



TECHNISCHE
UNIVERSITÄT
DRESDEN

Faculty of Computer Science Institute for System Architecture, Operating Systems Group

Memory

Henning Schild

Dresden, 2009-10-27

- 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
 - synchronization
- **Memory**



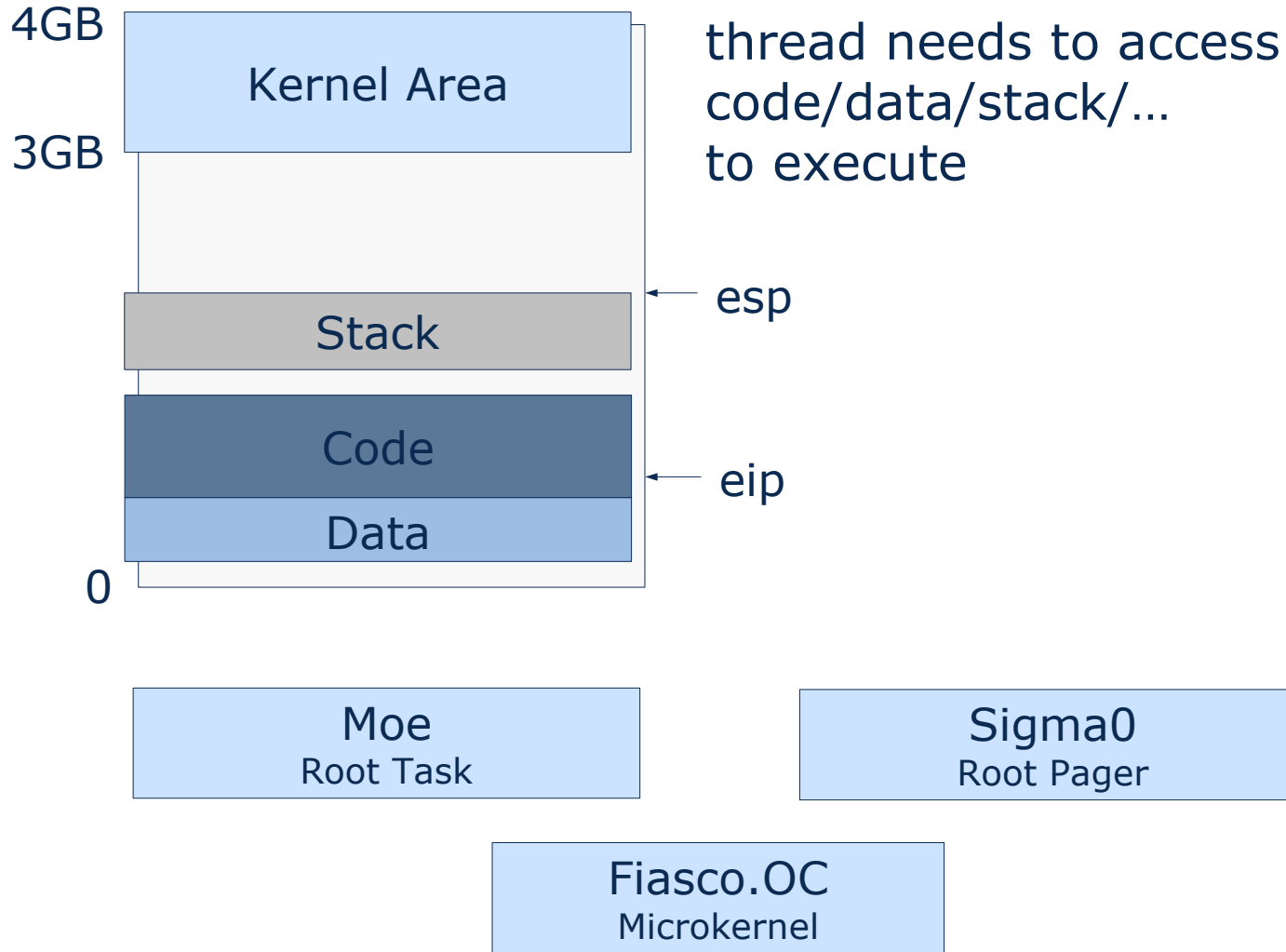
- Task creation
- Page-fault handling
- Flexpages
- Hierarchical pagers
- Region manager
- Dataspaces



Task Creation

- BIOS started, loaded and executed boot sector
 - boot loader (i.e. GRUB) loaded kernel and boot modules
 - bootstrap interpreted binary modules and set up kernel structures (Kernel Info Page)
 - kernel started by Bootstrap
 - Moe (root task) and Sigma0 (initial address space) started by kernel
- next step: start application tasks**

Address space layout



```
/* 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 */
```

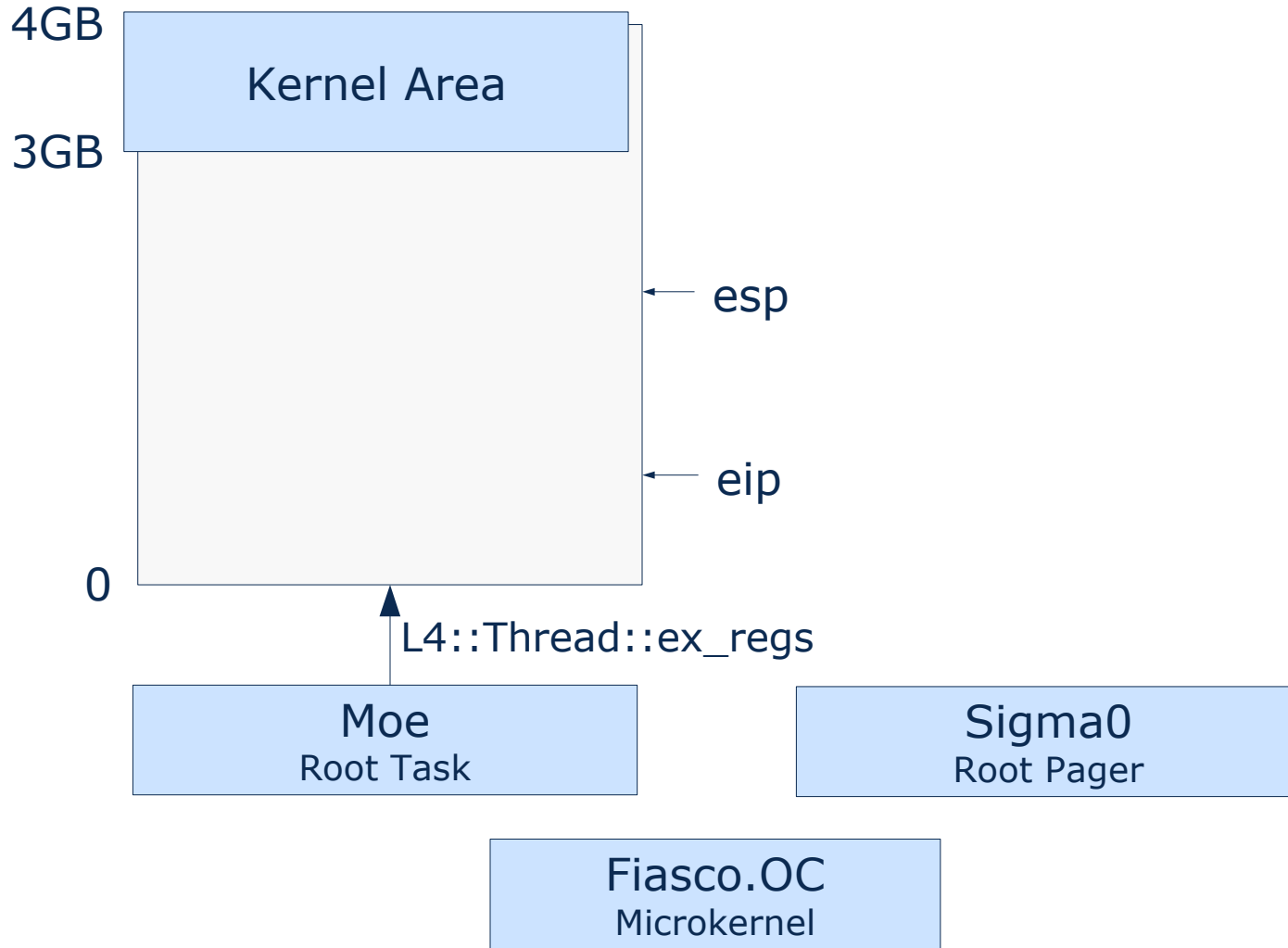
```
l4_addr_t sp, /* stack pointer */
```

```
l4_umword_t flags,
```

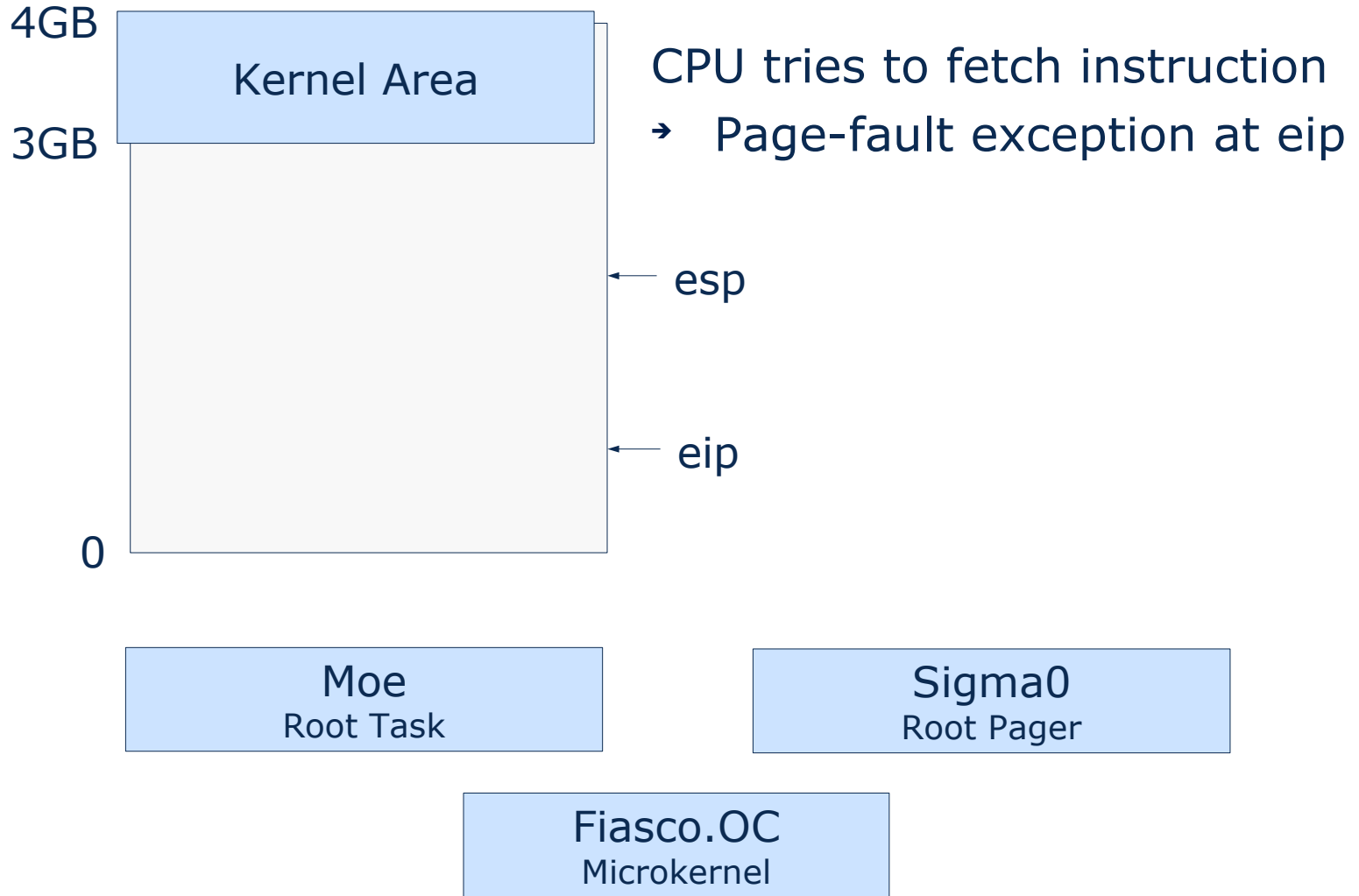
```
l4_utcb_t *utcb = l4_utcb()
```

```
)
```

