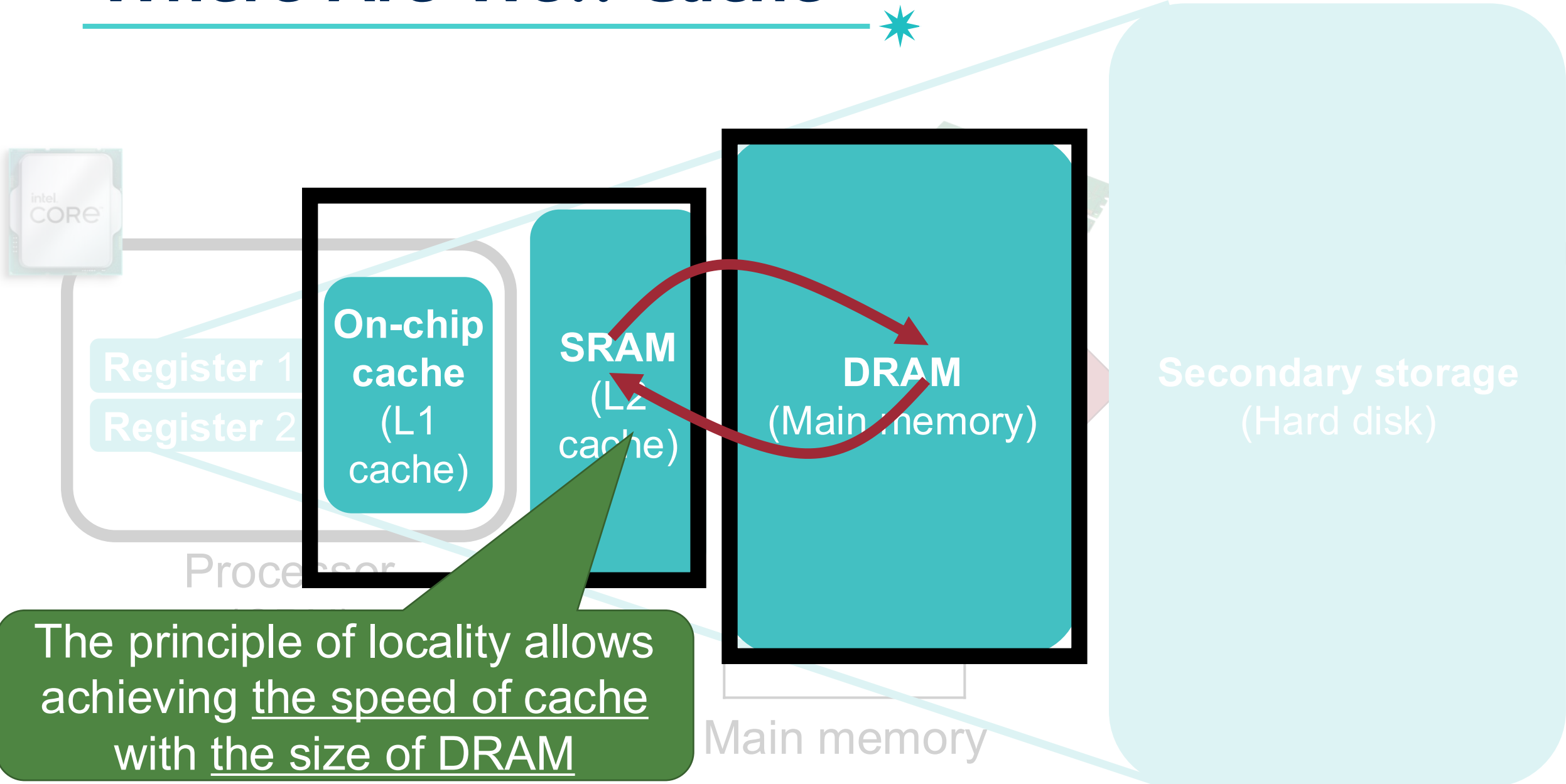


# CSE251: System Programming

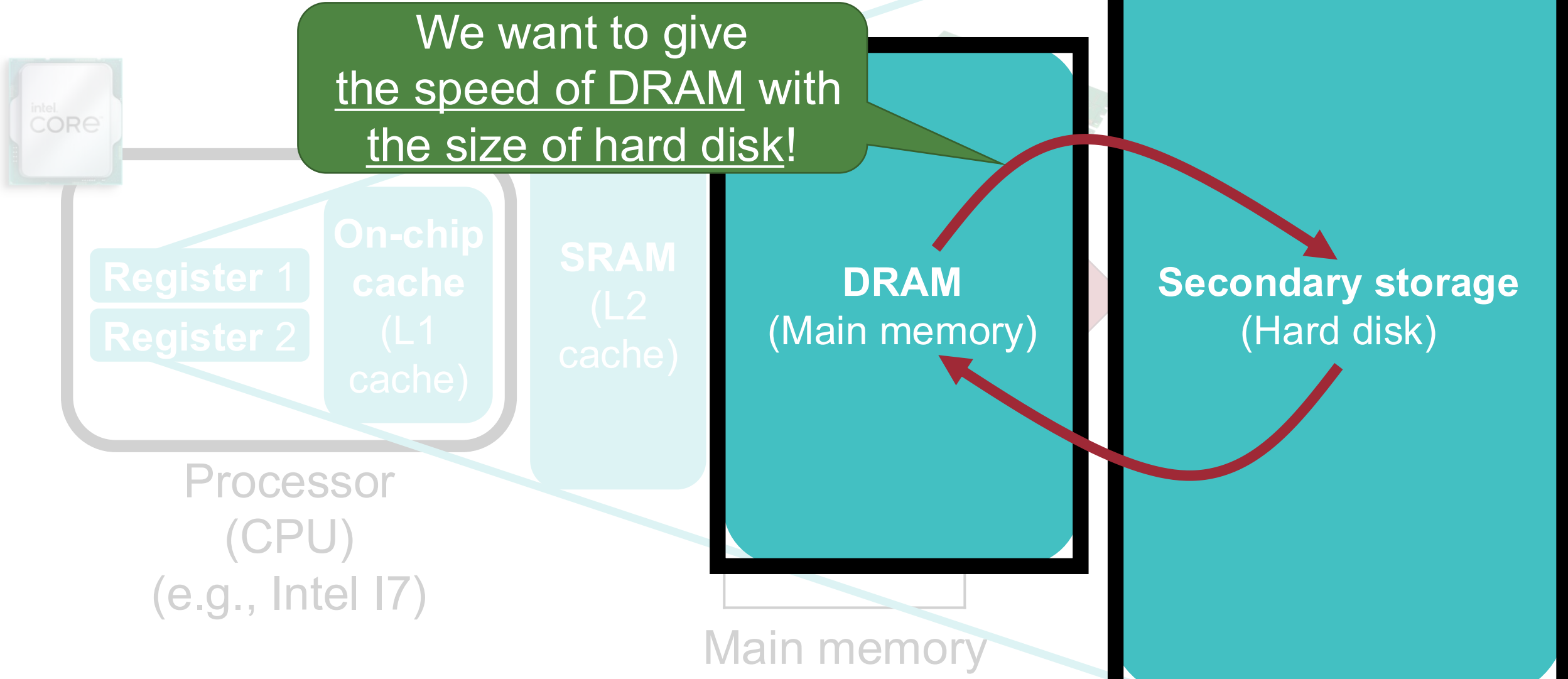
## 17. Dynamic Memory Allocation (1)

