

10 **3.1 Clock Synchronisation**

11 This section first presents some commonly-used definitions and notations
12 in the field of the clock synchronisation problem, and then gives an
13 overview of the known solutions that fits to the GUARDS constraints. At
14 last, the GUARDS synchronisation algorithm is detailed, including the
15 initial synchronisation.

16 **3.1.1 Definitions and Notations**

17 We consider a set of fully connected nodes. Each node has a physical
18 clock, and computes a logical clock time.

19 *Definition 1: Clock time*

20 Each node maintains a logical clock time $T = C(t)$, meaning that at real-
21 time t the clock time of the local node is T .

22 By convention, variables associated with clock time (local to a given
23 node) are in uppercase, whereas variables associated with real time
24 (measured in an assumed global Newtonian frame) are in lowercase.

25 *Definition 2 : drift rate*

1 A nonfaulty physical clock C has a maximum drift rate ρ if and only if for
2 any $t_2 > t_1$:

$$3 \quad 1 - \rho < \frac{C(t_2) - C(t_1)}{t_2 - t_1} < 1 + \rho$$

4 **Definition 3 : agreement**

5 The agreement condition is satisfied if and only if the skew δ between any
6 nonfaulty logical clocks C_i and C_j is bounded:

$$7 \quad |C_i(t) - C_j(t)| \leq \delta$$

8 **Definition 4 : accuracy**

9 The accuracy condition is satisfied if and only if any nonfaulty logical clock
10 C_i stay within a linear envelop of real time. That is, there exists a constant
11 $\gamma > 0$ and two constants a and b (that depend on initial conditions) such that:

$$12 \quad (1 - \gamma)t + a \leq C_i(t) \leq (1 + \gamma)t + b$$

13 **Definition 5 : synchronisation algorithm**

14 A synchronisation algorithm is an algorithm which satisfies the agreement
15 and accuracy conditions.

16

17 Some typical values of the parameters introduced so far or used in the
18 sequel of the section are given in Table 1.

19

| Variable | Content | Typical value |
|----------|---|-----------------------|
| ρ | Maximum drift rate of all physical clocks | 10^{-4} |
| f | Quartz frequency of physical clocks | 10 MHz |
| δ | Maximum skew between any two nonfaulty logical clocks | less than 100 μ s |

| Variable | Content | Typical value |
|------------|---|---------------|
| d | Upper bound on the transmission delay | 50 μ s |
| ϵ | Upper bound on the read error ² | 1 μ s |
| n | Number of clocks | at most 4 |
| m | Maximum number of (arbitrarily) faulty clocks | 1 |
| R | Resynchronisation interval | 500 ms |

1 **Table 1 : Typical values**

2 Note that the maximum skew between any two logical clocks is necessarily
3 greater than the skew between their physical clocks, and that this latter
4 skew can easily be obtained through the simple formula $\delta_{\text{physical}} = 2\rho R$,
5 where R is the interval between two resynchronisations. With the values
6 given in \$\$\$Table 1, this gives $\delta_{\text{physical}} = 2*10^{-4}*500 \text{ ms} = 100 \mu\text{s}$. Hence
7 the fact that the typical value of δ in \$\$\$Table 1 is greater than 100 μ s (see
8 section \$\$\$3.1.1.3 for a more precise estimation of δ).

