

Verification of OS-level Cache Management



Renato Mancuso



Sagar Chaki

OSPERT 2018

Goal + Approach

Colored Lockdown
for deterministic cache management

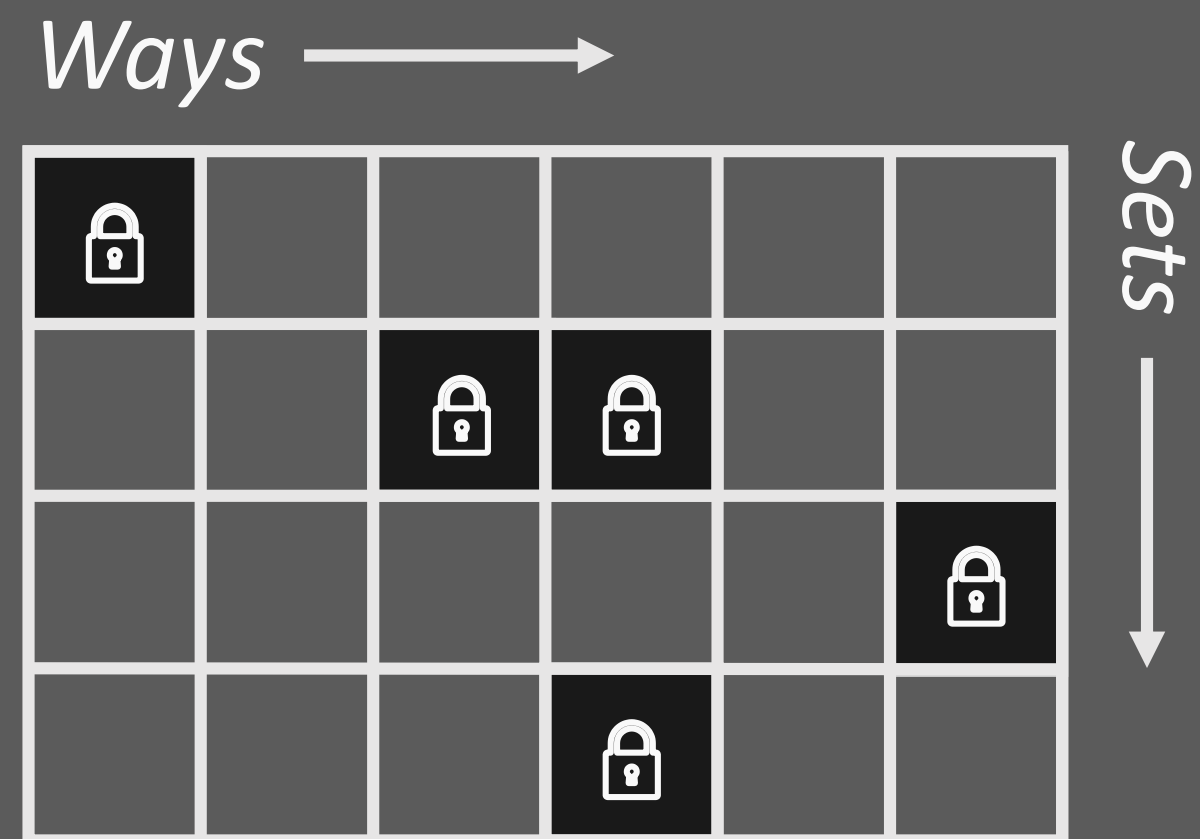
C source code

via **CBMC**

Linux Kernel
Module

Last Level Cache

Management Model



- ✓ Addresses all the sources of interference
- ✓ Converts the LLC cache in a deterministic object at the granularity of a single memory page
- ✓ Allows the use of legacy code
- ✓ Provides flexibility in cache assignment

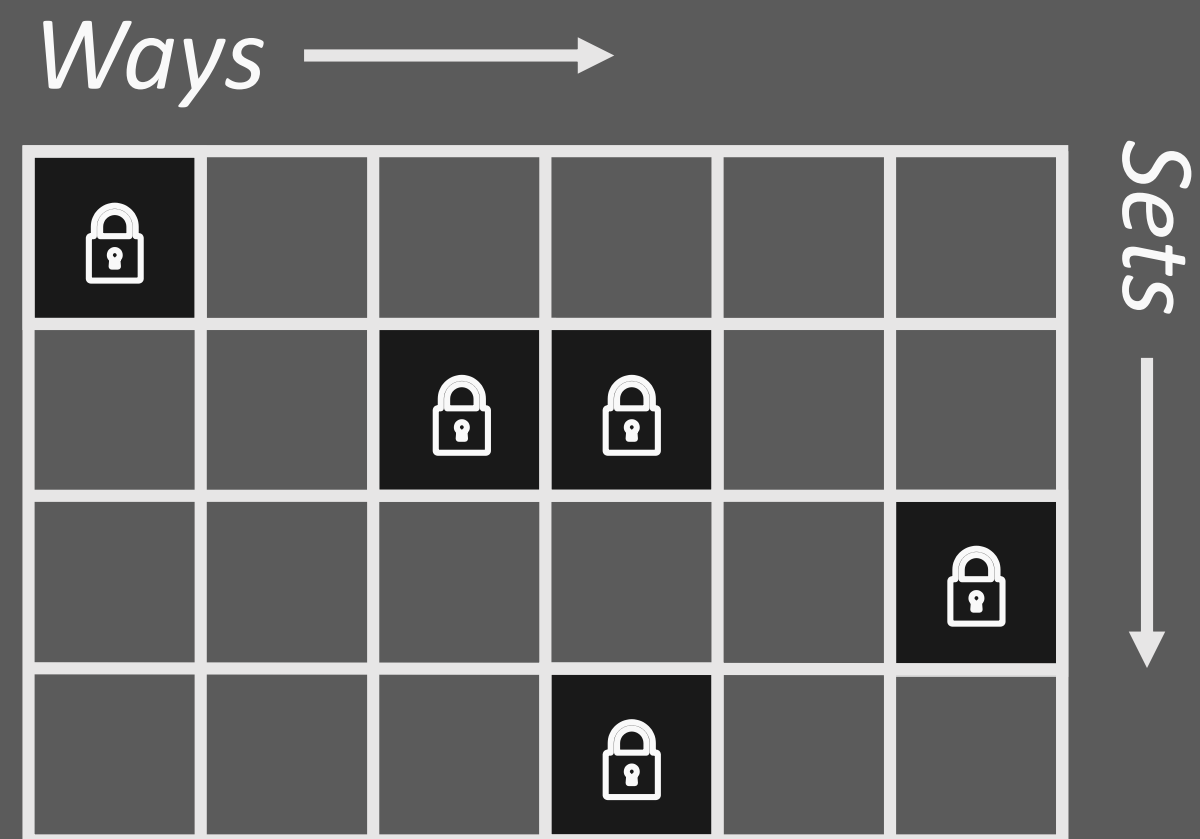


Background

Colored Lockdown

Last Level Cache

Management Model



- ✓ Addresses all the sources of interference
- ✓ Converts the LLC cache in a deterministic object at the granularity of a single memory page
- ✓ Allows the use of legacy code
- ✓ Provides flexibility in cache assignment



100% hits
on allocated pages

100% misses
on non-allocated pages

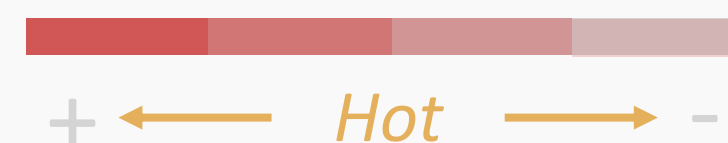
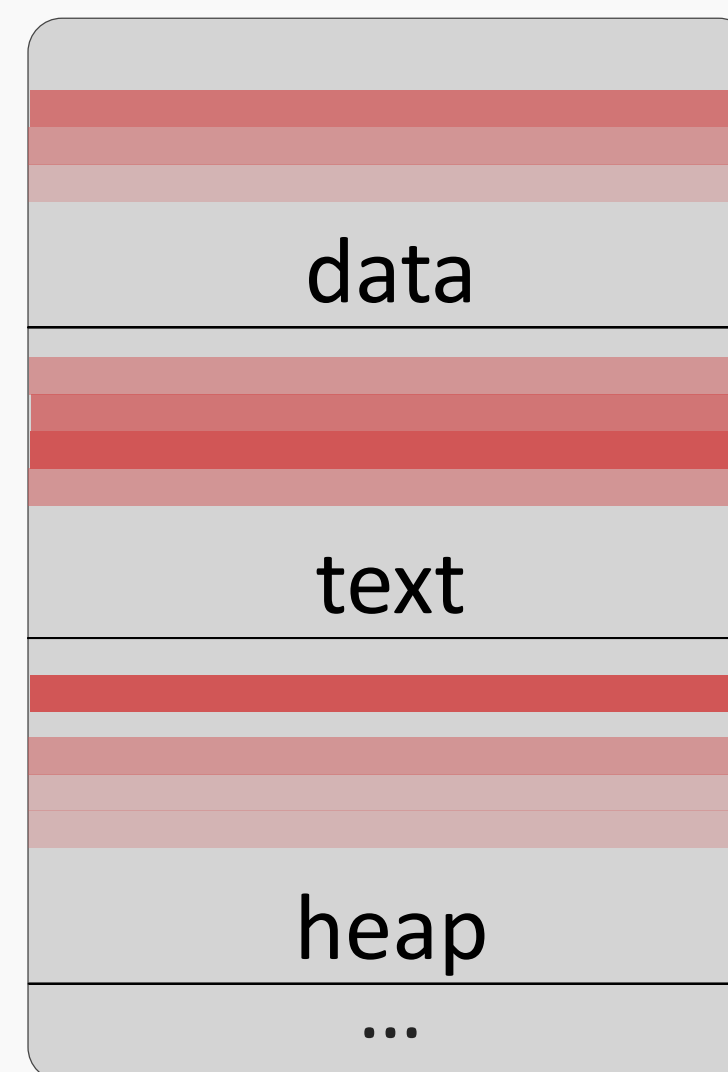
Profile-Driven Cache Allocation



PROBLEM

Caches are critical,
constrained resources.
Optimal allocation ?

