



TECHNISCHE
UNIVERSITÄT
DARMSTADT

 Parallel
 Programming

- Parallel Programming - Processes & Threads -

Prof. Dr. Felix Wolf
Technical University of Darmstadt

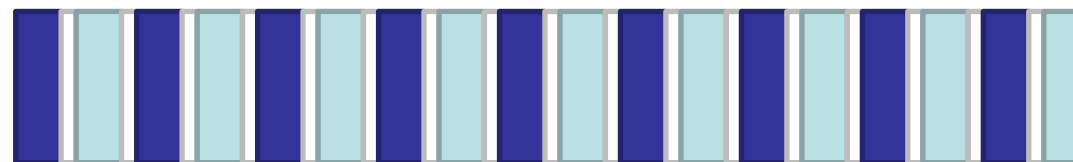
Outline

- Processes vs. threads
- On-chip multithreading

Time sharing

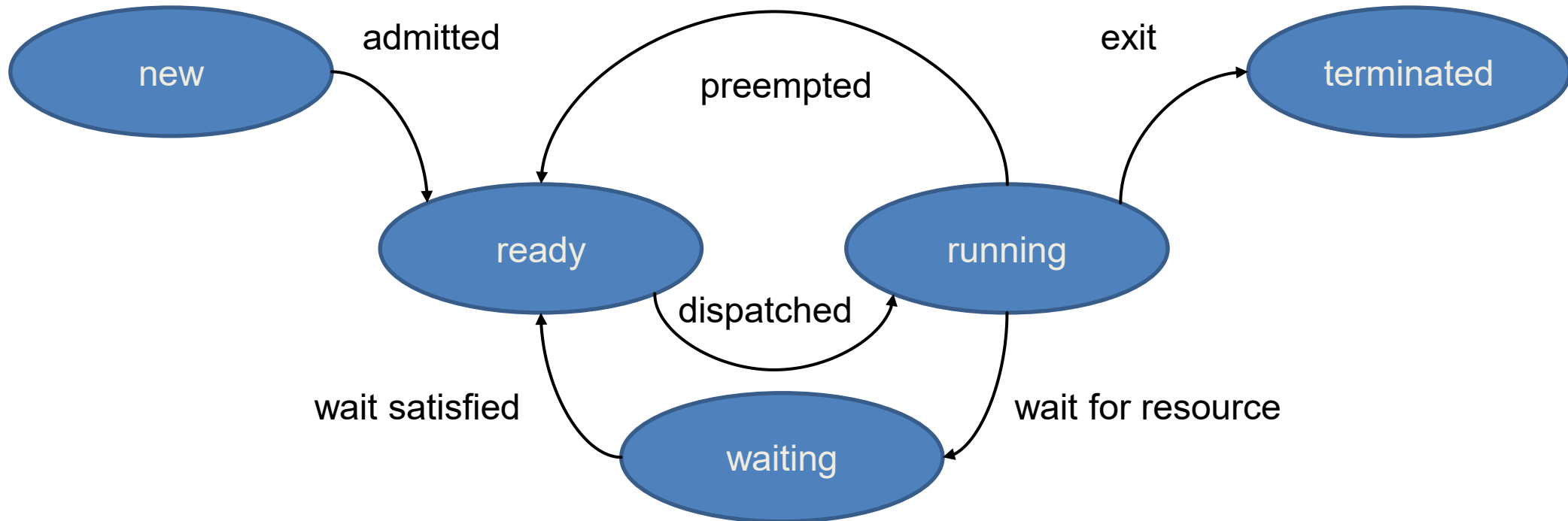
- Also called **multitasking**
- Logical extension of multiprogramming
- **CPU** executes multiple processes by switching among them
- Switches occur frequently enough so that users can interact with each program while it is running
- On a multiprocessor, processes can also run truly concurrently, taking advantage of additional processors

