# Memory

Björn Döbel

Dresden, 2010-10-26

# So far ...

- Introduction
  - monolithic vs. microkernels
  - microkernel history / L4 family
  - L4 concepts: tasks, threads and IPC
  - Fiasco.OC/TUDOS introduction

- Threads & Synchronization
  - address spaces / tasks
  - threads
    - TCB, kernel entry
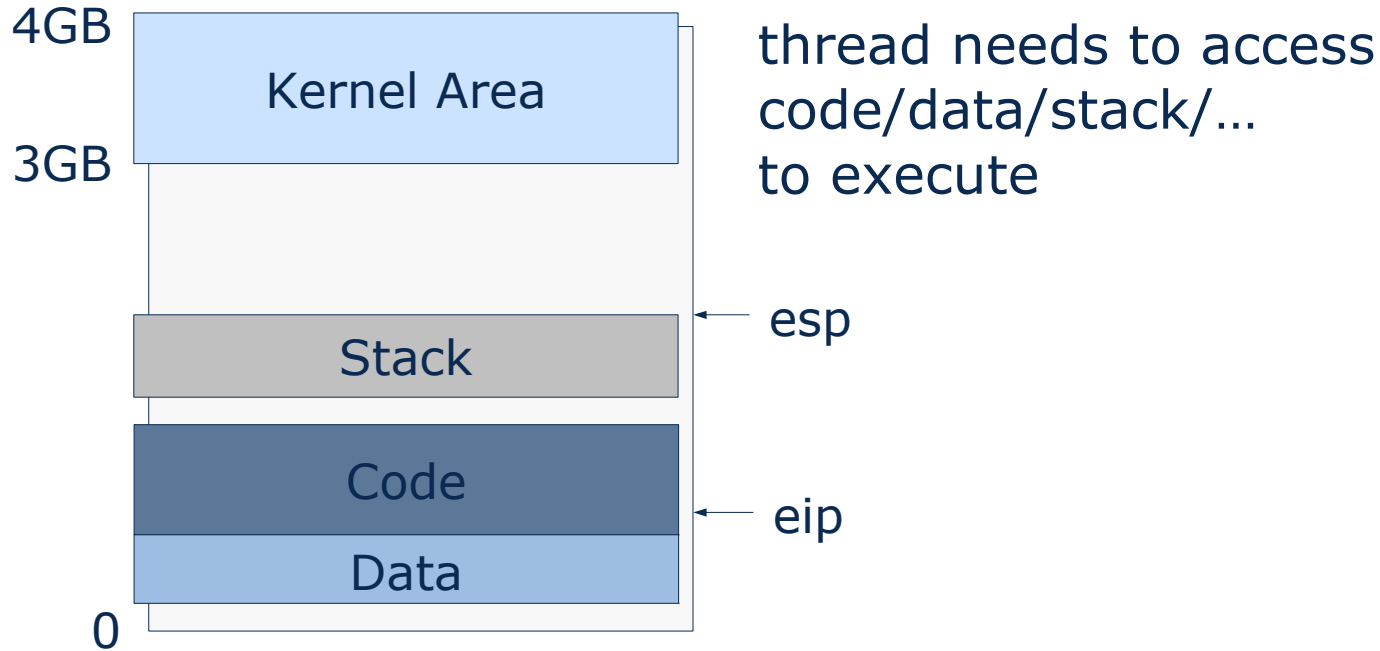  - scheduling

- **Memory**

- Task creation

- Page-fault handling

- Flexpages

- Hierarchical pagers

- Region manager

- Dataspaces

# Task Creation

- **BIOS** starts, then loads and executes boot sector

- **boot loader** (i.e. GRUB) loads kernel and multiboot modules

- **bootstrap** interprets binary modules and sets up kernel structures (Kernel Info Page)