L4::Thread::ex_regs



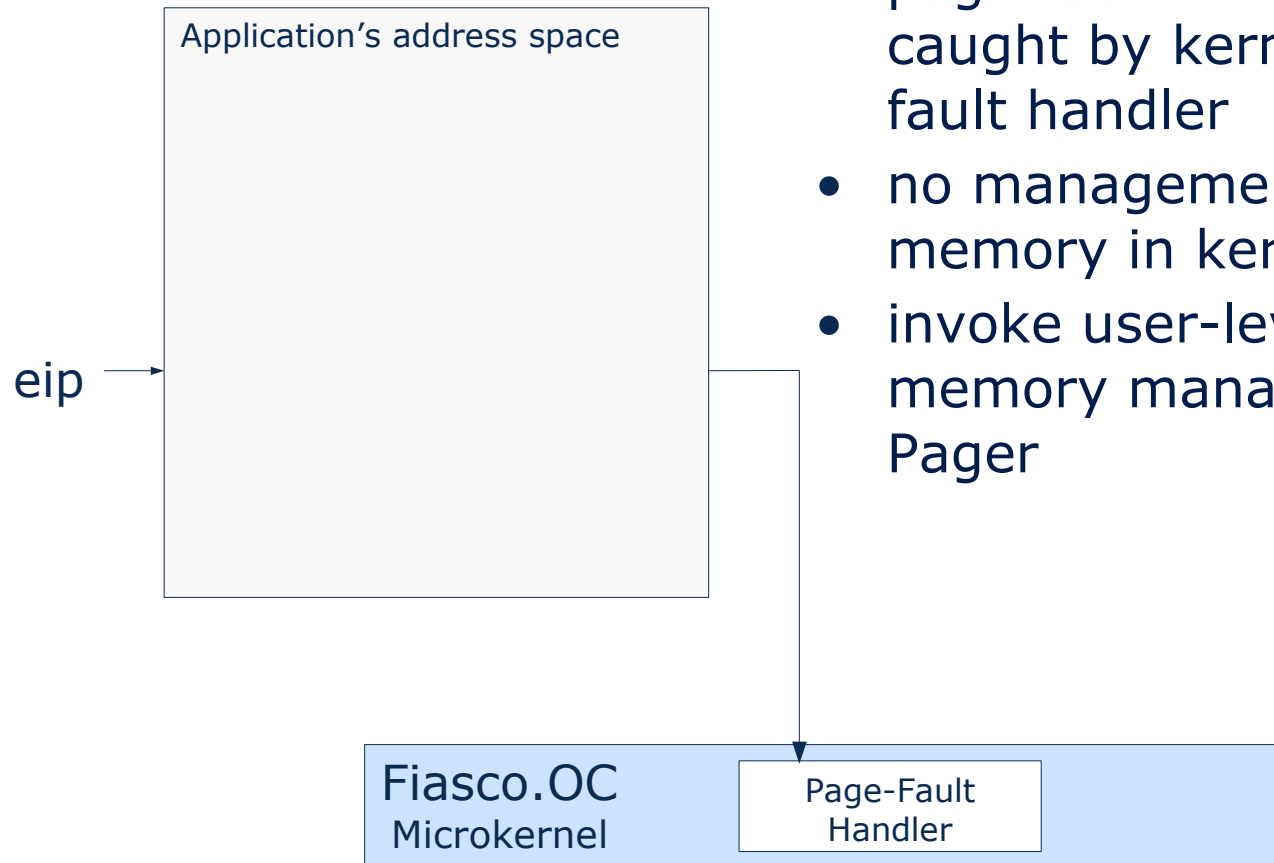
First thread executes ...





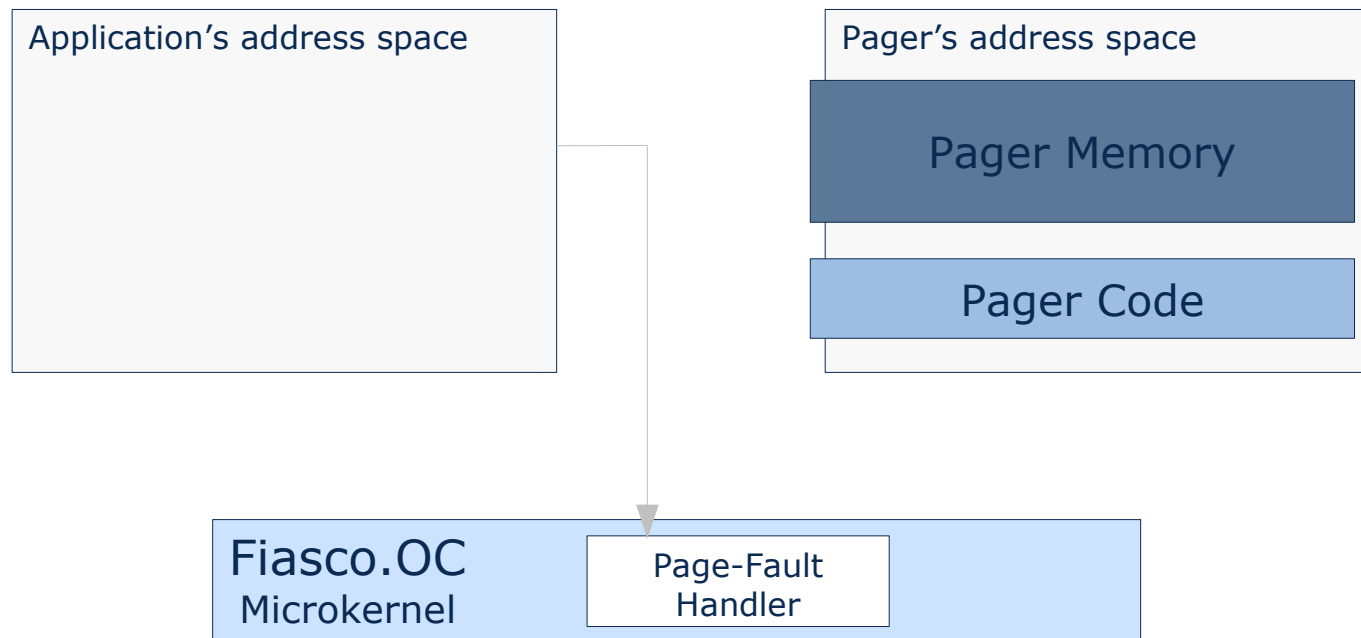
Page-Fault Handling

Page-fault handling



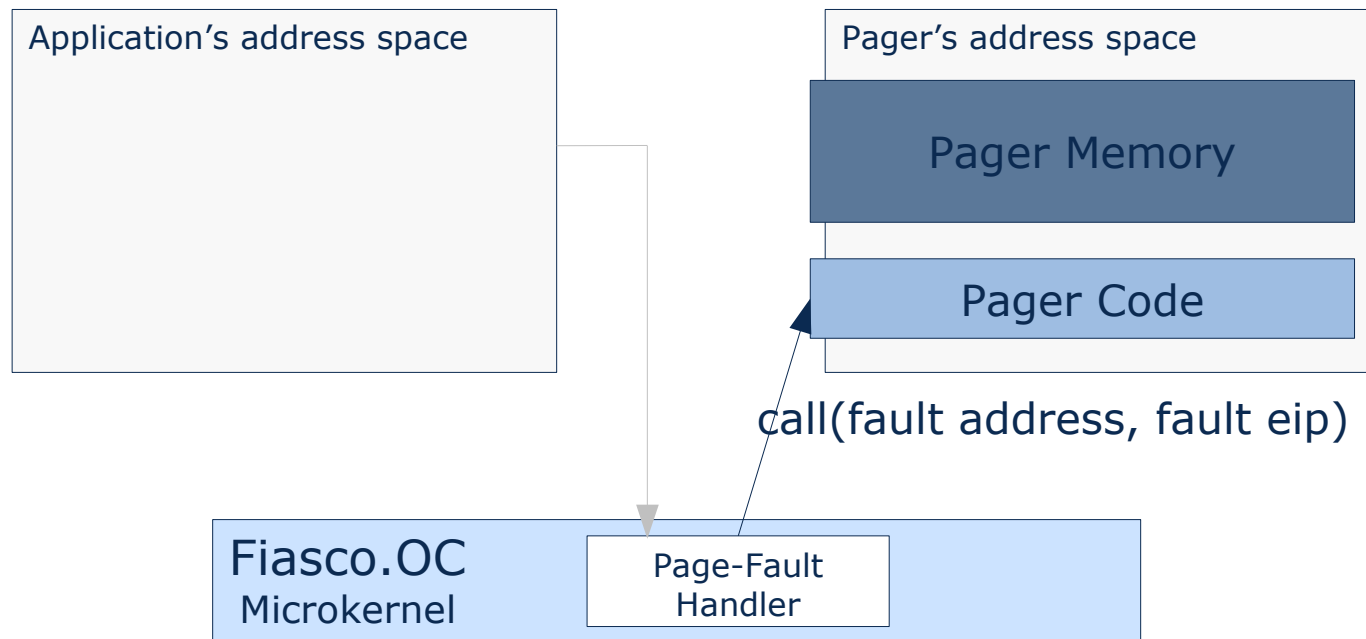
- page-fault exception is caught by kernel page-fault handler
- no management of user memory in kernel
- invoke user-level memory management ⇒ Pager