Seongil Wi

# Where Are We?: Cache



# Where Are We?: Virtual Memory



# Today's Topic: Dynamic Memory Allocation <sup>1</sup>

- Programmers use dynamic memory allocators (e.g., `malloc`) to acquire virtual memory **at runtime**

*Why do we need this?*



# Problems of Statically Declared Variables

5

```
int main(void) {
    int scores[100];
    int n;

    printf("Number of students: ");
    scanf("%d", &n);

    for (int i = 0; i < n; i++) {
        scanf("%d", &scores[i]);
    }

    return 0;
}
```

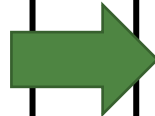
The array size is fixed before runtime

- If  $n > 100$ , the program may access memory out of bounds
- If  $n < 100$ , unused space is wasted

n is known only at runtime

# Motivation of Dynamic Memory Allocation

```
int main(void) {  
    int scores[100];  
    int n;  
  
    printf("Number of students: ");  
    scanf("%d", &n);  
  
    for (int i = 0; i < n; i++) {  
        scanf("%d", &scores[i]);  
    }  
  
    return 0;  
}
```



```
int main(void) {  
  
    int n;  
  
    printf("Number of students: ");  
    scanf("%d", &n);  
  
    int *scores = malloc(n*sizeof(int));  
    for (int i = 0; i < n; i++) {  
        scanf("%d", &scores[i]);  
    }  
  
    return 0;  
}
```

# Motivation of Dynamic Memory Allocation <sup>7</sup>

```
int main(void) {  
    int s;  
    int n;  
  
    printf("Number of students: ");  
    scanf("%d", &n);  
  
    for (int i = 0; i < n; i++) {  
        scanf("%d", &scores[i]);  
    }  
}
```

malloc allocates exactly  
enough space for n integers

```
int main(void) {  
    int n;  
  
    printf("Number of students: ");  
    scanf("%d", &n);  
  
    int *scores = malloc(n*sizeof(int));  
    for (int i = 0; i < n; i++) {  
        scanf("%d", &scores[i]);  
    }  
}
```

Dynamic memory allocation is needed when  
the required memory size is *known only at runtime*