- Fiasco.OC **kernel** started by bootstrap

- moe (**roottask**) and Sigma0 (initial address space) started by kernel

➔ next step: start application tasks

# Address space layout

4GB

Kernel Area

3GB

Stack ← esp

Code ← eip

Data

0

thread needs to access
code/data/stack/…
to execute

Moe
Root Task

Sigma0
Root Pager

Fiasco.OC
Microkernel

```
/* Create a new task. */
l4_msgtag_t
L4::Factory::create_task (Cap< Task > const &  task_cap,
                          l4_fpage_t const &     utcb_area,
                          l4_utcb_t              *utcb = l4_utcb()
)


/* Create a new thread. */
l4_msgtag_t
L4::Factory::create_thread (Cap< Thread > const &  target_cap,
                            l4_utcb_t                *utcb = l4_utcb()
)
```

```
/* Commit the given thread-attributes object. */
l4_msgtag_t
L4::Thread::control (Attr const &  attr)


/* Exchange basic thread registers. */
l4_msgtag_t
L4::Thread::ex_regs (l4_addr_t      ip,          /* instruction pointer */
                     l4_addr_t      sp,          /* stack pointer */
                     l4_umword_t    flags,
                     l4_utcb_t      *utcb = l4_utcb()
)
```
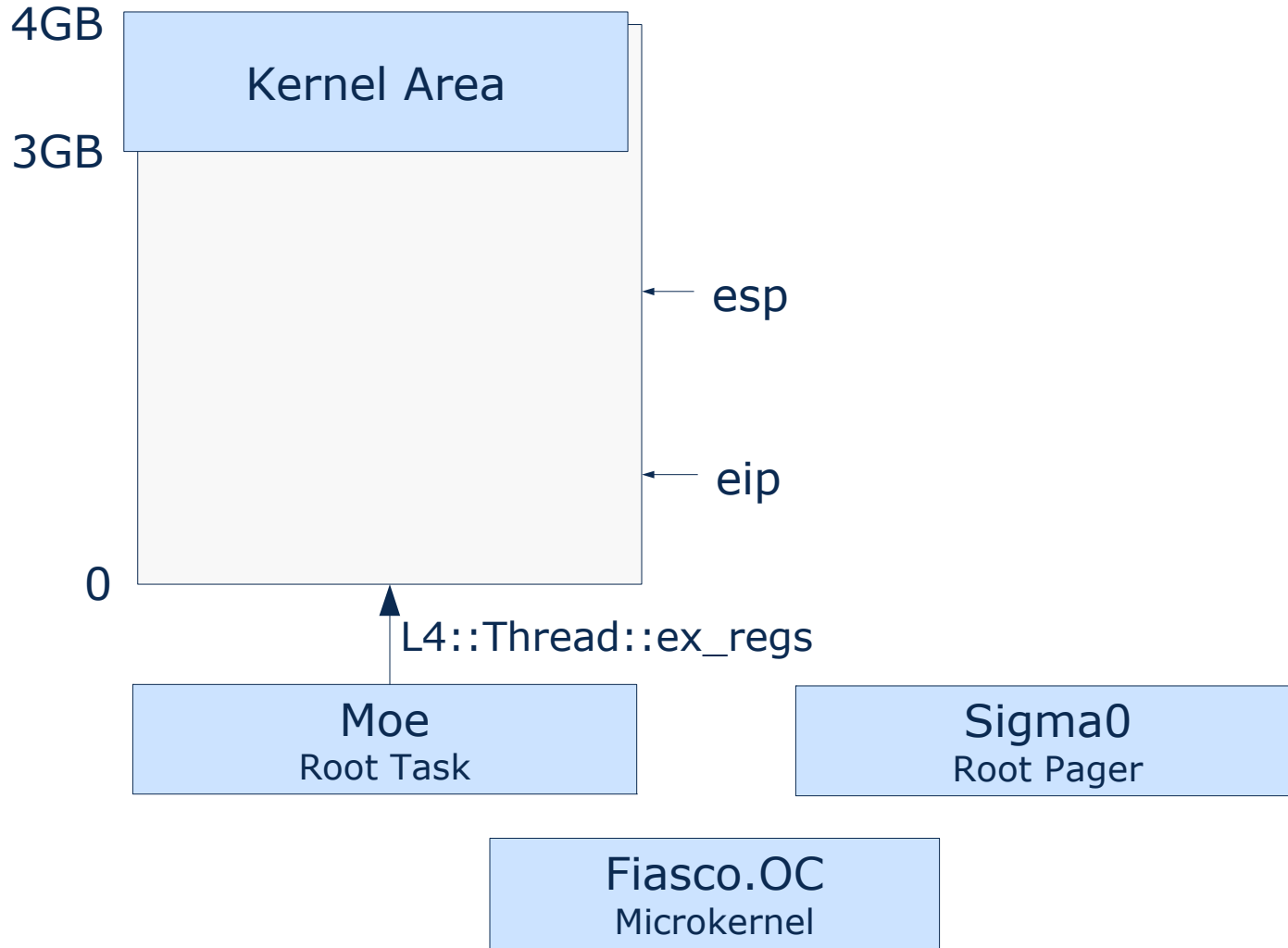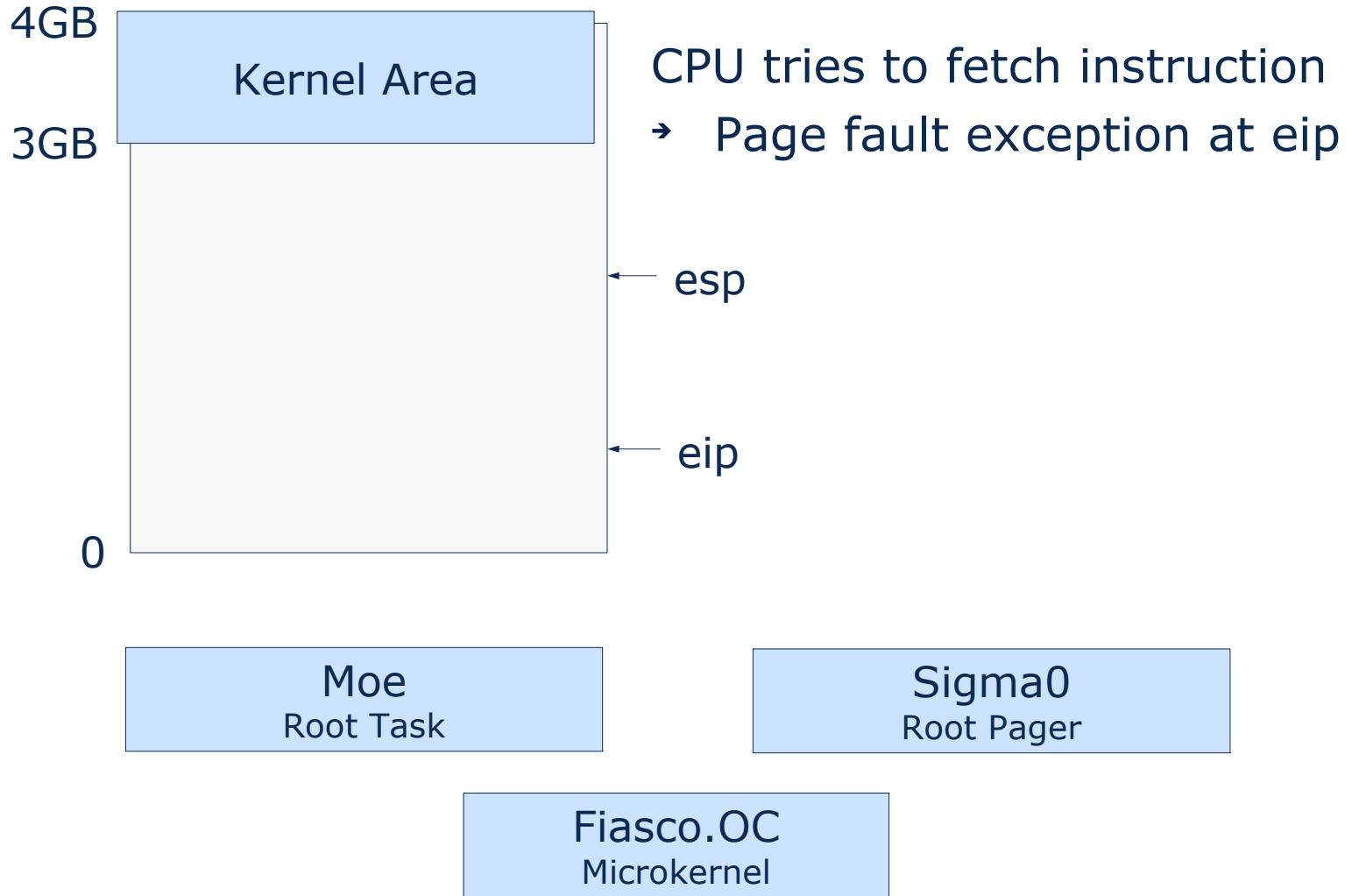
4GB

