

## System Software (2)

Peter Marwedel  
TU Dortmund, Informatik 12  
Germany

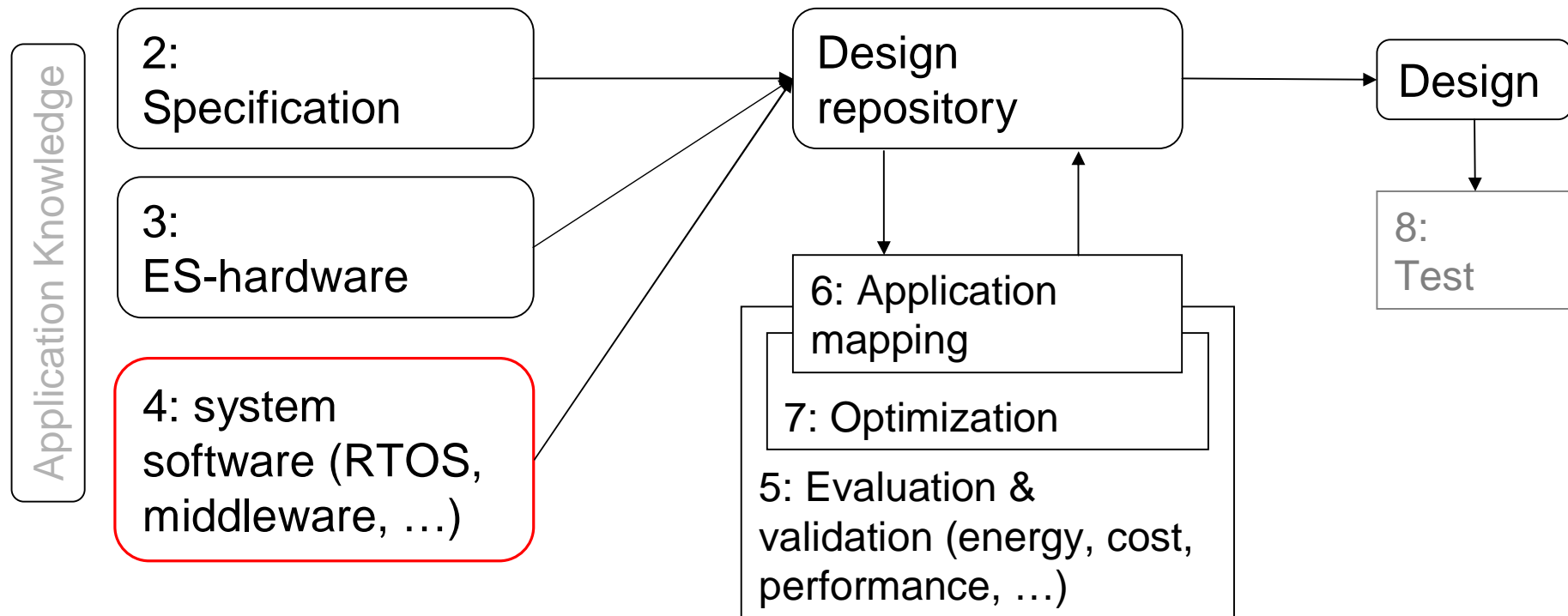


© Springer, 2010

2012年 11 月 28 日

These slides use Microsoft clip arts. Microsoft copyright restrictions apply.

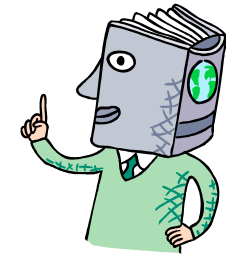
# Structure of this course



Numbers denote sequence of chapters

# Increasing design complexity + Stringent time-to-market requirements ➡ Reuse of components

Reuse requires knowledge from previous designs to be made available in the form of **intellectual property** (IP, for **SW & HW**).



- HW
- Operating systems
- ➡ ■ Middleware (Communication libraries, data bases, ...)
- ....

# Models of computation considered in this course

Communication/ local computations	Shared memory	Message passing	
		Synchronous	Asynchronous
Undefined components	Plain text, use cases   (Message) sequence charts		
Communicating finite state machines	StateCharts		SDL
Data flow	Scoreboarding + Tomasulo Algorithm (☞ Comp.Archict.)		Kahn networks, SDF
Petri nets		C/E nets, P/T nets, ...	
Discrete event (DE) model	VHDL*, Verilog*, SystemC*, ...	Only experimental systems, e.g. distributed DE in Ptolemy	
Imperative (Von Neumann) model	C, C++, Java [libraries]	C, C++, Java CSP, ADA	with libraries

\* Classification based on semantic model

---

# Pthreads

---

- **Shared memory model**
- Consists of standard API
  - Originally used for single processor
  - Locks ( mutex, read-write locks)

