

Synchronization of processes

The Coordinated Attack Problem

What is possible

□ It is possible for two processes to decide that they change their states **some time** in the future.



Typical communication protocol

Sender sends a message

• Now only the sender knows about the message.

- Receiver accepts message, sends an acknowledgement
 - Now both sender and receiver knows about the message.
 - Only the receiver knows that is has gone through.
- □ The sender get acknowledgement
 - Now both sender and receiver knows about the message.
 - Now both sender and receiver knows that is has gone through.
- Sender sends next message



6 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Arne Andreasson - Computer Science and Engineering

CHALMERS

Synchronization of processes

Consequences 2

□ It is impossible for the processes in a distributed system to know the **global** state of the system without stopping it.

Ordering the Events in a Distributed System

- Leslie Lamport, 1978.
- Presumptions:
 - A number of processes are communicating using message passing in an unreliable network. Each process is a sequence of events.
 - These events (in the same process) are totally ordered in time.

9 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Arne Andreasson - Computer Science and Engineering

- A message passing will have an unknown transmission time > 0.
- The sending and the receiving of messages will be seen as the events in a process.

() CHALMERS



Partial Ordering of Events

$" \rightarrow "$ (happens) before

Definition: " \rightarrow " on the set of events in a distributed system is the least relation that fulfills:

- (*i*) If *a* and *b* are two events in the same process and *a* happens before *b*, $a \rightarrow b$ holds.
- (*ii*) If a is the sending of a message in one process and b is the receiving of the same message in another process, then $a \rightarrow b$ holds.
- (*iii*) If both $a \rightarrow b$ and $b \rightarrow c$ holds, then $a \rightarrow c$ holds.
- (*iv*) If an event a "not happens before" an event b it is denoted as $a \rightarrow b$.
- (v) If $a \rightarrow b$ and $b \rightarrow a$ holds, then a and b in a way are "concurrent", this is denoted as a / / b:
- \square " \rightarrow " is an irreflexive partial order on the set of events in a system.

10 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Arne Andreasson - Computer Science and Engineering

```
CHALMERS
```

Vector Clocks

- □ A vector (array) is used to represent the partial ordering " \rightarrow ".
- Each process holds a clock vector *C_i* [1..*n*] *n* is the number of processes in the system.

It is defined as:

- $C_i[i]$ describes the time in process P_i .
- C_i [j] describes process P_i's knowledge about the time in process P_j. Updated each time P_i gets a message from P_i.
- Each event in P_i is time stamped with the entire vector C_i [1.n]. This represents P_i's global view of the event.
- This leads to the following rules for each process P_i:

(i) Each process P_i increments its local register C_i [i] for each event.

(ii) Each message is stamped with the actual full clock vector C_i [1..n].

(iii) When receiving a message a process should adjust its local clock value, C_i [1..n], comparing it to the message timestamp T [1..n] according to:
1 ≤ k ≤ n: C_i [k] := max (C_i [k],T [k]), then C_i [i] is incriminated.

11 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Ame Andreasson - Computer Science and Engineering

```
() CHALMERS
```



Implementing Partial Order using Vector Clocks

 \Box The partial ordering, " \rightarrow ", can be calculated by comparing the vector clock values.

Definitions:

(i) If $C_i \le C_j \iff \forall x: C_i [x] \le C_j [x]$ vector C_i is said to dominate vector C_i .

(ii) If $C_i < C_j \iff C_i \le C_j \land \exists x: C_i [x] < C_j [x]$ vector C_i is said to strictly dominate vector C_i .

(iii) If $C_i // C_j \iff \neg (C_i < C_j) \land \neg (C_j < C_i)$ vector C_i and vector C_i is said to be "parallel".

 \Box Then it holds for two events *a* and *b* with vector timestamps C_a and C_b :

$$a \to b \iff C_a < C_b$$

 $a \parallel b \iff C_a \parallel C_b$

□ Note that this is a partial ordering, there will be events *a* and *b* with vectors like (iii) above. Sometimes we say we have a conflict when this happens.

14 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Ame Andreasson - Computer Science and Engineering

```
CHALMERS
```

Logical Clocks

Leslie Lamport, 1978.

A Logical Clock is a function that assigns an ordering number C(a) to an event a in the system.