Process Address Space



Location of hot
region(s) is unknown

Absolute virtual
memory addresses
may change

Profile-Driven Cache Allocation



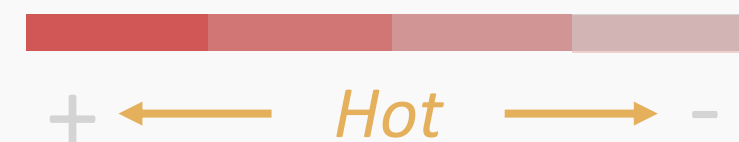
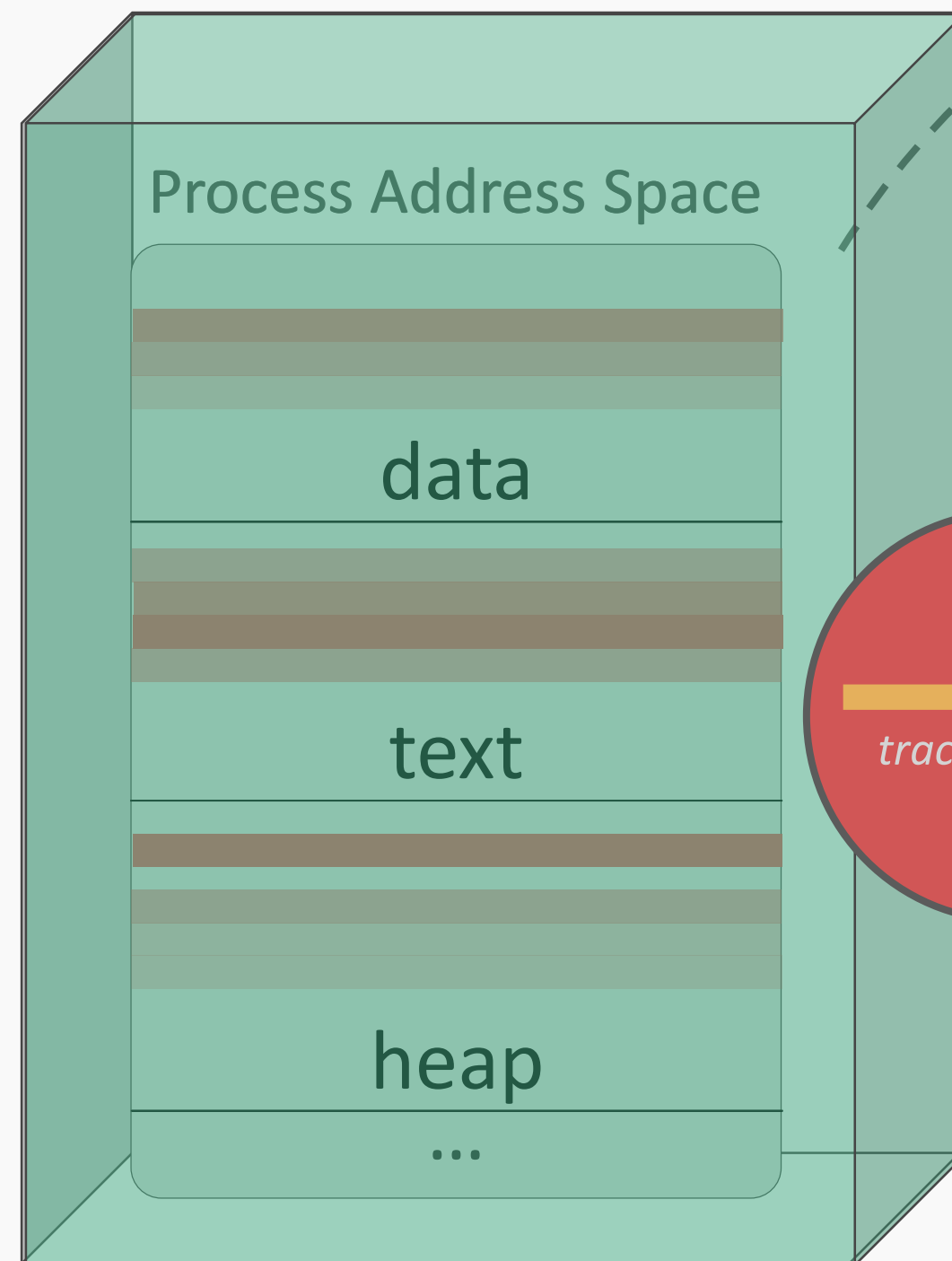
PROBLEM

Caches are critical, constrained resources. Optimal allocation ?



PROFILE MEMORY

Extract memory traces and produce memory usage profile.



Application Profile

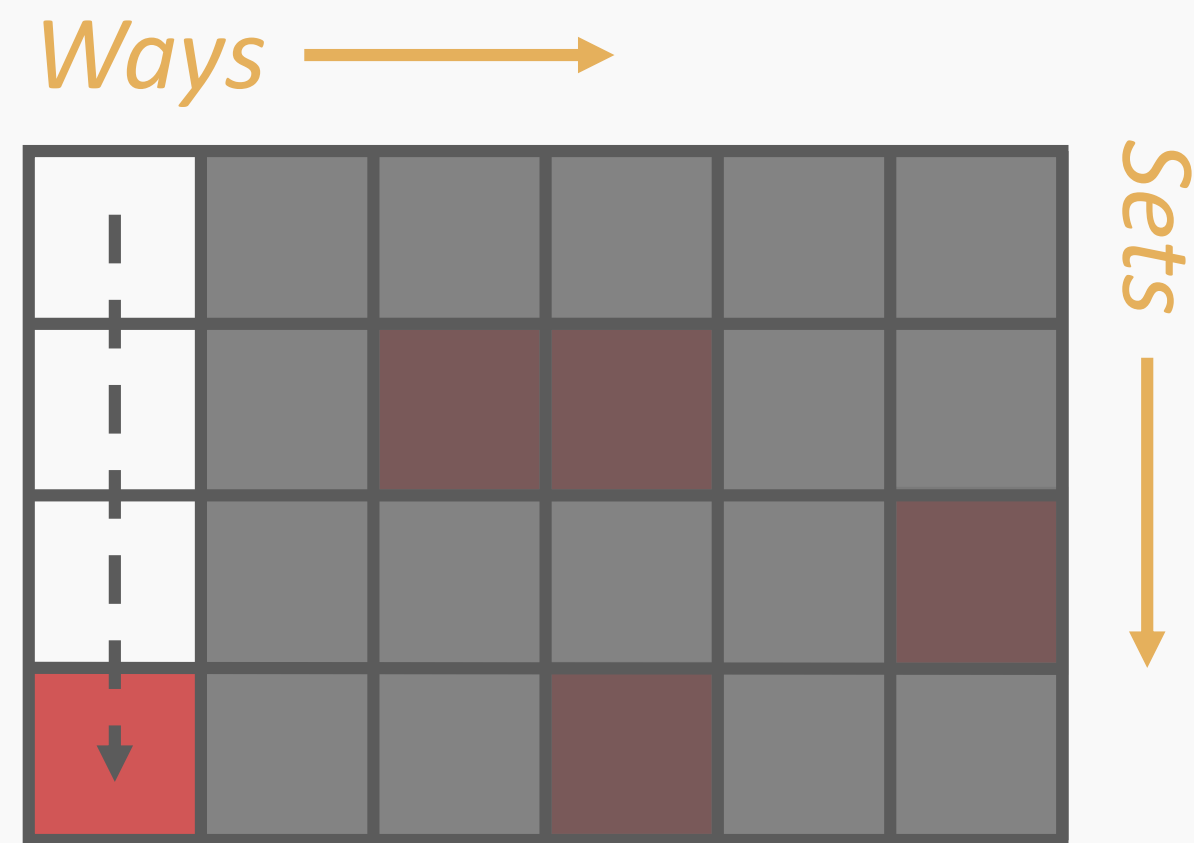
Page	Accesses
A	100 K
B	10 K
C	1 K
D	700
E	500
F	90
C	50
D	10
E	5