Single core

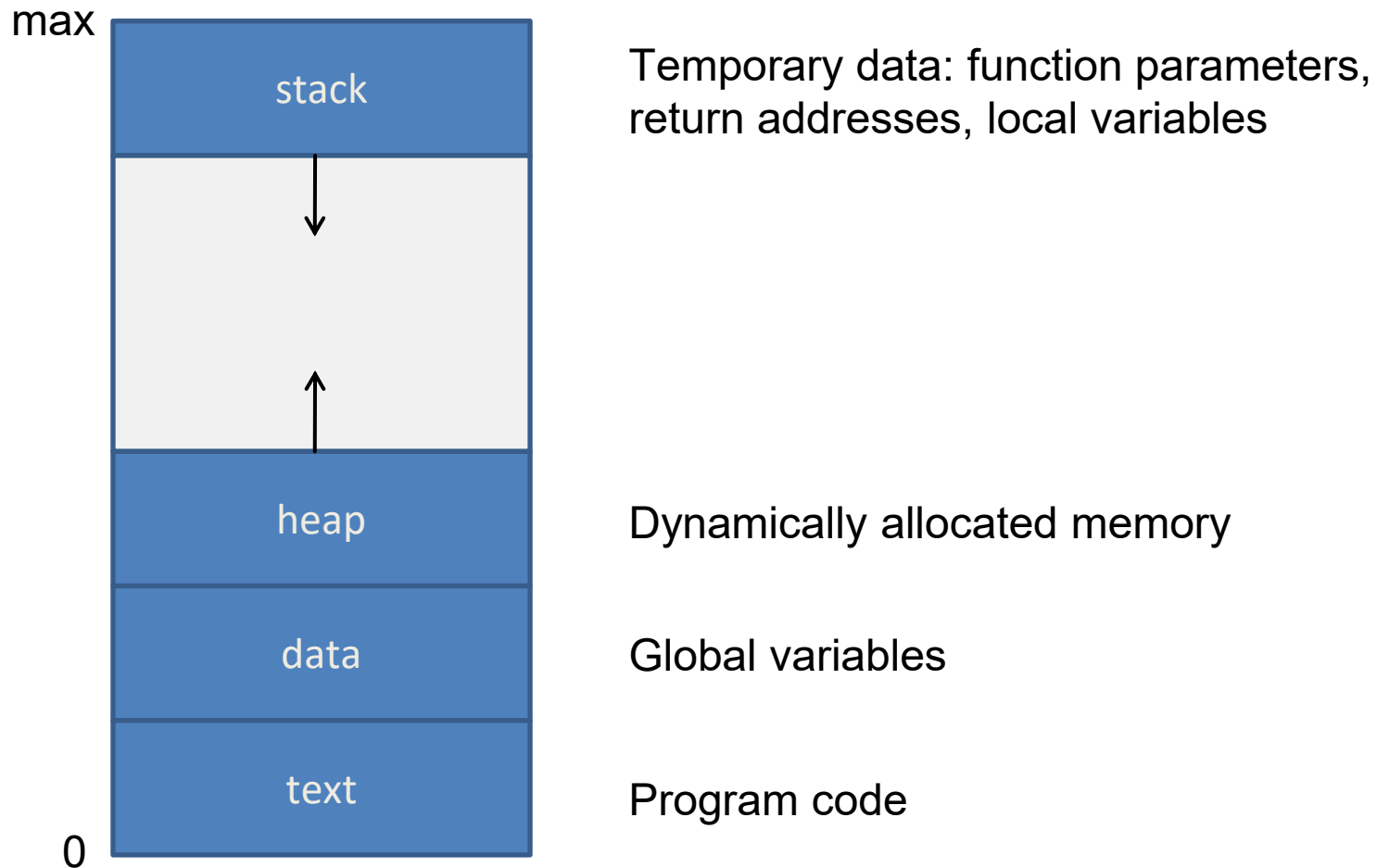


Process

- A process is a program in execution
- States of a process



Processes in memory



Thread

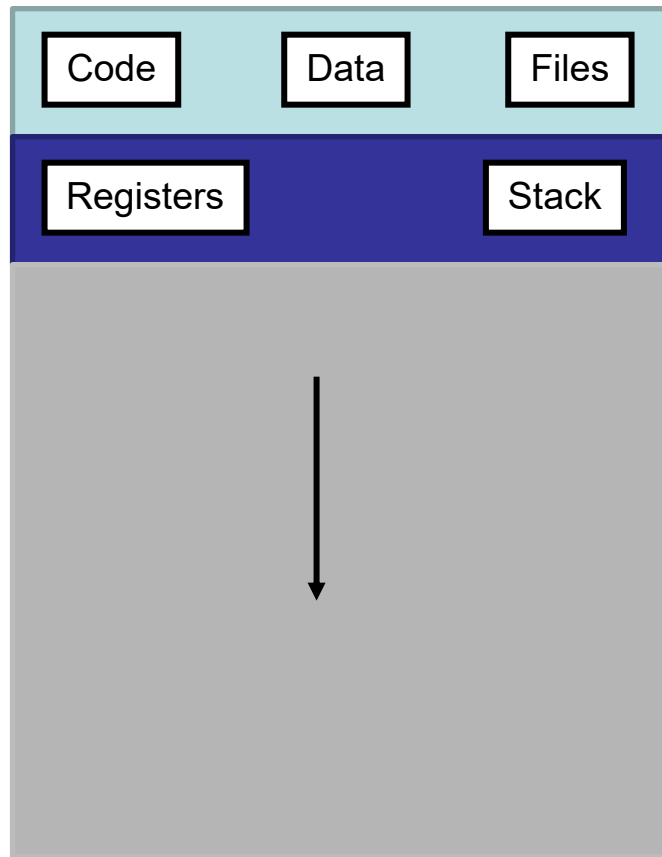
- Basic unit of CPU utilization
 - Flow of control within a process
- A thread includes
 - Thread ID
 - Program counter
 - Register set
 - Stack
- Shares resources with other threads belonging to the same process
 - Text (i.e., code) section
 - Data section (i.e., address space)
 - Other operating system resources
 - E.g., open files, signals

