# Microkernel Construction

## Introduction

SS2016

Hermann Härtig
Benjamin Engel
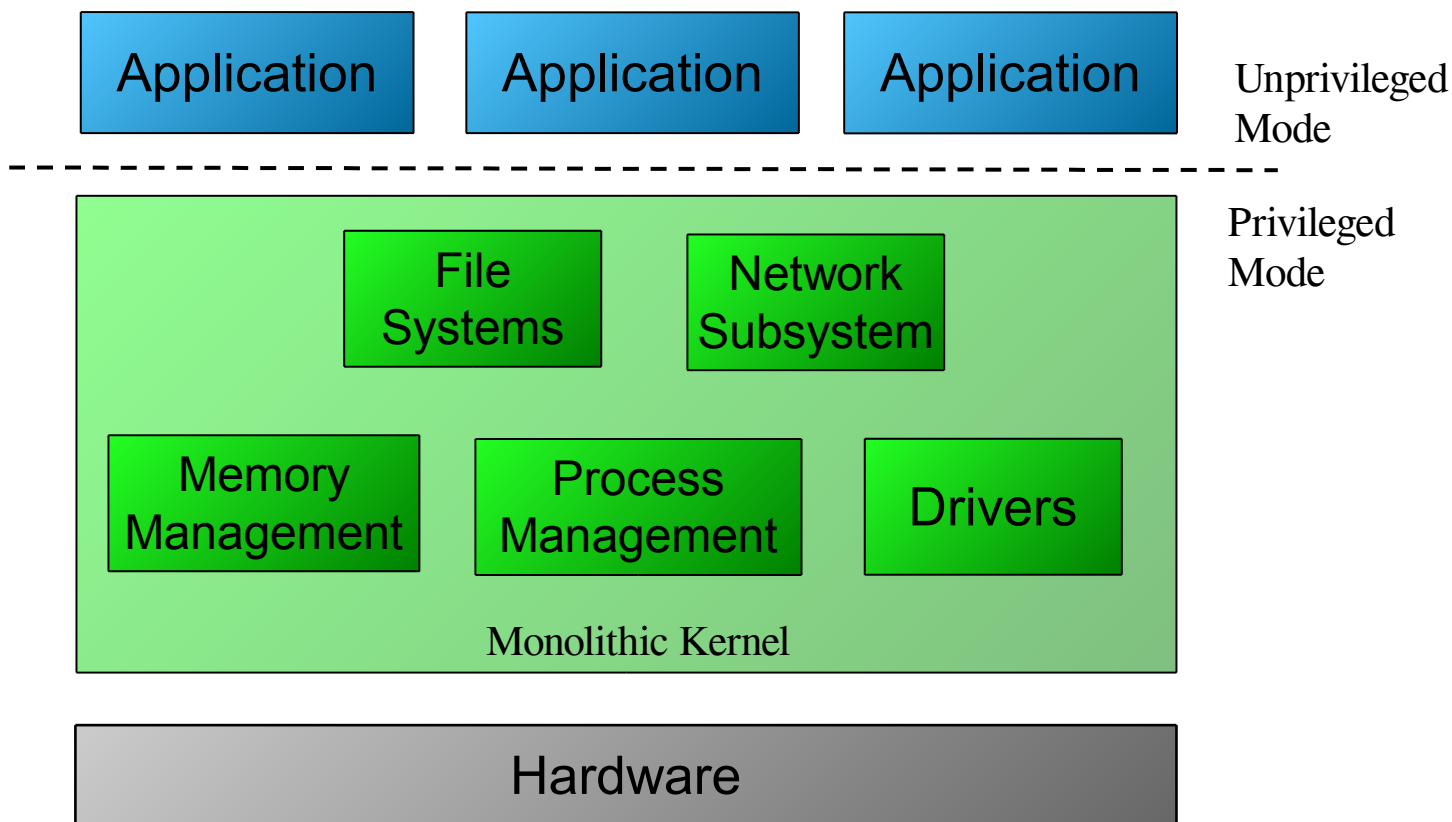**Tobias Stumpf**

TU Dresden
Operating
Systems Group

# Class Goals

- Provide deeper understanding of OS mechanisms

- Introduce L4 principles and concepts

- Make you become enthusiastic L4 hackers

- Propaganda for OS research at TU Dresden

Hermann
Härtig
Benjamin
Engel
**Tobias
Stumpf**

TU Dresden
Operating
Systems Group

# Administration

- Thursday, 4<sup>th</sup> DS, 2 SWS

- Slides: http://www.tudos.org → Teaching →
  Microkernel Construction

- Subscribe to our mailing list:
  http://www.tudos.org/mailman/listinfo/mkc2016

- In winter term:
  - Construction of Microkernel-based Systems (2 SWS)
  - Various Labs

Hermann
Härtig
Benjamin
Engel
**Tobias
Stumpf**

TU Dresden
Operating
Systems Group

# „Monolithic" Kernel System Design

| Application | Application | Application | Unprivileged Mode |
|---|---|---|---|

Privileged Mode

Monolithic Kernel:
- File Systems
- Network Subsystem
- Memory Management
- Process Management
- Drivers

Hardware

Hermann
Härtig
Benjamin
Engel
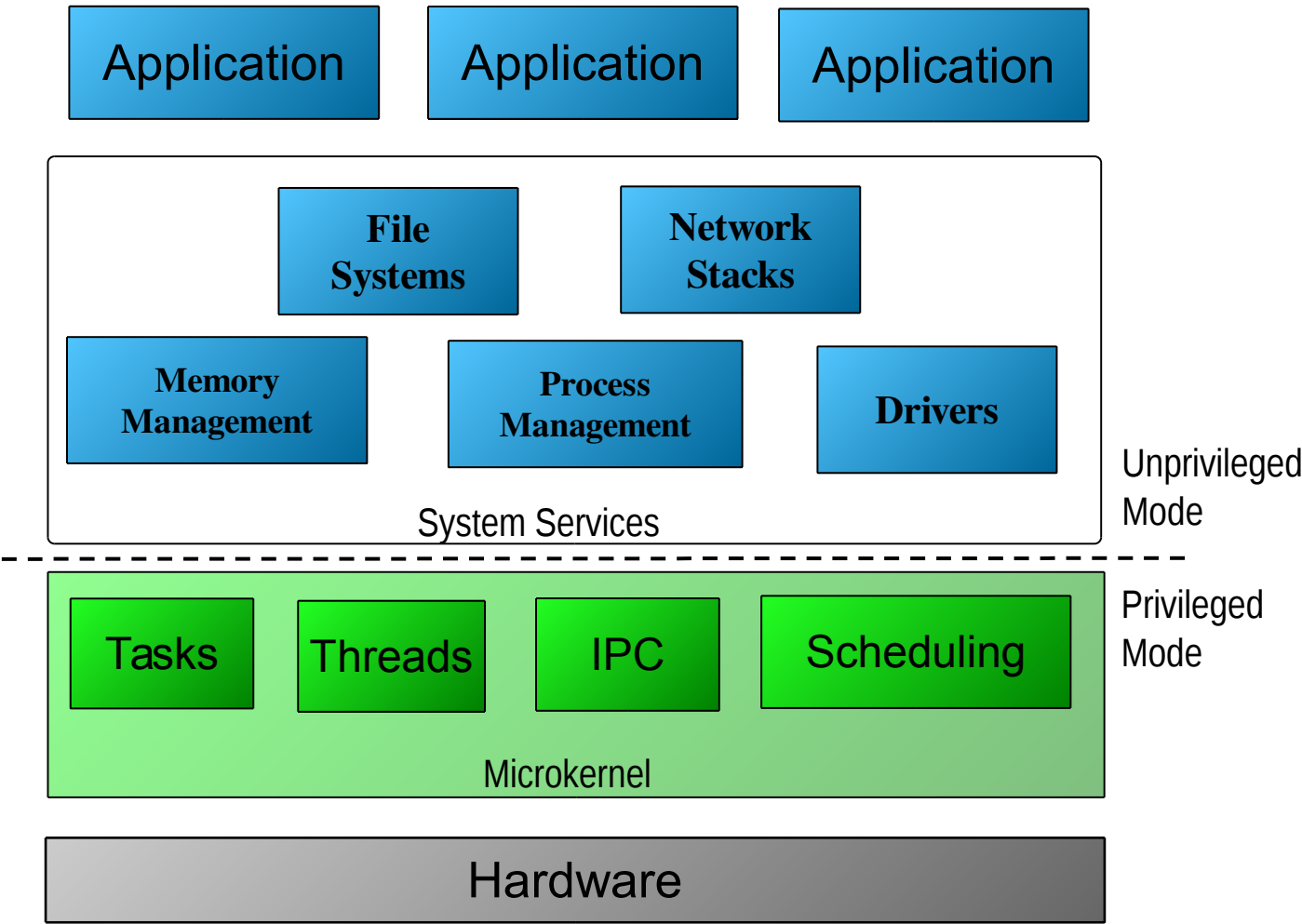**Tobias**
**Stumpf**