threshold

μ

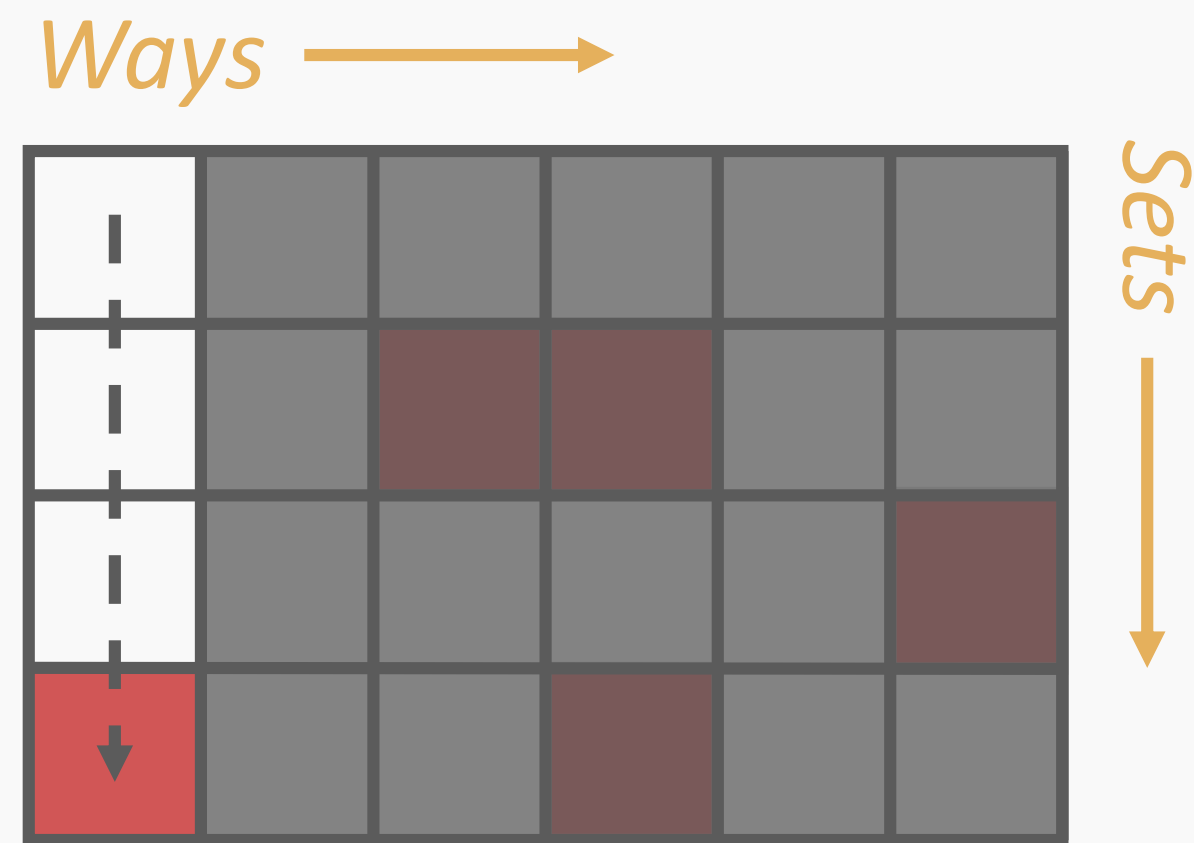
residual
cache misses

COLORING



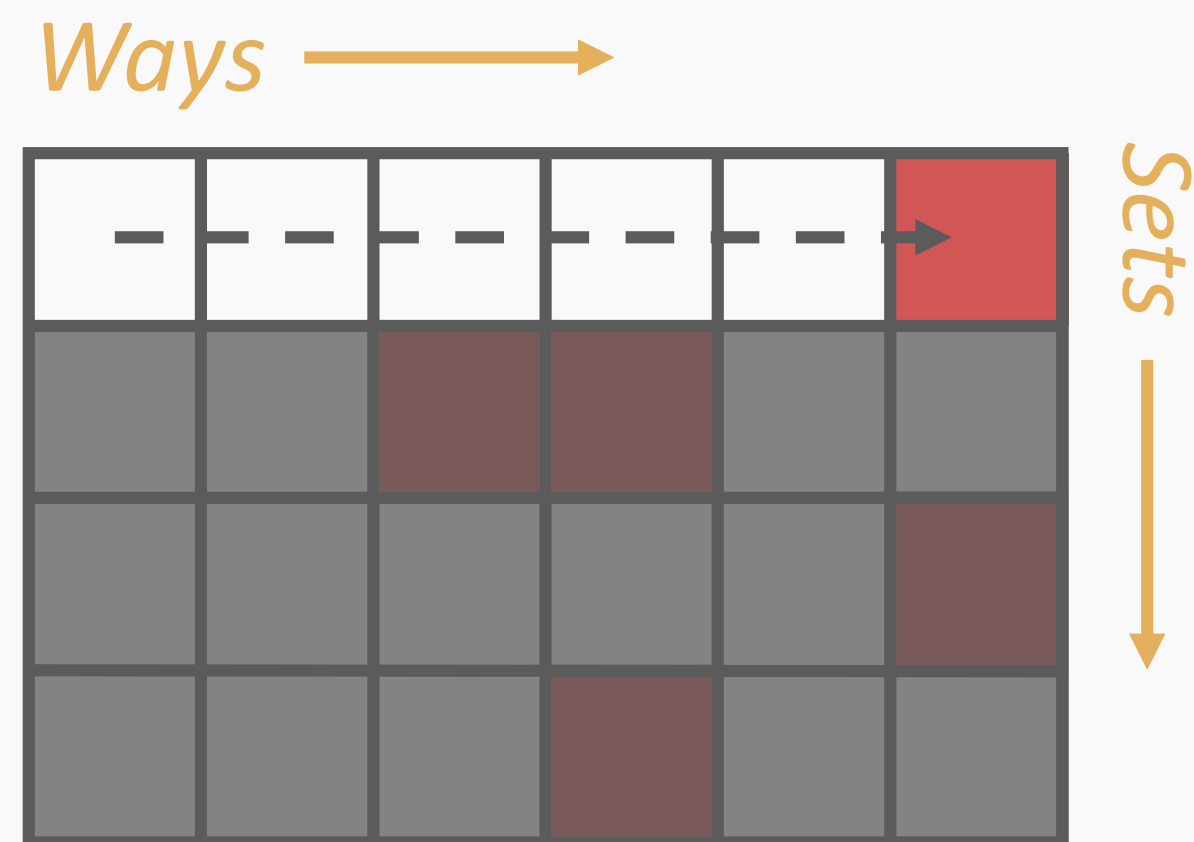
- Leverages on the **virtual** → **physical** translation layer
- Used to **move page mapping** across sets (up/down)
- Transparent to the **application**

COLORING



- Leverages on the **virtual** → **physical** translation layer
- Used to **move page mapping** across sets (up/down)
- Transparent to the **application**

LOCKDOWN



- Uses **architecture-specific** lockdown features
- Used to **allocate pages** on selected ways (left/right)
- Once allocated, pages trigger **cache hits** until deallocation

