# System Architecture

# 16 Memory Management

RAM, Design Space, Examples

January 12 2009
Winter Term 2008/09
Gerd Liefländer

# Recommended Reading

- Bacon, J.: Operating Systems (5)
- Bovet, D.: Understanding the Linux Kernel
- Knuth, D.: The Art of Computer Programming, Vol. 1, Ch. "Dynamic Storage Allocation"

- Nehmer, J.: Grundlagen modener BS, (4)
- Silberschatz, A.: Operating System Concepts (7)
- Stallings, W.: Operating Systems (7)
- Tanenbaum, A.: Modern Operating Systems (4)

# Agenda

- **Motivation**

- **Architecture of RAM Management**

- **Design Parameters**

- **Example Memory Managers**
  - Ring Buffer
  - Stack
  - Boundary Tag Systems
  - Buddy System
  - Slab Allocating
  - Heap Management

# Why Memory Management?

1. Each application needs RAM to run on a CPU, $\Rightarrow$ we have to establish

   - static address regions, e.g. code, global data
   - dynamic address regions, e.g. heap or stack

2. The kernel needs memory for its

   - residentpart
   - loadable kernel modules

3. Devices often only can use physically addressed memory for their buffers etc.

$\Rightarrow$ Every executable needs some RAM
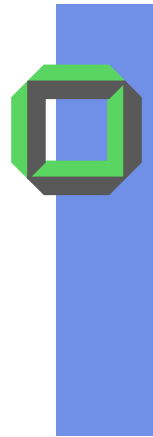
# *Why Memory Management?*

Entities of an "application" AS (address regions):

- Code ("text segment" in Unix jargon)

- Data ((un)initialized data segment)

- Heap

- Stack

Data types (entities) of the kernel AS*:

- Buffer

- TCB, page table, free list, bit map, …

*At the end of the course you should be able to enumerate at least 10 different kernel data types

# System Goals concerning Memory

- Increase (maximize)
  - memory utilization
  - Sometimes we must reserve some RAM capacity for
    - high priority applications or
    - system emergency functions

- Reduce (minimize)
  - application's response time
  - application's turnaround time
  - memory manager's overhead

# *What does a Memory Manager?*

- Keeping track of memory that is currently

  - allocated, i.e. in use or reserved for future use

    - Pinning parts of memory for specific tasks

  - free

- Looking for fitting free memory in case of a request

  - If found, allocate free memory according to some policy

- Free memory in case of a release

  - potentially look for free neighbors in order to reunify free neighbored memory pieces

# Memory Management (1)

- **Programmers want memory being**
  - large & fast & non volatile

- *Current technology* does not support all of these at once, but *future* technologies like MRAM* might do

- System architects offer a *memory hierarchy*
  - small high-speed caches (expensive)
  - medium sized fast main memory (RAM)
  - flash memory
  - Giga bytes of slow, cheap disk storage
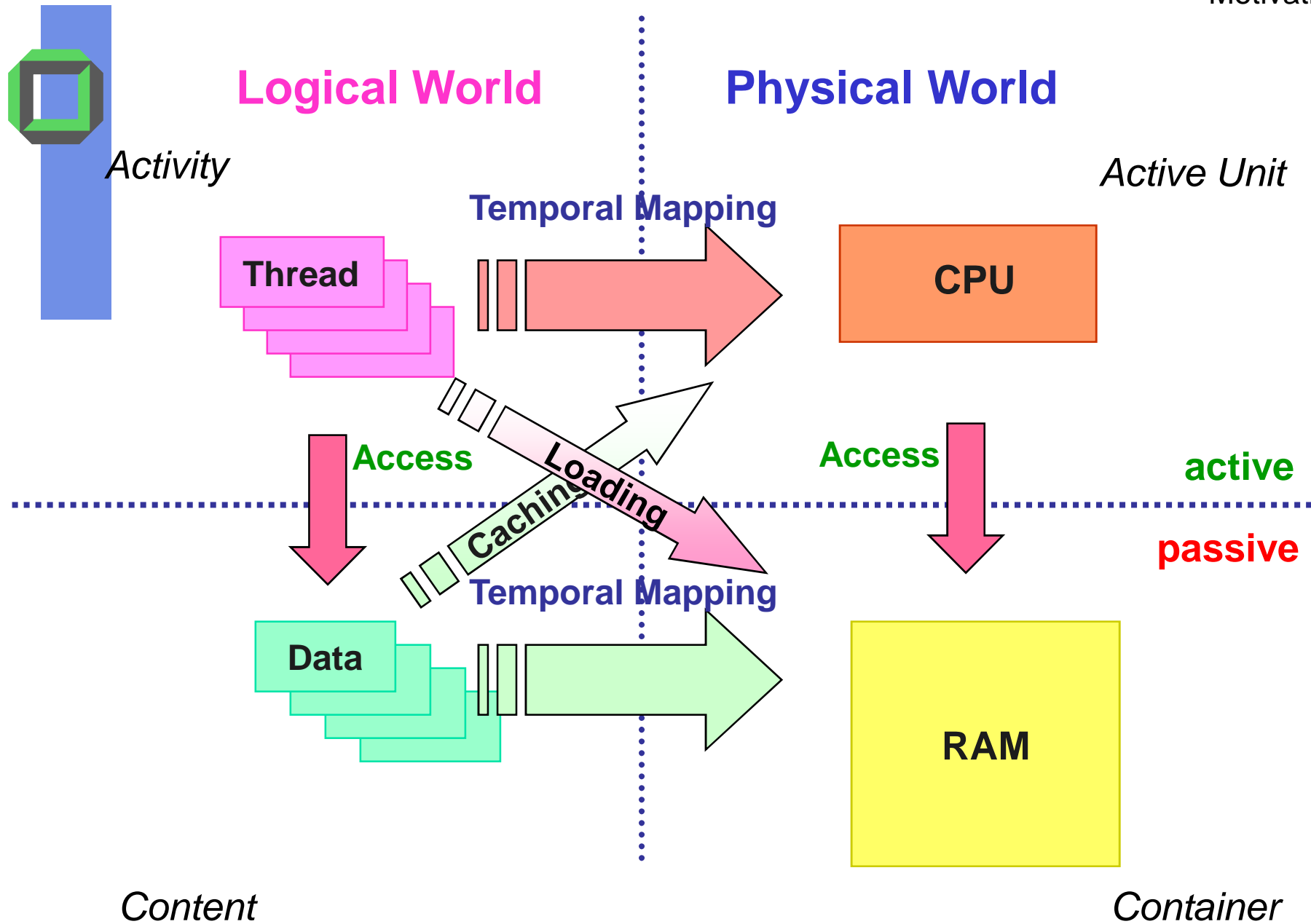  - Terra bytes of very slow archive memory

*MRAM = Magnetic RAM

# Memory Management (2)

Two <u>main goals:</u>

1. Manage memory efficiently

   - appropriate algorithms & data structures
   - program MMU (e.g. TLB)

2. Establish effective usability of memory

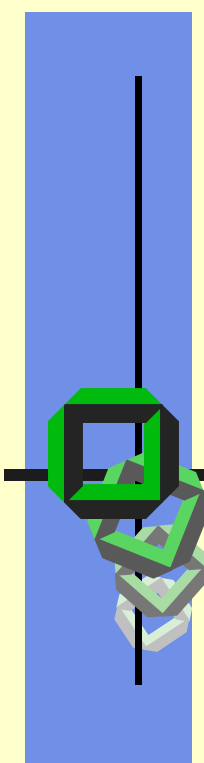   - maximize usage of (main) memory
   - support low cache footprint

**Logical World**  **Physical World**

*Activity*  *Active Unit*

**Temporal Mapping**

**Thread**  →  **CPU**

**Access**  *Loading*  **Access**  **active**

*Caching*  **passive**

**Temporal Mapping**

**Data**  →  **RAM**

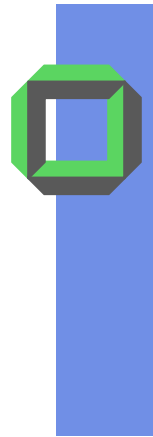*Content*  *Container*

# Why to bother about RAM?

**Memory is large <u>and</u> cheap, and if not, we'll use virtual memory.**

However, reality tells us:

1. Many computers don't use virtual memory at all

2. Modern programs tend to be memory greedy & virtual memory ≠ unlimited memory

3. Memory management has to be done anyway at
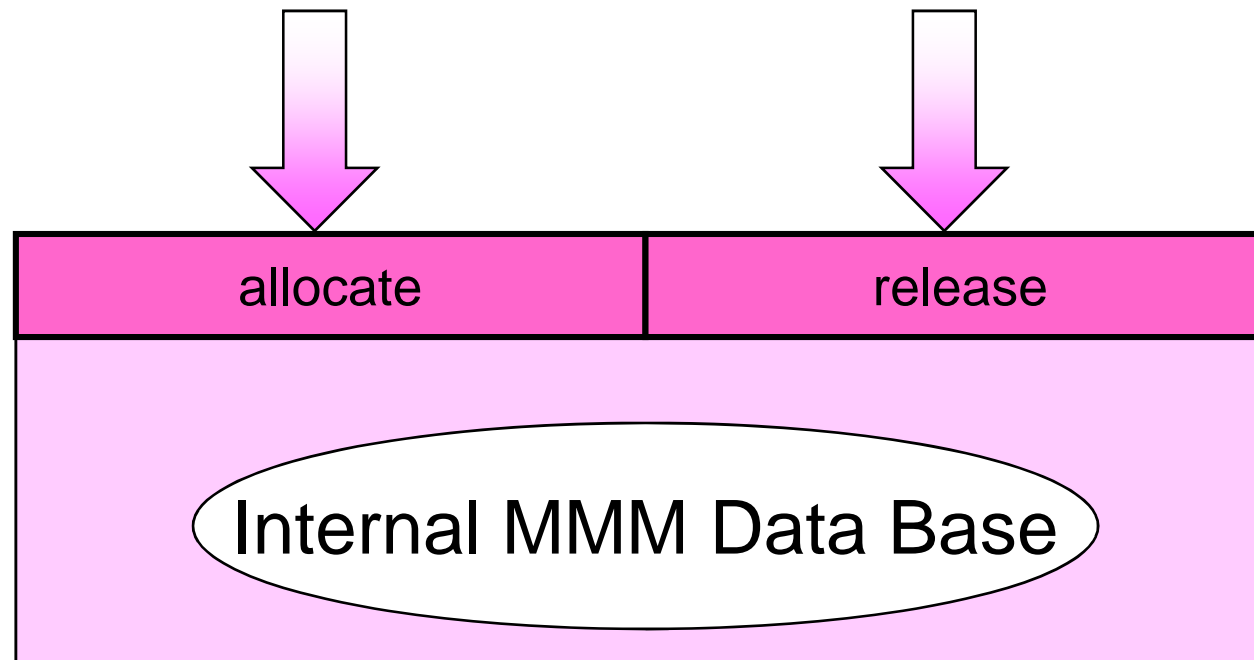   - system level
   - application level

# Architecture of RAM Management

# Memory Management Module

MMM-Interface



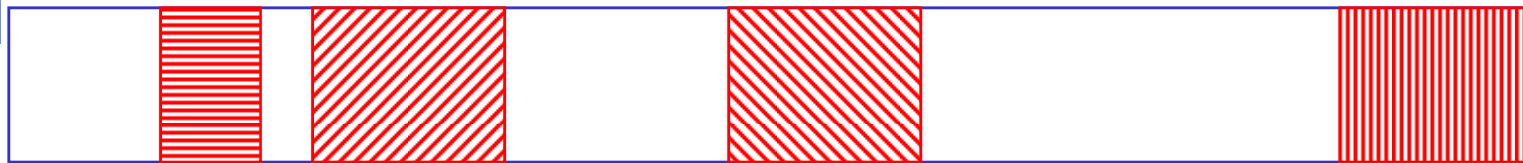| allocate | release |
|----------|---------|

Internal MMM Data Base

## Note:

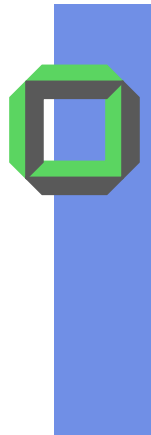∃ several related and/or unrelated memory managers

# MMM Data Base

MMMDB reflects what parts are allocated and what parts are free
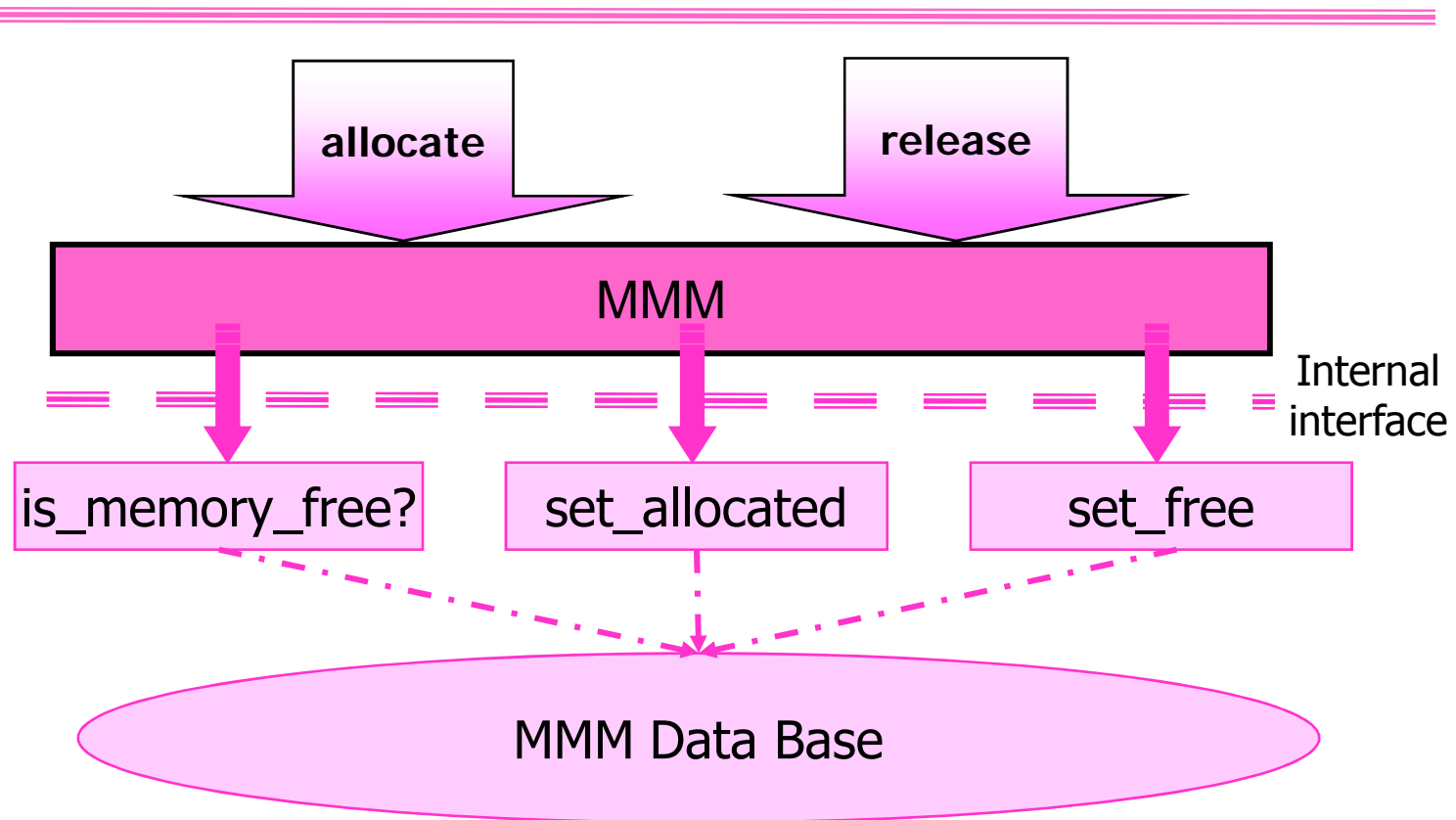


We must solve the following problems when handling memory requests

- **Efficiently** select a fitting free memory block
  - As fast as possible

- Try to utilize the memory efficiently
  - Avoid unusable memory leaks

- Meet additional constraints
  - Avoid unbounded waiting in front of MMM

# Smart System Architecture of a MMM

MMM-Interface

allocate release

**MMM**

Internal interface

is_memory_free? set_allocated set_free

MMM Data Base

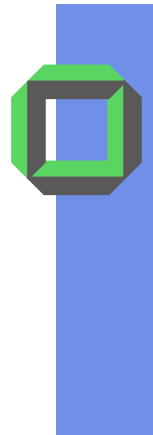*Main advantage of the above architecture?*

# Design Parameters
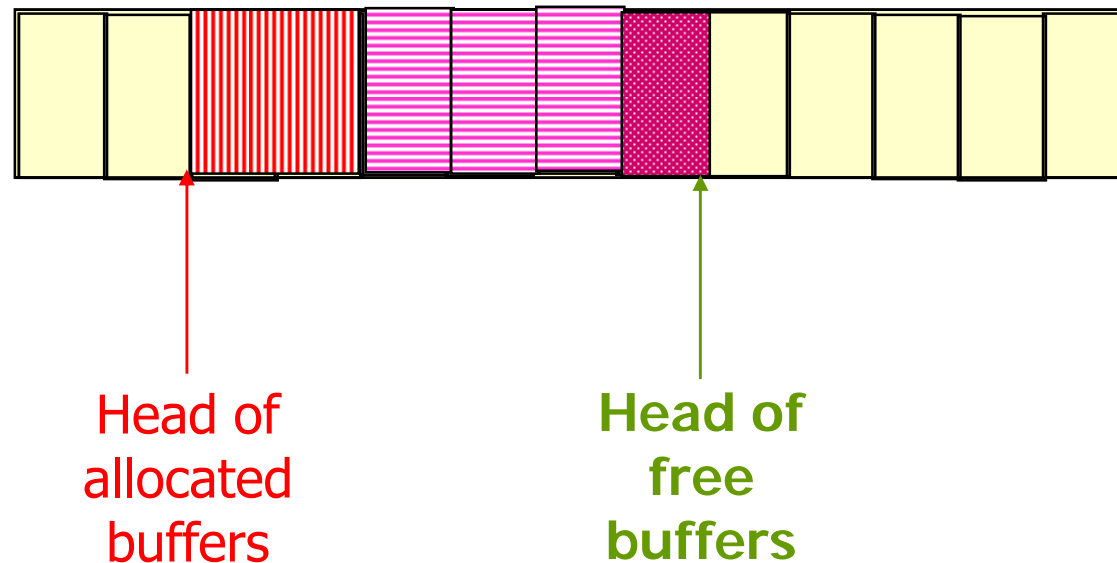
# Orthogonal Parameters

- Sequence of allocate/release-operations

- Size of memory blocks

- Data structures

- *Fragmentation (not design parameter but result of a design!!!!)*

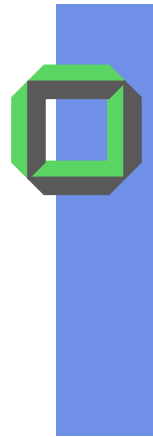- Allocation policy

- Reunification of released blocks

- …?