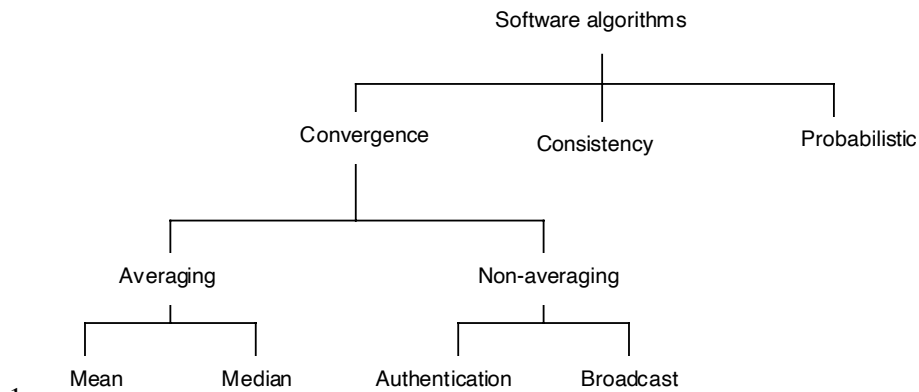
9 **3.1.2 Existing Algorithms**

10 A lot of clock synchronisation algorithms are described in the literature.
11 [Ramanathan et al. 1990] provides a good survey of these algorithms.
12 Restricting only to software implemented algorithms³ yields three main
13 kinds of algorithms (see \$\$\$Figure 1):

² The read error is mainly due to the transmission delay incertitude. Thus, $d-2\epsilon$ is the minimum message transit delay between any two nodes in the system.

³ In order to minimize the amount of specific hardware, in conformance to the GUARDS requirement to use as much COTS components as possible.

- 1 • Convergence algorithms, where nodes periodically resynchronise by
2 exchanging synchronisation messages.
 - 3 • Consistency algorithms, where nodes use an interactive consistency
4 algorithm to achieve agreement on clock values. The little amount of
5 literature on this kind of algorithm [Lamport & Melliar-Smith 1985]
6 does not allow us to use them within GUARDS without spending
7 significant efforts on carefully studying them.
 - 8 • Probabilistic algorithms [Cristian 1989], where nodes can make the
9 worst-case skew as small as desired, but with an increasing probability
10 of loss of synchronisation. This is clearly unacceptable in the field of
11 safety-critical applications.
- 12 Focusing on convergence algorithms, we see that there exists two main
13 kinds of such algorithms:
- 14 • Convergence-averaging algorithms, where each node resynchronises
15 according to clock values obtained through periodic one-round clock
16 exchanges. On each node, the other clocks can be taken into account
17 through a mean-like function [Lamport & Melliar-Smith 1985], or a
18 median-like function [Lundelius-Welch & Lynch 1988].
 - 19 • Convergence-nonaveraging algorithms, where each node periodically
20 seeks to be the system synchroniser. To deal with possible byzantine
21 behaviors, the exchanged messages can be authenticated or broadcast
22 [Srikanth & Toueg 1987; Dolev et al. 1995].



1

2

Figure 1 : Synchronisation algorithms

3.1.2.1 Convergence-averaging Algorithms

4 These synchronisation algorithms can tolerate m arbitrary faulty nodes in a
 5 (fully connected) network of n nodes, under the sufficient condition⁴ that
 6 $n \geq 3m+1$. This implies in particular that they *cannot* tolerate a byzantine
 7 faulty clock in a three node system (see \$\$\$Figure 2).

8 The basic idea of these algorithms is as follows. The
 9 resynchronisation is performed periodically: each node broadcasts a
 10 resynchronisation message when its local clock time has counted R
 11 seconds since the last resynchronisation period. At the same time, the node
 12 collects during a given waiting period (which is far shorter than R) the
 13 resynchronisation messages broadcast by other nodes: it records the arrival
 14 time (according to its local clock time) of each received resynchronisation
 15 message. After the waiting period, it uses a fault-tolerant averaging
 16 function on the arrival times to compute an averaged clock time. It then

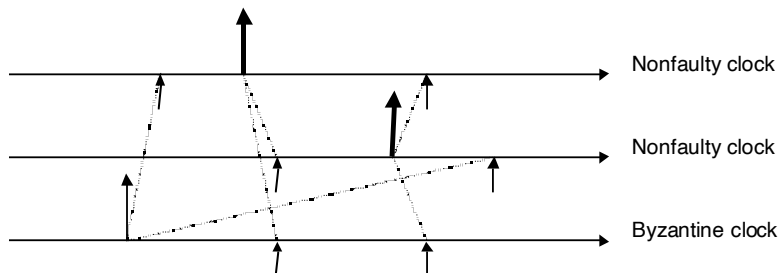
⁴ This condition is also necessary if authentication is not used.