TU Dresden
Operating
Systems Group

# Monolithic Kernel OS (Propaganda)

- **System components run in privileged mode**

➜ No protection between system components

 – Faulty driver can crash the whole system

 – More than 2/3 of today's OS code are drivers

➜ No need for good system design

 – Direct access to data structures

 – Undocumented and frequently changing interfaces

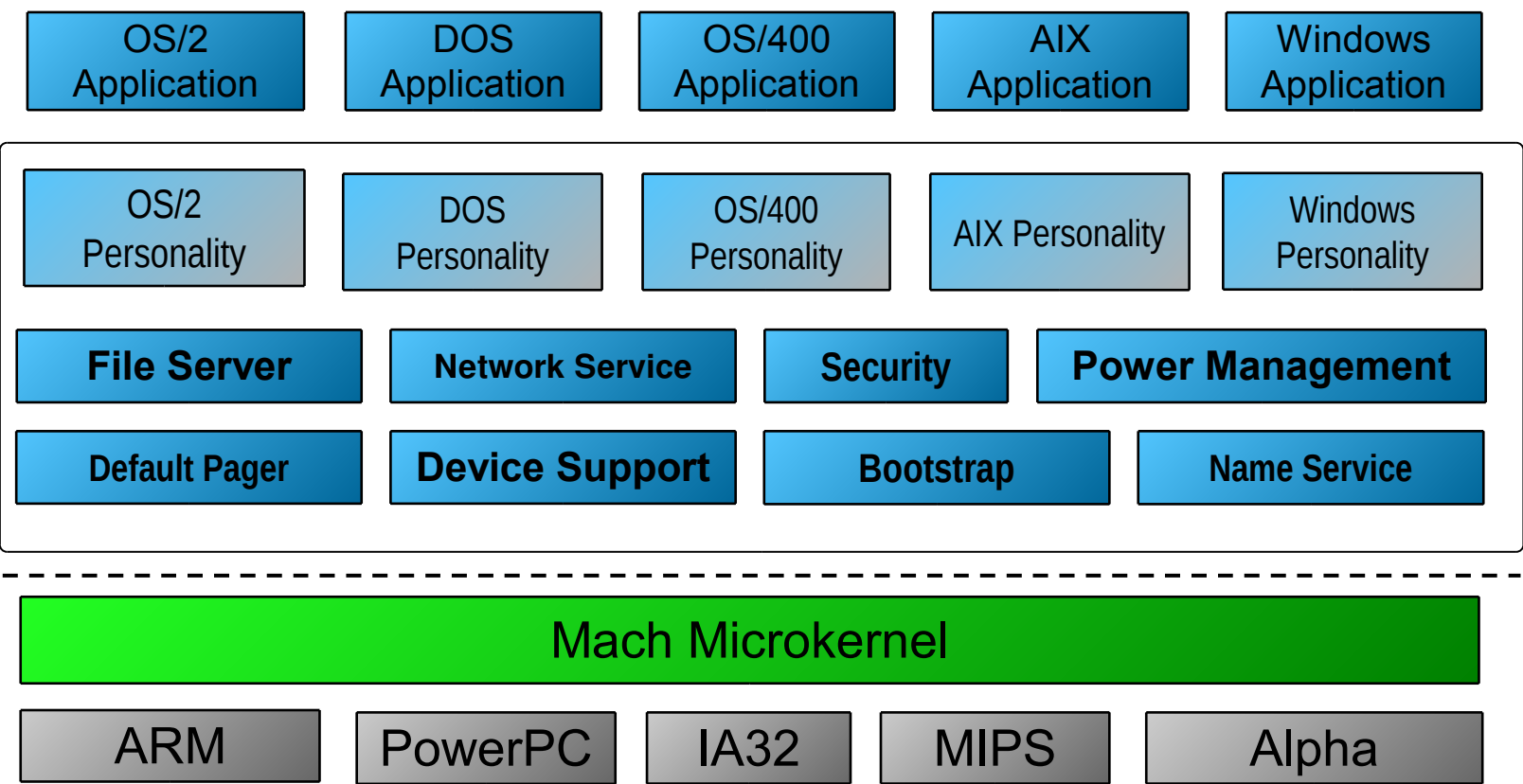➜ Big and inflexible

 – Difficult to replace system components

Why something different?