offline

online



offline

Profile
generation

Task tracing +
analysis

online



offline

Profile load

Color detection +
re-assignment

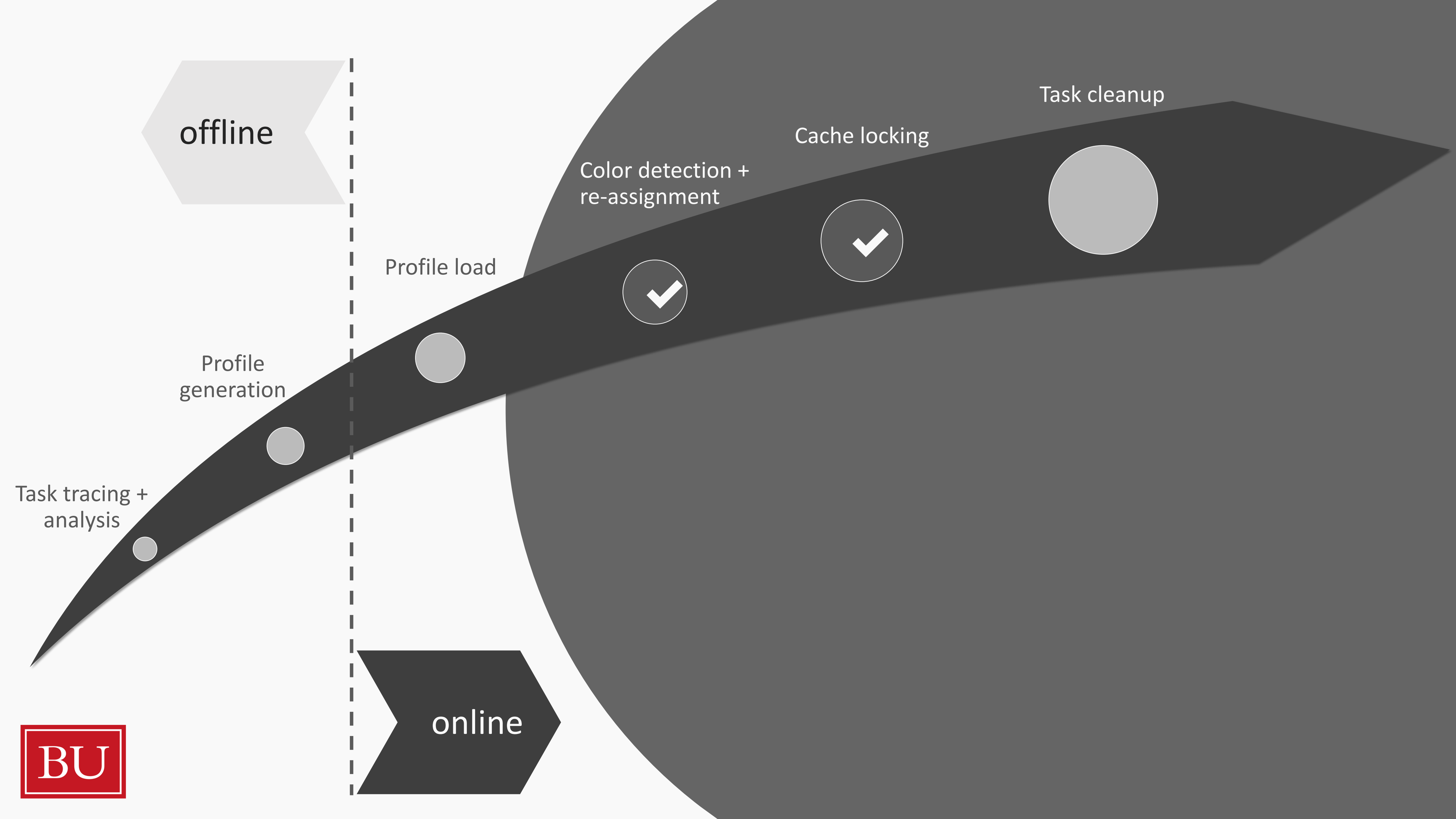
Cache locking

Task cleanup

Profile
generation

Task tracing +
analysis

online



offline

Profile load

Profile generation

Task tracing + analysis

Color detection + re-assignment

Cache locking

Task cleanup

1

2

3

4

VERIFIED PROPERTIES

online



If cache is large enough, allocation is performed

offline

Profile load

Profile generation

Task tracing + analysis

Color detection + re-assignment

Cache locking

Task cleanup

VERIFIED PROPERTIES

1

2

3

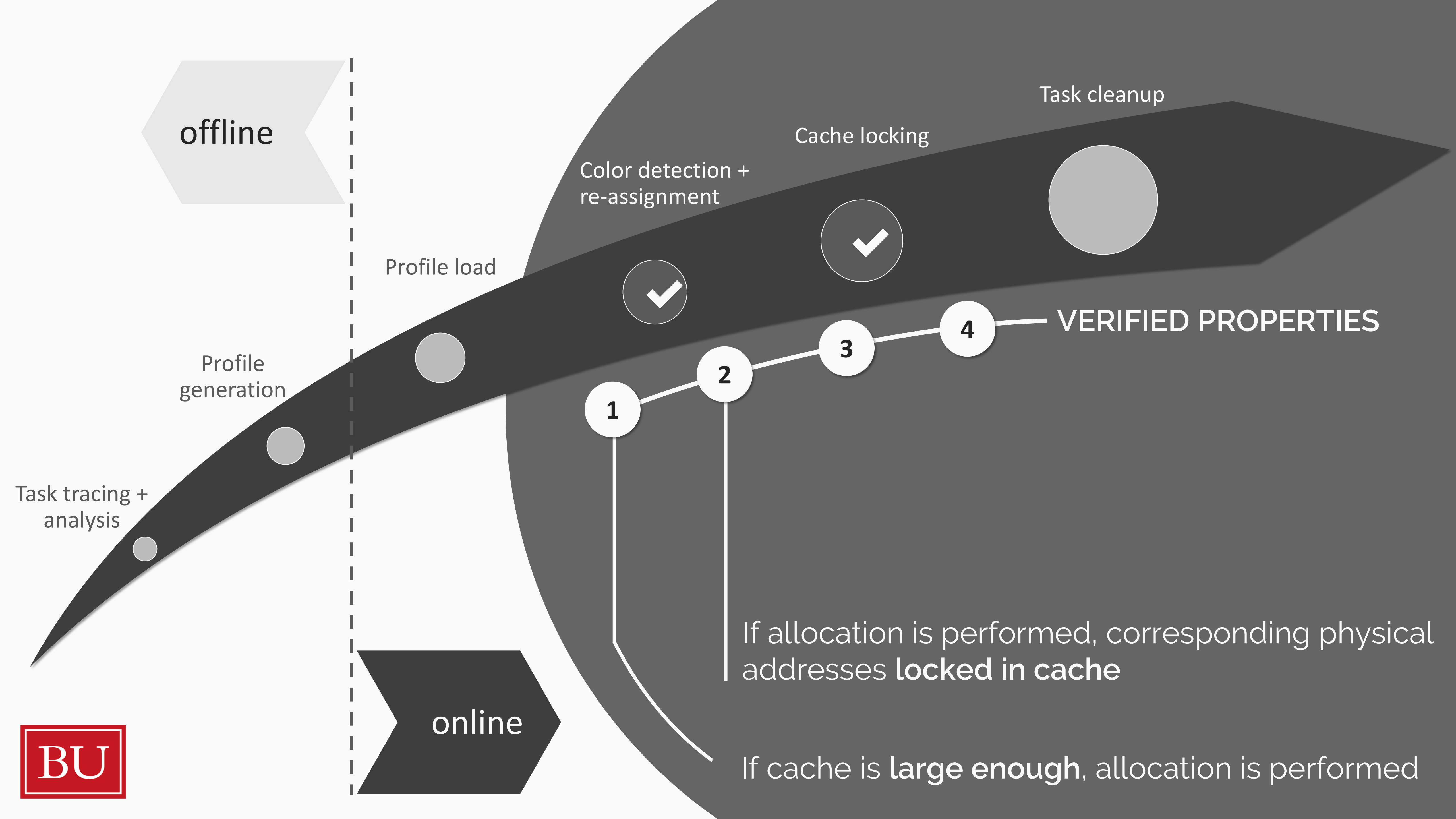
4

If allocation is performed, corresponding physical addresses **locked in cache**

If cache is **large enough**, allocation is performed

online

BU



offline

Profile load

Profile generation

Task tracing + analysis

Color detection + re-assignment

Cache locking

Task cleanup

VERIFIED PROPERTIES

1

2

3

4

No **extra** memory locked in cache

If allocation is performed, corresponding physical addresses **locked in cache**

If cache is **large enough**, allocation is performed

online

BU

offline

Profile load

Profile generation

Task tracing + analysis

Color detection + re-assignment

Cache locking

Task cleanup

VERIFIED PROPERTIES

All temporary **kernel resources** released at the end

No **extra** memory locked in cache

If allocation is performed, corresponding physical addresses **locked in cache**

If cache is **large enough**, allocation is performed

online

what was Verified

1 After locking issued, line is in cache

4

All temporary **kernel resources** released at the end

3

No **extra** memory locked in cache

2

If allocation is performed, corresponding physical addresses **locked in cache**

1

If cache is **large enough**, allocation is performed

what was Assumed

what was Verified

4

All temporary **kernel resources** released at the end

3

No **extra** memory locked in cache

2

If allocation is performed, corresponding physical addresses **locked in cache**

1

If cache is **large enough**, allocation is performed

1

After locking issued, line is in cache

2

Profile is correct

what was Assumed

what was Verified

4

All temporary **kernel resources** released at the end

3

No **extra** memory locked in cache

2

If allocation is performed, corresponding physical addresses **locked in cache**

1

If cache is **large enough**, allocation is performed

1

After locking issued, line is in cache

2

Profile is correct

3

Kernel descriptors correct

what was Assumed

what was Verified

4

All temporary **kernel resources** released at the end

3

No **extra** memory locked in cache

2

If allocation is performed, corresponding physical addresses **locked in cache**

1

If cache is **large enough**, allocation is performed

1

After locking issued, line is in cache

2

Profile is correct

3

Kernel descriptors correct

4

Kernel routines correct

what was Assumed

what was Verified

4

All temporary **kernel resources** released at the end

3

No **extra** memory locked in cache

2

If allocation is performed, corresponding physical addresses **locked in cache**

1

If cache is **large enough**, allocation is performed

1

After locking issued, line is in cache

2

Profile is correct

3

Kernel descriptors correct

4

Kernel routines correct

5

Initial status of lockdown bit is known

what was Assumed

Verification Boundaries

Verified
Logic

Verification Boundaries

User-Space

Template
Generic
Profile

Verified
Logic

Verification Boundaries

User-Space

Template
Generic
Profile

Verified
Logic

Hardware

Instantiate
memory
model

Abstract
cache
model

Verification Boundaries

User-Space

Hardware

Template
Generic
Profile

Instantiate
memory
model

Verified
Logic

Abstract
cache
model

Initialize
descriptors

Abstract &
replace
routines

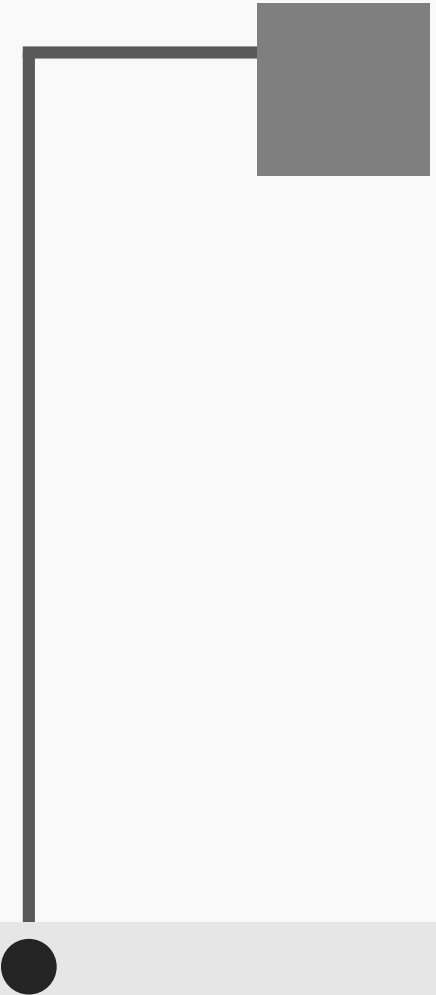
Rest of Kernel

Set Associative Cache

BU

Cache Model
& **Locking**

Set Associative Cache



A diagram showing a vertical line with a solid black circle at the bottom. A horizontal line extends to the right from the top of the vertical line, ending in a gray square. This represents a pointer to a cache line structure.

```
typedef struct {  
    void * addr;  
    char locked;  
} cache_line_t;
```

Cache Model
& Locking

Set Associative Cache

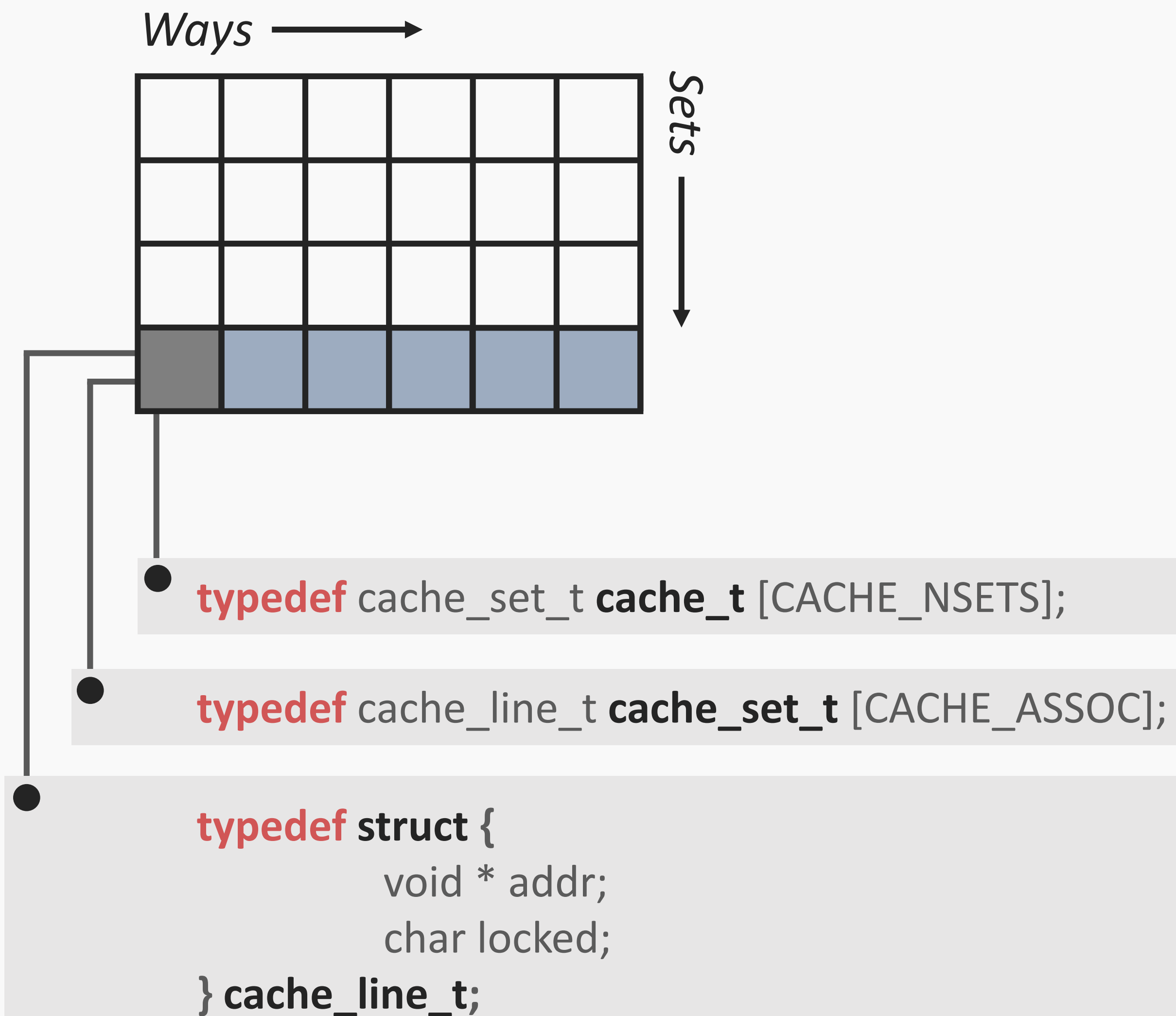


```
typedef cache_line_t cache_set_t [CACHE_ASSOC];
```

```
typedef struct {  
    void * addr;  
    char locked;  
} cache_line_t;
```

