

# #20



`std::unordered_set`  
`std::list`  
`std::deque`  
`std::vector`  
`std::set`

investigating the efficiency of C++ containers

#cpp\_friends, 14::30

# Efficiency of standard containers

The standard library provides *containers* to store collections of same type.

The following containers are *dynamic*, meaning that they can grow or shrink as needed — their size is not fixed.

- `std::list`
- `std::vector`
- `std::deque`
- `std::set`
- `std::unordered_set`

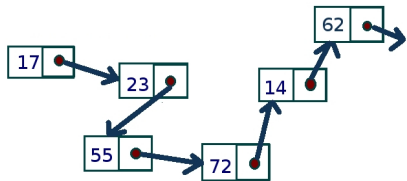
# Efficiency of standard containers

- `std::list`
- `std::vector`
- `std::deque`
- `std::set`
- `std::unordered_set`

Do you know how do they *store the data internally*?

Can you tell which one is *the most optimal* for a given task?

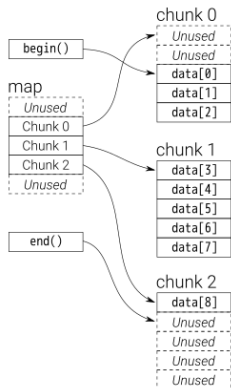
## Bit of theory — `std::list`



Operation	Complexity
<code>get(int index)</code>	$O(n)$
<code>push_back(T element)</code> <code>push_front(T element)</code>	$O(1)$
<code>insert(T element, int index)</code>	$O(1)$
<code>remove(int index)</code>	$O(1)$

image source: a detail of [https://commons.wikimedia.org/wiki/File:Linked\\_list\\_data\\_format.jpg](https://commons.wikimedia.org/wiki/File:Linked_list_data_format.jpg)

# Bit of theory — `std::deque`



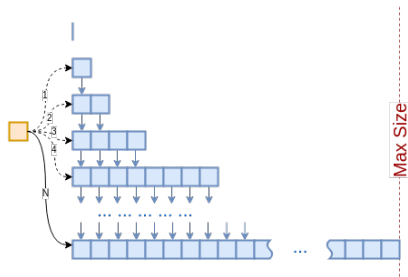
Operation	Complexity
<code>get(int index)</code>	$O(1)$
<code>push_back(T element)</code> <code>push_front(T element)</code>	$O(1)$
<code>insert(T element, int index)</code>	$O(n)$
<code>remove(int index)</code>	$O(n)$

$O(1)$  insert and remove at both ends of deque

Note: data is not stored contiguously.

image source: part of <https://stackoverflow.com/questions/6292332/what-really-is-a-deque-in-stl>

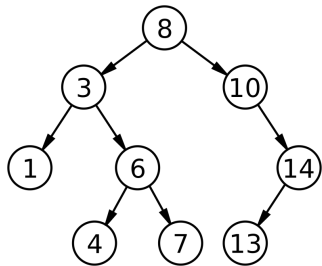
## Bit of theory — `std::vector`



Operation	Complexity
<code>get(int index)</code>	$O(1)$
<code>push_back(T element)</code>	$O(1)$ , may be worse
<code>push_front(T element)</code>	$O(n)$
<code>insert(T element, int index)</code>	$O(n)$
<code>remove_back</code>	$O(1)$
<code>remove(int index)</code>	$O(n)$

Note: data is **stored contiguously**.

image source: part of <https://stackoverflow.com/questions/52330010/what-does-stdvector-look-like-in-memory>



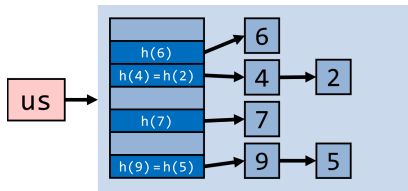
Operation	Complexity
<code>get(int index)</code>	<i>not provided</i>
<code>find(T element)</code>	$O(\log(n))$
<code>insert(T element)</code>	$O(\log(n))$
<code>remove(T element)</code>	$O(\log(n))$

Note [1]: sequential access is not supported.

Note [2]: uses a *binary search tree*.

image source: part of [https://en.wikipedia.org/wiki/Binary\\_search\\_tree](https://en.wikipedia.org/wiki/Binary_search_tree)

## Bit of theory — `std::unordered_set`



Operation	Complexity
<code>get(int index)</code>	<i>not provided</i>
<code>find(T element)</code>	$O(1)$ or $O(n)$
<code>insert(T element)</code>	$O(1)$ or $O(n)$
<code>remove(T element)</code>	$O(1)$ or $O(n)$

Note: all operation take average  $O(1)$  time.

It can go up to  $O(n)$  in worst case, but in general they perform well.

