20.2 — The stack and the heap

The memory that a program uses is typically divided into a few different areas, called segments:

- The code segment (also called a text segment), where the compiled program sits in memory. The code segment is typically read-only.
- The bss segment (also called the uninitialized data segment), where zero-initialized global and static variables are stored.
- The data segment (also called the initialized data segment), where initialized global and static variables are stored.
- The heap, where dynamically allocated variables are allocated from.
- The call stack, where function parameters, local variables, and other function-related information are stored.

For this lesson, we'll focus primarily on the heap and the stack, as that is where most of the interesting stuff takes place.

The heap segment

The heap segment (also known as the "free store") keeps track of memory used for dynamic memory allocation. We talked about the heap a bit already in lesson 19.1 -- Dynamic memory allocation with new and delete (https://www.learncpp.com/cpp-tutorial/dynamic-memory-allocation-with-new-and-delete/)², so this will be a recap.

In C++, when you use the new operator to allocate memory, this memory is allocated in the application's heap segment.

Assuming an int is 4 bytes:

```
1 | int* ptr { new int }; // new int allocates 4 bytes in the heap
2 | int* array { new int[10] }; // new int[10] allocates 40 bytes in the heap
```

The address of this memory is passed back by operator new, and can then be stored in a pointer. You do not have to worry about the mechanics behind the process of how free memory is located and allocated to the user. However, it is worth knowing that sequential memory requests may not result in sequential memory addresses being allocated!

```
1  int* ptr1 { new int };
2  int* ptr2 { new int };
3  // ptr1 and ptr2 may not have sequential addresses
```

When a dynamically allocated variable is deleted, the memory is "returned" to the heap and can then be reassigned as future allocation requests are received. Remember that deleting a pointer does not delete the variable, it just returns the memory at the associated address back to the operating system.

The heap has advantages and disadvantages:

- Allocating memory on the heap is comparatively slow.
- Allocated memory stays allocated until it is specifically deallocated (beware memory leaks) or the application ends (at which point the OS should clean it up).
- Dynamically allocated memory must be accessed through a pointer. Dereferencing a pointer is slower than accessing a variable directly.
- Because the heap is a big pool of memory, large arrays, structures, or classes can be allocated here.

The call stack

The **call stack** (usually referred to as "the stack") has a much more interesting role to play. The call stack keeps track of all the active functions (those that have been called but have not yet terminated) from the start of the program to the current point of execution, and handles allocation of all function parameters and local variables.

The call stack is implemented as a stack data structure. So before we can talk about how the call stack works, we need to understand what a stack data structure is.

The stack data structure

A **data structure** is a programming mechanism for organizing data so that it can be used efficiently. You've already seen several types of data structures, such as arrays and structs. Both of these data structures provide mechanisms for storing data and accessing that data in an efficient way. There are many additional data structures that are commonly used in programming, quite a few of which are implemented in the standard library, and a stack is one of those.

Consider a stack of plates in a cafeteria. Because each plate is heavy and they are stacked, you can really only do one of three things:

- 1. Look at the surface of the top plate
- 2. Take the top plate off the stack (exposing the one underneath, if it exists)
- 3. Put a new plate on top of the stack (hiding the one underneath, if it exists)

In computer programming, a stack is a container data structure that holds multiple variables (much like an array). However, whereas an array lets you access and modify elements in any order you wish (called **random access**), a stack is more limited. The operations that can be performed on a stack correspond to the three things mentioned above:

- 1. Look at the top item on the stack (usually done via a function called top(), but sometimes called peek())
- 2. Take the top item off of the stack (done via a function called pop())
- 3. Put a new item on top of the stack (done via a function called push())

A stack is a last-in, first-out (LIFO) structure. The last item pushed onto the stack will be the first item popped off. If you put a new plate on top of the stack, the first plate removed from the stack will be the plate you just pushed on last. Last on, first off. As items are pushed onto a stack, the stack grows larger -- as items are popped off, the stack grows smaller.