# Dynamic Memory Allocation

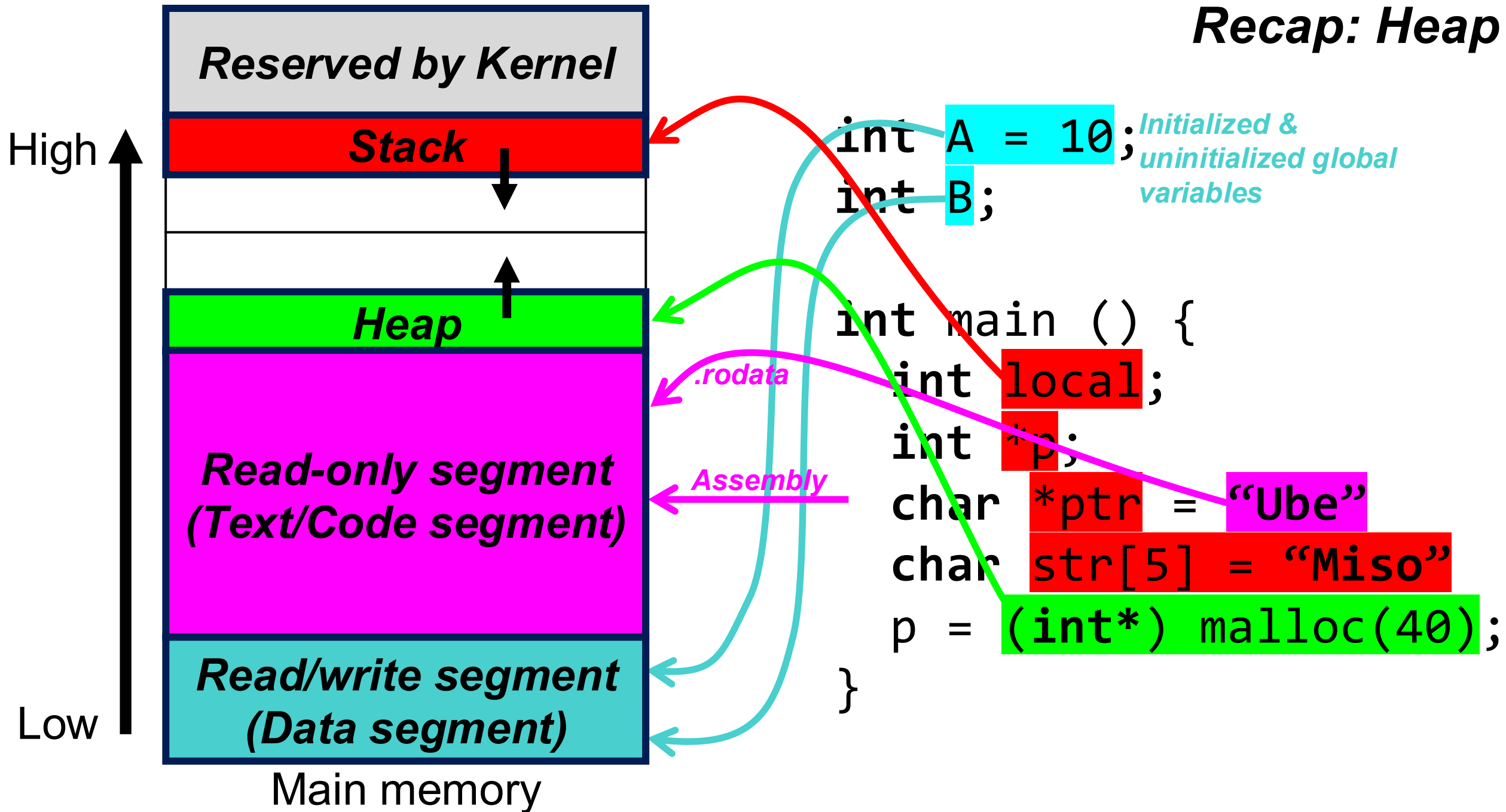
# Dynamic Memory Allocation

---



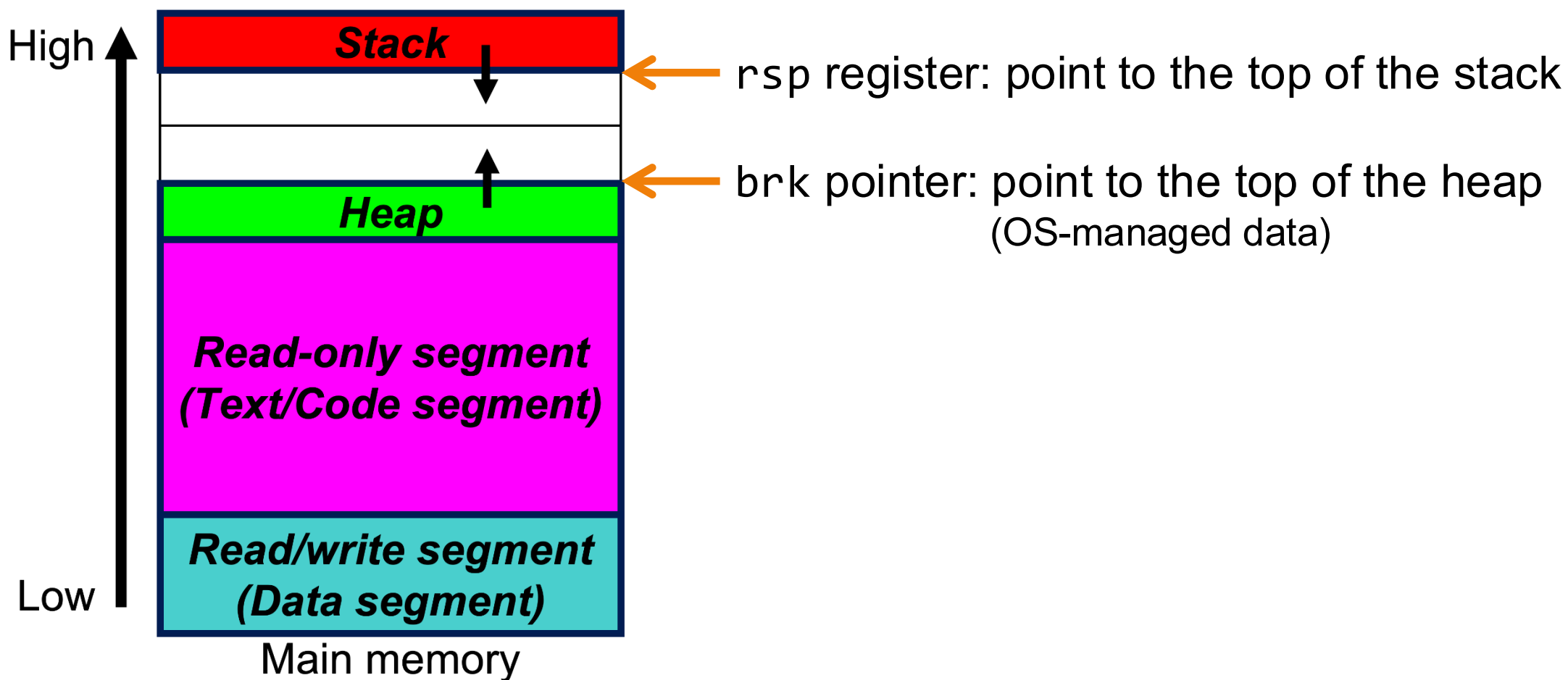
- Programmers use dynamic memory allocators (e.g., `malloc`) to acquire virtual memory **at runtime**
  - For data structures whose size is only known at runtime (It is very common to receive runtime data from user input via `stdin`, files, or the network)
- Dynamic memory allocators manage an area of process virtual memory known as the *heap*

# Recap: Heap



# Heap

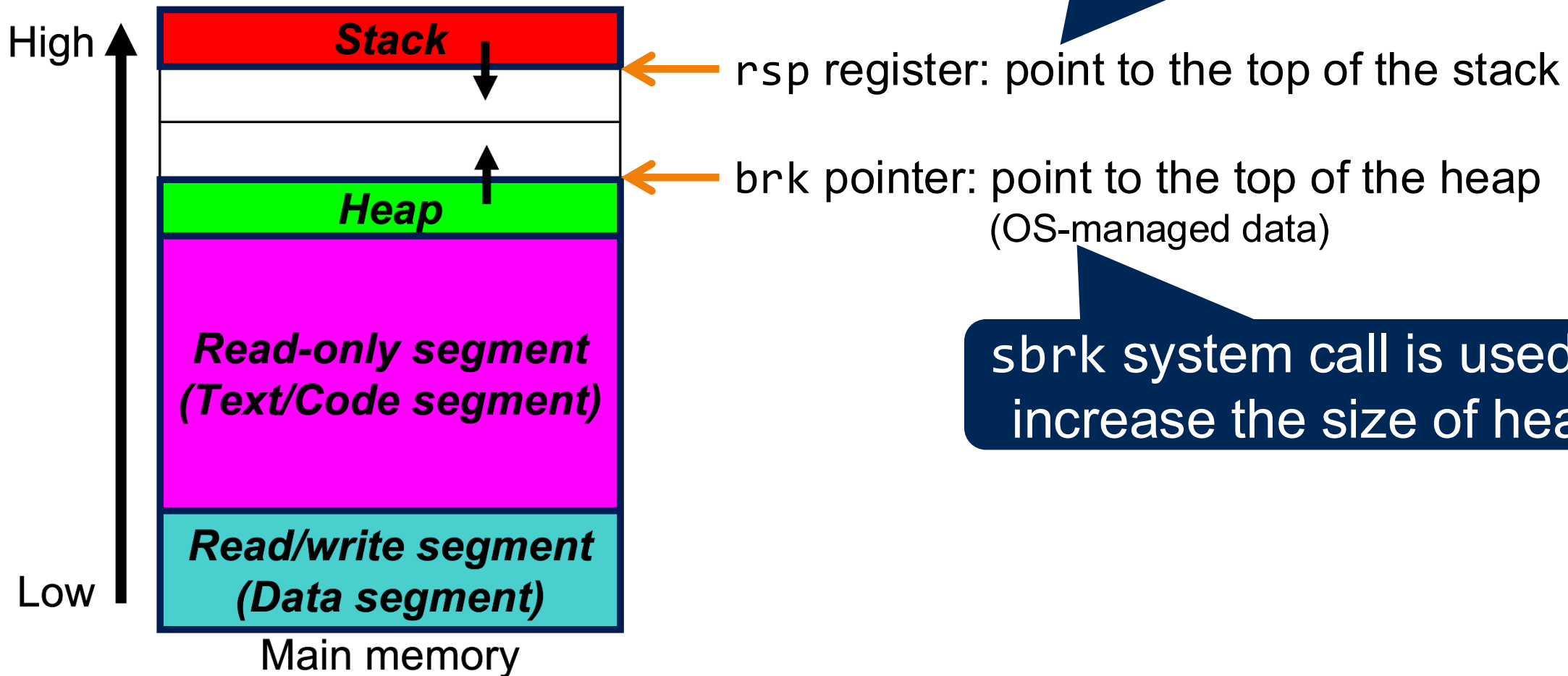
- Heap is a memory area where objects are dynamically allocated and freed