Kernel Area

3GB

← esp

← eip

0

L4::Thread::ex_regs

Moe
Root Task

Sigma0
Root Pager

Fiasco.OC
Microkernel

4GB

Kernel Area

3GB

0

CPU tries to fetch instruction
→ Page fault exception at eip

← esp

← eip

Moe
Root Task

Sigma0
Root Pager

Fiasco.OC
Microkernel

# Page Fault Handling

# Page-fault handling

Application's address space

eip →

Fiasco.OC
Microkernel
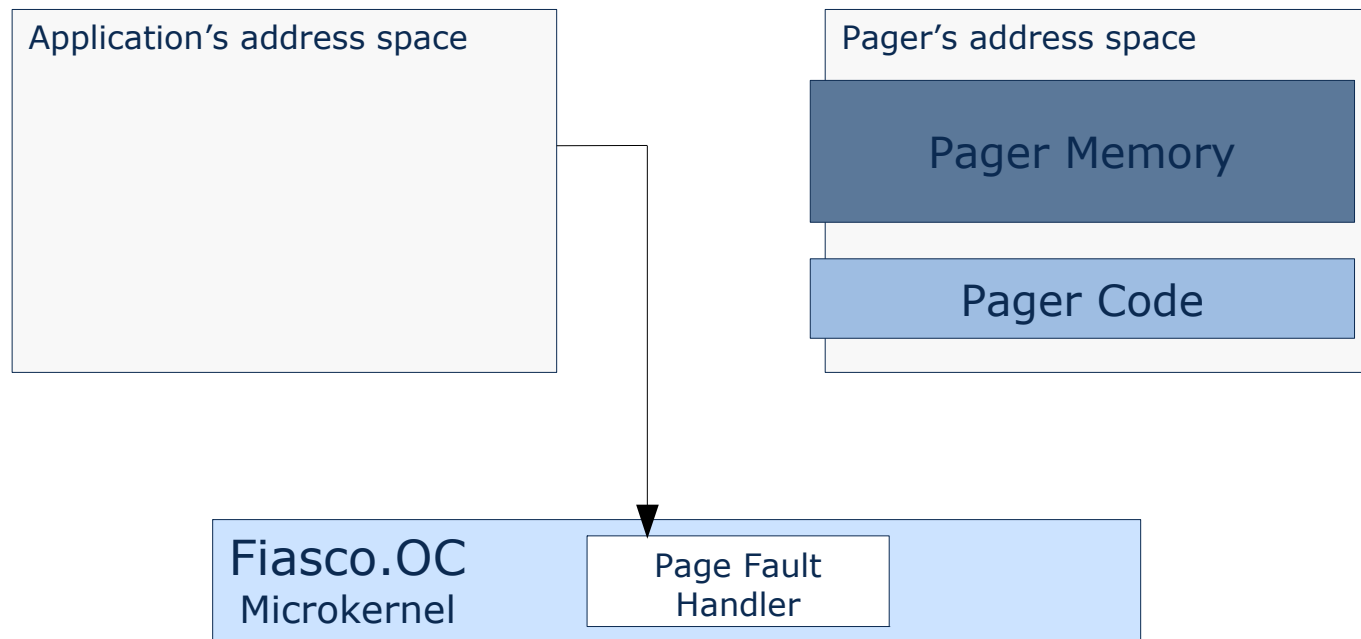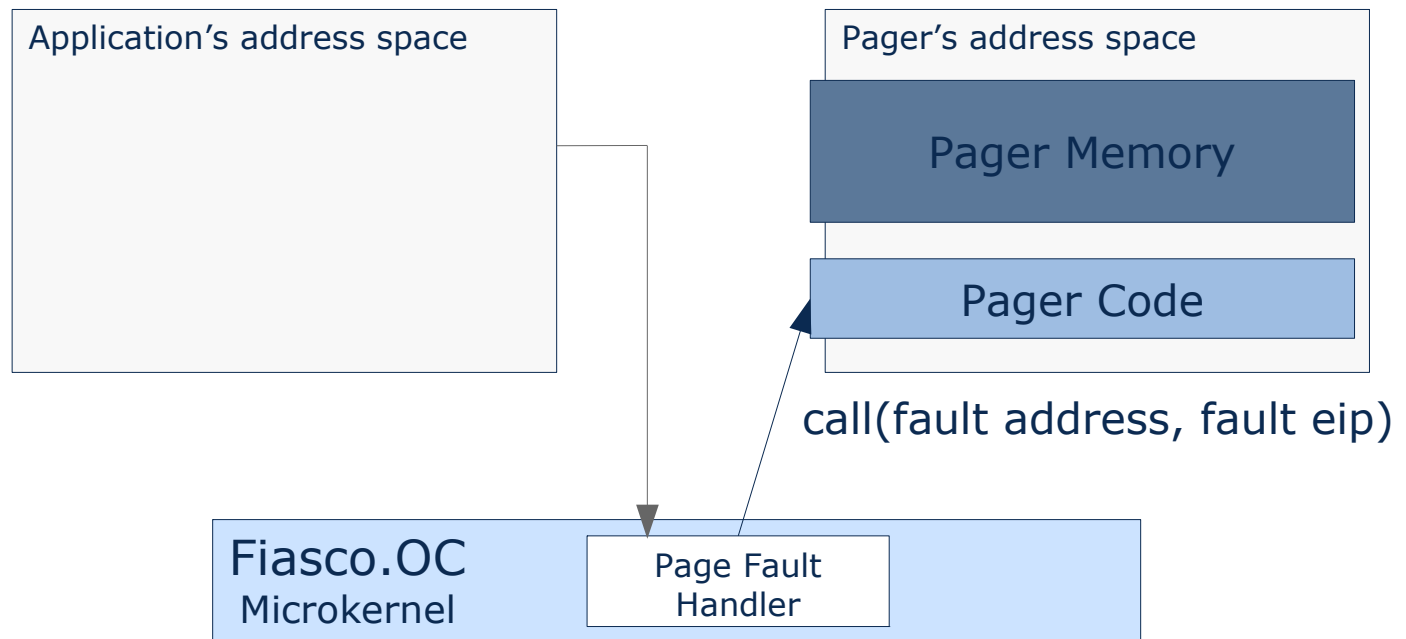
Page Fault
Handler