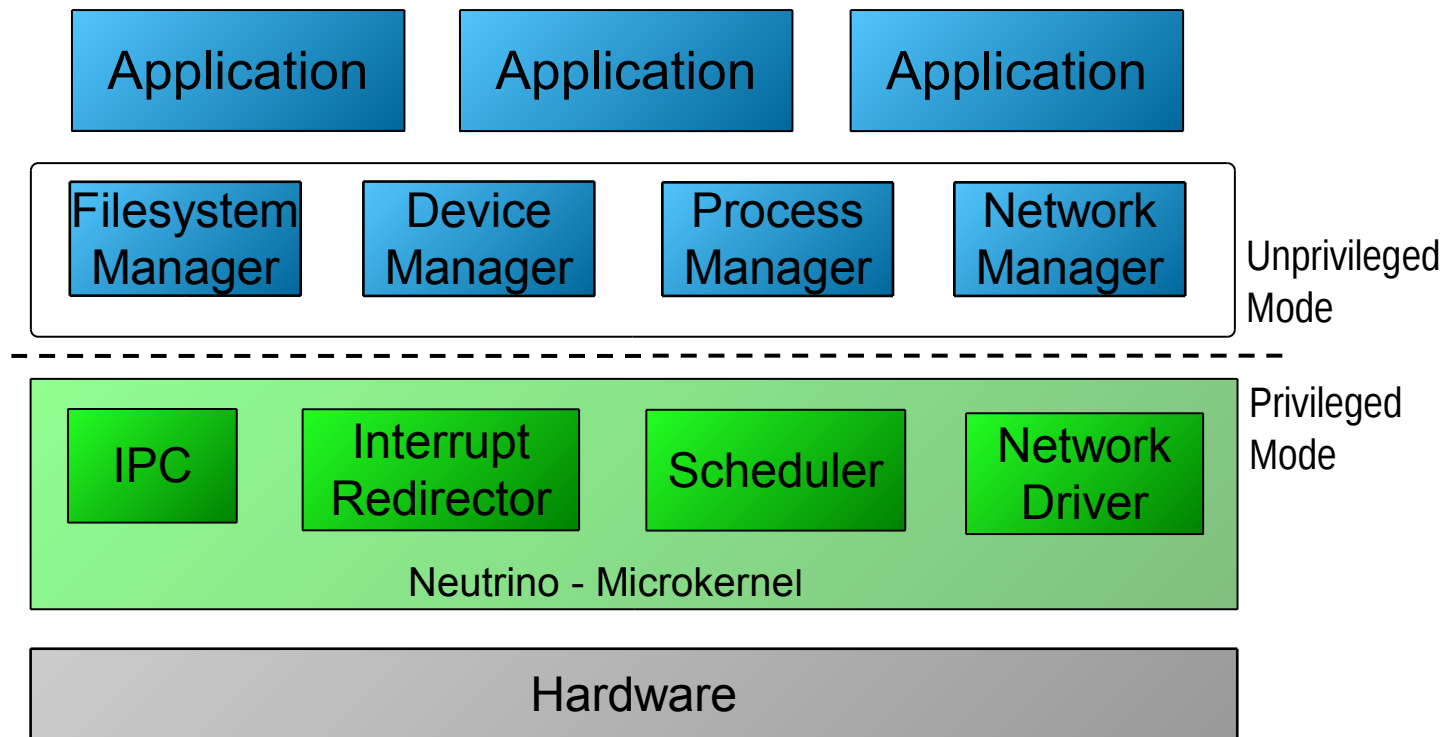- More and more difficult to manage increasing OS complexity

Hermann
Härtig
Benjamin
Engel
**Tobias
Stumpf**

TU Dresden
Operating
Systems Group

# Microkernel System Design



Hermann
Härtig
Benjamin
Engel
**Tobias
Stumpf**

TU Dresden
Operating
Systems Group

# Example – IBM Workplace OS / Mach

| OS/2 Application | DOS Application | OS/400 Application | AIX Application | Windows Application |
|---|---|---|---|---|

| OS/2 Personality | DOS Personality | OS/400 Personality | AIX Personality | Windows Personality |
|---|---|---|---|---|

**File Server**  **Network Service**  **Security**  **Power Management**

**Default Pager**  **Device Support**  **Bootstrap**  **Name Service**

## Mach Microkernel

| ARM | PowerPC | IA32 | MIPS | Alpha |
|---|---|---|---|---|

Hermann
Härtig
Benjamin
Engel
**Tobias
Stumpf**

TU Dresden
Operating
Systems Group

7

# Example – QNX / Neutrino

- Embedded systems
- Message passing system (IPC)
- Network transparency

Hermann
Härtig
Benjamin
Engel
**Tobias**
**Stumpf**

TU Dresden
Operating
Systems Group

# Example – Fiasco.OC

Hermann
Härtig
Benjamin
Engel
**Tobias**
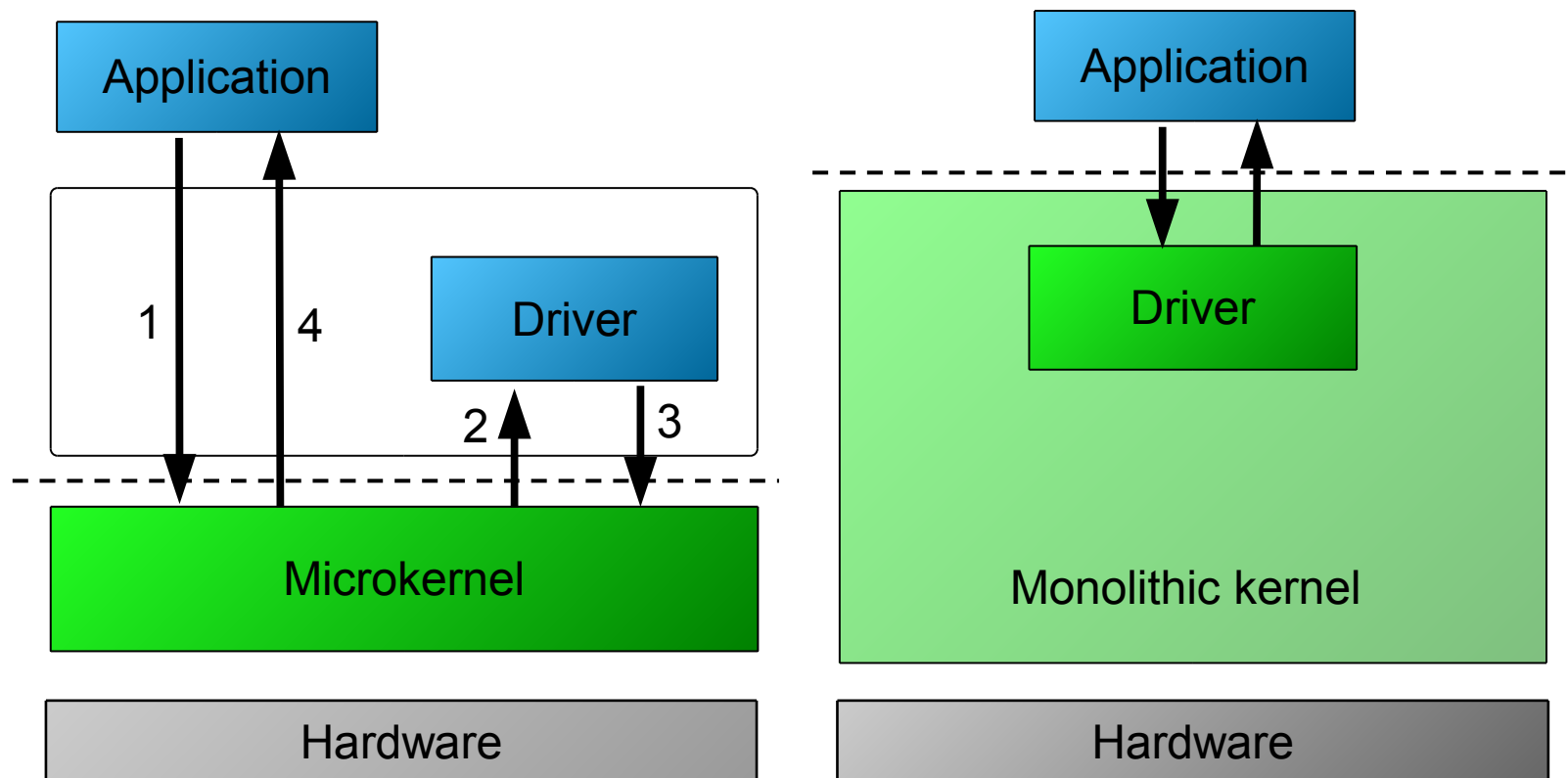**Stumpf**

TU Dresden
Operating
Systems Group

9

# Visions vs. Reality

- Flexibility and Customizable
  - Monolithic kernels are modular

- Maintainability and complexity
  - Monolithic kernel have layered architecture

- ✔ Robustness
  - Microkernels are superior due to isolated system components
  - Trusted code size (i386)
    - Fiasco kernel: about 20.000 loc
    - Linux kernel: > 1.000.000 loc (without drivers, arch, fs)