# PThreads Example

```
threads = (pthread_t *) malloc(n*sizeof(pthread_t));
pthread_attr_init(&pthread_custom_attr);
for (i=0;i<n; i++)
    pthread_create(&threads[i],
&pthread_custom_attr, task, ...);
for (i=0;i<n; i++) {
    pthread_mutex_lock(&mutex);
    <receive message>
    pthread_mutex_unlock(&mutex);
}
for (i=0;i<n; i++)
    pthread_join(threads[i], NULL);
```

```
void* task(void *arg) {
    ...
    pthread_mutex_lock(&mutex);
    <send message>
    pthread_mutex_unlock(&mutex);
    return NULL
}
```

---

# Pthreads

---

- Consists of standard API
  - Locks ( mutex, read-write locks)
  - Condition variables
  - Completely explicit synchronization
  - Synchronization is very hard to program correctly
- Typically supported by a mixture of hardware (shared memory) and software (thread management)
- Exact semantics depends on the memory consistency model
- Support for efficient producer/consumer parallelism relies on murky parts of the model
- Pthreads can be used as back-end for other programming models (e.g. OpenMP)

# OpenMP

---

Implementations target **shared memory** hardware

## Parallelism expressed using pragmas

- Parallel loops  
(`#pragma omp for {...}` ;focus: data parallelism)
- Parallel sections
- Reductions

## Explicit

- Expression of parallelism (mostly explicit)

## Implicit

- Computation partitioning
- Communication
- Synchronization
- Data distribution

Based on W. Verachtert (IMEC):  
*Introduction to Parallelism*,  
tutorial, DATE 2008

Lack of control over partitioning can cause problems



# Models of computation considered in this course

Communication/ local computations	Shared memory	Message passing	
		Synchronous	Asynchronous
Undefined components	Plain text, use cases   (Message) sequence charts		
Communicating finite state machines	StateCharts		SDL
Data flow	(Not useful) <sup>°</sup>		Kahn networks, SDF
Petri nets		C/E nets, P/T nets, ...	
Discrete event (DE) model	VHDL*, Verilog*, SystemC*, ...	Only experimental systems, e.g. distributed DE in Ptolemy	
Imperative (Von Neumann) model	C, C++, Java [libraries]	C, C++, Java CSP, ADA	with libraries

\* Classification based on semantic model

° Somewhat related: Scoreboarding + Tomasulo-Algorithm

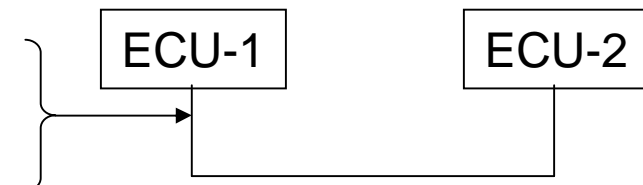
# OSEK/VDX COM

## OSEK/VDX COM

- is a special communication standard for the OSEK automotive OS Standard
- provides an “Interaction Layer” as an API for internal and external communication via a “Network Layer” and a “Data Link” layer (some requirements for these are specified)
- specifies the functionality, it is not an implementation.



