes time t_{max} , within that time a transaction either commits or aborts and chooses a new random number. The time to choose $e \cdot t_{max} \cdot c \cdot \ln n$ numbers is thus at most $e \cdot t_{max} \cdot d_T \cdot c \cdot \ln n = O(t_{max} \cdot d_T \cdot \ln n)$. Therefore, with high probability each transaction makes progress within a finite amount of time, i.e., our algorithm ensures wait freedom. Furthermore, the competitive ratio of our randomized contention manger for the previously considered dining philosophers problem is w.h.p. only $\mathcal{O}(\ln n)$, since any transaction only conflicts with two other transactions.

Chapter 18

Dominating Set

In this chapter we present another randomized algorithm that demonstrates the power of randomization to break symmetries. We study the problem of finding a small dominating set of the network graph. As it is the case for the MIS, an efficient dominating set algorithm can be used as a basic building block to solve a number of problems in distributed computing. For example, whenever we need to partition the network into a small number of local clusters, the computation of a small dominating set usually occurs in some way. A particularly important application of dominating sets is the construction of an efficient backbone for routing.

Definition 18.1 (Dominating Set). *Given an undirected graph $G = (V, E)$, a dominating set is a subset $S \subseteq V$ of its nodes such that for all nodes $v \in V$, either $v \in S$ or a neighbor u of v is in S .*

Remarks:

- It is well-known that computing a dominating set of minimal size is NP-hard. We therefore look for approximation algorithms, that is, algorithms which produce solutions which are optimal up to a certain factor.
- Note that every MIS (cf. Chapter 8) is a dominating set. In general, the size of every MIS can however be larger than the size of an optimal minimum dominating set by a factor of $\Omega(n)$. As an example, connect the centers of two stars by an edge. Every MIS contains all the leaves of at least one of the two stars whereas there is a dominating set of size 2.

All the dominating set algorithms that we study throughout this chapter operate in the following way. We start with $S = \emptyset$ and add nodes to S until S is a dominating set. To simplify presentation, we color nodes according to their state during the execution of an algorithm. We call nodes in S *black*, nodes which are covered (neighbors of nodes in S) *gray*, and all uncovered nodes *white*. By $W(v)$, we denote the set of white nodes among the direct neighbors of v , including v itself. We call $w(v) = |W(v)|$ the *span* of v .

18.1 Sequential Greedy Algorithm

Intuitively, to end up with a small dominating set S , nodes in S need to cover as many neighbors as possible. It is therefore natural to add nodes v with a large span $w(v)$ to S . This idea leads to a simple greedy algorithm:

Algorithm 63 Greedy Algorithm

```

1:  $S := \emptyset$ ;
2: while  $\exists$  white nodes do
3:    $v := \{v \mid w(v) = \max_u \{w(u)\}\}$ ;
4:    $S := S \cup v$ ;
5: end while
  
```

Theorem 18.2. *The Greedy Algorithm computes a $\ln \Delta$ -approximation, that is, for the computed dominating set S and an optimal dominating set S^* , we have*

$$\frac{|S|}{|S^*|} \leq \ln \Delta.$$

Proof. Each time, we choose a new node of the dominating set (each greedy step), we have cost 1. Instead of letting this node pay the whole cost, we distribute the cost equally among all newly covered nodes. Assume that node v , chosen in line 3 of the algorithm, is white itself and that its white neighbors are v_1, v_2, v_3 , and v_4 . In this case each of the 5 nodes v and v_1, \dots, v_4 get charged $1/5$. If v is chosen as a gray node, only the nodes v_1, \dots, v_4 get charged (they all get $1/4$).

Now, assume that we know an optimal dominating set S^* . By the definition of dominating sets, to each node which is not in S^* , we can assign a neighbor from S^* . By assigning each node to exactly one neighboring node of S^* , the graph is decomposed into stars, each having a dominator (node in S^*) as center and non-dominators as leaves. Clearly, the cost of an optimal dominating set is 1 for each such star. In the following, we show that the amortized cost (distributed costs) of the greedy algorithm is at most $\ln \Delta + 2$ for each star. This suffices to prove the theorem.