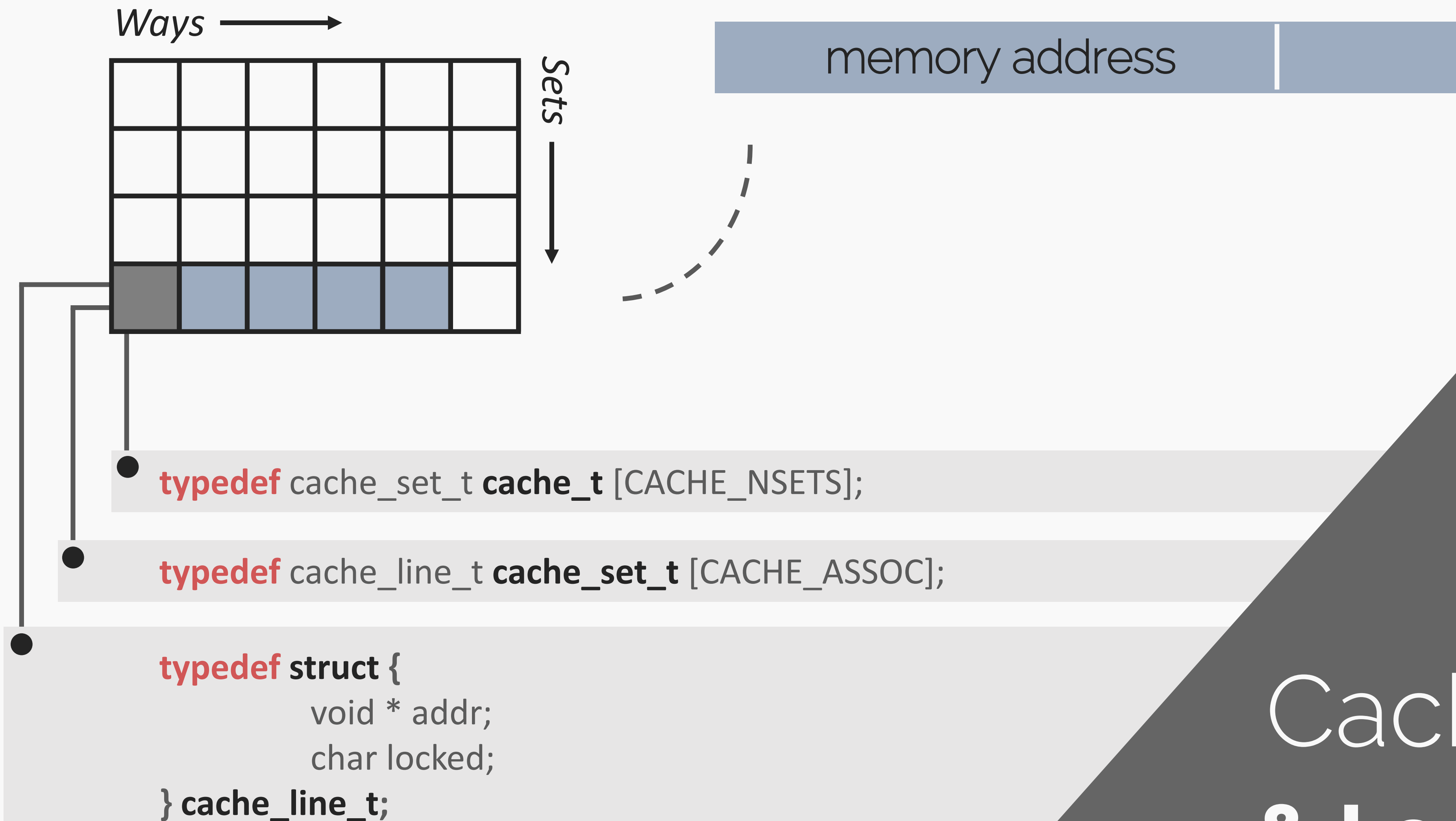
Cache Model
& Locking

Set Associative Cache



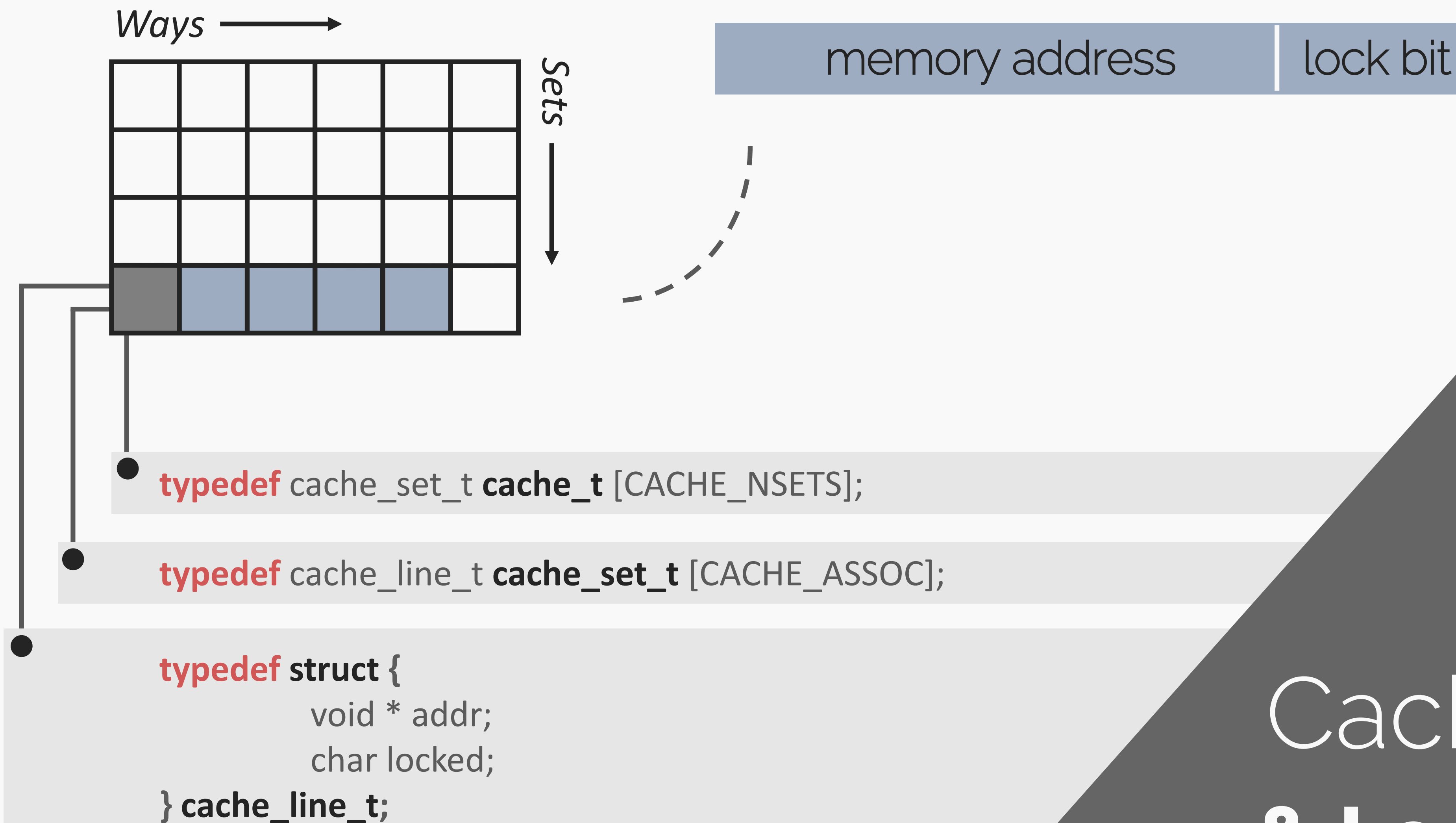
Cache Model
& Locking

Set Associative Cache



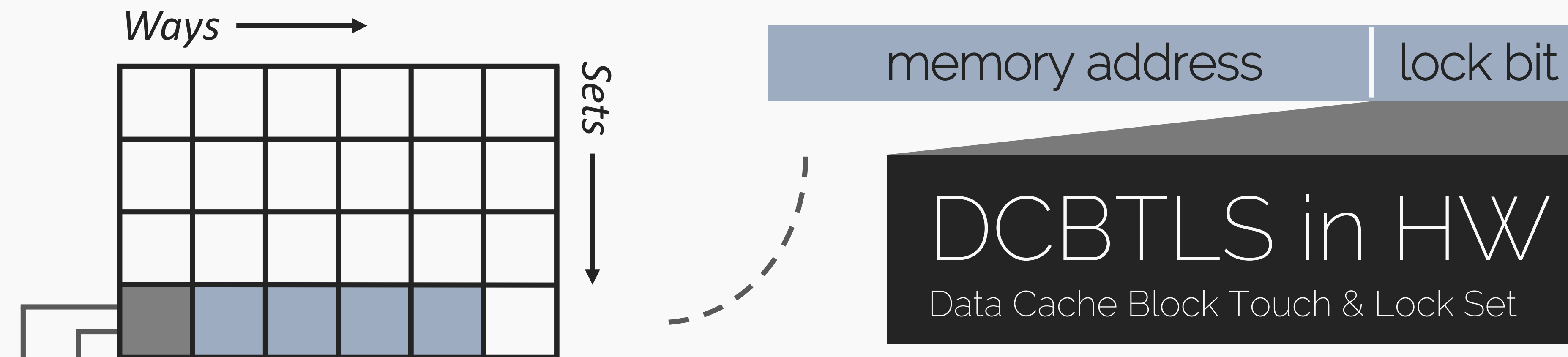
Cache Model
& Locking

Set Associative Cache



Cache Model
& Locking

Set Associative Cache



```
typedef cache_set_t cache_t [CACHE_NSETS];
```

```
typedef cache_line_t cache_set_t [CACHE_ASSOC];
```

```
typedef struct {  
    void * addr;  
    char locked;  
} cache_line_t;
```

Cache Model
& Locking

get_user_pages(...)

used to pin physical pages, returns page descriptors from virtual addresses

Abstraction of **Kernel Routine**

main parameters

(as of kernel 3.0-rc7)

— **struct** **task_struct** * tsk

— **struct** **mm_struct** * mm

— **unsigned long** start

— **unsigned long** nr_pages

— **int** write, **int** force

— **struct** **page** ** pages

get_user_pages(...)

used to pin physical pages, returns page descriptors from virtual addresses

main parameters

(as of kernel 3.0-rc7)

— **struct task_struct** * tsk

— **struct mm_struct** * mm

— **unsigned long** start

— **unsigned long** nr_pages

— **int** write, **int** force

— **struct page** ** pages

Abstraction of **Kernel Routine**

{

}

get_user_pages(...)

used to pin physical pages, returns page descriptors from virtual addresses

main parameters

(as of kernel 3.0-rc7)

— **struct task_struct** * tsk

— **struct mm_struct** * mm

— **unsigned long** start

— **unsigned long** nr_pages

— **int** write, **int** force

— **struct page** ** pages

Abstraction of **Kernel Routine**

```
{  
    struct page * page_ptr;
```

```
}
```

get_user_pages(...)

used to pin physical pages, returns page descriptors from virtual addresses

main parameters

(as of kernel 3.0-rc7)

— **struct task_struct** * tsk

— **struct mm_struct** * mm

— **unsigned long** start

— **unsigned long** nr_pages

— **int** write, **int** force

— **struct page** ** pages

Abstraction of **Kernel Routine**

```
{  
    struct page * page_ptr;  
    assert(nr_pages == 1);  

```

```
}
```

get_user_pages(...)

used to pin physical pages, returns page descriptors from virtual addresses

main parameters

(as of kernel 3.0-rc7)

— **struct task_struct** * tsk

— **struct mm_struct** * mm

— **unsigned long** start

— **unsigned long** nr_pages

— **int** write, **int** force

— **struct page** ** pages

Abstraction of **Kernel Routine**

```
{  
    struct page * page_ptr;  
    assert(nr_pages == 1);  
    assert(write == 0);
```

```
}
```


get_user_pages(...)

