

Software Systems

Lecture 2

Jonathan Dönszelmann

Vivian Roest

Delft University of Technology, The Netherlands

January 24, 2024

Previous Lectures

- Threads
- Spawning and Joining
- Channels
- Threadpools
- Asynchronous programming

Today

How to find and improve the performance of (existing) programs.

Today

How to find and improve the performance of (existing) programs.

Question:

Do you know any techniques already?

Today

How to find and improve the performance of (existing) programs.

- ① Running code less
- ② Make code that has to run faster

Today

How to find and improve the performance of (existing) programs.

- ① Running code less
- ② Make code that has to run faster

But first! How do we measure performance?

Measuring performance

- ① Benchmarking
- ② Profiling

Measuring performance: Benchmarking

```
time ./some/binary  
time cargo run
```

Warning:

Timing `cargo run` will also measure compile time. Make sure you compile first.

- 1 Very crude way measure performance
- 2 Not very accurate
- 3 Very useful to get an estimate!

Measuring performance: Benchmarking

Question:

Why is this not very accurate?

Measuring performance: Benchmarking

Remove influences of things we don't want to measure

Average results to remove random fluctuations

Measuring performance: Benchmarking

- Rust used to have this built-in
- Now Rust only has the framework¹, which external libraries can take advantage of. e.g. `criterion`
- `criterion` will:
 - Take samples out of many runs
 - Warm up
 - Show difference with old runs
 - Detect outliers
 - Give statistical significance (p-value)

Demo!

¹Technically it is still all there, but it's recommended to use an external library for it

Measuring performance: Profiling

Question:

Did you use profiling in Advanced Computing Systems?

Measuring performance: Profiling

- Measure execution times of different parts of an application
- On Linux: Perf
 - Hardware counters (`perf stat`)
 - Sampling
- On Windows: use WSL.
- On MacOS: Use DTrace

<https://perf.wiki.kernel.org/index.php/Tutorial>

Measuring performance: Profiling

- Measure execution times of different parts of an application
- On Linux: Perf
 - Hardware counters (`perf stat`)
 - Sampling
- On Windows: use WSL. On MacOS: Use DTrace

Demo!

Measuring performance: Profiling

Other useful tools:

- Valgrind
- Cachegrind
- Callgrind
- Dumb but useful: htop

Other kinds of performance

- for example: memory usage
- always a tradeoff

Question:

Can you think of other measures of performance?

Increasing performance

Split into three categories:

- Algorithmic improvement (being smart)
- Making code faster
- Running less code

Memoization

```
1  fn fibonacci(n: u64) -> u64 {
2      match n {
3          0 => 1,
4          1 => 1,
5          n => fibonacci(n-1) + fibonacci(n-2),
6      }
7  }
```

- fibonacci(5) needs fibonacci(4) and fibonacci(3)
- fibonacci(4) needs fibonacci(3) and fibonacci(2)

If we added some memory we could “cache” the results of fibonacci

Memoization

```
1 use memoize::memoize;
2
3 #[memoize]
4 fn fibonacci(n: u64) -> u64 {
5     match n {
6         0 => 1,
7         1 => 1,
8         n => fibonacci(n-1) + fibonacci(n-2),
9     }
10 }
```

- Adds a mapping from parameters to results
- “caches” the results of fibonacci to reduce calls

Lazy evaluation

- Iterators don't compute values immediately
- Only when `nth` is called do we run code

```
1  vec![1, 2, 3, 4].into_iter().map(|i: i32| i.pow(2)).nth(2)
```

- Slower in general
- Sometimes faster due to a smart compiler
- May be faster when few items end up needing processing

IO Buffering

- Every time `read_exact` is used, a system call is performed
- Goal: perform fewer system calls
- read as much as we can per system call

```
1 use std::fs::File;
2 use std::io::Read;
3
4 fn main() {
5     let mut file = File::open("very_large_file.txt").unwrap();
6     let mut buf = [0; 5];
7
8     while file.read_exact(&mut buf).is_ok() {
9         // process the buffer
10    }
11 }
```

IO Buffering

- Read as much as possible
- Read reads from internal buffer until empty
- When empty: perform another system call

```
1 use std::io::{BufReader, Read};
2
3 fn main() {
4     let mut file = File::open("very_large_file.txt").unwrap();
5     // only line changed
6     let mut reader = BufReader::new(file);
7
8     let mut buf = [0; 5];
9
10    while reader.read_exact(&mut buf).is_ok() {
11        // process the buffer
12    }
13 }
```