- Page fault exception is caught by kernel page-fault handler

- no management of user memory in kernel

- invoke user-level memory management ⇨ **Pager**
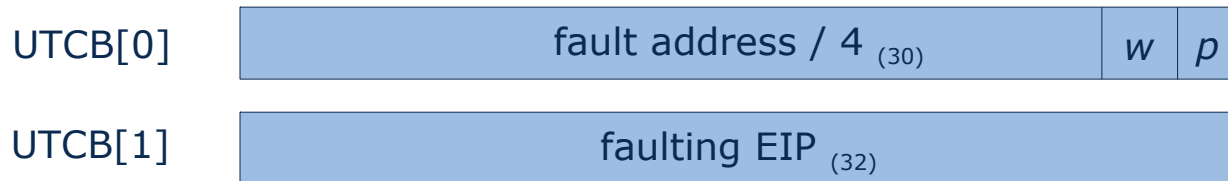
- thread which is invoked on page fault

- each thread has a (potentially different) pager assigned

- communication with pager thread using IPC

- kernel page fault handler sets up IPC to pager

- pager sees faulting thread as sender of IPC

UTCB[0]

| fault address / 4 $_{(30)}$ | $w$ | $p$ |
| --- | --- | --- |

UTCB[1]

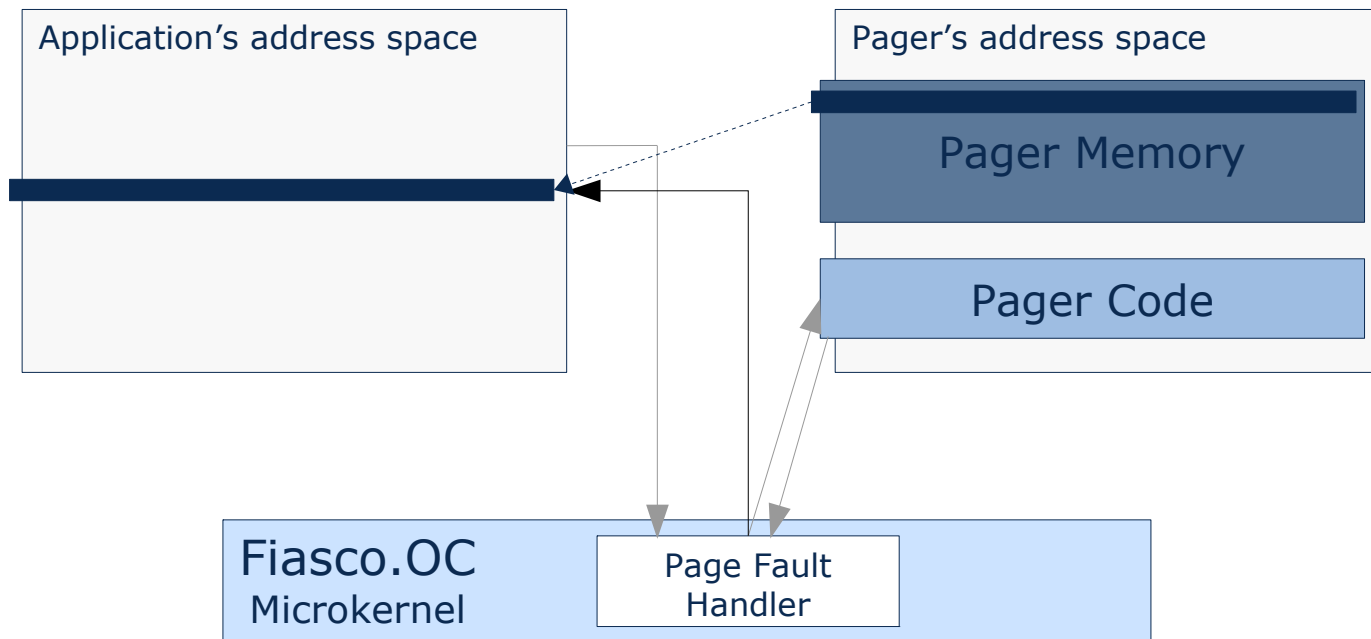| faulting EIP $_{(32)}$ |
| --- |

$w = 0$     read page fault
$w = 1$     write page fault
$p = 0$     no page present
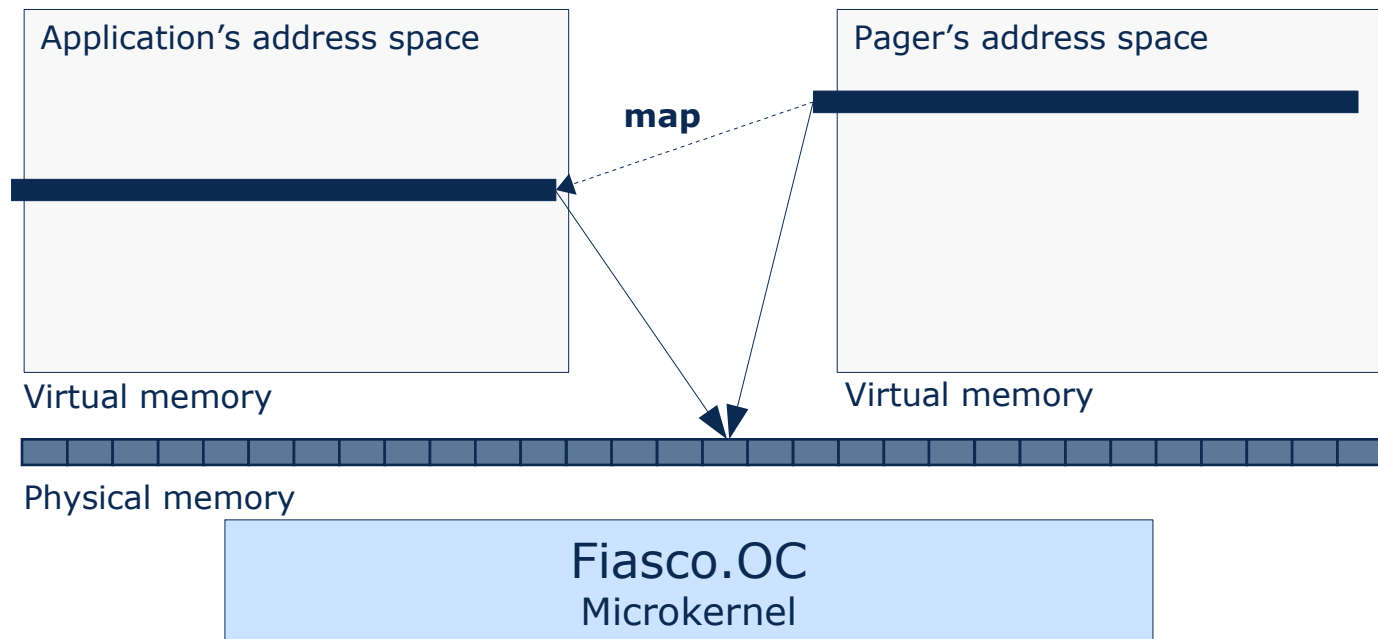$p = 1$     page present

- pager maps pages of it's own address space to the application's address space

- flexpage IPC enables these mappings

- map creates an entry in the receiver's address space pointing to the same page frame

- only valid pager address space entries can be mapped

- Special case: grant pages (flag: L4_FPAGE_GRANT)
➔ Removes mapping from sender's address space



Application's address space

Pager's address space

**grant**

Virtual memory

Virtual memory

Physical memory

Fiasco.OC
Microkernel