Single-threaded vs. multi-threaded

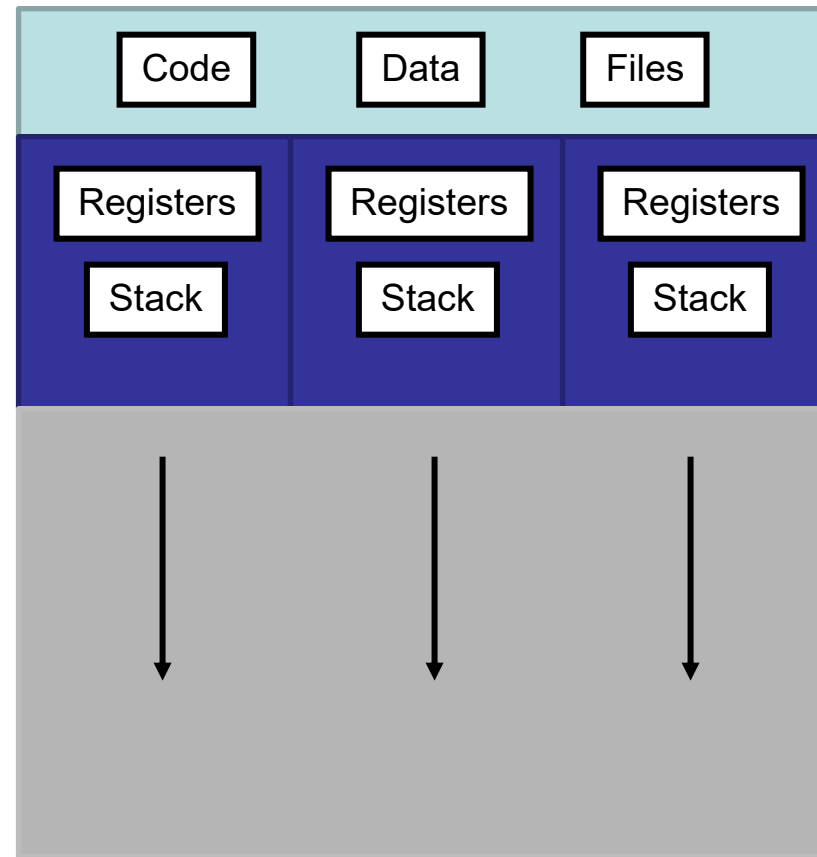


TECHNISCHE
UNIVERSITÄT
DARMSTADT

Parallel
Programming



Single-threaded



Multi-threaded

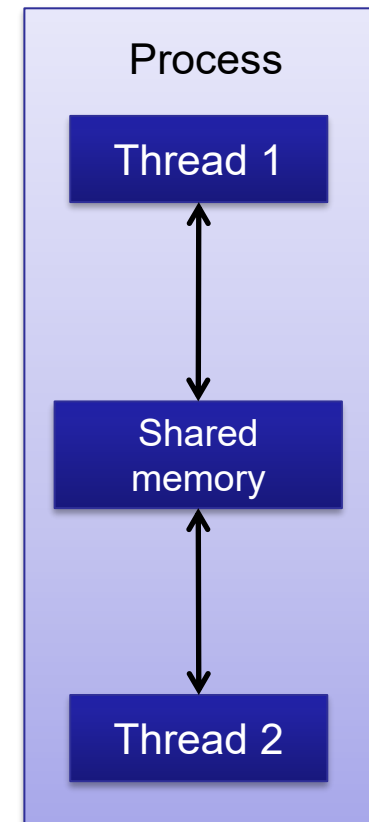
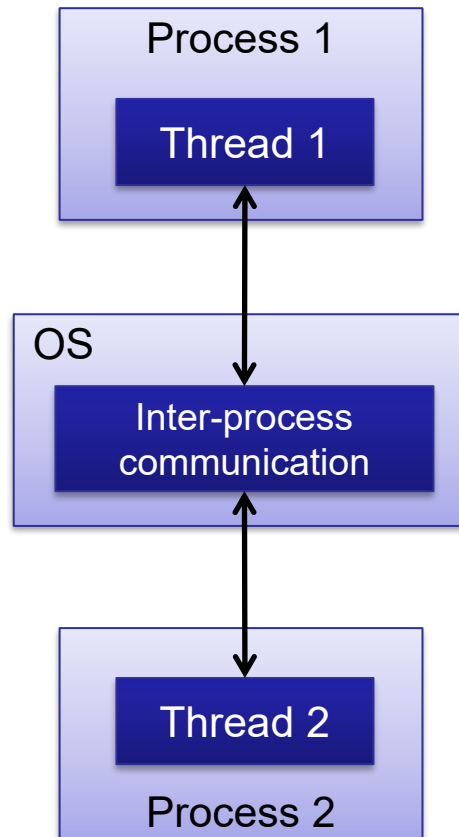