Consider a single star with center $v^* \in S^*$ before choosing a new node u in the greedy algorithm. The number of nodes that become dominated when adding u to the dominating set is $w(u)$. Thus, if some white node v in the star of v^* becomes gray or black, it gets charged $1/w(u)$. By the greedy condition, u is a node with maximal span and therefore $w(u) \geq w(v^*)$. Thus, v is charged at most $1/w(v^*)$. After becoming gray, nodes do not get charged any more. Therefore first node that is covered in the star of v^* gets charged at most $1/(d(v^*) + 1)$. Because $w(v^*) \geq d(v^*)$ when the second node is covered, the second node gets charged at most $1/d(v^*)$. In general, the i^{th} node that is covered in the star of v^* gets charged at most $1/(d(v^*) + i - 2)$. Thus, the total amortized cost in the star of v^* is at most

$$\frac{1}{d(v^*) + 1} + \frac{1}{d(v^*)} + \dots + \frac{1}{2} + \frac{1}{1} = H(d(v^*) + 1) \leq H(\Delta + 1) < \ln(\Delta) + 2$$

where Δ is the maximal degree of G and where $H(n) = \sum_{i=1}^n 1/i$ is the n^{th} number. \square

Remarks:

- One can show that unless $\text{NP} \subseteq \text{DTIME}(n^{O(\log \log n)})$, no polynomial-time algorithm can approximate the minimum dominating set problem better than $(1 - o(1)) \cdot \ln \Delta$. Thus, unless $\text{P} \approx \text{NP}$, the approximation ratio of the simple greedy algorithm is optimal (up to lower order terms).

18.2 Distributed Greedy Algorithm

For a distributed algorithm, we use the following observation. The span of a node can only be reduced if any of the nodes at distance at most 2 is included in the dominating set. Therefore, if the span of node v is greater than the span of any other node at distance at most 2 from v , the greedy algorithm chooses v before any of the nodes at distance at most 2. This leads to a very simple distributed version of the greedy algorithm. Every node v executes the following algorithm.

Algorithm 64 Distributed Greedy Algorithm (at node v):

```

1: while  $v$  has white neighbors do
2:   compute span  $w(v)$ ;
3:   send  $w(v)$  to nodes at distance at most 2;
4:   if  $w(v)$  largest within distance 2 (ties are broken by IDs) then
5:     join dominating set
6:   end if
7: end while

```

Theorem 18.3. *Algorithm 64 computes a dominating set of size at most $\ln \Delta + 2$ times the size of an optimal dominating set in $\mathcal{O}(n)$ rounds.*

Proof. The approximation quality follows directly from the above observation and the analysis of the *greedy algorithm*. The time complexity is at most linear because in every iteration of the while loop, at least one node is added to the dominating set and because one iteration of the while loop can be implemented in a constant number of rounds. \square

The approximation ratio of the above distributed algorithm is best possible (unless $\text{P} \approx \text{NP}$ or unless we allow local computations to be exponential). However, the time complexity is very bad. In fact, there really are graphs on which in each iteration of the while loop, only one node is added to the dominating set. As an example, consider a graph as in Figure 18.1. An optimal dominating set consists of all nodes on the center axis. The *distributed greedy algorithm* computes an optimal dominating set, however, the nodes are chosen sequentially from left to right. Hence, the running time of the algorithm on the graph of Figure 18.1 is $\Omega(\sqrt{n})$. Below, we will see that there are graphs on which Algorithm 64 even needs $\Omega(n)$ rounds.

The problem of the graph of Figure 18.1 is that there is a long path of descending degrees (spans). Every node has to wait for the neighbor to the left. Therefore, we want to change the algorithm in such a way that there are no long paths of descending spans. Allowing for an additional factor 2 in

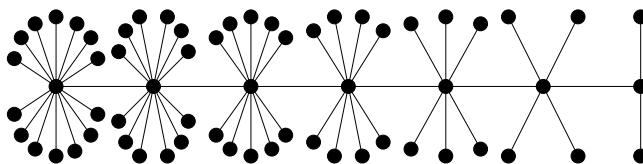


Figure 18.1: Distributed greedy algorithm: Bad example

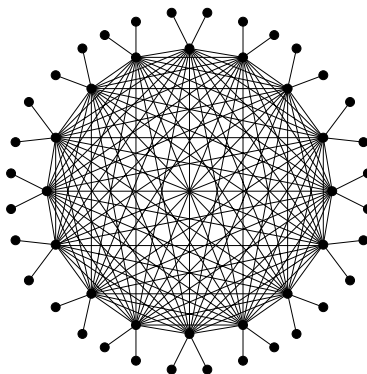


Figure 18.2: Distributed greedy algorithm with rounded spans: Bad example

the approximation ratio, we can round all spans to the next power of 2 and let the greedy algorithm take a node with a maximal rounded span. In this case, a path of strictly descending rounded spans has at most length $\log n$. For the distributed version, this means that nodes whose rounded span is maximal within distance 2 are added to the dominating set. Ties are again broken by unique node IDs. If node IDs are chosen at random, the time complexity for the graph of Figure 18.1 is reduced from $\Omega(\sqrt{n})$ to $\mathcal{O}(\log n)$.