Playground for an innovative system architect

# MM Design Parameters (1a)

- ## Sequence of allocate/release-operations
  - ~FIFO = queue



Head of
allocated
buffers

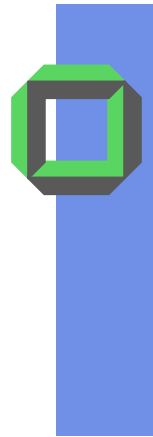**Head of
free
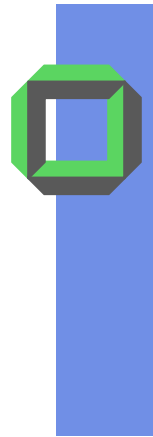buffers**

# MM Design Parameters (1b)

- Sequence of allocate/release-operations
  - FIFO, LIFO = stack

# MM Design Parameters (1c)

- Sequence of allocate/release-operations
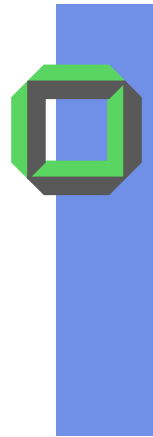  - FIFO, LIFO, arbitrary = in general

*Additionally we have to solve the problem of memory holes*

# MM Design Parameters (2)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary
- Size of memory blocks

# MM Design Parameters (2a)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary
- Size of memory blocks
  - Constant size = most buffering

# MM Design Parameters (2b)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary
- **Size of memory blocks**
  - Constant size, multiple of fixed size

# MM Design Parameters (2c)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary

- **Size of memory blocks**
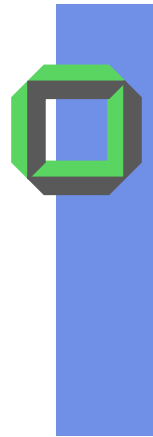  - Constant size, multiple size, fixed size = reservoir of frequently used blocks

# MM Design Parameters (2d)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary

- Size of memory blocks
  - Identical size, multiple size, fixed size, exponential size = buddy system

# MM Design Parameters (2e)

- Sequence of allocate/release-operations
    - FIFO, LIFO, arbitrary

- Size of Memory blocks
    - Identical size, multiple size fixed size, exponential size, arbitrary

# MM Design Parameters (3a)

- Sequence of allocate/release-operations
    - FIFO, LIFO, arbitrary
- Size of memory blocks
    - Identical size, multiple size, fixed size, exponential size, arbitrary
- **Data structures**
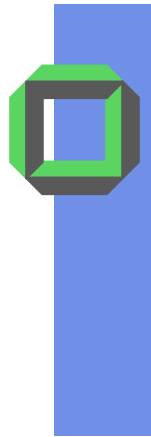    - Integrated within memory block(s) = mostly with larger blocks

# MM Design Parameters (3b)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary

- Size of memory blocks
  - Identical size, multiple size, fixed size, exponential size, arbitrary

- **Data structures**
  - Integrated, special memory block(s) = stack

# Result of Design (4)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary
- Size of memory blocks
  - Identical size, multple size, fixed size, exponential size, arbitrary
- Data Structures
  - Integrated, special memory block(s)
- **Fragmentation**
  - With(out) in(ex)ternal fragmentation

# MM Design Parameters (5a)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary

- Size of memory blocks
  - Identical size, multiple size, fixed size, exponential size, arbitrary

- Data structures
  - Integrated, special memory block(s)

- Fragmentation
  - With(out) in(ex)ternal fragmentation

- **Allocation policy**
  - First-Fit = take the **first fitting block** within an **ordered set of free blocks**

*How would you order the set of free blocks?*

# MM Design Parameters (5b)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary
- Size of memory blocks
  - Identical size, multiple size, fixed size, exponential sized, arbitrary
- Data structures
  - Integrated, special memory block(s)
- Fragmentation
  - With(out) in(ex)ternal fragmentation
- **Allocation policy**
  - First-, Next-Fit = take the next fitting block within an ordered set of free blocks

# MM Design Parameters (5c)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary

- Size of memory blocks
  - Identical size, multiple size, fixed size, exponential sized, arbitrary

- Data structures
  - Integrated, special memory block(s)

- Fragmentation
  - With(out) in(ex)ternal fragmentation

- **Allocation policy**
  - First-, Next-, BestFit = take the best fitting block within an ordered set of free blocks

# MM Design Parameters (5d)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary

- Size of memory blocks
  - Identical size, multiple size, fixed size, exponential sized, arbitrary

- Data structures
  - Integrated, special memory block(s)

- Fragmentation
  - With(out) in(ex)ternal fragmentation

- **Allocation policy**
  - First-, Next-, Best-, Worst-Fit = take the largest fitting block within an ordered set of free blocks

# MM Design Parameters (5e)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary

- Size of memory blocks
  - Identical size, multiple size, fixed size, exponential sized, arbitrary

- Data structures
  - Integrated, special memory block(s)

- Fragmentation
  - With(out) in(ex)ternal fragmentation

- **Allocation policy**
  - First-, Next-, Best-, Worst-, Nearest-Fit, = take the fitting block closest to the previous fitting block

# MM Design Parameters (6a)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary
- Size of memory blocks
  - Identical size, multiple size, fixed size, exponential size, arbitrary
- Data structures
  - Integrated, special memory block(s)
- Fragmentation
  - With(out) in(ex)ternal fragmentation
- **Allocation policy**
  - First-, Next-, Best-, Worst-, Nearest-Fit
- **Reunification of released blocks**
  - **Eager reunification** with neighbored free blocks (if any)

# MM Design Parameters (6b)

- Sequence of allocate/release-operations
  - FIFO, LIFO, arbitrary
- Size of memory blocks
  - Identical size, multiple size, fixed size, exponential size, arbitrary
- Data structures
  - Integrated, special memory block(s)
- Fragmentation
  - With(out) in(ex)ternal fragmentation
- Allocation policy
  - First-, Next-, Best-, Worst-, ..., Nearest-Fit
- **Reunification of Released blocks**
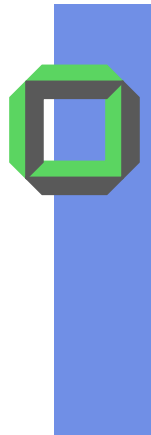  - Eager, **lazy reunification** = wait a while until you reunify

# Summary of Designing MM

- **Sequence of allocate/release-operations**
  - FIFO, LIFO, arbitrary

- **Size of memory blocks**
  - Identical size, multiple size, fixed size, exponential size, arbitrary

- **Data structures**
  - Integrated, special memory block(s)

- **Fragmentation**
  - With(out) internal or external fragmentation

- **Allocation policy**
  - First-, Next-, Best-, Worst-, …, Nearest-Fit

- **Reunification of Released blocks**
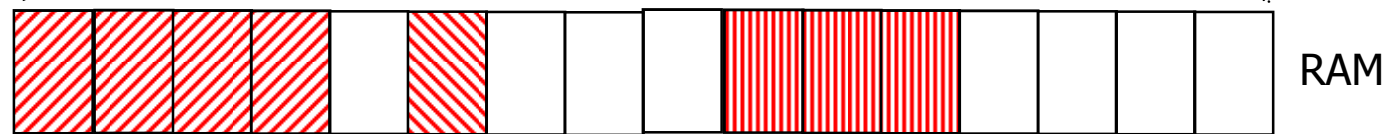  - Immediate, lazy reunification

# Data Structures

- **Bit Map (for fixed sized units, e.g. pages)**
  - Extra data structure
  - Integrated