For example, here's a short sequence showing how pushing and popping on a stack works:

Stack: empty
Push 1
Stack: 1
Push 2
Stack: 1 2
Push 3
Stack: 1 2 3
Pop
Stack: 1 2
Pop

Stack: 1

The plate analogy is a pretty good analogy as to how the call stack works, but we can make a better analogy. Consider a bunch of mailboxes, all stacked on top of each other. Each mailbox can only hold one item, and all mailboxes start out empty. Furthermore, each mailbox is nailed to the mailbox below it, so the number of mailboxes can not be changed. If we can't change the number of mailboxes, how do we get a stack-like behavior?

First, we use a marker (like a post-it note) to keep track of where the bottom-most empty mailbox is. In the beginning, this will be the lowest mailbox (on the bottom of the stack). When we push an item onto our mailbox stack, we put it in the mailbox that is marked (which is the first empty mailbox), and move the marker up one mailbox. When we pop an item off the stack, we move the marker down one mailbox (so it's pointed at the top non-empty mailbox) and remove the item from that mailbox. Anything below the marker is considered "on the stack". Anything at the marker or above the marker is not on the stack.

The call stack segment

The call stack segment holds the memory used for the call stack. When the application starts, the main() function is pushed on the call stack by the operating system. Then the program begins executing.

When a function call is encountered, the function is pushed onto the call stack. When the current function ends, that function is popped off the call stack (this process is sometimes called **unwinding the stack**). Thus, by looking at the functions that are currently on the call stack, we can see all of the functions that were called to get to the current point of execution.

Our mailbox analogy above is fairly analogous to how the call stack works. The stack itself is a fixed-size chunk of memory addresses. The mailboxes are memory addresses, and the "items" we're pushing and popping on the stack are called **stack frames**. A stack frame keeps track of all of the data associated with one function call. We'll talk more about stack frames in a bit. The "marker" is a register (a small piece of memory in the CPU) known as the stack pointer (sometimes abbreviated "SP"). The stack pointer keeps track of where the top of the call stack currently is.

We can make one further optimization: When we pop an item off the call stack, we only have to move the stack pointer down -- we don't have to clean up or zero the memory used by the popped stack frame (the equivalent of emptying the mailbox). This memory is no longer considered to be "on the stack" (the stack pointer will be at or below this address), so it won't be accessed. If we later push a new stack frame to this same memory, it will overwrite the old value we never cleaned up.

The call stack in action

Let's examine in more detail how the call stack works. Here is the sequence of steps that takes place when a function is called:

- 1. The program encounters a function call.
- 2. A stack frame is constructed and pushed on the stack. The stack frame consists of:
- The address of the instruction beyond the function call (called the **return address**). This is how the CPU remembers where to return to after the called function exits.
- All function arguments.
- Memory for any local variables
- Saved copies of any registers modified by the function that need to be restored when the function returns
- 3. The CPU jumps to the function's start point.
- 4. The instructions inside of the function begin executing.

When the function terminates, the following steps happen:

- 1. Registers are restored from the call stack
- 2. The stack frame is popped off the stack. This frees the memory for all local variables and arguments.
- 3. The return value is handled.
- 4. The CPU resumes execution at the return address.

Return values can be handled in a number of different ways, depending on the computer's architecture. Some architectures include the return value as part of the stack frame. Others use CPU registers.

Typically, it is not important to know all the details about how the call stack works. However, understanding that functions are effectively pushed on the stack when they are called and popped off (unwound) when they return gives you the fundamentals needed to understand recursion, as well as some other concepts that are useful when debugging.

A technical note: on some architectures, the call stack grows away from memory address 0. On others, it grows towards memory address 0. As a consequence, newly pushed stack frames may have a higher or a lower memory address than the previous ones.