Unfortunately, there still is a problem remaining. To see this, we consider Figure 18.2. The graph of Figure 18.2 consists of a clique with $n/3$ nodes and two leaves per node of the clique. An optimal dominating set consists of all the $n/3$ nodes of the clique. Because they all have distance 1 from each other, the described distributed algorithm only selects one in each while iteration (the one with the largest ID). Note that as soon as one of the nodes is in the dominating set, the span of all remaining nodes of the clique is 2. They do not have common neighbors and therefore there is no reason not to choose all of them in parallel. However, the time complexity of the simple algorithm is $\Omega(n)$. In order to improve this example, we need an algorithm that can choose many nodes simultaneously as long as these nodes do not interfere too much, even if they are neighbors. In Algorithm 65, $N(v)$ denotes the set of neighbors of v (including v itself) and $N_2(v) = \bigcup_{u \in N(V)} N(u)$ are the nodes at distance at most 2 of v . As before, $W(v) = \{u \in N(v) : u \text{ is white}\}$ and $w(v) = |W(v)|$. It is clear that if Algorithm 65 terminates, it computes a valid dominating set. We will now show that the computed dominating set is small and that the algorithm terminates quickly.

Theorem 18.4. *Algorithm 65 computes a dominating set of size at most $(6 \cdot \ln \Delta + 12) \cdot |S^*|$, where S^* is an optimal dominating set.*

Algorithm 65 Fast Distributed Dominating Set Algorithm (at node v):

```

1:  $W(v) := N(v)$ ;  $w(v) := |W(v)|$ ;
2: while  $W(v) \neq \emptyset$  do
3:    $\tilde{w}(v) := 2^{\lfloor \log_2 w(v) \rfloor}$ ; // round down to next power of 2
4:    $\hat{w}(v) := \max_{u \in N_2(v)} \tilde{w}(u)$ ;
5:   if  $\tilde{w}(v) = \hat{w}(v)$  then  $v.active := \mathbf{true}$  else  $v.active := \mathbf{false}$  end if;
6:   compute support  $s(v) := |\{u \in N(v) : u.active = \mathbf{true}\}|$ ;
7:    $\hat{s}(v) := \max_{u \in W(v)} s(u)$ ;
8:    $v.candidate := \mathbf{false}$ ;
9:   if  $v.active$  then
10:     $v.candidate := \mathbf{true}$  with probability  $1/\hat{s}(v)$ 
11:   end if;
12:   compute  $c(v) := |\{u \in W(v) : u.candidate = \mathbf{true}\}|$ ;
13:   if  $v.candidate$  and  $\sum_{u \in W(v)} c(u) \leq 3w(v)$  then
14:     node  $v$  joins dominating set
15:   end if
16:    $W(v) := \{u \in N(v) : u \text{ is white}\}$ ;  $w(v) := |W(v)|$ ;
17: end while

```

Proof. The proof is a bit more involved but analogous to the analysis of the approximation ratio of the greedy algorithm. Every time, we add a new node v to the dominating set, we distribute the cost among v (if it is still white) and its white neighbors. Consider an optimal dominating set S^* . As in the analysis of the greedy algorithm, we partition the graph into stars by assigning every node u not in S^* to a neighbor v^* in S^* . We want to show that the total distributed cost in the star of every $v^* \in S^*$ is at most $6H(\Delta + 1)$.

Consider a node v that is added to the dominating set by Algorithm 65. Let $W(v)$ be the set of white nodes in $N(v)$ when v becomes a dominator. For a node $u \in W(v)$ let $c(u)$ be the number of candidate nodes in $N(u)$. We define $C(v) = \sum_{u \in W(v)} c(u)$. Observe that $C(v) \leq 3w(v)$ because otherwise v would not join the dominating set in line 15. When adding v to the dominating set, every newly covered node $u \in W(v)$ is charged $3/(c(u)w(v))$. This compensates the cost 1 for adding v to the dominating set because

$$\sum_{u \in W(v)} \frac{3}{c(u)w(v)} \geq w(v) \cdot \frac{3}{w(v) \cdot \sum_{u \in W(v)} c(u)/w(v)} = \frac{3}{C(v)/w(v)} \geq 1.$$

The first inequality follows because it can be shown that for $\alpha_i > 0$, $\sum_{i=1}^k 1/\alpha_i \geq k/\bar{\alpha}$ where $\bar{\alpha} = \sum_{i=1}^k \alpha_i/k$.