- Clock condition: $\forall a,b: a \text{ and } b$ are events in the system: $a \rightarrow b \Rightarrow C(a) < C(b).$
- (N.B.! We don't require the reverse)
 - Each process P_i has an en instance C_i of the logical clock
 - If *a* is an event in process $P_i \implies C(a) = C_i(a)$
 - If *a* and *b* are events in process P_i and $a \rightarrow b \Rightarrow C_i(a) < C_i(b)$
 - the clock must step forward (tick) at least by 1 between two events in the same process.
 - If *a* is the sending of a message by process P_i and *b* is the receiving of the same message by process $P_j \Rightarrow C_i(a) < C_j(b)$
 - a message must be received at a clock value that is higher than the sending clock value



Implementing Logical Clocks

- \bigcirc Each process P_i increments its local register C_i after each event in the process.
- Each message contains a timestamp T_m .
 - If *a* is the event that process P_i sends a message *m*, then m will contain the timestamp $T_m = C_i(a)$.
- When a process receives a message it first compares its timestamp to the local logical clock register.
 - Before a process P_j registers the reception of a message m, P_j will set its local logical clock register $C_j := max(C_j, T_m + I)$ which will be the logical clock value for the message receiving event.

17 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Arne Andreasson - Computer Science and Engineering



- Since the processes do not know the physical time they have to use the Logical Clock values, i.e. they will have to follow the green lines.
 - For the processes the green lines will be thought of as vertical.
 - The events might be in different order than according to real time, but if for two events, a and b, that $a \rightarrow b$ holds then the logical clock value for a must be lower than for b.



Let \boldsymbol{a} be an event in process P_i and \boldsymbol{b} be an event in process P_i :

(*i*) if $C_i(a) < C_i(b)$ then $a \rightarrow_T b$.

(*ii*) if $C_i(a) = C_i(b)$ and $P_i \nmid P_i$ then $a \to_T b$.

 \square " \rightarrow_T " is a total ordering of the events in a system where processes communicate using message passing in a computer network:

$$- \neg (a \to_{\mathrm{T}} b) \Rightarrow (b \to_{\mathrm{T}} a)$$

$$- a \rightarrow b \Rightarrow a \rightarrow_{\mathrm{T}} b$$

() CHALMERS

 \Box " \langle " will introduce priorities on the processes

20 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Arne Andreasson - Computer Science and Engineering







- \bigcirc P; REOUEST:
 - $< T_i$; P_i; *REQUEST*> is sent to all other processes, T_i is the timestamp taken from C_i
 - $< T_i; P_i; REQUEST >$ is also put in P_i's local copy of the Request Queue sorted according to " \rightarrow_T "
 - P_i increments its local Logical Clock value
- P_i receives $REQUEST < T_i$; P_i; REQUEST >:
 - P_i adjusts its local copy of the Logical Clock according to the Logical Clock definition
 - $< T_i$; P_i ; REQUEST > is put in P_i 's local copy of the Request Queue sorted according to " \rightarrow_T "
 - P_i updates its TS-table
 - P_i increments its local Logical Clock value
 - P_i sends an acknowledgement $\langle T_i; P_i; ACK \rangle$ to P_i
 - P_i increments its local Logical Clock value
- P_i receives an acknowledgement $\langle T_i; P_i; ACK \rangle$ from P_i:

() CHALMERS

Distributed Resource Allocation Algorithm using Logical Clock

- □ Algorithm conditions:
 - The algorithm uses
 - Distributed Request Queue (empty at start)
 - · Logical Clock
 - each process administers:
 - a local copy of the Request Queue
 - · a local copy of the Logical Clock
 - a table containing the latest received timestamp from each of the other processes in the system (TS-table)
 - \bigcirc When placing a request at process P_i, it is associated with the current local logical clock value C_i in P_i.
 - all messages between two processes is delivered in the same order as they were sent and no message will disappear (FIFO). This will be the case when using a communication protocol such as TCP/IP.

22 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Arne Andreasson - Computer Science and Engineering

CHALMERS

- Pi adjusts its local copy of the Logical Clock according to the Logical Clock definition.
- P_i updates its TS-table with T_i from P_i
- · P_i increments its local Logical Clock value
- \bigcirc P_i is allowed access to the resource when:
 - $\langle T_i; P_i; REQUEST \rangle$ is number one in the (local) Request Queue
 - $T_i \rightarrow_T T_i$ for all T_i in the (local) TS-table
- \bigcirc P_i want to *RELEASE*:
 - $\langle T_i \rangle$; P; RELEASE> is sent to all other processes, T_i is the timestamp taken from actual C_i
 - <*T_i*;*P_i*;*REQUEST*> is erased from the (local) Request Queue
 - · P_i increments its local Logical Clock value
- \bigcirc P_i receives $\langle T_i; P_i; RELEASE \rangle$:
 - P_i adjusts its local copy of the Logical Clock according to the Logical Clock definition.
 - $<T_i$; P_i; *REQUEST*> is erased from the (local) Request Queue
 - P_i updates its TS-table with T_i from P_i
 - · P_i increments its local Logical Clock value

24 (48) - DISTRIBUTED SYSTEMS Process Synchronization - Sven Arne Andreasson - Computer Science and Engineering