- **Table/List (for arbitrary sized units, e.g. segments)**
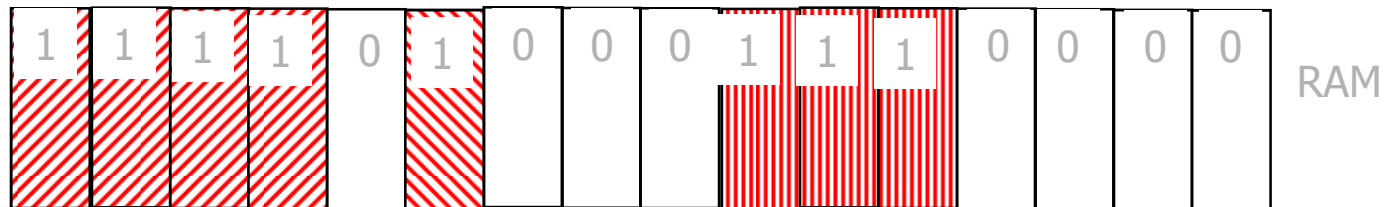  - Extra data structure
  - Integrated

# Bit Map Data Base

Extra Bit Map:

1111010001110000  BITMAP

RAM

Overhead per bit map:
The smaller the memory units, the larger the bit map

Integrated Bit Map:

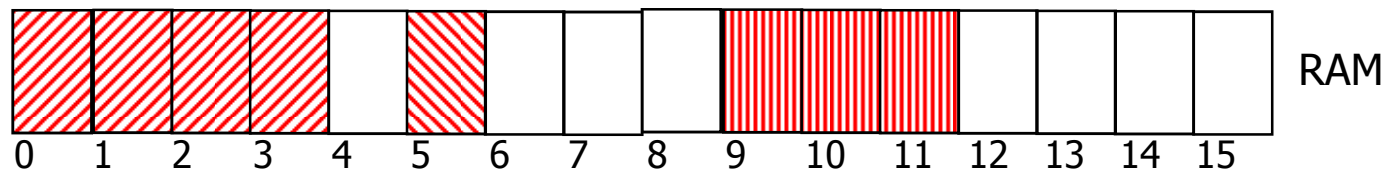1 1 1 1 0 1 0 0 0 1 1 1 0 0 0 0

RAM

# Table Data Base

Extra Tables (sorted according to addresses):

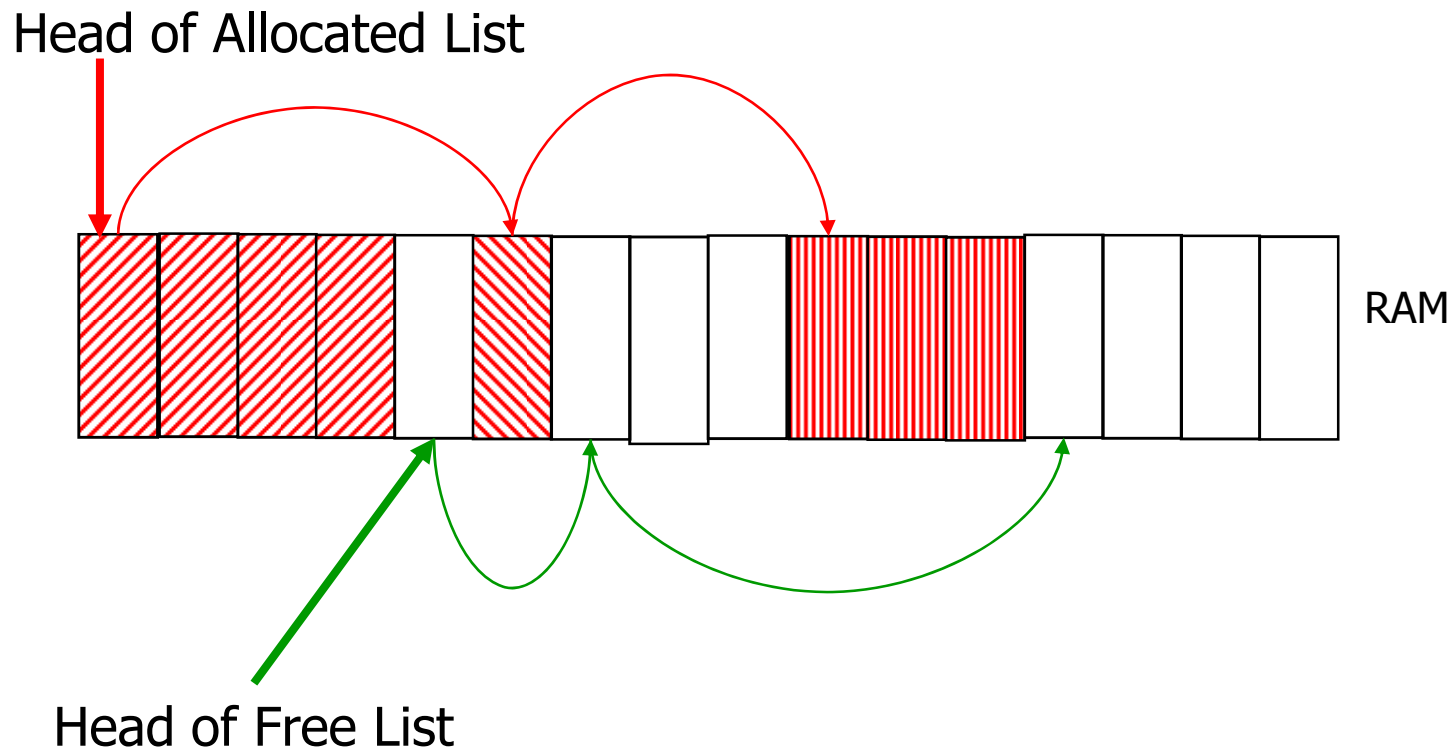| Free Table: | | | Allocate Table: | |
| --- | --- | --- | --- | --- |
| Address | Length | | Address | Length |
| 4 | 1 | | 0 | 4 |
| 6 | 3 | | 5 | 1 |
| 12 | 4 | | 9 | 3 |

RAM

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

∃ *other useful sorting criteria?*

# List oriented Data Base

Integrated List (sorted according to addresses):

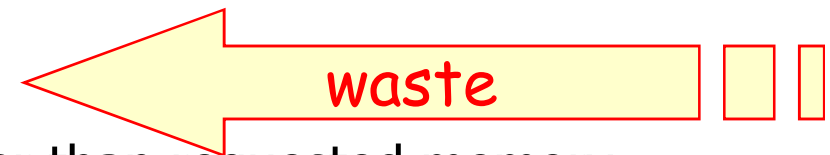Head of Allocated List

RAM

Head of Free List
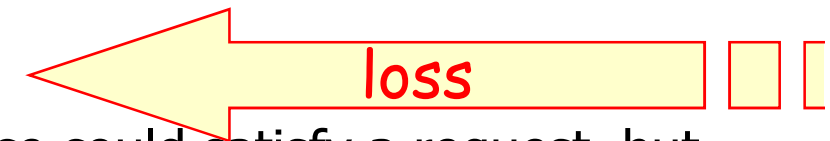
# Fragmentation

<u>Internal Fragmentation:</u>

**waste**

- allocated memory can be larger than requested memory
- memory management rounds up requested memory to the next manageable memory block unit

<u>External Fragmentation:</u>

**loss**

- Sum of total free memory space could satisfy a request, but free memory is scattered and is not contiguous
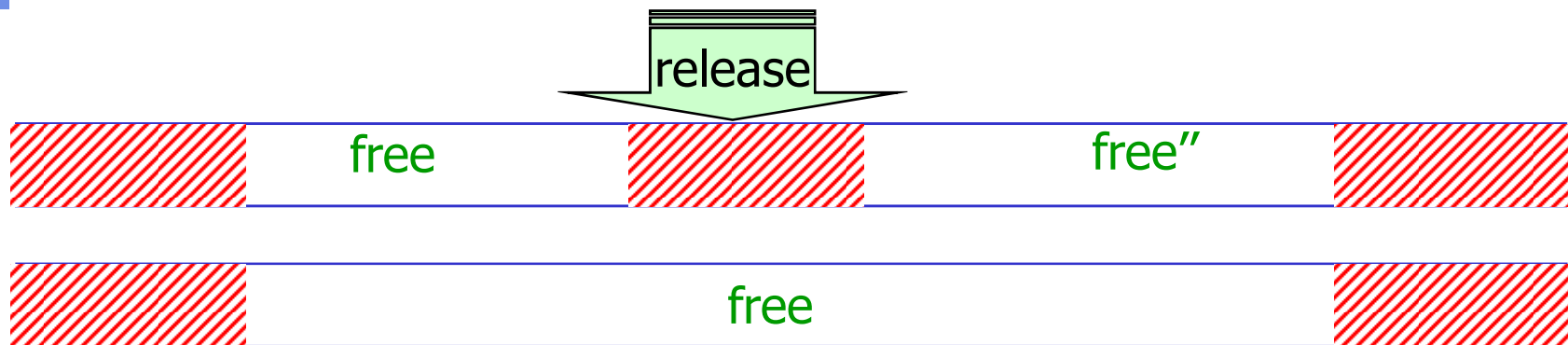
# Allocation Policies

- First-Fit

- Next-Fit (Rotating First-Fit)

