---

# Concurrency Control in Distributed Environment

*Ch 8 [BenA 06]*

Messages
Channels
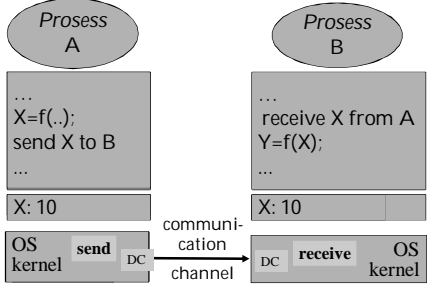Rendezvous
RPC and RMI

---

# Distributed System

- No shared memory
- Communication with messages
- Tightly coupled systems
  - Processes alive at the same time
- Persistent systems
  - Data stays even if processes die
- Fully distributed systems
  - Everything goes

---

# Communication with Messages (4)

Prosess A

```
…
X=f(..);
send X to B
…
```
X: 10

OS kernel  **send**  DC

communi-cation channel →

Prosess B

```
…
receive X from A
Y=f(X);
…
```
X: 10

DC  **receive**  OS kernel

- Sender, receiver
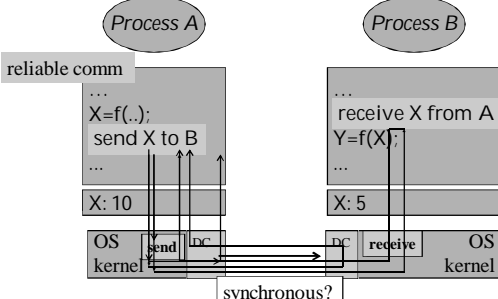- Synchronous/asynchronous communication

---

# Message Passing

- Synchronous communication
  - Atomic action
  - <u>Both</u> wait until communication complete
- Asynchronous communication          Usual case
  - Sender <u>continues</u> after giving the message to OS for delivery
  - <u>May</u> get an acknowledgement later on
    - Message received or not
- Addressing
  - Some address for receiver <u>process</u>          prosessi
    - Process name, id, node/name, …
  - Some address for the communication <u>channel</u>          kanava
    - Port number, channel name, …
  - Some address for <u>requested service</u>          palvelu
    - <u>Broker</u> will find out, sooner or later          meklari
      - After message has been sent?
    - Service address not known at service request time

---

# Synchronization levels (10)

Process A

reliable comm
```
…
X=f(..);
send X to B
…
```
X: 10

OS kernel  **send**  DC

Process B

```
…
receive X from A
Y=f(X);
…
```
X: 5

DC  **receive**  OS kernel

synchronous?

---

# Synchronization levels (1/5)

Process A

```
…
X=f(..);
send X to B
…
```
X: 10

OS kernel  **send**  DC

Process B

```
…
receive X from A
Y=f(X);
…
```
X: 5

DC  receive  OS kernel

---

## Synchronization levels (2/5)

Process A

Process B

asynchronous?

...
X=f(..);
send X to B
...

X: 10

receive X from A
Y=f(X);
...

X: 5

OS kernel | **send** | Data Com

DC | receive | OS kernel

14.3.2011                     Copyright Teemu Kerola 2011                     7

## Synchronization levels (3/5)

Process A

Process B

asynchronous?

...
X=f(..);
send X to B
...

X: 10

receive X from A
Y=f(X);
...

X: 5

OS kernel | **send** | DC
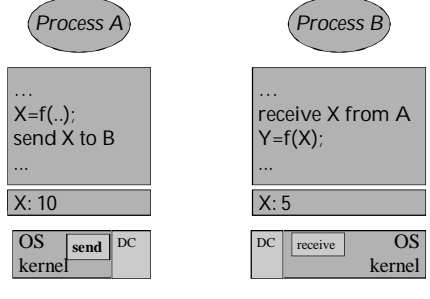
DC | receive | OS kernel

14.3.2011                     Copyright Teemu Kerola 2011                     8

## Synchronization levels (4/5)

Process A

Process B

reliable comm

...
X=f(..);
send X to B
...

X: 10

receive X from A
Y=f(X);
...