# Heap

- Heap is a memory area where objects are allocated and freed

push instruction is used to increase the size of stack



sbrk system call is used to increase the size of heap

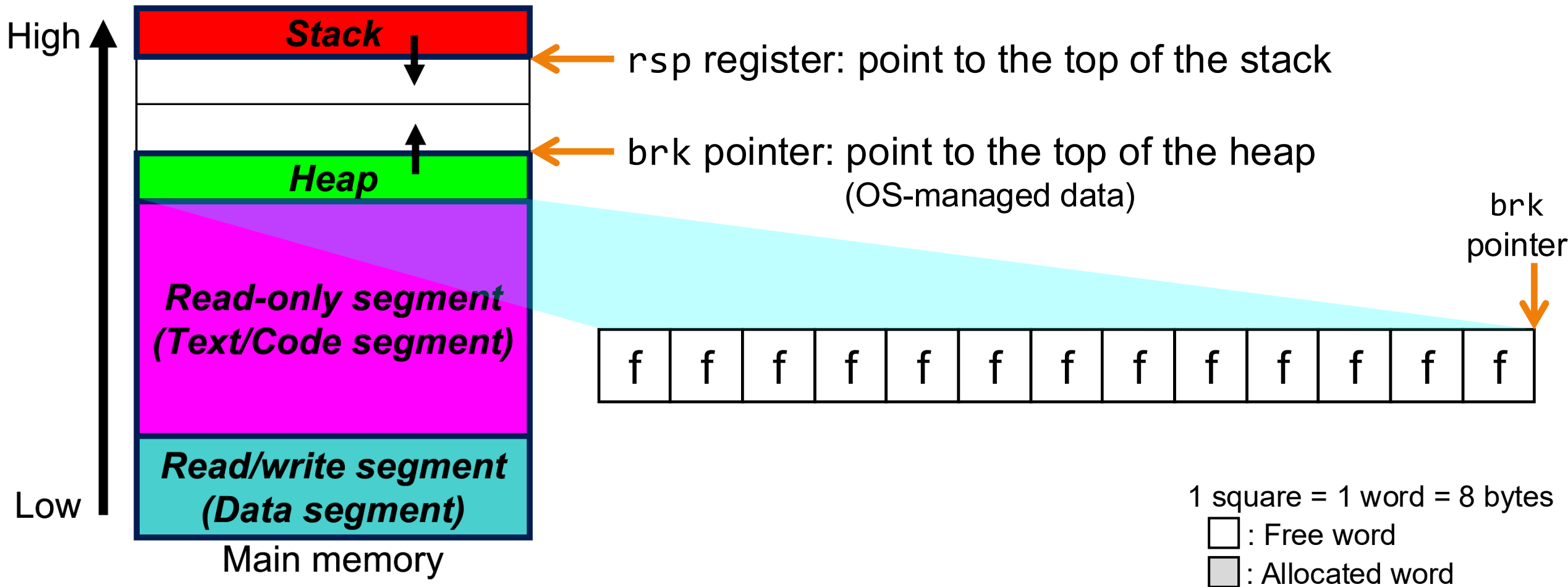
# Heap Allocator (a.k.a, Dynamic Memory Allocator, Heap Manager)

13

- Maintains heap as collection of variable sized **blocks**, which are either **allocated** or **free**

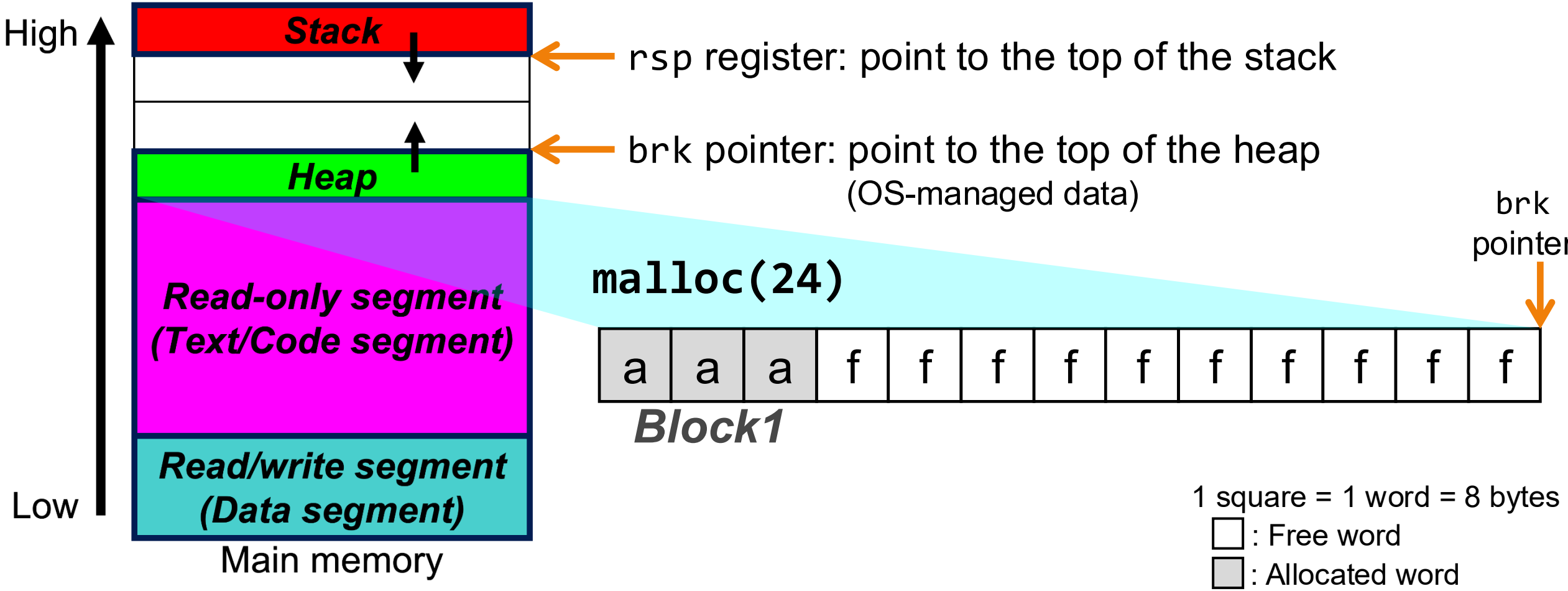
# Heap Allocator (a.k.a, Dynamic Memory Allocator, Heap Manager)

- Maintains heap as collection of variable sized **blocks**, which are either **allocated** or **free**



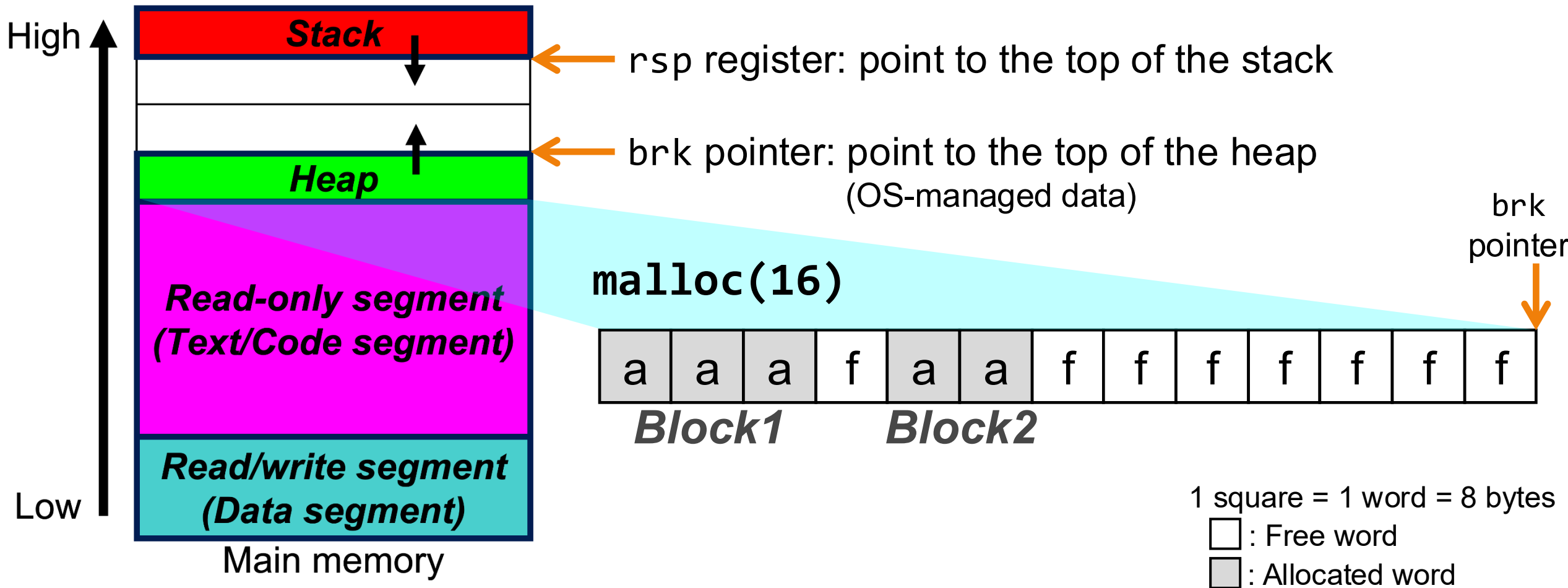
# Heap Allocator (a.k.a, Dynamic Memory Allocator, Heap Manager)