- ✗ Performance
  - Application performance degraded
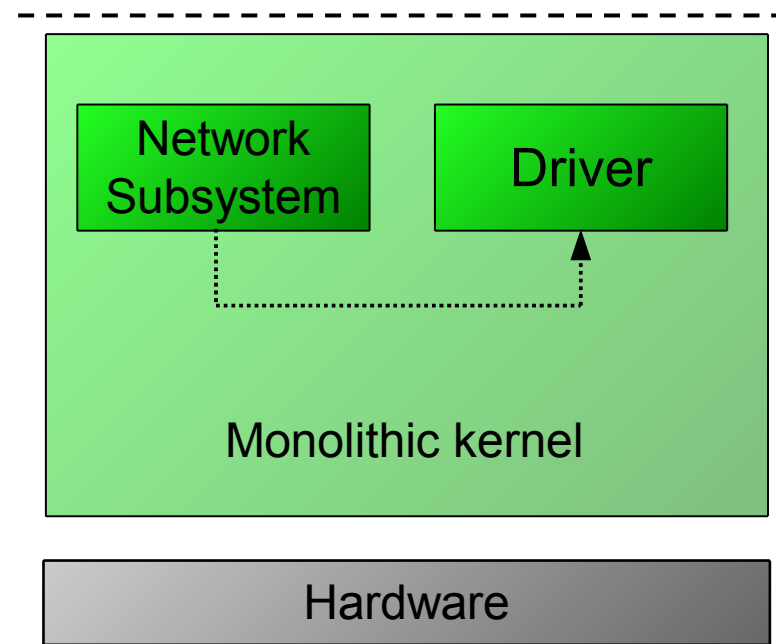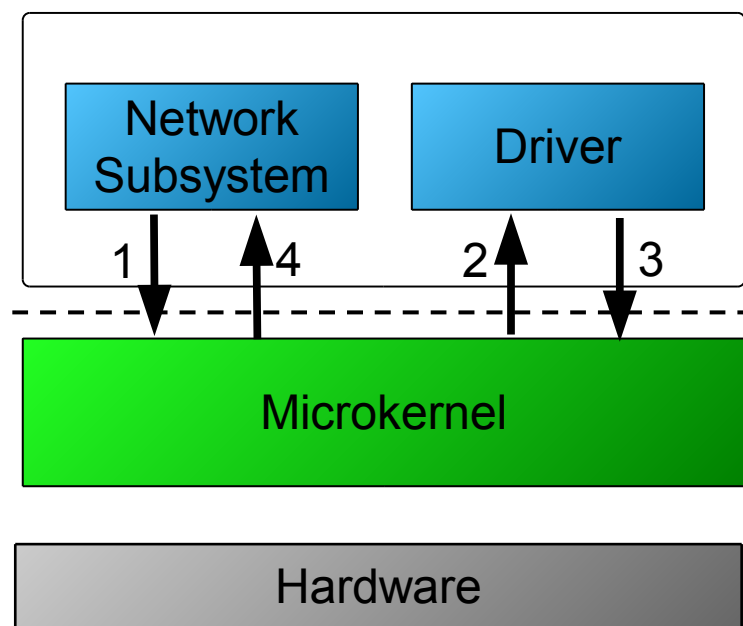  - Communication overhead (see next slides)

Hermann
Härtig
Benjamin
Engel
**Tobias
Stumpf**

TU Dresden
Operating
Systems Group

10

# Robustness vs. Performance (1)

- ## System calls
  - Monolithic kernel: 2 kernel entries/exits
  - Microkernel: 4 kernel entries/exits + 2 context switches



Hermann
Härtig
Benjamin
Engel
**Tobias**
**Stumpf**

TU Dresden
Operating
Systems Group

# Robustness vs. Performance (2)

- Calls between system services
  - Monolithic kernel: 1 function call
  - Microkernel: 4 kernel entries/exits + 2 context switches

Hermann
Härtig
Benjamin
Engel
**Tobias**
**Stumpf**

TU Dresden
Operating
Systems Group

# Challenges

- **Build functional powerful and fast microkernels**
    - Provide abstractions and mechanisms
    - Fast communication primitive (IPC)
    - Fast context switches and kernel entries/exits
- ➜ *Subject of this lecture*

- **Build efficient OS services**
    - Memory Management
    - Synchronization
    - Device Drivers
    - File Systems
    - Communication Interfaces
- ➜ *Subject of lecture "Construction of Microkernel-based systems" (in winter term)*

Hermann
Härtig
Benjamin
Engel
**Tobias
Stumpf**

TU Dresden
Operating
Systems Group

# L4 Microkernel Family

- Originally developed by Jochen Liedtke
  (GMD / IBM Research)

- Current development:
  - UNSW/NICTA/OKLABS: OKL4, SEL4, L4Verified
  - TU Dresden: Fiasco.OC, M3, Nova