X: 5

OS kernel | **send** | DC

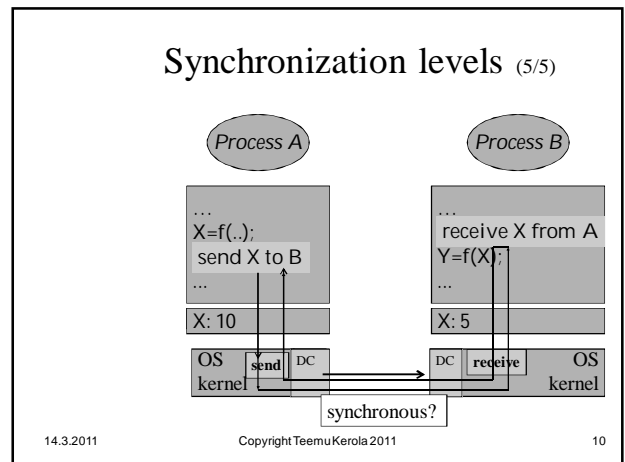DC | **receive** | OS kernel

14.3.2011                     Copyright Teemu Kerola 2011                     9

## Synchronization levels (5/5)

Process A

Process B

...
X=f(..);
send X to B
...

X: 10

receive X from A
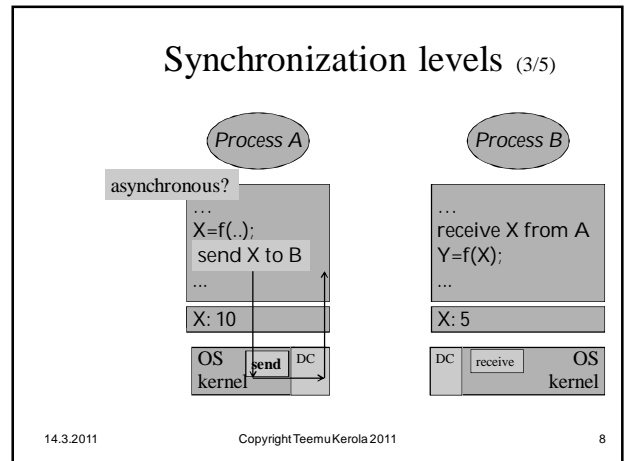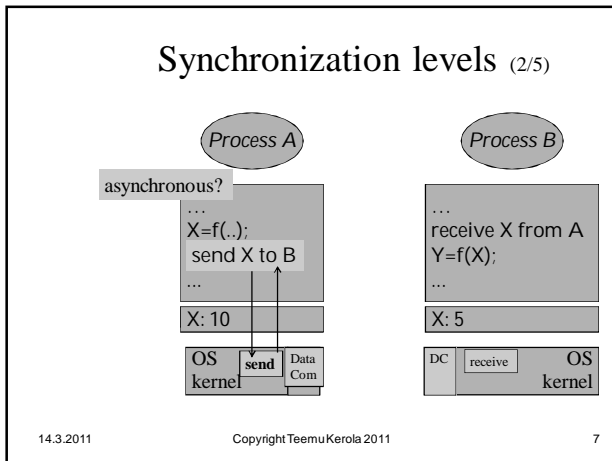Y=f(X);
...

X: 5

OS kernel | **send** | DC

DC | **receive** | OS kernel

synchronous?

14.3.2011                     Copyright Teemu Kerola 2011                     10

## Message Passing

- Symmetric communication
  - Cooperating processes at same level
  - Both know about each others address
  - Communication method for a fixed channel

- Asymmetric communication
  - Different status for communicating processes
  - Client-server model
    - Server address known, client address given in request

- Broadcast communication
  - Receiver not addressed directly
  - Message sent to everybody (in one node?)
  - Receivers may be limited in number
    - Just one?
    - Only the intended recipient(s) will act on it?

14.3.2011                     Copyright Teemu Kerola 2011                     11

## Wait Semantics

- Sender
  - Continue after OS has taken the message   `Usual case`
    - Non-blocking send
  - Continue after message reached receiver <u>node</u>
    - Blocking send
  - Continue after message reached receiver <u>process</u>
    - Blocking send
- Receiver
  - Continue only after message received   `Usual case`
    - Blocking receive
  - Continue even if no message received
    - Status indicated whether message received or not
    - Non-blocking receive

14.3.2011                     Copyright Teemu Kerola 2011                     12

## Message Passing

- Data flow
  - One-way
    - Synchronous may be one-way
    - Asynchronous is always one-way
  - Two-way
    - Synchronous may be two-way
    - Two asynchronous communications

  data flow vs. control flow!

