



THANK YOU to our sponsors	
Bronze	
🙏 Undo	
for supporting the NYC++ Me	etup
Adobe 3	© 2024 Adobe. All Rights Reserved





Concurrency is difficult because of mutation and efficiency.



Adobe Lightroom is a professional image editing tool targeting photographers.



Futures with continuations are a reasonable model for transforming synchronous code to asynchronous code and async code to non-blocking concurrent code. They don't handle all concurrency patterns - they are just part of the solution, but an important part we are focused on today.

My first thought was to use boost futures with continuations. I wrote a test piece of code equivalent to this. The second line crashed - with boost futures, `then()` is a consuming operation. If it had relied on C++11, it would have required a move. You could write a "split" operation (although, at the time, I didn't see how), but there were more issues that led to me starting what became the stlab concurrency library.



Let's back up a bit and look at one of the reasons concurrency is used in an interactive application—emphasis on interactive. Concurrency is a tool to improve performance with a fundamental tradeoff between latency and throughput. An interactive application strives for specific (real-time) latency guarantees. You want the application to feel fluid and responsive, regardless of how long the processing takes.



These are slides from an old product I worked on, Adobe Revel, this was built on Adobe Camera Raw, the same rendering engine that powers Adobe Lightroom and the Camera Raw component in Photoshop. The code is optimized for latency over throughput.



To provide smooth pan and zoom without seeing black the code keeps a GPU texture of the entire image at a low resolution.



With a high resolution image overlaying it to fill the current view port.





When you make an adjustment, like sliding a slider...





We render a draft of the foreground texture at a low quality.





Then, if there are no pending render requests, we render at a higher quality...



And then again at full quality.



And then we update the background layer.







When panning or zooming we animate the two layers in the viewport at full device frame rate (120Hz on an iPad Pro).





But we chase the animation, starting a render for a higher quality image at the, current, position we are animating too.



With the goal that by the time the motion stops, you are looking at full quality pixels. The human visual system is poor at discerning detail of objects in motion or objects that are chaning. The effect is it appears to be rendering high quality at full frame rate, even though it may only be capable of rendering full quality at a frame-rate of 1/4 or less.



The state machine that drives the rendering for display is implemented as a coroutine. The operations at each state are asynchronous. As the state machine transitions, it may _cancel_ the operations from prior states as they are no longer required, typically canceling a "stale" high-quality render request in favor of a lower-quality, faster request. But lower-quality renders may be a delta-render of a high-quality render, or there may be caches shared between renders that may be built. So when canceling, the question becomes, cancel what exactly?



Futures are a building block for asynchronous code. They provide a simple way to transform synchronous code into asynchronous code. Conceptually, a future is a token for the result of a function executed in a different context. Futures carry the dependencies—destructing the future cancels the task.



Continuations are just function composition. Again, dependencies are carried by the futures. If you don't want a result, then destructing it "unravels" the dependencies.



With joins (when all) and splits (multiple continuations on the same future), the graph becomes a DAG. Canceling an operation may cancel may in turn cancel other operations, but not those that are required to satisfy other futures.



We can define efficient cancelation. Cancelation is also asynchronous - the canceling thread should not await active operations. Of corse the stop message propagation takes time - and is inherently racy. We want to minimize the message propagation time and avoid serializing cancelation between active tasks.



Let's back up a bit and look at one of the reasons concurrency is used in an interactive application—emphasis on interactive. Concurrency is a tool to improve performance with a fundamental tradeoff between latency and throughput. An interactive application strives for specific (real-time) latency guarantees. You want the application to feel fluid and responsive, regardless of how long the processing takes.



Waiting for a task to complete... such as on a condition variable or a semaphore



The first reason is Amdahl's. This graph shows the effect of serialization in 10% increments up to a 16 core machine. If only 10% of our code is serialized, our performance with 16 cores is only slightly more than 6x and will never exceed 10x no matter how many processors we throw at the problem. Every synchronization point is expensive.



Even in a multithreaded environment, sync-waits can lead to deadlocks. Here is a visualization of a deadlock on a with a single thread of execution.

A task issues another task which is placed in the queue. Then stops to wait for the task to complete.

There is only one thread and one queue so the queued task will never complete, this is a deadlock.



With two threads, the first task may spawn a task, and wait for it to complete. That task is picked up by the second thread it spawns some additional work. And waits for that work to be completed. Deadlocked.

Any sync wait operation requires non-local reasoning about where the task that is being waited on is scheduled, what tasks that task may spawn, and how many threads are available in the system.

This leads to the "function color problem" - If a function is async, every operation that depends on it is also async. There is no way to rejoin.

To use sync wait safely, you need to know that thread resources are available to execute the item being waited on. Sync wait breaks local reasoning.

History

Adobe

- I gave a talk on how I solved these problems, Better Code: Concurrency
- I cleaned up my code and packaged it as a little library on GitHub, stlab/libraries/concurrency

36

© 2024 Adobe. All Rights Reserved

- Written almost entirely in the hotel bar at C++Now
- Felix Petriconi found it useful and contributed most of the current code
- The library is now in broad use at Adobe and other companies



stlab futures are regular, they can be copied, you can attach multiple continuations, destruction is cancelation. RAII applies to the processor (arguably even more important a resource than memory!).