- Special case: grant pages (flag: L4_FPAGE_GRANT)

➔ Removes mapping from sender's address space

  ➔ ATTENTION: aliases remain

- Removes entries to a page frame (fpage is specified in invoker's address space)

➔ Kernel tracks mappings in a database

# Flexpages

- In general, flexpages represent areas within an address space

- Flexpages in Fiasco.OC are used to describe:
  - Memory pages
  - I/O ports
  - Capabilities

- Today only flexpages for memory pages are described.

# Flexpage IPC in detail

- Size-aligned
- Size $2^{size}$, smallest is hardware page
- Source and target area of a map IPC are described by flexpages



Sender's Address Space

l4_ipc_receive(…)

Receiver's Address Space

| page address$_{(20)}$ | size$_{(6)}$ | ~$_{(2)}$ | rights$_{(4)}$ |

Receive Window

| send base$_{(32)}$ | | | |
| page address$_{(20)}$ | size$_{(6)}$ | ~$_{(2)}$ | rights$_{(4)}$ |

l4_ipc_send(…)

# Flexpage offset

- send flexpage is smaller than the receive window
  - target position is derived from send flexpage alignment and send base

l4_ipc_send(…)

l4_ipc_receive(…)

| send base$_{(32)}$ | | | |
|---|---|---|---|
| page address$_{(20)}$ | size$_{(6)}$ | ~$_{(2)}$ | rights$_{(4)}$ |

| page address$_{(20)}$ | size$_{(6)}$ | ~$_{(2)}$ | rights$_{(4)}$ |
|---|---|---|---|