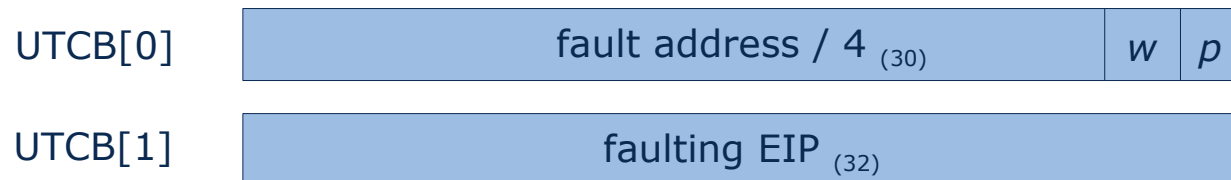
- thread which is invoked on page-fault
- each thread has a (potentially different) pager assigned



Pager invocation

- communication with pager thread by IPC
- kernel page-fault handler sets up IPC to pager
- pager sees faulting thread as sender of IPC

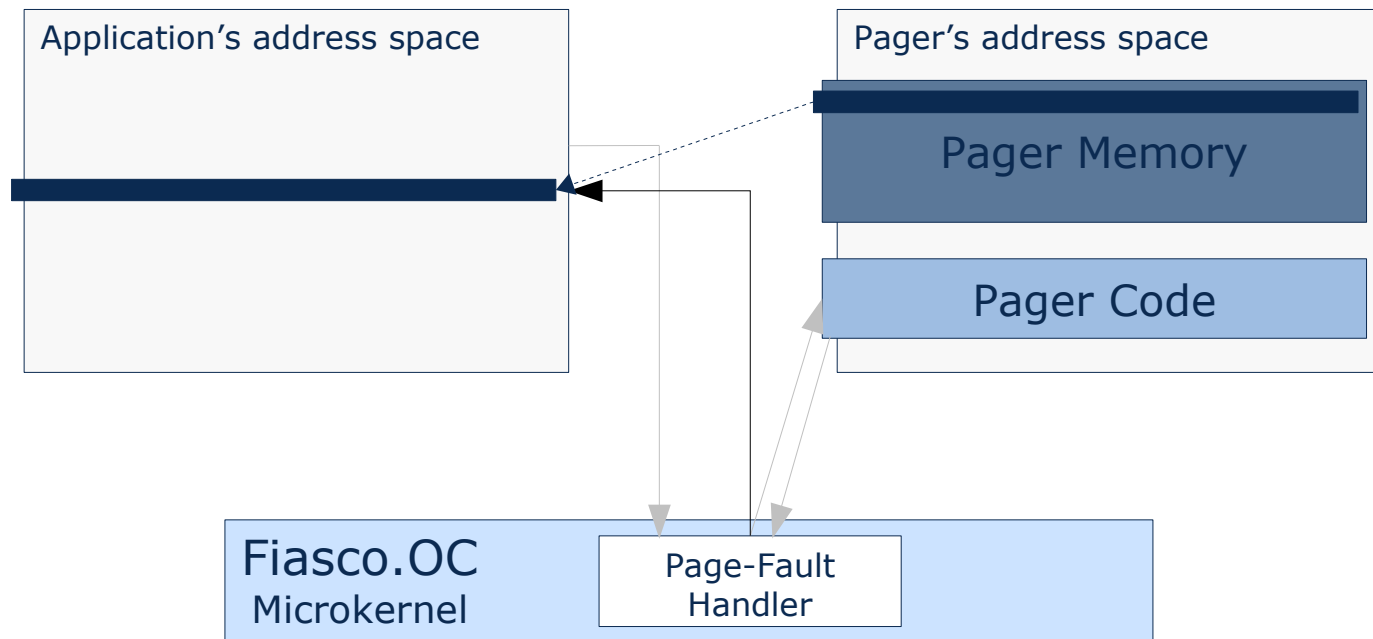




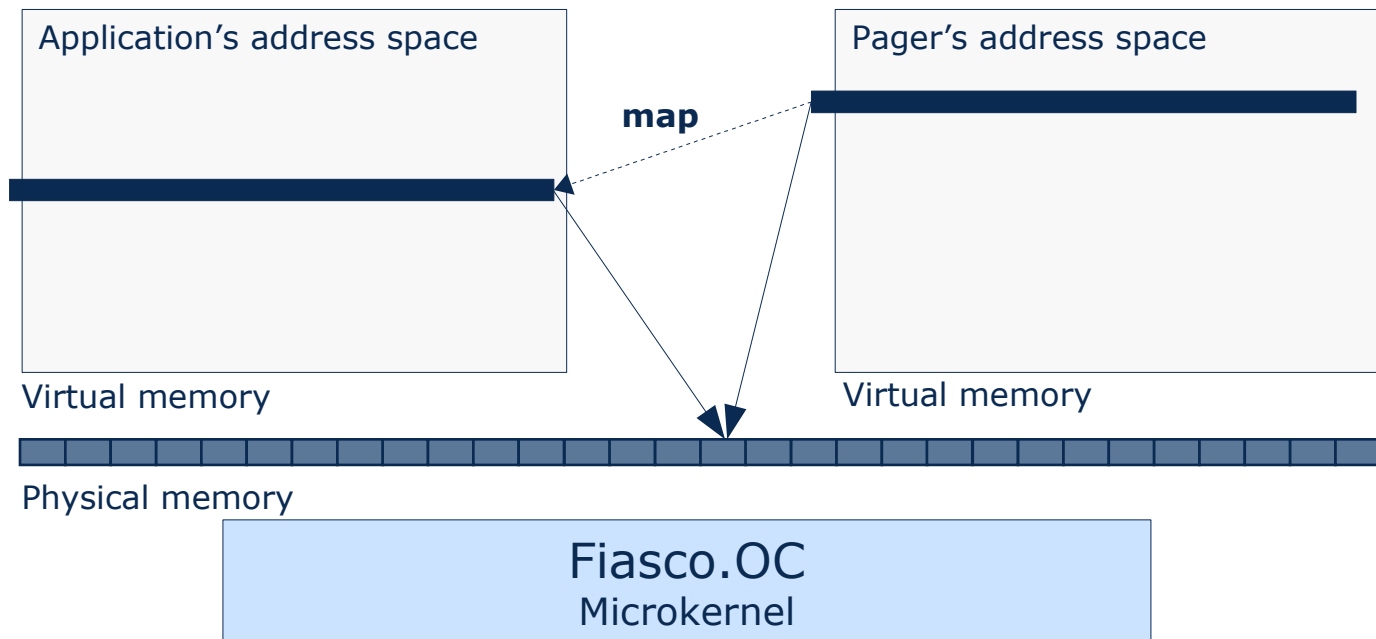
$w = 0$ read page fault
 $w = 1$ write page fault
 $p = 0$ no page present
 $p = 1$ page present

Pager's reply

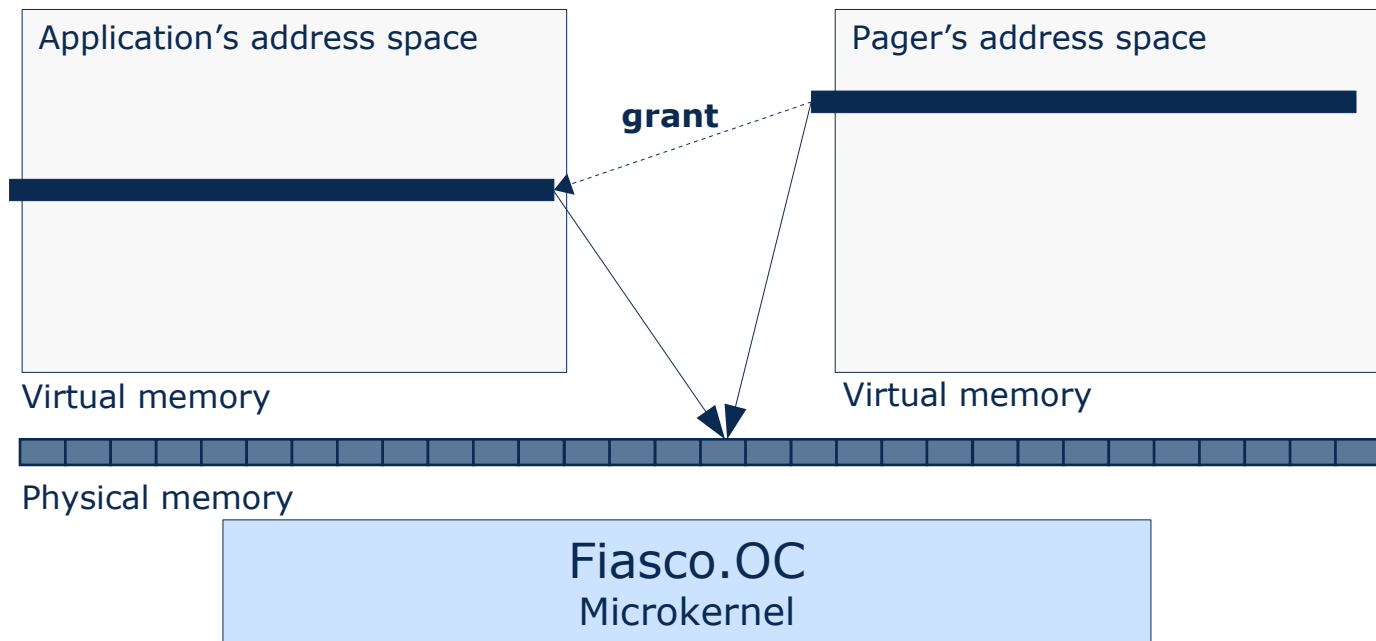
- pager maps pages of it's own address space to the application's address space
- flexpage IPC enables these mappings