Processes vs. threads

Communication



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Parallel
Programming



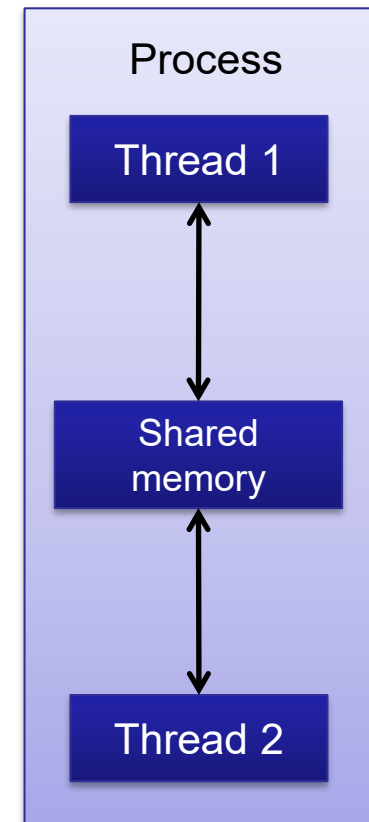
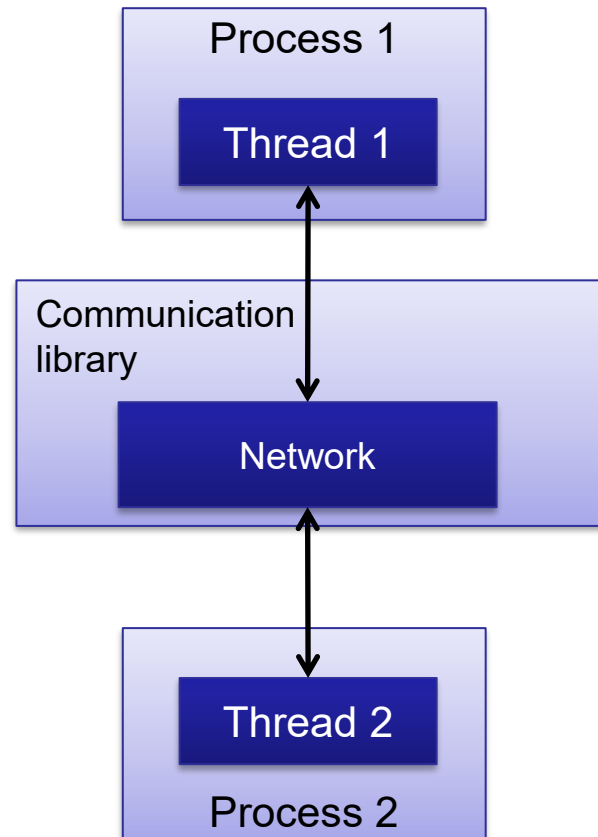
Processes vs. threads