Now consider a node $v^* \in S^*$ and assume that a white node $u \in W(v^*)$ turns gray or black in iteration t of the while loop. We have seen that u is charged $3/(c(u)w(v))$ for every node $v \in N(u)$ that joins the dominating set in iteration t . Since a node can only join the dominating set if its span is largest up to a factor of two within two hops, we have $w(v) \geq w(v^*)/2$ for every node $v \in N(u)$ that joins the dominating set in iteration t . Because there are at most $c(u)$ such nodes, the charge of u is at most $6/w(v^*)$. Analogously to the sequential greedy algorithm, we now get that the total cost in the star of a node $v^* \in S^*$ is at

most

$$\sum_{i=1}^{|N(v^*)|} \frac{6}{i} \leq 6 \cdot H(|N(v^*)|) \leq 6 \cdot H(\Delta + 1) = 6 \cdot \ln \Delta + 12.$$

□

To bound the time complexity of the algorithm, we first need to prove the following lemma.

Lemma 18.5. *Consider an iteration of the while loop. Assume that a node u is white and that $2s(u) \geq \max_{v \in \mathcal{C}(u)} \hat{s}(v)$ where $\mathcal{C}(u) = \{v \in N(u) : v.\text{candidate} = \mathbf{true}\}$. Then, the probability that u becomes dominated (turns gray or black) in the considered while loop iteration is larger than $1/9$.*

Proof. Let $D(u)$ be the event that u becomes dominated in the considered while loop iteration, i.e., $D(u)$ is the event that u changes its color from white to gray or black. Thus, we need to prove that $\Pr[D(u)] > 1/9$. We can write this probability as

$$\Pr[D(u)] = \Pr[c(u) > 0] \cdot \Pr[D(u)|c(u) > 0] + \underbrace{\Pr[c(u) = 0] \cdot \Pr[D(u)|c(u) = 0]}_{=0}.$$

It is therefore sufficient to lower bound the probabilities $\Pr[c(u) > 0]$ and $\Pr[D(u)|c(u) > 0]$. We have $2s(u) \geq \max_{v \in \mathcal{C}(u)} \hat{s}(v)$. Therefore, in line 10, each of the $s(u)$ active nodes $v \in N(u)$ becomes a candidate node with probability $1/\hat{s}(v) \geq 1/(2s(u))$. The probability that at least one of the $s(u)$ active nodes in $N(u)$ becomes a candidate therefore is

$$\Pr[c(u) > 0] > 1 - \left(1 - \frac{1}{2s(u)}\right)^{s(u)} > 1 - \frac{1}{\sqrt{e}} > \frac{1}{3}.$$

We used that for $x \geq 1$, $(1 - 1/x)^x < 1/e$. We next want to bound the probability $\Pr[D(u)|c(u) > 0]$ that at least one of the $c(u)$ candidates in $N(u)$ joins the dominating set. We have

$$\Pr[D(u)|c(u) > 0] \geq \min_{v \in N(u)} \Pr[v \text{ joins dominating set} | v.\text{candidate} = \mathbf{true}].$$

Consider some node v and let $C(v) = \sum_{v' \in W(v)} c(v')$. If v is a candidate, it joins the dominating set if $C(v) \leq 3w(v)$. We are thus interested in the probability $\Pr[C(v) \leq 3w(v) | v.\text{candidate} = \mathbf{true}]$. Assume that v is a candidate. Let $c'(v') = c(v') - 1$ be the number of candidates in $N(v') \setminus \{v\}$. For a node $v' \in W(v)$, $c'(v')$ is upper bounded by a binomial random variable $\text{Bin}(s(v') - 1, 1/s(v'))$ with expectation $(s(v') - 1)/s(v')$. We therefore have

$$\mathbb{E}[c(v') | v.\text{candidate} = \mathbf{true}] = 1 + \mathbb{E}[c'(v')] = 1 + \frac{s(v') - 1}{s(v')} < 2.$$