- Primitives
  - One message at a time
  - Need addresses for communicating processes
  - Operating system level service
  - Usually not programming language level construct
    - Too primitive: need to know node id, process id, port number,…

14.3.2011    Copyright Teemu Kerola 2011    13

14.3.2011    Copyright Teemu Kerola 2011    14

## Channels

- History of languages utilizing channels
  - Guarded Commands    vartioidut komennot
    - Dijkstra, 1975

    Edsger Dijkstra

  - Communicating Sequential Processes    kommunikoivat sarjalliset prosessit
    - CSP, Hoare, 1978

  - Occam
    - David May et al 1983    C.A.R. Hoare
    - Hoare as consultant
    - Inmos Transputer

    David May

14.3.2011    Copyright Teemu Kerola 2011    15

## Guarded Commands (Dijkstra)

- Way to describe predicate transformer semantics
- Communication not really specified    predikaatti-muunnos-semantiikka
- Guarded command    C → S
  - Condition or guard
  - Statement
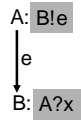
greatest common divisor

```
x, y = X, Y    -- statement (unguarded)
do   -- loop command, loop terminates when x = y
  x ≠ y →
    if   -- conditional command (itself guarded)
      x > y → x := x-y    -- guarded statement in the if
      y > x → y := y-x
    fi
od
print x ; -- another statement, also unguarded
```

guard

vartioitu lauseke

can be also input/output statement

http://en.wikipedia.org

14.3.2011    Copyright Teemu Kerola 2011    16

## Communicating Sequential Processes – CSP (Hoare)

- Language for modeling and analyzing the behavior of concurrent communicating systems
- A known group of processes A, B, …
- Communication:
  - output statement: B!e
    - evaluate e, send the value of e to B
  - input statement: A?x
    - receive the value from A to x
  - input, output: blocking statements
  - output & input: "distributed assignment"
    - Communicate value from one process to a variable in some other process

A: B!e

e

B: A?x

14.3.2011    Copyright Teemu Kerola 2011    17

## CSP communication

- Input/output statements
  - Destination!port ($e_1$, …, $e_n$) ;
  - Source?port ($x_1$, …, $x_n$) ;
- Binding
  - Communication with **named processes**
  - Matching types for communication
- Example:    **Copy** ( West => Copy => East )

**West:**
```
do true ->
  Copy!c;
  …
od
```

**Copy:**
```
do true ->
  West?c;
  East!c ;
od
```

**East:**
```
do true ->
  Copy?c;
  …
od
```

14.3.2011    Copyright Teemu Kerola 2011    18

## OCCAM Language

- Communication through **named channels**
  - Globally defined
    - Somewhere, in advance
  - Each channel has <u>one</u> sender and <u>one</u> receiver
    - both <u>processes</u> in some <u>nodes</u>
- Transputer
  - Multicomputer
    - E.g., 100 node Hathi-2 in ÅA
  - Automatic message routing for channels
  - Programmed with OCCAM

X := 5
c ! X

c ? Y
Y := Y*2

PAR
 SEQ
  X := 5
  c ! X
 SEQ
  c ? Y
  Y := Y*2

FPU
CPU
RAM
Links
Memory Interface

**IMS T800**

14.3.2011          Copyright Teemu Kerola 2011          19

## OCCAM Example

(Andrews, p 331)

West | East
c1
c2
dummy | EAsks

Copy

```
PROC Copy (CHAN OF BYTE West, EAsks, East)
  BYTE c1, c2, dummy;  -- buffer size = 2
  SEQ
    West ? c1              -- West has 1st byte
    WHILE TRUE
      ALT
        West ? c2          -- West has new byte
          SEQ
            East ! c1      -- send previous byte
            c1 := c2       -- copy to buffer c1
        EAsks ? dummy      -- East wants a byte
          SEQ
            East ! c1      -- send previous byte
            West ? c1      -- receive next one
```

block here, until other end ready ("guards")

- How to bind processes to nodes? 8 vs. 100 nodes?
- How to bind channels to processes, physical system?
  - 4 physical ports (N, S, E, W) in each processor    Discuss

14.3.2011          Copyright Teemu Kerola 2011          20

## Inmos Transputer



http://www.cs.bris.ac.uk/~dave/transputer.html

- B0042
- 2D array
- 10 boards 420 cpu's
- 30 boards 1260 cpu's

14.3.2011          Copyright Teemu Kerola 2011          21

## Channels

- Communication through <u>named channels</u>
  - Typed, global to processes
  - Programming language concept
  - <u>Any one</u> can read/write (usually limited in practice)

many readers/writers? same process writes and reads?

