

# Chapter 10

## 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 3) 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.

### 10.1 Basics

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

**Definition 10.1** (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  $i^{\text{th}}$  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 10.2.** *If all generated clock pulses are valid according to Definition 10.1, the above method provides an asynchronous algorithm that behaves exactly the same way as the given synchronous algorithm.*

*Proof.* When the  $i^{\text{th}}$  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 10.1, 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  $T_{\text{tot}}$  and  $M_{\text{tot}}$  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 10.3** (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 10.4.** *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 10.1 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.

## 10.2 The Local Synchronizer $\alpha$

---

**Algorithm 40** Synchronizer  $\alpha$  (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  $\alpha$  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  $\alpha$ .

**Theorem 10.5.** *The time and message complexities of synchronizer  $\alpha$  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  $\alpha$  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  $\alpha$  was presented in a framework, mostly set up to have a common standard to discuss different synchronizers. Without the framework, synchronizer  $\alpha$  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  $\alpha$  allows for simple and fast synchronization, it produces awfully many messages. Can we do better? Yes.

## 10.3 The Global Synchronizer $\beta$

---

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

**Theorem 10.6.** *The time and message complexities of synchronizer  $\beta$  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 3, 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 ( $\alpha$ ) and a message-efficient synchronizer ( $\beta$ ), 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:  $\gamma$ .

## 10.4 The Hybrid Synchronizer $\gamma$

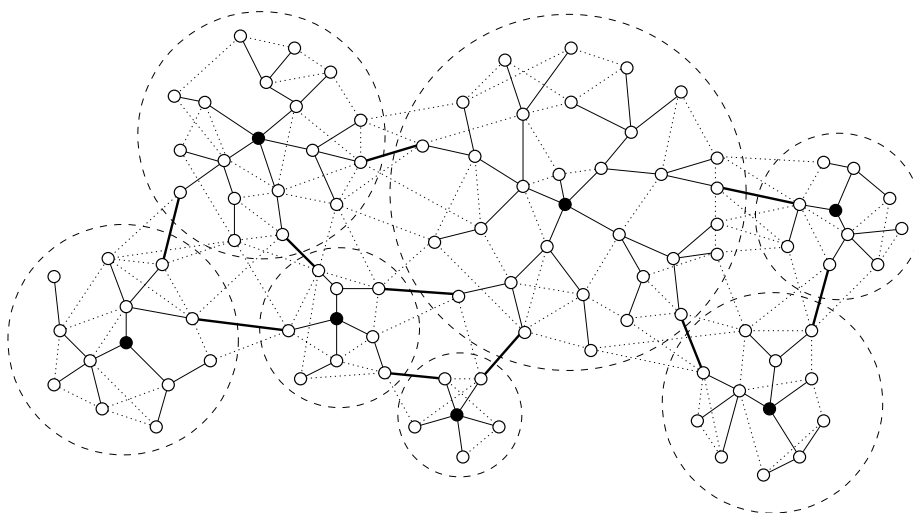


Figure 10.1: A cluster partition of a network: The dashed cycles specify the clusters, cluster leaders are black, the solid edges are the edges of the intracluster trees, and the bold solid edges are the intercluster edges

Synchronizer  $\gamma$  can be seen as a combination of synchronizers  $\alpha$  and  $\beta$ . 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 *intracluster 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 10.1 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  $\gamma$  works in two phases. In a first phase, synchronizer  $\beta$  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  $\alpha$  between clusters. A detailed description is given by Algorithm 42.

---

**Algorithm 42** Synchronizer  $\gamma$  (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 10.7.** *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  $\gamma$  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  $4k + 1 = O(k)$ .  $\square$

## 10.5 Network Partition

We will now look at the initialization phase of synchronizer  $\gamma$ . Algorithm 43 describes how to construct a partition into clusters that can be used for synchronizer  $\gamma$ . 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  $\rho$ .

---

**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  $\alpha$  and  $\beta$  as special cases of synchronizer  $\gamma$ .

**Theorem 10.8.** *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  $\rho$ . 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 10.9.** *Using  $\rho = 2$ , Algorithm 43 computes a clustering with cluster radius at most  $\log_2 n$  and with at most  $n$  intercluster edges.*

**Corollary 10.10.** *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  $\gamma$ .

**Corollary 10.11.** *The time and message complexities of synchronizer  $\gamma$  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 3, 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  $\gamma$  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  $\alpha$ ,  $\beta$ , and  $\gamma$  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  $\gamma$  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.