By linearity of expectation, we hence obtain

$$\begin{aligned} \mathbb{E}[C(v) | v.\text{candidate} = \mathbf{true}] &= \sum_{v' \in W(v)} \mathbb{E}[c(v') | v.\text{candidate} = \mathbf{true}] \\ &< 2w(v). \end{aligned}$$

We can now use Markov's inequality to bound the probability that $C(v)$ becomes too large:

$$\Pr[C(v) > 3w(v) | v.\text{candidate} = \mathbf{true}] < \frac{2}{3}.$$

Combining everything, we get

$$\begin{aligned} & \Pr[v \text{ joins dom. set} | v.\text{candidate} = \mathbf{true}] \\ &= \Pr[C(v) \leq 3w(v) | v.\text{candidate} = \mathbf{true}] > \frac{1}{3} \end{aligned}$$

and hence

$$\Pr[D(u)] = \Pr[c(u) > 0] \cdot \Pr[D(u) | c(u) > 0] > \frac{1}{3} \cdot \frac{1}{3} = \frac{1}{9}.$$

□

Theorem 18.6. *In expectation, Algorithm 65 terminates in $\mathcal{O}(\log^2 \Delta \cdot \log n)$ rounds.*

Proof. First observe that every iteration of the while loop can be executed in a constant number of rounds. Consider the state after t iterations of the while loop. Let $\tilde{w}_{\max}(t) = \max_{v \in V} \tilde{w}(v)$ be the maximal span rounded down to the next power of 2 after t iterations. Further, let $s_{\max}(t)$ be the maximal support $s(v)$ of any node v for which there is a node $u \in N(v)$ with $w(u) \geq \tilde{w}_{\max}(t)$ after t while loop iterations. Observe that all nodes v with $w(v) \geq \tilde{w}_{\max}(t)$ are active in iteration $t+1$ and that as long as the maximal rounded span $\tilde{w}_{\max}(t)$ does not change, $s_{\max}(t)$ can only get smaller with increasing t . Consider the pair $(\tilde{w}_{\max}, s_{\max})$ and define a relation \prec such that $(w', s') \prec (w, s)$ iff $w' < w$ or $w = w'$ and $s' \leq s/2$. From the above observations, it follows that

$$(\tilde{w}_{\max}(t), s_{\max}(t)) \prec (\tilde{w}_{\max}(t'), s_{\max}(t')) \implies t > t'. \quad (18.1)$$

For a given time t , let $T(t)$ be the first time for which

$$(\tilde{w}_{\max}(T(t)), s_{\max}(T(t))) \prec (\tilde{w}_{\max}(t), s_{\max}(t)).$$

We first want to show that for all t ,

$$\mathbb{E}[T(t) - t] = O(\log n). \quad (18.2)$$

Let us look at the state after t while loop iterations. By Lemma 18.5, every white node u with support $s(u) \geq s_{\max}(t)/2$ will be dominated after the following while loop iteration with probability larger than $1/9$. Consider a node u that satisfies the following three conditions:

- (1) u is white
- (2) $\exists v \in N(u) : w(v) \geq \tilde{w}_{\max}(t)$
- (3) $s(u) \geq s_{\max}(t)/2$.

As long as u satisfies all three conditions, the probability that u becomes dominated is larger than $1/9$ in every while loop iteration. Hence, after $t + \tau$ iterations (from the beginning), u is dominated or does not satisfy (2) or (3) with probability larger than $(8/9)^\tau$. Choosing $\tau = \log_{9/8}(2n)$, this probability becomes $1/(2n)$. There are at most n nodes u satisfying Conditions (1) – (3). Therefore, applying union bound, we obtain that with probability more than $1/2$, there is no white node u satisfying Conditions (1) – (3) at time $t + \log_{9/8}(2n)$. Equivalently, with probability more than $1/2$, $T(t) \leq t + \log_{9/8}(2n)$. Analogously, we obtain that with probability more than $1/2^k$, $T(t) \leq t + k \log_{9/8}(2n)$. We then have

$$\begin{aligned} \mathbb{E}[T(t) - t] &= \sum_{\tau=1}^{\infty} \Pr[T(t) - t = \tau] \cdot \tau \\ &\leq \sum_{k=1}^{\infty} \left(\frac{1}{2^k} - \frac{1}{2^{k+1}} \right) \cdot k \log_{9/8}(2n) = \log_{9/8}(2n) \end{aligned}$$

and thus Equation (18.2) holds.

