

Lecture 13: March 27

*Lecturer: Prashant Shenoy**Scribe: Ritvika Pillai, Mohini Jain*

13.1 Overview

The topic of the lecture is “Time ordering and clock synchronization.” This lecture covered the following topics:

Clock Synchronization: Cristian’s algorithm, Berkeley algorithm, NTP, GPS

Logical Clocks: Event ordering

13.2 Clock Synchronization

13.2.1 The motivation of clock synchronization

Just like usual clocks, systems have a clock that tells time to applications running on the system. Centralized machines have just 1 clock, but in the case of distributed systems each machine has its own clock. All these clocks might not be in sync and may drift over time. Hence the time will be dependent on which local clock is being checked.

For example, you modify files and save them on one machine A, and use another machine B to compile the files modified. If machine B has a faster clock than machine A, you may not correctly compile the files modified because the time of compiling files on machine B may be later than the time of editing files on machine A according to local timestamp on different machines, thus leading to errors.

13.2.2 How physical clocks and time work

1) One approach is to use astronomical metrics (solar day) to tell time. For example, solar noon is the time that sun is directly overhead. “Noon” on our clocks is different from solar noon. Noon depends on time zone while solar noon is a fixed physical fact. We typically use the notion of solar day to tell there are 24 hours between the time that sun is directly overhead on a particular location. Although this method was used for centuries, it is not accurate since it based on the length of a day.

2) Atomic clocks use the properties of atoms to measure time. Atomic clocks are the most accurate clocks and other clocks derive from this time. Typically, you will have some centralized atomic clock broadcast its time. The receivers, which may use less accurate mechanisms, are then synchronized with the atomic clock. For example, cell-phone clock also uses atomic clock to synchronize time with cell-phone broadcast tower. Some satellites also broadcast time based on atomic clocks.

3) Coordinated universal time (UTC) is based on noon in Greenwich (UK). All time zones are offset by UTC. Most of the atomic clocks broadcast UTC time regardless of timezone using wireless channels, satellites, FM radios, etc. Receivers listen to this and set local time based on the broadcast time.

4) Mechanical clocks are less accurate—the accuracy is roughly one part per million. Computers typically use mechanical clocks. This small amount of inaccuracy results in clock drift because the property of the physical mechanism (usually quartz) can change with environmental properties such as temperature or humidity. To avoid clock drift, we need to synchronize machines with a master or with one another.

13.2.3 Drift tolerance and frequency of synchronization

Actual clocks have drift and hence need synchronization once in a while. Just like the drift developed over time in a mechanical clock due to the quartz crystal inside it, computer clocks also develop drift over time since they also use a chip which depends on a physical material to tell time. This drift needs to be calculated to determine how frequently synchronization is needed.

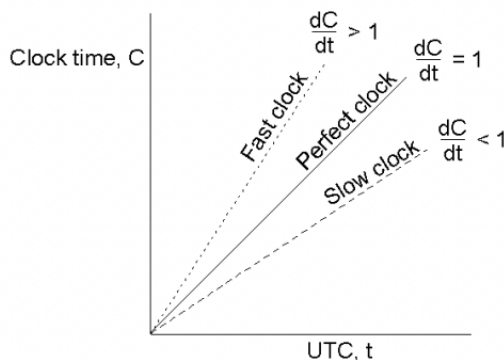


Figure 13.1: Clock drift relative to a perfect clock.

In Figure 12.1, Here, t is UTC time, C is clock time, and slope, dC/dt , is the rate of advancement of the clock. ρ indicates the inaccuracy of the clock. Consider the following cases: a. If the $dC/dt = 1$, the real time and clock advances proportionally and are in sync. b. If the $dC/dt < 1$, real time advances by 1 second and the clock will advance by $(1 - \rho)$ second, i.e., the clock runs slower. c. If the $dC/dt > 1$, real time advances by 1 second then clock will advance by $(1 + \rho)$ second, i.e., the clock runs faster. To limit the error in clock to δ , we need to synchronize every $\delta/2\rho$ seconds.

Question: Do we assume the drifts of two machines in a distributed system to be the same? If so, why do we use 2ρ ?

Answer: Yes, we assume that both the machines in the system are using a similar chip that produces a drift of ρ . Since we want to limit the drift of the entire distributed system and not just one machine in it, relative total drift between machines is considered for synchronization. So, if one machine is slower by $-\rho$ and the other is ahead by $+\rho$, the total difference becomes 2ρ (relative to one another).