used to pin physical pages, returns page descriptors from virtual addresses

main parameters

(as of kernel 3.0-rc7)

— **struct task_struct** * tsk

— **struct mm_struct** * mm

— **unsigned long** start

— **unsigned long** nr_pages

— **int** write, **int** force

— **struct page** ** pages

Abstraction of **Kernel Routine**

```
{  
    struct page * page_ptr;  
    assert(nr_pages == 1);  
    assert(write == 0);  
    assert(force == 0);
```

```
}
```

get_user_pages(...)

used to pin physical pages, returns page descriptors from virtual addresses

main parameters

(as of kernel 3.0-rc7)

— **struct task_struct** * tsk

— **struct mm_struct** * mm

— **unsigned long** start

— **unsigned long** nr_pages

— **int** write, **int** force

— **struct page** ** pages

Abstraction of **Kernel Routine**

```
{  
    struct page * page_ptr;  
    assert(nr_pages == 1);  
    assert(write == 0);  
    assert(force == 0);  
    assert(tsk == current);
```

```
}
```

get_user_pages(...)

used to pin physical pages, returns page descriptors from virtual addresses

main parameters

(as of kernel 3.0-rc7)

— **struct task_struct** * tsk

— **struct mm_struct** * mm

— **unsigned long** start

— **unsigned long** nr_pages

— **int** write, **int** force

— **struct page** ** pages

Abstraction of **Kernel Routine**

```
{  
    struct page * page_ptr;  
    assert(nr_pages == 1);  
    assert(write == 0);  
    assert(force == 0);  
    assert(tsk == current);  
    assert(mm == tsk->mm);
```

```
}
```

get_user_pages(...)

used to pin physical pages, returns page descriptors from virtual addresses

main parameters

(as of kernel 3.0-rc7)

— **struct task_struct** * tsk

— **struct mm_struct** * mm

— **unsigned long** start

— **unsigned long** nr_pages

— **int** write, **int** force

— **struct page** ** pages

Abstraction of Kernel Routine

```
{
    struct page * page_ptr;
    assert(nr_pages == 1);
    assert(write == 0);
    assert(force == 0);
    assert(tsk == current);
    assert(mm == tsk->mm);
    page_ptr = __CPROVER_ui_void_ptr(...);
}
```

get_user_pages(...)

used to pin physical pages, returns page descriptors from virtual addresses

main parameters

(as of kernel 3.0-rc7)

— **struct task_struct** * tsk

— **struct mm_struct** * mm

— **unsigned long** start

— **unsigned long** nr_pages

— **int** write, **int** force

— **struct page** ** pages

Abstraction of Kernel Routine

```
{
    struct page * page_ptr;
    assert(nr_pages == 1);
    assert(write == 0);
    assert(force == 0);
    assert(tsk == current);
    assert(mm == tsk->mm);
    page_ptr = __CPROVER_ui_void_ptr(...);
    __CPROVER_assume(page_ptr >= mem_map &&
                     page_ptr < (mem_map + MAX_PAGES));
}
```

get_user_pages(...)

used to pin physical pages, returns page descriptors from virtual addresses

main parameters

(as of kernel 3.0-rc7)

— **struct task_struct** * tsk

— **struct mm_struct** * mm

— **unsigned long** start

— **unsigned long** nr_pages

— **int** write, **int** force

— **struct page** ** pages

Abstraction of Kernel Routine

```
{
    struct page * page_ptr;
    assert(nr_pages == 1);
    assert(write == 0);
    assert(force == 0);
    assert(tsk == current);
    assert(mm == tsk->mm);
    page_ptr = __CPROVER_ui_void_ptr(...);
    __CPROVER_assume(page_ptr >= mem_map &&
                     page_ptr < (mem_map + MAX_PAGES));
    __CPROVER_assume(ALIGNED_TO(page_ptr,
                                sizeof(struct page)));
}
```


get_user_pages(...)

used to pin physical pages, returns page descriptors from virtual addresses

main parameters

(as of kernel 3.0-rc7)

— **struct task_struct** * tsk

— **struct mm_struct** * mm

— **unsigned long** start

— **unsigned long** nr_pages

— **int** write, **int** force

— **struct page** ** pages

Abstraction of Kernel Routine

```
{
    struct page * page_ptr;
    assert(nr_pages == 1);
    assert(write == 0);
    assert(force == 0);
    assert(tsk == current);
    assert(mm == tsk->mm);
    page_ptr = __CPROVER_ui_void_ptr(...);
    __CPROVER_assume(page_ptr >= mem_map &&
                     page_ptr < (mem_map + MAX_PAGES));
    __CPROVER_assume(ALIGNED_TO(page_ptr,
                                sizeof(struct page)));
    *pages = page_ptr;
}
```

get_user_pages(...)

used to pin physical pages, returns page descriptors from virtual addresses

main parameters

(as of kernel 3.0-rc7)

— **struct task_struct** * tsk

— **struct mm_struct** * mm

— **unsigned long** start

— **unsigned long** nr_pages

— **int** write, **int** force

— **struct page** ** pages

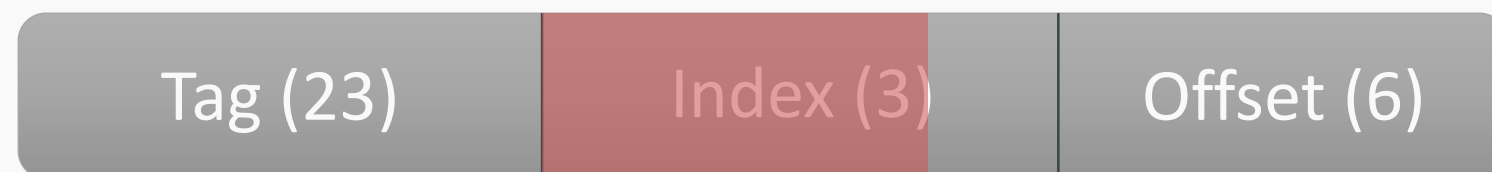
Abstraction of Kernel Routine

```
{
    struct page * page_ptr;
    assert(nr_pages == 1);
    assert(write == 0);
    assert(force == 0);
    assert(tsk == current);
    assert(mm == tsk->mm);
    page_ptr = __CPROVER_ui_void_ptr(...);
    __CPROVER_assume(page_ptr >= mem_map &&
                     page_ptr < (mem_map + MAX_PAGES));
    __CPROVER_assume(ALIGNED_TO(page_ptr,
                                 sizeof(struct page)));
    *pages = page_ptr;
    return 1;
}
```

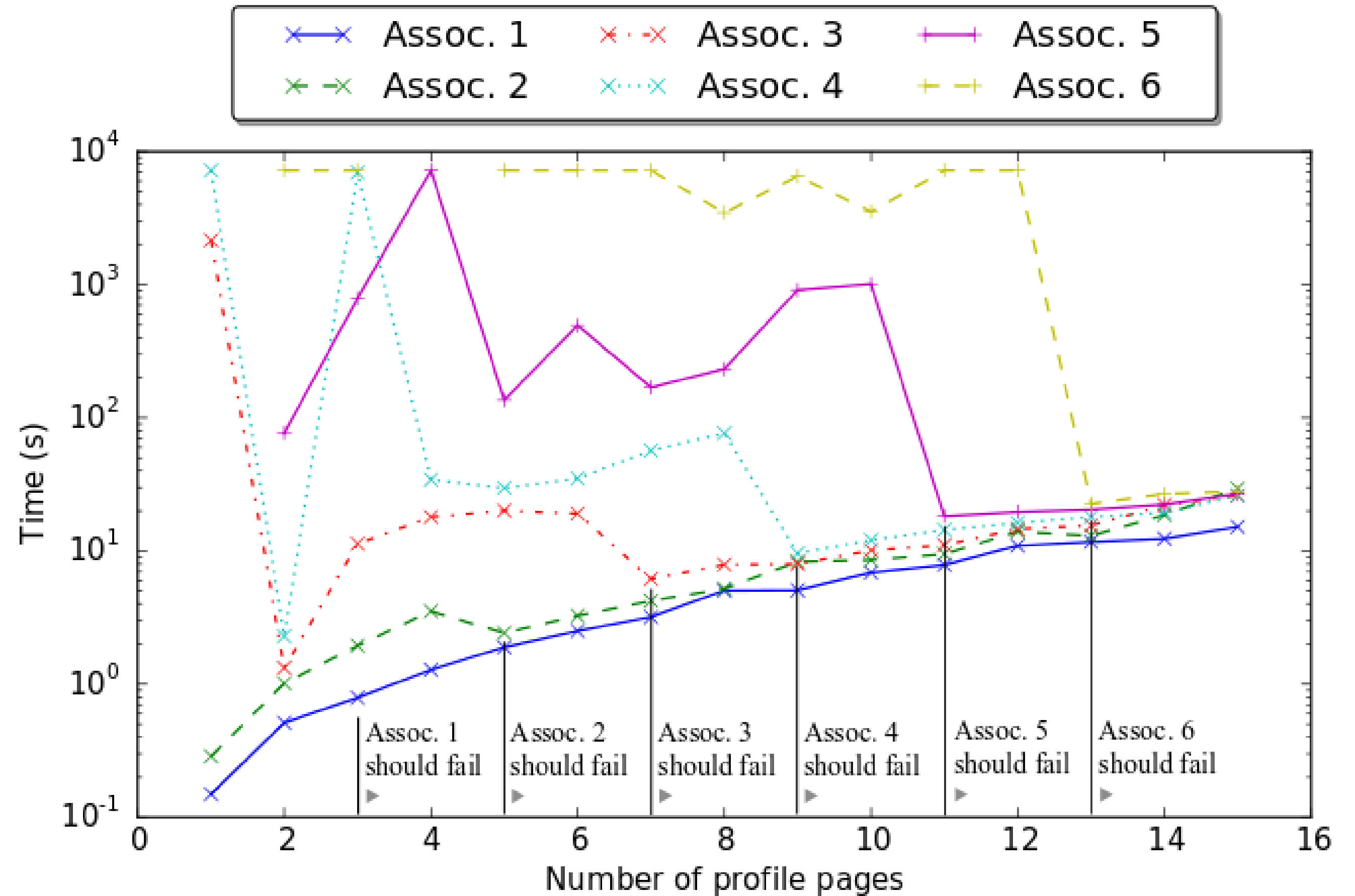
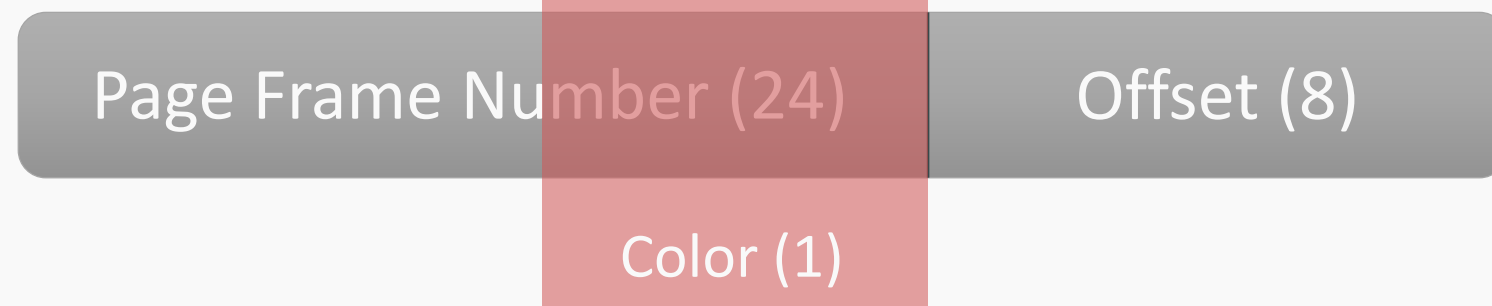
Evaluation