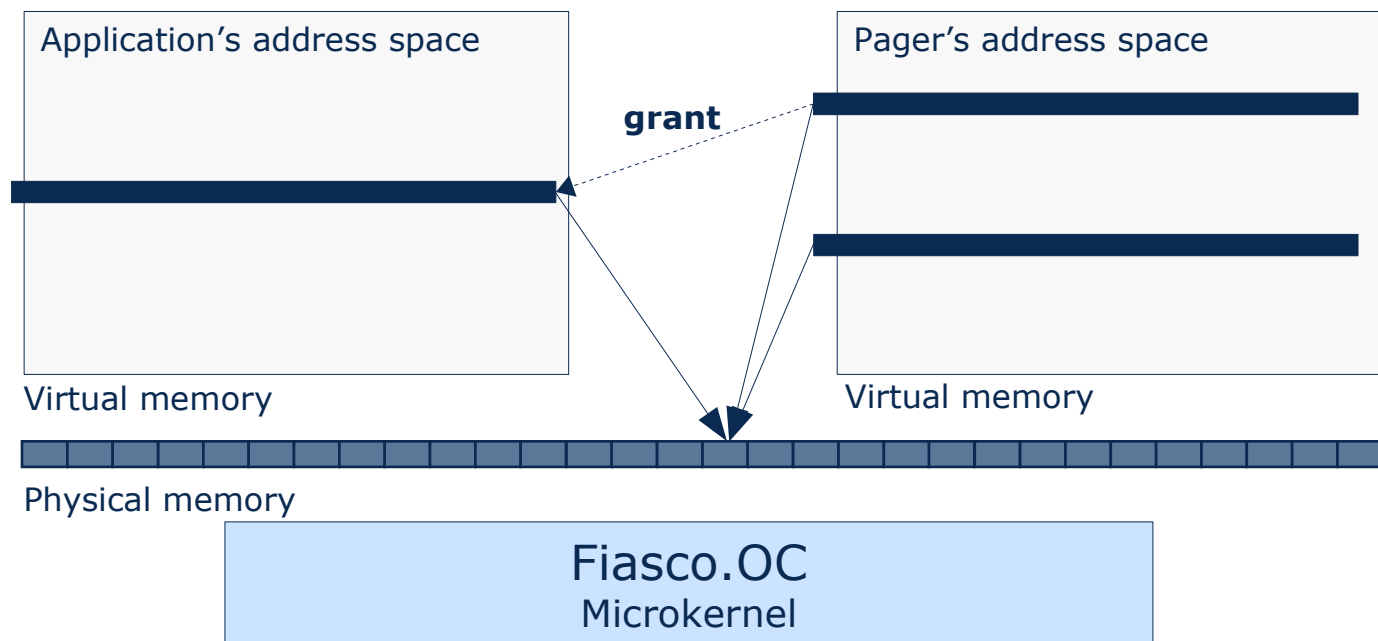
- map creates an entry in the receiver's address space pointing to the same page frame
- only valid pager address space entries can be mapped



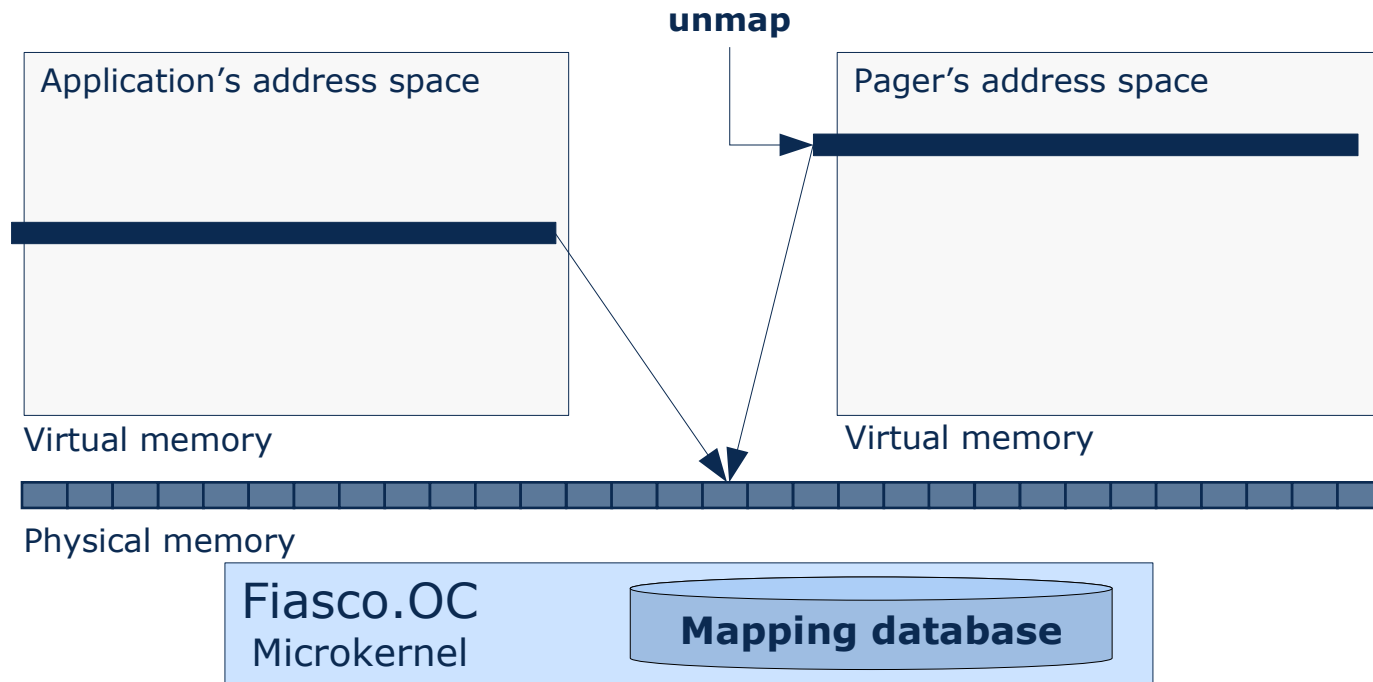
- Special case: grant pages (flag: L4_FPAGE_GRANT)
- Removes mapping from sender's address space



- Special case: grant pages (flag: L4_FPAGE_GRANT)
- Removes mapping from sender's address space
 - **ATTENTION: aliases may remain**



- Unmap removes entries to a page frame (fpage is specified in invoker's address space)
- Kernel needs to track mappings in a database



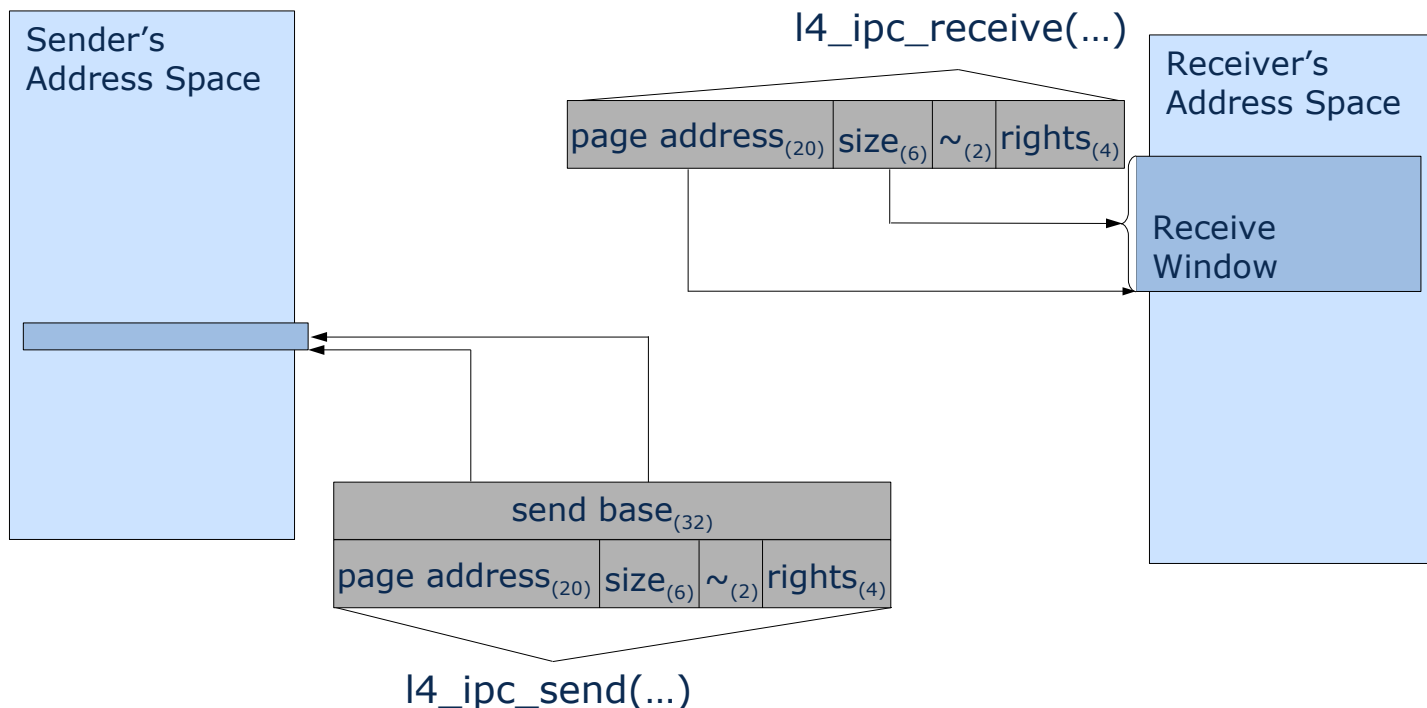


Flexpages

- in general, flexpages represent areas within a name space
- flexpages in Fiasco.OC are used to describe:
 - Memory pages
 - I/O ports
 - Capabilities
- in the following, only flexpages for memory pages are described

Flexpage IPC in detail

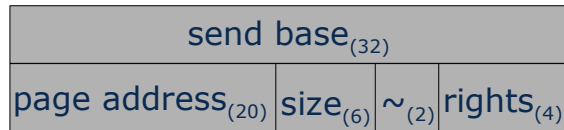
- flexpages are size aligned
- flexpage size 2^{size} , smallest is hardware page
- source and target area of a map IPC are described by flexpages



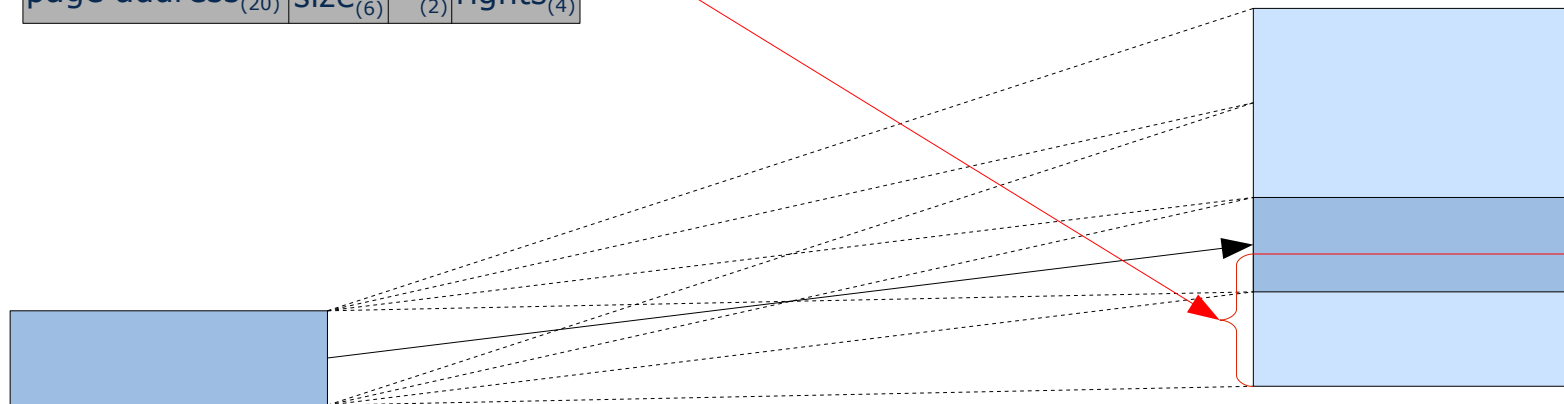
Flexpage offset

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

l4_ipc_send(...)