A quick and dirty call stack example

Consider the following simple application:

```
int foo(int x)
1
3
         // b
         return x;
5 | } // foo is popped off the call stack here
7
     int main()
 8
 9
         // a
         foo(5); // foo is pushed on the call stack here
 10
11
 12
 13
         return 0;
     }
 14
```

The call stack looks like the following at the labeled points:

```
main()
```

b:

```
foo() (including parameter x)
main()
```

c:

```
main()
```

Stack overflow

The stack has a limited size, and consequently can only hold a limited amount of information. On Visual Studio for Windows, the default stack size is 1MB. With g++/Clang for Unix variants, it can be as large as 8MB. If the program tries to put too much information on the stack, stack overflow will result. **Stack overflow** happens when all the memory in the stack has been allocated -- in that case, further allocations begin overflowing into other sections of memory.

Stack overflow is generally the result of allocating too many variables on the stack, and/or making too many nested function calls (where function A calls function B calls function C calls function D etc...) On modern operating systems, overflowing the stack will generally cause your OS to issue an access violation and terminate the program.

Here is an example program that will likely cause a stack overflow. You can run it on your system and watch it crash:

```
#include <iostream>

int main()

int stack[10000000];

std::cout << "hi" << stack[0]; // we'll use stack[0] here so the compiler won't

optimize the array away

return 0;
}</pre>
```

This program tries to allocate a huge (likely 40MB) array on the stack. Because the stack is not large enough to handle this array, the array allocation overflows into portions of memory the program is not allowed to use.

On Windows (Visual Studio), this program produces the result:

```
HelloWorld.exe (process 15916) exited with code -1073741571.
```

-1073741571 is c0000005 in hex, which is the Windows OS code for an access violation. Note that "hi" is never printed because the program is terminated prior to that point.

Here's another program that will cause a stack overflow for a different reason:

```
1 | // h/t to reader yellowEmu for the idea of adding a counter
     #include <iostream>
     int g_counter{ 0 };
 4
 5
 6
     void eatStack()
7
 8
         std::cout << ++g_counter << ' ';</pre>
9
10
         // We use a conditional here to avoid compiler warnings about infinite recursion
11
         if (g_counter > 0)
             eatStack(); // note that eatStack() calls itself
12
13
         // Needed to prevent compiler from doing tail-call optimization
14
15
         std::cout << "hi";
     }
16
17
18
     int main()
19
         eatStack();
20
21
22
         return 0;
23
```

In the above program, a stack frame is pushed on the stack every time function eatStack() is called. Since eatStack() calls itself (and never returns to the caller), eventually the stack will run out of memory and cause an overflow.

Author's note

When run on the author's Windows 10 machine (from within the Visual Studio Community IDE), eatStack() crashed after 4848 calls in debug mode, and 128,679 calls in release mode.

Related content

We talk more about functions that call themselves in upcoming lesson <u>20.3 -- Recursion</u> (https://www.learncpp.com/cpp-tutorial/recursion/)³.

The stack has advantages and disadvantages:

- Allocating memory on the stack is comparatively fast.
- Memory allocated on the stack stays in scope as long as it is on the stack. It is destroyed when it is popped off the stack.
- All memory allocated on the stack is known at compile time. Consequently, this memory can be accessed directly through a variable.
- Because the stack is relatively small, it is generally not a good idea to do anything that eats up lots of stack space. This includes allocating or copying large arrays or other memory-intensive structures.

Author's note

<u>This comment (https://www.learncpp.com/cpp-tutorial/introduction-to-objects-and-variables/#comment-560618)</u> has some additional (simplified) information about how variables on the stack are laid out and receive actual memory addresses at runtime.



3



Back to table of contents

5



<u>Previous lesson</u>

20.1 <u>Function Pointers</u>

6

7





Notify me about replies:

POST COMMENT



★ Find a mistake? Leave a comment above!

Avatars from https://gravatar.com/⁸ are connected to your provided email address.

