# Asychronous Rust

## Performant Software Systems with Rust — Lecture 13

Baochun Li, Professor

Department of Electrical and Computer Engineering

University of Toronto

# Entering **Asynchronous Rust**

# Asychronous Programming

- Multi-threaded programming → **asynchronous** programming

  - operations may not finish sequentially in the order they were started

- CPU-bound vs. I/O-bound operations: video export vs. downloading from the Internet

# Blocking System Calls

- Many operating system APIs — in the form of **system calls** are **blocking**

- We can certainly spawn new threads to use these blocking system calls

  - but watch out on the **overhead** of spawning threads

  - it is not a good idea to have too many threads

# Non-Blocking Calls

- Compared to a `spawn` followed by a `join`, it would be great to have something like this:

```
1  let data = fetch_data_from(url).await;
2
3  println!("{data}");
```

# Threads → Stackless Coroutines

- Migration to **stackless coroutines**

- Concurrency happens entirely (or mostly) within your code

- An **async runtime** — which is just another crate — manages async tasks

- Your code yields control from time to time, say using the `await` keyword

- Rust supports **async** programming with `Future`s and the `async` / `await` syntax

# Revisiting the Idea of Parallelism vs. Concurrency

- **Analogy**: a team splitting up work for the course project
  - assign each member multiple tasks, assign each member one task, or a mix of both
- **Concurrency**: assigning each member multiple tasks, can switch but only make progress one at a time — a **single CPU core**
- **Parallelism**: each member takes one task, making progress at the same time — **multiple CPU cores**

# Futures

- A **future**: is a value that may not be ready now, but will become ready at some point in the future

  - called a **task** or a **promise** in other programming languages

- Rust provides a `Future` trait as a building block from the Rust standard library (`std::future::Future`)

- Futures are simply types that implement the `Future` trait

- Each future holds its own information about the progress that has been made and what "ready" means

# The Async Syntax

- You can apply the `async` keyword to blocks and functions to specify that they can be **interrupted** and **resumed**

- Within an `async` block or `async` function, you can use the `await` keywoard to **await** an **future** — waiting for it to become "ready"

- Anywhere you `await` a future, that `async` block or function can **pause** and **resume**

- The process of checking with a future to see if its value is available yet is called **polling**

# Hello, world! in Async Rust

```rust
1   // Define an async function.
2   async fn say_hello() {
3       println!("Hello, world!");
4   }
5
6   // Boilerplate which lets us write `async fn main`, needs
7   // `tokio = { version = "1", features = ["full"] }` in Cargo.toml
8   #[tokio::main]
9   async fn main() {
10      // Call an async function and await its result.
11      say_hello().await;
12  }
```

# The Async Runtime

- The async runtime schedules tasks to be run on CPU cores

- When one task suspends (say by calling `recv()` or `await`), another must be picked to run

- When this task resumes (say the message it is waiting for arrives), it should be ready to run

# The Async Runtime

- Rust doesn't provide such a runtime itself — leaving this to crates

- Much better choice than other languages, where the runtime may also manage memory (such as garbage collection in Go), or even becomes a full virtual machine

- A simple runtime (also called **executor**) can be single-threaded, can be written in less than 500 lines of code

# Choice of Async Runtimes

- `tokio` is a full-featured, multi-threaded runtime that we often use

  - But other runtimes are also excellent and can have less overhead

  - such as `async-std` and `smol`

  - Honourable mention: `bytedance/monoio`

    - a more efficient thread-per-core runtime without work-stealing task queues, so that local data does not need to be `Sync` and `Send`

# async and await

```
1  // An async function, but it doesn't need to wait for anything.
2  async fn add(a: u32, b: u32) → u32 {
3    a + b
4  }
5
6  async fn wait_to_add(a: u32, b: u32) → u32 {
7    sleep(1000).await;
8    a + b
9  }
```

Calling an async function directly — without await —
returns a **future**, which is a struct or enum that implements
the Future trait and represents deferred computation

# A Real-World Example —
`page_title()`

# A Real-World Example — main()

```rust
1 #[tokio::main]
2 async fn main() {
3     let url = "https://www.rust-lang.org";
4
5     match page_title(url).await {
6         Some(title) ⇒ println!("Title of {} is: {}", url, title),
7         None ⇒ eprintln!("Could not fetch/parse title for {}", url),
8     }
9 }
```