- Maintains heap as collection of variable sized **blocks**, which are either **allocated** or **free**



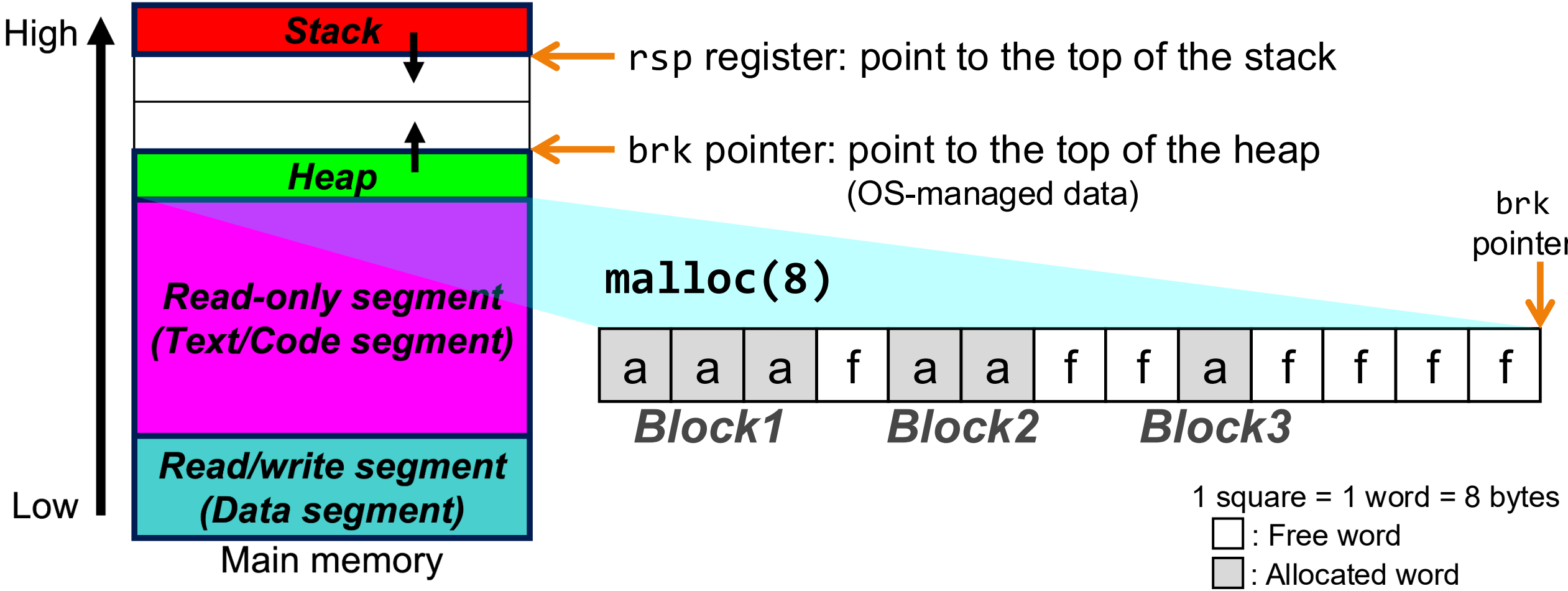
# Heap Allocator (a.k.a, Dynamic Memory Allocator, Heap Manager)

- Maintains heap as collection of variable sized **blocks**, which are either **allocated** or **free**



# Heap Allocator (a.k.a, Dynamic Memory Allocator, Heap Manager)

- Maintains heap as collection of variable sized **blocks**, which are either **allocated** or **free**



# Heap Allocator (a.k.a, Dynamic Memory Allocator, Heap Manager)

18

- Maintains heap as collection of variable sized **blocks**, which are either **allocated** or **free**
- Types of allocators:
  - **Explicit allocator**: program allocates and frees space
    - e.g., malloc and free in C
  - **Implicit allocator**: program allocates, **but does not** free space
    - e.g., new and garbage collection in Java
- Will discuss simple explicit memory allocation today

# The malloc Package

---



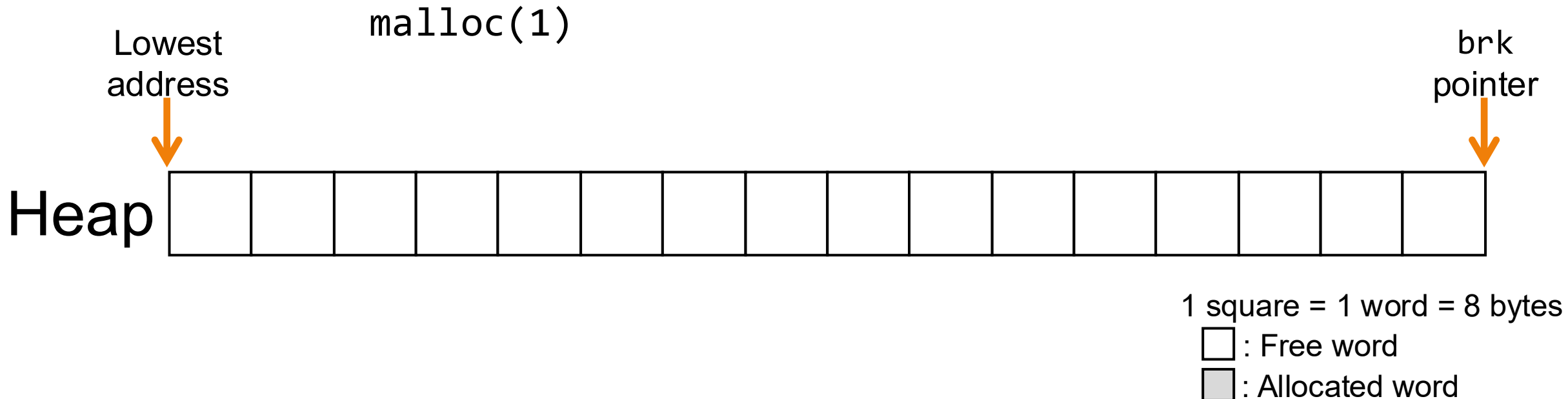
```
#include <stdlib.h>
void *malloc(size_t size)
```