1 corrects its own clock time by setting it to the averaged clock time before
2 the next resynchronisation period.

3 In [Lamport & Melliar-Smith 1985], the fault-tolerant averaging
4 function F is the mean of all values (with values differing from the local
5 value by more than δ replaced by the local value).

6 In [Lundelius-Welch & Lynch 1988], the function F is a median-like
7 function. It consists in rejecting the m highest and m lowest values, and
8 then performing the average of all remaining values.

9 With three clocks, the above fault synchronisation function F
10 consists simply in taking the median value. So a byzantine faulty clock can
11 send synchronisation messages to two other nodes in such a way that they
12 find themselves as being always the median: the byzantine faulty clocks
13 has only to send synchronisation messages so that the earliest (respectively
14 latest) nonfaulty clock receive the Byzantine synchronisation message
15 before (respectively after) its own synchronisation emission. Thus each
16 nonfaulty node keeps its own value and the two nonfaulty clocks drifts
17 apart progressively.



18

1 When a node broadcasts its synchronisation message, it is assumed that it
2 receives it immediately (because for the local node, there is no transmission
3 delay for the reception of its own synchronisation message). So, the test of
4 reception of $(m+1)$ signed messages includes the possible locally emitted
5 synchronisation message.

6 \$\$\$Figure 3 gives an example of tolerance of a byzantine faulty clock
7 in a three node system (m is equal to 1). The faulty clock (clock number 3
8 on the figure) broadcasts its (signed) resynchronisation message⁵ so that
9 the fastest nonfaulty clock (clock number 1 on the figure) receives this
10 message before the end of its own resynchronisation interval and the
11 slowest nonfaulty clock (clock number 2 on the figure) receives it after the
12 end of its own resynchronisation interval. We have seen that such an
13 inconsistent broadcast can desynchronise the two nonfaulty clocks under a
14 convergence-averaging synchronisation algorithm. Let us analyze what
15 happens here:

- 16 • The fastest nonfaulty clock receives the resynchronisation message of
17 clock 3. A short time after, this non-faulty clock reaches the end of its
18 own resynchronisation interval. According to the algorithm, it has
19 received at this time $m+1$ resynchronisation messages (the one coming
20 from clock 3 and its own). So it starts its new local clock (horizontal
21 bold arrow on the figure) and relays to the two other nodes two (signed)
22 resynchronisation messages.
- 23 • The slowest nonfaulty clock receives the two signed resynchronisation
24 messages (just emitted by the clock 1), and thus according to the
25 algorithm restarts its new local clock (horizontal bold arrow on the
26 figure) and relays to the two other nodes these two (signed)
27 resynchronisation messages. A short time after, this clock 2 reaches the
28 end of its own resynchronisation interval and then broadcast its own
29 resynchronisation message. A short time after, this clock 2 receives the

5 The signature of a message m by a node i is denoted by $m : i$ on the figure.

1 resynchronisation message from clock 3 which it may identify as being
2 faulty.

3 We can see that the clock 2 starts its new local clock at most d seconds
4 (maximum transmission time) after clock 1 has started its own.