- send flexpage is larger than receive window
  - target position is derived from receive flexpage alignment and send base
➔ send base depends on information about the receiver

l4_ipc_send(…)

| send base$_{(32)}$ | | | |
|---|---|---|---|
| page address$_{(20)}$ | size$_{(6)}$ | ~$_{(2)}$ | rights$_{(4)}$ |

l4_ipc_receive(…)

| page address$_{(20)}$ | size$_{(6)}$ | ~$_{(2)}$ | rights$_{(4)}$ |
|---|---|---|---|

# Page fault in detail

- kernel page fault handler sets receive window to whole address space
➜ pager can map more than just one page, where the page fault happened to the client



| Application's Address Space | send base ⇨ pf address | | | | Pager's Address Space |
|---|---|---|---|---|---|
| | page address$_{(20)}$ | size$_{(6)}$ | ~$_{(2)}$ | rights$_{(4)}$ | |

whole space

Page Fault Receive Window

Fiasco.OC Microkernel

Page Fault Handler

# Root task's pager

- pages are mapped as they are needed
- → *demand paging*

Application's address space

esp →  Stack

eip →  Code
       Data

Root Task's address space

Memory

App Code Section
App Data Section

Roottask

**Fiasco.OC**
Microkernel

# Hierarchical Pagers

- initial pager can only implement basic memory management

- no knowlegde about application requirements
  - different requirements at the same time

- missing services for advanced memory management
  - e.g. no disk driver for swapping

- build more advanced pagers on top of the initial one

➔ pager hierarchy

# Pager hierarchy



Application

Application

Pager 3

Disk Driver

Pager 1

Pager 2

Phys. Memory
1-to-1 mapped

Initial Address Space

Fiasco.OC
Microkernel

- L$^4$Linux implements Linux paging policy
- RT pager implements real-time paging policy
  (e.g. no swapping)

- pager has to specify send base
- pager needs to know client's address space layout
  - no problems with only one pager (e.g. L$^4$Linux)
- possible conflicts if more than one pager manages an address space:



➔ virtual memory must be managed independent of pagers

# Region Mapper

- per address space map that keeps track which part of the address space is managed by which pager

| Address space | Region Map |
|---|---|

**Address space**

Pager 1

Pager 2

Pager 2

**Region Map**

| | |
|---|---|
| <start, end> | **Pager 1** |
| <start, end> | **Pager 2** |
| <start, end> | **Pager 2** |

- intermediate pager that identifies which pager should handle a page fault
- can reside in the application's address space
➔ region mapper has to be pager of all thread of a task

| VM Region | Pager |
|---|---|
| <start, end> | Pager 1 |
| <start, end> | Pager 2 |
| <start, end> | Pager 2 |

Application

Region Mapper

Pager 1

Pager 2

Pager 2

call( region mapper,
pf addr,
pf eip, ...)

Pager 1

Pager Memory

Pager Code

Fiasco.OC
Microkernel

Page-Fault
Handler

# Region mapper

- region mapper calls the pager that is responsible
- receive window gets restricted to the area managed by that pager
- ➔ no interference between different pagers

| Application | | VM Region | Pager | | Pager 1 |
|---|---|---|---|---|---|
| Region Mapper | | <start, end> | Pager 1 | | |
| | | <start, end> | Pager 2 | | Pager Memory |
| Pager 1 | | <start, end> | Pager 2 | | |
| Pager 2 | | | | | Pager Code |
| Pager 2 | | | | | |

**Call(pager1, pf addr,
pf eip,
receive window,...)**

**Fiasco.OC**
Microkernel

Page-Fault
Handler

- memory management in terms of pages so far

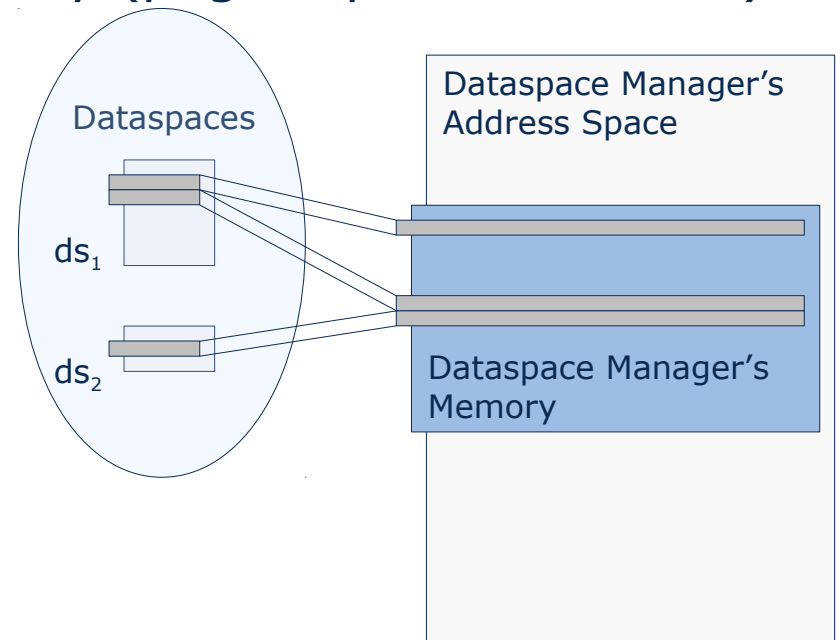- application's view to memory
  - code / data sections
  - memory mapped files
  - anonymous memory (heaps, stacks, …)
  - network / file system buffers
  - …

