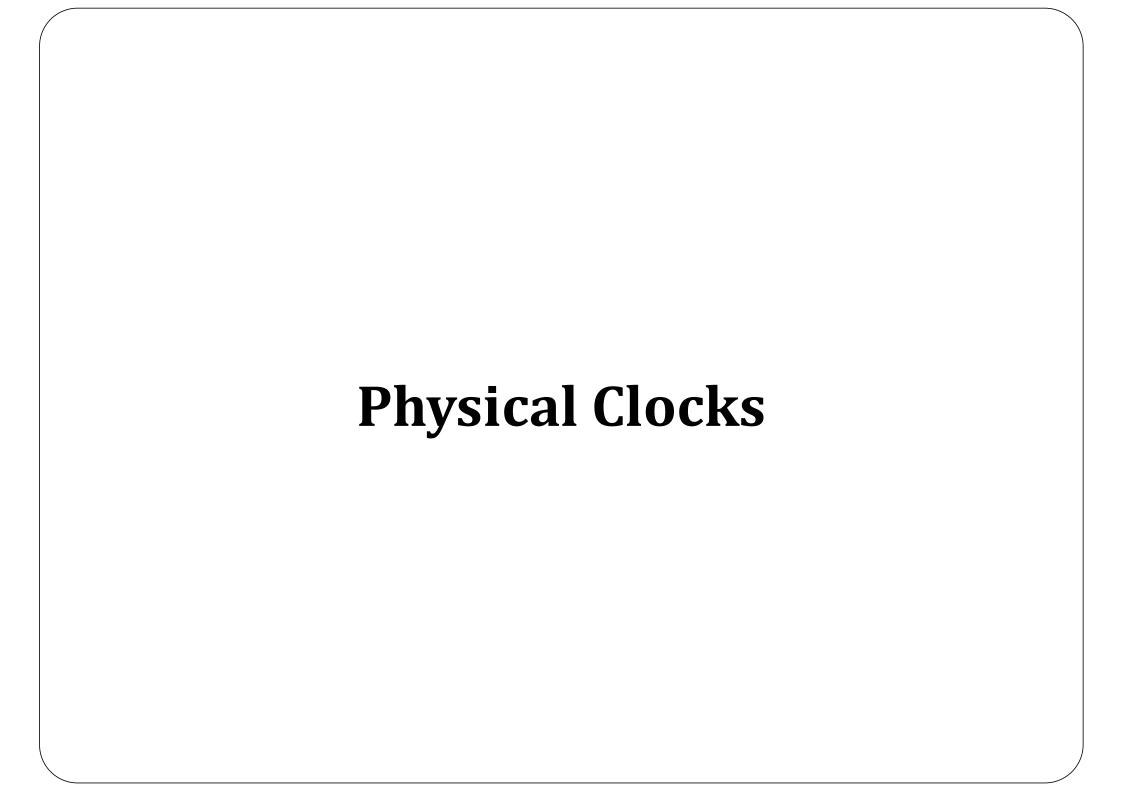# Clocks in Distributed Systems

- Needed to
  - Order two or more events happening at same or different nodes (Ex: Consistent ordering of updates at different replicas, ordering of multicast messages sent in a group)
  - Decide if two events happened between some fixed duration of each other (Ex: Replay of stolen messages in distributed authentication protocols like Kerberos)
  - Start events at different nodes together at the same time (Ex: tracking in sensor networks, sleep/wakeup scheduling)

- Easy if a globally synchronized clock is available, but
  - Perfectly synchronized clocks are impossible to achieve
  - But perfect synchronization may not be needed always; synchronization within bounds may be enough
    - Degree of synchronization needed depends on application
      - Kerberos requires synchronization of the order of minutes
      - Tracking applications may require synchronization of the order of seconds
  - Still not sufficient for ordering events always
    - Suppose each node timestamps events at the node by its local clock
    - Suppose synchronization is accurate within bound $\partial$
    - May not be able to order events whose timestamps differ by less than $\partial$

- Two approaches for building clocks
  - **Physical Clocks**
    - Each machine has its own local clock
    - Clock synchronization algorithms run periodically to keep them synchronized with each other within some bounds
    - Useful for giving a consistent view of "current time" across all nodes within some bounds, but cannot order events always
  - **Logical Clocks**
    - Use the notion of **causality** to order events
    - Can what happened in one event affect what happens in another?
      - Because if not, ordering them is not important
    - Useful for ordering events, but not for giving a consistent view of "current time" across all nodes

# Physical Clocks

# Physical Clocks

- Each node has a local clock used by it to timestamp events at the node

- Local clocks of different nodes may vary

- Need to keep them synchronized (Clock Synchronization Problem)

- Perfect synchronization not possible because of inability to estimate network delays exactly

# Clock Synchronization

- Internal Synchronization
  - Requires the clocks of the nodes to be synchronized to within a pre-specified bound
  - However, the clock times may not be synchronized to any external time reference, and can vary arbitrarily from any such reference
- External Synchronization
  - Requires the clocks to be synchronized to within a pre-specified bound of an external reference clock

# How Computer Clocks Work

- Computer clocks are based on crystals that oscillate at a certain frequency
- Every H oscillations, the timer chip interrupts once (clock tick)
  - Resolution: time between two interrupts
- The interrupt handler increments a counter that keeps track of no. of ticks from a reference in the past (epoch)
- Knowing no. of ticks per second, we can calculate year, month, day, time of day etc.