- Successful:
  - Returns a pointer to a memory block of at least size bytes aligned to a 16-byte boundary (on x86-64)

# The malloc Package

```
#include <stdlib.h>
void *malloc(size_t size)
```

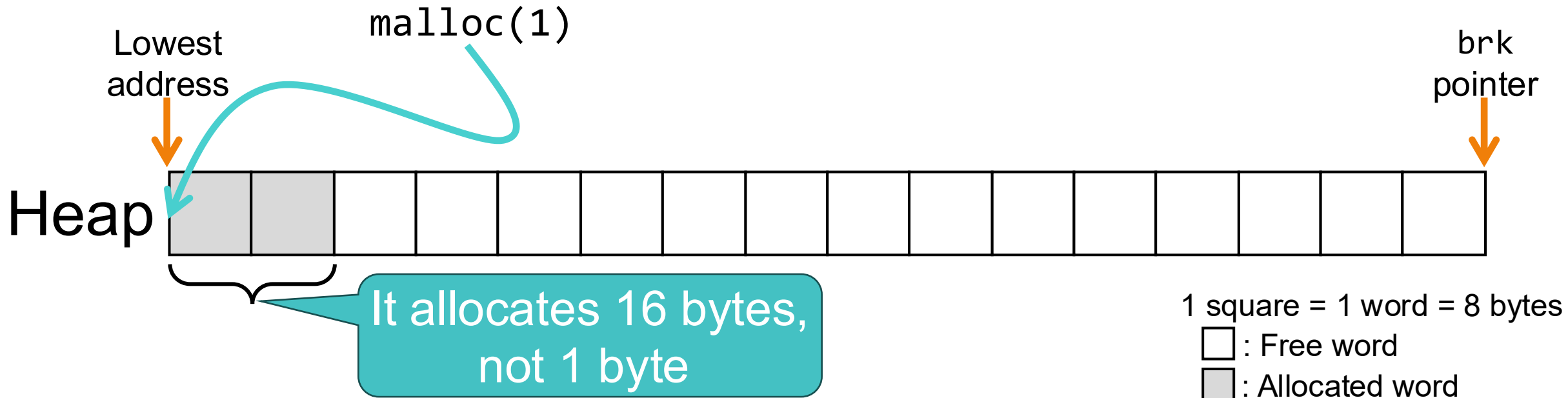
- Successful:
  - Returns a pointer to a memory block of at least size bytes aligned to a 16-byte boundary (on x86-64)



# The malloc Package

```
#include <stdlib.h>
void *malloc(size_t size)
```

- Successful:
  - Returns a pointer to a memory block of at least size bytes aligned to a 16-byte boundary (on x86-64)



# The malloc Package

---



```
#include <stdlib.h>
void *malloc(size_t size)
```

- **Successful:**
  - Returns a pointer to a memory block of at least size bytes aligned to a 16-byte boundary (on x86-64)
- **Unsuccessful:**
  - returns NULL (0) and sets errno

```
void free(void *p)
```

- Returns the block pointed at by p to pool of available memory
- p must come from a previous call to malloc, calloc, or realloc

# Other Functions

---



- `calloc`: Version of `malloc` that initializes allocated block to zero
- `realloc`: Changes the size of a previously allocated block
- `sbrk`: Used internally by allocators to grow or shrink the heap

# Malloc Example

---



```
#include <stdio.h>
#include <stdlib.h>
```

```
void foo(long n) {
    long i, *p;
```

Allocate a block of n longs

```
    p = (long *) malloc(n * sizeof(long));
```

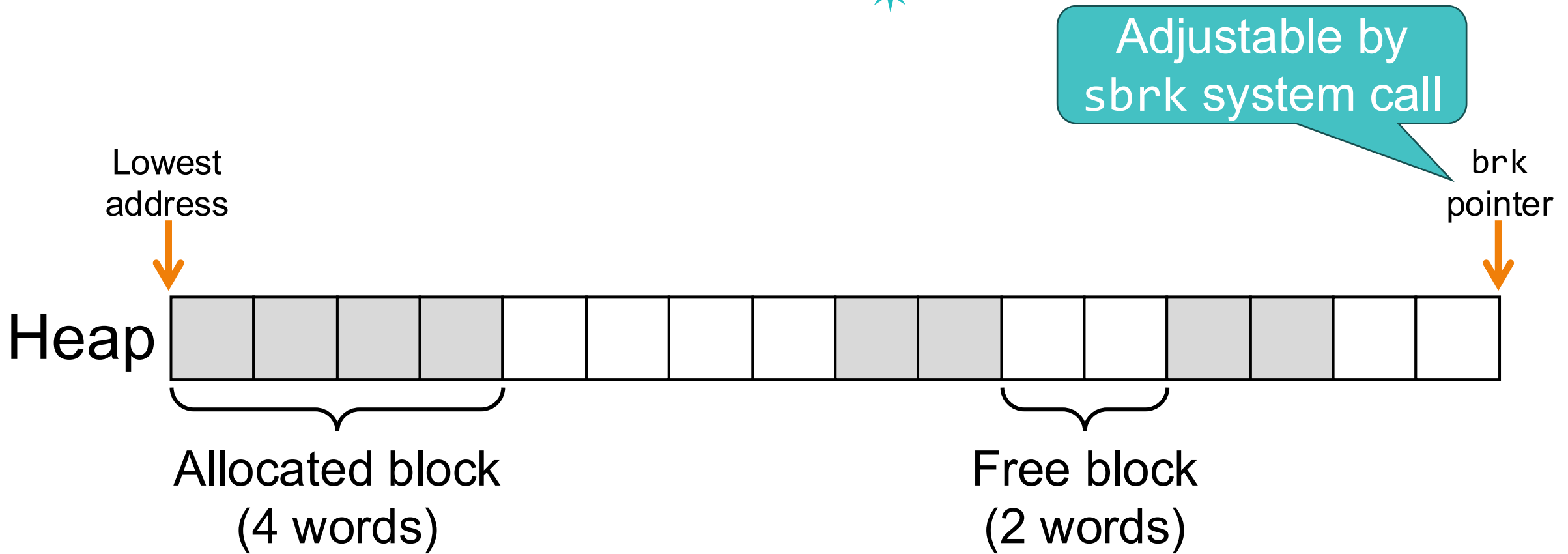
```
    for (i=0; i<n; i++)
        p[i] = i;
```

Return allocated block to the heap

```
    free(p);
```

```
}
```

# Heap Visualization Convention



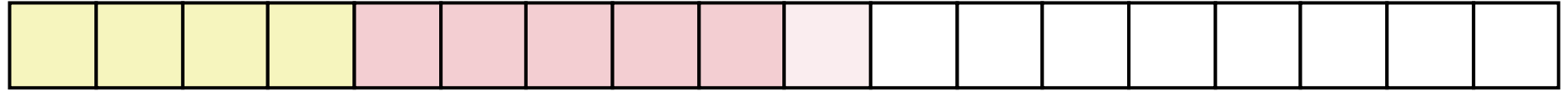
1 square = 1 word = 8 bytes  
□ : Free word  
■ : Allocated word