Physical Clocks Totally ordered Multicast Why? Let S stand for all events in the system, □ The algorithm can be used for sending totally ordered Multicast Messages. not only messages in the network but also other interactions □ Instead of sending Requirement messages, Multicast messages with logical clock values can be sent and \square " \rightarrow_{a} " is defined as happens before in S ordered according to the total order in a message queue. How avoid abnormal behavior? The messages are delivered to the receiving process following the total order definition. (i) require that all interaction should involve exchange of timestamps (ii) use synchronized physical clocks Strong clock condition: $\forall a, b \in S: a \rightarrow_a b \Rightarrow C(a) < C(b)$ 25 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Arne Andreasson - Computer Science and Engineering 26 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Arne Andreasson - Computer Science and Engineering () CHALMERS CHALMERS **Clock Conditions** Setting clock values using messages between nearby processes • Each clock must have high precision, i.e. have correct speed: $\exists \kappa \ll l: \forall i: \left| \frac{d}{dt} c_i(t) - 1 \right| < \kappa$ • We use the following notation: μ : the minimum time to send a message between two processes Typical value for a modern clock might be $\kappa < 10^{-6}$. v_m : the real time to send a message m • The value difference between two clocks should not be too large, ξ_m : the den unpredictable delay of a message m i.e. the clocks should be synchronized: then $\forall i,j: |C_i(t) - C_i(t)| < \varepsilon$ where ε is the maximal difference between any two clocks. $\xi_m = v_m - \mu$. for message m The clocks must show the "same" time. \Box Let μ denote the shortest time for a process to influence another process

For a correct behavior of the system it is required that $\epsilon/(1-\kappa) \le \mu$, since then $C_i(t+\mu) - C_j(t) > 0$ This requirement will assure that every message will experience that time is moving forward

• consequently $\varepsilon \leq \mu(1-\kappa)$ must hold

One solution would be to send clock values between the processes to set clock values as close as possible

Clock adjusting algorithm

- Leslie Lamport, 1978.
 - If process P_i doesn't get a message at time t it lets its local clock C_i step forward using its internal time keeper.
 - If process P_i sends a message m at time t, this message will contain a timestamp $T_m = C_i(t)$
 - If process P_i receives a message m at time t containing a timestamp T_m P_i should set its local clock according to:
 C_i := max(C_i, T_m + μ_m)
- A network with diameter D (the maximal number of the minimum number of node jumps between any two nodes in the network).
- \Box Each node sends a message to each neighbor at least every τ time units.
 - $\label{eq:constraint} \begin{array}{l} \bullet \\ \bullet \\ \bullet \\ \bullet \\ D(2\kappa\tau + \xi) \end{array}$ Then the clock difference will be maximized by
 - The problem here is that theoretically this is unlimited since ξ is unlimited.
- □ The fastest clock will push the other clocks forward.

29 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Ame Andreasson - Computer Science and Engineering

() CHALMERS

() CHALMERS



Instead of taking the highest value as the new value, the mean value of the clocks is chosen.. Since no clock is allowed to be changed backwards they change their speeds in order to reach the next calculated value, the aiming value.

The blue line shows what is supposed to be the accurate time of the system.



- Most messages have short transmission time.
- \Box but at least a minimal time, μ





Probabilistic Clock Synchronizing

- □ Flaviu Cristian, 1989.
- □ To set the clients clock a time request message is sent to the clock server.
 - \bigcirc Total waiting time from sending request until receiving reply is 2 * U
 - The received time value is *T*. It is assumed that this time was read by the server after half the waiting time *U*. Then the estimated time when receiving the reply will be T + U.
- **Calculaing the error**:
 - The server can only have read this time between the sending of the request and the receiving of the reply. Thus the error can be estimated to be $\pm U$.
 - Since there is a minimum transmission time, μ , for the messages this error estimation can be refined to $\pm (U \mu)$
 - **O** But we have to read our own clock to get *U*. Since the clock might have a speed error κ we have to add $\kappa * U$ to the error estimathion
- Thus the time can be estimated as $T + U \pm (U \mu + \kappa * U)$
- 33 (48) DISTRIBUTED SYSTEMS Process-Process Synchronization Sven Arne Andreasson Computer Science and Engineering

CHALMERS

Snapshot Algorithm

Chandy and Lamport, 1985.

- Calculates a "useful" approximate global state of a distributed system, although it might never have existed.
- Algorithm conditions:
 - the nodes are part of a directed graph, real or virtual
 - the graph is strongly connected, i.e. there is a path from any node to any other node
 - the algorithm messages are only sent according to the directed graph
 - FIFO secure transmissions on the communication links (e.g. TCP/IP)
 - the process can continue working and change their states while the snapshot takes place
- Calculates an estimated global state of the system consistent with the real state.
 - this might never have existed
 - but this state is a possible state that doesn't conflict with the real states and might be used for system control.