- Support for hardware architectures:
  - **X86, ARM, ...**

Hermann
Härtig
Benjamin
Engel
**Tobias
Stumpf**

TU Dresden
Operating
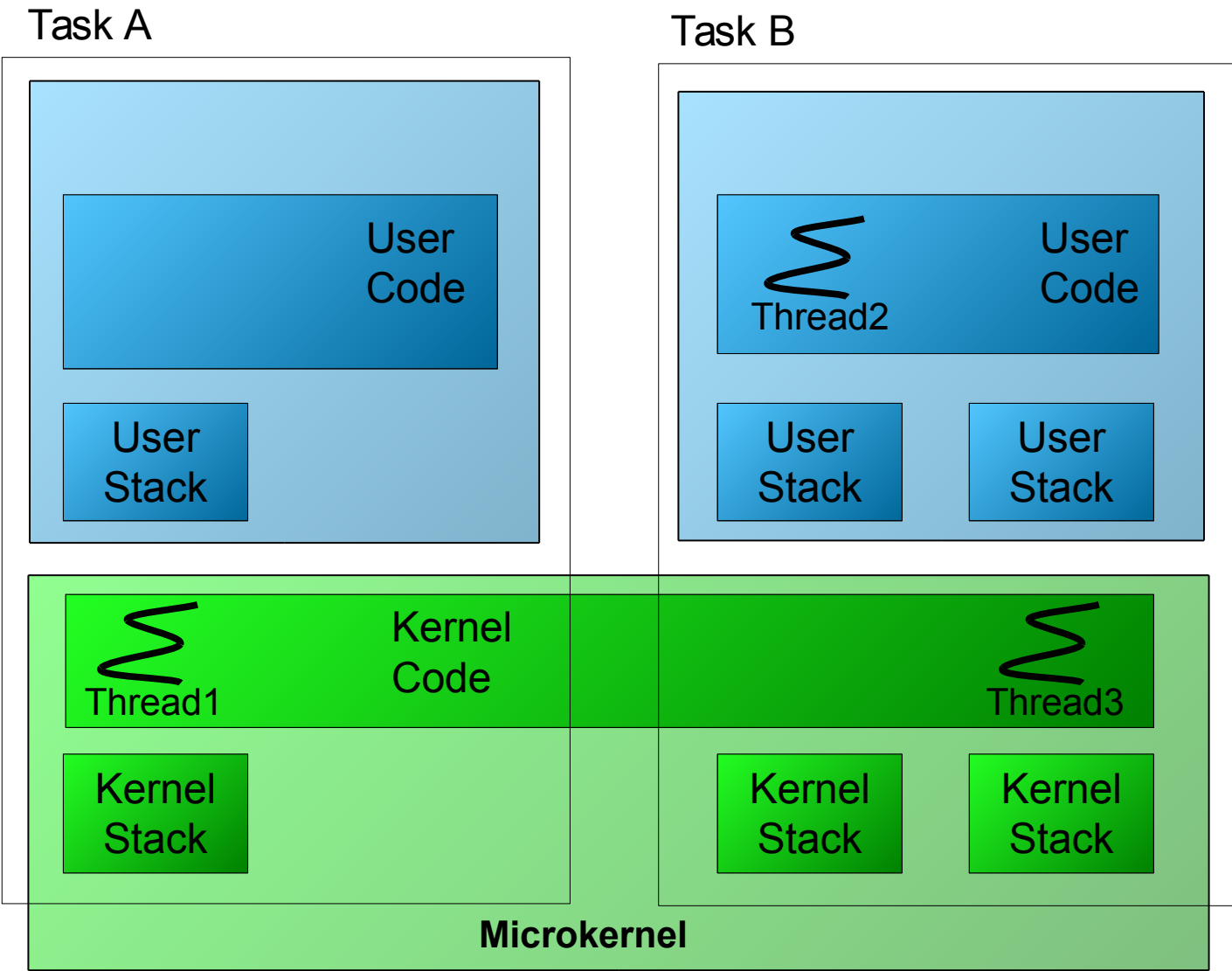Systems Group

# More Microkernels (Incomplete list)

- ## Commercial kernels
  - Singularity @ Microsoft Research
  - K42 @ IBM Research
  - velOSity/INTEGRITY @ Green Hills Software
  - Chorus/ChorusOS @ Sun Microsystems
  - PikeOS @ SYSGO AG

- ## Research kernels
  - EROS/CoyotOS @ John Hopkins University
  - Minix @ FU Amsterdam
  - Amoeba @ FU Amsterdam
  - Pebble @ Bell Labs
  - Grasshopper @ University of Sterling
  - Flux/Fluke @ University of Utah
  - Pistachio @ KIT

Hermann
Härtig
Benjamin
Engel
**Tobias
Stumpf**

TU Dresden
Operating
Systems Group

# L4 - Concepts

- Jochen Liedtke: "A microkernel does no real work"
  - Kernel provides only inevitable mechanisms
  - No policies implemented in the kernel
- Abstractions
  - Tasks with address spaces
  - Threads executing programs/code
- Mechanisms
  - Resource access control
  - Scheduling
  - Communication (IPC)

Hermann
Härtig
Benjamin
Engel
**Tobias
Stumpf**

TU Dresden
Operating
Systems Group

# Threads and Tasks

Hermann
Härtig
Benjamin
Engel
**Tobias**
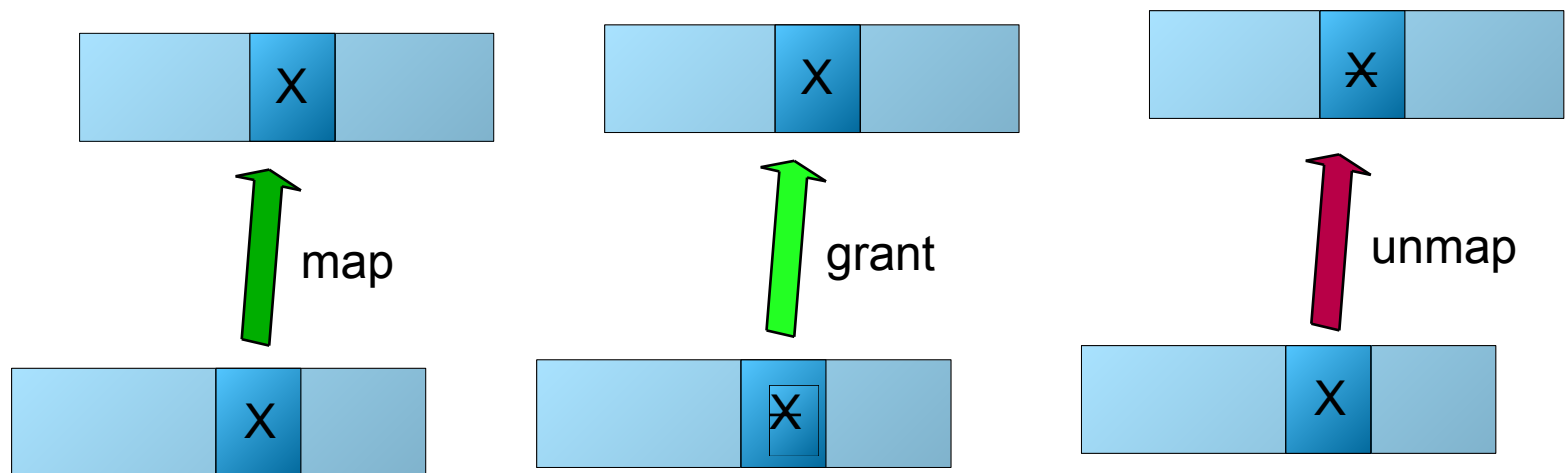**Stumpf**