Memory Instance #1

From cache controller's perspective



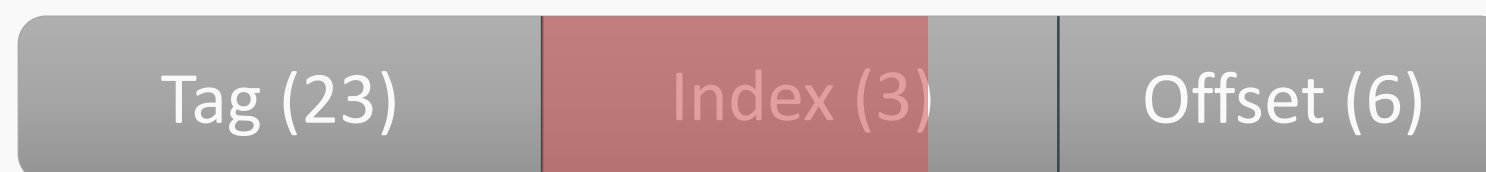
From OS's perspective



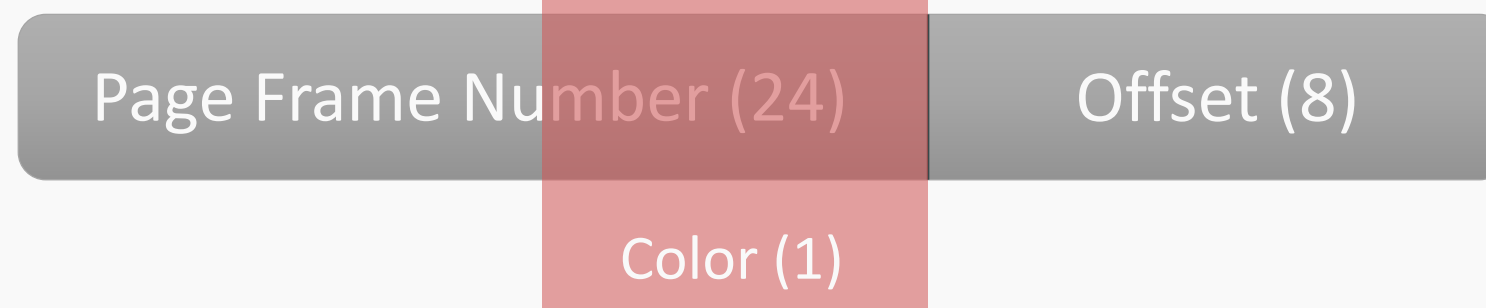
Evaluation

Memory Instance #1

From cache controller's perspective

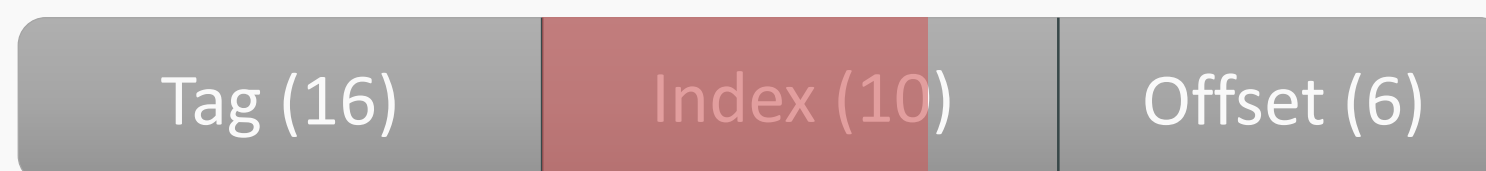


From OS's perspective

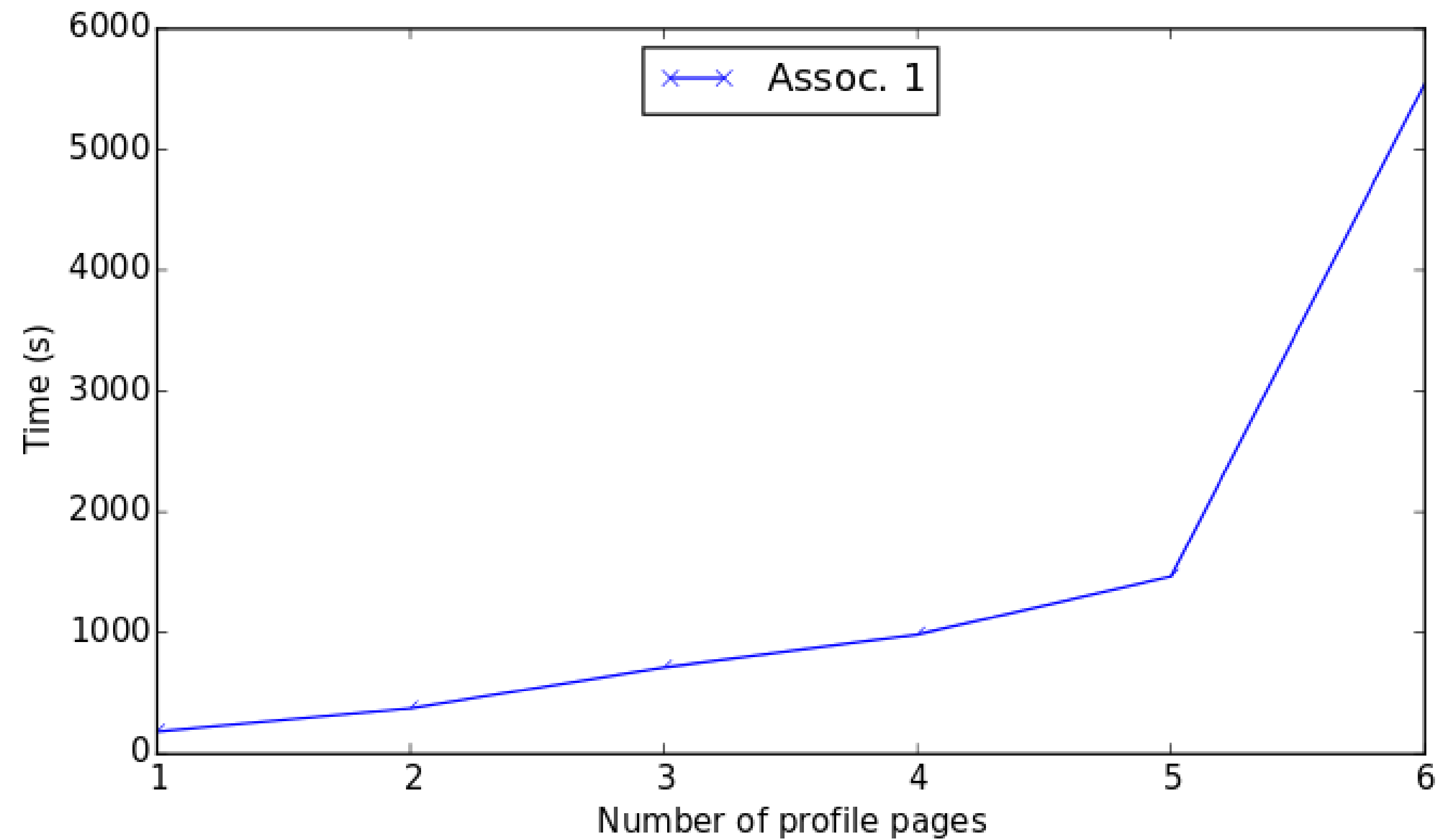
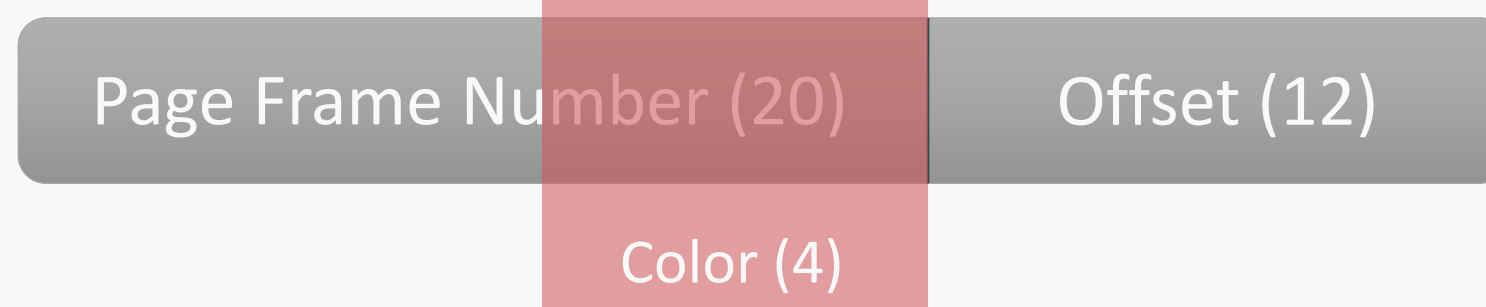