- Best-Fit

- Nearest-Fit

Analyze pros and cons of each of them

# Reunification Policies

- Eager reunification (when releasing memory)



- Lazy reunification

# Additional Requirements

- With lazy reunification MM neighboring blocks can be free. However, none of these can satisfy the current memory request:

  ⇒ Garbage Collection[1]

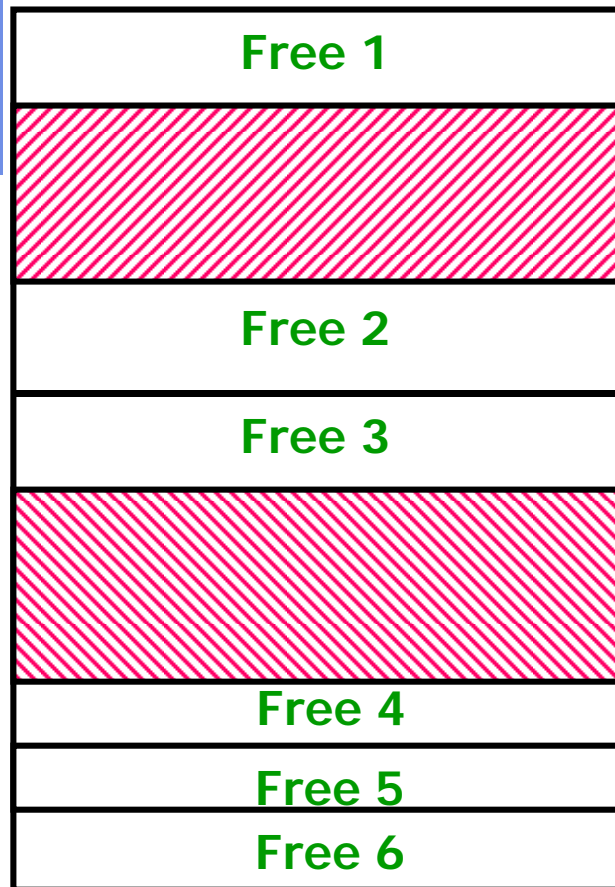- With an arbitrary allocation scheme (e.g. pure segmentation) we might get external fragmentation with a lot of scattered free blocks (Swiss cheese)

  ⇒ Compaction

[1]Garbage collection = releasing memory of no longer referenced objects

# Garbage Collection (1)

| Free 1 |
|---|
| (allocated) |
| Free 2 |
| Free 3 |
| (allocated) |
| Free 4 |
| Free 5 |
| Free 6 |

None of the 5 free blocks is large enough, but there are some neighboring free blocks

$\Rightarrow$

**New block to allocate**

Can be reunified to become a larger free block

# Garbage Collection (2)

Free 1′

Free 2′

Free 3′

Now 2 of the 3 free blocks
are large enough

either

or

**New block to allocate**

# Compaction (1)

| |
|---|
| Free 1 |
| (block) |
| Free 2 |
| (block) |
| Free 3 |

Observation:
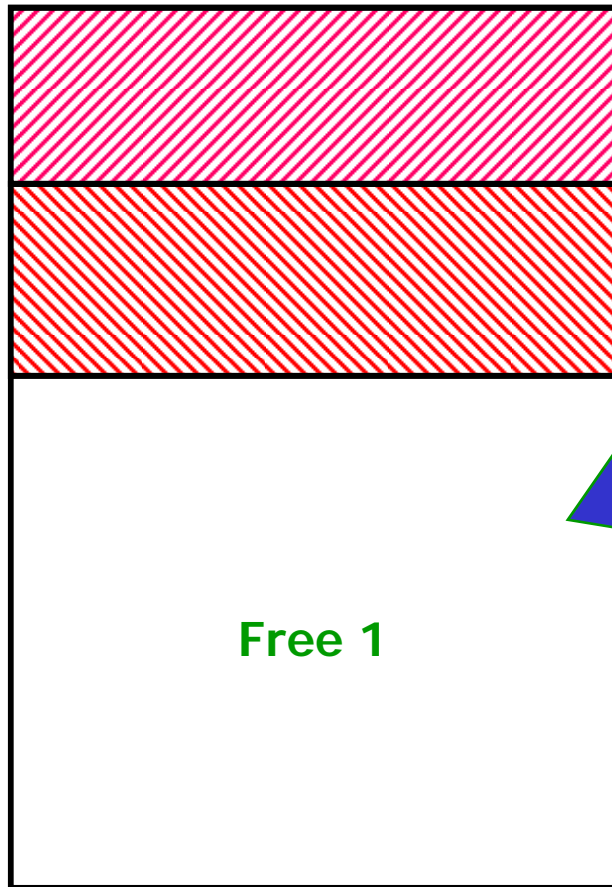None of the 3 free blocks is large enough, but ∃ enough free memory

New block to allocate

Idea:
Move the allocated blocks towards a chosen memory boundary

# Compaction (2)

In textbooks garbage collection is often mixed up with compaction

**New block to allocate**

**Free 1**

# Examples: Memory Manager

- Ringbuffer
- Stack

Very specific design parameters

- Boundary Tag System ("Randkennzeichnungsverfahren")
  - Operations in arbitrary order
  - Arbitrary sized blocks
  - Integrated management data structures
  - Allocation according to xyz-Fit (best-fit is possible)
  - External fragmentation
  - Immediate reunification

- Buddy System ("Halbierungsverfahren")
  - Operations in arbitrary order
  - Allocated blocks of $2^0$, $2^1$, $2^2$, $2^3$, ...
  - Explicit management data structures
  - Allocation according to "Best-Fit"
  - Internal and external fragmentation
  - Immediate (or lazy) reunification

- Linux Slab Allocator ("Stückchenzuteiler")

Very specific design parameters

# Boundary Tag

allocated

length

**Not Free**

| free | a | l | allocated | a | l | free |

free

**Not Allocated**

| allocated | f | l | fp | bp | free | f | l | allocated |

pair of pointers
for the free list

# Reunification (1)

Release this Portion !!

f l1 fpbp | free | f l1 | a l2 | allocated | a l2 | f l3 fpbp | free | f l3

*What to do?*

Because length and status field are of fixed size ⇒
we can easily get the length of the block to be released.

Furthermore we can look over both boundaries to get
necessary information about the neighboring blocks.

# Reunification (2)

Release this block !!



| f | l1 | fp | bp | free | | f | l1 | a | l2 | allocated | | a | l2 | f | l3 | fp | bp | free | | f | l3 |

| f | l1 | fp | bp | free | | f | l1 | f | l2 | fp | bp | releasing | | f | l2 | f | l3 | fp | bp | free | | f | l3 |

We can detect whether both (one or no) neighbor(s) are(is) free
$\Rightarrow$

# Reunification (3)

Release this block !!



Finally we have to adjust the new pointers in the free list and the resulting length field (in this case: $l1' = l1 + l2 + l3$)