TU Dresden
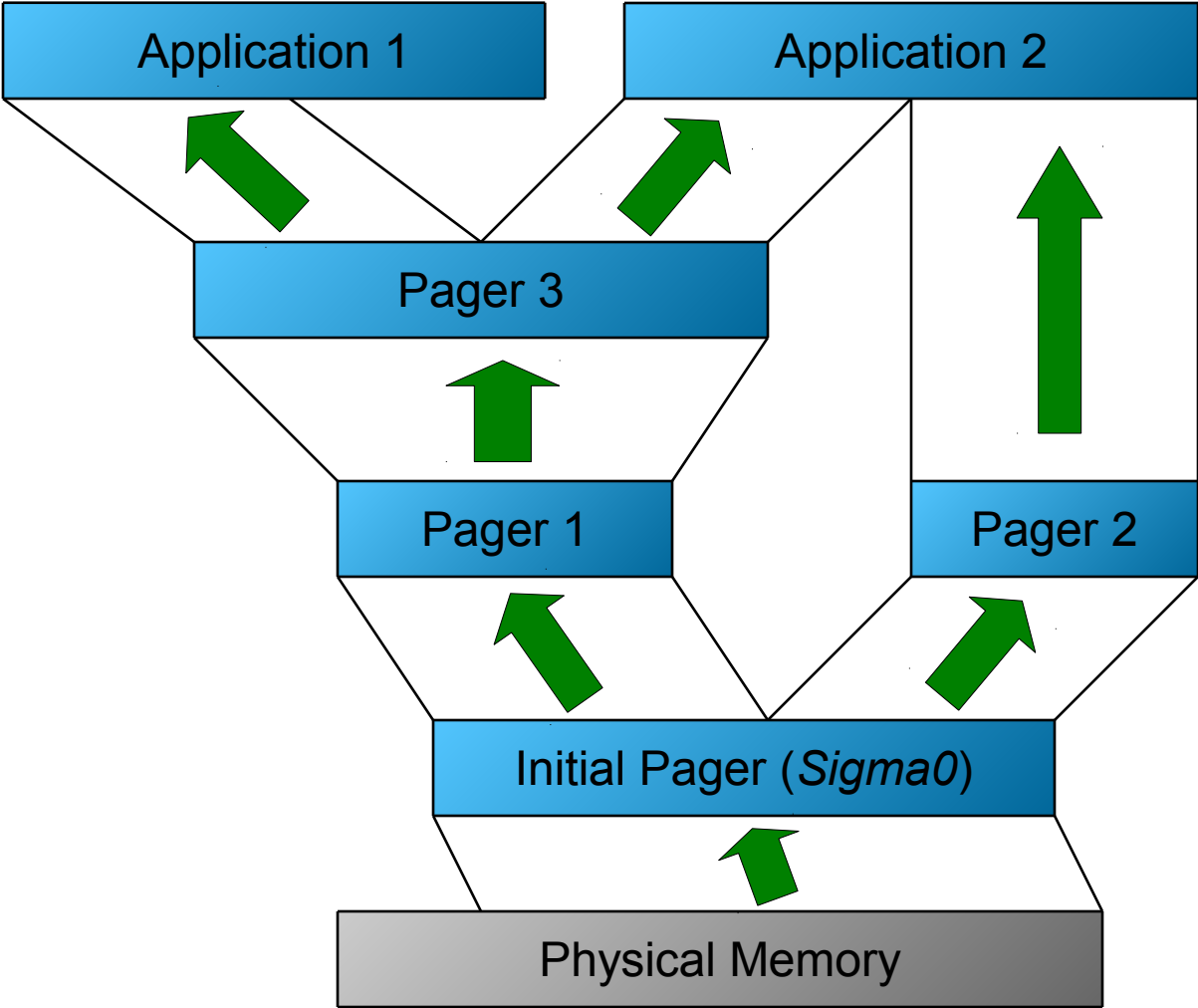Operating
Systems Group

# Threads

- ## Represent unit of execution

  - Execute user code (application)

  - Execute kernel code (system calls, page faults, interrupts, exceptions)

- ## Subject to scheduling

  - Quasi-parallel execution on one CPU

  - Parallel execution on multiple CPUs

  - Voluntarily switch to another thread possible

  - Preemptive scheduling by the kernel according to certain parameters

- ## Associated with an address space

  - Executes code in one task at one point in time

    - Migration allows threads move to another task

  - Several threads can execute in one task

Hermann
Härtig
Benjamin
Engel
**Tobias
Stumpf**

TU Dresden
Operating
Systems Group

# Tasks

- Represent domain of protection and isolation

- Container for code, data and resources

- Address space: capabilities + memory pages

- Three management operations:

  - Map: share page with other address space

  - Grant: give page to other address space
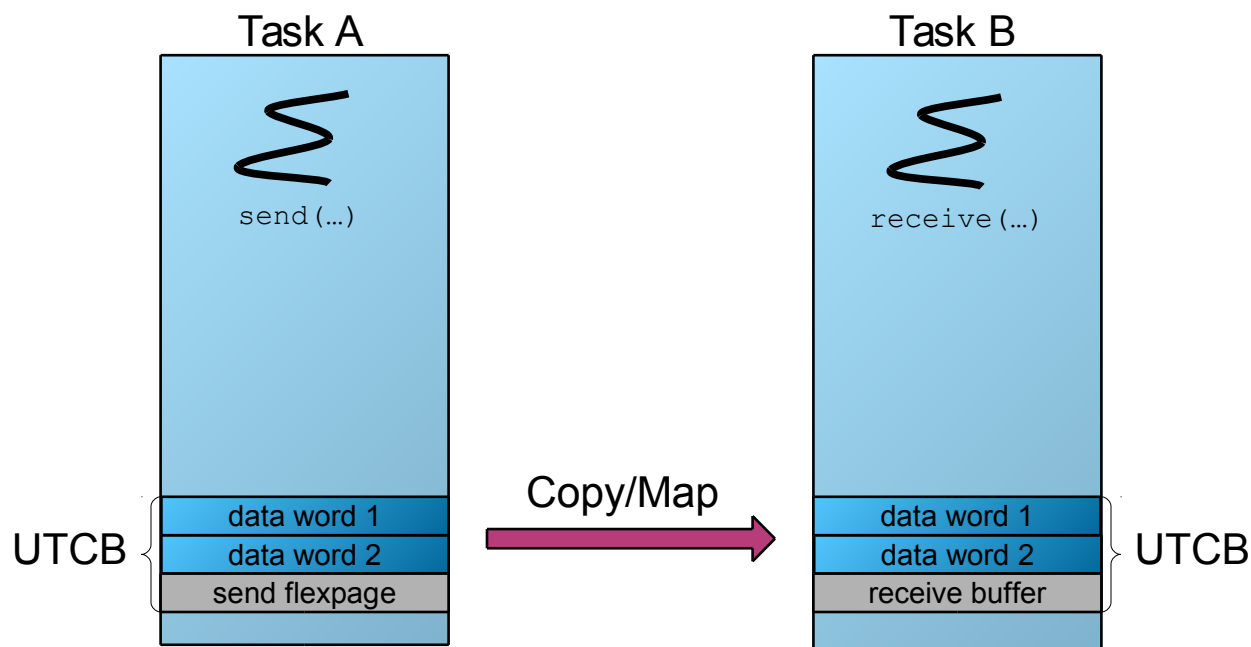
  - Unmap: revoke previously mapped page

Hermann
Härtig
Benjamin
Engel
**Tobias
Stumpf**

TU Dresden
Operating
Systems Group

map

grant

unmap

# Recursive Address Spaces

Hermann
Härtig
Benjamin
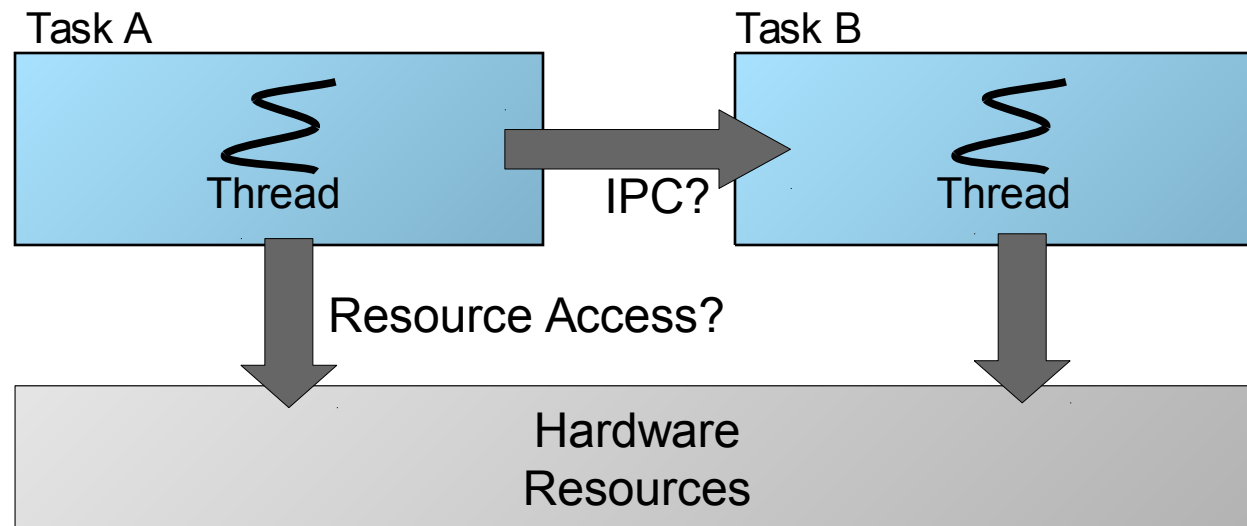Engel
**Tobias**
**Stumpf**