Let $t_0 = 0$ and $t_i = T(t_{i-1})$ for $i = 1, \dots, k$, where $t_k = \min_t \tilde{w}_{\max}(t) = 0$. Because $\tilde{w}_{\max}(t) = 0$ implies that $w(v) = 0$ for all $v \in V$ and that we therefore have computed a dominating set, by Equations (18.1) and (18.2) (and linearity of expectation), the expected number of rounds until Algorithm 65 terminates is $\mathcal{O}(k \cdot \log n)$. Since $\tilde{w}_{\max}(t)$ can only have $\lfloor \log \Delta \rfloor$ different values and because for a fixed value of $\tilde{w}_{\max}(t)$, the number of times $s_{\max}(t)$ can be decreased by a factor of 2 is at most $\log \Delta$ times, we have $k \leq \log^2 \Delta$. \square

Remarks:

- It is not hard to show that Algorithm 65 even terminates in $\mathcal{O}(\log^2 \Delta \cdot \log n)$ rounds with probability $1 - 1/n^c$ for an arbitrary constant c .
- Using the median of the supports of the neighbors instead of the maximum in line 8 results in an algorithm with time complexity $\mathcal{O}(\log \Delta \cdot \log n)$. With another algorithm, this can even be slightly improved to $\mathcal{O}(\log^2 \Delta)$.
- One can show that $\Omega(\log \Delta / \log \log \Delta)$ rounds are necessary to obtain an $\mathcal{O}(\log \Delta)$ -approximation.
- It is not known whether there is a fast deterministic approximation algorithm. This is an interesting and important open problem. The best deterministic algorithm known to achieve an $\mathcal{O}(\log \Delta)$ -approximation has time complexity $2^{\mathcal{O}(\sqrt{\log n})}$.

Chapter 19

Routing

19.1 Array

(Routing is important for any distributed system. This chapter is only an introduction into routing; we will see other facets of routing in a next chapter.)

Definition 19.1 (Routing). *We are given a graph and a set of routing requests, each defined by a source and a destination node.*

Definition 19.2 (One-to-one, Permutation). *In a one-to-one routing problem, each node is the source of at most one packet and each node is the destination of at most one packet. In a permutation routing problem, each node is the source of exactly one packet and each node is the destination of exactly one packet.*

Remarks:

- Permutation routing is a special case of one-to-one routing.

Definition 19.3 (Store and Forward Routing). *The network is synchronous. In each step, at most two packets (one in each direction) can be sent over each link.*

Remarks:

- If two packets want to follow the same link, then one is queued (stored) at the sending node. This is known as contention.

Algorithm 66 Greedy on Array

An array is a linked list of n nodes; that is, node i is connected with nodes $i - 1$ and $i + 1$, for $i = 2, \dots, n - 1$. With the greedy algorithm, each node injects its packet at time 0. At each step, each packet that still needs to move rightward or leftward does so.

Theorem 19.4 (Analysis). *The greedy algorithm terminates in $n - 1$ steps.*

Proof. By induction two packets will never contend for the same link. Then each packet arrives at its destination in d steps, where d is the distance between source and destination. \square

Remarks:

- Unfortunately, only the array (or the ring) allows such a simple contention-free analysis. Already in a tree (with nodes of degree 3 or more) there might be two packets arriving at the same step at the same node, both want to leave on the same link, and one needs to be queued. In a “Mercedes-Benz” graph $\Omega(n)$ packets might need to be queued. We will study this problem in the next section.
- There are many strategies for scheduling packets contending for the same edge (e.g. “farthest goes first”); these queuing strategies have a substantial impact on the performance of the algorithm.

19.2 Mesh

Algorithm 67 Greedy on Mesh

A mesh (a.k.a. grid, matrix) is a two-dimensional array with m columns and m rows ($n = m^2$). Packets are routed to their correct column (on the row in greedy array style), and then to their correct row. The farthest packet will be given priority.

Theorem 19.5 (Analysis). *In one-to-one routing, the greedy algorithm terminates in $2m - 2$ steps.*

Proof. First note that packets in the first phase of the algorithm do not interfere with packets in the second phase of the algorithm. With Theorem 19.4 each packet arrives at its correct column in $m - 1$ steps. (Some packets may arrive at their turning node earlier, and already start the second phase; we will not need this in the analysis.) We need the following Lemma for the second phase of the algorithm.