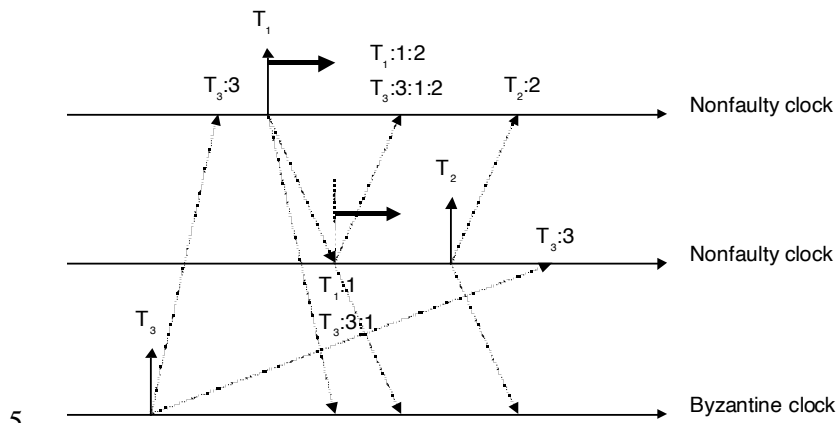


Figure 3 : Tolerance of a byzantine clock in a 3 node system

7 3.1.2.3 Maximum Skew Estimation

8 One of the most important characteristics of a synchronisation algorithm is
9 its ability to achieve a small maximum skew between non-faulty clocks.
10 Therefore, it is important that we estimate and compare the maximum
11 skew of [Lundelius-Welch & Lynch 1988] (hereafter denoted by LL) and
12 [Srikanth & Toueg 1987] (hereafter denoted by ST).

13 Theoretical results provide us with the two following formulas:

- 1 • for the LL convergence-averaging algorithm (theorem 16, page 21)⁶:
 2 $\delta = \beta + \varepsilon + \rho(7\beta+3d+4\varepsilon) + 4\rho^2(2+\rho)(\beta+d)$, with $\beta = 4(\varepsilon+\rho R)$
- 3 • for the ST convergence-nonaveraging algorithm (lemmas 6 and 8, page
 4 630-631): $\delta = [R(1+\rho) + d][\rho(2+\rho)/(1+\rho)] + d(1+\rho)$
- 5 Using the typical values given in \$\$\$Table 1, we have $\varepsilon = 10^{-6}$ s, $\rho R =$
 6 $5 \cdot 10^{-5}$ s, $d = 5 \cdot 10^{-5}$ s, and we can see that terms including factors like ρ^2 ,
 7 $\rho\varepsilon$, ρd are negligible. By dropping these higher order terms, we obtain:
- 8 • for the LL algorithm: $\delta \sim 5\varepsilon + 4\rho R$
- 9 • for the ST algorithm: $\delta \sim d + 2\rho R$
- 10 We can see that the maximum skew δ is of the same order of magnitude in
 11 the two algorithms (a few hundreds of microseconds). However, several
 12 remarks must be made:
- 13 • We have taken a rather pessimistic maximum drift rate of 10^{-4} for
 14 physical clocks. For example, ρ is estimated at 10^{-6} by [Ramanathan et
 15 al. 1990], leading to $\rho R = 5 \cdot 10^{-7}$ s. If we take such a smaller value for ρ ,
 16 then LL becomes better than ST.
- 17 • We have taken a rather optimistic maximum transmission delay d ,
 18 taking into account that GUARDS channels are expected to be near each
 19 other (typically a few decimeters). If channels are distributed over
 20 longer distances, then this could lead to an increase of d , thus making
 21 again LL better than ST.

6 The factorisation of this fomula follows actually those given in [Ramanathan et al. 1990]. But note that the exact reformulation of [Ramanathan et al. 1990] is incorrect (omission of the scalar 3 in front of d in the third factor).

- 1 • From a practical viewpoint, it is not clear whether the transmission
2 delay d is the same for LL and for ST. Indeed, LL is a single-round
3 algorithm where synchronisation messages are deterministically and
4 synchronously broadcast. Thus d corresponds to the maximum time
5 required for a message to be prepared by a node, broadcast to all other
6 nodes and processed by these other nodes. For ST the issue is that a
7 node may have to broadcast its synchronisation message (first round),
8 and then to relay immediately after the synchronisation messages it
9 received (second round), precisely because it received the $(m+1)^{\text{th}}$
10 signed synchronisation message just after having sent its own
11 synchronisation message. The maximum transmission delay has to take
12 into account this unfavorable case (which is likely to happen precisely
13 when clocks are well synchronised), thus leading to a larger value.