TU Dresden
Operating
Systems Group

# Messages: Copy Data

- UTCB (user-level thread control block)
- Always mapped, no page faults
- Transfer data and map capabilities

Hermann
Härtig
Benjamin
Engel
**Tobias**
**Stumpf**

TU Dresden
Operating
Systems Group

Task A

send(…)

Task B

receive(…)

UTCB

| data word 1 |
| data word 2 |
| send flexpage |

Copy/Map

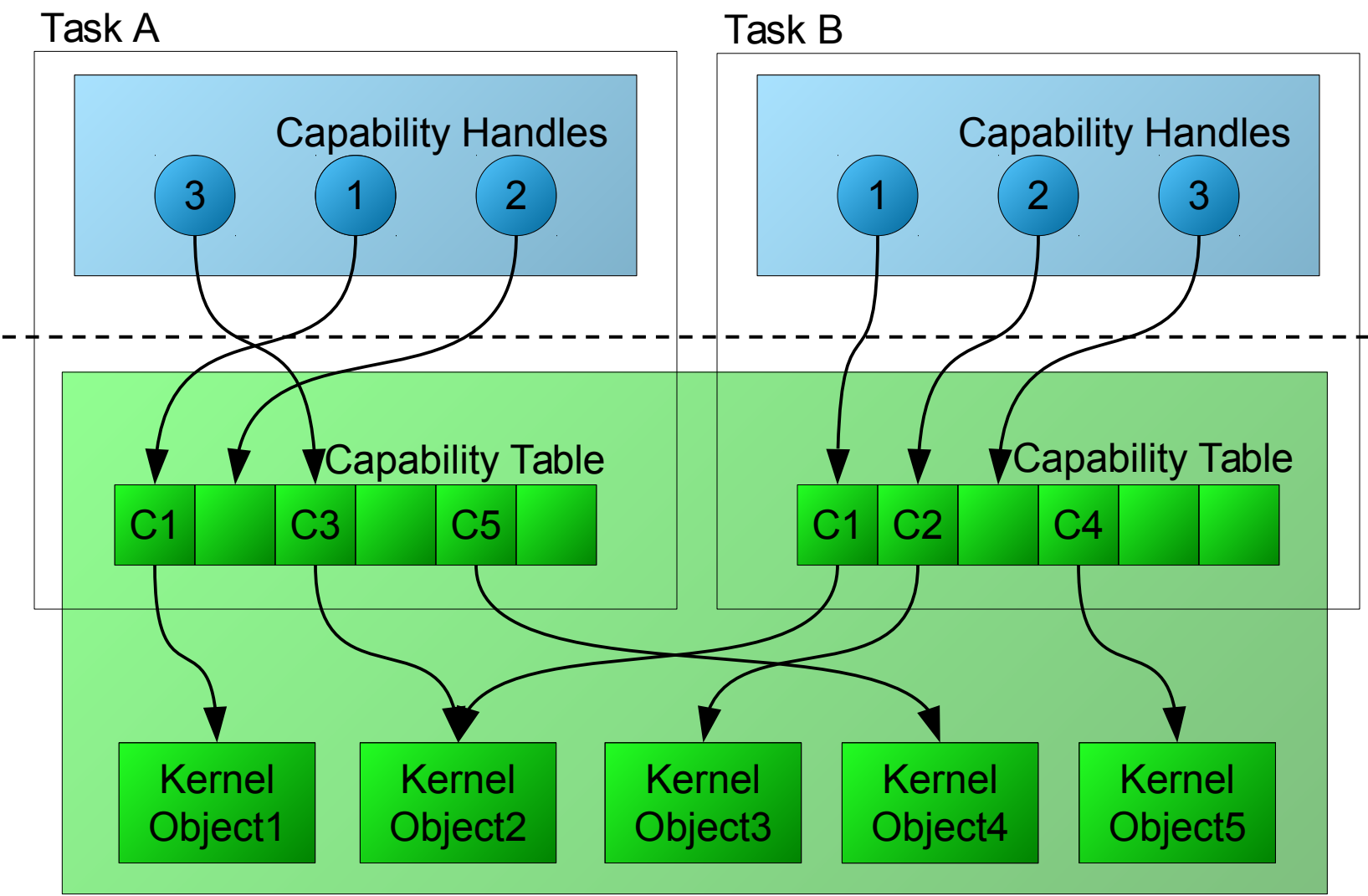| data word 1 |
| data word 2 |
| receive buffer |

UTCB

# Communication and Resource Control

- Need to control who can send data to whom
  - Security and isolation
  - Access to resources

- Approaches
  - IPC-redirection/introspection
  - Central vs. distributed policy and mechanism
  - ACL-based vs. capability-based

Hermann
Härtig
Benjamin
Engel
**Tobias**
**Stumpf**

TU Dresden
Operating
Systems Group

Task A

Task B

Thread

IPC?

Thread

Resource Access?

Hardware
Resources

# Capabilities

Hermann
Härtig
Benjamin
Engel
**Tobias**
**Stumpf**

TU Dresden
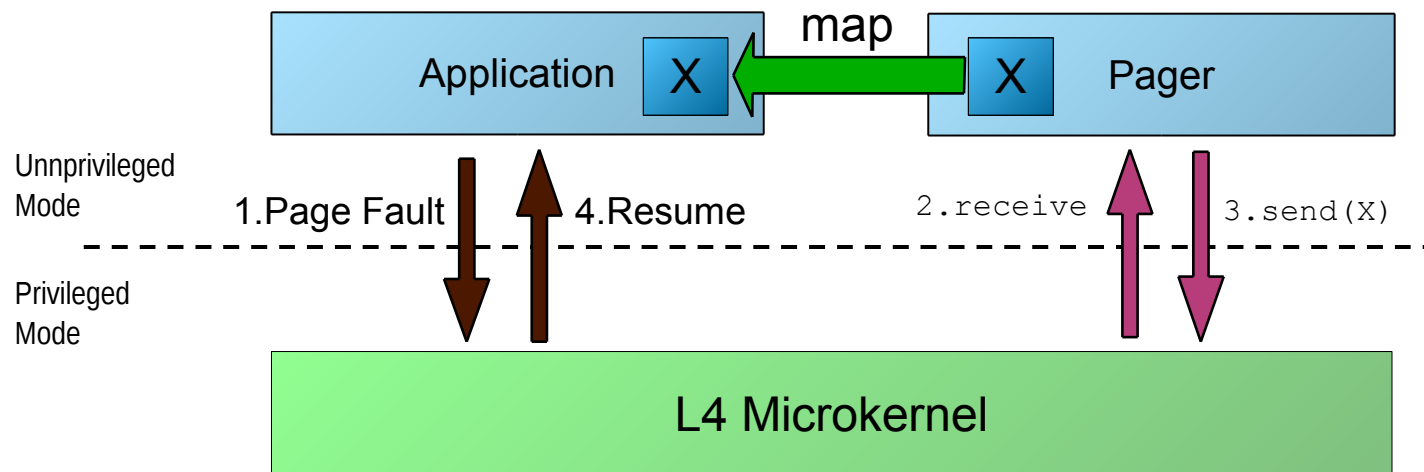Operating
Systems Group

# Capabilities - Details

- Kernel objects represent resources and communication channels

- Capability

  – Reference to kernel object

  – Associated with access rights

  – Can be mapped from task to another task

- Capability table is task-local data structure inside the kernel

  – Similar to page table

  – Valid entries contain capabilities

- Capability handle is index number to reference entry into capability table

  – Similar to file handle (in POSIX)