- Pipe or mailbox
- <u>Synchronous,</u> one-way (?)
- How to tie in with many nodes?
  - Not really thought through! Easy with shared memory!

### Algorithm 8.1: Producer-consumer (channels)

| channel of <u>integer</u> ch | | |
|---|---|---|
| **producer** | | **consumer** |
| integer x | | integer y |
| loop forever | | loop forever |
| p1: | $x \leftarrow produce$ | q1: | $ch \Rightarrow y$ |
| p2: | $ch \Leftarrow x$ | buffer size? | q2: | $consume(y)$ |

14.3.2011          Copyright Teemu Kerola 2011          22

## Filtering Problem

(inC) → compress aaaa→4a → (pipe) → output k'th ch →\n → (outC)
...aaaabbb...                                      ...4a\n3b...

- Compress many (at most MAX) similar characters to pairs …
  "compress"
  - {nr of chars, char}
- … <u>and</u> place newline (\n) after every K'th character in the compressed string
  "output"
- Why is it called "Conway's problem"?
  - "Classic <u>coroutine</u> example"   vuorottaisrutiinit

Conway, M. "Design of a separable transition-diagram compiler," *CACM* 6, 1963, pages 396–408.

14.3.2011          Copyright Teemu Kerola 2011          23

Filtering Problem with Channels

### Algorithm 8.2: Conway's problem

constant integer MAX ← 9
constant integer K ← 4
channel of integer inC, pipe, outC

| **compress** | | **output** | |
|---|---|---|---|
| char c, previous ← 0 | | char c | |
| integer n ← 0 | | integer m ← 0 | |
| $inC \Rightarrow previous$ | | | |
| loop forever | no last char? | loop forever | |
| p1: | $inC \Rightarrow c$ | q1: | $pipe \Rightarrow c$ |
| p2: | if (c = previous) and (n < MAX − 1) | q2: | $outC \Leftarrow c$ |
| p3: | $n \leftarrow n + 1$ | q3: | $m \leftarrow m + 1$ |
| p4: | else if n > 0 | q4: | if m >= K |
| p5: | $pipe \Leftarrow intToChar(n+1)$ | q5: | $outC \Leftarrow newline$ |
| p6: | $n \leftarrow 0$ | q6: | $m \leftarrow 0$ |
| p7: | $pipe \Leftarrow previous$ | q7: | |
| p8: | $previous \leftarrow c$ | | |

14.3.2011          Copyright Teemu Kerola 2011          24

## Matrix Multiplication with Channels

$$\begin{matrix} 1\ 2\ 3 \\ 4\ 5\ 6 \\ 7\ 8\ 9 \end{matrix} \quad X \quad \begin{matrix} 1\ 0\ 2 \\ 0\ 1\ 2 \\ 1\ 0\ 0 \end{matrix} \quad = \quad \begin{matrix} 4\ \ 2\ \ 6 \\ 10\ 5\ 18 \\ 16\ 8\ 30 \end{matrix}$$

- $16 = (7\ 8\ 9) \bullet (1\ 0\ 1)$
- $30 = (7\ 8\ 9) \bullet (2\ 2\ 0)$      $\underline{7*2+\ \ 8*2+\ \ 9*0+}\ 0\ = 30$
- Process for every <u>multiply-add</u>

channels        other processes

Result   |30|   * = 7 +   16   * = 8 +   0   * = 9 +   0   Zero

14.3.2011          Copyright Teemu Kerola 2011          25

---

Process Array for Matrix Multiplication

27 processes
24 channels

column 1

contains 1 row, sends it down one element at a time

West-bound multiply-add, South-bound copy North

contains 1 value, makes three multiply-adds, forwards values down

How to initialize everything?

How to synchronize everything?