# Buddy System (1)

Memory to be managed

$2^{20}$ Byte = 1 MB free

Suppose a client requests 100 KB

# Buddy System (2)

Dividing free memory until an appropriate free block

| 1 MB free | | | |
|---|---|---|---|

| | | 512 KB | |
|---|---|---|---|

| | 256 KB | 512 KB | |
|---|---|---|---|

| | 128 KB | 256 KB | 512 KB |
|---|---|---|---|

| 100 KB | 128 KB | 256 KB | 512 KB |
|---|---|---|---|

Internal fragmentation

# Buddy System (3)

| 1 MB free |
|---|

| 100 KB | | 128 KB | 256 KB | 512 KB |
|---|---|---|---|---|

**Suppose another 200 KB block is requested**

# Buddy System (4)

| 1 MB free | | | |
|---|---|---|---|

| 100 KB | | 128 KB | 256 KB | 512 KB |
|---|---|---|---|---|

| 100 KB | | 128 KB | 200 KB | | 512 KB |
|---|---|---|---|---|---|

**Suppose another 200 KB block is requested**

# Buddy System (5)

| 1 MB free | | | |
|---|---|---|---|

| 100 KB | 128 KB | 256 KB | 512 KB |
|---|---|---|---|

| 100 KB | 128 KB | 200 KB | 512 KB |
|---|---|---|---|

| 100 KB | 128 KB | 200 KB | 200 KB | 256 KB |
|---|---|---|---|---|

**Suppose another 130 KB block is requested**

# Buddy System (6)

| 1 MB free |
|---|

| 100 KB | | 128 KB | 256 KB | 512 KB |
|---|---|---|---|---|

| 100 KB | | 128 KB | 200 KB | | 512 KB |
|---|---|---|---|---|---|

| 100 KB | | 128 KB | 200 KB | | 200 KB | | 256 KB |
|---|---|---|---|---|---|---|---|

| 100 KB | | 128 KB | 200 KB | | 200 KB | | 130 KB | |
|---|---|---|---|---|---|---|---|---|

**Suppose we have to release the first 200 KB**

# Buddy System (7)

| 1 MB free |
|---|

| 100 KB | 128 KB | 256 KB | 512 KB |
|---|---|---|---|

| 100 KB | 128 KB | 200 KB | | 512 KB |
|---|---|---|---|---|

| 100 KB | 128 KB | 200 KB | | 200 KB | | 256 KB |
|---|---|---|---|---|---|---|

| 100 KB | 128 KB | 200 KB | | 200 KB | | 130 KB | |
|---|---|---|---|---|---|---|---|

| 100 KB | 128 KB | 256 KB | 200 KB | | 130 KB | |
|---|---|---|---|---|---|---|

**Suppose we have to release the other 200 KB**

# Buddy System (8)

| 1 MB free |
|---|

| 100 KB | | 128 KB | 256 KB | 512 KB |
|---|---|---|---|---|

| 100 KB | | 128 KB | 200 KB | | 512 KB |
|---|---|---|---|---|---|

| 100 KB | | 128 KB | 200 KB | | 200 KB | | 256 KB |
|---|---|---|---|---|---|---|---|

| 100 KB | | 128 KB | 200 KB | | 200 KB | | 130 KB | |
|---|---|---|---|---|---|---|---|---|

| 100 KB | | 128 KB | 256 KB | 200 KB | | 130 KB | |
|---|---|---|---|---|---|---|---|

| 100 KB | | 128 KB | 256 KB | 256 KB | 130 KB | |
|---|---|---|---|---|---|---|

*Can we reunify these two blocks forming a 512 KB free block?*

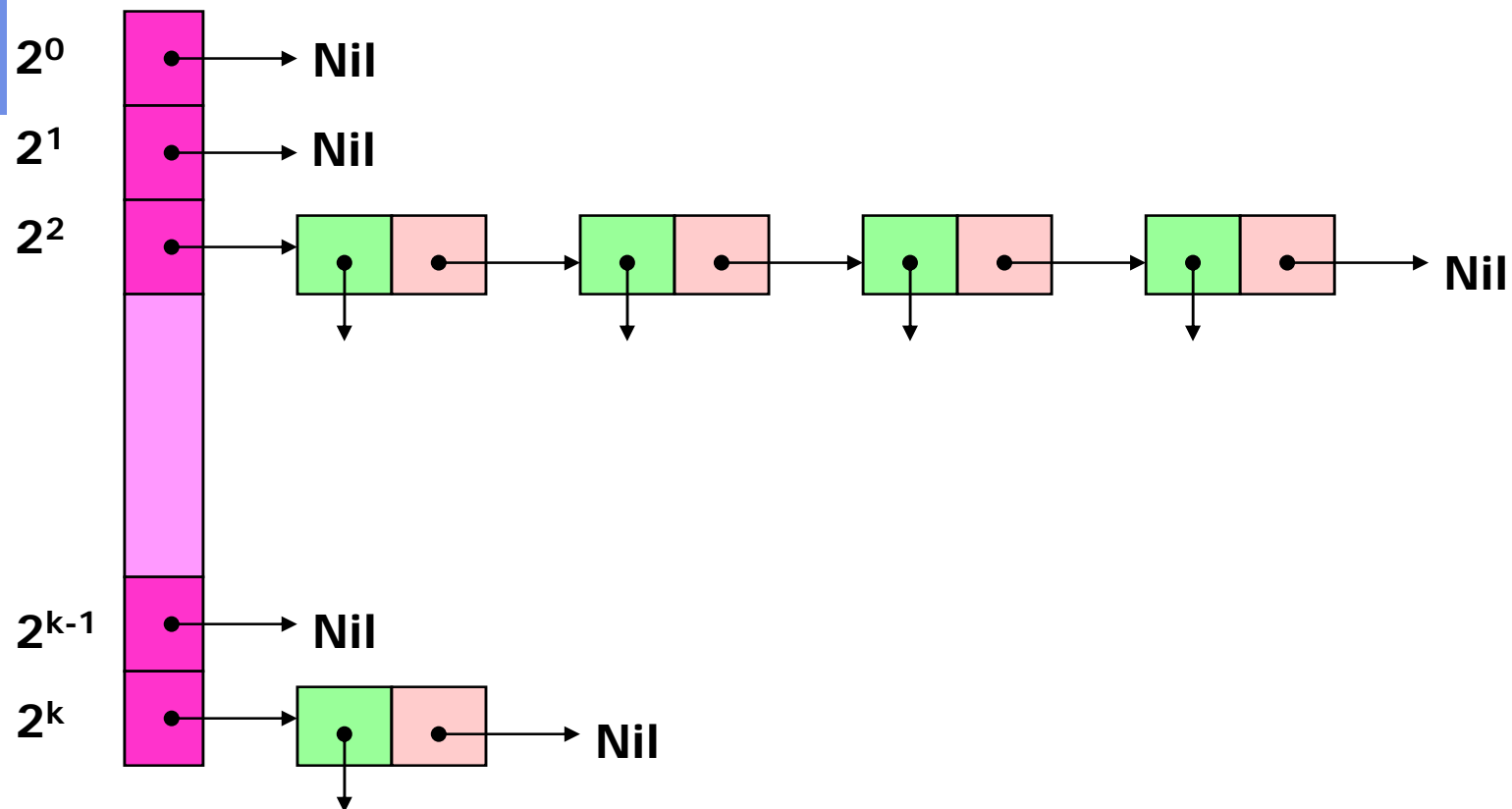Not at all, they are not buddies! Only buddies belong with each other!

# Data Structures for Buddy System



$2^0$    → Nil

$2^1$    → Nil

$2^2$    → Nil

$2^{k-1}$    → Nil

$2^k$    → Nil

Array of heads pointing to free blocks of a certain size $2^k$

# Operations of Buddy System (1)

Allocating block of size s:

- Round up s to next power of 2, say $2^i$
  ($\Rightarrow$ internal fragmentation)

- Access head of the list for the $2^i$ pieces

- If list is not empty get first element of list

- If list is empty (recursively) do:
  - Access head of list for the $2^{i+1}$ pieces
  - If list isn't empty get first element of list
  - cut element in halves
  - take lower half, insert upper half into list for $2^i$ pieces …

# Operations of Buddy System (2)

## Releasing a block of size $2^i$

- Determine its buddy

- if buddy is (partly) allocated

  insert the block to be released

  into the list of pieces of size $2^i$

- if buddy is free reunify both buddies

*Question:*
*How can we efficiently determine the appropriate buddy?*

# Determining the Buddy

## Address calculation for a piece of size $2^{12}$

| 15 | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|----|----|----|----|----|----|---|---|---|---|---|---|---|---|---|---|
| X | Y | Z | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

**address of a piece of size $2^{12}$**

_Question:_ What's the address of its buddy?

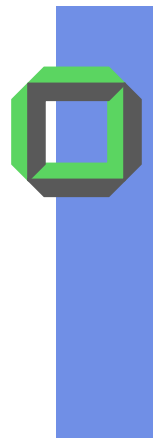| 15 | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|----|----|----|----|----|----|---|---|---|---|---|---|---|---|---|---|
| X | Y | Z | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

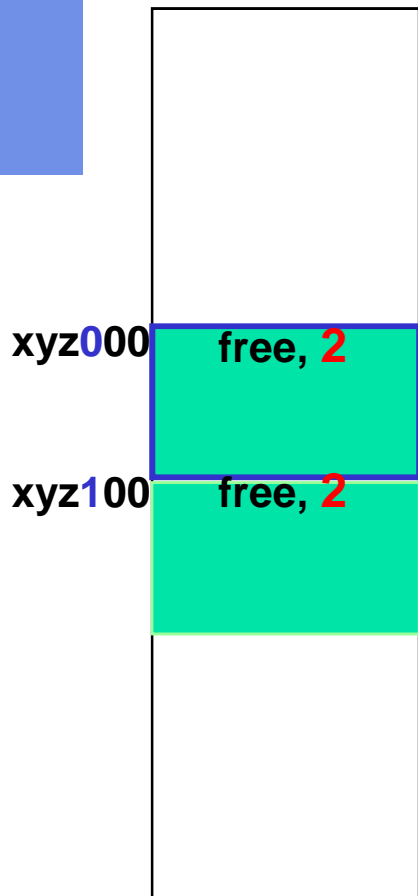Remark: Other pieces of size $2^{12}$ differ only in the leading address bits.