Note: also works for (network) sockets and other file-like objects

Lock contention

Lock before slow operation

```
1  let result =  
    → Arc::new(Mutex::new(0));  
2  
3  for i in 0..3 {  
4      let r = result.clone();  
5      spawn(move || {  
6          let mut guard = r.lock();  
7  
8          *guard += slow(i);  
9      });  
10 }
```

Lock after slow operation

```
1  let result =  
    → Arc::new(Mutex::new(0));  
2  
3  for i in 0..3 {  
4      let r = result.clone();  
5      spawn(move || {  
6          let answer = slow(i);  
7  
8          *r.lock() += answer;  
9      });  
10 }
```

- Try to make critical sections smaller
- Hint: use scope blocks to be explicit about critical sections

Moving code outside of a loop

```
1  fn example(a: usize, b: u64) {
2      for _ in 0..a {
3          some_other_function(b + 3)
4      }
5  }
```

Compute $b + 3$ first:

```
1  fn example(a: usize, b: u64) {
2      let c = b + 3;
3      for _ in 0..a {
4          some_other_function(c)
5      }
6  }
```

Compilers are good at this: <https://godbolt.org/z/n58GjhsP5>

Memory allocation

- Memory allocation and deallocation takes time. For example:
 - `Vec::new()`, `HashMap::new()`, `String::new()`
 - Resizing any of the above
 - `Box::new()`
 - Cloning any of the above
 - `Arc::new()`, `Rc::new()`
 - Dropping any of the above
- Time depends on allocator you use (this can be changed!)

Question:

What can we do to reduce allocation?

Memory allocation

- Static allocation
- Preallocating to the right size
- Moving allocations outside an inner loop
- Arena allocation:
 - Small special-purpose allocator
 - (usually) no individual freeing supported
 - free all at once

Memory Allocation

```
1  struct Doggo {
2      cuteness: u64,
3      scritches_required: bool,
4  }
5
6  // Create a new arena to allocate into.
7  let bump = Bump::new();
8
9  // Allocate values into the arena.
10 let scooter = bump.alloc(Doggo {
11     cuteness: u64::MAX_VALUE,
12     scritches_required: true,
13 });
14
15 // Mutable references to the just-allocated value are returned.
16 assert!(scooter.scritches_required);
17 scooter.cuteness += 1;
```

From: <https://docs.rs/bumpalo/latest/bumpalo/#example>

Memory Allocation

Question:

How is this different to a Vec?

Making code faster

Question:

Do you know any techniques already to create faster code?

Making code faster

Question:

Do you know any techniques already to create faster code?

Inlining

```
1 // compiler may
2 // ignore this
3 #[inline]
4 fn example_1() {
5     ...
6 }
7
8 // not this
9 #[inline(always)]
10 fn example_2() {
11     ...
12 }
13
14 fn main() {
15     example_1();
16     example_2();
17 }
```

- Calling functions works with call instructions
- This represents some overhead
- What if we pasted the body of a function at the callsite?
- Larger code size, (sometimes) faster code

Compiler Options

- The compiler has multiple optimization levels
- 0, 1, 2 and 3 for “speed”
- ‘s’ and ‘z’ for “size” (code size)
- The default is very low → fast(er) compile times
- To select better options: `cargo run --release`

Changing Compiler Options

Cargo.toml

```
1  [package]
2  ...
3
4  [dependencies]
5  ...
6
7  # for cargo run
8  [profile.dev]
9  opt-level = 1
10
11 # for cargo run --release
12 [profile.release]
13 opt-level = 3
```

- `opt-level` can be 0-3, 's', 'z'
- 3 not always the best: experiment!
- Optimization levels are a tradeoff, 3 may be fast but the code size could be huge
- <https://godbolt.org/z/nTnef7888>

Link-Time optimization

Cargo.toml

```
1  [package]
2  ...
3
4  [dependencies]
5  ...
6
7  # for cargo run
8  [profile.release]
9  opt-level = 3
10 lto = true
```

- Each crate compiled separately
- No optimization between crates
- With lto there is

Target CPU

- ① Modern CPUs sometimes have specialized hardware
- ② Not every CPU has the same hardware
- ③ Programs cannot assume they can use this hardware
- ④ With `target-cpu` you choose a specific cpu