In more precise terms, `unordered_set`'s operations are **amortized  $O(1)$**

image source: [https://hackingcpp.com/cpp/std/unordered\\_set\\_thumb.svg](https://hackingcpp.com/cpp/std/unordered_set_thumb.svg)



# Benchmark — storage of data

```
template <typename T>
void test_storing() {
    std::cout << "Testing " << typeid(T).name() << " ... ";
    auto start {std::chrono::system_clock::now()};
    T container {};
    for (uint64_t i {0}; i < TEST_SIZE; ++i) {
        container.push_back(i);
    }
    auto finish {std::chrono::system_clock::now()};
    auto elapsed {std::chrono::duration_cast<std::chrono::milliseconds>(finish - start)};
    std::cout << "done, took " << elapsed.count() << " ms\n";
}
```

Trial	<b>vector</b>	<b>deque</b>	<b>list</b>	<b>set</b>	<b>unordered_set</b>
1.	70 ms	53 ms	309 ms	2745 ms	761 ms
2.	70 ms	54 ms	316 ms	2667 ms	730 ms
3.	71 ms	55 ms	332 ms	2700 ms	769 ms

Tested on an x86\_64 machine with Ryzen 5600X processor. Full code: [storing.cpp](#).

# Benchmark — accessing data randomly

```
for (std::size_t k {0}; k < NUM_RANDOM_READS; k++) {  
    read_value = *(container.cbegin() + INDEX_FIRST_QUARTER);  
    read_value = *(container.cbegin() + INDEX_SECOND_QUARTER);  
    read_value = *(container.cbegin() + INDEX_THIRD_QUARTER);  
    read_value = *(container.cbegin() + INDEX_FOURTH_QUARTER);  
}
```

Trial	<b>vector</b>	<b>deque</b>	<b>list</b>	<b>set</b>	<b>unordered_set</b>
1.	150 ms	304 ms	2449 ms	1743 ms	907 ms
2.	146 ms	295 ms	2437 ms	1754 ms	902 ms
3.	148 ms	308 ms	2446 ms	1729 ms	891 ms

Tested on an x86\_64 machine with Ryzen 5600X processor. Full code: [reading\\_random.cpp](#).

# Benchmark — accessing data sequentially

```
for (std::size_t k {0}; k < NUM_SEQUENTIAL_READS; k++) {  
    // begin s-th sequential read  
    for (std::size_t r {0}; r < CONTAINER_SIZE; r++) {  
        // read r-th element of the container (a deque or a vector)  
        read_value = *(container.cbegin() + r);  
    }  
}
```

Trial	<b>vector</b>	<b>deque</b>	<b>list</b>	<b>set</b>	<b>unordered_set</b>
1.	11 ms	22 ms	1357 ms	143 ms	45 ms
2.	9 ms	21 ms	1353 ms	145 ms	45 ms
3.	10 ms	21 ms	1360 ms	143 ms	45 ms

Tested on an x86\_64 machine with Ryzen 5600X processor. Full code: [reading\\_sequential.cpp](#).

# Benchmark — `std::vector` or `std::unordered_set`?

```
26 void test_vector() {
27     std::cout << "Testing std::vector ... " << std::flush;
28     auto start {std::chrono::system_clock::now()};
29     std::vector<tested_type> vec{};
30     for (std::size_t i {0}; i < CONTAINER_SIZE; ++i) {
31         vec.push_back(dist10G(rng));
32     }
33     std::sort(std::begin(vec), std::end(vec));
34     auto finish {std::chrono::system_clock::now()};
35     auto elapsed {std::chrono::duration_cast<timing_granularity>(finish - start)};
36     std::cout << "done, took " << elapsed.count() << " units\n";
37 }

39 void test_unordered_set() {
40     std::cout << "Testing std::unordered_set ... " << std::flush;
41     auto start {std::chrono::system_clock::now()};
42     std::unordered_set<tested_type> myset{};
43     for (std::size_t i {0}; i < CONTAINER_SIZE; ++i) {
44         myset.insert(dist10G(rng));
45     }
46     auto finish {std::chrono::system_clock::now()};
47     auto elapsed {std::chrono::duration_cast<timing_granularity>(finish - start)};
48     std::cout << "done, took " << elapsed.count() << " units\n";
49 }
50 }
```

Container size	<code>vector</code>	<code>unordered_set</code>
10	15 $\mu s$	7 $\mu s$
100	23 $\mu s$	22 $\mu s$
1000	251 $\mu s$	241 $\mu s$
10K	3834 $\mu s$	2659 $\mu s$
1M	374 ms	426 ms
10M	4152 ms	6109 ms

Tested on an x86\_64 machine with Ryzen 5600X processor. Full code: `vector_or_hashset.cpp`.

# Why is that?

modern CPUs: every possible trick to boost performance

- accessing contiguous memory is *significantly* faster than accessing random addresses
- `std::list`: usually very fragmented (just a few insert/delete operations suffice)
- `std::deque`: fragmented chunks of contiguous memory
- `std::vector`: average theoretical performance, great practical performance
  - — especially for operations on the entire sequence
- `std::vector`: the only container that stores data *sequentially*

# Key takeaways

- know the pros and cons of various standard containers
- know your machine too — caching, paging, speculative execution  
*can provide a significant performance gain — or performance penalty*
- linked list is still useful — e.g. great for caching data
- `std::unordered_set`: almost a silver bullet
- when in doubt — measure, then fix, then measure again
- `std::vector` is almost always good enough

Thank you!