13.3 Centralized clock synchronization algorithms

13.3.1 Cristian's Algorithm

In Cristian's Algorithm, is a master machine called *time server* which is the authoritative clock for telling time. It is in sync with the atomic clock via a UTC receiver. Other machines in the system synchronize

with the time server.

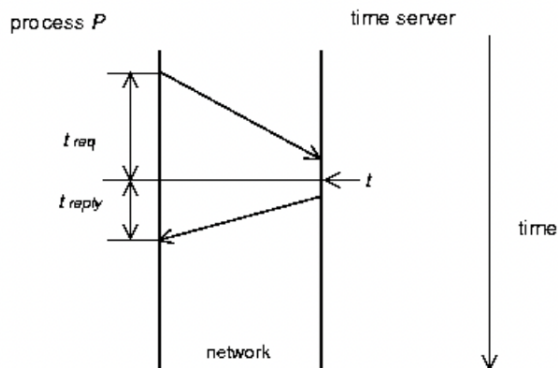


Figure 13.2: Cristian's Algorithm.

Machine P sends message to the time server to check the current time. After taking some time (t_{req}) to propagate, the request reaches the time server and will then be processed. The time server then returns the current time (t) and machine P uses this time to reset its clock. The machine P will set its time as $(t + t_{reply})$ and not just t . This is done so as to take the propagation delay from server to machine P into account. We can use $(t_{req} + t_{reply})/2$ as an estimation of t_{reply} . The better the estimate, the better the synchronization.

13.3.2 Berkeley Algorithm

This algorithm doesn't use a time server. Instead, clocks are synchronized with one another in a group, and no machine in this group synchronize with external atomic clock. We use leader election to select a "master" in a group to run clock synchronization while others are "slaves." This master clock is known as the coordinator. Each machine sends their local time to the coordinator. The coordinator then calculates an average of these times. Based on the value of average, the time of all clocks are adjusted. For example, three machines reply with their clock values as time difference of 0, -10, +25 at 3:00, then the master will tell all those machines to set their clock at 3 : 00 + 5(5 = (0 - 10 + 25)/3). This is a relative clock synchronization algorithm, not an absolute synchronization algorithm. The propagation times are estimated in the same way as Cristian's algorithm.

Question: Does the drift rate change over time?

Answer: It may change over time, but the manufacturer of the clocks guarantees that the drift will be between 0 and ρ (both included).

Question: How often does the time daemon send synchronization requests?

Answer: Synchronization here works the same way as in Cristian's Algorithm, i.e., if we want the clocks to not be off by more than δ , the clocks will need to be synchronized every $\delta/2\rho$ seconds. The only difference is synchronization is relative unlike Cristian's Algorithm where a time server is used to tell the correct time.

13.4 Distributed clock synchronization approaches

Both Cristian's and Berkeley are centralized algorithms. Apart from these, there are also decentralized algorithms using resynchronized intervals. In a decentralized version of Berkeley, the role of coordinator is

eliminated. Instead, all machines broadcast their times to all other machines at the start of the interval. At every machine, suppose n clock values are received within the interval. Then at the end of period S , their average is calculated which is then used to set their local time. The latency caused by the network while communicating is also adjusted during this time at the receiver end. For the outliers, machines can throw away few highest and lowest values to avoid negative influence of extremely fast or slow clocks relative to the average time.

There are two decentralized approaches in use today. One approach is using NTP which is used by most computers. It uses a time server and advanced techniques to deal with network propagation delays. The accuracy is typically between 1 and 50ms. The other approach is `rdate`, which synchronizes a machine with a specified machine. In many cases, you can run `rdate` with the argument of the name of server and just synchronize clock with that server.

13.4.1 Network Time Protocol (NTP)

NTP is widely used standard which based on Cristian's algorithm. In NTP clock synchronization, you also want to find out network propagation delay (dT_{res}). NTP clock synchronization uses a hierarchical protocol and unlike Cristian's algorithm, it does not let the clock be set backward. Since the fast clock cannot go backward, it is synchronized by slowing it down. Letting a clock go backward can have many negative consequences (such as two files having the same timestamp). This is the reason why NTP is widely used compared to Cristian's algorithm.