Source    Source    Source

2 0 1        2 1 0        0 0 1

Result  4,2,6   1   3,2,4   2   3,0,0   3   0,0,0   Zero

2 0 1        2 1 0        0 0 1

Result  10,5,18   4   6,5,10   5   6,0,0   6   0,0,0   Zero

2 0 1        2 1 0        0 0 1

Result  16,8,30   * = 7 + N   9,8,16   8   9,0,0   9   0,0,0   Zero

2 0 1        2 1 0        0 0 1

Sink    Sink    Sink

14.3.2011          Copyright Teemu Kerola 2011          26

---

### Algorithm 8.3: Multiplier process with channels

integer FirstElement
channel of integer North, East, South, West
integer Sum, integer SecondElement          Relative names?

loop forever
p1:    North ⇒ SecondElement          ← wait 1st for this (*)
p2:    East ⇒ Sum          ← and <u>then</u> for this
p3:    Sum ← Sum + FirstElement · SecondElement
p4:    South ⇐ SecondElement
p5:    West ⇐ Sum

2 0 1

16,8,30   7   9,8,16

2 0 1

- How to map processes to nodes?
- How to map channels to processes?
  - North channel of one process the South channel of some other
- North-South data flow has priority (*)
  - Waiting even when data-flow East-West available
  - Node on East may be blocked unnecessarily

Discuss

14.3.2011          Copyright Teemu Kerola 2011          27

---

### Algorithm 8.4: Multiplier with channels and selective input

integer FirstElement
channel of integer North, East, South, West
integer Sum, integer SecondElement

loop forever
  either
p1:        North ⇒ SecondElement          If message from North available, do this
p2:        East ⇒ Sum
  or
p3:        East ⇒ Sum          If message from East available, do this
p4:        North ⇒ SecondElement
p5:    South ⇐ SecondElement          sequential block
p6:    Sum ← Sum + FirstElement · SecondElement
p7:    West ⇐ Sum

- Guarded statement
  - Execute <u>one</u> selective input statement
    - Nondeterministic selection (if both available)
    - p2 follows p1, it does not compete with p3

Discuss

14.3.2011          Copyright Teemu Kerola 2011          28

---

## Dining Philosophers with Channels

- Each <u>fork i</u> is a process, forks[i] is a channel
- Each <u>philosopher i</u> is a process

### Algorithm 8.5: Dining philosophers with channels

channel of boolean forks[5]

| philosopher i | fork i |
|---|---|
| boolean <u>dummy</u> | boolean <u>dummy</u> |
| loop forever | loop forever |
| p1:  think | q1:  forks[i] ⇐ true |
| p2:  forks[i] ⇒ dummy | q2:  forks[i] ⇒ dummy |
| p3:  forks[i⊕1] ⇒ dummy | q3: |
| p4:  eat | q4:     mutex? |
| p5:  forks[i] ⇐ true    (would *false* | q5:     deadlock? |
| p6:  forks[i⊕1] ⇐ true   be ok?) | q6: |

- Would it be enough to initialize each *forks[i] <= true* ?
  - Do you really need *forks[i] => dummy* in fork i? Why?

14.3.2011          Copyright Teemu Kerola 2011     Discuss     29

---

14.3.2011          Copyright Teemu Kerola 2011          30

## Rendezvous (1978, Abrial & Andrews)

- Synchronization with communication
  - No channels, usage similar to procedure calls
  - One (*accepting*) process waits for one of the (*calling*) processes
    - One request in service at a time   `asymmetric`
  - Calling process must know id of the accepting process
  - Accepting process does not need to know the id of calling process
  - May involve parameters and return value
- Good for client-server synchronization
  - Clients are calling processes   `server.service(parm, result)`
  - Server is accepting process   `accept service(p, r)`
  - Server is active process
  - Language construct, no mapping for real system nodes

14.3.2011          Copyright Teemu Kerola 2011          31

---

**Algorithm 8.6: Rendezvous**

| client | server |
|---|---|
| integer parm, result | integer p, r |
| loop forever | loop forever |
| p1:  parm ← . . . | q1: |
| p2:  server.service(parm, result) | q2:    accept service(p, r) |
| p3:  use(result) | q3:        r ← do the service(p) |

- Can have many similar clients
- Implementation with messages (e.g.)
  - Service request in one message
    - Arguments must be marshalled (make them suitable for transmission)
  - Wait until reply received
  - Reply result in another message