Lemma 19.6 (Many-to-One on Array, Lemma 1.5 in Leighton Section 1.7). *We are given an array with n nodes. Each node is a destination for at most one packet (but may be the source of many). If edge contention is resolved by farthest-to-go (FTG), the algorithm terminates in $n - 1$ steps.*

Leighton Section 1.7 Lemma 1.5. Leftward moving packets and rightward moving packets never interfere; so we can restrict ourselves to rightward moving packets. We name the packets with their destination node. Since the queuing strategy is FTG, packet i can only be stopped by packets $j > i$. Note that a packet i may be contending with the same packet j several times. However, packet i will either find its destination “among” the higher packets, or directly after the last of the higher packets. More formally, after k steps, packets $j, j + 1, \dots, n$ do not need links $1, \dots, l$ anymore, with $k = n - j + l$. Proof by induction: Packet n has the highest priority: After k steps it has escaped

the first k links. Packet $n - 1$ can therefore use link l in step $l + 1$, and so on. Packet i not needing link i in step $k = n$ means that packet i has arrived at its destination node i in step $n - 1$ or earlier. \square

Lemma 19.6 completes the proof. \square

Remarks:

- A $2m - 2$ time bound is the best we can hope for, since the distance between the two farthest nodes in the mesh is exactly $2m - 2$.
- One thing still bugs us: The greedy algorithm might need queues in the order of m . And queues are expensive! In the next section, we try to bring the queue size down!

19.3 Routing in the Mesh with Small Queues

(First we look at a slightly simpler problem.)

Definition 19.7 (Random Destination Routing). *In a random destination routing problem, each node is the source of at most one packet with destination chosen uniformly at random.*

Remarks:

- Random destination routing is not one-to-one routing. In the worst case, a node can be destination for all n packets, but this case is very unlikely (with probability $1/n^{n-1}$)
- We study algorithm 19.2, but this time in the random destination model. Studying the random destination model will give us a deeper understanding of routing... and distributed computing in general!

Theorem 19.8 (Random destination analysis of algorithm 19.2). *If destinations are chosen at random the maximum queue size is $O(\log n / \log \log n)$ with high probability. (With high probability means with probability at least $1 - O(1/n)$.)*

Proof. We can restrict ourselves to column edges because there will not be any contention at row edges. Let us consider the queue for a north-bound column edge. In each step, there might be three packets arriving (from south, east, west). Since each arriving south packet will be forwarded north (or consumed when the node is the destination), the queue size can only grow from east or west packets – packets that are “turning” at the node. Hence the queue size of a node is always bounded by the number of packets turning at the node. A packet only turns at a node u when it is originated at a node in the same row as u (there are only m nodes in the row). Packets have random destinations, so the probability to turn for each of these packets is $1/m$ only. Thus the probability P that r or more packets turn in some particular node u is at most

$$P \leq \binom{m}{r} \left(\frac{1}{m}\right)^r$$

(The factor $(1 - 1/m)^{m-r}$ is not present because the event “exactly r ” includes the event “more than r ” already.) Using

$$\binom{x}{y} < \left(\frac{xe}{y}\right)^y, \text{ for } 0 < y < x$$

we directly get

$$P < \left(\frac{me}{r}\right)^r \left(\frac{1}{m}\right)^r = \left(\frac{e}{r}\right)^r$$

Hence most queues do not grow larger than $\mathcal{O}(1)$. Also, when we choose $r := \frac{e \log n}{\log \log n}$ we can show $P = o(1/n^2)$. The probability that any of the $4n$ queues ever exceeds r is less than $1 - (1 - P)^{4n} = o(1/n)$. \square

Remarks:

- OK. We got a bound on the queue size. Now what about time complexity?!? The same analysis as for one-to-one routing applies. The probability that a node sees “many” packets in phase 2 is small... it can be shown that the algorithm terminates in $\mathcal{O}(m)$ time with high probability.
- In fact, maximum queue sizes are likely to be a lot less than logarithmic. The reason is the following: Though $\Theta(\log n / \log \log n)$ packets might turn at some node, these turning packets are likely to be spread in time. Early arriving packets might use gaps and do not conflict with late arriving packets. With a much more elaborate method (using the so-called “wide-channel” model) one can show that there will never be more than four(!) packets in any queue (with high probability only, of course).
- Unfortunately, the above analysis only works for random destination problems. Question: Can we devise an algorithm that uses small queues only but for any one-to-one routing problem? Answer: Yes, we can! In the simplest form we can use a clever trick invented by Leslie Valiant: Instead of routing the packets directly on their row-column path, we route each packet to a randomly chosen intermediate node (on the row-column path), and from there to the destination (again on the row-column path). Valiant’s trick routes all packets in $\mathcal{O}(m)$ time (with high probability) and only needs queues of size $\mathcal{O}(\log n)$. Instead of choosing a random intermediate node one can choose a random node that is more or less in the direction of the destination, solving any one-to-one routing problem in $2m + \mathcal{O}(\log n)$ time with only constant-size queues. You don’t wanna know the details...
- What about no queues at all?!?