© P. Marwedel, 2011



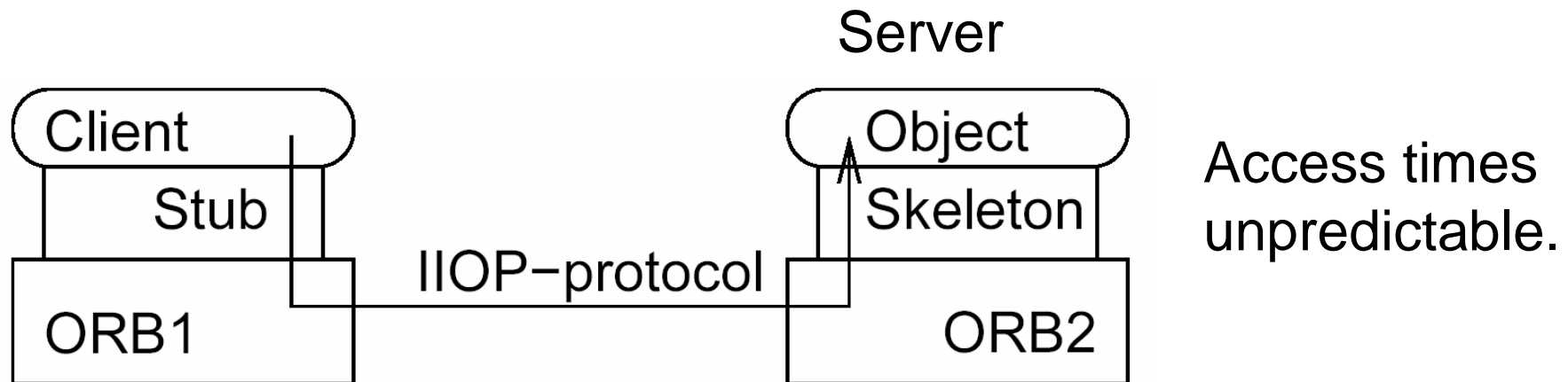
# CORBA

## (Common Object Request Broker Architecture)

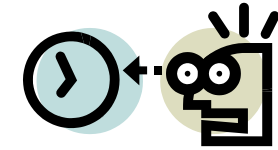
Software package for access to remote objects;

Information sent to Object Request Broker (ORB) via local stub.

ORB determines location to be accessed and sends information via the IIOP I/O protocol.

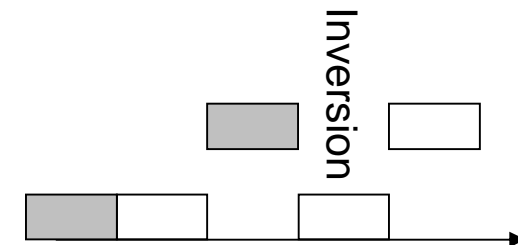
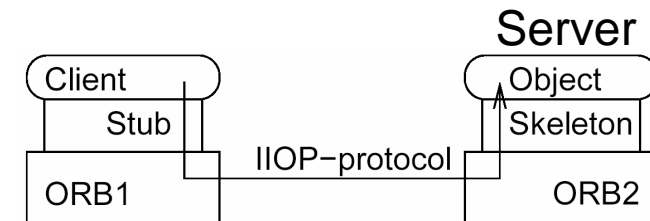


# Real-time (RT-) CORBA



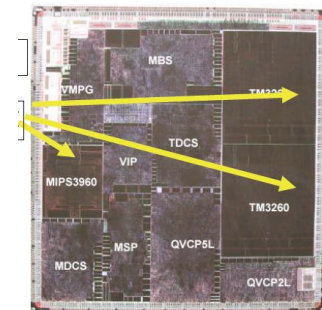
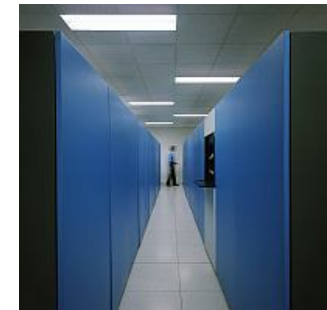
## RT-CORBA

- provides *end-to-end predictability of timeliness in a fixed priority system*.
- *respects thread priorities between client and server for resolving resource contention,*
- provides thread priority management,
- provides priority inheritance,
- bounds latencies of operation invocations,
- provides pools of preexisting threads.



# Message passing interface (MPI)

- Asynchronous/synchronous **message passing**
- Designed for high-performance computing
- Comprehensive, popular library
- Available on a variety of platforms
- Mostly for homogeneous multiprocessing
- Considered for MPSoC programs for ES;
- Includes many copy operations to memory (memory speed ~ communication speed for MPSoCs); Appropriate MPSoC programming tools missing.



[http://www.mhpcc.edu/training/workshop/mipi/MAIN.html#Getting\\_Started](http://www.mhpcc.edu/training/workshop/mipi/MAIN.html#Getting_Started)

# MPI (1)

---

## Sample blocking library call (for C):

- `MPI_Send(buffer, count, type, dest, tag, comm)` where
  - *buffer*: Address of data to be sent
  - *count*: number of data elements to be sent
  - *type*: data type of data to be sent (e.g. `MPI_CHAR`, `MPI_SHORT`, `MPI_INT`, ...)
  - *dest*: process id of target process
  - *tag*: message id (for sorting incoming messages)
  - *comm*: communication context = set of processes for which destination field is valid
  - function result indicates success

[http://www.mhpcc.edu/training/workshop/mpi/MAIN.html#Getting\\_Started](http://www.mhpcc.edu/training/workshop/mpi/MAIN.html#Getting_Started)

---

# MPI (2)

---

## Sample non-blocking library call (for C):

- `MPI_Isend(buffer, count, type, dest, tag, comm, request)`  
where
  - *buffer ... comm*: same as above
  - *request*: unique "request number". "handle" can be used (in a WAIT type routine) to determine completion

[http://www.mhpcc.edu/training/workshop/mpi/MAIN.html#Getting\\_Started](http://www.mhpcc.edu/training/workshop/mpi/MAIN.html#Getting_Started)