## 10.6 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  $\beta$  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 10.12.** *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 10.12, 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  $\alpha$  pays most attention to the local synchronization, we take a look at a protocol inspired by the synchronizer  $\alpha$ . A pseudo-code representation for the clock synchronization protocol  $\alpha$  is given in Algorithm 44.

---

**Algorithm 44** Clock synchronization  $\alpha$  (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 10.13.** *The clock synchronization protocol  $\alpha$  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  $\alpha$  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 Notes

The idea behind synchronizers is quite intuitive and as such, synchronizers  $\alpha$  and  $\beta$  were implicitly used in various asynchronous algorithms [Gal76, Cha79, CL85] before being proposed as separate entities. The general idea of applying synchronizers to run synchronous algorithms in asynchronous networks was first introduced by Awerbuch [Awe85a]. His work also formally introduced the synchronizers  $\alpha$  and  $\beta$ , whereas other constructions were presented in [AP90, PU87].

Naturally, as synchronizers are motivated by practical difficulties with local clocks, there are plenty of real life applications. Studies regarding applications can be found in, e.g., [SM86, Awe85b, LTC89, AP90, PU87]. Synchronizers in the presence of network failures have been discussed in [AP88, HS94].

It has been known for a long time that the global clock skew is  $\Theta(D)$  [LL84, ST87]. The problem of synchronizing the clocks of nearby nodes was introduced by Fan and Lynch in [LF04]; they proved a surprising lower bound of  $\Omega(\log D / \log \log D)$  for the local skew. The first algorithm providing a non-trivial local skew of  $\mathcal{O}(\sqrt{D})$  was given in [LW06]. Later, matching upper and lower bounds of  $\Theta(\log D)$  were given in [LLW10]. The problem has also been studied in a dynamic setting [KLO09, KLLO10].

Clock synchronization is a well-studied problem in practice, for instance regarding the global clock skew in sensor networks, e.g. [EGE02, GKS03, MKSL04, PSJ04]. One more recent line of work is focussing on the problem of minimizing the local clock skew [BvRW07, SW09, LSW09, FW10, FZTS11].

## Bibliography

- [AP88] Baruch Awerbuch and David Peleg. Adapting to Asynchronous Dynamic Networks with Polylogarithmic Overhead. In *24th ACM Symposium on Foundations of Computer Science (FOCS)*, pages 206–220, 1988.

- [AP90] Baruch Awerbuch and David Peleg. Network Synchronization with Polylogarithmic Overhead. In *Proceedings of the 31st IEEE Symposium on Foundations of Computer Science (FOCS)*, 1990.
- [Awe85a] Baruch Awerbuch. Complexity of Network Synchronization. *Journal of the ACM (JACM)*, 32(4):804–823, October 1985.
- [Awe85b] Baruch Awerbuch. Reducing Complexities of the Distributed Max-flow and Breadth-first-search Algorithms by Means of Network Synchronization. *Networks*, 15:425–437, 1985.
- [BvRW07] Nicolas Burri, Pascal von Rickenbach, and Roger Wattenhofer. Dozer: Ultra-Low Power Data Gathering in Sensor Networks. In *International Conference on Information Processing in Sensor Networks (IPSN)*, Cambridge, Massachusetts, USA, April 2007.
- [Cha79] E.J.H. Chang. *Decentralized Algorithms in Distributed Systems*. PhD thesis, University of Toronto, 1979.
- [CL85] K. Mani Chandy and Leslie Lamport. Distributed Snapshots: Determining Global States of Distributed Systems. *ACM Transactions on Computer Systems*, 1:63–75, 1985.
- [EGE02] Jeremy Elson, Lewis Girod, and Deborah Estrin. Fine-grained Network Time Synchronization Using Reference Broadcasts. *ACM SIGOPS Operating Systems Review*, 36:147–163, 2002.
- [FW10] Roland Flury and Roger Wattenhofer. Slotted Programming for Sensor Networks. In *International Conference on Information Processing in Sensor Networks (IPSN)*, Stockholm, Sweden, April 2010.
- [FZTS11] Federico Ferrari, Marco Zimmerling, Lothar Thiele, and Olga Saukh. Efficient Network Flooding and Time Synchronization with Glossy. In *Proceedings of the 10th International Conference on Information Processing in Sensor Networks (IPSN)*, pages 73–84, 2011.
- [Gal76] Robert Gallager. Distributed Minimum Hop Algorithms. Technical report, Lab. for Information and Decision Systems, 1976.
- [GKS03] Saurabh Ganeriwal, Ram Kumar, and Mani B. Srivastava. Timing-sync Protocol for Sensor Networks. In *Proceedings of the 1st international conference on Embedded Networked Sensor Systems (SenSys)*, 2003.
- [HS94] M. Harrington and A. K. Somani. Synchronizing Hypercube Networks in the Presence of Faults. *IEEE Transactions on Computers*, 43(10):1175–1183, 1994.
- [KLLO10] Fabian Kuhn, Christoph Lenzen, Thomas Locher, and Rotem Oshman. Optimal Gradient Clock Synchronization in Dynamic Networks. In *29th Symposium on Principles of Distributed Computing (PODC)*, Zurich, Switzerland, July 2010.