340 COMMENTS Newest ▼



KLAP

(1) July 7, 2025 5:10 am PDT

When a dynamically allocated variable is deleted, the memory is "returned" to the heap and can then be reassigned as future allocation requests are received. Remember that deleting a pointer does not delete the variable, it just returns the memory at the associated address back to the operating system. The second phrase confused me. I thought when you delete a dynamically allocated object the pointer remain and variable is deleted. That's why you can't de-reference the pointer after using delete.







PooperShooter

① May 18, 2025 6:22 pm PDT

"However, it is worth knowing that sequential memory requests may not result in sequential memory addresses being allocated!"

Is this due to segmentation? So like if the free blocks of mem in the free list aren't compacted and they are spread out then the OS may not be able to meet the request with just 1 chunk of memory and it might have to allocate a portion in one block and the other portion in another block of free mem? If I am not understanding correctly please help fix what I got wrong, would be much appreciated.



Reply



Kania

① April 5, 2025 8:14 am PDT

Lmaao my device so bad it crashed after 200 calls







March 28, 2025 1:01 am PDT

Hello Alex.

You've described that arrays let us access elements in any order, but call stack is limited and let us access the elements by pop instruction. On the other hand, you've mentioned call stack stores local variables. If the mechanics of arrays differ from the call stack, then how a local array variable is stored in stack?









(1) March 6, 2025 1:10 pm PST

Hi Alex!

"-1073741571 is c0000005 in hex". I'm converting the number and it's giving me c000000fd in hex, i did a little research and skimming through this: https://learn.microsoft.com/en-

us/openspecs/windows_protocols/ms-erref/596a1078-e883-4972-9bbc-49e60bebca55? redirectedfrom=MSDN i found the "STATUS_STACK_OVERFLOW" error which is effectively mapped as 0xC00000FD in hex.

Is this correct though? Anyways thanks for all the hardwork, keep it up!







October 13, 2024 12:27 am PDT

Hi Alex.

How do the constexpr functions execute if they do execute during compile time. My guess is that since compiler is a program itself, it must have memory allocated for call stack and hence the constexpr function will be called on that stack.

Is this right or am i completely wrong on this.



Reply



Alex Author

Reply to ssd ⁹ O October 16, 2024 2:17 pm PDT

How a compiler internally implements the execution of constexpr functions is an implementationspecific detail. Abstractly, the compiler basically acts as an interpreter. Since constexpr functions can call other functions, I imagine the compiler mimics a call stack (probably in heap-allocated memory) for purpose of evaluation.







Joyboy

① October 1, 2024 11:52 am PDT

If inline functions just expand inline wherever they are called , then does this mean they dont enter the stack? if so is this one of the reasons why many small functions are made instead a big one so as to make the compiler expand them inline so that they dont saturate the call stack?









Alex Author

Reply to Joyboy 10 © October 1, 2024 3:46 pm PDT

Inline functions aren't always expanded inline, but when they are this avoids having to do stack operations, which is where most of the overhead of making a function call comes from.

Small functions are preferred because they tend to be more modular and easier to understand. The fact that they are more likely to be inlined is secondary benefit, as function calls are typically fast on modern architectures.



Reply



ssd

Q Reply to Alex ¹¹ **(** October 13, 2024 12:17 am PDT

But at the end, if inline functions do expand inline, the expansion will be happen on the function they are called meaning that the local variables and function parameters(related to inline) will be still be allocated to stack frame in which they are called in.

For example, consider two functions function1(inline) and function2(calling the function1).so when the stack frame for function 2 is allocated the local variables and function parameters related to function1 are allocated in fucntion2.

Meaning although local variables and function parameters are allocated in call stack we have avoided a whole another function getting pushed and popped in the call stack which is very

advantageous in terms of performance.

But i guess this expansion is done during build(compile and link) so the built time might increase if we are using too much inline(because there must be some cost related to expansion) but we do gain some performance gain in runtime.

I hope this is what expansion of inline looks like in stack.

↑ Reply

Joyboy

Reply to Alex 11 ① October 2, 2024 4:42 am PDT

i see . ty for for the lesson and the response.

