

# 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 2) 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 10.5** 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 10.5 formally describes the synchronizer  $\alpha$ .

**Theorem 10.6.** *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 10.7** 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 10.7.

**Theorem 10.8.** *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 the improved variant of the GHS algorithm (Algorithm 2.18) mentioned in Chapter 2, 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 $\gamma$

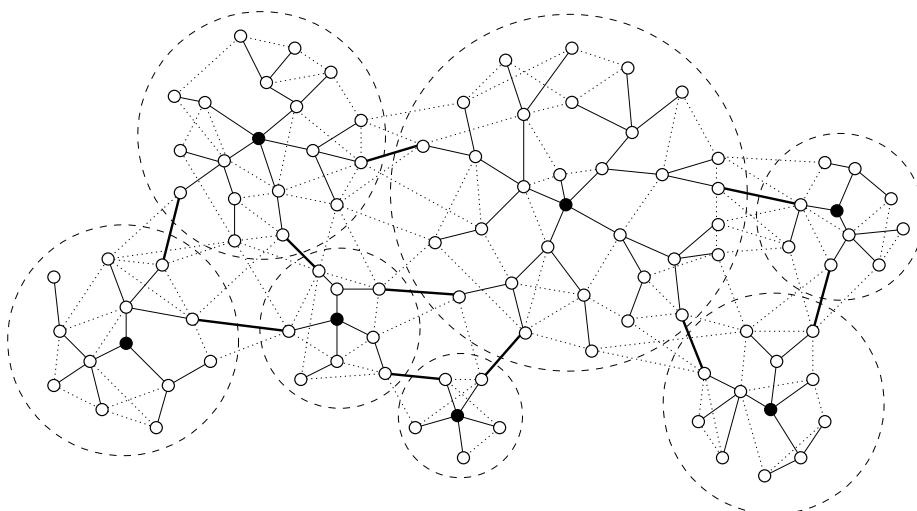


Figure 10.9: A cluster partition of a network: The dashed cycles specify the clusters, cluster leaders are black, the solid edges are the edges of the intracluster trees, and the bold solid edges are the intercluster edges

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 10.9 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 10.10.

---

**Algorithm 10.10** 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 10.11.** *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 10.12 describes how to construct a partition into clusters that can be used for synchronizer  $\gamma$ . In Algorithm 10.12,  $B(v, r)$  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 10.12** 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 10.13.** *Algorithm 10.12 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 10.14.** *Using  $\rho = 2$ , Algorithm 10.12 computes a clustering with cluster radius at most  $\log_2 n$  and with at most  $n$  intercluster edges.*

**Corollary 10.15.** *Using  $\rho = n^{1/k}$ , Algorithm 10.12 computes a clustering with cluster radius at most  $k$  and at most  $\mathcal{O}(n^{1+1/k})$  intercluster edges.*

**Remarks:**

- Algorithm 10.12 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 10.12 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 10.16.** *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:**

- In Chapter 2, 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. In such scenarios, it is possible to achieve 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.