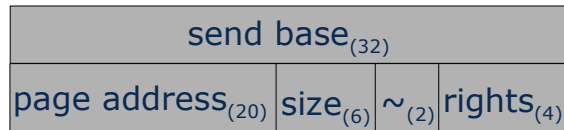
l4_ipc_receive(...)



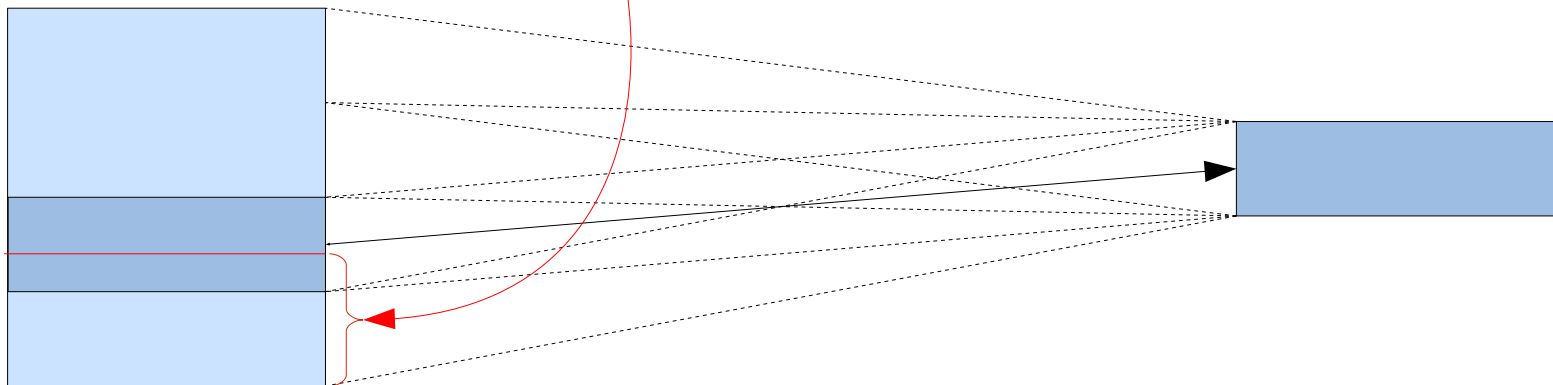
Flexpage offset

- 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(...)

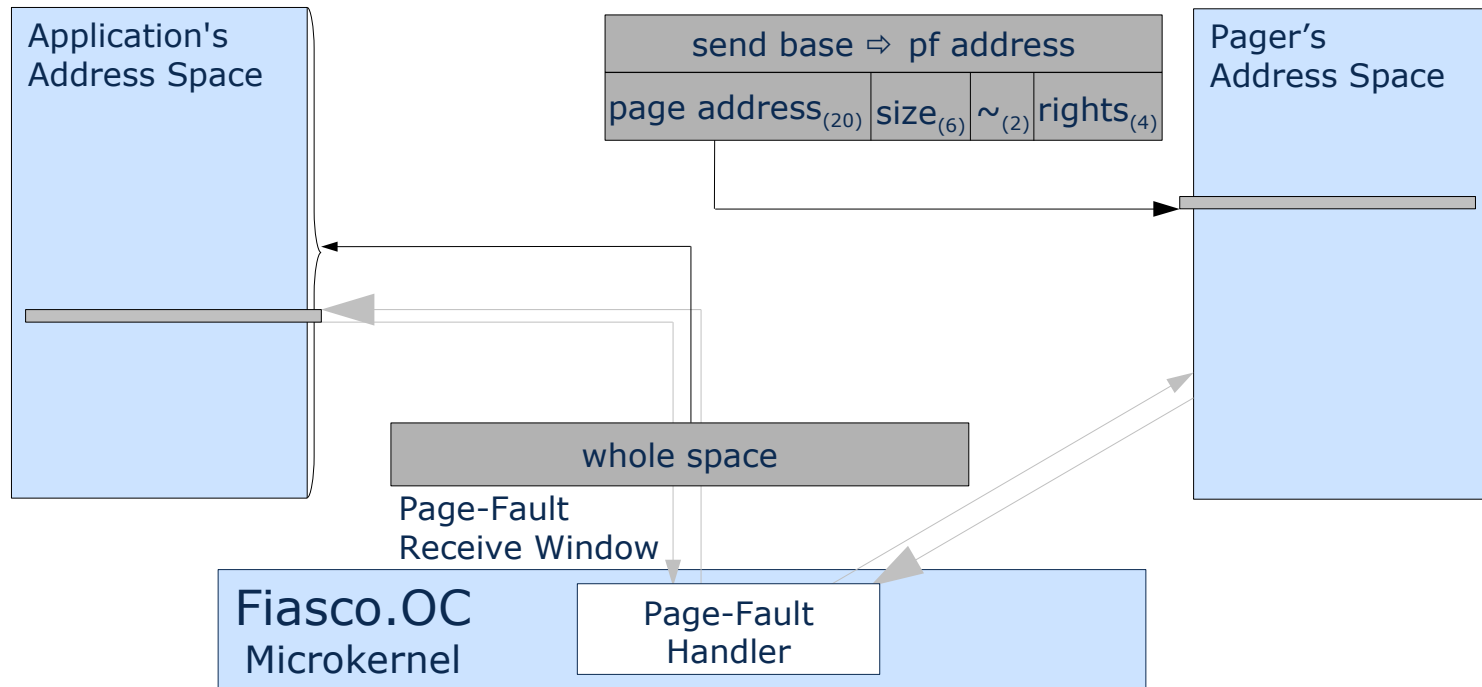


l4_ipc_receive(...)