1 0 → Reply



Filipe

① July 18, 2024 9:37 pm PDT

Hi, there's a part that nobody explains if the heap has unlimited memory than it must mean it's allocated outside the "heap region" that is shown to belong to the program memory layout.

So if this is true than why people keep showing that same heap memory region beloging to the program?





Alex Author

In modern OSes, all processes are given their own virtual address space. From the process's point of view, these appear to be real addresses. However, the virtual address space for the process is managed by the memory management unit of the OS, and that includes mapping the virtual address of a process to address of physical memory (or secondary storage).

So when you request heap memory, the address that the process/C++ sees appears to be part of the heap for the process, but in reality it's a virtual address that's mapped to somewhere else.

2 → Reply



Filipe

Thank you very much for yor explanation.

■ 0 Reply



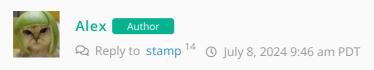
stamp

(Line of the control of the control

```
1 // release mode
     void func(){
3
         int i1{ 1 };
 4
          int i2{ 2 };
 5
         std::cout << &i1 << '\n'; // 00000065D86FFCE0
  6
         std::cout << &i2 << '\n'; // 00000065D86FFCF0
7
     }
 8
9
     int main(int argc, char* argv[]){
 10
          func();
11
         int i1{ 3 };
 12
          int i2{ 4 };
         std::cout << &i1 << '\n'; // 00000065D86FFCF8
 13
 14
          std::cout << &i2 << '\n'; // 00000065D86FFCB0
15 | }
```

This code running in release mode on my machine all addresses are not sequential, but in debug mode they are. Course mentions that it will happens to the heap, is it same on the stack memory?

1 0 → Reply



i1 and i2 are allocated on the stack. There is no requirement that variables are allocated consecutively on the stack. Perhaps your implementation is doing something different for some reason.

↑ Reply



Swaminathan

① June 22, 2024 2:51 am PDT

```
1 | int* ptr { new int }; // ptr is assigned 4 bytes in the heap
2 | int* array { new int[10] }; // array is assigned 40 bytes in the heap
```

I thought the pointer itself will be assigned in the stack and only the memory for the object being pointed to will be in heap..is my assumption wrong that in your examples, 4 bytes will be assigned in the stack one each for ptr and array, and they will point to the heap where they point to 4 and 40 bytes respectively??

1 0 → Reply



Alex Author

You are correct. I updated the wording in the comment to make it less ambiguous.

0 Reply

Links

- 1. https://www.learncpp.com/author/Alex/
- 2. https://www.learncpp.com/cpp-tutorial/dynamic-memory-allocation-with-new-and-delete/
- 3. https://www.learncpp.com/cpp-tutorial/recursion/
- 4. https://www.learncpp.com/cpp-tutorial/introduction-to-objects-and-variables/#comment-560618
- 5. https://www.learncpp.com/
- 6. https://www.learncpp.com/cpp-tutorial/function-pointers/
- 7. https://www.learncpp.com/the-stack-and-the-heap/
- 8. https://gravatar.com/
- 9. https://www.learncpp.com/cpp-tutorial/the-stack-and-the-heap/#comment-603021
- 10. https://www.learncpp.com/cpp-tutorial/the-stack-and-the-heap/#comment-602580
- 11. https://www.learncpp.com/cpp-tutorial/the-stack-and-the-heap/#comment-602603
- 12. https://www.learncpp.com/cpp-tutorial/the-stack-and-the-heap/#comment-599823
- 13. https://www.learncpp.com/cpp-tutorial/the-stack-and-the-heap/#comment-599897
- 14. https://www.learncpp.com/cpp-tutorial/the-stack-and-the-heap/#comment-599208
- 15. https://www.learncpp.com/cpp-tutorial/the-stack-and-the-heap/#comment-598693
- 16. https://g.ezoic.net/privacy/learncpp.com