# Desugared async function

```
 1  fn page_title(url: &str) → impl Future<Output = Option<String>> + '_ {
 2      async move {
 3          let response = reqwest::get(url).await.ok()?;
 4          let response_text = response.text().await.ok()?;
 5          let document = Html::parse_document(&response_text);
 6          let selector = Selector::parse("title").ok()?;
 7
 8          document
 9              .select(&selector)
10              .next()
11              .map(|title| title.inner_html())
12      }
13  }
```

# Desugaring `async` function

- Uses the `impl Trait` syntax — Traits as parameters
- Returns a `Future` with an associated type of `Output`
  - type is `Option<String>`
  - the same as the original return type from the `async fn` version of `page_title`

# Desugaring async function

# What Is + `'_ Doing Here?

```rust
1 fn page_title(url: &str) → impl Future<Output = Option<String>> + '_ {
2     async move {
3         // captures `url` inside this closure
4     }
5 }
```

# What Is `+ `_` Doing Here?

- The returned future captures `url`, which is a `&str`
  - The future borrows from `url`
  - So the future cannot `outlive` the reference it holds
- Therefore its lifetime is tied to the lifetime of `url`
- `+ `_` implies "This `impl Future` may contain borrows, and its lifetime is some anonymous lifetime that's at most as long as the references it captures (like `url`)."

**Good news**: Only needed in older versions of the Rust compiler (before v1.7), no longer needed now.

# Desugaring `#[tokio::main]`

```rust
1  #[tokio::main]
2  async fn main() {
3      let url = "https://www.rust-lang.org";
4
5      match page_title(url).await {
6          Some(title) ⇒ println!("Title of {} is: {}", url, title),
7          None ⇒ eprintln!("Could not fetch/parse title for {}", url),
8      }
9  }
```

# Desugaring `#[tokio::main]`

```rust
 1 fn main() {
 2     // Create a multi-threaded Tokio runtime
 3     let rt = Runtime::new().expect("failed to create Tokio runtime");
 4     let url = "https://www.rust-lang.org";
 5
 6     rt.block_on(async {
 7         match page_title(url).await {
 8             Some(title) ⇒ println!("Title of {} is: {}", url, title),
 9             None ⇒ eprintln!("Could not fetch or parse title for {}", url)
10         }
11     });
12 }
```

Performant Software Systems with Rust

# Deep dive: What's a Future in Rust anyway?

# A value that represents a **<span style="color:yellow">computation</span>** that may not be finished yet.

# Future: **Formal Definition**

```rust
1 pub trait Future {
2     type Output;
3
4     fn poll(
5         self: Pin<&mut Self>,
6         cx: &mut Context<'_>,
7     ) → Poll<Self::Output>;
8 }
```

# Key Points

- type `Output`: the type the future will eventually produce (like `u32`, `Result<T, E>`, etc.)

- `poll(...)`: asks the future "are you ready yet?"

  - Returns `Poll::Pending` if it's not done

  - Returns `Poll::Ready(value)` when it has the final result

- You almost never write `poll` by hand in normal code — the `async/await` syntax and executors do that for you

# Futures are Lazy

- Just like **iterators** and unlike **threads**

  - A `std::thread::spawn` starts running immediately on another OS thread

  - A `Future` does nothing until an executor polls it

```
1  let fut = async { println!("hi"); }; // Nothing printed yet
2  // Only when we .await it (or poll it) will it actually run.
```

# Executors and Tasks: The Runtime

- A runtime (like `Tokio`, `async-std`, etc.) provides an executor that

    - Keeps a queue of futures (often called `tasks`)

    - Repeatedly calls `poll` on them

- Uses a `Waker` to be notified when I/O or timers are ready, so it can poll again

# What Does `self: Pin<&mut Self>` Mean?

```rust
1 pub trait Future {
2     type Output;
3
4     fn poll(
5         self: Pin<&mut Self>,
6         cx: &mut Context<'_>,
7     ) -> Poll<Self::Output>;
8 }
```

# What Does `self: Pin<&mut Self>` Mean?

```
1  future.poll(cx);
```

- `poll` takes `self` by `Pin<&mut Self>` instead of by `&mut Self`

- `Pin<T>` is a wrapper that says "Through this pointer, you are not allowed to move the value `T` in memory."