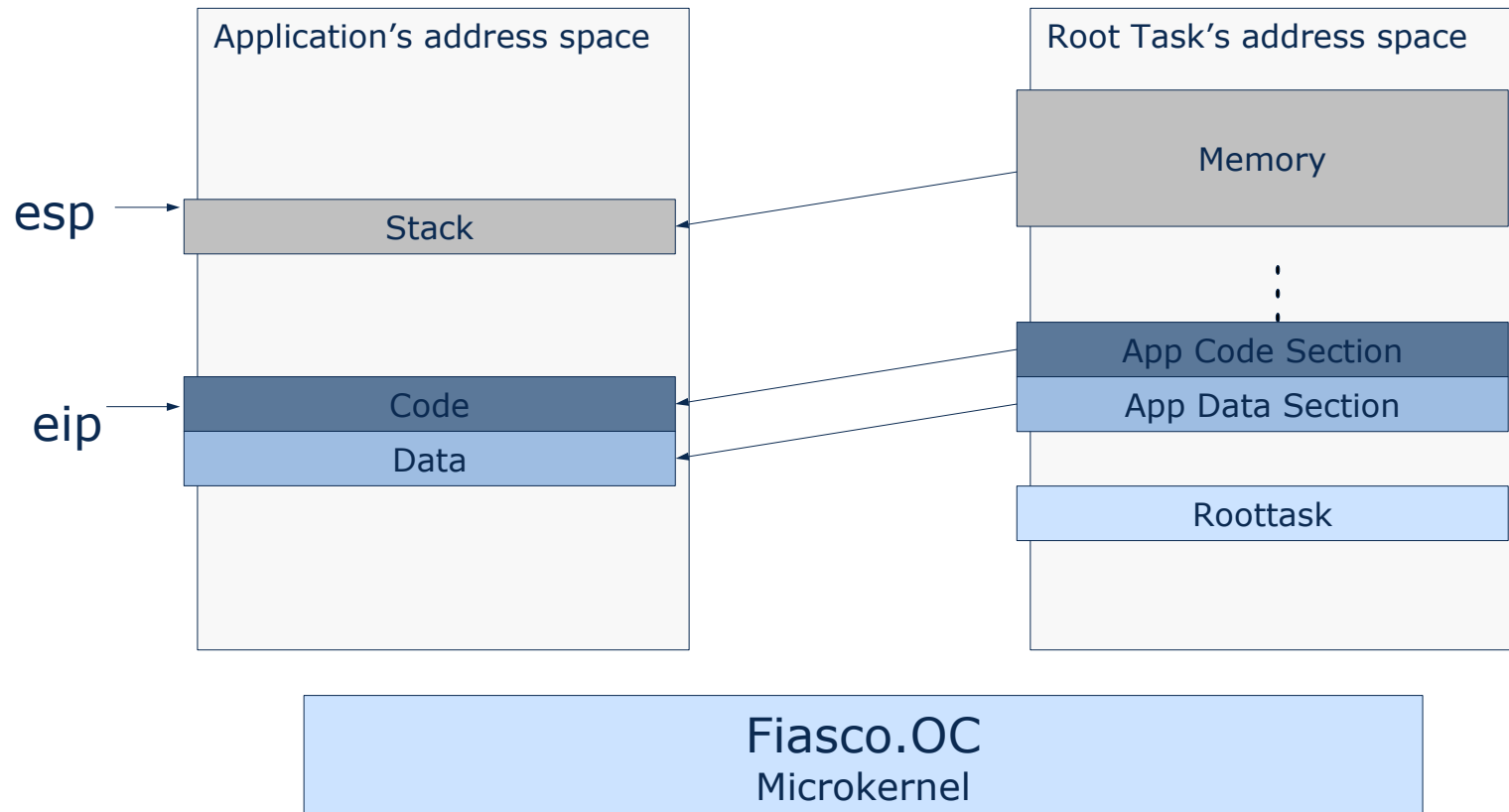


Page-fault in detail

- 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



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

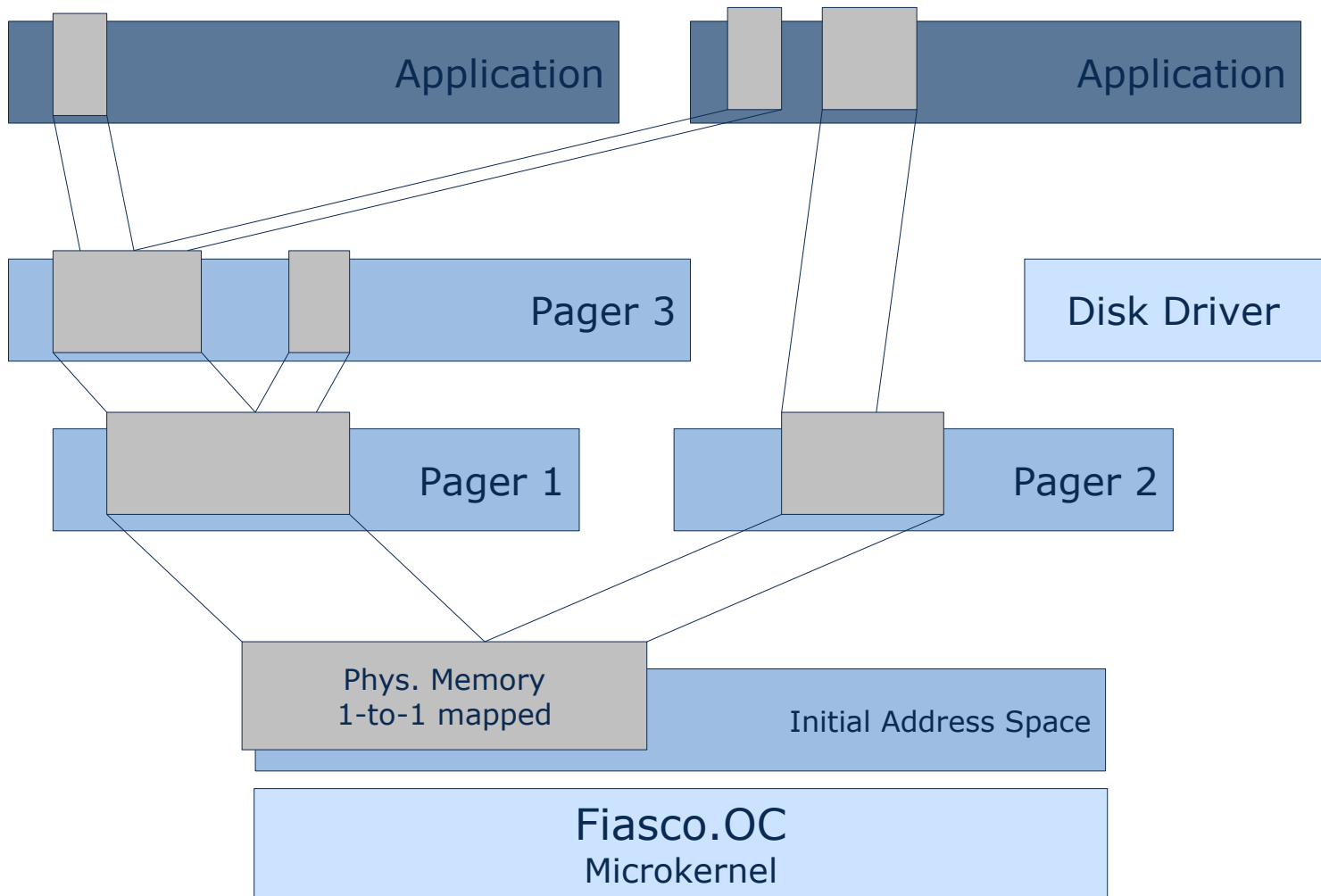




Hierarchical Pagers

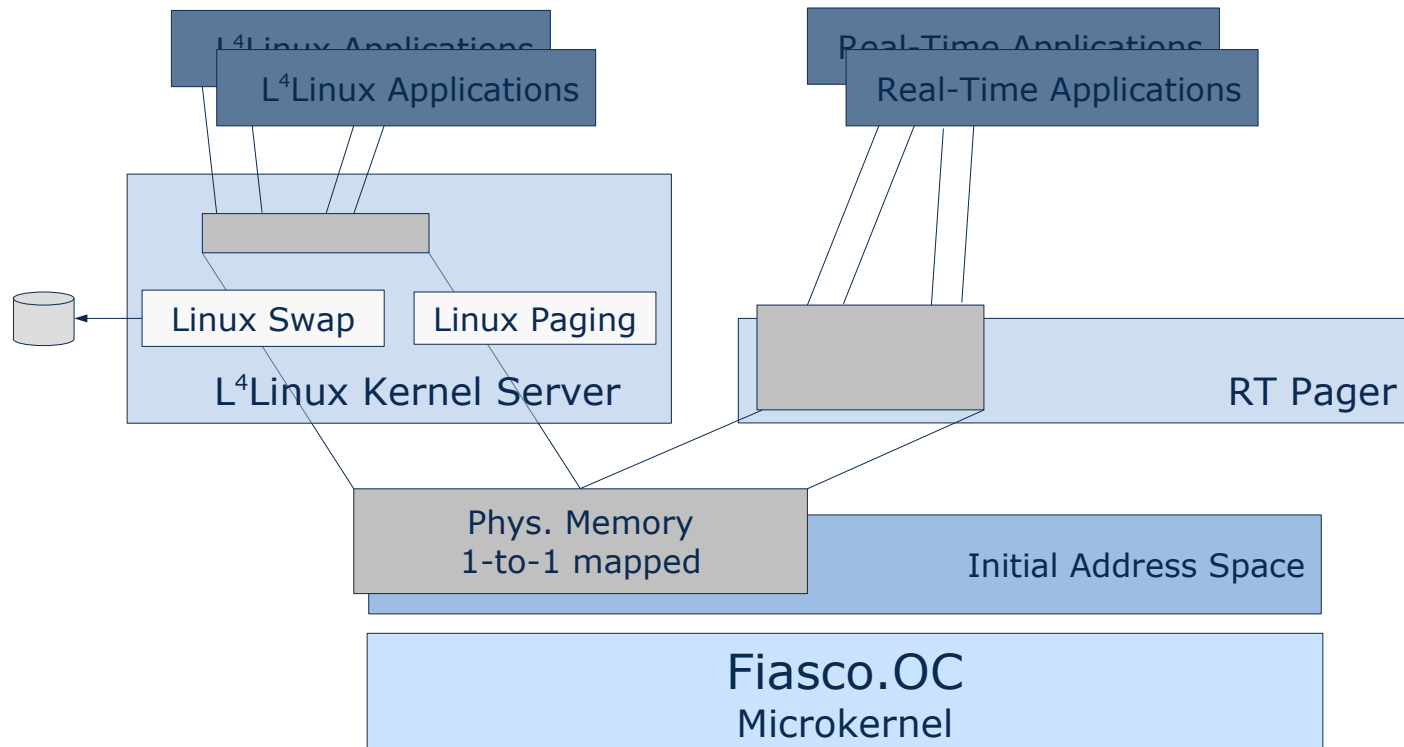
- initial pager can only implement basic memory management
 - no knowledge about application requirements
 - different requirements at the same time
 - missing services for advanced memory management
 - e.g. no disk driver for swapping
 - build more advanced pagers on top of the initial one
- pager hierarchy

Pager hierarchy



Real world example

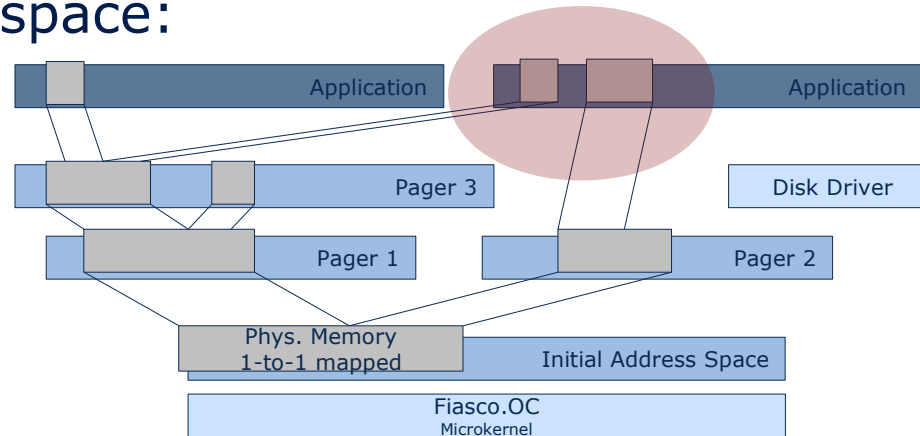
- L⁴Linux implements Linux paging policy
- RT pager implements real-time paging policy (e.g. no swapping)