Memory Instance #2

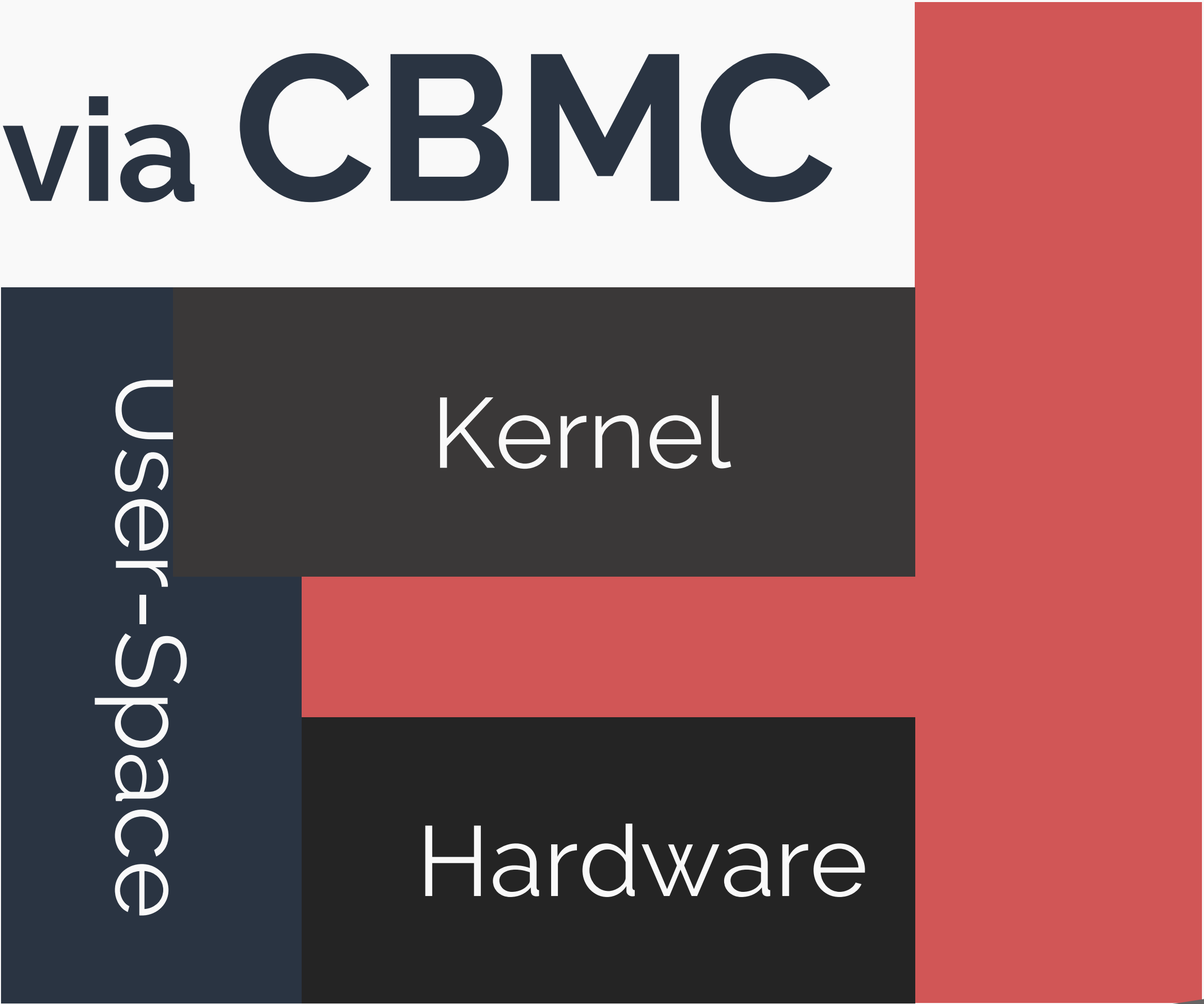
From cache controller's perspective



From OS's perspective

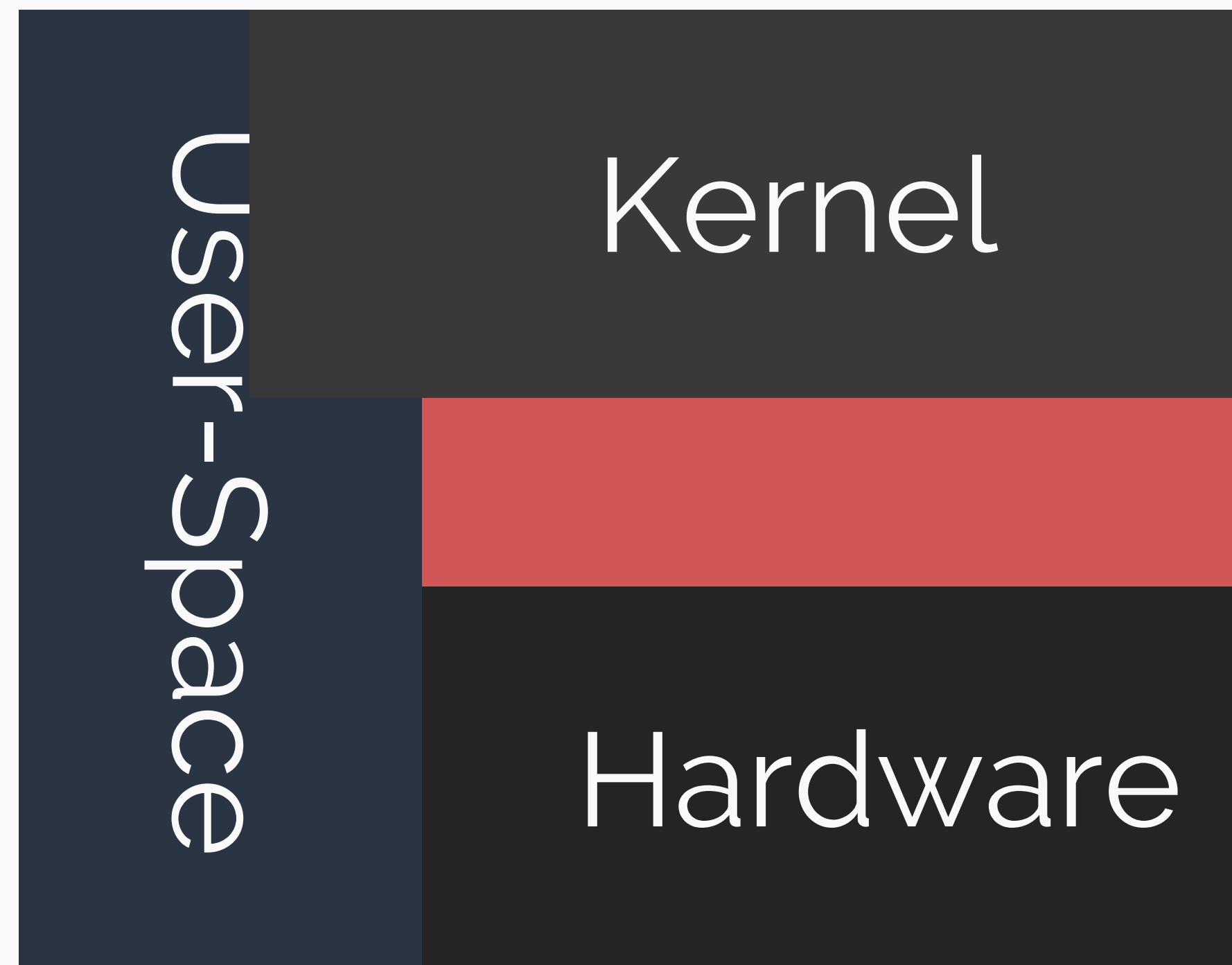


Colored Lockdown via CBMC



Summary

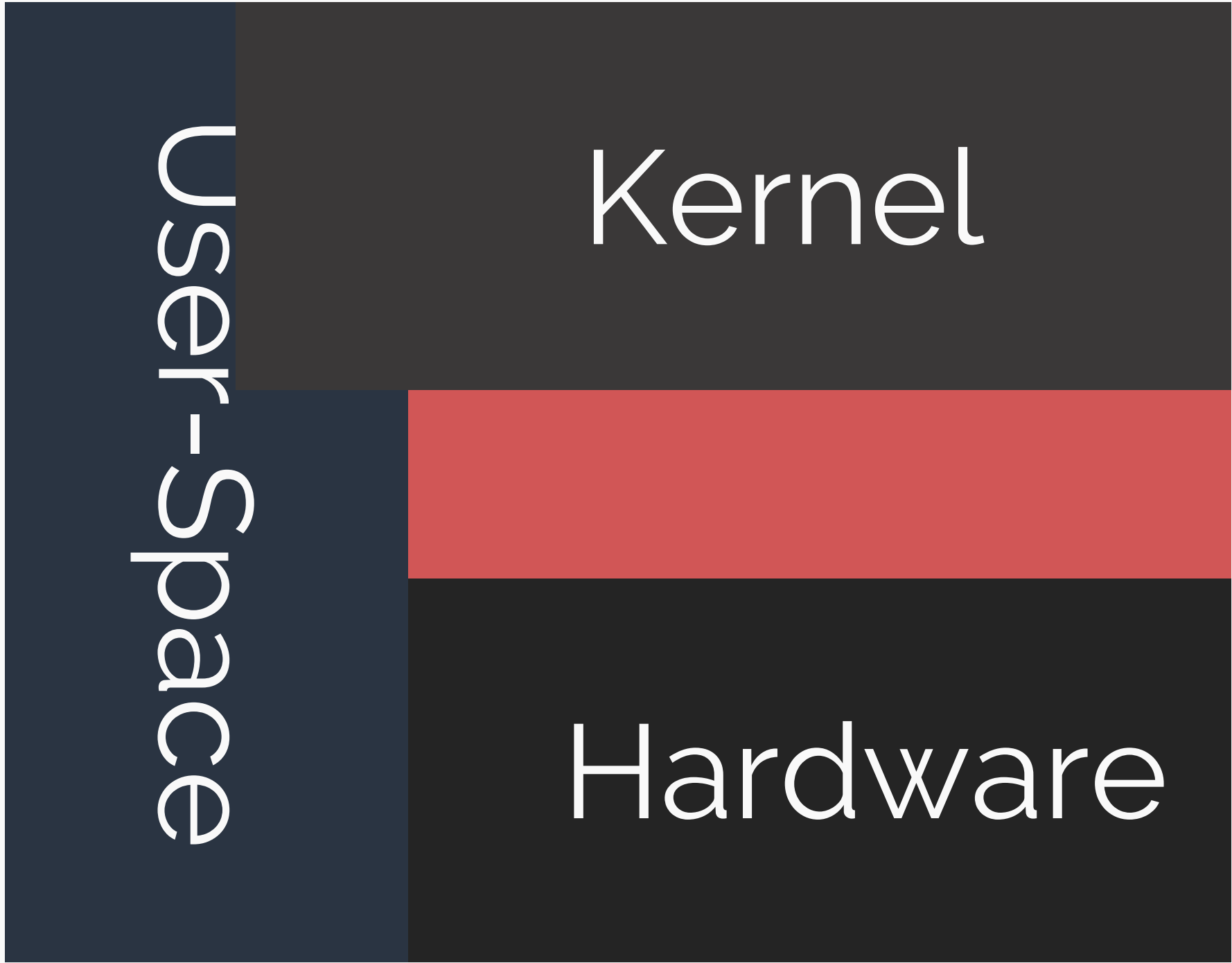
Colored Lockdown via CBMC



successful with
realistic memory model

Summary

Colored Lockdown via CBMC



scalability needed
+ + + + + + +

Summary



Thanks.

rmancuso@bu.edu

BU