19.4 Hot-Potato Routing

Definition 19.9 (Hot-Potato Routing). *Like the store-and-forward model the hot-potato model is synchronous and at most two packets (one in each direction) can be sent over a link. However, contending packets cannot be stored; instead all but one contending packet must be sent over a “wrong link” (known as deflection) immediately, since the hot-potato model does not allow queuing.*

Remarks:

- Don't burn your fingers with "hot-potato" packets. If you get one you better forward it directly!
- A node with degree δ receives up to δ packets at the beginning of each step – since the node has δ links, it can forward all of them, but unfortunately not all in the right direction.
- Hot-potato routing is easier to implement, especially on light-based networks, where you don't want to convert photons into electrons and then back again. There are a couple of parallel machines that use the hot-potato paradigm to simplify and speed up routing.
- How bad does hot-potato routing get (in the random or the one-to-one model)? How bad can greedy hot-potato routing (greedy: whenever there is no contention you must send a packet into the right direction) get in a worst case?

Algorithm 68 Greedy Hot-Potato Routing on a Mesh

Packets move greedy towards their destination (any good link is fine if there is more than one). If a packet gets deflected, it gets excited with probability p (we set $p = \Theta(1/m)$). An excited packet has higher priority. When being excited, a packet tries to reach the destination on the row-column path. If two excited packets contend, then the one that wants to exit the opposite link is given priority. If an excited packet fails to take its desired link it becomes normal again.

Theorem 19.10 (Analysis). *A packet will reach its destination in $\mathcal{O}(m)$ expected time.*

Sketch, full proof in Busch et al., SODA 2000. An excited packet can only be deflected at its start node (after becoming excited), and when trying to turn. In both cases, the probability to fail is only constant since other excited packets need to be at the same node at exactly the right instant. Thus the probability that an excited packet finds its destination is constant, and therefore a packet needs to "try" (to become excited) only constantly often. Since a packet tries every p 'th time it gets deflected, it only gets deflected $\mathcal{O}(1/p) = \mathcal{O}(m)$ times in expectation. Since each time it does not get deflected, it gets closer to its destination, it will arrive at the destination in $\mathcal{O}(m)$ expected time. \square

Remarks:

- It seems that at least in expectation having no memory at all does not harm the time bounds much.
- It is conjectured that one-to-one routing can be shown to have time complexity $\mathcal{O}(m)$ for this greedy hot-potato routing algorithm. However, the best known bound needs an additional logarithmic factor.

19.5 More Models

Routing comes in many flavors. We mention some of them in this section for the sake of completeness.

Store-and-forward and hot-potato routing are variants of packet-switching. In the circuit-switching model, the entire path from source to destination must be locked such that a stream of packets can be transmitted.

A packet-switching variant where more than one packet needs to be sent from source to destination in a stream is known as wormhole routing.

Static routing is when all the packets to be routed are injected at time 0. Instead, in dynamic routing, nodes may inject new packets constantly (at a certain rate). Not much is known for dynamic routing.

Instead of having a single source and a single destination for each packet as in one-to-one routing, researchers have studied many-to-one routing, where a node may be destination for many sources. The problem of many-to-one routing is that there might be congested areas in the network (areas with nodes that are destinations of many packets). Packets that can be routed around such a congested area should do that, or they increase the congestion even more. Such an algorithm was studied by Busch et al. at STOC 2000.

Also one-to-many routing (multicasting) was considered, where a source needs to send the same packet to many destinations. In one-to-many routing, packets can be duplicated whenever needed.

Nobody knows the topology of the Internet (and it is certainly not an array or a mesh!). The problem is to find short paths without storing huge routing tables at each node. There are several forms of routing (e.g. compact routing, interval routing) that study the trade-off between routing table size and quality of routing.

Also, researchers started studying the effects of mixing various queuing strategies in one network. This area of research is known as adversarial queuing theory.

And last not least there are several special networks. A mobile ad-hoc network, for example, consists of mobile nodes equipped with a wireless communication device. In such a network nodes can only communicate when they are within transmission range. Since the network is mobile (dynamic), and since the nodes are considered to be simple, a variety of new problems arise.