➔ abstraction to map this view to low-level memory management

# Dataspaces

- Dataspace: *unstructured data container*
- abstraction for anything that contains data:
  - Files
  - Anonymous memory
  - I/O adapter memory
  - …
- Dataspaces are implemented by *Dataspace Managers*
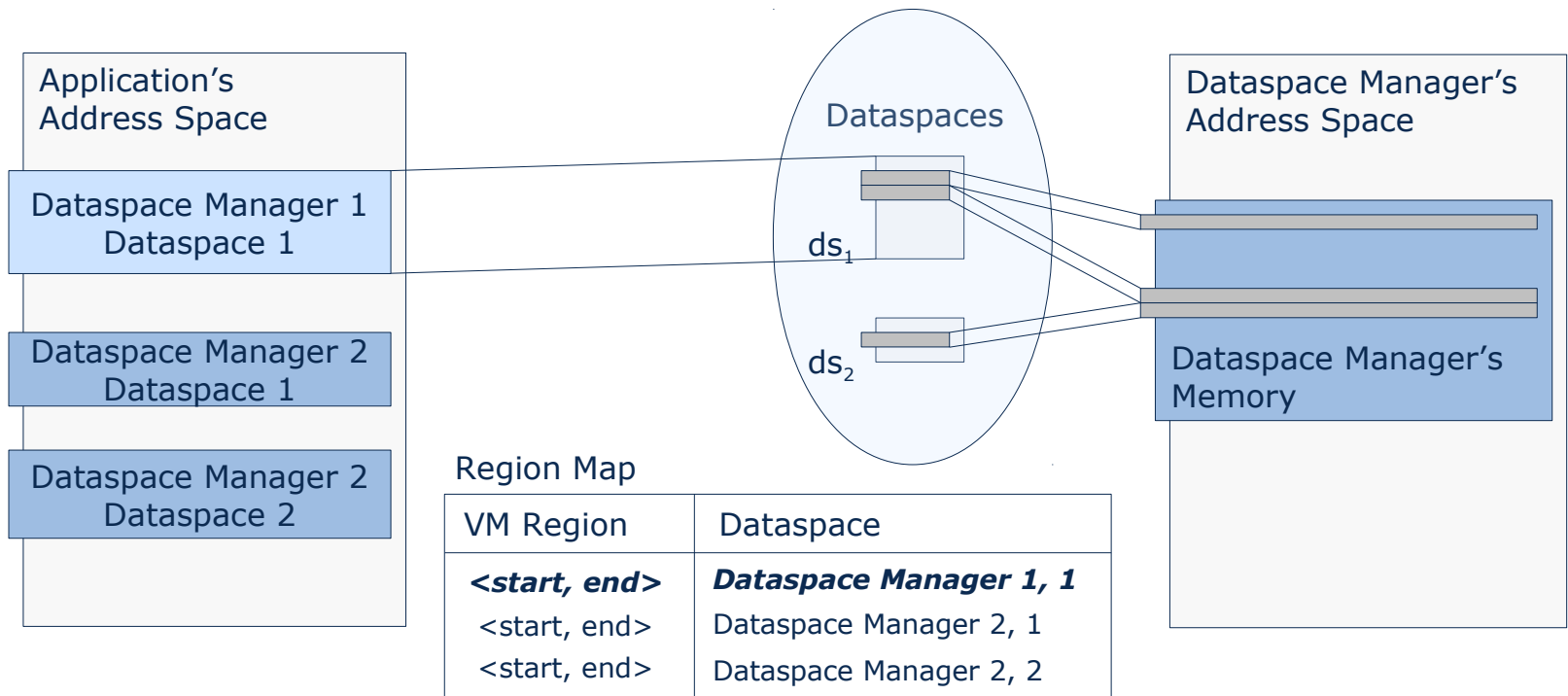- Dataspaces can be attached to regions of an address space

- a Dataspace Manager determines the semantic of a dataspace
- each Dataspace Manager is the pager for its dataspaces
➔ implements the paging policy (page replacement etc.)



Dataspaces

$ds_1$

$ds_2$

Dataspace Manager's Address Space

Dataspace Manager's Memory

# Dataspaces & region mapper

- region map keeps track which dataspaces are attached to which virtual memory regions

- region mapper translates page faults to dataspace offsets



Application's Address Space

Dataspace Manager 1
Dataspace 1

Dataspace Manager 2
Dataspace 1

Dataspace Manager 2
Dataspace 2

Dataspaces

$ds_1$

$ds_2$

Dataspace Manager's Address Space

Dataspace Manager's Memory

Region Map

| VM Region | Dataspace |
|---|---|
| **<start, end>** | **Dataspace Manager 1, 1** |
| <start, end> | Dataspace Manager 2, 1 |
| <start, end> | Dataspace Manager 2, 2 |