---

# Evaluation

---

## Explicit

- Computation partitioning
- Communication
- Data distribution

## Implicit

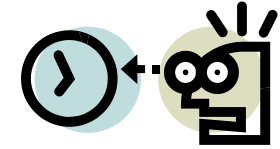
- Synchronization (implied by communic., explicit possible)
- Expression of parallelism (implied)
- Communication mapping

## Properties

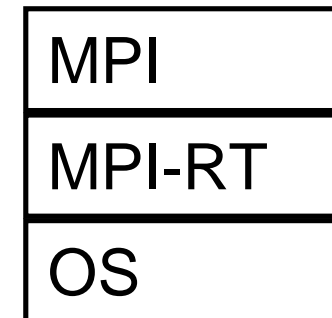
- Most things are explicit
- Lots of work for the user (*“assembly lang. for parallel prog.”*)
- doesn't scale well when # of processors is changed heavily



# RT-issues for MPI



- MPI/RT: a real-time version of MPI [MPI/RT forum, 2001].
- MPI-RT does not cover issues such as thread creation and termination.
- MPI/RT is conceived as a potential layer between the operating system and standard (non real-time) MPI.



# Universal Plug-and-Play (UPnP)

- Extension of the plug-and-play concept
- *Enable emergence of easily connected devices & simplify implementation of networks @ home & corporate environments!*
- Examples: Discover printers, storage space, control switches in homes & offices
- **Exchanging data**, no code (reduces security hazards)
- Agreement on data formats & protocols
- Classes of predefined devices (printer, mediaserver etc.)
- <http://upnp.org>



© P. Marwedel, 2012

# Devices Profile for Web Services (DPWS)

- More general than UPnP
- ... *DPWS defines a minimal set of implementation constraints to enable secure Web Service messaging, discovery, description, and eventing on resource-constrained devices.*  
...
- *DPWS specifies a set of built-in services:*
  - *Discovery services ...*
  - *Metadata exchange services...*
  - *Publish/subscribe eventing services...*
- *Lightweight protocol, supporting dynamic discovery, ... its application to automation environments is clear.*



# Network Communication Protocols

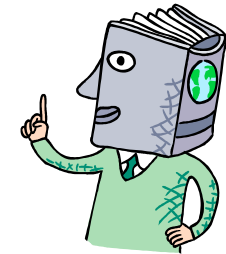
## - e.g. JXTA -

- *Open source peer-to-peer protocol specification.*
- *Defined as a set of XML messages that allow any device connected to a network to exchange messages and collaborate independently of the network topology.*
- *.. Can be implemented in any modern computer language.*
- *JXTA peers create a virtual overlay network, allowing a peer to interact with other peers even when some of the peers and resources are behind firewalls and NATs or use different network transports.*



# Increasing design complexity + Stringent time-to-market requirements ➡ Reuse of components

Reuse requires knowledge from previous designs to be made available in the form of **intellectual property** (IP, for **SW & HW**).

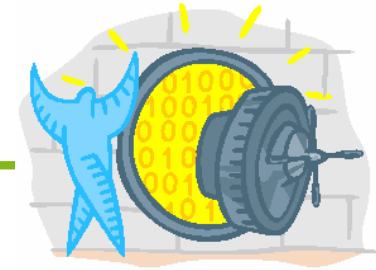


- HW
- Operating systems
- Middleware (Communication libraries, data bases, ...)
- ....



# Data bases

---



**Goal:** store and retrieve persistent information

Transaction= sequence of read and write operations

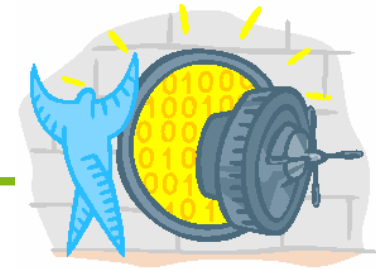
Changes not final until they are committed

Requested (“ACID”) properties of transactions



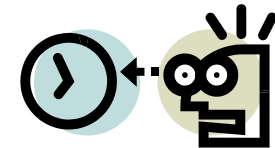
1. **Atomic:** state information as if transaction is either completed or had no effect at all.
2. **Consistent:** Set of values retrieved from several accesses to the data base must be possible in the world modeled.
3. **Isolation:** No user should see intermediate states of transactions
4. **Durability:** results of transactions should be persistent.

# Real-time data bases



Problems with implementing real-time data bases:

1. transactions may be aborted various times before they are finally committed.
2. For hard discs, the access times to discs are hardly predictable.



Possible solutions:

1. Main memory data bases
2. Relax ACID requirements

---

# Summary

---

- Communication middleware
  - Pthreads
  - OpenMP
  - OSEK/VDX COM
  - CORBA
  - MPI
  - JXTA
  - DPWS
- RT-Data bases (brief)