Compile time:

```
RUSTFLAGS="-C target-cpu=native" cargo run
```

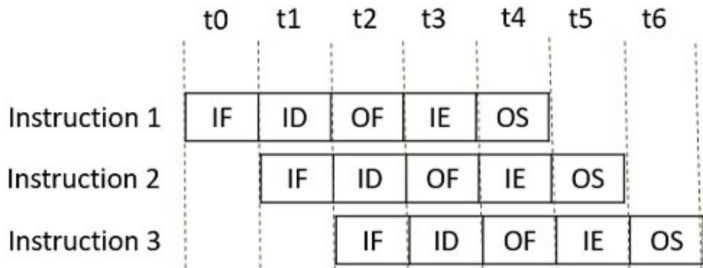
Target CPU (runtime)

```
1  #[inline(always)]
2  fn foo_impl() { ... }
3
4  // This generates a stub for CPUs that support SSE4:
5  #[target_feature(enable = "sse4")]
6  unsafe fn foo_sse4() {
7      // inlining here will recompile
8      // foo_impl for sse4
9      foo_impl()
10 }
11
12 // This generates a stub for CPUs that support AVX:
13 #[target_feature(enable = "avx")]
14 unsafe fn foo_avx() {
15     foo_impl()
16 }
```

Branch prediction

```
1  if a > b {  
2      do_x();  
3  } else {  
4      do_y();  
5  }
```

- Condition only available late due to pipelining
- Predict the outcome of the condition
- Start executing most-likely branch
- Wrong prediction → performance penalty



Cold paths

```
1  #[cold]
2  fn rarely_executed_function() { }
```

- Mark rarely used functions
- Generated code will favor optimizing other paths

Better branch prediction?

Question:

How do we know what to inline and what to mark as cold?

PGO: better branch prediction

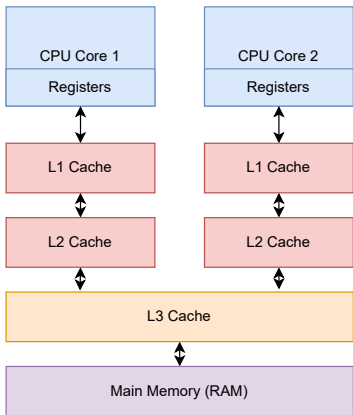
- Manually figuring out what to inline / mark as cold is hard
 - Profiling can help! → “Profile-Guided Optimization”
 - See lecture notes how to do this²
- ① Build special ‘Instrumented’ binary
 - ② Gather statistics while running
 - ③ Build a better program

Note:

If the actual workload is substantially different from the instrumented run then the program could perform worse!

²<https://cese.pages.eui.tudelft.nl/software-systems/part-1/lecture-notes/lecture-2.html>

Caching



- Closer caches are smaller and faster
- Smaller code and smaller data may mean more of it fits in cache
- Consecutive data is usually better
- Optimized code may be larger
- Benchmark! (cachegrind)

Rust: Zero cost abstractions

Rust provides abstractions:

- Iterators
- Traits & Generics
- Built-in collections
- Closures (lambda functions)

Question:

Does this mean Rust is slower than C?

Rust: Zero cost abstractions

Question:

Does this mean Rust is slower than C?

- Not necessarily!
- Abstractions have no cost if you don't use them (unlike Python)
- If you do use them, they are close to what you could manually make
- The compiler is very smart!

Rust: Static dispatch

```
1  struct A; struct B;
2  trait X {
3      fn something(&self);
4  }
5  impl X for A {
6      fn something(&self) { ... }
7  }
8  impl X for B {
9      fn something(&self) { ... }
10 }
11
12 fn run_something<T: X>(x: T) {
13     x.something();
14 }
15
16 run_something(A);
17 run_something(B);
```