# Reunification of Buddies

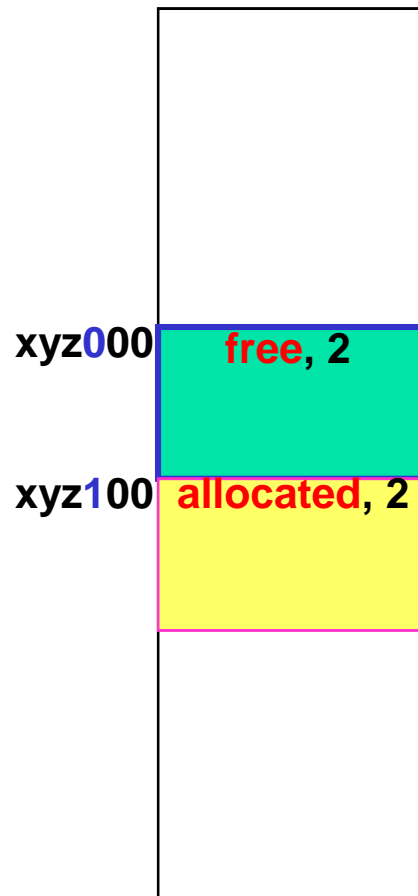Each block of size $2^i$ of buddy system contains:

- State of the block (allocated or free)

- Length of the block

## *Is there an alternative?*

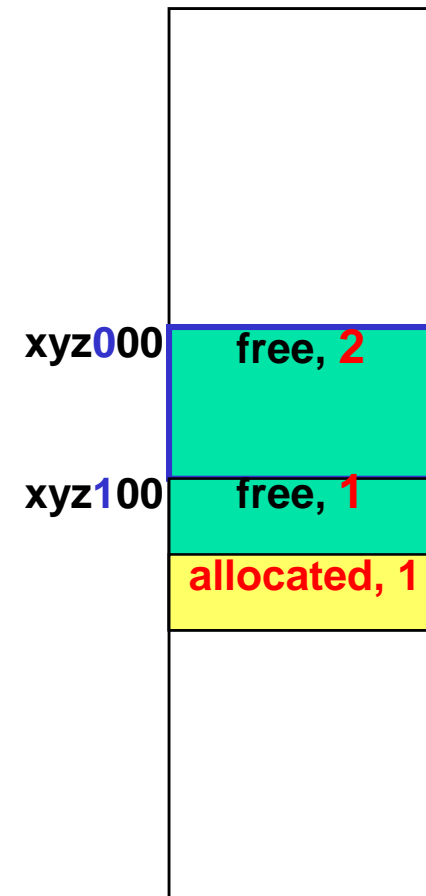- Analyze the implementation of Linux

- Compare both implementations

# Scenarios for Reunifications

xyz**0**00 | free, **2**

xyz**1**00 | free, **2**

**reunification**
**buddy free**

xyz**0**00 | free, **2**

xyz**1**00 | allocated, **2**

**no reunification**
**buddy allocated**

xyz**0**00 | free, **2**

xyz**1**00 | free, **1**

allocated, **1**

**no reunification**
**buddy partly allocated**

# Summary of Buddy System

On average, the internal fragmentation is about 25%
- each memory block is at least 50% occupied

Works efficiently if the size M of RAM is a power of 2

# Slab Allocating (in Linux)

- Inside a kernel there are a few data types, e.g.
  - TCB
  - LAS descriptor
  - Page table
  - File handle …

- Kernel tend to request these data types over and over again, e.g. to establish a new thread/task

- Slab allocation tries to support reusage of previously used data types $\Rightarrow$ install object type caches

- Following slides:
  - Steffen Wolfer, Proseminar WS 2003 Linux Internals
  - see also: Sven Krohlas,　　　"　　　WS 2004

# Slab Allocating (1)

New approach (Sun Microsystems, 1994)

- regard a memory area as an „object container"

- collect objects of the same type in specific "logical kernel caches"

- divide caches into slabs

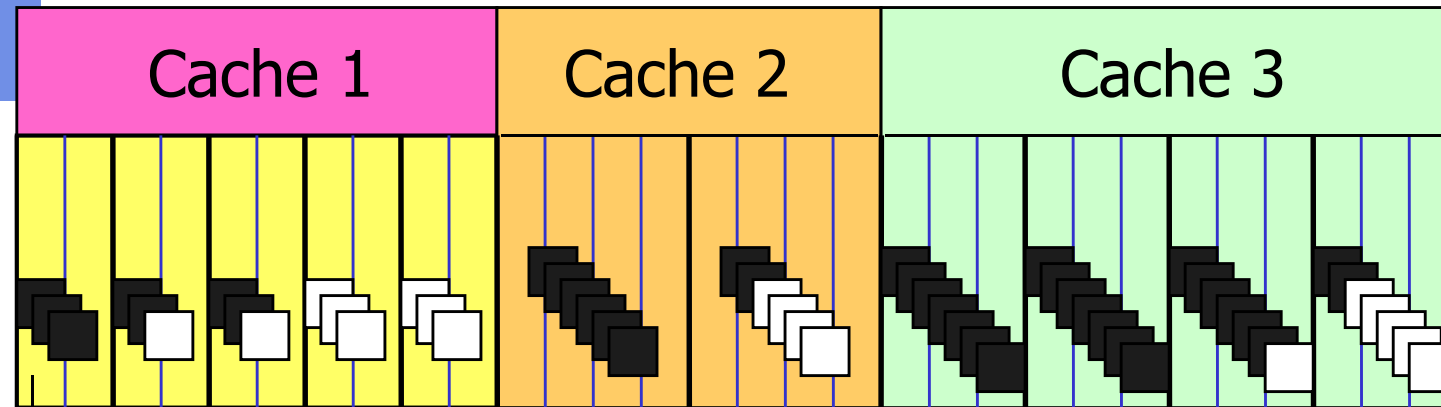- slab is part of a specific cache

- slab size = s*page frames, s ≥ 1

# Slab Allocating (2)

- Reuse already initialized objects

- Simplify complex allocation

- Take page frames from the buddy system

  - Give back only if buddy system needs them

# Slab Allocating (3)



Cache 1 | Cache 2 | Cache 3

Page frame

Slabs

■ Allocated object

□ Free object

# Slab Allocating (4)

Table for the caches:

| Cache Descriptor |
| --- |
| slabs_full |
| slabs_partial |
| slabs_free |
| next |
| num |
| ... |

- doubled linked lists with slab descriptors for
  - full
  - partially full
  - empty slabs

- pointer to next cache descriptor

- number of objects per slab, size of object, flags, ...
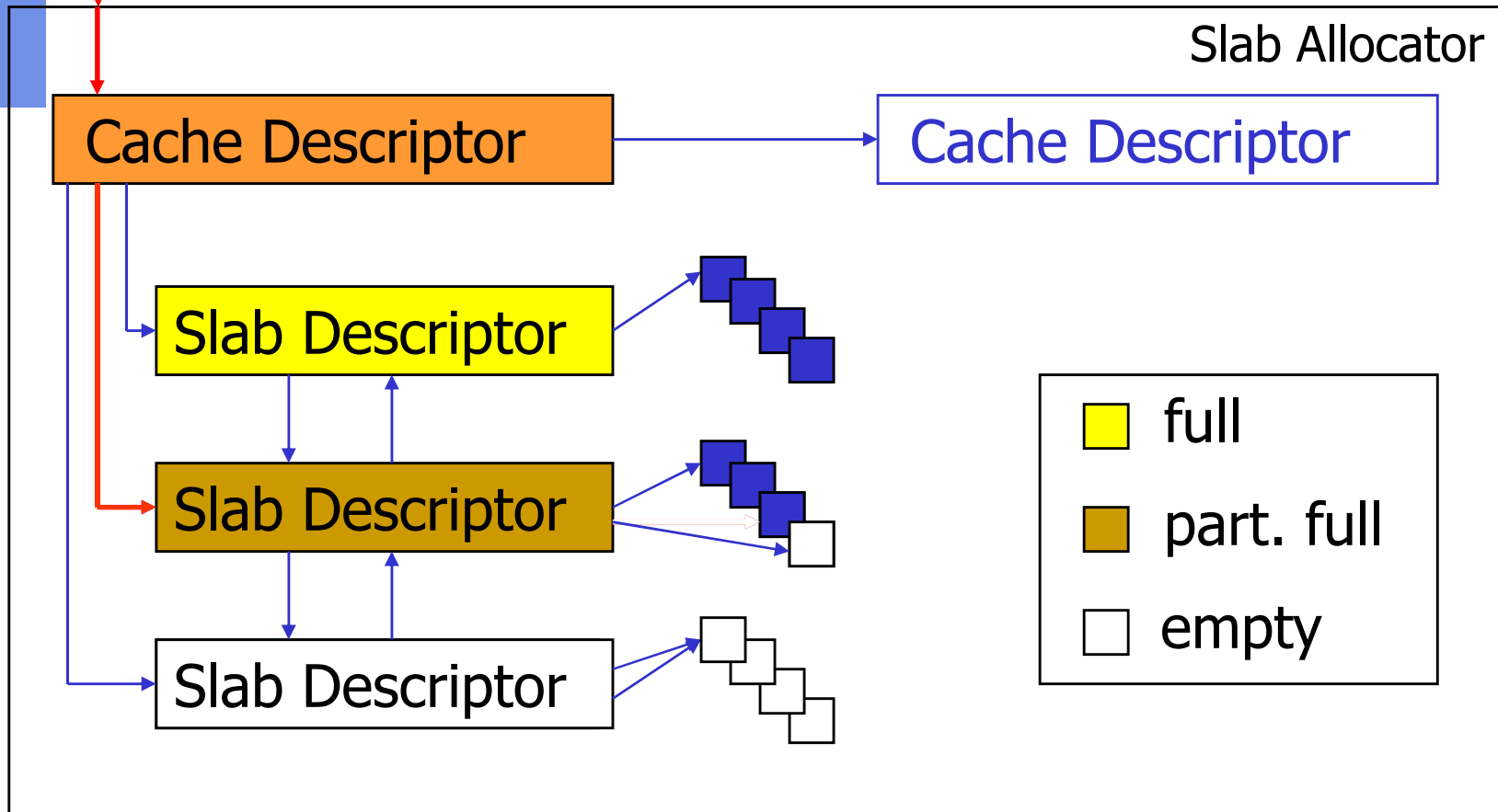
# Slab Allocating (5)