# Why Clocks Differ: Clock Drift

- Period of crystal oscillation varies slightly due to temperature, humidity, ageing,…
- If it oscillates faster, more ticks per real second, so clock runs faster; similar for slower clocks
- For machine p, when correct reference time is t, let machine clock show time as $C = C_p(t)$
- Ideally, $C_p(t) = t$ for all p, t
- In practice, $1 - \rho \leq dC/dt \leq 1 + \rho$
- $\rho$ = max. clock drift rate, usually around $10^{-5}$ for cheap oscillators
- Drift results in skew between clocks (difference in clock values of two machines)

# Resynchronization

- Periodic resynchronization needed to offset skew
- If two clocks are drifting in opposite directions, max. skew after time t is $2\rho t$
- If application requires that clock skew $< \delta$, then resynchronization period

$$r < \delta / (2 \rho)$$

- Usually $\rho$ and $\delta$ are known
  - $\rho$ given by crystal manufacturer
  - $\delta$ specified from application requirement

# Cristian's Algorithm

- One node acts as the time server

- All other nodes sends a message periodically (within resync. period r) to the time server asking for current time

- Time server replies with its time to the client node

- Client node sets its clock to the reply

- Problems:
  - How to estimate the delay incurred by the server's reply in reaching the client?
  - What if time server time is less than client's current time?

- Handling message delay: try to estimate the time the message with the timer server's time took to reach the client
  - Measure round trip time and halve it
  - Make multiple measurements of round trip time, discard too high values, take average of rest
  - Make multiple measurements and take minimum
  - Use knowledge of processing time at server if known to eliminate it from delay estimation (How to know?)
- Handling fast clocks
  - Do not set clock backwards; slow it down over a period of time to bring in tune with server's clock
    - Ex: increase the software clock every two interrupts instead of one

- Can be used for external synchronization if the time server is synchronized with external clock reference
  - Requires a special node with a time source
- What if the time server fails?
  - Usually a problem, as it is assumed that the time server is special (synchronized with external clock or at least with a more reliable clock)
- Works well in small LANs, not scalable to large number of nodes over WANs
  - Load on the central server will be high, affecting its processing time, in turn affecting synchronization error
  - Delay variance increases in larger networks

# Berkeley Algorithm

- Centralized as in Cristian's, but the time server is active
- Time server asks for time of other nodes at periodic intervals
- Other nodes reply with their time
- Time server averages the times and sends the adjustments (difference from local clock) needed to each machine
  - Adjustments may be different for different machines
  - Why do we send adjustments, and not the new absolute clock value?
- Nodes sets their time (advances immediately or slows down slowly) to the new time

- Time server can handle  faulty clocks by eliminating client clock values that are too low or too high

- What if the time server fails?

  - Just elect another node as the time server (Leader Election Problem)

  - Note that the actual time of the central server does not matter, enough for it to tick at around the same rate as other clocks to compute average correctly (why?)

- Cannot be used for external synchronization

- Works well in small LANs only for the same reason as Cristian's

# External Synchronization with Real Time

- Clocks must be synchronized with real time

- But what is "real time" anyway?

# Measurement of Time

- Astronomical Time
  - Traditionally used
  - Based on earth's rotation around its axis and around the sun
  - Solar day : interval between two consecutive transits of the sun
  - Solar second : 1/86,400 of a solar day
  - Period of earth's rotation varies, so solar second is not stable
  - Mean solar second : average length of large no of solar days, then divide by 86,400

- Atomic Time
  - Based on the transitions of Cesium 133 atom
  - 1 sec. = time for 9,192,631,770 transitions
  - Many labs worldwide maintain a network of atomic clocks
  - International Atomic Time (TAI) : mean no. of ticks of the clocks since Jan 1, 1958
  - Highly stable
  - But slightly off-sync with mean solar day (since solar day is getting longer)
  - A leap second inserted occasionally to bring it in sync.
  - Resulting clock is called UTC – Universal Coordinated Time

- UTC time is broadcast from different sources around the world, ex.
  - National Institute of Standards & Technology (NIST) – runs WWV radio station, anyone with a proper receiver can tune in
  - United States Naval Observatory (USNO) – supplies time to all defense sources
  - National Physical Laboratory in UK
  - Satellites
  - Many others
  - Accuracies can vary (< 1 milliseconds to a few milliseconds)

# Synchronizing with UTC Time

- Put an atomic clock in each node!!
  - Too costly
  - Most often the accuracy is not needed, so the cost is not worth it
- Put a GPS receiver at each node
  - Still costly
  - GPS does not work well indoor
- Can use a Cristian-like algorithm with the time server sync'ed to a UTC source
  - Not scalable for internet-scale synchronization
- Solution: Use a hierarchical approach

# NTP : Network Time Protocol

- Protocol for time synchronization in the internet
- Hierarchical architecture
  - Stratum 0: reference clocks (atomic clocks or receivers for time broadcast by national time standards or satellites, ex. GPS)
  - Stratum 1: primary servers with reference clocks
    - Most accurate
  - Stratum 2, 3,… servers synchronize to primary servers in a hierarchical manner (stratum 2 servers sync. with stratum 1, stratum 3 with stratum 2 etc.)
    - Lower stratum no. means more accurate
    - More servers at higher stratum no.

- Different communication modes
  - Multicast (usually within LAN servers)
    - One or more servers periodically multicasts their time to other servers
  - Symmetric (usually within multiple geographically close servers)
    - Two servers directly exchange timing information
  - Client server (to higher stratum servers)
    - Cristian-like algorithm
- Communicates over UDP
- Reliability ensured by synchronizing with redundant servers
- Accuracy ensured by combining and filtering multiple time values from multiple servers
- Sync. possible to within tens of milliseconds for most machines
  - But just a best-effort service, no guarantees