13.4.2 Global Positioning System (GPS)

GPS is a technology that allows any device to figure out its location. It requires clock synchronization to accurately figure out where the device is located. For example, a phone has a GPS chip that listens to satellite broadcasts. Uses the principle of triangulation to know where you are with respect to the known position of the satellite. These known positions are called landmarks. GPS achieves high accuracy because it is synchronized with satellites which use atomic clocks without a hierarchical protocol. It is assumed that the satellites are in perfect synchronization with an atomic clock.

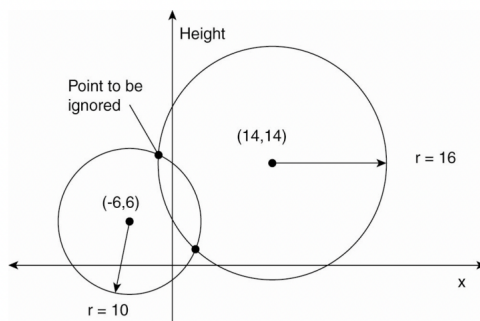


Figure 13.3: Global Positioning System (GPS)

2D space: Let the 2 landmarks be $(14,14)$ and $(-6,6)$. A device somewhere in this space will measure its distance with respect to these landmarks. Say, it is 16 units from the first landmark and 10 units from the second, this means that the device is on the intersection of the 2 circles because it has to satisfy both the distance constraints. If a 3rd landmark is added, then the exact position of the device can be known.

3D space: We assume GPS landmark A with its position (x_1, y_1, z_1) and its timestamp t_1 , and GPS receiver B (e.g. a car) with its unknown position (x, y, z) and the timestamp t receiving broadcast t_1 message from a GPS landmark. Then the distance between A and B is $di = \sqrt{(x - x_1)^2 + (y - y_1)^2 + (z - z_1)^2}$, and di also equals $c(t_2 - t_1)$, where c is the speed of light. If we assume the receiver has a drift time dr from landmark A, then we can use Equation (1) to show that $c(t_2 + dr - t_1) = \sqrt{(x - x_1)^2 + (y - y_1)^2 + (z - z_1)^2}$. From Equation (1), we can see that there are 4 unknowns, x, y, z and dr , thus we need minimum 4 satellites to compute the location of a GPS receiver as well as its time value. If we get 4 satellites, then we can get multiple solutions of the location of the receiver. If we have 6 or 8 satellites, we can quickly narrow the solutions. Therefore, GPS does clock synchronization as well as computing the receiver's location, the satellites just keep broadcasting the time.

13.5 Logical clock

The above approaches use timestamps to reason the order of events. If the time difference between two events is smaller than the accuracy, then we cannot say which event happens first, thus problems may be caused. In some cases, if processes need to know the order in which the events occurred instead of the exact time, then logical clocks should be used. Hence absolute time isn't important and clock synchronization isn't needed.

13.5.1 Event Ordering

In logical clocks, there is no global clock and local clocks may run faster or slower. The ordering of events needs to be figured out in such a situation. There are some key ideas of logical clocks proposed by the scientist Lamport: we can use send/receive messages exchanged between processes/machines to order events, and if 2 processes never communicate with each other, then in such cases we don't need to find order.

The happened-before relation: We use the fundamental property that the send event occurs before the receive event. The relation is transitive, i.e, if A occurs before B and B occurs before C, then A occurs before C.

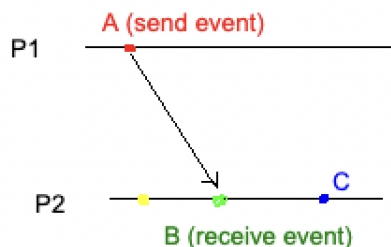


Figure 13.4: The “happened-before” relation. Note that we cannot say anything about the relation with the yellow event.

Each processor has a logical clock which gets incremented whenever an event occurs. Suppose when process i sends a message to process j , it piggybacks its local timestamp (say $LC_i=3$) along with the message. The receiver takes this timestamp and its local timestamp (say $LC_j=4$). Then the maximum of both these values is calculated and incremented by 1, i.e, $\max(LC_i, LC_j) + 1$. This makes sure that the timestamp assigned to the receiver event is higher than the sending event. This technique was invented by Leslie Lamport.

The above algorithm only solves half the problem as it gives only forward property and not the reverse property (i.e. if $timestampA < timestampB$ then A has occurred before B). Hence further changes are needed in this approach.

Question: Do data centers or large servers use logical clocks?

Answer: Although logical clocks are a good idea and they are being used in some techniques, in real world applications, as of now only real clocks (atomic clock synchronization) are being used.