Table for Slabs:

| Slab Descriptor |
| --- |
| inuse |
| s_mem |
| free |
| list |
| ... |

- number of currently allocated objects

- pointer to 1. Object

- pointer to 1. free Object

- pointer to list of slab descriptors
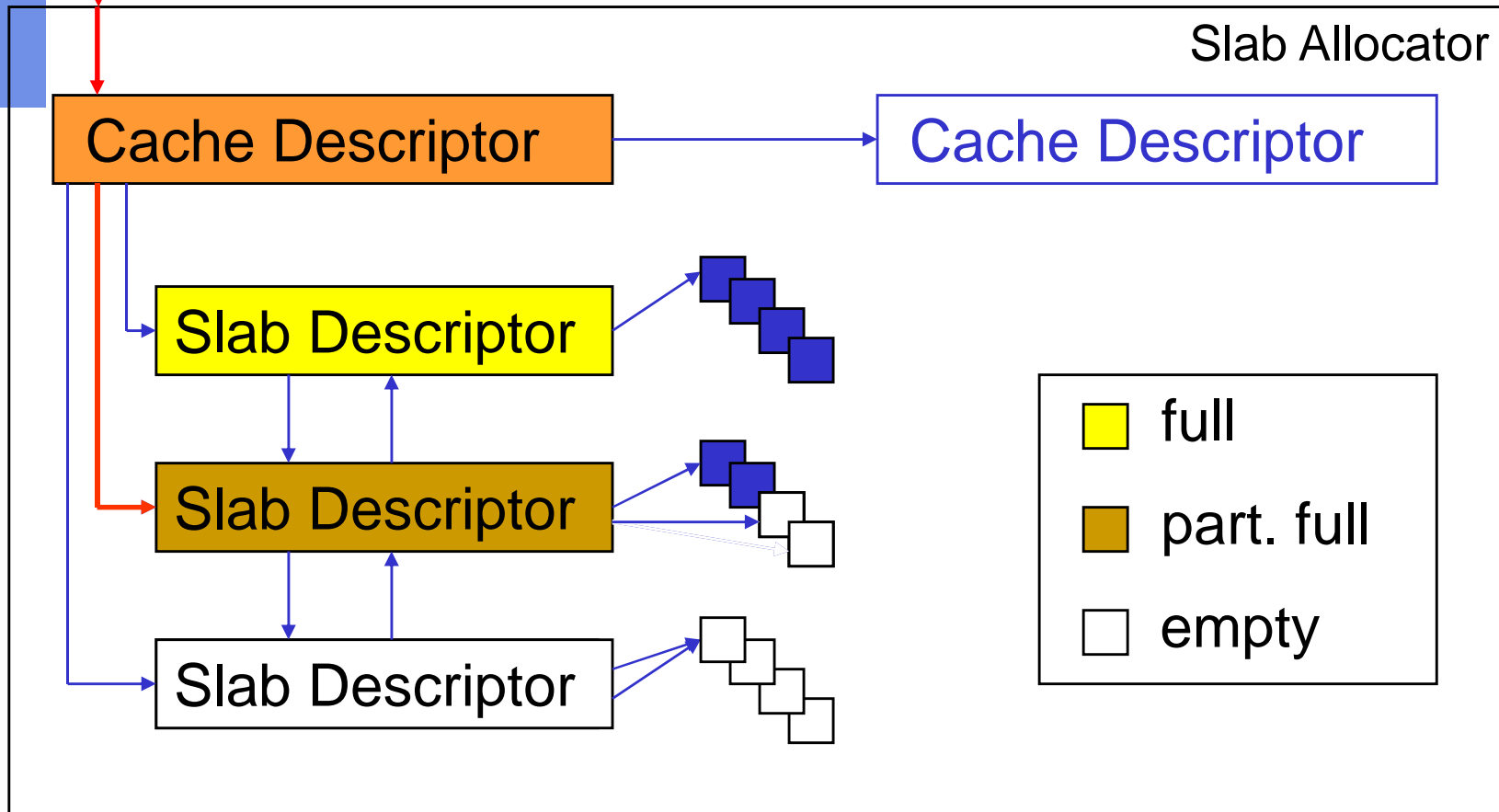
# Slab Allocating (6)

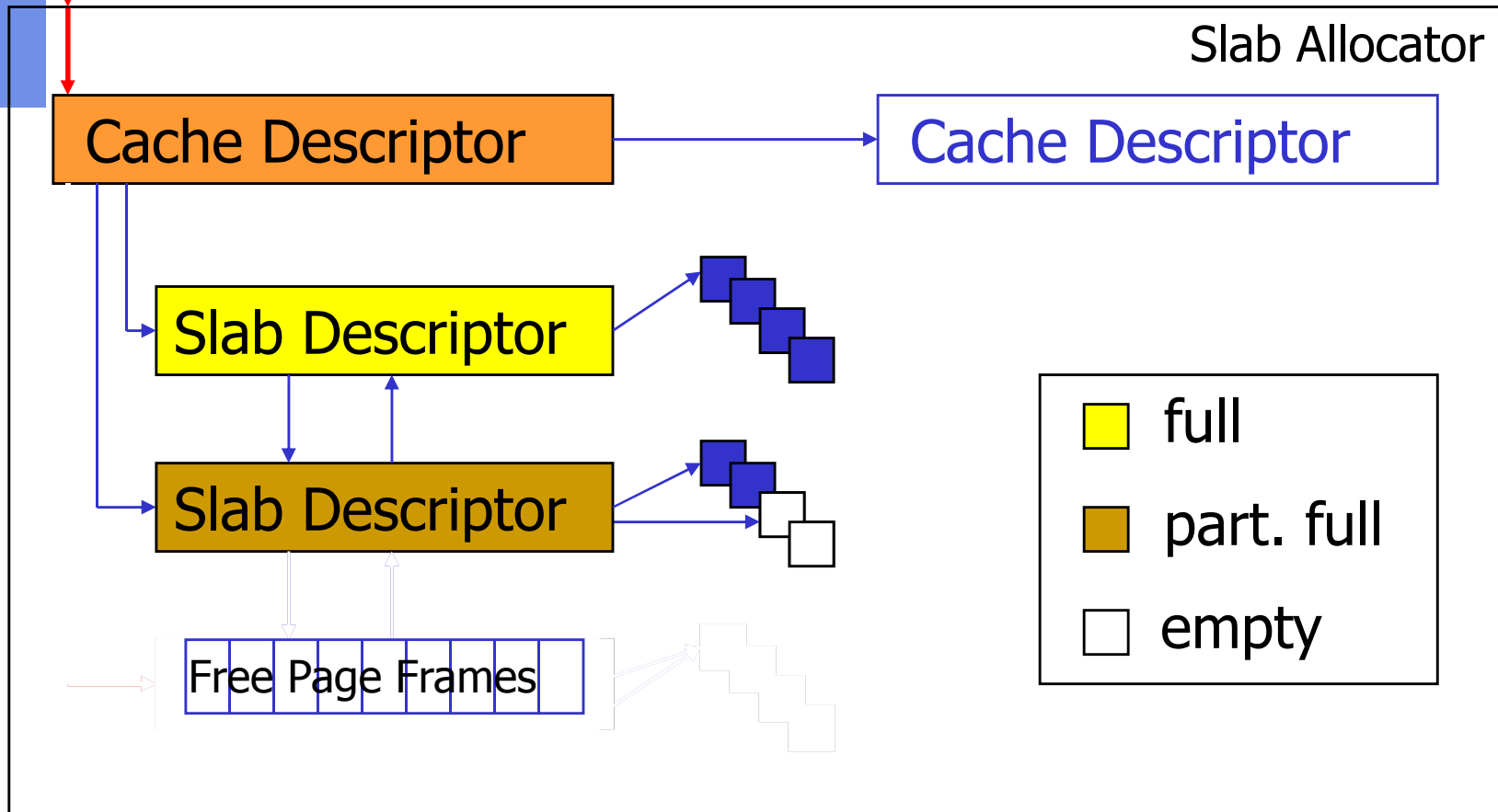Allocate memory: kmem_cache_alloc()

# Slab Allocating (7)

Release memory: kmem_cache_free()



Slab Allocator

Cache Descriptor

Cache Descriptor

Slab Descriptor

Slab Descriptor

Slab Descriptor

| | |
|---|---|
| ▢ | full |
| ▢ | part. full |
| ▢ | empty |

# Slab Allocating (8)

Free slabs: kmem_cache_destroy()



Slab Allocator

Cache Descriptor → Cache Descriptor

Slab Descriptor

Slab Descriptor

Free Page Frames

- full
- part. full
- empty

# Summary of Slab Allocation

- *Who is using slab allocation?*

- *What kernel data types are mapped to slabs?*

- *How much space in total is managed by slabs?*

- *Do slabs have to be mapped to contiguous memory?*
- *...*