# Allocation Example

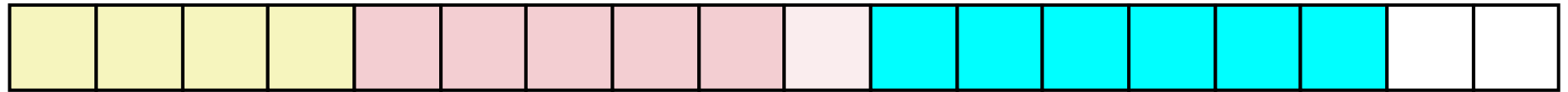
p1 = malloc(32)



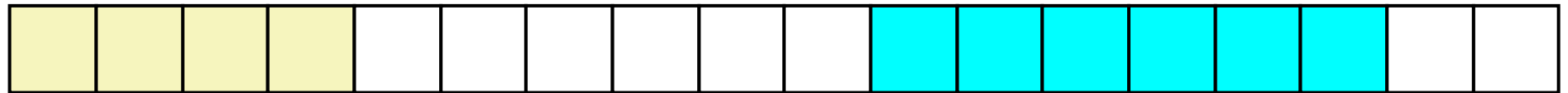
p2 = malloc(40)



p3 = malloc(48)



free(p2)



p4 = malloc(16)

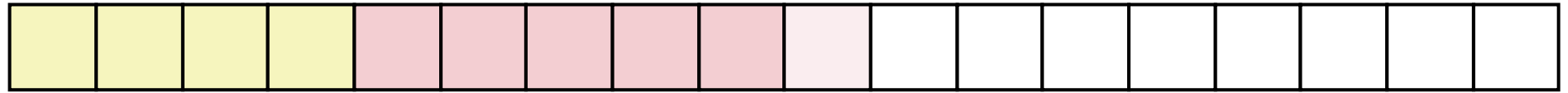


# Allocation Example

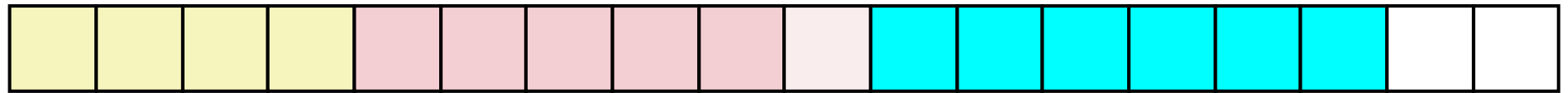
p1 = malloc(32)



p2 = malloc(40)



p3 = malloc(48)



free(p2)



Depending on the allocator's placement policy

p4 = malloc(16)



# (Traditional) Constraints of Explicit Allocator

28

- Programs can issue arbitrary sequence of `malloc` and `free` requests
- Allocators handle only the current request as best as it can (It does not know what requests will come in the future)
  - Must respond *immediately* to `malloc` requests
  - Must allocate blocks from free memory
- Must align blocks so they satisfy all alignment requirements
  - 16-byte (x86-64) alignment on 64-bit systems

# Design Goals of Heap Allocator

---



- Given some sequence of malloc and free requests:
  - $R_0, R_1, \dots, R_k, \dots, R_{n-1}$
- Goals: maximize throughput and peak memory utilization

# Throughput

---



- Measures the number of completed requests per unit time
- Example:
  - 5,000 malloc and 5,000 free calls in 10 seconds
  - Throughput = 1,000 operations/second

# Peak Memory Utilization



- Measures the allocator's space efficiency
- **Definition #1:** Aggregated payload ( $P_k$ )
  - The sum of the allocated payloads after request  $R_k$
- **Definition #2:** Current heap size ( $H_k$ )
  - The heap size after request  $R_k$  (heap only grows when allocator uses sbrk)
- **Peak memory utilization**  $U_k = \max(P_0, P_1, \dots, P_k) / H_k$ 
  - Peak memory utilization after request  $R_k$

# Peak Memory Utilization



- Measures the allocator's space efficiency

- **Definition #1:** Aggregated payload ( $P_k$ )

- The sum of the allocated payloads after request  $R_k$ .

- **Definition #2:** Current

- The heap size after request  $R_k$ .

Peak aggregate payload: the maximum aggregate payload at any point in the sequence up to request

- **Peak memory utilization**  $U_k = \max(P_0, P_1, \dots, P_k) / H_k$

- Peak memory utilization after request  $R_k$

# Peak Memory Utilization

$R_0$  `a = malloc(16)`

$P_0 = 16$  bytes

$R_1$  `b = malloc(32)`

$P_1 = 48$  bytes

$R_2$  `free(a)`

$P_2 = 32$  bytes

$R_3$  `c = malloc(48)`

$P_3 = 80$  bytes

$R_4$  `free(c)`

$P_4 = 32$  bytes

- Peak aggregate payload =  $\max(P_0, P_1, P_2, P_3, P_4) = 80$  bytes
- Current heap size  $H_4 = 128$  bytes (Assumption)
- Peak memory utilization  $U_4 = 80 / 128 = 62.5\%$

# Design Goals of Heap Allocator

- Given some sequence of malloc and free requests:

–  $R_0, R_1, \dots, R_k, \dots, R_{n-1}$

- Goals: maximize throughput and peak memory utilization

*Measures the number of completed requests per unit time*

*Measures the allocator's space efficiency*

***These goals are often conflicting***

# Any Problems?



```
p1 = malloc(32)
```



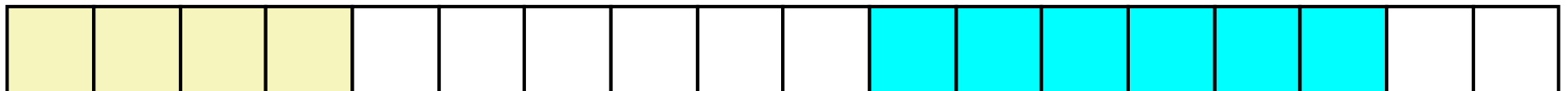
```
p2 = malloc(40)
```



```
p3 = malloc(48)
```



```
free(p2)
```



```
p4 = malloc(64)
```



# Any Problems?



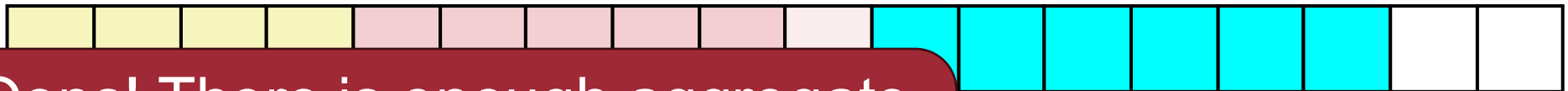
```
p1 = malloc(32)
```



```
p2 = malloc(40)
```



```
p3 = malloc(48)
```



```
free(p2)
```

Oops! There is enough aggregate heap memory, but no single free block is large enough



```
p4 = malloc(64)
```