Consistent State

- \Box Assume that S₁ is the real global state when the algorithm starts and S₂ is the real state when the algorithm terminates.
- \Box A calculated state S_c is consistent with the real state if
 - \circ S_c = S₁, or
 - \bigcirc S_c = S₂, or
 - \bigcirc S₁ \rightarrow S_c \rightarrow S₂ is a possible state sequence.

() CHALMERS

^{35 (48) -} DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Arne Andreasson - Computer Science and Engineering

Snapshot-algorithm

- uses Marker Messages (MM)
- □ Local Calculated State (LS)
 - A node initiates the algorithm
 - records its local state, LS
 - puts itself in recording state
 - sends an MM-message on each out-going link
 - When a node gets its first *MM*-message on an in-going link
 - records its local state, LS
 - puts itself in recording state
 - marks the in-going link as ready
 - sends an MM-message on each out-going link
 - When a node in recording state gets an *MM*-message on an in-going link
 - marks the in-going link as ready

37 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Arne Andreasson - Computer Science and Engineering

CHALMERS

() CHALMERS

Example: Counting Capabilities

- A Capability is a permission (grant) ticket to be used in a system for access to different resources.
 - the CORBA ior-string
 - the right to use a certain program (license)
 - the right to read/change a file
 - a Control Token
- □ In a distributed system a capability can be sent in a message between processes.
- □ In this example we assume that the number of capabilities in the system should be constant (e.g. we have paid a license for a program that allows n users in the same time.
- A capability might get lost or one process might not delete its own copy when sending it to another process.
 - We need an algorithm to count the capabilities in the system at one moment
 - We don't want to stop the system to do this
 - Then we can use the Snapshot Algorithm



• its recorded state is - sent to one node for assembling the global calculated state

• When a node in recording state gets another message on an in-going link

- or sent to all other nodes so all can calculate the global state
- it leaves recording state
- A calculating node that has got the local calculated states, *LS*, from all other nodes can put them together to the Calculated Global State

• if the in-going link is not marked as ready the recorded state, LS, will be recalculated according to

□ The Calculated Global State might not have been existing for the system.

38 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Arne Andreasson - Computer Science and Engineering

() CHALMERS

- \square Nodes N₁, N₂, N₃, N₄, N₅ and N₆
- \Box Capabilities K₁, K₂ and K₃.
- \square Node N₁ starts the algorithm by recording K₁ and then sending out MM on outgoing links.
- \Box Capabilities K_2 and K_2 are being transferred among nodes.



(CHALMERS

- \Box When node N₃ and node N₄ received the MM-messages they had no capabilities. They both record empty.
- \Box Capability K₃ has reached node N₆.



- The MM-messages have now reached nodes N₂, N₄ and N₅.
- Since node N₄ now has got MM-messages on all in links it is ready with its part of the algorithm. For the global state N₄ contributes with K₂. Nodes N₂ and N₅ record empty
- Nodes N₂ and N₅ have not got an MM before so they send an MM-message on each outgoing link. Node N₅ has got MM-message on all in links and is ready.
- □ Capability K₁ has reached node N₃. Since N₃ already has got an MM-message on the corresponding link it does not record K₁ as a part of its contribution to the algorithm.



43 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Arne Andreasson - Computer Science and Engineering

- □ The nodes N₃ and N₄ sends MM-messages on their outgoing links.
- Capability K₂ then reached N₄ and is recorded since N₄ has not got an MM-message on the corresponding link yet.
- Also node N_1 has sent capability K_1 in a message in its way to N_3



- **C**apabilities K_3 and K_2 have been sent.
- Since the nodes N₂ and N₅ have not got an MM-message before they send MM-messages on all their outgoing links.
- \Box K₂ and K₃ have reached N₅ and N₃

Since N_5 has already got an MM-message on the corresponding link K_2 is not recorded by N_5 . K_3 is on the other side recorded by node N_3 since N_3 has not got an MM-message on the corresponding link.



44 (48) - DISTRIBUTED SYSTEMS Process-Process Synchronization - Sven Ame Andreasson - Computer Science and Engineering

() CHALMERS

- \square The MM-message from N₂ has reached N₁ and the MM-message from N₅ has reached N₆.
- \Box Node N₁ then is ready with its part of the algorithm.
- \Box Node N₆ gets its first MM-message and records its state which is empty.



- □ N₆ then sends MM-messages on its outgoing links.
- \Box K₃ has arrived to N₄.



- \Box K₂ has been sent from N₅ to N₆.
- \Box The MM-messages have reached N₂ and N₃. Now all nodes have got MM-messages on all in links and the local calculations are ready.



- **u** By letting one or all nodes combine the local recorded states a global state will be calculated.
- □ This state might never existed in reality but can still be used for the distributed control.
 - In our example it gives the correct number of capabilities in the system although it does not tell where they actually are.
- □ The calculated state:



