# Memory and Data Locality

(The C Memory Model)

https://cs.pomona.edu/classes/cs140/

# Outline

Topics and Learning Objectives

- Motivate our discussion on memory

- Discuss memory access timing

- Discuss the C memory model

- Discuss locality

Exercise

- Sets and Lists (tangentially related)

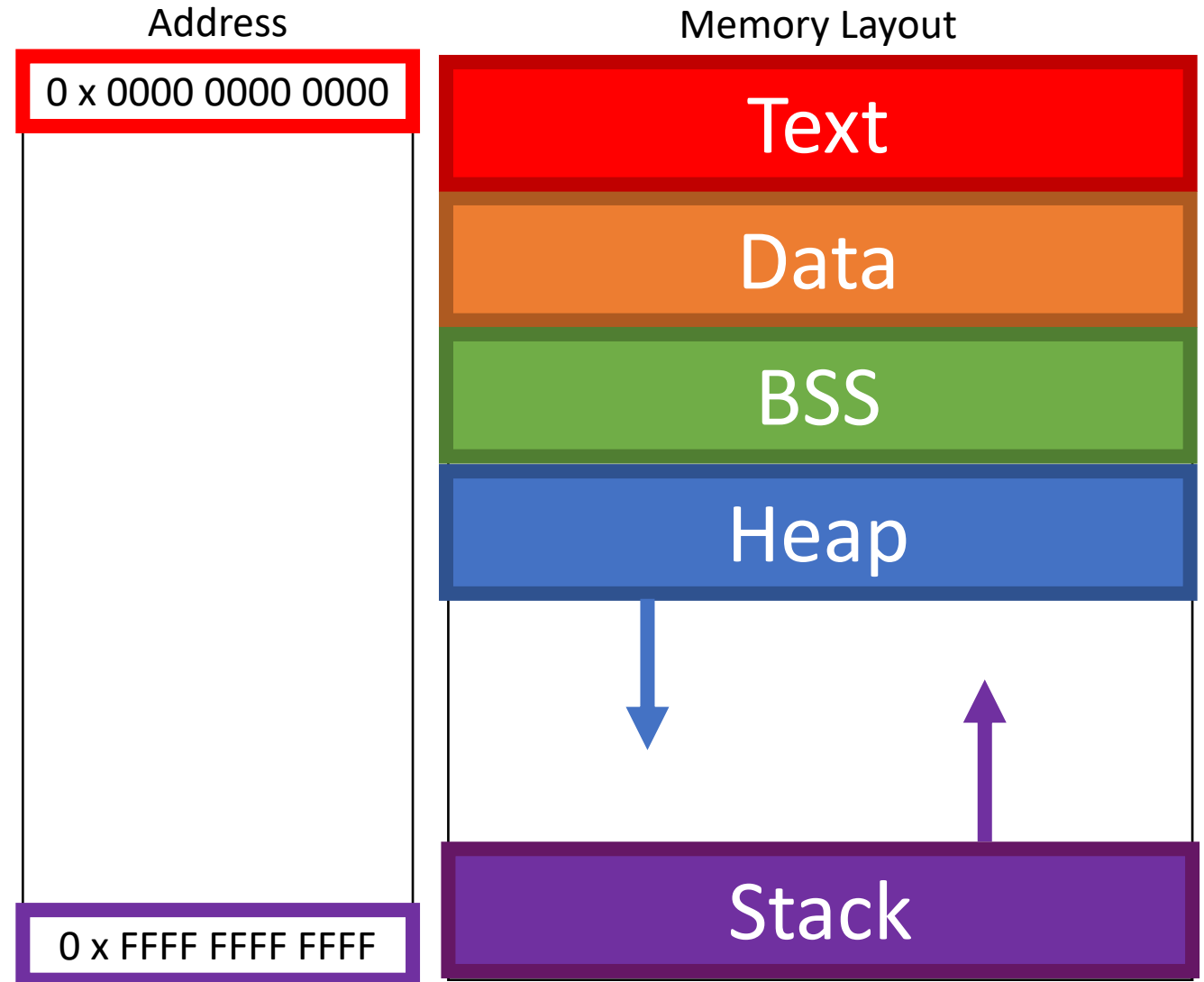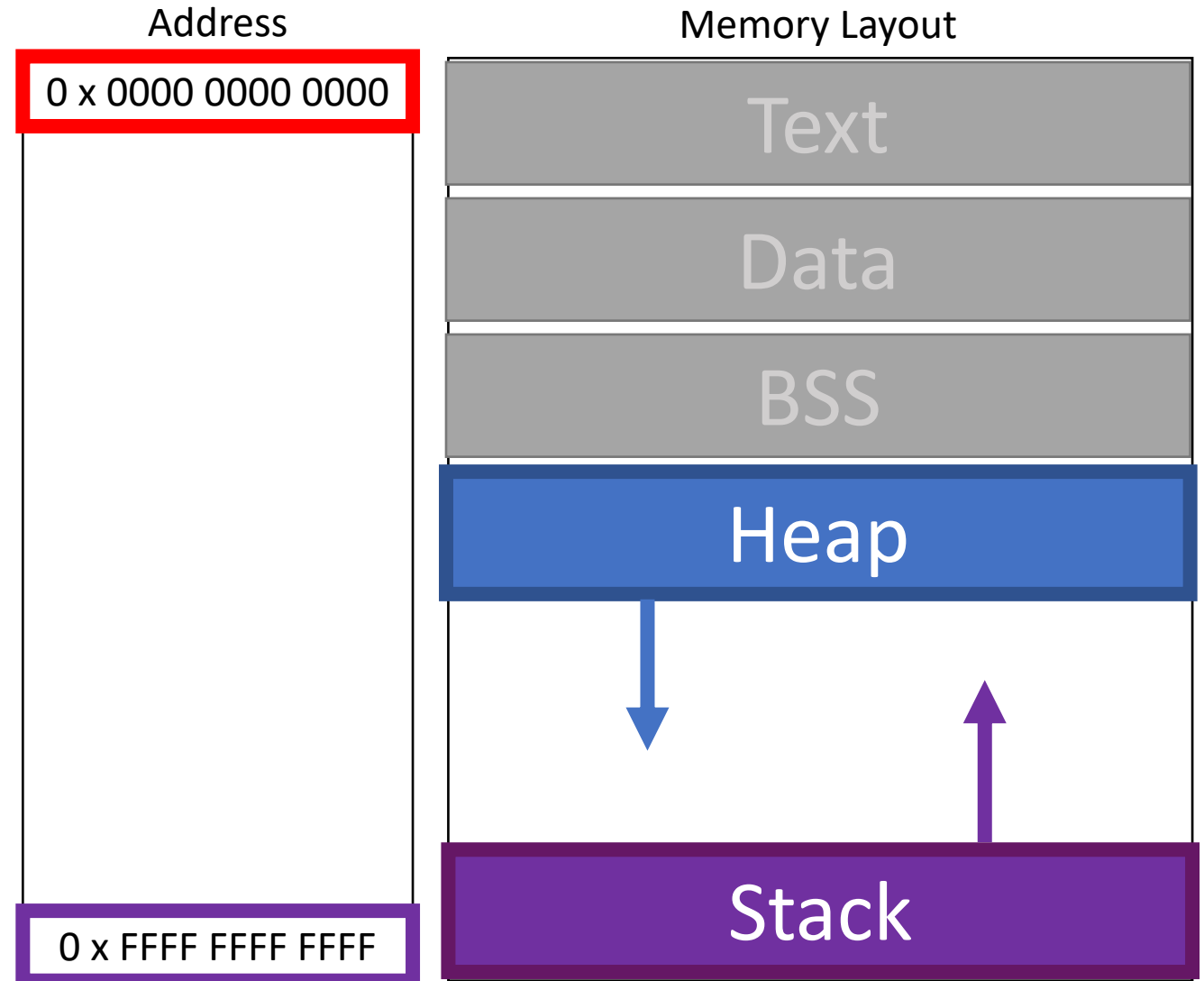# Exercise

## Process Memory

Text: contains your compiled code

Data: contains **initialized** static and global variables

BSS: contains uninitialized static and global variables

Heap: contains dynamically allocated memory

Stack: contains local variables

| Address | Memory Layout |
|---|---|
| 0 x 0000 0000 0000 | Text |
| | Data |
| | BSS |
| | Heap |
| | |
| 0 x FFFF FFFF FFFF | Stack |

4

# Process Memory

Text: contains your compiled code

Data: contains **initialized** static and global variables

BSS: contains uninitialized static and global variables
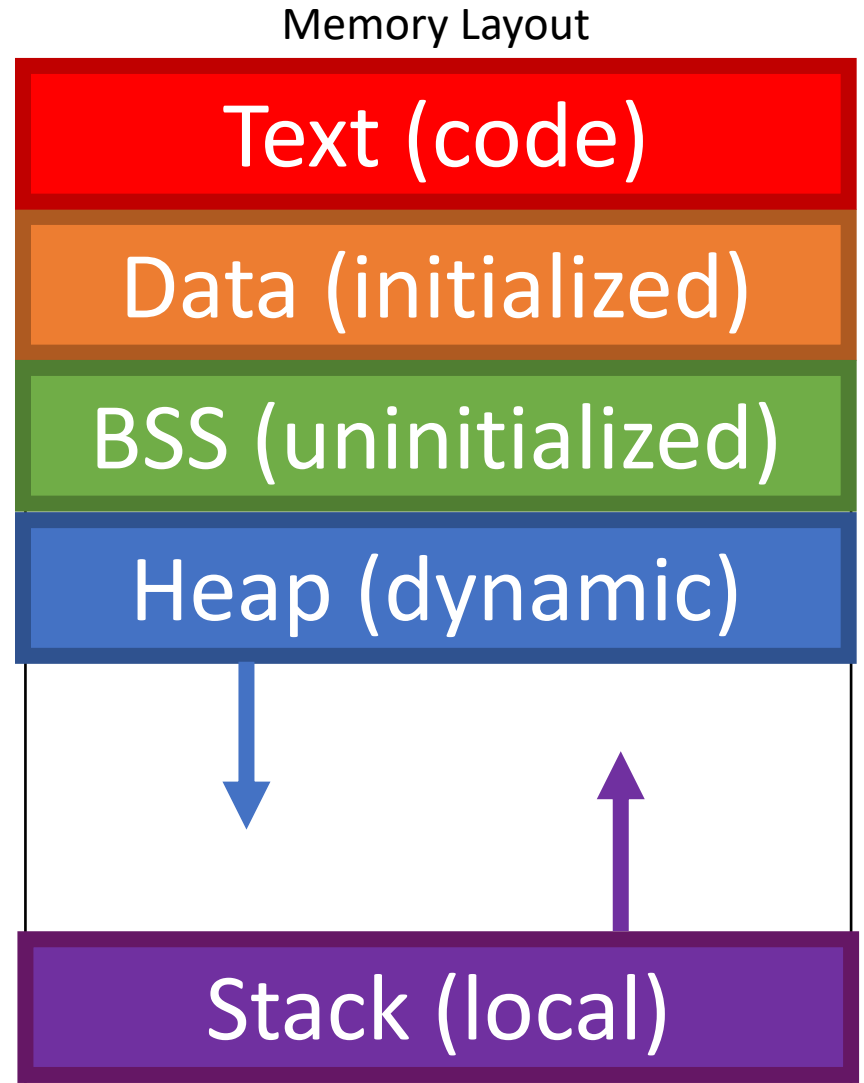
Heap: contains dynamically allocated memory

Stack: contains local variables

Address

0 x 0000 0000 0000

0 x FFFF FFFF FFFF

Memory Layout

Text

Data

BSS

Heap

Stack

Read through the code for a few moments.

```c
char* string = "hello";
int iSize;

char* f(void)
{
    char* p;
    iSize = 8;
    p = malloc(iSize);
    return p;
}
```

Memory Layout

| Text (code) |
| --- |
| Data (initialized) |
| BSS (uninitialized) |
| Heap (dynamic) |

| Stack (local) |

What goes in Text?

```
char* string = "hello";
int iSize;

char* f(void)
{
    char* p;
    iSize = 8;
    p = malloc(iSize);
    return p;
}
```

Memory Layout



Text (code)

Data (initialized)

BSS (uninitialized)

Heap (dynamic)

Stack (local)

```
char* string = "hello";
int iSize;

char* f(void)
{
    char* p;
    iSize = 8;
    p = malloc(iSize);
    return p;
}
```

Memory Layout

Text (code)

Data (initialized)

BSS (uninitialized)

Heap (dynamic)

Stack (local)

```
char* string = "hello";
int iSize;

char* f(void)
{
    char* p;
    iSize = 8;
    p = malloc(iSize);
    return p;
}
```

Memory Layout

Text (code)

Data (initialized)

BSS (uninitialized)

Heap (dynamic)

Stack (local)

What goes in Data?

```
char* string = "hello";
int iSize;

char* f(void)
{
    char* p;
    iSize = 8;
    p = malloc(iSize);
    return p;
}
```

Memory Layout

| Text (code) |
| Data (initialized) |
| BSS (uninitialized) |
| Heap (dynamic) |
| Stack (local) |

What goes in BSS?

```
char* string = "hello";
int iSize;

char* f(void)
{
    char* p;
    iSize = 8;
    p = malloc(iSize);
    return p;
}
```

Memory Layout

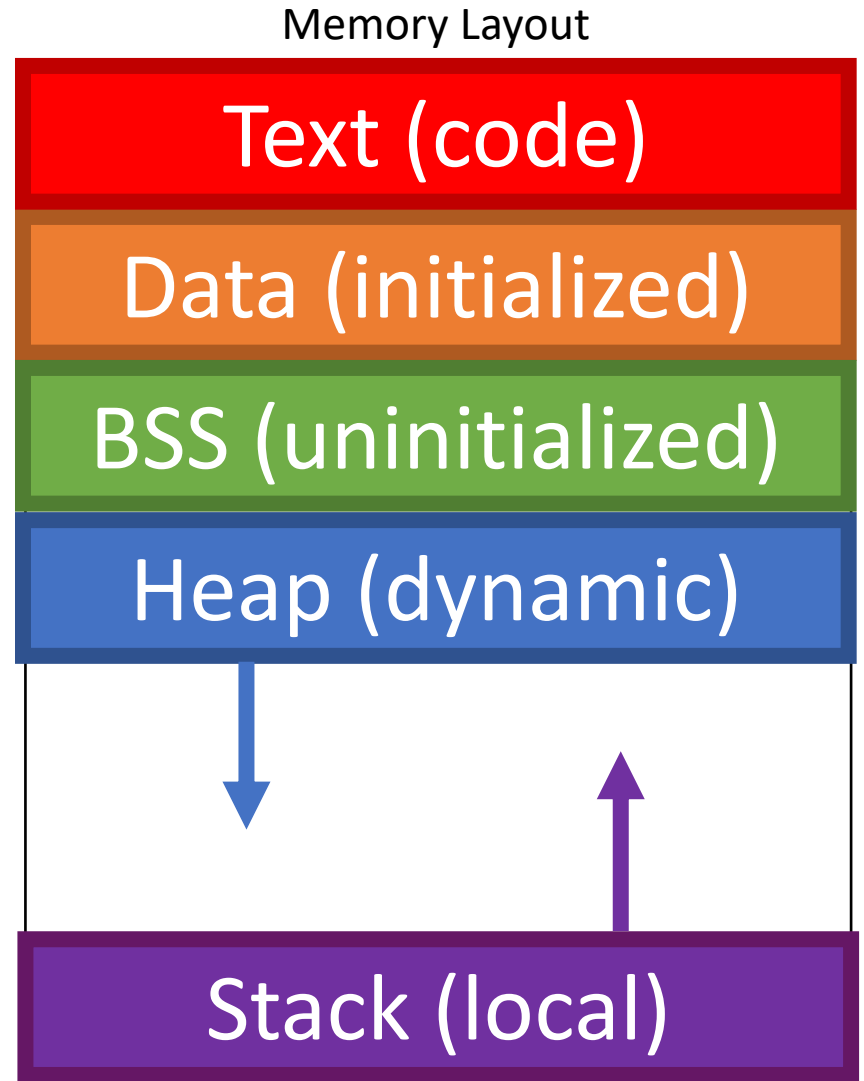| Text (code) |
| Data (initialized) |
| BSS (uninitialized) |
| Heap (dynamic) |
| Stack (local) |

11

```c
char* string = "hello";
int iSize;

char* f(void)
{
    char* p;
    iSize = 8;
    p = malloc(iSize);
    return p;
}
```

Memory Layout



Text (code)

Data (initialized)

BSS (uninitialized)

Heap (dynamic)

Stack (local)

What goes in Heap?

```
char* string = "hello";
int iSize;

char* f(void)
{
    char* p;
    iSize = 8;
    p = malloc(iSize);
    return p;
}
```

Memory Layout

Text (code)

Data (initialized)

BSS (uninitialized)

Heap (dynamic)

Stack (local)

What goes in Heap?

```c
char* string = "hello";
int iSize;

char* f(void)
{
    char* p;
    iSize = 8;
    p = malloc(iSize);
    return p;
}
```

Memory Layout

Text (code)

Data (initialized)

BSS (uninitialized)

Heap (dynamic)

Stack (local)

```
char* string = "hello";
int iSize;

char* f(void)
{
        char* p;
        iSize = 8;
        p = malloc(iSize);
        return p;
}
```

Memory Layout

| Text (code) |
| Data (initialized) |
| BSS (uninitialized) |
| Heap (dynamic) |
| Stack (local) |

What goes in Stack?

```c
char* string = "hello";
int iSize;

char* f(void)
{
    char* p;
    iSize = 8;
    p = malloc(iSize);
    return p;
}
```

Memory Layout

Text (code)

Data (initialized)

BSS (uninitialized)

Heap (dynamic)

Stack (local)

```c
char* string = "hello";
int iSize;

char* f(void)
{
    char* p;
    iSize = 8;
    p = malloc(iSize);
    return p;
}
```

Memory Layout

| Text (code) |
| Data (initialized) |
| BSS (uninitialized) |
| Heap (dynamic) |
| Stack (local) |

17

# Stack and Heap Resources

- https://www.cs.princeton.edu/courses/archive/fall07/cos217/lectures/06MemoryAllocation-3x1.pdf

- https://doc.rust-lang.org/1.6.0/book/the-stack-and-the-heap.html

- https://www.usna.edu/Users/cs/aviv/classes/ic221/s16/lec/08/lec.html

- https://manybutfinite.com/post/anatomy-of-a-program-in-memory/

- https://stackoverflow.com/questions/79923/what-and-where-are-the-stack-and-heap

# Principle of Locality

- <u>Locality</u>: programs tend to reuse data and instructions near those they have used recently

- <u>Temporal locality</u>: recently referenced items are likely to be referenced again in the near future

- <u>Spatial locality</u>: items with nearby addresses tend to be referenced close together in time

*Discontiguous data structures are the root of all evil. This is simply a fact. If you don't believe me, I'll try to convince you. Specifically, please say "no" to linked lists. OK. Please. Please say no to linked lists. There is almost nothing more harmful you can do to the performance of an actual modern microprocessor than to use a linked list data structure.*

Chandler Carruth (Engineer at Google) @ CppCon 2014

# JUST SAY NO TO LINKED LISTS

*Discontiguous data structures are the root of all evil. This is simply a fact. If you don't believe me, I'll try to convince you. Specifically, please say "no" to linked lists. OK. Please. Please say no to linked lists. There is almost nothing more harmful you can do to the performance of an actual modern microprocessor than to use a linked list data structure.*

Chandler Carruth (Engineer at Google) @ CppCon 2014

# DISCONTIGUOUS DATA STRUCTURES ARE THE ROOT OF ALL (PERFORMANCE) EVIL

# Intel Core i7 cache hierarchy (2014)

Notice that the amount of space gets bigger as you go down the hierarchy

Processor package

Core 0

Core 3

16 (rax, rbx...)    Regs          Regs

32kB each    L1 d-cache | L1 i-cache    ...    L1 d-cache | L1 i-cache

256kB    L2 unified cache    L2 unified cache

2-20MB    L3 unified cache (shared by all cores)

L4* 128 MB (Iris Pro)

Main memory

# (Simplified) Cache View

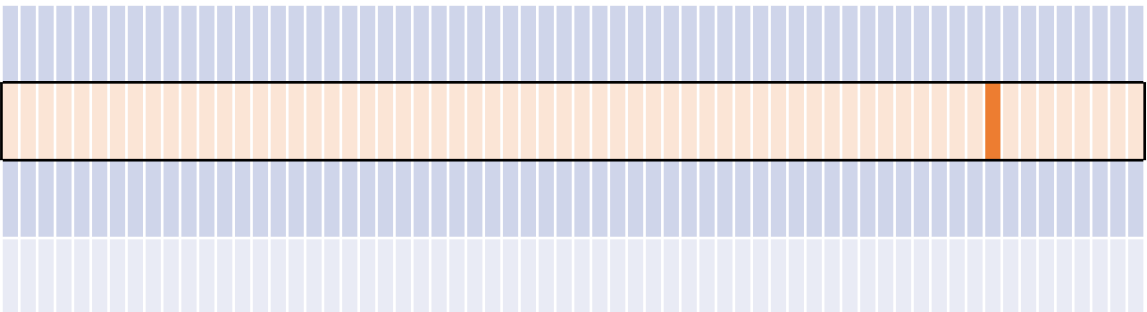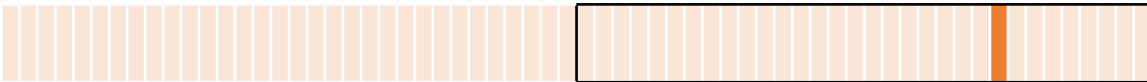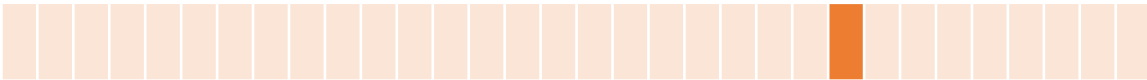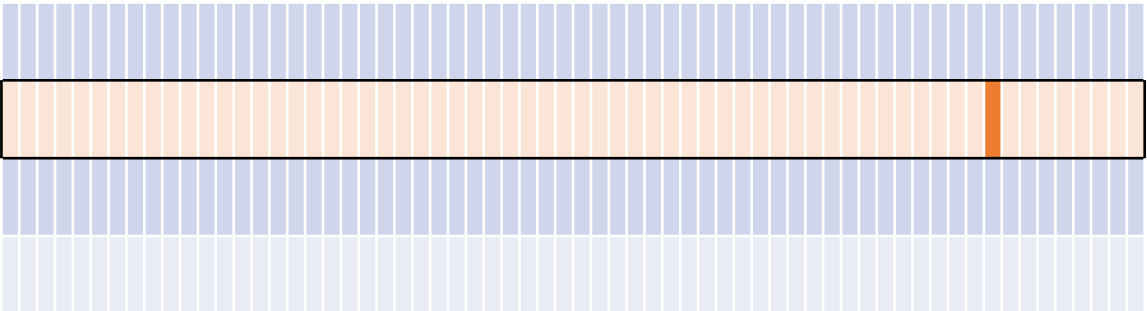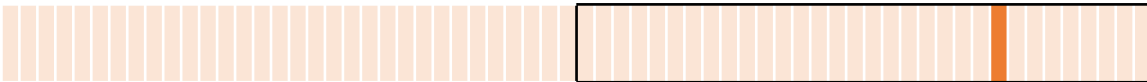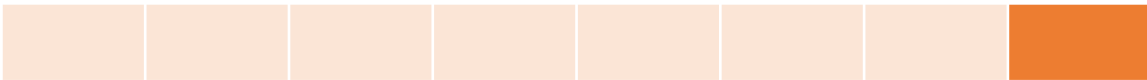Faster   Bigger

Register

L1 Cache

L2 Cache

L3 Cache

Memory

1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1 Cache

4. On miss, looks in L2 Cache

5. On miss, looks in L3 Cache

6. On miss, finds it in memory
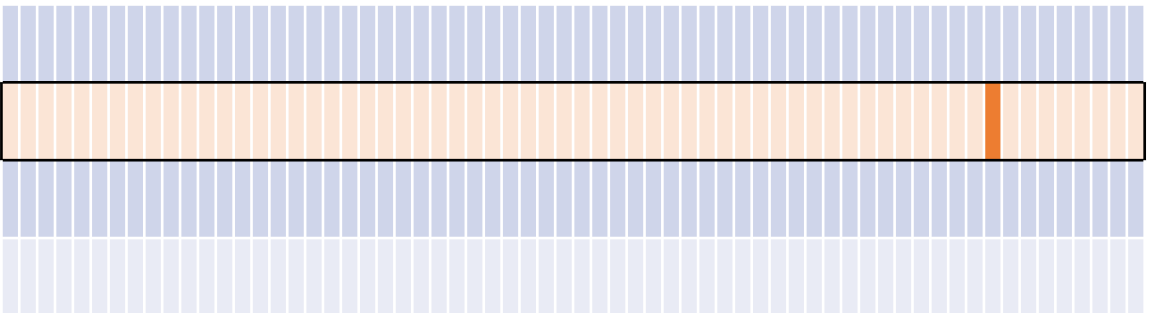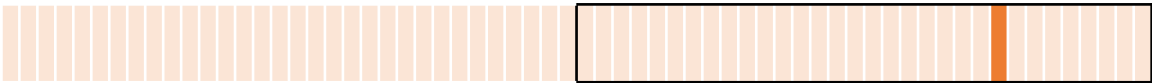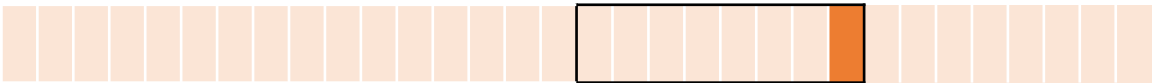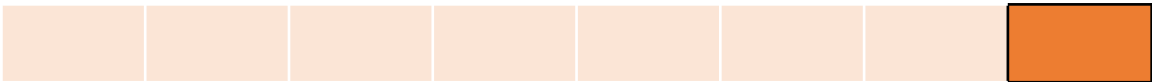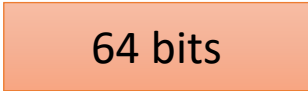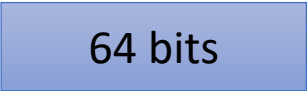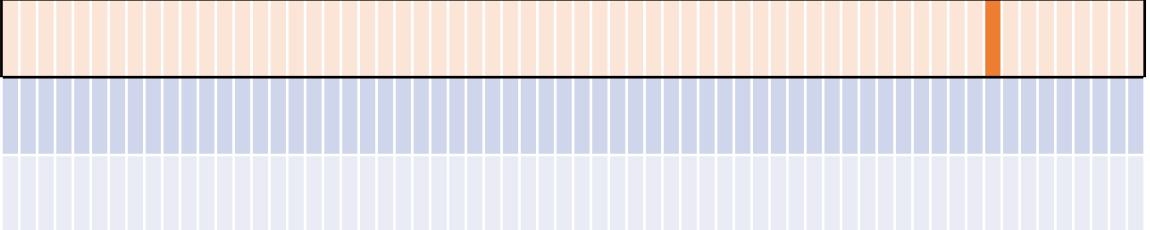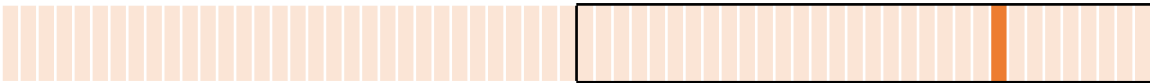
(in "virtual" memory)

| Entity | Time (nanoseconds) | Note |
|---|---:|---|
| One cycle on a 3 GHz processor | 1 | |
| L1 cache reference | 0.5 | |
| Branch mis-predict | 5 | |
| L2 cache reference | 7 | 14x L1 cache |
| Mutex lock/unlock | 25 | |
| Main memory reference | 100 | 20x L2; 200x L1 |
| Compress 1k bytes with Snappy | 3,000 | |
| Send 1K bytes over 1 Gbps network | 10,000 | |
| Read 4K randomly from SSD | 150,000 | |
| Read 1 MB sequentially from main memory | 250,000 | |
| Round trip with the same datacenter | 500,000 | |
| Read 1 MB sequentially from SSD | 1,000,000 | 4x main memory |
| Disk seek | 10,000,000 | 20x datacenter round trip |
| Read 1 MB sequentially from Disk | 20,000,000 | 80x main memory; 20x SSD |
| Send package CA -> Netherlands -> CA | 150,000,000 | |

Efficiency with Algorithms, Performance with Data Structures

# Memory Layout

| |
|:---:|
| **Text** |
| **Data** |
| **BSS** |
| **Heap** |
| |
| **Stack** |

# Intel Core i7 cache hierarchy (2014)

Processor package

Core 0

16 (rax, rbx…) — Regs

32kB each — L1 d-cache / L1 i-cache

256kB — L2 unified cache

Core 3

Regs

L1 d-cache / L1 i-cache

L2 unified cache

…

2-20MB — L3 unified cache (shared by all cores)

L4* 128 MB (Iris Pro)

Main memory

27

# Intel Haswell    Latency (cycles)



Bar chart showing latency in cycles: L1 ≈ 4, L2 ≈ 11, L3 ≈ 35, Memory ≈ 228.

CPU

64 bits          64 bits
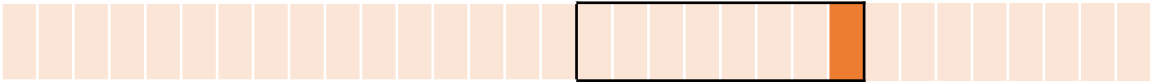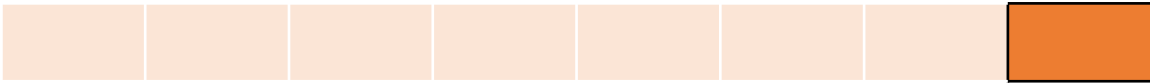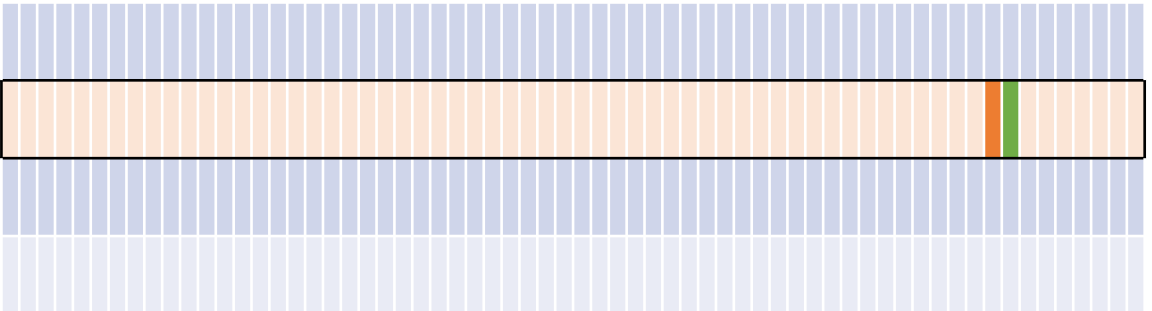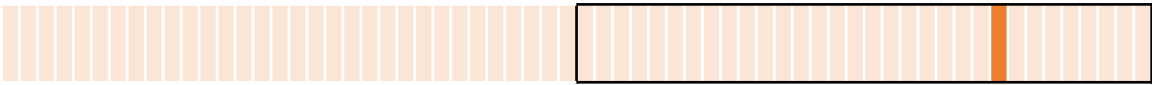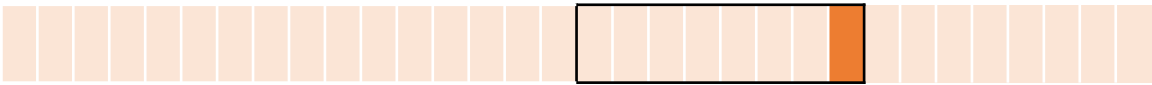
1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory
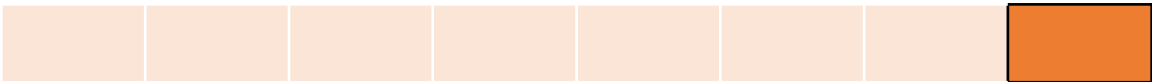
(in "virtual" memory)

29

CPU

64 bits    64 bits

1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory

(in "virtual" memory)
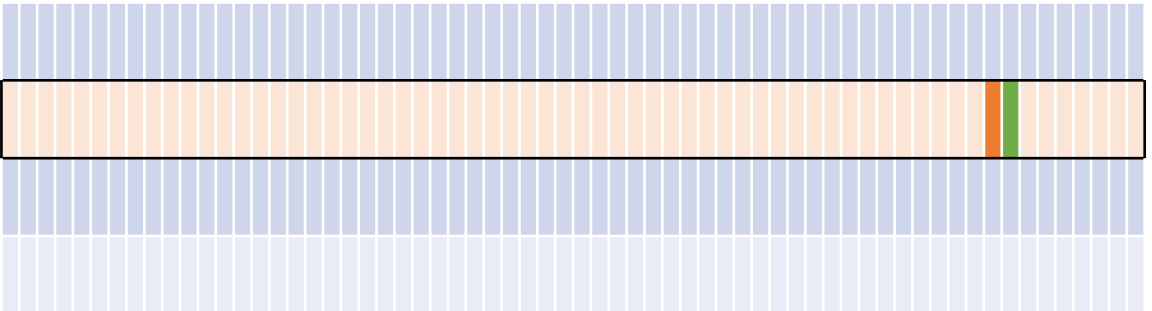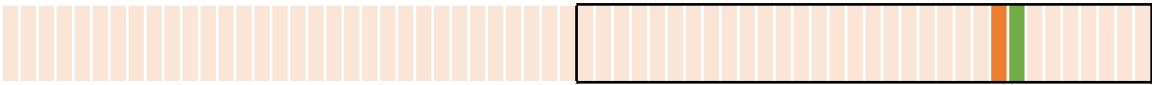
CPU

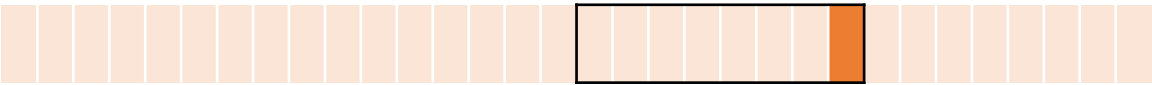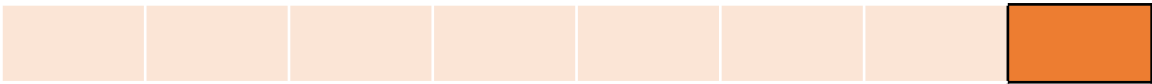| 64 bits | | 64 bits |

1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On <span style="color:red">miss</span>, looks in L2

5. On <span style="color:red">miss</span>, looks in L3

6. On <span style="color:red">miss</span>, <span style="color:green">finds</span> it in memory
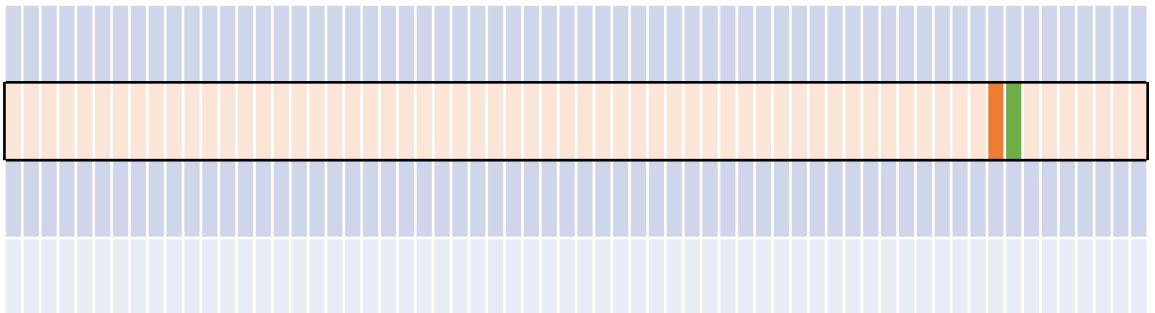
(in "virtual" memory)

1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory
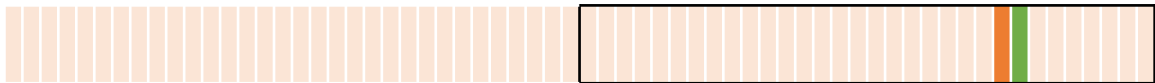
(in "virtual" memory)

CPU

64 bits     64 bits
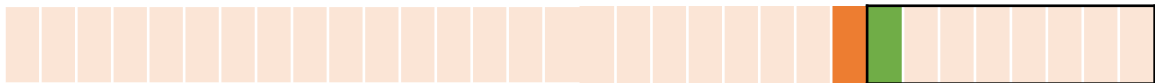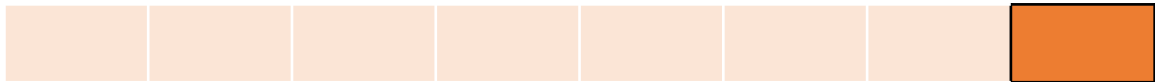
1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory

(in "virtual" memory)

CPU

64 bits          64 bits
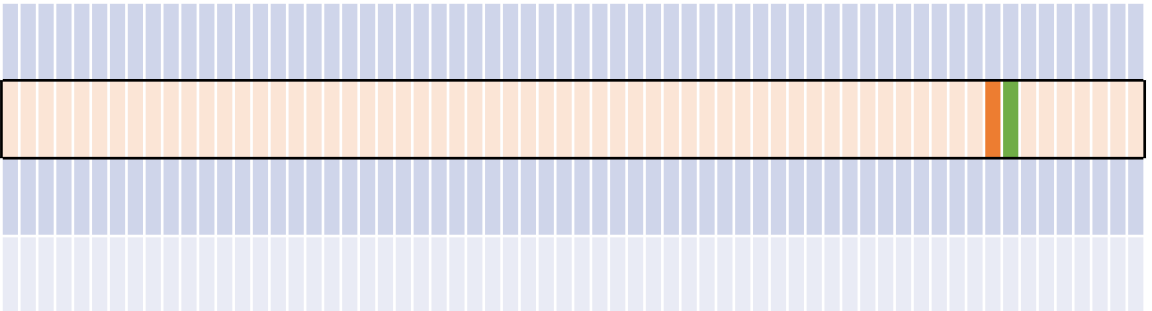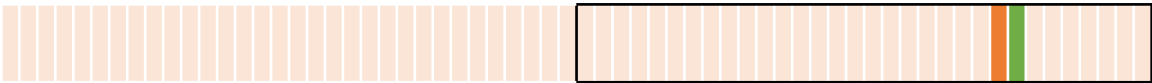
1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory

(in "virtual" memory)

34

CPU

64 bits          64 bits
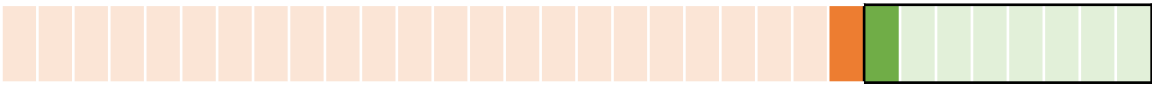
1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory

(in "virtual" memory)
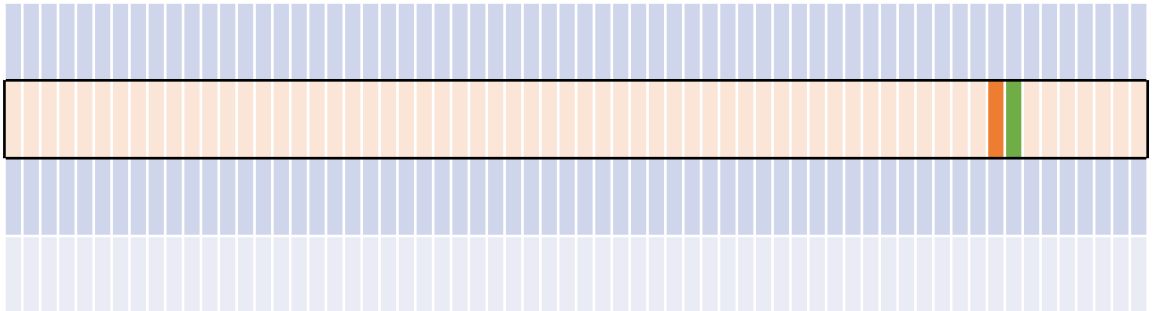
35

CPU

64 bits     64 bits

1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory

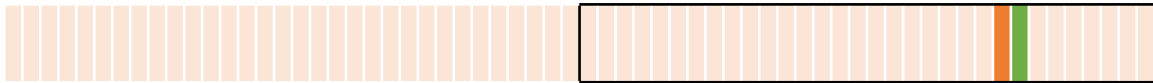(in "virtual" memory)
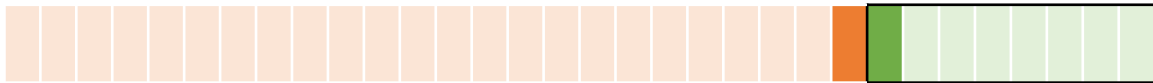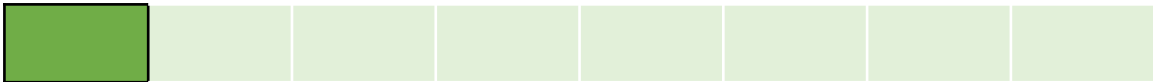
CPU

64 bits    64 bits

1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory

(in "virtual" memory)