There is more to the library than just futures with continuations. There is a high performance dynamic task stealing executor (that back-ends to the platform executor, i.e. Apple's libdispatch, or Windows thread pool, if available), a "main" scheduler to schedule on the main event loop (if the platform has one), an implementation of channels (needs a major update), and a number of smaller utilities...



The programming model presented by futures with continuations is deceptively simple -But every continuation comes at a cost.



A small object allocation/deallocation is 200-500 cycles. Every atomic operation is 15-30 cycles.

The overhead of an async continuation vs sequential function composition can easily be 100-1000x.

We only want to attach continuations when the operation is significantly large or a specific execution context is required.

(Side note: The last line on this graph is a thread context switch. This is why thread pools exist and the importance of not overcommitting.)



Futures and continuations work well when the grain size (how much work we have to do) in the continuation is relatively large. The more we reduce the overhead the greater the performance benefit - not just because we reduce the cost per call but because we enable additional parallelism to benefit.



But this doesn't compose well, an async operation now needs to take a callback so it can compose the operation prior to starting it. The simple syntax of futures with continuations means they get used for simple function composition - at a high (very measurable) cost. Even this function composition has a cost - at each level we are adding another stack from. But what if instead of starting the async operations, we simply described them - enter sender/receivers.



An async operation needs to take a callback so it can compose the operation prior to starting it. The complexity grows rapidly.

The simple syntax of futures with continuations means they get used for simple function composition - at a high (very measurable) cost. Even this function composition has a cost - at each level we are adding another stack from. But what if instead of starting the async operations, we simply described them - enter sender/receivers. The complexity grows rapidly.



This is the problem that C++ sender/receivers solve. Instead of launching async code and attaching continuation, we build a _program_ (a function) that describes the async operation and then start that. This _program_ is built flat, without heap allocations, with no need for synchronization because the connections are made before it starts.

`then()` is compose logically compose, it doesn't return g(f(x)) it returns the logical function g composed with f. This composition pattern is obscured with a complex sender/receiver interface. We're building up functions that describe an async computation graph, and then we start it. By appending application and composition operations, we can build a computation tree.

In sender/receiver terminology, `then()` returns a sender adapter closure.

In std::exec, "transfer" transfers execution to, for example, a thread, serial queue, GPU, etc.



The programming model presented by futures with continuations is deceptively simple -But every continuation comes at a cost.



Let's start with flat composition - we will drop the `then`.

We're building a functional model to describe the execution of a program



Links are just functions - stored as a flat sequence. The links form a single logical function.



With sender receiver we can pass the error from one `then()` to the next, but what do we do if this is just a simple compose? If we were using std::expected, it doesn't change this, instead it changes how the application happens.

We only need a sequence of functions to describe the computation



The programming model presented by futures with continuations is deceptively simple -But every continuation comes at a cost.



This is binding s to f and g - but we don't actually do the bind directly, we store the structure



Segments are an applicator and a sequence of links. A segment logically forms a single function.

```
Application - Segments
auto expector = [](auto f, auto... args) -> expected<...> {
    try { f(args...); }
    catch { return current_exception(); }
};
auto badd = [](int x, int y){
    if (x == 42) throw runtime_error("bad x");
    return x + y;
}
(on(expector) | badd | [](int x){ return x * 2; })(12, 5);
expected{34};
```

As an example, we can write an apply function that maps exception to std::expected and then use it to compose functions and return expected on an error.

If we pass 42 for x will result in an expected carrying the bad x exception.



The programming model presented by futures with continuations is deceptively simple -But every continuation comes at a cost.



If s sends f | g to another thread, how do you get the result back? Are we back to a future?



Chains are two-dimensional data structures. We can append functions to the end of the last segment or segments to the end of the chain. They are stored as a tuple of tuples of functions - flat with no heap allocations.



The chain is executed by appending each segment to the end of the prior segment, so the result of one segment is passed as an argument to the next segment. A chain is logically a single function.



The programming model presented by futures with continuations is deceptively simple -But every continuation comes at a cost.



Instead of wrapping individual operations



p is a promise associated with a future. Having access to the structure allows operations to make transformations. Errors only need to be caught and propagated from the segment level (same for cancelation).



* How general purpose? I believe chains with an S & K combinator are Turing complete - still working on the proof.

** Kind of - they don't have to be used for asynchronous descriptions, but they carry the complexity of signaling cancelations and exceptions at every step.



The current implementation with sender/receiver split is broken with respect to cancelation. let_value() provides a form of split for computation, but split must be contained with a rejoin. I've managed to spark interest in trying to find a solution on the sender/receiver forum



* Maybe. I'm only two weeks into working on chains, and each step has been challenging. I may just rediscover the complexity of sender/receivers. A structured form for the complete "program" may have advantages for scheduling (especially on GPUs). But also some of the same problems - the entire structure of the program is carried in the type system. I suspect I'll hit implementation limits, however, my types are less complex then sender-receiver types. Eric is currently exploring a form of type erasure through lambda expressions that is questionable.

... And I still have many open questions. Declaring signatures up-front would simplify the code and produce better compiler diagnostics, but it would also limit "generic" code, which would have to build programs on instantiation. I'm 2 weeks in where sender/receivers are nearly 10 years in... At the very least the process of thinking through and building chains has greatly improved my understanding of Sender/Receivers and where the complexity comes from.