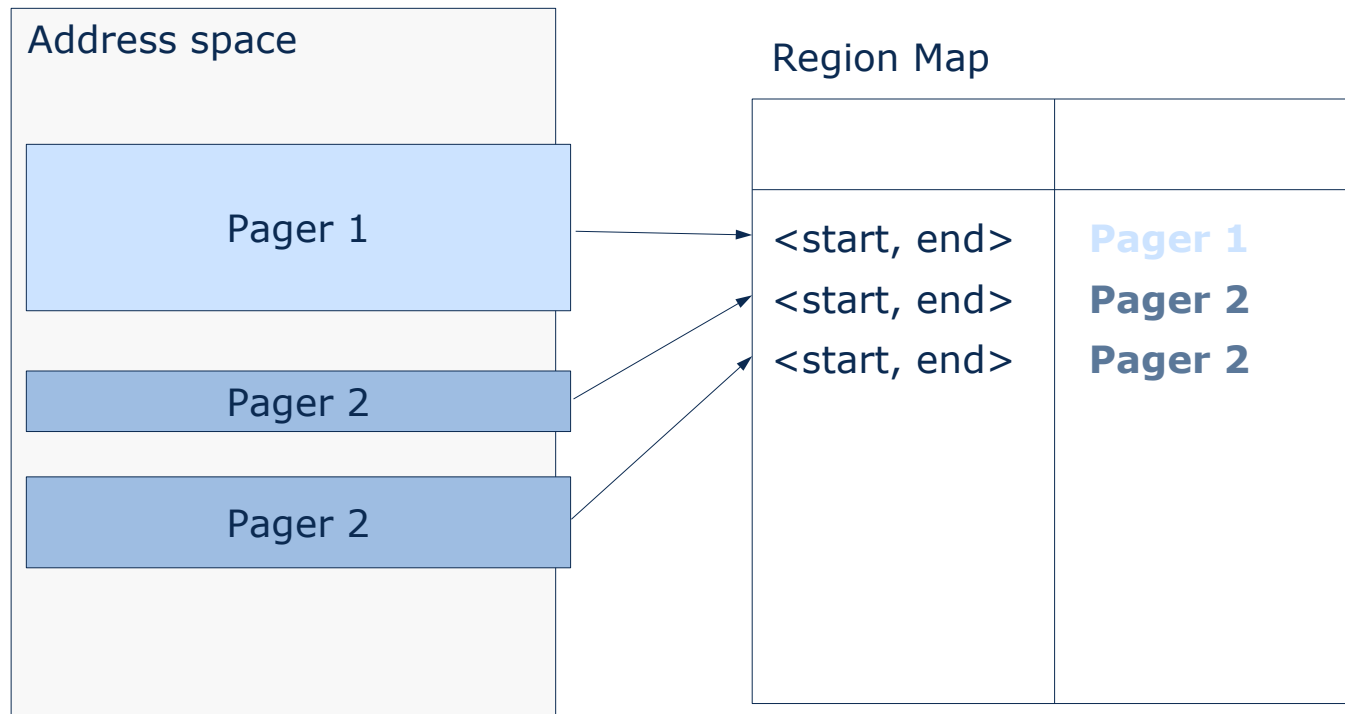
Region Mapper

- pager has to specify send base
- pager needs to know client's address space layout
 - no problems with only one pager (e.g. L⁴Linux)
- possible conflicts if more than one pager manages an address space:



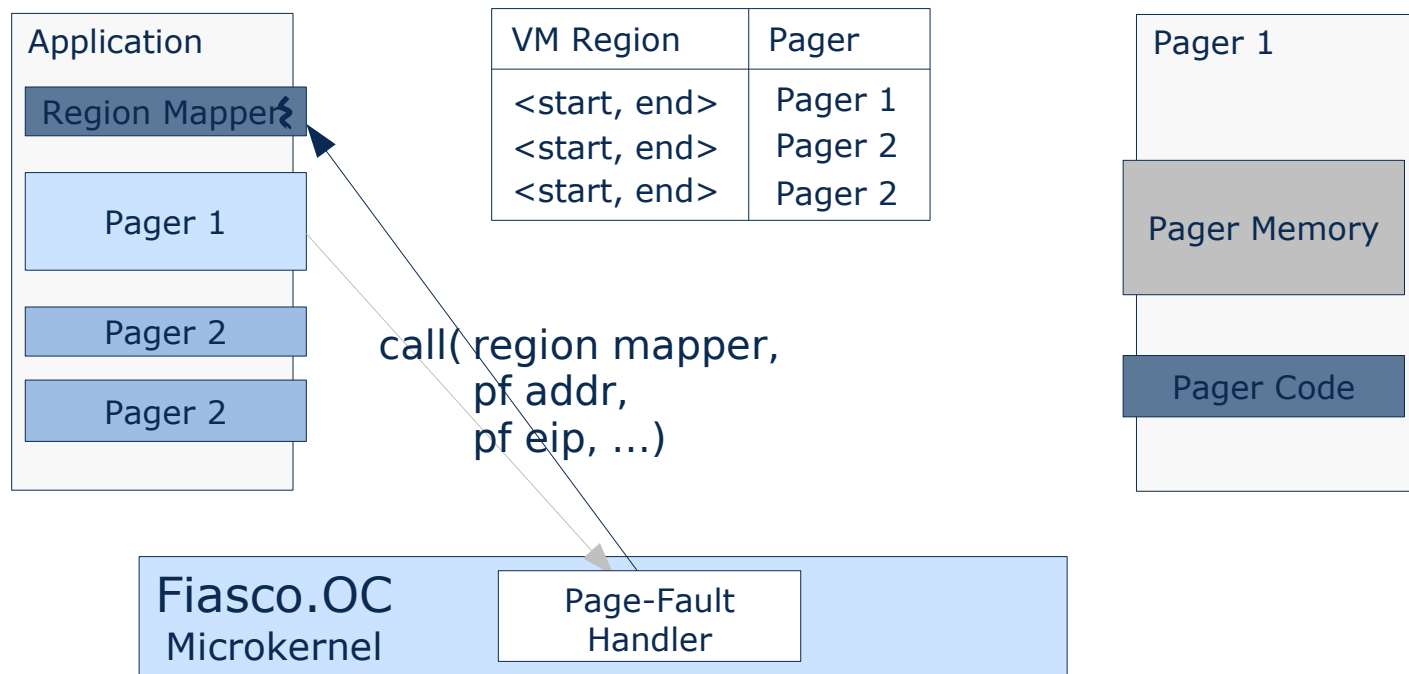
→ virtual memory must be managed independent of pagers

