CprE 450/550x Distributed Systems and Middleware

#### Synchronization

Yong Guan 3216 Coover Tel: (515) 294-8378 Email: guan@ee.iastate.edu April 6 & 8, 2004

#### Readings for Today's Lecture

- References
- > Chapter 5 of "Distributed Systems: Principles and Paradigms"
- Chapter 14 of Coulouris: "Distributed Systems"

#### Clock Synchronization

- In a centralized system, time is unambiguous
- In a distributed system, achieving agreement on time is not trivial.
- Example: UNIX makefile
- A change to one source file only requires one file to be recompiled, not all the files
- How make works?
  - Examine the times at which all the source and object files were last modified.
  - In a distributed system in which there is no global agreement on time, how?
- . Is it possible to synchronize all the clocks in a distributed system?



#### **Physical Clocks**

- Almost all computers have a circuit for keeping track of time.
- Computer Timer is a machined quartz crystal When kept under tension, quartz crystal oscillates at a well-defined frequency, depending on the kind of crystal, how it is cut, and the amount of tension.
  - Two registers: a counter and a holding register
  - Each oscillation decrements the counter by one, when it gets to 0, an interrupt is generated and the counter is reset from the holding register.
  - Each interrupt is called a clock tick.
  - When the system is booted initially, date and time are required to be entered and deposited in CMOS RAM. Each clock tick increases the time stored in CMOS RAM by one such that software clock can be maintained.

## **Physical Clocks**

- It doesn't matter if the clock is off by a small of amount for a single computer with a single clock.
- For multiple CPUs with their own clocks, things change: Though the frequency at which a crystal oscillator runs is fairly stable, it is impossible to guarantee the crystals on different computers run at the same frequency.
  - The crystals will run at slightly different rates, which result in the clocks out-of-sync. The time value difference is called clock sk
  - Programs depending on the time associate with files, objects, messages may fail due to these clock skew
- How do we synchronize the clocks with real-world clocks?



#### **Physical Clocks**

- With the invention of atomic clock in 1948, measuring time becomes more accurately by counting transitions of the cesium 133 atom.
- Physicists took over the job of timekeeping from astronomers
- A second is defined as the time it takes the cesium 133 atom to make exactly 9, 192,631,770 transitions. This number makes an atomic second equal to the mean solar second.
- BI H averages the number of clock ticks from 50 laboratories in the world to produce I nternational Atomic Time (TAI).
- TAI=the mean number of ticks of the cesium 133 clocks since midnight on Jan. 1, 1958 divided by 9, 192,631,770

#### Physical Clocks

- 86,400 TAI seconds is 3 msec less than a mean solar day.
- Over the years, noon would become earlier and earlier.
- BI H introduce leap seconds whenever the difference between TAI and solar time grows to 800 msec.
- Universal Coordinated Time (UTC) (replaced Greenwich Mean Time, which is astronomical time)
- NIST operates a shortwave radio station with call letters WWV from Fort Colins, CO.
- WWV broadcasts a short pulse at the start of each UTC second. +-1msec (+-10msce due to atmosphere fluctuations).
- Similar services, UK's MSF, GEOS (earth satellite), etc.



## Clock Synchronization Algorithms

 Each machine is assumed to have a timer that causes an interrupt *H* time a second. When the timer goes off, the interrupt handler adds one to a software clock.

C: value of the clock  $C_p(t)\text{: The value of the clock at machine p at UTC time }t.$ 

- I deally, C<sub>p</sub>(t)=t for all p and t. i.e., dC/dt=1
- In practice, the relative error obtainable with modern timer chips is 10<sup>-5</sup>.
- Maximum drift rate r, where 1-r <=dC/dt <=1+r.</p>







## Averaging algorithm

- Dividing time into fixed-length re-sync intervals.
- At the beginning, each machine broadcasts its own time.
- After a machine broadcasts its time, it starts a local timer to collect all other broadcasts that arrive during some time interval S.
- Then,
  - Average the values from all the other machines
  - Discard the m highest and m lowest values, and average the remaining ones.

  - NTP (Network Time Protocol) - Can be further improved

## Multiple External Time Sources

WWV, GEOS receivers

## Use of Synchronized Clocks

At-most-once message dilivery

- Message seq number, what about system crashes and reboots?
- ConnID+timestamp

























## Global State (1)

- Knowing the global state in distributed systems is useful on many occasions.
- The global state consists of the local state of each process, together with the messages-in-transit.
- Distributed Snapshot (Chandy and Lamport'85)









## **Election Algorithms**

- Many distributed algorithms requires one process in the system acts as a leader (coordinator, initiator).
- It does not matter which process it is, but one of them has to do it.
- The goal of election algorithm is to ensure that when an election starts, it concludes with all processes agreeing on who the new coordinator is to be.















ompariso	on		
Algorithm	Messages per entry/exit	Delay before entry (in message times)	Problems
Centralized	3	2	Coordinator crash
Distributed	2 ( n – 1 )	2 (n - 1)	Crash of any process
Token ring	1 to ∞	0 to n – 1	Lost token, process crash
A compa	rison of three	e mutual exclusio	on algorithms.





he Transaction I	Model (2)	
Primitive	Description	
BEGIN_TRANSACTION	Make the start of a transaction	
END_TRANSACTION	Terminate the transaction and try to commit	
ABORT_TRANSACTION	Kill the transaction and restore the old values	
READ	Read data from a file, a table, or otherwise	
WRITE	Write data to a file, a table, or otherwise	
Examples of	primitives for transactions.	







				5
Writeahead Log				
x = 0; y = 0; BEGIN TRANSACTION:	Log	Log	Log	
x = x + 1; y = y + 2 x = y * y;	[x = 0 / 1]	[x = 0 / 1] [y = 0/2]	[x = 0 / 1] [y = 0/2] [x = 1/4]	
END_TRANSACTION; (a)	(b)	(c)	(d)	
a) A transaction b) – d) The log be	efore each state	ment is execut	ted	























• Theorem: If *S* is any schedule of two-phase transactions, then *S* is serializable.

Two-Phase Locking (cont)• Theorem: If S is any schedule of two-phase transactions,<br/>then S is serializable.• Proof:<br/>Suppose not. Then the serialization graph G for S has a<br/>cycle,<br/> $T_{II} \rightarrow T_{I2} \rightarrow \dots \rightarrow T_{Ip} \rightarrow T_{II}$ Therefore, a lock by  $T_{II}$  follows an unlock by  $T_{II}$ ,<br/>contradicting the assumption that  $T_{II}$  is two-phase.



# Optimistic Concurrency Control

"Forgiveness is easier to get than permission"

- Basic idea:
  - Process transaction without attention to serializability.
  - Keep track of accessed data items.
  - At commit point, check for conflicts with other transactions.
  - Abort if conflicts occurred.

#### Problem:

 would have to keep track of conflict graph and only allow additional access to take place if it does not cause a cycle in the graph.