37

CPU

64 bits          64 bits
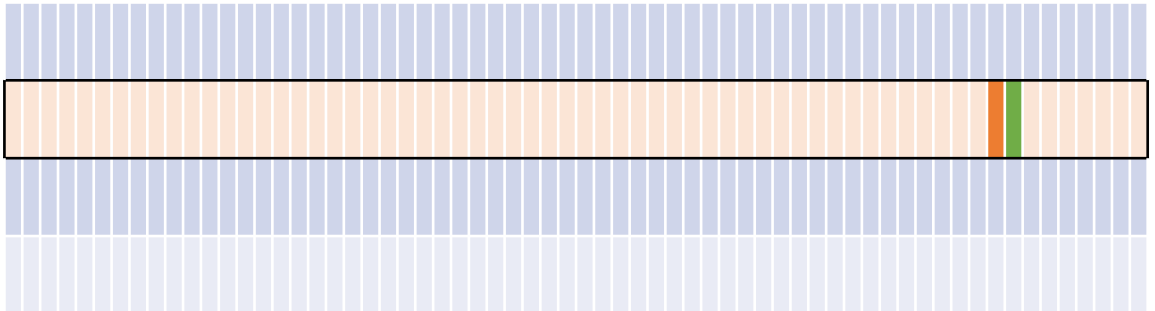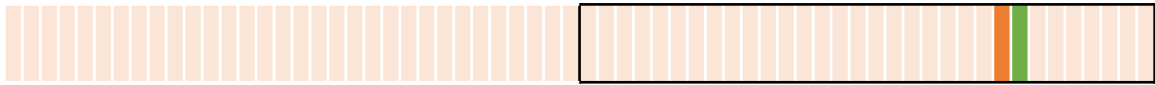
1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory
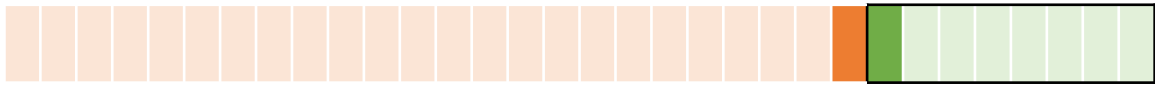
(in "virtual" memory)

38

CPU

64 bits    64 bits

1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory
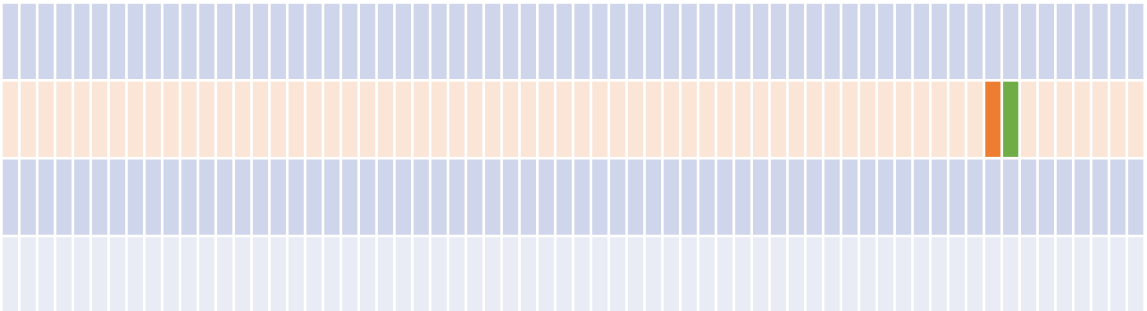
(in "virtual" memory)

39

CPU

64 bits    64 bits

1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory
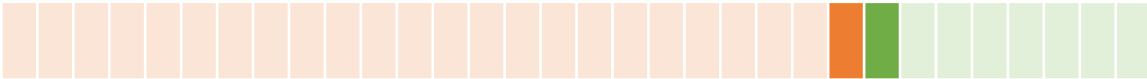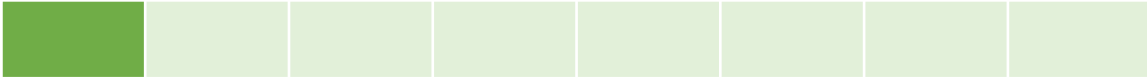
(in "virtual" memory)

40

CPU

64 bits    64 bits

1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory
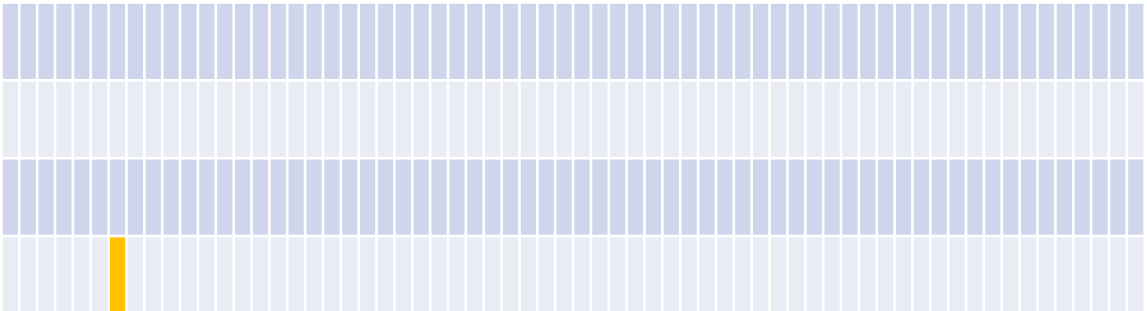
(in "virtual" memory)

41

CPU

64 bits    64 bits

1.  CPU needs a piece of memory

2.  Needs to load it into a register

3.  Looks in L1

4.  On miss, looks in L2

5.  On miss, looks in L3

6.  On miss, finds it in memory

(in "virtual" memory)

42

CPU

64 bits   64 bits

1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory
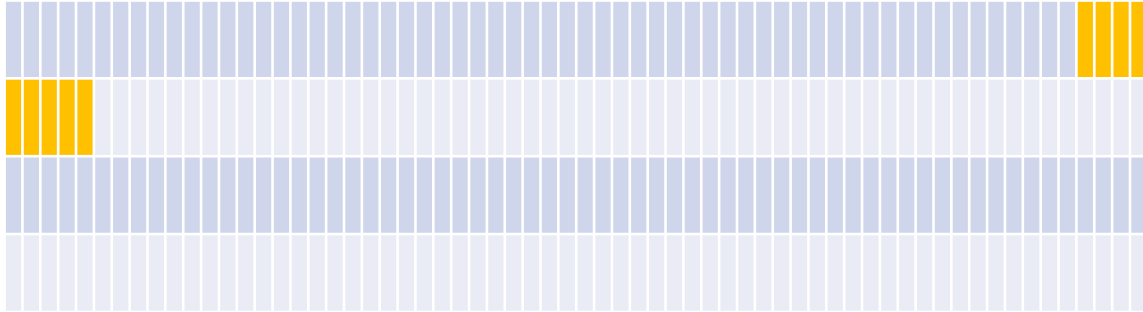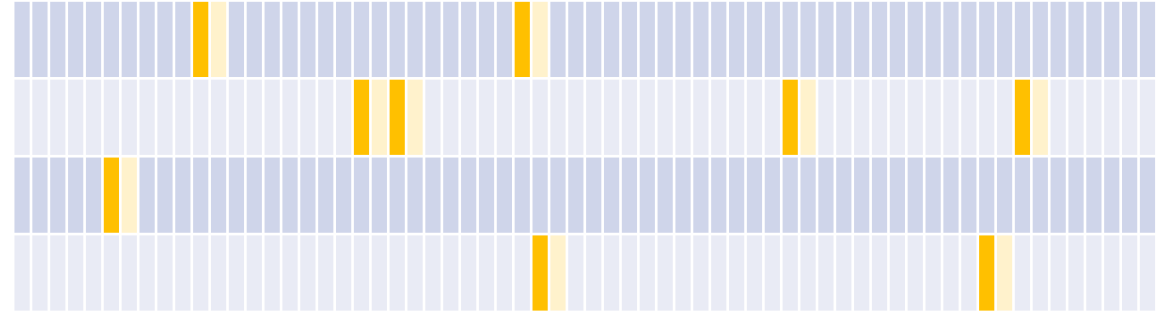
(in "virtual" memory)

43

CPU

64 bits          64 bits

1. CPU needs a piece of memory

2. Needs to load it into a register

3. Looks in L1

4. On miss, looks in L2

5. On miss, looks in L3

6. On miss, finds it in memory

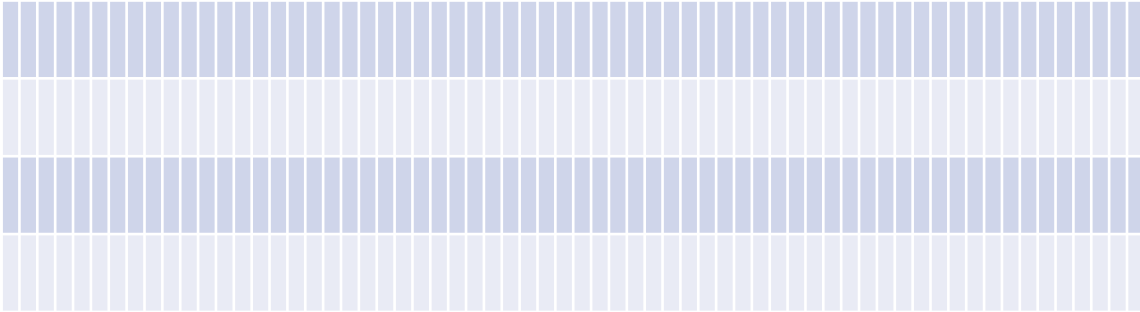(in "virtual" memory)

44

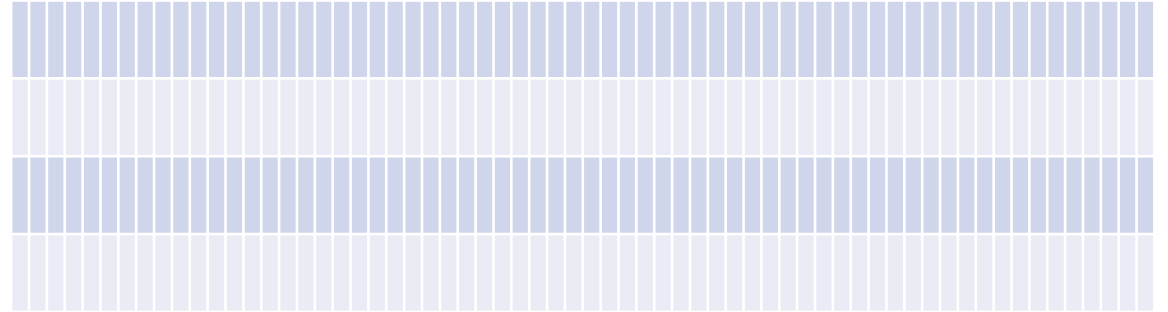**std::vector**                    **vs std::list**



Which has better locality?
Meaning, which will work better with cache?
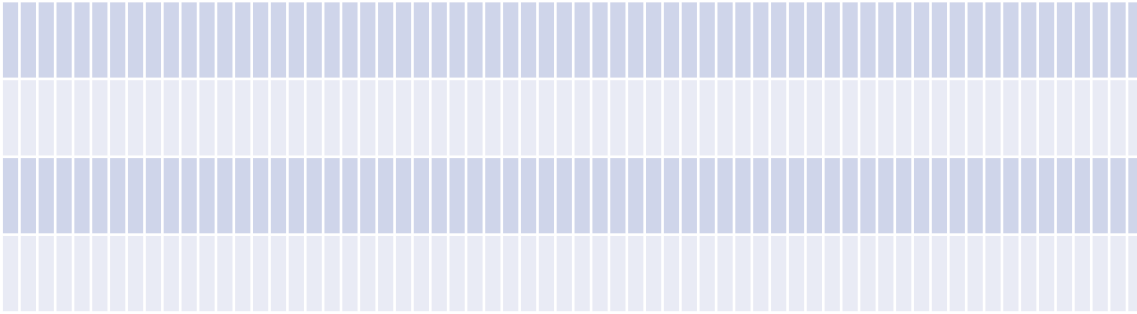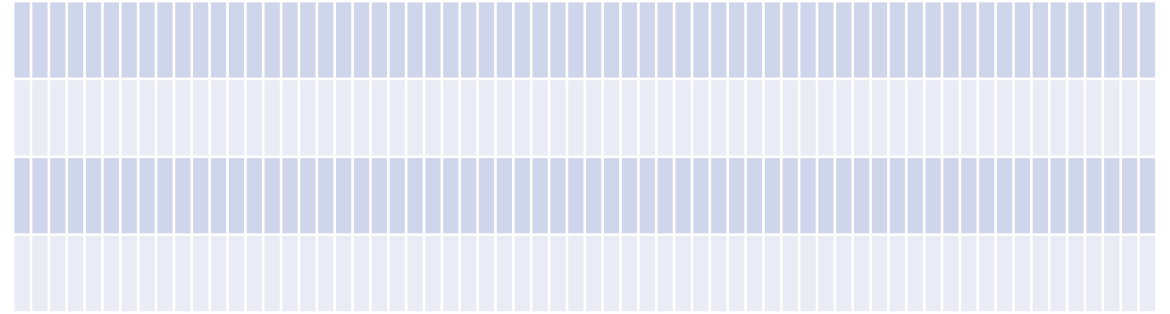
**Adjacency Matrix**

**Adjacency List**

What does this picture look like for:
1. An Adjacency Matrix
2. An Adjacency List

**Sorted Array**

**Search Tree**

**Heap**
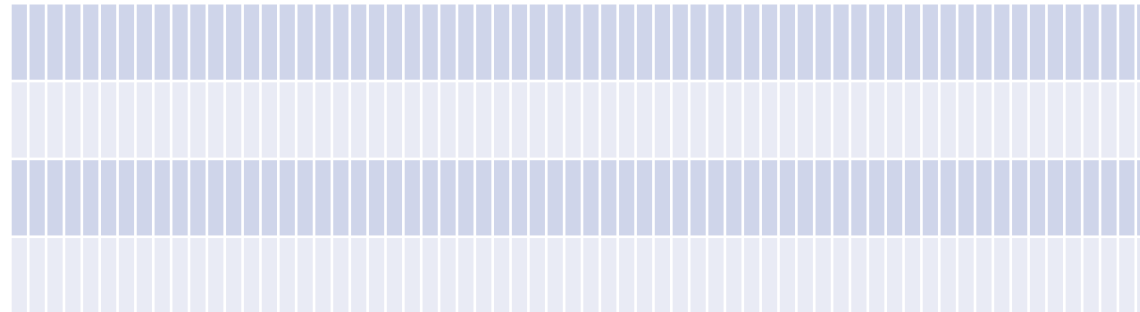
What does extract min look like?

What does this picture look like for:
1. A Sorted Array
2. A Binary Search Tree
3. A Heap

47