14.3.2011          Copyright Teemu Kerola 2011          32

---

## Guards in Rendezvous

- Additional constraint for accepting given service call
- Accept service call, if
  - Someone requests it and
  - Guard for that request type is true
    - Guard is based on local state
- If many such requests (with open guards) available, select one randomly
- Complete one request at a time
  - Implicit mutex

14.3.2011          Copyright Teemu Kerola 2011          33

---

## Ada Rendezvous

```
                                begin
                                loop
                                   select
Bounded Buffer in Ada              when Count < Index'Last =>
Export public ops defined             accept Append(I: in Integer ) do
before task body                         B(In_Ptr) := I;
task body Buffer is                   end Append;
  B: Buffer_Array;                     Count := Count − 1; In_Ptr := In_Ptr + 1;
  In_Ptr, Out_Ptr, Count: Index := 0;  or
                                   when Count > 0 =>
  ...                                accept Take(I: out Integer ) do
  Buffer.Append (456);                  I := B(Out_Ptr);
  Buffer.Append (333);              end Take;
  ...                               Count := Count − 1; Out_Ptr := Out_Ptr − 1;;

  ...                               or
  Buffer.Take(x);                   terminate;      Terminates when no
  Buffer.Take(y);             end select;           rendezvous processes
  ...                            end loop;          available? Tricky!
                                                    How to know?
How is buffer mutex problem solved?  end Buffer;    No concurrent operations!
```

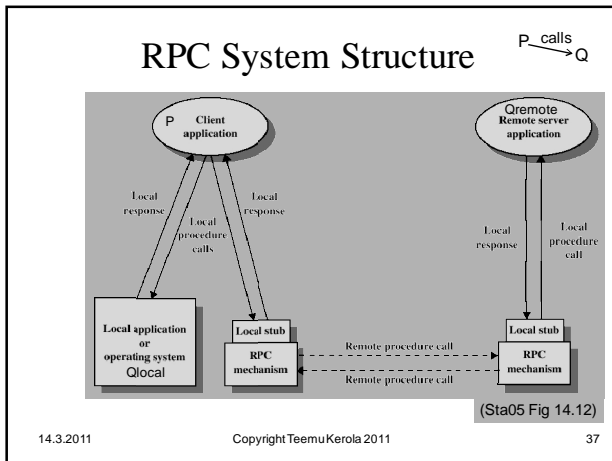14.3.2011          Copyright Teemu Kerola 2011          34

---

14.3.2011          Copyright Teemu Kerola 2011          35

---

## Remote Procedure Call

- Common operating system service for client-server model synchronization
  - Implemented with messages
  - Parameter marshalling
    - Semantics remain, implementation may change
  - Mutex problem
    - Combines monitor and synchronized messages?
      - Automatic mutex for service
    - Multiple calls active simultaneously?   `Usual case`
      - Mutex problems solved within called service
  - Semantics similar to ordinary procedure call
    - But no global environment (e.g., shared array)
  - Two-way synchronized communication channel
    - Client waits until service completed (usually)

14.3.2011          Copyright Teemu Kerola 2011          36

## RPC System Structure

P_calls
→Q



Local application or operating system Qlocal

Local stub

RPC mechanism

Local stub

RPC mechanism

Remote procedure call

(Sta05 Fig 14.12)

14.3.2011          Copyright Teemu Kerola 2011          37

## RPC Module

```
module mname
    op opname (formals)  [returns result]   Export public ops
body
    variable declarations;
    initialization code;
    proc opname (formal identifiers) returns result identifier
        declarations of local variables;
        statements
    end
    local procedures and processes;
end mname
```

Call:    `call mname.opname (arguments)`

14.3.2011          Copyright Teemu Kerola 2011          38

## RPC Example: Time Server

```
module TimeServer
    op get_time() returns int;   # retrieve time of day
    op delay(int interval);      # delay interval ticks
body
    int tod = 0;          # the time of day
    sem m = 1;            # mutual exclusion semaphore
    sem d[n] = ([n] 0);   # private delay semaphores
    queue of (int waketime, int process_id) napQ;
    ## when m == 1, tod < waketime for delayed processes

    proc get_time() returns time {
        time = tod;
    }

    proc delay(interval) {    # assume interval > 0
        int waketime = tod + interval;
        P(m);
        insert (waketime, myid) at appropriate place on napQ;
        V(m);
        P(d[myid]);   # wait to be awakened
    }
```

mutex

(And00 Fig 8.1)