- [KLO09] Fabian Kuhn, Thomas Locher, and Rotem Oshman. Gradient Clock Synchronization in Dynamic Networks. In *21st ACM Symposium on Parallelism in Algorithms and Architectures (SPAA)*, Calgary, Canada, August 2009.
- [LF04] Nancy Lynch and Rui Fan. Gradient Clock Synchronization. In *Proceedings of the 23rd Annual ACM Symposium on Principles of Distributed Computing (PODC)*, 2004.
- [LL84] Jennifer Lundelius and Nancy Lynch. An Upper and Lower Bound for Clock Synchronization. *Information and Control*, 62:190–204, 1984.
- [LLW10] Christoph Lenzen, Thomas Locher, and Roger Wattenhofer. Tight Bounds for Clock Synchronization. In *Journal of the ACM, Volume 57, Number 2*, January 2010.
- [LSW09] Christoph Lenzen, Philipp Sommer, and Roger Wattenhofer. Optimal Clock Synchronization in Networks. In *7th ACM Conference on Embedded Networked Sensor Systems (SenSys)*, Berkeley, California, USA, November 2009.
- [LTC89] K. B. Lakshmanan, K. Thulasiraman, and M. A. Comeau. An Efficient Distributed Protocol for Finding Shortest Paths in Networks with Negative Weights. *IEEE Trans. Softw. Eng.*, 15:639–644, 1989.
- [LW06] Thomas Locher and Roger Wattenhofer. Oblivious Gradient Clock Synchronization. In *20th International Symposium on Distributed Computing (DISC)*, Stockholm, Sweden, September 2006.
- [MKSL04] Miklós Maróti, Branislav Kusy, Gyula Simon, and Ákos Lédeczi. The Flooding Time Synchronization Protocol. In *Proceedings of the 2nd international Conference on Embedded Networked Sensor Systems, SenSys '04*, 2004.
- [PSJ04] Santashil PalChaudhuri, Amit Kumar Saha, and David B. Johnson. Adaptive Clock Synchronization in Sensor Networks. In *Proceedings of the 3rd International Symposium on Information Processing in Sensor Networks, IPSN '04*, 2004.
- [PU87] David Peleg and Jeffrey D. Ullman. An Optimal Synchronizer for the Hypercube. In *Proceedings of the sixth annual ACM Symposium on Principles of Distributed Computing, PODC '87*, pages 77–85, 1987.
- [SM86] Baruch Shieber and Shlomo Moran. Slowing Sequential Algorithms for Obtaining Fast Distributed and Parallel Algorithms: Maximum Matchings. In *Proceedings of the fifth annual ACM Symposium on Principles of Distributed Computing, PODC '86*, pages 282–292, 1986.
- [ST87] T. K. Srikanth and S. Toueg. Optimal Clock Synchronization. *Journal of the ACM*, 34:626–645, 1987.

- [SW09] Philipp Sommer and Roger Wattenhofer. Gradient Clock Synchronization in Wireless Sensor Networks. In *8th ACM/IEEE International Conference on Information Processing in Sensor Networks (IPSN)*, San Francisco, USA, April 2009.