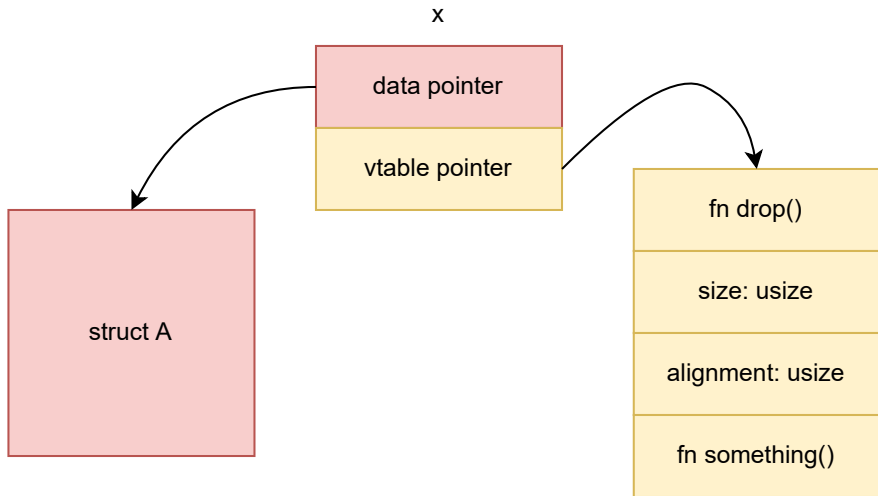
- Static Dispatch
- Covered in Software Fundamentals (Lecture 5)
- Code is duplicated for different types

Rust: Dynamic dispatch

```
1  struct A; struct B;
2  trait X {
3      fn something(&self);
4  }
5  impl X for A {
6      fn something(&self) { ... }
7  }
8  impl X for B {
9      fn something(&self) { ... }
10 }
11
12 fn run_something(x: &dyn X) {
13     x.something();
14 }
15
16 run_something(&A);
17 run_something(&B);
```

- Dynamic Dispatch
- Only one version of code
- Vtables!

Rust: Dynamic dispatch



Vtable: A description of a type

Rust: Struct layout

```
1  struct A {  
2      a: u8;  
3      b: u32;  
4  }  
5  
6  #[repr(packed)]  
7  struct A {  
8      a: u8;  
9      b: u32;  
10 }  
11  
12 #[repr(C, align(2))]  
13 struct A {  
14     a: u8;  
15     b: u32;  
16 }
```

- By default, structs have padding
- Packed makes structs smaller
- Alignment can make accessing fields slower

Rust: Memory management

```
1  bitflags! {  
2      struct Flags: u8 {  
3          const A = 0b00000001;  
4          const B = 0b00000010;  
5          const C = 0b00000100;  
6          const D = 0b00001000;  
7          const E = 0b00010000;  
8          const F = 0b00100000;  
9          const G = 0b01000000;  
10         const H = 0b10000000;  
11     }  
12 }
```

- A boolean is a byte
- Sometimes, packing tighter helps performance too

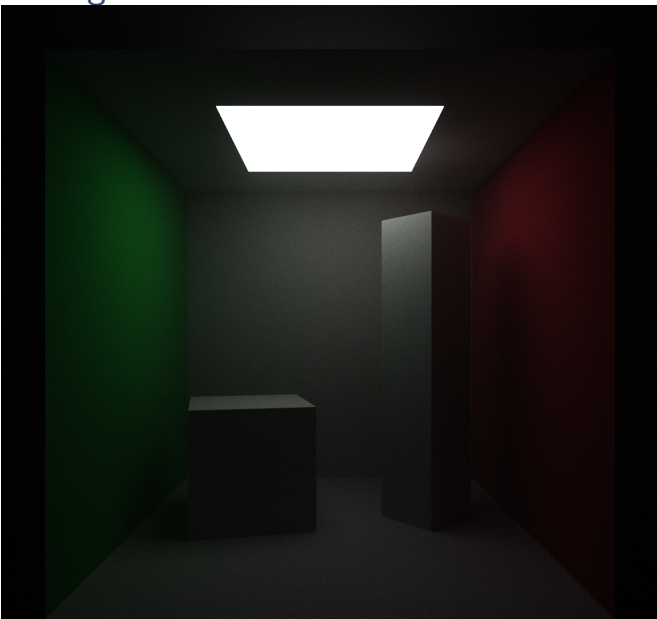
Rust: Zero copy design

- Cloning takes time
- If we could keep our data in one place:
 - An arena
 - A buffer (for example: packets from network)
- Don't pass around data, pass around offsets and lengths: slices
- Serde: Serializing and Deserializing (<https://docs.rs/serde>)

Much more!

- No silver bullet: Everything is a tradeoff
- Benchmark & Profile to know for sure!
- Lots of other ways:
 - Algorithmic improvements
 - Better/different hash functions
 - Hardware specific optimizations
 - Using dedicated hardware (ASIC, GPU, etc.)
 - Probabilistic datastructures
 - Operating system configuration

Assignment 2: Performance



- Raytracing
- Deliberately slow
- Apply techniques from lecture