# Pin<&mut T>

- Logically a mutable reference, but with an extra guarantee:

  - you can mutate T

  - but you must not move T to a different memory location via this reference

# Why Do We Care About Moving?

# **What Does** `self: Pin<&mut Self>` **Enforce?**

- Before calling `poll`, the executor must have pinned the future

- Inside `poll`, the future is treated as **immovable** (through this handle)

# What Does an Executor Do?

```rust
1  use std::pin::Pin;
2  use std::future::Future;
3  use std::task::Context;
4
5  fn drive<F: Future + Unpin>(fut: &mut F, cx: &mut Context<'_>) {
6      // Pin the &mut F temporarily
7      let pinned = Pin::new(fut);
8      match Future::poll(pinned, cx) {
9          Poll::Pending => { /* ... */ }
10         Poll::Ready(output) => { /* ... */ }
11     }
12 }
```

# The Unpin Marker Trait

If the future is not Unpin (e.g., most async fn futures), the executor must pin it in a stable place first:

```
1  let fut = my_async_fn();
2  let mut fut = Box::pin(fut); // heap-allocate and pin
3
4  // later:
5  let waker = /* ... */;
6  let mut cx = Context::from_waker(&waker);
7  let _ = fut.as_mut().poll(&mut cx);
```

# The Unpin Marker Trait

```
1  pub trait Unpin {}
```

"Is it still safe to move this thing around in memory, even if it has been pinned?"

- **Yes**: the type is Unpin

- **No**: the type is not Unpin (written !Unpin)

# The Unpin Marker Trait

- If `Self: Unpin`, then `Pin<&mut Self>` is the same as `&mut Self` in practice

- Most "normal" types are `Unpin`

- Compiler-generated `async fn` futures are usually not `Unpin`, so they need pinning

- So `self: Pin<&mut Self>` is a general signature:

  - Works for both `Unpin` and `!Unpin` futures

  - Enforces "must be pinned before polling" at the type level

# What's a `Waker`?

- A `Waker` is basically the doorbell for a future

- A `Waker` is a handle created by the executor that a future can store and later call `wake()` on, to say:

- **"Hey runtime, I might be ready now — please poll me again soon."**

# Where Does Waker Live?

```rust
1 pub trait Future {
2     type Output;
3
4     fn poll(
5         self: Pin<&mut Self>,
6         cx: &mut Context<'_>,
7     ) -> Poll<Self::Output>;
8 }
```

The Context carries a Waker:

```rust
1 impl<'a> Context<'a> {
2     pub fn waker(&self) -> &Waker { ... }
3 }
```

# poll()

```rust
1 fn poll(self: Pin<&mut Self>, cx: &mut Context<'_>) → Poll<Self::Output> {
2     let waker = cx.waker();
3     // ... maybe store it somewhere ...
4 }
```

# **Why Do We Need a `Waker`?**

- Without `Waker`, if a future isn't ready yet, the executor would have to:

  - keep polling it in a loop (busy-waiting), or

  - poll every future all the time "just in case"

- That's terrible for CPU!

# **Why Do We Need a `Waker`?**

- The contract in async Rust is:

  - When `poll` returns `Poll::Pending`, the future is not ready yet

  - If it ever becomes ready to make progress, it must ensure that someone calls `waker.wake()`

  - The executor, when `wake()` is called, will schedule that future to be polled again

# Why Do We Need a Waker?

- So Waker is how a future tells the executor:

  - **"Some external event happened (I/O ready, timer fired, etc.). Please come back and poll me."**

# Spawning Tasks

```rust
1  use tokio::{spawn, time::{sleep, Duration}};
2
3  async fn say_hello() {
4      // Wait for a while before printing
5      sleep(Duration::from_millis(100)).await;
6      println!("hello");
7  }
8
9  async fn say_world() {
10     sleep(Duration::from_millis(100)).await;
11     println!("world!");
12 }
13
14 #[tokio::main]
15 async fn main() {
16     spawn(say_hello());
17     spawn(say_world());
18
```

# Joining Tasks

```rust
1  #[tokio::main]
2  async fn main() {
3      let handle1 = spawn(say_hello());
4      let handle2 = spawn(say_world());
5
6      // Wait for both tasks to finish
7      // spawn returns `JoinHandle`, which implements `Future`
8      let _ = handle1.await;
9      let _ = handle2.await;
10 }
```

# Required Additional Reading

The Rust Programming Language, Chapter 17

The Rustonomicon, Chapter 8

Asynchronous Programming in Rust, Chapter 1-4

The State of Async Rust: Runtimes

Tokio: An Asynchronous Runtime