Communication (2)



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Parallel
Programming



Processes vs. threads

Communication (3)

Processes

- Communication explicit
- Often requires replication of data
- Address spaces protected
- Parallelization usually implies profound redesign
- Writing debuggers is harder
- More scalable

Threads

- Convenient communication via shared variables
- More space efficient – sharing of code and data
- Context switch cheaper
- Incremental parallelization easier
- Harder to debug – race conditions



On-chip multithreading

- All modern, pipelined CPUs suffer from the following problem
 - When a memory reference misses L1 or L2 caches, it takes a long time until the requested word is loaded into the cache
- On-chip multithreading allows the CPU to manage multiple threads to mask these stalls
- If one thread is stalled, the CPU can run another thread and keep the hardware busy
- Four hardware threads per core often sufficient to hide latency

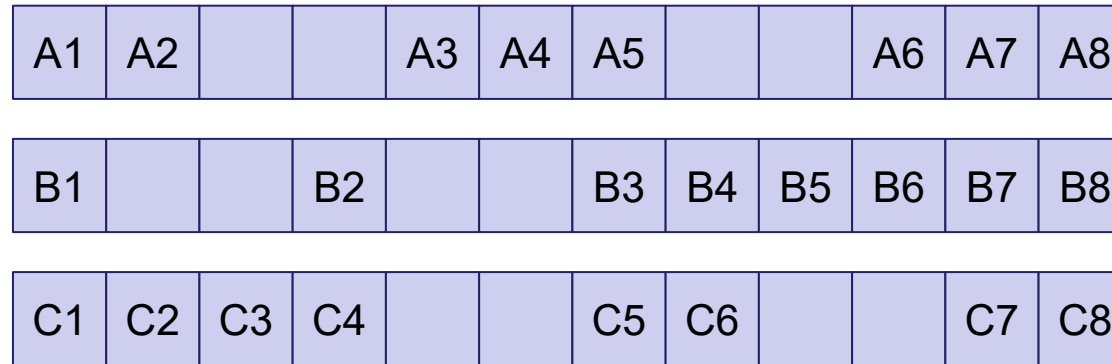
Fine-grained vs. coarse-grained multithreading



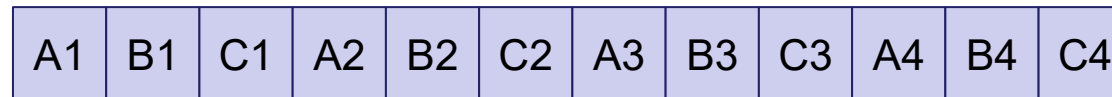
TECHNISCHE
UNIVERSITÄT
DARMSTADT

Parallel
Programming

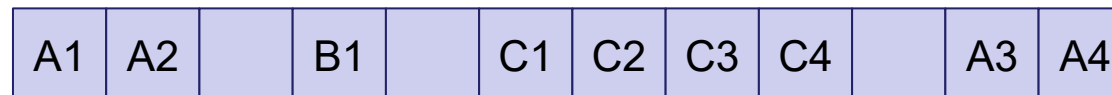
3 threads



Fine grained – threads
run round-robin



Coarse grained – switch
occurs only upon stall



Switching occurs in
regular intervals, e.g.,
at the boundary of an
instruction cycle.

Five ways of improving CPU performance

- Increase the clock speed
- Put two cores on a chip
- Add functional units
- Make pipeline longer
- Use on-chip multithreading

On-chip multithreading support can improve performance over-proportionally in comparison to the required extra chip area



Hyperthreading in the Intel Core i7

- Two threads (or processes) can run at once on the same core
- Looks from far like a dual processor in which both CPUs share a common cache and main memory
- However, many hardware resources are shared between threads
- Advantage – enables true concurrency within the same core
- Disadvantage – resource contention (e.g., cache) may lower throughput

Summary

- A process is a program in execution
- Threads are light-weight processes that can share code and a global address space
 - Advantages – responsiveness, utilization of multiprocessors
 - Disadvantages – synchronization overhead, programming complexity
- Hardware threads can hide latency
- Hyperthreading can boost performance