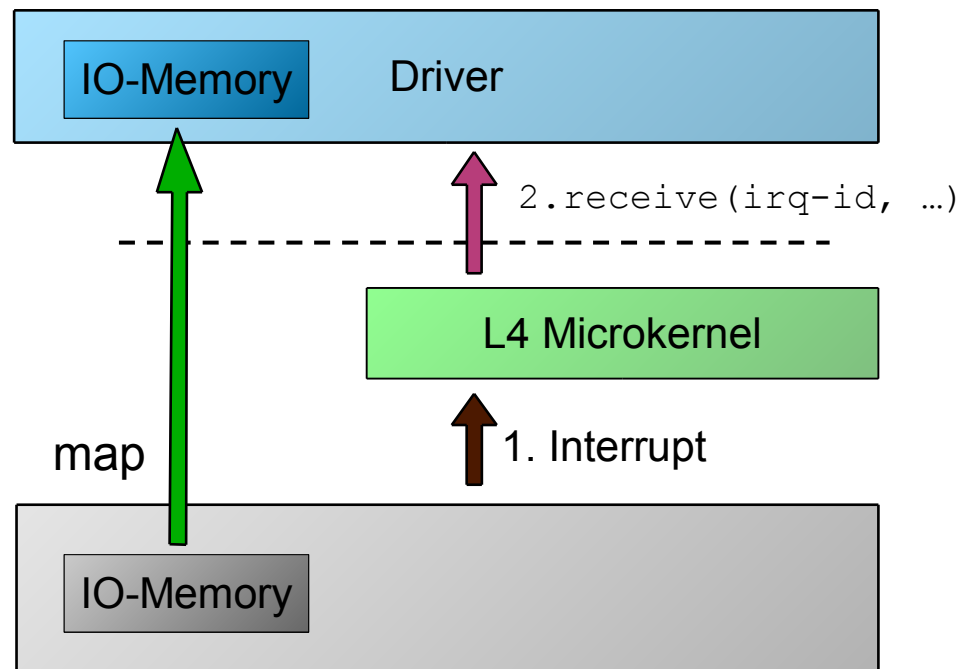- Mapping capabilities establishes a new valid entry into the capability table

Hermann
Härtig
Benjamin
Engel
Tobias
Stumpf

TU Dresden
Operating
Systems Group

# Page Faults and Pagers

- ## Page Faults are mapped to IPC
  - Pager is special thread that receives page faults
  - Page fault IPC cannot trigger another page fault
- ## Kernel receives the flexpage from pager and inserts mapping into page table of application
- ## Other faults normally terminate threads

Hermann
Härtig
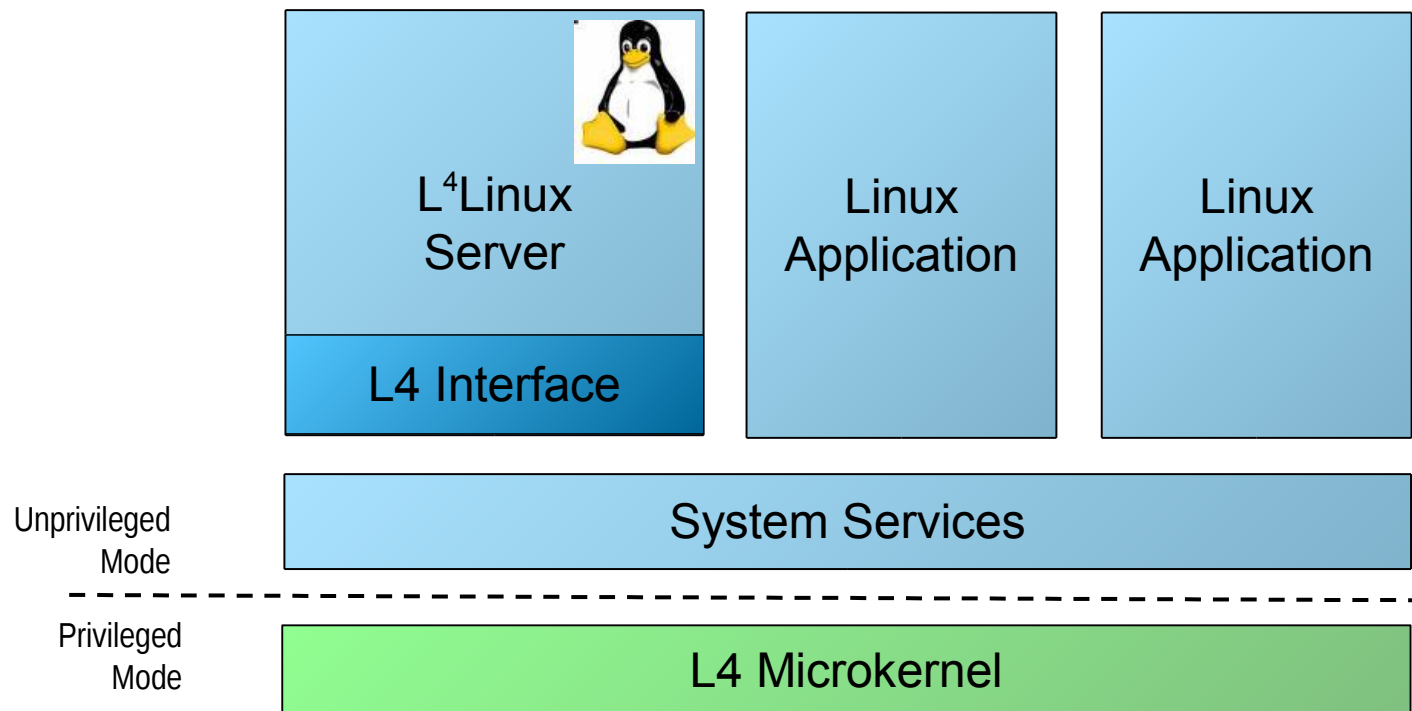Benjamin
Engel
**Tobias**
**Stumpf**

TU Dresden
Operating
Systems Group

map

Application    X          X    Pager

Unnprivileged
Mode

1.Page Fault        4.Resume        `2.receive`        `3.send(X)`

Privileged
Mode

L4 Microkernel

26

# Device Drivers

- Hardware interrupts: mapped to IPC
- I/O memory & I/O ports: mapped via flexpages



Hermann
Härtig
Benjamin
Engel
**Tobias**
**Stumpf**

TU Dresden
Operating
Systems Group

# L4 Applications - L⁴Linux

- Paravirtualized Linux kernel and native Linux applications run as user-level L4 tasks
- System calls / page faults are mapped to L4 IPC

Hermann
Härtig
Benjamin
Engel
**Tobias**
**Stumpf**

TU Dresden
Operating
Systems Group

# Lecture Outline

- **Introduction**
- Address spaces, threads, thread switching
- Kernel entry and exit
- IPC
- Address space management
- Capabilities
- Synchronization
- Case Studies: Escape, M3, Genode
- Hands-on experience

Hermann
Härtig
Benjamin
Engel
**Tobias
Stumpf**

TU Dresden
Operating
Systems Group

29