- region mapper propagates fault to dataspace manager's fault handler
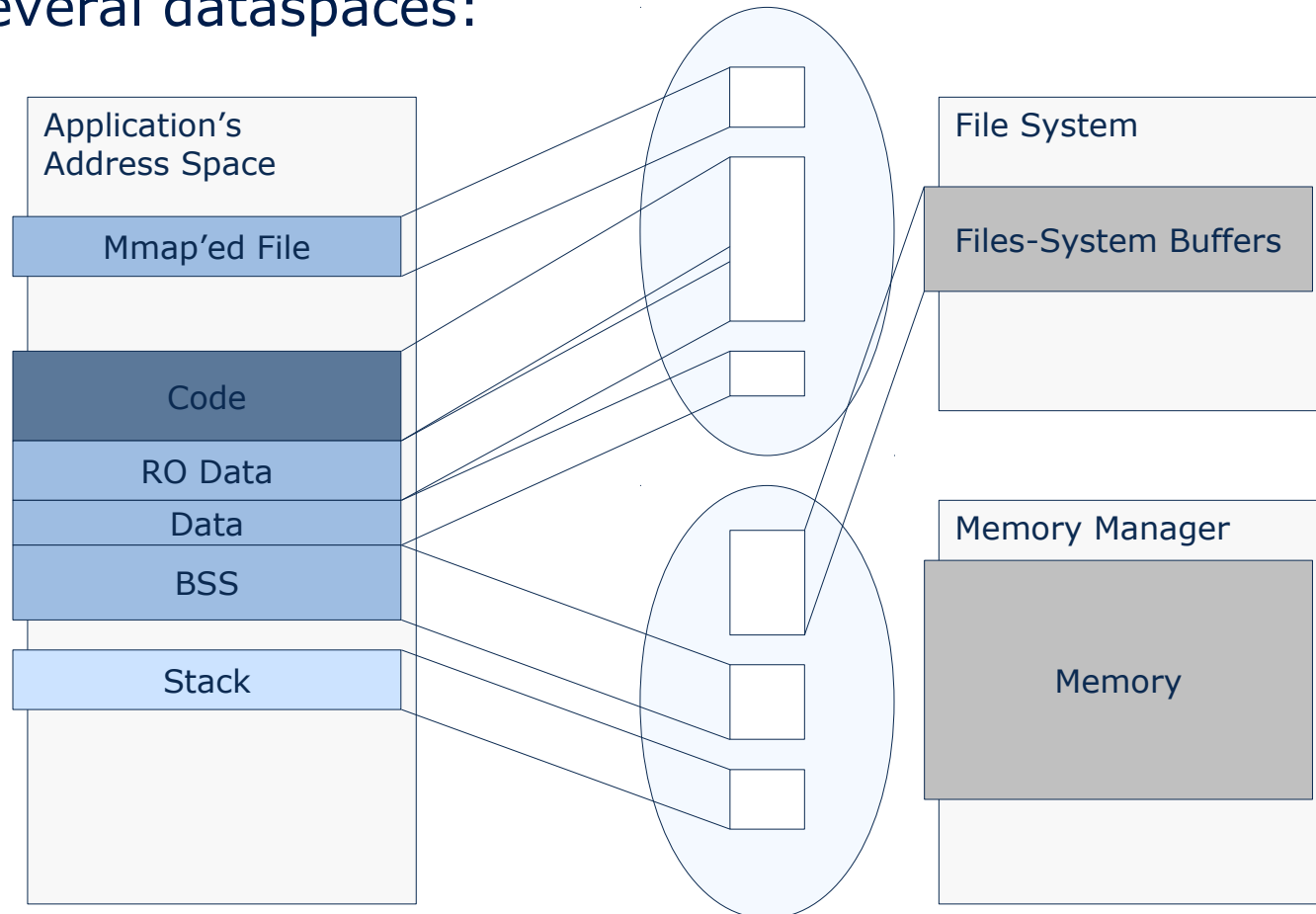➔ dataspace fault (ds_manager_id, ds_id, offset)

| Application's Address Space | Dataspaces | Dataspace Manager's Address Space |
|---|---|---|
| Dataspace Manager 1 Dataspace 1 | $ds_1$ | |
| Dataspace Manager 2 Dataspace 1 | $ds_2$ | Dataspace Manager's Memory |
| Dataspace Manager 2 Dataspace 2 | | |

Region Map

| VM Region | Dataspace |
|---|---|
| *<start, end>* | *Dataspace Manager 1, 1* |
| <start, end> | Dataspace Manager 2, 1 |
| <start, end> | Dataspace Manager 2, 2 |

- allocate / free dataspaces
  - create / destroy dataspace
  - semantic depends on dataspace type:
    - anonymous memory: open (size)
    - file: open (filename, mode, …)
    - …
- attach / detach dataspace
  - create / remove entry in region map
  - ➔ Makes dataspace contents accessible to application
- propagate capability
  - grant access rights to other applications
  - ➔ very easy shared memory implementation

- application address spaces are constructed from several dataspaces:

- page Allocation Algorithms
  - list-based algorithms, bitmaps, trees, …
- page Replacement Algorithms
  - Least-Recently-Used (LRU)
  - Working Sets
  - Clock
  - …

➔ both page allocation and page replacement are implemented by dataspace managers

➔ can have different strategies for the dataspaces of an application

- memory sharing important for
  - shared libraries
  - data transfer between system components
  - …
- different types of sharing
  - full sharing, all clients see modifications
  - ➔ easy to implement, pager / dataspace manager grants access rights to pages / dataspaces
  - lazy copying of dataspaces
  - ➔ copy-on-write

# Summary

- closer look on tasks/threads:
  - creation
  - page-fault handling
- flexpages
  - memory pages, I/O ports, Capabilities
  - structure
  - offset computation
- pager hierarchy
- region mapper & dataspaces

# More information

- Flexpages

  H.Härtig, J.Wolter, J.Liedtke: "*Flexible sized page objects*" , http://os.inf.tu-dresden.de/papers_ps/flexpages.pdf

- Dataspaces

  Mohit Aron, Yoonho Park, Trent Jaeger, Jochen Liedtke, Kevin Elphinstone, Luke Deller: "*The SawMill Framework for VM Diversity*", ftp://ftp.cse.unsw.edu.au/pub/users/disy/ papers/Aron_PJLED_01.ps.gz

- next lecture (2.11.) on 'Communication':
  - IPC flavors
  - communication control

- next exercise (2.11.) on:
  - Practical exercise, computer pool