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

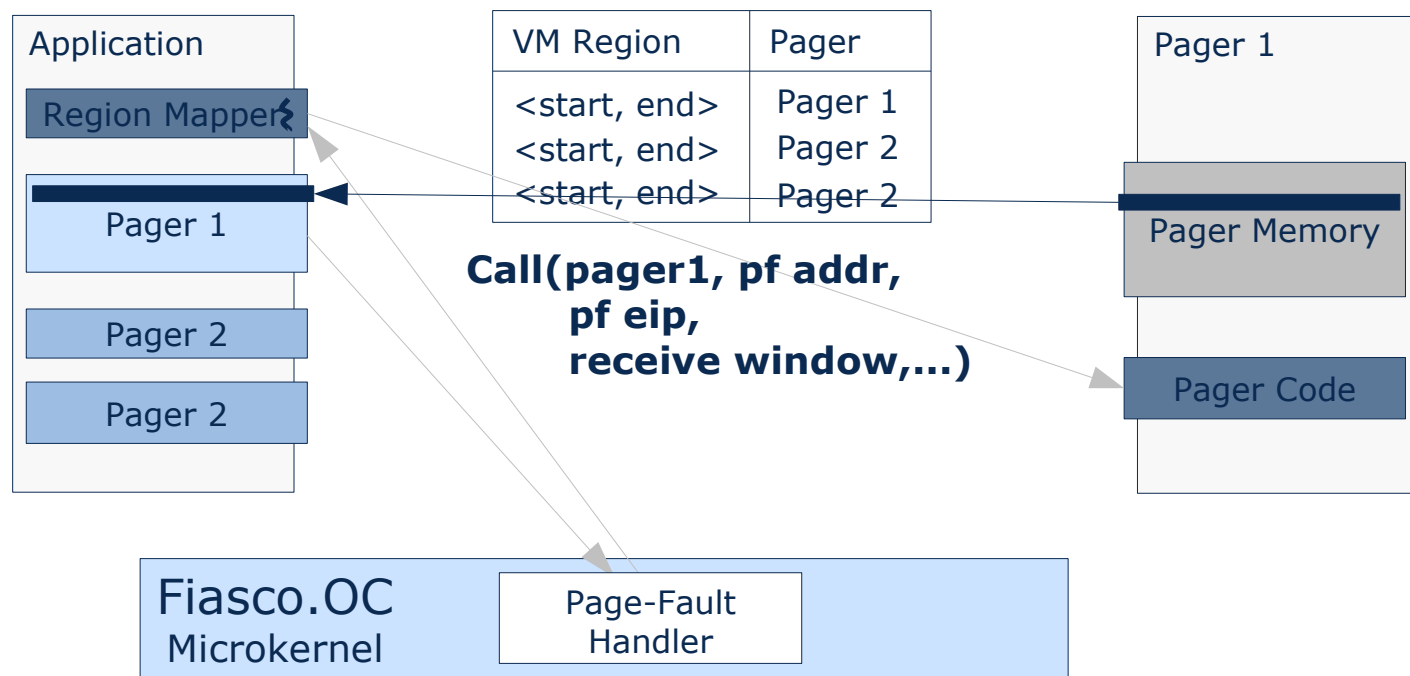


Region mapper

- 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



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



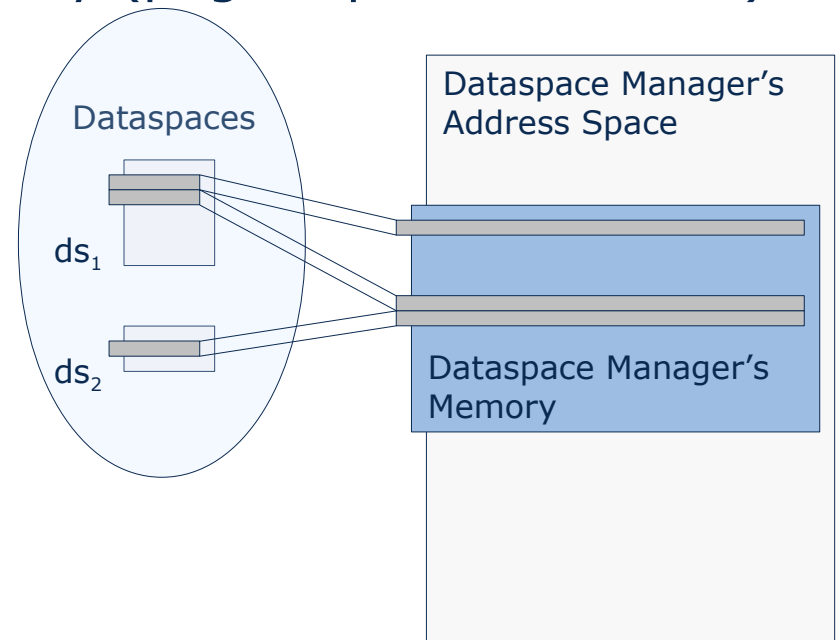


Dataspaces

- 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

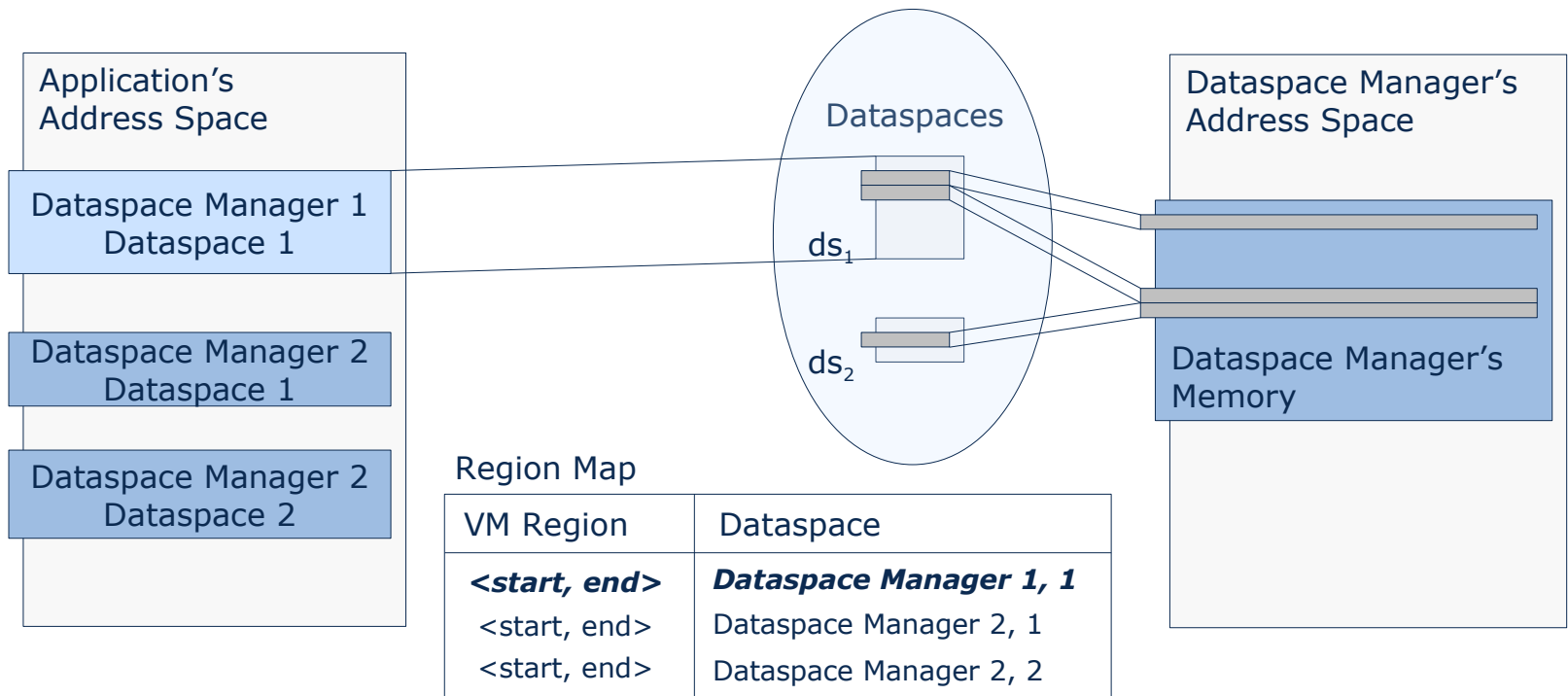
- Dataspace: *unstructured data container*
- abstraction for anything that contains data:
 - Files
 - Anonymous memory
 - I/O adapter memory
 - ...
- Dataspaces are implemented by *Dataspace Managers*
- Dataspaces can be attached to regions of an address space

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



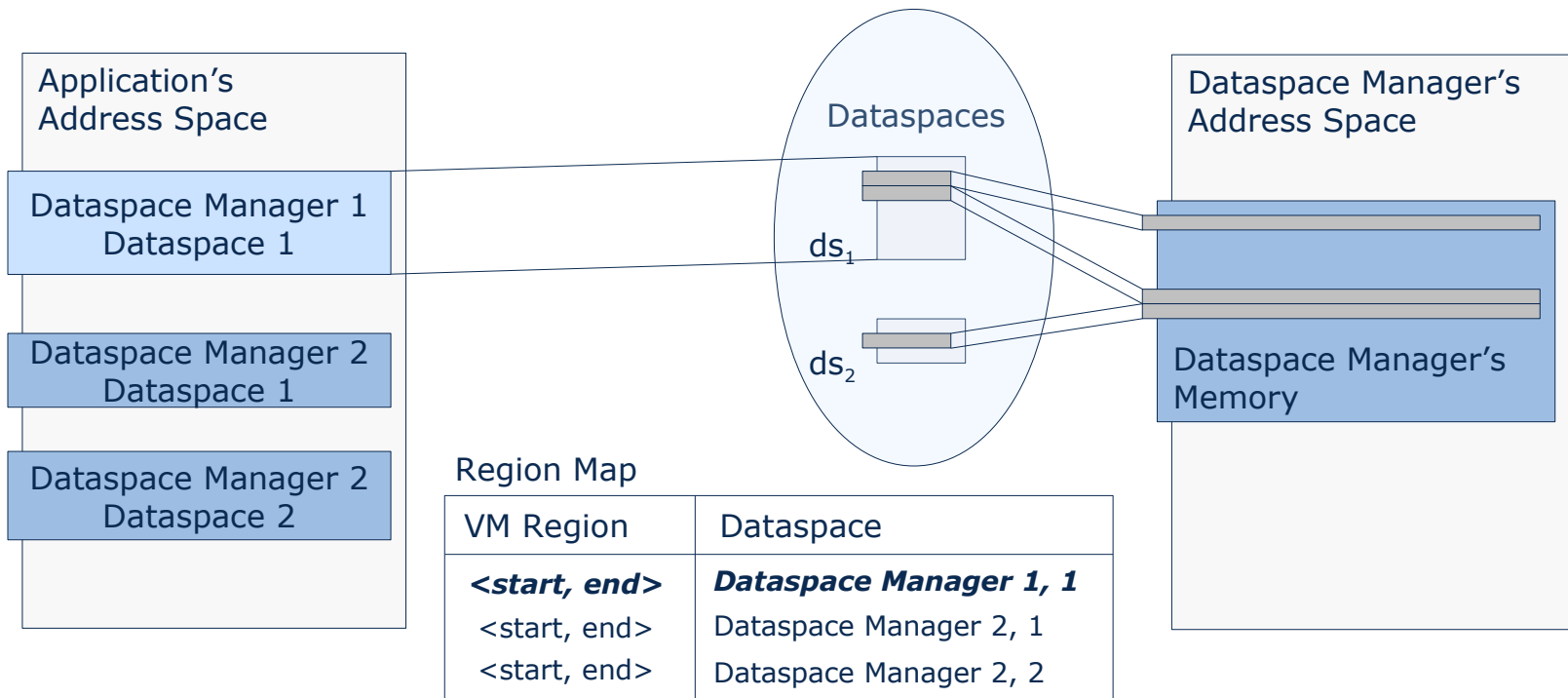
Dataspaces & region mapper

- region map keeps track which dataspaces are attached to which virtual memory regions
- region mapper translates page faults to dataspace offsets



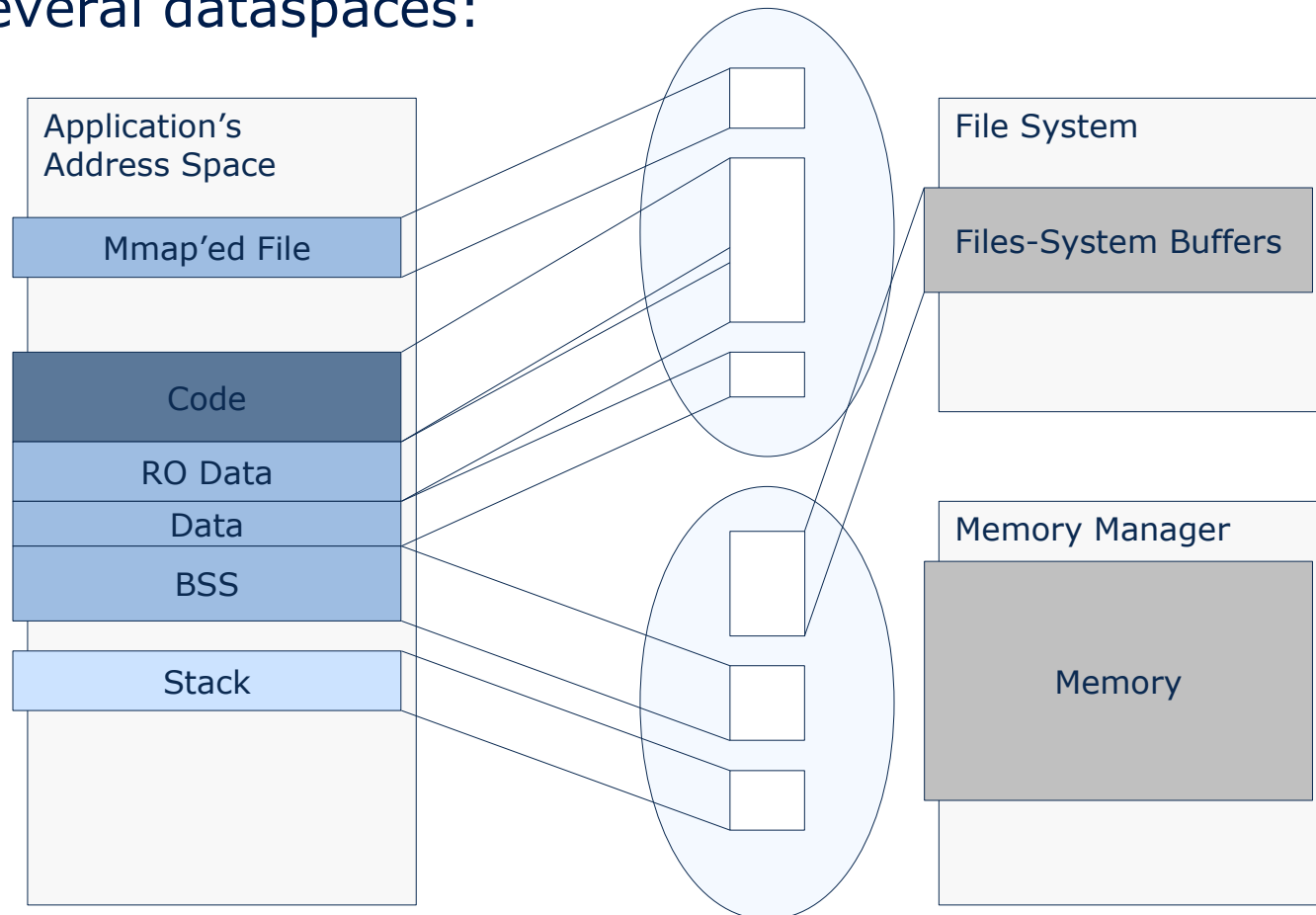
Dataspaces & region mapper

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



- allocate / free dataspace
 - 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

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

- 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](ftp://ftp.cse.unsw.edu.au/pub/users/disy/papers/Aron_PJLED_01.ps.gz)

- next lecture (3.11.) on 'Communication':
 - IPC flavors
 - communication control
- next exercise (10.11.) on:
 - paper reading exercise
 - Brinch Hansen: "The nucleus of a multiprogramming system"
 - **read the paper thoroughly**
 - **be prepared to summarize and discuss it**