(process Clock{} on next slide)

14.3.2011          Copyright Teemu Kerola 2011          39

```
process Clock {
    start hardware timer;
    while (true) {
        wait for interrupt, then restart hardware timer;
        tod = tod+1;
        P(m);
        while (tod >= smallest waketime on napQ) {
            remove (waketime, id) from napQ;
            V(d[id]);   # awaken process id
        }
        V(m);
    }
}
end TimeServer
```

- Internal process
  - Keeps the time
  - Wakes up delayed clients
- Service RPC's:     time = TimeServer.get_time();
                      TimeServer.delay(10);

Discuss

14.3.2011          Copyright Teemu Kerola 2011          40

---

Linux machine>> man rpc

RPC(3)                                              RPC(3)

NAME
    rpc - library routines for remote procedure calls

SYNOPSIS AND DESCRIPTION
    These routines allow C programs to make procedure calls on other machines across the network. First, the client calls a procedure to send a data packet to the server. Upon receipt of the packet, the server calls a dispatch routine to perform the requested service, and then sends back a reply. Finally, the procedure call returns to the client.

```
callrpc(host, prognum, versnum, procnum, inproc, in, outproc, out)
    char *host;           remote process
    u_long prognum, versnum, procnum;
    char *in, *out;                        decode/encode
    xdrproc_t inproc, outproc;             parameters/results
```

14.3.2011          Copyright Teemu Kerola 2011          41

## Remote Method Invocation (RMI)

```
package example.hello;
                                rmi server
import java.rmi.Remote;
import java.rmi.RemoteException;

public interface Hello extends Remote {
    String sayHello() throws RemoteException;
}
```
http://java.sun.com/j2se/1.5.0/docs/guide/rmi/hello/hello-world.html

- Java RPC
- Start rmiregistry        rmiregistry &      start rmiregistry
  - Stub lookup (default at port 1099)
- Start rmi server
  - Server runs until explicitly terminated by user
    `java -classpath classDir example.hello.Server &`
    `start java -classpath classDir example.hello.Server`

14.3.2011          Copyright Teemu Kerola 2011          42

```
package example.hello;
import java.rmi.registry.Registry;                               rmi server
import java.rmi.registry.LocateRegistry;
import java.rmi.RemoteException;
import java.rmi.server.UnicastRemoteObject;
public class Server implements Hello {
    public Server() {}
    public String sayHello() {
        return "Hello, world!"; }
    public static void main(String args[]) {
        try { Server obj = new Server();
            Hello stub = (Hello) UnicastRemoteObject.exportObject(obj, 0);
                    // Bind the remote object's stub in the registry
            Registry registry = LocateRegistry.getRegistry();
            registry.bind("Hello", stub);
            System.err.println("Server ready");
        } catch (Exception e) {
            System.err.println("Server exception: " + e.toString());
            e.printStackTrace();
        }
    }
}
        Output: Server ready
```

14.3.2011                    Copyright Teemu Kerola 2011                    43

```
package example.hello;
import java.rmi.registry.LocateRegistry;                         rmi client
import java.rmi.registry.Registry;

public class Client {
    private Client() {}
    public static void main(String[] args) {
        String host = (args.length < 1) ? null : args[0];
        try {
            Registry registry = LocateRegistry.getRegistry(host);
            Hello stub = (Hello) registry.lookup("Hello");
            String response = stub.sayHello();
            System.out.println("response: " + response);
        } catch (Exception e) {
            System.err.println("Client exception: " + e.toString());
            e.printStackTrace();
        }
    }
}
    Output:  response: Hello, world!
```

14.3.2011              Copyright Teemu Kerola 2011              Discussion 6    44

# Summary

- Distributed communication with messages
  - Synchronization and communication
  - Computation time + communication time = ?
- Higher level concepts
  - Guarded commands (theoretical background)
  - CSP (idea) & Occam (application)
  - Named Channels (ok without shared memory?)
  - Rendezvous
  - RPC & RMI (Java)

14.3.2011                    Copyright Teemu Kerola 2011                    45