# Any Problems?



```
p1 = malloc(32)
```



```
p2 = malloc(40)
```



```
p3 = malloc(48)
```

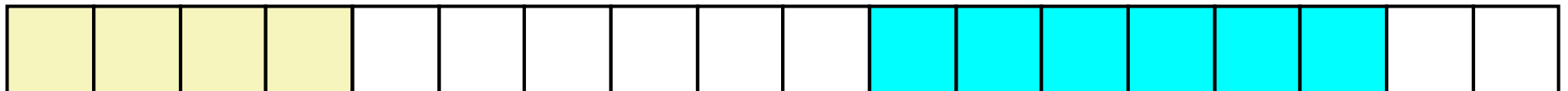


```
free(p2)
```

Fragmentation (단편화, 파편화)



```
p4 = malloc(64)
```



# Fragmentation (단편화, 파편화)

---

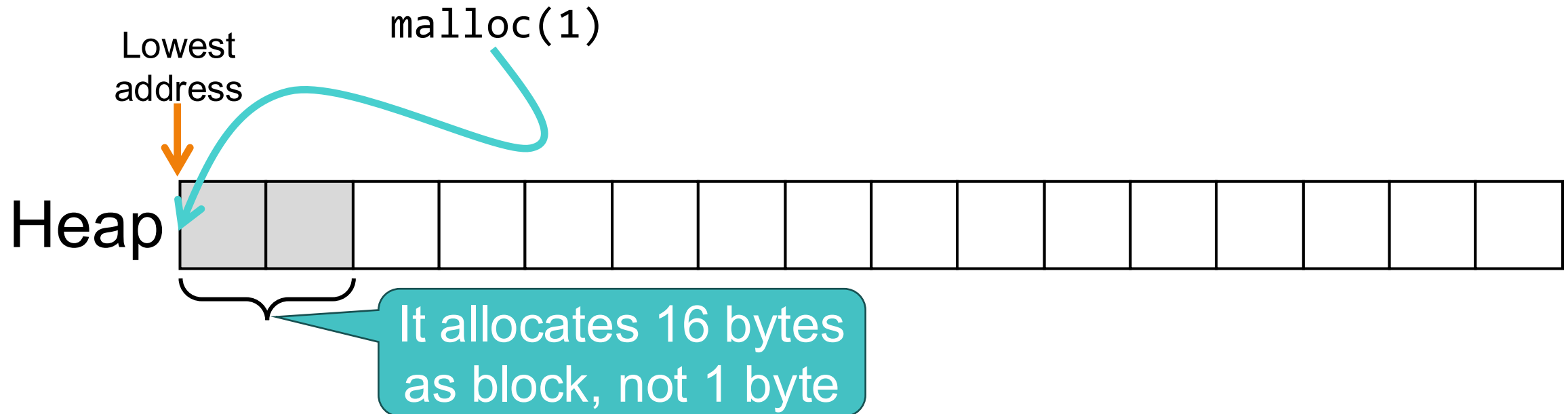


- Poor memory utilization caused by *fragmentation*
  - *Internal* fragmentation
  - *External* fragmentation

# Internal Fragmentation

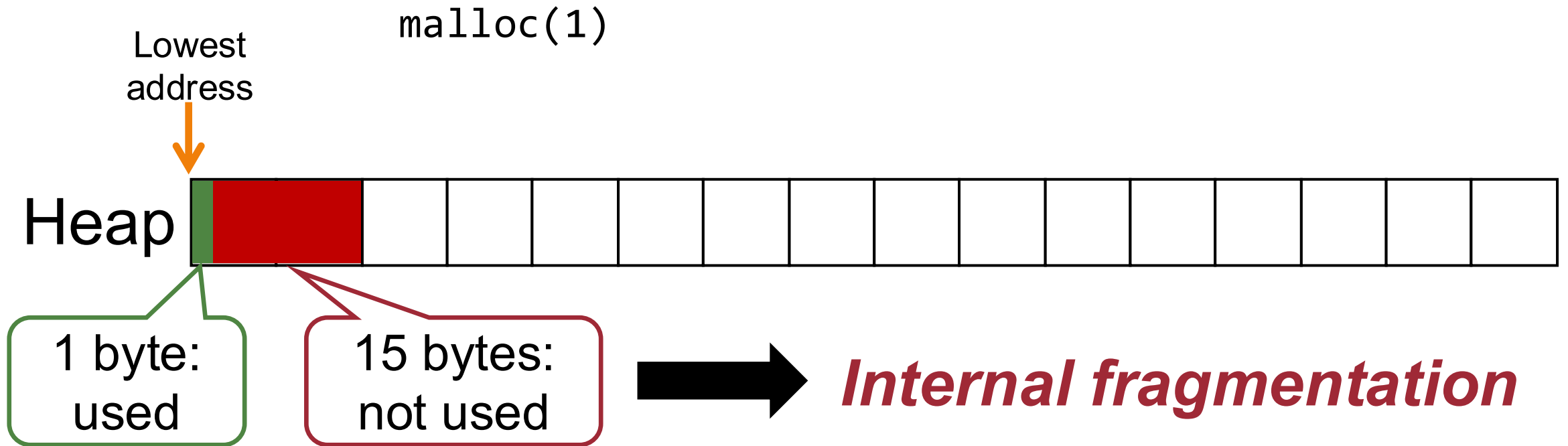


Occurs if payload is smaller than block size



# Internal Fragmentation

Occurs if payload is smaller than block size



# Internal Fragmentation

---

Occurs if payload is smaller than block size

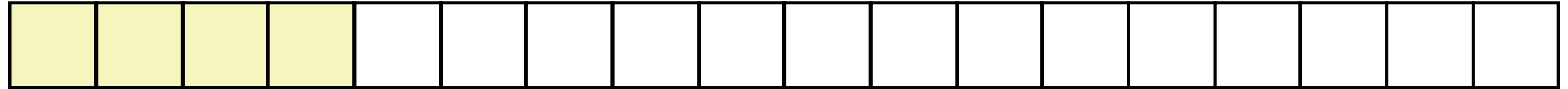
- Caused by
  - Overhead of maintaining heap data structure (e.g., header, footer ← we will cover about it later)
  - Padding for alignment purposes
  - Explicit policies (e.g., to return a big block to satisfy a small request)
- Depends only on the pattern of *previous* requests
  - Thus, easy to measure

# External Fragmentation



- Occurs when there is enough aggregate heap memory, but no single free block is large enough

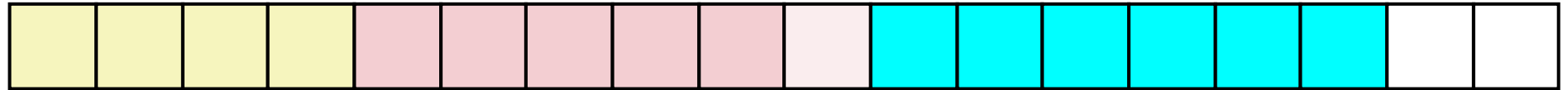
`p1 = malloc(32)`



`p2 = malloc(40)`



`p3 = malloc(48)`

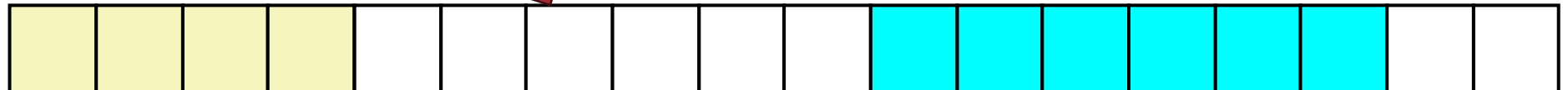


`free(p2)`

External fragmentation



`p4 = malloc(64)`



# External Fragmentation

---



- Occurs when there is enough aggregate heap memory, but no single free block is large enough
- Depends on the pattern of *future* requests
  - Thus, difficult to measure

**Question?**